



Rejecta Mathematica

Volume 1, Number 1 — July 2009

ISSN 1948-8351

math.rejecta.org

Rejecta Mathematica is an open access, online journal that publishes papers that have been rejected from peer-reviewed journals in the mathematical sciences. Each paper is accompanied by an open letter from its authors discussing the original review process and stating the case for its value to the research community.

Letter from the Editors	1
<hr/>	
Articles	
Subspaces that Minimize the Condition Number of a Matrix	4
<i>Siddharth Joshi and Stephen Boyd</i>	
Automatic Countings	10
<i>Doron Zeilberger</i>	
Alexander Duality for Monomial Ideals and their Resolutions	18
<i>Ezra Miller</i>	
Meaning-Imposers versus Meaning-Derivers	58
<i>Gary Harper</i>	
WInHD: Wavelet-based Inverse Halftoning via Deconvolution	84
<i>Ramesh Neelamani and Richard Baraniuk</i>	
Mass Matrix Transforms in Qubit Field Theory	104
<i>Marni Sheppeard</i>	

Michael Wakin — Christopher Rozell — Mark Davenport — Jason Laska
editors@rejecta.org

2009 Rejecta Publications, Inc.

This work is published under the Creative Commons Attribution-NonCommercial License.



License <http://creativecommons.org/licenses/by-nc/2.5/legalcode>

Human-Readable Summary <http://creativecommons.org/licenses/by-nc/2.5/>



Welcome to the inaugural issue of *Rejecta Mathematica*! Thank you for joining us for what we hope will be a unique and interesting experiment. For those unfamiliar with our mission, *Rejecta Mathematica* is an open access, online journal that publishes only papers that have been rejected from peer-reviewed journals in the mathematical sciences. In addition, every paper appearing in *Rejecta Mathematica* includes an open letter from its authors discussing the paper's original review process, disclosing any known flaws in the paper, and stating the case for the paper's value to the community.

Since starting this endeavor, the questions we've been asked most often are "Why are you doing this?" and "Is it a joke?" While we are not above admitting that we have had a few good laughs in this process, we hope that this issue will serve as definitive proof that *Rejecta Mathematica* is not a joke. Despite the central role that peer review (and even rejection) must play in the scientific process [1], we believe there are several reasons why this project can make a positive and valuable contribution to the mathematical sciences research community.

First, there is ample evidence that in the traditional review process, significant (even Nobel prize-winning) research is occasionally overlooked and groundbreaking work is sometimes actively shunned [2–4]. Perhaps this is most dramatically illustrated in the fact that at least "36 future Nobel Laureates encountered resistance on [the] part of scientific journal editors or referees to manuscripts that dealt with discoveries that on [a] later date would assure them the Nobel Prize" [5]. While it would be presumptuous for us to assume that we can spot significant work that others may have missed, we can provide a venue to introduce rejected work to the community and increase the chances that its value will be appreciated sooner rather than later.

Second, there is also evidence that a research community can derive value from a centralized repository of rejected papers, even when (and perhaps especially when) the results are either incorrect or not significant enough to warrant consideration for a major international prize. *Rejecta Mathematica* can benefit authors looking for feedback on their work, wanting to warn the community against false starts (i.e., the classic "null results" that never see the light of day, only to be repeated by others) [6, 7], or wanting to illuminate the occasional vagaries of the peer review process to enhance accountability and scientific integrity [8]. Our journal can also benefit readers who want access to "minor results" that may be useful but not publishable in isolation. Indeed, *Rejecta Mathematica* has existed in folklore for many years as a fictitious place to send papers that were never to see the light of day, and the concept of a formal repository for rejected papers hoping to be discovered and revived (called *Rejuvenatable Mathematics*) has also been proposed [9].

While such a project as *Rejecta Mathematica* would have been impracticable in the pre-internet age, the flood of resources available today begs another oft-posed question: "Why do we need a new journal? Isn't this what a preprint server (like the arXiv), a blog, or a personal website is for?" We believe that a central collection of articles that have been selected for their potential interest to the community will increase their visibility beyond what could be achieved through a general preprint server or personal website. We also believe that the commentary and advocacy by the authors will increase the value of the papers beyond what



would exist from the appearance of the paper alone. Finally, we believe that the availability of thoughtful technical discussion (via *Rejecta Mathematica* “correspondences” following up on previously published articles) has the potential to generate more valuable interaction than the immediate commentary generally available on a blog. There is no doubt, however, that blogs and online archives can also play a significant role in advocating for rejected papers.

Finally, we would be remiss not to mention that being researchers ourselves, at some level we simply wanted to conduct an experiment. What started as a fleeting idea around the lunch table (discussing one of our own rejected papers) turned into the type of inquiry that fuels even the most serious of studies: if we build *Rejecta Mathematica* and ask for papers, what will happen? Will we get any papers, and if so, will they all be the delusional output of mathematical cranks? (This has been a common conjecture.)

Other questions concern our editorial policies. Should we simply publish every article we receive, and if not, how should we evaluate the submissions? After careful consideration, we have settled on an editorial process that includes no technical peer review (hence our slogan “Caveat Emptor”). Rather, we will rely on the technical review provided by the journal from which the paper was originally rejected and focus instead on selecting papers based on their apparent potential interest to researchers in the mathematical sciences. Admittedly, and perhaps necessarily in a journal of this scope, the concept of “potential interest” encompasses a broad set of loosely defined criteria. Ultimately, we will try to choose papers that allow some opportunity for learning. For example, we do not see much value to the community in publishing papers that were rejected solely for their incomprehensibility.

The open letter plays a major part in our decision process, as we view its role in a *Rejecta Mathematica* article as being at least as important as the technical content of the research paper. The open letters are where the authors can both tell the history of the paper and convey the lessons learned from the rejection. Undoubtedly, many open letters will provide a frank commentary on the peer-review process. Some may even be controversial. At the very least, they should help others benefit from the (technical and nontechnical) mistakes of their peers. To address the original question, there have indeed been papers rejected from *Rejecta Mathematica*.

We are delighted to say that the content of this first issue runs the gamut of genres included in our mission: minor or traditionally unpublishable results, non-traditional ideas and proof techniques, misunderstood genius, results based on questionable assumptions, and controversial papers and open letters. We are also pleased that the papers span several areas of the mathematical sciences, including pure mathematics, applied mathematics, theoretical physics, and engineering. We hope that you enjoy the issue with as much good humor and intellectual stimulation as we have encountered in putting it together. We welcome feedback, future submissions, and support for the *Rejecta Mathematica* mission through our website: math.rejecta.org.

Michael Wakin — Christopher Rozell — Mark Davenport — Jason Laska
editors@rejecta.org



References

- [1] F. C. Fang, “On rejection,” *Infection and Immunity*, vol. 76, no. 5, pp. 1802–1803, 2008. [Online]. Available: <http://iai.asm.org>
- [2] “Coping with peer rejection,” *Nature*, vol. 425, no. 6959, p. 645, 2003.
- [3] B. Barber, “Resistance by scientists to scientific discovery: This source of resistance has yet to be given the scrutiny accorded religious and ideological sources,” *Science*, vol. 134, no. 3479, pp. 596–602, 1961.
- [4] J. M. Campanario, “Have referees rejected some of the most-cited articles of all times?” *Journal of the American Society for Information Science*, vol. 47, no. 4, pp. 302–310, 1996.
- [5] ——. [Online]. Available: <http://www2.uah.es/jmc/nobel/nobel.html>
- [6] J. Stallings, “How not to prove the Poincare conjecture,” *Topology Seminar Wisconsin, 1965, Annals of Mathematics Studies*, vol. 60, pp. 83–88, 1966. [Online]. Available: <http://math.berkeley.edu/~stall/notPC.pdf>
- [7] *Journal of Negative Results in BioMedicine*. [Online]. Available: <http://www.jnrbm.com/>
- [8] *Philica*. [Online]. Available: <http://philica.com/>
- [9] A. Magid, “Theorems that should never have been proven,” *Notices of the AMS*, vol. 44, no. 7, 1997.



An open letter concerning Subspaces that Minimize the Condition Number of a Matrix

Siddharth Joshi

Stephen Boyd

This article poses and answers the following question: How do you choose a subspace of given dimension that minimizes the condition number of a given matrix on that subspace? Part of the answer is a bit surprising (at least to us): When the subspace dimension is no more than half the size of the matrix, a subspace can be found on which the matrix has condition number one.

We think our paper makes it clear that we consider our result simple, but interesting and not obvious. We certainly make no claims as to its depth, or its potential applications. It is not in the literature, and does not follow in any direct or simple way from existing results. In other words, it is, as far as we know, new.

The manuscript was rejected by two journals. The first rejection was based on the reviewers and editor noting that someone had written a paper that seemed to cover similar material. But a cursory reading of that paper, and ours, show that while the other paper shared a few key words with ours, the results were in no way related. On the positive side, one reviewer suggested a simplification of our proof, which we gladly used in our revision, which was also rejected.

We then submitted the article to another journal. In this case, the editor apparently did not even understand the result, which is stated very clearly, in completely standard, and elementary, mathematical language. Moreover, he insisted that we describe an application, so we added a simple application involving an ellipsoid intersected with a subspace. It was rejected.

Subspaces that Minimize the Condition Number of a Matrix

Siddharth Joshi

Stephen Boyd*

Abstract

We define the condition number of a nonsingular matrix on a subspace, and consider the problem of finding a subspace of given dimension that minimizes the condition number of a given matrix. We give a general solution to this problem, and show in particular that when the given dimension is less than half the dimension of the matrix, a subspace can be found on which the condition number of the matrix is one.

1 The problem

Suppose $A \in \mathbf{R}^{n \times n}$ and $\mathcal{V} \subseteq \mathbf{R}^n$ is a subspace with $\dim \mathcal{V} = k \geq 1$. We define the *maximum gain* (*minimum gain*) of A on \mathcal{V} , as

$$G_{\max} = \sup_{x \in \mathcal{V}, x \neq 0} \frac{\|Ax\|}{\|x\|}, \quad G_{\min} = \inf_{x \in \mathcal{V}, x \neq 0} \frac{\|Ax\|}{\|x\|},$$

respectively, where $\|\cdot\|$ denotes the Euclidean norm. When A is nonsingular, we define its *condition number on the subspace* \mathcal{V} as

$$\kappa_{\mathcal{V}}(A) = G_{\max}/G_{\min}.$$

The condition number of A on any one-dimensional subspace is 1, and its condition number on $\mathcal{V} = \mathbf{R}^n$ is the (usual) condition number of A , which we denote $\kappa(A)$. The condition number on any subspace is between 1 and $\kappa(A)$. If $\kappa_{\mathcal{V}}(A) = 1$, we say that A is isotropic on \mathcal{V} , since its gain $\|Ax\|/\|x\|$ is the same for any nonzero vector $x \in \mathcal{V}$.

In this note we address the following problem: Given a nonsingular matrix $A \in \mathbf{R}^{n \times n}$, and $k \in \{1, \dots, n\}$, find a subspace $\mathcal{V} \subseteq \mathbf{R}^n$ of dimension k which minimizes $\kappa_{\mathcal{V}}(A)$. The number $\kappa_{\mathcal{V}}(A)$ is a measure of the anisotropy of the linear function induced by A , restricted to the subspace \mathcal{V} , so our problem is to find a subspace of dimension k on which A is maximally isotropic.

*The authors are with the department of Electrical Engineering at Stanford University. Email addresses: Siddharth Joshi: sidj@stanford.edu, Stephen Boyd: boyd@stanford.edu.

We will show that the minimum possible condition number of A , on a subspace of dimension k , is given by

$$\inf_{\mathcal{V} : \dim \mathcal{V} = k} \kappa_{\mathcal{V}}(A) = \max \left(\frac{\sigma_{n-k+1}}{\sigma_k}, 1 \right) = \begin{cases} 1 & k \leq \lceil n/2 \rceil, \\ \sigma_{n-k+1}/\sigma_k & k > \lceil n/2 \rceil, \end{cases} \quad (1)$$

where $\sigma_1 \geq \dots \geq \sigma_n > 0$ are the singular values of A . (The infimum is over all subspaces of \mathbf{R}^n of dimension k .) This means, in particular, that for $k \leq \lceil n/2 \rceil$, we can find a subspace of dimension k on which A is isotropic.

There are many classical results that identify a subspace of a given dimension that minimizes or maximizes some quantity that depends on the subspace and matrix. For example, the Courant-Fischer theorem tells us that the minimum value of G_{\max} , over all subspaces of dimension k , is σ_{n-k+1} , and the maximum value of G_{\min} , over all subspaces of dimension k , is σ_k . For these and similar results, see, e.g., [3, §4.2] or [1]. Also, the idea of condition number of a matrix restricted to a particular subspace can be seen in [2].

We can give a geometric application (or interpretation) of our problem. We are given an ellipsoid $\mathcal{E} = \{z \mid \|Az\| \leq 1\}$ in \mathbf{R}^n , where $A \in \mathbf{R}^{n \times n}$ is nonsingular. Our goal is to find a k dimensional subspace \mathcal{V} so that the ellipsoid $\mathcal{V} \cap \mathcal{E}$ is as spherical as possible, *i.e.*, has minimum eccentricity. (The eccentricity of $\mathcal{V} \cap \mathcal{E}$ is defined as the ratio of its maximum semi-axis length to its minimum semi-axis length, which is exactly $\kappa_{\mathcal{V}}(A)$.) The solution is to choose \mathcal{V} that minimizes the condition number of A on \mathcal{V} . Our result (1) can be interpreted in this geometric setting. For example, if $k < \lceil n/2 \rceil$, we can always find a subspace of dimension k for which $\mathcal{V} \cap \mathcal{E}$ is perfectly spherical, *i.e.*, a ball. As a very simple special case, we see that for any ellipsoid in \mathbf{R}^3 , there is a plane that intersects it in a ball. Our general result (1) can be considered a generalization of this simple fact.

2 The solution

Suppose Q and Z are $n \times n$ orthogonal matrices, *i.e.*, $Q^T Q = Z^T Z = I$. Then we have

$$\kappa_{\mathcal{V}}(QAZ) = \kappa_{\mathcal{W}}(A),$$

where $\mathcal{W} = Z\mathcal{V} = \{Zv \mid v \in \mathcal{V}\}$. It follows that

$$\inf_{\mathcal{V} : \dim \mathcal{V} = k} \kappa_{\mathcal{V}}(A) = \inf_{\mathcal{V} : \dim \mathcal{V} = k} \kappa_{\mathcal{V}}(QAZ),$$

since the first orthogonal matrix Q has no effect, and the second orthogonal matrix Z simply changes the parametrization of subspaces of dimension k .

Now let $A = U\Sigma V^T$ be a singular value decomposition of A , *i.e.*, U and V are orthogonal, and $\Sigma = \mathbf{diag}(\sigma_1, \dots, \sigma_n)$. Our observation above, with $Q = U^T$, $Z = V$, shows that

$$\inf_{\mathcal{V} : \dim \mathcal{V} = k} \kappa_{\mathcal{V}}(A) = \inf_{\mathcal{V} : \dim \mathcal{V} = k} \kappa_{\mathcal{V}}(\Sigma).$$

So we can just as well solve the problem for the diagonal matrix Σ . (To reconstruct a subspace of dimension k on which A has least condition number, we find a subspace of dimension k for which Σ has least condition number, and multiply it by V .)

Now our problem is to find a subspace \mathcal{V} of dimension k which minimizes $\kappa_{\mathcal{V}}(\Sigma)$. We will show that

$$\inf_{\mathcal{V} : \dim \mathcal{V} = k} \kappa_{\mathcal{V}}(\Sigma) = \max \left(\frac{\sigma_{n-k+1}}{\sigma_k}, 1 \right) = \begin{cases} 1 & k \leq \lceil n/2 \rceil, \\ \sigma_{n-k+1}/\sigma_k & k > \lceil n/2 \rceil, \end{cases} \quad (2)$$

Let $\{e_1, \dots, e_n\}$ be the standard basis for \mathbf{R}^n , *i.e.*, for $i = 1, \dots, n$, $e_{ij} = 0$ if $i \neq j$ and $e_{ij} = 1$ otherwise.

We first give a simple result. Suppose $i < j$, and let σ satisfy $\sigma_i \geq \sigma \geq \sigma_j$. Then there is a unit vector $z \in \mathbf{span}\{e_i, e_j\}$ for which $\|\Sigma z\| = \sigma$. This can be seen several ways. For example, we can rotate a unit vector z from e_i towards e_j . The norm $\|\Sigma z\|$ varies continuously from σ_i to σ_j , and therefore has the value σ at some rotation angle. We can easily construct such a z . If $\sigma_i = \sigma_j$, we can take $z = e_i$ or $z = e_j$. If $\sigma_i > \sigma_j$, we can take

$$z = \frac{(\sigma^2 - \sigma_j^2)^{1/2} e_i + (\sigma_i^2 - \sigma^2)^{1/2} e_j}{(\sigma_i^2 - \sigma_j^2)^{1/2}}.$$

It is easily verified that $\|z\| = 1$ and $\|\Sigma z\| = \sigma$.

2.1 Case 1: $k \leq \lceil n/2 \rceil$

To establish (2), we will construct a subspace \mathcal{V}^* of dimension k , with $\kappa_{\mathcal{V}^*}(\Sigma) = 1$. We will construct an orthonormal basis $\{z_0, z_1, \dots, z_{k-1}\}$ for \mathcal{V}^* . We start with $z_0 = e_{\lceil n/2 \rceil}$. Note that $\|\Sigma z_0\| = \sigma_{\lceil n/2 \rceil}$.

Next, we choose a unit vector $z_1 \in \mathbf{span}\{e_{\lceil n/2 \rceil - 1}, e_{\lceil (n+1)/2 \rceil + 1}\}$ that satisfies $\|\Sigma z_1\| = \sigma_{\lceil n/2 \rceil}$. We can do this using our simple result above, noting that

$$\sigma_{\lceil n/2 \rceil - 1} \geq \sigma_{\lceil n/2 \rceil} \geq \sigma_{\lceil (n+1)/2 \rceil + 1}.$$

We note that $z_1 \perp z_0$ and $\Sigma z_1 \perp \Sigma z_0$.

We continue the construction, taking z_2 as any unit vector

$$z_2 \in \mathbf{span}\{e_{\lceil n/2 \rceil - 2}, e_{\lceil (n+1)/2 \rceil + 2}\}$$

that satisfies $\|\Sigma z_2\| = \sigma_{\lceil n/2 \rceil}$. This continues, until we have unit vectors z_0, \dots, z_{k-1} . These vectors are mutually orthogonal, since each one is in the span of two standard basis vectors, and these pairs of standard basis vectors are disjoint. Since Σ is a diagonal matrix, the vectors $\Sigma z_0, \dots, \Sigma z_{k-1}$ are mutually orthogonal.

We now show that $\kappa_{\mathcal{V}^*}(\Sigma) = 1$. For any nonzero vector $b \in \mathcal{V}^*$, the gain of Σ in the direction of b , $\|\Sigma b\|/\|b\| = \sigma_{\lceil n/2 \rceil}$, because the gain of Σ in the direction of any unit vector in the orthonormal basis $\{z_0, \dots, z_{k-1}\}$ of \mathcal{V}^* is $\sigma_{\lceil n/2 \rceil}$. Thus $G_{\max} = G_{\min} = \sigma_{\lceil n/2 \rceil}$, and therefore $\kappa_{\mathcal{V}^*}(\Sigma) = 1$.

2.2 Case II: $k > \lceil n/2 \rceil$

To establish (2), we first construct a subspace \mathcal{V}^* of dimension k , with $\kappa_{\mathcal{V}^*}(\Sigma) = \sigma_{n-k+1}/\sigma_k$, and then show that for any subspace \mathcal{V} of dimension k , $\kappa_{\mathcal{V}}(\Sigma) \geq \kappa_{\mathcal{V}^*}(\Sigma)$.

We will construct an orthonormal basis for \mathcal{V}^* . We start with the $2k - n$ vectors $\{e_{n-k+1}, e_{n-k}, \dots, e_{k-1}, e_k\}$. We will choose $n - k$ unit vectors, z_1, \dots, z_{n-k} , such that

$$\{z_1, \dots, z_{n-k}, e_{n-k+1}, \dots, e_{k-1}, e_k\}$$

forms an orthonormal basis for \mathcal{V}^* . The $n - k$ unit vectors z_1, \dots, z_{n-k} will be chosen in $\text{span}\{e_1, \dots, e_{n-k}, e_{k+1}, \dots, e_n\}$, and will therefore be orthogonal to $\{e_{n-k+1}, \dots, e_k\}$.

Choose a unit vector $z_1 \in \text{span}\{e_1, e_n\}$, satisfying $\|\Sigma z_1\| = \sigma_k$. We can do this using the simple result given earlier, since $\sigma_1 \geq \sigma_k \geq \sigma_n$. We note that $z_1 \perp e_j$, and $\Sigma z_1 \perp \Sigma e_j$, $j = n - k + 1, \dots, k$.

We continue the construction, choosing a unit vector $z_2 \in \text{span}\{e_2, e_{n-1}\}$, satisfying $\|\Sigma z_2\| = \sigma_k$. This continues, until we have chosen a unit vector z_{n-k} in $\text{span}\{e_{n-k}, e_{k+1}\}$, satisfying $\|\Sigma z_{n-k}\| = \sigma_k$.

The vectors z_1, \dots, z_{n-k} are mutually orthogonal, since each one is in the span of two standard basis vectors, and these pairs of standard basis vectors are disjoint. Also $z_i \perp e_j$ for $i = 1, \dots, n - k$ and $j = n - k + 1, \dots, k$, since each vector z_i is in the span of two standard basis vectors which are not in the set $\{e_{n-k+1}, \dots, e_k\}$. Thus $\{z_1, \dots, z_{n-k}, e_{n-k+1}, e_{n-k}, \dots, e_{k-1}, e_k\}$ forms an orthonormal basis for \mathcal{V}^* . Similarly, since Σ is a diagonal matrix, the vectors $\Sigma z_1, \dots, \Sigma z_{n-k}, \Sigma e_{n-k+1}, \dots, \Sigma e_k$ are mutually orthogonal.

We now show $\kappa_{\mathcal{V}^*}(\Sigma) = \sigma_{n-k+1}/\sigma_k$. Let b any nonzero vector in \mathcal{V}^* , say,

$$b = \beta_1 z_1 + \dots + \beta_{n-k} z_{n-k} + \beta_{n-k+1} e_{n-k+1} + \dots + \beta_k e_k.$$

The gain of Σ in the direction b is

$$\begin{aligned} \frac{\|\Sigma b\|}{\|b\|} &= \left(\frac{\sum_{i=1}^{n-k} \beta_i^2 \|\Sigma z_i\|^2 + \sum_{j=n-k+1}^k \beta_j^2 \|\Sigma e_j\|^2}{\sum_{i=1}^{n-k} \beta_i^2 \|z_i\|^2 + \sum_{j=n-k+1}^k \beta_j^2 \|e_j\|^2} \right)^{1/2} \\ &= \left(\frac{\sum_{i=1}^{n-k} \beta_i^2 \sigma_k^2 + \sum_{j=n-k+1}^k \beta_j^2 \sigma_j^2}{\sum_{i=1}^n \beta_i^2} \right)^{1/2}, \end{aligned}$$

and therefore $\sigma_{n-k+1} \geq \|\Sigma b\|/\|b\| \geq \sigma_k$. For $b = e_{n-k+1}$, $\|\Sigma b\|/\|b\| = \sigma_{n-k+1}$, so $G_{\max} = \sigma_{n-k+1}$; for $b = e_k$, we have $\|\Sigma b\|/\|b\| = \sigma_k$, so $G_{\min} = \sigma_k$. It follows that $\kappa_{\mathcal{V}^*}(\Sigma) = \sigma_{n-k+1}/\sigma_k$.

Now we will show that for any subspace \mathcal{V} of dimension k , $\kappa_{\mathcal{V}}(\Sigma) \geq \sigma_{n-k+1}/\sigma_k$. By the Courant-Fischer theorem, for any subspace \mathcal{V} of dimension k , $G_{\max} \geq \sigma_{n-k+1}$ and $G_{\min} \leq \sigma_k$. It follows that $\kappa_{\mathcal{V}}(\Sigma) = G_{\max}/G_{\min} \geq \sigma_{n-k+1}/\sigma_k$. This establishes (2), and therefore (1).

References

- [1] D. S. Bernstein. *Matrix Mathematics: Theory, facts, and formulas, with application to linear systems theory*. Princeton University Press, 2005.

- [2] T. F. Chan and D. E. Foulser. Effectively well-conditioned linear systems. *SIAM Journal on Scientific and Statistical Computing*, 9(6):963–968, November 1988.
- [3] R. A. Horn and C. A. Johnson. *Matrix Analysis*. Cambridge University Press, 1985.



A note from the Rejecta Mathematica editorial board regarding
Automatic CounTilings by Doron Zeilberger

The following paper was submitted by Doron Zeilberger in response to an invitation to contribute a paper to the inaugural issue of *Rejecta Mathematica*. In lieu of a traditional open letter, we would like to refer the reader to Prof. Zeilberger's website:

<http://www.math.rutgers.edu/~zeilberg/mamarim/mamarimhtml/tilings.html>

As a brief summary, this paper considers the problem of computing the number of ways in which a $k \times n$ rectangle can be covered by a given set of tiles. The paper in fact describes a Maple program – available at the website above – which will tackle the problem for all n , given any k and any set of allowable tiles. What was once a problem that would have to be tackled on a case-by-case basis is approached with a unified treatment that relies on a computer to discover the appropriate “structure theorems.” As Prof. Zeilberger claims on his website, “what is so nice about it is that everything is done by machine: the combinatorics, the algebra, and the analysis.”

The rejection history of this paper is well-documented on Prof. Zeilberger's website. A primary source of dispute in the original referee review was whether Prof. Zeilberger's algorithm alone constituted a significant, original mathematical contribution. Prof. Zeilberger includes the original review, which he describes as “narrow-minded and ignorant”, offers his own reply, and expresses his opinion about the “flawed” editorial policies of the rejecting journal.

For full context, software code, and examples, the interested reader is invited to read both the paper and the website.

Rejecta Mathematica Editors

Automatic CounTilings

Doron ZEILBERGER¹

I Have a Dream

One day it would be possible to write in English, or in an English-like super-high-level programming language, the following command:

Write a Maple program that inputs an arbitrary positive integer k , a symbol n , and an arbitrary set T of tiles, as well as a symbol t , and outputs the rational function of t whose Maclaurin series's coefficient of t^n is the number of different tilings of a k by n rectangle by the tiles of T .

But Yesterday's Dream Came To Pass

This day hasn't arrived yet, so I had to spend a few weeks writing such a program myself. This is the human-generated Maple package TILINGS accompanying this article.

But, this is already a great step forward. Once I wrote the program, my beloved servant, Shalosh B. Ekhad, can solve (potentially) infinitely many enumeration problems, completely rigorously, for enumerating tilings of rectangles of *arbitrary* width and using an *arbitrary* set of tiles (the tiles even do not have to be connected).

Traditionally, a human would have had to tackle each specific width (k) and each specific set of tiles (T), (even the set consisting of the two dimers, the vertical and horizontal 1 by 2 rectangles) *one at a time*. He or she would have had to figure out the combinatorics of the situation, establish a "structure theorem", then deduce from this a *set of equations*, and finally, solve this set of equations (this very last part, starting in the seventies, was possibly aided by Maple or Mathematica). Of course, the human will only be able to do it for very small k , and very simple sets of tiles, since very soon the structure theorem, and consequently the set of equations, would be too hard to *derive* by hand, let alone solve.

Once the meta-program above would be realized, I would become superfluous. But even at the present level, the human-made Maple program TILINGS can generate, in principle, infinitely many articles, and Ph.D. theses, proving general, rigorous and *interesting* results. Except for very simple cases, of course, the complexity of the proofs and results are beyond mere humans.

The Simplest Not-Completely-Trivial Example

Let $a(n)$ be the number of ways of tiling a $2 \times n$ rectangle using the horizontal and vertical dominoes,

¹ Department of Mathematics, Rutgers University (New Brunswick), Hill Center-Busch Campus, 110 Frelinghuysen Rd., Piscataway, NJ 08854-8019, USA. [zeilberg at math dot rutgers dot edu](mailto:zeilberg@math.rutgers.edu), <http://www.math.rutgers.edu/~zeilberg>. First version: Jan. 20, 2006. Accompanied by Maple packages TILINGS and RecTILINGS downloadable from Zeilberger's website. Supported in part by the NSF.

i.e. the set of tiles is:

$$XX \quad , \quad \begin{array}{c} X \\ X \end{array} \quad .$$

Let $\mathcal{A}(n)$ be the set of tilings of a $2 \times n$ rectangular board by the above two tiles. Of course, if $n = 0$, then $\mathcal{A}(n)$ consists of one tiling only: the **empty tiling**, but if $n > 0$ then there are two cases to consider, regarding the bottom-left cell.

Case I: It participates in the vertical tile, in which case the board looks like

$$\begin{array}{c} 10000 \dots 0 \\ 10000 \dots 0 \end{array} \quad ,$$

where 1 means an occupied cell and 0 means a still-untiled cell.

Case II: It participates in the horizontal tile, in which case the board looks like

$$\begin{array}{c} 00000 \dots 0 \\ 11000 \dots 0 \end{array} \quad .$$

Now the first case is clearly “isomorphic” to $\mathcal{A}(n-1)$, and the isomorphism is obtained by “chop the 1’s” (i.e. remove the vertical tile), reducing the problem to that of tiling a $2 \times (n-1)$ board. But Case II requires us to introduce a new **auxiliary** set, let’s call it $\mathcal{B}(n)$, that of tiling a “jagged” rectangle with the two leftmost cells of the bottom row removed.

So now we forget about the original problem, and try to find a “structure theorem” for $\mathcal{B}(n)$. Of course, if $n = 0, 1$, then $\mathcal{B}(n)$ is the empty set, but if $n \geq 2$ then, we consider the leftmost unoccupied cell of the top row. Now there is only one case to consider: it is tiled by a horizontal tile, getting a board of the form

$$\begin{array}{c} 11000 \dots 0 \\ 11000 \dots 0 \end{array} \quad ,$$

which is trivially “isomorphic” to $\mathcal{A}(n-2)$.

We now have two “structure theorems”:

$$\mathcal{A}(n) \equiv \mathcal{A}(n-1) \cup \mathcal{B}(n) \quad (n > 0) \quad ,$$

$$\mathcal{B}(n) \equiv \mathcal{A}(n-2) \quad . \quad \quad \quad (\text{MarkovianScheme})$$

Now taking cardinalities, calling $a(n) := |\mathcal{A}(n)|$, $b(n) := |\mathcal{B}(n)|$, we have the system of **linear recurrences**

$$a(n) = a(n-1) + b(n) \quad (n > 0) \quad ,$$

$$b(n) = a(n-2) \quad ,$$

with the obvious initial conditions $a(-2) = 0, a(-1) = 0, a(0) = 1$.

Doing the usual generatingfunctionology, we get the generating functions

$$f(t) := \sum_{n=0}^{\infty} a(n)t^n = \frac{1}{1-t-t^2} \quad ,$$

$$g(t) := \sum_{n=0}^{\infty} b(n)t^n = \frac{t^2}{1-t-t^2} \quad .$$

Equivalently, we can use **Polya's Picture Writing**[P] and use **weight-enumerators** to deduce the system

$$f(t) = 1 + tf(t) + g(t) \quad ,$$

$$g(t) = t^2 f(t) \quad ,$$

and solve the system of **two linear equations** in the **two unknowns** $f(t), g(t)$. But we can do better still: keep track of the number of occurrences of each tile. Introducing the variables h and v for a horizontal and vertical tile respectively, and defining the **weight** of a tiling to be $h^{\#horizontal_tiles} v^{\#vertical_tiles} t^n$, where n is its length, and defining $F(t, h, v)$ to be the sum of all the weights of all tilings, and analogously $G(t, h, v)$ for the set \mathcal{B} , we get the system

$$F(t, h, v) = 1 + tvF(t, h, v) + hG(t, h, v) \quad ,$$

$$G(t, h, v) = t^2 h F(t, h, v) \quad ,$$

and solving this system gives:

$$F(t, h, v) := \frac{1}{1 - vt - h^2 t^2} \quad ,$$

$$G(t, h, v) := \frac{t^2 h}{1 - vt - h^2 t^2} \quad .$$

Note that $f(t)$, $g(t)$, and $F(t, h, v)$, $G(t, h, v)$ are **rational functions**. In particular, taking the **partial-fraction** decomposition of $f(t)$, over the reals, one easily gets the **asymptotics** for $a(n)$, namely $(\sqrt{5}/\phi) \cdot \phi^n$, where ϕ is the Golden Ratio. Also by differentiating w.r.t. h and v and then plugging-in $h = 1$ and $v = 1$ one gets the generating function for “the total number of horizontal tiles” and “the total number of vertical tiles” respectively, from which once again, we can deduce asymptotics and get the asymptotic averages by dividing by the already-known asymptotic “number of tilings”. Ditto for the higher moments and correlation (of course in this toy problem the correlation is tautologically -1).

But there is more! Now that we have quick ways for computing $a(n)$ and $b(n)$, up to very large n , either via the generating function or directly by (*MarkovianScheme*), we can use the latter, with the help of **Wilf's methodology**[W] to do **sequencing**, **ranking**, and **random selection**. Let's just consider the last, that of generating **uniformly at random**, such a tiling of the $2 \times n$ rectangular board. Use a loaded coin, with probabilities $a(n-1)/a(n)$ and $b(n)/a(n)$ to decide whether to cover the leftmost-bottommost cell with a vertical or horizontal tile respectively. In the former case, continue recursively. In the latter case, use a very biased coin with probabilities

$b(n)/a(n-2) = 1$ and 0 to decide whether to tile the leftmost 0 at the top row, with a horizontal tile or vertical tile respectively, and then continue recursively.

The General Case

The same reasoning applies to an **arbitrary** (but fixed, i.e. numeric) width k and an *arbitrary* (but of course finite) set of (finite) “tiles”. Here a tile is any finite set of lattice points, that does not need to be connected, and a tiling is a covering by translations of the tiles.

The logic is the same as above, but instead of only one ‘auxiliary’ set $\mathcal{B}(n)$ (that was trivially isomorphic to $\mathcal{A}(n-2)$) we get many more such ‘stepping-stones’. These auxiliary sets are tilings of jagged rectangular boards that can be represented as, for example, (here the width, k , is 4) :

$$\begin{array}{r} 0101\dots 0 \\ 1001\dots 0 \\ 0000\dots 0 \\ 1110\dots 0 \end{array} \quad ,$$

where \dots stands for 0’s (i.e. still unoccupied cells). We can code such jagged boards on the computer just as above as matrices of 0’s and 1’s leaving the \dots implicit, and making it right-justified (in addition to left-justified), and hence there should be at least one 1 in the rightmost column. We also assume that the leftmost column has at least one 0, or else it is equivalent to a smaller board with that leftmost column removed (giving the removed column its due weight by multiplying by t when we do the equations). For each such jagged board we (or rather the computer) locates the **fundamental free cell** which is the lowest 0 in the leftmost column. We then consider which of the tiles in the set of tiles can cover that cell, and in what relative position of the considered tile. For each such scenario, the computer places that tile, getting a different jagged board, that may be a previously-encountered one or a brand-new one. For each new board, we do it again. The computer checks at each stage whether each and every jagged-board encountered so-far already has a structure theorem expressing it a union of isomorphic copies of other jagged boards, and if not, keeps creating new structure theorems, that in turn usually bring in new kinds of jagged boards. Eventually there won’t be any new ones (pigeonhole!) and then the process halts.

Next, to get generating functions, the computer (all by itself) takes weight-enumerators, getting a (usually very large) system of linear equations for the set of unknown weight-enumerators for these jagged boards, that includes, in its midst, the original, non-jagged rectangular board. Then, Maple **solves** this system, giving in particular the generating function for the original set, the rational function, (in t , or in t as well as in the tile-variables) whose Maclaurin series’s coefficient of t^n is the number of (or weight-enumerating polynomial for) tilings of a $k \times n$ rectangular board by the given set of tiles.

Once the generating function is known, we can process it, using Maple’s built-in partial-fraction (over the reals) (**parfrac**) and differentiation, **diff**, to get asymptotic size, and asymptotic averages, variance, covariance, and higher moments for the random variables “number of occurrences

of a given tile-type”.

Disclaimer: This last part involves floating-point calculations and may not be entirely rigorous. In principle, of course, one can stay within algebraic numbers, but the output would not be insightful to human eyes.

The Maple Package TILINGS

All this (and more!) is implemented in the Maple package TILINGS, that can be downloaded from

<http://www.math.rutgers.edu/~zeilberg/tokhniot/TILINGS> .

Once you downloaded it into a directory in your computer, calling it TILINGS (and **not** TILINGS.txt), go into Maple, and type:

```
read TILINGS:
```

and then follow the instructions given there. In particular, for a list of the main functions, type `ezra()`; and for help with a specific function, type `ezra(FunctionName)`;

Also available is a precursor Maple package RecTILINGS that only handles rectangular tiles. It can be downloaded from:

<http://www.math.rutgers.edu/~zeilberg/tokhniot/RecTILINGS> .

A Very Quick Overview

Full details are given on-line, but for the benefit of the lazy reader who does not use Maple, let me just mention that the main functions are:

`GFt(k,Tiles,t)` and `GFtH(k,Tiles,t,H)` .

The former computes the generating function (in t only) and the latter the weight-enumerator (in t and in the variables $H[tile]$, one for each tile).

`Rt(a,b)` is the package’s shorthand for a $b \times a$ rectangular tile. For example, the toy problem we did above by hand is reproduced by:

```
GFt(2,[Rt(1,2),Rt(2,1)],t); and
```

```
GFtH(2,[Rt(1,2),Rt(2,1)],t,H); .
```

`Sidra(k,ListOfTiles,N)`; gives you the first $N + 1$ terms of the enumerating sequence (that can automatically be sent to Sloane’s database), while `SidraW(k,H,ListOfTiles,N)`; gives you the first $N + 1$ terms in the sequence of weight-enumerators, that are polynomials in $H[tile]$ ’s.

`Astat(k,ListOfTiles,n)`; gives you the asymptotic statistics: average number of occurrences for

each tile, the variance, and the asymptotic covariance matrix for these random variables (in the order of appearance in the list `ListOfTiles`). `AstatV` is a verbose version of `Astat`.

Finally, `RandTiling(ListOfTiles,k0,n0)`; gives you a (uniformly) random tiling of a $k0 \times n0$ rectangular board, using the given tiles (e.g. try: `RandTiling([Rt(1,2),Rt(2,1)],4,10)`), and

`RandTilingNice(ListOfTiles,k0,n0)`; gives you the same in nice, matrix form.

`BuildTree` and `BuildTreeH` actually gives you the Markovian Enumeration Scheme, complete with generating functions, and `Sipur` and `SipurArokh` tell you everything you always wanted to know about the tilings.

Sample Input and Output

The webpage of this article,

<http://www.math.rutgers.edu/~zeilberg/mamarim/mamarimhtml/tilings.html> ,

contains several sample input and output files. Readers can generate their own output for their favorite set of tiles and for their favorite width.

An “Almost” Automatic Proof of Kasteleyn’s Formula

In the very special case of two dimer tiles (1×2 and 2×1), Kasteleyn and Fisher&Temperley ([K][FT], see [AS] for a beautiful recent survey) gave their deservedly celebrated beautiful formula for the weight-enumerator of an $m \times n$ rectangle for *arbitrary* (i.e. symbolic) m and n . To wit, if m and n are both even positive integers, and z and z' the variables for the horizontal and vertical dimers respectively (what we called h and v above), then

$$Z_{m,n}(z, z') = 2^{mn/2} \prod_{r=1}^{m/2} \prod_{s=1}^{n/2} \left[z^2 \cos^2 \frac{r\pi}{m+1} + z'^2 \cos^2 \frac{s\pi}{n+1} \right] , \quad (K)$$

with a similar formula when n is odd. At this time of writing, computers can’t prove this formula in general. But for any *specific* m (but general(!) n), it is now a routine verification, thanks to our Maple package `TILINGS`. Of course, in practice, it can only be done for small values of m ($m \leq 10$ on our computer), but with bigger and future computers, one would be able to go further. The proof is a beautiful example of (rigorous!) generalization from finitely many cases, combined with ‘general-nonsense’ linear algebra *handwaving*, that is nevertheless fully rigorous.

Let’s call the right side of (K) $W_{m,n}(z, z')$. We have to prove, for our given m , that $W_{m,n}(z, z') = Z_{m,n}(z, z')$ for *all* n . For a fixed even m it is readily seen that $W_{m,n}(z, z')$ is expressible as a product of $m/2$ different dilations of $U_n(z)$, the Tchebycheff polynomials of the second kind. It is well-known (and easily seen) that $U_n(z)$ satisfies a (homogeneous) *linear recurrence equation with constant* (i.e. *not depending on* n) *coefficients* of order 2, and hence so does any of its dilations $U_n(\beta z)$, and the product of $m/2$ of these creatures consequently satisfies a linear recurrence equation with constant coefficients of order $2^{m/2}$. Hence its generating function, in t , is a rational function

whose denominator has degree $2^{m/2}$ and numerator degree $\leq 2^{m/2} - 1$. We have to prove that this generating function coincides with the ‘real thing’, the rational function outputted by procedure GFtH of TILINGS, that also turns out to have the same degrees. But to prove that two rational functions whose numerator-degree is p and denominator-degree is q are identical, we only have to check, by elementary linear algebra, that the first $p + q + 1$ terms in their Maclaurin series coincide (in our case $2^{m/2+1}$), and this is a routine check, done in our package TILINGS by typing `ProveKasteleyn(m);`, for the general case (with z and z'), and `ProveKasteleyn1(m);`, for the straight-enumeration case ($z = z' = 1$).

By hindsight, Kasteleyn and Fisher&Temperley were extremely lucky, since in the very special case of the two dimer tiles, they were able to use graph theory and Pfaffians. The general case seems, at present, out of reach. Even the monomer-dimer problem is still wide open, let alone an “explicit” expression for the weight-enumerators for tilings of an $m \times n$ rectangle using other sets of tiles. Our Maple package TILINGS is a *research tool* that can automatically discover and rigorously prove results for **specific** m but general n . Let’s hope that the human can use it to discover new **ansatzes** by which to conjecture and hopefully prove some “explicit” form for the monomer-dimer and more general sets of tiles, where *both* m and n are **symbolic**. From this, one should be able to extract, at least, an “explicit” expression for the so-called *thermodynamic limit*, or at the very least, determine rigorously the *critical exponent*.

I hope to explore these speculations in a subsequent article.

References

- [AS] Frederico Ardila and Richard S. Stanley, *Tilings*, preprint, www.arxiv.org/math.CO/0501170, 2005.
- [FT] M. Fisher and H. Temperley, *Dimer Problems in Statistical Mechanics-an exact result*, Philos. Mag. **6** (1961), 1061-1063.
- [K] P. W. Kasteleyn, *The statistics of dimers on a lattice: I. The number of dimer arrangements in a quadratic lattice*, Physica **27** (1961), 1209-1225.
- [P] G. Polya, *On Picture Writing*, Amer. Math. Monthly **63** (1956), 689-697.
- [W] Herbert S. Wilf, *A Unified Setting for Sequencing, Ranking, and Selection Algorithms for Combinatorial Objects*, Adv. in Math. **24** (1977), 281-291.



An open letter concerning
*Alexander duality for monomial ideals and
their resolutions*

Dear Reader,

This article was submitted to *Journal of Pure and Applied Algebra* on December 15, 1998, and it was rejected with a very short report about eight months later, the cited reason being that it was too long for its content. By the time I received that overdue rejection, I was nearly done writing a sequel,

Ezra Miller, *The Alexander duality functors and local duality with monomial support*, *Journal of Algebra* **231** (2000), 180–234.

which contained more general results. The sequel has been well-cited, but the current article was already on the arXiv (math.AC/9812095), and according to Google Scholar it has also been well-cited. In fact, this article has been cited more than most of my others—as much or more, for example, than my articles in *Journal of the American Mathematical Society* and *Duke Mathematical Journal*. It seemed a shame that what is apparently a useful article should languish in eternal semipublication, so I submitted it to *Rejecta Mathematica*.

Why is this article useful? It is more concrete than its sequel: more examples, more illustrations, and fewer functors. The article contains no known errors and no known uncited rederivations of earlier work; in fact, subsequent work (by other authors as well as in its sequel) has confirmed the results herein by independent methods many times over. The article is unchanged from the version submitted to *Journal of Pure and Applied Algebra*.

Ezra Miller
25 February 2008

Alexander Duality for Monomial Ideals and Their Resolutions

Ezra Miller

Abstract

Alexander duality has, in the past, made its way into commutative algebra through Stanley-Reisner rings of simplicial complexes. This has the disadvantage that one is limited to square-free monomial ideals. The notion of Alexander duality is generalized here to arbitrary monomial ideals. It is shown how this duality is naturally expressed by Bass numbers, in their relations to the Betti numbers of a monomial ideal and its Alexander dual. The effect of Alexander duality on free resolutions is studied in the context of cellular resolutions. Relative cohomological constructions on cellular complexes are shown to relate cellular free resolutions of a monomial ideal to free resolutions of its Alexander dual ideal.

Introduction

Alexander duality in its most basic form is a relation between the homology of a simplicial complex Γ and the cohomology of another simplicial complex Γ^\vee , called the *dual* of Γ . Recently there has been much interest in the consequences of this relation when applied to the monomial ideals which are the Stanley-Reisner ideals I_Γ and I_{Γ^\vee} for the given simplicial complex and its Alexander dual. This has the limitation that Stanley-Reisner ideals are always squarefree. The first aim of this paper is to define Alexander duality for arbitrary monomial ideals and then generalize some of the relations between I_Γ and I_{Γ^\vee} . A second goal is to demonstrate that Bass numbers are the proper vessels for the translation of Alexander duality into commutative algebra. The final goal is to reveal the connections between Alexander duality and the recent work on cellular resolutions.

There are two “minimal” ways of describing an arbitrary monomial ideal: via the minimal generators or via the (unique) irredundant irreducible decomposition. Given a monomial ideal I , Definition 1.5 describes a method for producing another monomial ideal I^\vee whose minimal generators correspond to the irredundant irreducible components of I . Miraculously, this is enough to guarantee that the minimal generators of I correspond to the irreducible components of I^\vee . It is particularly easy to verify that this reversal of roles takes place for the squarefree ideals $I = I_\Gamma$ and $I^\vee = I_{\Gamma^\vee}$ above (Proposition 1.10). A connection with linkage and canonical modules is described in Theorem 2.1.

One can also deal with Alexander duality as a combinatorial phenomenon, thinking of Γ as an order ideal in the lattice of subsets of $\{1, \dots, n\}$. The Alexander dual Γ^\vee is then given by the complement of the order ideal, which gives an order ideal in the opposite lattice. For squarefree monomial ideals all is well since the only monomials we care about are represented precisely by the lattice of subsets of $\{1, \dots, n\}$. For general monomial ideals we instead consider the larger lattice \mathbb{Z}^n , by which we mean the poset with its natural partial order \preceq . Then a monomial ideal I can be regarded as a dual order ideal in \mathbb{Z}^n , and I^\vee is constructed (roughly) from the complementary set of lattice points, which is an order ideal—see Definition 2.9. It is Theorem 2.13 which proves the equivalence of the two definitions.

Bass numbers first assert themselves in Section 3. Their relations to Betti numbers for monomial modules (Corollary 3.6 and Theorem 3.12) are derived as consequences of graded local duality and Alexander duality (in its avatar as lattice duality in \mathbb{Z}^n). The Bass-Betti relations are then massaged to equate the localized Bass numbers of I (Definition 4.8) with the Betti numbers of I^\vee in the first of the two central results of this paper, Theorem 4.10. Theorem 2.13 is then recovered as a special case of this main result, which also finds an application to inequalities between the Betti numbers of dual ideals (Theorem 4.13) generalizing those for squarefree ideals in [2].

The extension of Alexander duality to resolutions is accomplished in Sections 5 and 6. A new canonical and geometric resolution, the *cohull resolution* is constructed in Definition 5.15. It should be thought of as Alexander dual to the *hull resolution* of [4] (which is similarly canonical and geometric). Roughly speaking, the cohull resolution is constructed from the irreducible components instead of the minimal generators. The cohull resolution owes its existence to the second central result of the paper, Theorem 5.8, which is a more general result on duality for cellular resolutions. Its proof, which is resolutely algebraic, is the content of Section 6. The idea is to deform an ideal into its dual step by step via Definition 6.1 and keep track of the deformations on cellular resolutions (Theorem 6.9). The final step, taken in Theorem 6.11, is to check the effect of the deformations on the homology of the resolutions.

Acknowledgements. The author would like to express his thanks to Dave Bayer, David Eisenbud, Serkan Hosten, Sorin Popescu, Stefan Schmidt, Frank Sottile, Bernd Sturmfels, and Kohji Yanagawa for their helpful comments and discussions.

1 Definitions and basic properties

For notation, let S be the \mathbb{Z}^n -graded k -algebra $k[x_1, \dots, x_n] \subseteq T := S[x_1^{-1}, \dots, x_n^{-1}]$, where k is a field and $n \geq 2$. If $A \subseteq T$ is any subset, $\langle a \mid a \in A \rangle$ will denote the S -submodule generated by the elements in A , and it may also be regarded as an ideal if $A \subseteq S$. The maximal \mathbb{Z}^n -graded ideal $\langle x_1, \dots, x_n \rangle$ of S will be denoted by \mathfrak{m} . Each (Laurent) monomial in T is specified uniquely by a single vector $\mathbf{a} = (a_1, \dots, a_n) = \sum_i a_i \mathbf{e}_i \in \mathbb{Z}^n$, while each irreducible monomial ideal is specified uniquely by a vector $\mathbf{b} = (b_1, \dots, b_n) \in \mathbb{N}^n$, so the notation

$$x^{\mathbf{a}} = x_1^{a_1} \cdots x_n^{a_n} \quad \text{and} \quad \mathfrak{m}^{\mathbf{b}} = \langle x_i^{b_i} \mid b_i \geq 1 \rangle$$

will be used to highlight the similarity. The \mathbb{Z}^n -graded prime ideals, which are precisely the monomial prime ideals, are indexed by faces of the $(n-1)$ -simplex $\Delta := 2^{\{1, \dots, n\}}$ with vertices $1, \dots, n$. Identifying a face $F \in \Delta$ with its characteristic vector in \mathbb{Z}^n , the monomial prime corresponding to F may be written with the above notation as \mathfrak{m}^F . Note, in particular, that $\mathfrak{m}^{\mathbf{b}}$ need not be an artinian ideal, just as $x^{\mathbf{a}}$ need not have full support. In fact, $\mathfrak{m}^{\mathbf{b}}$ is $\mathfrak{m}^{\sqrt{\mathbf{b}}}$ -primary, where $\sqrt{\mathbf{b}} \in \Delta$ is the face representing the support of \mathbf{b} ; that is, $\sqrt{\mathbf{b}}$ has i^{th} coordinate 1 if $b_i \geq 1$ and 0 otherwise. With this notation, taking radicals can be expressed as $\sqrt{\mathfrak{m}^{\mathbf{b}}} = \mathfrak{m}^{\sqrt{\mathbf{b}}}$.

All modules N and homomorphisms of such will be \mathbb{Z}^n -graded, so that $N = \bigoplus_{\mathbf{a} \in \mathbb{Z}^n} N_{\mathbf{a}}$. In addition, any module that is isomorphic to a submodule of T as a \mathbb{Z}^n -graded module will, if it is convenient, be freely identified with that submodule of T . For instance, the principal ideal generated by $x_1 \cdots x_n$ can be identified with the module $S[-\mathbf{1}]$, where $\mathbf{1} = (1, \dots, 1) \in \mathbb{Z}^n$ and $N[\mathbf{a}]_{\mathbf{b}} = N_{\mathbf{a}+\mathbf{b}}$ for $\mathbf{a}, \mathbf{b} \in \mathbb{Z}^n$. In this paper, ideals will all be proper monomial ideals, and the symbol I will always denote such an ideal. The vector \mathbf{a}_I will denote the exponent on the least common multiple of the minimal generators of I .

Before making the definition of Alexander dual ideal, the next few results make sure that the exponents used to define the set $\text{Irr}(I)$ of irredundant irreducible components of I are $\preceq \mathbf{a}_I$. For the next two results, let Λ denote the set of irreducible ideals containing I .

Lemma 1.1 *If $\mathfrak{m}^{\mathbf{b}} \in \text{Irr}(I)$ then $\mathfrak{m}^{\mathbf{b}}$ is minimal (under inclusion) in Λ .*

Proof: Suppose $\mathfrak{m}^{\mathbf{b}} \neq \mathfrak{m}^{\mathbf{c}}$ and that $\mathfrak{m}^{\mathbf{b}} \supseteq \mathfrak{m}^{\mathbf{c}} \in \Lambda$. If now $I = \mathfrak{m}^{\mathbf{b}} \cap I'$ for some ideal I' then also $I = \mathfrak{m}^{\mathbf{c}} \cap I'$, whence $\mathfrak{m}^{\mathbf{b}} \notin \text{Irr}(I)$. \square

Proposition 1.2 *If $\mathfrak{m}^{\mathbf{b}} \in \text{Irr}(I)$ then for each $i \in \sqrt{\mathbf{b}}$ there is a minimal generator $x^{\mathbf{c}}$ of I with $b_i = c_i$.*

Proof: Suppose $\mathfrak{m}^{\mathbf{b}} \in \text{Irr}(I)$ but the conclusion does not hold. Then given any minimal generator $x^{\mathbf{c}}$ of I , either $b_{i'} \leq c_{i'}$ for some $i \neq i' \in \sqrt{\mathbf{b}}$, or else $b_i < c_i$. In either case, $x^{\mathbf{c}} \in \mathfrak{m}^{\mathbf{b}+\mathbf{e}_i}$, where \mathbf{e}_i is the i^{th} unit vector in \mathbb{Z}^n . Then $\mathfrak{m}^{\mathbf{b}+\mathbf{e}_i} \supseteq I$, contradicting the minimality of $\mathfrak{m}^{\mathbf{b}}$ in Λ . \square

Corollary 1.3 *For any $\mathfrak{m}^{\mathbf{b}} \in \text{Irr}(I)$ we have $\mathbf{b} \preceq \mathbf{a}_I$.* \square

The following notation will be very convenient in the definition and handling of Alexander duality. For any vector $\mathbf{a} \in \mathbb{Z}^n$ and any face $F \in \Delta$, let $\mathbf{a} \cdot F$ denote the restriction of \mathbf{a} to F :

$$(\mathbf{a} \cdot F)_i = \begin{cases} a_i & \text{if } i \in F \\ 0 & \text{otherwise} \end{cases}.$$

This operation may also be thought of as the coordinatewise product of \mathbf{a} and F . If, in addition, $\mathbf{0} \preceq \mathbf{b} \preceq \mathbf{a}$, define $\mathbf{b}^{\mathbf{a}}$ to be the vector whose i^{th} coordinate is $a_i + 1 - b_i$ if $b_i \geq 1$ and 0 otherwise; more compactly,

$$\mathbf{b}^{\mathbf{a}} = (\mathbf{a} + \mathbf{1} - \mathbf{b}) \cdot \sqrt{\mathbf{b}} = (\mathbf{a} + \mathbf{1}) \cdot \sqrt{\mathbf{b}} - \mathbf{b},$$

where $\sqrt{\mathbf{b}}$ is the support of \mathbf{b} , as above. The next result is a first indication of the utility of $\mathbf{b}^{\mathbf{a}}$ when applied to irreducible ideals $\mathfrak{m}^{\mathbf{b}}$.

Proposition 1.4 *If $\mathbf{0} \preceq \mathbf{b}, \mathbf{c} \preceq \mathbf{a}$ then $\mathfrak{m}^{\mathbf{b}} \supseteq \mathfrak{m}^{\mathbf{c}}$ if and only if $\mathbf{b}^{\mathbf{a}} \succeq \mathbf{c}^{\mathbf{a}}$.*

Proof: The condition $\mathfrak{m}^{\mathbf{b}} \supseteq \mathfrak{m}^{\mathbf{c}}$ is equivalent to the combination of (i) $\sqrt{\mathbf{b}} \succeq \sqrt{\mathbf{c}}$ and (ii) $\mathbf{b} \cdot \sqrt{\mathbf{c}} \preceq \mathbf{c}$. Now consider the inequalities in the following chain:

$$\mathbf{b}^{\mathbf{a}} = (\mathbf{a} + \mathbf{1} - \mathbf{b}) \cdot \sqrt{\mathbf{b}} \succeq (\mathbf{a} + \mathbf{1} - \mathbf{b}) \cdot \sqrt{\mathbf{c}} \succeq (\mathbf{a} + \mathbf{1} - \mathbf{c}) \cdot \sqrt{\mathbf{c}} = \mathbf{c}^{\mathbf{a}}.$$

The left inequality is equivalent to (i) since $\mathbf{a} + \mathbf{1} - \mathbf{b}$ has full support, and the right inequality is equivalent to (ii) since $\mathbf{c} \cdot \sqrt{\mathbf{c}} = \mathbf{c}$. It remains only to show that $\mathbf{b}^{\mathbf{a}} \succeq \mathbf{c}^{\mathbf{a}}$ implies both inequalities, and this can be checked coordinatewise. If $c_i = 0$, then both inequalities become trivial; if $c_i > 0$ then $b_i > 0$, and the left inequality becomes an equality while the right inequality becomes $(\mathbf{b}^{\mathbf{a}})_i = a_i + 1 - b_i \geq a_i + 1 - c_i = (\mathbf{c}^{\mathbf{a}})_i$. \square

Corollary 1.3 clears the way for the main definition of this paper:

Definition 1.5 (Alexander duality) *Given an ideal I and $\mathbf{a} \succeq \mathbf{a}_I$, the Alexander dual ideal $I^{\mathbf{a}}$ with respect to \mathbf{a} is defined by*

$$I^{\mathbf{a}} = \langle x^{\mathbf{b}^{\mathbf{a}}} \mid \mathfrak{m}^{\mathbf{b}} \in \text{Irr}(I) \rangle.$$

For the special case when $\mathbf{a} = \mathbf{a}_I$, let $I^{\vee} = I^{\mathbf{a}_I}$.

Remark 1.6 (i) We will never have occasion to take an Alexander dual of the ideal \mathfrak{m} , so $\mathfrak{m}^{\mathbf{a}}$ will retain its original definition.

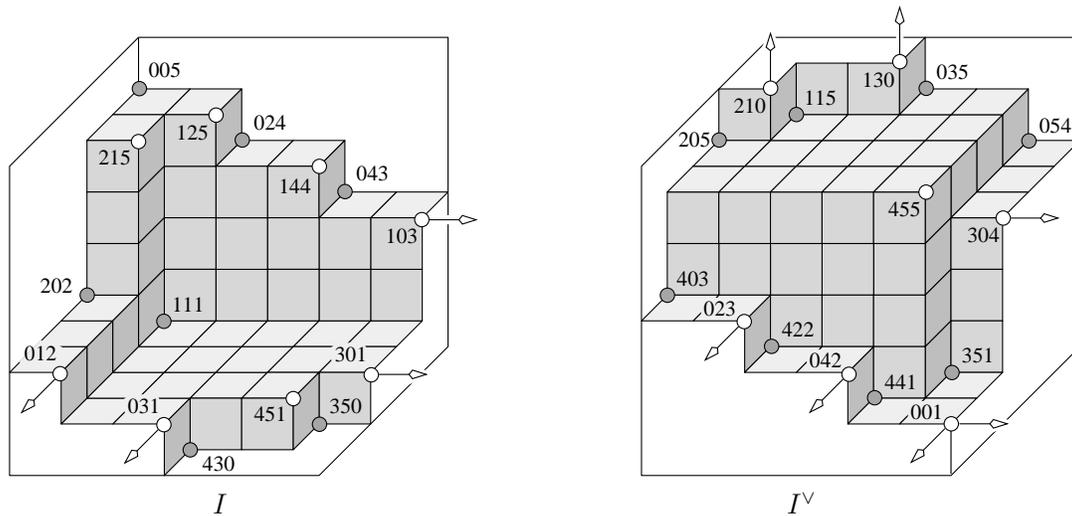
(ii) The dual $I^{\mathbf{a}}$ with respect to any $\mathbf{a} \succeq \mathbf{a}_I$ depends only on $\mathbf{a} \cdot \sqrt{\mathbf{a}_I}$. This is because \mathbf{b} and $\mathbf{a} \cdot \sqrt{\mathbf{b}}$ determine $\mathbf{b}^{\mathbf{a}}$, and $\mathbf{a} \cdot \sqrt{\mathbf{b}} = (\mathbf{a} \cdot \sqrt{\mathbf{a}_I}) \cdot \sqrt{\mathbf{b}}$ for all of the relevant \mathbf{b} by Corollary 1.3. In particular, $I^{\vee} = I^{\mathbf{1}}$ if I is squarefree.

(iii) I^{\vee} is not gotten by taking the depolarization of the Alexander dual of the polarization of I (see [14], Chapter II for polarization). For instance, when $I = \langle x^2, xy, y^2 \rangle$, the polarization is $I_{\text{polar}} = \langle x_1x_2, x_1y_1, y_1y_2 \rangle$, whose canonical Alexander dual is $I_{\text{polar}}^{\vee} = \langle x_1y_1, x_1y_2, x_2y_1 \rangle$. Removing the subscripts on x and y then yields the principal ideal $\langle xy \rangle$, whereas $I^{\vee} = \langle xy^2, x^2y \rangle$.

Proposition 1.7 *The set of generators for $I^{\mathbf{a}}$ given by the definition is minimal. More generally, suppose $\mathbf{a} \succeq \mathbf{a}_I$ and Λ is a collection of integer vectors $\preceq \mathbf{a}$ such that $I = \bigcap_{\mathbf{b} \in \Lambda} \mathfrak{m}^{\mathbf{b}}$. Then $I^{\mathbf{a}} = \langle x^{\mathbf{b}^{\mathbf{a}}} \mid \mathbf{b} \in \Lambda \rangle$, and the intersection determined by Λ is irredundant if and only if the set of generators for $I^{\mathbf{a}}$ is minimal.*

Proof: This follows from Corollary 1.3 and Proposition 1.4. \square

Example 1.8 Let $n = 3$, so that $S = k[x, y, z]$. Figure 1 lists the minimal generators and irredundant irreducible components of an ideal $I \subseteq S$ and its dual I^{\vee} with respect to \mathbf{a}_I . The (truncated) “staircase diagrams” representing the monomials not in these ideals are also rendered in Figure 1. In fact, the staircase diagram for I^{\vee} is gotten by literally turning the staircase diagram for I upside-down (the reader is encouraged to try this). Notice that the support of a minimal generator of I is equal to the support of the corresponding irreducible component of I^{\vee} . \square



$$\begin{aligned}
 I &= \langle z^5, x^2z^2, x^4y^3, x^3y^5, y^4z^3, y^2z^4, xyz \rangle \\
 &= \langle x^2, y, z^5 \rangle \cap \langle y, z^2 \rangle \cap \langle y^3, z \rangle \cap \langle x^4, y^5, z \rangle \cap \langle x^3, z \rangle \cap \langle x, z^3 \rangle \cap \langle x, y^4, z^4 \rangle \cap \langle x, y^2, z^5 \rangle
 \end{aligned}$$

$$\mathbf{a} := \mathbf{a}_I = (4, 5, 5)$$

$$\begin{aligned}
 I^\vee &= \langle z \rangle \cap \langle x^3, z^4 \rangle \cap \langle x, y^3 \rangle \cap \langle x^2, y \rangle \cap \langle y^2, z^3 \rangle \cap \langle y^4, z^2 \rangle \cap \langle x^4, y^5, z^5 \rangle \\
 &= \langle x^3y^5z, y^5z^4, y^3z^5, xyz^5, x^2z^5, x^4z^3, x^4y^2z^2, x^4y^4z \rangle.
 \end{aligned}$$

Figure 1: The truncated staircase diagrams, minimal generators, and irredundant irreducible components for I and I^\vee . Black lattice points are generators, and white lattice points indicate irreducible components. The numbers are to be interpreted as vectors, e.g. $205 = (2, 0, 5)$. The arrows attached to a white lattice point indicate the directions in which the component continues to infinity; it should be noted that a white point has a zero in some coordinate precisely when it has an arrow pointing in the corresponding direction.

Example 1.9 Let Σ_n denote the symmetric group on $\{1, \dots, n\}$ and $\mathbf{c} = (1, 2, \dots, n) \in \mathbb{N}^n$. The ideal $I = \langle x^{\sigma(\mathbf{c})} \mid \sigma \in \Sigma_n \rangle$ is the *permutahedron ideal* determined by \mathbf{c} , introduced in [4], Example 1.9. The results of Example 5.22 below imply that the canonical Alexander dual is the *forest ideal*, which is generated by $2^n - 1$ monomials: $I^\vee = \langle (x^F)^{n-|F|+1} \mid \emptyset \neq F \in \Delta \rangle$. For instance, when $n = 3$,

$$\begin{aligned}
 I &= \langle xy^2z^3, xy^3z^2, x^2yz^3, x^2y^3z, x^3yz^2, x^3y^2z \rangle \\
 I^\vee &= \langle xyz, x^2y^2, x^2z^2, y^2z^2, x^3, y^3, z^3 \rangle.
 \end{aligned}$$

The quotient of S by the forest ideal has the same dimension (over k) as the algebra \mathcal{A}_n generated by the Chern 2-forms of the tautological line bundles over a flag manifold (see [10] and [13]). More precisely, the standard monomials of I^\vee , which are known to be in bijection with the forests on n labelled vertices, are shown in [10] to be a k -basis of \mathcal{A}_n . The minimal free resolution of I^\vee is obtained in Example 5.22, below. \square

Recall that for a simplicial complex $\Gamma \subseteq \Delta$ the *Stanley-Reisner ideal* I_Γ of Γ is defined by the nonfaces of Γ :

$$I_\Gamma = \langle x^F \mid F \notin \Gamma \rangle,$$

and the *Alexander dual simplicial complex* Γ^\vee consists of the complements of the nonfaces of Γ :

$$\Gamma^\vee = \{F \in \Delta \mid \bar{F} \notin \Gamma\},$$

where $\bar{F} = \{1, \dots, n\} \setminus F$. Recall also that I_Γ may be equivalently described as

$$I_\Gamma = \bigcap_{\bar{F} \in \Gamma} \mathfrak{m}^{\bar{F}},$$

since $\mathfrak{m}^F \supseteq I \Leftrightarrow F$ has at least one vertex in each nonface of $\Gamma \Leftrightarrow \bar{F}$ is missing at least one vertex from each nonface of $\Gamma \Leftrightarrow \bar{F}$ is a face of Γ . Applying Definition 1.5 to the latter characterization of I_Γ yields:

Proposition 1.10 *For a simplicial complex $\Gamma \subseteq \Delta$ we have $I_\Gamma^\vee = I_{\Gamma^\vee}$.*

Proof: Observe that $\mathbf{b}^1 = \mathbf{b}$ if $\mathbf{b} \in \{0, 1\}^n$, and use Proposition 1.7 along with Remark 1.6(ii). We get $I_\Gamma^\vee = \langle x^F \mid \bar{F} \in \Gamma \rangle = \langle x^F \mid F \notin \Gamma^\vee \rangle = I_{\Gamma^\vee}$. \square

Thus, as promised, Definition 1.5 generalizes to arbitrary monomial ideals the definition of Alexander duality for squarefree monomial ideals. The connection with the squarefree case is never lost, however, because the general definition does the same thing to the zero-set of I as the squarefree definition does:

Proposition 1.11 *Taking Alexander duals commutes with taking radicals: $\sqrt{I^\vee} = \sqrt{I}^\vee$.*

Proof: Since $\mathbf{0} \preceq \mathbf{b} \preceq \mathbf{a}_I$ whenever $\mathfrak{m}^{\mathbf{b}} \in \text{Irr}(I)$, the equality $\sqrt{\mathbf{b}} = \sqrt{\mathbf{b}^{\mathbf{a}_I}}$ follows from the definitions. Thus,

$$\begin{aligned} \sqrt{I^\vee} &= \langle x^{\sqrt{\mathbf{b}}} \mid \mathfrak{m}^{\mathbf{b}} \in \text{Irr}(I) \rangle \\ &= \langle x^F \mid \mathfrak{m}^F \text{ is minimal among primes containing } I \rangle \\ &= \sqrt{I}^\vee, \end{aligned}$$

the last equality using again the facts mentioned in the first line of the proof of Proposition 1.10. \square

The notion of Alexander duality sheds light on the interconnections between some of the developments in [3], [4], and [15] concerning cellular resolutions and (co)generic monomial ideals. To begin with, consider the following condition on a set of vectors $\{\mathbf{b}^j = (b_1^j, \dots, b_n^j) \in \mathbb{N}^n\}_{j=1}^r$:

$$b_i^j \geq 1 \Rightarrow b_i^j \neq b_i^{j'} \text{ for all } j' \neq j.$$

A *generic* ideal, as defined in [3], is an ideal whose minimal generators have exponent vectors satisfying the above condition; similarly, a *cogeneric* ideal, as defined in [15], is an ideal whose irredundant irreducible components have exponent vectors satisfying the above condition. Using Definition 1.5 the following is immediate (for any $\mathbf{a} \succeq \mathbf{a}_I$).

Proposition 1.12 *$I^{\mathbf{a}}$ is generic if and only if I is cogeneric.* □

Example 1.13 The ideal I in Example 1.8 is generic, while I^{\vee} is cogeneric. □

The connections between the minimal resolutions of such ideals and cellular resolutions will be explored in Section 5.

Recall that the *Castelnuovo-Mumford regularity* and *initial degree* of a \mathbb{Z} -graded S -module L defined respectively by

$$\operatorname{reg}(L) := \max\{j \in \mathbb{Z} \mid \operatorname{Tor}_i(L, k)_{i+j} \neq 0\} \quad \text{and} \quad \operatorname{indeg}(L) := \min\{j \in \mathbb{Z} \mid L_j \neq 0\},$$

where L_j is the j^{th} component of L . The question was raised in [8], Question 10 whether there is a duality for possibly nonradical monomial ideals with the “amazing properties”

- $\operatorname{reg}(I) - \operatorname{indeg}(I) = \dim(S/I^{\vee}) - \operatorname{depth}(S/I^{\vee})$
- I is componentwise linear if and only if S/I^{\vee} is sequentially Cohen-Macaulay

obeyed by Alexander duals in the squarefree case. Here, I is considered in its \mathbb{Z} -grading. Having defined a duality operation in this paper, some comments are obviously warranted.

First of all, it is unrealistic to expect the first property to extend to the arbitrary (nonradical) case since the right-hand side of the equation is bounded while the left-hand side is not, in general. For instance, if $d \in \mathbb{N}$ then $\operatorname{reg}(\mathfrak{m}^{d-1}) - \operatorname{indeg}(\mathfrak{m}^{d-1}) = n(d-1) - d$ while $(\mathfrak{m}^{d-1})^{\vee} = \langle x_1 \cdots x_n \rangle$ is Cohen-Macaulay. Nevertheless, there may be some class of ideals which behaves nicely under some kind of duality, not necessarily as defined here. As to whether or not such a class of ideals exists for the Alexander duality as defined here, such an investigation has not yet been made.

Unfortunately, the second property also fails for I and $I^{\mathbf{a}}$, for somewhat trivial reasons: almost every ideal has an artinian Alexander dual. Specifically, if I is arbitrary and $x = x_1 \cdots x_n$, then $S/(xI)^{\mathbf{a}}$ is artinian (for any $\mathbf{a} \succeq \mathbf{a}_I$), and hence Cohen-Macaulay. But the minimal free resolution of xI is just the shift by $\mathbf{1}$ of the minimal resolution of I . Thus every minimal resolution, be it componentwise linear or not, appears as the resolution of an ideal whose dual is a Cohen-Macaulay ideal; i.e. $S/I^{\mathbf{a}}$ Cohen-Macaulay $\not\Rightarrow I$ componentwise linear.

One might still hope that the implication “ I has a linear resolution $\Rightarrow S/I^{\mathbf{a}}$ is sequentially Cohen-Macaulay” would hold, but even this fails, as the example below shows. The fundamental problem with the nonsquarefree case is that the \mathbb{Z} -degree of an element is not determined by the support of its \mathbb{Z}^n -graded degree, as it is with squarefree monomials. Thus an ideal might have

a linear resolution while its generators have support sets of varying sizes, wreaking havoc with the equidimensionality required for the Cohen-Macaulayness of the dual. Even so, it would be very interesting to know what is the property Alexander dual to “sequentially Cohen-Macaulay”; perhaps this property could relax the requirements of componentwise linearity in a nice way.

Example 1.14 Let $I' = \langle ab, bc, cd \rangle \subseteq S = k[a, b, c, d]$ be the ideal of the “stick twisted cubic” simplicial complex spanned by the edges $\{b, d\}$, $\{b, c\}$, and $\{a, c\}$. It is readily checked that I' has a linear resolution: indeed, $(I')^\vee$ is the ideal of another stick twisted cubic, which is Cohen-Macaulay because the stick twisted cubic is connected and has dimension 1, so [6], Theorem 3 applies. Let

$$\begin{aligned} I &= \mathfrak{m}I' = \langle a^2b, abc, acd, ab^2, b^2c, bcd, abc, bc^2, c^2d, abd, bcd, cd^2 \rangle \\ I^\vee &= \langle b^2d^2, b^2c^2, a^2c^2, abc^2d^2, a^2bcd^2, a^2b^2cd \rangle \end{aligned}$$

with $\mathbf{a}_I = (2, 2, 2, 2)$. Then I has a linear resolution by [8], Lemma 1, and we show that S/I^\vee is not sequentially Cohen-Macaulay.

Recall that for a module N to be sequentially Cohen-Macaulay, we require that there exist a filtration $0 = N_0 \subset N_1 \subset \cdots \subset N_r = N$ such that N_i/N_{i-1} is Cohen-Macaulay for all $i \leq r$ and $\dim(N_{i+1}/N_i) > \dim(N_i/N_{i-1})$ for all $i < r$. It follows from the equidimensionality of N/N_{r-1} and the strict reduction of dimension in successive quotients that N_{r-1} is the top dimensional piece of N ; i.e. N_{r-1} is the intersection of all primary components (of 0 in N) which have dimension $\dim(N)$. Thus it suffices to check that S/I_{top}^\vee is not Cohen-Macaulay, where $I_{\text{top}}^\vee = \langle b^2d^2, b^2cd, abcd, b^2c^2, abc^2, a^2c^2 \rangle$ is the intersection of all primary components of I^\vee which have dimension $2 = \dim(S/I^\vee)$. \square

2 Alternate characterizations of the Alexander dual ideal

Definition 1.5 is quite satisfactory for the consequences just derived from it, but it can sometimes be inconvenient to work with. For instance, it is not obvious from the definition that $(I^{\mathbf{a}})^{\mathbf{a}} = I$, which is fundamental—see Corollary 2.14. For this and other applications, we set out now to find other characterizations of the Alexander dual ideal in Theorem 2.1 and in Definition 2.9 with Theorem 2.13. Along the way, an algebraic analogue of combinatorial lattice duality in \mathbb{Z}^n is defined in Definition 2.3.

First, a result relating Alexander duality to linkage (see [17], Appendix A.9 for a brief introduction to linkage, and references):

Theorem 2.1 *If $\mathbf{a} \succeq \mathbf{a}_I$ then $(\mathfrak{m}^{\mathbf{a}+1} : I^{\mathbf{a}}) = I + \mathfrak{m}^{\mathbf{a}+1}$.*

Proof: Let $\text{Min}(I^{\mathbf{a}})$ denote the exponents on the minimal generators of $I^{\mathbf{a}}$. Then $(\mathfrak{m}^{\mathbf{a}+1} : I^{\mathbf{a}}) = \bigcap_{\mathbf{b} \in \text{Min}(I^{\mathbf{a}})} (\mathfrak{m}^{\mathbf{a}+1} : x^{\mathbf{b}})$. But $x^{\mathbf{c}} \cdot x^{\mathbf{b}} \in \mathfrak{m}^{\mathbf{a}+1} \Leftrightarrow \mathbf{b} + \mathbf{c} \not\leq \mathbf{a} \Leftrightarrow \mathbf{c} \not\leq \mathbf{a} - \mathbf{b} \Leftrightarrow x^{\mathbf{c}} \in \mathfrak{m}^{\mathbf{a}+1-\mathbf{b}}$. Thus,

taking all intersections over $\mathbf{b} \in \text{Min}(I^{\mathbf{a}})$,

$$\bigcap (\mathfrak{m}^{\mathbf{a}+1} : x^{\mathbf{b}}) = \bigcap \mathfrak{m}^{\mathbf{a}+1-\mathbf{b}} = \bigcap (\mathfrak{m}^{\mathbf{b}^{\mathbf{a}}} + \mathfrak{m}^{\mathbf{a}+1}) = \left(\bigcap \mathfrak{m}^{\mathbf{b}^{\mathbf{a}}} \right) + \mathfrak{m}^{\mathbf{a}+1} = I + \mathfrak{m}^{\mathbf{a}+1}$$

since $(\mathbf{b}^{\mathbf{a}})^{\mathbf{a}} = \mathbf{b}$ for all $\mathbf{b} \preceq \mathbf{a}$. □

Remark 2.2 Using Corollary 2.14, below, this theorem provides a useful way to compute the Alexander dual ideal, given a set of generators. Indeed, the generators for $I^{\mathbf{a}}$ are simply those generators of $(\mathfrak{m}^{\mathbf{a}+1} : I)$ whose exponents are $\preceq \mathbf{a}$. Using Definition 1.5 (and Corollary 2.14 again), this can also be construed as a method for computing irreducible components of I given a generating set for I , or vice versa.

Denoting the \mathbb{Z}^n -graded Hom functor by $\underline{\text{Hom}}$, the next duality that comes into play is the k -dual $N^\wedge := \underline{\text{Hom}}_k(N, k)$, which is a \mathbb{Z}^n -graded S -module with the grading $(N^\wedge)_{\mathbf{c}} = \text{Hom}_k(N_{-\mathbf{c}}, k)$. It is a simple but very important observation that $T^\wedge \cong T$ as \mathbb{Z}^n -graded modules. This can be exploited: let $M \subseteq T$ be a submodule (the \mathbb{Z}^n -graded submodules of T are precisely the *monomial modules* of [4]). Taking the k -dual of the surjection $T \rightarrow T/M$ yields an injection $(T/M)^\wedge \rightarrow T^\wedge \cong T$. This makes $(T/M)^\wedge$ into a submodule of T which we call the T -dual of M and denote by M^T . If one thinks of the module M as a set of lattice points in \mathbb{Z}^n , then M^T can be thought of as the negatives of the lattice points in the complement of M ; i.e. we can make the equivalent

Definition 2.3 The T -dual M^T of a monomial module $M \subseteq T$ is defined by $x^{-\mathbf{b}} \in M^T \Leftrightarrow x^{\mathbf{b}} \notin M$.

The equivalence with the earlier formulation can be seen simply by examining which \mathbb{Z}^n -graded pieces of M and M^T have dimension 1 over k and which have dimension 0. Observe the striking similarity of Definition 2.3 with definition of the dual simplicial complex: $\overline{F} \in \Gamma^\vee \Leftrightarrow F \notin \Gamma$. Here are some properties of the T -dual which will be used later (possibly without explicit reference). Note the similarity of (i)–(iii) to the laws governing complements, unions, and intersections.

Proposition 2.4 Let M and N be submodules of T . Then

$$\begin{array}{ll} (i) & (M^T)^T = M \\ (ii) & M \subseteq N \Leftrightarrow N^T \subseteq M^T \\ (iii) & (M + N)^T = M^T \cap N^T \\ (iv) & M[\mathbf{a}]^T = M^T[-\mathbf{a}] \end{array} \quad \begin{array}{ll} (v) & T/M^T = M^\wedge \\ (vi) & (N/M)^\wedge = M^T/N^T \text{ if } M \subseteq N \\ (vii) & (N/N \cap M)^\wedge = M^T/M^T \cap N^T \end{array}$$

Proof: Statements (i)–(iv) follow from Definition 2.3, and (v) follows either from the definition and (i) or as a special case of (vi). To prove (vi) observe that $N/M = \ker(T/M \rightarrow T/N)$ so that $(N/M)^\wedge = \text{coker}((T/N)^\wedge \rightarrow (T/M)^\wedge)$ and use the definition of T -dual. Finally, (vii) is just (vi) and (iii) applied to $(N + M)/M = N/N \cap M$. □

Definition 2.5 Given a monomial ideal $I \subseteq S$ define the Čech hull of I in T :

$$\tilde{I} := \langle x^{\mathbf{b}} \mid \mathbf{b} \in \mathbb{Z}^n \text{ and } x^{\mathbf{b}^+} \in I \rangle,$$

where $\mathbf{b}^+ \in \mathbb{N}^n$ is, as usual, the join (componentwise maximum) of \mathbf{b} and $\mathbf{0}$ in the order lattice \mathbb{Z}^n .

Proposition 2.6 Taking the Čech hull commutes with finite intersections and sums. Furthermore,

- (i) \tilde{I} is the largest monomial submodule of T whose intersection with S is equal to I .
- (ii) \tilde{I} can be generated by (possibly infinitely many) monomials in T of degree $\preceq \mathbf{a}_I$.
- (iii) \tilde{I}^T is generated in degrees $\preceq \mathbf{0}$.

Proof: The first statement follows from (i) and the definitions.

(i) It is clear from the definition that \tilde{I} contains I ; and if $x^{\mathbf{b}} \in \tilde{I} \cap S$ then $\mathbf{b}^+ = \mathbf{b}$ whence $x^{\mathbf{b}} \in I$. Thus $\tilde{I} \cap S = I$. On the other hand, if M is a monomial submodule of T satisfying $M \cap S = I$ and $x^{\mathbf{b}} \in M$, then $x^{\mathbf{b}^-} \cdot x^{\mathbf{b}} = x^{\mathbf{b}^+} \in M \cap S = I$, where $\mathbf{b}^- := \mathbf{b}^+ - \mathbf{b}$. Thus $M \subseteq \tilde{I}$.

(ii) If $x^{\mathbf{b}} \in \tilde{I}$ then $\mathbf{c} \preceq \mathbf{b}^+$ for some minimal generator $x^{\mathbf{c}}$ of I , whence $x^{\mathbf{c}-\mathbf{b}^-}$ is in \tilde{I} , divides $x^{\mathbf{b}}$, and has exponent $\preceq \mathbf{a}_I$.

(iii) The following statement is precisely the T -dual to statement (i): \tilde{I}^T is the smallest submodule whose sum with $\tilde{\mathbf{m}}$ is equal to I^T . As $\tilde{\mathbf{m}}$ already contains all degrees $\not\preceq \mathbf{0}$, minimality of \tilde{I}^T implies that it is generated in degrees $\preceq \mathbf{0}$. \square

Example 2.7 (i) Recall that for $F \in \Delta$, the complement $\{1, \dots, n\} \setminus F$ is denoted by \bar{F} . Using this, the localization $S[x^{-\bar{F}}]$ is achieved by inverting the variables x_i for $i \notin F$. Now let $\mathbf{b} \succ \mathbf{0}$ and $F = \sqrt{\mathbf{b}}$. Then

$$\left(\widetilde{\mathbf{m}^{\mathbf{b}}}\right)^T = \left(S[x^{-\bar{F}}]\right)[\mathbf{b} - F].$$

To see this, first observe that if $\mathbf{c} \in \mathbb{N}^n$ then $x^{\mathbf{c}} \notin \mathbf{m}^{\mathbf{b}} \Leftrightarrow \mathbf{c} \cdot F \preceq \mathbf{b} - F$. Therefore, if $\mathbf{c} \in \mathbb{Z}^n$ then $x^{\mathbf{c}} \notin \widetilde{\mathbf{m}^{\mathbf{b}}} \Leftrightarrow \mathbf{c}^+ \preceq \mathbf{b} - F \Leftrightarrow \mathbf{c} \cdot F \preceq \mathbf{b} - F$. This last condition is equivalent to $-\mathbf{c} \cdot F \succeq F - \mathbf{b}$, and this occurs if and only if $x^{-\mathbf{c}} \in \left(S[x^{-\bar{F}}]\right)[\mathbf{b} - F]$.

(ii) For a special case, it follows that when $\mathbf{b} \succeq \mathbf{0}$, $\widetilde{\mathbf{m}^{\mathbf{b}+1}} = S[\mathbf{b}]^T$. \square

Remark 2.8 Example 2.7(ii) is the reason for the name *Čech hull*: when $\mathbf{b} = \mathbf{0}$, we find that $\tilde{\mathbf{m}}$ is the kernel of the last map in the Čech complex on x_1, \dots, x_n .

Definition 2.9 For any monomial ideal I and $\mathbf{a} \succeq \mathbf{a}_I$, define

$$I^{[\mathbf{a}]} := \tilde{I}^T[-\mathbf{a}] \cap S.$$

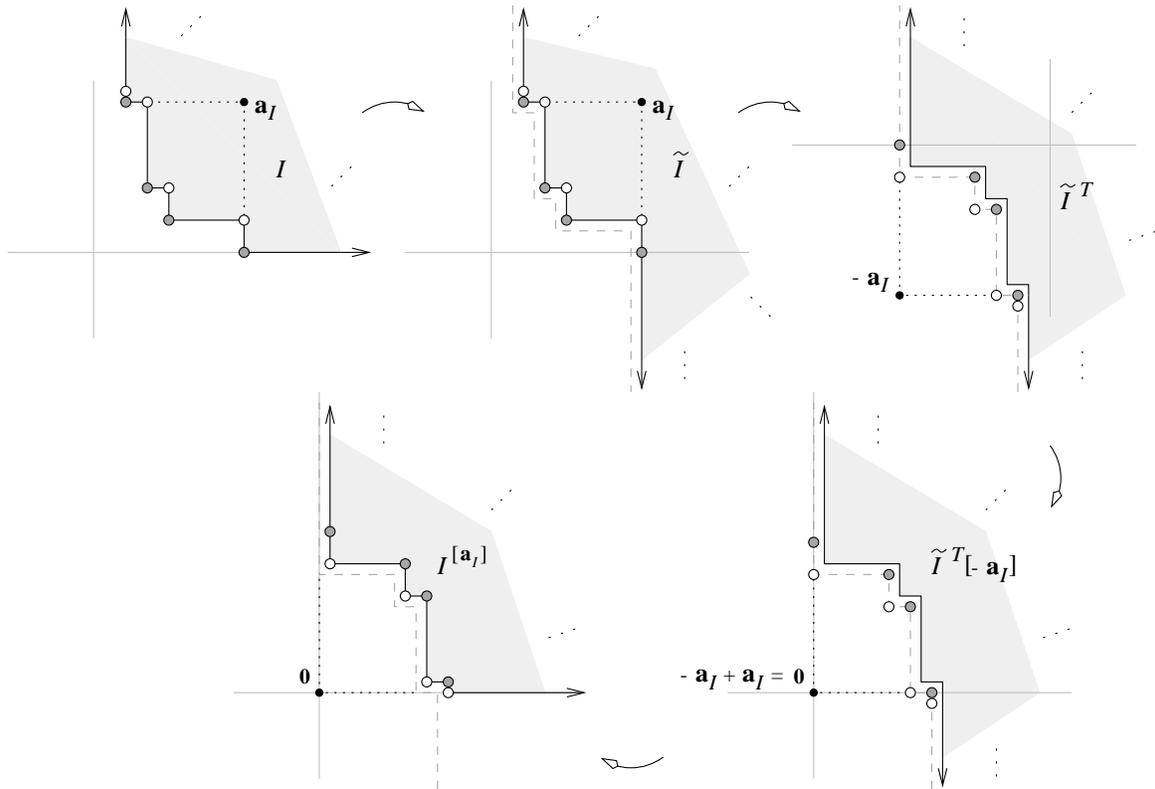


Figure 2

Example 2.10 Figure 2 is a schematic diagram depicting the transformation in stages from I to $I^{[a_I]}$. The black and white dots shift by $\mathbf{1}$ from the penultimate stage to the last; they are left in place (with respect to the dark black dot and the dark dotted lines) for the rest of the transformation. This shift is the reason for the $\mathbf{1}$ in the definition of \mathbf{b}^a , and it occurs because the flip-flop from \tilde{I} to \tilde{I}^T leaves a space of $\mathbf{1}$. The crux of this whole theory is that the “boundaries” of \tilde{I} and \tilde{I}^T have the same shape, but reversed, thus switching the roles of the black and white dots. This schematic may be helpful in parsing the proof of Theorem 2.13, below. \square

Lemma 2.11 $(I^{[a]})^\sim = \tilde{I}^T[-\mathbf{a}]$.

Proof: Let $M = \tilde{I}^T[-\mathbf{a}]$. By Proposition 2.6(i), $(M \cap S)^\sim \supseteq M$ since their intersections with S are equal by definition. Thus $((M \cap S)^\sim)^T \subseteq M^T$, with equality in degrees $\preceq \mathbf{0}$. But $M^T = \tilde{I}^{[a]}$ is generated in negative degrees by Proposition 2.6(ii), so that in fact $((M \cap S)^\sim)^T = M^T$. Taking T -duals of this equality gives the desired result. \square

The upshot is that I may be reconstructed from $I^{[a]}$ via the same construction which produces $I^{[a]}$ from I :

Proposition 2.12 $\mathbf{a}_{I^{[a]}} \preceq \mathbf{a}$ and $I = (I^{[a]})^{[a]}$.

Proof: By Proposition 2.6(iii) $\tilde{I}^T[-\mathbf{a}]$ is generated in degrees $\preceq \mathbf{a}$, so Lemma 2.11 implies that the same holds for $(I^{[a]})^\sim$. It is trivial to check that if any monomial module $M \subseteq T$ is generated in degrees $\preceq \mathbf{a}$ then so is $M \cap S$, because $\mathbf{a} \succeq \mathbf{0}$. Thus $\mathbf{a}_{I^{[a]}} \preceq \mathbf{a}$, and, in particular, $(I^{[a]})^{[a]}$ is well-defined. Now

$$\begin{aligned} (I^{[a]})^{[a]} &= ((I^{[a]})^\sim)^T[-\mathbf{a}] \cap S && \text{by definition} \\ &= (\tilde{I}^T[-\mathbf{a}])^T[-\mathbf{a}] \cap S && \text{by the previous lemma} \\ &= \tilde{I} \cap S && \text{by Proposition 2.4(iv) and (i)} \\ &= I. && \square \end{aligned}$$

The real cause for introducing $I^{[a]}$ is the next result, which may not be so unexpected at this point. It would seem that Theorem 2.13 makes the notation $I^{[a]}$ superfluous, and it does; nevertheless, the notation will be retained for emphasis, to indicate that Sections 3 and 4 (and, in particular, Theorem 4.10) are logically independent from Theorem 2.13.

Theorem 2.13 $I^{\mathbf{a}} = I^{[a]}$.

Proof: To simplify notation, declare that $\mathbf{b} \in \text{Irr}(I)$ if $\mathbf{m}^{\mathbf{b}} \in \text{Irr}(I)$. For each $\mathbf{b} \in \text{Irr}(I)$, let $S^{\mathbf{b}} = S[x^{-\sqrt{\mathbf{b}}}]$ be the localization of S at the prime $\mathbf{m}^{\sqrt{\mathbf{b}}}$. Then for each $\mathbf{b} \in \text{Irr}(I)$ and any $\mathbf{c} \in \mathbb{N}^n$ we have the following two facts:

- (i) $S^{\mathbf{b}}[-\mathbf{c}] \cong S^{\mathbf{b}}[-\mathbf{c} \cdot \sqrt{\mathbf{b}}]$ since multiplication by $x^{\mathbf{c} \cdot \sqrt{\mathbf{b}}}$ is a \mathbb{Z}^n -graded automorphism of $S^{\mathbf{b}}[-\mathbf{c}]$.
- (ii) $S \cap S^{\mathbf{b}}[-\mathbf{c} \cdot \sqrt{\mathbf{b}}] = S[-\mathbf{c} \cdot \sqrt{\mathbf{b}}]$. Indeed, this is equivalent to $(\langle x^{\mathbf{c} \cdot \sqrt{\mathbf{b}}} \rangle \cdot S^{\mathbf{b}}) \cap S = \langle x^{\mathbf{c} \cdot \sqrt{\mathbf{b}}} \rangle$, which holds because $\langle x^{\mathbf{c} \cdot \sqrt{\mathbf{b}}} \rangle \subseteq S$ is saturated with respect to $\langle x^{\sqrt{\mathbf{b}}} \rangle$; i.e. $(\langle x^{\mathbf{c} \cdot \sqrt{\mathbf{b}}} \rangle : x^{\sqrt{\mathbf{b}}}) = \langle x^{\mathbf{c} \cdot \sqrt{\mathbf{b}}} \rangle$.

Creating $I^{[a]}$ from I in stages yields

$$\begin{aligned} \tilde{I} &= \bigcap \tilde{\mathbf{m}}^{\mathbf{b}} && \text{by Proposition 2.6} \\ \Rightarrow \tilde{I}^T &= \sum (\tilde{\mathbf{m}}^{\mathbf{b}})^T && \text{by Proposition 2.4(iii)} \\ &= \sum S^{\mathbf{b}}[\mathbf{b} - \sqrt{\mathbf{b}}] && \text{by Example 2.7(i)} \\ \Rightarrow \tilde{I}^T[-\mathbf{a}] &= \sum S^{\mathbf{b}}[-\mathbf{b}^{\mathbf{a}}] && \text{by (i) above, with } \mathbf{c} = \mathbf{a} + \sqrt{\mathbf{b}} - \mathbf{b} \\ \Rightarrow S \cap \tilde{I}^T[-\mathbf{a}] &= \sum S[-\mathbf{b}^{\mathbf{a}}] && \text{by (ii) above, with } \mathbf{c} = \mathbf{b}^{\mathbf{a}} \end{aligned}$$

where the intersection and all of the summations are taken over all $\mathbf{b} \in \text{Irr}(I)$. The last summation above is equal to $I^{\mathbf{a}}$ since each summand $S[-\mathbf{b}^{\mathbf{a}}]$ is just a principal ideal $\langle x^{\mathbf{b}^{\mathbf{a}}} \rangle$. \square

Corollary 2.14 $(I^{\mathbf{a}})^{\mathbf{a}} = I$. Furthermore, $(\mathbf{b}^{\mathbf{a}})^{\mathbf{a}} = \mathbf{b}$, so that

$$I^{\mathbf{a}} = \bigcap \{ \mathbf{m}^{\mathbf{b}^{\mathbf{a}}} \mid x^{\mathbf{b}} \text{ is a minimal generator of } I \}. \quad \square$$

Remark 2.15 In general, one has $(I^\vee)^\vee \neq I$. However, in the special case when $I = \sqrt{I}$, it will always happen that $(I^\vee)^\vee = I$. This follows from an application of Corollary 2.14 to Remark 1.6(ii). The difference $\mathbf{a}_I - \mathbf{a}_{I^\vee}$ measures the extent to which $(I^\vee)^\vee \neq I$ fails, in the sense that $(I^\vee)^\vee = I[\mathbf{a}_I - \mathbf{a}_{I^\vee}] \cap S$. However $((I^\vee)^\vee)^\vee = I^\vee$, so that an ideal which is already an Alexander dual is maximal in some sense. It is unclear what the invariant $\mathbf{a}_I - \mathbf{a}_{I^\vee}$ means, in general, although the passage from I to $(I^\vee)^\vee$ can sometimes be thought of as a “tightening” that may resolve some amount of nonminimality in the hull resolution of [4]—see Example 5.27. See also Remark 5.9(ii) below for another occurrence of the invariant $\mathbf{a}_I - \mathbf{a}_{I^\vee}$.

The reader interested in cellular resolutions may wish to skip directly to Section 5, whose only logical dependence on Sections 3 and 4 is Proposition 3.11.

3 Bass numbers versus Betti numbers

Algebraically, Alexander duality is best expressed in terms of relations between Betti and Bass numbers (Definition 3.1), as evidenced by this section and the next. The principle behind this is that the T -duality of Section 2, which can be thought of as lattice duality in \mathbb{Z}^n , can also be interpreted (Corollary 3.6) as a manifestation of the same process that interchanges flat and injective modules (in the appropriate category). In Theorem 3.12 this results in equalities between Bass and Betti numbers of I . Though perhaps not so interesting a statement in its own right, Proposition 3.11 is the workhorse for the remainder of the paper—it is the *reason* everything else is true. It encapsulates simultaneously the relations between all of the dualities that enter into this paper: k - and T -duality, Alexander duality, linkage, local duality, and Matlis duality.

Definition 3.1 *The derived functors of the \mathbb{Z}^n -graded functor $\underline{\text{Hom}}$ will be called $\underline{\text{Ext}}$, and the left derived functor of \otimes , which is also \mathbb{Z}^n -graded, will be called $\underline{\text{Tor}}$. For a module N define*

$$\begin{aligned}\mu_{i,\mathbf{b}}(N) &= \dim_k \left(\underline{\text{Ext}}_S^i(k, N)_{\mathbf{b}} \right) \\ \beta_{i,\mathbf{b}}(N) &= \dim_k \left(\underline{\text{Tor}}_i^S(k, N)_{\mathbf{b}} \right),\end{aligned}$$

the i^{th} Bass and Betti numbers of N in degree \mathbf{b} .

Remark 3.2 (i) In order to compute these derived functors in the category \mathcal{M} of \mathbb{Z}^n -graded S -modules (see Proposition 3.3), we need to know that \mathcal{M} has enough injective and projective modules, just as in the nongraded case. Of course, there are always free modules, so this takes care of the projectives; for injectives one can easily modify the proof of [5], Theorem 3.6.2 to fit the \mathbb{Z}^n -graded case.

(ii) If M is finitely generated then $\underline{\text{Ext}}^i(M, -) = \text{Ext}^i(M, -)$. In particular, summing the Betti or Bass numbers over all \mathbf{b} (or all \mathbf{b} with fixed \mathbb{Z} -degree) gives the same result as computing directly in the nongraded (or \mathbb{Z} -graded) case.

In what follows, we will need the notion of a *flat resolution in \mathcal{M}* . This is defined exactly like a free resolution, except that the resolving modules are required to be flat instead of free, where *flat* means acyclic for $\underline{\text{Tor}}$ [18], Section 2.4. Recall that free and flat are equivalent for finitely generated S -modules; this is a simple consequence of the grading and Nakayama's lemma. However, non-finitely generated flat modules, such as localizations of S , may fail to be free, or even projective.

Proposition 3.3 (i) $\underline{\text{Ext}}^i(M, N)$ can be calculated as the homology of the complexes obtained either by applying $\underline{\text{Hom}}(-, N)$ to a projective resolution of M in \mathcal{M} or by applying $\underline{\text{Hom}}(M, -)$ to an injective resolution of N in \mathcal{M} .

(ii) $\underline{\text{Tor}}_i(M, N)$ can be calculated as the homology of the complexes obtained by either tensoring with N a flat resolution of M in \mathcal{M} or by tensoring with M a flat resolution of N in \mathcal{M} .

Proof: (i) Remark 3.2(i) above provides enough injectives to use [18], Definition 2.5.1, Example 2.5.3, and Exercise 2.7.4.

(ii) [18], Theorem 2.7.2 and Exercise 2.4.3. □

Lemma 3.4 $N^\wedge = \underline{\text{Hom}}_S(N, S^\wedge)$.

Proof: [5], Proposition 3.6.16(c), whose proof holds just as easily in the \mathbb{Z}^n -graded case. □

The next theorem is the starting point for the comparison of Betti and Bass numbers. Its corollary, which carries out the lattice complementation, is fundamental to the rest of the results in this section.

Proposition 3.5 For any module N , $\mu_{i,\mathbf{b}}(N) = \beta_{i,-\mathbf{b}}(N^\wedge)$.

Proof: A module J is injective if and only if J^\wedge is flat, because

$$(1) \quad \underline{\text{Hom}}(-, J) = \underline{\text{Hom}}(-, \underline{\text{Hom}}(J^\wedge, S^\wedge)) = \underline{\text{Hom}}(- \otimes J^\wedge, S^\wedge).$$

Indeed, the first term being an exact functor means that J is injective, while the last term being an exact functor means that J^\wedge is flat, since $\underline{\text{Hom}}(-, S^\wedge)$ is *a priori* a faithful exact functor. It follows that a complex $J : 0 \rightarrow J^0 \rightarrow J^1 \rightarrow \dots$ is an injective resolution of N in \mathcal{M} if and only if $(J)^\wedge$ is a flat resolution of N^\wedge . Substituting k for $(-)$ in Equation (1) and applying Proposition 3.3 we get

$$(2) \quad \underline{\text{Ext}}^i(k, N) \cong \underline{\text{Tor}}_i(k, N^\wedge)^\wedge$$

from which the result follows at once. □

Corollary 3.6 $\mu_{i,\mathbf{b}}(T/M) = \beta_{i,-\mathbf{b}}(M^T)$ for any monomial module $M \subseteq T$. □

The next few results are preliminary to the theorems relating the Betti numbers of I to the Bass numbers of I (Theorem 3.12) and the Bass numbers of $I^{[a]}$ (Theorem 4.10).

Proposition 3.7 *Let I be an ideal. Then*

- (i) $\beta_{i,\mathbf{b}}(\tilde{I}) = 0$ if $\mathbf{b} \not\preceq \mathbf{a}_I$.
- (ii) $\beta_{i,\mathbf{b}}(\tilde{I}) = 0$ if $\mathbf{b} \not\preceq \mathbf{1}$.
- (iii) $\beta_{i,\mathbf{b}}(\tilde{I}) = \beta_{i,\mathbf{b}}(I)$ if $\mathbf{1} \preceq \mathbf{b}$.

Proof: Given any submodule $M \subseteq T$, define for each $\mathbf{b} \in \mathbb{Z}^n$ the following simplicial subcomplex of Δ :

$$K_{\mathbf{b}}(M) = \{F \in \Delta \mid x^{\mathbf{b}-F} \in M\}.$$

It is a result of [9] and [12] (and extended to the case $M \subseteq T$ by [4]) that

$$\beta_{i,\mathbf{b}}(M) = \dim_k \tilde{H}_i(K_{\mathbf{b}}(M); k),$$

the dimension of the i^{th} simplicial homology of $K_{\mathbf{b}}(M)$ with coefficients in k . To prove (i) and (ii) it suffices to show that $K_{\mathbf{b}}(\tilde{I})$ is a cone (and therefore acyclic) unless $\mathbf{1} \preceq \mathbf{b} \preceq \mathbf{a}_I$. If $\mathbf{a}_I = (a_1, \dots, a_n)$ and $b_i \geq a_i + 1$, then it follows from Proposition 2.6(ii) that $K_{\mathbf{b}}(\tilde{I})$ is a cone with vertex $\{i\}$, proving (i). That $K_{\mathbf{b}}(\tilde{I})$ is a cone with vertex $\{i\}$ if $b_i \leq 0$ follows directly from the definition of Čech hull, proving (ii). Finally, (iii) holds because $K_{\mathbf{b}}(\tilde{I}) = K_{\mathbf{b}}(I)$ whenever $\mathbf{b} \succeq \mathbf{1}$. \square

Lemma 3.8 *Let $M \subseteq T$. Then $\beta_{i,\mathbf{b}}(M) = \beta_{i,\mathbf{b}}(M/M \cap \widetilde{\mathfrak{m}^{\mathbf{a}+1}})$ if $\mathbf{b} \preceq \mathbf{a}$.*

Proof: It follows from Example 2.7(ii) that $(M \cap \widetilde{\mathfrak{m}^{\mathbf{a}+1}})_{\mathbf{b}} = 0$ if $\mathbf{b} \preceq \mathbf{a}$, so the Taylor resolution of it (see [16] for the original or [4], Proposition 1.5 for a treatment including submodules of T) forces $\beta_{i,\mathbf{b}}(M \cap \widetilde{\mathfrak{m}^{\mathbf{a}+1}}) = 0$ for all $\mathbf{b} \preceq \mathbf{a}$. Applying $\underline{\text{Tor}}$ to the exact sequence

$$0 \rightarrow M \cap \widetilde{\mathfrak{m}^{\mathbf{a}+1}} \rightarrow M \rightarrow M/M \cap \widetilde{\mathfrak{m}^{\mathbf{a}+1}}$$

yields the result. \square

Lemma 3.9 *If $i < n$ then $\underline{\text{Ext}}^i(k, S/I) \cong \underline{\text{Ext}}^i(k, T/I)$, and in the remaining case $i = n$ we have $\underline{\text{Ext}}^n(k, S/I) = k[\mathbf{1}]$.*

Proof: One can first calculate $\underline{\text{Ext}}^i(k, S) = \begin{cases} k[\mathbf{1}] & \text{if } i = n \\ 0 & \text{otherwise} \end{cases}$ from the Koszul complex and $\underline{\text{Ext}}^i(k, T) = 0$ for all i because T is injective in the category \mathcal{M} . Using the long exact sequence of $\underline{\text{Ext}}$ from $0 \rightarrow S \rightarrow T \rightarrow T/S \rightarrow 0$ one finds that $\underline{\text{Ext}}^i(k, S) \cong \underline{\text{Ext}}^{i-1}(k, T/S)$.

From the above calculations and the long exact sequence of $\underline{\text{Ext}}$ arising from

$$0 \rightarrow S/I \rightarrow T/I \rightarrow T/S \rightarrow 0$$

the lemma will follow if we can show that the map

$$\underline{\text{Ext}}^{n-1}(k, T/S) \rightarrow \underline{\text{Ext}}^n(k, S/I)$$

is an isomorphism. But S is a regular ring, so $\underline{\text{Ext}}^n(k, S/I)$ is nonzero *a priori* because of [5], Proposition 3.1.14 and [5], Theorem 3.1.17, so it is enough to prove that the 1-dimensional vector space $\underline{\text{Ext}}^{n-1}(k, T/S) \cong \underline{\text{Ext}}^n(k, S) \cong k[\mathbf{1}]$ maps surjectively, i.e. that $\underline{\text{Ext}}^n(k, T/I) = 0$. Now $\underline{\text{Ext}}^n(k, T/I) \cong \underline{\text{Ext}}^{n+1}(k, I)$ because of the exact sequence

$$0 \rightarrow I \rightarrow T \rightarrow T/I \rightarrow 0,$$

and $\underline{\text{Ext}}^{n+1}(k, I) = 0$ because of the same [5] reference as above. \square

The next main result, Theorem 3.12, is really a rephrasing of an observation made in the proof of [9], Theorem 5.2. While it is possible, by quoting the self-duality of the Koszul complex, to extend the result to include all S -modules, the proof here demonstrates effectively the interaction of Alexander duality with other kinds of duality. Aside from the intrinsic interest in its proof, Theorem 3.12 will find an application in the proof of Theorem 4.10. Two preliminary results are needed, the first of which will also be used in the proof of Proposition 4.6.

Lemma 3.10 *With $J = I + \mathfrak{m}^{a+1}$ we have $\tilde{J}^T = I^{[a]}[\mathbf{a}]$. The same is true if I and $I^{[a]}$ are reversed.*

Proof: The last statement is because of Proposition 2.12. By Example 2.7(ii) and Proposition 2.4,

$$I^{[a]}[\mathbf{a}] = \tilde{I}^T \cap S[\mathbf{a}] = (\tilde{I} + S[\mathbf{a}]^T)^T = \tilde{J}^T. \quad \square$$

The reader knowledgeable about linkage will recognize a hint of Theorem 2.1 in the next proposition. Only the special case $\mathbf{b} = \mathbf{0}$ is required in this section. However, the more general result is a major component in the proof of Theorem 6.11.

Proposition 3.11 *Let $\mathbf{a} \succeq \mathbf{a}_I$, $J = I^{[a]} + \mathfrak{m}^{a+1}$, and $\mathbf{b} \in \mathbb{N}^n$. Then*

$$\underline{\text{Ext}}_S^n(S[\mathbf{b}]/S[\mathbf{b}] \cap \tilde{J}, S) = (I/I \cap \mathfrak{m}^{a+b+1})[\mathbf{a} + \mathbf{1}].$$

In particular, taking $\mathbf{b} = \mathbf{0}$ yields $\underline{\text{Ext}}^n(S/J, S) \cong (I/I \cap \mathfrak{m}^{a+1})[\mathbf{a} + \mathbf{1}]$.

Proof: The module T/\tilde{J} is the k -dual of the finitely generated module $I[\mathbf{a}]$ by Lemma 3.10, and is hence artinian by Matlis duality, cf. [7], Theorem 2.1.4. Thus our module $S[\mathbf{b}]/S[\mathbf{b}] \cap \tilde{J} \subseteq T/\tilde{J}$ is also artinian, and (obviously) finitely generated, as well. Since the canonical module of S is $S[-\mathbf{1}]$ by [7], Corollary 2.2.6, local duality (in the form of [7], Theorem 2.2.2) applied to the zeroeth local

cohomology module implies the first equality of the following:

$$\begin{aligned}
\underline{\text{Ext}}_S^n(S[\mathbf{b}]/S[\mathbf{b}] \cap \tilde{J}, S) &= (S[\mathbf{b}]/S[\mathbf{b}] \cap \tilde{J})^\wedge[\mathbf{1}] && \text{by local duality} \\
&= (\tilde{J}^T / \tilde{J}^T \cap S[\mathbf{b}]^T)[\mathbf{1}] && \text{by Proposition 2.4(vii)} \\
&= (I/I \cap S[\mathbf{a} + \mathbf{b}]^T)[\mathbf{a} + \mathbf{1}] && \text{by Lemma 3.10 and shifting by } [-\mathbf{a}][\mathbf{a}] \\
&= (I/I \cap \mathfrak{m}^{\mathbf{a}+\mathbf{b}+1})[\mathbf{a} + \mathbf{1}] && \text{by Example 2.7(ii).} \quad \square
\end{aligned}$$

Given an artinian ideal J , the list of Betti numbers for the canonical module $\underline{\text{Ext}}^n(S/J, S[-\mathbf{1}])$ of S/J is essentially the reverse of the list of Betti numbers for J ; see, for instance, [5], Corollary 3.3.9. On the other hand, there is the lattice-complementation view of Alexander duality, which emerges in Corollary 3.6 as a relation between the Betti numbers of a monomial module and the Bass numbers of its T -dual. These two dualities can be combined to compare the Betti numbers of I to the Bass numbers of the same ideal I :

Theorem 3.12 For all $i \in \mathbb{Z}$ and $\mathbf{b} \in \mathbb{Z}^n$,

$$\beta_{n-i, \mathbf{b}}(S/I) = \mu_{i, \mathbf{b}-\mathbf{1}}(S/I).$$

Proof: The case $i = n$ follows from the calculations of Lemma 3.9, so assume from now on that $i \leq n - 1$. In particular, we can calculate the Bass numbers from T/I instead of S/I by Lemma 3.9. Let $\mathbf{a} = \mathbf{a}_I + \mathbf{1}$. All of the Betti numbers are zero in degrees $\mathbf{b} \not\leq \mathbf{a}$ by Proposition 3.7(i) and (iii). As for the Bass numbers, we can use the fact that, with $J := I^{[\mathbf{a}]} + \mathfrak{m}^{\mathbf{a}+1}$, we have $I^T = \tilde{J}[\mathbf{a}]$ by Lemma 3.10. It follows that $\mu_{i, \mathbf{b}-\mathbf{1}}(T/I) = \beta_{i, \mathbf{1}-\mathbf{b}}(\tilde{J}[\mathbf{a}]) = \beta_{i, \mathbf{a}+1-\mathbf{b}}(\tilde{J})$ by Corollary 3.6, and then Proposition 3.7(ii) implies that these numbers are zero if $\mathbf{b} \not\leq \mathbf{a}$.

From now on, assume $\mathbf{b} \leq \mathbf{a}$ and $0 \leq i \leq n - 1$. Let $J = I^{[\mathbf{a}]} + \mathfrak{m}^{\mathbf{a}+1}$ and calculate

$$\begin{aligned}
\mu_{i, \mathbf{b}-\mathbf{1}}(S/I) &= \mu_{i, \mathbf{b}-\mathbf{1}}(T/I) && \text{by Lemma 3.9 and } i \leq n - 1 \\
&= \beta_{i, \mathbf{a}+1-\mathbf{b}}(\tilde{J}) && \text{by Corollary 3.6, since } I^T = \tilde{J}[\mathbf{a}] \\
&= \beta_{i, \mathbf{a}+1-\mathbf{b}}(J) && \text{by Proposition 3.7(iii) and } \mathbf{b} \leq \mathbf{a} \\
&= \beta_{i+1, \mathbf{a}+1-\mathbf{b}}(S/J) && \text{since } i \geq 0 \\
&= \beta_{n-i-1, \mathbf{b}-\mathbf{1}-\mathbf{a}}(\underline{\text{Ext}}^n(S/J, S)) && \text{since } S \text{ is Gorenstein and } S/J \text{ is artinian} \\
&= \beta_{n-i-1, \mathbf{b}}((I/I \cap \mathfrak{m}^{\mathbf{a}+1})) && \text{by Proposition 3.11} \\
&= \beta_{n-i-1, \mathbf{b}}(I) && \text{by Lemma 3.8 and } \mathbf{b} \leq \mathbf{a} \\
&= \beta_{n-i, \mathbf{b}}(S/I) && \text{since } i \leq n - 1
\end{aligned}$$

proving the theorem. □

4 Localization and restriction

This section aims to reveal the equality (Theorem 4.10) between Betti numbers of I and localized Bass numbers (Definition 4.8) of $I^{[\mathbf{a}]}$. This equality generalizes Theorem 2.13. As a consequence of the equality, an inequality between Betti numbers of I and $I^{[\mathbf{a}]}$ is obtained in Theorem 4.13, generalizing to arbitrary monomial ideals an inequality of [2] which was proven for radical ideals.

The next proposition should be thought of as the nonlocalized precursor to Theorem 4.10(i).

Proposition 4.1 *Let I be an ideal and $\mathbf{a} \succeq \mathbf{a}_I$. If $\mathbf{1} \preceq \mathbf{b} \preceq \mathbf{a}$ then $\beta_{i,\mathbf{b}}(I) = \mu_{i,\mathbf{b}\mathbf{a}-\mathbf{1}}(S/I^{[\mathbf{a}]})$.*

Proof: If, to start with, $\mathbf{b} \preceq \mathbf{a}$, then

$$\begin{aligned} \beta_{i,\mathbf{b}}(M) &= \mu_{i,-\mathbf{b}}\left((M/M \cap \widetilde{\mathbf{m}^{\mathbf{a}+\mathbf{1}}})^\wedge\right) \text{ by Lemma 3.8 and Proposition 3.5} \\ &= \mu_{i,-\mathbf{b}}\left(S[\mathbf{a}]/S[\mathbf{a}] \cap M^T\right) \text{ by Proposition 2.4(vii) and Example 2.7(ii)} \\ &= \mu_{i,\mathbf{a}-\mathbf{b}}\left(S/M^T[-\mathbf{a}] \cap S\right). \end{aligned}$$

Substituting $M = \widetilde{I}$ we get $\beta_{i,\mathbf{b}}(\widetilde{I}) = \mu_{i,\mathbf{a}-\mathbf{b}}(S/I^{[\mathbf{a}]})$ if $\mathbf{b} \preceq \mathbf{a}$, and when the assumption $\mathbf{1} \preceq \mathbf{b}$ is added, $\mathbf{a} - \mathbf{b} = \mathbf{b}^{\mathbf{a}} - \mathbf{1}$ and the result is a consequence of Proposition 3.7(iii). \square

Theorem 4.10 is the combination of the previous proposition with localization and restriction of scalars. The following definitions will provide concise notation for these operations, which will be needed also for the definition of Bass numbers at primes other than \mathbf{m} (Definition 4.8). Recall that $\overline{F} = \{1, \dots, n\} \setminus F = \mathbf{1} - F$.

Definition 4.2 *Let Δ be the $(n-1)$ -simplex on the vertices $\{1, \dots, n\}$ and $F \in \Delta$. Define*

$$\begin{aligned} (i) \quad N(-\overline{F}) &:= S[x^{-\overline{F}}] \otimes_S N \text{ for arbitrary modules } N \\ (ii) \quad S_{[F]} &:= k[x_i \mid i \in F] \text{ a } \mathbb{Z}^F\text{-graded } k\text{-subalgebra of } S \\ (iii) \quad N_{[F]} &:= \bigoplus_{\mathbf{b} \in \mathbb{Z}^F} N_{\mathbf{b}} \text{ a } \mathbb{Z}^F\text{-graded } S_{[F]}\text{-module} \\ (iv) \quad N_{(F)} &:= N(-\overline{F})_{[F]} \end{aligned}$$

The operations on N listed above are all exact and commute with sums. They should be thought of as: (i) homogeneous localization at \mathbf{m}^F , (iii) taking the “degree zero part” of N with respect to F , and (iv) taking the “degree zero part of the homogeneous localization at \mathbf{m}^F ” as in algebraic geometry. In (ii) and (iii), the copy of \mathbb{Z}^F may be thought of as sitting inside \mathbb{Z}^n in the obvious way: as the space spanned by the basis vectors in the support of F . Thinking of \mathbb{Z}^F this way can cause notational problems, however. For instance, any \mathbb{Z}^n -graded S -module N can be thought of as a \mathbb{Z}^F -graded $S_{[F]}$ -module which in degree $\mathbf{b} \in \mathbb{Z}^F$ is

$$\bigoplus_{\mathbf{c}^F = \mathbf{b}} N_{\mathbf{c}} = \bigoplus_{\mathbf{b}' \in \mathbb{Z}^F} N_{\mathbf{b} + \mathbf{b}'},$$

where $\mathbf{c} \cdot F$ denotes the restriction to F as in Section 1. Note that the right-hand side gives this vector space the structure of a $\mathbb{Z}^{\bar{F}}$ -graded $S_{[\bar{F}]}$ -module. The convention will be the following:

If N is a \mathbb{Z}^F -graded $S_{[F]}$ -module and $\mathbf{b} \in \mathbb{Z}^F$, the graded piece of N in degree \mathbf{b} will be denoted $N_{\mathbf{b},F}$. That way, if N happens also to be a \mathbb{Z}^n -graded S -module, the usual notation $N_{\mathbf{b}}$ can retain its old meaning as the degree \mathbf{b} part in the \mathbb{Z}^n -grading.

Even if $\mathbf{b} \notin \mathbb{Z}^F$ it will sometimes be convenient to use $N_{\mathbf{b},F}$ to denote the $\mathbf{b} \cdot F$ graded piece in the \mathbb{Z}^F -grading; i.e. with $\mathbf{c} = \mathbf{b} \cdot F \in \mathbb{Z}^F$, we set $N_{\mathbf{b},F} := N_{\mathbf{c},F}$. The next Lemma follows from the definitions and the convention above. In each of (i)–(v), the objects are \mathbb{Z}^F -graded $S_{[F]}$ -modules, but in (i), the objects may also be considered as $\mathbb{Z}^{\bar{F}}$ -graded $S_{[\bar{F}]}$ -modules or even $\mathbb{Z}^{\bar{F}} \times \mathbb{Z}^F = \mathbb{Z}^n$ -graded $S_{[\bar{F}]} \otimes_k S_{[F]} = S$ -modules.

Lemma 4.3 For any $F \in \Delta$,

- (i) $M(-\bar{F}) = T_{[\bar{F}]} \otimes_k M_{(F)} = S(-\bar{F}) \otimes_{S_{(F)}} M_{(F)}$
- (ii) $M_{[F]} = M_{\mathbf{0},\bar{F}}$
- (iii) $M[\mathbf{a}]_{[F]} = M_{\mathbf{a},\bar{F}}[\mathbf{a} \cdot F]$
- (iv) $(\widetilde{I})_{[F]} = \widetilde{I}_{[F]}$
- (v) $(M^T)_{[F]} = (M_{[F]})^{T_{[F]}}$

where the right-hand sides of (iv) and (v) are, respectively, the Čech hull and T -dual over $S_{[F]}$. \square

For submodules $M \subseteq T$ the various gradings allow for convenient characterizations of localization as in Definition 4.2(iv). They use the fact that for any $\mathbf{b} \in \mathbb{Z}^n$, $M_{\mathbf{b},\bar{F}}$ is naturally a submodule of $T_{[F]} = T_{(F)}$.

Proposition 4.4 Let M be a monomial module.

$$(i) \quad M_{(F)} = \bigcup_{\mathbf{b} \in \mathbb{Z}^{\bar{F}}} M_{\mathbf{b},\bar{F}}.$$

If M can be generated in degrees \mathbf{c} satisfying $\mathbf{c} \cdot \bar{F} \preceq \mathbf{a} \cdot \bar{F}$ then

$$(ii) \quad M_{(F)} = M_{\mathbf{a},\bar{F}}.$$

Proof: (i) Observe that $M \subseteq M(-\bar{F}) \subseteq T$ because everything is torsion-free. Thus, if $\mathbf{b} \in \mathbb{Z}^{\bar{F}}$, then multiplication by $x^{-\mathbf{b}}$ induces an inclusion $M_{\mathbf{b},\bar{F}} \rightarrow M_{(F)}$. For the other inclusion, note that any monomial in $M_{(F)}$ can be written as $x^{\mathbf{b}} \cdot x^{\mathbf{c}}$ for some $x^{\mathbf{c}} \in M$ and $\mathbf{b} = -(\mathbf{c} \cdot \bar{F}) \in \mathbb{Z}^{\bar{F}}$.

(ii) The collection $\{M_{\mathbf{b},\bar{F}}\}_{\mathbf{b} \in \mathbb{Z}^{\bar{F}}}$ of $S_{[F]}$ -submodules of $T_{[F]}$ is partially ordered by inclusion because M is a module. The union in (i) stabilizes after $\mathbf{a} \cdot \bar{F}$ if M is generated in degrees \mathbf{c} satisfying $\mathbf{c} \cdot \bar{F} \preceq \mathbf{a} \cdot \bar{F}$. \square

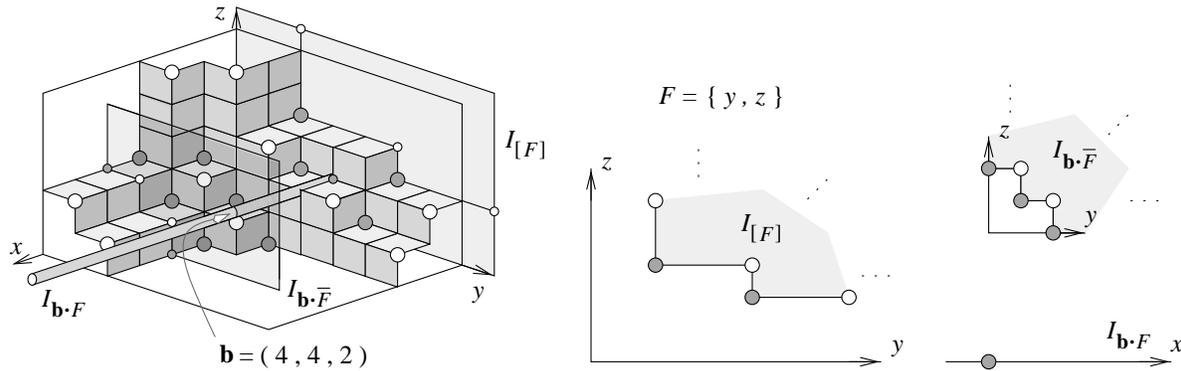


Figure 3

Example 4.5 Figure 3 illustrates some parts of Definition 4.2 and Lemma 4.3 in a specific case. For notation, $x, y,$ and z are identified with $1, 2,$ and $3 \in \{1, \dots, 3\} = \Delta$. The face F is $\{y, z\} = \{2, 3\}$, while $\mathbf{b} = (4, 4, 2)$. The small colored dots represent generators or irreducible components in the restricted ideals. It is not true that $\mathbf{b} \succeq \mathbf{a}_I$, so Proposition 4.4 does not apply; nevertheless, $I_{\mathbf{b} \cdot \bar{F}} = I_{(F)}$ for these $\mathbf{b}, I,$ and F . Figure 3 can also be used as a test case for Proposition 4.6.

Proposition 4.6 $(I^{[\mathbf{a}]})_{(F)} = (I_{[F]})^{[\mathbf{a} \cdot F]}$ as ideals in $S_{(F)} = S_{[F]}$. In words, dualizing and then localizing is the same as restricting and then dualizing.

Proof: It is enough to show that $(I^{[\mathbf{a}]})_{(F)}[\mathbf{a} \cdot F] = (I_{[F]})^{[\mathbf{a} \cdot F]}[\mathbf{a} \cdot F]$. Now

$$\begin{aligned} (I^{[\mathbf{a}]})_{(F)}[\mathbf{a} \cdot F] &= (I^{[\mathbf{a}]})_{\mathbf{a} \cdot \bar{F}}[\mathbf{a} \cdot F] && \text{by Proposition 4.4(ii) and Proposition 2.12} \\ &= (I^{[\mathbf{a}]}[\mathbf{a}])_{[F]} && \text{by Lemma 4.3(iii)} \\ &= \left(\left((I + \mathfrak{m}^{\mathbf{a}+1})^\sim \right)^T \right)_{[F]} && \text{by Lemma 3.10,} \end{aligned}$$

and one can use the rules 4.3(v) and then 4.3(iv) for interchanging the various operations to get the last line to equal

$$\left((I_{[F]} + \mathfrak{m}^{\mathbf{a} \cdot F + F}_{[F]})^\sim \right)^{T_{[F]}}$$

where $(-)^{T_{[F]}}$ is T -duality over $S_{[F]}$ as in Lemma 4.3(v). Another application of Lemma 3.10 (over $S_{[F]}$ this time) gives the desired result. \square

Proposition 4.7 Let $I \subseteq S$ and $\mathbf{b} \in \mathbb{Z}^F$. Then $\beta_{i, \mathbf{b}}(I) = \beta_{i, \mathbf{b} \cdot F}(I_{[F]})$.

Proof: Let \mathbb{F} be the Taylor resolution of I (see the proof of Lemma 3.8 for references). Then $\mathbb{F}_{[F]}$ is the Taylor resolution of $I_{[F]}$. Furthermore, $(k \otimes_S \mathbb{F})_{[F]} = k \otimes_{S_{[F]}} \mathbb{F}_{[F]}$ because if $\mathbf{b} \in \mathbb{N}^n$ then

$$\left(k \otimes_S S[-\mathbf{b}] \right)_{[F]} = k \otimes_{S_{[F]}} S[-\mathbf{b}]_{[F]} = \begin{cases} k[-\mathbf{b}] & \text{if } \mathbf{b} \in \mathbb{Z}^F \\ 0 & \text{if } \mathbf{b} \notin \mathbb{Z}^F \end{cases} .$$

Thus the Betti numbers in question are calculated from the same complex of k -vector spaces. \square

Definition 4.8 (Bass numbers for arbitrary monomial primes) *Given a module N and a degree $\mathbf{b} \in \mathbb{Z}^F$, the i^{th} Bass number of N with respect to F (or the prime ideal \mathfrak{m}^F) in degree \mathbf{b} is defined as*

$$\mu_{i,\mathbf{b}}(F, N) := \dim_k \left(\underline{\text{Ext}}_{S(F)}^i(k, N_{(F)\mathbf{b}}) \right).$$

Remark 4.9 When $F = \mathbf{1}$ this definition agrees with the Bass numbers of Definition 3.1.

Now comes the main result of this section. It can be thought of as a far-reaching generalization of Theorem 2.13, which is a special case, pending the appropriate interpretation of Bass numbers—see Proposition 4.12 and the second proof of Theorem 2.13 that follows it. In part (i) of the next theorem, the case where \mathbf{b} has full support is Proposition 4.1.

Theorem 4.10 *If $\mathbf{0} \neq F \preceq \mathbf{b} \preceq \mathbf{a} \cdot F$ then for all $i \in \mathbb{Z}$ we have*

$$\begin{aligned} (i) \quad & \beta_{i,\mathbf{b}}(I) = \mu_{i,\mathbf{b}\mathbf{a}-F}(F, S/I^{[\mathbf{a}]}) \\ (ii) \quad & \mu_{n-i-1,\mathbf{b}-\mathbf{1}}(S/I) = \mu_{i,\mathbf{b}\mathbf{a}-F}(F, S/I^{[\mathbf{a}]}) \\ (iii) \quad & \beta_{i,\mathbf{b}}(I) = \beta_{|F|-i-1,\mathbf{b}\mathbf{a}}(I^{[\mathbf{a}]_{(F)}}). \end{aligned}$$

In any of these formulas, I and $I^{[\mathbf{a}]}$ can be switched, and the same goes for \mathbf{b} and $\mathbf{b}^{\mathbf{a}}$.

Proof: Statements (ii) and (iii) follow easily from (i), in view of Theorem 3.12 and the fact that $\beta_{i,\mathbf{b}}(I) = \beta_{i+1,\mathbf{b}}(S/I)$ when $\mathbf{b} \neq \mathbf{0}$. To prove (i), note that $\mathbf{b}^{\mathbf{a}} = (\mathbf{b} \cdot F)^{\mathbf{a}\cdot F}$, so

$$\begin{aligned} \beta_{i,\mathbf{b}}(I) &= \beta_{i,\mathbf{b}\cdot F}(I_{[F]}) && \text{by Proposition 4.7} \\ &= \mu_{i,\mathbf{b}\mathbf{a}-F}(S_{[F]}/I_{[F]}^{[\mathbf{a}\cdot F]}) && \text{by Proposition 4.1} \\ &= \mu_{i,\mathbf{b}\mathbf{a}-F}(S_{(F)}/I_{(F)}^{[\mathbf{a}]}) && \text{by Proposition 4.6} \\ &= \mu_{i,\mathbf{b}\mathbf{a}-F}(F, S/I^{[\mathbf{a}]}) && \text{by definition} \end{aligned}$$

since $(-)_{(F)}$ is exact. Note that the Bass number in the penultimate line is with respect to the maximal ideal of $S_{(F)}$. The last statement of the theorem is true because $(\mathbf{b}^{\mathbf{a}})^{\mathbf{a}} = \mathbf{b}$ and $(I^{[\mathbf{a}]})^{[\mathbf{a}]} = I$, and because imposing the condition on \mathbf{b} is equivalent to imposing the same condition on $\mathbf{b}^{\mathbf{a}}$. \square

Remark 4.11 Part (i) of the theorem can be thought of as the generalization to arbitrary monomial ideals of the formulas in [6], Proposition 1 and [2], Theorem 2.4, using [9], Theorem 5.2 and the fact that links come from localization ([9], Proposition 5.6).

As a consequence of the theorem, the list of Betti numbers of $I^{\mathbf{a}}$ will be independent of \mathbf{a} , though the \mathbb{Z}^n -degrees in which they occur will vary with \mathbf{a} . Indeed, the list of Betti numbers of $I^{\mathbf{a}}$ is just the list of (localized) Bass numbers of I by part (i) of the theorem. Thus the collection of ideals that are dual to I are very closely related homologically. This will be highlighted again in Section 5 in terms of various geometrically defined resolutions.

Before the above remark, the symbol $I^{\mathbf{a}}$ had not appeared in this section (or the last) without brackets on the \mathbf{a} ; that is, none of the results have been logically dependent on Definition 1.5 or Theorem 2.13. Therefore, Theorem 4.10 can be used to give a second proof of Theorem 2.13. In fact, this “second proof” was discovered before the more elementary proof in Section 2. The next proposition is what allows the irreducible decomposition to be read off of the zeroeth Bass numbers just as the minimal generators are read off the zeroeth Betti numbers.

Proposition 4.12 *Given an ideal $I \subseteq S$ the following are equivalent for $\mathbf{b} \in \mathbb{Z}^F$:*

- (i) $\mathfrak{m}^{\mathbf{b}}$ is an irredundant irreducible component of I .
- (ii) $\mu_{0, \mathbf{b}-F}(F, S/I) = 1$.
- (iii) $\mu_{0, \mathbf{b}-F}(F, S/I) \neq 0$.

Proof: Let $I = \bigcap_j Q_j$ be the (unique) irredundant decomposition of I into irreducible ideals Q_j . Then we have an injection $0 \rightarrow S/I \rightarrow \bigoplus_j S/Q_j$ which, by the proofs of [17], Propositions 3.16 and 3.17, induces an *isomorphism*

$$(3) \quad \underline{\mathrm{Hom}}_S(S/\mathfrak{m}^F, S/I)(-\overline{F}) \rightarrow \bigoplus_j \underline{\mathrm{Hom}}_S(S/\mathfrak{m}^F, S/Q_j)(-\overline{F});$$

this is because the functor $\Delta_{\mathfrak{p}}(\cdot)_{\mathfrak{p}}$ in the [17] reference is easily seen to be $\mathrm{Hom}(R/\mathfrak{p}, \cdot)_{\mathfrak{p}}$ (so we can take $\mathfrak{p} = \mathfrak{m}^F$). Using Lemma 4.3(i) we can move the localization into and out of the $\underline{\mathrm{Hom}}$: for any finitely generated S -modules L and N ,

$$\begin{aligned} \underline{\mathrm{Hom}}_S(L, N)(-\overline{F}) &\cong \underline{\mathrm{Hom}}_{S(-\overline{F})}(L(-\overline{F}), N(-\overline{F})) \\ &\cong S(-\overline{F}) \otimes_{S(F)} \underline{\mathrm{Hom}}_{S(F)}(L_{(F)}, N_{(F)}) \\ &\cong T_{[\overline{F}]} \otimes_k \underline{\mathrm{Hom}}_{S(F)}(L_{(F)}, N_{(F)}). \end{aligned}$$

Treating these as $\mathbb{Z}^{\overline{F}}$ -graded $S_{[\overline{F}]}$ -modules and taking the degree $\mathbf{0} \cdot \overline{F}$ part in the last line yields $\underline{\mathrm{Hom}}_{S(F)}(L_{(F)}, N_{(F)})$. Applying this procedure to Equation (3) reveals an isomorphism

$$\underline{\mathrm{Hom}}_{S(F)}(k, (S/I)_{(F)}) \cong \bigoplus_j \underline{\mathrm{Hom}}_{S(F)}(k, (S/Q_j)_{(F)}).$$

Since we can calculate

$$\mu_{0, \mathbf{b}-F}(F, S/Q_j) = \begin{cases} 1 & \text{if } Q_j = \mathfrak{m}^{\mathbf{b}} \\ 0 & \text{otherwise} \end{cases}$$

the proposition follows from the definition of Bass numbers. \square

Second proof of Theorem 2.13: Every generator $x^{\mathbf{b}}$ of I corresponds to a nontrivial zeroeth Betti number of I which satisfies the condition $F \preceq \mathbf{b} \preceq \mathbf{a} \cdot F$ for $F = \sqrt{\mathbf{b}}$ because $I \subseteq S$ and $\mathbf{a} \succeq \mathbf{a}_I$. After applying Theorem 4.10(i) and the previous proposition, we can conclude that each generator of I does indicate the presence of an appropriate irreducible component of $I^{[\mathbf{a}]}$. To show that each nontrivial zeroeth Bass number of $I^{[\mathbf{a}]}$ comes from some Betti number of I , we demonstrate that if $\mathbf{b} \in \mathbb{Z}^F$ and $\mu_{0, \mathbf{b}-F}(F, S/I) \neq 0$ then $F \preceq \mathbf{b} \preceq \mathbf{a} \cdot F$. Localizing at \mathfrak{m}^F , we may assume that $F = \mathbf{1}$. Clearly $\mathbf{b} \succeq \mathbf{1}$ since $\mathfrak{m}^{\mathbf{b}}$ is \mathfrak{m} -primary, so the desired result falls out of Theorem 3.12 and Proposition 3.7. \square

Next on the agenda is the generalization to arbitrary monomial ideals of an inequality of [2] for squarefree ideals. The topological argument involving links employed there is preempted here by a simple algebraic observation involving localization (which gives links in the squarefree case, see [9], Proposition 5.6).

Theorem 4.13 *If $\mathbf{a} \succeq \mathbf{a}_I$ and $F \preceq \mathbf{b} \preceq \mathbf{a} \cdot F$ then*

$$\beta_{i, \mathbf{b}}(I) \leq \sum_{\mathbf{c} \cdot F = \mathbf{b}^{\mathbf{a}}} \beta_{|F|-i-1, \mathbf{c}}(I^{\mathbf{a}}).$$

Proof: Let \mathbb{F} be a minimal free resolution of $I^{\mathbf{a}}$. Localizing at \mathfrak{m}^F we obtain a free resolution $\mathbb{F}_{(F)}$ of $I^{\mathbf{a}_{(F)}}$ over $S_{(F)}$. The generators of $\mathbb{F}_{(F)}$ as a free $S_{(F)}$ -module are in bijective correspondence with the generators of \mathbb{F} itself. Specifically, for any $\mathbf{b}' \in \mathbb{Z}^F$ we find that $S[\mathbf{c}]_{(F)} = S_{(F)}[\mathbf{b}' \cdot F]$ if and only if $\mathbf{c} \cdot F = \mathbf{b}'$. Thus the number of summands of $\mathbb{F}_{(F)}$ in homological degree $|F| - i - 1$ and \mathbb{Z}^F -degree $\mathbf{b}^{\mathbf{a}}$ is equal to

$$\sum_{\mathbf{c} \cdot F = \mathbf{b}^{\mathbf{a}}} \beta_{|F|-i-1, \mathbf{c}}(I^{\mathbf{a}})$$

since \mathbb{F} is minimal. On the other hand, the number of such summands is clearly $\geq \beta_{|F|-i-1, \mathbf{b}^{\mathbf{a}}}(I^{\mathbf{a}_{(F)}})$, with equality if and only if $\mathbb{F}_{(F)}$ is minimal. Since this last number is equal to $\beta_{i, \mathbf{b}}(I)$ by Theorem 4.10, we are done. \square

Corollary 4.14 (Bayer-Charalambous-Popescu) *If I is squarefree then*

$$\beta_{i, \mathbf{b}}(I) \leq \sum_{\mathbf{b} \preceq \mathbf{c} \preceq \mathbf{1}} \beta_{|\mathbf{b}|-i-1, \mathbf{c}}(I^V)$$

for $0 \leq i \leq n - 1$ and $\mathbf{0} \preceq \mathbf{b} \preceq \mathbf{1}$.

Proof: This is a special case of the theorem once it is noted that (i) $\beta_{|\mathbf{b}|-i-1, \mathbf{c}}(I^V) = 0$ unless $\mathbf{0} \preceq \mathbf{c} \preceq \mathbf{1}$, and (ii) $\mathbf{0} \preceq \mathbf{c}$ and $\mathbf{c} \cdot \sqrt{\mathbf{b}} = \mathbf{b}$ imply $\mathbf{c} \succeq \mathbf{b}$. \square

5 Duality for cellular complexes: the cohull resolution

This section explores the effect of Alexander duality on various geometrically defined free resolutions, in the spirit of [3], [4], and [15]. First, the concept of a geometrically defined resolution is broadened past *cellular resolutions* to include *relative cocellular resolutions* (Definition 5.3). The key result (Theorem 5.8) is presented, though the majority of its proof occupies Section 6. As an application, it is shown how irreducible decompositions can be specified by cellular resolutions (Theorem 5.12). The culmination of these ideas is a new canonical geometric resolution for monomial ideals (Definition 5.15). It is called the *cohull resolution*, and is defined by applying Alexander duality to the hull resolution of [4]. As a special case, the co-Scarf resolution of a cogeneric monomial ideal of [15] is seen to be the cohull resolution (Theorem 5.23), and is thus Alexander dual to the Scarf resolution of a generic monomial ideal [3]. A number of examples are presented, including permutahedron and forest ideals.

Conventions regarding grading and chain complexes:

A chain complex of S -modules

$$\mathbb{F} : \cdots \rightarrow N_{i+1} \rightarrow N_i \rightarrow N_{i-1} \rightarrow \cdots, \quad N_i \text{ in homological degree } i,$$

is viewed as a (homologically) \mathbb{Z} -graded S -module $\bigoplus N_i$ with a differential ∂ of degree -1 . If “[\mathbf{a}]” is attached to \mathbb{F} then each summand is to be shifted in its \mathbb{Z}^n -grading to the left by \mathbf{a} , while “(j)” indicates that the homological degrees are to be shifted down by j , yielding the notation

$$\mathbb{F}[\mathbf{a}](j) : \cdots \rightarrow N_{i+1}[\mathbf{a}] \rightarrow N_i[\mathbf{a}] \rightarrow N_{i-1}[\mathbf{a}] \rightarrow \cdots, \quad N_i \text{ in homological degree } i - j.$$

Here, $N[\mathbf{a}]_{\mathbf{b}} = N_{\mathbf{a}+\mathbf{b}}$ for any S -module N by definition. Taking the S -dual $\mathbb{F}^* := \underline{\text{Hom}}(\mathbb{F}, S)$ changes ∂ to its transpose δ , and makes homological degrees into cohomological degrees, which are the negatives of homological degrees:

$$\begin{aligned} \mathbb{F}^* : \cdots \leftarrow N_{i+1}^* \leftarrow N_i^* \leftarrow N_{i-1}^* \leftarrow \cdots, \quad N_i^* \text{ in homological degree } -i \\ = \text{cohomological degree } i. \end{aligned}$$

Labelled cell complexes provide compact vessels for recording the monomial entries in certain \mathbb{Z}^n -graded free resolutions of an ideal. [4] introduces this notion in the context of monomial modules, but attention is restricted to boundary operators of the cell complex. The definitions below extend the concept to include coboundary operators, as well. For the reader’s convenience, the definition of a labelled regular cell complex and the cellular free complex it determines is recalled briefly below, although the reader is urged to consult [4], Section 1 for more details.

Let $\Lambda \subseteq \mathbb{Z}^n$ be a set of vectors, and let X be a *regular cell complex* whose vertices are indexed by the elements of Λ . For $\mathbf{c}, \mathbf{c}' \in \mathbb{Z}^n$, define the *join* $\mathbf{c} \vee \mathbf{c}'$ to be the componentwise maximum,

i.e. $\mathbf{c} \vee \mathbf{c}'$ is the smallest vector which is greater than or equal to both \mathbf{c} and \mathbf{c}' in each coordinate: $(\mathbf{c} \vee \mathbf{c}')_i = \max(c_i, c'_i)$. Given a face $F \in X$, define the *label* \mathbf{a}_F of F to be the join $\bigvee_{v \in F} \mathbf{a}_v$ of the labels on the vertices in F , where $\mathbf{a}_v \in \Lambda$ is the element corresponding to v . To avoid confusion, the symbol $|X|$ will be used to denote the unlabelled underlying cell complex of the labelled cell complex X .

We assume that $|X|$ comes equipped with an *incidence function* $\varepsilon(F, F') \in \{1, 0, -1\}$ defined on pairs of faces, which is used to define the boundary map in the oriented chain complex of $|X|$ (with coefficients in k). For each $F \in X$, let SF be the free S -module with one generator F in degree \mathbf{a}_F . The *cellular complex* \mathbb{F}_X is the homologically and \mathbb{Z}^n -graded chain complex of S -modules

$$\mathbb{F}_X = \bigoplus_{F \in X, F \neq \emptyset} SF \quad \text{with differential} \quad \partial F = \sum_{G \in X, G \neq \emptyset} \varepsilon(F, G) \frac{m_F}{m_G} \cdot G,$$

where $m_F := x^{\mathbf{a}_F}$. The homological degree of the basis vector $F \in SF$ is the dimension of $F \in |X|$. If \mathbb{F}_X is acyclic, it will be said that X *supports a free resolution* of the module $\langle x^{\mathbf{a}_v} \mid v \in X \text{ is a vertex} \rangle$.

Remark 5.1 In [4] it is assumed that the elements of Λ are pairwise incomparable (as elements in the poset \mathbb{Z}^n), but Λ is not assumed to be finite. Here, however, Λ will always be finite, but pairwise incomparability is not assumed. It is easily verified that all of the results in [4], Section 1 remain true under these hypotheses.

Definition 5.2 (Relative Cellular Complexes) A relative cellular complex $\mathbb{F}_{(X,Y)}$ is the quotient of a cellular complex \mathbb{F}_X supported on a labelled regular cell complex X by a subcomplex \mathbb{F}_Y for some regular cell subcomplex $Y \subseteq X$, with the labelling on Y induced by the labelling on X .

Definition 5.3 (Relative Cocellular Complexes) A relative cocellular complex $\mathbb{F}^{(X,Y)}$ is obtained by taking $\mathbb{F}_{(X,Y)}^*$ for a pair (X, Y) of labelled relative regular cell complexes. If Y is empty, $\mathbb{F}^{(X,Y)}$ may be denoted \mathbb{F}^X and called a cocellular complex supported on X .

Remark 5.4 The relative cocellular complex $\mathbb{F}^{(X,Y)}$ can be viewed as the homogenization of the relative cochain complex of the pair (X, Y) , as long as the label on a dual face F^* is the negative $-\mathbf{a}_F$ of the label on the face F . The coboundary can then be written as $\delta G^* = \sum_{(F \in X, F \neq \emptyset)} \varepsilon(F, G) \frac{m_F}{m_G} \cdot F^*$.

Definition 5.5 Given a labelled regular cell complex X and a vector $\mathbf{b} \in \mathbb{Z}^n$, define the following two labelled subcomplexes of X :

- (i) $X_B(\mathbf{b}) := \{F \in X \mid \mathbf{a}_F \preceq \mathbf{b}\}$, the positively bounded subcomplex of X with respect to \mathbf{b} .
- (ii) $X_U(\mathbf{b}) := \{F \in X \mid \mathbf{a}_F \not\preceq \mathbf{b}\}$, the negatively unbounded subcomplex of X with respect to \mathbf{b} .

Finally, let $X_U := X_U(\mathbf{1})$ be simply the negatively unbounded subcomplex of X .

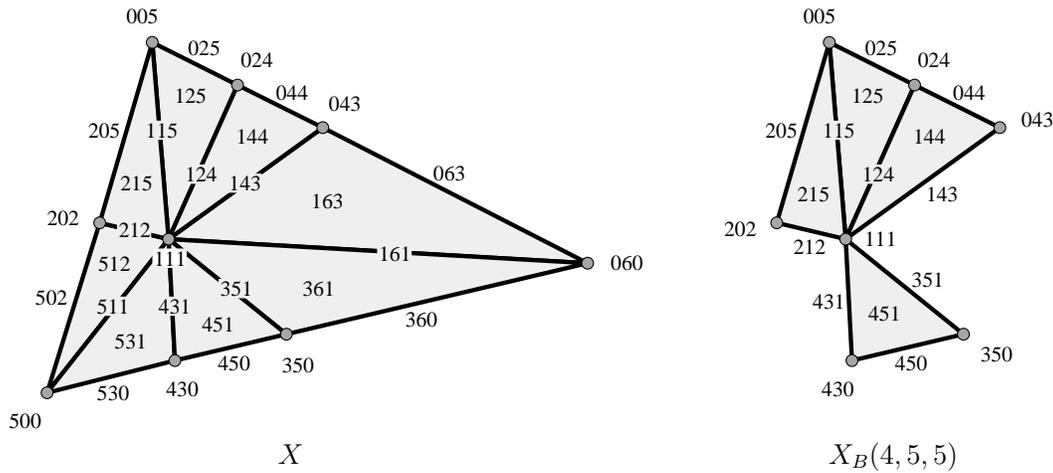


Figure 4

Example 5.6 Let I be as in Example 1.8. The labelled complex X in Figure 4 is the Scarf complex [3] of $I + \mathfrak{m}^{(5,6,6)}$ (see also Example 5.14, below). Hence \mathbb{F}_X is a minimal free resolution by [3], Theorem 3.2. In this case, $(5, 6, 6) = \mathbf{a}_I + \mathbf{1}$, but z^5 is already in I . The label “215” in the diagrams is short for $(2,1,5)$. The subcomplex $X_B(4, 5, 5)$, which is the Scarf complex of the ideal I itself, is also depicted in Figure 4 (see Proposition 5.7, below). The subcomplex X_U is depicted in Figure 5 along with a representation of the labelled relative cellular complex (X, X_U) and the relative cocellular complex $\mathbb{F}^{(X, X_U)}$ of free S -modules determined by it. For this, the edges have been oriented towards the center and the faces counterclockwise. The left copy of S^8 represents the 2-cells in clockwise order starting from 361, the right copy of S^8 represents the edges clockwise starting from 161, and the copy of S represents the lone vertex. The other vertices and edges are not considered since they lie in the subcomplex X_U . It is not a coincidence that the negatively unbounded subcomplex of X is the topological boundary of X —this will always happen for the Scarf complex of a generic artinian monomial ideal, cf. Theorem 5.18. \square

Recall that \mathbf{a}_I is the exponent on the least common multiple of the minimal generators for I . Suppose that we have a cellular resolution \mathbb{F}_X of the ideal $I + \mathfrak{m}^{\mathbf{a}+1}$ with $\mathbf{a} \succeq \mathbf{a}_I$.

Proposition 5.7 $\mathbb{F}_{X_B(\mathbf{b})}$ is a cellular resolution of I for any \mathbf{b} such that $\mathbf{a}_I \preceq \mathbf{b} \preceq \mathbf{a}$.

Proof: With the conditions on \mathbf{b} , the ideal I is generated by all monomials in $I + \mathfrak{m}^{\mathbf{a}+1}$ whose exponent is $\preceq \mathbf{b}$, so the result is a direct consequence of [4], Corollary 1.3. \square

Duality for cellular resolutions says that if the cellular resolution \mathbb{F}_X of $I + \mathfrak{m}^{\mathbf{a}+1}$ has minimal length, a resolution for the Alexander dual $I^{\mathbf{a}}$ with respect to \mathbf{a} can also be recovered from X :

Theorem 5.8 If the cellular resolution \mathbb{F}_X of $I + \mathfrak{m}^{\mathbf{a}+1}$ has length $n - 1$ then $\mathbb{F}^{(X, X_U)}[-\mathbf{a} - \mathbf{1}](1 - n)$ is a relative cocellular resolution of $I^{\mathbf{a}}$. Furthermore, this dual resolution is minimal if \mathbb{F}_X is.

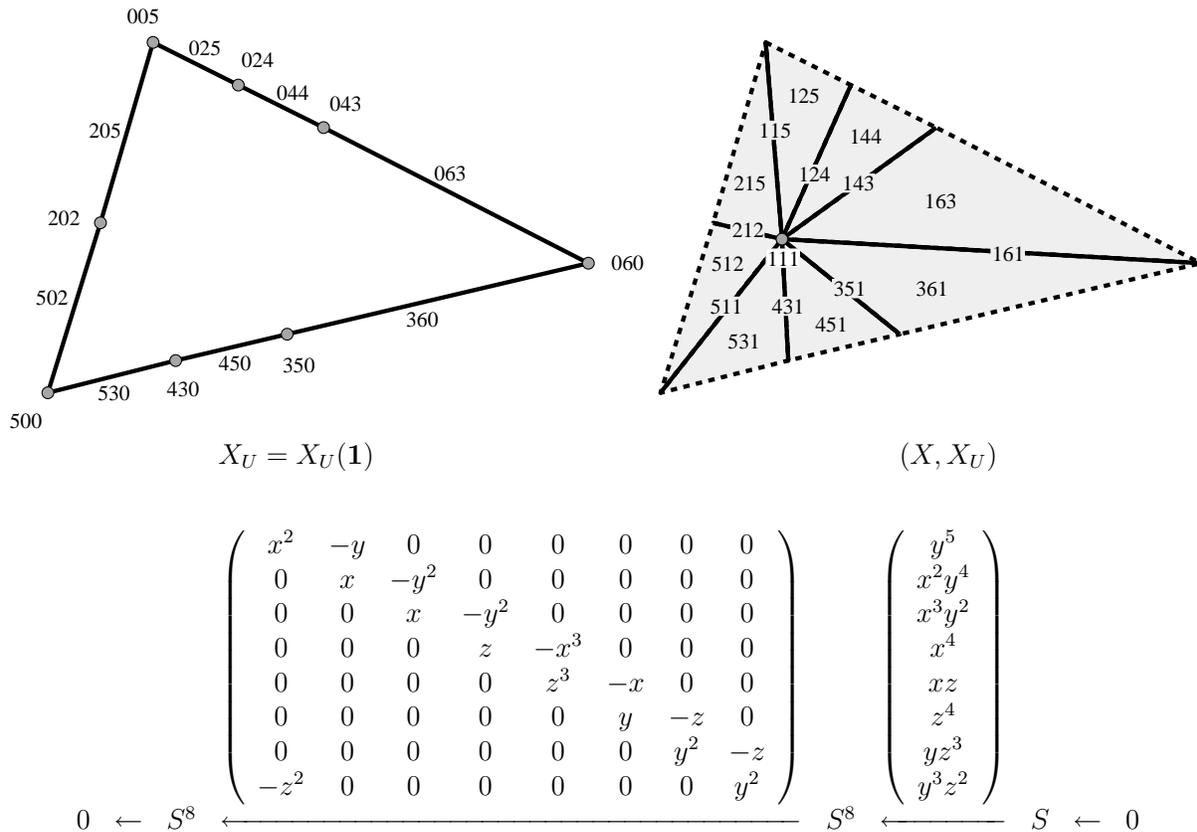


Figure 5

Proof: The first statement will be a direct consequence of Theorem 6.11, below; the necessary assumption here that \mathbb{F}_X has length $n - 1$ is what makes $\mathbb{F}^{(X, X_U)}[-\mathbf{a} - \mathbf{1}](1 - n)$ a *resolution* instead of just a free complex—that is, there are no terms in negative homological degrees. The construction of $\mathbb{F}^{(X, X_U)}$ from \mathbb{F}_X preserves minimality because the matrices defining the differential of the former are submatrices of the transposes of those defining the latter, and we need only check that these entries are in \mathfrak{m} . \square

Remark 5.9 (i) The hypothesis of the theorem requires that X have dimension $(n - 1)$, so that \mathbb{F}_X has minimal *length*, but it does not require that \mathbb{F}_X actually be a minimal resolution.

(ii) It can be shown that X_U may be replaced in the theorem by $X_U(\mathbf{b} + \mathbf{1})$ for any \mathbf{b} satisfying $\mathbf{0} \preceq \mathbf{b} \preceq \mathbf{a}_I - \mathbf{a}_{I^c}$. Here, again, is the mysterious invariant from the remark after Corollary 2.14. In most cases of interest, though, $X_U = X_U(\mathbf{b})$ for all such \mathbf{b} .

Example 5.10 The free complex in Figure 5 is the minimal free resolution of the ideal I^\vee from Example 1.8. The reader may check, for instance, that the product of the large matrix in Figure 5 with the list of generators for I^\vee (which may be treated as a matrix with one row) is zero. Note that the homological and \mathbb{Z}^n -graded shifts promised by Theorem 5.8 aren't visible from the matrices. \square

Theorem 5.8 affords a generalization of [3], Theorem 8.3 on reading irreducible decompositions off of cellular resolutions. We will need the following.

Lemma 5.11 *If the labelled cell complex X supports a minimal free resolution of an artinian ideal $J \subseteq S$ then X is pure of dimension $n - 1$.*

Proof: Any facet has dimension > 0 , so suppose that F is a facet of dimension $d > 0$. Denote by F^* the basis element of the cocellular complex \mathbb{F}^X . The modules $\underline{\text{Ext}}^d(J, S)$ can be calculated as the cohomology of \mathbb{F}^X by definition, and the coboundary $\delta(F^*)$ is zero because F is a facet. Moreover, the image of δ is contained in $\mathfrak{m}\mathbb{F}^X$ by minimality of \mathbb{F}_X , whence F^* is not itself a coboundary. Thus F^* represents a nonzero element of $\underline{\text{Ext}}^d(J, S) \cong \underline{\text{Ext}}^{d+1}(S/J, S)$. It follows that $d = n - 1$ because S/J has only one nonzero such $\underline{\text{Ext}}$ module [5], Proposition 3.3.3(b)(i). \square

For the statement of the next theorem, the following notation is convenient. Suppose $\mathbf{a} \succeq \mathbf{a}_I$ and define, for any $\mathbf{1} \preceq \mathbf{b} \preceq \mathbf{a} + \mathbf{1}$, the *bounded part* $\mathbf{b}_B := (\mathbf{a} + \mathbf{1} - \mathbf{b})^{\mathbf{a}}$ of \mathbf{b} with respect to \mathbf{a} to be the vector whose i^{th} coordinate is b_i if $b_i \leq a_i$ and zero if $b_i = a_i + 1$.

Theorem 5.12 *Let \mathbb{F}_X be a minimal cellular resolution of $I + \mathfrak{m}^{\mathbf{a}+1}$. Then the facets of X are in bijection with the irredundant irreducible components of I , and the intersection $\bigcap_F \mathfrak{m}^{(\mathbf{a}_F)_B}$ over all facets $F \in X$ is an irredundant irreducible decomposition of I .*

Proof: It follows from Lemma 5.11 that under these conditions X must be pure of dimension $n - 1$. Using this, it suffices to show that the label on any facet is $\succeq \mathbf{1}$, for then each facet corresponds to a minimal generator of $I^{\mathbf{a}}$ by Theorem 5.8 and we are done by Proposition 1.7. Suppose, then, that \mathbf{a}_F is ≤ 0 in some coordinate for some facet F ; say $(\mathbf{a}_F)_n = 0$. For $t \gg 0$ consider $Y := X_B(t, t, \dots, t, 0)$, which gives a resolution of $J := (I + \mathfrak{m}^{\mathbf{a}+1}) \cap k[x_1, \dots, x_{n-1}]$ by [4], Corollary 1.3. The resolution \mathbb{F}_Y is minimal because \mathbb{F}_X is, and Y has dimension $n - 1$ because $F \in Y$. On the other hand, J is an artinian ideal of $k[x_1, \dots, x_{n-1}]$, which contradicts Lemma 5.11 (with n replaced by $n - 1$). \square

The major consequence of Theorem 5.8 is the construction of the cohull resolution (Definition 5.15) from the hull resolution [4], Section 2. Therefore, we recall here the definition of the hull complex. Let $t > (n + 1)!$ and define $t^{\mathbf{b}} := (t^{b_1}, \dots, t^{b_n})$. The convex hull of the points $\{t^{\mathbf{b}} \mid x^{\mathbf{b}} \in I\}$ is a polyhedron P_t whose face poset is independent of t . It is shown in [4] that the vertices of P_t are given by those $t^{\mathbf{b}}$ such that $x^{\mathbf{b}}$ is a minimal generator of I . The hull complex $\text{hull}(I)$ is defined to be the bounded faces of P_t , but it may also be described as those faces of P_t admitting a strictly positive inner normal. The hull complex is labelled via the labels on its vertices.

Theorem 5.13 (Bayer-Sturmfels) *The free complex $\mathbb{F}_{\text{hull}(I)}$ is a cellular resolution of I . \square*

Example 5.14 Let Λ be the set of exponents on the minimal generators of a generic monomial ideal I , and let X be the labelled simplex with vertices in Λ . The *Scarf complex* of I is the labelled subcomplex $\Delta_I \subseteq X$ determined by

$$\Delta_I = \{F \in X \mid \mathbf{a}_F = \mathbf{a}_G \Rightarrow F = G\}.$$

It is minimal and coincides with the hull resolution of I by [4], Theorem 2.9. See Example 5.6. \square

Definition 5.15 (The cohull resolution) *The cohull resolution $\text{cohull}_{\mathbf{a}}(I)$ of an ideal I with respect to $\mathbf{a} \succeq \mathbf{a}_I$ is defined to be the relative cocellular resolution $\mathbb{F}^{(X, X_U)}[-\mathbf{a} - \mathbf{1}](1 - n)$, where X is the hull complex of $I^{\mathbf{a}} + \mathfrak{m}^{\mathbf{a}+1}$. The canonical cohull resolution, or simply the cohull resolution $\text{cohull}(I)$ of I is obtained by taking $\mathbf{a} = \mathbf{a}_I$.*

The cohull resolution, like the hull resolution, is a possibly nonminimal resolution that preserves some of the symmetry (in the generators and irreducible components) of an ideal.

There are some geometric properties of hull resolutions of artinian ideals that make cohull resolutions a little more tangible. Suppose, for instance, that J is an artinian monomial ideal, with $x_1^{d_1}, \dots, x_n^{d_n}$ among its minimal generators. Choose $t > (n+1)!$, and let v_1, \dots, v_n be the vertices of the polyhedron P_t determined by these minimal generators. The vertices $\{v_i\}$ of P_t span an affine hyperplane which will be denoted by H .

Fix a strictly positive inner normal φ_G for each $G \in \text{hull}(J)$. Recall that P_t is contained in the (closed) polyhedron $\mathbf{1} + \mathbb{R}_+^n$ (since monomials in S have no negative exponents). Each face $G \in \text{hull}(J)$ spans an affine space which does not contain the vector $\mathbf{1} \in \mathbb{R}^n$ because the hyperplane containing G and normal to φ_G does not contain $\mathbf{1}$. Therefore the projection π from the point $\mathbf{1}$ to the hyperplane H induces a homeomorphism $\text{hull}(J) \rightarrow \pi(\text{hull}(J))$. In fact,

Proposition 5.16 *If J is artinian, $\pi(\text{hull}(J))$ is a regular polytopal subdivision of the simplex $H \cap P_t$.*

Proof: That $H \cap P_t \subset \mathbf{1} + \mathbb{R}_+^n$ is a simplex follows because it is convex and contains v_1, \dots, v_n . Now π induces a map of the boundary $\partial P_t \rightarrow H \cap P_t$ which is obviously surjective. Suppose that $\pi(\mathbf{w})$ is in the interior of $H \cap P_t$ for some $\mathbf{w} \in \partial P_t$. It is enough to show that if a nonzero support functional φ attains its minimum on P_t at \mathbf{w} then φ is strictly positive. All coordinates of φ are ≥ 0 a priori because it attains a minimum on P_t ; but if the i^{th} coordinate of φ is zero then $\langle \varphi, v_i \rangle < \langle \varphi, \mathbf{w} \rangle$ and φ cannot be minimized at \mathbf{w} . \square

Remark 5.17 *This generalizes the result [3], Corollary 4.5 for generic artinian monomial ideals, in view of [4], Theorem 2.9. Regular subdivisions here are as in [19], Definition 5.3.*

We arrive at the following characterization for artinian hull complexes:

Theorem 5.18 *If X is the hull complex of an artinian monomial ideal, then $|X|$ is a simplex and the negatively unbounded complex X_U is the topological boundary of X .*

Proof: By the previous proposition, it suffices to show that a face G of the hull complex of any (not necessarily artinian) ideal has a label without full support if and only if it is contained in the topological boundary of the shifted positive orthant $\mathbf{1} + \mathbb{R}_+^n$. But this holds because the i^{th} coordinate of \mathbf{a}_G is zero if and only if every vertex of G (and hence every point in G) has i^{th} coordinate 1. \square

Although cohull resolutions are relative cocellular by definition, they can frequently be viewed as cellular resolutions, as well. In fact, with a slight weakening of the notion of labelled cell complex, all cohull resolutions are weakly cellular. To be precise, define a *weakly labelled cell complex* to be the same as a labelled cell complex, except that instead of requiring that the label \mathbf{a}_F be equal to the join $\bigvee_{v \in F} \mathbf{a}_v$, we require only that $\mathbf{a}_F \preceq \bigvee_{v \in F} \mathbf{a}_v$ whenever $\dim F > 0$. A free complex or resolution is called *weakly cellular* if it is supported on a weakly labelled cell complex.

Theorem 5.19 *The cohull resolution of I with respect to \mathbf{a} is weakly cellular for any $\mathbf{a} \succeq \mathbf{a}_I$.*

Proof: Let $J = I + \mathfrak{m}^{\mathbf{a}+1}$ and assume the notation from after Definition 5.15. Define Q_t to be the intersection of P_t with the closed half-space containing the origin and determined by the hyperplane H . Then Q_t is a polytope which may also be described as the convex hull of (all of) the vertices of P_t . Furthermore, the bounded faces of P_t are simply those faces of Q_t which admit a strictly positive inner normal. Thus $X := \text{hull}(J)$ is a subcomplex of the boundary complex of Q_t , as is the boundary ∂X .

Let $Y \subset \partial Q_t$ be the subcomplex generated by the facets of Q_t whose inner normal is *not* strictly positive. Denote chain and relative cochain complexes over k by $\mathcal{C}(-)$ and $\mathcal{C}^*(-, -)$. Then $Y \cap X = \partial X$ and the $\mathcal{C}^*(Q_t, Y) = \mathcal{C}^*(X, \partial X)$. For elementary reasons, $\mathcal{C}^*(Q_t, Y) \cong \mathcal{C}^*(X^\vee)$ for some subcomplex X^\vee of the polar polytope Q_t^\vee (use, for instance, the methods of [19], Sections 2.2–2.3). Note that the isomorphisms will exist regardless of the incidence functions in question, by [5], Theorem 6.2.2. That X^\vee is weakly labelled follows from the isomorphism $\mathcal{C}^*(X^\vee) \cong \mathcal{C}^*(X, \partial X)$ and the remark after Definition 5.3. Indeed, the condition $F \supseteq G \Rightarrow \mathbf{a}_F \succeq \mathbf{a}_G$ for faces $F, G \in (X, \partial X)$ is equivalent to the condition $F^\vee \subseteq G^\vee \Rightarrow -\mathbf{a}_F \preceq -\mathbf{a}_G$ for faces of X^\vee , and this need only be applied when F is a facet containing G and F^\vee is a vertex of G^\vee . \square

Proposition 5.20 *If a weakly cellular resolution is minimal, it is cellular. In particular, if a cohull resolution is minimal, it is cellular.*

Proof: Let $(\widetilde{\mathbb{F}}, \widetilde{\partial})$ denote the augmented complex $\mathbb{F}_X \rightarrow I \rightarrow 0$, where X is a weakly labelled complex supporting a free resolution of I . We show that if $G \in X$ then $\mathbf{a}_G \succ \bigvee_{v \in G} \mathbf{a}_v$ implies \mathbb{F} is not

minimal. This is vacuous if $\dim G = 0$, so assume $\dim G$ has minimal dimension ≥ 1 , and suppose that $\mathbf{a}_G - \mathbf{e}_i \succeq \bigvee_{v \in G} \mathbf{a}_v$. Then $\tilde{\partial}(G) = x_i y$ for some $y \in \tilde{\mathbb{F}}$ because $\dim G$ is minimal. It follows that $x_i \tilde{\partial}(y) = \tilde{\partial}(x_i y) = 0$, whence $\tilde{\partial}(y) = 0$ because $\tilde{\mathbb{F}}$ is torsion-free. Thus $\tilde{\partial}(G) \in x_i \ker(\tilde{\partial}) \subseteq \mathbf{m} \cdot \ker(\tilde{\partial})$ does not represent a minimal generator of $\ker(\tilde{\partial})$ by Nakayama's Lemma for graded modules. \square

Remark 5.21 For cohull resolutions the proposition is probably true without the hypothesis of minimality, but a proof (which would likely be geometric instead of algebraic) has not been found. In particular, all cohull resolutions in the examples below are cellular. Cellularity of the cohull resolution is equivalent to the following more concrete statement: the label on any interior face of the hull complex of an artinian ideal is the greatest common divisor of the labels on the facets that contain it.

Example 5.22 (continuation of Example 1.9) The minimal resolution of the permutahedron ideal I of Example 1.9 is, by [4], Example 1.9, the hull resolution, which is supported on a permutahedron. The minimal resolution of $I + \mathbf{m}^{(n+1)\mathbf{1}}$ is also the hull resolution, and is supported on the complex X which may be described as follows.

There are two kinds of faces of X . The first kind are those that make up the boundary ∂X ; these are indexed by the proper nonempty faces $F \in \Delta$ and have vertices $t^{(n+1)\mathbf{e}_i} \in P_t$ for $i \in F$ (recall from Section 1 that \mathbf{e}_i denotes the i^{th} basis vector of \mathbb{Z}^n and $\Delta = \{1, \dots, n\}$ is the $(n-1)$ -simplex). On the other hand, the interior p -faces of X are in bijection with the chains

$$(4) \quad \emptyset \prec F_1 \prec F_2 \prec \dots \prec F_{n-p}$$

of faces of Δ , where F_{n-p} might (or might not) equal Δ . Note that the interior faces of X for which $F_{n-p} = \Delta$ are faces of the permutahedron itself.

More generally, an interior p -face G given by (4) for which $F_{n-p} \neq \Delta$ is affinely spanned by the permutahedral $(p-1)$ -face $G' : \emptyset \prec F_1 \prec \dots \prec F_{n-p} \prec \Delta$ and the ‘‘artinian’’ vertices $\{t^{(n+1)\mathbf{e}_i} \mid i \notin F_{n-p}\}$ of P_t . In fact, a functional which attains its minimum (in P_t) on G may be produced directly. For this purpose, define for any $F \in \Delta$ the functional F^\dagger on \mathbb{R}^n to be the transpose of F ; i.e. $\langle F^\dagger, \mathbf{e}_i \rangle = 1$ if $i \in F$ and zero otherwise. Then the functional $\varphi_\epsilon := \mathbf{1}^\dagger + \epsilon \sum_{j=1}^{n-p} F_j^\dagger$ attains its minimum (in P_t) on G' for all $0 < \epsilon \ll 1$. But for $\epsilon \gg 0$ we have $\langle \varphi_\epsilon, t^{(n+1)\mathbf{e}_i} \rangle < \langle \varphi_\epsilon, G' \rangle$ whenever $i \notin F_{n-p}$. Thus we can choose the unique ϵ that makes $\langle \varphi_\epsilon, t^{(n+1)\mathbf{e}_i} \rangle = \langle \varphi_\epsilon, G' \rangle$ for all $i \notin F_{n-p}$, so that φ_ϵ attains its minimum on G .

It is easy to check that the labels on the faces of X are distinct, whence \mathbb{F}_X is the minimal resolution of $I + \mathbf{m}^{(n+1)\mathbf{1}}$ by [4], Remark 1.4. In particular, the irredundant irreducible components of I are in bijection with the facets of X by Theorem 5.12, and the generators of the forest ideal I^\vee are given by $x^{(n+1)\mathbf{1} - \mathbf{a}_G}$ for facets $G \in X$. This recovers the generators for I^\vee in Example 1.9.

Retaining earlier notation, the face G has dimension $1 + \dim(G')$. Thus the p -faces of X are in bijection with the collection of p - and $(p-1)$ -faces of the permutahedron. In fact, the (unlabelled)

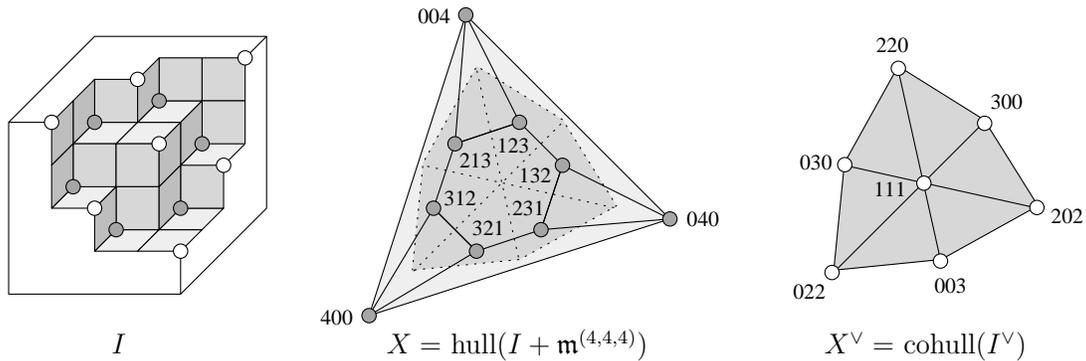


Figure 6: I and I^\vee are the permutahedron and forest ideals when $n = 3$. The complex X is the (labelled) regular polytopal subdivision of the simplex promised by Proposition 5.16. Overlaid on this figure is the dual complex X^\vee (without its labelling). At right, X^\vee is shown with its labelling, which is \mathbb{Z}^n -shifted as per Theorem 5.8. Turn the picture over for the staircase of I^\vee .

pair $(|X|, |\partial X|)$ has the same faces as the pair $(\partial(v * Y), v)$ consisting of the boundary of the cone over the permutahedron Y rel the apex of the cone. The cellular complex X^\vee supporting the cohull resolution of the forest ideal I^\vee is therefore easy to describe. Let Y be the permutahedron in \mathbb{R}^n and Y^\vee its polar. Then X^\vee is the cone over ∂Y^\vee from the barycenter of Y^\vee . The vertices G^\vee of X^\vee , which are labelled by the generators of I^\vee , almost all correspond to the facets G' of Y (whose labellings are as above). Only the apex of the cone is an exception, corresponding instead to the interior of Y . The case $n = 3$ is depicted in Figure 6; it should be noted that the equality $Y = Y^\vee$ is only because Y is 2-dimensional, not some more general self-duality.

Now $\text{cohull}(I^\vee)$ is a cellular resolution of $I^\vee = I^\vee + \mathfrak{m}^{\mathbf{a}_I + \mathbf{1}}$, so we can dualize this cellular resolution using Theorem 5.8 again. This yields a minimal relative cocellular resolution of I , which is seen to be cellular and (coincidentally?) equal to $\text{hull}(I)$. \square

Recall from Section 1 that an ideal is *cogeneric* if it is Alexander dual to a generic ideal. The minimal resolution of such an ideal was introduced in [15], where it was dubbed the *co-Scarf resolution*. The next theorem, along with the proof of Theorem 5.19 above, explains why the construction in [15] involved a subcomplex of the boundary of the simple polytope dual to the simplicial polytope of which the Scarf complex is a subcomplex. The theorem is a direct consequence of Theorem 5.8, Example 5.14, and Proposition 5.20.

Theorem 5.23 *Any cohull resolution of a cogeneric monomial ideal is minimal and cellular.* \square

Remark 5.24 That the co-Scarf resolution is cellular as opposed to weakly cellular was assumed in [4], Example 1.8 but overlooked in [15].

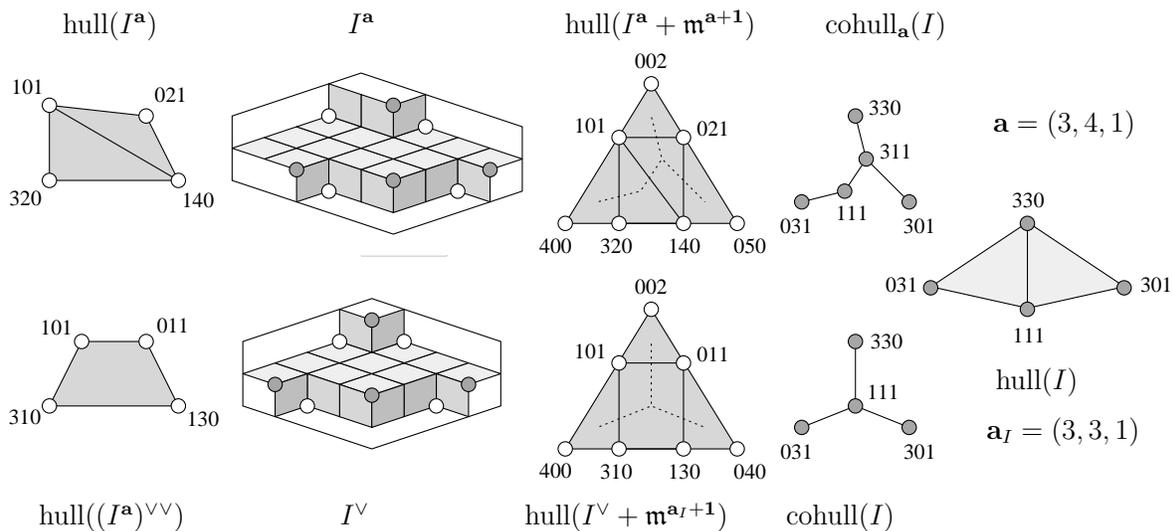


Figure 7

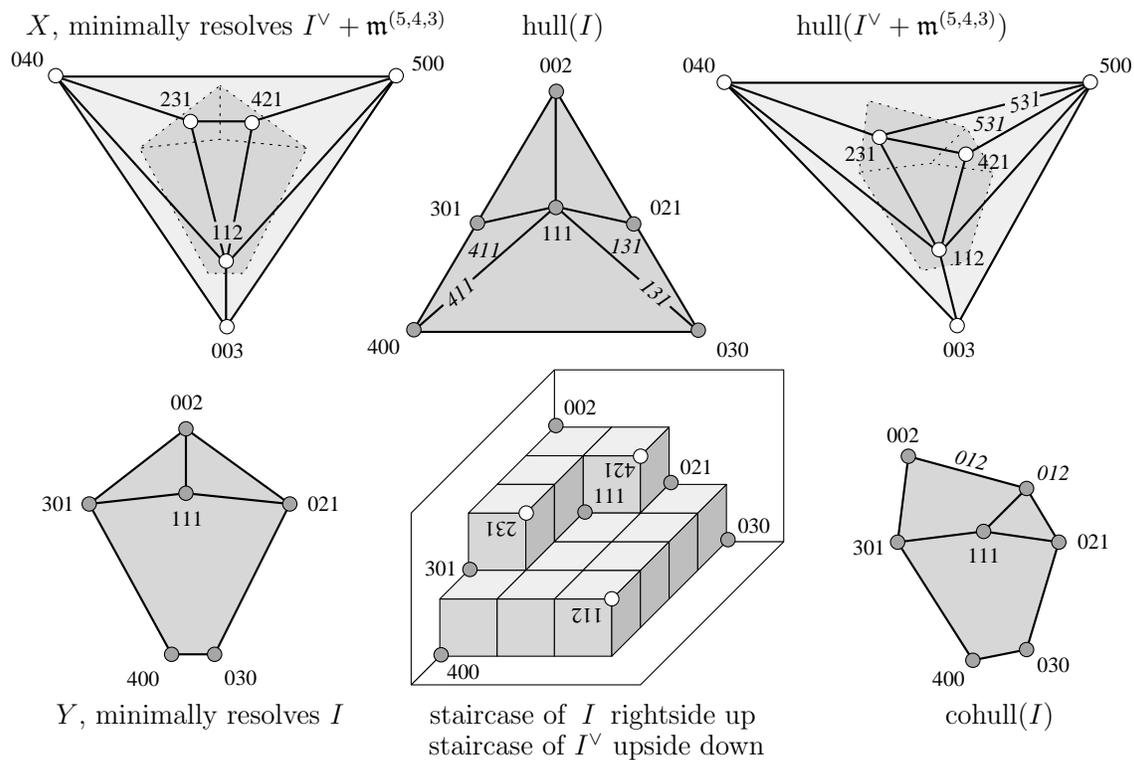


Figure 8

Example 5.25 It is possible for the hull and cohull resolutions to coincide for a given ideal I . For instance, this occurs if $I = \mathfrak{m}$; or if I is simultaneously generic and cogeneric (which turns out to be pretty hard to accomplish!); or if I is the permutahedron ideal in 3 variables. Conjecturally, the hull and cohull resolutions should coincide for permutahedron ideals of all dimensions. \square

Example 5.26 Of course, it is also possible for the hull and cohull resolutions to be very different. For instance, the cohull resolution of the ideal I^\vee from Examples 1.8 and 5.10 is the co-Scarfi resolution, which is cellular and supported on an octagon with only one maximal face (dualize the picture in Figure 5). On the other hand, $\text{hull}(I^\vee)$ is a triangulation of the same octagon. \square

Example 5.27 The canonical cohull resolution can differ from a noncanonical cohull resolution. For instance, let $I = \langle x^3z, xyz, y^3z, x^3y^3 \rangle$ and $\mathbf{a} = (3, 4, 1)$, so $I^{\mathbf{a}} = \langle xz, x^3y^2, xy^4, y^2z \rangle$ and $I^\vee = \langle xz, x^3y, xy^3, yz \rangle$. Since $\text{hull}(I)$ is not minimal, we look elsewhere for the minimal resolution of I . But $\text{hull}(I^{\mathbf{a}} + \mathfrak{m}^{\mathbf{a}+1})$ is not minimal, and the failure of minimality occurs in such a way that $\text{cohull}_{\mathbf{a}}(I)$ is also not minimal. On the other hand, the offending nonminimal edge is not present in $\text{hull}(I^\vee + \mathfrak{m}^{\mathbf{a}+1})$, and this resolution is minimal. It follows that $\text{cohull}(I)$ is minimal. Note how the passage from $I^{\mathbf{a}}$ to $(I^{\mathbf{a}})^{\vee\vee} = I^\vee$ “tightens” the hull resolution of $I^{\mathbf{a}}$ to make the nonminimal edge disappear in $\text{hull}((I^{\mathbf{a}})^{\vee\vee})$, cf. the remark after Corollary 2.14.

The labelled complexes supporting these resolutions are all depicted in Figure 7, where the resolutions with black vertices are drawn “upside down” to make their superimposition on the staircase diagram for I easier to visualize. Observe that a staircase diagram for I can be obtained by turning over the staircase diagram for either $I^{\mathbf{a}}$ or I^\vee , although these result in different “bounding boxes” for I . Note that all of the complexes, particularly the cohull complexes, are labelled and not just weakly labelled. \square

Example 5.28 Finally, an example to illustrate that not all cellular resolutions come directly from hull and cohull resolutions, so that the algebraic techniques to prove exactness in Section 6 prove a stronger duality for resolutions than a geometric treatment such as that in [4] or [15] could provide. All of the labelled cellular complexes from this example are depicted in Figure 8.

Let $I = \langle z^2, x^3z, x^4, y^3, y^2z, xyz \rangle$, so that $I^\vee = \langle xyz^2, x^2y^3z, x^4y^2z \rangle$. Then $\text{hull}(I)$ and $\text{cohull}(I)$ are not minimal (the offending cells have italic labels); moreover, $\text{cohull}_{\mathbf{a}}(I) = \text{cohull}(I)$ for all $\mathbf{a} \succeq \mathbf{a}_I = (4, 3, 2)$. Nonetheless, the minimal resolution \mathbb{F}_X of $I^\vee + \mathfrak{m}^{(5,4,3)}$ is cellular, so Theorem 5.8 applies, yielding a minimal relative cocellular resolution $\mathbb{F}^{(X, X^\vee)}[-(5, 4, 3)](-2)$ for I . In fact, this relative cocellular resolution is cellular, supported on the labelled cell complex Y . \square

6 Deformations and limits of resolutions

The final item on the agenda is the proof of Theorem 5.8. To that end, the goal of this section is Theorem 6.11, which is actually a little more general than Theorem 5.8. It can be viewed as the

result of applying a limiting process to a collection of pairs of linked artinian monomial ideals that are deformations of a given pair. The entire section is a setup to apply a limit to Proposition 3.11, and is another manifestation of the kinship of Alexander duality and other types of duality for Gorenstein rings. The maps $f_{\mathbf{b}}$ in the following definition accomplish the deformations.

Definition 6.1 Define the map $f_{\mathbf{b}}: \mathbb{Z}^n \rightarrow \mathbb{Z}^n$ for $\mathbf{b} \succeq \mathbf{0}$ by the coordinatewise formula

$$f_{\mathbf{b}}(\mathbf{c})_i = \begin{cases} c_i - b_i & \text{if } c_i \leq 0 \\ c_i & \text{if } c_i \geq 1 \end{cases}$$

To avoid messy exponents we also let $f_{\mathbf{b}}(x^{\mathbf{c}}) = x^{f_{\mathbf{b}}(\mathbf{c})}$. Whenever the symbol $f_{\mathbf{b}}$ is written, it will be assumed that $\mathbf{b} \succeq \mathbf{0}$.

Proposition 6.2 If $I \subseteq S$ is any monomial ideal then $\langle f_{\mathbf{b}}(I) \rangle = S[\mathbf{b}] \cap \tilde{I}$.

Proof: It is clear from the definition that $f_{\mathbf{b}}(\mathbf{c}) \succeq -\mathbf{b}$ if $\mathbf{c} \succeq \mathbf{0}$, whence $\langle f_{\mathbf{b}}(I) \rangle \subseteq S[\mathbf{b}]$. Since also $f_{\mathbf{b}}(\mathbf{c})^+ = f_{\mathbf{b}}(\mathbf{c}^+)$, we conclude that $\langle f_{\mathbf{b}}(I) \rangle \subseteq \tilde{I}$ as well. For the reverse inclusion, assume $x^{\mathbf{c}} \in S[\mathbf{b}] \cap \tilde{I}$. Then $f_{\mathbf{b}}(x^{\mathbf{c}^+}) \in f_{\mathbf{b}}(I)$ and divides $x^{\mathbf{c}}$ because $f_{\mathbf{b}}(\mathbf{c}^+) \preceq \mathbf{c}$ whenever $\mathbf{c} \succeq -\mathbf{b}$, a fact which is easily seen from the definition. \square

Recall from Section 5 that the join of $\mathbf{c}, \mathbf{c}' \in \mathbb{Z}^n$ is the componentwise maximum.

Lemma 6.3 The map $f_{\mathbf{b}}$ preserves joins; that is, $f_{\mathbf{b}}(\mathbf{c} \vee \mathbf{c}') = f_{\mathbf{b}}(\mathbf{c}) \vee f_{\mathbf{b}}(\mathbf{c}')$. \square

This lemma is important because of the next proposition, originally due to D. Bayer. Let X be a labelled cell complex, and suppose $f: \mathbb{Z}^n \rightarrow \mathbb{Z}^n$ is a map respecting joins. Denote by $f(X)$ the labelled cell complex which is obtained by applying f to the labels on the faces of X . Thus $G \in f(X)$ is labelled by $f(\mathbf{a}_G)$ whenever $G \in X$ is labelled by \mathbf{a}_G .

Proposition 6.4 Let \mathbb{F}_X be a cellular resolution of a finitely generated module $M \subseteq T$. If $f: \mathbb{Z}^n \rightarrow \mathbb{Z}^n$ preserves joins then $\mathbb{F}_{f(X)}$ is a resolution of $\langle f(M) \rangle$.

Proof: Note that because f respects joins the effect of f is determined by its effect on the vertex labels. Similarly, $\langle f(M) \rangle = \langle f(x^{\mathbf{b}}) \mid \mathbf{b} \text{ is a vertex label of } X \rangle$. Thus one only needs to check that $\mathbb{F}_{f(X)}$ is acyclic. It suffices to check that $X_B(\mathbf{b})$ is acyclic for all $\mathbf{b} \in \mathbb{Z}^n$, by the acyclicity criterion of [4], Proposition 1.2.

Suppose, then, that α is a cycle of the reduced chain complex of $|f(X)_B(\mathbf{b})|$. Then α also represents a cycle of $|X|$. Let \mathbf{c} be the join of the labels on the faces in the support of α , considered as faces of X . Since f preserves joins, $f(\mathbf{c}) \preceq \mathbf{b}$ and $|X_B(\mathbf{c})| \subseteq |f(X)_B(\mathbf{b})|$. Now α is a boundary in the reduced chain complex of $|X_B(\mathbf{c})|$ by [4], Proposition 1.2, and it follows that α is also a boundary in the reduced chain complex of $|X_B(\mathbf{b})|$, completing the proof. \square

Corollary 6.5 *If \mathbb{F}_X is a cellular resolution of I then $\mathbb{F}_{f_{\mathbf{b}}(X)}$ is a cellular resolution of $S[\mathbf{b}] \cap \tilde{I}$. \square*

Keeping the notation of the corollary we can augment $\mathbb{F}_{f_{\mathbf{b}}(X)}$ to a resolution of $S[\mathbf{b}]/S[\mathbf{b}] \cap \tilde{I}$, homologically shifted down 1, by adding a summand $S[\mathbf{b}]$ in homological degree -1 . We denote this augmented resolution by $\mathbb{F}_{\mathbf{b}}(X)$, and we let $\mathbb{F}^{\mathbf{b}}(X) := \mathbb{F}_{\mathbf{b}}(X)^*$, with differential $\delta_{\mathbf{b}}$. The generator of the summand $S[-f_{\mathbf{b}}(\mathbf{a}_F)] \subseteq \mathbb{F}_{\mathbf{b}}(X)$ corresponding to the face F will be denoted by $F_{\mathbf{b}}$, while the generator of $S[f_{\mathbf{b}}(\mathbf{a}_F)] = S[-f_{\mathbf{b}}(\mathbf{a}_F)]^* \subseteq \mathbb{F}^{\mathbf{b}}(X)$ will be denoted by $F^{\mathbf{b}}$. Keep in mind that $F^{\mathbf{b}}$ is in \mathbb{Z}^n -graded degree $-f_{\mathbf{b}}(\mathbf{a}_F)$.

We will soon be defining maps between the $\mathbb{F}^{\mathbf{b}}(X)$ for various \mathbf{b} , and the following lemma, particularly part (ii), will be the tool used to prove that these maps are well-defined, commute with the differentials, and form an inverse system.

Lemma 6.6 *If $\mathbf{b} \succeq \mathbf{b}' \succeq \mathbf{0}$ then*

$$\begin{aligned} (i) \quad & f_{\mathbf{b}} = f_{\mathbf{b}-\mathbf{b}'} \circ f_{\mathbf{b}'}, \\ (ii) \quad & f_{\mathbf{b}'}(\mathbf{c}) - f_{\mathbf{b}}(\mathbf{c}) = \mathbf{c} - f_{\mathbf{b}-\mathbf{b}'}(\mathbf{c}). \end{aligned}$$

Proof: Plug and chug, using the equality $f_{\mathbf{b}}(\mathbf{c})^+ = \mathbf{c}^+$ for (i). \square

Lemma 6.7 *For every $\mathbf{b} \succeq \mathbf{b}' \succeq \mathbf{0}$ we have an injection of chain complexes $\varphi_{\mathbf{b},\mathbf{b}'}: \mathbb{F}^{\mathbf{b}}(X) \hookrightarrow \mathbb{F}^{\mathbf{b}'}(X)$ sending $F^{\mathbf{b}}$ to $\frac{m_F}{f_{\mathbf{b}-\mathbf{b}'}(m_F)} F^{\mathbf{b}'}$.*

Proof: There are two aspects to the proof: (i) the given map is an injection of homologically graded modules which (as a map of \mathbb{Z}^n -graded modules) has degree $\mathbf{0}$, and (ii) the injections commute with the differentials. The first follows from the equality $-f_{\mathbf{b}}(\mathbf{a}_F) = -f_{\mathbf{b}'}(\mathbf{a}_F) + \mathbf{a}_F - f_{\mathbf{b}-\mathbf{b}'}(\mathbf{a}_F)$ which is easily seen to be equivalent to Lemma 6.6(ii) when $\mathbf{c} = \mathbf{a}_F$. The second is a longer calculation directly from the definition of the differentials $\delta_{\mathbf{b}}$ and $\delta_{\mathbf{b}'}$ of the chain complexes $\mathbb{F}_{\mathbf{b}}$ and $\mathbb{F}_{\mathbf{b}'}$.

The definitions imply that $\delta_{\mathbf{b}}$ is just the transpose of the differential from the cellular free complex as defined in [4]. Thus, $\delta_{\mathbf{b}}(F^{\mathbf{b}}) = \sum_{G \in X} \varepsilon(G, F) \frac{f_{\mathbf{b}}(m_G)}{f_{\mathbf{b}}(m_F)} G^{\mathbf{b}}$, where ε is the incidence function defining the differential of X . Note that $\varepsilon(G, F)$ is nonzero only if $G \supseteq F$. We have

$$\begin{aligned} \delta_{\mathbf{b}'} \circ \varphi_{\mathbf{b},\mathbf{b}'}(F^{\mathbf{b}}) &= \sum_{G \in X} \varepsilon(G, F) \frac{f_{\mathbf{b}'}(m_G)}{f_{\mathbf{b}'}(m_F)} \cdot \frac{m_F}{f_{\mathbf{b}-\mathbf{b}'}(m_F)} G^{\mathbf{b}'} \\ &= \sum_{G \in X} \varepsilon(G, F) \frac{m_G}{f_{\mathbf{b}-\mathbf{b}'}(m_G)} \cdot \frac{f_{\mathbf{b}}(m_G)}{f_{\mathbf{b}}(m_F)} G^{\mathbf{b}'} \\ &= \varphi_{\mathbf{b},\mathbf{b}'} \circ \delta_{\mathbf{b}}(F^{\mathbf{b}}), \end{aligned}$$

where the transition from the first line to the second is accomplished by two applications of Lemma 6.6(ii). \square

Lemma 6.8 *If $\mathbf{b} \succeq \mathbf{b}' \succeq \mathbf{b}'' \succeq \mathbf{0}$ then $\varphi_{\mathbf{b}, \mathbf{b}''} = \varphi_{\mathbf{b}', \mathbf{b}''} \circ \varphi_{\mathbf{b}, \mathbf{b}'}$.*

Proof: We need only check the equality as maps of modules. The proof again uses property (ii) from Lemma 6.6, and it involves manipulations similar to those in the proof of Lemma 6.7. \square

These lemmata show that we have an inverse system of complexes of free modules, so it is natural now to take the inverse limit. With $\mathbb{F}^t(X) := \mathbb{F}^{t-1}(X)$ we can simplify a little since the inverse systems $\{\mathbb{F}^{\mathbf{b}}(X)\}_{\mathbf{b} \succeq \mathbf{0}}$ and $\{\mathbb{F}^t(X)\}_{t \in \mathbb{N}}$ are cofinal, so that their limits are the same. We take this opportunity to note that our inverse limits, when taken in the category of \mathbb{Z}^n -graded objects and degree zero maps, will be denoted by \varprojlim_t , and that S is complete in this category. Recall that, for our inverse system $\{\mathbb{F}^t(X)\}_{t \in \mathbb{N}}$ of chain complexes, for instance, this is defined as

$$\varprojlim_t \mathbb{F}^t(X) = \bigoplus_{\mathbf{c} \in \mathbb{Z}^n} \varprojlim_t \mathbb{F}^t(X)_{\mathbf{c}},$$

where the inverse limits on the right are in the category of chain complexes of k -vector spaces.

At each stage in the inverse system, $f_{\mathbf{b}}$ moves the labels on X_U away from the first orthant, in negative directions, turning any zeros into arbitrarily large negative integers (hence the name “negatively unbounded” for the subcomplex X_U of X). Then S -duality makes the negative integers positive. Thus the maps $f_{\mathbf{b}}$, combined as they are with S -duality in the definition of $\mathbb{F}^{\mathbf{b}}$, create irreducible components of $I^{\mathbf{a}}$ from those generators of I which do not have full support by pushing the zeros out to (positive) infinity. In the limit, the vertices defining those generators disappear. This provides the intuition for the next result.

Theorem 6.9 $\mathbb{F}^{(X, X_U)} = \varprojlim_t \mathbb{F}^t(X)$.

Proof: The first observations are that $\mathbb{F}^{(X, X_U)}$ is a subcomplex of \mathbb{F}^t for all t , and that the maps $\varphi_{t, t'} := \varphi_{t-1, t'-1}$ defining the inverse system restrict to the identity on $\mathbb{F}^{(X, X_U)}$. This is because of the way $f_t := f_{t-1}$ is defined:

$$(5) \quad f_{t-t'}(m_F) = m_F \iff t = t' \text{ or } F \notin X_U$$

because $f_{\mathbf{b}}(\mathbf{c})_i = c_i$ for all i precisely when $\mathbf{c} \succeq \mathbf{1}$. Thus we have, for all $t \geq 0$, exact sequences

$$(6) \quad 0 \rightarrow \mathbb{F}^{(X, X_U)} \rightarrow \mathbb{F}^t(X) \rightarrow \mathbb{F}^t(X_U) \rightarrow 0$$

giving rise to a corresponding exact sequence of inverse systems. To be more precise, the maps $\{\varphi_{t, t'}\}$ from the inverse system $\{\mathbb{F}^t(X)\}_{t \in \mathbb{N}}$ induce maps $\{\psi_{t, t'}: \mathbb{F}^t(X_U) \rightarrow \mathbb{F}^{t'}(X_U)\}_{t \geq t'}$ which make $\{\mathbb{F}^t(X_U)\}_{t \in \mathbb{N}}$ into an inverse system.

It is readily seen that the maps $\psi_{t, t'}$ are injections, so that $\varprojlim_t \mathbb{F}^t(X_U) = \bigcap_t \psi_{t, 0}(\mathbb{F}^t(X_U))$. Furthermore, statement (5) implies that $\psi_{t, t'}(\mathbb{F}^t(X_U)) \subseteq \mathfrak{m}\mathbb{F}^{t'}(X_U)$ if $t > t'$. It follows from the

Krull intersection theorem that the inverse limit is zero. Since the inverse limit is always left exact our exact sequence of inverse systems arising from (6) yields the desired isomorphism. \square

So we can write $\mathbb{F}^{(X, X_U)}$ as an inverse limit. What have we gained? In the category of \mathbb{Z}^n -graded objects in which each graded piece has finite dimension over k (e.g. if the objects are chain complexes which are finitely generated as S -modules), the functor \varprojlim is exact, at least in the case where the inverse systems are indexed by \mathbb{N} —see [18], Exercise 3.5.2. With this in mind the following corollary is a simple consequence of [18], Theorem 3.5.8.

Corollary 6.10 *To compute homology we have $H_i(\mathbb{F}^{(X, X_U)}) = \varprojlim_t H_i(\mathbb{F}^t(X))$.* \square

Until this point in this section, the labelled cell complex X has been arbitrary. Now, however, we suppose that X supports a cellular free resolution of the ideal $I + \mathfrak{m}^{\mathbf{a}+1}$, with $\mathbf{a} \succeq \mathbf{a}_I$. We will see shortly that for any t the only nonvanishing homology of $\mathbb{F}^t(X)$ is in homological degree $1 - n$, so the previous corollary implies that the same holds for $\mathbb{F}^{(X, X_U)}$. Now \mathbb{F}_X has length at least $n - 1$ (i.e. $\dim X \geq n - 1$) because it gives a free resolution of an artinian ideal; if we are so lucky that \mathbb{F}_X has length exactly $n - 1$, then the summand of $\mathbb{F}^{(X, X_U)}$ in homological degree $1 - n$ will be the last nonzero term. In other words, $\mathbb{F}^{(X, X_U)}$ will be a *free resolution* of some S -module. This is what makes Theorem 5.8 a special case of the next result. Even if we aren't so lucky with the length of \mathbb{F}_X , at least it will be split exact in homological degrees $> n - 1$ (so that $\mathbb{F}^{(X, X_U)}$ is split exact in homological degrees $< 1 - n$), and we can still determine what the nonzero homology module is:

Theorem 6.11 *Under the above conditions, $H_i(\mathbb{F}^{(X, X_U)}) = 0$ if $i \neq 1 - n$, and $H_{1-n}(\mathbb{F}^{(X, X_U)}) = I^{[\mathbf{a}]}[\mathbf{a} + \mathbf{1}]$.*

Proof: Let $J = I + \mathfrak{m}^{\mathbf{a}+1}$. For any $\mathbf{b} \succeq \mathbf{0}$ Corollary 6.5 implies that $\mathbb{F}_{\mathbf{b}}(X)$ is a free resolution of the module $S[\mathbf{b}]/S[\mathbf{b}] \cap \tilde{J}$, homologically shifted down by 1. Thus $\mathbb{F}^{\mathbf{b}}(X)$, which is the S -dual of $\mathbb{F}_{\mathbf{b}}(X)$, is a complex whose homology in degree $i - 1$ is $\underline{\text{Ext}}_S^i(S[\mathbf{b}]/S[\mathbf{b}] \cap \tilde{J}, S)$. Now $S[\mathbf{b}]/S[\mathbf{b}] \cap \tilde{J} \subseteq T/\tilde{J}$ is artinian since $J = I + \mathfrak{m}^{\mathbf{a}+1}$ is, and it is noetherian because $S[\mathbf{b}]$ is. Hence the Ext module in question is, by [5], Theorem 3.3.10(c), nonzero only for $i = n$. Moreover, Proposition 3.11 produces the equality

$$\underline{\text{Ext}}_S^n(S[\mathbf{b}]/S[\mathbf{b}] \cap \tilde{J}, S) = (I/I \cap \mathfrak{m}^{\mathbf{a}+\mathbf{b}+1})[\mathbf{a} + \mathbf{1}].$$

Taking the $\varprojlim_{\mathbf{b}}$ of this last line and applying Corollary 6.10 along with the completeness of S proves the theorem. \square

References

- [1] *D. Bayer*, Monomial ideals and duality, Lecture notes, Berkeley 1995–96, available by anonymous ftp at math.columbia.edu/pub/bayer/monomials_duality/monomials.ps.

- [2] *D. Bayer, H. Charalambous, and S. Popescu*, Extremal Betti numbers and applications to monomial ideals, preprint (math.AG/9804052), 1998.
- [3] *D. Bayer, I. Peeva, and B. Sturmfels*, Monomial resolutions, *Math. Res. Letters* **5** (1998), 31-46.
- [4] *D. Bayer and B. Sturmfels*, Cellular resolutions of monomial modules *J. reine angew. Math.* **502** (1998), 123–140.
- [5] *W. Bruns and J. Herzog*, Cohen-Macaulay rings, Cambridge University Press, 1993.
- [6] *J. Eagon and V. Reiner*, Resolutions of Stanley-Reisner rings and Alexander duality, *J. Pure and Appl. Algebra*, to appear.
- [7] *S. Goto and K. Watanabe*, On graded rings, II (\mathbb{Z}^n -graded rings), *Tokyo J. Math.* **1** (1978), 237–261.
- [8] *J. Herzog, V. Reiner, and V. Welker*, Componentwise linear ideals and Golod rings, preprint, 1997.
- [9] *M. Hochster*, Cohen-Macaulay rings, combinatorics, and simplicial complexes, *Ring theory II* (B. R. McDonald and R. Morris, eds.), *Lect. notes in Pure and Appl. Math.*, No. 26, Dekker, New York, (1977), 171–223.
- [10] *A. Postnikov, B. Shapiro, and M. Shapiro*, Ring of Chern forms on flag manifolds and forests, Extended abstract (paper in preparation).
- [11] *G. A. Reisner*, Cohen-Macaulay quotients of polynomial rings, *Adv. in Math.* **21** (1976), 30–49.
- [12] *I. Z. Rosenknop*, Polynomial ideals that are generated by monomials (Russian), *Moscow. Oblast. Ped. Inst. Uw cen Zap.* **282** (1970), 151–159.
- [13] *B. Shapiro and M. Shapiro*, On algebra generated by Bott-Chern 2-forms on SL_n/B , *C. R. Acad. Sci. Paris Sér. I Math.* **326** ser. 1 (1998), 75-80.
- [14] *J. Stückrad and W. Vogel*, Buchsbaum rings and applications, Springer, 1986.
- [15] *B. Sturmfels*, The co-Scarf resolution, to appear in *Commutative Algebra and Algebraic Geometry*, Proceedings Hanoi 1996 (D. Eisenbud and N.V. Trung, eds.), Springer Verlag.
- [16] *D. Taylor*, Ideals generated by monomials in an R -sequence, thesis, Chicago University, 1966.
- [17] *W. Vasconcelos*, Computational methods in commutative algebra and algebraic geometry, Springer, 1998.
- [18] *C. Weibel*, An introduction to homological algebra, Cambridge University Press, 1994.
- [19] *G. Ziegler*, Lectures on Polytopes, *Graduate Texts in Mathematics* **152**, Springer, 1995.

current address: Department of Mathematics, Duke University, Durham, NC 27708

e-mail: ezra@math.duke.edu



On Rejection

...for centuries of an elegant geometric language

Gary Harper

My paper, *Meaning-Imposers versus Meaning-Derivers*, was first rejected by the esteemed *American Journal of Physics* in January, 2008. “Read our editorial policy” the editor wrote, “our readers are not interested in a new interpretation of Geometric Algebra.” (Quoting from memory, since I had rejected his rejection letter.) New interpretation, sir?!—the whole point of the paper was to avoid *imposing* an interpretation, thereby *deriving* “the keystone of the entire structure of mathematics” to echo the hero of the paper.

The editor had the grace to say that he considered himself an associate member of my “organization” (his quotes, referring to my *Institute for Nagging Doubt* organization, unquoted, no less serious than the *Rejecta Mathematica* organization, unquoted.) I happen to admire this particular editor, enough not to embarrass him by name, and would be proud to have him as a full member, if he would actually read my paper. (Okay, he had a point: it contained mostly mathematics, not physics, even tho it did have *torque* in it, and *linear force*, *angular velocity*, not to mention *balance point*.)

The second rejection occurred in February from the august *American Mathematical Monthly*. This time my paper received a good editorial review, as I know from the kind rejection letter and from concurrent hits on my website where the ideas are leisurely developed. “I have reviewed your submission in detail with our editorial board” the rejector wrote, “and we have reluctantly concluded that it does not have broad appeal to our diverse audience. We have only a limited amount of space available each month, and are forced by the enormous volume of submission to reject many fine papers.” (Quoting from memory.)

Where did I go wrong? Perhaps I should not have called professional mathematicians, in the first paragraph, “meaning-imposers” who generate “inconsistencies and confusions”. Perhaps I should not have asked “Where did I go wrong?” on page 61; or, “Have I made another blunder?” on page 62. Perhaps I should not have given equal time to the blunders of my Geometric Algebra heroes, especially the living ones; or should not have written, “*points have not yet become full-fledged geometric objects*, like scalars!” on page 68. Maybe I used too many exclamation marks.

Please tell me, sirs, what to change because this paper is dead serious. It is an attempt to introduce the reader to the very expressive geometric language that germinated in the fertile young mind of Hermann Gunther Grassmann in the early 1800s. Altho we have recently understood a good half of his language, the other half, which is just as good—or perhaps better since it is the foundation—remains unknown because it seems strange. But it really isn’t, what *really* is strange is the perverse historical trajectory that makes it seem strange.

Just by reading this paper you have no hope of becoming articulate with Grassmann’s full language. For that you will of course have to also read his two books and play with the ideas. And good luck with that—it took me a good ten years to really understand his fundamentals, and then another good ten years to make them cohere in my mind. So, if this paper succeeds in its purpose, you will have twenty good years in front of you.

Meaning-Imposers versus Meaning-Derivers

Gary Harper

The Geometric Algebra community has evolved into a large segment of meaning-imposers and a tiny segment of meaning-derivars. Meaning-imposers begin with an abstract mathematical formalism, viewed as a gift from heaven regardless of how it had actually been achieved historically; and then *impose* on it whatever geometric meaning seems convenient or appropriate. Such meaning, as many readers know, is called an *interpretation* informally, or a *model* formally. Meaning-derivars, as few readers may know, begin with geometric meaning, viewed as the primitive starting point, and then *derive* everything else from that, *including the mathematical formalism itself*.

Meaning imposition has been undeniably fruitful, but it generates subtle inconsistencies and confusions that have stalled us in the purely *free* Geometric Algebra. Hence we cannot articulate *bound* things, like points for example, except by imposing clever but clumsy artifices on the free language from *outside* it. Nearly two centuries ago Hermann Grassmann showed how to articulate bound things from *inside* the language,¹ but he suffered from the distinct inconsistencies and confusions of a creator who has not had time to polish his creation. By starting over and carefully re-deriving his *full* language from seminal geometric concepts, we can dispel the fog and gain a more expressive language.

Here are the seminal ideas: (1) the concept of geometric points, (2) what it means to summarize points, and (3) what it means to extend something from a point, to wit: (1*) Points have fixed distances among themselves. (2*) Summarizing points is like summarizing anything: order doesn't matter; grouping doesn't matter; a point summarized with Nothing is just the point itself. Finally, Grassmann's gem, somewhat polished: (3*) extending something from a point sweeps it from there directly back to its original position, filling in as it returns, which increments dimension. Hence, to begin at the beginning, extending a *point* from another point produces a directed line segment that has a dimension one higher than that of a point.

(So, clearly, Grassmann was the founding meaning-deriver; but he fell under the seductive spell of mathematical abstraction, and became a resolute meaning-*remover*. Since geometric meaning had already generated his symbolism, Grassmann never could have become a bona fide meaning-imposer. Such persons arrived later, after Grassmann's resolutely abstract symbolism had been cleaned up and unified by William Kingdon Clifford.²)

At this juncture, meaning-imposers will ask what the three primitive concepts are, if not preliminary meaning imposition. Point well taken—we meaning-derivars are closet meaning-imposers; but we are *timid* ones who impose meaning only at the very beginning, before any symbolism has been established, and not just when it seems convenient or appropriate. If meaning ever comes to seem *overwhelmingly* convenient or appropriate, we go right back to the closet and start all over again, convinced that we, in our naivety, have neglected something important *that will change the symbolism*.

Well then, a meaning-imposer may say, I, the meaning-deriver, will shortly need to return to the closet if I expect to have free vectors in my language. There is just no way that *roving* directed line segments can ever be *derived* from *fixed* points!—the roving idea has to be *imposed*. Again, point well taken—it does seem implausible that securely *bound* things, unassisted, could ever produce something *free*. But let us just see if it might be true.

The primitive idea of point summary immediately generates some symbolism; and it looks exactly like the rules for elementary-school *addition*, applied unfamiliarly to points, with *zero* taking the role of Nothing. But it also looks like the rules for *logical or*, again applied unfamiliarly to points, with *false* taking the role of Nothing. So, which will it be?

One can't be sure immediately because, as mentioned, point summary is still unfamiliar, even these several centuries after Grassmann (and Möbius) introduced it. Indeed, it is commonly understood to be, for example, “non-geometric”; and it “makes no intrinsic sense”.³ Here is where an eager young meaning-deriver will probably have to return to the closet and start all over again because she, in her naivety, will decide that summarizing two points, \mathbf{a} and \mathbf{b} say, produces the *midpoint*, \mathbf{m} , between them. What could be more natural?—this immediately establishes summary indifference to order, the *commutative law*. But it also implies that summarizing point \mathbf{a} , say, with itself simply reproduces point \mathbf{a} . This is clearly a kind of geometric *logic*, devoid of numbers; not arithmetic, wallowing in numbers.

With that satisfying thought, the young meaning-deriver begins to investigate how this relates to the primitive idea of fixed points. Whoa! *That* idea applied to the midpoint idea would invalidate summary indifference to grouping,⁴ the *associative law*. Back to the closet: midpoint \mathbf{m} summarizes *two* points, so it should have *twice* the significance of either *one* of them alone. So this is not *geometric logic*, it is *geometric addition*, and the symbolism now becomes $\mathbf{a} + \mathbf{b} = 2\mathbf{m}$ and $\mathbf{a} + \mathbf{a} = 2\mathbf{a}$, rather than naked \mathbf{a} , as before. This new symbolism now provides enough information to validate the associative law, and all the other laws of summary. (Terminology: \mathbf{m} is a *location*; 2 is its *weight*, which is a kind of *magnitude* like length, area, and volume. The weight of a sum point is the sum of its summand weights. A naked location like \mathbf{a} is called a *simple point*, or a *unit point* since $\mathbf{a} = 1\mathbf{a}$.)

To summarize, *points must be weighted* for the first two primitive concepts to be validated together—points must wallow in numbers. This is your very first *derived meaning*; and it may appear insignificant until you consider that it seems manifestly contrary to hoary Euclidean convention, which denies points “magnitude”.⁵ In retrospect, it is clear that this ancient convention must be misleading, at best: distances among points is *all about* numbers simply because distance is an ordered continuum; and summary of points should somehow cause points to inherit that continuum.

Your second derived meaning may not seem so insignificant: *a sum point always lies on the line thru its two summands*. This arises directly from previous equation $\mathbf{a} + \mathbf{b} = 2\mathbf{m}$, where \mathbf{m} is still the *midpoint* (to keep the commutative law valid); hence it necessarily *lies on the line thru \mathbf{a} and \mathbf{b}* , tho it now has a weight of 2. You can use this equation to approximate any sum of two weighted points as closely as you wish by adjusting weights to express *midpoints of midpoints*, iterated; and such an approximation becomes exact in the limit. *Midpoints of midpoints* necessarily generate a result lying on the line thru the two summands.

What other derived concepts do the first two primitive concepts mandate? When you play with these concepts as Grassmann did,⁶ you soon discover that they require a sum of two weighted points, $aa + bb$ say (having scalar weights a and b), to obey this simple rule:

$$\text{weight-distance}(a) = \text{weight-distance}(b)$$

This means that the weight of a times its distance to the sum point equals the weight of b times its distance to the sum point. This is your third derived meaning, and it is quite significant because it tells *exactly* how summary of points causes them to inherit the distance continuum. The rule is just as valid for negative weights as for positive ones provided you carefully distinguish signs as follows: give a summand-to-sum distance that crosses the other summand the opposite sign to a distance that doesn't.

The weight-distance rule induces the following intuitive concept of point summary: a sum point is physically a *balance point* so it always lies nearer the heavier summand (the one with greatest absolute value). When this idea is applied to a sum of points having opposite signs, the opposite-sign distinction kicks in, requiring the sum point to lie on the line *from* the lighter point (in absolute value) *thru and beyond* the heavier point. Which prepares you for...

Some magic: what is $a - b$?

Whatever the *location* of this sum, call it l , its weight, $1 - 1$, is 0; so this sum has this form: $0l$. Since anything multiplied by zero is just zero, the sum $a - b$ is clearly zero. Right?

There is a quick way to test this: give the entire sum a non-presupposing name, “ v ” say, meaning that $v = a - b$ (“ T ” was presupposing since it was a *location*). Now see if indifferent v acts like zero: Add v to the second nearest thing in sight, namely point b . When you apply the primitive rules of point summary, you get point a . Hmmm... Okay then, *subtract* v from the nearest thing in sight, point a . You get point b .

That is not how zero acts!—zero doesn't *change* things when it is added or subtracted with them. This v thing is *changing* points under addition and subtraction; in fact it is *moving them around*. Where did I go wrong reasoning that $a - b$ must be zero? I went wrong in assuming there was some location, “ T ” I called it, for its zero weight to multiply. There's not. The point sum $a - b$ really *has no location*. It really *has no weight*. And yet it is *not zero*. It is truly bizarre to the modern mind, which has come to shun point summary, altho many minds born in the 1800s were comfortable with it. Let's reacquaint ourselves with their old friend:

You can sneak up on $a - b$ by approximating it with non-zero weights that approach zero. For example, start by giving a a weight of $1/2$, and then successively halve a 's actual-minus-approximate weight like this: $1/2$, $3/4$, $7/8$, $15/16$, etc. This will successively halve the approximate *sum* weight. At each weight *halving*, the weight-distance rule will scoot the approximate sum location *twice* as far away along the line thru a and b .

This removes some of the mystery: as the sum point *weight* goes to zero, its *location* goes to infinity, in *lock-step*; so the diminishing weight and the receding location effectively cancel each other. Which is why $a - b$ is not zero: it is actually a peculiar kind of *zero times infinity*. The satisfying conclusion is that $a - b$ is a point at infinity. Right?

This is certainly a modern concept, quite familiar from Projective Geometry, which is redolent with classically imposed points at infinity. But I just said that the result of $a - b$ really *has no*

location. How can it not have a location if it *resides at infinity*? Have I made another blunder? Or is this just an innocuous problem with our language?

There is an easy way to test this: start \mathbf{a} with a weight of $1\frac{1}{2}$ (rather than $1/2$) and then sneak up on $\mathbf{a} - \mathbf{b}$, as before. As before, halving the approximate sum weight scoots the approximate sum location twice as far away. *But it does so in the opposite direction*. This is again mandated by the weight–distance rule, which carefully notices, during the approximation, which summand is lighter, and which is heavier. In consequence, since both approximations approach $\mathbf{a} - \mathbf{b}$ in the limit, it appears that *this sum is infinitely distant from itself!*

However, if this sum *really* has no location, then the problem disappears because such a sum cannot be *any distance* from *anything*, let alone from itself. But if it “resides at infinity” then there is a problem with our language, and it definitely is *not* innocuous. It generates the subtle confusion that Geometric Algebra directly articulates “points at infinity”. The full Geometric Algebra does not. It cannot. *It can articulate only finite representations*. Moreover, *such a finite representation cannot be a single non-decomposable thing*—it is *intrinsically* composite. (Fore-shadowing query: Indifferent v therefore cannot be a sum, as naively assumed, so what is it?) To peek ahead, there is no “at” at infinity; rather there are “ats” at finity.

You have just seen that v , under addition and subtraction, can move the two points that compose it, \mathbf{a} and \mathbf{b} . Look closely: from elementary-school rules of addition, $v + \mathbf{b} = (\mathbf{a} - \mathbf{b}) + \mathbf{b} = \mathbf{a} + (\mathbf{b} - \mathbf{b}) = \mathbf{a} + 0 = \mathbf{a}$. So point \mathbf{b} has effectively been carried from one end of v to the other end. And the reason is clear: under addition, \mathbf{b} annihilates one of v 's endpoints, *poof*, leaving the other endpoint as residue. It seems natural to call the *poofing* endpoint the *tail*, the residual endpoint the *head*; and say that $v + \mathbf{b}$ carries point \mathbf{b} from v 's tail to v 's head. Altho this nomenclature seems natural, One wonders how generally useful it might be since this obviously works only because v is being added with a copy of its *own* endpoint. Right?— v doesn't carry *other* points around under addition, does it?

Well, let's just see: given an *arbitrary* point \mathbf{r} , what is $v + \mathbf{r}$? To ask this question in the fresh young symbolism, solve this equation: $v + \mathbf{r} = x$, where x is unknown, *utterly* unknown as indicated by its generic font.

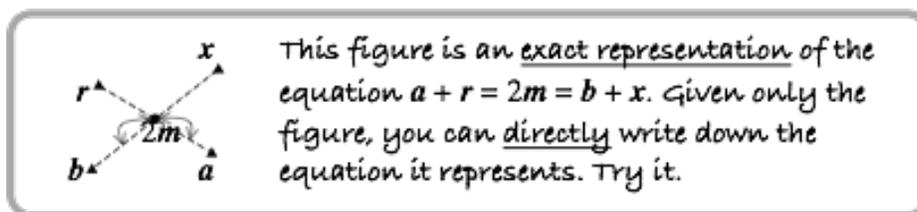
What is going to happen next is so important to your understanding of the full Geometric Algebra that I am going to present it in complete detail, with **powerful** emphasis on the crucial part. If you hope to acquire a more expressive geometric language then you will have to wrestle with this until you understand it completely. How will you know whether you've understood it? If my experience is any indication, *you will become amazed*. If you don't, then you may be suffering from traditional meaning-imposing habits. To help overcome that, remember that we are articulating *fixed points*, and nothing else. We began *bound*; we are *bound* now; and it looks like we will *stay bound* because we are too timid to cavalierly *impose* any kind of geometric freedom. If freedom arises, it will be *entirely derived* from things that are *entirely bound*. Who would ever bet on that?

Okay, expand equation $v + \mathbf{r} = x$, giving $\mathbf{a} - \mathbf{b} + \mathbf{r} = x$. Now pull the purely-positive-equation trick by putting \mathbf{b} on x 's side of the equation: $\mathbf{a} + \mathbf{r} = \mathbf{b} + x$. The left side has the sum of two simple points. The right side has the sum of a simple point and *something*, namely x . For the right side to equal the left side, this *something* must also be a simple point (do the weight calculation), so denote it in point font, x . Hence, utterly unknown x has become somewhat known *simple point* x . So, apparently v really does move arbitrary points around since that question will be an-

swered by x 's eventual location. To find this location, we need to visualize the transformed equation: $a + r = b + x$.

This equation involves *addition*, *equality* and *simple points*. These are the elements that have to be displayed geometrically. *Addition* of two points can be indicated by a dashed line connecting them. The *equals sign* is too imprecise about location to be useful on a geometric figure. Instead, a skinny curved line with tiny arrows on each end will be used. Call this the *geometric equals sign*. Its two tiny heads will just touch the things that are equal. *Simple points* are so useful that they should be distinguished from generic weighted points; let's use a little triangle for them and a little dot for generic points.

With these conventions, the transformed equation becomes geometrically obvious: two little triangles connected by a dashed line denote addition of two simple points, so their sum, $2m$, lies at their *midpoint*. There are two of these additions connected by equality, so they share the *same* midpoint sum. Here is a picture:



Visualizing simple point sums.

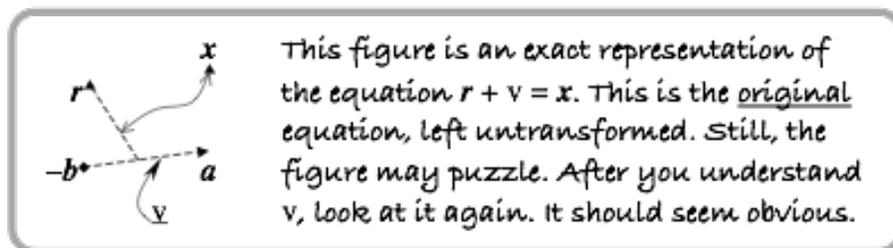
This kind of figure has seminal importance so let's dignify it as an *X-diagram*. It answers the question about v 's ability to move points other than the two it comprises: x is the solution to the original equation $v + r = x$. In words, v added to r moves that over to x . Since the location of r is entirely *arbitrary*, v moves *any* simple point in a similar way. Here is how to be sure you understand this completely:

Sketch weighted points a and $-b$. Connect them with a dashed line to indicate that they are being added together. This makes $a - b$ a kind of *cohesive bundle* (hint), deserving its own name, v ; and deserving more intimate notation: $a-b$. Have a friend sketch a point r somewhere—*anywhere*. Now add your little bundle to it like this: **Wham!**—equate the summary result to x so you have a concrete result to work with. **Bang!**—unbundle v and swap $-b$ to x 's side of the equation (by erasing the dashed-line addition and the minus sign). *This gives two simple point sums that equal each other.* **Pow!**—do the sum you immediately know, namely $a + r$. This gives midpoint $2m$, which is *also* the sum you didn't immediately know, namely $b + x$. So now you know it too, and you therefore know where x is. (For graphical precision, you should, of course, sketch the dashed-line additions as you do each sum, thereby making them neatly X d together right in the middle.) Next, have your assistant sketch a different point somewhere, *anywhere*, **pow!** Another: **pow!**... If you can do ten distinct $v+r$ sums in a row, correctly, without batting an eye, then you understand this. *Please* understand this—it is really quite simple; but the main reason we still suffer from geometric inconsistencies and confusion is nearly universal chronic ignorance of its various unexpected consequences.

Having understood this, you may think that it does not seem amazing. But it might seem surprising, or at least *peculiar*: Recall that v was able to move a copy of its own endpoint under ad-

dition by *poofing* it (to speak technically), leaving the other endpoint as residue. But here v is moving an arbitrary point under addition by a kind of *scissoring mechanism*, the X-diagram, in which nothing is being *poofed*. And yet...

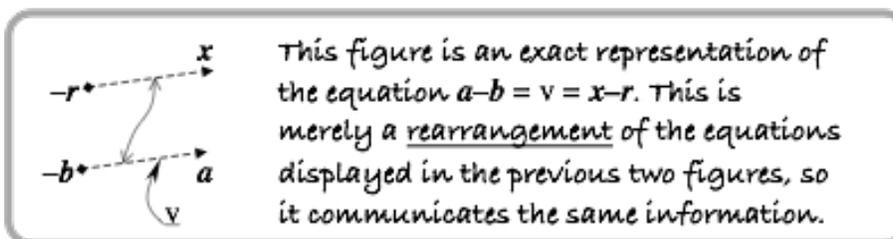
And yet v is moving point r exactly as tho v could move parallel to itself to place its tail over r , and then do the addition; in which case v would be *poofing*, exactly as before. Dust off your high-school geometry and gaze at the X scissoring mechanism until you understand this in your bones. (You see the abm triangle, congruent to the rxm triangle, don't you?—take it from there.) To help understand this, here is the previous figure, *exactly as before*, geometrically, with the *exact same equation*, except that it has been left untransformed:



A direct visualization of $r + v = x$.

This figure may not seem as geometrically obvious as the X-diagram. Nevertheless, it directly represents the *same* equation, left untransformed. To understand that, superimpose the previous X-diagram on top of it—corresponding points will match exactly. The X-diagram should seem obvious if you understand how simple points add. When you understand how v acts, this figure should also seem obvious. Notice that the sum $2m$ has disappeared from it because this point was not present in the original equation—it served as an illuminating *centered pivot point* in the transformed equation that you may now discard as a conceptual crutch. Notice also that, altho simple points are still depicted with tiny triangles, a *negative* simple point has been depicted with a tiny square. This makes it easy to determine the head and tail of v , as you see. Speaking of v , let me say again, for emphasis: in this figure it looks as tho v is able to move parallel to itself to engage the *poof* method of addition. If so, any puzzlement about the figure would evaporate.

It turns out that v actually is able to move parallel to itself; and a different transformation of the original equation will directly display this. In this new equation, don't unbundle v . Instead move r over to the x side of the equation: $v = x - r$. **Kazaam!**—the right side of this equation has *exactly the same form as* v does, namely a simple point minus another simple point. Here is the figure representing this re-transformed equation:



geometric magic

To see that it communicates the same information, superimpose the X-diagram on top of it—corresponding points will again match exactly. In consequence, $a-b$ equals $x-r$. Hence, because the location of r is completely arbitrary, v , in effect, is able to move *anywhere* parallel to itself. In short, *geometric freedom has been entirely derived from things that are entirely bound. That, I submit, is amazing.* Grassmann's protracted account of discovery seems to indicate that he found it amazing too.⁷

Apparently the strange *un-point-like* point v has a fixed *separation* and a fixed *direction* but no particular location. But exactly *what* has separation-and-direction? And exactly *what* has no location? Are these the *same thing*? No: the *summands* of v , considered as a separated-directed whole, can reside *anywhere* because the *sum itself* resides *nowhere*.

To be precise: v is a peculiar addition of oppositely weighted points whose *sum loses, in lock-step, both magnitude and location, which, in the limit, makes its summands gain both separation and direction.* It's magic—without magic, to borrow John Wheeler's aphorism. This raises some...

Perplexities

- Roving v acts much like a conventional vector (so the name “ v ” was deviously presupposing), except that *it is not a line segment.* Right?...
- I mean, addition of weighted points always produces another *point*, at least *formally*, doesn't it?—addition never *changes dimension*, does it? (Everyone knows that a point is zero-dimensional, a line segment one-dimensional.)
- Speaking of changing things, the previous magic *changed focus* in the limit from a *sum* to its *summands*. Is this distinction important *formally*?
- For example, is this what makes v *intrinsically composite*?
- Is that why the full Geometric Algebra always articulates things *at finity*?
- If so, do the *conventional rules* of Geometric Algebra make the sum–summand distinction properly?
- If not, should they? Could they?

Some of these perplexities may be superfluous. After all, extending a point from another point produces a *directed line segment*, which would be *exactly a conventional vector* if it turns out to be as mobile, and as mobilizing, as v turned out to be; and that now seems likely, doesn't it? But that would raise another perplexity: what exactly, then, is the distinction between *subtracting* a point from another point, and *extending* that point from the other point? To find out requires specifying the...

Relationship between extension and addition

The relationship, like all healthy ones, is built on *mutual respect*: (4a) Addition respects extension enough not to change the properties that extension hath wrought. (4b) Extension respects addition enough to treat addition's *result* as a genuine summary of its *arguments*. These new concepts require a trip back to the closet to start all over again. Fortunately they merely augment the first three concepts, so, once we understand how they *change the symbolism*, we can just pick up where we left off.

Addition does not change the properties that extension hath wrought. This is just an unfamiliar specialization of the concept of summary. One naturally expects a summary *not to change any properties* of what it summarizes; to do otherwise would make a mockery of the concept. Refusal to change properties is really more central to the idea of summary than indifference to order, or indifference to grouping; but it remains unfamiliar because of our present narrow experience with summary of *homogeneous* things, ordinary numbers. Extension generates *inhomogeneous* things, things having *different dimensions*; and *that* is the property that addition does not change. Which removes one perplexity: (4a*) *Addition does not change dimension.* (This is a *semantic* formality that is surprisingly tricky to symbolize properly, as you shall see.)

And that, in turn, begins to remove the perplexity about v being intrinsically composite. When addition is presented with summands having different dimensions, it can't summarize them to anything simpler because it can't give them a common dimension. Hence, it merely bundles them into a summary list, with the plus sign serving as conjoining punctuation. This bundle is *intrinsically composite* because (1) its contents cannot reduce to just one thing, and (2) it can always be decomposed and re-bundled differently, using addition's associative law.

It is obvious to any student of Geometric Algebra that a sum of things having different dimensions is intrinsically composite: these things are obviously *too distinct* to merge in summary. But it is almost always a surprise that sometimes—*oftentimes*—even things of the *same dimension* are that distinct. This surprise can be blamed on a historical mishap: we have become stuck in the purely free Geometric Algebra. In that language, all *readily imaginable* things of the same dimension can always sum to a single thing simply because imaginable space just happens to be a perfect cage for free things.

If you go *just* beyond imaginable space, however, you bump into free things of the same dimension that are *too distinct* to sum to a single thing. Free bivectors in free 4-space, for example, are that distinct if the planes thru them intersect in just one point. This possibility arises naturally from the extra dimension (and has obvious expression in the *full* language), but it seems so bizarre to most students that they dismiss the idea of intrinsically composite same-dimensioned sums as too esoteric to worry about. Even Grassmann may have had that attitude during his early “geometry” phase,⁸ as he dismissively called it.

If he did, he certainly revised his opinion after he encountered bound things late in his explorations. Grassmann began his language like all students today begin it, purely free; and perhaps humankind's roving spatial experience makes this approach natural. But his incredible curiosity and creativity eventually introduced him to *bound* points via *free* vectors!⁷ This is *exactly backward* logically; and it is truly, *astonishingly*, extraordinary as witnessed by the fact that in nearly two centuries of ignorance about Grassmann's bound language, no one else has made the trip backward and installed points within the formalities like Grassmann did. When he did that, he quickly discovered that there are *readily imaginable* same-dimensioned bound things that are too distinct to sum to a single thing.⁹ They are not esoteric at all—in fact they are more common than same-dimensioned things that *can* sum to a single thing. (To peek ahead, they aren't simple points, are they?—they always sum to a *single* thing.) To really understand this, you need to know more about extension.

And that requires notation. Extension had initially been denoted in about as many different ways as there were authors writing about it—Grassmann himself used several distinct notations—but it has recently stabilized on Cartan's wedge, \wedge , meaning *extended to*. Unfortunately,

that has two serious problems in the *full* Geometric Algebra, which must, above all, articulate points well since they are the generative elements:

(1) $a \wedge b$ would generate a directed line segment with tail at a , head at b ; which is opposite to $a - b$, which generates separated-directed points with tail at b , head at a . This inconsistency would be confusing of course, but the worst of it is that these two expressions have an elegant relationship (coming up), *fundamental to the full language*, that would be obscured if they did not have their heads and tails in the same order. This really begs for extension to be *from* rather than *to*, which somewhat polishes Grassmann's gem. Consistency with point subtraction prompted Hamilton to adopt a similar convention.¹⁰

(2) From item 1, extension is clearly *directed*, so it really should have a *directed* symbol, rather than one with bilateral symmetry like \wedge . How about \blacktriangleleft ? This clearly indicates *from*, and its *filled-in* form indicates *extension*. Hence $a \blacktriangleleft b$ is " a extended from b ", like $a - b$ is " a subtracted from b "; and these two expressions have their ducks aligned. As a bonus, this distinct notation should help clarify the transition from the conventional purely free language to Grassmann's full language for those readers crossing that bridge.

With notation established, we can pick up where we left off: *Extension respects addition enough to treat addition's result as a genuine summary of its arguments*. Which is to say, extension with a point is indifferent to whether it operates on addition's *arguments*, or on addition's *result*. Here is how this augments the symbolism:

$$(4b.1^*) \quad (A \blacktriangleleft c) + (B \blacktriangleleft c) = (A + B) \blacktriangleleft c \quad \text{and} \quad (c \blacktriangleleft A) + (c \blacktriangleleft B) = c \blacktriangleleft (A + B)$$

Notice that there are two rules, commuted, because extension is *directed*, so extending *from* c on the left is generally different from extending *to* c (reading backward) on the right. Mathematicians call these rules *distributive laws*, which focuses on *syntax*. This may seem appropriate since these rules *are* part of the syntax, as explained shortly; but they, like all rules in this paper, were motivated by primitive semantics, so this paper will call them extension's *respect for summary* to emphasize their meaningful origin.

When you apply these rules multiple times to *scalar-weighted* points via a valid limiting process you get, for *scalar* c :

$$(4b.2^*) \quad c(a \blacktriangleleft b) = (ca) \blacktriangleleft b = a \blacktriangleleft (cb)$$

This will be called extension's *respect for multiple summary*, again focusing on geometric meaning (and it can also be generalized to generic A and B). To indicate that extension has both kinds of respect, let's say that it has *strong respect for summary*. Strong respect for summary makes the language versatile and expressive by decoupling extension from addition, and from addition's infinitesimal multiple limit, scalar multiplication. This prepares you to start...

Extending things

Here is where timid meaning derivation begins to really pay off. A meaning-deriver has to start with the primitive concept for extension, the third concept, which requires extension from a point to *increment dimension*. That, in turn, requires establishing the *primitive dimension*, the dimension of a *point*.

Meaning-imposers long ago agreed that a point is zero-dimensional, but that poses a serious conundrum for meaning-derivers: Since extension from a point *increments dimension*, shouldn't

the dimension of a point therefore establish the *dimensional increment* that gives everything else a dimension? This is simply a natural requirement for the dimension of an extension result to be the sum of its argument dimensions. If so, then *points must be one-dimensional*. This would imply that line segments are really *two-dimensional*; patches of plane are *three-dimensional*; and so on. This seems silly—we have known for millennia that lines are one-dimensional, planes are two-dimensional, and on up. Nevertheless, in a last-gasp nag, the meaning-deriver asks, *what about on down?*

If points were one-dimensional, then scalars would be zero-dimensional. Suddenly it is the meaning-imposers who have a serious conundrum: They have recently reached universal agreement that scalars are indeed zero-dimensional. If *points* were zero-dimensional, as *also* agreed, then, by dimensional decrement scalars would have a dimension of minus one. (Unless points *are* scalars. Well, are they? If meaning-imposing habits incline you to think so, please ponder the elegant relationship (coming up) between points and scalars before deciding.) Here you have an example of the inconsistency that meaning imposition generates. Exposed like this, *zero is not minus one*, it doesn't seem subtle, does it?

Why haven't we noticed this problem for the last several thousand years? First, only recently have scalars acquired a dimension, when they were belatedly recognized to be full-fledged *geometric objects* like lines, planes and so on. When scalars interacted with geometric things, it was seen that they must have a dimension of zero because they do not change the dimension of what they multiply. Strangely, second, points have yet to be emancipated like that—*points have not yet become full-fledged geometric objects*, like scalars! Meaning-imposers have so far *refused to allow points into the formalities* alongside scalars, vectors, etc; except as outcasts, undesirables who are denied full geometric rights. It is tempting to blame this on Euclid, who refused to grant points “magnitude”, which effectively exiled them to the interpretation where, third, they have been neglected, orphaned from their geometric family, and underfed to the extent that they literally have no weight at all. *Exile to the interpretation*—let David Hilbert describe that:

“One should always be able to say, instead of ‘points, lines, and planes’, ‘tables, chairs, and beer mugs’.”¹¹ Well, lines long ago managed to escape from Hilbert's beer hall by dressing up as vectors, able to participate in black-tie formalities. Planes have recently pulled off the same formal getaway by dressing up as bivectors; but points are still stuck in the pub in their underwear. Since they are, who really cares what their dimension is? Apparently it is very much like the dimension of a table, or perhaps a beer mug?—who cares? Meaning-derivars care, and they want to get the orphan point out of the unruly interpretation and into the ruly formalities alongside its geometric kin: scalars, vectors, bivectors and so on. Transition into the formalities has been a paradigm for mathematical progress for thousands of years; but unexpectedly, for points it will require, *gazook, meaning inside the language, formal semantics.*

Which, for distinction, requires *formal syntax*. Some of this syntax has already been presented: it is just the collection of *conventional rules* for Geometric Algebra—the equations, like the commutative, associative and distributive laws, that serve as axioms. These equations establish the valid sentences in the language.

All the rest is semantics, which traditionally—*dogmatically*—has resided almost entirely within the mind of the person composing the sentences. That turns out to be woefully inadequate for the *full* language, where *bound* points generate *free* things. The important *formal* distinction between bound and free requires *formal* semantics because the syntax intentionally ignores the

distinction, for good reason. Moreover, such semantics rest, in an unanticipated way, on *formal* dimension, which, because it cannot be defined by equations, is also part of the semantics.

Formal dimension presents a rare opportunity to please everyone. To distinguish it from the previous decidedly *informal* dimension, give it a distinguished name: *extent*, which means *number of points required in an extension*. This will please the meaning-imposer since a line segment obviously requires *two* points, so it has extent two; a patch of plane requires *three* points, so it has extent three; and on up. Certainly, a meaning-deriver is pleased because this gets the *foundational* dimensions right: a point requires *one* point in the trivial do-nothing extension, so it has *extent one*; a scalar requires (dare I say) *zero* points, so it has *extent zero*. The meaning-imposer might be doubly pleased to discover that formal extent, in its intrinsically *separated* form, automatically articulates conventional dimension. Hence, conventional geometric dimension is not wrong, it is just a *special kind* of dimension.

To be specific, addition in the full language makes a distinction it could not have made with points absent, namely the distinction between the *separated* extent of free things, and the *filled-in* extent of bound things. This distinction is definitely part of the semantics because the syntax—the *conventional rules* of Geometric Algebra—simply cannot make it. To begin understanding that, investigate filled-in extent from the beginning:

Extending a point from another point produces a directed line segment that has a dimension one higher than that of a point. Start formalizing this by expressing it in the young symbolism: $a \blacktriangleleft b$.

Now proceed to formalize dimension by making extent an operator that accepts an argument, so that, for example, $\text{extent}(a)$ produces $\{1\}$ since a is a simple point. Curly braces indicate a *list* of extents, necessary because $\text{extent}()$'s argument might be intrinsically composite. For example, $\text{extent}(b + a + a \blacktriangleleft b)$ produces $\{0, 1, 2\}$ if simple point a has a different location than b . (If these points had the same location, the extension would produce Nothing with extent $\{2\}$ (a line segment with no length); in which case $\text{extent}(b + a + a \blacktriangleleft b)$ would have been $\{0, 1\}$. Such potential disappearance is just one of the reasons $\text{extent}()$ is semantic—discovering disappearance requires computation.)

So, extension from a point increments extent, as required; and what does *addition* of two points do?—what, for example, is $\text{extent}(a + b)$? You already know: since $a + b$ generates a single thing, and since addition does not change dimension, this extent must be $\{1\}$. Which brings up a subtle but very important point: because $a + b$ generates a point with a weight of 2, the $\text{extent}()$ operator clearly ignores weight; so, in general, for *any* generic weighted point xx , $\text{extent}(xx)$ produces $\{1\}$.

In consequence, *weight is not Euclidian “magnitude”*. When Euclid asserted that a point “has no magnitude” he meant that it has no *spatial extent* like a line does, like a plane does, like a volume does... Euclid was asserting, in the technical language of the full Geometric Algebra, that a point has no extent greater than one. This is true: it has precisely extent $\{1\}$, and *this has nothing whatsoever to do with the point’s weight*, which specifies its potential scaling relations with its geometric kinfolk. These distinct concepts have been confused for millennia because there was no terminology that clearly distinguished them. The full Geometric Algebra remedies that by quantifying Euclidian “magnitude” with *extent*, and scaling relations with *weight*, *length*, *area*, *volume*... (Each of which is a formal kind of *magnitude*, un-scare-quoted—see how deviously confusing the vernacular is?)

Now for the conventional-vector question: how does directed line segment $\mathbf{a}\blacktriangleleft\mathbf{b}$ move around? Meaning-imposers long ago asserted that directed line segments are free to move anywhere parallel to themselves; and that seems to have been wildly successful. Nevertheless, a meaning-deriver is not that bold; in fact he is *so* timid that he won't let $\mathbf{a}\blacktriangleleft\mathbf{b}$ move *at all*, unless the primitive semantics allow it. Fortunately, the primitive semantics have already generated things that can move points around; so the meaning-deriver can try moving the endpoints of $\mathbf{a}\blacktriangleleft\mathbf{b}$ to see what happens, like this:

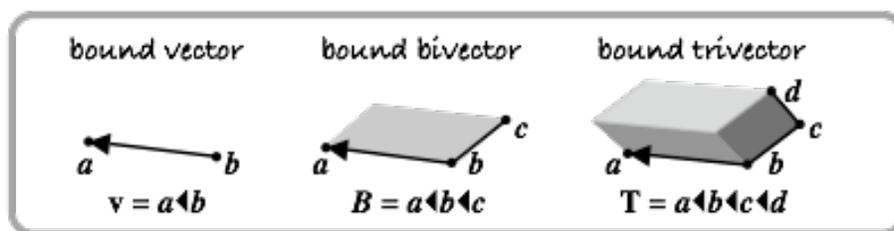
Generate a roving separated-directed pair of points, add it to both \mathbf{a} and \mathbf{b} , and then extend them, thereby translating $\mathbf{a}\blacktriangleleft\mathbf{b}$ parallel to itself. *This parallel-translated version of $\mathbf{a}\blacktriangleleft\mathbf{b}$ will almost never equal $\mathbf{a}\blacktriangleleft\mathbf{b}$* for a simple reason: To be equal to $\mathbf{a}\blacktriangleleft\mathbf{b}$, it would, for starters, obviously have to be expressible entirely in terms of points \mathbf{a} and \mathbf{b} . This is generally not possible because the translator itself is generally not so expressible.

Well then, suppose the translator *were* so expressible. Then it would be a scaled version of $\mathbf{a}-\mathbf{b}$, which translates $\mathbf{a}\blacktriangleleft\mathbf{b}$ somewhere along the line thru itself. In this case, however, extension utterly ignores the translation, thereby making the translated version of $\mathbf{a}\blacktriangleleft\mathbf{b}$ equal to it. This is mathematical poetry arising from extension's strong respect for summary, which, in particular, requires that a point extended from itself vanishes. (For pleasure and education, you might compose this simple poetry yourself.)

In consequence, $\mathbf{a}\blacktriangleleft\mathbf{b}$ is *not* a conventional vector, even tho is looks like one (since it is a directed line segment). Rather it is a *bound vector*, bound to the line thru itself, which will be called its *confining space*. Contrariwise, $\mathbf{a}-\mathbf{b}$ actually *is* a conventional *free vector*, even tho it does *not* look like one (it is *not* a directed line segment). It does not look like a conventional vector because it has been *unconventionally disciplined to treat points as bona fide geometric objects*.

This long-overdue discipline prepares you for a hint about the elegant relationship between bound and free: What is a free vector extended from a simple point? For example, what is $\mathbf{a}-\mathbf{b}$ extended from \mathbf{b} ? Extended from \mathbf{a} ? Extended from \mathbf{c} , not on $\mathbf{a}\blacktriangleleft\mathbf{b}$'s confining line? You can easily do the math for the first two questions by appealing to respect for summary (and you'll get the same answer); but to really understand the elegant relationship, you need to make acquaintance with all readily imaginable bound things, and then watch how they generate their free counterparts.

First, extend bound vector $\mathbf{a}\blacktriangleleft\mathbf{b}$, call it \mathbf{v} , from point \mathbf{c} not on $\mathbf{a}\blacktriangleleft\mathbf{b}$'s confining line. This sweeps \mathbf{v} directly back from \mathbf{c} , filling in as it returns, thereby generating *bound bivector* $\mathbf{a}\blacktriangleleft\mathbf{b}\blacktriangleleft\mathbf{c}$, call it \mathbf{B} . This bivector is bound to the *plane* thru itself for the same reason a bound vector is bound to the *line* thru itself: a parallel-translated version of $\mathbf{a}\blacktriangleleft\mathbf{b}\blacktriangleleft\mathbf{c}$ cannot equal $\mathbf{a}\blacktriangleleft\mathbf{b}\blacktriangleleft\mathbf{c}$ unless the translating free vector is expressible entirely in terms of generating points \mathbf{a} , \mathbf{b} , \mathbf{c} . Next, extend bound bivector \mathbf{B} from point \mathbf{d} not on \mathbf{B} 's confining plane. This generates *bound trivector* $\mathbf{a}\blacktriangleleft\mathbf{b}\blacktriangleleft\mathbf{c}\blacktriangleleft\mathbf{d}$, call it \mathbf{T} , which is bound to the linear space thru *itself*. This space happens to model extent- $\{4\}$ physical space (count the points required in the extension), so a bound trivector is a *ceiling* for that space. Here is a picture:



Readily imaginable bound things

Now for the free counterparts to these bound things. You have already seen the free counterpart to bound vector $\mathbf{v} = \mathbf{a}\blacktriangleleft\mathbf{b}$, namely free vector $\mathbf{v} = \mathbf{a}-\mathbf{b}$. Note these two important properties:

(1) Bound vector \mathbf{v} is free vector \mathbf{v} extended from a simple point on the confining space thru the bound vector, *exactly*. You just discovered this if you took the previous hint about the relationship between bound and free. This establishes the elegant relationship between $\mathbf{a}-\mathbf{b}$ and $\mathbf{a}\blacktriangleleft\mathbf{b}$, which removes the perplexity about the exact distinction between them. (Is it now clear why these corresponding vectors should be articulated in the same order?) To dignify the relationship, call free vector $\mathbf{v} = \mathbf{a}-\mathbf{b}$ the *free part* of bound vector $\mathbf{v} = \mathbf{a}\blacktriangleleft\mathbf{b}$.

(2) The free vector is composed of *separate, but otherwise exactly opposite bound things added together*.

The emphasized phrases are *universal attributes* of the free-bound relationship, so it will be useful to ponder them briefly before examining that relationship in detail.

First, to extend free vector \mathbf{v} from a simple point, the simplest strategy is to place \mathbf{v} 's tail right over the point before extending. *Poof*, the tail-on-point part of the extension will vanish because a point extended from itself vanishes. This leaves the head of \mathbf{v} extended from the point, which is just bound vector \mathbf{v} . This is the *poof* method of point extension, even more wonderful than the *poof* method of point addition because it will apply to things of even higher extent.

Second, ponder what it will mean for *separate but otherwise exactly opposite bound things* of higher extent to be *added together*. As with primitive things, it will mean that sum magnitude diminishes to zero as sum location recedes to infinity, which will, in the limit, shift focus from sum to summands. There is a transparent way to demonstrate this: successively extend by the independent free vectors hidden within these higher-extent things. This will automatically produce roving things having separate but otherwise exactly opposite bound ends because free vectors have those properties. As a bonus, it will show that even tho bound generates free, free does not generate bound, which is one reason we are still stuck in the free language. (Being stuck there impels us to persistently try to represent bound with free, typically points with vectors, which is *inherently contradictory* because free *cannot* generate bound.)

Here are the free vectors hidden in the previous readily imaginable bound things: $\mathbf{v} = \mathbf{a}-\mathbf{b}$, $\mathbf{w} = \mathbf{b}-\mathbf{c}$ and $\mathbf{x} = \mathbf{c}-\mathbf{d}$ (gaze at the previous figure).

To begin, extend \mathbf{v} from \mathbf{w} : $\mathbf{v}\blacktriangleleft\mathbf{w} = \mathbf{v}\blacktriangleleft(\mathbf{b}-\mathbf{c}) = \mathbf{v}\blacktriangleleft\mathbf{b} - \mathbf{v}\blacktriangleleft\mathbf{c}$. It would be enlightening to descend further toward points, but there is no real need to do so because you know, from the *poof* method of free-vector extension, that this is a pair of separate but otherwise exactly opposite bound vectors added together. You also know, from freedom of the \mathbf{v} and \mathbf{w} arguments, that there are countless other exactly opposite bound vector pairs equivalent to this one, differing only in loca-

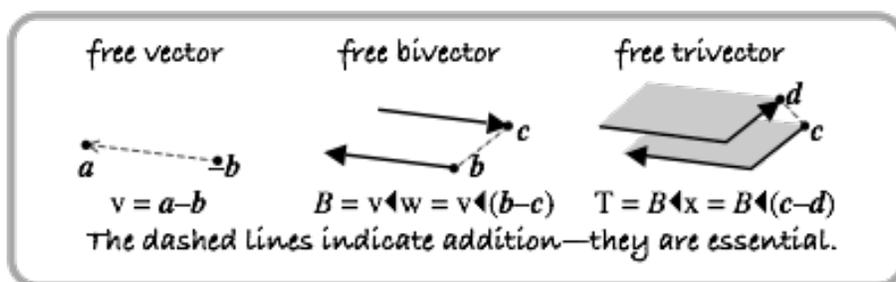
tion, all having the same separation and direction. (This is *area* separation, as dimensionally distinguished from the *length* separation of a free vector; just as *area* magnitude is dimensionally distinguished from *length* magnitude, and so on.)

Here is an intuitively appealing way to understand why these variously located $v\langle(b-c)$ extensions are all equal: incrementally approximate them in unison by sneaking up on free vector $b-c$ as before, but extend at each step. The various bound vector results will dwindle away in lock step as they recede to infinity until they all seem to merge together as a tiny directed dash on the horizon. In the limit this dash loses (length) magnitude and location, which causes its various summand pairs to gain (area) separation and direction. At that limit, the focus necessarily shifts from non-existent sum to existent summands—these summand pairs have suddenly become a *free* bivector, call it B .

Does free bivector B have an elegant relationship with bound bivector B ? Here is where the *poof* method of point extension really shines: To extend free B from a simple point on bound B 's confining plane, take advantage of B 's freedom to place one of its bound vector ends over the point, then extend from there. *Poof*, the summand-on-point part of the extension vanishes because it produces no area (technically: it produces Nothing, *zero*, having extent $\{3\}$). This leaves only the other summand extended from the point—a **bound** bivector that is just a filled-in version of the *free* one with an incremented extent. This is indeed bound bivector B . To dignify this elegant relationship, call free B the *free part* of bound B . Query: do you still get bound B if you put free B 's *other* end on the point before extending? What if you don't place *either* end on the point?

Now here is a curiosity: to get bound B from free B , you had to use something **bound**, namely a simple point. However, the game we are now playing is to extend by *free* vectors hidden within bound things; and so far that has generated something *free*, namely B .

To continue this game, extend free B from the last hidden free vector, x : $B\langle x = B\langle(c-d) = B\langle c - B\langle d$. You could descend further toward points, but there is no need to do so because you know, from a paragraph ago, that this is a pair of separate but otherwise exactly opposite bound bivectors added together. You also know, from freedom of the B and x arguments, that there are countless other exactly opposite bound bivector pairs equivalent to this one, differing only in location, all having the same volume separation and direction. Here is an intuitively appealing way to understand why these variously located $B\langle(c-d)$ extensions are all equal: incrementally approximate them in unison by sneaking up on free vector $c-d$ as before, but extend at each step. The various bound bivector results will dwindle away in lock step as they recede to infinity until they all seem to merge together as a tiny directed patch on the horizon. In the limit this patch loses (area) magnitude and location, which causes its various summand pairs to gain (volume) separation and direction. At that limit the focus necessarily shifts from non-existent sum to existent summands—these summand pairs have suddenly become a *free* trivector, call it T . Here is a picture of it with the free consorts that led to it:



Extension generates free things from free vectors.

By now there should be no need to explain the relationship between free T and bound T by describing how the former is the free part of the latter. In general, it should now be clear that generic bound \mathbf{a} is just its free part a extended from simple point a lying within its confining space. This fills in free a and binds it thru point a . In algebra, $\mathbf{a} = a \wedge a$. This is truly elegant because it *universally* relates free with bound, a relationship that applies *even to the latecomers*, scalars and points. It is an algebraic fact that weighted point aa equals $a \wedge a$.¹² Which is to say, *a weighted point is a bound scalar*, and the typeface emphasizes this, as you may have noticed. So, even tho scalars aren't points, since they have lower dimension, nevertheless, by extension they give magnitude to simple points; and that is the deep reason points wallow in numbers in the full Geometric Algebra. Hence scalars are really, *really*, full-fledged geometric objects.

The previous derivations showed that when extension's arguments are free, its result is also free. This holds for addition as well, and here is a peek at the reason: When two elementary free things are so distinct that their sum cannot coalesce to a single thing, then that composite sum of *free* things is naturally declared *free* by fiat. On the other hand, when such a sum *can* coalesce to a single thing, it does so by algebraically pre-shaping and positioning the two free summands so that, when added, an end of one *cancels* an end of the other, *poof* (imagine adding, for example, two free bivectors from the previous figure). The canceling pre-shaping ensures that the two surviving ends are exactly opposite (and separate, else they didn't survive)—again something free. So, if you begin with free vectors as your primitive elements, then you will be stuck in the free sub-language of the full Geometric Algebra. Therefore, don't imagine that you can represent points with the conventional free language alone—that will set you up for inconsistency and confusion.

The free sub-language has the lovely property that its elements can always be juxtaposed, which allows you to not only *extend* to higher dimensions, but to also *retract* to lower dimensions. Extension and retraction have complementary symmetries that, together, provide full information about geometric relationships. It was Clifford's genius to conjoin them into a very informative, widely celebrated *Clifford product*.² In Geometric Algebra this is called the *geometric product*, and it makes the free sub-language extremely versatile and expressive—it's a wonderful place to get stuck in.

How not to distinguish free from bound.

When you examine any contemporary book on Geometric Algebra, you discover that the vectors, bivectors and trivectors within it *are all depicted filled-in*, as tho they were bound. And yet they are allowed to roam around, *as tho they were not filled in*. How can these books get away

with such blatant inconsistency? *By refusing to allow points into the formalities except as outcasts*, that's how. The precise formal distinction between bound and free (coming up) disappears with points left in the interpretation. This renders the inconsistency so subtle it doesn't get noticed.

It's not as tho there were a malicious conspiracy to exclude points; it's more subtle than that. Geometers, despite mathematics' renowned proud rejection of meaning, approach their intrinsically meaningful subject with deeply held preconceptions that are fertile and mostly correct. Points within the symbolism would crumble these preconceptions around the edges like this: *Points would need the same dimension as free vectors to add properly with them. This would require free vectors not to be filled-in* for Grassmann's gem to be able to assign dimension consistently. Points? One-dimensional? Don't kid me. Free vectors? Not filled-in? Ha!—how could such things have *fixed* length and direction?

When reason and logical consistency nudge comfortable misconception, misconception typically remains complacent and unmoved; so points remain in their underwear breathing beer fumes. Except in one glorious yet sobering case: the curious case of Hermann Gunther Grassmann.

When he happened on points late in his investigations, he quickly realized that they must have the same *order* as free vectors, namely one. This of course had generative consequences for everything on up, so he gave them *orders* too, corresponding to the formal dimensions developed in this paper. His supple accommodation of points within the formalities has yet to be matched. That's the glorious part.

Here's the sobering part: his *order* was for him not *dimension*, but rather a way to get things to interact properly, with no other meaning. Listen to this:

There are seven types of spatial magnitude, divided into four orders:

- | | |
|-----------|---|
| 1st order | 1. Simple or multiple points |
| | 2. Straight lines of definite length and direction |
| 2nd order | 3. Definite parts of definite infinite straight lines |
| | 4. Plane areas of definite magnitude and direction |
| 3rd order | 5. Definite parts of definite infinite planes |
| | 6. Definite volumes |
| 4th order | 7. Definite volumes |

Volumes appear twice here, once as magnitudes of third order, once as magnitudes of fourth order, according as they are regarded as products of three straight lines of definite direction and length or as products of four points.¹³

Do you recognize these things? Here's a hint: the first item in each pair is bound, the second item is free with the same extent (neglecting separation, which *must not* be neglected, as explained shortly). So, in the first order there are points and free vectors, in the second order there are bound vectors and free bivectors, and so on. Since each pair has the same numerical extent, you see that Grassmann's *order* indeed corresponds to my *formal dimension*; but *his* dimension is, again, decidedly *informal*. This is most evident from his perplexing comment on "volumes",

in which he explicates the generative distinction between a free trivector and a bound one (in that order), *without making a dimensional distinction between them*.¹⁴

Why did he fail to do that? Remember, he arrived at points via *roving arrows*, the very same imagery we still have of free vectors. Even tho he explicitly discarded this imagery in favor of abstract algebra; nevertheless, free-vectors-as-roving-arrows must have become for him an inviolable concept, given how incredibly fruitful it had been. When he happened on points he was already comfortable with meaningless abstraction. Indeed, by then he embraced it; so he left *order* as an abstract formality that merely oiled the gears in his algebraic machinery. To have interpreted *order* geometrically would have required him to remove the shafts of his roving arrows—his “straight lines of definite direction and length”—leaving only arrow-heads and arrow-tails possessing mysteriously fixed separation and direction. But he clearly had no inclination to interpret *order* geometrically; and almost certainly no inclination to dismantle his fertile preconceptions. So, here again, comfortable misconception remained complacent and unmoved—Grassmann was human after all. That’s my guess.

Finite intrinsically composite semantic formalities

Intrinsically composite, as mentioned, is hard to imagine for same-extent free things, but easy to do so for *bound* things. Vectors bound to skewed non-intersecting lines, for example, do not have a common-enough extension factor to sum to a *single* thing; so their sum has extent $\{2, 2\}$. These vectors can, however, sum to *two* things in many different ways, and the most perspicuous is a free bivector perpendicular to a bound vector. In physical terms, the free bivector can articulate an angular velocity while the bound vector articulates a velocity along some line. Or these things can articulate a torque combined with a linear force. In short, addition of skewed lines generates an expressive *screw algebra*, reinvented by just about everyone who has really understood Grassmann.¹⁵

Are there any other imaginable same-extent bound sums that are intrinsically composite? How about the simple point sum $a+b$? It certainly is not intrinsically composite because it reduces to single midpoint $2m$. In fact, the humdrum sum of most weighted points reduces to a single thing. Well then, how about the *magic* sum, free vector $a-b$? It *is* intrinsically composite because (1) it cannot reduce to just one thing, and (2) it can always be decomposed by being un-bundled and re-bundled differently (that’s how the *poof* interactions work). To put this intuitively and generalize it, *separate but otherwise exactly opposite bound things are too distinct to sum to a single thing*. This should not come as a surprise—exactly opposite *is* quite distinct. Therefore, a free vector has extent $\{1, 1\}$; a free bivector has extent $\{2, 2\}$; and so on.

This notation transparently displays the *separated* extent of free things, but it fails to make a further crucial distinction. For example, it would give a free bivector the same extent as the sum of skewed vectors bound to non-intersecting lines, extent $\{2, 2\}$. A free bivector is a sum of *parallel* vectors bound to non-intersecting lines, which generally produces a single extent $\{2\}$ result (as previously demonstrated by approximation); except when the bound vectors are *exactly opposite*, in which case you get *exactly opposite* extent $\{2, 2\}$. *Exactly opposite* addition is what distinguishes free from bound; and such addition refuses to reduce to a single result by attempting to assign *contrary* properties to that result: zero magnitude and infinite “location”. The infinite “location” cannot be computed because it does not exist, but the zero magnitude is straightforward to compute. It is the *semantic formality* that shifts focus from sum to summands, from

infinite to finite, from bound to free. Which motivates a peculiar yet precise definition of a *free thing*: a non-zero thing with zero magnitude. That is the crucial distinction.

Let's put it to use: To be non-zero, a free thing, a bivector for example, must have formal *separation* (also readily computable), transparently annotated by *composite* extent $\{2, 2\}$. To further indicate that this is *yin-yang* composite in a cohesive *exactly opposite* way, One could call it extent $\{2, 2\}$ -with-zero-magnitude. This is an accurate but clumsy way of specifying that it is a free bivector. Since it is *free*, why not instead distinguish it with *free non-bold* notation?—extent $\{2\}$. Hence, extent $\{2\}$ means cohesive extent $\{2, 2\}$ -with-zero-magnitude. Similarly extent $\{1\}$ means cohesive extent $\{1, 1\}$ -with-zero-magnitude, and so on.

You can think of cohesive free extent as addition's respectful way of "extending". When addition is presented with separate but otherwise exactly opposite bound summands, it leaves them "extended" not by *incrementing* extent, but rather by *separating* it in a formal way. So, a free vector has formal separated extent $\{1\}$, a free bivector has separated extent $\{2\}$, and so on, just the numerical dimensions meaning-imposers have been declaring all along. Which is to say, free extent is conventional geometric dimension (called *grade* in the conventional free language), now well distinguished.

Distinguished by *separation*—that's what a tag of zero magnitude means: *this thing is not filled in*. Extension from a simple point fills it in, suddenly giving it *non-zero magnitude* equal to its just-departed separation. This magnitude becomes annotated with incremented **bold** singleton extent. Hence, separation describes a *pair of opposite summands*, something free; magnitude describes *one result*, something bound. The lowest extent magnitude is weight.

Weight: a scalar extended from a simple point generates weight—*non-zero magnitude* annotated with incremented bold singleton extent $\{1\}$. So, algebraically, a scalar is free, meaning that it is a non-zero thing¹⁶ tagged with a formal magnitude of zero, like all free things in the full Geometric Algebra. Such lowest-extent "separation" is the *value* of the scalar, which gets "filled in"—acquires locus—by extension with a simple point.

All this expressive formal distinction arises from finally letting points enter the symbolism as full-fledged geometric objects. To see that this emancipation is well worthwhile, examine the machinations necessary to keep points out while trying nonetheless to gain their expressive power. You need to impose...

Models

The three most popular models are the *vector space model*, the *homogeneous model*, and the *conformal model*. They are easiest to understand by the way they represent the plane. (Technically, the collection of primitive semantics is a model itself, *Grassmann's point model*; but it, unlike these, is not a clever artifice *imposed* on the symbolism, rather it is the DNA from which the symbolism is *derived*: Grassmann's point model—seed for a growing symbolism; conventional models—straightjackets for an inert symbolism.)

The vector space model of the plane begins with a formal algebra of two free vectors, whose inherent limitations are traditionally overcome informally. First, since these vectors are *free*, where are you going to put them? Answer: *implicitly* anchor them to a point, the *origin*. With their tails firmly fixed in one place, their heads can represent points—you get free vectors *and* points! Free vectors?—but you just bound them! No, let them roam around when you need them

to. But then you can't use them to represent points! No, just attach them to the origin when you need points. And so on. This has worked surprisingly well because, even tho the various fleeting distinctions all reside *outside* the symbolism, they nonetheless reside *inside* a human mind, which is superb with fleeting distinctions.

Fleeting distinctions won't do for a *model*, however, so the modeling community has decided that free vectors shall be *explicitly* anchored to the origin. This allows the fertile vector space idea to be unambiguously implemented on a computer. It has the ironic consequence that all the *free* elements in the language are effectively *bound* thru the origin, which has become *semi-formal* since it now has explicit representation in the software, even tho it has none in the algebra proper. Modeling enthusiasts don't mind this self-imposed handicap because they have more spiffy models that overcome it.

There is a different way to overcome the handicap that should be clear by now: having moved the origin from the informalities into the semi-formalities, why not continue this advance by moving it right into the formalities? As previously explained, this effectively moves some of the semantics into the symbolism. Here are the advantages of a meaningful symbolism: nebulous distinctions *outside* the language become precise distinctions *inside* the language, which now lets free things move parallel to themselves but requires bound things to stay in their confining spaces. Moreover, with the origin formal, everyone will have to implement it in the same way, as specified by the syntax of Geometric Algebra. (With the origin semi-formal, this is left to the digression of the implementer—need I say more?) The natural and expressive origin-in-the-formalities solution is obvious only in retrospect because comfortable misconception has rendered it almost inconceivable.

Consequently, modelers overcome the handicaps inherent in the vector space model in a different way, by moving to the *homogeneous model*. They always describe this by saying that you must move up an *extra* dimension above the plane. **Not!**—a *healthy plane already requires three dimensions*. You just saw that a plane with only two dimensions is crippled. By formally introducing the origin to heal it, *you increased its dimension by one*; but this is not an *extra* dimension, it is a *missing* dimension!

The origin increases dimension by one because it is just as variable as a free vector—that's what enables it to make a formal getaway from Hilbert's beer hall. Speaking abstractly, *dimension* just counts the number of variables available. With two free vectors, you have two variables available corresponding to the separation ("length") of each vector. But *you don't have points yet*—you don't *really* have a plane. To get points, you need a point to refer your free vectors to; and if that is done *outside* the symbolism, as it always has been in the last century, then you mangle your free vectors in the alleyway, as just described. By formally introducing a point for your free vectors to collaborate with, your symbolism suddenly acquires *the missing variable*, weight, which generates weighted points thruout the plane—this is now a *genuine* plane; it is not pointless anymore. It has abstract dimension 3 corresponding to formal extent $\{\mathbf{3}\}$.

Consequently, "move up an extra dimension" really means "use another free vector to stand in for the missing variable that a point would supply, had not preconception abandoned it in the gutter." Stating it baldly (badly?) like this makes achieving it obvious: let the separation of the "extra" free vector basis element correspond to the weight of the missing origin. Hence, distance above the plane corresponds to point weight; so unit points—*locations*—can be represented by anchored free vectors whose heads lie one unit above the original plane. The new unit-separated

plane becomes a model-with-“points”. (Modelers call them *points*, un-quoted, but their “points” are always free vectors masquerading as points. Such things aren’t real points—they are undesirables, tacitly denied full geometric rights. This is easy to demonstrate: real points with *full* rights would, for consistency, induce free vectors with *separation*, absent in every model except the generative one, Grassmann’s point model.)

The homogeneous model is fun to play with because it shows, in an unexpected way, how perpendicular distance can precisely represent point weight. By applying some mind-boggling dimension hopping, you can use this model to articulate both free and bound things, thereby overcoming the shortcomings of the vector space model. Of course such contortions make sense only if you are *absolutely determined* to keep real points out of your formalities.

Semi-formally anchored free vectors give the homogeneous model its own peculiar handicaps. To bypass my point sympathies, let a model enthusiast expose them:

... the geometric algebra approach exposes some weaknesses in the homogeneous model. It turns out that we cannot really define a useful inner product in the representation space \mathbb{R}^{n+1} that represents the metric aspects of the original space \mathbb{R}^n well; we can only revert to the inner product of \mathbb{R}^n . As a consequence, we also have no compelling geometric product and our geometric algebra of \mathbb{R}^{n+1} is impoverished ...¹⁷

Not to worry—there is another model that overcomes this fresh impoverishment, the *conformal model*, “which requires *two* extra dimensions”,¹⁷ meaning *one* dimension above the defective homogeneous model, which constitutes one *genuine* dimension above Grassmann’s formal point space. The genuine extra dimension is given negative distances, thereby causing the augmented space to curve in such a way that extension in it can be projected down to *rounds* in the original space.¹⁸ So, in a conformal representation of physical space, the extension of three points generates the unique circle thru them; the extension of four points generates the unique sphere thru them. “Points” themselves are rounds with zero radius (null vectors). Clever, huh? It gets even better: by including a special “point at infinity”, ∞ , you can generate *flats*, rounds with infinite radius. Moreover...

Our model also solves another problem that perplexed Grassmann throughout his life. He was finally forced to conclude that it is impossible to define a geometrically meaningful inner product between points. The solution eluded him because it requires the concept of indefinite metric that accompanies the concept of null vector. Our model supplies an inner product $a \cdot b$ that directly represents the Euclidean distance between the points. This is a boon to distance geometry, because it greatly facilitates computation of distances among many points.¹⁹

Altho it is true that the bondage of Grassmann’s points a and b naturally precludes an inner product for them (since they cannot be juxtaposed), it is not true that this precludes finding the distance between them. Conjure up free vector $a - b$, and then use Grassmann’s inner product—there is no need to hop on the conformal pony, lovely tho it may be, to access its high-dimensional inner product.

There is no question that models are lovely, with beautiful, fruitful mathematics generated by incredibly curious, creative mathematicians whom I deeply admire. But models typically solve problems they have inflicted on themselves by leaving the origin in the semi-formalities, where it cannot interact properly with its geometric kin. Even worse, they solve problems in an indirect, obscure and inefficient way that Grassmann’s full language can solve in a direct, transparent and efficient way. The *formal* point-generated distinction between bound and free (obviously lacking

in the previous purely free models) enables this. This distinction, coupled with Clifford's unification of the free sublanguage, gives you an exceptionally expressive way to articulate geometric concepts: hop in the free sublanguage when you need its services; hop in the bound part when you need things in certain places; stay in the free sublanguage as much as possible because it is most versatile and expressive.

To illustrate, simple subtraction of points \mathbf{a} and \mathbf{b} moved these *bound* things into the *free* language where distance calculations are available. Simple subtraction can also generate rounds by moving points into the free language. For example, to generate the circle thru three points, subtract the points pairwise to form three free vectors, then apply symmetry to find the center point. Finally, do a direct, transparent and efficient fixed-radius computation. Or just apply symmetry *directly* to generate peripheral points iteratively—this is even more direct, transparent and efficient. (As for *flats*, why not generate them with ordinary low-dimensional extension? This avoids the superfluous imposed “point at infinity” and is (need I say?) direct, transparent and efficient.)

Simple subtraction of separate but otherwise identical bound things can always be used to move them into the versatile free sublanguage. This is seldom convenient for anything but points, however, and seldom necessary either because the elegant relationship between bound and free offers an easier way to hop into the free sublanguage: extract free parts.

Extracting free parts is such a crucial bridge from the full Geometric Algebra into its free sublanguage that I like to consider it a primitive operation, on par with extension, retraction and their unification, the geometric product. This requires a pithy notation for extracting free parts; and it also requires the elegant relationship, $\mathbf{a} = a\mathbf{e}_1$, to be added to the symbolism as an axiom. (For the purpose of generating free parts, I'm guessing it really is an axiom: If you don't want it as an axiom, then you have to isolate free a on the right to directly generate free parts; and good luck with that. Remember, point \mathbf{a} cannot participate in a retraction (an inner product), nor a geometric product, so how are you going to *un-extend* it to the scalar unit to isolate free a ? If this intrigues you, study how Whitehead did it by crippling his language with tacit context.²⁰)

As a practical matter, free parts discard locus information so they are easy to compute. *Especially* easy if you keep your basis as free as possible by allowing just one point in it, the origin. With this discipline, the origin is the sole source of bondage; so extracting a free part amounts to extricating a (generally translated) origin.

Finally, hopping into the bound part of Geometric Algebra from its free sublanguage is trivial in two ways: (1) Extend from a simple point. Since this point can be smoothly moved, any bound thing can be smoothly moved. (2) Decompose the free thing using addition's associative law. Suddenly you have two *relatively bound* things, one of which gets associated with *something else*, thereby transferring the relative bondage to it. The screw algebra illustrates this well, as the following section explains.

The free–bound distinction is intrinsically semantic.

Can the conventional rules of Geometric Algebra, the axioms, make the free–bound distinction? If so, they would have to distinguish between a sum and its summands. But they can't—as far as the axioms are concerned, *a sum and its added summands are literally equal*. That is a great boon because it allows free and bound to be articulated together, and intermingled. For example, altho the semantics make a clear distinction between magnitude and separation, the axioms can't because they cannot distinguish a sum from its summands. Instead, the axioms simply

articulate magnitude and separation *simultaneously*, indifferently; and automatically switch from one to the other as the situation dictates. To begin understanding this, scale a free vector as you sneak up on it. During the approximation you will be scaling a diminishing *weight*. At the limit, however, you will suddenly be scaling a *separation*—a startling revelation for me. (At that limit you will be *simultaneously* scaling a zero weight, which will of course remain zero.) The axioms' indifferent automatic switching allows free things to be decomposed into bound things when need be, or vice versa. In short, the axioms do their syntactic duty, which is: let you express any valid sentence in Geometric Algebra, and let you transform that into more informative sentences.

The free-bound distinction *requires* distinguishing between a sum and its summands; so if syntax can't do it, semantics will have to. You have just seen that this is done by a formal zero-nonzero distinction. To illustrate just how adamantly semantic such distinction is, let's resolve the final dangling perplexity:

“Vectors bound to skewed non-intersecting lines, for example, cannot sum to a *single* thing; consequently their sum has extent $\{\mathbf{2}, \mathbf{2}\}$. These vectors can, however, sum to *two* things in many different ways, and the most perspicuous is a free bivector perpendicular to a bound vector.” ???

How can that be? A free bivector plus a bound vector seems to have extent $\{\mathbf{2}, \mathbf{2}\}$, which would expand into extent $\{\{\mathbf{2}, \mathbf{2}\}\text{-with-zero-magnitude}, \mathbf{2}\}$. Can these *three* things (when fully decomposed) possibly reduce to extent $\{\mathbf{2}, \mathbf{2}\}$? Yes—the perspicuous sum is a convenient and illuminating *canonical* form, not a *minimal* form. A minimal form has extent $\{\mathbf{2}, \mathbf{2}\}$ and this is easy to see: move the free bivector so that the tail of one of its ends is right on the tail of the bound vector. Conjoining these two vectors like this gives them a common extension factor that engages extension's respect for summary. This collapses the two vectors to one, leaving two skewed bound vectors. (Here you see *relative* bondage, transparently exposed. Reverse the procedure orthogonally to get the canonical form.)

So you see, for dimension to be well defined, addition must present the extent() operator with a minimal form. This is inherently computational—intrinsically semantic. Magnitude is an important part of this computation since it distinguishes free from bound, an essential distinction for a minimal form.

(If you want addition to give *you* a canonical form, you will have to ask for it—that is how semantics works; and it is just one more reason an expressive geometric language requires semantics. Whether your request should be formal or semi-formal is a question we haven't pondered adequately because we have shunned semantics.)

The computer people know that, in pathological cases, computation cannot distinguish between zero and darn-close-to-zero. With respect to magnitude this means that, in pathological cases, the full Geometric Algebra cannot distinguish between free and bound—it will not produce a genuine minimal form. This does not, however, invalidate the distinction that magnitude makes; it just requires extra care to do well, as with all formal semantics. Keeping free things bundled goes a long way toward minimizing this problem—to repeat, stay as free as possible and take care to represent free bundles by individual names, which keeps their bundles intact during computation. This forces their zero magnitude to *stay* zero, unambiguously. The extreme way to do that is to remain in the comfortable conventional free language; but then you are back where you started—bigoted against points.

Finally, it is not as if mathematics has been immaculately devoid of formal semantics, though many mathematicians are reluctant to admit it. What, for example, is a *metric* but a precise assignment of *meaning* to distance? That is just as semantic, and just as formal, as the distinction between free and bound; in fact it helps establish that distinction in the full Geometric Algebra. Mathematicians should come out of the closet about semantics. Computer scientists outed long ago and they feel liberated now.

References

Clif.1881

William Kingdon Clifford, *Mathematical Papers*, edited by Robert Tucker, reprinted 1968, Chelsea. “Clifford was above all and before all a geometer ... if he had lived, we might have known something.” We did.

Dors.2007

Leo Dorst, Daniel Fontijne, Stephen Mann, *Geometric Algebra for Computer Science*, Morgan Kaufmann. An eloquent presentation of the meaning-imposer’s approach to Geometric Algebra.

Gras.18??

Hermann Gunther Grassmann, *A New Branch of Mathematics, The Ausdehnungslehre of 1844, and Other works*, Open Court. Translated by Lloyd Kannenberg, 1995. Effectively three books in one: the 1878 edition of the *Ausdehnungslehre* of 1844, the 1847 *Geometrisch Analyse*, and a rich sampler of Grassmann’s articles on mathematics and physics.

Gras.1862

Hermann Grassmann, *Extension Theory*, American Mathematical Society. Translated by Lloyd Kannenberg, 2000. Grassmann’s attempt to appeal to mathematicians after his 1844 opus had failed to do so.

Harp.201?

Gary Harper, *Playing with Geometric Algebra—Stalking a coherent language*. Forthcoming. My own modest contribution, the hardest thing I’ve ever tried. Foundational chapters may be downloaded from gary-harper.com/ Some of the prose makes me wince now.

Hest.198?-20??

David Hestenes, modelingNTS.la.ASU.edu/ Collects most of Hestenes’s papers on mathematics and physics. Hestenes is the prolific lion of the free Geometric Algebra and has become an enthusiast of geometric models.

¹Gras.1844, Part Two, “Elementary Magnitudes”, devoted to the bound language. Part One, “Extensive Magnitudes” is devoted to the free language.

²Clif. “Applications of Grassmann’s Extensive Algebra”, 1878, and “On the classification of Geometric Algebras”, 1876.

³Trust me on this—to embarrass these *particular* authors would unfairly single them out from among the *many* others who hold similar opinions, saying, for example, that you can’t add London to Paris.

⁴Harp. “Speaking of Space” p5.

⁵John Playfair, (Euclid's) *Elements of Geometry*, J. B. Lippincott, 1857. p8: "A point is that which has position, but not magnitude." This is Playfair's elucidation of Euclid's "that which has no parts, or which has no magnitude."

⁶Grass.1844 p162, Grass.1847 p326. He developed point sums obliquely in terms of center of gravity, and displacements from an arbitrary origin. For a direct midpoints-of-midpoints development, see Harp, "Adding Points".

⁷Grass.1844, p154–161. This discovery prompted Grassmann to abandon the notation he had used for displacements in the first half of his book. For commentary on the mathematics, see Harp, "Speaking of Space", p43.

⁸Grass.1844 p11: "I found that the analysis I had discovered did not touch only on the subject of geometry, as it seemed before. Rather, I soon realized that I had come upon the domain of a new science, of which geometry itself is only a special application."

⁹Grass.1844 p184–185, 192–198.

¹⁰William Rowan Hamilton, *On Symbolical Geometry*, p2. Available from maths.tcd.ie/pub/HistMath/People/Hamilton/

¹¹Otto Blumenthal, *Lebensgeschichte*, Berlin, 1935 in David Hilbert, "Gesammelte Abhandlungen", p403.

¹²Set $\hat{a} = a\mathbf{a}$, then $\hat{a} = a \blacktriangleleft \mathbf{a}$ in parallel with $\mathbf{a} = a \blacktriangleleft \hat{a}$. In this *uniform notation*, \hat{a} denotes *both* magnitude and location, as bold generic \mathbf{a} does. Similarly, non-bold \hat{a} is the *free part* of bold \hat{a} , namely its weight a .

¹³Grass.1845 p289. Written by request for clarification from his editor.

¹⁴Altho a bound trivector has extent $\{4\}$, a free trivector has extent $\{3\}$, meaning extent $\{\mathbf{3}, \mathbf{3}\}$ -with-zero-magnitude.

¹⁵Clifford was working on a screw algebra he called *biquaternions* when he encountered Grassmann's ideas; and he realized Grassmann could unify his screw algebra. This is clear from the main subsection of Clifford's previously cited "Applications of Grassmann's Extensive Algebra"², which is titled "*On the Relation of Grassmann's Method to Quaternions and Biquaternions; and on the Generalization of these Systems*". See also Clif.1873 p181, "Preliminary Sketch of Biquaternions".

¹⁶In the full Geometric Algebra zero cannot be a scalar or any other kind of number. See Harp. "Speaking of Space" p45.

¹⁷Dors.2007 p246.

¹⁸Hence, you can import the conformal model into Grassmann's language of imaginable space by augmenting that language with an extra dimension having Minkowski metric. This would require *explicit* projection back into imaginable space, which is certainly more expressive than the *implicit* projection of the conformal model, but more cumbersome when imaginable space is the only space of interest. But in that narrow case, why even bother with the conformal ploy?—Grassmann's origin-in-the-formalities already articulates imaginable space directly, transparently, efficiently.

¹⁹Hest.2001 *Unified Algebraic Framework for Classical Geometry* (UAFCG.html). I am guessing that Grassmann would have been as appalled by the conformal model as he was by Hamilton's vector algebra: it's not *his baby*, and it demonstrates for the n^{th} time that we still haven't understood the bound part of *his baby*.

²⁰Alfred North Whitehead, *Universal Algebra*, Cambridge, 1898, p516. Whitehead called extracting a free part the "operation of taking the vector". He achieved it by setting the free ceiling, the unit trivector, equal to scalar one, which introduces a tacit context that precludes moving to other dimensions. Even worse, it obscures important dimensional distinctions. "Vector" had its Latin meaning, "*carrier*", for Whitehead; and because free things are able to *carry* bound things of the same extent, *vector* became synonymous with *free* for him, which makes for perplexing reading.



An open letter concerning
WInHD: Wavelet-based Inverse Halftoning via Deconvolution

Ramesh Neelamani and Richard Baraniuk

Birth: The niche problem of inverse halftoning error-diffused halftones has been addressed by a number of solid researchers using several practical and effective methods. However, due to the non-linearity of the halftoning process and the complexities of the human visual system, the methods proposed to date have been ad hoc.

At first glance, we thought that we had little chance of coming up with even a mediocre solution to the nonlinear inverse halftoning problem. We pursued lines of research from photon-limited imaging and Polya trees, but those approaches lead nowhere. One day, Rob Nowak found some literature on an intriguing linear approximation to halftoning. We were pleasantly surprised when a wavelet-thresholding based estimator based on this linear approximation produced competitive results (not only in terms of the workhorse mean-squared-error (MSE) metric but also in terms of a standard visual quality metric). We called our algorithm Wavelet-based Inverse Halftoning via Deconvolution (WInHD).

We thought that WInHD would be a “slam-dunk” paper that would certainly interest the image processing community, because in addition to presenting competitive results near the state-of-the-art, our insights also reduced the inverse halftoning problem to a well-understood deconvolution problem. Furthermore, assuming that the linear approximation was accurate and that the model noise was Gaussian, we were able to derive and analyze bounds on WInHD’s MSE performance as the image resolution increased.

With high optimism, we submitted a paper to a top-tier image processing journal.

Death: But alas, our enthusiasm was deflated due to the following review points, which we disagree with.

- The linear approximation was deemed questionable. Any claims about optimality were deemed to be overstated.
- Our results were deemed to be visually inferior. The metrics used used to evaluate our simulation results did not conform to the quality of the images as perceived. We were urged to seek input from the experts in the field and then publish the results of the survey.

The combination of lukewarm reviews and diverging author interests meant that the paper had to be abandoned.

After-life: With its publication in *Rejecta Mathematica*, we would like to honestly address some of the issues raised in our paper’s day of reckoning.

We believe that while the reviewers raised several valid points, the paper contained several contributions that would benefit the image processing community. Addressing the linear approximation point, we agree that a linear approximation to the halftoning process is not suitable for all purposes. However, our view is that the surprising results obtained using such a model make our paper more, not less, interesting. We do concede that the optimality claims made in the paper need to be taken with a this linear approximation in mind. However, the limitations of our analysis



have been clearly stated in the paper (it was termed as conditional optimality in the paper, but perhaps our analysis required some bigger and bolder disclaimers).

On the visual quality issue, beauty indeed lies in the eyes of the beholder! Like a majority of image processing practitioners, we agree that the MSE may be inadequate to measure the visual quality of an image. However, in our paper, we employed all of the metrics that were accessible in the literature (that is, we did not cherry-pick them) to substantiate that our method provided “superior visual” performance (arguably a strong term to use, but certainly not obviously wrong). Surveys can certainly be an effective approach to analyzing an image processing result. But, while useful, conducting surveys for every image processing paper borders on onerous. As an alternative, we published our code so that our results were reproducible and so that our method could be tested on anyone’s images of choice.

The tussle about the visual quality improvement afforded by WInHD seems to have no easy resolution in sight. However an even larger question emerges. Is it really necessary for follow-on papers to always significantly improve upon previous results? Should a paper’s publishability be so heavily reliant on the improved results that it produces? How about insights that may open some closed doors?

WInHD: Wavelet-based Inverse Halftoning via Deconvolution

Ramesh Neelamani and Richard Baraniuk*

Abstract

We propose the *Wavelet-based Inverse Halftoning via Deconvolution* (WInHD) algorithm to perform inverse halftoning of error-diffused halftones. WInHD is motivated by our realization that inverse halftoning can be formulated as a deconvolution problem under Kite et al.'s linear approximation model for error diffusion halftoning. Under the linear model, the error-diffused halftone comprises the original gray-scale image blurred by a convolution operator and colored noise; the convolution operator and noise coloring are determined by the error diffusion technique. WInHD performs inverse halftoning by first inverting the model-specified convolution operator and then attenuating the residual noise using scalar wavelet-domain shrinkage. Since WInHD is model-based, it is easily adapted to different error diffusion halftoning techniques. Using simulations, we verify that WInHD is competitive with state-of-the-art inverse halftoning techniques in the mean-squared-error sense and that it also provides good visual performance. We also derive and analyze bounds on WInHD's mean-squared-error performance as the image resolution increases.

Key words: inverse halftoning, error diffusion, deconvolution, wavelets, wavelet-vaguelette.

*Contact author: R. Neelamani. neelsh@gmail.com.

1 Introduction

Digital halftoning is a common technique used to render a sampled gray-scale image using only black or white dots [1] (see Figures 3(a) and (b)); the rendered bi-level image is referred to as a halftone. *Inverse halftoning* is the process of retrieving a gray-scale image from a given halftone. Applications of inverse halftoning include rehalftoning, halftone resizing, halftone tone correction, and facsimile image compression [2, 3]. In this paper, we focus on inverse halftoning images that are halftoned using popular error diffusion techniques such as those of Floyd et al. [4], and Jarvis et al. [5] (hereby referred to as Floyd and Jarvis respectively).

Error-diffused halftoning is non-linear because it uses a quantizer to generate halftones. Recently, Kite et al. proposed an accurate linear approximation model for error diffusion halftoning (see Figure 4) [6, 7]. Under this model, the halftone $y(n_1, n_2)$ is expressed in terms of the original gray-scale image $x(n_1, n_2)$ and additive white noise $\gamma(n_1, n_2)$ as (see Figure 1)

$$\begin{aligned} y(n_1, n_2) &= \mathcal{P}x(n_1, n_2) + \mathcal{Q}\gamma(n_1, n_2) \\ &= (p * x)(n_1, n_2) + (q * \gamma)(n_1, n_2), \end{aligned} \quad (1)$$

with $*$ denoting convolution and (n_1, n_2) indexing the pixels. The \mathcal{P} and \mathcal{Q} are the linear time-invariant (LTI) systems with respective impulse responses $p(n_1, n_2)$ and $q(n_1, n_2)$ determined by the error diffusion technique.

From (1), we infer that inverse halftoning can be posed as the classical *deconvolution* problem because the gray-scale image $x(n_1, n_2)$ can be obtained from the halftone $y(n_1, n_2)$ by deconvolving the filter \mathcal{P} in the presence of the colored noise $\mathcal{Q}\gamma(n_1, n_2)$. Conventionally, deconvolution is performed in the Fourier domain. The Wiener deconvolution filter, for example, would estimate $x(n_1, n_2)$ by inverting \mathcal{P} and *regularizing* the resulting noise with scalar Fourier shrinkage. As we will see, inverse halftoning using a Gaussian low-pass filter (GLPF) [8] can be interpreted as a naive Fourier deconvolution approach to inverse halftoning.

Unfortunately, all Fourier-based deconvolution techniques induce ringing and blurring artifacts due to the fact that the energy of edge discontinuities spreads over many Fourier coefficients. As a result of this uneconomical representation, the desirable edge Fourier coefficients are easily confounded with those due to the noise [9–11].

In contrast, the wavelet transform provides an economical representation for images with sharp edges [12]. This economy makes edge wavelet coefficients easy to distinguish from those due to the noise and has led to powerful image estimation algorithms based on scalar wavelet shrinkage [11, 13].

The wavelet transform was first exploited in inverse halftoning by J. Luo et al. [14]. Xiong et al. extended this algorithm using non-orthogonal, redundant wavelets to obtain improved results for error-diffused halftones [15]. Both these algorithms rely on a variety of steps such as clipping and edge-adapted noise attenuation in the wavelet subbands to exploit different empirical observations. However, these steps and their implications are not theoretically well-justified.

To simultaneously exploit the economy of wavelet representations and the interplay between inverse halftoning and deconvolution, we propose the *Wavelet-based Inverse Halftoning via Deconvolution* (WInHD) algorithm (see Figure 2) [16]. WInHD provides robust estimates by first inverting the convolution operator \mathcal{P} determined by the linear model (1) for error diffusion and then effectively attenuating the residual colored noise using wavelet-domain scalar shrinkage operations [13, 17]. Since WInHD is model-based, it easily adapts to different error diffusion halftoning techniques. See Figure 3 for simulation results.

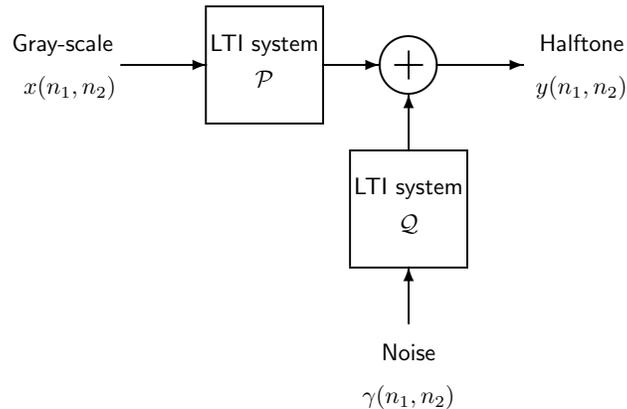


Figure 1: *Linear approximation for error diffusion halftoning.* Under the linear model of [6, 7], the error-diffused halftone $y(n_1, n_2)$ comprises the original gray-scale image $x(n_1, n_2)$ passed through an LTI system \mathcal{P} and white noise $\gamma(n_1, n_2)$ colored by an LTI system \mathcal{Q} . The systems \mathcal{P} and \mathcal{Q} are determined by the error diffusion technique.

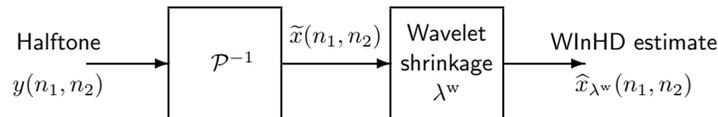


Figure 2: *Wavelet-based Inverse Halftoning via Deconvolution (WInHD).* WInHD inverts the convolution operator \mathcal{P} to obtain a noisy estimate $\tilde{x}(n_1, n_2)$ of the gray-scale image. Subsequent scalar shrinkage with λ^w in the wavelet domain (for example, level-dependent hard thresholding) effectively attenuates the residual noise corrupting $\tilde{x}(n_1, n_2)$ to yield the WInHD estimate $\hat{x}_{\lambda^w}(n_1, n_2)$.

Unlike previous inverse halftoning algorithms, we can analyze the theoretical performance of WInHD under certain conditions. For images in a Besov smoothness space, we derive the minimum rate at which the WInHD estimate's mean-squared-error (MSE) decays as the resolution increases; that is, as number of pixels in the gray-scale image tends to infinity. We assume that the linear model for error diffusion (1) is exact and that the noise $\gamma(n_1, n_2)$ is Gaussian. Further, if the gray-scale image $x(n_1, n_2)$ contains some additive noise (say, scanner noise) before halftoning that is Gaussian, then we show that the MSE decay rate enjoyed by WInHD in estimating the noise-free $x(n_1, n_2)$ is optimal; that is, no other inverse halftoning algorithm can have a better error decay rate for every Besov space image as the number of image pixels increases.

Section 2 describes Kite et al.'s linear model for error diffusion halftoning from [6, 7]. Section 3 clarifies the equivalence between inverse halftoning and deconvolution and also analyzes Fourier-domain inverse halftoning. Section 4 presents a brief overview of wavelets. Section 5 discusses the proposed WInHD algorithm and its theoretical performance. Section 6 illustrates the experimental performance of WInHD. Section 7 provides conclusions and future directions. A technical proof in Appendix A completes the paper.



Figure 3: (a) Original Lena image (512×512 pixels). (b) Floyd halftone. (c) Multiscale gradient-based estimate [18], PSNR = 31.3 dB. (d) WInHD yields competitive PSNR performance (32.1 dB) and visual performance. (All documents including the above images undergo halftoning during printing. To minimize the halftoning effect, the images have been reproduced at the maximum size possible.) See Figure 8 for image close-ups.

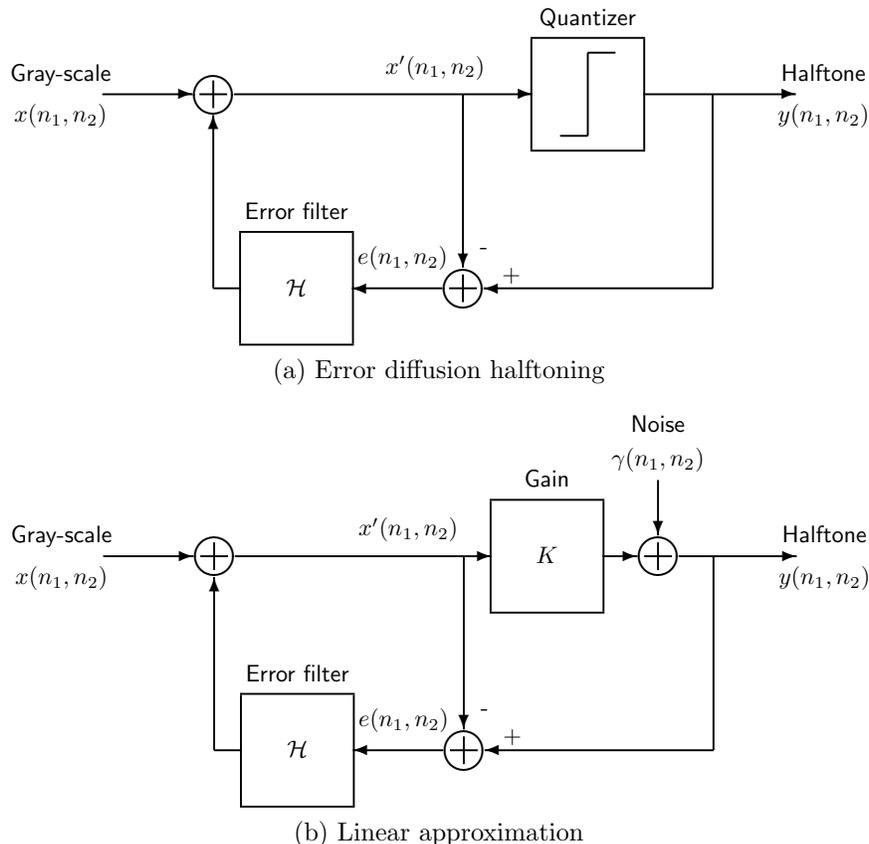


Figure 4: (a) Error diffusion halftoning. The gray-scale image pixels $x(n_1, n_2)$ are quantized to yield $y(n_1, n_2)$ and the quantization error $e(n_1, n_2)$ is diffused over a causal neighborhood by the error filter \mathcal{H} . (b) The linear model approximates the quantizer with gain K and additive white noise $\gamma(n_1, n_2)$ [6].

2 Linear Model for Error Diffusion

In this section, we describe the non-linear error diffusion halftoning and the linear approximation proposed in [6, 7].

Digital halftoning is a process that converts a given gray-scale digital image (for example, each pixel value $\in [0, 1, \dots, 255]$) into a bi-level image (for example, each pixel value = 0 or 255) [1]. Error diffusion halftoning is one popular approach to perform digital halftoning. The idea is to take the error from quantizing a gray-scale pixel to a bi-level pixel and diffuse this quantization error over a causal neighborhood. The error diffusion ensures that the spatially-localized average pixel values of the halftone and original gray-scale image are similar; therefore, the halftone visually resembles the gray-scale image. Figure 4(a) illustrates the block diagram for error diffusion halftoning. The $x(n_1, n_2)$ denote the pixels of the input gray-scale image and $y(n_1, n_2)$ denote the bi-level pixels of the output halftone. The $x'(n_1, n_2)$, which yields $y(n_1, n_2)$ after quantization, is obtained by diffusing the quantization error $e(n_1, n_2)$ over a causal neighborhood of $x(n_1, n_2)$ using the error filter \mathcal{H} . The quantizer makes error-diffused halftoning a non-linear technique. Error diffusion

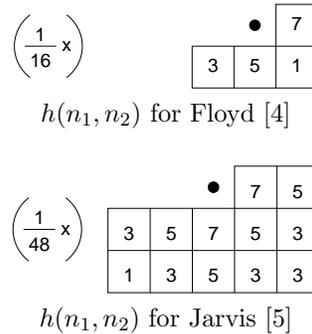


Figure 5: Error filters $h(n_1, n_2)$ for Floyd [4] and Jarvis [5] error diffusion. The quantization error at the black dot is diffused over a causal neighborhood according the displayed weights.

techniques such as Floyd [4] and Jarvis [5] are characterized by their choice of \mathcal{H} 's impulse response $h(n_1, n_2)$ (see Figure 5).

Recently, Kite et al. proposed an accurate linear model for error diffusion halftoning [6, 7]. This model accurately predicts the “blue noise” (high-frequency noise) and edge sharpening effects observed in various error-diffused halftones. As shown in Figure 4(b), this model approximates the effects of quantization using a gain K followed by the addition of white noise $\gamma(n_1, n_2)$. The halftone $y(n_1, n_2)$ can then be written in terms of the gray-scale image $x(n_1, n_2)$ and the additive white noise $\gamma(n_1, n_2)$ as in (1); the error diffusion technique determines the 2-dimensional (2-D) frequency responses of the LTI filters \mathcal{P} and \mathcal{Q} as

$$P(f_1, f_2) := \frac{K}{1 + (K - 1)H(f_1, f_2)}, \tag{2}$$

$$Q(f_1, f_2) := \frac{1 - H(f_1, f_2)}{1 + (K - 1)H(f_1, f_2)} \tag{3}$$

with $P(f_1, f_2)$, $Q(f_1, f_2)$, and $H(f_1, f_2)$ denoting the respective 2-D Fourier transforms of $p(n_1, n_2)$, $q(n_1, n_2)$, and $h(n_1, n_2)$. For any given error diffusion technique, Kite et al. found that the gain K is almost constant for different images. However, the K varied with the error diffusion technique [6]; for example, $K = 2.03$ for Floyd, while $K = 4.45$ for Jarvis. Figure 6 (a) and (b) illustrate the radially-averaged frequency response magnitudes of the filters \mathcal{P} and \mathcal{Q} for Floyd and Jarvis respectively; these responses are obtained by averaging over an annulus of constant radius in the 2-D frequency domain [1]. In [7], Kite et al. further generalized the linear model of (1) by using different gains K_s and K_n in the signal transfer function $P(f_1, f_2)$ and the noise transfer function $Q(f_1, f_2)$ respectively: $P(f_1, f_2) := \frac{K_s}{1+(K_s-1)H(f_1, f_2)}$ and $Q := \frac{1-H(f_1, f_2)}{1+(K_n-1)H(f_1, f_2)}$. In this paper, we assume a single gain factor K for both the signal and noise transfer functions as proposed in [6].

3 Inverse Halftoning \approx Deconvolution

Given a halftone $y(n_1, n_2)$ (see Figure 4(a)), inverse halftoning aims to estimate the gray-scale image $x(n_1, n_2)$. In the classical deconvolution problem, given the blurred and noisy observation $y(n_1, n_2)$ as in (1) with known LTI filters responses $p(n_1, n_2)$ and $q(n_1, n_2)$, we seek to estimate

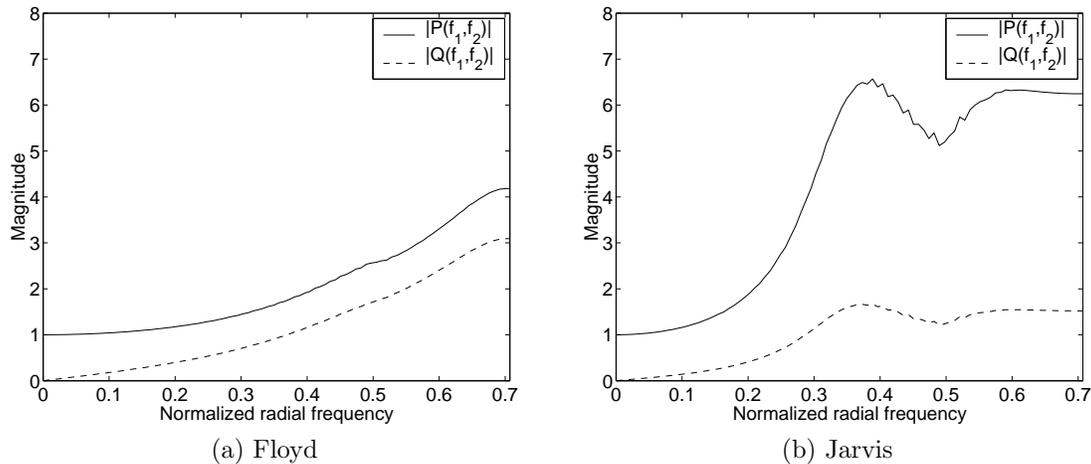


Figure 6: Plots (a) and (b) respectively illustrate the radially-averaged frequency response magnitudes $|P(f_1, f_2)|$ (solid line) and $|Q(f_1, f_2)|$ (dotted line) for Floyd and Jarvis. The high-pass $|P(f_1, f_2)|$ explains the sharpened edges, while the high-pass $|Q(f_1, f_2)|$ explains the “blue noise” behavior in the halftones.

$x(n_1, n_2)$. Thus, under the linear model of [6, 7], inverse halftoning can be posed as a deconvolution problem.

3.1 Deconvolution

Due to the interplay between inverse halftoning and deconvolution, the well-studied deconvolution literature [19–21] can be exploited to understand inverse halftoning as well. Deconvolution algorithms conceptually consist of the following steps:

1. *Operator inversion*

Invert the convolution operator \mathcal{P} to obtain a noisy estimate $\tilde{x}(n_1, n_2)$ of the input signal¹

$$\tilde{x}(n_1, n_2) := \mathcal{P}^{-1}y(n_1, n_2) = x(n_1, n_2) + \mathcal{P}^{-1}\mathcal{Q}\gamma(n_1, n_2). \quad (4)$$

2. *Transform-domain shrinkage*

Attenuate the colored noise $\mathcal{P}^{-1}\mathcal{Q}\gamma(n_1, n_2)$ by expressing $\tilde{x}(n_1, n_2)$ in terms of a chosen orthonormal basis $\{b_k\}_{k=0}^{N-1}$ and shrinking the k -th component with a scalar λ_k , $0 \leq \lambda_k \leq 1$ [22]

$$\hat{x}_\lambda := \sum_k \langle \tilde{x}, b_k \rangle \lambda_k b_k = \sum_k (\langle x, b_k \rangle + \langle \mathcal{P}^{-1}\mathcal{Q}\gamma, b_k \rangle) \lambda_k b_k \quad (5)$$

to obtain the deconvolution estimate \hat{x}_λ .

The $\sum_k \langle x, b_k \rangle \lambda_k b_k$ in (5) denotes the *retained part* of the signal $x(n_1, n_2)$ that shrinkage preserves from (4), while $\sum_k \langle \mathcal{P}^{-1}\mathcal{Q}\gamma, b_k \rangle \lambda_k b_k$ denotes the *leaked part* of the colored noise

¹For non-invertible \mathcal{P} , we replace \mathcal{P}^{-1} by its pseudo-inverse and $x(n_1, n_2)$ by its orthogonal projection onto the range of \mathcal{P} in (4).

$\mathcal{P}^{-1}\mathcal{Q}\gamma(n_1, n_2)$ that shrinkage fails to attenuate. Clearly, we should set $\lambda_k \approx 1$ if the variance $\sigma_k^2 := \mathbb{E}(|\langle \mathcal{P}^{-1}\mathcal{Q}\gamma, b_k \rangle|^2)$ of the k -th colored noise component is small relative to the energy $|\langle x, b_k \rangle|^2$ of the corresponding signal component and set $\lambda_k \approx 0$ otherwise. The shrinkage by λ_k can also be interpreted as a form of *regularization* for the deconvolution inverse problem [20].

The choice of transform domain to perform the shrinkage in deconvolution (see Step 2 above) critically influences the MSE of the deconvolution estimate. An important fact is that for a given transform domain, even with the best possible λ_k 's, the estimate \hat{x}_λ 's MSE is lower-bounded within a factor of 2 by [9–11]

$$\sum_k \min(|\langle x, b_k \rangle|^2, \sigma_k^2). \quad (6)$$

From (6), \tilde{x}_λ can have small MSE only when most of the signal energy ($= \sum_k |\langle x, b_k \rangle|^2$) and colored noise energy ($= \sum_k \sigma_k^2$) is captured by just a few transform-domain coefficients — we term such a representation *economical* — and when the energy-capturing coefficients for the signal and noise are different. Otherwise, the \tilde{x}_λ is either excessively noisy due to leaked noise components or distorted due to lost signal components.

Traditionally, the Fourier domain (with sinusoidal b_k) is used to estimate $x(n_1, n_2)$ from $\tilde{x}(n_1, n_2)$ because it represents the colored noise $\mathcal{P}^{-1}\mathcal{Q}\gamma(n_1, n_2)$ in (4) most economically. That is, among orthonormal transforms, the Fourier transform captures the maximum colored noise energy using a fixed number of coefficients because it diagonalizes convolution operators [23]. Fourier-based deconvolution performs both the operator inversion and the shrinkage simultaneously in the Fourier domain as

$$\hat{X}_{\lambda^f} := Y(f_1, f_2) \frac{1}{P(f_1, f_2)} \lambda_{f_1, f_2}^f \quad (7)$$

with shrinkage

$$\lambda_{f_1, f_2}^f := \frac{|P(f_1, f_2)|^2}{|P(f_1, f_2)|^2 + \Upsilon(f_1, f_2)|Q(f_1, f_2)|^2} \quad (8)$$

at the different frequencies. The $Y(f_1, f_2)$ and $\hat{X}_{\lambda^f}(f_1, f_2)$ denote the 2-D Fourier transforms of $y(n_1, n_2)$ and the deconvolution estimate $\hat{x}_{\lambda^f}(n_1, n_2)$ respectively. The $\Upsilon(f_1, f_2)$ in (8) is called the *regularization term* and is set appropriately during deconvolution [20]. For example, using the signal to noise ratio to set $\Upsilon(f_1, f_2) = \frac{\mathbb{E}(|\Gamma(f_1, f_2)|^2)}{|X(f_1, f_2)|^2}$ in (7) yields the Wiener deconvolution estimate [24]; the $\Gamma(f_1, f_2)$ and $X(f_1, f_2)$ denote the respective Fourier transforms of $\gamma(n_1, n_2)$ and $x(n_1, n_2)$. The $\frac{1}{P(f_1, f_2)} \lambda_{f_1, f_2}^f$ in (7) constitutes the frequency response of the so-called *deconvolution filter*.

Fourier-based deconvolution suffers from the drawback that its estimates for images with sharp edges are unsatisfactory either due to excessive noise or due to distortions such as blurring or ringing. Since the energy due to the edge discontinuities spreads over many image Fourier coefficients, as dictated by the MSE bound in (6), any estimate obtained via Fourier-domain shrinkage suffers from a large MSE.

3.2 Inverse half-toning via Gaussian low-pass filtering (GLPF)

Conventionally, inverse half-toning has been performed using a finite impulse response (FIR) Gaussian filter with coefficients $g(n_1, n_2) \propto \exp[-(n_1^2 + n_2^2)/(2\sigma_g^2)]$, where $-4 \leq n_1, n_2 \leq 4$, and σ_g determines the bandwidth [8]. We can interpret inverse half-toning using GLPF as a naive Fourier-domain deconvolution approach to inverse half-toning. This is substantiated by our observation that the deconvolution filter $\frac{1}{P(f_1, f_2)} \lambda_{f_1, f_2}^f$ (see (7) and (8)) constructed with the linear model filters \mathcal{P}

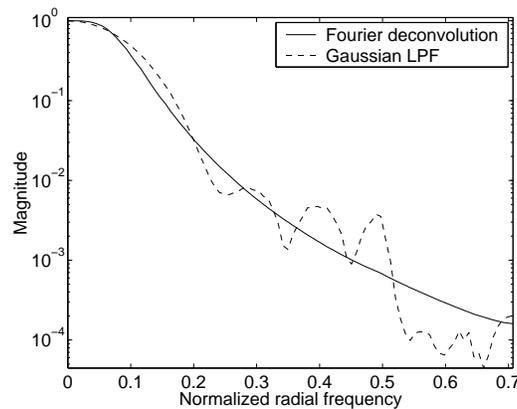


Figure 7: Comparison of radially-averaged frequency response magnitudes of the FIR GLPF (dashed line) used for inverse halftoning in [8] with the response of the deconvolution filter (solid line) constructed with filters \mathcal{P} and \mathcal{Q} for Floyd and with $\Upsilon(f_1, f_2) \propto \frac{1}{f_1^2 + f_2^2}$ (see (7) and (8)). Ripples in the GLPF frequency response result because the filter is truncated in space.

and \mathcal{Q} for Floyd and with regularization $\Upsilon(f_1, f_2) \propto \frac{1}{f_1^2 + f_2^2}$ has a frequency response that closely matches the frequency response of the GLPF (see Figure 7) [8]. The corresponding inverse halftone estimates obtained using simulations are also nearly identical. Predictably, GLPF estimates suffer from the same drawbacks that afflict any Fourier-based deconvolution estimate — excessive noise (when σ_g is small) or significant blurring (when σ_g is large). Exploiting the insights provided by the deconvolution perspective, we can infer that unsatisfactory GLPF estimates result because the Fourier domain does not economically represent images with edges.

4 Background on Wavelets

In contrast to Fourier representations, wavelets provide economical representations for a diverse class of signals including images with edges [11, 12].

4.1 Wavelet transform

The 2-D discrete wavelet transform (DWT) represents a spatially-continuous image $x(t_1, t_2) \in L^2([0, 1]^2)$ in terms of shifted versions of a low-pass scaling function ϕ and shifted and dilated versions of prototype bandpass wavelet functions $\{\psi^{LH}, \psi^{HL}, \psi^{HH}\}$ [11, 25]. For special choices of ϕ and ψ 's, the shifted and dilated functions $\phi_{j,k_1,k_2}(t_1, t_2) := 2^j \phi(2^j t_1 - k_1, 2^j t_2 - k_2)$, and $\psi_{j,k_1,k_2}^b := 2^j \psi^b(2^j t_1 - k_1, 2^j t_2 - k_2)$ with $b \in \mathcal{B} := \{LH, HL, HH\}$, where the LH , HL , and HH denote the *subbands* of the wavelet decomposition, form an orthonormal basis. The j parameter corresponds to the *scale* of the analysis, while the k_1, k_2 parameters correspond to the *location*. A finite-resolution approximation $x^J(t_1, t_2)$ to $x(t_1, t_2)$ is given by

$$x^J(t_1, t_2) = \sum_{k_1, k_2 \in \mathbb{Z}} s_{j_0, k_1, k_2} \phi_{j_0, k_1, k_2}(t_1, t_2) + \sum_{b \in \mathcal{B}} \sum_{j=j_0}^J \sum_{k_1, k_2 \in \mathbb{Z}} w_{j, k_1, k_2}^b \psi_{j, k_1, k_2}^b(t_1, t_2),$$

with scaling coefficients $s_{j_0, k_1, k_2} := \langle x, \phi_{j_0, k_1, k_2} \rangle$ and wavelet coefficients $w_{j, k_1, k_2}^b := \langle x, \psi_{j, k_1, k_2}^b \rangle$. The parameter J controls the resolution of the wavelet reconstruction $x^J(t_1, t_2)$ of $x(t_1, t_2)$; in fact, the L_2 error $\|x^J - x\|_2 \rightarrow 0$ as $J \rightarrow \infty$.

The DWT can be extended to transform sampled images as well. Consider, for example, a sampled image obtained by sampling $x(t_1, t_2)$ uniformly as

$$x(n_1, n_2) = N \int_{\frac{n_2}{\sqrt{N}}}^{\frac{n_2+1}{\sqrt{N}}} \int_{\frac{n_1}{\sqrt{N}}}^{\frac{n_1+1}{\sqrt{N}}} x(t_1, t_2) dt_1 dt_2, \quad 0 \leq n_1, n_2 \leq \sqrt{N} - 1. \quad (9)$$

For such N -pixel images, the N wavelet coefficients can be efficiently computed in $O(N)$ operations using a filter bank consisting of low-pass filters, high-pass filters, and decimators [11].

Purely for notational convenience, we henceforth refer to the location parameters k_1, k_2 by k and do not explicitly specify the different wavelet subbands: w_{j, k_1, k_2}^b and ψ_{j, k_1, k_2}^b for $b \in \mathcal{B} := \{LH, HL, HH\}$ will be referred to simply as $w_{j, k}$ and $\psi_{j, k}$. Further, we discuss the processing of only the wavelet coefficients, but all steps are replicated on the scaling coefficients as well.

4.2 Economy of wavelet representations

Wavelets provide economical representations for images in smoothness spaces such as Besov spaces [9, 12]. Roughly speaking, a Besov space $B_{p, q}^s$ contains functions with “ s derivatives in L_p ,” with q measuring finer smoothness distinctions [12]. Besov spaces with different s , p , and q characterize many classes of functions; for example, $B_{1, \infty}^1$ contains piece-wise polynomial images [26]. If a continuous-space image $x(t_1, t_2) \in B_{p, q}^s$, $s > \frac{2}{p} - 1$, $1 \leq p, q \leq \infty$, then the DWT coefficients computed using the image samples (see (9)) satisfies (for all N)

$$\frac{1}{\sqrt{N}} \left(\sum_j 2^{jq(s+1-\frac{2}{p})} \left(\sum_k |w_{j, k}|^p \right)^{\frac{q}{p}} \right)^{\frac{1}{q}} < \infty, \quad (10)$$

assuming the underlying wavelet basis functions are sufficiently smooth [10, 17, 27].² From (10), we can infer that the wavelet coefficients of Besov space images decay exponentially fast with increasing scale j . Further, among all orthogonal transforms, the wavelet transform captures the maximum (within a constant factor) signal energy using a fixed number of coefficients for the worst-case Besov space signal [9].

4.3 Wavelet-based signal estimation

The wavelet transform’s economical representations have been exploited in many fields [11]. For example, wavelets provide an effective solution to the problem of estimating signal samples $x(n_1, n_2)$ from additive white Gaussian noise (AWGN) corrupted observations [17, 27, 28]

$$\tilde{x}(n_1, n_2) = x(n_1, n_2) + \gamma(n_1, n_2), \quad (11)$$

²The traditional Besov space characterizing equation in [10, 17, 27] assumes L_2 -normalized wavelet coefficients $w_{j, k}$; that is, $\sum_{j, k} |w_{j, k}|^2 = \|x(t_1, t_2)\|_2^2$. Because the $w_{j, k}$ used in (10) are computed using signal samples $x(n_1, n_2)$ that satisfy $\sum_{j, k} |w_{j, k}|^2 = \sum_{n_1, n_2} |x(n_1, n_2)|^2 \approx N \|x(t_1, t_2)\|_2^2$, a normalization factor of \sqrt{N} is required.

with $\gamma(n_1, n_2)$ denoting AWGN samples of variance σ^2 . Such a setup is similar to estimating $x(n_1, n_2)$ from (4) but with $\mathcal{P}^{-1}\mathcal{Q}$ equal to identity. Simple shrinkage in the wavelet domain with scalars λ^w can provide excellent estimates of $x(n_1, n_2)$. This shrinkage is illustrated by (5) with wavelet basis functions as b_k 's and with identity $\mathcal{P}^{-1}\mathcal{Q}$. For example, *hard thresholding* shrinks the wavelet coefficients of $\tilde{x}(n_1, n_2)$ with scalars

$$\lambda_{j,k}^w = \begin{cases} 1, & \text{if } |\tilde{w}_{j,k}| > \rho_j \sigma_j, \\ 0, & \text{if } |\tilde{w}_{j,k}| \leq \rho_j \sigma_j, \end{cases} \quad (12)$$

with $\tilde{w}_{j,k} := \langle \tilde{x}, \psi_{j,k} \rangle$, σ_j^2 the noise variance at wavelet scale j , and ρ_j the scale-dependent threshold factor (for examples, see [11, p. 442]). When the pixels $x(n_1, n_2)$ arise from a continuous-space image $x(t_1, t_2) \in B_{p,q}^s$ with $s > \frac{2}{p} - 1$ and $1 \leq p, q \leq \infty$, hard thresholding (with judiciously chosen ρ_j [28]) provides estimates whose MSE-per-pixel decays at least as fast as $N^{-\frac{s}{s+1}}$ as $N \rightarrow \infty$ [17, 27]. Further, no estimator can achieve a better error decay rate for every $x(t_1, t_2) \in B_{p,q}^s$. If the threshold factor ρ_j is chosen to be scale-independent, then the MSE decay rate is decelerated by an additional $\log N$ factor.

In practice, the *Wavelet-domain Wiener Filter* (WWF) improves on the MSE performance of hard thresholding by employing Wiener estimation on each wavelet coefficient [29, 30]. WWF chooses

$$\lambda_{j,k}^w = \frac{|w_{j,k}|^2}{|w_{j,k}|^2 + \sigma_j^2}. \quad (13)$$

However, the coefficients $w_{j,k}$ required to construct the $\lambda_{j,k}^w$ are unknown. Hence, a ‘‘pilot’’ estimate of the unknown signal is first computed using hard thresholding. Then, using λ^w constructed with the pilot estimate’s wavelet coefficients in (13), WWF shrinkage is performed. Sufficiently different wavelet basis functions must be used in the two steps [29, 30].

5 Wavelet-based Inverse HalfToning Via Deconvolution (WInHD)

To simultaneously exploit the economy of wavelet representations and our realization about the interplay between inverse halfToning and deconvolution, we propose the WInHD algorithm [16]. WInHD adopts the wavelet-based deconvolution approach of [10] to perform inverse halfToning.

5.1 WInHD algorithm

WInHD employs scalar shrinkage in the wavelet domain to perform inverse halfToning as follows (see Figure 2):.

1. *Operator inversion*
As in (4), obtain a noisy estimate $\tilde{x}(n_1, n_2)$ of the input image by inverting \mathcal{P} .
2. *Wavelet-domain shrinkage*
Employ scalar shrinkage in the wavelet domain to attenuate the noise $\mathcal{P}^{-1}\mathcal{Q}\gamma(n_1, n_2)$ in $\tilde{x}(n_1, n_2)$ and obtain the WInHD estimate $\hat{x}_{\lambda^w}(n_1, n_2)$ as follows:
 - (a) Compute the DWT of the noisy \tilde{x} to obtain $\tilde{w}_{j,k} := \langle \tilde{x}, \psi_{j,k} \rangle$.

- (b) Shrink the noisy $\tilde{w}_{j,k}$ with scalars $\lambda_{j,k}^w$ (using (12) or (13)) to obtain $\hat{w}_{j,k;\lambda^w} := \tilde{w}_{j,k} \lambda_{j,k}^w$. The colored noise variance at each scale j determining the $\lambda_{j,k}^w$ is given by $\sigma_j^2 := \mathbb{E} \left(|\langle \mathcal{P}^{-1} \mathcal{Q} \gamma, \psi_{j,k} \rangle|^2 \right)$.
- (c) Compute the inverse DWT with the shrunk $\hat{w}_{j,k;\lambda^w}$ to obtain the WInHD estimate $\hat{x}_{\lambda^w}(n_1, n_2)$.

For error diffusion systems, \mathcal{P}^{-1} is an FIR filter. Hence, the noisy estimate $\tilde{x}(n_1, n_2)$ obtained in Step 1 using \mathcal{P}^{-1} is well-defined. The subsequent wavelet-domain shrinkage in Step 2 effectively extracts the few dominant wavelet components of the desired gray-scale image $x(n_1, n_2)$ from the noisy $\tilde{x}(n_1, n_2)$ because the residual noise $\mathcal{P}^{-1} \mathcal{Q} \gamma(n_1, n_2)$ corrupting the wavelet components is not excessive.

WInHD can be easily adapted to different error diffusion techniques simply by choosing the gain K recommended by [6] and the error filter response $h(n_1, n_2)$ for the target error diffusion technique. K and $h(n_1, n_2)$ determine the filters \mathcal{P} and \mathcal{Q} (see (2) and (3)) required to perform WInHD. In contrast, the gradient-based inverse half-toning method [18] adapts to a given error diffusion technique by employing a set of smoothing filters that need to be designed carefully.

5.2 Asymptotic performance of WInHD

With advances in technology, the spatial resolution of digital images (controlled by the number of pixels N) has been steadily increasing. Hence any inverse half-toning algorithm should not only perform well at a fixed resolution but should also guarantee good performances at higher spatial resolutions. In this section, under some assumed conditions, we deduce the rate at which the per-pixel MSE for WInHD decays as number of pixels $N \rightarrow \infty$.

Invoking established results in wavelet-based image estimation in Gaussian noise, we prove the following proposition in Appendix A about the asymptotic performance of WInHD.

Proposition 1 *Let $x(n_1, n_2)$ be a N -pixel gray-scale image obtained as in (9) by uniformly sampling a continuous-space image $x(t_1, t_2) \in B_{p,q}^s$ with $t_1, t_2 \in [0, 1)$, $s > \frac{2}{p} - 1$, and $1 \leq p, q, \leq \infty$. Let $p(n_1, n_2)$ and $q(n_1, n_2)$ denote known filter impulse responses that are invariant with N and with Fourier transform magnitudes $|P(f_1, f_2)| \geq \epsilon > 0$ and $|Q(f_1, f_2)| < \infty$. Let $y(n_1, n_2)$ be observations obtained as in (1) with $\gamma(n_1, n_2)$ zero-mean AWGN samples with variance σ^2 . Then, the per-pixel MSE of the WInHD estimate $\hat{x}(n_1, n_2)$ obtained from $y(n_1, n_2)$ using hard thresholding behaves as*

$$\frac{1}{N} \mathbb{E} \left(\sum_{n_1, n_2} |\hat{x}(n_1, n_2) - x(n_1, n_2)|^2 \right) \leq C N^{-\frac{s}{s+1}}, \quad N \rightarrow \infty, \quad (14)$$

with constant $C > 0$ independent of N .

The above proposition affirms that the per-pixel MSE of the WInHD estimate decays as $N^{-\frac{s}{s+1}}$ with increasing spatial resolution ($N \rightarrow \infty$) under the mild assumptions discussed below.

The central assumption in Proposition 1 is that the linear model (1) for error diffusion is accurate. This is well-substantiated in [6, 7]. The conditions $|P(f_1, f_2)| \geq \epsilon > 0$ and $|Q(f_1, f_2)| < \infty$ respectively ensure that \mathcal{P} is invertible and that the variance of the colored noise $\mathcal{Q} \gamma(n_1, n_2)$ is bounded. We have verified that for common error diffusion half-toning techniques such as Floyd and Jarvis, the filters \mathcal{P} and \mathcal{Q} recommended by the linear model of Kite et al. satisfy these

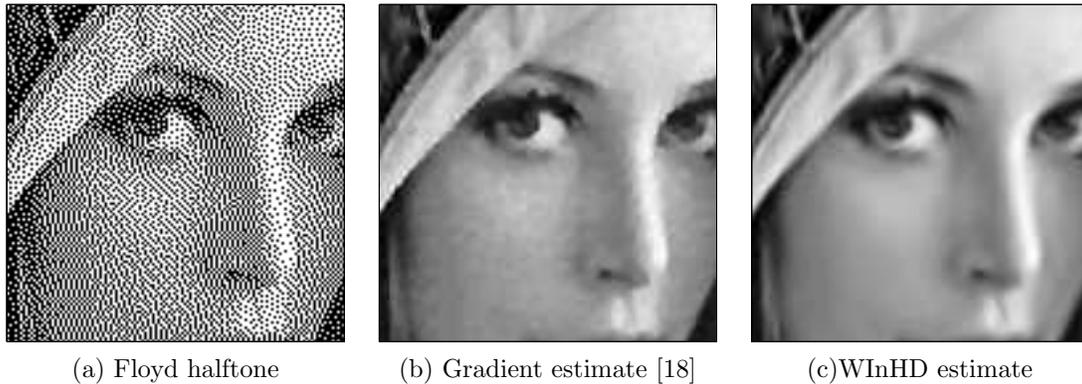


Figure 8: Close-ups (128×128 pixels) of (a) Floyd halftone, (b) Gradient estimate [18], and (c) WInHD estimate.

conditions (see Figure 6). The final assumption is that the noise $\gamma(n_1, n_2)$ is Gaussian; this is required to invoke the established results on the asymptotics of wavelet-based estimators [17]. However, recently, wavelet-domain thresholding has been shown to be optimal for many other noise distributions as well [31, 32]. Hence the noise Gaussianity assumption in Proposition 1 could be relaxed.

Often, gray-scale digital images are corrupted with some noise before being subjected to halftoning. For example, sensor noise corrupts images captured by charged coupled device (CCD) based digital cameras. In such cases as well, WInHD can effectively estimate the noise-free gray-scale image with an MSE decay rate of $N^{\frac{-s}{s+1}}$ as in Proposition 1. Further, WInHD's MSE decay rate can be shown to be optimal. The noise-free gray-scale image and resulting halftone can be related using the linear model of [6, 7] as

$$y(n_1, n_2) = \mathcal{P} [x(n_1, n_2) + \beta(n_1, n_2)] + \mathcal{Q}\gamma(n_1, n_2), \quad (15)$$

with $\beta(n_1, n_2)$ denoting the noise corrupting the gray-scale image before halftoning. If the $\beta(n_1, n_2)$ is AWGN with non-zero variance, then we can easily infer that the residual noise after inverting \mathcal{P} in Step 1 of WInHD can be analyzed like white noise because its variance is bounded but non-zero [10]. Hence we can invoke well-established results on the performance of wavelet-based signal estimation in the presence of white noise [17, 27, 28] to conclude that no estimator can achieve a better error decay rate than WInHD for every Besov space image. Thus, WInHD is an optimal estimator for inverse halftoning error-diffused halftones of noisy images.

6 Results

We illustrate WInHD's performance using 512×512 -pixel *Lena* and *Peppers* test images halftoned using the Floyd algorithm [4] (see Figure 3 and 8). All WInHD estimates and software are available at www.dsp.rice.edu/software. We set the gain $K = 2.03$, as calculated for Floyd in [6, 7], and use the Floyd error filter response $h(n_1, n_2)$ (see Figure 5) to characterize the impulse responses $p(n_1, n_2)$ and $q(n_1, n_2)$. Inverting the operator \mathcal{P} (Step 2) requires $O(N)$ operations and memory

Table 1: PSNR and computational complexity of inverse halfToning algorithms (N pixels).

Inverse halfToning algorithm	PSNR (dB)		Computational complexity
	<i>Lena</i>	<i>Peppers</i>	
Gaussian [8]	28.6	27.6	$O(N)$
Kernel [2]	32.0	30.2	$O(N)$
Gradient [18]	31.3	31.4	$O(N)$
Wavelet denoising [15]	31.7	30.7	$O(N \log N)$
WInHD	32.1	31.2	$O(N)$

for a N -pixel image since \mathcal{P}^{-1} is FIR. To perform the wavelet-domain shrinkage (Step 2), we choose the WWF because it yields better estimates compared to schemes such as hard thresholding.

Estimates obtained by shrinking DWT coefficients are not shift-invariant; that is, translations of $y(n_1, n_2)$ will result in different estimates. Hence, we exploit the *complex wavelet transform* (CWT) instead of the usual DWT to perform the WWF. The CWT expands images in terms of shifted and dilated versions of *complex-valued* basis functions instead of the real-valued basis functions used by the DWT [33, 34]; the expansion coefficients are also complex-valued. Wavelet-domain shrinkage using WWF on the CWT coefficient magnitudes yields significantly improved near shift-invariant estimates with just $O(N)$ operations and memory. (The redundant, shift-invariant DWT can also be used instead of the CWT to obtain shift-invariant estimates [11], but the resulting WInHD algorithm requires $O(N \log N)$ operations and memory.) The standard deviation of the noise $\gamma(n_1, n_2)$, which is required during wavelet shrinkage, is calculated using the standard deviation of $y(n_1, n_2)$'s finest scale CWT coefficients.

Figures 3 and 8 compares the WInHD estimate with the multiscale gradient-based estimate [18] for the Lena image. We quantify the WInHD's performance by measuring the peak signal-to-noise ratio $\text{PSNR} := 20 \log_{10} \frac{512 \times 255}{\|\hat{x} - x\|_2}$ (for 512×512 -pixel images with gray levels $\in [0, 1, \dots, 255]$) with $\hat{x}(n_1, n_2)$ the estimate. Table 1 summarizes the PSNR performance and computational complexity of WInHD compared to published results for inverse halfToning with Gaussian filtering [8], kernel estimation [2], gradient estimation [18], and wavelet denoising with edge-detection [15]. We can see that WInHD is competitive with the best published results.

The WInHD estimate yields competitive visual performance as well. We quantify visual performance using two metrics: weighted SNR (WSNR) [35, 36] and the *Universal Image Quality Index* (UIQI) [37]. Both metrics were computed using the halfToning toolbox of [38]. The WSNR is obtained by weighting the SNR in the frequency domain according to a linear model of the human visual system [35, 36]. The WSNR numbers in Table 2 are calculated at a spatial Nyquist frequency of 60 cycles/degree. The recently proposed UIQI metric of Wang et al. effectively models image distortion with a combination of correlation loss, luminance distortion, and contrast distortion [37]; UIQI $\in [-1, 1]$ with larger values implying better image quality. For the Lena image, WInHD's performance in terms of both the visual metrics is competitive with the gradient estimate's performance (see Table 2).

Table 2: Visual metrics for inverse halftoned estimates of Lena.

Algorithm	WSNR (dB)	UIQI
Gradient [18]	34.0	0.62
WInHD	35.9	0.62

7 Conclusions

Using the linear error diffusion model of [6, 7], we have demonstrated that inverse halftoning can be posed as a deconvolution problem in the presence of colored noise. Exploiting this new perspective, we have proposed the simple *Wavelet-based Inverse Halftoning via Deconvolution* (WInHD) algorithm based on wavelet-based deconvolution to perform inverse halftoning. Since WInHD is model-based, it is easily tunable to the different error diffusion halftoning techniques. WInHD yields state-of-the-art performance in the MSE sense and visually.

WInHD also enjoys desirable theoretical properties under certain mild conditions. For images in a Besov space, WInHD estimate's MSE is guaranteed to decay rapidly as the spatial resolution of the input gray-scale image increases. Further, if the gray-scale image lies in a Besov space and is noisy before halftoning, then WInHD's MSE decay rate cannot be improved upon by any estimator.

We have assumed *a priori* knowledge of the error diffusion filter in this paper. However, the error diffusion filter is not always known. Under such circumstances, the error diffusion filter coefficients could be estimated by integrating adaptive techniques such as the one proposed by Wong [39] into our algorithm. However, this remains a topic of future study.

To facilitate efficient hardware implementation, in addition to requiring minimal memory and computations, an inverse halftoning algorithm should also be compatible with fixed-point digital signal processors. For example, the gradient-based algorithm [18] is optimized for hardware implementation while still obtaining good inverse halftoning results. Since our focus in this paper has been primarily theoretical, we have not specifically addressed any hardware optimization issues. The design of a hardware-compatible inverse halftoning algorithm based on WInHD is a topic of interesting future study.

A Decay Rate of WInHD's MSE

We deduce the asymptotic performance of WInHD as claimed in Proposition 1.

Instead of analyzing the problem of estimating $x(n_1, n_2)$ from $y(n_1, n_2)$, we can equivalently analyze the estimation of $x(n_1, n_2)$ from the noisy observation $\tilde{x}(n_1, n_2)$ obtained after inverting \mathcal{P} (see (4)). The reduction is equivalent because $P(f_1, f_2)$ is known and invertible (since $|P(f_1, f_2)| \geq \epsilon > 0$).³

The frequency components of the colored noise $\mathcal{P}^{-1}\mathcal{Q}\gamma(n_1, n_2)$ corrupting the $\tilde{x}(n_1, n_2)$ in (4) is given by $\frac{Q(f_1, f_2)\Gamma(f_1, f_2)}{P(f_1, f_2)}$. These frequency components are independent and Gaussian because the Fourier transform diagonalizes convolution operators. Since $|P(f_1, f_2)|$ is strictly non-zero and $|Q(f_1, f_2)|$ is bounded, the variance of $\frac{Q(f_1, f_2)\Gamma(f_1, f_2)}{P(f_1, f_2)}$ is uniformly bounded — say with variance ζ^2

³Since the filter \mathcal{P}^{-1} is FIR for error diffusion systems, boundary effects are negligible asymptotically because only a finite number of boundary pixels are corrupted.

— at all frequencies.

Because the estimation error due to wavelet-domain hard thresholding is monotone with respect to noise variance [10], the error in estimating $x(n_1, n_2)$ from (4) using wavelet-domain hard thresholding is less than the error in estimating $x(n_1, n_2)$ observed in white noise as in (11) but with variance ζ^2 . Hence the per-pixel MSE in estimating $x(n_1, n_2)$ from (4) can be bounded with the decay rate $N^{\frac{-s}{s+1}}$ established for the white noise setup (see Section 4.3) to yield (14) with a constant $C > 0$ independent of N [27, 28]. \square

Acknowledgments

We thank Dr. Brian Evans for his many constructive comments.

References

- [1] R. Ulichney, *Digital Halftoning*. Cambridge, MA: MIT Press, 1987.
- [2] P. Wong, “Inverse halftoning and kernel estimation for error diffusion,” *IEEE Trans. Image Processing*, vol. 6, pp. 486–498, Apr. 1995.
- [3] M. Y. Ting and E. A. Riskin, “Error-diffused image compression using a binary-to-gray scale decoder and predictive pruned tree-structured vector quantization,” *IEEE Trans. Image Processing*, vol. 3, pp. 854–857, Nov. 1994.
- [4] R. W. Floyd and L. Stienberg, “An adaptive algorithm for spatial grayscale,” *Proc. Soc. Image Display*, vol. 17, no. 2, pp. 75–77, 1976.
- [5] J. Jarvis, C. Judice, and W. Ninke, “A survey of techniques for the display of continuous tone pictures on bilevel displays,” *Comput. Graph and Image Process.*, vol. 5, pp. 13–40, 1976.
- [6] T. D. Kite, B. L. Evans, A. C. Bovik, and T. L. Sculley, “Digital halftoning as 2-D delta-sigma modulation,” *Proc. IEEE ICIP '97*, vol. 1, pp. 799–802, Oct. 26–29 1997.
- [7] T. D. Kite, B. L. Evans, A. C. Bovik, and T. L. Sculley, “Modeling and quality assessment of halftoning by error diffusion,” *IEEE Trans. Image Processing*, vol. 9, pp. 909–922, May 2000.
- [8] S. Hein and A. Zakhor, “Halftone to continuous-tone conversion of error-diffusion coded images,” *IEEE Trans. Image Processing*, vol. 4, pp. 208–216, Feb. 1995.
- [9] D. L. Donoho, “Unconditional bases are optimal bases for data compression and for statistical estimation,” *Appl. Comput. Harmon. Anal.*, vol. 1, pp. 100–115, Dec. 1993.
- [10] D. L. Donoho, “Nonlinear solution of linear inverse problems by Wavelet-Vaguelette Decomposition,” *Appl. Comput. Harmon. Anal.*, vol. 2, pp. 101–126, 1995.
- [11] S. Mallat, *A Wavelet Tour of Signal Processing*. Academic Press, 1998.
- [12] R. A. DeVore, B. Jawerth, and B. J. Lucier, “Image compression through wavelet transform coding,” *IEEE Trans. Inform. Theory*, vol. 38, pp. 719–746, Mar. 1992.

- [13] D. L. Donoho and I. M. Johnstone, "Ideal spatial adaptation via wavelet shrinkage," *Biometrika*, vol. 81, pp. 425–455, 1994.
- [14] J. Luo, R. de Queiroz, and Z. Fan, "A robust technique for image descreening based on the wavelet transform," *IEEE Trans. Signal Processing*, vol. 46, pp. 1179–1184, Apr. 1998.
- [15] Z. Xiong, M. T. Orchard, and K. Ramchandran, "Inverse half-toning using wavelets," *IEEE Trans. Signal Processing*, vol. 8, pp. 1479–1482, Oct. 1999.
- [16] R. Neelamani, R. D. Nowak, and R. G. Baraniuk, "Model-based inverse half-toning with Wavelet-Vaguelette Deconvolution," in *Proc. IEEE ICIP '00*, (Vancouver, Canada), pp. 973–976, Sept. 2000.
- [17] D. L. Donoho, "De-noising by soft-thresholding," *IEEE Trans. Inform. Theory*, vol. 41, pp. 613–627, May 1995.
- [18] T. D. Kite, N. Damera-Venkata, B. L. Evans, and A. C. Bovik, "A fast, high-quality inverse half-toning algorithm for error diffused halftones," *IEEE Trans. Image Processing*, vol. 9, pp. 1583–1592, Sept. 2000.
- [19] A. K. Jain, *Fundamentals of Digital Image Processing*. Englewood Cliffs, NJ: Prentice-Hall, 1989.
- [20] A. K. Katsaggelos (Ed.), *Digital Image Restoration*. New York: Springer-Verlag, 1991.
- [21] R. Neelamani, H. Choi, and R. G. Baraniuk, "Wavelet-based deconvolution using optimally regularized inversion for ill-conditioned systems," in *Wavelet Applications in Signal and Image Processing VII, Proc. SPIE*, vol. 3813, pp. 58–72, July 1999.
- [22] W. James and C. Stein, "Estimation with quadratic loss," in *Proc. Fourth Berkeley Symp. Math. Statist. Probab.*, vol. 1, pp. 361–380, Univ. California Press, 1961.
- [23] G. Davis and A. Nosratinia, "Wavelet-based image coding: An overview," in *Appl. Comput. Control Signals Circuits* (B. N. Datta, ed.), vol. 1, Birkhauser, 1999.
- [24] K. R. Castleman, *Digital Image Processing*. New Jersey: Prentice Hall, 1996.
- [25] C. S. Burrus, R. A. Gopinath, and H. Guo, *Introduction to Wavelets and Wavelet Transforms: A Primer*. Prentice-Hall, 1998.
- [26] J. Kalifa and S. Mallat, "Thresholding estimators for linear inverse problems," *Ann. Statist.*, vol. 31, Feb. 2003.
- [27] D. L. Donoho and I. M. Johnstone, "Asymptotic minimaxity of wavelet estimators with sampled data," *Statist. Sinica*, vol. 9, no. 1, pp. 1–32, 1999.
- [28] D. L. Donoho and I. Johnstone, "Minimax estimation by wavelet shrinkage," *Ann. Statist.*, vol. 26, pp. 879–921, 1998.
- [29] S. Ghael, A. M. Sayeed, and R. G. Baraniuk, "Improved wavelet denoising via empirical Wiener filtering," in *Wavelet Applications in Signal and Image Processing V, Proc. SPIE*, vol. 3169, pp. 389–399, Oct. 1997.

- [30] H. Choi and R. G. Baraniuk, "Analysis of wavelet domain Wiener filters," in *IEEE Int. Symp. Time-Frequency and Time-Scale Analysis*, (Pittsburgh), Oct. 1998.
- [31] R. Averkamp and C. Houdre, "Wavelet thresholding for non (necessarily) Gaussian noise: Functionality," *Ann. Statist.*, vol. 31, Feb. 2003.
- [32] H.-Y. Gao, "Choice of thresholds for wavelet shrinkage estimate of the spectrum," *Journal of Time Series Analysis*, vol. 18, pp. 231–251, 1997.
- [33] N. G. Kingsbury, "Complex wavelets for shift invariant analysis and filtering of signals," *Appl. Comput. Harmon. Anal.*, vol. 10, pp. 234–253, May 2001.
- [34] J. K. Romberg, H. Choi, R. G. Baraniuk, and N. G. Kingsbury, "A hidden Markov tree model for the complex wavelet transform," *IEEE Trans. Signal Processing*, 2002. Submitted.
- [35] T. Mitsa and K. Varkur, "Evaluation of contrast sensitivity functions for the formulation of quality measures incorporated in halftoning algorithms," in *Proc. IEEE ICASSP '93*, vol. 5, pp. 301–304, 1993.
- [36] T. D. Kite, N. Damera-Venkata, B. L. Evans, and A. C. Bovik, "Image quality assessment based on a degradation model," *IEEE Trans. Image Processing*, vol. 9, pp. 636–650, Apr. 2000.
- [37] Z. Wang and A. C. Bovik, "A universal image quality index," *IEEE Signal Processing Lett.*, vol. 9, pp. 81–84, Mar. 2002.
- [38] V. Monga, N. Damera-Venkata, and B. L. Evans., *Halftoning Toolbox for MATLAB*. 2002. Available: www.ece.utexas.edu/~bevans/projects/halftoning/toolbox.
- [39] P. W. Wong, "Inverse halftoning and kernel estimation for error diffusion," *IEEE Trans. Image Processing*, vol. 4, pp. 486–498, Apr. 1995.



An open letter concerning
Mass matrix transforms in qubit field theory

Marni Sheppeard

To whomever ...

This small paper reports on the initial observation that Carl Brannen's mass operators are naturally expressed as discrete Fourier series, common in the theory of quantum computation. Our obsession with these simple matrices has generated a great deal of criticism. Lubos Motl, the string theorist, called us F-ing Crackpots on my blog, Arcadian Functor, and all attempts to have this paper endorsed for the preprint arxiv failed. Although it is to be recognized that the abundance of errors in much of my writing is regrettable, in my experience these errors are never corrected by the people who think that this work is trivial and wrong, because it would be beneath them to consider it seriously.

The difficulty here is that our motivation for studying these mass operators lies not in standard particle physics, nor in standard theories of gravity, neither of which have anything whatsoever to say about the rest masses of fundamental particles. Unfortunately, a basic idea in quantum field theory is that certain parameters, such as mass, vary continuously, in a very complex way, from singular raw values. An unwavering belief in this idea leads people to conclude that simple formulae for rest masses cannot exist. Of those knowledgeable people willing to consider that such formulae may exist, most appear to believe that one should make no attempt to publish papers about it until one has constructed a complete theory of quantum gravity.

Since Carl Brannen, and many others, have now traveled a fair distance down this road, it seems that a rather impressive theory will actually exist before a single paper is published in a highly regarded peer reviewed journal. Fortunately, thanks to the Internet Age, a rapidly growing number of people are now working on this subject.

Mass matrix transforms in qubit field theory

M. D. Sheppard

Department of Physics and Astronomy

University of Canterbury,

Christchurch, New Zealand

Circulant mass matrices for triples of charged and neutral leptons have been studied in the context of qubit quantum field theory. This note describes the discrete Fourier transform behind such matrices, and discusses a category theoretic interpretation of these operators.

PACS numbers: 03.67.-a, 03.67.Lx, 04.60.-m

INTRODUCTION

Using a measurement algebra approach to QFT, Brannen [1] recently recovered the Koide [2] formula

$$(\sqrt{m_e} + \sqrt{m_\mu} + \sqrt{m_\tau})^2 = \frac{3}{2}(m_e + m_\mu + m_\tau) \quad (1)$$

for charged lepton masses in the form of a 3×3 circulant complex matrix, whose eigenvalues squared give the lepton masses to experimental precision. This analysis was extended to a set of three neutrinos, and the mass ratio predictions agree with preliminary neutrino oscillation data.

Here it is observed that the discrete Fourier transform [3] provides a further interpretation of the mass matrices, both as a duality between operators and eigenvalues and also as a link to the theory of quantum computation [4].

It is expected that other triples of Standard Model particles, namely baryons and mesons, will also be associated with 3×3 matrix operators of the same kind in accord with their preon structure [1] and the association of spatial directions to the number of particle generations, given by the three primitive idempotents of the measurement algebra.

FOURIER TRANSFORMS AND MASS MATRICES

A *circulant* matrix is built from its first row by adding cyclic permutations. In particular, a 3×3 circulant takes the form

$$\begin{pmatrix} A & B & C \\ C & A & B \\ B & C & A \end{pmatrix} \quad (2)$$

where A , B and C will be complex numbers. Note that any such circulant is a combination of the three permutations (123), (231) and (312). For real eigenvalues λ_k it is essential that A be real and $C = \bar{B}$. Thus a mass matrix [1] takes the form

$$C = \eta \begin{pmatrix} 1 & re^{i\theta} & re^{-i\theta} \\ re^{-i\theta} & 1 & re^{i\theta} \\ re^{i\theta} & re^{-i\theta} & 1 \end{pmatrix} \quad (3)$$

for real η , r and θ . In terms of these parameters, the eigenvalues are given by

$$\lambda_k = \eta(1 + 2r\cos(\theta + \frac{2\pi k}{3}))$$

The Koide formula (1) follows when $r^2 = \frac{1}{2}$ and this choice may be applied also to the neutrino matrix.

In the $n \times n$ case, the discrete Fourier transform [3][4] interchanges the set of eigenvalues λ_k (assumed distinct) and matrix entries $A_1, A_2, A_3, \dots, A_n$ via

$$\begin{aligned} \lambda_k &= \sum_j e^{\frac{2\pi ijk}{n}} A_j \\ A_j &= \frac{1}{n} \sum_k e^{-\frac{2\pi ijk}{n}} \lambda_k \end{aligned} \quad (4)$$

Viewing the eigenvalues as a diagonal matrix, the transform interchanges the bases of projection operators and cyclic permutations. For real eigenvalues (m_1, m_2, m_3) with $m_i = \lambda_i^2$ in the above, and letting $\omega = e^{\frac{2\pi i}{3}}$, the transform takes the diagonal matrix to the circulant

$$\begin{pmatrix} m_1 + m_2 + m_3 & m_1\omega + m_2\omega^2 + m_3 & m_1\omega^2 + m_2\omega + m_3 \\ m_1\omega^2 + m_2\omega + m_3 & m_1 + m_2 + m_3 & m_1\omega + m_2\omega^2 + m_3 \\ m_1\omega + m_2\omega^2 + m_3 & m_1\omega^2 + m_2\omega + m_3 & m_1 + m_2 + m_3 \end{pmatrix}$$

which must be the square of (3) since the square of a circulant is a circulant. Thus a choice of scale is specified by $\eta = \frac{1}{3}(m_1 + m_2 + m_3)$.

A 3×3 matrix is viewed as a function on the discrete torus $\mathbb{Z}_3 \times \mathbb{Z}_3$, which has a quantum description and a convolution product for matrices [3]. Letting $D_{ij} = \delta_{ij}\omega^i$ there is a Weyl rule

$$D \circ (312) = \omega(312) \circ D$$

where the phase $\frac{2\pi}{3}$ is proportional to \hbar^{-1} . This associates Planck's constant with a hierarchy \mathbb{N} determined by the size of the matrix, but the continuum limit is obtained via $\hbar \rightarrow \infty$ rather than $\hbar \rightarrow 0$.

If masses are to be thought of as quantum numbers, then why are their values so awkward in comparison to, say, spin? For 2×2 circulants with entries A and B , the eigenvectors are $(1, 1)$ and $(1, -1)$ with eigenvalues $(A + B)$ and $(A - B)$ respectively. For example, for the Pauli swap matrix σ_x , with $A = 0$, the spin eigenvalues are ± 1 . Complexity in the eigenvalue set only arises in dimension three or higher.

Degenerate eigenvalues $\frac{\lambda_k}{\eta} \in \{1 - r, 1 - r, 1 + 2r\}$ occur when $\theta = 0$ and all matrix entries are real. Although this pattern does not describe the leptons, we observe that it is the typical composition of masses for baryon constituents. Since such mass operators arise in a preon model that unifies particle structure, it is expected that all standard model bound states and resonances may be arranged into mass triples.

In quantum computation [4] a Fourier transform is also defined in this way, acting on a set of n basis states. For example, an N qubit computation uses $n = 2^N$ basis states. The transform is unitary and it may be built from unitary gates, namely the Hadamard gate $H = \frac{1}{\sqrt{2}}(\sigma_x + \sigma_z)$ and the series

$$B_k = \begin{pmatrix} 1 & 0 \\ 0 & e^{\frac{2\pi i}{2^k}} \end{pmatrix}$$

By analogy, a mass computation with 3^N basis states uses ternary digits, so the gates B_k would be replaced by gates

$$T_k = \begin{pmatrix} 1 & 0 & 0 \\ 0 & e^{\frac{2\pi i}{3^k}} & 0 \\ 0 & 0 & e^{\frac{4\pi i}{3^k}} \end{pmatrix} \quad (5)$$

which are also unitary. In general, the Fourier operator entries F_{ij} are given by ω^{ij} , and the theory of *mutually unbiased bases* generalises the Pauli operator algebra in all prime power dimensions.

A basic time evolution operator exists for each dimension n . Note, however, that unlike in conventional constructions, this local evolution is not in any way associated with an emergent

cosmic clock, the latter being more closely related to the scale \hbar , given here by the matrix dimension. That is, this approach does not assume a globally defined time for a nonsensical universal observer.

DISCUSSION

The mass matrices arise from a one dimensional discrete transform, which itself involves commutative variables. However, it is seen that phase space variables satisfy the Weyl algebra of the quantum plane. Is there a noncommutative transform that extends this analysis to nonclassical underlying spaces? This is relevant to the question of extending the perturbative rest mass computations [1] to nonperturbative regimes.

Kapranov [5] has recently considered path spaces approximated by cubical paths, each of which is represented by a noncommutative monomial in the spatial directions. In dimension $d > 1$ a noncommutative Fourier transform relates measures on the space of paths to functions of the noncommuting variables. The basic idea is that a path integral is just a map from a noncommutative ring to a suitable commutative subring. In this way, particle masses [1] could arise as path integral invariants.

Taking T-duality seriously, one also expects to deal with nonassociativity. From a category theoretic point of view, both noncommutative and nonassociative structures can be dealt with in a unified framework. The cohomological element of interest here is the parity cube axiom, which describes the now familiar pentagon law on five of its faces. In a sufficiently lax algebraic setting, such as for tetracategories, the sixth face may break this law, providing the deformation parameter that turns a pentagon into a hexagon representing the permutation group S_3 [6].

The generation count by primitive idempotents [1] is confirmed by the string theoretic index theorem argument applied to the Riemann moduli space of the six punctured sphere, which has an orbifold Euler characteristic [7] of -6. The six punctures are associated to the six faces of a cube via a dual vertex, which is thickened to a sphere. Note that cohomological integrals for such moduli spaces commonly appear in QFT computations as multiple zeta values and polylogarithms.

For helpful discussions I thank Carl Brannen, Michael Rios, Matti Pitkanen, Tony Smith and Louise Riofrio.

[1] C. A. Brannen, <http://brannenworks.com/dmaa.pdf>.

[2] Y. Koide, *Lett. Nuovo Cim.* **34**, 201 (1982).

- [3] R. Aldrovandi, *Special matrices of mathematical physics* (World Scientific, 2001).
- [4] M. A. Nielsen and I. L. Chuang, *Quantum computation and quantum information* (Cambridge, 2000).
- [5] M. Kapranov, math.QA/0612411.
- [6] M. Batanin, math.CT/0301221.
- [7] M. Mulase and M. Penkava, math-ph/9811024.