
GENERALIZED PERSISTENCE ALGORITHM FOR DECOMPOSING MULTIPARAMETER PERSISTENCE MODULES

Tamal K. Dey
tamaldey@purdue.edu

Cheng Xin
xinc@purdue.edu

Department of Computer Science
Purdue University

ABSTRACT

The classical persistence algorithm computes the unique decomposition of a persistence module implicitly given by an input simplicial filtration. Based on matrix reduction, this algorithm is a cornerstone of the emergent area of topological data analysis. Its input is a simplicial filtration defined over the integers \mathbb{Z} giving rise to a 1-parameter persistence module. It has been recognized that multiparameter version of persistence modules given by simplicial filtrations over d -dimensional integer grids \mathbb{Z}^d is equally or perhaps more important in data science applications. However, in the multiparameter setting, one of the main challenges is that topological summaries based on algebraic structure such as decompositions and bottleneck distances cannot be as efficiently computed as in the 1-parameter case because there is no known extension of the persistence algorithm to multiparameter persistence modules. We present an efficient algorithm to compute the unique decomposition of a finitely presented persistence module M defined over the multiparameter \mathbb{Z}^d . The algorithm first assumes that the module is presented with a set of N generators and relations that are *distinctly graded*. Based on a generalized matrix reduction technique it runs in $O(N^{2\omega+1})$ time where $\omega < 2.373$ is the exponent for matrix multiplication. This is much better than the well known algorithm called Meataxe which runs in $\tilde{O}(N^{6(d+1)})$ time on such an input. In practice, persistence modules are usually induced by simplicial filtrations. With such an input consisting of n simplices, our algorithm runs in $O(n^{(d-1)(2\omega+1)})$ time for $d \geq 2$. For the special case of zero dimensional homology, it runs in time $O(n^{2\omega+1})$.

1 Introduction

Persistence modules defined over a single parameter such as \mathbb{Z} or \mathbb{R} have become a central object of study in topological data analysis (TDA). It is an indexed set of vector spaces connected by linear maps most commonly arising from applying a homology functor to a simplicial filtration—another well known construct in TDA. Under some mild conditions [47], such a module decomposes uniquely into interval modules called *bars*. These bars or its equivalent persistence diagrams encode the input module completely. Naturally, computing these bars from an input persistence module efficiently becomes an important endeavor in TDA. Starting with the persistence algorithm [29], a number of improvements and extensions have been proposed for computing the bar decompositions in the single parameter case. However, the problem in the multi parameter case has not received as much attention. Other than some specific cases [9, 13, 20, 26], the only known algorithm for the purpose can be derived from the so-called Meataxe algorithm which applies to much more general modules than we consider in TDA at the expense of high computational cost. Sacrificing this generality and still encompassing a large class of modules that appear in TDA, we can design a much more efficient algorithm. Specifically, we present an algorithm that can decompose a finitely presented module (unique decomposition is guaranteed by Krull-Schmidt theorem [2]) with a time complexity that is far better than the Meataxe algorithm though we lose the generality as the module needs to be *distinctly graded*, that is, no two generators and no two relations of the module have the same grade. If this condition is not satisfied, a simple observation implies that the algorithm still produces an output that can be viewed as a decomposition of a module close under the interleaving distance.

For measuring algorithmic efficiency, it is imperative to specify how the input module is presented. Assuming an index set of size m and vector spaces of dimension $O(s)$, a one-parameter persistence module can be presented by a set of $O(s) \times O(s)$ matrices each representing a linear map $M_i \rightarrow M_{i+1}$ between two consecutive vector spaces M_i and M_{i+1} . This input format is costly as it takes $O(ms^2)$ space ($O(s^2)$ -size matrix for each index) and also does not appear to offer any benefit in time complexity for computing the bars. An alternative presentation is obtained by considering the persistence module as a graded module over a polynomial ring $\mathbb{k}[t]$ and presenting it with the so-called *generators* $\{g_i\}$ of the module and *relations* $\{\sum_i \alpha_i g_i = 0 \mid \alpha_i \in \mathbb{k}[t]\}$ among them. A presentation matrix encoding the relations in terms of the generators characterizes the module completely. Then, a matrix reduction algorithm akin to the persistence algorithm [22] provides the desired decomposition. Figure 1 illustrates the advantage of this presentation over the other costly presentation. In practice, when the 1-parameter persistence module is given by an implicit simplicial filtration, one can apply the matrix reduction algorithm directly on a boundary matrix rather than first computing a presentation matrix from it and then decomposing it. If there are $O(N)$ generators (creator simplices) and relations (destructor simplices), the algorithm runs in $O(N^3)$ time with simple matrix reductions and in $O(N^\omega)$ time with more sophisticated matrix multiplication techniques where $\omega < 2.373$ is the exponent for matrix multiplication.

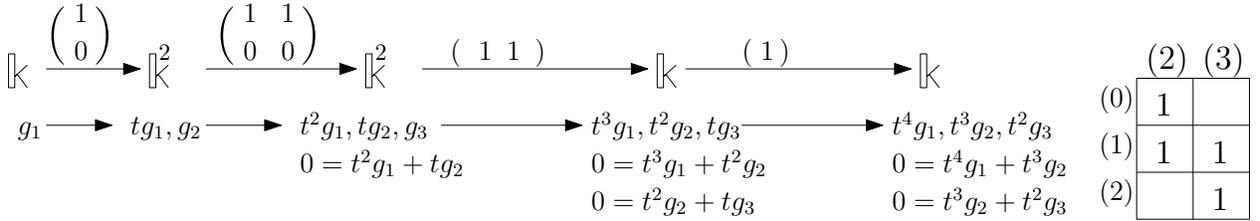


Figure 1: Costly presentation (top) vs. graded presentation (bottom, right). The top chain can be summarized by three generators g_1, g_2, g_3 at grades (0), (1), (2) respectively, and two relations $0 = t^2g_1 + tg_2, 0 = t^2g_2 + tg_3$ at grades (2), (3) respectively. The grades of generators and relations are given by the first times they appear in the chain. Finally, this information can be summarized succinctly by the presentation matrix on the right.

The Meataxe algorithm for multiparameter persistence modules follows the costly approach analogous to the one in the one-parameter case that expects the presentation of each individual linear map explicitly. In particular, it expects the input d -parameter module M over a finite subset of \mathbb{Z}^d to be given as a large matrix in $\mathbb{k}^{D \times D}$ with entries in a fixed field $\mathbb{k} = \mathbb{F}_q$, where D is the sum of dimensions of vector spaces over all points in \mathbb{Z}^d supporting M . The time complexity of the Meataxe algorithm is $O(D^6 \log q)$ [33]. In general, D might be quite large. It is not clear what is the most efficient way to transform an input that specifies generators and relations (or a simplicial filtration) to a representation matrix required by the Meataxe algorithm. A naive approach is to consider the minimal subgrid in \mathbb{Z}^d that supports the non-trivial maps. In the worst-case, with N being the total number of generators and relations, one has to consider $O(\binom{N}{d}) = O(N^d)$ grid points in \mathbb{Z}^d each with a vector space of dimension $O(N)$. Therefore, $D = O(N^{d+1})$ giving a worst-case time complexity of $O(N^{6(d+1)} \log q)$. Even allowing approximation, the algorithm [34] runs in $O(N^{3(d+1)} \log q)$ time.

We take the alternate approach where the module is treated as a finitely presented graded module over multivariate polynomial ring $R = \mathbb{k}[t_1, \dots, t_d]$ [23] and presented with a set of generators and relations graded appropriately. The fact that the persistence modules in TDA can be modeled as a graded module studied in algebraic geometry and commutative algebra [30, 42] was recognized in [16, 17, 37] and further studied in [38, 40]. Given a presentation matrix encoding relations with generators, our algorithm computes a diagonalization of the matrix giving a presentation of each module called *indecomposable* which the input module decomposes into. These indecomposables are the higher dimensional analogues of the bars. Compared to the one-parameter case, we have to cross two main barriers for computing the indecomposables. First, we need to allow row operations along with column operations for reducing the input matrix. In one-parameter case, row operations become redundant because column operations already produce the bars. Second, differently from the one-parameter case, we cannot allow all left-to-right column or bottom-to-top row operations for the matrix reduction because the parameter space $\mathbb{Z}^d, d > 1$, unlike \mathbb{Z} only induces a partial order on these operations. We show how these two difficulties can be overcome by an incremental approach combined with a linearization trick. Given a presentation matrix with a total of N generators and relations that are graded distinctly, our algorithm runs in $O(N^{2\omega+1})$ time. Surprisingly, the complexity does not depend on the parameter d . with time complexity $O(N^{2\omega+1})$ where $\omega < 2.373$ is the matrix multiplication exponent.

In practice, we are often given a simplicial filtration instead of a presentation matrix relating the generators of the induced persistence module. In this case, one has to compute presentation matrices from the input filtration consisting

of n simplices. For 2-parameter persistence modules, we can compute a presentation matrix of size $O(n) \times O(n)$ using the algorithm of Lesnick and Wright [40] in $O(n^3)$ time whereas for d -parameter persistence modules, we can adapt an algorithm of Skryzalin [45] to compute the presentation in $O(n^{d+1})$ time. For $d \geq 2$, this algorithm produces a presentation matrix of dimension $O(n^{d-1}) \times O(n^{d-1})$. Therefore, with $N = O(n^{d-1})$, the decomposition algorithm takes $O(n^{(d-1)(2\omega+1)})$ time. Combining the costs for computing a presentation and its decomposition, the time complexity of our algorithm becomes $O(n^{(d-1)(2\omega+1)})$ for $d \geq 2$. The time complexity of the Meataxe algorithm remains the same, $O(n^{6(d+1)} \log q)$, with a simplicial filtration of n simplices because the highest dimension of the vector space at each grid point is $O(n)$. Our algorithm is better than the Meataxe algorithm in this case too.

As a generalization of the traditional persistence algorithm, it is expected that our algorithm can be interpreted as computing invariants such as persistence diagrams [21] or barcodes [48]. A roadblock to this goal is that d -parameter persistence modules do not have complete discrete invariants for $d \geq 2$ [17, 38]. Consequently, one needs to invent other invariants suitable for multiparameter persistence modules. A natural way to generalize the invariant in traditional persistent homology would be to consider the decomposition and take the discrete invariants in each indecomposable component. This gives us invariants which are no longer complete but still contain rich information.

We offer two interpretations of the output of our algorithm as two different invariants: *persistent graded Betti numbers* as a generalization of persistence diagrams and *blockcodes* as a generalization of barcodes. The persistent graded Betti numbers are linked to the graded Betti numbers studied in commutative algebra, which is introduced in TDA for multiparameter persistence modules in the work of [17] and [37]. The bigraded Betti numbers are further studied in [40]. By constructing the free resolution of a persistence module, we can compute its graded Betti numbers and then decompose them according to each indecomposable module, which results into the persistent graded Betti numbers. For each indecomposable, we apply the dimension function which is also known as the Hilbert function in commutative algebra to summarize the graded Betti numbers for each indecomposable module. This constitutes a blockcode for the indecomposable module of the persistence module. The blockcode is a good vehicle for visualizing lower dimensional persistence modules such as 2- or 3-parameter persistence modules.

1.1 Other related work

Since it is known that there is no complete discrete invariant for multiparameter persistence, researchers have proposed various reasonable summaries that can be computed in practice. Among them the rank invariant proposed by Carlsson et al. [16, 17] is a popular one. Cerri et al. [19] propose multiparameter persistent Betti number as a stable invariant. Lesnick and Wright introduce the computational tool of fibered barcode in [39, 40] as an interactive vehicle to visualize the one-parameter restriction of biparameter persistence modules.

Another related line of work focuses on defining distances and their stabilities in the space of multiparameter persistence modules. The interleaving distance [3, 4, 8, 38], and multi-matching distance [19, 18] are some of the work to mention a few. The relation between interleaving distance and bottleneck distance is studied in [8, 5, 11]. On the computational front, Dey and Xin showed that the bottleneck distance can be computed in polynomial time for the special cases of interval decomposable modules [26] though the general problem is proved to be NP-hard [3, 4]. A recent work of Kerber et al. shows that the matching distance [36] can be computed efficiently in polynomial time.

1.2 Outline

The rest of the paper is organized as follows. In Section 2, we introduce some background materials on persistence modules in the language of graded modules. In Section 3, we introduce the presentation of a persistence module and its presentation matrix which assists in computing the decomposition of persistence modules. The 1-1 correspondence between the decompositions of the persistence module and its presentation is a fundamental fact which is presented as our first main theorem. Based on this correspondence, we observe that two main computational problems need to be solved, (i) computing the decomposition of the presentation matrix, and (ii) constructing a valid presentation. In Section 4, we handle the first problem by designing an algorithm for computing a decomposition of the presentation matrix. We observe that this problem can be transformed to a what we call generalized matrix reduction problem. Based on that, we propose an algorithm to solve this problem and prove the correctness of our algorithm, and illustrate it with an example. In Section 5, we introduce the strategies for the second problem of computing presentations and analyze the total time complexity for computing presentations together with matrix reduction. In Section 6, we give two interpretations of the results of our decomposition of persistence modules as two different invariants, persistent graded Betti numbers as a generalization of persistence diagrams and blockcodes as a generalization of barcodes. In Section 7, we conclude with suggesting some future direction.

2 Persistence modules

We want to study the *total decomposition* of a persistence module arising from a simplicial filtration in the multi-parameter setting. We first present some preliminary concepts from commutative algebra that lay the foundation of this work. For more details on multiparameter persistent homology and commutative algebra, we refer the readers to [10, 17, 24, 42]. Mainly, we need the concept of *graded modules* because as in [17] we treat the familiar persistence modules in topological data analysis as graded modules. Let $R = \mathbb{k}[t_1, \dots, t_d]$ be the d -variate Polynomial ring for some $d \in \mathbb{Z}_+$ with \mathbb{k} being a field. Throughout this paper, we assume coefficients are in \mathbb{k} . Hence homology groups are vector spaces.

Definition 2.1. A \mathbb{Z}^d -graded R -module (graded module in brief) is an R -module M that is a direct sum of \mathbb{k} -vector spaces $M_{\mathbf{u}}$ indexed by $\mathbf{u} \in \mathbb{Z}^d$, i.e. $M = \bigoplus_{\mathbf{u}} M_{\mathbf{u}}$, such that the ring action satisfies that $\forall i, \forall \mathbf{u} \in \mathbb{Z}^d, t_i \cdot M_{\mathbf{u}} \subseteq M_{\mathbf{u}+e_i}$, where $\{e_i\}_{i=1}^d$ is the standard basis in \mathbb{Z}^d .

Another interpretation of graded module is that, for each $\mathbf{u} \in \mathbb{Z}^d$, the action of t_i on $M_{\mathbf{u}}$ determines a linear map $t_i \bullet : M_{\mathbf{u}} \rightarrow M_{\mathbf{u}+e_i}$ by $(t_i \bullet)(m) = t_i \cdot m$. So, we can also describe a graded module equivalently as a collection of vectors spaces $\{M_{\mathbf{u}}\}_{\mathbf{u} \in \mathbb{Z}^d}$ with a collection of linear maps $\{t_i \bullet : M_{\mathbf{u}} \rightarrow M_{\mathbf{u}+e_i}, \forall i, \forall \mathbf{u}\}$ where the commutative property $(t_j \bullet) \circ (t_i \bullet) = (t_i \bullet) \circ (t_j \bullet)$ holds. The commutative diagram in Figure 2 shows a graded module for $d = 2$, also called a bigraded module.

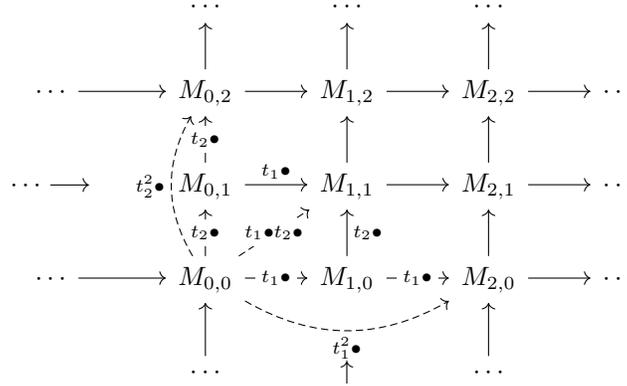


Figure 2: A graded 2-parameter module. All sub-diagrams of maps and compositions of maps are commutative.

We call a graded module M *finitely generated* if there exists a finite set of elements $\{g_1, \dots, g_n\} \subseteq M$ such that each element $m \in M$ can be written as an R -linear combination of these elements, i.e. $m = \sum_{i=1}^n \alpha_i g_i$ with $\alpha_i \in R$. We call this set $\{g_i\}$ a *generating set* of M . In this paper, we assume that all modules are finitely generated. Such modules always admit a minimal generating set.

Definition 2.2. A graded module morphism, called *morphism* in short, between two modules M and N is defined as an R -linear map $f : M \rightarrow N$ preserving grades: $f(M_{\mathbf{u}}) \subseteq N_{\mathbf{u}}, \forall \mathbf{u} \in \mathbb{Z}^d$. Equivalently, it can also be described as a collection of linear maps $\{f_{\mathbf{u}} : M_{\mathbf{u}} \rightarrow N_{\mathbf{u}}\}$ which gives the following commutative diagram for each \mathbf{u} and i :

$$\begin{array}{ccc} M_{\mathbf{u}} & \xrightarrow{t_i} & M_{\mathbf{u}+e_i} \\ f_{\mathbf{u}} \downarrow & & \downarrow f_{\mathbf{u}+e_i} \\ N_{\mathbf{u}} & \xrightarrow{t_i} & N_{\mathbf{u}+e_i} \end{array}$$

Two graded modules M, N are isomorphic if there exist two morphisms $f : M \rightarrow N$ and $g : N \rightarrow M$ such that $g \circ f$ and $f \circ g$ are identity maps.

For a graded module M , define a shifted graded module $M_{\rightarrow \mathbf{u}}$ for some $\mathbf{u} \in \mathbb{Z}^d$ by requiring $(M_{\rightarrow \mathbf{u}})_{\mathbf{v}} = M_{\mathbf{v}-\mathbf{u}}$ for each \mathbf{v} .

Definition 2.3 (Free module). We say a graded module is *free* if it is isomorphic to the direct sum of a collection of $R_{\rightarrow \mathbf{u}_j}$'s for some \mathbf{u}_j 's in \mathbb{Z}^d , denoted as $\bigoplus_j R_{\rightarrow \mathbf{u}_j}$.

Definition 2.4 (homogeneous). We say an element $m \in M$ is *homogeneous* if $m \in M_{\mathbf{u}}$ for some $\mathbf{u} \in \mathbb{Z}^d$. We denote $\text{gr}(m) = \mathbf{u}$ as the *grade* of such homogeneous element. To emphasize the grade of a homogeneous element, we also write $m^{\text{gr}(m)} := m$.

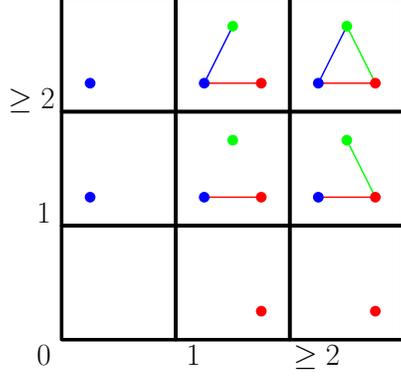


Figure 3: The working example on a 2-parameter simplicial filtrations. Each square box indicates what is the current (filtered) simplicial complex at the grade of the box. It has one connected component in 0th homology groups at grades except $(0, 0)$ and $(1, 1)$, and has two connected components at grade $(1, 1)$.

A minimal generating set of a free module is called a *basis*. We usually further require that all the elements in a basis, also called *generators*, are homogeneous. For a free module $F \simeq \bigoplus_j R_{\rightarrow \mathbf{u}_j}$, $\{e_j : j = 1, 2, \dots\}$ is a homogeneous basis of F , where e_j indicates the multiplicative identity in $R_{\rightarrow \mathbf{u}_j}$. The generating set $\{e_j : j = 1, 2, \dots\}$ is often referred to as the standard basis of $\bigoplus_j R_{\rightarrow \mathbf{u}_j} = \langle \{e_j : j = 1, 2, \dots\} \rangle$.

A *d-parameter persistence module* is a graded R -module obtained by applying the homology functor with some field \mathbb{k} on a d -parameter simplicial filtration defined below. In this paper, it can be treated as a synonym for a \mathbb{Z}^d -graded R -module. Formally, a (d -parameter) *simplicial filtration* is a family of simplicial complexes $\{X_{\mathbf{u}}\}_{\mathbf{u} \in \mathbb{Z}^d}$ such that for each grade $\mathbf{u} \in \mathbb{Z}^d$ and each $i = 1, \dots, d$, $X_{\mathbf{u}} \subseteq X_{\mathbf{u}+e_i}$. We obtain a simplicial chain complex $(C_{\bullet}(X_{\mathbf{u}}), \partial_{\bullet})$ for each $X_{\mathbf{u}}$ in this simplicial filtration. For each chain complex $C_{\bullet}(X_{\mathbf{u}})$, we have the cycle spaces $Z_p(X_{\mathbf{u}})$'s and boundary spaces $B_p(X_{\mathbf{u}})$'s as kernels and images of boundary maps ∂_p 's respectively, and the homology group $H_p(X_{\mathbf{u}}) = Z_p(X_{\mathbf{u}})/B_p(X_{\mathbf{u}})$ as the cokernel of the inclusion maps $B_p(X_{\mathbf{u}}) \hookrightarrow Z_p(X_{\mathbf{u}})$. Taking $M_{\mathbf{u}} = H_p(X_{\mathbf{u}})$ and the linear maps $H_p(X_{\mathbf{u}}) \rightarrow H_p(X_{\mathbf{v}})$ induced by inclusions $X_{\mathbf{u}} \subseteq X_{\mathbf{v}}$ define a d -parameter persistence module. More details of this construction is given later in Section 5.

For illustration purpose, we describe a working example of a 2-parameter persistence module induced from a 2-parameter simplicial filtration shown in Figure 3. We will use this example throughout to show its induced persistence module and computational results of our algorithm. Later in the context, when we mention an example without reference, we refer to this working example.

Example 1. [Working example] In practice, the most common simplicial filtration is obtained from the sublevel sets $\{X_{\mathbf{u}} := f^{-1}(-\infty, \mathbf{u}]\}_{\mathbf{u} \in \mathbb{R}^n}$ of a given (one-critical) filtration function $f : X \rightarrow \mathbb{Z}^d$ on a topological space represented by a simplicial complex X .

For example, let the space X be a simplicial 1-complex with 0-simplices consisting of three vertices, blue vertex v_b , red vertex v_r , and green vertex v_g , connected by three edges, blue edge e_b , red edge e_r , and green edge e_g as 1-simplices. Assign a filtration function $f : X \rightarrow \mathbb{Z}^2$ as follows:

$$\begin{aligned} f(v_b) &= (0, 1), f(v_r) = (1, 0), f(v_g) = (1, 1) \\ f(e_b) &= (1, 2), f(e_r) = (1, 1), f(e_g) = (2, 1) \end{aligned}$$

Based on this filtration function, the subcomplex $X_{\mathbf{u}}$ for each $\mathbf{u} \in \mathbb{Z}^2$ is illustrated in Figure 3. Take vertices as basis of each $C_0(X_{\mathbf{u}})$ and edges as basis of $C_1(X_{\mathbf{u}})$. Recall that to emphasize the grades, we denote $v_*^{\mathbf{u}} \in C_0(X_{\mathbf{u}})$ to be the basic element in the vector space $C_0(X_{\mathbf{u}})$. All these $v_*^{\mathbf{u}}$ are homogeneous element in the graded module $C_0(X)$. For each vertex $v_* \in \{v_b, v_r, v_g\}$, there is a unique smallest grade $\text{gr}(v_*) \triangleq f(v_*)$ such that $v_*^{\text{gr}(v_*)}$ is a homogeneous basic element in $C_0(X_{\text{gr}(v_*)})$ and $\mathbf{u}' \not\geq \text{gr}(v_*) \implies v_* \notin C_0(X_{\mathbf{u}'})$. We call this grade $\text{gr}(v_*)$ the birth time of v_* . Then for all $\mathbf{u} \geq \text{gr}(v_*)$, by the definition of scalar multiplication of graded modules, $\mathbf{t}^{\mathbf{u}-\text{gr}(v_*)} v_*^{\text{gr}(v_*)} = v_*^{\mathbf{u}} \in C_0(X_{\mathbf{u}})$ is the image of $v_*^{\text{gr}(v_*)}$ under the inclusion map $C_0(X_{\text{gr}(v_*)}) \hookrightarrow C_0(X_{\mathbf{u}})$, which is the homogeneous basis element of $C_0(X_{\mathbf{u}})$ corresponding to the vertex v_* . Sometime, we omit the upper index by writing $v_* = v_*^{f(v_*)}$ for brevity. With these conventions, we can see that for each $\mathbf{u} \in \mathbb{Z}^2$, the vector space $C_0(X_{\mathbf{u}})$ is generated by all $v_*^{\mathbf{u}} = \mathbf{t}^{\mathbf{u}-\text{gr}(v_*)} v_*^{\text{gr}(v_*)}$ such that v_* is born before or at \mathbf{u} , which means $\text{gr}(v_*) \leq \mathbf{u}$. In fact, $C_0(X)$ is a free module

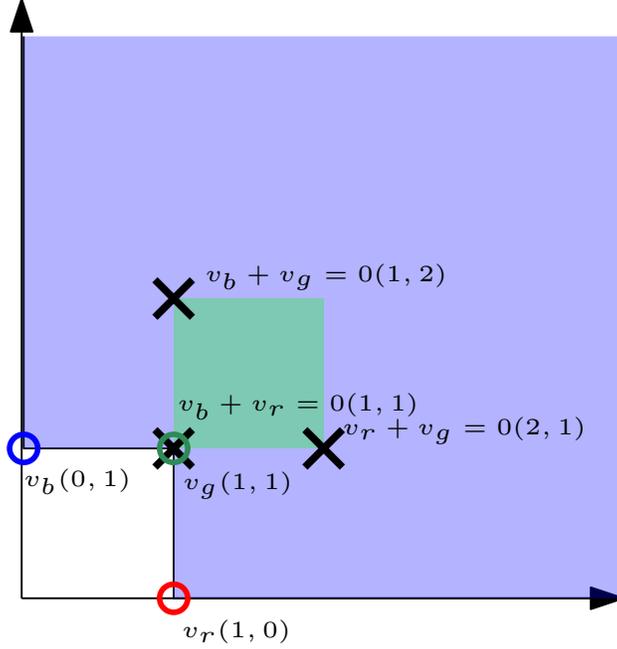


Figure 4: Persistence module whose presentation matrix is $[\partial_1]$ described in the working example.

and $\{v_b^{\text{gr}(v_b)}, v_r^{\text{gr}(v_r)}, v_g^{\text{gr}(v_g)}\}$ forms a basis of $C_0(X)$. That means, any element of $C_0(X)$ can be written as a R -linear combination of these $v_*^{\text{gr}(v_*)}$'s and all these $v_*^{\text{gr}(v_*)}$'s are linearly independent.

Similarly, for each $e_* \in \{e_b, e_r, e_g\}$, we have the earliest basic element $e_*^{\text{gr}(e_*)}$ of $C_1(X_{\text{gr}(e_*)})$ for $\text{gr}(e_*) \triangleq f(e_*)$. The set $\{e_b^{\text{gr}(e_b)}, e_r^{\text{gr}(e_r)}, e_g^{\text{gr}(e_g)}\}$ forms a basis of the free module $C_1(X)$. Furthermore, by the commutative property of morphisms, we have for each $\mathbf{u} \geq \text{gr}(e_*)$,

$$\partial_{1,\mathbf{u}}(e_*^{\mathbf{u}}) = \partial_{1,\mathbf{u}}(\mathbf{t}^{\mathbf{u}-\text{gr}(e_*)} e_*^{\text{gr}(e_*)}) = \mathbf{t}^{\mathbf{u}-\text{gr}(e_*)} \circ \partial_{1,\text{gr}(e_*)}(e_*) = \mathbf{t}^{\mathbf{u}-\text{gr}(e_*)} \circ \partial_1(e_*^{\text{gr}(e_*)}) \quad (1)$$

In fact, as a morphism between two free modules, ∂_1 is fully determined by $\partial_1(e_*)$. Consider, for example, the red edge e_r connecting v_b and v_r . With the field chosen to be \mathbb{F}_2 , one has $\partial_1(e_r^{\text{gr}(e_r)}) = \partial_1(e_r^{(1,1)}) = v_b^{(1,1)} + v_r^{(1,1)} = \mathbf{t}^{(1,0)}v_b^{(0,1)} + \mathbf{t}^{(0,1)}v_r^{(1,0)}$. Similar for $\partial_1(e_b^{(1,2)})$ and $\partial_1(e_g^{(2,1)})$. Therefore, ∂_1 can be represented as a matrix with entries in $R = \mathbb{k}[\mathbf{t}]$ as follows:

$$[\partial_1] \begin{matrix} e_r^{(1,1)} & e_b^{(1,2)} & e_g^{(2,1)} \\ v_b^{(0,1)} \\ v_r^{(1,0)} \\ v_g^{(1,1)} \end{matrix} \begin{pmatrix} \mathbf{t}^{(1,0)} & \mathbf{t}^{(1,1)} & 0 \\ \mathbf{t}^{(0,1)} & 0 & \mathbf{t}^{(1,1)} \\ 0 & \mathbf{t}^{(0,1)} & \mathbf{t}^{(1,0)} \end{pmatrix}$$

Now consider the 0^{th} persistence homology module induced from boundary morphism $\partial_1 : C_1(X) \rightarrow C_0(X)$. Note that the 0^{th} homology is a space of connected components. For each \mathbf{u} , $H_0(X_{\mathbf{u}}) = \frac{Z_0(X_{\mathbf{u}})}{B_0(X_{\mathbf{u}})} = \frac{C_0}{\text{im } \partial_{1,\mathbf{u}}}$. With bases of C_0 and C_1 chosen above, we have the 0^{th} persistence homology on grades from $(0,0)$ (bottom-left corner) to $(2,2)$ (top-right corner) described as following diagram (also illustrated in Figure 4):

$$\begin{array}{ccccc}
\mathbb{k} & \xrightarrow{1} & \mathbb{k} & \xrightarrow{1} & \mathbb{k} \\
\uparrow & & \uparrow & & \uparrow \\
1 & & [1,1] & & 1 \\
\downarrow & & \downarrow & & \downarrow \\
\mathbb{k} & \xrightarrow{[1,0]^T} & \mathbb{k}^2 & \xrightarrow{[1,1]} & \mathbb{k} \\
\uparrow & & \uparrow & & \uparrow \\
0 & & [1,0]^T & & 1 \\
\downarrow & & \downarrow & & \downarrow \\
0 & \xrightarrow{0} & \mathbb{k} & \xrightarrow{1} & \mathbb{k}
\end{array}$$

In what follows, we take the liberty of omitting X and p if they are clear from the context. Thus, we may denote $Z_p(X)$, $B_p(X)$, and $H_p(X)$ as Z , B and H respectively.

Definition 2.5 (decomposition). For a finitely generated module M , we call $M \simeq \bigoplus M^i$ a *decomposition* of M for some collection of modules $\{M^i\}$. We say a module M is *indecomposable* if $M \simeq M^1 \oplus M^2 \implies M^1 = 0$ or $M^2 = 0$. By the Krull-Schmidt theorem [2], there exists an essentially unique (up to permutation and isomorphism) decomposition $M \simeq \bigoplus M^i$ with every M^i being indecomposable. We call it the *total* decomposition of M .

For example, the free module R is generated by $\langle e_1^{(0,0)} \rangle$ and the free module $R_{\rightarrow(0,1)} \oplus R_{\rightarrow(1,0)}$ is generated by $\langle e_1^{(0,1)}, e_2^{(1,0)} \rangle$. A free module M generated by $\langle e_j^{\mathbf{u}_j} : j = 1, 2, \dots \rangle$ has a (total) decomposition $M \simeq \bigoplus_j R_{\rightarrow \mathbf{u}_j}$.

Definition 2.6. Two morphisms $f : M \rightarrow N$ and $f' : M' \rightarrow N'$ are isomorphic, denoted as $f \simeq f'$, if there exist isomorphisms $g : M \rightarrow M'$ and $h : N \rightarrow N'$ such that the following diagram commutes:

$$\begin{array}{ccc}
M & \xrightarrow{f} & N \\
\cong \downarrow & & \downarrow \cong \\
M' & \xrightarrow{f'} & N'
\end{array}$$

Essentially, like isomorphic modules, two isomorphic morphisms can be considered the same. For two morphisms $f_1 : M^1 \rightarrow N^1$ and $f_2 : M^2 \rightarrow N^2$, there exists a canonical morphism $g : M^1 \oplus M^2 \rightarrow N^1 \oplus N^2$, $g(m_1, m_2) = (f_1(m_1), f_2(m_2))$, which is essentially uniquely determined by f_1 and f_2 and is denoted as $f_1 \oplus f_2$. We denote a trivial module by bold $\mathbf{0}$, and a trivial morphism by $\mathbf{0}$. Analogous to the decomposition of a module, we can also define a decomposition of a morphism.

Definition 2.7. A morphism f is indecomposable if $f \simeq f_1 \oplus f_2 \implies f_1$ or f_2 is the trivial morphism $\mathbf{0} : \mathbf{0} \rightarrow \mathbf{0}$. We call $f \simeq \bigoplus f_i$ a decomposition of f . If each f_i is indecomposable, we call it a *total* decomposition of f .

Like decompositions of modules, the total decompositions of a morphism is also essentially unique.

3 Presentation and its decomposition

To study total decompositions of persistence modules as graded modules, we borrow the idea of *presentations* of graded modules and build a bridge between decompositions of persistence modules and corresponding presentations. The later ones can be transformed to a matrix reduction problem with possibly nontrivial constraints which we will introduce in Section 4. Our first main result is that there are 1-1 correspondences between persistence modules, presentations, and presentation matrices. Recall that, by assumption, all modules are finitely generated. A graded module hence a persistence module accommodates a description called its *presentation* that aids finding its decomposition.

Definition 3.1. A presentation of a graded module H is an exact sequence $F^1 \xrightarrow{f} F^0 \twoheadrightarrow H \rightarrow 0$. We call f a presentation map. We say a graded module H is *finitely presented* if there exists a presentation of H with both F^1 and F^0 being finitely generated.

It follows from the definition that a presentation of H is determined by the presentation map f where $\text{coker } f \simeq H$.

Remark 3.1. Presentations of a given graded module are not unique. However, there exists an essentially unique (up to isomorphism) presentation f of a graded module in the sense that any presentation f' of that module can be written as $f' \simeq f \oplus f''$ with $\text{coker } f'' = 0$. We call this unique presentation *the minimal presentation*. See more details of the construction and properties of minimal presentation in Appendix A.1.

Fixed bases of nonzero free modules F^1 and F^0 provide a matrix form $[f]$ of the presentation map f , which we call a *presentation matrix* of H . It has entries in R . In the special case that H is a free module with F^1 being a zero module, we define the presentation matrix $[f]$ of H to be a *null column matrix* with matrix size $\ell \times 0$ for some $\ell \in \mathbb{N}$. An important property of a persistence module H is that a decomposition of its presentation f corresponds to a decomposition of H itself. The decomposition of f can be computed by *diagonalizing* its presentation matrix $[f]$. Informally, a diagonalization of a matrix \mathbf{A} is an equivalent matrix \mathbf{A}' in the following form (see formal Definition 4.2 later):

$$\mathbf{A}' = \begin{bmatrix} \mathbf{A}_1 & 0 & \cdots & 0 \\ 0 & \mathbf{A}_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{A}_k \end{bmatrix}$$

All nonzero entries are in \mathbf{A}_i 's and we write $\mathbf{A} \simeq \bigoplus \mathbf{A}_i$. It is not hard to see that for a map $f \simeq \bigoplus f_i$, there is a corresponding diagonalization $[f] \simeq \bigoplus [f_i]$. With these definitions and the fact that persistence modules are graded modules, we have the following theorem that motivates our decomposition algorithm (proof in Appendix A).

Theorem 3.1. *There are 1-1 correspondences between the following three structures arising from a minimal presentation map $f : F^1 \rightarrow F^0$ of a persistence module H , and its presentation matrix $[f]$:*

1. A decomposition of the persistence module $H \simeq \bigoplus H^i$;
2. A decomposition of the presentation map $f \simeq \bigoplus f_i$
3. A diagonalization of the presentation matrix $[f] \simeq \bigoplus [f_i]$

Remark 3.2. In practice, we might be given a presentation which is not necessarily minimal. One way to handle this case is to compute the minimal presentation of the given presentation first. For 2-parameter modules, this can be done by the algorithm in [40]. The other choice is to compute the decomposition of the given presentation (not necessarily minimal) directly, which is sufficient to get the decomposition of the module thanks to the following proposition (proof at the end of Appendix A).

Proposition 3.2. *Let f be any presentation of a graded module H .*

1. For a decomposition of $H \simeq \bigoplus H^i$, there exists a decomposition of $f \simeq \bigoplus f^i$ so that $\text{coker } f^i = H^i, \forall i$.
2. The total decomposition of H follows from the total decomposition of f .

Remark 3.3. By Remark 3.1, any presentation f can be written as $f \simeq f^* \oplus f'$ with f^* being the minimal presentation and $\text{coker } f' = 0$. Furthermore, f' can be written as $f' \simeq g \oplus h$ where g is an identity map and h is a zero map. The corresponding matrix form is $[f'] \simeq [f^*] \oplus [g] \oplus [h]$ with $[g]$ being an identity submatrix and $[h]$ being a collection of zero column vectors. Therefore, one can easily read these trivial parts from the result of matrix diagonalization if it is total, meaning that none of its constituents (\mathbf{A}_i s) can be decomposed (diagonalized) further (Definition 4.2). See the following diagram for an illustration.

$$[f] = \left(\begin{array}{c|c|c} f^* & g & h \\ \hline [f^*] & & \\ \hline & 1 & \\ & & 1 \\ & & & 1 \end{array} \right)$$

It follows from Theorem 3.1 that we have to address the problem of *total diagonalization* of a presentation matrix $[f]$, $f : F^1 \rightarrow F^0$ for computing a total decomposition of the module H it represents. Each row r_i and column c_j of $[f]$ represent a generator g_i and a relation s_j respectively where $[f]_{ij} = \alpha_{ij}$ if $s_j = \sum_i \alpha_{ij} g_i$ for $\alpha_{ij} \in R$. We label the rows and columns with the grades of the elements they represent, that is, $\text{gr}(r_i) := \text{gr}(g_i)$ and $\text{gr}(c_j) := \text{gr}(s_j)$. Furthermore, we can simplify $[f]$ by observing that α_{ij} has the form $\alpha_{ij} = k \cdot t_1^{u_1} t_2^{u_2}$ where $k \in \mathbb{k}$ and $\mathbf{u} = (u_1, u_2) = \text{gr}(c_j) - \text{gr}(c_i)$. For all homogeneous transformations of bases, only the value of k changes in α_{ij} . Therefore, we can replace α_{ij} with the value of k which is either 0 or 1 when $\mathbb{k} = \mathbb{F}_2$. See the matrices in Example 2 given in Section 4.1.

With this change, the matrix $[f]$ becomes a matrix over the field \mathbb{F}_2 with the following operations for transformations as summarized below:

1. $[f]_{ij} = 1$ if and only if $\alpha_{ij} = 1$ in the relation $s_j = \sum_i \alpha_{ij} g_i$.
2. A column c_i can be added to column c_j , denoted as $c_i \rightarrow c_j$, only if $i \neq j$ and $\text{gr}(c_i) \leq \text{gr}(c_j)$. A row r_i can be added to row r_j denoted $r_i \rightarrow r_j$ only if $i \neq j$ and $\text{gr}(r_i) \leq \text{gr}(r_j)$.

4 Computing decomposition

In this section, we present an algorithm for computing a total decomposition of a distinctly graded module, which means that no two columns and no two rows in the presentation matrix have the same grades. All modules are assumed to be finitely presented and we take $\mathbb{k} = \mathbb{F}_2$ for simplicity though our method works for any finite field. We have observed that a total decomposition of a module can be achieved by computing a total decomposition of its presentation f . This in turn means a total diagonalization of the presentation matrix $[f]$. Here we formally define some notations about the diagonalization.

Given an $\ell \times m$ matrix $\mathbf{A} = [\mathbf{A}_{i,j}]$, with row indices $\text{Row}(\mathbf{A}) = [\ell] := \{1, 2, \dots, \ell\}$ and column indices $\text{Col}(\mathbf{A}) = [m] := \{1, 2, \dots, m\}$, we define an *index block* B of \mathbf{A} as a pair $[\text{Row}(B), \text{Col}(B)]$ with $\text{Row}(B) \subseteq \text{Row}(\mathbf{A}), \text{Col}(B) \subseteq \text{Col}(\mathbf{A})$. We say an index pair (i, j) is in B if $i \in \text{Row}(B)$ and $j \in \text{Col}(B)$, denoted as $(i, j) \in B$. We denote a *block* of \mathbf{A} on B as the matrix restricted to the index block B , i.e. $[\mathbf{A}_{i,j}]_{(i,j) \in B}$, denoted as $\mathbf{A}|_B$. We call B the index of the block $\mathbf{A}|_B$. We abuse the notations $\text{Row}(\mathbf{A}|_B) := \text{Row}(B)$ and $\text{Col}(\mathbf{A}|_B) := \text{Col}(B)$. For example, the i th row $r_i = \mathbf{A}_{i,*} = \mathbf{A}|_{[\{i\}, \text{Col}(\mathbf{A})]}$ and the j th column $c_j = \mathbf{A}_{*,j} = \mathbf{A}|_{[\text{Row}(\mathbf{A}), \{j\}]}$ are blocks with indices $[\{i\}, \text{Col}(\mathbf{A})]$ and $[\text{Row}(\mathbf{A}), \{j\}]$ respectively. Specifically, $[\emptyset, \{j\}]$ represents an index block of a single column j and $[\{i\}, \emptyset]$ represents an index block of a single row i . We call $[\emptyset, \emptyset]$ the empty index block.

A morphism can have presentation matrices in different equivalent forms depending on the bases chosen.

Definition 4.1. We say a matrix \mathbf{A}' is equivalent to \mathbf{A} , denoted as $\mathbf{A}' \sim \mathbf{A}$, if they present the same morphism.

Definition 4.2. A matrix $\mathbf{A}' \sim \mathbf{A}$ is called a *diagonalization* of \mathbf{A} with a set of nonempty index blocks $\mathcal{B} = \{B_1, B_2, \dots, B_k\}$ if rows and columns of \mathbf{A} are partitioned into these blocks, i.e., $\text{Row}(\mathbf{A}) = \coprod_i \text{Row}(B_i)$ and $\text{Col}(\mathbf{A}) = \coprod_i \text{Col}(B_i)$, and all the nonzero entries of \mathbf{A}' have indices in some B_i . We write $\mathbf{A}' = \bigoplus_{B_i \in \mathcal{B}} \mathbf{A}'|_{B_i}$. We say $\mathbf{A}' = \bigoplus_{B_i \in \mathcal{B}} \mathbf{A}'|_{B_i}$ is *total* if no block in this diagonalization can be decomposed further into smaller nonempty blocks. That means, for each block $\mathbf{A}'|_{B_i}$, there is no nontrivial diagonalization. Specifically, when \mathbf{A} is a null column matrix (the presentation matrix of a free module), we say \mathbf{A} is itself a total diagonalization with index blocks $\{[\{i\}, \emptyset] \mid i \in \text{Row}(\mathbf{A})\}$.

Note that each nonempty matrix \mathbf{A} has a trivial diagonalization with the set of index blocks being $\{(\text{Row}(\mathbf{A}), \text{Col}(\mathbf{A}))\}$. Guaranteed by Krull-Schmidt theorem [2], all total diagonalizations are unique up to permutations of their rows and columns, and equivalent transformation within each block. The total diagonalization of \mathbf{A} is denoted generically as \mathbf{A}^* . All total diagonalizations of \mathbf{A} have the same set of index blocks, denoted as \mathcal{B}^* , unique up to permutations of rows and columns.

4.1 Simplification of presentation matrix

First we want to transform the diagonalization problem to an equivalent problem that involves matrices with a simpler form. The idea is to simplify the presentation matrix to have entries only in \mathbb{k} . There is a correspondence between diagonalizations of the original presentation matrix and certain constrained diagonalizations of the corresponding transformed \mathbb{k} -matrix under this subset of basic operations.

Inspired by the ideas from [17], we first make some observations about the homogeneous property of presentation maps and presentation matrices. Equivalent matrices actually represent isomorphic presentations $f' \simeq f$ that admit commutative diagram,

$$\begin{array}{ccc} F^1 & \xrightarrow{f} & F^0 \\ \downarrow \simeq^{h^1} & & \downarrow \simeq^{h^0} \\ F^1 & \xrightarrow{f'} & F^0 \end{array}$$

where h^1 and h^0 are endomorphisms on F^1 and F^0 respectively. The endomorphisms are realized by basis changes between corresponding presentation matrices $[f] \simeq [f']$. Since all morphisms between graded modules are required

to be homogeneous by definition, we can use homogeneous bases (all the basis elements chosen are homogeneous elements¹) for F^0 and F^1 to represent matrices, say $F^0 = \langle g_1, \dots, g_n \rangle$ and $F^1 = \langle s_1, \dots, s_m \rangle$. Therefore, for simplicity we can consider only equivalent presentation matrices under homogeneous basis changes. Each entry $[f]_{i,j}$ is also homogeneous. That means, $[f]_{i,j} = t^{\mathbf{u}}$ with $\mathbf{u} = \text{gr}(s_j) - \text{gr}(g_i)$. We call such presentation matrix *homogeneous presentation matrix*.

For example, given $F^0 = \langle g_1^{(1,1)}, g_2^{(2,2)} \rangle$, the basis change $g_2^{(2,2)} \leftarrow g_2^{(2,2)} + g_1^{(1,1)}$ is not homogeneous since $g_2^{(2,2)} + g_1^{(1,1)}$ is no longer a homogeneous element. However, $g_2^{(2,2)} \leftarrow g_2^{(2,2)} + t^{(1,1)}g_1^{(1,1)}$ is a homogeneous change with $\text{gr}(g_2^{(2,2)} + t^{(1,1)}g_1^{(1,1)}) = \text{gr}(g_2^{(2,2)}) = (2, 2)$, which results in a new homogeneous basis, $\{g_1^{(1,1)}, g_2^{(2,2)} + t^{(1,1)}g_1^{(1,1)}\}$. Homogeneous basis changes always result in homogeneous bases.

Notation. Let $[f]$ be a homogeneous presentation matrix of $f : F^1 \rightarrow F^0$ with bases $F^0 = \langle g_1, \dots, g_n \rangle$ and $F^1 = \langle s_1, \dots, s_m \rangle$. We extend the notation of grading to every row r_i and every column c_j from the basis elements g_i and s_j they represent respectively, that is, $\text{gr}(r_i) := \text{gr}(g_i)$ and $\text{gr}(c_j) := \text{gr}(s_j)$. We define a partial order \leq_{gr} on rows $\{r_i\}$ by asserting $r_i \leq_{\text{gr}} r_j$ iff $\text{gr}(r_i) \leq \text{gr}(r_j)$. Similarly, we define a partial order on columns $\{c_j\}$.

For such a homogeneous presentation matrix $[f]$, we have the following observations:

1. $\text{gr}([f]_{i,j}) = \text{gr}(c_j) - \text{gr}(r_i)$
2. a nonzero entry $[f]_{i,j}$ can only be zeroed out by column operations from columns $c_k \leq_{\text{gr}} c_j$ or by row operations from rows $r_\ell \geq_{\text{gr}} r_i$.

Observation (2) indicates which subset of column and row operations is sufficient to zero out the entry $[f]_{i,j}$. We restate the diagonalization problem as follows:

Given an $n \times m$ homogeneous presentation matrix $\mathbf{A} = [f]$ consisting of entries in $\mathbb{k}[t_1, \dots, t_d]$ with grading on rows and columns, find a total diagonalization of \mathbf{A} under the following admissible row and column operations:

- multiply a row or column by nonzero $\alpha \in \mathbb{k}$; (For $\mathbb{k} = \mathbb{F}_2$, we can ignore these operations).
- for two rows r_i, r_j with $j \neq i$ and $r_j \leq_{\text{gr}} r_i$, set $r_j \leftarrow r_j + t^{\mathbf{u}} \cdot r_i$ where $\mathbf{u} = \text{gr}(r_i) - \text{gr}(r_j)$
- for two columns c_i, c_j with $j \neq i$ and $c_i \leq_{\text{gr}} c_j$, set $c_j \leftarrow c_j + t^{\mathbf{v}} \cdot c_i$ where $\mathbf{v} = \text{gr}(c_j) - \text{gr}(c_i)$

The above operations realize all possible homogeneous basis changes. That means, any homogeneous presentation matrix can be realized by a combination of the above operations.

In fact, the values of nonzero entries in the matrix are redundant under the homogeneous property $\text{gr}(\mathbf{A}_{i,j}) = \text{gr}(c_j) - \text{gr}(r_i)$ given by observation (1). So, we can further simplify the matrix by replacing all the nonzero entries with their \mathbb{k} -coefficients. For example, we can replace $2 \cdot t^{\mathbf{u}}$ with 2. What really matters are the partial orders defined by the grading of rows and columns. With our assumption of $\mathbb{k} = \mathbb{F}_2$, all nonzero entries are replaced with 1. Based on above observations, we further simplify the diagonalization problem to be the one as follows.

Given a \mathbb{k} -valued matrix \mathbf{A} with a partial order on rows and columns, find a total diagonalization $\mathbf{A}^* \sim \mathbf{A}$ with the following admissible operations:

- multiply a row or column by nonzero $\alpha \in \mathbb{k}$; (For $\mathbb{k} = \mathbb{F}_2$, we can ignore these operations).
- Add c_i to c_j only if $j \neq i$ and $\text{gr}(c_i) \leq \text{gr}(c_j)$; denoted as $c_i \rightarrow c_j$.
- Add r_k to r_l only if $l \neq k$ and $\text{gr}(r_\ell) \leq \text{gr}(r_k)$; denoted as $r_k \rightarrow r_l$.

The assumption of $\mathbb{k} = \mathbb{F}_2$ allows us to ignore the first set of multiplication operations on the binary matrix obtained after transformation. Also, with the assumption of distinct grading, the second two sets of admissible operations become:

- Add column c_i to column c_j , denoted as $c_i \rightarrow c_j$, only if $\text{gr}(c_i) < \text{gr}(c_j)$.
- Add row r_i to row r_j , denoted $r_i \rightarrow r_j$, only if $\text{gr}(r_i) > \text{gr}(r_j)$.

We denote the set of all admissible column and row operations as

$$\begin{aligned} \text{Colop} &= \{(i, j) \mid c_i \rightarrow c_j \text{ is an admissible column operation}\}, \\ \text{Rowop} &= \{(k, l) \mid r_k \rightarrow r_l \text{ is an admissible row operation}\}. \end{aligned}$$

¹Recall that an element $m \in M$ is homogeneous with grade $\text{gr}(m) = \mathbf{u}$ for some $\mathbf{u} \in \mathbb{Z}^d$ if $m \in M_{\mathbf{u}}$.

Under the assumption that no two columns nor rows have same grades, Colop and Rowop are closed under transitive relation.

Proposition 4.1. $(i, j), (j, k) \in \text{Colop (Rowop)} \implies (i, k) \in \text{Colop (Rowop)}$.

Given a solution of the diagonalization problem in the simplified form, one can reconstruct a solution of the original problem on the presentation matrix by reversing the above process of simplification. We will illustrate it by running our algorithm on the working example 1 at the end of this section. The matrix reduction we employ for diagonalization may be viewed as a *generalized matrix reduction* because the matrix is reduced under constrained operations Colop and Rowop which might be a nontrivial subset of all basic operations.

Remark 4.1. There are two extreme but trivial cases: (i) there are no \leq_{gr} -comparable pair of rows and columns. In this case, $\text{Colop} = \text{Rowop} = \emptyset$ and the original matrix is a trivial solution. (ii) All pairs of rows and all pairs of columns are \leq_{gr} -comparable. Or equivalently, both Colop and Rowop are totally ordered. In this case, one can apply traditional matrix reduction algorithm to reduce the matrix to a diagonal matrix with all nonzero blocks being 1×1 minors. This is also the case for traditional 1-parameter persistence module if one further applies row reduction after column reduction. Note that row reductions are not necessary for reading out persistence information because it essentially does not change the persistence information. However, in multiparameter cases, both column and row reductions are necessary to obtain a diagonalization from which the persistence information can be read. From this view-point, our algorithm can be thought of as a generalization of the traditional persistence algorithm.

Example 2. Consider our working example 1. One can see later in Section 5 (Case 1) that the presentation matrix of this example can be chosen to be the same as the matrix of the boundary morphism $\partial_1 : C_1 \rightarrow C_0$. With fixed bases $C_0 = \langle v_b^{(0,1)}, v_r^{(1,0)}, v_g^{(1,1)} \rangle$ and $C_1 = \langle e_r^{(1,1)}, e_b^{(1,2)}, e_g^{(2,1)} \rangle$, this presentation matrix $[\partial_1]$ and the corresponding binary matrix \mathbf{A} can be written as follows (recall that superscripts indicate the grades) :

$$\begin{array}{c} [\partial_1] \\ v_b^{(0,1)} \\ v_r^{(1,0)} \\ v_g^{(1,1)} \end{array} \begin{pmatrix} e_r^{(1,1)} & e_b^{(1,2)} & e_g^{(2,1)} \\ \mathbf{t}^{(1,0)} & \mathbf{t}^{(1,1)} & 0 \\ \mathbf{t}^{(0,1)} & 0 & \mathbf{t}^{(1,1)} \\ 0 & \mathbf{t}^{(0,1)} & \mathbf{t}^{(1,0)} \end{pmatrix} \longrightarrow \mathbf{A} \begin{array}{c} c_1^{(1,1)} \\ c_2^{(1,2)} \\ c_3^{(2,1)} \end{array} \begin{pmatrix} r_1^{(0,1)} \\ r_2^{(1,0)} \\ r_3^{(1,1)} \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix}$$

Four admissible operations are: $r_3 \rightarrow r_1, r_3 \rightarrow r_2, c_1 \rightarrow c_2, c_1 \rightarrow c_3$.

4.2 Total diagonalization algorithm

Recall that we assume distinct grading, i.e., no two columns nor two rows have same grades. We make some comments on the output of our algorithm without this assumption later in the conclusion.

Let \mathbf{A} be the presentation matrix whose total diagonalization we are looking for. We order the rows and columns of the matrix \mathbf{A} according to any topological order that extends the partial order on the grades to a total order, e.g., dictionary order. We fix the indices $\text{Row}(\mathbf{A}) = \{1, 2, \dots, \ell\}$ and $\text{Col}(\mathbf{A}) = \{1, 2, \dots, m\}$ according to this order. With this ordering, observe that, for each admissible column operation $c_i \rightarrow c_j$, we have $i < j$, and for each admissible row operation $r_l \rightarrow r_k$, we have $l > k$.

For any column c_t , let $\mathbf{A}_{\leq t} := \mathbf{A}|_C$ denote the left submatrix on $C = [\text{Row}(\mathbf{A}), \{j \in \text{Col}(\mathbf{A}) \mid j \leq t\}]$ and $\mathbf{A}_{< t}$ denote its stricter version obtained by excluding column c_t from $\mathbf{A}_{\leq t}$. Our algorithm starts with the finest decomposition and coarsens it as necessary. The main steps of our algorithm runs as follows (see Figure 5 for an illustration):

0. **Initialization:** Initialize the collection of index blocks $\mathcal{B}^{(0)} := \{B_i^{(0)} := [\{i\}, \emptyset] \mid i \in \text{Row}(\mathbf{A})\}$, for the total diagonalization of null column matrix $\mathbf{A}_{\leq 0}$.
1. **Main iteration:** Process \mathbf{A} from left to right incrementally by introducing a column c_t and considering left submatrices $\mathbf{A}_{\leq t}$ for $t = 1, 2, \dots, m$. We update and maintain the collection of index blocks $\mathcal{B}^{(t)} \leftarrow \{B_i^{(t)}\}$ for the current submatrix $\mathbf{A}_{\leq t}$ in each iteration by using column and block updates stated below. Here we use upper index $(\cdot)^{(t)}$ to emphasize the iteration t .
2. **Sub-column update:** Partition the column c_t into sub-columns $c_t|_{\text{Row}B_i^{(t-1)}} := \mathbf{A}|_{[\text{Row}(B_i^{(t-1)}), \{t\}]}$, one for the set of rows $\text{Row}(B_i^{(t-1)})$ for each block from the previous iteration. We process each such sub-column $c_t|_{\text{Row}B_i^{(t-1)}}$ one by one, checking whether there exists a sequence of admissible operations that are able to reduce the sub-column to zero while *preserving the prior*, according to the definition below.

can be visited and modified anywhere by every subroutines called. That means, every time we update values in \mathbf{A} by some operations, these operations are applied to the latest \mathbf{A} .

Algorithm 1: TOTDIAGONALIZE(\mathbf{A})

Input: \mathbf{A} = input matrix treated as a global variable whose columns and rows are totally ordered respecting some fixed partial order given by the grading.

Result: a total diagonalization \mathbf{A}^* with index blocks \mathcal{B}^*

```

1  $\mathcal{B}^{(0)} \leftarrow \{B_i^{(0)} := [\{i\}, \emptyset] \mid i \in \text{Row}(\mathbf{A})\};$ 
2 for  $t \leftarrow 1$  to  $m := |\text{Col}(\mathbf{A})|$  do
3    $B_0^{(t)} \leftarrow [\emptyset, \{t\}];$ 
4   for each  $B_i^{(t-1)} \in \mathcal{B}^{(t-1)}$  do
5      $T := [\text{Row}(B_i^{(t-1)}), \text{Col}(\mathbf{A}_{\leq t}) \setminus \text{Col}(B_i^{(t-1)})];$ 
6     if BLOCKREDUCE( $T$ ) == false then
7        $B_i^{(t)} \leftarrow B_i^{(t-1)} \oplus B_0^{(t)};$  // update  $B_i$  by appending  $t$ 
8     end
9     else
10       $B_i^{(t)} \leftarrow B_i^{(t-1)};$  //  $B_i$  remains unchanged
11    end
12  end
13   $\mathcal{B}^{(t)} \leftarrow \{B_i^{(t)}\}$  with all  $B_i^{(t)}$  containing  $t$  merged as one block.
  //  $\mathbf{A}$  and  $c_t$  are updated in BLOCKREDUCE whenever it returns true
14 end
15 return  $(\mathbf{A}, \mathcal{B}^{(m)});$ 

```

Remark 4.2. Our algorithm does not require the input presentation matrix to be minimal. As indicated in Remark 3.3, the trivial parts result in either identity blocks or single column blocks like $[\emptyset, \{j\}]$. A single column block corresponds to a zero morphism and is not merged with any other blocks. Therefore, c_j is a zero column. For a single row block $[\{i\}, \emptyset]$ which is not merged with any other blocks, r_i is a zero row vector. It represents a free indecomposable submodule in the total decomposition of the input persistence module.

We first prove the correctness of TOTDIAGONALIZE assuming that BLOCKREDUCE routine works as claimed, namely, it checks if a sub-column of the current column c_t can be zeroed out while preserving the prior, that is, without changing the left submatrix from the previous iteration and also the other sub-columns of c_t that have already been zeroed out.

Proposition 4.2. *At the end of each iteration t , $\mathbf{A}_{\leq t}$ is a total diagonalization.*

Proof. We prove it by induction on t . For the base case $t = 0$, it follows trivially by definition. Now assume $\mathbf{A}^{(t-1)}$ is the matrix we get at the end of iteration $(t-1)$ with $\mathbf{A}_{\leq t-1}^{(t-1)}$ totally diagonalized. That means, $\mathbf{A}_{\leq t-1}^{(t-1)} = \mathbf{A}_{\leq t-1}^*$ where $\mathbf{A} = \mathbf{A}^{(0)}$ is the original given matrix. For contradiction, assume at the end of iteration t , the matrix we get, $\mathbf{A}^{(t)}$, has left submatrix $\mathbf{A}_{\leq t}^{(t)}$ which is not totally diagonalized. That means, some index block $B \in \mathcal{B}^{(t)}$ can be decomposed further. Observe that such B must contain t because all other index blocks (not containing t) in $\mathcal{B}^{(t)}$ are also in $\mathcal{B}^{(t-1)}$ which cannot be decomposed further by our inductive assumption. We denote this index block containing t as B_t . Let \mathbf{A}' be the equivalent matrix of $\mathbf{A}^{(t)}$ such that $\mathbf{A}'_{\leq t}$ is a total diagonalization with index blocks \mathcal{B}' . Let F be an equivalent transformation from $\mathbf{A}^{(t)}$ to \mathbf{A}' , which decomposes B_t into at least two distinct index blocks of \mathcal{B}' , say B_0, B_1, \dots . Only one of them contains t , say B_0 . Then B_1 consists of only indices that are from $\mathbf{A}_{\leq t-1}$, which means B_1 equals some index block $B_i \in \mathcal{B}^{(t-1)}$. Therefore, the transformation F gives a sequence of admissible operations which can reduce the sub-column $c_t|_{\text{Row}(B_i)}$ to zero in $\mathbf{A}^{(t)}$. Note that we just use F to decompose the block of B_t . Therefore, we can choose a sequence of admissible operations which only involves indices of B_t . This gives us a sequence of admissible operations that does not change other sub-columns $c_t|_{\text{Row}(B_j)}$ for $B_j \neq B_t$. Starting with this sequence of admissible operations, we construct another sequence of admissible operations which further keeps $\mathbf{A}_{\leq t-1}^{(t)}$ unchanged to reach the contradiction. Note that $\mathbf{A}_{\leq t-1}^{(t)} = \mathbf{A}_{\leq t-1}^{(t-1)}$

Observe that all index blocks of \mathcal{B}' other than B_0 are also index blocks in $\mathcal{B}^{(t-1)}$, i.e. $\mathcal{B}' \setminus \{B_0\} \subseteq \mathcal{B}^{(t-1)}$. B_0 can be written as $B_0 = \bigoplus_{B_j \in \mathcal{B}^{(t-1)} \setminus \mathcal{B}'} B_j \oplus [\emptyset, \{t\}]$. Let B_a be the merge of index blocks that are in $\mathbf{A}^{(t-1)}$ and also in \mathbf{A}' and B_b be the merge of the rest of the index blocks of $\mathbf{A}^{(t-1)}$, i.e., $B_a = \bigoplus_{B_j \in \mathcal{B}' \cap \mathcal{B}^{(t-1)}} B_j$ and

$B_b = \bigoplus_{B_j \in \mathcal{B}^{(t-1)} \setminus \mathcal{B}'} B_j$. Then $\{B_a, B_b\}$ can be viewed as a coarser decomposition on $\mathbf{A}_{\leq t-1}^{(t-1)}$ and also on $\mathbf{A}'_{\leq t-1}$. By taking restrictions, we have $\mathbf{A}'|_{B_a} \sim \mathbf{A}^{(t-1)}|_{B_a}$ with equivalent transformation F_a and $\mathbf{A}'|_{B_b} \sim \mathbf{A}^{(t-1)}|_{B_b}$ with equivalent transformation F_b . Then F_a gives a sequence of admissible operations with indices in B_a and F_b gives a sequence of admissible operations with indices in B_b . By applying these operations on \mathbf{A}' , we can transform $\mathbf{A}'_{\leq t-1}$ to $\mathbf{A}_{\leq t-1}^{(t-1)}$ with sub-column $[\text{Row}(\mathbf{A}) \setminus \text{Row}(B_0), \{t\}]$ unchanged, which consists of the sub-columns that have already been reduced to zero. Combining all admissible operations from the three transformations F, F_a and F_b together, we get a sequence of admissible operations that reduce sub-column $[\text{Row}(B_i), \{t\}]$ to zero without changing $\mathbf{A}_{\leq t}^{(t)}$ and also those sub-columns which have already been reduced. But, then BLOCKREDUCE would have returned ‘true’ signaling that B_i should not be merged with any other block required to form the block B_t reaching a contradiction. \square

Now we design the BLOCKREDUCE subroutine as required. With the requirement of prior preservation, observe that reducing the sub-column $c_t|_{\text{Row}B}$ for some $B \in \mathcal{B}^{(t-1)}$ is the same as reducing $T = [\text{Row}(B), (\text{Col}(\mathbf{A}_{\leq t}) \setminus \text{Col}(B))]$ called the *target block* (see Figure 5 on the right). The main idea of BLOCKREDUCE is to consider a specific subset of admissible operations called *independent operations*. Within $\mathbf{A}_{\leq t}$, these operations only change entries in T and this change is independent of their order of application. Our BLOCKREDUCE is designed to search for a sequence of admissible operations within this subset and reduce T with it, if it exists. Clearly, the prior is preserved with these operations. The only thing we need to ensure is that searching within the set of independent operations is sufficient. That means, if there exists a sequence of admissible operations that can reduce T to 0 and meanwhile preserves the prior, then we can always find one such sequence with only independent operations. This is what we show next.

Consider the following matrices for each admissible operation. For each admissible column operation $c_i \rightarrow c_j$, let

$$\mathbf{Y}^{i,j} := \mathbf{A} \cdot [\delta_{i,j}]$$

where $[\delta_{i,j}]$ is the $m \times m$ square matrix with only one non-zero entry at (i, j) . Observe that $\mathbf{A} \cdot [\delta_{i,j}]$ is a matrix with the only nonzero column at j with entries copied from c_i in \mathbf{A} . Similarly, for each admissible row operation $r_l \rightarrow r_k$, let $[\delta_{k,l}]$ be the $\ell \times \ell$ matrix with only non-zero entry at (k, l) . let

$$\mathbf{X}^{k,l} := [\delta_{k,l}] \cdot \mathbf{A}$$

Application of a column operation $c_i \rightarrow c_j$ can be viewed as updating \mathbf{A} to $\mathbf{A} \cdot (\mathbf{I} + [\delta_{i,j}]) = \mathbf{A} + \mathbf{Y}^{i,j}$. Similar observation holds for row operations as well. For a target block $T = [\text{Row}(B), \text{Col}(\mathbf{A}_{\leq t}) \setminus \text{Col}(B)]$ defined on some $B \in \mathcal{B}^{(t-1)}$, we say an admissible column (row) operation, $c_i \rightarrow c_j$ ($r_l \rightarrow r_k$ resp.) is *independent on T* if $i \notin \text{Col}(T), j \in \text{Col}(T)$ ($l \notin \text{Row}(T), k \in \text{Row}(T)$ resp.). Briefly, we just call such operations *independent operations* if T is clear from the context.

We have two observations about independent operations that are important. The first one follows from the definition that $T = [\text{Row}(B), \text{Col}(\mathbf{A}_{\leq t}) \setminus \text{Col}(B)]$.

Observation 4.3. Within $\mathbf{A}_{\leq t}$, an independent column or row operation only changes entries on T .

Observation 4.4. For any independent column operation $c_i \rightarrow c_j$ and row operation $r_l \rightarrow r_k$, we have $[\delta_{k,l}] \cdot \mathbf{A} \cdot [\delta_{i,j}] = 0$. Or, equivalently

$$(\mathbf{I}_{\ell \times \ell} + [\delta_{k,l}]) \cdot \mathbf{A} \cdot (\mathbf{I}_{m \times m} + [\delta_{i,j}]) = \mathbf{A} + [\delta_{k,l}] \mathbf{A} + \mathbf{A} [\delta_{i,j}] = \mathbf{A} + \mathbf{X}^{k,l} + \mathbf{Y}^{i,j} \quad (2)$$

Proof. $[\delta_{k,l}] \cdot \mathbf{A} \cdot [\delta_{i,j}] = \mathbf{A}_{l,i} [\delta_{k,j}]$ (see Fig 6 for an illustration). By definitions of independence and T , we have $l \notin \text{Row}(B), i \in \text{Col}(B)$. That means they are row index and column index from different blocks. Therefore, $\mathbf{A}_{l,i} = 0$. \square

The following proposition reveals why we are after the independent operations.

Proposition 4.5. *The target block $\mathbf{A}|_T$ can be reduced to 0 while preserving the prior if and only if $\mathbf{A}|_T$ can be written as a linear combination of independent operations. That is,*

$$\mathbf{A}|_T = \sum_{\substack{l \notin \text{Row}(T) \\ k \in \text{Row}(T)}} \alpha_{k,l} \mathbf{X}^{k,l}|_T + \sum_{\substack{i \notin \text{Col}(T) \\ j \in \text{Col}(T)}} \beta_{i,j} \mathbf{Y}^{i,j}|_T \quad (3)$$

where $\alpha_{k,l}$'s and $\beta_{i,j}$'s are coefficient in $\mathbb{k} = \mathbb{F}_2$.

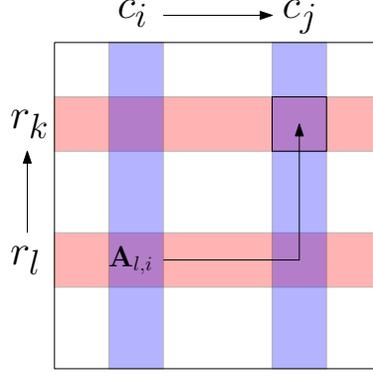


Figure 6: $[\delta_{k,l}]\mathbf{A}[\delta_{i,j}]$ is a matrix with the only nonzero entry at (k, j) being a copy of $\mathbf{A}_{l,i}$.

Proof. The full proof is in Appendix B, here we give some intuitive explanation. Reducing the target block $\mathbf{A}|_T$ to 0 is equivalent to finding matrices \mathbf{P} and \mathbf{Q} encoding sequences of admissible row operations and admissible column operations respectively so that $\mathbf{PAQ}|_T = 0$. For \Leftarrow direction, we can build $\mathbf{P} = \mathbf{I} + \sum \alpha_{k,l}[\delta_{k,l}]$ and $\mathbf{Q} = \mathbf{I} + \sum \beta_{i,j}[\delta_{i,j}]$ with binary coefficients $\alpha_{k,l}$'s and $\beta_{i,j}$'s given in Equation 3. Then using Observations 4.3 and 4.4, we show \mathbf{PAQ} indeed reduces $\mathbf{A}|_T$ to 0 with the prior being preserved.

For \Rightarrow direction, as long as we show that the existence of a transformation reducing $\mathbf{A}|_T$ to 0 implies the existence of a transformation reducing $\mathbf{A}|_T$ to 0 by independent operations, we are done. This is formally captured as Proposition B.1 and proved in Appendix B. \square

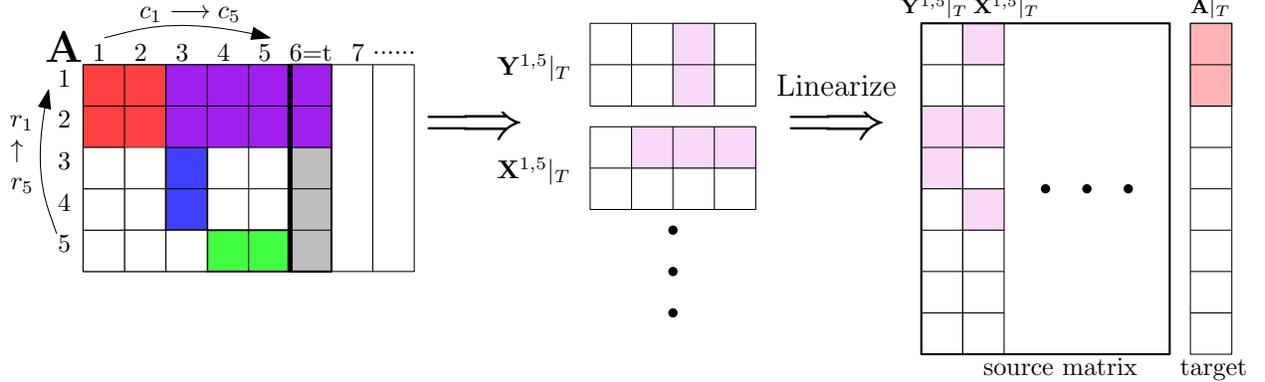


Figure 7: $t = 6$ is the current column. For the current block, the red one, our target block T is the purple one. We build $\mathbf{X}^{k,l}|_T$'s for admissible row operations from blue or green into purple. For example, $r_5 \rightarrow r_1$ illustrated on the left. Also, build $\mathbf{Y}^{i,j}|_T$'s for admissible column operations from the red block to the purple block. For example, $c_1 \rightarrow c_5$ on top-left. The middle picture shows $\mathbf{X}^{1,5}|_T$ and $\mathbf{Y}^{1,5}|_T$ for these two operations. After linearizing, the corresponding vectors are added into the source matrix, which is finally used to reduce the target $\mathbf{A}|_T$.

We can view $\mathbf{A}|_T$, $\mathbf{Y}^{i,j}|_T$, $\mathbf{X}^{k,l}|_T$ as binary vectors in the same $|T|$ -dimensional space. Proposition 4.5 tells us that it is sufficient to check if $\mathbf{A}|_T$ can be a linear combination of the vectors corresponding to a set of independent operations. So, we first *linearize* each of the matrices $\mathbf{Y}^{i,j}|_T$'s, $\mathbf{X}^{k,l}|_T$'s, and $\mathbf{A}|_T$ to a column vector as described later (see Figure 7). Then, we check if $\mathbf{A}|_T$ is in the span of $\mathbf{Y}^{i,j}|_T$'s and $\mathbf{X}^{k,l}|_T$'s. This is done by collecting all vectors $\mathbf{X}^{i,j}|_T$'s and $\mathbf{Y}^{k,l}|_T$'s into a matrix S called the *source matrix* (Figure 7(right)) and then reducing the vector $c := \mathbf{A}|_T$ with S by some standard matrix reduction algorithm with left-to-right column additions, which is the subroutine called COLREDUCE in BLOCKREDUCE described below. If $c = \mathbf{A}|_T$ can be reduced to 0, we apply the corresponding independent operations to update \mathbf{A} . Observe that all column operations used in reducing $\mathbf{A}|_T$ together only change the sub-column $c_t|_{\text{Row}B}$ while row operations may change \mathbf{A} to the right of the column t .

Here we provide a short description and the pseudo-code of the subroutine COLREDUCE.

For a column c_j , we use $\text{Low}(c_j)$ to indicate the lowest row number such that c_j has 1 in that row. Let $\text{Low}(c_j) = -1$ if c_j is a zero column. We call a matrix $\mathbf{S}' \sim \mathbf{S}$ lowest-conflict-free for \mathbf{S} if for each row index $i = \text{Low}(c_j) \neq -1$ there is no $j' \neq j$ so that $\text{Low}(c_{j'}) = i$. Notice that \mathbf{S}' is not necessarily unique. However, all the claims do not depend on the choice of \mathbf{S}' . The algorithm $\text{COLREDUCE}(\mathbf{S}, c)$ transforms the matrix $[\mathbf{S}|c]$ to a lowest-conflict-free matrix and as a result reduces the column c . We say this procedure *reduces c with \mathbf{S}* . Note that this algorithm is the traditional persistence algorithm.

Algorithm 2: $\text{COLREDUCE}(\mathbf{S}, c)$

Input: \mathbf{S} =source matrix, c =target column to reduce.

Result: return the reduced target column

```

1  $\mathbf{S}' \leftarrow [\mathbf{S}|c]$ ;
2 for  $i \leftarrow 1$  to  $|\text{Col}(\mathbf{S})|$  do // Transform  $[\mathbf{S}|c]$  to be lowest-conflict-free
3    $\ell \leftarrow \text{Low}(c_i)$ ;
4   if  $\ell \neq -1$  then
5     for  $j \leftarrow 1$  to  $i - 1$  do
6       if  $\text{Low}(c_j) == \ell$  then
7          $c_i \leftarrow c_j + c_i$ ;
8         go to 3
9       end
10    end
11  end
12 end
13 return  $c$ 

```

The following fact is well known and is the basis of the classical matrix based persistence algorithm.

Fact 1. There exists a set of column operations adding a column only to its right such that the matrix $[\mathbf{S}|c]$ is reduced to $[\mathbf{S}'|0]$ if and only if $\text{COLREDUCE}(\mathbf{S}, c)$ returns a zero vector.

Now we describe the linearization used in routine BLOCKREDUCE as presented in Algorithm 3: BLOCKREDUCE . We fix a linear order \leq_{Lin} on the set of matrix indices, $\text{Row}(\mathbf{A}) \times \text{Col}(\mathbf{A})$, as follows: $(i, j) \leq_{\text{Lin}} (i', j')$ if $j > j'$ or $j = j', i < i'$. Explicitly, we linearly order the indices as:

$$((1, m), (2, m), \dots, (\ell, n), (1, m - 1), (2, m - 1), \dots).$$

For any index block B , let $\text{Lin}(\mathbf{A}|_B)$ be the vector of dimension $|\text{Col}(B)| \cdot |\text{Row}(B)|$ obtained by linearizing $\mathbf{A}|_B$ to a vector in the above linear order on the indices.

Algorithm 3: $\text{BLOCKREDUCE}(T)$

Data: \mathbf{A} =global variable of the given matrix.

Input: T =index of target block to be reduced; t =index of current column

Result: Return a boolean to indicate whether $\mathbf{A}|_T$ can be reduced. Reduce block $\mathbf{A}|_T$ if possible.

```

1 Compute  $c := \text{Lin}(\mathbf{A}|_T)$  and initialize empty matrix  $\mathbf{S}$ ;
2 for each admissible column operation  $c_i \rightarrow c_j$  with  $i \notin \text{Col}(T), j \in \text{Col}(T)$ , do
3   compute  $\mathbf{Y}^{i,j}|_T := (\mathbf{A} \cdot [\delta_{i,j}])|_T$  and  $y^{i,j} = \text{Lin}(\mathbf{Y}^{i,j}|_T)$ ; update  $\mathbf{S} \leftarrow [\mathbf{S}|y^{i,j}]$ ;
4 end
5 for each admissible row operation  $r_l \rightarrow r_k$  with  $l \notin \text{Row}(T), k \in \text{Row}(T)$  do
6   compute  $\mathbf{X}^{k,l}|_T := ([\delta_{k,l}] \cdot \mathbf{A})|_T$  and  $x^{k,l} = \text{Lin}(\mathbf{X}^{k,l}|_T)$ ; update  $\mathbf{S} \leftarrow [\mathbf{S}|x^{k,l}]$ ;
7 end
8  $\text{COLREDUCE}(\mathbf{S}, c)$  returns indices of  $y^{i,j}$ 's and  $x^{k,l}$ 's used to reduce  $c$  if possible;
9 For every such index of  $y^{i,j}$  or  $x^{k,l}$  apply  $c_i \rightarrow c_j$  or  $r_l \rightarrow r_k$  to transform  $\mathbf{A}$ ;
10 return  $\mathbf{A}|_T == 0$ ;

```

Proposition 4.6. *The target block on T can be reduced to zero in \mathbf{A} while preserving the prior if and only if $\text{BLOCKREDUCE}(T)$ returns true.*

Time complexity. First we analyze the time complexity of TOTDIAGONALIZE assuming that the input matrix has size $\ell \times m$. Clearly, $\max\{\ell, m\} = O(N)$ where N is the total number of generators and relations. For each of $O(N)$ columns, we attempt to zero out every sub-column with row indices coinciding with each block B of the previously determined $O(N)$ blocks. Let B has N_B rows. Then, the block T in step 5 has N_B rows and $O(N)$ columns.

To zero-out a sub-column, we create a source matrix out of T which has size $O(NN_B) \times O(N^2)$ because each of $O(\binom{N}{2})$ possible operations is converted to a column of size $O(NN_B)$ in the source matrix. The source matrix \mathbf{S} with the target vector c can be reduced with an efficient algorithm [12, 35] in $O(a + N^2(NN_B)^{\omega-1})$ time where a is the total number of nonzero elements in $[\mathbf{S}|c]$ and $\omega \in [2, 2.373]$ is the exponent for matrix multiplication. We have $a = O(NN_B \cdot N^2) = O(N^3N_B)$. Therefore, for each block B we spend $O(N^3N_B + N^2(NN_B)^{\omega-1})$ time in step 6. Then, observing $\sum_{B \in \mathcal{B}} N_B = N$, for each column we spend a total time of

$$\sum_{B \in \mathcal{B}} O(N^3N_B + N^2(NN_B)^{\omega-1}) = O(N^4 + N^{\omega+1} \sum_{B \in \mathcal{B}} N_B^{\omega-1}) = O(N^4 + N^{2\omega}) = O(N^{2\omega}) \quad (4)$$

Therefore, counting for all of the $O(N)$ columns, the total time for decomposition takes $O(N^{2\omega+1})$ time.

4.3 Running TOTDIAGONALIZE on the working example 1

Example 3. Consider the binary matrix after simplification as illustrated in Example 2.

$$\mathbf{A} \begin{matrix} & c_1^{(1,1)} & c_2^{(1,2)} & c_3^{(2,1)} \\ r_1^{(0,1)} & \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix} \\ r_2^{(1,0)} \\ r_3^{(1,1)} \end{matrix}$$

It has 4 admissible operations: $r_3 \rightarrow r_1, r_3 \rightarrow r_2, c_1 \rightarrow c_2, c_1 \rightarrow c_3$.

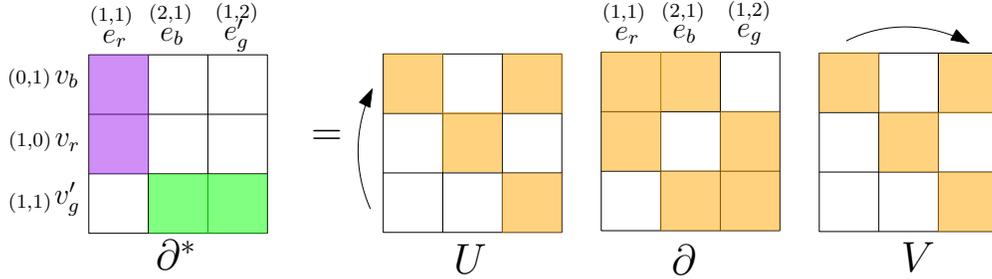


Figure 8: Diagonalizing the binary matrix given in Example 2: It is equivalent to multiplying the original matrix ∂ with a left matrix U that represents the row operation and a right matrix V that represents the column operations.

Before the first iteration, \mathcal{B} is initialized to be $\mathcal{B} = \{B_1 = (\{1\}, \emptyset), B_2 = (\{2\}, \emptyset), B_3 = (\{3\}, \emptyset)\}$. In the first iteration when $t = 1$, we have block $B_0 = (\emptyset, \{1\})$ for column c_1 . For $B_1 = (\{1\}, \emptyset)$, the target block we hope to zero out is $T = (\{1\}, \{1\})$. So we call **BLOCKREDUCE**(T) to check if $\mathbf{A}|_T$ can be zeroed out and update the entries on T according to the results of **BLOCKREDUCE**(T). There is only one admissible operation from outside of T into it, namely, $r_3 \rightarrow r_1$. The target vector $c = \text{Lin}(\mathbf{A}|_T)$ and the source matrix $\mathbf{S} = \{\text{Lin}([\delta_{1,3}]\mathbf{A})|_T\}$ are:

$$\mathbf{S} \quad \text{Lin}([\delta_{1,3}]\mathbf{A})|_T \quad \Bigg| \quad c = \text{Lin}(\mathbf{A}|_T) \\ \left[\begin{array}{c|c} 0 & 1 \end{array} \right]$$

The result of **COLREDUCE**(\mathbf{S}, c) stays the same as its input. That means we cannot reduce c at all. Therefore, **BLOCKREDUCE**(T, t) returns FALSE and nothing is updated in the original matrix.

It is not surprising that the matrix remains the same because the only admissible operation that can affect T does not change any entries in T at all. So there is nothing one can do to reduce it, which results in merging $B_1 \oplus B_0 = (\{1\}, \{1\})$. Similarly, for B_2 with $T = (\{2\}, \{1\})$, the only admissible operation $r_3 \rightarrow r_2$ does not change anything in T . Therefore, the matrix does not change and B_2 is merged with $B_1 \oplus B_0$, which results in the block $(\{1, 2\}, \{1\})$. For B_3 with $T = (\{3\}, \{1\})$, there is no admissible operation. So the matrix does not change. But $\mathbf{A}|_T = \mathbf{A}|_{(\{3\}, \{1\})} = 0$. That means **BLOCKREDUCE** returns TRUE. Therefore, we do not merge B_3 . In summary, B_0, B_1, B_2 are merged to be one block $(\{1, 2\}, \{1\})$ in the first iteration. So after the first iteration, there are two index blocks in $\mathcal{B}^{(1)}$: $(\{1, 2\}, \{1\})$ and $(\{3\}, \emptyset)$.

In the second iteration $t = 2$, we process the second column c_2 . Now $B_1 = (\{1, 2\}, \{1\}), B_2 = (\{3\}, \emptyset)$ and $B_0 = (\emptyset, \{2\})$. For the block $B_1 = (\{1, 2\}, \{1\})$, the target block we hope to zero out is $T = (\{1, 2\}, \{2\})$. There are

three admissible operations from outside of T into T , $r_3 \rightarrow r_1, r_3 \rightarrow r_2, c_1 \rightarrow c_2$. **BLOCKREDUCE**(T) constructs the target vector $c = \text{Lin}(\mathbf{A}|_T)$ and the source matrix $\mathbf{S} = \{\text{Lin}([\delta_{1,3}]\mathbf{A})|_T, \text{Lin}([\delta_{2,3}]\mathbf{A})|_T, \text{Lin}([\mathbf{A}\delta_{1,2}])|_T\}$ illustrated as follows:

$$\mathbf{S} \quad \begin{array}{c|ccc|c} \text{Lin}([\delta_{1,3}]\mathbf{A})|_T & \text{Lin}([\delta_{2,3}]\mathbf{A})|_T & \text{Lin}([\mathbf{A}\delta_{1,2}])|_T & c = \text{Lin}(\mathbf{A}|_T) \\ \hline \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} & & & \begin{bmatrix} 1 \\ 0 \end{bmatrix} \end{array}$$

The result of **COLREDUCE**(\mathbf{S}, c) is

$$\begin{array}{c|ccc|c} \mathbf{S} & & & c \\ \hline \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} & & & \begin{bmatrix} 0 \\ 0 \end{bmatrix} \end{array}$$

So the **BLOCKREDUCE** updates $\mathbf{A}|_T$ to get the following updated matrix:

$$\mathbf{A}' \quad \begin{array}{c|ccc} c_1^{(1,1)} & c_2^{(1,2)} & c_3^{(2,1)} \\ \hline r_1^{(0,1)} + r_3^{(1,1)} & \begin{pmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix} \\ r_2^{(1,0)} \\ r_3^{(1,1)} \end{array}$$

and return **TRUE** since $\mathbf{A}'|_T = 0$. Therefore, we do not merge B_1 . We continue to check for the block $B_2 = (\{3\}, \emptyset)$ and $T = (\{3\}, \{1, 2\})$, whether $\mathbf{A}'|_T$ can be reduced to zero. There is no admissible operation for this block at all. Therefore, the matrix stays the same and **BLOCKREDUCE** returns **FALSE**. We merge $B_2 \oplus B_0 = (\{3\}, \{2\})$.

Continuing the process for the last column c_3 in the third iteration $t = 3$, we see that $B_1 = (\{1, 2\}, \{1\}), B_2 = (\{3\}, \{2\})$ and $B_0 = (\emptyset, \{3\})$. For the block $B_1 = (\{1, 2\}, \{1\})$, the target block we hope to zero out is $T = (\{1, 2\}, \{2, 3\})$. There are four admissible operations from outside of T into T , $r_3 \rightarrow r_1, r_3 \rightarrow r_2, c_1 \rightarrow c_2, c_1 \rightarrow c_3$. **BLOCKREDUCE**(T) constructs the target vector $c = \text{Lin}(\mathbf{A}|_T)$ and the source matrix $\mathbf{S} = \{\text{Lin}([\delta_{1,3}]\mathbf{A})|_T, \text{Lin}([\delta_{2,3}]\mathbf{A})|_T, \text{Lin}([\mathbf{A}\delta_{1,2}])|_T, \text{Lin}([\mathbf{A}\delta_{1,3}])|_T\}$ illustrated as follows:

$$\mathbf{S} \quad \begin{array}{c|cccc|c} \text{Lin}([\delta_{1,3}]\mathbf{A})|_T & \text{Lin}([\delta_{2,3}]\mathbf{A})|_T & \text{Lin}([\mathbf{A}\delta_{1,2}])|_T & \text{Lin}([\mathbf{A}\delta_{1,3}])|_T & c = \text{Lin}(\mathbf{A}|_T) \\ \hline \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix} & & & & \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix} \end{array}$$

The result of **COLREDUCE**(\mathbf{S}, c) is

$$\begin{array}{c|cccc|c} \mathbf{S} & & & & c \\ \hline \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} & & & & \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \end{array}$$

So the **BLOCKREDUCE** updates $\mathbf{A}|_T$ to get the following updated matrix:

$$\mathbf{A}' \quad \begin{array}{c|ccc} c_1^{(1,1)} & c_2^{(1,2)} + c_1^{(1,1)} & c_3^{(2,1)} \\ \hline r_1^{(0,1)} & \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 1 \end{pmatrix} \\ r_2^{(1,0)} + r_3^{(1,1)} \\ r_3^{(1,1)} \end{array}$$

and returns **TRUE** since $\mathbf{A}'|_T = 0$. Therefore, we do not merge B_1 with any other block. We continue to check for the block $B_2 = (\{3\}, \{2\})$ and $T = (\{3\}, \{1, 3\})$, whether $\mathbf{A}'|_T$ can be reduced to zero. There is no admissible operation for this block at all. Therefore, the matrix stays the same and **BLOCKREDUCE** returns **FALSE**. We merge $B_2 \oplus B_0 = (\{3\}, \{2, 3\})$.

Finally the algorithm returns the matrix \mathbf{A}' shown above as the final result. It is the correct total diagonalization with two index blocks in $\mathcal{B}^{\mathbf{A}'}$: $B_1 = (\{1, 2\}, \{1\})$ and $B_2 = (\{3\}, \{2, 3\})$. An examination of **COLREDUCE**(\mathbf{S}, c) in all three iterations over columns reveals that the entire matrix \mathbf{A} is updated by operations $r_3 \rightarrow r_2$ and $c_1 \rightarrow c_2$.

We can further transform it back to the original form of the presentation matrix $[\partial_1]$. Observe that a row addition $r_i \leftarrow r_i + r_j$ reverts to a basis change in the opposite direction.

$$\begin{aligned}
& [\partial_1] \begin{pmatrix} e_r^{(1,1)} & e_b^{(1,2)} & e_g^{(2,1)} \\ v_b^{(0,1)} & \mathbf{t}^{(1,0)} & \mathbf{t}^{(1,1)} & 0 \\ v_r^{(1,0)} & \mathbf{t}^{(0,1)} & 0 & \mathbf{t}^{(1,1)} \\ v_g^{(1,1)} & 0 & \mathbf{t}^{(0,1)} & \mathbf{t}^{(1,0)} \end{pmatrix} \\
\Rightarrow & [\partial_1]^* \begin{pmatrix} e_r^{(1,1)} & e_b^{(1,2)} + \mathbf{t}^{(0,1)} e_r^{(1,1)} & e_g^{(2,1)} \\ v_b^{(0,1)} & \mathbf{t}^{(1,0)} & 0 & 0 \\ v_r^{(1,0)} & \mathbf{t}^{(0,1)} & 0 & 0 \\ v_g^{(1,1)} + \mathbf{t}^{(0,1)} v_r^{(1,0)} & 0 & \mathbf{t}^{(0,1)} & \mathbf{t}^{(1,0)} \end{pmatrix}
\end{aligned}$$

5 Computing presentations

Now that we know how to decompose a presentation by diagonalizing its matrix form, we describe how to construct and compute these matrices in this section. In practice, as described in Example 1, a persistence module is given implicitly with a simplicial filtration from which a graded module of simplicial chain complex can be inferred as we discussed before. We always assume that the simplicial filtration is 1-critical, which means that each simplex has a unique earliest birth time. For the case which is not 1-critical, called multi-critical, one may utilize the *mapping telescope*, a standard algebraic construction [31], which transforms a multi-critical filtration to a 1-critical one. However, notice that this transformation increases the input size depending on the multiplicity of the incomparable birth times of the simplices. For 1-critical filtrations, each module C_p is free. With a fixed basis for each free module C_p , a concrete matrix $[\partial_p]$ for each boundary morphism ∂_p based on the chosen bases can be constructed.

With this input, we discuss our strategies for different cases that depend on two parameters, d , the number of parameters of filtration function, and p , the dimension of the homology groups in the persistence modules.

We already stated that a simplicial filtration induces a persistence module. Here, we first give the details of this construction. Recall that a (d -parameter) *simplicial filtration* is a family of simplicial complexes $\{X_{\mathbf{u}}\}_{\mathbf{u} \in \mathbb{Z}^d}$ such that for each grade $\mathbf{u} \in \mathbb{Z}^d$ and each $i = 1, \dots, d$, $X_{\mathbf{u}} \subseteq X_{\mathbf{u} + \mathbf{e}_i}$. A *d -parameter persistence module* is a graded R -module where the vector spaces $M_{\mathbf{u}}$ are homology groups and linear maps among them are induced by a d -parameter simplicial filtration.

We obtain a simplicial chain complex $(C_{\bullet}(X_{\mathbf{u}}), \partial_{\bullet})$ for each $X_{\mathbf{u}}$ in this simplicial filtration. For each comparable pairs in the grading $\mathbf{u} \leq \mathbf{v} \in \mathbb{Z}^d$, a family of inclusion maps $C_{\bullet}(X_{\mathbf{u}}) \hookrightarrow C_{\bullet}(X_{\mathbf{v}})$ is induced by the canonical inclusion $X_{\mathbf{u}} \hookrightarrow X_{\mathbf{v}}$ giving rise to the following diagram:

$$\begin{array}{ccccccc}
C_{\bullet}(X_{\mathbf{u}}) : & \cdots & \xrightarrow{\partial_{p+2}} & C_{p+1}(X_{\mathbf{u}}) & \xrightarrow{\partial_{p+1}} & C_p(X_{\mathbf{u}}) & \xrightarrow{\partial_p} & C_{p-1}(X_{\mathbf{u}}) & \xrightarrow{\partial_{p-1}} & \cdots \\
\downarrow & & & \downarrow & & \downarrow & & \downarrow & & \\
C_{\bullet}(X_{\mathbf{v}}) : & \cdots & \xrightarrow{\partial_{p+2}} & C_{p+1}(X_{\mathbf{v}}) & \xrightarrow{\partial_{p+1}} & C_p(X_{\mathbf{v}}) & \xrightarrow{\partial_p} & C_{p-1}(X_{\mathbf{v}}) & \xrightarrow{\partial_{p-1}} & \cdots
\end{array}$$

For each chain complex $C_{\bullet}(X_{\mathbf{u}})$, we have the cycle spaces $Z_p(X_{\mathbf{u}})$'s and boundary spaces $B_p(X_{\mathbf{u}})$'s as kernels and images of boundary maps ∂_p 's respectively, and the homology group $H_p(X_{\mathbf{u}}) = Z_p(X_{\mathbf{u}})/B_p(X_{\mathbf{u}})$ as the cokernel of the inclusion maps $B_p(X_{\mathbf{u}}) \hookrightarrow Z_p(X_{\mathbf{u}})$. In line with category theory we use the notations im , ker , coker for indicating both the modules of kernel, image, cokernel and the corresponding morphisms uniquely determined by their constructions². We obtain the following commutative diagram:

$$\begin{array}{ccc}
& B_p(X_{\mathbf{u}}) \hookrightarrow Z_p(X_{\mathbf{u}}) \xrightarrow{\text{coker}} H_p(X_{\mathbf{u}}) & \\
& \text{im } \partial_{p+1} \nearrow & \searrow \text{ker } \partial_p \\
\cdots & C_{p+1}(X_{\mathbf{u}}) \xrightarrow{\partial_{p+1}} C_p(X_{\mathbf{u}}) \cdots &
\end{array}$$

²e.g. $\text{ker } \partial_p$ denotes the inclusion of Z_p into C_p

In the language of graded modules, for each p , the family of vector spaces and linear maps (inclusions) $(\{C_p(X_{\mathbf{u}}}\}_{\mathbf{u} \in \mathbb{Z}^d}, \{C_p(X_{\mathbf{u}}) \hookrightarrow C_p(X_{\mathbf{v}})\}_{\mathbf{u} \leq \mathbf{v}})$ can be summarized as a \mathbb{Z}^d -graded R -module:

$$C_p(X) := \bigoplus_{\mathbf{u} \in \mathbb{Z}^d} C_p(X_{\mathbf{u}}), \text{ with the ring action } t_i \cdot C_p(X_{\mathbf{u}}) : C_p(X_{\mathbf{u}}) \hookrightarrow C_p(X_{\mathbf{u}+e_i}) \quad \forall i, \forall \mathbf{u}.$$

That is, the ring R acts as the linear maps (inclusions) between pairs of vector spaces in $C_p(X_{\cdot})$ with comparable grades. It is not too hard to check that this $C_p(X_{\cdot})$ is indeed a graded module. Each p -chain in a chain space $C_p(X_{\mathbf{u}})$ is a homogeneous element with grade \mathbf{u} . Then we have a chain complex of graded modules $(C_*(X), \partial_*)$ where $\partial_* : C_*(X) \rightarrow C_{*-1}(X)$ is the boundary morphism given by $\partial_* \triangleq \bigoplus_{\mathbf{u} \in \mathbb{Z}^d} \partial_{*,\mathbf{u}}$ with $\partial_{*,\mathbf{u}} : C_*(X_{\mathbf{u}}) \rightarrow C_{*-1}(X_{\mathbf{u}})$ being the boundary map on $C_*(X_{\mathbf{u}})$.

The kernel and image of a graded module morphism are also graded modules as submodules of domain and codomain respectively whereas the cokernel is a quotient module of the codomain. They can also be defined grade-wise in the expected way:

$$\text{For } f : M \rightarrow N, (\ker f)_{\mathbf{u}} = \ker f_{\mathbf{u}}, (\text{im } f)_{\mathbf{u}} = \text{im } f_{\mathbf{u}}, (\text{coker } f)_{\mathbf{u}} = \text{coker } f_{\mathbf{u}}.$$

All the linear maps are naturally induced from the original linear maps in M and N . In our chain complex cases, the kernel and image of the boundary morphism $\partial_p : C_p(X) \rightarrow C_{p-1}(X)$ is the family of cycle spaces $Z_p(X)$ and family of boundary spaces $B_{p-1}(X)$ respectively with linear maps induced by inclusions. Also, from the inclusion induced morphism $B_p(X) \hookrightarrow Z_p(X)$, we have the cokernel module $H_p(X)$, consisting of homology groups $\bigoplus_{\mathbf{u} \in \mathbb{Z}^d} H_p(X_{\mathbf{u}})$ and linear maps induced from inclusion maps $X_{\mathbf{u}} \hookrightarrow X_{\mathbf{v}}$ for each comparable pairs $\mathbf{u} \leq \mathbf{v}$. This $H_p(X)$ is an example of *persistence module* M we mentioned in the beginning of this section, which we will study. It is called a persistence module M because not only does it encode the information of homology groups by each graded component $M_{\mathbf{u}}$, but, roughly speaking, also tracks birth, death, merging and persistence of the homological cycles through all admissible linear maps $M_{\mathbf{u}} \rightarrow M_{\mathbf{v}}, \forall \mathbf{u} \leq \mathbf{v}$. Classical persistence modules arising from a filtration of a simplicial complex over \mathbb{Z} is an example of a 1-parameter persistence module where the action $t_1 \cdot M_{\mathbf{u}} \subseteq M_{\mathbf{u}+e_1}$ signifies the linear map $M_{\mathbf{u}} \rightarrow M_{\mathbf{v}}$ between homology groups induced by the inclusion of the complex at \mathbf{u} into the complex at $\mathbf{v} = \mathbf{u} + e_1$.

In our case, we have chain complex of graded modules and induced homology groups which can be succinctly described by the following diagram:

$$\begin{array}{ccccccc} & & B_p(X) & \hookrightarrow & Z_p(X) & \twoheadrightarrow & H_p(X) & & B_{p-1}(X) & \hookrightarrow & Z_{p-1}(X) & \twoheadrightarrow & H_{p-1}(X) \\ & & \nearrow \text{im } \partial_{p+1} & & \searrow \ker(\partial_p) & & \nearrow \text{im } \partial_p & & \searrow \ker \partial_{p-1} & & & & & \\ \cdots & C_{p+1}(X) & \xrightarrow{\partial_{p+1}} & C_p(X) & \xrightarrow{\partial_p} & C_{p-1}(X) & \cdots & & & & & & & & \end{array}$$

Now we show how to compute presentations of persistence modules.

Note that a presentation gives an exact sequence $F^1 \rightarrow F^0 \twoheadrightarrow H \rightarrow 0$. To reveal further details of a presentation of H , we recognize that it respects the following commutative diagram,

$$\begin{array}{ccccc} & & Y^1 & & \\ & & \nearrow \text{im } f^1 & & \searrow \ker f^0 \\ F^1 & \xrightarrow{f^1} & F^0 & \xrightarrow{f^0 = \text{coker } f^1} & H \end{array}$$

where $Y^1 \hookrightarrow F^0$ is the kernel of f^0 . With this diagram being commutative, all maps in this diagram are essentially determined by the presentation map f^1 . We call the surjective map $f^0 : F^0 \twoheadrightarrow H$ *generating map*, and $Y^1 = \ker f^0$ the *1st syzygy module* of H .

We introduce the following useful properties of graded modules which are used in the justifications later. They are similar to Proposition (1.3) in Chapter 6 of [24].

Fact 2. Let M be a persistence module.

1. Choosing a homogeneous element in M with grade \mathbf{u} is equivalent to choosing a morphism $R_{\rightarrow \mathbf{u}} \rightarrow M$.
2. Choosing a set of homogeneous elements in M with grades $\mathbf{u}_1, \dots, \mathbf{u}_n$ is equivalent to choosing a morphism $\bigoplus_{i=1}^n R_{\rightarrow \mathbf{u}_i} \rightarrow M$.
3. Choosing a generating set of M consisting of n homogeneous elements with grades $\mathbf{u}_1, \dots, \mathbf{u}_n$ is equivalent to choosing a surjective morphism $\bigoplus_{i=1}^n R_{\rightarrow \mathbf{u}_i} \twoheadrightarrow M$.

4. If $M \simeq \bigoplus_i R_{\rightarrow \mathbf{u}_i}$ is a free module, choosing a basis of M is equivalent to choosing an isomorphism $\bigoplus R_{\rightarrow \mathbf{u}_i} \rightarrow M$.

5.1 Multiparameter filtration, zero-dimensional homology

In this case $p = 0$ and $d > 0$. This special case corresponds to determining clusters in the multiparameter setting. Importance of clusters obtained by classical one-parameter persistence has already been recognized in the literature [15, 41]. Our algorithm computes such clusters in a multiparameter setting. In this case, we obtain a presentation matrix straightforwardly with the observation that the module Z_0 of cycle spaces coincides with the module C_0 of chain spaces.

- Presentation: $C_1 \xrightarrow{\partial_1} C_0 \xrightarrow{\text{coker } \partial_1} H_0$
- Presentation matrix = $[\partial_1]$ is given as part of the input.

Justification. For $p = 0$, the cycle module $Z_0 = C_0$ is a free module. So we have the presentation of H as follows:

$$\begin{array}{ccccc} & & B_0 & & \\ & \nearrow \text{im } \partial_1 & \searrow & & \\ C_1 & \xrightarrow{\partial_1} & C_0 & \xrightarrow{\text{coker } \partial_1} & H_0 \end{array}$$

It is easy to check that $\partial_1 : C_1 \rightarrow C_0$ is a presentation of H_0 since both C_1 and C_0 are free modules. With standard basis of chain modules C_p 's, we have a presentation matrix $[\partial_1]$ as the valid input to our decomposition algorithm.

The 0^{th} homology in our working example 1 corresponds to this case. The presentation matrix is the same as the matrix of boundary morphism ∂_1 .

For convenience, we introduce a compact description of a presentation $f^1 : F^1 \rightarrow F^0$ of a module H . We write $H = \langle g_1, \dots, g_n : s_1, \dots, s_m \rangle$ where $\{g_i\}$ is a chosen basis of F^0 and $\{s_j\}$ is a chosen generating set of $\text{im } f^1 \subseteq F^0$ of F^0 . In the working example 1, we can write $H_0 = \langle v_b^{(0,1)}, v_r^{(1,0)}, v_g^{(1,1)} : \partial_1(e_r^{(1,1)}), \partial_1(e_b^{(1,2)}), \partial_1(e_g^{(2,1)}) \rangle$.

5.2 2-parameter filtration, multi-dimensional homology

In this case, $d = 2$ and $p \geq 0$. Lesnick and Wright [40] have presented an algorithm to compute a presentation, in fact a minimal presentation, for this case. We restate some of their observations for completeness here. When $d = 2$, by Hilbert Syzygy Theorem [32], the kernel of a morphism between two free graded modules is always free. This implies that the canonical surjective map $Z_p \rightarrow H_p$ from free module Z_p can be naturally chosen as a generating map in the presentation of H_p . In this case we have:

- Presentation: $C_{p+1} \xrightarrow{\bar{\partial}_{p+1}} Z_p \xrightarrow{\text{coker } \bar{\partial}_{p+1}} H_p$ where $\bar{\partial}_{p+1}$ is the induced map from the diagram:

$$\begin{array}{ccccc} & & B_p & \hookrightarrow & Z_p & \twoheadrightarrow & H_p \\ & \nearrow \text{im } \partial_{p+1} & \nearrow \bar{\partial}_{p+1} & & \searrow \ker \partial_p & & \\ C_{p+1} & \xrightarrow{\partial_{p+1}} & C_p & & & & \end{array}$$

- Presentation matrix = $[\bar{\partial}_{p+1}]$ is constructed as follows:
 1. Compute a basis $G(Z_p)$ for the free module Z_p where $G(Z_p)$ is presented as a set of generators in the basis of C_p . This can be done by an algorithm in [40]. Take $G(Z_p)$ as the row basis of the presentation matrix $[\bar{\partial}_{p+1}]$.
 2. Present $\text{im } \partial_{p+1}$ in the basis of $G(Z_p)$ to get the presentation matrix $[\bar{\partial}_{p+1}]$ of the induced map as follows. Originally, $\text{im } \partial_{p+1}$ is presented in the basis of C_p through the given matrix $[\partial_{p+1}]$. One needs to rewrite each column of $[\partial_{p+1}]$ in the basis $G(Z_p)$ computed in the previous step. This can be done as follows. Let $[G(Z_p)]$ denote the matrix presenting basis elements in $G(Z_p)$ in the basis of C_p . Let c be any column vector in $[\partial_{p+1}]$. We reduce c to zero vector by the matrix $[G(Z_p)]$ and note the columns that are added to c . These columns provide the necessary presentation of c in the basis $G(Z_p)$. This reduction can be done through the traditional persistent algorithm [28].

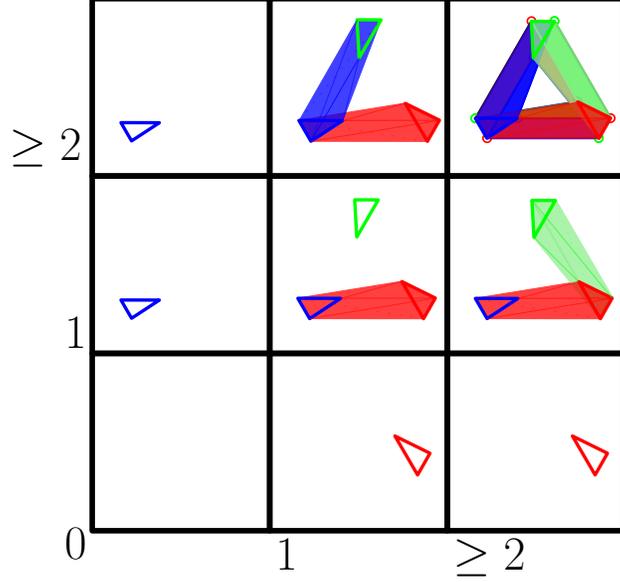


Figure 9: An example of 2-parameter simplicial filtrations. Each square box indicates what is the current (filtered) simplicial complex at the grade of the box. This example has one nontrivial cycle in 1st homology groups at grades except $(0, 0)$, $(1, 1)$, $(2, 2)$, and has two nontrivial cycles at grades $(1, 1)$ and $(2, 2)$. Note that all tunnels connecting triangles are hollow.

Justification. Unlike $p = 0$ case, for $p > 0$, we just know Z_p is a (proper) submodule of C_p , which means that Z_p is not necessarily equal to the free module C_p . However, fortunately for $d = 2$, the module Z_p is free, and we have an efficient algorithm to compute a basis of Z_p as the kernel of the boundary map $\partial_p : C_p \rightarrow C_{p-1}$. Then, we can construct the following presentation of H_p :

$$\begin{array}{ccccccc}
 & & & B_p & & & \\
 & & \nearrow & & \searrow & & \\
 & \text{im } \partial_{p+1} & & & & & \\
 \longrightarrow & C_{p+1} & \xrightarrow{\bar{\partial}_{p+1}} & Z_p & \xrightarrow{\text{coker } \bar{\partial}_{p+1}} & H_p & \longrightarrow 0
 \end{array}$$

Here the $\bar{\partial}_{p+1}$ is an induced map from ∂_{p+1} . With a fixed basis on Z_p and standard basis of C_{p+1} , we rewrite the presentation matrix $[\partial_{p+1}]$ to get $[\bar{\partial}_{p+1}]$, which constitutes a valid input to our decomposition algorithm.

Example 4. Consider the simplicial complex described in Figure 9. This is a hollow torus consisting of three empty triangles on three corners and each pair of triangles is connected by a hollow tunnel. This example is quite similar to the working example if we view the red, blue, green triangles as three generators in the H_1 persistence homology and three tunnels as relations connecting them. Then, we get an almost same presentation except that at grade $(2, 2)$, the triangular torus introduces a new cycle which is different from any previous generators. For fixed bases of Z_1 and B_1 , we can build the presentation matrix of $\bar{\partial}_2$. After doing some basic reduction, it can be shown that this presentation matrix is equivalent to:

$$[\bar{\partial}_2] \begin{pmatrix} s_r^{(1,1)} & s_b^{(1,2)} & s_g^{(2,1)} \\ g_b^{(0,1)} & \mathbf{t}^{(1,0)} & \mathbf{t}^{(1,1)} & 0 \\ g_r^{(1,0)} & \mathbf{t}^{(0,1)} & 0 & \mathbf{t}^{(1,1)} \\ g_g^{(1,1)} & 0 & \mathbf{t}^{(0,1)} & \mathbf{t}^{(1,0)} \\ g_\infty^{(2,2)} & 0 & 0 & 0 \end{pmatrix}$$

where $g_r^{(0,1)}$, $g_b^{(1,0)}$, $g_r^{(1,1)}$ represent the three triangles at the corners and $g_\infty^{(2,2)}$ represents the new cycle generated by the torus; images of $s_r^{(1,1)}$, $s_b^{(1,2)}$, $s_g^{(2,1)}$ under $\bar{\partial}_2$ represent the boundaries of three tunnels.

5.3 Multiparameter filtration, multi-dimensional homology

Now we consider the most general case where $p > 0$ and $d > 0$. The issue is that now Z_p is not free. So, it cannot be chosen as the 0th free module F^0 in the presentation of H_p . In what follows, we drop the index p from all modules for simplicity. We propose the following procedure to construct the presentation of H_p . Here we use lower indices for morphisms f_0 and f_1 between free modules in presentations instead of upper indices as in f^0 and f^1 in order to write the inverse f_i^{-1} of a map f_i more clearly.

- Presentation is constructed as follows:

1. Construct a minimal presentation of Z with 1st syzygy module Y^1 :

$$\begin{array}{ccccc} & & Y^1 & & \\ & \nearrow & \swarrow & & \\ F^1 & \xrightarrow{f_1} & F^0 & \xrightarrow{f_0} & Z \end{array}$$

2. With the short exact sequence $B \hookrightarrow Z \xrightarrow{\pi} H$, construct the presentation of H :

$$\begin{array}{ccccc} & & f_0^{-1}(B) & & \\ & \nearrow & \swarrow & & \\ F^1 \oplus C & \xrightarrow{\ker(\pi \circ f_0) \circ (\text{im } f_1 + \text{im } \partial)} & F^0 & \xrightarrow{\pi \circ f_0} & H \end{array}$$

where $\pi \circ f_0$ is the composition of surjective morphisms $F^0 \xrightarrow{f_0} Z \xrightarrow{\pi} H$; the inclusion map $f_0^{-1}(B) \hookrightarrow F^0$ is given by the kernel map $\ker(\pi \circ f_0)$; the surjective map $\text{im } f_1 + \text{im } \partial : F^1 \oplus C \rightarrow f_0^{-1}(B)$ is induced by the following diagram:

$$\begin{array}{ccccccc} 0 & \longrightarrow & F^1 & \hookrightarrow & F^1 \oplus C & \twoheadrightarrow & C \longrightarrow 0 \\ & & \text{im } f_1 \downarrow & & \exists \text{im } f_1 + \text{im } \partial \downarrow & & \downarrow \text{im } \partial \\ 0 & \longrightarrow & Y^1 & \hookrightarrow & f_0^{-1}(B) & \twoheadrightarrow & B \longrightarrow 0 \end{array}$$

where $\text{im } \partial : C \rightarrow B$ is the canonical surjective map induced from boundary map ∂ .

And finally, the presentation map $F^1 \oplus C \rightarrow F^0$ is just the composition $\ker(\pi \circ f_0) \circ (\text{im } f_1 + \text{im } \partial)$.

Presentation matrix = $[\ker(\pi \circ f_0) \circ (\text{im } f_1 + \text{im } \partial)]$ is constructed as follows:

1. Construct a presentation matrix $[\bar{\partial}]$ the same way as in the previous case.
2. Compute for Y^1 a minimal generating set $G(Y^1)$ in the basis of $G(Z)$. Let $[G(Y^1)]$ be the resulting matrix. Combine $[\bar{\partial}]$ with $[G(Y^1)]$ from right to get a larger matrix $[G(Y^1) \mid \bar{\partial}]$.

Justification. First, we take a presentation of Z ,

$$\begin{array}{ccccc} & & Y^1 & & \\ & \nearrow & \swarrow & & \\ F^1 & \xrightarrow{f_1} & F^0 & \xrightarrow{f_0} & Z \end{array}$$

Here Y^1 is the 1st syzygy module of Z . Combining it with the short exact sequence $B \hookrightarrow Z \rightarrow H$, we have,

$$\begin{array}{ccccc} f_0^{-1}(B) & \hookrightarrow & F^0 & & \\ \downarrow & & \downarrow f_0 & \searrow \bar{f}_0 = \pi \circ f_0 & \\ B & \hookrightarrow & Z & \xrightarrow{\pi} & H \end{array}$$

The map $\bar{f}_0 = \pi \circ f_0$ is a composition of surjections and thus is a surjection from a free module F^0 to H , which is a valid candidate for the 0th free module of a presentation of H . Observe that the 1st syzygy module of H , $\ker \bar{f}_0 = \ker(\pi \circ f_0) = f_0^{-1}(\ker \pi) = f_0^{-1}(B)$, and that $f_0^{-1}(B)$ can be constructed as the pullback of the maps from B, F^0 to Z . The left square commutative diagram preserves the inclusion and surjection in parallel.

Now the only thing left is to find a surjection from a free module to $f_0^{-1}(B)$. First, by the property of pullback, we know that $\ker f_0 = \ker(f_0^{-1}(B) \rightarrow B)$ in a commutative way. It implies that the following diagram commutes.

$$\begin{array}{ccc}
Y^1 & & \\
\downarrow \ker g & \searrow \ker f_0 & \\
f_0^{-1}(B) & \hookrightarrow & F^0 \\
\downarrow g & & \downarrow f_0 \\
B & \hookrightarrow & Z
\end{array}$$

Now focus on the left vertical line of the above commutative diagram. We have a short exact sequence $Y^1 \hookrightarrow f_0^{-1}(B) \twoheadrightarrow B$. By the horseshoe lemma (see lemma 2.2.8 in [46] for details), we can build the generating set of $f_0^{-1}(B)$ as illustrated in the following diagram:

$$\begin{array}{ccccccc}
0 & \longrightarrow & F^1 & \hookrightarrow & F^1 \oplus C & \twoheadrightarrow & C \longrightarrow 0 \\
& & \downarrow \text{im } f_1 & & \downarrow \exists \text{im } f_1 + \text{im } \partial & & \downarrow \text{im } \partial \\
0 & \longrightarrow & Y^1 & \hookrightarrow & f_0^{-1}(B) & \twoheadrightarrow & B \longrightarrow 0
\end{array}$$

The left projection $F^1 \twoheadrightarrow Y^1$ comes from the previous presentation of Z . The $C \twoheadrightarrow B$ is the image map induced from boundary map $\partial : C_{p+1} \rightarrow C_p$. We take the direct sum of $F^1 \oplus C$ and the horseshoe lemma indicates that there exists a projection $F^1 \oplus C \twoheadrightarrow f_0^{-1}(B)$ making the whole diagram commute. So finally, we have the valid presentation of $F^1 \oplus C \rightarrow F^0 \twoheadrightarrow H$.

Now we identify a generating set of $f_0^{-1}(B)$ that helps us constructing a matrix for the presentation of H . From the surjection $F^1 \oplus C \rightarrow f_0^{-1}(B)$ in the above commutative diagram, one can see that the combination of generators from $B = \text{im } \partial$ and $Y^1 = \text{im } f_1$ forms a generating set of $f_0^{-1}(B)$. The generators from $B = \text{im } \partial$ can be computed as in the previous case, which results in the matrix $[\bar{\partial}]$. The generators $G(Y^1)$ from $Y^1 = \text{im } f_1$ are obtained as a result of computing the presentation of Z , which can be done by an algorithm presented in Skryzalin's thesis [45]. Combining these two together, we get the presentation matrix $[\bar{\partial}, G(Y^1)]$ of H as desired. So, now we have the solutions for all general cases.

The above construction of presentation matrix can be understood as follows. The issue caused by non-free Z is that, if we use the same presentation matrix as we did in the previous case with free Z , we may lose some relations coming from the inner relations of a generating set of Z . We fix this problem by adding these inner relations into the presentation matrix.

Figure 10 shows a simple example of a filtration of simplicial complex whose persistence module H for $p = 1$ is a quotient module of non-free module Z . The module H is generated by three 1-cycles presented as $g_1^{(0,1,1)}, g_2^{(1,0,1)}, g_3^{(1,1,0)}$. But when they appear together in $(1, 1, 1)$, there is a relation between these three: $\mathbf{t}^{(1,0,0)}g_1^{(0,1,1)} + \mathbf{t}^{(0,1,0)}g_2^{(1,0,1)} + \mathbf{t}^{(0,0,1)}g_3^{(1,1,0)} = 0$. Although $\text{im } \partial_1 = 0$, we still have a nontrivial relation from Z . So, we have $H = \langle g_1^{(0,1,1)}, g_2^{(1,0,1)}, g_3^{(1,1,0)} : s^{(1,1,1)} = \mathbf{t}^{(1,0,0)}g_1^{(0,1,1)} + \mathbf{t}^{(0,1,0)}g_2^{(1,0,1)} + \mathbf{t}^{(0,0,1)}g_3^{(1,1,0)} \rangle$. The presentation matrix turns out to be the following:

$$\begin{array}{c}
s^{(1,1,1)} \\
g_1^{(0,1,1)} \\
g_2^{(1,0,1)} \\
g_3^{(1,1,0)}
\end{array}
\begin{pmatrix}
\mathbf{t}^{(1,0,0)} \\
\mathbf{t}^{(0,1,0)} \\
\mathbf{t}^{(0,0,1)}
\end{pmatrix}$$

5.4 Time complexity

Now we consider the time complexity for computing presentation and decomposition together. Let n be the size of the input filtration, that is, total number of simplices obtained by counting new simplices added to the filtration (at most one new simplex at a grid point of \mathbb{Z}^d). We consider three different cases as before:

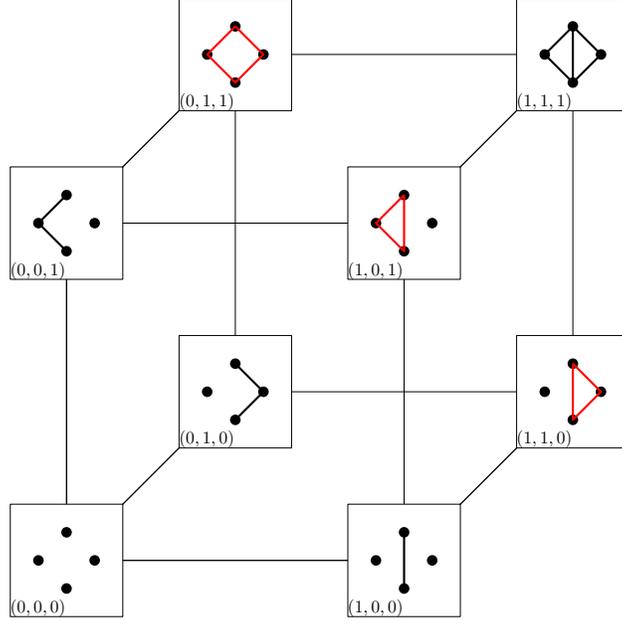


Figure 10: An example of a filtration of simplicial complex for $d = 3$ with non-free Z when $p = 1$. The three red circles are three generators in Z_1 . However, at grading $(1, 1, 1)$, the earliest time these three red circles exist simultaneously, there is a relation among these three generators.

Multiparameter, 0-dimensional homology: In this case, the presentation matrix $[\partial_1]$ where $\partial_1 : C_1 \rightarrow C_0$ has size $O(n) \times O(n)$. So, we can take $N = O(n)$ in Eqn. (4) for the time complexity analysis of decomposition. Therefore, the total time complexity for this case is $O(n^{2\omega+1})$.

2-parameter, multi-dimensional homology: In this case, as described in section 5.2, first we compute a basis $G(Z_p)$ that is presented in the basis of C_p . This is done by the algorithm of Lesnick and Wright [40] (henceforth called LW-algorithm) which runs in $O(n^3)$ time. Using $[G(Z_p)]$, we compute the presentation matrix $[\partial_{p+1}]$ as described in section 5.2. This can be done in $O(n^3)$ time assuming that $G(Z_p)$ has at most n elements. The presentation matrix is decomposed with TOTDIAGONALIZE as in the previous case. However, to claim that it runs in $O(n^{2\omega+1})$ time, one needs to ensure that the basis $G(Z_p)$ has $O(n)$ elements. This follows from the fact that Z_p is a free submodule of the free module C_p which has rank n .

Multiparameter, multi-dimensional homology: In this case, we need to compute a generating set $G(Z)$ for $Z = Z_p$ and then a generating set $G(Y^1)$ for the 1st Syzygy module Y^1 . Both of these generating sets can be computed by the algorithm of Skryzalin [45](Theorem 2.6.4). The algorithm considers $O(n^{d-2})$ slices of 2-parameter modules and generates the basis for each of them in $O(n^{d+1})$ total time. This also implies that the size of the basis the algorithm computes is at most $O(n \cdot n^{d-2}) = O(n^{d-1})$.

Next, we compute a generating set $G(Y^1)$ for the syzygy module Y^1 . Recall that $Y^1 \xrightarrow{\ker f_0} Z_p$ where $F_0 \xrightarrow{f_0} Z_p$. Taking the generating set computed in the previous step as a basis for F_0 and observing that $\ker f_0 = \ker \bar{f}_0$ where $F_0 \xrightarrow{f_0} Z_p \xrightarrow{i} C_p$ and $\bar{f}_0 = i \circ f_0$ we can compute a generating set $G(Y^1)$ in terms of a basis of $G(Z_p)$ using the algorithm of Skryzalin again. Again, each of the $O(n^{d-2})$ 2-parameter slices will generate at most $O(n)$ basis elements giving a total of $O(n^{d-1})$ basis elements.

The matrix $[G(Y^1)]$ appended with the matrix $[\bar{\delta}_p]$, becomes the presentation matrix for H_p of size $O(n^{d-1}) \times O(n^{d-1})$. The decomposition algorithm on such a matrix takes at most $O(n^{(d-1)(2\omega+1)})$.

In summary, we have the following time complexity:

- d -parameter 0-dimensional case: $O(n^{2\omega+1})$.
- d -parameter multi-dimensional case(general case): $O(n^{d+1}) + O(n^{(d-1)(2\omega+1)}) = O(n^{(d-1)(2\omega+1)})$.

6 Persistent graded Betti numbers and blockcodes

For 1-parameter persistence modules, the traditional persistence algorithm computes a complete invariant called the persistence diagram [28] which also has an alternative representation called barcodes [17]. As a generalization of the traditional persistence algorithm, it is expected that the result of our algorithm should also lead to similar invariants. We propose two interpretations of our result as two different invariants, *persistent graded Betti numbers* as a generalization of persistence diagrams and *blockcodes* as a generalization of barcodes.

Both of them depend on the ideas of free resolution and graded Betti numbers which are well studied in commutative algebra and are first introduced in TDA by Knudson [37]. A brief introduction to free resolutions and their construction are given in Appendix A. Here, we focus more on the two invariants mentioned above. In a nutshell, a free resolution is an extension of free presentation. Consider a free presentation of M as depicted below.

$$\begin{array}{ccccccc}
 & & & Y^1 & & & \\
 & & & \nearrow & \searrow & & \\
 & & \text{im } f^1 & & \text{ker } f^0 & & \\
 F^1 & \xrightarrow{f^1} & & & & \xrightarrow{f^0 = \text{coker } f^1} & M
 \end{array}$$

If the presentation map f^1 has nontrivial kernel, we can find a nontrivial map $f^2 : F^2 \rightarrow F^1$ with $\text{im } f^2 = \text{ker } f^1$, which implies $\text{coker } f^2 \cong \text{im } f^1 = \text{ker } f^0 = Y^1$. Therefore, f^2 is in fact a presentation map of the first syzygy module Y^1 of M . We can keep doing this to get f^3, f^4, \dots by constructing presentation maps on higher order syzygy modules Y^2, Y^3, \dots of M , which results in a diagram depicted below, which is called a free resolution of M .

$$\begin{array}{ccccccccccc}
 & & & Y^3 & & & Y^2 & & & Y^1 & & \\
 & & & \nearrow & \searrow & & \nearrow & \searrow & & \nearrow & \searrow & \\
 & & \text{im } f^3 & & \text{ker } f^2 & & \text{im } f^2 & & \text{ker } f^1 & & \text{im } f^1 & & \text{ker } f^0 & & \\
 \dots & \longrightarrow & F^3 & \xrightarrow{f^3} & F^2 & \xrightarrow{f^2} & F^1 & \xrightarrow{f^1} & F^0 & \xrightarrow{f^0 = \text{coker } f^1} & M
 \end{array}$$

Free resolution is not unique. However, there exists an essentially unique minimal free resolution in the sense that any free resolution can be obtained by summing the minimal free resolution with a free resolution of a trivial module. For a graded module M , consider the multiset consisting of the grades of homogeneous basis elements for each F^j in the minimal free resolution of M . We record the multiplicity of each grade $\mathbf{u} \in \mathbb{Z}^d$ in this multiset, denoted as $\beta_{j,\mathbf{u}}^M$. Then, the mapping $\beta_{(-,-)}^M : \mathbb{Z}_{\geq 0} \times \mathbb{Z}^d \rightarrow \mathbb{Z}_{\geq 0}$ can be viewed as an invariant of graded module M , which is called the graded Betti numbers of M . By applying the decomposition of module $M \simeq \bigoplus M^i$, we have for each indecomposable M^i , the refined graded Betti numbers $\beta^{M^i} = \{\beta_{j,\mathbf{u}}^{M^i} \mid j \in \mathbb{N}, \mathbf{u} \in \mathbb{Z}^d\}$. We call the set $\mathcal{PB}(M) := \{\beta^{M^i}\}$ *persistent graded Betti numbers* of M . For the working example 1, the persistent graded Betti numbers are given in two tables listed in Table 1.

One way to summarize the information of graded Betti numbers is to use the Hilbert function, which is also called dimension function [26] in TDA defined as:

$$\text{dm}M : \mathbb{Z}^d \rightarrow \mathbb{Z}_{\geq 0} \quad \text{dm}M(\mathbf{u}) = \dim(M_{\mathbf{u}})$$

Fact 3. There is a relation between the graded Betti numbers and dimension function of a persistence module as follows:

$$\forall \mathbf{u} \in \mathbb{Z}^d, \quad \text{dm}M(\mathbf{u}) = \sum_{\mathbf{v} \leq \mathbf{u}} \sum_j (-1)^j \beta_{j,\mathbf{v}}$$

Then for each indecomposable M^i , we have the dimension function $\text{dm}M^i$. We call the set of dimension functions $\mathcal{B}_{\text{dm}}(M) := \{\text{dm}M^i\}$ the *blockcode* of M .

For our working example, the dimension functions of indecomposable summands M^1 and M^2 are:

$$\text{dm}M^1(\mathbf{u}) = \begin{cases} 1 & \text{if } \mathbf{u} \geq (1, 0) \text{ or } \mathbf{u} \geq (0, 1) \\ 0 & \text{otherwise} \end{cases} \quad \text{dm}M^2(\mathbf{u}) = \begin{cases} 1 & \text{if } \mathbf{u} = (1, 1) \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

They can be visualized as in Figure 11.

The information which can be read out from graded Betti numbers and dimension functions are similar. We take the dimension functions of our working example as an example. For $\text{dm}M^1$, two connected components are born at the two left-bottom corners of the purple region. They are merged together immediately when they meet at grade $(1, 1)$.

β^{M^1}	(1,0)	(0,1)	(1,1)	(2,1)	(1,2)	(2,2)	...
β_0	1	1					
β_1			1				
$\beta_{\geq 2}$							
β^{M^2}	(1,0)	(0,1)	(1,1)	(2,1)	(1,2)	(2,2)	...
β_0			1				
β_1				1	1		
β_2						1	
$\beta_{\geq 3}$							

Table 1: Persistence grades $\mathcal{PB}(M) = \{\beta^{M^1}, \beta^{M^2}\}$. All nonzero entries are listed in this table. Blank boxes indicate 0 entries.

After that, they persist forever as one connected component. For $\text{dm}M^2$, one connected component born at the left-bottom corner of the square green region. Later at the grades of left-top corner and right-bottom corner of the green region, it is merged with some other connected component with smaller grades of birth. Therefore, it only persists within this green region.

Remark. In general, both persistent graded Betti numbers and blockcodes are not sufficient to classify multiparameter persistence modules, which means they are not complete invariants. As indicated in [14], there is no complete discrete invariant for multiparameter persistence modules. However, interestingly, these two invariants are indeed complete invariants for interval decomposable modules like this example, which are recently studied in [5, 8, 26].

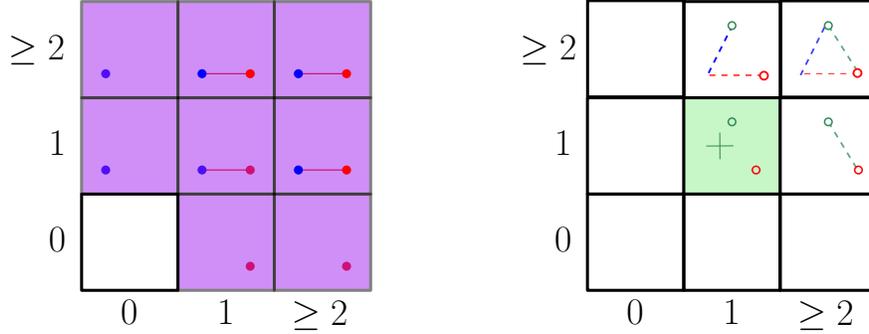


Figure 11: $\text{dm}M^1$ and $\text{dm}M^2$. Each colored square represents an 1-dimensional vector space \mathbb{k} and each white square represents a 0-dimensional vector space. In the left picture, M^1 is generated by $v_b^{0,1}, v_r^{1,0}$ which are drawn as a blue dot and a red dot respectively. They are merged at $(1, 1)$ by the red edge e_r . In the right picture, M^2 is generated by $v_g^{(1,1)} + t^{(0,1)}v_r^{1,0}$ which is represented by the combination of the green circle and the red circle together at $(1, 1)$. After this point $(1, 1)$, the generator is mod out to be zero by relation of e_g starting at $(2, 1)$, represented by the green dashed line segment, and by relation of $e_b + t^{(0,1)}e_r$ starting at $(1, 2)$, represented by the blue dashed line segment connected with the red dashed line segment.

6.1 Analogy with 1-parameter persistence modules

In this section, we draw an analogy between the well known invariants, persistent diagrams and barcodes, in 1-parameter persistence modules and the invariants which we called the persistent graded Betti numbers and blockcodes respectively.

We first give an illustration of the decomposition of an 1-parameter persistence module with a simple example.

Consider the 0^{th} persistence module induced by the 1-parameter simplicial filtration shown in Figure 12. The 0^{th} homology group encodes the connected components. From the simplicial filtration, first we can see that the number of connected components from grades 1 to 5 are $(1, 2, 2, 1, 1)$. This corresponds to the dimensions of homology vector space at each grade, which is also called the dimension function of the persistence module. Three vertices in blue, red, and green constitute three generators denoted as g_1, g_2, g_3 for homology groups introduced at grades 1, 2, and 3 respectively. In the filtration, g_2 is merged with g_1 at grade 3, and g_3 is merged with g_2 (hence also g_1) at grade 4.

The persistence algorithm computes the decomposition of this persistence module which results in a decomposition consisting of three indecomposable components. Each one of them corresponds to one generator. The persistence diagram summarizes the result as three pairings of grades: $(2, 3)$, $(3, 4)$, and $(1, \infty)$. The explanations are:

$(2,3)$: The generator g_2 born at grade 2 is merged with a generator born earlier than it at grade 3.

$(3,4)$: The generator g_3 born at grade 3 is merged with a generator born earlier than it at grade 4.

$(1, \infty)$: The generator g_1 born at grade 1 is never merged with some other generator.

The barcode represents the graph of dimension functions of each component in the decomposition. From the barcode of the example illustrated in figure 12, we can track directly when each generator gets born, merges (dies), and persists during its lifetime.

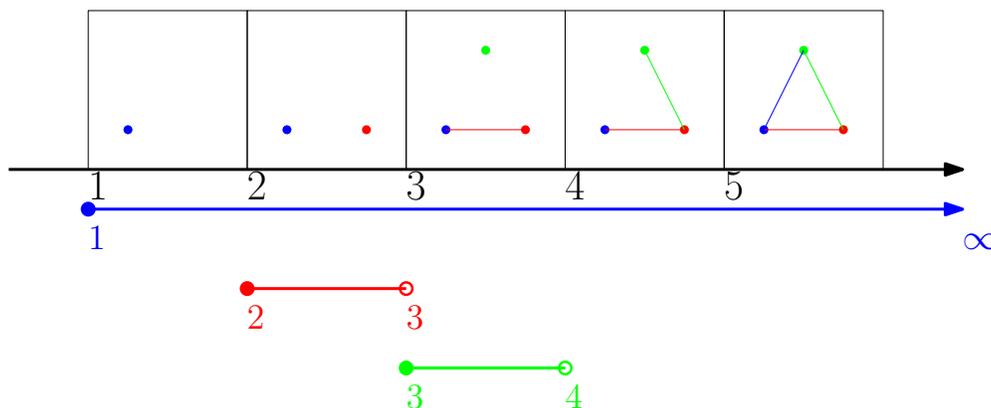


Figure 12: An example of 1-parameter simplicial filtration and its barcode.

For multiparameter persistence, we aim to compute a summary which encodes similar information as in the 1-parameter case. Consider the simplicial bi-filtration for the working example 1. Similar to our example in the 1-parameter case, we have three generators $g_{0,1}, g_{1,0}, g_{1,1}$ which are born at grades $(0, 1), (1, 0), (1, 1)$ respectively.

As shown in 4.3, the decomposition consists of two indecomposable component. One corresponds to $g_{0,1}, g_{1,0}$, and the other corresponds to $g_{1,1}$. Roughly speaking, the reason we cannot decompose $g_{0,1}$ and $g_{1,0}$ further is that their birth time are incomparable based on standard partial order of grades in \mathbb{Z}^2 . When they merge together at grade $(1, 1)$, neither one of them could be claimed to be merged with the other one having an earlier grade. So we have to keep them together in the same indecomposable component. However, for $g_{1,1}$, when it is merged with $g_{0,1}$ and $g_{1,0}$ at grades $(2, 1)$ and $(1, 2)$ respectively, both $g_{0,1}$ and $g_{1,0}$ have earlier grades than $g_{1,1}$.

Note that this explanation of decomposability based on the comparability of grades of generators does not always work as in 1-parameter case. That is why the decomposition in multiparameter case is much more complicated. But it is interesting to ask when this rule works in multiparameter case.

If we check the blockcode illustrated in Figure 11, we can see that for the first component, two generators are born at the two left-bottom corner of the purple region, which are grades $(0, 1), (1, 0)$. They are merged immediately at grade $(1, 1)$. After that, none of them is merged with anything else. Therefore, the merged generator persists forever. For the second component, one generator gets born at the left-bottom corner of the green region, which is grade $(1, 1)$. It persists in the green region. It is stopped by something else at grades $(1, 2)$ and $(2, 1)$. Therefore, it cannot persist beyond the green region.

7 Concluding remarks

In this paper, we propose an algorithm that generalizes the traditional persistence algorithm to the general case of multiparameter persistence. Even if its utility was clear, its design was illusive. The results of this algorithm are interpreted as invariants we call persistent graded Betti numbers and blockcode, which can be viewed as generalizations of the persistent diagram and the barcode computed with the traditional persistence algorithm. Specifically, our algorithm can be applied to determine whether a persistence module is interval decomposable or block decomposable,

which plays important roles in the computation of bottleneck distances and interleaving distances in multiparameter cases [1, 5, 8, 26].

The assumption that no two columns nor rows have same grades is necessary for our current algorithm. If we consider the persistence modules induced from filtration functions in the space of all real valued functions, our assumption of distinct grading is a generic property meaning that almost all induced persistence modules satisfy the assumption. However, it is still possible that in practice the induced persistence module does not satisfy the assumption. Without this assumption, our algorithm computes a (not necessarily total) decomposition which represents a total decomposition of some persistence module M' which can be viewed as a perturbed version of the original persistence module M by an arbitrarily small amount (considering $\mathbb{Z}^d \subseteq \mathbb{R}^d$). That means, the interleaving distance between this M' and M is arbitrarily small. The computed decomposition $M' = \oplus M'_j$ depends on the order with which we break the ties. How useful is this proposed strategy in practice? This question essentially relates to the question of "stability" of decomposition structures for which there is no satisfactory answer yet in the literature. Currently there is no universal definition about the stability of the decomposition structure of persistence modules. A simple way to address it is to find a matching with minimal bottleneck distance between two decomposition structures of two persistence modules with some cost function chosen for each paired indecomposable components. The most common cost function used so far is the interleaving distance. The stability-like property under this setting is an active area of recent research. There are some results of stability-like property on some special cases of multiparameter persistence modules, such as rectangle decomposable, triangle decomposable, block decomposable module and some other special interval decomposable modules [5]. For general multiparameter persistence modules, we know that the bottleneck distance, defined as we described above, can be much larger than interleaving distance [8]. One possible solution is that, based on the stability-like property of rectangle decomposable modules and special interval decomposable modules, one can approximate the original persistence modules with rectangle or interval decomposable modules with a similar decomposition structure which may provide some stability like property [27].

We believe the two invariants that we discussed are interesting summaries containing rich information about the multiparameter persistence modules. It motivates some interesting questions for future work. What kind of new meaningful pseudo-metrics on the space of persistence modules can be constructed and computed based on these invariants, and what are the relations between the new pseudo-metrics and the existing pseudo-metrics like interleaving distance, bottleneck distance, multi-matching distance, and so on? How stable will these pseudo-metrics be?

The time complexity of our algorithm is more than $O(n^4)$ in the 2-parameter case. An interesting question is if one can apply approximation techniques such as those in [25, 44] to design an approximation algorithm with time complexity $o(n^4)$. We also believe that most of the techniques for speeding up computation in the traditional persistence algorithm, like those in [6, 7], can be applied to our algorithm.

Acknowledgments.

This research is supported partially by the NSF grants CCF-1740761, CCF-2049010 and DMS-1547357. We acknowledge the influence of the BIRS Oaxaca workshop on Multiparameter Persistence which partially seeded this work.

References

- [1] H. Asashiba, M. Buchet, E. G. Escolar, K. Nakashima, and M. Yoshiwaki. On Interval Decomposability of 2D Persistence Modules. *arXiv e-prints arXiv:1812.05261*, Dec 2018.
- [2] M. Atiyah. On the krull-schmidt theorem with application to sheaves. *Bulletin de la Société Mathématique de France*, 84:307–317, 1956.
- [3] H. B. Bjerkevik and M. B. Botnan. Computational Complexity of the Interleaving Distance. *arXiv e-prints arXiv:1712.04281*, Dec 2017.
- [4] H. B. Bjerkevik, M. B. Botnan, and M. Kerber. Computing the interleaving distance is NP-hard. *Found. Comput. Math*, pages 1237–1271, 2020.
- [5] M. B. Botnan and M. Lesnick. Algebraic stability of zigzag persistence modules. *Algebraic & Geometric Topology*, 18:3133–3204, 2018.
- [6] U. Bauer, M. Kerber, and J. Reininghaus. Distributed computation of persistent homology. *algorithm engineering and experimentation*, pages 31–38, 2014.
- [7] U. Bauer, M. Kerber, J. Reininghaus, and H. Wagner. Phat - persistent homology algorithms toolbox. *Journal of Symbolic Computation*, 78:76–90, 2017.

-
- [8] H. Bjerkevik. Stability of higher-dimensional interval decomposable persistence modules. *arXiv e-print arXiv:1609.02086*, 2016.
- [9] M. B. Botnan, V. Lebovici, and S. Y. Oudot. On rectangle-decomposable 2-parameter persistence modules. In *36th International Symposium on Computational Geometry, SoCG 2020*, volume 164 of *LIPICs*, pages 22:1–22:16. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2020.
- [10] W. Bruns and H. J. Herzog. *Cohen-Macaulay Rings*. Cambridge University Press, 1998.
- [11] M. Buchet and E. G. Escolar. Every 1D Persistence Module is a Restriction of Some Indecomposable 2D Persistence Module. *arXiv e-prints*, page arXiv:1902.07405, Feb 2019.
- [12] J. R. Bunch and J. E. Hopcroft. Triangular factorization and inversion by fast matrix multiplication. *Math. Comput.*, 28:231–236, 1974.
- [13] C. Cai, W. Kim, F. Mémoli, and Y. Wang. Elder-rule-staircodes for augmented metric spaces. In *36th International Symposium on Computational Geometry (SoCG 2020)*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2020.
- [14] G. Carlsson. Topology and data. *Bulletin of the American Mathematical Society*, 46(2):255–308, 2009.
- [15] G. Carlsson and F. Mémoli. Persistent clustering and a theorem of J. Kleinberg. *arXiv e-print arXiv:0808.2241*, 2008.
- [16] G. Carlsson, G. Singh, and A. Zomorodian. Computing multidimensional persistence. In *International Symposium on Algorithms and Computation*, pages 730–739. Springer, 2009.
- [17] G. Carlsson and A. Zomorodian. The theory of multidimensional persistence. *Discrete & Computational Geometry*, 42(1):71–93, Jul 2009.
- [18] A. Cerri, M. Ethier, and P. Frosini. On the geometrical properties of the coherent matching distance in 2D persistent homology. *arXiv e-prints arXiv:1801.06636*, Jan 2018.
- [19] A. Cerri, B. D. Fabio, M. Ferri, P. Frosini, and C. Landi. Betti numbers in multidimensional persistent homology are stable functions. *Mathematical Methods in the Applied Sciences*, 36(12):1543–1557.
- [20] J. Cochoy and S. Y. Oudot. Decomposition of exact pfd persistence bimodules. *Discret. Comput. Geom.*, 63(2):255–293, 2020.
- [21] D. Cohen-Steiner, H. Edelsbrunner, and J. Harer. Stability of persistence diagrams. *Discrete & Computational Geometry*, 37(1):103–120, Jan 2007.
- [22] D. Cohen-Steiner, H. Edelsbrunner, and D. Morozov. Vines and vineyards by updating persistence in linear time. In *Proc. 22nd Annu. Sympos. Comput. Geom.*, pages 119–126, 2006.
- [23] R. Corbet and M. Kerber. The representation theorem of persistence revisited and generalized. *Journal of Applied and Computational Topology*, 2(1):1–31, Oct 2018.
- [24] D. A. Cox, J. Little, and D. O’Shea. *Using algebraic geometry*, volume 185. Springer Science & Business Media, 2006.
- [25] T. K. Dey, D. Shi, and Y. Wang. Simba: An efficient tool for approximating Rips-filtration persistence via simplicial batch-collapse. *European Symposium on Algorithms*, pages 35:1–16, 2016.
- [26] T. K. Dey and C. Xin. Computing bottleneck distance for 2-d interval decomposable modules. In *34th International Symposium on Computational Geometry, SoCG 2018, June 11-14, 2018, Budapest, Hungary*, pages 32:1–32:15, 2018.
- [27] T. K. Dey and C. Xin. Rectangular approximation and stability of 2-parameter persistence modules. *arXiv e-print arXiv:2108.07429*, 2021.
- [28] H. Edelsbrunner and J. Harer. *Computational Topology: An Introduction*. Applied Mathematics. American Mathematical Society, 2010.
- [29] H. Edelsbrunner, D. Letscher, and A. Zomorodian. Topological persistence and simplification. In *Proceedings 41st Annual Symposium on Foundations of Computer Science*, pages 454–463. IEEE, 2000.
- [30] D. Eisenbud. *The Geometry of Syzygies: A Second Course in Algebraic Geometry and Commutative Algebra*, volume 229. Springer Science & Business Media, 2005.
- [31] A. Hatcher. *Algebraic topology*. Cambridge Univ. Press, Cambridge, 2000.
- [32] D. Hilbert. Ueber die theorie der algebraischen formen. *Mathematische annalen*, 36(4):473–534, 1890.
- [33] D. F. Holt. The meataxe as a tool in computational group theory. *LONDON MATHEMATICAL SOCIETY LECTURE NOTE SERIES*, pages 74–81, 1998.

-
- [34] D. F. Holt and S. Rees. Testing modules for irreducibility. *Journal of the Australian Mathematical Society*, 57(1):1–16, 1994.
- [35] O. H. Ibarra, S. Moran, and R. Hui. A generalization of the fast lup matrix decomposition algorithm and applications. *J. Algorithms*, 3(1):45 – 56, 1982.
- [36] M. Kerber, M. Lesnick, and S. Oudot. Exact Computation of the Matching Distance on 2-Parameter Persistence Modules. In *35th International Symposium on Computational Geometry (SoCG 2019)*, volume 129 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 46:1–46:15, 2019.
- [37] K. P. Knudson. A refinement of multi-dimensional persistence. *arXiv e-print arXiv:0706.2608*, 2007.
- [38] M. Lesnick. The theory of the interleaving distance on multidimensional persistence modules. *Foundations of Computational Mathematics*, 15(3):613–650, 2015.
- [39] M. Lesnick and M. Wright. Interactive Visualization of 2-D Persistence Modules. *arXiv e-prints arXiv:1512.00180*, Dec 2015.
- [40] M. Lesnick and M. Wright. Computing Minimal Presentations and Betti Numbers of 2-Parameter Persistent Homology. *arXiv e-prints arXiv:1902.05708*, Feb 2019.
- [41] S. Liu, D. Maljovec, B. Wang, P. T. Bremer, and V. Pascucci. Visualizing high-dimensional data: Advances in the past decade. *IEEE Transactions on Visualization and Computer Graphics*, 23(3):1249–1268, 2017.
- [42] E. Miller and B. Sturmfels. *Combinatorial Commutative Algebra*. 2004.
- [43] T. Römer. On minimal graded free resolutions. *Illinois J. Math*, 45(2):1361–1376, 2001.
- [44] D. R. Sheehy. Linear-size approximations to the vietoris–rips filtration. *Discrete & Computational Geometry*, 49(4):778–796, 2013.
- [45] J. Skryzalin. Numeric invariants from multidimensional persistence, 2016. PhD thesis, Stanford University.
- [46] C. A. Weibel. *An introduction to homological algebra*. Number 38. Cambridge university press, 1995.
- [47] S. Y. Oudot. *Persistence Theory: From Quiver Representations to Data Analysis*, volume 209. American Mathematical Society, 2015.
- [48] A. Zomorodian and G. Carlsson. Computing persistent homology. *Discrete & Computational Geometry*, 33(2):249–274, 2005.

Appendices

A Free resolution and graded Betti numbers

Here we introduce free resolutions and graded Betti numbers of graded modules. Based on these tools, we give a proof of our Theorem 3.1.

Definition A.1. For a graded module M , a free resolution $\mathcal{F} \rightarrow M$ is an exact sequence:

$$\dots \longrightarrow F^2 \xrightarrow{f^2} F^1 \xrightarrow{f^1} F^0 \xrightarrow{f^0} M \longrightarrow 0 \quad \text{where each } F^i \text{ is a free graded } R\text{-module.}$$

We say two free resolutions \mathcal{F}, \mathcal{G} of M are isomorphic, denoted as $\mathcal{F} \simeq \mathcal{G}$, if there exists a collection of isomorphisms $\{h^i : F^i \rightarrow G^i\}_{i=0,1,\dots}$ which commutes with f^i 's and g^i 's. That is, for all $i = 0, 1, \dots$, $g^i \circ h^i = h^{i-1} \circ f^i$ where h^{-1} is the identity map on M . See the following commutative diagram as an illustration.

$$\begin{array}{ccccccc} \dots & \longrightarrow & F^1 & \xrightarrow{f^1} & F^0 & \xrightarrow{f^0} & M \longrightarrow 0 \\ & & \downarrow h^1 & & \downarrow h^0 & & \downarrow 1 \\ \dots & \longrightarrow & G^1 & \xrightarrow{g^1} & G^0 & \xrightarrow{g^0} & M \longrightarrow 0 \end{array}$$

For two free resolutions $\mathcal{F} \rightarrow M$ and $\mathcal{G} \rightarrow N$, by taking direct sums of free modules $F^i \oplus G^i$ and morphisms $f^i \oplus g^i$, we get a free resolution of $M \oplus N$, denoted as $\mathcal{F} \oplus \mathcal{G}$.

Note that a presentation of M can be viewed as the tail part $F^1 \xrightarrow{f^1} F^0 \xrightarrow{f^0} M \longrightarrow 0$ of a free resolution $\mathcal{F} \rightarrow M$. Free resolutions and presentations are not unique. But there exists a unique minimal free resolution in the following sense:

Fact 4. For a graded module M , there exists a unique free resolution such that $\forall i \geq 0, \text{im } f_{i+1} \subseteq \mathfrak{m}F_i$, where $\mathfrak{m} = (x_1, \dots, x_d)$ is the unique maximal ideal of the graded ring $R = \mathbb{k}[x_1, \dots, x_d]$.

Definition A.2. In a minimal free resolution $\mathcal{F} \rightarrow M$, the tail part $F^1 \xrightarrow{f^1} F^0 \xrightarrow{f^0} M \rightarrow 0$ is called the *minimal presentation* of M and f^1 is called the *minimal presentation map* of M .

Here we briefly state the construction of the unique free resolution without formal proof. More details can be found in [10, 43]:

Construction A.1. Choose a minimal set of homogeneous generators g_1, \dots, g_n of M . Let $F^0 = \bigoplus_{i=1}^n R_{\rightarrow \text{gr}(g_i)}$ with standard basis $e_1^{\text{gr}(g_1)}, \dots, e_n^{\text{gr}(g_n)}$ of F^0 . The homogeneous R -map $f^0 : F^0 \rightarrow M$ is determined by $f^0(e_i) = g_i$.

Now the 1st syzygy module of M , $S_1 \xrightarrow{\ker f^0} F^0$, is again a finitely generated graded R -module. We choose a minimal set of homogeneous generators s_1, \dots, s_m of S_1 and let $F^1 = \bigoplus_{j=1}^m R_{\rightarrow \text{gr}(s_j)}$ with standard basis $e_1^{\text{gr}(s_1)}, \dots, e_m^{\text{gr}(s_m)}$ of F^1 . The homogeneous R -map $f^1 : F^1 \rightarrow F^0$ is determined by $f^1(e_j) = s_j$. By repeating this procedure for $S_2 = \ker f^1$ and moving backward further, one gets a graded free resolution of M .

Fact 5. Any free resolution of M can be obtained (up to isomorphism) from the minimal free resolution by summing it with free resolutions of trivial modules, each with the following form

$$\dots \longrightarrow 0 \longrightarrow F^{i+1} \xrightarrow{f^{i+1}=\mathbb{1}} F^i \longrightarrow 0 \longrightarrow \dots \longrightarrow N = 0 \longrightarrow 0$$

Note that the only nontrivial morphism $F^{i+1} \xrightarrow{f^{i+1}=\mathbb{1}} F^i$ is the identity map $\mathbb{1}$.

From the above constructions, it is not hard to see that this unique free resolution is a minimal one in the sense that each free module F^j has smallest possible size of basis.

For this unique free resolution, for each j , we can write $F^j \simeq \bigoplus_{\mathbf{u} \in \mathbb{Z}^d} \bigoplus^{\beta_{j,\mathbf{u}}^M} R_{\rightarrow \mathbf{u}}$ (the notation $\bigoplus^{\beta_{j,\mathbf{u}}^M} R_{\rightarrow \mathbf{u}}$ means the direct sum of $\beta_{j,\mathbf{u}}^M$ copies of $R_{\rightarrow \mathbf{u}}$). The set $\{\beta_{j,\mathbf{u}}^M \mid j \in \mathbb{N}, \mathbf{u} \in \mathbb{Z}^d\}$ is called *the graded Betti numbers* of M . When M is clear, we might omit the upper index in Betti number. For example, the graded Betti number of the persistence module for our working example 1 is listed as Table 2.

β^M	(1,0)	(0,1)	(1,1)	(2,1)	(1,2)	(2,2)	\dots
β_0	1	1	1				
β_1			1	1	1		
β_2						1	
$\beta_{\geq 3}$							

Table 2: All the nonzero graded Betti numbers $\beta_{i,\mathbf{u}}$ are listed in the table. Empty items are all zeros.

Note that the graded Betti number of a module is uniquely determined by the unique minimal free resolution. On the other hand, if a free resolution $\mathcal{G} \rightarrow M$ with $G^j \simeq \bigoplus_{\mathbf{u} \in \mathbb{Z}^d} \bigoplus^{\gamma_{j,\mathbf{u}}^M} R_{\rightarrow \mathbf{u}}$ satisfies $\gamma_{j,\mathbf{u}}^M = \beta_{j,\mathbf{u}}^M$ everywhere, then $\mathcal{G} \simeq \mathcal{F}$ is also a minimal free resolution of M .

Fact 6. $\beta_{*,*}^{M \oplus N} = \beta_{*,*}^M + \beta_{*,*}^N$

Proposition A.1. Given a graded module M with a decomposition $M \simeq M^1 \oplus M^2$, let $\mathcal{F} \rightarrow M$ be the minimal resolution of M , and $\mathcal{G} \rightarrow M^1$ and $\mathcal{H} \rightarrow M^2$ be the minimal resolution of M^1 and M^2 respectively, then $\mathcal{F} \simeq \mathcal{G} \oplus \mathcal{H}$.

Proof. $\mathcal{G} \oplus \mathcal{H} \rightarrow M$ is a free resolution. We need to show it is a minimal free resolution. By previous argument, we just need to show that the graded Betti numbers of $\mathcal{G} \oplus \mathcal{H} \rightarrow M^1 \oplus M^2$ coincide with graded Betti numbers of $\mathcal{F} \rightarrow M$. This is true by the fact 6. \square

Note that the free resolution is an extension of free presentation. So the above proposition applies to free presentation, which immediately results in the following Corollary.

Corollary A.2. Given a graded module M with a decomposition $M \simeq M^1 \oplus M^2$, let f be its minimal presentation map, and g, h be the minimal presentation maps of M^1, M^2 respectively, then $f \simeq g \oplus h$.

We also have the following fact relating morphisms:

Fact 7. $\ker(f^1 \oplus f^2) = \ker f^1 \oplus \ker f^2$; $\text{coker}(f^1 \oplus f^2) = \text{coker} f^1 \oplus \text{coker} f^2$.

Based on the above statements, now we can prove Theorem 3.1

proof of Theorem 3.1. With the obvious correspondence $[f_i] \leftrightarrow [f]_i$, (2 \leftrightarrow 3) easily follows from our arguments about matrix diagonalization in the main context.

(1 \rightarrow 2) Given $H \simeq \bigoplus H^i$ with the minimal presentation maps f of H : For each H^i , there exists a minimal presentation map f_i . By Corollary A.2, we have $f \simeq \bigoplus f_i$.

(2 \rightarrow 1) Given $f \simeq \bigoplus f_i$: Since $H = \text{coker} f = \text{coker}(\bigoplus f_i) = \bigoplus \text{coker} f_i$, let $H^i = \text{coker} f_i$, we have the decomposition $H \simeq \bigoplus H^i$.

It follows that the above two constructions together give the desired 1-1 correspondence. \square

proof of Proposition 3.2. We start with (2). Consider the total decomposition $f \simeq \bigoplus f^i$. By Remark 3.1, any presentation is isomorphic to a direct sum of the minimal presentation and some trivial presentations. Let $f \simeq g \oplus h$ with g being the minimal presentation. So $\text{coker} h = 0$. Let $g \simeq \bigoplus g^j$ and $h \simeq \bigoplus h^k$ be the total decomposition of g and h . Note that $\forall k, \text{coker} h^k = 0$. Now we have $\text{coker} f \simeq \bigoplus \text{coker} f^i$ with $\text{coker} f^i$ being either $\text{coker} g^j$ or 0, by the essentially uniqueness of total decomposition. With $H \simeq \bigoplus \text{coker} g^j$ being a total decomposition of H by Remark 3.2, and $\bigoplus \text{coker} f^i = \bigoplus \text{coker} g^j \oplus 0$, we can say that $H \simeq \bigoplus \text{coker} f^i$ is also a total decomposition.

Now for (1). For any decomposition $H \simeq \bigoplus H^i$, it is not hard to see that each H^i can be written as a direct sum of a subset of H_*^j 's with $H \simeq \bigoplus H_*^j$ being the total decomposition of H . One just need to combine the f^i 's correspondingly in the total decomposition of $f \simeq \bigoplus f^i$ to get the desired decomposition of f . \square

B Missing proofs in Section 4

Proposition (4.5). The target block $\mathbf{A}|_T$ can be reduced to 0 while preserving the prior if and only if $\mathbf{A}|_T$ can be written as a linear combination of independent operations. That is,

$$\mathbf{A}|_T = \sum_{\substack{l \notin \text{Row}(T) \\ k \in \text{Row}(T)}} \alpha_{k,l} \mathbf{X}^{k,l}|_T + \sum_{\substack{i \notin \text{Col}(T) \\ j \in \text{Col}(T)}} \beta_{i,j} \mathbf{Y}^{i,j}|_T \quad (6)$$

where $\alpha_{k,l}$'s and $\beta_{i,j}$'s are coefficient in $\mathbb{k} = \mathbb{F}_2$.

Proof. Everything in the statement of the proposition is restricted to T . For simplicity of notations, we omit the lower script $\leq t$ by assuming $\mathbf{A}_{\leq t} = \mathbf{A}$, i.e., $t = m$ is the last column index. It can be verified that this omission does not affect the proof. The simple reason is that because of the admissible rules of column operations, entries beyond column t carried by any admissible operations will never affect entries in $\mathbf{A}_{\leq t}$.

Recall that $\mathbf{Y}^{i,j} = \mathbf{A} \cdot [\delta_{i,j}]$ for some $(i, j) \in \text{Colop}$ and $\mathbf{X}^{k,l} = [\delta_{k,l}] \cdot \mathbf{A}$ for some $(l, k) \in \text{Rowop}$ where

$$\text{Colop} = \{(i, j) \mid c_i \rightarrow c_j \text{ is an admissible column operation}\} \subseteq \text{Col}(\mathbf{A}) \times \text{Col}(\mathbf{A}) \text{ and}$$

$$\text{Rowop} = \{(l, k) \mid r_l \rightarrow r_k \text{ is an admissible row operation}\} \subseteq \text{Row}(\mathbf{A}) \times \text{Row}(\mathbf{A})$$

Let \mathbf{I} be the identity matrix. We say a matrix \mathbf{P} is an admissible left multiplication matrix if $\mathbf{P} = \mathbf{I} + \sum_{\text{Rowop}} \alpha_{k,l} [\delta_{k,l}]$ for some $(l, k) \in \text{Rowop}$, $\alpha_{k,l} \in \{0, 1\}$. Similarly, we say a matrix \mathbf{Q} is an admissible right multiplication matrix if $\mathbf{Q} = \mathbf{I} + \sum_{\text{Colop}} \beta_{i,j} [\delta_{i,j}]$ for some $(i, j) \in \text{Colop}$, $\beta_{i,j} \in \{0, 1\}$. In short, we just say \mathbf{P} and \mathbf{Q} are admissible.

It is not difficult to observe the following properties of admissible matrices:

Fact 8. Matrix $\mathbf{A}' \sim \mathbf{A}$ is an equivalent matrix transformed from \mathbf{A} by a sequence of admissible operations if and only if $\mathbf{A}' = \mathbf{PAQ}$ for some admissible \mathbf{P} and \mathbf{Q} .

Fact 9. Admissible matrices are closed under multiplication and taking inverse.

Fact 10. For any admissible \mathbf{P} , let $S \subseteq \text{Row}(\mathbf{P})$ be any subset of row indices. Then $\mathbf{P}|_{S \times S}$ is invertible.

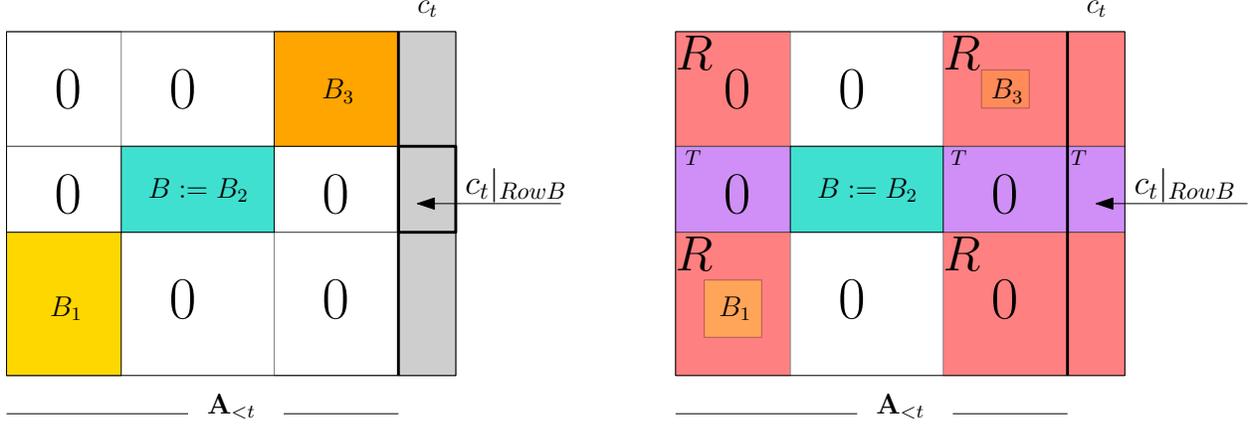


Figure 13: (Left) \mathbf{A} at iteration t during reduction of the sub-column $c_t|_{\text{Row}(B)}$ for the block $B = B_2$. (Right) Target block T shown in magenta includes the sub-column of c_t . It does not include $B := B_2$. All rows external to T have zeros in the columns external to T . All columns external to T have zeros in the rows external to T . Red regions combined form R .

For the last fact, observe that the matrix $\mathbf{P}|_{S \times S}$ can be embedded as a block of an admissible matrix \mathbf{P}' constructed by making all off-diagonal entries of \mathbf{P} whose indices are not in $S \times S$ to be zero. The matrix \mathbf{P}' is obviously admissible. So by the second fact, it is invertible. Also, \mathbf{P}' can be written in block diagonal form with two blocks $\mathbf{P}'|_{S \times S}$ and $\mathbf{P}'|_{\bar{S} \times \bar{S}} = \mathbf{I}$ where $\bar{S} = \text{Row}(\mathbf{P}') - S$. Therefore, if \mathbf{P}' is invertible, so is $\mathbf{P}|_{S \times S} = \mathbf{P}'|_{S \times S}$.

We write the matrix \mathbf{A} in the following block forms with respect to B and T with necessary reordering of rows and columns (See Figure 13 for a simple illustration without reordering rows and columns):

$$\mathbf{A} = \left[\begin{array}{c|c} R & 0 \\ \hline T & B \end{array} \right]$$

Here we abuse the notations of block and index block to make the expression more legible. In the above block forms of \mathbf{A} , for example, T represents the entries of \mathbf{A} on the index block T , that is the block $\mathbf{A}|_T$, which is the target block we want to reduce. Note that

$$\begin{aligned} R &= [\text{Row}(\mathbf{A}) \setminus \text{Row}(B), \text{Col}(\mathbf{A}_{<t}) \setminus \text{Col}(B)] \\ &= \left[\bigoplus_{B_i \neq B} B_i \right] \cup [\text{Row}(\mathbf{A}) \setminus \text{Row}(B), \{t\}] \end{aligned}$$

which is the block obtained by merging all other previous index blocks together with the sub-column of t excluding entries on $\text{Row}(B)$. The right top block is zero since it belongs to the intersections of rows and columns from different blocks.

Observe that, the target block T can be reduced to 0 in \mathbf{A} with prior preserved if and only if

$$\mathbf{P}\mathbf{A}\mathbf{Q} := \mathbf{P} \cdot \left[\begin{array}{c|c} R & 0 \\ \hline T & B \end{array} \right] \cdot \mathbf{Q} = \left[\begin{array}{c|c} R & 0 \\ \hline 0 & B \end{array} \right] \quad (7)$$

for some admissible \mathbf{P} and \mathbf{Q} .

For \Leftarrow direction, consider $\mathbf{P} = \mathbf{I} + \sum \alpha_{k,l}[\delta_{k,l}]$ and $\mathbf{Q} = \mathbf{I} + \sum \beta_{i,j}[\delta_{i,j}]$ with binary coefficients $\alpha_{k,l}$'s and $\beta_{i,j}$'s given in Equation 6. Then, we have

$$\mathbf{P}\mathbf{A}\mathbf{Q} = (\mathbf{I} + \sum \alpha_{k,l}[\delta_{k,l}])\mathbf{A}(\mathbf{I} + \sum \beta_{i,j}[\delta_{i,j}]) \quad (8)$$

$$= \mathbf{A} + \sum \alpha_{k,l}[\delta_{k,l}]\mathbf{A} + \sum \beta_{i,j}\mathbf{A}[\delta_{i,j}] + \sum \sum \alpha_{k,l}\beta_{i,j}[\delta_{k,l}]\mathbf{A}[\delta_{i,j}] \quad (9)$$

$$= \mathbf{A} + \sum \alpha_{k,l}[\delta_{k,l}]\mathbf{A} + \sum \beta_{i,j}\mathbf{A}[\delta_{i,j}] \quad (10)$$

$$= \mathbf{A} + \sum \alpha_{k,l}\mathbf{X}^{k,l} + \sum \beta_{i,j}\mathbf{Y}^{i,j} \quad (11)$$

The third Equation (10) follows from Observations 4.4. After restriction to T , by the assumption that $\sum \alpha_{k,l} \mathbf{X}^{k,l} + \sum \beta_{i,j} \mathbf{Y}^{i,j} = \mathbf{A}|_T$, we get $\mathbf{P}\mathbf{A}\mathbf{Q}|_T = 0$. By the definition of independent operations and Observation 4.3, one can see that our \mathbf{P}, \mathbf{Q} solves Equation 7.

For \Rightarrow , we will show that if the above equation is solvable, then there always exist solutions \mathbf{P}' and \mathbf{Q}' in a simpler forms as stated in the following proposition.

Proposition B.1. *Equation (7) is solvable for some admissible \mathbf{P} and \mathbf{Q} if and only if it is solvable for some admissible \mathbf{P}' and \mathbf{Q}' in the following form:*

$$\mathbf{P}' = \left[\begin{array}{c|c} I & 0 \\ \hline U & I \end{array} \right] \text{ and } \mathbf{Q}' = \left[\begin{array}{c|c} I & 0 \\ \hline V & I \end{array} \right] \quad (12)$$

Before we prove Proposition B.1, we show how one can prove the \Rightarrow direction in Proposition 4.5 from it. Based on the equivalent condition Equation 7 and Proposition B.1, we can write \mathbf{P}' and \mathbf{Q}' in formula 12 as

$$\mathbf{P}' = \mathbf{I} + \sum_{\substack{(l,k) \in \\ \text{Rowop}_{R \rightarrow T}}} \alpha_{k,l} [\delta_{k,l}] \quad \mathbf{Q}' = \mathbf{I} + \sum_{\substack{(i,j) \in \\ \text{Colop}_{B \rightarrow T}}} \beta_{i,j} [\delta_{i,j}]$$

where $\text{Rowop}_{R \rightarrow T} = \{(l, k) \in \text{Rowop} \mid (l, k) \in \text{Row}(R) \times \text{Row}(T)\}$ and $\text{Colop}_{B \rightarrow T} = \{(i, j) \in \text{Colop} \mid (i, j) \in \text{Col}(B) \times \text{Col}(T)\}$, and $\alpha_{k,l}, \beta_{i,j} \in \{0, 1\}$. Then, similar to Equation 11, we get

$$\begin{aligned} \mathbf{P}'\mathbf{A}\mathbf{Q}' &= (\mathbf{I} + \sum_{\substack{(l,k) \in \\ \text{Rowop}_{R \rightarrow T}} \alpha_{k,l} [\delta_{k,l}]) \cdot \mathbf{A} \cdot (\mathbf{I} + \sum_{\substack{(i,j) \in \\ \text{Colop}_{B \rightarrow T}} \beta_{i,j} [\delta_{i,j}]) \\ &= \mathbf{A} + \sum_{\substack{(l,k) \in \\ \text{Rowop}_{R \rightarrow T}} \alpha_{k,l} [\delta_{k,l}] \cdot \mathbf{A} + \sum_{\substack{(i,j) \in \\ \text{Colop}_{B \rightarrow T}} \beta_{i,j} \mathbf{A} \cdot [\delta_{i,j}] + \sum_{\substack{(l,k) \in \\ \text{Rowop}_{R \rightarrow T}}} \sum_{\substack{(i,j) \in \\ \text{Colop}_{B \rightarrow T}}} \alpha_{k,l} \beta_{i,j} [\delta_{k,l}] \cdot \mathbf{A} \cdot [\delta_{i,j}] \\ &= \mathbf{A} + \sum_{\substack{(l,k) \in \\ \text{Rowop}_{R \rightarrow T}} \alpha_{k,l} [\delta_{k,l}] \cdot \mathbf{A} + \sum_{\substack{(i,j) \in \\ \text{Colop}_{B \rightarrow T}} \beta_{i,j} \mathbf{A} \cdot [\delta_{i,j}] \\ &= \mathbf{A} + \sum_{\substack{(l,k) \in \\ \text{Rowop}_{R \rightarrow T}} \alpha_{k,l} \mathbf{X}^{k,l} + \sum_{\substack{(i,j) \in \\ \text{Colop}_{B \rightarrow T}} \beta_{i,j} \mathbf{Y}^{i,j} \end{aligned}$$

By restriction on T we have

$$\mathbf{P}'\mathbf{A}\mathbf{Q}'|_T = \mathbf{A}|_T + \sum_{\substack{(l,k) \in \\ \text{Rowop}_{R \rightarrow T}}} \alpha_{k,l} \mathbf{X}^{k,l}|_T + \sum_{\substack{(i,j) \in \\ \text{Colop}_{B \rightarrow T}}} \beta_{i,j} \mathbf{Y}^{i,j}|_T \quad (13)$$

With $\mathbf{P}'\mathbf{A}\mathbf{Q}'|_T = 0$ by our assumption, we get

$$\mathbf{A}|_T = \sum_{\substack{(l,k) \in \\ \text{Rowop}_{R \rightarrow T}}} \alpha_{k,l} \mathbf{X}^{k,l}|_T + \sum_{\substack{(i,j) \in \\ \text{Colop}_{B \rightarrow T}}} \beta_{i,j} \mathbf{Y}^{i,j}|_T$$

This is exactly what we want

$$\mathbf{A}|_T = \sum_{\substack{l \notin \text{Row}(T) \\ k \in \text{Row}(T)}} \alpha_{k,l} \mathbf{X}^{k,l}|_T + \sum_{\substack{i \notin \text{Col}(T) \\ j \in \text{Col}(T)}} \beta_{i,j} \mathbf{Y}^{i,j}|_T \quad (14)$$

□

Now we give the proof of Proposition B.1.

Proof of Proposition B.1. The \Leftarrow direction is trivial. For the \Rightarrow direction, we want to show that, if Equation (7) is solvable for some admissible \mathbf{P} and \mathbf{Q} , then there exists admissible \mathbf{P}' and \mathbf{Q}' so that

$$\mathbf{P}' = \left[\begin{array}{c|c} I & 0 \\ \hline U & I \end{array} \right], \mathbf{Q}' = \left[\begin{array}{c|c} I & 0 \\ \hline V & I \end{array} \right], \text{ and } \mathbf{P}' \cdot \left[\begin{array}{c|c} R & 0 \\ \hline T & B \end{array} \right] \cdot \mathbf{Q}' = \left[\begin{array}{c|c} R & 0 \\ \hline UR + BV + T & B \end{array} \right] = \left[\begin{array}{c|c} R & 0 \\ \hline 0 & B \end{array} \right]$$

We write \mathbf{P} and \mathbf{Q} in corresponding block forms as follows:

$$\mathbf{P} = \left[\begin{array}{c|c} P_1 & P_2 \\ \hline P_3 & P_4 \end{array} \right] \text{ and } \mathbf{Q} = \left[\begin{array}{c|c} Q_1 & Q_2 \\ \hline Q_3 & Q_4 \end{array} \right] \quad (15)$$

From Equation (7) one can get a set of equations

$$P_1 R Q_2 + P_2 B Q_4 = 0 \quad (16)$$

$$P_1 R Q_1 + P_2 B Q_3 = R \quad (17)$$

$$P_3 R Q_2 + P_4 B Q_4 = B \quad (18)$$

$$P_3 R Q_1 + P_4 B Q_3 = T \quad (19)$$

From Fact 10, we know that P_1, P_4, Q_1, Q_4 are invertible. By left multiplication with P_1^{-1} and right multiplication with Q_4^{-1} on both sides of Equation (16), one can get :

$$P_1^{-1} P_1 R Q_2 Q_4^{-1} + P_1^{-1} P_2 B Q_4 Q_4^{-1} = R Q_2 Q_4^{-1} + P_1^{-1} P_2 B = 0 \implies -R Q_2 Q_4^{-1} = P_1^{-1} P_2 B \quad (20)$$

Similarly, by left multiplication with P_1^{-1} on both sides of Equation (17) and by right multiplication with Q_4^{-1} on both sides of Equation (18), one can get the following equations:

$$P_1 R Q_1 + P_2 B Q_3 = R \implies R Q_1 = P_1^{-1} R - P_1^{-1} P_2 B Q_3 \quad (21)$$

$$P_3 R Q_2 + P_4 B Q_4 = B \implies P_4 B = B Q_4^{-1} - P_3 R Q_2 Q_4^{-1} \quad (22)$$

Now from Equation 19, we have:

$$T = P_3 R Q_1 + P_4 B Q_3$$

Equation 21 and 22 \rightarrow

$$\begin{aligned} T &= P_3 (P_1^{-1} R - P_1^{-1} P_2 B Q_3) + (B Q_4^{-1} - P_3 R Q_2 Q_4^{-1}) Q_3 \\ &= P_3 P_1^{-1} R + B Q_4^{-1} Q_3 - P_3 P_1^{-1} P_2 B Q_3 - P_3 R Q_2 Q_4^{-1} Q_3 \end{aligned}$$

Equation 20 \rightarrow

$$\begin{aligned} T &= P_3 P_1^{-1} R + B Q_4^{-1} Q_3 - P_3 P_1^{-1} P_2 B Q_3 + P_3 P_1^{-1} P_2 B Q_3 \\ &= P_3 P_1^{-1} R + B Q_4^{-1} Q_3 \end{aligned}$$

Letting $U = P_3 P_1^{-1}$ and $V = Q_4^{-1} Q_3$, we get the desired equation. Now we just need to show that \mathbf{P}', \mathbf{Q}' are both admissible. We prove it for \mathbf{Q}' . Similar proof holds for \mathbf{P}' . We want to show that for any $(i, j) \in \text{Row}(V) \times \text{Col}(V)$, if $\mathbf{Q}'_{i,j} = 1$, then $(i, j) \in \text{Colop}$. From equality, $V = Q_4^{-1} Q_3$, which implies $\mathbf{Q}'_{i,j} = \sum_k (Q_4^{-1})_{i,k} \cdot (Q_3)_{k,j} = 1$, we know that $(Q_4^{-1})_{i,k} = (Q_3)_{k,j} = 1$ for some k . Since Q_4^{-1} and Q_3 are both blocks in the admissible matrix \mathbf{Q} , by the definition of admissible left multiplication matrix, we have $(i, k), (k, j) \in \text{Colop}$. Note that Colop is closed under transitive relation by Proposition 4.1. So we have $(i, j) \in \text{Colop}$.

□