# POLITECNICO DI TORINO Repository ISTITUZIONALE

### Waring's theorem for binary powers

Original

Waring's theorem for binary powers / Kane Daniel, M.; Sanna, Carlo; Shallit, Jeffrey. - In: COMBINATORICA. - ISSN 0209-9683. - STAMPA. - 39:6(2019), pp. 1335-1350. [10.1007/s00493-019-3933-3]

Availability: This version is available at: 11583/2819915 since: 2020-05-05T17:54:36Z

Publisher: Springer

Published DOI:10.1007/s00493-019-3933-3

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)

# Waring's Theorem for Binary Powers

Daniel M. Kane<sup>\*</sup> Mathematics University of California, San Diego 9500 Gilman Drive #0404 La Jolla, CA 92093-0404 USA dakane@math.ucsd.edu

Carlo Sanna<sup>†</sup> Dipartimento di Matematica "Giuseppe Peano" Università degli Studi di Torino Via Carlo Alberto 10 10123 Torino Italy carlo.sanna.dev@gmail.com

Jeffrey Shallit<sup>‡</sup> School of Computer Science University of Waterloo Waterloo, ON N2L 3G1 Canada shallit@uwaterloo.ca

February 15, 2019

#### Abstract

A natural number is a *binary* k *th power* if its binary representation consists of k consecutive identical blocks. We prove, using tools from combinatorics, linear algebra, and number theory, an analogue of Waring's theorem for sums of binary k *th* 

<sup>\*</sup>D. M. Kane was supported by NSF Award CCF-1553288 (CAREER) and a Sloan Research Fellowship. \*C. Sanna is a member of the INdAM group GNSAGA.

<sup>&</sup>lt;sup>‡</sup>J. Shallit was supported by NSERC Discovery Grants #105829/2013 and 2018-04118.

powers. More precisely, we show that for each integer  $k \geq 2$ , there exists an effectively computable natural number n such that every sufficiently large multiple of  $E_k := \gcd(2^k - 1, k)$  is the sum of at most n binary k'th powers. (The hypothesis of being a multiple of  $E_k$  cannot be omitted, since we show that the gcd of the binary k'th powers is  $E_k$ .) Furthermore, we show that  $n = 2^{O(k^3)}$ . Analogous results hold for arbitrary integer bases b > 2.

# 1 Introduction

Let  $\mathbb{N} = \{0, 1, 2, ...\}$  be the natural numbers and let  $S \subseteq \mathbb{N}$ . The principal problem of additive number theory is to determine whether every integer N (resp., every sufficiently large integer N) can be represented as the sum of some *constant* number of elements of S, not necessarily distinct, where the constant does not depend on N. For a superb introduction to this topic, see [14].

Probably the most famous theorem of additive number theory is Lagrange's theorem from 1770: every natural number is the sum of four squares [10]. Waring's problem (see, e.g., [21, 22]), first stated by Edward Waring in 1770, is to determine g(k) such that every natural number is the sum of g(k) k'th powers. (A priori, it is not even clear that  $g(k) < \infty$ , but this was proven by Hilbert in 1909 [7].) From Lagrange's theorem we know that g(2) = 4. For other results concerning sums of squares, see, e.g., [6, 13].

If every natural number is the sum of k elements of S, we say that S forms a *basis* of order k. If every sufficiently large natural number is the sum of k elements of S, we say that S forms an *asymptotic basis* of order k.

In this paper, we consider a variation on Waring's theorem, where the ordinary notion of integer power is replaced by a related notion inspired from formal language theory. There is other recent work along the same lines. For example, Banks [1] recently proved that every natural number is the sum of at most 49 numbers whose base-10 expansion is a palindrome, and Cilleruelo, Luca, and Baxter [4] improved this result to 3 summands for all bases  $b \geq 5$ .

Our main result is Theorem 1.1 below, which we prove using arguments from combinatorics, linear algebra, and number theory; it concerns sums of binary k'th powers. We say that a natural number N is a *base-b* k'th power if its base-b representation consists of kconsecutive identical blocks. For example, 3549 in base 2 is

#### 1101 1101 1101,

so 3549 is a base-2 (or binary) cube. Throughout this paper, we consider only *canonical* base-*b* expansions (that is, those without leading zeros). Hence a number N > 0 is a base-*b* k'th power if and only if

$$N = a \cdot c_k^b(n),$$

for some  $n \geq 1$ , where

$$c_k^b(n) := \frac{b^{kn} - 1}{b^n - 1} = 1 + b^n + \dots + b^{(k-1)n}$$

and

$$b^{n-1} \le a < b^n. \tag{1}$$

The latter condition is needed to ensure that the base-*b* k'th power is formed by the concatenation of blocks that begin with a nonzero digit. Such a number consists of k consecutive blocks of digits, each of length n. For example,  $3549 = 13 \cdot c_3^2(4)$ .

The binary squares

 $0, 3, 10, 15, 36, 45, 54, 63, 136, 153, 170, 187, 204, 221, 238, 255, 528, 561, 594, 627, \ldots$ 

form sequence <u>A020330</u> in Sloane's *On-Line Encyclopedia of Integer Sequences* [20]. The binary cubes

 $0, 7, 42, 63, 292, 365, 438, 511, 2184, 2457, 2730, 3003, 3276, 3549, 3822, 4095, 16912, \ldots$ 

form sequence  $\underline{A297405}$ .

We define

 $\mathcal{S}_k^b := \{ n \ge 0 : n \text{ is a base-} b \ k' \text{th power} \} = \{ 0 \} \cup \{ a \cdot c_k^b(n) : n \ge 1, \ b^{n-1} \le a < b^n \}.$ 

The set  $\mathcal{S}_k^b$  is an interesting and natural set to study because its counting function is  $\Omega(N^{1/k})$ , just like the ordinary k'th powers. It has also appeared in a number of recent papers (e.g., [2]). However, there are two significant differences between the ordinary k'th powers and the base-b k'th powers.

The first difference is that 1 is not a base-b k'th power for k > 1. Thus, the base-b k'th powers cannot, in general, form a basis of finite order, but only an asymptotic basis.

A more significant difference is that the gcd of the ordinary k'th powers is always equal to 1, while the gcd  $E_k$  of the base-b k'th powers may, in some cases, be greater than one. This is quantified in Section 2. Thus, it is not reasonable to expect that every sufficiently large natural number can be the sum of a fixed number of base-b k'th powers; only those that are also a multiple of  $E_k$  can be so represented.

Our main result is the following:

**Theorem 1.1.** For every integer  $k \ge 1$  there exists a natural number n such that every sufficiently large multiple of  $E_k = \gcd(2^k - 1, k)$  is representable as the sum of n binary k'th powers. Furthermore, if W(k) is the least such n, then  $W(k) = 2^{O(k^3)}$ .

*Remark* 1.2. The fact that W(2) = 4 was proved in [12].

It may be worth noting that the methods that we use for the binary k'th powers have almost nothing in common with the deep number-theoretic tools (such as the circle method, [14]) that have been developed to handle the ordinary version of Waring's theorem. Indeed, it is not even clear that those tools could be adapted for use in our problem.

# 2 The greatest common divisor of $S_k^b$

We need the following classic lemma, sometimes called the "lifting-the-exponent" or LTE lemma [3]. Let  $\nu_p(n)$  denote the *p*-adic valuation of *n* (the exponent of the highest power of *p* dividing *n*).

**Lemma 2.1.** If p is a prime number and  $c \neq 1$  is an integer such that  $p \mid c - 1$ , then

$$\nu_p\left(\frac{c^n-1}{c-1}\right) \ge \nu_p(n),$$

for all positive integers n.

Now we prove several formulas for the greatest common divisor of the elements of  $\mathcal{S}_k^b$ .

**Theorem 2.2.** For  $k \ge 1$  define

$$A_k = \gcd(\mathcal{S}_k^b),$$
  

$$B_k = \gcd(c_k^b(1), c_k^b(2), \ldots),$$
  

$$C_k = \gcd(c_k^b(1), c_k^b(2), \ldots, c_k^b(k)),$$
  

$$D_k = \gcd(c_k^b(1), c_k^b(k)),$$
  

$$E_k = \gcd\left(\frac{b^k - 1}{b - 1}, k\right).$$

Then  $A_k = B_k = C_k = D_k = E_k$ .

*Proof.*  $A_k = B_k$ : If d divides  $B_k$ , then it clearly also divides all numbers of the form  $a \cdot c_k^b(n)$  with  $b^{n-1} \leq a < b^n$  and hence  $A_k$ .

On the other hand if d divides  $A_k$ , then it divides  $c_k^b(1)$ . Furthermore, d divides  $b^{n-1} \cdot c_k^b(n)$ and  $(b^{n-1}+1)c_k^b(n)$  (both of which are members of  $\mathcal{S}_k^b$  provided  $n \geq 2$ ). So it must divide their difference, which is just  $c_k^b(n)$ . So d divides  $B_k$ .

 $B_k = C_k$ : Note that d divides  $B_k$  if and only if it divides  $c_k^b(1)$  and also  $c_k^b(n) \mod c_k^b(1)$  for all  $n \ge 1$ . Now it is well known that, for  $b \ge 2$  and integers  $n, k \ge 1$ , we have

$$b^n \equiv b^{n \bmod k} \pmod{b^k - 1}.$$

(See, for example, [9, Ex. 4.3.2.6 and 4.5.3.31].) Hence

$$\begin{aligned} c_k^b(n) &= 1 + b^n + \dots + b^{(k-1)n} \equiv 1 + b^{n \mod k} + \dots + b^{(k-1)n \mod k} \pmod{b^k - 1} \\ &\equiv 1 + b^a + \dots + b^{(k-1)a} \pmod{b^k - 1} \\ &\equiv 1 + b^a + \dots + b^{(k-1)a} \pmod{c_k^b(1)} \\ &\equiv c_k^b(a) \pmod{c_k^b(1)}, \end{aligned}$$

where  $a = n \mod k$ . Thus any divisor of  $C_k$  is also a divisor of  $B_k$ . The converse is clear.

 $D_k = E_k$ : It suffices to observe that

$$c_k^b(k) = 1 + b^k + \dots + b^{(k-1)k}$$
$$\equiv \overbrace{1+1+\dots+1}^k \pmod{b^k-1}$$
$$\equiv k \pmod{b^k-1}$$
$$\equiv k \pmod{\frac{b^k-1}{b-1}}$$
$$\equiv k \pmod{c_k^b(1)}.$$

 $B_k = E_k$ : Every divisor of  $B_k$  clearly divides  $D_k$ , and above we saw  $D_k = E_k$ . We now show that every prime divisor of  $E_k$  divides  $B_k$  to at least the same order, thus showing that every divisor of  $E_k$  divides  $B_k$ .

Fix an integer  $\ell \geq 1$  and let p be a prime factor of  $E_k$ . On the one hand, if  $p \mid b^{\ell} - 1$ , then by Lemma 2.1 we get that

$$\nu_p(c_k^b(\ell)) = \nu_p\left(\frac{b^{k\ell} - 1}{b^\ell - 1}\right) \ge \nu_p(k) \ge \nu_p(E_k),$$

since  $E_k \mid k$ . Hence  $p^{\nu_p(E_k)} \mid c_k^b(\ell)$ . On the other hand, if  $p \nmid b^\ell - 1$ , then  $p^{\nu_p(E_k)}$  divides  $c_k^b(\ell) = \frac{b^{k\ell} - 1}{b^\ell - 1}$  simply because  $p^{\nu_p(E_k)}$  divides the numerator but does not divide the denominator. In both cases, we have that  $p^{\nu_p(E_k)} \mid c_k^b(\ell)$ , and since this is true for all prime divisors of  $E_k$ , we get that  $E_k \mid c_k^b(\ell)$ , as desired.

Remark 2.3. For b = 2, the sequence  $E_k$  is sequence <u>A014491</u> in Sloane's *Encyclopedia*. We make some additional remarks about the values of  $E_k$  in Section 5.

In the remainder of the paper, for concreteness, we focus on the case b = 2. We set  $c_k(n) := c_k^2(n)$  and  $S_k := S_k^2$ . However, everything we say also applies more generally to bases b > 2.

# **3** Waring's theorem for binary k'th powers: proof outline and tools

In this section, we give an outline of the proof of Theorem 1.1. All of the mentioned constants depend only on k.

Given a number N, a multiple of  $E_k$ , that we wish to represent as a sum of binary k'th powers, we first choose a suitable power of 2, say  $x = 2^n$ , and think of N as a degree-kpolynomial p evaluated at x. For example, we can represent N in base  $2^n$ ; the "digits" of this representation then correspond to the coefficients of p.

Similarly, the integers  $c_k(n), c_k(n+1), \ldots, c_k(n+k-1)$  can also be viewed as polynomials in  $x = 2^n$ . By linear algebra, there is a unique way to rewrite p as a linear combination of  $c_k(n), c_k(n+1), \ldots, c_k(n+k-1)$ , and this linear transformation can be represented by a matrix M that depends only on k, and is independent of n.

At first glance, such a linear combination would seem to provide a suitable representation of N in terms of binary k'th powers, but there are three problems to overcome:

- (a) the coefficients of  $c_k(i)$ ,  $n \le i < n + k$ , could be much too large;
- (b) the coefficients could be too small (by Eq. (1), the coefficient of  $c_k(i)$  needs to be at least  $2^{i-1}$ ), or even negative;
- (c) the coefficients might not be integers.

Issue (a) can be handled by choosing n such that  $2^n \approx N^{1/k}$ . This guarantees that the resulting coefficients of the  $c_k(n)$  are at most a constant factor larger than  $2^n$ . Using Lemma 3.1 below, the coefficients can be "split" into at most a constant number of coefficients lying in the desired range.

Issue (b) is handled by not working with N, but rather with Y := N - D, where D is a suitably chosen linear combination of  $c_k(n), c_k(n+1), \ldots, c_k(n+k-1)$  with large positive integer coefficients. Any negative coefficients arising in the expression for Y can now be offset by adding the large positive coefficients corresponding to D, giving us coefficients for the representation of N that are positive and lie in a suitable range.

Issue (c) is handled by finding  $d_k$ , the common denominator of the rational numbers involved, and working with  $\lfloor Y/d_k \rfloor$  instead of Y. Once a representation is found, multiplying by  $d_k$  gives us a representation with integer coefficients for a number Y' close to Y. The difference is sufficiently small that it can be handled. This completes the sketch of our construction. It is carried out in more detail in the rest of the paper.

## **3.1** Expressing multiples of $c_k(n)$ as a sum of binary k'th powers

As we have seen in Eq. (1), a positive integer of the form  $a \cdot c_k(n)$  is a binary k'th power if  $2^{n-1} \leq a < 2^n$ . But how about larger multiples of  $c_k(n)$ ? The following lemma will be useful.

**Lemma 3.1.** Let  $a \ge 2^{n-1}$ . Then  $a \cdot c_k(n)$  is the sum of at most  $\lceil \frac{a}{2^n-1} \rceil$  binary k'th powers.

*Proof.* Clearly the claim is true for  $2^{n-1} \leq a < 2^n$ . Otherwise, define  $b := \lceil \frac{a}{2^n-1} \rceil$  and  $c := (2^n - 1)b - a$ , so that  $0 \leq c < 2^n - 1$ . Then  $a = (b - 2)(2^n - 1) + d_1 + d_2$ , where  $d_1 = \lfloor (2^n - 1) - \frac{c}{2} \rfloor$  and  $d_2 = \lceil (2^n - 1) - \frac{c}{2} \rceil$ . A routine calculation now shows that  $2^{n-1} \leq d_1 \leq d_2 < 2^n$ , and so  $a \cdot c_k(n)$  is the sum of b binary k'th powers.

### 3.2 Change of basis and the Vandermonde matrix

In what follows, matrices and vectors are always indexed starting at 0. Recall that a Vandermonde matrix

$$V(a_0, a_1, \ldots, a_{k-1})$$

is a  $k \times k$  matrix where the entry in the *i*'th row and *j*'th column, for  $0 \le i, j < k$ , is defined to be  $a_i^j$ . The matrix is invertible if and only if the  $a_i$  are distinct.

Recall that  $c_k(n) = 1 + 2^n + 2^{2n} + \dots + 2^{(k-1)n}$ . For  $k \ge 1$  and  $n \ge 0$  we have

$$\begin{bmatrix} c_k(n) \\ c_k(n+1) \\ \vdots \\ c_k(n+k-1) \end{bmatrix} = M_k \begin{bmatrix} 1 \\ 2^n \\ \vdots \\ 2^{(k-1)n} \end{bmatrix},$$
(2)

where  $M_k = V(1, 2, 4, ..., 2^{k-1})$ . For example,

$$M_4 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 4 & 8 \\ 1 & 4 & 16 & 64 \\ 1 & 8 & 64 & 512 \end{bmatrix}$$

Let a natural number Y be represented as an  $\mathbb{N}$ -linear combination

$$Y = a_0 + a_1 2^n + \dots + a_{k-1} 2^{(k-1)n}$$

Then, multiplying Eq. (2) on the left by

$$[b_0 \quad b_1 \quad \cdots \quad b_{k-1}] := [a_0 \quad a_1 \quad \cdots \quad a_{k-1}]M_k^{-1}, \tag{3}$$

we get the following expression for Y as a  $\mathbb{Q}$ -linear combination of binary k'th powers:

$$Y = b_0 c_k(n) + b_1 c_k(n+1) + \dots + b_{k-1} c_k(n+k-1).$$
(4)

It remains to estimate the size of the coefficients  $b_i$ , as well as the sizes of their denominators.

The Vandermonde matrix is well studied (e.g., [16, pp. 43, 105]). We recall one basic fact about it.

**Lemma 3.2.** The determinant of  $V(a_0, a_1, \ldots, a_{k-1})$  is

$$\prod_{0 \le i < j < k} (a_j - a_i).$$

We now define  $d_k$  to be the determinant of  $M_k$ , and  $\ell_k$  to be the largest of the absolute values of the entries of  $M_k^{-1}$ . Note that, by Lemma 3.2,  $d_k$  is positive. Also, Laplace's formula tells us that  $M_k^{-1} = M'_k d_k^{-1}$ , where  $M'_k$  is the adjugate (classical adjoint)  $M'_k$  of  $M_k$ . Furthermore, since  $M_k$  has integer entries, so does  $M'_k$ .

**Proposition 3.3.** We have  $0 < d_k < 2^{k^3/3}$  for  $k \ge 1$ .

*Proof.* By the formula of Lemma 3.2 we know that

$$d_k = \prod_{0 \le i < j < k} (2^j - 2^i) < \prod_{0 \le i < j < k} 2^j = 2^{k^3/3 - k^2/2 + k/6} < 2^{k^3/3}$$

for  $k \geq 1$ .

The sequence  $(d_k)$  is sequence <u>A203303</u> in the OEIS [20].

Our next result demonstrates that  $\ell_k$ , the absolute value of the largest entry in  $M_k^{-1}$ , is bounded above by a constant.

**Proposition 3.4.** We have  $\ell_k < 34$ .

*Proof.* As is well known (see, e.g., [8, Exercise 1.2.3.40], the *i*'th column in the inverse of the Vandermonde matrix  $V(a_0, a_1, \ldots, a_{k-1})$  consists of the coefficients of the polynomial

$$p_i(x) := \prod_{\substack{0 \le j < k \\ j \ne i}} \frac{x - a_j}{a_i - a_j}.$$

We also observe that if

$$(x - b_1)(x - b_2) \cdots (x - b_n) = x^n + c_{n-1}x^{n-1} + \dots + c_1x + c_0,$$

is a polynomial with real roots, then the absolute value of every coefficient  $c_i$  is bounded by

$$|c_0| + \dots + |c_{n-1}| \le \prod_{1 \le i \le n} (1 + |b_i|).$$

Putting these two facts together, we see that all of the entries in the *i*'th column of  $V(a_0, a_1, \ldots, a_{k-1})^{-1}$  are, in absolute value, bounded by

$$P_k(i) := \frac{\prod_{\substack{0 \le j < k \\ j \ne i}} (1 + |a_j|)}{\prod_{\substack{0 \le j < k \\ j \ne i}} |a_j - a_i|}$$

Now let's specialize to  $a_{\ell} = 2^{\ell}$ . We get

$$P_k(i) := \frac{\prod_{\substack{0 \le j < k \\ j \ne i}} (2^j + 1)}{\prod_{\substack{0 \le j < k \\ j \ne i}} |2^j - 2^i|} \le \frac{\prod_{0 \le j < k} (2^j + 1)}{\prod_{\substack{0 \le j < k \\ j \ne i}} |2^j - 2^i|}$$

To finish the proof of the upper bound, it remains to find a lower bound for the denominator

$$Q_k(i) := \prod_{\substack{0 \le j < k \\ j \ne i}} |2^j - 2^i|$$

We claim, for  $k \geq 2$ , that

$$Q_k(0) \ge Q_k(1) \tag{5}$$

and

$$Q_k(1) \le Q_k(2) \le \dots \le Q_k(k-1).$$
(6)

To see (5), note that  $Q_k(0) = \prod_{2 \le j < k} (2^j - 1)$  and  $Q_k(1) = \prod_{2 \le j < k} (2^j - 2)$ . On the other hand, by telescoping cancellation we see, for  $1 \le i \le k - 2$ , that

$$\frac{Q_k(i)}{Q_k(i+1)} = \frac{2^{k-1} - 2^i}{(2^{i+1} - 1)2^{k-2}} < \frac{2^{k-1}}{3 \cdot 2^{k-2}} = \frac{2}{3},$$

which proves (6). Hence  $Q_k(i)$  is minimized at i = 1. Now

$$\ell_k \le \max_{0 \le i < k} \frac{\prod_{0 \le j < k} (2^j + 1)}{Q_k(i)} \le \frac{\prod_{0 \le j < k} (2^j + 1)}{Q_k(1)}$$
$$= \frac{\prod_{0 \le j < k} (2^j + 1)}{\prod_{2 \le j < k} (2^j - 2)} < 2 \cdot 3 \cdot \prod_{j \ge 2} \frac{2^j + 1}{2^j - 2} \doteq 33.023951743 \dots < 34,$$

where the last product has been estimated with a personal computer, using the inequalities

$$1 < \prod_{j \ge N} \frac{2^j + 1}{2^j - 2} = \prod_{j \ge N} \left( 1 + \frac{3}{2^j - 2} \right) < \exp\left(\sum_{j \ge N} \frac{3}{2^j - 2}\right) < \exp\left(\frac{3}{2^{N-2}}\right),$$
  
= 50.

with N = 50.

Remark 3.5. The tightest upper bound seems to be  $\ell_k < 5.194119929183 \cdots$  for all k, but we do not prove this here; see [18].

### 3.3 The Frobenius number

Let S be a set and x be a real number. By xS we mean the set  $\{xs : s \in S\}$ .

Let  $S \subseteq \mathbb{N}$  with gcd(S) = 1. The Frobenius number of S, written F(S), is the largest integer that cannot be represented as a non-negative integer linear combination of elements of S. See, for example, [17].

As we have seen,  $gcd(\mathcal{S}_k) = E_k = gcd(k, 2^k - 1)$ . Thus  $gcd(E_k^{-1}\mathcal{S}_k) = 1$ . Define  $F_k$  to be the Frobenius number of the set  $E_k^{-1}\mathcal{S}_k$ . In this section we give a weak upper bound for  $F_k$ .

**Lemma 3.6.** For  $k \ge 2$  we have  $F_k \le 2^{k^2+k}$ .

Proof. Consider  $T = \{g_1, g_2, g_3\}$  where  $g_1 = 2^k - 1$ ,  $g_2 = (2^k - 2)\frac{2^{k^2} - 1}{2^k - 1}$ , and  $g_3 = (2^k - 1)\frac{2^{k^2} - 1}{2^k - 1}$ . We have  $T \subseteq S_k$ . Let d be the greatest common divisor of T. Then d divides  $g_3 - g_2 = \frac{2^{k^2} - 1}{2^k - 1}$  and  $g_1 = 2^k - 1$ . So d divides  $D_k$ . On the other hand, clearly,  $A_k$  divides d, while from Theorem 2.2 we know that  $A_k = D_k = E_k$ . Hence,  $d = E_k$ .

Clearly  $F(E_k^{-1}S_k) \leq F(E_k^{-1}T)$ . Furthermore, since  $g_1 | g_3$ , it follows that  $F(E_k^{-1}T) = F(\{E_k^{-1}g_1, E_k^{-1}g_2\})$ . By a well-known result (see, e.g., [17, Theorem 2.1.1, p. 31]), we have  $F(\{a, b\}) = ab - a - b$ , and the desired claim follows.

*Remark* 3.7. We compute explicitly that  $F_2 = 17$ ,  $F_3 = 723$ ,  $F_4 = 52753$ ,  $F_5 = 49790415$ , and  $F_6 = 126629$ . This is sequence <u>A298306</u> in the OEIS [20].

## 4 The complete proof

We are now ready to fill in the details of the proof of our main result, Theorem 1.1. We recall the definitions of the following quantities that will figure in the proof:

- $c_k(n) = 1 + 2^n + \dots + 2^{(k-1)n};$
- $E_k = \gcd(k, 2^k 1)$  is the greatest common divisor of the set  $S_k$  of binary k'th powers;
- $F_k$  is the Frobenius number of the set  $E_k^{-1}\mathcal{S}_k$ ;
- $d_k$  is the determinant of the Vandermonde matrix  $M_k = V(1, 2, ..., 2^{k-1});$
- $\ell_k$  is the largest of the absolute values of the entries of  $M_k^{-1}$

Proof of Theorem 1.1. The result is clear for k = 1, so let us assume  $k \ge 2$ . Set  $Z := (F_k + 1)E_k$  and  $c := k^2 \ell_k d_k 2^{k^2 - k + 1}$ . We construct a representation for every  $N > Z + c2^k$  that is also a multiple of  $E_k$ .

Define X := N - Z. By above X is positive. Choose n as large as possible so that  $X > c2^{kn}$ . By above  $n \ge 1$ , and by our choice of n we have  $X \le c2^{k(n+1)}$ .

First we explain how to write  $X = T + X_3$ , where

- (a)  $X_3 < c_k(n)$ ; and
- (b) T is an N-linear combination of  $c_k(n), \ldots, c_k(n+k-1)$  with all coefficients sufficiently large.

To do so, first define  $Q := c_k(n) + \cdots + c_k(n+k-1)$ . Note that

$$Q \leq kc_k(n+k-1)$$
  
=  $k(1+2^{n+k-1}+2^{2(n+k-1)}+\dots+2^{(k-1)(n+k-1)})$   
 $\leq k(1+2+2^2+\dots+2^{(k-1)(n+k-1)})$   
 $\leq k2^{(k-1)(n+k-1)+1}$   
=  $k2^{(k-1)n}2^{k^2-2k+2}$ .

It now follows that

$$\frac{X}{Q} > \frac{c2^{kn}}{k2^{(k-1)n}2^{k^2-2k+2}} = k\ell_k d_k 2^{n+k-1}.$$

Hence, if we define  $R := \lfloor X/Q \rfloor$ , then

$$R \ge k\ell_k d_k 2^{n+k-1}.\tag{7}$$

We have now obtained RQ (a good approximation of X), which is an N-linear combination of  $c_k(n), \ldots, c_k(n+k-1)$  with every coefficient equal to R, where  $0 \le X - RQ < Q$ . Set  $X_2 := X - RQ$ . We now improve this approximation of X using a greedy algorithm, as follows: from  $X_2$  we remove as many copies as possible of  $c_k(n + k - 1)$  while leaving the remainder nonnegative, then similarly as many copies as possible of  $c_k(n + k - 2)$  from what is left, and so forth, down to  $c_k(n)$ . More precisely, for each index  $i = k - 1, k - 2, \ldots, 0$  (in that order) set

$$r_i = \left\lfloor \frac{X_2 - \sum_{i < j < k} r_j c_k(n+j)}{c_k(n+i)} \right\rfloor,\,$$

and then define  $X_3 := X_2 - D$ , where

$$D = r_0 c_k(n) + r_1 c_k(n+1) + \dots + r_{k-1} c_k(n+k-1)$$

By the way we chose the  $r_i$ , we have

$$0 \le r_{k-1} < 2 \tag{8}$$

$$0 \le r_i < \frac{c_k(n+i+1)}{c_k(n+i)} < 2^{k-1} \text{ for } 0 \le i \le k-2.$$
(9)

Furthermore  $0 \le X_3 < c_k(n)$ . If T = RQ + D, then (a) and (b) above are now satisfied. Next, define  $Y := \lfloor X_3/d_k \rfloor$ . Since  $0 \le Y \le X_3 < c_k(n)$ , we can express Y in base  $2^n$  as

$$Y = a_0 + a_1 2^n + \dots + a_{k-1} 2^{(k-1)n},$$

where each  $a_i$  is an integer satisfying  $0 \le a_i < 2^n$ .

Applying the transformation discussed above in Section 3.2 to Y, we obtain the  $\mathbb{Q}$ -linear combination

$$Y = b_0 c_k(n) + b_1 c_k(n+1) + \dots + b_{k-1} c_k(n+k-1).$$

From Eqs. (3) and (4) we know that

$$|b_i| \le k\ell_k \cdot 2^n \tag{10}$$

for  $0 \leq i < k$ , and, furthermore, the denominator of each  $b_i$  divides  $d_k$ . Hence  $d_k Y$  is an integer that is a  $\mathbb{Z}$ -linear combination of the  $c_k(n), \ldots, c_k(n+k-1)$ .

Set  $X_4 := X_3 - d_k Y$ . Clearly  $0 \le X_4 < d_k$ . Putting this all together, we have

$$N = X + Z = X_2 + RQ + Z = X_3 + D + RQ + Z = X_4 + (D + d_kY + RQ) + Z.$$

From above we have that  $D+d_kY+RQ$  is a  $\mathbb{Z}$ -linear combination of  $c_k(n), \ldots, c_k(n+k-1)$ , say

$$D + d_k Y + RQ = s_0 c_k(n) + \dots + s_{k-1} c_k(n+k-1),$$

where  $s_i = r_i + d_k b_i + R$ . We now obtain upper and lower bounds on the  $s_i$ .

We have

$$s_i \ge d_k b_i + R \ge d_k (-k\ell_k 2^n) + k\ell_k d_k 2^{n+k-1} \ge k\ell_k d_k (2^{n+k-1} - 2^n) \ge 2^{n+k-1},$$

where we have used Eqs. (7) and (10) and the fact that  $k \ge 2$ . This gives the lower bound, and shows that no  $s_i$  is too small.

For the upper bound, note that

$$r_i \le 2^{k-1} \tag{11}$$

by Eqs. (8) and (9), that

$$d_k b_i \le k \ell_k d_k 2^n \tag{12}$$

by Eq. (10), and

$$R \le \frac{X}{Q} + 1 \le \frac{c2^{k(n+1)}}{Q} + 1 \le \frac{k^2\ell_k d_k 2^{k^2 - k + 1} 2^{k(n+1)}}{2^{(k-1)(n+k-1)}} + 1 \le k^2\ell_k d_k 2^{2k+n} + 1.$$
(13)

Putting together Eqs. (11), (12), and (13), and using Propositions 3.3 and 3.4, we get  $s_i = 2^{n+O(k^3)}$ . Using Lemma 3.1, we see that each  $s_i c_k(n+i)$  is the sum of at most  $2^{O(k^3)}$  binary k'th powers, and hence  $D + d_k Y + RQ$  is the sum of at most  $k 2^{O(k^3)} = 2^{O(k^3)}$  binary k'th powers.

Now by construction N, D,  $d_kY$ , RQ, and Z are all integer multiples of  $E_k$ , so  $X_4$  is also a multiple of  $E_k$ . Furthermore  $X_4 + Z > (F_k + 1)E_k$ , so  $X_4 + Z$  can be represented as a non-negative integer linear combination of binary k'th powers. On the other hand  $X_4 + Z \leq (F_k + 1)E_k + d_k = 2^{O(k^3)}$ . It therefore follows that N is the sum of at most  $2^{O(k^3)}$ binary k'th powers.

Remark 4.1. A more explicit version of our bound on W(k) is effectively computable from our proof.

## 5 Final remarks

Everything we have done in this paper is equally applicable to expansions in bases b > 2.

The bound  $W(k) = 2^{O(k^3)}$  we obtained in this paper is rather weak, and can certainly be improved. We leave this as work for the future. For example, we have

**Conjecture 5.1.** Every natural number > 147615 is the sum of at most nine binary cubes. The total number of exceptions is 4921.

*Remark* 5.2. We have verified this claim up to  $2^{27}$ .

There is another approach to Waring's theorem for binary powers that could potentially give much better bounds for W(k). For sets  $S, T \subseteq \mathbb{N}$  define the sumset S + T as follows:

$$S + T = \{s + t : s \in S, t \in T\}.$$

We make the following conjecture:

**Conjecture 5.3.** Writing  $C_n$  for the set  $\{a \cdot c_k^2(n) : 2^{n-1} \le a < 2^n\}$  of cardinality  $2^{n-1}$  (i.e., the *kn*-bit binary *k*'th powers), for  $n, k \ge 1$ , all the elements in the sumset

$$C_n + C_{n+1} + \dots + C_{n+k-1},$$

are actually represented uniquely as a sum of k elements, one chosen from each of the summands.

If this conjecture were true — we have proved it for  $1 \le k \le 3$  — it would prove that the sumset

$$\overbrace{\mathcal{S}_k+\dots+\mathcal{S}_k}^k$$

has positive density, and hence, by a result of Nathanson [15, Theorem 11.7, p. 366] (building on earlier work of Schirelmann), that  $S_k$  forms an asymptotic additive basis. From this we could obtain better bounds on W(k).

In the light of our results, it seems natural to ask about the set  $\mathcal{T}_1^b$  of positive integers k such that  $gcd(\mathcal{S}_k^b) = 1$ . Indeed, we have that the elements of  $\mathcal{T}_1^b$  are exactly the integers k such that  $\mathcal{S}_k^b$  forms an asymptotic additive basis for  $\mathbb{N}$ . It turn out that  $\mathcal{T}_1^b$  has a natural density, and even more can be said: since  $\left(\frac{b^k-1}{b-1}\right)_{k\geq 1}$  is a Lucas sequence, we can employ the same methods of [19] to prove the following result:

**Theorem 5.4.** For all integers  $g \ge 1$ ,  $b \ge 2$ , the set  $\mathcal{T}_g^b$  of positive integers k such that  $gcd(\mathcal{S}_k^b) = g$  has a natural density, given by

$$\mathbf{d}(\mathcal{T}_g^b) = \sum_{\substack{d \ge 1\\ \gcd(b,d)=1}} \frac{\mu(d)}{L_b(dg)}$$

where  $\mu$  is the Möbius function and  $L_b(x) := \operatorname{lcm}(x, \operatorname{ord}_x(b))$ , where  $\operatorname{ord}_x(b)$  is the multiplicative order of b, modulo x. In particular, the series converges absolutely.

Furthermore,  $\mathbf{d}(\mathcal{T}_g^b) > 0$  if and only if  $\mathcal{T}_g^b \neq \emptyset$  if and only if  $g = \gcd\left(L_b(g), \frac{b^{L_b(g)}-1}{b-1}\right)$ .

Also, employing the methods of [11], the counting function of the set  $\{g \ge 1 : \mathcal{T}_g^b \neq \emptyset\}$  can be shown to be  $\gg x/\log x$  and at most o(x), as  $x \to +\infty$ . Note only that, in doing so, where in [11] results of Cubre and Rouse [5] on the density of the set of primes p such that the rank of appearance of p in the Fibonacci sequence is divisible by a fixed positive integer m are used, one should instead use results on the density of the set of primes p such that ord<sub>p</sub>(b) is divisible by m — for example, those given by Wiertelak [23].

## Acknowledgments

We are grateful to Igor Pak for introducing the first and third authors to each other. We thank the referees for their careful reading of the paper.

# References

- W. D. Banks, Every natural number is the sum of forty-nine palindromes, *INTEGERS* — *Electronic J. Combinat. Number Theory*, **16** (2016), Paper #A3.
- [2] A. Bridy, R. J. Lemke-Oliver, A. Shallit, and J. Shallit, The generalized Nagell-Ljunggren problem: powers with repetitive representations, *Experimental Math.* 0 (2018), 1–12. https://doi.org/10.1080/10586458.2017.1419391
- [3] R. D. Carmichael, On the numerical factors of certain arithmetic forms, Amer. Math. Monthly 16 (1909), 153–159.
- [4] J. Cilleruelo, F. Luca, and L. Baxter, Every positive integer is a sum of three palindromes, Math. Comp. 87 (2018), 3023-3055. https://doi.org/10.1090/mcom/3221
- [5] P. Cubre and J. Rouse, Divisibility properties of the Fibonacci entry point, Proc. Amer. Math. Soc. 142 (11) (2014), 3771–3785.
- [6] E. Grosswald, Representations of Integers as Sums of Squares, Springer-Verlag, 1985.
- [7] D. Hilbert, Beweis für die Darstellbarkeit der ganzen Zahlen durch eine feste Anzahl n-ter Potenzen (Waringsches Problem), Math. Annalen 67 (1909), 281–300.
- [8] D. E. Knuth, The Art of Computer Programming, Vol. 1, Fundamental Algorithms, Third Edition, Addison-Wesley, 1997.
- [9] D. E. Knuth, *The Art of Computer Programming*, Vol. 2, Seminumerical Algorithms, Second Edition, Addison-Wesley, 1981.
- [10] J. L. Lagrange, Démonstration d'un théorème d'arithmétique, Nouv. Mém. Acad. Roy. Sc. de Berlin, 1770, pp. 123–133. Also in Oeuvres de Lagrange, 3 (1869), pp. 189–201.
- [11] P. Leonetti and C. Sanna, On the greatest common divisor of n and the nth Fibonacci number, Rocky Mountain J. Math., 48 (2018), 1191–1199.
- [12] P. Madhusudan, D. Nowotka, A. Rajasekaran, and J. Shallit, Lagrange's theorem for binary squares, in I. Potapov, P. Spirakis, and J. Worrell, eds., 43rd International Symposium on Mathematical Foundations of Computer Science (MFCS 2018), Leibniz International Proceedings in Informatics (LIPIcs), 117 (2018), 18:1–18:14. https:// doi.org/10.4230/LIPIcs.MFCS.2018.18.
- [13] C. J. Moreno and S. S. Wagstaff, Jr., Sums of Squares of Integers. Chapman and Hall/CRC, 2005.
- [14] M. B. Nathanson, Additive Number Theory: The Classical Bases, Springer, 1996.
- [15] M. B. Nathanson, *Elementary Methods in Number Theory*, Springer, 2000.

- [16] G. Pólya and G. Szegö, Problems and Theorems in Analysis II, Springer-Verlag, 1976.
- [17] J. L. Ramírez-Alfonsín, The Diophantine Frobenius Problem, Oxford University Press, 2006.
- [18] C. Sanna, J. Shallit, and S. Zhang, Largest entry in the inverse of a Vandermonde matrix, manuscript in preparation, February 2019.
- [19] C. Sanna and E. Tron, The density of numbers n having a prescribed G.C.D. with the nth Fibonacci number, Indag. Math. 29 (2018), 972–980.
- [20] N. J. A. Sloane et al., The On-Line Encyclopedia of Integer Sequences, 2017. Available at https://oeis.org.
- [21] C. Small, Waring's problem. *Math. Mag.* **50** (1977), 12–16.
- [22] R. C. Vaughan and T. Wooley, Waring's problem: a survey. In M. A. Bennett, B. C. Berndt, N. Boston, H. G. Diamond, A. J. Hildebrand, and W. Philipp, editors, *Number Theory for the Millennium. III*, pp. 301–340. A. K. Peters, 2002.
- [23] K. Wiertelak, On the density of some sets of primes p, for which  $n \mid \operatorname{ord}_p a$ , Funct. Approx. Comment. Math. 28 (2000), 237–241.