# The Complexity of Generating Functions for Integer Points in Polyhedra and Beyond

Alexander Barvinok*

**Abstract.** Motivated by the formula for the sum of the geometric series, we consider various classes of sets $S \subset \mathbb{Z}^d$ of integer points for which an a priori "long" Laurent series or polynomial $\sum_{m \in S} \mathbf{x}^m$ can be written as a "short" rational function $f(S; \mathbf{x})$. Examples include the sets of integer points in rational polyhedra, integer semigroups, and Hilbert bases of rational cones, among others. We discuss applications to efficient counting and optimization and open questions.

**Mathematics Subject Classification (2000).** Primary 05A15; Secondary 68W30, 11P21, 52C07, 11H06.

**Keywords.** lattice point, rational polytope, generating function, rational function, Laurent polynomial, integer semigroup, Hilbert basis, efficient counting, computational complexity

## 1. Introduction

Our inspiration comes from a formula for the sum of a finite geometric series:

$$\sum_{m=0}^{n} x^m = \frac{1 - x^{n+1}}{1 - x}. \tag{1.1}$$

We look at the formula from several points of view.

Geometrically, the left hand side of (1.1) represents the sum over all integer points in a one-dimensional polytope. Namely, with every integer point $m$ we associate a monomial $x^m$ and then consider the sum over all integer points in the interval $[0, n]$.

From the computational complexity point of view, the left hand side of (1.1) is a "long" polynomial whereas the right hand side of (1.1) is a "short" rational function. More precisely, to write an integer $m$ we need about $\log m$ digits or bits. Consequently, to write the left hand side of (1.1), we need about $n \log n$ bits. On

the other hand, to write the right hand side of (1.1) we need only about $\log n$ bits. Thus the left hand side is exponentially longer than the right hand side.

Finally, let us read (1.1) from right to left. We can ask how to extract various facts about the set $S$ of integer points in the interval $[0, n]$ from the rational function encoding. For example, to compute the number $|S|$ of points we substitute $x = 1$ into the right hand side of (1.1). Although $x = 1$ is a pole of the rational function, we can compute the desired value by applying l'Hospital's rule.

Let $\mathbb{R}^d$ be Euclidean space with the standard basis $e_1, \ldots, e_d$, so a point $x \in \mathbb{R}^d$ is identified with the $d$-tuple $x = (\xi_1, \ldots, \xi_d)$ of its coordinates, and let $\mathbb{Z}^d \subset \mathbb{R}^d$ be the standard integer lattice, that is the set of points with integer coordinates. With every integer point $m = (\mu_1, \ldots, \mu_d)$ we associate the Laurent monomial

$$\mathbf{x}^m = x_1^{\mu_1} \ldots x_d^{\mu_d}$$

in $d$ complex variables $\mathbf{x} = (x_1, \ldots, x_d)$. We agree that $x_i^0 = 1$.

Let $S \subset \mathbb{Z}^d$ be a finite set and let us consider the sum

$$f(S; \mathbf{x}) = \sum_{m \in S} \mathbf{x}^m.$$

Thus $f(S; \mathbf{x})$ is a Laurent polynomial that is the generating function of the set $S$. We are interested in the following general questions:

• For which sets $S \subset \mathbb{Z}^d$ a potentially long Laurent polynomial $f(S; \mathbf{x})$ can be written as a short rational function?

• What information about the set $S$ can be extracted from $f(S; \mathbf{x})$ given as a short rational function?

The paper is organized as follows.

In Section 2, we discuss necessary preliminaries from the theory of computational complexity, define what "long" and "short" means and show that if $S$ is the set of integer points in a rational polyhedron $P \subset \mathbb{R}^d$ then the generating function $f(S; \mathbf{x})$ can be computed in polynomial time as a short rational function, provided the dimension $d$ of the ambient space is fixed in advance. We discuss applications to efficient counting and optimization and practical implementations of the algorithms.

In Section 3, we discuss what information can we extract from a set $S \subset \mathbb{Z}^d$ *defined* by its generating function $f(S; \mathbf{x})$ written as a rational function. In particular, we show that if $S_1, S_2 \subset \mathbb{Z}^d$ are two finite sets defined by their rational generating functions $f(S_1; \mathbf{x})$ and $f(S_2; \mathbf{x})$, then the generating function $f(S; \mathbf{x})$ of their intersection $S = S_1 \cap S_2$ can be computed in polynomial time as a rational function.

In Section 4, we show that if $S \subset \mathbb{Z}_+$ is an integer semigroup with a fixed number $d$ of generators, then $f(S; x)$ can be computed in polynomial time as a short rational function. This result is obtained as a corollary of a more general

result that the *projection* of the set of integer points in a rational polytope admits a polynomial time computable rational generating function. We mention some other examples such as Hilbert bases of rational cones.

In Section 5, we consider the results of Sections 2 and 4 in the general context of Presburger arithmetic. We argue that the "natural" class of sets $S \subset \mathbb{Z}^d$ with short rational generating functions $f(S; \mathbf{x})$ would have been the class of sets defined by formulas of Presburger arithmetic where all combinatorial parameters (the number of variables and Boolean operations) are fixed and only numerical constants are allowed to vary. As the paper is being written, this is still a conjecture.

In Section 6, we try to identify the natural boundaries of the developed theory. We also discuss the emerging picture of what happens if the dimension $d$ of the ambient space is allowed to grow.

## 2. Rational Polyhedra

Formula (1.1) admits an extension to general rational polyhedra.

**Definition 2.1.** The set $P \subset \mathbb{R}^d$ of solutions to a system of finitely many linear inequalities is called a *polyhedron*:

$$P = \left\{ (\xi_1, \ldots, \xi_d) : \quad \sum_{j=1}^{d} \alpha_{ij} \xi_j \leq \beta_i, \quad i = 1, \ldots, n \right\}. \qquad (2.1)$$

Here $\alpha_{ij}$ and $\beta_i$ are real numbers. A bounded polyhedron is called a *polytope*. A polyhedron $P$ is called *rational* if in (2.1) one can choose all $\alpha_{ij}$ and $\beta_i$ integer.

To state an analogue of formula (1.1) we need to discuss the notion of the input size. As we remarked earlier, to write an integer $a$ we need roughly $\lceil \log_2 |a| \rceil + 1$ bits. Consequently, to define a rational polyhedron $P \subset \mathbb{R}^d$ by the inequalities (2.1) we need about

$$\mathcal{L} = n(d+1) + \sum_{i,j} \lceil \log_2 |\alpha_{ij}| \rceil + \sum_{i} \lceil \log_2 |\beta_i| \rceil \qquad (2.2)$$

bits. The number $\mathcal{L}$ is called the *input size* of representation (2.1) of $P$.

We are interested in the computational complexity of formulas and algorithms. In particular, we are interested in *polynomial time* algorithms, that is, in the algorithms whose running time is at most $\mathcal{L}^{O(1)}$, where $\mathcal{L}$ is the input size. In what follows, often the dimension $d$ of the ambient space will be fixed in advance and the algorithms will run in polynomial time *for any fixed dimension $d$*. In other words, the running time of such an algorithm is at most $\mathcal{L}^{\phi(d)}$ for some function $\phi$. We use [28] as a general reference in the area of computational complexity and algorithms.

Let $P \subset \mathbb{R}^d$ be a rational polyhedron with a vertex (equivalently, a non-empty polyhedron without lines), possibly unbounded, and let $S = P \cap \mathbb{Z}^d$ be the set of integer points in $P$.

To simplify notation, we denote the generating function

$$f(S; \mathbf{x}) = \sum_{m \in S} \mathbf{x}^m,$$

where $S = P \cap \mathbb{Z}^d$, just by $f(P, \mathbf{x})$.

It is not hard to show that there exists a non-empty open set $U \subset \mathbb{C}^d$ such that for all $\mathbf{x} \in U$ the series

$$f(P, \mathbf{x}) = \sum_{m \in P \cap \mathbb{Z}^d} \mathbf{x}^m$$

converges absolutely and uniformly on compact subsets of $U$ to a rational function in $\mathbf{x}$. It turns out that this rational function can be efficiently computed as long as the dimension $d$ of the ambient space is fixed in advance.

The following result was proved, essentially, in [3] although the formal statement and better complexity bounds did not appear until [4].

**Theorem 2.2.** *Let us fix $d$. Then there exists a polynomial time algorithm, which, for a rational polyhedron $P \subset \mathbb{R}^d$ without lines defined by inequalities (2.1) computes the generating function*

$$f(P, \mathbf{x}) = \sum_{m \in P \cap \mathbb{Z}^d} \mathbf{x}^m$$

*in the form*

$$f(P, \mathbf{x}) = \sum_{i \in I} \epsilon_i \frac{\mathbf{x}^{v_i}}{(1 - \mathbf{x}^{u_{i1}}) \cdots (1 - \mathbf{x}^{u_{id}})}, \tag{2.3}$$

*where $\epsilon_i \in \{-1, 1\}$, $v_i, u_{ij} \in \mathbb{Z}^d$, and $u_{ij} \neq 0$ for all $i, j$.*

The complexity of the algorithm is $\mathcal{L}^{O(d)}$, where $\mathcal{L}$ is the input size of $P$ defined by (2.2). In particular, the number $|I|$ of terms in (2.3) is $\mathcal{L}^{O(d)}$, which is why we call (2.3) a *short rational function*.

Rational cones play the crucial role in the proof of Theorem 2.2.

**2.1. Rational cones.** A non-empty rational polyhedron $K$ is called a *rational cone* if for every $x \in K$ and $\lambda \geq 0$ we have $\lambda x \in K$. We are interested in *pointed* rational cones, that is, cones not containing lines (equivalently, cones for which 0 is the vertex). A basic example of a pointed rational cone is provided by the non-negative orthant $\mathbb{R}_+^d$ consisting of the points with non-negative coordinates. The generating function for the set of integer points in $\mathbb{R}_+^d$ is a multiple geometric series

$$f(\mathbb{R}_+^d, \mathbf{x}) = \sum_{m \in \mathbb{Z}_+^d} \mathbf{x}^m = \prod_{i=1}^d \frac{1}{1 - x_i}.$$

A *unimodular cone* $K$ is the set of non-negative linear combinations of a given basis $u_1, \ldots, u_d$ of the lattice $\mathbb{Z}^d$. Up to an integral change of coordinates, a unimodular

cone $K$ looks like the non-negative orthant $\mathbb{R}^d_+$. Consequently, the generating function for the set of integer points in $K$ is a multiple geometric series

$$f(K, \mathbf{x}) = \sum_{m \in K \cap \mathbb{Z}^d} \mathbf{x}^m = \prod_{i=1}^d \frac{1}{1 - \mathbf{x}^{u_i}}.$$

It is well-known that any rational cone $K$ can be subdivided into unimodular cones, cf., for example, Section 2.6 of [16]. However, even for $d = 2$, the number of the unimodular cones may have to be exponentially large in the input size: consider the cone $K \subset \mathbb{R}^2$ spanned by $(1, 0)$ and $(1, n)$ for a positive integer $n$. Nevertheless, where exists a computationally efficient procedure for constructing a more general *decomposition* of a rational cone into unimodular cones.

**Definition 2.3.** For a set $A \subset \mathbb{R}^d$, let $[A] : \mathbb{R}^d \longrightarrow \mathbb{R}$ be the indicator of $A$ defined by

$$[A](x) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{if } x \notin A. \end{cases}$$

Let $\mathcal{P}(\mathbb{Q}^d)$ be the vector space (over $\mathbb{C}$) spanned by the indicators $[P]$ of rational polyhedra $P \subset \mathbb{R}^d$. We call $\mathcal{P}(\mathbb{Q}^d)$ the *algebra of rational polyhedra*. Vector space $\mathcal{P}(\mathbb{Q}^d)$ possesses an interesting and useful algebra structure, cf. [26], which we don't discuss here.

The idea is to write the indicator $[K]$ of a given rational cone $K \subset \mathbb{R}^d$ as a linear combination of indicators of unimodular cones. For $d = 2$ such an efficient procedure has long been known via the *continued fractions* method, cf., for example, [22]. We give a simple example below.

Suppose that $K \subset \mathbb{R}^2$ is the cone spanned by vectors $(1, 0)$ and $(31, 164)$. Writing the continued fraction expansion, we obtain

$$\frac{164}{31} = 5 + \cfrac{1}{3 + \cfrac{1}{2 + \cfrac{1}{4}}} \quad ,$$

so we write $164/31 = [5; 3, 2, 4]$. Next, we compute the *convergents*

$$[5; 3, 2] = 5 + \cfrac{1}{3 + \cfrac{1}{2}} = \frac{37}{7}, \quad [5; 3] = 5 + \frac{1}{3} = \frac{16}{3}, \quad \text{and} \quad [5] = \frac{5}{1}$$

and notice that

$$[K] = [K_0] - [K_1] + [K_2] - [K_3] + [K_4],$$

where $K_0$ is spanned by $(1, 0)$ and $(0, 1)$, $K_1$ is spanned by $(0, 1)$ and $(1, 5)$, $K_2$ is spanned by $(1, 5)$ and $(3, 16)$, $K_3$ is spanned by $(3, 16)$ and $(7, 37)$, and $K_4$ is

spanned by $(7, 37)$ and $(31, 164)$. Since $K_i$ turn out to be unimodular for $i = 0, 1, 2, 3, 4$, we get the short rational function expression

$$f(K, \mathbf{x}) = \frac{1}{(1 - x_1)(1 - x_2)} - \frac{1}{(1 - x_2)(1 - x_1 x_2^5)} + \frac{1}{(1 - x_1 x_2^5)(1 - x_1^3 x_2^{16})}$$
$$- \frac{1}{(1 - x_1^3 x_2^{16})(1 - x_1^7 x_2^{37})} + \frac{1}{(1 - x_1^7 x_2^{37})(1 - x_1^{31} x_2^{164})}.$$

A polynomial time algorithm for computing a unimodular cone decomposition in any (fixed in advance) dimension $d$ was suggested in [3]. Using triangulations, it is not hard to reduce the case of an arbitrary rational cone to that of a *simple rational cone* $K \subset \mathbb{R}^d$

$$K = \left\{ \sum_{i=1}^d \lambda_i u_i : \quad \lambda_i \geq 0 \right\}$$

spanned by linearly independent vectors $u_1, \ldots, u_d \in \mathbb{Z}^d$, which may not, however, constitute a basis of the lattice $\mathbb{Z}^d$. As a measure of how far is $K$ from being unimodular, we introduce the *index* $\text{ind}(K)$ of $K$ as the index of the sublattice generated by $u_1, \ldots, u_d$ in the ambient lattice $\mathbb{Z}^d$. Thus $\text{ind}(K)$ is a positive integer and $\text{ind}(K) = 1$ if and only if $K$ is a unimodular cone.

Let us consider the parallelepiped

$$\Pi = \left\{ \sum_{i=1}^d \lambda_i u_i : \quad |\lambda_i| \leq \text{ind}^{-1/d}(K) \quad \text{for} \quad i = 1, \ldots, d \right\}.$$

Then $\Pi$ is a convex body symmetric about the origin and $\text{vol } \Pi = 2^d$. Therefore, by the Minkowski Theorem there is a non-zero point $w \in \Pi \cap \mathbb{Z}^d$, cf., for example, Section VII.3 of [5]. Moreover, such a point $w$ can be constructed in polynomial time as long as the dimension $d$ is fixed, cf. Section 6.7 of [17]. Replacing $w$ by $-w$ if needed, we can also ensure that $w$ lies in the same halfspace as $u_1, \ldots, u_d$. Let $K_i$ be the cone spanned by $u_1, \ldots, u_d$ with the vector $u_i$ replaced by $w$ and let $\epsilon_i = 1$ or $\epsilon_i = -1$ depending on whether this replacement preserves or reverses the orientation of the set $u_1, \ldots, u_d$ (we choose $\epsilon_i = 0$ if we obtain a linearly dependent set). Then we observe that

$$[K] = \sum_{i=1}^d \epsilon_i [K_i] \quad \pm \quad \text{indicators of lower-dimensional cones}$$

$$\text{and} \tag{2.4}$$

$$\text{ind}(K_i) \leq \text{ind}^{(d-1)/d}(K) \quad \text{if } \dim K_i = d.$$

As we iterate the above procedure, on the $n$th step, we obtain a decomposition of the cone $K$ as a linear combination of at most $d^n$ cones $K_i$ (not counting smaller-dimensional cones) with

$$\text{ind}(K_i) \leq (\text{ind}(K))^{\left(\frac{d-1}{d}\right)^n}.$$

To ensure that all $K_i$ are unimodular, we can choose $n = O\big(d \log \log \text{ind}(K)\big)$, which results in a polynomial time algorithm for a fixed $d$.

To prove a weaker version of Theorem 2.2 (with $d$ replaced by $d + 2$ in (2.3) and $\mathcal{L}^{O(d^2)}$ complexity) one can note that a rational polyhedron $P \subset \mathbb{R}^d$ without lines can be represented as the section of a pointed rational cone $K \subset \mathbb{R}^{d+1}$ by the affine hyperplane $\xi_{d+1} = 1$. Consequently, we have

$$f(P, \mathbf{x}) = \frac{\partial}{\partial x_{d+1}} f\left(K, (\mathbf{x}, x_{d+1})\right)\Big|_{x_{d+1}=0}. \tag{2.5}$$

### 2.2. Using identities in the algebra of polyhedra. 
The following remarkable result was proved by A.G. Khovanskii and A.V. Pukhlikov [23], and, independently, by J. Lawrence [25].

**Theorem 2.4.** *Let $\mathcal{P}\left(\mathbb{Q}^d\right)$ be the vector space spanned by the indicators of rational polyhedra and let $\mathbb{C}(\mathbf{x})$ be the vector space of rational functions in $d$ complex variables $\mathbf{x} = (x_1, \ldots, x_d)$. There exists a linear transformation $\mathcal{F} : \mathcal{P}\left(\mathbb{Q}^d\right) \longrightarrow \mathbb{C}(\mathbf{x})$ such that*

1. *If $P \subset \mathbb{R}^d$ is a rational polyhedron with a vertex then $\mathcal{F}\big([P]\big) = f(P, \mathbf{x})$, where $f(P, \mathbf{x})$ is the rational function defined as the sum of the series*

$$\sum_{m \in P \cap \mathbb{Z}^d} \mathbf{x}^m$$

*when the series converges absolutely.*

2. *If $P \subset \mathbb{R}^d$ is a rational polyhedron without vertices then $\mathcal{F}\big([P]\big) = 0$.*

*Proof.* Let us fix a decomposition

$$\mathbb{R}^d = \sum_{i \in I} \alpha_i [Q_i] \tag{2.6}$$

for some rational polyhedra $Q_i$ with vertices and some numbers $\alpha_i$. Multiplying (2.6) by $[P]$, we get

$$[P] = \sum_{i \in I} \alpha_i [P \cap Q_i], \tag{2.7}$$

from which we deduce that $\mathcal{P}\left(\mathbb{Q}^d\right)$ is spanned by indicators of rational polyhedra with vertices.

Suppose that we have a linear relation

$$\sum_{j \in J} \beta_j [P_j] = 0 \tag{2.8}$$

for some polyhedra $P_j$ with vertices. Multiplying (2.8) by $[Q_i]$, we get

$$\sum_{j \in J} \beta_j [P_j \cap Q_i] = 0.$$

Since $Q_i$ has a vertex and $P_j \cap Q_i \subset Q_i$, there exists a non-empty open set $U_i \subset \mathbb{C}^d$ such that for all $\mathbf{x} \in U_i$ all the series defining $f(P_j \cap Q_i, \mathbf{x})$ converge absolutely and uniformly on compact subsets of $U_i$. Therefore, we must have

$$\sum_{j \in J} \beta_j f(P_j \cap Q_i, \mathbf{x}) = 0 \quad \text{for all} \quad i \in I.$$

Similarly, from (2.7) we get

$$f(P_j, \mathbf{x}) = \sum_{i \in I} \alpha_i f(P_j \cap Q_i, \mathbf{x}) \quad \text{for all} \quad j \in J.$$

Combining the last two equations, we conclude that

$$\sum_{j \in j} \beta_j f(P_j, \mathbf{x}) = \sum_{i \in I, j \in J} \alpha_i \beta_j f(P_j \cap Q_i, \mathbf{x}) = 0. \tag{2.9}$$

Thus a linear dependence (2.8) among indicators of rational polyhedra $P_j$ with vertices implies the corresponding linear dependence (2.9) among the generating functions $f(P_j, \mathbf{x})$. Therefore, the correspondence

$$[P] \longmapsto f(P, \mathbf{x})$$

extends to a linear transformation $\mathcal{F} : \mathcal{P}(\mathbb{Q}^d) \longrightarrow \mathbb{C}(\mathbf{x})$. It remains to show that $\mathcal{F}([P]) = 0$ if $P$ is a rational polyhedron with a line.

We observe that if $P' = P + u$ is a translation of $P$ by a lattice vector $u$, we must have $f(P', \mathbf{x}) = \mathbf{x}^u f(P, \mathbf{x})$ for all rational polyhedra $P$ with vertices. By linearity, we must have $\mathcal{F}([P + u]) = \mathbf{x}^u \mathcal{F}([P])$ for all rational polyhedra $P$. However, if $P$ contains a line then there is a vector $u \in \mathbb{Z}^d \setminus \{0\}$ such that $P + u = P$. Therefore, we must have $\mathcal{F}([P]) = 0$ for $P$ with a line. $\qquad \square$

Theorem 2.4 provides a powerful tool for computing the generating function of the set of integer points in a rational polyhedron. The following "duality trick" going back to the seminal paper of M. Brion [11] turns out to be particularly useful.

Let $\langle \cdot, \cdot \rangle$ be the standard scalar product in $\mathbb{R}^d$ and let $K \subset \mathbb{R}^d$ be a cone. The cone

$$K^* = \left\{ x \in \mathbb{R}^d : \quad \langle x, y \rangle \geq 0 \quad \text{for all} \quad y \in K \right\}$$

is called the *dual* to $K$. It is easy to see that if $K$ is rational (resp. unimodular) cone then $K^*$ is a rational (resp. unimodular) cone, and that if $K$ contains a line (resp. lies in a proper subspace of $\mathbb{R}^d$) then $K^*$ lies in a proper subspace of $\mathbb{R}^d$ (resp. contains a line). A standard duality argument implies that $(K^*)^* = K$ for

closed convex cones $K$. A less obvious observation is that duality preserves linear relations among indicators of closed convex cones:

$$\sum_{i \in I} \alpha_i [K_i] = 0 \qquad \text{implies} \qquad \sum_{i \in I} \alpha_i [K_i^*] = 0,$$

see, for example, Section IV.1 of [5] for a proof.

Now, to compute the generating function $f(K, \mathbf{x})$ one can do the following. First, we compute the dual cone $K^*$, and, iterating (2.4), we compute unimodular cones $K_i$ and numbers $\epsilon_i \in \{-1, 1\}$ such that

$$[K^*] \equiv \sum_{i \in I} \epsilon_i [K_i] \quad \text{modulo indicators of lower-dimensional cones.}$$

Then, dualizing again, we get

$$[K] \equiv \sum_{i \in I} \epsilon_i [K_i^*] \quad \text{modulo indicators of cones with lines.} \qquad (2.10)$$

In view of Theorem 2.4, cones with lines can be ignored as far as generating functions are concerned. This gives us

$$f(K, \mathbf{x}) = \sum_{i \in I} \epsilon_i f(K_i^*, \mathbf{x}).$$

Since $K_i^*$ are unimodular cones, this completes computation of $f(K, \mathbf{x})$. This trick allows us to reduce the complexity of the algorithm in Theorem 2.2 from $\mathcal{L}^{O(d^2)}$ to $\mathcal{L}^{O(d)}$, where $\mathcal{L}$ is the size of the input.

Another important identity is Brion's Theorem [11], which expresses the generating function of the set of integer points in $P$ as the sum of generating functions for the sets of integer points in the tangent (supporting) cones at the vertices of $P$. Namely, for a vertex $v$ of a polyhedron $P$ let us define the tangent cone $K_v$ as

$$K_v = \Big\{ x : \quad \epsilon x + (1 - \epsilon) v \in P \quad \text{for all sufficiently small} \quad \epsilon > 0 \Big\}.$$

We note that $K_v$ is not a cone per se but rather a translation of the cone $K_v - v$.

**Theorem 2.5** (Brion's Theorem). *For a rational polyhedron $P$ we have*

$$f(P, \mathbf{x}) = \sum_v f(K_v, \mathbf{x}),$$

*where the sum is taken over all vertices of $P$ and the identity is understood as the identity among rational functions.*

Discovered by M. Brion [11], Theorem 2.5 started an avalanche of research. The original proof of Theorem 2.5 was based on algebro-geometric methods. Later,

elementary proofs were discovered in [23] and [25]. One can deduce Theorem 2.5 from Theorem 2.4 and an elementary identity

$$[P] \equiv \sum_v [K_v] \quad \text{modulo indicators of polyhedra with lines,}$$

cf. Section VIII.4 of [5].

Theorem 2.5 together with the unimodular decomposition of Section 2.1 and the duality trick provide the proof of Theorem 2.2 as stated. Another advantage of using Theorem 2.5 is that it allows us to understand how the generating function $f(P, \mathbf{x})$ changes as the facets of $P$ move parallel to themselves so that the combinatorial structure of $P$ does not change. In this case, the tangent cones $K_v$ get translated by vectors linearly depending on the displacements of the facets of $P$. Writing $K_v$ as combinations of translated unimodular cones $K_i + v$ as in (2.10), we notice that as far as lattice points are concerned, a *rational* translation $K_i + v$ of a *unimodular* cone $K_i$ is equivalent to a certain *integer* translation $K_i + u$:

$$\left(K_i + v\right) \cap \mathbb{Z}^d = \left(K_i + u\right) \cap \mathbb{Z}^d \quad \text{for some} \quad u \in \mathbb{Z}^d$$

and hence we have

$$f\left(K_i + v, \mathbf{x}\right) = f\left(K_i + u, \mathbf{x}\right) = \mathbf{x}^u f\left(K_i, \mathbf{x}\right).$$

If $K = \mathbb{R}_+^d$ then $u$ is obtained from $v$ by rounding up the coordinates to the nearest integer. The case of a general unimodular cone differs by a unimodular linear transformation, see [4] for details.

**2.3. Implementation.** The algorithm of Theorem 2.2 appears to be practical. First, it was implemented by J. De Loera et al. [12], who wrote the `LattE` (Lattice point enumeration) software package. The authors of `LattE` discovered that often the most practically efficient way to handle computations is to represent a polyhedron $P$ as a hyperplane section of a higher-dimensional cone as in (2.5) and then use the "dualized" decomposition (2.10). The package allows one to compute the number of integer points in a given rational polytope. Formally speaking, to compute the number $|P \cap \mathbb{Z}^d|$ of integer points in a given rational polytope $P$, we should substitute $\mathbf{x} = (1, \ldots, 1)$ into the rational function $f(P, \mathbf{x})$. However, we need to be careful since this particular value is a pole of every fraction in (2.3). Nevertheless, the substitution can be done efficiently, see Section 3.1 and [3], [4], [7], and [12] for details.

In addition, `LattE` allows one to compute the Ehrhart (quasi)-polynomial of a given rational polytope $P$, that is, to find a formula for the number of integer points in the dilated polytope $nP$, where $n$ is a positive integer, see also Section 6.1.

Testing whether a given rational polyhedron $P$ contains an integer point, or, equivalently, whether $f(P, \mathbf{x}) \not\equiv 0$ is a non-trivial problem related to the general *integer programming problem* of optimizing a given linear function on the set $P \cap \mathbb{Z}^d$. `LattE` package contains also an implementation of an integer programming algorithm based on rational functions $f(P, \mathbf{x})$.

Another implementation, called `barvinok`, was written by S. Verdoolaege, see [36]. Among other features, the implementation allows one to obtain closed explicit formulas for the number of integer points in a parametric polytope as a function of displacement parameters when the facets of the polytope move parallel to themselves, see Theorem 2.5 and the subsequent discussion.

There is an extensive literature devoted to the lattice point enumeration in polytopes, whether from algorithmic, structural, or application points of view. For the classical Ehrhart theory in the context of enumerative combinatorics, see [34] and [9] for a clever simplification of the proofs of the main results of the theory. For an approach featuring Dedekind sums and other analytic tools, see [8]. It does not seem to be possible to survey all the literature in the paper. In addition to already mentioned papers, we provide only a few references among many good papers which appeared after the survey [4].

Efficient counting in special situations with applications to computational questions in representation theory and network flows is discussed in [2]. For a recent advance connecting lattice point counting with algebraic geometry, see [29]. For a computationally efficient version of the Euler-Maclaurin formula, satisfying, in addition, some natural "local" conditions, see [10].

## 3.  Operations on Sets and Generating Functions

Motivated in part by Theorem 2.2, let us consider sets $S \subset \mathbb{Z}^d$ *defined* by their generating functions

$$f(S; \mathbf{x}) = \sum_{m \in S} \mathbf{x}^m$$

written as rational functions in the form

$$f(S; \mathbf{x}) = \sum_{i \in I} \epsilon_i \frac{\mathbf{x}^{a_i}}{(1 - \mathbf{x}^{b_{i1}}) \dots (1 - \mathbf{x}^{b_{ik}})}. \tag{3.1}$$

Here $I$ is a finite set of indices, $\epsilon_i \in \mathbb{Q}$, $a_i, b_{ij} \in \mathbb{Z}^d$, and $b_{ij} \neq 0$ for all $i, j$. To avoid ambiguity, we assume that either $S$ is finite, or, if $S$ is infinite, then there is a non-empty open set $U \subset \mathbb{C}^d$ such that the series defining $f(S; \mathbf{x})$ converges absolutely and uniformly on compact subsets of $U$ and for every fraction in (3.1) there is the Laurent series (multiple geometric series) expansion

$$\frac{\mathbf{x}^{a_i}}{(1 - \mathbf{x}^{b_{i1}}) \dots (1 - \mathbf{x}^{b_{ik}})} = \sum_{(\mu_1, \dots, \mu_k) \in \mathbb{Z}_+^k} \mathbf{x}^{a_i + \mu_1 b_{i1} + \dots + \mu_k b_{ik}}$$

in $U$.

To indicate the computational complexity level of our set $S$, we consider the two parameters *fixed* in formula (3.1): the number $d$ of variables and the number $k$ of binomials in the denominator of each fraction. Note that if we happen to have a smaller number of binomials in some fraction, we can formally "pad" it to $k$ by

multiplying both the numerator and denominator of the fraction by some artificial binomials. Since $k$ is fixed, that would increase the length of the formula by a constant factor.

Next, we discuss what information about the set $S$ can be extracted from $f(S; \mathbf{x})$ given in the form of (3.1).

### 3.1. Monomial substitutions and differentiation. 
One piece of information we can get is the cardinality $|S|$ of a finite set $S$. To compute $|S|$, we would like to substitute $\mathbf{x} = (1, \ldots, 1)$ in (3.1), but this should be done carefully since this particular value of $\mathbf{x}$ is the pole of every single fraction in (3.1). The procedure is introduced in [3].

We choose a sufficiently generic vector $c \in \mathbb{Z}^d$, $c = (\gamma_1, \ldots, \gamma_d)$, so that $\langle c, b_{ij} \rangle \neq 0$ for all $i, j$. For a $\tau \in \mathbb{C}$, let

$$\mathbf{x}(\tau) = \left( e^{\tau \gamma_1}, \ldots, e^{\tau \gamma_d} \right).$$

Thus we want to compute

$$\lim_{\tau \longrightarrow 0} f\big(S; \mathbf{x}(\tau)\big).$$

Let us compute

$$\alpha_i = \langle c, a_i \rangle \quad \text{and} \quad \beta_{ij} = \langle c, b_{ij} \rangle.$$

Then

$$f\big(S; \mathbf{x}(\tau)\big) = \sum_{i \in I} \epsilon_i \frac{e^{\alpha_i \tau}}{(1 - e^{\beta_{i1}\tau}) \ldots (1 - e^{\beta_{ik}\tau})}. \tag{3.2}$$

Next, we note that $f\big(S; \mathbf{x}(\tau)\big)$ is a meromorphic function in $\tau$ and that we want to compute the constant term of its Laurent expansion in the neighborhood of $\tau = 0$. To do that, we deal with every fraction separately. We write each fraction of (3.2) as

$$\frac{e^{\alpha_i \tau}}{(1 - e^{\beta_{i1}\tau}) \ldots (1 - e^{\beta_{ik}\tau})} = \tau^{-k} e^{\alpha_i \tau} \prod_{j=1}^{k} g_{ij}(\tau), \qquad \text{where}$$

$$g_{ij}(\tau) = \frac{\tau}{1 - e^{\beta_{ij}\tau}}.$$

Now, each $g_{ij}(\tau)$ is an analytic function of $\tau$ and we compute is Taylor series expansion $p_{ij}(\tau)$ up to the $\tau^{k+1}$ term:

$$\frac{\tau}{1 - e^{\beta_{ij}\tau}} \equiv p_{ij}(\tau) \mod \tau^{k+1}.$$

Similarly, we compute a polynomial $q_i(\tau)$ such that

$$e^{\alpha_i \tau} \equiv q_i(\tau) \mod \tau^{k+1}.$$

Finally, successively multiplying polynomials $\mod \tau^{k+1}$ we compute the polynomial $h_i(\tau)$ with $\deg h_i \leq k$ such that

$$q_i p_{i1} \cdots p_{ik} \equiv h_i \mod \tau^{k+1}.$$

Letting

$$h(\tau) = \sum_{i \in I} h_i(\tau),$$

we conclude that the coefficient of $\tau^k$ in $h(\tau)$ is the desired value of (3.2) at $\tau = 0$ and hence is the value $f(S; \mathbf{x})$ at $\mathbf{x} = (1, \ldots, 1)$. We note that the procedure has a polynomial time complexity even if both $k$ and $d$ are allowed to vary and if we allow different numbers $k_i \leq k$ of binomials in different fractions of (3.1).

A more general operation which can be computed in polynomial time is that of a *monomial substitution*. Let $f(\mathbf{x})$ be an expression of the type (3.1). Let $\mathbf{z} = (z_1, \ldots, z_n)$ be a new set of variables, let $l_1, \ldots, l_d \in \mathbb{Z}^n$ be vectors, and let $\phi : \mathbb{C}^n \longrightarrow \mathbb{C}^d$ be the transformation defined by

$$(z_1, \ldots, z_n) \longmapsto (x_1, \ldots, x_d) \quad \text{where} \quad x_i = \mathbf{z}^{l_i}.$$

If the image $\phi(\mathbb{C}^n)$ does not lie in the set of poles of $f$, one can define a rational function $g(\mathbf{z}) = f(\phi(\mathbf{z}))$. Function $g$ can be computed in polynomial time in the form

$$g(\mathbf{z}) = \sum_{i \in I'} \delta_i \frac{\mathbf{z}^{q_i}}{\left(1 - \mathbf{z}^{b_{i1}}\right) \cdots \left(1 - \mathbf{z}^{b_{ik_i}}\right)},$$

where $\delta_i \in \mathbb{Q}$, $q_i, b_{ij} \in \mathbb{Z}^n$, $b_{ij} \neq 0$ for all $i, j$ and $k_i \leq k$ for all $i \in I'$.

The case of $l_1 = \ldots = l_d = 0$ corresponds to the case of $\mathbf{x} = (1, \ldots, 1)$ considered above. As above, the general case of a monomial substitution is handled by a one-parametric perturbation and computation with univariate polynomials. Details can be found in [7] (the assumption that $k$ is fixed in advance is not needed there).

The operation of monomial substitution has the following geometric interpretation. Let $T : \mathbb{R}^d \longrightarrow \mathbb{R}^n$ be the linear transformation whose matrix in the standard bases consists of the integer column vectors $l_1, \ldots, l_d$. Let $S \subset \mathbb{Z}^d$ be a set and suppose that for all $m \in T(S)$ the set $T^{-1}(m) \cap S$ is finite. The monomial substitution $x_i = \mathbf{z}^{l_i}$ into the generating function $f(S; \mathbf{x})$ produces the weighted generating function $g(\mathbf{z})$ of the image $T(S) \subset \mathbb{Z}^n$, where each monomial $\mathbf{z}^m$ for $m \in T(S)$ is counted with multiplicity $|T^{-1}(m) \cap S|$.

Another useful operation is that of differentiation. Let $p$ be a $d$-variate polynomial. We can write

$$\sum_{m \in S} p(m) \mathbf{x}^m = p\left(x_1 \frac{\partial}{\partial x_1}, \ldots, x_d \frac{\partial}{\partial x_d}\right) f(S; \mathbf{x}).$$

As long as $k$ is fixed in advance, the result can be computed in polynomial time in the form

$$\sum_{i \in I'} \delta_i \frac{\mathbf{x}^{q_i}}{\left(1 - \mathbf{x}^{b_{i1}}\right)^{\gamma_{i1}} \cdots \left(1 - \mathbf{x}^{b_{ik}}\right)^{\gamma_{ik}}},$$

where $\delta_i \in \mathbb{Q}$, $a_i, b_{ij} \in \mathbb{Z}^d$, $b_{ij} \neq 0$, and $\gamma_{ij}$ are non-negative integers such that $\gamma_{i1} + \ldots + \gamma_{ik} \leq k + \deg p$ for all $i$, see [6].

This observation is used in [6], see also [10] and [13].

One corollary of Theorem 2.2 is that we can efficiently perform set-theoretic operations (intersection, union, difference) of finite sets defined by (3.1). The following result is proved in [7].

**Theorem 3.1.** *Let us fix positive integers $d$ and $k$. Then there exists a polynomial time algorithm, which, for any two finite sets $S_1, S_2 \subset \mathbb{Z}^d$ given by their rational generating functions*

$$f(S_1; \mathbf{x}) = \sum_{i \in I_1} \alpha_i \frac{\mathbf{x}^{p_i}}{(1 - \mathbf{x}^{a_{i1}}) \ldots (1 - \mathbf{x}^{a_{ik}})} \tag{3.3}$$

*and*

$$f(S_2; \mathbf{x}) = \sum_{i \in I_2} \beta_i \frac{\mathbf{x}^{q_i}}{(1 - \mathbf{x}^{b_{i1}}) \ldots (1 - \mathbf{x}^{b_{ik}})} \tag{3.4}$$

*computes the generating function $f(S; \mathbf{x})$ of their intersection $S = S_1 \cap S_2$ in the form*

$$f(S; \mathbf{x}) = \sum_{i \in I} \gamma_i \frac{\mathbf{x}^{u_i}}{(1 - \mathbf{x}^{v_{i1}}) \ldots (1 - \mathbf{x}^{v_{is}})},$$

*where $s \leq 2k$.*

*Proof.* The idea of the proof is to *linearize* the operation of intersection of sets. Suppose we have two Laurent series

$$g_1(\mathbf{x}) = \sum_{m \in \mathbb{Z}^d} \rho_{1m} \mathbf{x}^m \quad \text{and} \quad g_2(\mathbf{x}) = \sum_{m \in \mathbb{Z}^d} \rho_{2m} \mathbf{x}^m.$$

Let us define their *Hadamard product* $g_1(\mathbf{x}) \star g_2(\mathbf{x})$ as

$$g(\mathbf{x}) = \sum_{m \in \mathbb{Z}^d} \rho_m \mathbf{x}^m \quad \text{where} \quad \rho_m = \rho_{1m} \rho_{2m}.$$

Then, clearly,

$$f(S_1 \cap S_2; \mathbf{x}) = f(S_1; \mathbf{x}) \star f(S_2; \mathbf{x}).$$

Without loss of generality, we assume that there is a non-empty open set $U \subset \mathbb{C}^d$ such that for all $\mathbf{x} \in U$ and every fraction of (3.3) and (3.4) we have the multiple geometric series expansions:

$$\frac{\mathbf{x}^{p_i}}{(1 - \mathbf{x}^{a_{i1}}) \ldots (1 - \mathbf{x}^{a_{ik}})} = \sum_{(\mu_1, \ldots, \mu_k) \in \mathbb{Z}_+^k} \mathbf{x}^{p_i + \mu_1 a_{i1} + \ldots + \mu_k a_{ik}} \tag{3.5}$$

and

$$\frac{\mathbf{x}^{q_i}}{(1 - \mathbf{x}^{b_{i1}}) \ldots (1 - \mathbf{x}^{b_{ik}})} = \sum_{(\nu_1, \ldots, \nu_k) \in \mathbb{Z}_+^k} \mathbf{x}^{q_i + \nu_1 b_{i1} + \ldots + \nu_k b_{ik}}. \tag{3.6}$$

As usual, we assume that for all $\mathbf{x} \in U$ the convergence in (3.5) and (3.6) is absolute and uniform on all compact subsets of $U$. To ensure that such a set $U$

indeed exists, we choose a sufficiently generic linear function $\ell : \mathbb{R}^d \longrightarrow \mathbb{R}$ and make sure that $\ell\left(a_{ij}\right), \ell\left(b_{ij}\right) > 0$ for all $i, j$ by reversing, if necessary, the direction of $a_{ij}$ and $b_{ij}$ via the identity

$$\frac{1}{1 - \mathbf{x}^a} = -\frac{\mathbf{x}^{-a}}{1 - \mathbf{x}^{-a}}.$$

Here we use that $S_1$ and $S_2$ are finite so that $f(S_1; \mathbf{x})$ and $f(S_2; \mathbf{x})$ are, in fact, Laurent polynomials.

Since the Hadamard product is a bilinear operation on series, to compute $f(S_1; \mathbf{x}) \star f(S_2; \mathbf{x})$ it suffices to compute the Hadamard product of every pair of series (3.5) and (3.6).

In the space $\mathbb{R}^{2k}$ of two $k$-tuples $(x, y)$, where $x = (\xi_1, \ldots, \xi_k)$ and $y = (\eta_1, \ldots, \eta_k)$, let us introduce the polyhedron

$$Q_i = \left\{ (x, y) : \begin{array}{l} \xi_1, \ldots, \xi_k; \eta_1, \ldots, \eta_k \geq 0 \\ p_i + \xi_1 a_{i1} + \ldots + \xi_k a_{ik} = q_i + \eta_1 b_{i1} + \ldots + \eta_k b_{ik} \end{array} \right\} \quad (3.7)$$

and let $\mathbb{Z}^{2k} \subset \mathbb{R}^{2k}$ be the standard integer lattice.

Since the Hadamard product is bilinear and for monomials we have

$$\mathbf{x}^{m_1} \star \mathbf{x}^{m_2} = \begin{cases} \mathbf{x}^m & \text{if } m_1 = m_2 = m \\ 0 & \text{if } m_1 \neq m_2, \end{cases}$$

the Hadamard product of the series (3.5) and (3.6) can be expressed as the sum

$$\sum_{(m,n)\in Q_i \cap \mathbb{Z}^{2k}} \mathbf{x}^{p_i + \mu_1 a_{i1} + \ldots + \mu_k a_{ik}}, \qquad \text{where} \tag{3.8}$$
$$m = (\mu_1, \ldots, \mu_k) \quad \text{and} \quad n = (\nu_1, \ldots, \nu_k).$$

On the other hand, (3.8) is obtained from the generating function $f(Q_i, \mathbf{z})$ with $\mathbf{z} = (z_1, \ldots, z_{2k})$ by the monomial substitution

$$z_i = \mathbf{x}^{a_i} \quad \text{for} \quad i = 1, \ldots, k \qquad \text{and} \qquad z_i = 1 \quad \text{for} \quad i = k+1, \ldots, 2k \tag{3.9}$$

and multiplication by $\mathbf{x}^{p_i}$.

We use Theorem 2.2 to compute $f(Q_i, \mathbf{z})$. The monomial substitution (3.9) can also be computed in polynomial time, cf. Section 3.1. $\qquad \square$

Therefore, one can compute the generating functions of the union and difference:

$$f\left(S_1 \cup S_2; \mathbf{x}\right) = f(S_1; \mathbf{x}) + f(S_2; \mathbf{x}) - f\left(S_1 \cap S_2; \mathbf{x}\right) \quad \text{and}$$
$$f\left(S_1 \setminus S_2; \mathbf{x}\right) = f(S_1; \mathbf{x}) - f\left(S_1 \cap S_2; \mathbf{x}\right).$$

Theorem 3.1 allows us to work with generating functions (3.1) directly as with data structures bypassing any more explicit descriptions of sets $S$ in question. Of

course, there is a price to pay: with every set-theoretic operation, the complexity level of the set, the number $k$ of binomials in the denominator of each fraction in (3.1), doubles. From the definition (3.7) of $Q_i$ we can notice that in a sufficiently general position we will have $\dim Q_i = 2k - d$, so we would be able to choose $s = 2k - d$ in Theorem 3.1. Theorem 3.1 admits an extension to infinite sets $S_1$ and $S_2$ provided there is a non-empty open set $U \subset \mathbb{C}^d$ such that the multiple geometric series expansions (3.5) and (3.6) hold for all fractions in (3.3) and (3.4). K. Woods [38] used the construction of the Hadamard product to show that in any fixed dimension there is a polynomial time algorithm to check if a given integer is a period of the Ehrhart quasi-polynomial of a given rational polytope.

# 4. Beyond Polyhedra: Projections

There are other interesting sets admitting short rational generating functions (3.1). We start with examples.

**4.1. Integer semigroups.** Let $S$ be the semigroup generated by positive coprime integers $a_1$ and $a_2$, that is, the set of all non-negative integer combinations of $a_1$ and $a_2$:

$$S = \left\{ \mu_1 a_1 + \mu_2 a_2 : \quad \mu_1, \mu_2 \in \mathbb{Z}_+ \right\}.$$

It is not hard to show that

$$f(S; x) = \frac{1 - x^{a_1 a_2}}{(1 - x^{a_1})(1 - x^{a_2})}$$

(the series defining $f(S; x)$ converges for all $|x| < 1$).

Let $S$ be the semigroup generated by positive coprime integers $a_1, a_2$, and $a_3$,

$$S = \left\{ \mu_1 a_1 + \mu_2 a_2 + \mu_3 a_3 : \quad \mu_1, \mu_2, \mu_3 \in \mathbb{Z}_+ \right\}.$$

Then there exist positive integers $p_1, p_2, p_3, p_4$, and $p_5$, not necessarily distinct, such that

$$f(S; x) = \frac{1 - x^{p_1} - x^{p_2} - x^{p_3} + x^{p_4} + x^{p_5}}{(1 - x^{a_1})(1 - x^{a_2})(1 - x^{a_3})}.$$

This interesting result was rediscovered a number of times. It was explicitly stated by M. Morales [27]; the proof wasn't published though. Independently, the proof was rediscovered by G. Denham [14]. Both proofs are algebraic and based on the interpretation of $f(S; x)$ as the Hilbert series of a graded ring $\mathbb{C}[t^{a_1}, t^{a_2}, t^{a_3}]$. In this special case (a Cohen-Macaulay ring of codimension 2), the Hilbert series can be computed via the Hilbert-Burch Theorem, cf. also [18]. Meanwhile, a combinatorial proof of a somewhat weaker result (up to 12 monomials in the numerator) independently appeared in [35].

The pattern breaks down for semigroups with $d \geq 4$ generators, meaning that if we choose the denominator of $f(S; x)$ in the form $(1 - x^{a_1}) \cdots (1 - x^{a_d})$, the

number of monomials in the numerator does not remain constant for a particular value of $d$, and, moreover, grows exponentially with the input size of $a_1, \ldots, a_d$. As shown in [35], for $d = 4$ the number of the monomials in the numerator can grow as fast as $\min^{1/2}\{a_1, a_2, a_3, a_4\}$, whereas the input size is only about $\log(a_1 a_2 a_3 a_4)$.

Nevertheless, the generating function $f(S; x)$ admits a short rational function representation for any number $d$ of generators fixed in advance. The following result was proved in [7].

**Theorem 4.1.** *Let us fix $d$. Then there exists a positive integer $s = s(d)$ and a polynomial time algorithm, which, given positive integers $a_1, \ldots, a_d$, computes the generating function $f(S; x)$ of the semigroup*

$$S = \left\{ \sum_{i=1}^{d} \mu_i a_i : \qquad \mu_1, \ldots, \mu_d \in \mathbb{Z}_+ \right\}$$

*generated by $a_1, \ldots, a_d$ in the form*

$$f(S; x) = \sum_{i \in I} \alpha_i \frac{x^{p_i}}{(1 - x^{b_{i1}}) \cdots (1 - x^{b_{is}})}, \tag{4.1}$$

*where $\alpha_i \in \mathbb{Q}$, $p_i, b_{ij} \in \mathbb{Z}$ and $b_{ij} \neq 0$ for all $i, j$.*

In particular, for any fixed $d$, the number $|I|$ of fractions in (4.1) is bounded by a polynomial in the input size, that is, in $\log(a_1 \cdots a_d)$.

Theorem 4.1 is obtained as a corollary of a more general result that the *projection* of the set of integer points in a rational polytope of a fixed dimension admits a short rational generating function [7].

**Theorem 4.2.** *Let us fix $d$. Then there exists a number $s = s(d)$ and a polynomial time algorithm, which, given a rational polytope $P$ and a linear transformation $T : \mathbb{R}^d \longrightarrow \mathbb{R}^k$ such that $T(\mathbb{Z}^d) \subset \mathbb{Z}^k$, computes the generating function $f(S; \mathbf{x})$ for $S = T(P \cap \mathbb{Z}^d)$, $S \subset \mathbb{Z}^k$, in the form*

$$f(S; \mathbf{x}) = \sum_{i \in I} \frac{\mathbf{x}^{p_i}}{(1 - \mathbf{x}^{b_{i1}}) \cdots (1 - \mathbf{x}^{b_{is}})}, \tag{4.2}$$

*where $\alpha_i \in \mathbb{Q}$, $p_i, b_{ij} \in \mathbb{Z}^k$ and $b_{ij} \neq 0$ for all $i, j$.*

One can observe that Theorem 4.1 is a corollary of Theorem 4.2. Indeed, let $T : \mathbb{R}^d \longrightarrow \mathbb{R}$ be the linear transformation defined by

$$T(\xi_1, \ldots, \xi_d) = a_1 \xi_1 + \ldots + a_d \xi_d.$$

Then the semigroup $S$ generated by $a_1, \ldots, a_d$ is the image $S = T(\mathbb{Z}_+^d)$ of the set $\mathbb{Z}_+^d$ of integer points in the rational polyhedron $\mathbb{R}_+^d \subset \mathbb{R}^d$. The polyhedron $\mathbb{R}_+^d$ is unbounded, so Theorem 4.2 cannot be applied immediately. However, it is not hard to show that $S \subset \mathbb{Z}_+$ stabilizes after a while (if $a_1, \ldots, a_d$ are coprime then

$S$ includes all sufficiently large positive integers). Thus only the initial interval of $S$ is of interest, to get which we replace $\mathbb{R}_+^d$ by a sufficiently large simplex

$$P = \left\{ (\xi_1, \ldots, \xi_d) : \quad \sum_{i=1}^{d} \xi_i \leq t \quad \text{and} \quad \xi_i \geq 0 \quad \text{for} \quad i = 1, \ldots, d \right\},$$

see [7] for details.

We sketch the proof of Theorem 4.2 below.

Without loss of generality we assume that $\dim \ker T = d - k$. The proof then proceeds by induction on $d - k$. If $d = k$ we are in the situation of Theorem 2.2. We note that for any $k$ and $d$, if the restriction $T : P \cap \mathbb{Z}^d \longrightarrow S$ is one-to-one, we can compute the generating function $f(S; \mathbf{x})$ from that of the set $P \cap \mathbb{Z}^d$ using an appropriate monomial substitution, cf. Section 3.1. Otherwise, the monomial substitution will account for each point $m \in S$ with the multiplicity equal to the number of the points in $P \cap \mathbb{Z}^d$ mapped onto $m$. Thus our goal is to eliminate multiplicities.

The case of $d = k + 1$ illuminates some of the ideas used in the proof for an arbitrary $d - k$. Suppose that

$$T : \mathbb{R}^{k+1} \longrightarrow \mathbb{R}^k, \quad (\xi_1, \ldots, \xi_{k+1}) \longmapsto (\xi_1, \ldots, \xi_k)$$

is the projection (this is a sufficiently general case). Let $\widehat{S} = P \cap \mathbb{Z}^{k+1}$ and let us consider the restriction $T : \widehat{S} \longrightarrow S$. Then, for every point $m \in S$, the preimage $T^{-1}(m) \subset \widehat{S}$ is the set of integer points in the interval $T^{-1}(m) \cap P$ which all agree in their first $k$ coordinates and disagree in the last coordinate. Let $e_{k+1}$ be the last basis vector and let us consider

$$Y = \widehat{S} \setminus \left( \widehat{S} + e_{k+1} \right).$$

In words: we subtract from $\widehat{S}$ its translation by 1 in the last coordinate.

Then the restriction $T : Y \longrightarrow S$ is one-to-one since the preimage $T^{-1}(m) \subset Y$ consists of the single point in $T^{-1}(m) \subset \widehat{S}$ with the smallest last coordinate. Now, $\widehat{S}$ is the set of integer points in a rational polytope and we compute its generating function using Theorem 2.2. Then we compute the generating function of $Y$ using Theorem 3.1. Finally, we obtain $f(S; \mathbf{x})$ by substituting $x_{k+1} = 1$ in the generating function $f\big(Y; (\mathbf{x}, x_{k+1})\big)$, cf. Section 3.1.

Let us consider the case of general $k$ and $d$. Let $pr : \mathbb{Z}^{k+1} \longrightarrow \mathbb{Z}^k$ be the natural projection, $pr(\mu_1, \ldots, \mu_{k+1}) = (\mu_1, \ldots, \mu_k)$. Let $\widehat{T} : \mathbb{Z}^d \longrightarrow \mathbb{Z}^{k+1}$ be a linear transformation which is a lifting of $T$ so that $pr\left(\widehat{T}(m)\right) = T(m)$ for all $m \in \mathbb{Z}^d$. We define $\widehat{S} = \widehat{T}(S)$, $\widehat{S} \subset \mathbb{Z}^{k+1}$, and consider the restriction

$$pr : \widehat{S} \longrightarrow S.$$

For every $m \in S$ the preimage $pr^{-1}(m) \subset \widehat{S}$ consists of the points which differ in their last coordinate only. Suppose that we managed to construct $\widehat{T}$ in such a way

that the set $pr^{-1}(m) \subset \widehat{S}$ has *small gaps*, meaning that there exists a constant $l = l(d)$ such that if there are two points in $pr^{-1}(m)$ whose $(k+1)$st coordinates differ by more than $l$, there must be a point in $pr^{-1}(m)$ lying strictly between them.

In this case, we compute $f(S; \mathbf{x})$ as follows. Let us define

$$Y = \widehat{S} \setminus \bigcup_{j=1}^{l} \left( \widehat{S} + je_{k+1} \right).$$

In words: we subtract from $\widehat{S}$ its $l$ translates by $1, \ldots, l$ in the last coordinate. Because of the small gap property, the restriction $pr : Y \longrightarrow S$ is one-to-one: now, the preimage $pr^{-1}(m) \subset Y$ consists of the single point in $pr^{-1}(m) \subset \widehat{S}$ with the smallest last coordinate. Using the induction hypothesis, we compute the generating function of $\widehat{S}$. Then, applying Theorem 3.1 $l$ times, we compute the generating function of $Y$. Finally, $f(S; \mathbf{x})$ is obtained from $f\big(Y; (\mathbf{x}, x_{k+1})\big)$ by the substitution $x_{k+1} = 1$, see Section 3.1.

In general, we cannot construct a lifting $\widehat{T}$ with the small gap property but the next best thing is possible. Namely, we can construct in polynomial time a decomposition $\mathbb{R}^k = \bigcup_i Q_i$ of $\mathbb{R}^k$ into a union of non-overlapping rational polyhedra $Q_i$ such that for each piece $S_i = S \cap Q_i$ a lifting $\widehat{T}_i$ with the small gap property indeed exists. The generating functions $f(S_i; \mathbf{x})$ are computed as above and then patched together into a single generating function $f(S; \mathbf{x})$. The construction of such polyhedra $Q_i$ and liftings $\widehat{T}_i$ is based on the results of [21] and [20]. The main tool is the following *Flatness Theorem*, see, for example, Section 6.7 of [17] or Section VII.8 of [5].

**Theorem 4.3** (Flatness Theorem). *For each dimension $d$ there exists a constant $\omega(d)$ with the following property: if $V$ is a $d$-dimensional real vector space, $\Lambda \subset V$ is a lattice of rank $d$, $\Lambda^* \subset V^*$ is the reciprocal lattice, and $K \subset V$ is a convex compact set with non-empty interior such that $K \cap \Lambda = \emptyset$ then there is an $\ell \in \Lambda^* \setminus \{0\}$ such that*

$$\max_{x \in K} \ell(x) - \min_{x \in K} \ell(x) \leq \omega(d). \tag{4.3}$$

In words: a lattice-free convex body is flat in some lattice direction. The number in the left hand side of (4.3) is called the *width of $K$ with respect to $\ell$* and denoted width$(K, \ell)$. The infimum of width$(K, \ell)$ over all $\ell \in \Lambda^*$ is called the *lattice width of $K$* and denoted width$(K)$. A simple and crucial observation relating the lattice width and the small gap property is that if for $\ell \in \Lambda^*$ we have width$(K, \ell) \leq \gamma$ width$(K)$ then the gaps between the consecutive integers in the set $\ell(K \cap \Lambda)$ do not exceed $\gamma \omega(d)$.

We go back to finish the sketch of the proof of Theorem 4.2. Let $\Lambda = \mathbb{Z}^k \cap \ker(T)$ be the lattice in $\ker(T)$. For $y \in \mathbb{R}^d$, let $P_y = P \cap T^{-1}(x)$ be the fiber of the polytope $P$ over $x$. We will measure the lattice width of $P_y$ with respect to $\Lambda$. The results of [21] and [20] allow us to construct a polyhedral decomposition $\mathbb{R}^k = \bigcup_i Q_i$ and vectors $\ell_i \in \Lambda^*$ such that for all $y \in Q_i$ we have either width $(P_y, \ell_i) \leq 2$ width$(P_y)$

or width$(P_y, \ell_i) \leq 1$. We then define

$$\widehat{T}_i(x) = \big(T(x), \ell_i(x)\big) \quad \text{if} \quad T(x) \in Q_i.$$

This completes the sketch of proof of Theorem 4.2.

**4.2. Applications.** Theorem 4.1 implies polynomial time solvability of a variety of problems about integer semigroups. Suppose that the generators $a_1, \ldots, a_d$ are coprime. As is known, all sufficiently large integers lie in the semigroup $S$ generated by $a_1, \ldots, a_d$. In the situation when the number $d$ of generators is fixed, R. Kannan [20] constructed a polynomial time algorithm to compute the largest integer not in $S$. Theorem 4.1 implies that one can compute in polynomial time the number of positive integers not in $S$, the number of integers in $S$ belonging to a particular interval, etc.

Unlike the algorithm of Theorem 2.2, the algorithms of Theorems 4.1 and 4.2 seem to be unimplementable at the moment. Indeed, the way Theorem 4.2 is proved gives $s = d^{\Omega(d)}$ at best and, similarly, in Theorem 4.1. It is not clear at the moment whether a smaller value of $s$ is possible.

In Theorem 4.1, apart from $d = 1, 2, 3$, the value of $d = 4$ seems to indicate a possibility of a "special treatment". The approach of [33] combined with the continued fraction method, see Section 2.1, may lead to a practically efficient algorithm to compute $f(S; x)$.

Theorem 4.2 implies that some other interesting sets admit short rational generating functions. One class of such sets consists of the Hilbert bases of rational cones. Let $K \subset \mathbb{R}^d$ be a pointed rational cone. The set $S \subset K \cap \mathbb{Z}^d$, $0 \notin S$, is called the (minimal) *Hilbert basis* of the semigroup $K \cap \mathbb{Z}^d$ if every point in $K \cap \mathbb{Z}^d$ can be represented as a sum of some points in $S$ and if no point in $S$ is a sum of other points in $S$. In other words, $S$ consists of the points in $K \cap \mathbb{Z}^d$ that cannot be written as a sum of non-zero points in $K \cap \mathbb{Z}^d$. Theorem 4.2 implies that as long as the dimension $d$ remains fixed, given a rational cone $K$, the generating function $f(S; \mathbf{x})$ can be computed in polynomial time as a short rational function of the type (3.1). Consequently, the number $|S|$ of points in the Hilbert basis of $K \cap \mathbb{Z}^d$ can be computed in polynomial time.

To deduce this result from Theorem 4.2, let $Q \subset K$ be a rational polyhedron containing all integer points in $K$ except 0 (to get $Q$ from $K$, we cut the vertex of $K$ by a hyperplane), let $P = Q \times Q \subset \mathbb{R}^d \oplus \mathbb{R}^d = \mathbb{R}^{2d}$ and let $T$ be the projection $P \longmapsto K$, $T(x, y) = x + y$. Then the Hilbert basis $S$ is the complement in $Q \cap \mathbb{Z}^d$ of the image $T\big(P \cap \mathbb{Z}^{2d}\big)$. The obstacle that the polyhedron $Q$ is not bounded, so Theorem 4.2 cannot be applied immediately, can be easily fixed since only the "initial part" of the semigroup $K \cap \mathbb{Z}^d$ is of interest, see [7].

Another class of sets allowing short rational generating functions via Theorem 4.2 are the *test sets* in integer programming, see [30].

It should be noted that the short rational function description provides only very general characterization of the set. For example, many of the fine properties of test sets [30] do not seem to be picked up by rational generating functions and some empirically observed phenomena are still waiting for their explanation. For

structural results (without complexity estimates) regarding $f(S; \mathbf{x})$, where $S$ is the projection of the set of integer points in a rational polyhedron, see [24].

# 5. Beyond Projections: Presburger Arithmetic

Let us consider formulas we can construct by using integer variables, operations of addition, subtraction, and multiplication by an integer constant (but not multiplication of two integer variables), comparison ($<$, $>$, $=$), Boolean operations ("and", "or", "not"), and quantifiers ($\forall$, $\exists$). The realm of such formulas is *Presburger arithmetic*. Thus the set $P \cap \mathbb{Z}^d$ of integer points in a rational polyhedron can be described by a quantifier-free formula of Presburger arithmetic: the set $P \cap \mathbb{Z}^d$ consists of $d$-tuples of integer variables that satisfy a number of linear constraints with constant integer coefficients. Similarly, the projection $T\left(P \cap \mathbb{Z}^d\right)$ of the set of integer points in a polyhedron is described by a formula of Presburger arithmetic with existential quantifiers only (no quantifier alternations).

With a little work, Theorem 2.2 can be extended as follows. Let us fix the number $d$ of variables. Then there exists a polynomial time algorithm, which, given a quantifier-free formula $F$ of Presburger arithmetic, computes the generating function $f(S; \mathbf{x})$ of the set $S \subset \mathbb{Z}^d$ defined by $F$ as a rational function (2.3). Some routine precautions regarding convergence of the series defining $f(S; \mathbf{x})$, if $S$ is infinite, should be taken. The general case of a set defined by a quantifier-free formula $F$ reduces to that of the set integer points in a rational polyhedron by some more or less straightforward "cutting and pasting" of polyhedra. Since the dimension $d$ of the ambient space is fixed, this cutting and pasting can be performed in polynomial time.

Theorem 4.2 can be extended as follows. Let us fix the number of variables *and* the number of Boolean operations used. Then there exists a polynomial time algorithm, which, given a formula $F$ of Presburger arithmetic without quantifier alternations, computes the generating function $f(S; \mathbf{x})$ of the finite set $S \subset \mathbb{Z}^k$ defined by $F$ as a rational function (4.2). Note that here we have to fix not only the number of variables but also the number of Boolean operations. For example, unless **P=NP** one cannot hope to compute the generating function of the projection of the set of integer points in a union of rational polytopes if the number of polytopes is allowed to vary, cf. Section 5.3 of [37] and [31].

One can ask whether the results can be extended even further. Let us fix the number of variables and the number of Boolean operations, making numerical constants essentially the only parameters of the formula. Is there a polynomial time algorithm which computes the generating function (3.1) of the (finite) set $S$ of points described by such a formula? This indeed seems very plausible, see the discussion in Chapter V of [37]. Intuitively, such sets should have some "hidden periodicity" and short rational generating functions should reveal that periodicity. Besides, it seems hard to prove that a particular finite, but large, set $S \subset \mathbb{Z}^d$ does not admit a short rational generating function: if a particular candidate expression for $f(S; \mathbf{x})$ is not short, one can argue that we haven't searched hard enough and

that there is another, better candidate.

We mention that the result of R. Kannan [19] establishes polynomial time solvability of decision problems for formulas with not more than one quantifier alternation. If the number of variables is not fixed, the complexity of decision problems in Presburger arithmetic is double exponential by the result of M. Fischer and M. Rabin [15].

# 6.  Concluding Remarks

One can ask whether some of the technique discussed in this paper can be extended to lattice points satisfying some non-linear constraints. The answer seems to be "no". For example, lattice points in the standard Euclidean ball exhibit phenomena explained not by rational but rather by theta functions. Let

$$B_n = \left\{ (\xi_1, \xi_2, \xi_3, \xi_4) : \quad \xi_1^2 + \xi_2^2 + \xi_3^2 + \xi_4^2 \leq n \right\}$$

be the Euclidean ball of radius $\sqrt{n}$. Jacobi's formula asserts that the number $|B_n \cap \mathbb{Z}^4| - |B_{n-1} \cap \mathbb{Z}^4|$ of integer points on the sphere of radius $\sqrt{n}$ is equal to

$$8 \sum_{4 \nmid r \mid n} r$$

(in words: eight times the sum of divisors of $n$ that are not divisible by four). One can show then [1] that if one can count points in a 4-dimensional ball efficiently (in polynomial time), one can factor integers efficiently (in randomized polynomial time).

We note also that lattice points in *irrational* polyhedra exhibit a very interesting behavior, see [32].

**6.1.  Large dimensions.**  Almost everywhere in this paper we assumed that the dimension $d$ of the ambient space is fixed in advance. But what if the dimension is allowed to grow? Given a rational polyhedron $P \subset \mathbb{R}^d$, it is an NP-hard problem to determine whether $P \cap \mathbb{Z}^d = \emptyset$ (even when $P$ is a rational simplex). Thus there is little hope to compute the generating function $f(P, \mathbf{x})$ in polynomial time. However, it appears that some interesting "residues" or "shadows" of $f(P, \mathbf{x})$ can be efficiently computed even when the dimension $d$ is allowed to grow, cf. [10] and [6].

The number $e(P) = |P \cap \mathbb{Z}^d|$ of integer points in a rational polyhedron is an example of a lattice invariant *valuation*, see [26]. That is, the map $P \longmapsto e(P)$ extends to a linear functional on the space spanned by the indicators $[P]$ of rational polyhedra, cf. Definition 2.3, and the linear functional is invariant under lattice shifts: $e(P) = e(P + u)$, $u \in \mathbb{Z}^d$. One can ask if there is another lattice invariant valuation $\nu$ on rational polytopes which is efficiently computable in interesting cases and which, in some sense, approximates the counting valuation $e(P)$. For example, the volume $\mathrm{vol}\, P$ may serve as the "0th" approximation to $e(P)$.

With every lattice invariant valuation $\nu$ one can associate the expression

$$\nu(nP) = \sum_{i=0}^{d} \nu_i(P;n) n^i, \tag{6.1}$$

where $nP$ is a dilation of $P$ by an integer factor $n$ and the coefficients $\nu_i(P;n)$ are quasi-periodic: $\nu_i(P; n + t) = \nu_i(P;n)$ provided $tP$ is a polytope with integer vertices, cf. [26]. In the case of the counting valuation $e$, the expression (6.1) is called the *Ehrhart quasi-polynomial* of $P$ and $e_d(P;n) = \text{vol}\,P$. As the $k$th approximation to the counting valuation $e$ we consider a lattice invariant valuation $\nu$ which agrees with $e$ in the $k + 1$ highest terms:

$$\nu_i(P;n) = e_i(P;n) \quad \text{for} \quad i = d, d - 1, \ldots, d - k.$$

A natural goal is to construct such a valuation $\nu$, which is computable in polynomial time (at least, in some interesting cases) for any $k$ fixed in advance.

Abstractly speaking, to define the counting valuation $e$, we have to choose a finite-dimensional real vector space $V$ and a lattice $\Lambda \subset V$. Then we define $e(P) = |P \cap \Lambda|$ for every polytope $P \subset V$ such that the vertices of $tP$ belong to $\Lambda$ for some integer $t$. Apparently, to make a canonical choice of $\nu$, we have to fix some additional structure in $V$. In [6] a canonical valuation $\nu$ is constructed for rational polytopes whose facets are parallel to hyperplanes from a given finite collection of hyperplanes. Valuation $\nu$ agrees with $e$ in the $k + 1$ highest terms and for any fixed $k$ valuation $\nu$ is polynomially computable on polytopes with the number facets exceeding the dimension $d$ by not more than a constant fixed in advance (in particular, on rational simplices). In [10] a different canonical valuation $\mu$ is constructed provided a scalar product on $V$ is chosen. Valuation $\mu$ also agrees with $e$ on the $k + 1$ highest terms and polynomially computable on the same class of polytopes.

# References

[1] Bach, E., Miller, G., Shallit, J., Sums of divisors, perfect numbers and factoring, *SIAM J. Comput.* **15** (1986), 1143–1154.

[2] Baldoni-Silva, W., De Loera, J. A., Vergne, M. Counting integer flows in networks, *Found. Comput. Math.* **4** (2004), 277–314.

[3] Barvinok, A.I., A polynomial time algorithm for counting integral points in polyhedra when the dimension is fixed, *Math. Oper. Res.* **19** (1994), 769–779.

[4] Barvinok, A., Pommersheim, J.E., An algorithmic theory of lattice points in polyhedra. In *New Perspectives in Algebraic Combinatorics (Berkeley, CA, 1996–97)*. Math. Sci. Res. Inst. Publ., **38**, Cambridge Univ. Press, Cambridge, 1999, 91–147.

[5] Barvinok, A., *A Course in Convexity*, Graduate Studies in Mathematics, **54**. American Mathematical Society, Providence, RI, 2002.

[6] Barvinok, A., Computing the Ehrhart quasi-polynomial of a rational simplex, *Mathematics of Computation*, to appear.

[7]   Barvinok, A., Woods, K., Short rational generating functions for lattice point problems, *J. Amer. Math. Soc.* **16** (2003), 957–979.

[8]   Beck, M., Robins, S., *Computing the Continuous Discretely. Integer-point Enumeration in Polyhedra*, Undergraduate Texts in Mathematics, Springer-Verlag, Berlin, to appear.

[9]   Beck, M., Sottile, F., *Irrational proofs for three theorems of Stanley*, preprint arXiv math.CO/0501359, 2005.

[10]  Berline, N., Vergne, M., *Local Euler-Maclaurin formula for polytopes*, preprint arXiv math.CO/0507256, 2005.

[11]  Brion, M. Points entiers dans les polyédres convexes, *Ann. Sci. cole Norm. Sup. (4)* **21** (1988), 653–663.

[12]  De Loera, J. A., Hemmecke, R., Tauzer, J., Yoshida, R., Effective lattice point counting in rational convex polytopes, *J. Symbolic Comput.* **38** (2004), 1273–1302; see also `http://www.math.ucdavis.edu/~latte/`

[13]  De Loera, J.A., Hemmecke, R., Köppe, M., Weismantel, R., Integer polynomial optimization in fixed dimension, *Math. Oper. Res.*, to appear.

[14]  Denham, G., Short generating functions for some semigroup algebras., *Electron. J. Combin.* **10** (2003), Research Paper 36, 7 pp. (electronic).

[15]  Fischer, M. J., Rabin, M. O., Super-exponential complexity of Presburger arithmetic. In *Complexity of Computation (Proc. SIAM-AMS Sympos., New York, 1973)*, pp. 27–41. SIAM-AMS Proc., Vol. VII, Amer. Math. Soc., Providence, R.I., 1974.

[16]  Fulton, W., *Introduction to Toric Varieties*, Annals of Mathematics Studies, **131**. Princeton University Press, Princeton, 1993.

[17]  Grötschel, M., Lovász, L., Schrijver, A., *Geometric Algorithms and Combinatorial Optimization. Second edition*, Algorithms and Combinatorics, **2**. Springer-Verlag, Berlin, 1993.

[18]  Herzog, J., Generators and relations of abelian semigroups and semigroup rings, *Manuscripta Math.* **3** (1970), 175–193.

[19]  Kannan, R., Test sets for integer programs, $\forall\exists$ sentences. In *Polyhedral Combinatorics (Morristown, NJ, 1989)* 39–47, DIMACS Ser. Discrete Math. Theoret. Comput. Sci., 1, Amer. Math. Soc., Providence, RI, 1990.

[20]  Kannan, R., Lattice translates of a polytope and the Frobenius problem, *Combinatorica* **12** (1992), 161–177.

[21]  Kannan, R., Lovász, L., Scarf, H. E., The shapes of polyhedra, *Math. Oper. Res.* **15** (1990), 364–380.

[22]  Khinchin, A. Ya., *Continued Fractions*, The University of Chicago Press, Chicago, Ill.-London, 1964.

[23]  Khovanskii, A. G., Pukhlikov, A. V., The Riemann-Roch theorem for integrals and sums of quasipolynomials on virtual polytopes, (Russian) *Algebra i Analiz* **4** (1992), no. 4, 188–216; translation in *St. Petersburg Math. J.* **4** (1993), no. 4, 789–812.

[24]  Khovanskii, A. G., Sums of finite sets, orbits of commutative semigroups and Hilbert functions, (Russian) *Funktsional. Anal. i Prilozhen.* **29** (1995), no. 2, 36–50, 95; translation in *Funct. Anal. Appl.* **29** (1995), no. 2, 102–112.

[25] Lawrence, J., Rational-function-valued valuations on polyhedra. In *Discrete and Computational Geometry (New Brunswick, NJ, 1989/1990)*. DIMACS Ser. Discrete Math. Theoret. Comput. Sci., **6**, Amer. Math. Soc., Providence, RI, 1991, 199–208.

[26] McMullen, P., Valuations and dissections. In *Handbook of Convex Geometry, Vol. A, B*. North-Holland, Amsterdam, 1993, 933-988.

[27] M. Morales, Syzygies of monomial curves and a linear diophantine problem of Frobenius, *preprint*, Max-Planck-Institut für Mathematik, Bonn, 1986.

[28] Papadimitriou, C. H., *Computational Complexity*, Addison-Wesley Publishing Company, Reading, MA, 1994.

[29] Pommersheim, J., Thomas, H., Cycles representing the Todd class of a toric variety, *J. Amer. Math. Soc.* **17** (2004), 983–994.

[30] Scarf, H. E., Test sets for integer programs. In *Lectures on Mathematical Programming (ISMP97) (Lausanne, 1997)*. Math. Programming **79** (1997), no. 1-3, Ser. B, 355–368.

[31] Schöning, U., Complexity of Presburger arithmetic with fixed quantifier dimension, *Theory Comput. Syst.* **30** (1997), 423–428.

[32] Skriganov, M. M., Ergodic theory on $SL(n)$, Diophantine approximations and anomalies in the lattice point problem, *Invent. Math.* **132** (1998), 1–72.

[33] Shallcross, D., Neighbors of the origin for four by three matrices, *Math. Oper. Res.* **17** (1992), 608–614.

[34] Stanley, R. P., *Enumerative Combinatorics. Vol. 1. Corrected reprint of the 1986 original*, Cambridge Studies in Advanced Mathematics, **49**. Cambridge University Press, Cambridge, 1997.

[35] Székely L. A., Wormald, N. C., Generating functions for the Frobenius problem with 2 and 3 generators, *Math. Chronicle* **15** (1986), 49–57.

[36] Verdoolaege, S., Woods, K., Bruynooghe M., Cools R., Computation and manipulation of enumerators of integer projections of parametric polytopes. Preprint Katholieke Universiteit Leuven, Dept. of Computer Science, Report CW 392, 2005; see also `http://www.kotnet.org/~skimo/barvinok/`

[37] Woods, K. M., Rational generating functions and lattice point sets. Diss. Univ. of Michigan, 2004.

[38] Woods, K., Computing the period of an Ehrhart quasipolynomial, *The Electron J. Combin.* **12** (2005), Research paper 34, 12 pp.

Department of Mathematics, University of Michigan, Ann Arbor, MI 48109-1043

E-mail: barvinok@umich.edu