

# DECOMPOSING A FACTORIAL INTO LARGE FACTORS

BORIS ALEXEEV, EVAN CONWAY, ANDREW V. SUTHERLAND, TERENCE TAO, MARKUS UHR,  
AND KEVIN VENTULLO

ABSTRACT. Let  $t(N)$  denote the largest number such that  $N!$  can be expressed as the product of  $N$  integers greater than or equal to  $t(N)$ . The bound  $t(N)/N = 1/e - o(1)$  was apparently established in unpublished work of Erdős, Selfridge, and Straus; but the proof is lost. Here we obtain the more precise asymptotic

$$\frac{t(N)}{N} = \frac{1}{e} - \frac{c_0}{\log N} + O\left(\frac{1}{\log^{1+c} N}\right)$$

for an explicit constant  $c_0 = 0.30441901 \dots$  and some absolute constant  $c > 0$ , answering a question of Erdős and Graham. For the upper bound, a further lower order term in the asymptotic expansion is also obtained. With numerical assistance, we obtain highly precise computations of  $t(N)$  for wide ranges of  $N$ , establishing several explicit conjectures of Guy and Selfridge on this sequence. For instance, we show that  $t(N) \geq N/3$  for  $N \geq 43632$ , with the threshold shown to be best possible.

## 1. INTRODUCTION

Given a natural number  $M$ , define a *factorization* of  $M$  to be a finite multiset  $\mathcal{B}$  of natural numbers such that the product

$$\prod \mathcal{B} := \prod_{a \in \mathcal{B}} a$$

(where the elements are counted with multiplicity) is equal to  $M$ ; more generally, define a *subfactorization* of  $M$  to be a finite multiset  $\mathcal{B}$  such that  $\prod \mathcal{B}$  divides  $M$ . Given a threshold  $t$ , we say that a multiset  $\mathcal{B}$  is  *$t$ -admissible* if  $a \geq t$  for all  $a \in \mathcal{B}$ . For a given natural number  $N$ , we then define  $t(N)$  to be the largest  $t$  for which there exists a  $t$ -admissible factorization  $\mathcal{B}$  of  $N!$  of cardinality  $|\mathcal{B}| = N$ .

**Example 1.1.** The multiset

$$\{3, 3, 3, 3, 4, 4, 5, 7, 8\}$$

is a 3-admissible factorization of

$$\prod \{3, 3, 3, 3, 4, 4, 5, 7, 8\} = 3^4 \times 4^2 \times 5 \times 7 \times 8 = 9!$$

of cardinality

$$|\{3, 3, 3, 3, 4, 4, 5, 7, 8\}| = 9,$$

hence  $t(9) \geq 3$ . One can check that no 4-admissible factorization of  $9!$  of this cardinality exists, hence  $t(9) = 3$ .

It is easy to see that  $t(N)$  is non-decreasing in  $N$  (any cardinality  $N$  factorization of  $N!$  can be extended to a cardinality  $N + 1$  factorization of  $(N + 1)!$  by adding  $N + 1$  to the multiset). The first few elements of the sequence  $t(N)$  are

$$1, 1, 1, 2, 2, 2, 2, 2, 3, 3, 3, 3, 3, 4, \dots$$

(OEIS A034258). The values of  $t(N)$  for  $N \leq 79$  were computed in [10], and the values for  $N \leq 200$  can be extracted from OEIS A034259, which describes the inverse sequence to  $t$ . As part of our work, we extend this sequence to  $N \leq 10^4$ ; see [17] and Figure 5.

When the factorial  $N!$  is replaced with an arbitrary number the problem of determining  $t(N)$  essentially becomes the bin covering problem, which is known to be NP-hard; see e.g., [2]. However, as we shall see in this paper, the special structure of the factorial (and in particular, the profusion of factors at the “tiny primes” 2, 3) make it more tractable to estimate  $t(N)$  with far higher precision than if one were factoring an arbitrary number. For instance, we can show that

$$0 \leq t(9 \times 10^8) - 316\,560\,601 \leq 113;$$

see Table 2.

**Remark 1.2.** One can equivalently define  $t(N)$  as the greatest  $t$  for which there exists a  $t$ -admissible *subfactorization* of  $N!$  of cardinality *at least*  $N$ . This is because every such subfactorization can be converted into a  $t$ -admissible factorization of cardinality exactly  $N$  by first deleting elements from the subfactorization to make the cardinality  $N$ , and then multiplying one of the elements of the subfactorization by a natural number to upgrade the subfactorization to a factorization. This “relaxed” formulation of the problem turns out to be more convenient both for theoretical analysis of  $t(N)$  and for numerical computations.

By combining the obvious lower bound

$$\prod B \geq t^{|B|} \tag{1.1}$$

for any  $t$ -admissible multiset  $B$  with Stirling’s formula (2.4), we obtain the trivial upper bound

$$\frac{t(N)}{N} \leq \frac{(N!)^{1/N}}{N} = \frac{1}{e} + O\left(\frac{\log N}{N}\right); \tag{1.2}$$

see Figure 1. In [9, p.75] it was reported that an unpublished work of Erdős, Selfridge, and Straus established the asymptotic

$$\frac{t(N)}{N} = \frac{1}{e} + o(1) \tag{1.3}$$

(first conjectured in [7]) and asked if one could show the bound

$$\frac{t(N)}{N} \leq \frac{1}{e} - \frac{c}{\log N} \tag{1.4}$$

for some constant  $c > 0$  (problem #391 in <https://www.erdosproblems.com>; see also [10, Section B22, p. 122–123]); it was also noted that similar results were obtained in [1] if one restricted the  $a_i$  to be prime powers. However, as later reported in [8], Erdős “believed that Straus had written up our proof [of (1.3)]. Unfortunately Straus suddenly died and no trace was ever found of his notes. Furthermore, we never could reconstruct our proof, so our assertion now can be called only a conjecture”. In [10] the lower bound  $\frac{t(N)}{N} \geq \frac{1}{4}$  was established for



FIGURE 1. The function  $t(N)/N$  (blue) for  $N \leq 200$ , using the data from OEIS A034258, as well as the trivial upper bound  $(N!)^{1/N}/N$  (green), the improved upper bound from Lemma 5.1 (pink), which is asymptotic to (1.5) (purple), and the function  $\lfloor 2N/7 \rfloor / N$  (brown), which we show to be a lower bound for  $N \neq 56$ . Theorem 1.3(iv) implies that  $t(N)/N$  is asymptotic to (1.5) (purple), which in turn converges to  $1/e$  (orange), although it appears that (1.8) (gray) eventually becomes a sharper approximation. The threshold  $1/3$  (red) is permanently crossed at  $N = 43632$ .

sufficiently large  $N$ , by rearranging powers of 2 and 3 in the obvious factorization  $1 \times 2 \times \dots \times N$  of  $N!$ . A variant lower bound of the asymptotic shape  $\frac{t(N)}{N} \geq \frac{3}{16} - o(1)$  obtained by rearranging only powers of 2, and which is superior for medium values of  $N$ , can also be found in [10]. The following conjectures in [10] were also made:

- (1) One has  $t(N) \leq N/e$  for  $N \neq 1, 2, 4$ .
- (2) One has  $t(N) \geq \lfloor 2N/7 \rfloor$  for  $N \neq 56$ .
- (3) One has  $t(N) \geq N/3$  for  $N \geq 3 \times 10^5$ . (It was also asked if the threshold  $3 \times 10^5$  could be lowered.)

In this paper we answer all of these questions.

**Theorem 1.3** (Main theorem). *Let  $N$  be a natural number.*

- (i) *If  $N \neq 1, 2, 4$ , then  $t(N) \leq N/e$ .*
- (ii) *If  $N \neq 56$ , then  $t(N) \geq \lfloor 2N/7 \rfloor$ .*
- (iii) *If  $N \geq 43632$ , then  $t(N) \geq N/3$ . The threshold 43632 is best possible.*

FIGURE 2. A continuation of Figure 1 to the region  $80 \leq N \leq 599$ .

(iv) For large  $N$ , one has

$$\frac{t(N)}{N} = \frac{1}{e} - \frac{c_0}{\log N} + O\left(\frac{1}{\log^{1+c} N}\right) \quad (1.5)$$

for some constant  $c > 0$ , where  $c_0$  is the explicit constant

$$\begin{aligned} c_0 &:= \frac{1}{e} \int_0^1 f_e(x) dx \\ &= 0.30441901 \dots \end{aligned} \quad (1.6)$$

and for any  $\alpha > 0$ ,  $f_\alpha : (0, \infty) \rightarrow \mathbb{R}$  denotes the piecewise smooth function

$$f_\alpha(x) := \left\lfloor \frac{1}{x} \right\rfloor \log \frac{\lceil 1/\alpha x \rceil}{1/\alpha x}. \quad (1.7)$$

In particular, (1.3) and (1.4) hold. In fact the upper bound can be sharpened to

$$\frac{t(N)}{N} \leq \frac{1}{e} - \frac{c_0}{\log N} - \frac{c_1 + o(1)}{\log^2 N} \quad (1.8)$$

for an explicit constant  $c_1 = 0.75554808 \dots$ ; see Proposition 5.2.

For future reference, we observe the simple bounds

$$\begin{aligned} 0 \leq f_\alpha(x) &\leq \frac{1}{x} \log \frac{1/\alpha x + 1}{1/\alpha x} \\ &= \frac{1}{x} \log(1 + \alpha x) \\ &\leq \alpha \end{aligned} \quad (1.9)$$



FIGURE 3. The piecewise continuous function  $x \mapsto \frac{1}{e}f_e(x)$ , together with its mean value  $c_0 = 0.30441901 \dots$  and the upper bound  $\frac{\log(1+ex)}{ex}$ . The function exhibits an oscillatory singularity at  $x = 0$  similar to  $\sin \frac{1}{x}$  (but it is always nonnegative and bounded). Informally, the function  $f_e$  quantifies the difficulty that large primes in the factorization of  $N!$  have in becoming only slightly larger than  $N/e$  after multiplying by a natural number.

for all  $x > 0$ ; in particular,  $f_e$  is a bounded function. It however has an oscillating singularity at  $x = 0$ ; see Figure 3.

In Appendix C we give some details on the numerical computation of the constant  $c_0$ .

**Remark 1.4.** In a previous version [16] of this manuscript, the weaker bounds

$$\frac{1}{e} - \frac{O(1)}{\log N} \leq \frac{t(N)}{N} \leq \frac{1}{e} - \frac{c_0 + o(1)}{\log N}$$

were established, which were enough to recover (1.3), (1.4), and Theorem 1.3(i). Numerically, the upper bound in (1.8) appears to be a rather good approximation, and we conjecture that it is a lower bound as well.

As one might expect, the proof of Theorem 1.3 proceeds by a combination of both theoretical analysis and numerical calculations. Our main tools to obtain upper and lower bounds on  $t(N)$  can be summarized as follows (and in Table 1):

- In Section 3, we discuss *greedy algorithms* to construct subfactorizations, that provide quickly computable, though suboptimal, lower bounds on  $t(N)$  for small, medium, and moderately large values;

- In Section 4, we present a *linear programming* (and *integer programming*) method that provides quite accurate upper and lower bounds on  $t(N)$  for small and medium values of  $N$ , and which we apply in Section 5 to establish a general upper bound (Lemma 5.1) on  $t(N)$  that can be used to obtain Theorem 1.3(i);
- In Section 6, we extend the *rearrangement approach* from [11] to give a computer-assisted proof that Theorem 1.3(iii) holds for sufficiently large  $N$ , as well as an analytic proof of (1.3).
- In Section 7, we introduce an *accounting equation* linking the “ $t$ -excess” of a subfactorization with its “ $p$ -surpluses” at various primes, which provides an alternate proof of Lemma 5.1, and also is the starting point for the modified factorization technique discussed below;
- In Section 8, we give *modified approximate factorization* strategy, which provides lower bounds on  $t(N)$ , that become asymptotically quite efficient.

Claim	Range of $N$	Method used
Theorem 1.3(ii)	$N \leq 3 \times 10^5$	Greedy
Theorem 1.3(ii), (iii)	$67425 \leq N \leq 10^{12}$	Greedy
Theorem 1.3(i)	$N \leq 10^4$	Linear programming
Theorem 1.3(iii)	$43632 \leq N \leq 8 \times 10^4$	Integer programming
Theorem 1.3(i)	$N > 80$	Lemma 5.1
Theorem 1.3(iv) (upper)	$N$ sufficiently large	Lemma 5.1
Theorem 1.3(ii), (iii)	$N$ sufficiently large	Rearrangement
Theorem 1.3(ii), (iii)	$N \geq 10^{11}$	Modified approximate factorization
Theorem 1.3(iv) (lower)	$N$ sufficiently large	Modified approximate factorization

TABLE 1. The techniques in this paper can establish the various components of Theorem 1.3, for various overlapping ranges of  $N$ .

The final approach is significantly more complicated than the other four, but gives the most efficient lower bounds in the asymptotic limit  $N \rightarrow \infty$ . The key idea is to start with an approximate factorization

$$N! \approx \left( \prod_{j \in I} j \right)^A$$

for some relatively small natural number  $A$  (e.g.,  $A = \lfloor \log^2 N \rfloor$ ) and a suitable set  $I$  of natural numbers greater than or equal to  $t$ ; there is some freedom to select parameters here, and we will take  $I$  to be the natural numbers in  $(t, t(1 + \sigma)]$  that are 3-rough (coprime to 6), where  $t$  is the target lower bound for  $t(N)$  we wish to establish, and  $\sigma := \frac{3N}{tA}$  is chosen to bring the number of terms in the approximate factorization close to  $N$ . With this choice of  $I$ , this product contains approximately the right number of copies of  $p$  for medium-sized primes  $p$ ; but it has the “wrong” number of copies of large primes, and is also constructed to avoid the “tiny” primes  $p = 2, 3$ . One then performs a number of alterations to this approximate factorization to correct for the “surpluses” or “deficits” at various primes  $p > 3$ , using the supply of available tiny primes  $p = 2, 3$  as a sort of “liquidity pool” to efficiently reallocate primes in the factorization. A key point will be that the incommensurability of  $\log 2$  and  $\log 3$  (i.e., the irrationality of  $\log 3 / \log 2$ ) means that the 3-smooth numbers (numbers of the

form  $2^n 3^m$ ) are asymptotically dense (in logarithmic scale), allowing for other factors to be exchanged for 3-smooth factors with little loss<sup>1</sup>.

**1.1. Author contributions and data.** This project was initially conceived as a single-author manuscript by Terence Tao, but since the release of the initial preprint [16], grew to become a collaborative project organized via the Github repository [17], which also contains the supporting code and data for the project. The contributions of the individual authors, according to the CRediT categories<sup>2</sup>, are as follows:

- Boris Alexeev: Formal Analysis, Investigation, Software.
- Evan Conway: ...
- Andrew Sutherland: ...
- Terence Tao: Conceptualization, Formal Analysis, Methodology, Project Administration, Visualization, Writing – original draft, Writing – review & editing.
- Markus Uhr: Formal Analysis, Software.
- Kevin Ventullo: Software

**1.2. Acknowledgments.** TT is supported by NSF grant DMS-2347850. We thank Thomas Bloom for the web site <https://www.erdosproblems.com>, where TT learned of this problem, as well as Bryna Kra and Ivan Pan for corrections.

## 2. NOTATION AND BASIC ESTIMATES

We use the usual asymptotic notation  $X = O(Y)$ ,  $X \ll Y$ , or  $Y \gg X$  to denote an inequality of the form  $|X| \leq CY$  for some absolute constant  $C$ ; if we need this constant to depend on additional parameters, we will indicate this by subscripts, thus for instance  $O_M(Y)$  denotes a quantity bounded in magnitude by  $C_M Y$  for some  $C_M$  depending on  $M$ . We also write  $X \asymp Y$  for  $X \ll Y \ll X$ . For effective estimates, we will use the more precise notation  $O_{\leq}(Y)$  to denote any quantity whose magnitude is bounded by exactly at most  $Y$ . We also use  $\bar{O}_{\leq}(Y)^+$  to denote a quantity of size  $O_{\leq}(Y)$  that is also non-negative, that is to say it lies in the interval  $[0, Y]$ . We also use  $o(X)$  to denote any quantity bounded in magnitude by  $c(N)X$ , for some  $c(N)$  that goes to zero as  $N \rightarrow \infty$ .

If  $S$  is a statement, we use  $1_S$  to denote its indicator, thus  $1_S = 1$  when  $S$  is true and  $1_S = 0$  when  $S$  is false. If  $x$  is a real number, we use  $\lfloor x \rfloor$  to denote the greatest integer less than or equal to  $x$ , and  $\lceil x \rceil$  to be the least integer greater than or equal to  $x$ .

Throughout this paper, the symbol  $p$  (or  $p_0, p_1$ , etc.) is always understood to be restricted to be prime. We use  $(a, b)$  to denote the greatest common divisor of  $a$  and  $b$ ,  $a|b$  to denote the assertion that  $a$  divides  $b$ , and  $\pi(x) = \sum_{p \leq x} 1$  to denote the usual prime counting function.

<sup>1</sup>The weaker results alluded to in Remark 1.4 only used the prime 2 as a supply of “liquidity”, and thus encountered inefficiencies due to the inability to “make change” when approximating another factor by a power of two.

<sup>2</sup><https://credit.niso.org/>

We use  $v_p(a/b) = v_p(a) - v_p(b)$  to denote the  $p$ -adic valuation of a positive rational number  $a/b$ , that is to say the number of times  $p$  divides the numerator  $a$ , minus the number of times  $p$  divides the denominator  $b$ . For instance,  $v_2(32/27) = 5$  and  $v_3(32/27) = -3$ . If one applies a logarithm to the fundamental theorem of arithmetic, one obtains the identity

$$\sum_p v_p(r) \log p = \log r \quad (2.1)$$

for any positive rational  $r$ .

For a natural number  $n$ , we can write

$$v_p(n) = \sum_{j=1}^{\infty} 1_{p^j | n}. \quad (2.2)$$

Upon taking partial sums, we recover Legendre's formula

$$v_p(N!) = \sum_{j=1}^{\infty} \left\lfloor \frac{N}{p^j} \right\rfloor = \frac{N - s_p(N)}{p - 1} \quad (2.3)$$

where  $s_p(N)$  is the sum of the digits of  $N$  in the base  $p$  expansion.

Given a putative factorization  $B$  of  $N!$ , we refer to the quantity  $v_p\left(\frac{N!}{\prod B}\right)$  as the  $p$ -surplus of  $B$  with respect to the target  $N!$ , and similarly refer to the negative  $-v_p\left(\frac{N!}{\prod B}\right) = v_p\left(\frac{\prod B}{N!}\right)$  of this surplus as the  $p$ -deficit, with the multiset being  $p$ -balanced if the  $p$ -surplus (or  $p$ -deficit) is zero. Thus, a factorization of  $N!$  is achieved if and only if one is balanced at every prime  $p$ , whereas a subfactorization is achieved if one is either in balance or surplus at every prime  $p$ .

Let  $M(N, t)$  denote the maximal cardinality of a  $t$ -admissible subfactorization of  $N!$ ; thus, by Remark 1.2,  $t(N) \geq t$  if and only if  $M(N, t) \geq N$ .

To bound the factorial, we have the explicit Stirling approximation [14]

$$\log N! = N \log N - N + \log \sqrt{2\pi N} + O_{\leq}^+\left(\frac{1}{12N}\right), \quad (2.4)$$

valid for all natural numbers  $N$ .

**2.1. Approximation by 3-smooth numbers.** The primes 2, 3 will play a special role<sup>3</sup> in this paper and will be referred to as *tiny primes*. Call a natural number *3-smooth* if it is the product of tiny primes, i.e., it is of the form  $2^n 3^m$  for some natural numbers  $n, m$ , and *3-rough* if it is not divisible by any tiny prime, that is to say it is coprime to 6. Given a positive real number  $x$ , we use  $\lceil x \rceil^{(2,3)}$  to denote the smallest 3-smooth number greater than or equal to  $x$ . For instance,  $\lceil 5 \rceil^{(2,3)} = 6$  and  $\lceil 10 \rceil^{(2,3)} = 12$ .

<sup>3</sup>One could also run analogous arguments with other sets of tiny primes; for instance, the initial version [16] of this paper only utilized the prime 2 in this fashion.





FIGURE 4. The function  $\log \frac{[x]^{(2,3)}}{x}$ , compared against  $\kappa_x$ .

It will be convenient to introduce a variant of this quantity that is close to a power<sup>4</sup> of 12. If  $1 \leq L \leq x$  is an additional real parameter, we define

$$[x]_L^{(2,3)} := 12^a [x/12^a]^{(2,3)} \quad (2.5)$$

for any real  $x \geq L \geq 1$ , where  $a := \left\lfloor \frac{x/L}{\log 12} \right\rfloor$  is the largest integer such that  $12^a \leq x/L$ .

For any  $L \geq 1$ , let  $\kappa_L$  be the least quantity such that

$$x \leq [x]_L^{(2,3)} \leq \exp(\kappa_L)x \quad (2.6)$$

holds for all  $x \geq L$ ; see Figure 4. In Appendix A we establish the following facts:

**Lemma 2.1** (Approximation by 3-smooth numbers).

- (i) We have  $\kappa_{4.5} = \log \frac{4}{3} = 0.28768 \dots$  and  $\kappa_{40.5} = \log \frac{32}{27} = 0.16989 \dots$ .
- (ii) For large  $L$ , one has  $\kappa_L \ll \log^{-c} L$  for some absolute constant  $c > 0$ .
- (iii) If  $1 \leq L \leq x$  are real numbers, then

$$x \leq [x]_L^{(2,3)} \leq \exp(\kappa_L)x \quad (2.7)$$

<sup>4</sup>The significance of the base 12 is that the 3-smooth portion  $2^{v_2(N!)} 3^{v_3(N!)}$  of  $N!$ , which serves as our “liquidity pool”, is approximately  $2^N 3^{N/2} = \sqrt{12}^N$ ; see (2.3) below. This makes  $\log \sqrt{12}$  a natural “unit of currency” in which to conduct various factor exchanges, with various integer linear combinations of  $\log 2$  and  $\log 3$  usable as “small change” to approximate quantities that are not integer multiples of  $\log \sqrt{12} = \log 2 + \frac{1}{2} \log 3$ .

and for any  $0 \leq \gamma < 1$  we have

$$\frac{v_2(\lceil x \rceil_L^{(2,3)}) - 2\gamma v_3(\lceil x \rceil_L^{(2,3)})}{1 - \gamma} \leq \frac{\log x + \kappa_{L,\gamma}^{(2)}}{\log \sqrt{12}} \quad (2.8)$$

and

$$\frac{2v_3(\lceil x \rceil_L^{(2,3)}) - \gamma v_2(\lceil x \rceil_L^{(2,3)})}{1 - \gamma} \leq \frac{\log x + \kappa_{L,\gamma}^{(3)}}{\log \sqrt{12}} \quad (2.9)$$

where

$$\kappa_{L,\gamma}^{(2)} := \left( \frac{\log \sqrt{12}}{(1 - \gamma) \log 2} - 1 \right) \log(12L) + \frac{\kappa_L \log \sqrt{12}}{(1 - \gamma) \log 2} \quad (2.10)$$

$$\kappa_{L,\gamma}^{(3)} := \left( \frac{\log \sqrt{12}}{(1 - \gamma) \log \sqrt{3}} - 1 \right) \log(12L) + \frac{\kappa_L \log \sqrt{12}}{(1 - \gamma) \log \sqrt{3}}. \quad (2.11)$$

We remark that when  $x$  is a power of 12, the left-hand sides of (2.8), (2.9) are both equal to  $\frac{\log x}{\log \sqrt{12}}$ ; thus the estimates (2.8), (2.9) are quite efficient asymptotically.

**2.2. Sums over primes.** We recall the effective prime number theorem from [6, Corollary 5.2], which asserts that

$$\pi(x) \geq \frac{x}{\log x} + \frac{x}{\log^2 x} \quad (2.12)$$

for  $x \geq 599$  and

$$\pi(x) \leq \frac{x}{\log x} + \frac{1.2762x}{\log^2 x} \quad (2.13)$$

for  $x > 1$ .

We will also need to control sums of somewhat oscillatory functions over primes, for which the bounds in (2.12), (2.13) are of insufficient strength. Let  $y < x$  be real numbers. Given a function  $b : (y, x] \rightarrow \mathbb{R}$ , its *total variation*  $\|b\|_{\text{TV}(y,x]}$  is defined as the supremum of the quantities  $\sum_{j=0}^{J-1} |b(x_{j+1}) - b(x_j)|$  for  $y < x_0 \leq \dots \leq x_J \leq x$ , and the *augmented total variation*  $\|b\|_{\text{TV}^*(y,x]}$  is defined as

$$\|b\|_{\text{TV}^*(y,x]} := |b(y^+)| + |b(x)| + \|b\|_{\text{TV}(y,x]},$$

$b(y^+) := \lim_{t \rightarrow y^+} b(t)$  denotes the right limit of  $b$  at  $y$  (which exists if  $b$  is of finite total variation). Equivalently,  $\|b\|_{\text{TV}^*(y,x]}$  is the total variation of  $b$  if extended by zero outside of  $(y, x]$ . The indicator function  $1_{(y,x]}$  clearly has an augmented total variation of 2.

We will use this augmented total variation to control sums over primes. More precisely, in Appendix B we will show

**Lemma 2.2** (Effective bounds for oscillatory sums over primes). *Let  $1423 \leq y \leq x$ , and let  $b : (y, x] \rightarrow \mathbb{R}$  be of bounded total variation. Then we have the bound*

$$\sum_{y < p \leq x} b(p) \log p = \int_y^x \left( 1 - \frac{2}{\sqrt{t}} \right) b(t) dt + O_{\leq}(\|b\|_{\text{TV}^*(y,x]} E(x)) \quad (2.14)$$

where the error function  $E(x)$  is defined as

$$E(x) := 0.95\sqrt{x} + 3.83 \times 10^{-9}x. \quad (2.15)$$

In particular one has

$$\pi(x) - \pi(y) = \int_y^x \left(1 - \frac{2}{\sqrt{t}}\right) \frac{dt}{\log t} + O_{\leq} \left(2 \frac{E(x)}{\log y}\right); \quad (2.16)$$

estimating

$$1 - \frac{2}{\sqrt{y}} \leq 1 - \frac{2}{\sqrt{t}} \leq 1$$

and using the convexity of  $t \mapsto \frac{1}{\log t}$ , we obtain the upper bound

$$\pi(x) - \pi(y) \leq \frac{x-y}{2 \log y} + \frac{x-y}{2 \log x} + 2 \frac{E(x)}{\log y} \quad (2.17)$$

and the lower bound

$$\pi(x) - \pi(y) \geq \left(1 - \frac{2}{\sqrt{y}}\right) \frac{x-y}{\log \frac{x+y}{2}} - 2 \frac{E(x)}{\log y}. \quad (2.18)$$

For non-negative  $b$  and the trivial inequalities

$$\frac{b(p) \log p}{\log x} \leq b(p) \leq \frac{b(p) \log p}{\log y}$$

we similarly obtain the upper bound

$$\sum_{y < p \leq x} b(p) \leq \frac{1}{\log y} \int_y^x b(t) dt + \|b\|_{\text{TV}^*(y,x]} \frac{E(x)}{\log y} \quad (2.19)$$

and the lower bound

$$\sum_{y < p \leq x} b(p) \leq \frac{1 - \frac{2}{\sqrt{y}}}{\log x} \int_y^x b(t) dt - \|b\|_{\text{TV}^*(y,x]} \frac{E(x)}{\log x}. \quad (2.20)$$

One can also replace all occurrences of  $E(x)$  here by the classical error term  $O(x \exp(-c \sqrt{\log x}))$  for some absolute constant  $c > 0$  (in which case the  $\frac{2}{\sqrt{t}}$  type terms can be absorbed into the error term).

We remark that the accuracy in (2.14), (2.16) in particular is on par with what would be provided by the Riemann hypothesis, as long as  $x$  is not too large (e.g.,  $x \leq 10^{18}$ ). The other estimates in this lemma are not quite as precise, but still adequate for our applications. The error term  $E(x)$  can be improved somewhat for large  $x$  (see (B.3)), but this simplified version will suffice for our analysis (in particular, the contribution of the second term in (2.15) will be negligible for our applications). We make the easy remark that  $E(x)$  is non-decreasing in  $x$ , while  $E(x)/x$  is non-increasing.

### 3. GREEDY ALGORITHMS

Recall that  $t(N)$  can be interpreted as the largest  $t$  for which one has  $M(N, t) \geq N$ . Because of this, any algorithm that can produce lower bounds on  $M(N, t)$ , can also produce lower bounds on  $t(N)$ , as follows:

- Step 1. Use some heuristics to start with an initial proposed lower bound  $t$  for  $t(N)$ .
- Step 2. Use the provided lower bound algorithm for  $M(N, t)$  to test if  $M(N, t) \geq N$ .
- Step 3. If this algorithm succeeds (in a reasonable amount of time), then either HALT or increase  $t$  by some amount (possibly guided by the extent to which  $M(N, t)$  exceeds  $N$ ), and return to Step 2.
- Step 4. If instead the algorithm fails (or times out), then either HALT (with an error) or else decrease  $t$  by some amount (again possibly guided by the extent to which  $M(N, T)$  falls short of  $N$ ) and return to Step 2.

One can similarly use an algorithm that can produce upper bounds for  $M(N, t)$  to produce upper bounds for  $t(N)$ .

The algorithm described above is imprecisely specified, because it requires one to make some implementation decisions about how to select the parameter  $t$  at various steps of the algorithm, and also when to abandon the  $M(N, t)$  algorithm in case it takes an excessive amount of time to run. In particular, having some accurate heuristics (or “hints”) about what the correct value of  $t(N)$  should be (possibly based on the outcomes of previous stages of the algorithm) can greatly accelerate its performance. But regardless of this variability in speed, the  $M(N, t)$  algorithm will in practice produce a certificate (e.g., an explicit subfactorization of  $N!$ ) that can be quickly and independently verified by a separate computer program to confirm the upper or lower bound on  $t(N)$ . So the output of such imprecisely specified algorithms can at least be independently confirmed, if not reproduced exactly. In particular, the lack of reproducibility is not a major concern when verifying a specific bound on  $t(N)$ , such as  $t(N) \geq N/3$ , so long as independently verifiable proof certificates (such as an  $N/3$ -admissible subfactorization of  $N!$  of length  $N$ ) are generated.

We therefore turn to the question of how to algorithmically obtain good upper and lower bounds on  $M(N, t)$ . In this section we will discuss greedy methods to obtain lower bounds on this quantity; in the next section we will discuss how linear programming and integer programming methods can also be used to obtain both upper and lower bounds on  $M(N, t)$ .

The following simple greedy algorithm is a fast method to give reasonably good lower bounds on  $M(N, t)$ :

- (0) Initialize  $\mathcal{B}$  to be the empty multiset.
- (1) If  $\mathcal{B}$  is not a factorization, locate the largest prime  $p$  which is currently in surplus:  $v_p(N! / \prod \mathcal{B}) > 0$ .
- (2) If  $N! / \prod \mathcal{B}$  contains a multiple of  $p$  that is greater than or equal to  $t$ , locate the smallest such multiple, add it to  $\mathcal{B}$ , and return to Step 1. Otherwise, HALT the algorithm.

This procedure clearly halts in finite time to produce a  $t$ -admissible subfactorization of  $N!$ , with the length of this subfactorization giving a lower bound on  $M(N, t)$  (which can then lead to lower bounds on  $t(N)$  as discussed above). For instance, applying this procedure with  $N = 9$ ,  $t = 3$  produces the 3-admissible subfactorization

$$\{7 \times 1, 5 \times 1, 3 \times 1, 3 \times 1, 3 \times 1, 3 \times 1, 2 \times 2, 2 \times 2, 2 \times 2\}$$

which recovers the bound  $M(9, 3) \geq 9$  (and hence  $t(9) \geq 3$ ) from Example 1.1, albeit with a slightly different subfactorization, in which the 8 is replaced by 4.

This procedure is efficient for small  $N$ , for instance attaining the exact value of  $t(N)$  for all  $N \leq 79$ , though it begins to degrade for larger  $N$ ; see Figure 9. The performance is also respectable (though not optimal) for medium  $N$ ; for instance, when  $N = 3 \times 10^5$  and  $t = N/3$ , it establishes the lower bound  $M(N, t) \geq N + 372$ , which is close to the exact value of  $M(N, t) = N + 455$  which we could establish by the linear programming methods of the next section (see (4.10)).

**discuss modifications to the algorithm to make it perform both faster and more accurately; talk about hints**

In order to get this algorithm to validate all  $8 \times 10^4 \leq N \leq 10^{11}$  on commodity hardware in a reasonable amount of time, two major modifications were implemented.

First, when trying to improve an inequality  $t(N) \geq \lambda N$  for all  $N$  in some range, one can avoid running the algorithm for every single  $N$  in that range by instead proving a stronger inequality on a sparse subset. Namely, if one can show that  $t(N_0) \geq (\lambda + \epsilon)N_0$ , it follows that  $t(N) \geq \lambda N$  for all  $N \in \left[N_0, (1 + \frac{\epsilon}{\lambda})N_0\right)$ . If one can fix  $\epsilon$ , this reduces a brute force check of every value in a range of length  $L$  to checking just  $\log_{1+\frac{\epsilon}{\lambda}}(L)$  values. From the estimates in Section 1, one would expect to be able to take  $\epsilon = \frac{1}{e} - \lambda$  asymptotically; in practice the algorithm uses slightly smaller values.

The other major modification is related to Step (2) in the above algorithm, where one is searching for the least value of  $c$  such that  $cp \geq t$  and  $c$  can be constructed from the remaining factors. In practice, the algorithm pre-computes and store information about all such candidates; absent any further heuristics, this amounts to storing information about all integers less than  $N$ , which becomes prohibitively expensive for large  $N$ . The key observation is that any such  $c$  must satisfy a further arithmetic condition, namely that  $\frac{cL(c)}{S(c)} < t$ . By only storing information about  $c$  which satisfy this condition, the memory footprint is reduced by several orders of magnitude.

By using the greedy method, Theorem 1.3(ii) can be verified for  $N \leq 3 \times 10^5$ , and Theorem 1.3(iii) can be verified for  $67425 \leq N \leq 10^{12}$  **provide more details and links to code**. Thus, to resolve these claims, it remains to only establish Theorem 1.3(iii) in the regime  $43632 \leq N < 67425$  and  $N > 10^{12}$ , and also to show that this claim fails for  $N = 43631$ .

## 4. LINEAR PROGRAMMING

It turns out that linear programming and integer programming methods are quite effective at upper and lower bounding  $M(N, t)$ . The starting point is the following integer program interpretation of  $M(N, t)$ . For any  $t, N$ , let  $J_{t,N}$  be the collection of all  $j \geq t$  that divide  $N!$ , and which do not have any proper factor  $j' < j$  that is also greater than or equal to  $t$ . For instance,

$$J_{4,5} = \{4, 5, 6, 9\}.$$

**Proposition 4.1** (Integer programming description of  $t(N)$ ). *For any  $N, t \geq 1$ ,  $M(N, t)$  is the maximum value of*

$$\sum_{j \in J_{t,N}} m_j \tag{4.1}$$

where the  $m_j$  are non-negative integers subject to the constraints

$$\sum_{j \in J_{t,N}} m_j v_p(j) \leq v_p(N!) \tag{4.2}$$

for all primes  $p \leq N$ .

*Proof.* If  $m_j, j \in J_{t,N}$  are non-negative numbers obeying (4.2), then clearly

$$\prod_{j \geq t} j^{m_j} \tag{4.3}$$

is a  $t$ -admissible subfactorization of  $N!$ , so that  $M(N, t)$  is greater than or equal to (4.1). Conversely, suppose that  $M(N, t) \geq M$ , thus we have a  $t$ -admissible subfactorization of  $N!$  into  $M$  factors. Clearly, each of these factors  $j$  is at least  $t$ , and divides  $N!$ . If one of these factors  $j$  has a proper factor  $j' < j$  that is greater than or equal to  $t$ , then we can replace the factor  $j$  by the factor  $j'$  in the subfactorization, and still obtain a  $t$ -admissible subfactorization of  $N!$ . Iterating this, we may assume without loss of generality that all the factors  $t$  lie in  $J_{t,N}$ . We can then express this subfactorization as a product (4.3), and by computing  $p$ -valuations we conclude the constraints (4.2). The claim follows.  $\square$

This integer program formulation can be used, when combined with standard packages such as Gurobi or Mojo **give more details here**, to compute  $M(N, t)$  (and hence  $t(N)$ ) precisely for any specific  $N, t$  with  $N$  as large as  $10^4$ , though in practice it is better to first use faster methods (which we discuss below) to control these quantities first, using integer programming as a last resort when these faster methods fail to achieve the desired result.

For larger  $N$ , the sets  $J_{t,N}$  become somewhat large, and the integer program becomes computationally expensive. For the purposes of lower bounding  $M(N, t)$ , one can arbitrarily replace  $J_{t,N}$  with a smaller set (effectively setting  $m_j = 0$  for all  $j$  outside this set) to speed up the integer program; empirically we have found that the set  $\{j : t \leq j \leq N\}$  is a good choice, as it appears to give the same bounds while being significantly faster.

For upper bounds, we can relax the integer program to a linear program. Let  $M_{\mathbb{R}}(N, t)$  denote the maximum value of (4.1) where the  $m_j, j \in J_{t,N}$  are now non-negative *real* numbers

obeying (4.2). Clearly we have the upper bound

$$M(N, t) \leq M_{\mathbb{R}}(N, t)$$

which can be improved slightly to

$$M(N, t) \leq \lfloor M_{\mathbb{R}}(N, t) \rfloor \quad (4.4)$$

since  $M(N, t)$  is an integer. We refer to these bounds as the *linear programming upper bounds*.

The quantity  $M_{\mathbb{R}}(N, t)$  can be computed by standard linear programming methods; in particular, upper bounds on  $M_{\mathbb{R}}(t, N)$  can be obtained by solving a dual linear program involving some weights  $w_p, p \leq N$  that obey constraints for each  $j \in J_{t,N}$ . In fact we can restrict attention to those constraints with  $j$  in the range  $t \leq j \leq N$ :

**Proposition 4.2** (Dual description of  $M_{\mathbb{R}}(N, t)$ ). *For any  $N, t \geq 1$ ,  $M_{\mathbb{R}}(N, t)$  is the minimum value of*

$$\sum_{p \leq N} w_p v_p(N!) \quad (4.5)$$

where  $w_p$  are non-negative reals for primes  $p \leq N$  subject to the constraints that the  $w_p$  are weakly decreasing, thus

$$w_{p_2} \leq w_{p_1} \quad (4.6)$$

whenever  $p_2 \geq p_1$ , and

$$\sum_{p \leq N} w_p v_p(j) \geq 1 \quad (4.7)$$

for all  $t \leq j \leq N$ . In particular, if

$$\sum_{p \leq N} w_p v_p(N!) < N \quad (4.8)$$

then  $t(N) < t$ .

*Proof.* Suppose first that  $w_p$  are non-negative reals obeying (4.7) for all  $t \leq j \leq N$ . We claim that (4.7) in fact holds for all  $j \geq t$ , not just for  $t \leq j \leq N$ . Indeed, if this were not the case, consider the first  $j \geq t$  where (4.7) fails. Take a prime  $p$  dividing  $j$  and replace it by a prime in the interval  $[p/2, p)$  which exists by Bertrand's postulate (or remove  $p$  entirely, if  $p = 2$ ); this creates a new  $j'$  in  $[j/2, j)$  which is still at least  $t$ . By the weakly decreasing hypothesis on  $w_p$ , we have

$$\sum_p w_p v_p(j) \geq \sum_p w_p v_p(j')$$

and hence by the minimality of  $j$  we have

$$\sum_p w_p v_p(j) > 1,$$

a contradiction.

Now let  $m_j, j \in J_{t,N}$  be non-negative reals obeying (4.2). Multiplying each constraint in (4.2) by  $w_p$  and summing, we conclude from (4.7) that

$$\sum_{j \in J_{t,N}} m_j \leq \sum_{j \in J_{t,N}} m_j \sum_{p \leq N} w_p v_p(j) \leq \sum_{p \leq N} v_p(N!)$$

and hence (4.5) is an upper bound for  $M_{\mathbb{R}}(N, t)$ .

In the converse direction, we need to locate weakly decreasing non-negative weights  $w_p$  obeying (4.7) for  $t \leq j \leq N$  for which

$$\sum_{p \leq N} w_p v_p(N!) = M_{\mathbb{R}}(N, t). \quad (4.9)$$

To do this, we first make the technical observation that in the definition of  $M_{\mathbb{R}}(N, t)$ , we can enlarge the index set  $J_{t,N}$  to the larger set  $J'_{t,N}$  of natural numbers  $j \geq t$  that divide  $N!$ . This follows by repeating the proof of Proposition 4.1: if  $m_j$  were non-zero for some  $j \geq t$  dividing  $N!$  that had a proper factor  $j' \geq t$ , then one could transfer the mass of  $m_j$  to  $m_{j'}$  (i.e., replace  $m_{j'}$  with  $m_{j'} + m_j$  and then set  $m_j$  to zero) without affecting (4.2).

If we then invoke the duality theorem of linear programming, we can find weights  $w_p \geq 0$  for  $p \leq N$  obeying (4.9) as well as (4.7) for all  $j \in J'_{t,N}$  (not just  $j \in J_{t,N}$ ). To conclude the proof, it suffices to show that the  $w_p$  are weakly decreasing. Suppose for contradiction that there are primes  $p_1 < p_2 \leq N$  such that  $w_{p_1} > w_{p_2}$ . Let  $\varepsilon > 0$  be a sufficiently small quantity, and define the modification  $\tilde{w}_p$  to  $w_p$  by decreasing  $w_{p_1}$  by  $\varepsilon$  and leaving all other  $w_p$  unchanged. This decreases the left-hand side of (4.9), so to get a contradiction with the already-obtained lower bound, it suffices to show that

$$\sum_p \tilde{w}_p v_p(j) \geq 1$$

for all  $j \geq t$  dividing  $N!$ . If  $j$  has a proper factor  $j'$  that is still at least  $t$ , the condition for  $j$  would follow from that of  $j'$ , so we may restrict attention to the case where  $j$  has no proper factor greater than or equal to  $t$ . We can assume that  $j$  is divisible by  $p_1$ , otherwise the claim follows from (4.7). For  $\varepsilon$  small enough, one has

$$\sum_p \tilde{w}_p v_p(j) \geq \sum_p w_p v_p\left(\frac{p_2}{p_1}j\right);$$

since  $\frac{p_2}{p_1}j \geq j \geq t$ , we are done unless  $\frac{p_2}{p_1}j$  is not divisible by  $N!$ . This only occurs when  $v_{p_2}(j) = v_{p_2}(N!)$ , but then

$$\frac{j}{p_1} \geq p_2^{v_{p_2}(N!)} \geq p_2^{\lfloor N/p_2 \rfloor}.$$

We claim that the right-hand side is at least  $N/2$ . This is clear for  $p_2 \geq N/2$ , and also for  $\sqrt{N/2} \leq p_2 < N/2$  since  $\lfloor N/p_2 \rfloor \geq 2$  in this case. For  $p_2 < \sqrt{N/2}$  one has

$$p_2^{\lfloor N/p_2 \rfloor} \geq 3^{\lfloor \sqrt{2N} \rfloor} \geq \frac{N}{2}$$

for all  $N$  (here we use that  $3^k \geq \frac{(k+1)^2}{4}$  for  $k \geq 1$ ). Thus in all cases we have  $j/p_1 \geq N/2 \geq t$ , contradicting the hypothesis that  $j$  has no proper factor that is at least  $t$ .  $\square$

Proposition 4.2 allows for a fast method to compute  $M_{\mathbb{R}}(N, t)$  by a linear program. In practice, we have found that even if we drop the explicit constraint (4.6) that the  $w_p$  are weakly decreasing (or equivalently, if we return to the original primal problem of optimizing (4.1) for real  $m_j \geq 0$  obeying (4.2), but now with  $j$  restricted to  $t \leq j \leq N$ ), the optimal weights  $w_p$  produced by the resulting linear program will be weakly decreasing anyway, although we



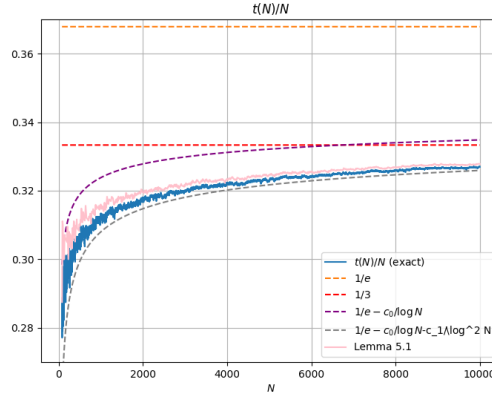


FIGURE 5. Exact values of  $t(N)/N$  for  $80 \leq N \leq 10^4$ , obtained via integer programming. The upper bound from Lemma 5.1 is surprisingly sharp, as is the refined asymptotic  $1/e - c_0/\log N - c_1/\log^2 N$ , though the cruder asymptotics  $1/e$  or  $1/e - c_0/\log N$  are significantly poorer approximations.

could not prove this empirically observed fact rigorously. For instance, when  $N = 3 \times 10^5$  and  $t = N/3$ , this linear program produces non-decreasing weights which certify that

$$M_{\mathbb{R}}(N, t) = N + 445.83398 \dots$$

and hence by (4.4)

$$M(N, t) \leq N + 445 \quad (4.10)$$

for this choice of  $N, t$ . In fact, as discussed later in this section, we know that equality holds in this particular case. For  $N \leq 10^4$ , we found that the linear programming upper bound on  $t(N)$  is in fact tight except when  $N = 155, 765, 1528, 1618, 1619, 2574, 2935, 3265, 5122, 5680, 9633$ , in which case integer programming was needed to precisely compute  $t(N)$ . The values of  $t(N)$  thus computed are plotted in Figure 5.

With integer programming, we could also establish<sup>5</sup>  $t(N) \geq N/3$  for all  $43632 \leq N \leq 8 \times 10^4$ . In particular, when combined with the greedy algorithm computations from the previous section, this resolves Theorem 1.3(ii), (iii) except in the asymptotic range  $N > 10^{12}$ , where it suffices to establish the lower bound  $t(N) \geq N/3$ .

For  $N \geq 10^4$ , the integer programming method to lower bound  $M(N, t)$  becomes slow. We found two faster methods to give slightly weaker lower bounds on this quantity, which we call the “floor+residuals” method, and the “smooth factorization” method.

The “floor+residuals” method proceeds by first running the primal linear program to find the real  $m_j \geq 0$  for  $t \leq j \leq N$  that maximize (4.1) subject to (4.2). The integer parts  $\lfloor m_j \rfloor$  will then of course also obey (4.2) and thus form a subfactorization; but this subfactorization is somewhat inefficient because there can be a  $p$ -surplus of  $v_p(N!) - \sum_{j \geq t} \lfloor m_j \rfloor v_p(j)$  at various primes  $p \leq N$ . We then apply the greedy algorithm of the previous section to fashion as many

<sup>5</sup>Explicit factorizations in this range can be found at <https://github.com/teorth/erdos-guy-selfridge/tree/main/Data/factorizations>.

factors greater than or equal to  $t$  from these residual primes, to obtain our final subfactorization that provides a lower bound on  $M(N, t)$ .

The floor+residuals method is fast and highly accurate for small and medium  $N$  (e.g.,  $N \leq 3 \times 10^5$ ). For instance:

- The method computes  $t(N)$  exactly for all  $N \leq 600$ , with the sole exception of  $N = 155$ ; see Figure 9. When  $N = 155$ , the floor+residuals method provides a subfactorization that certifies  $t(155) \geq 45$ , while the linear programming upper bound (4.4) gives  $t(155) \leq 46$ . Integer programming can then be deployed to confirm  $t(155) = 45$ .
- The method also verifies this lower bound for  $N = 41006$ , while the linear programming upper bound (4.4) shows that  $t(N) \geq N/3$  fails for all smaller  $N$  except for  $N = 1, 2, 3, 4, 5, 6, 9$ ; see Figure 7.
- The method establishes the matching lower bound

$$M(N, t) \geq N + 455 \tag{4.11}$$

to (4.10) when  $N = 3 \times 10^5$  and  $t = N/3$ .

For larger  $N$  (e.g.,  $10^4 \leq N \leq 9 \times 10^8$ ), the floor+residuals method becomes slow due to the large number  $\pi(N)$  of variables  $w_p$  that are involved in the linear program. We developed a *smooth factorization lower bound* method<sup>6</sup> to handle this range, by first using the greedy approach from the previous section to allocate all factors involving  $p \geq \sqrt{N}$ , and then using a version of the floor+residuals method to handle the smaller primes  $p < \sqrt{N}$  (with  $j$  now restricted to “smooth” numbers - numbers whose prime factors are less than  $\sqrt{N}$ ). Thus, the linear program now involves only  $\pi(\sqrt{N})$  variables  $w_p$ , and runs considerably faster in ranges such as  $10^4 < N \leq 9 \times 10^8$ . The lower bounds obtained by this method remain quite close to the linear programming upper bound (4.4), and outperforms the greedy algorithm; see Table 2 and Figure 6.

## 5. SOME UPPER BOUNDS

It is easy to check using (2.1) that the weights  $w_p := \frac{\log p}{\log t}$  will obey the conditions (4.8), (4.7) as long as  $0 > \log N! - N \log t$ . This recovers the trivial upper bound (1.2). By adjusting these weights at large primes, one can improve this bound as follows:

**Lemma 5.1** (Upper bound criterion). *Suppose that  $1 \leq t \leq N$  are such that*

$$\sum_{\substack{t \\ \lfloor \sqrt{t} \rfloor < p \leq N}} f_{N/t}(p/N) > \log N! - N \log t, \tag{5.1}$$

where  $f_{N/t}$  was defined in (1.7). Then  $t(N) < t$ .

<sup>6</sup>Details of this algorithm may be found at <https://github.com/teorth/erdos-guy-selfridge/tree/main/src/mojo>.

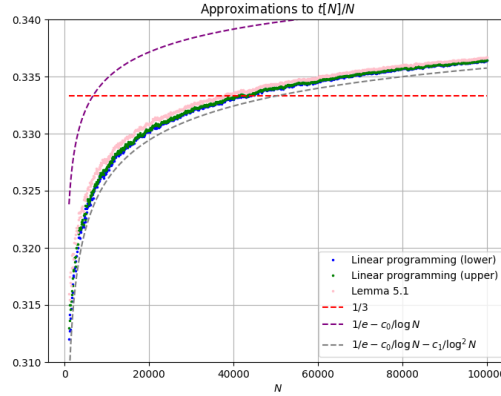


FIGURE 6. Upper and lower bounds  $t(N)/N$  obtained by linear programming upper bound and smooth factorization lower bound for  $10^3 \leq N \leq 10^5$  that are multiples of 100. The refined asymptotic  $1/e - c_0/\log N - c_1/\log^2 N$  is now a slight underestimate, hinting at further terms in the asymptotic expansion. For another view of the situation near the crossover point  $N = 43632$  for Theorem 1.3(iii), see Figure 7.

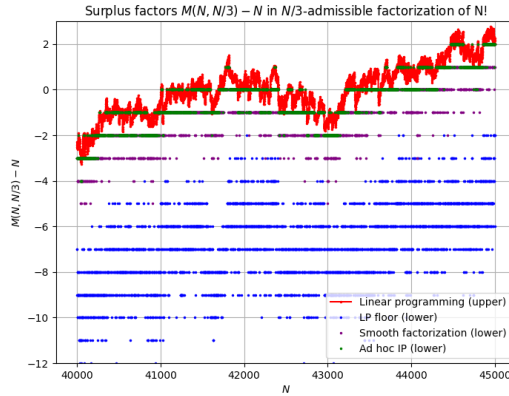


FIGURE 7. Bounds on  $M(N, N/3) - N$  for  $4 \times 10^4 \leq N \leq 4.5 \times 10^4$ . The linear programming upper bound (red) can be rounded down to the nearest integer, as per (4.4). Integer programming (green), implemented in an *ad hoc* fashion, provides lower bounds that almost always matches the upper bound (rounded down) at this range of  $N$ ; in particular it verifies  $t(N) \geq N/3$  for  $43632 \leq N \leq 4.5 \times 10^4$ . The smooth factorization method (purple) gives slightly worse bounds than the (slower) *ad hoc* integer programming methods, but slightly better than the (faster) linear algebra floor method (blue) **TODO: update this to floor+residuals.**

*Proof.* We introduce the weights

$$w_p := \begin{cases} \frac{\log p}{\log t} & p \leq \frac{t}{\lfloor \sqrt{t} \rfloor} \\ \frac{\log p}{\log t} - \frac{\log \lceil t/p \rceil}{\log t} = 1 - \frac{\log \lceil t/p \rceil}{\log t} & p > \frac{t}{\lfloor \sqrt{t} \rfloor}. \end{cases}$$

$N$	$t(N)_-$	$t(N)_+$	Lemma 5.1	$\frac{N}{e} - \frac{c_0 N}{\log N} - \frac{c_1 N}{\log^2 N}$
$1 \times 10^5$	33 642	$t(N)_- + 4$	$t(N)_- + 26$	$t(N)_- - 69$
$2 \times 10^5$	67 703	$t(N)_- + 1$	$t(N)_- + 36$	$t(N)_- - 130$
$3 \times 10^5$	101 903	$t(N)_- + 3$	$t(N)_- + 42$	$t(N)_- - 206$
$4 \times 10^5$	136 143	$t(N)_- + 6$	$t(N)_- + 43$	$t(N)_- - 248$
$5 \times 10^5$	170 456	$t(N)_- + 3$	$t(N)_- + 46$	$t(N)_- - 310$
$6 \times 10^5$	204 811	$t(N)_- + 4$	$t(N)_- + 47$	$t(N)_- - 373$
$7 \times 10^5$	239 187	$t(N)_- + 9$	$t(N)_- + 54$	$t(N)_- - 425$
$8 \times 10^5$	273 604	$t(N)_- + 6$	$t(N)_- + 64$	$t(N)_- - 490$
$9 \times 10^5$	308 029	$t(N)_- + 13$	$t(N)_- + 70$	$t(N)_- - 539$
$1 \times 10^6$	342 505	$t(N)_- + 3$	$t(N)_- + 62$	$t(N)_- - 619$
$2 \times 10^6$	687 796	$t(N)_- + 4$	$t(N)_- + 87$	$t(N)_- - 1180$
$3 \times 10^6$	1 033 949	$t(N)_- + 11$	$t(N)_- + 107$	$t(N)_- - 1736$
$4 \times 10^6$	1 380 625	$t(N)_- + 12$	$t(N)_- + 122$	$t(N)_- - 2286$
$5 \times 10^6$	1 727 605	$t(N)_- + 4$	$t(N)_- + 126$	$t(N)_- - 2763$
$6 \times 10^6$	2 074 962	$t(N)_- + 21$	$t(N)_- + 152$	$t(N)_- - 3326$
$7 \times 10^6$	2 422 486	$t(N)_- + 22$	$t(N)_- + 165$	$t(N)_- - 3819$
$8 \times 10^6$	2 770 212	$t(N)_- + 29$	$t(N)_- + 177$	$t(N)_- - 4316$
$9 \times 10^6$	3 118 129	$t(N)_- + 24$	$t(N)_- + 173$	$t(N)_- - 4834$
$1 \times 10^7$	3 466 235	$t(N)_- + 12$	$t(N)_- + 179$	$t(N)_- - 5392$
$2 \times 10^7$	6 952 243	$t(N)_- + 18$	$t(N)_- + 234$	$t(N)_- - 10 284$
$3 \times 10^7$	10 444 441	$t(N)_- + 13$	$t(N)_- + 253$	$t(N)_- - 14 975$
$4 \times 10^7$	13 940 484	$t(N)_- + 64$	$t(N)_- + 354$	$t(N)_- - 19 582$
$5 \times 10^7$	17 439 282	$t(N)_- + 33$	$t(N)_- + 356$	$t(N)_- - 24 124$
$6 \times 10^7$	20 940 210	$t(N)_- + 23$	$t(N)_- + 381$	$t(N)_- - 28 610$
$7 \times 10^7$	24 442 818	$t(N)_- + 37$	$t(N)_- + 415$	$t(N)_- - 32 996$
$8 \times 10^7$	27 946 958	$t(N)_- + 43$	$t(N)_- + 445$	$t(N)_- - 37 417$
$9 \times 10^7$	31 452 431	$t(N)_- + 23$	$t(N)_- + 428$	$t(N)_- - 41 882$
$1 \times 10^8$	34 958 725	$t(N)_- + 48$	$t(N)_- + 482$	$t(N)_- - 46 039$
$2 \times 10^8$	70 064 782	$t(N)_- + 45$	$t(N)_- + 644$	$t(N)_- - 87 837$
$3 \times 10^8$	105 218 403	$t(N)_- + 41$	$t(N)_- + 752$	$t(N)_- - 128 227$
$4 \times 10^8$	140 401 212	$t(N)_- + 80$	$t(N)_- + 887$	$t(N)_- - 167 495$
$5 \times 10^8$	175 605 266	$t(N)_- + 98$	$t(N)_- + 972$	$t(N)_- - 206 175$
$6 \times 10^8$	210 825 848	$t(N)_- + 68$	$t(N)_- + 1058$	$t(N)_- - 244 391$
$7 \times 10^8$	246 059 851	$t(N)_- + 89$	$t(N)_- + 1147$	$t(N)_- - 282 167$
$8 \times 10^8$	281 305 291	$t(N)_- + 1440$	$t(N)_- + 1158$	$t(N)_- - 319 700$
$9 \times 10^8$	316 560 601	$t(N)_- + 113$	$t(N)_- + 1238$	$t(N)_- - 357 029$

TABLE 2. For sample values of  $N \in [10^5, 9 \times 10^8]$ , the (remarkably precise) lower and upper bounds  $t(N)_- \leq t(N) \leq t(N)_+$  obtained by smooth factorization and linear programming respectively, the (slightly weaker) upper bound on  $t(N)$  from Lemma 5.1, and the conjectural approximation  $\frac{N}{e} - \frac{c_0 N}{\log N} - \frac{c_1 N}{\log^2 N}$  (rounded to the nearest integer).

Clearly the  $w_p$  are non-negative. It will suffice to verify the conditions (4.7), (4.8). If  $j \in J_{t,N}$  contains no prime factor  $p > \frac{t}{\lfloor \sqrt{t} \rfloor}$ , then from (2.1) we have

$$\sum_p w_p v_p(j) = \frac{\sum_p v_p(j) \log p}{\log t} = \frac{\log j}{\log t} \geq 1.$$

$N$	$t(N)_-$	Fast greedy	Time (s)	Exhaustive greedy	Time (s)
$1 \times 10^6$	342 505	$t(N)_- - 4863$	0.001	$t(N)_- - 3723$	0.118
$2 \times 10^6$	687 796	$t(N)_- - 7710$	0.001	$t(N)_- - 6874$	0.296
$3 \times 10^6$	1 033 949	$t(N)_- - 10 793$	0.002	$t(N)_- - 9791$	0.731
$4 \times 10^6$	1 380 625	$t(N)_- - 14 637$	0.003	$t(N)_- - 11 502$	1.654
$5 \times 10^6$	1 727 605	$t(N)_- - 20 847$	0.003	$t(N)_- - 15 778$	2.686
$6 \times 10^6$	2 074 962	$t(N)_- - 25 872$	0.002	$t(N)_- - 21 181$	4.337
$7 \times 10^6$	2 422 486	$t(N)_- - 30 334$	0.003	$t(N)_- - 22 348$	5.348
$8 \times 10^6$	2 770 212	$t(N)_- - 32 103$	0.003	$t(N)_- - 27 647$	6.474
$9 \times 10^6$	3 118 129	$t(N)_- - 30 721$	0.004	$t(N)_- - 30 069$	7.078
$1 \times 10^7$	3 466 235	$t(N)_- - 38 594$	0.007	$t(N)_- - 33 151$	10.576
$2 \times 10^7$	6 952 243	$t(N)_- - 80 533$	0.008	$t(N)_- - 67 701$	33.397
$3 \times 10^7$	10 444 441	$t(N)_- - 112 826$	0.013	$t(N)_- - 98 056$	87.062
$4 \times 10^7$	13 940 484	$t(N)_- - 109 009$	0.016	$t(N)_- - 104 535$	110.231
$5 \times 10^7$	17 439 282	$t(N)_- - 187 980$	0.026	$t(N)_- - 147 795$	204.029
$6 \times 10^7$	20 940 210	$t(N)_- - 246 858$	0.027	$t(N)_- - 228 610$	346.622
$7 \times 10^7$	24 442 818	$t(N)_- - 289 249$	0.018	$t(N)_- - 233 968$	505.893
$8 \times 10^7$	27 946 958	$t(N)_- - 245 135$	0.036	$t(N)_- - 231 144$	457.513
$9 \times 10^7$	31 452 431	$t(N)_- - 335 660$	0.042	$t(N)_- - 275 753$	657.308
$1 \times 10^8$	34 958 725	$t(N)_- - 337 459$	0.082	$t(N)_- - 292 563$	671.201
$2 \times 10^8$	70 064 782	$t(N)_- - 691 712$	0.144	$t(N)_- - 623 293$	2035.745
$3 \times 10^8$	105 218 403	$t(N)_- - 956 030$	0.092	$t(N)_- - 903 647$	3479.215
$4 \times 10^8$	140 401 212	$t(N)_- - 1 332 772$	0.111	$t(N)_- - 1 142 677$	5336.544
$5 \times 10^8$	175 605 266	$t(N)_- - 1 659 922$	0.303	$t(N)_- - 1 477 924$	7285.802
$6 \times 10^8$	210 825 848	$t(N)_- - 1 786 095$	0.322	$t(N)_- - 1 590 641$	8735.219
$7 \times 10^8$	246 059 851	$t(N)_- - 1 991 829$	0.258	$t(N)_- - 1 920 970$	10462.708

TABLE 3. For sample values of  $N \in [10^6, 7 \times 10^8]$ , the performance of the fast greedy algorithm (using heuristics to guess optimal choices of  $t$ ), and more thorough run of the algorithm (that exhaustively searches over  $t$  for which the algorithm works), compared against the lower bound  $t(N)_-$  obtained from the linear programming method. **need some details on what kind of machine these runs were made on, to make some sense of the run times.**

If  $j \in J_{t,N}$  is of the form  $j = mp_1$  where  $p_1 > \frac{t}{\lfloor \sqrt{t} \rfloor}$  and  $m$  contains no prime factor exceeding  $\frac{t}{\lfloor \sqrt{t} \rfloor}$ , then  $m \geq \lceil t/p_1 \rceil$ , and we have

$$\begin{aligned}
\sum_p w_p v_p(j) &= \frac{\sum_p v_p(j) \log p}{\log t} - \frac{\log \frac{\lceil t/p_1 \rceil}{t/p_1}}{\log t} \\
&= \frac{\log(mp_1)}{\log t} - \frac{\log \frac{m}{t/p_1}}{\log t} \\
&= 1.
\end{aligned}$$

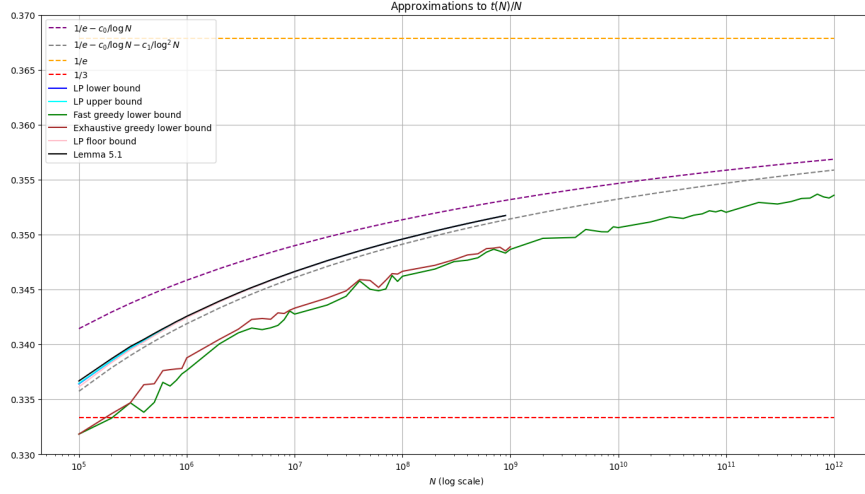


FIGURE 8. Lower bounds on  $t(N)/N$  coming from the linear programming and greedy methods in the range  $10^5 \leq N \leq 10^{12}$ , as well as upper bounds coming from linear programming and Lemma 5.1. Even if one simply takes floors from the linear program and discards residuals, the resulting lower bound is asymptotically almost indistinguishable from the upper bound. **would be nice to display data from other greedy methods for comparison**

Finally, if  $j \in J_{t,N}$  is divisible by two primes  $p_1, p_2 > t \lfloor \sqrt{t} \rfloor$  (possibly equal), then

$$\begin{aligned} \sum_p w_p v_p(j) &\geq 1 - \frac{\log \lceil t/p_1 \rceil}{\log t} + 1 - \frac{\log \lceil t/p_1 \rceil}{\log t} \\ &\geq 1 - \frac{\log \sqrt{t}}{\log t} + 1 - \frac{\log \sqrt{t}}{\log t} \\ &= 1. \end{aligned}$$

Thus we have verified (4.7) for all  $j \in J_{t,N}$ . Finally, from (2.1), (2.3), (5.1) we have

$$\begin{aligned} \sum_p w_p v_p(N!) &= \frac{\sum_p v_p(N!) \log p}{\log t} - \sum_{p > \frac{t}{\lfloor \sqrt{t} \rfloor}} \frac{v_p(N!) \log \frac{t/p}{t/p}}{\log t} \\ &\geq \frac{\log N!}{\log t} - \frac{\sum_{p > \frac{t}{\lfloor \sqrt{t} \rfloor}} \lfloor \frac{N}{p} \rfloor \log \frac{t/p}{t/p}}{\log t} \\ &= \frac{\log N!}{\log t} - \frac{\sum_{p > \frac{t}{\lfloor \sqrt{t} \rfloor}} f_{N/t}(p/N)}{\log t} \\ &< N, \end{aligned}$$

giving (4.5). The claim follows.  $\square$

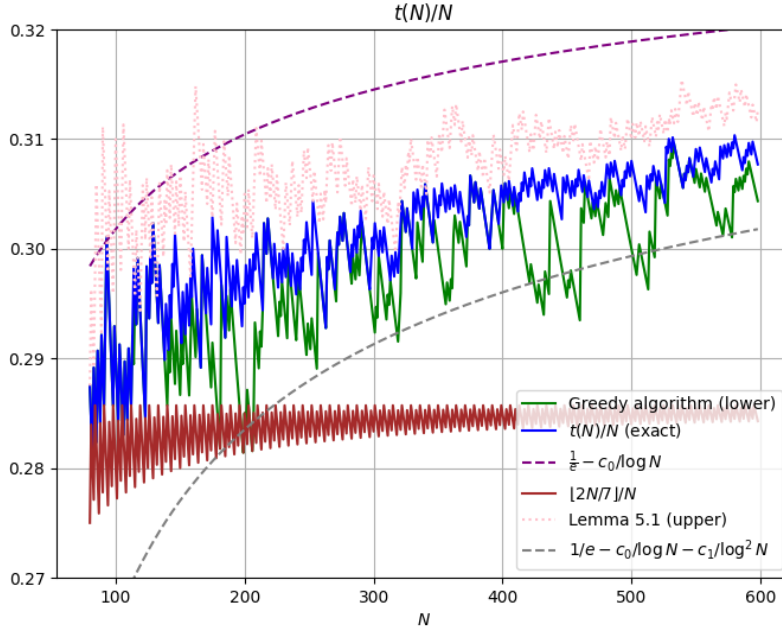


FIGURE 9. An enlarged version of Figure 2, displaying the lower bound from the greedy algorithm and the upper bound from Lemma 5.1. The linear programming upper bound and floor+residual bounds are exact in this region, except for  $N = 155$  in which the upper bound is off by one.

In practice, Lemma 5.1 gives quite good upper bounds on  $N$ , especially when  $N$  is large, although for medium  $N$  the linear programming method is superior: see Figure 1, Figure 2, Figure 9.

We can now prove the upper bound portion of Theorem 1.3(iv):

**Proposition 5.2.** *For large  $N$ , one has*

$$\frac{t(N)}{N} \leq \frac{1}{e} - \frac{c_0}{\log N} - \frac{c_1 + o(1)}{\log^2 N}$$

where

$$c'_1 := \frac{1}{e} \int_0^1 f_e(x) \log \frac{1}{x} dx = 0.3702015 \dots \quad (5.2)$$

$$c''_1 := \sum_{k=1}^{\infty} \frac{1}{k} \log \left( \frac{e}{k} \left\lceil \frac{k}{e} \right\rceil \right) \approx 1.6796 \quad (5.3)$$

$$c_1 := c'_1 + c_0 c''_1 - e c_0^2 / 2 \approx 0.75554808. \quad (5.4)$$

We discuss the numerical evaluation of these constants in Appendix C.

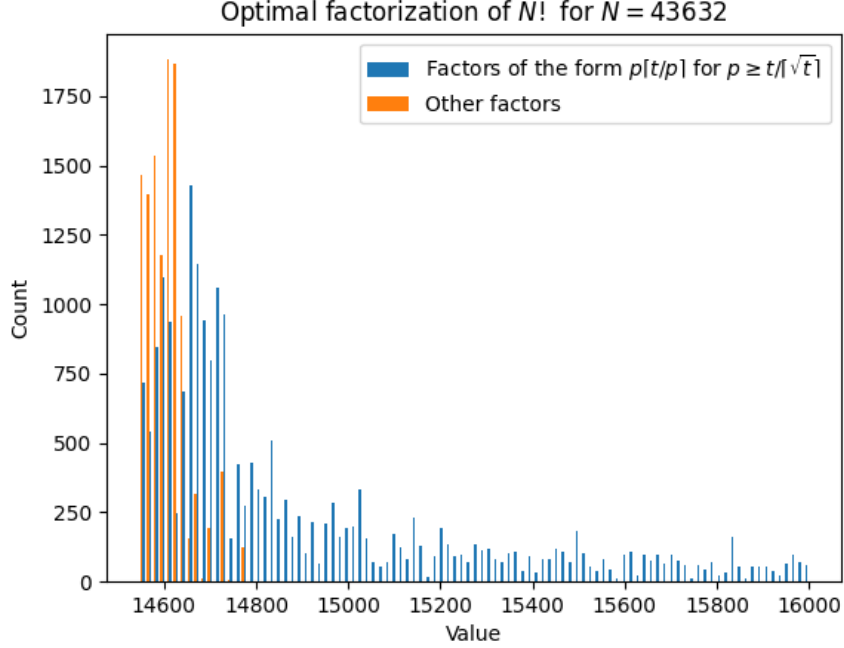


FIGURE 10. A histogram of an optimal factorization of  $N!$  for  $N = 43632$ , demonstrating that  $t(N) = N/3 + 1 = 14545$ . Of the  $N$  factors, most (32 147) of them are of the form  $p[t/p]$  for  $p > \frac{t}{\lfloor \sqrt{t} \rfloor}$ , thus directly contributing to the left-hand side of Lemma 5.1. (The factors arising from very large primes are lie to the right of the displayed graph, as a long tail of the histogram.) The remaining 11 485 factors stay close to  $t(N)$ , with the largest being 14 941.

Numerically, this bound is a reasonably good approximation for medium-sized  $N$ , see Figure 5, Figure 6, although it may be possible to improve the approximation further with additional terms. Based on these numerics it seems natural to conjecture that one in fact has

$$\frac{t(N)}{N} = \frac{1}{e} - \frac{c_0}{\log N} - \frac{c_1 + o(1)}{\log^2 N}$$

as  $N \rightarrow \infty$ .

*Proof.* We apply Lemma 5.1 with

$$t := \frac{1}{e} - \frac{c_0}{\log N} - \frac{c_1 - \varepsilon}{\log^2 N}$$

for a given small constant  $\varepsilon > 0$ . From Taylor expansion of the logarithm and the Stirling approximation (2.4) one sees that

$$\log N! - N \log t = ec_0 \frac{N}{\log N} + (ec_1 - \frac{1}{2}e^2 c_0^2 - e\varepsilon + o(1)) \frac{N}{\log^2 N}$$



so it will suffice to establish the lower bound

$$\sum_{\frac{t}{\lfloor \sqrt{t} \rfloor} < p \leq N} f_{N/t}(p/N) \geq ec_0 \frac{N}{\log N} + (ec_1 - \frac{1}{2}e^2c_0^2 - e\epsilon + o(1)) \frac{N}{\log^2 N} \quad (5.5)$$

for  $N$  sufficiently large depending on  $\epsilon$ .

For  $N$  large enough, we have  $\frac{t}{\lfloor \sqrt{t} \rfloor} \leq \frac{N}{\log^3 N}$ . On the interval  $[1/\log^3 N, 1]$ , the piecewise smooth function  $f_{N/t}$  is bounded by  $O(1)$  thanks to (1.9), and has an (augmented) total variation of  $O(\log^3 N)$ ; the same is then true for the rescaled function  $x \mapsto f_{N/t}(x/N)$  on  $[N/\log^3 N, N]$ . This implies that  $x \mapsto \frac{1}{\log x} f_{N/t}(x/N)$  has an (augmented) total variation of  $O(\log^2 N)$ . By Lemma 2.2 (with classical error term), we conclude that the left-hand side of (5.5) is at least

$$\int_{N/\log^3 N}^N f_{N/t}(x/N) \frac{dx}{\log x} + O\left(N \exp(-c\sqrt{\log N})\right)$$

for some  $c > 0$ . Performing a change of variable, we reduce to showing that

$$\int_{1/\log^3 N}^1 f_{N/t}(x) \frac{\log N}{\log(Nx)} dx \geq ec_0 + \frac{ec_1 - \frac{1}{2}e^2c_0^2 - e\epsilon + o(1)}{\log N}.$$

By Taylor expansion, we have

$$\frac{\log N}{\log(Nx)} = 1 + \frac{\log \frac{1}{x}}{\log N} + o\left(\frac{1}{\log N}\right)$$

and from dominated convergence we have

$$\int_{1/\log^3 N}^1 f_{N/t}(x) \log \frac{1}{x} dx = ec'_1 + o(1)$$

and hence by definition of  $c_1$ , we reduce to showing that

$$\int_{1/\log^3 N}^1 f_{N/t}(x) dx \geq ec_0 + \frac{ec_0c''_1 - e^2c_0^2 - e\epsilon + o(1)}{\log N}.$$

By performing a rescaling by  $N/et = 1 + \frac{ec_0+o(1)}{\log N}$ , the left-hand side may be written as

$$\left(1 - \frac{ec_0 + o(1)}{\log N}\right) \int_{N/et \log^3 N}^{N/et} \left\lfloor \frac{N/et}{x} \right\rfloor \log \left( ex \left\lceil \frac{1}{ex} \right\rceil \right)$$

so it will suffice to show that

$$\int_{N/et \log^3 N}^{N/et} \left\lfloor \frac{N/et}{x} \right\rfloor \log \left( ex \left\lceil \frac{1}{ex} \right\rceil \right) dx \geq ec_0 + \frac{ec_0c''_1 - e\epsilon + o(1)}{\log N}.$$

From (1.6), (1.9) we have

$$\int_{1/\log^2 N}^1 \left\lfloor \frac{1}{x} \right\rfloor \log \left( ex \left\lceil \frac{1}{ex} \right\rceil \right) = ec_0 - \frac{o(1)}{\log N}$$

it suffices to show that

$$\int_{1/\log^2 N}^{N/et} \left( \left\lfloor \frac{N/et}{x} \right\rfloor - \left\lfloor \frac{1}{x} \right\rfloor \right) \log \left( ex \left\lceil \frac{1}{ex} \right\rceil \right) dx \geq \frac{ec_0c''_1 - e\epsilon + o(1)}{\log N}.$$

Let  $K$  be sufficiently large depending on  $\varepsilon$ , then for  $N$  sufficiently large depending on  $K$  we can lower bound the left-hand side by

$$\sum_{k=1}^K \int_{1/k}^{N/etk} \log \left( ex \left\lceil \frac{1}{ex} \right\rceil \right) dx;$$

since  $\frac{N}{etk} = \frac{1}{k} + \frac{ec_0}{k \log N}$ , we can lower bound this (using the irrationality of  $e$ ) by

$$\frac{ec_0 + o(1)}{\log N} \sum_{k=1}^K \frac{1}{k} \log \left( \frac{e}{k} \left\lceil \frac{k}{e} \right\rceil \right)$$

for sufficiently large  $N$ . Since the sum here can be made arbitrarily close to  $c_0''$  by increasing  $K$ , we obtain the claim.  $\square$

We can now establish Theorem 1.3(i):

**Proposition 5.3.** *One has  $t(N)/N < 1/e$  for  $N \neq 1, 2, 4$ .*

*Proof.* From existing data on  $t(N)$  (or the linear programming method) one can verify this claim for  $N < 80$  (see Figure 1), so we assume that  $N \geq 80$ .

Applying Lemma 5.1 and (2.4), it suffices to show that

$$\sum_{p \geq \frac{N/e}{\lfloor \sqrt{N/e} \rfloor}} f_e(p/N) > \frac{1}{2} \log(2\pi N) + \frac{1}{12N}. \quad (5.6)$$

This may be easily verified numerically in the range  $80 \leq N \leq 5000$  (see Figure 11). We will discard the  $\lfloor \sqrt{N/e} \rfloor$  denominator, and reduce to showing

$$\sum_{N/e < p \leq N} f_e(p/N) > \frac{1}{2} \log(2\pi N) + \frac{1}{12N} \quad (5.7)$$

for  $N > 5000$ . On  $[1/e, 1]$ , one can compute

$$\|f_e\|_{\text{TV}^*(1/e, 1]} = 4 - 2 \log 2$$

so by Lemma 2.2 (noting that  $5000/e > 1423$ ) we have

$$\sum_{N/e < p \leq N} f_e(p/N) \geq \frac{N \left(1 - \frac{2}{\sqrt{N/e}}\right)}{\log N} \int_{1/e}^1 f_e(x) dx - (4 - 2 \log 2) \frac{E(N)}{\log N}$$

and so it suffices to show that

$$\left(1 - \frac{2}{\sqrt{N/e}}\right) \int_{1/e}^1 f_e(x) dx \geq (4 - 2 \log 2) \frac{E(N)}{N} + \frac{\log(2\pi N) \log N}{2N} + \frac{\log N}{12N^2}.$$

The right-hand side is increasing in  $N$  and the left-hand side is decreasing for  $N \geq 5000$ , so it suffices to verify this claim for  $N = 5000$ ; but this is a routine calculation (with plenty of room to spare; cf., Figure 11).  $\square$

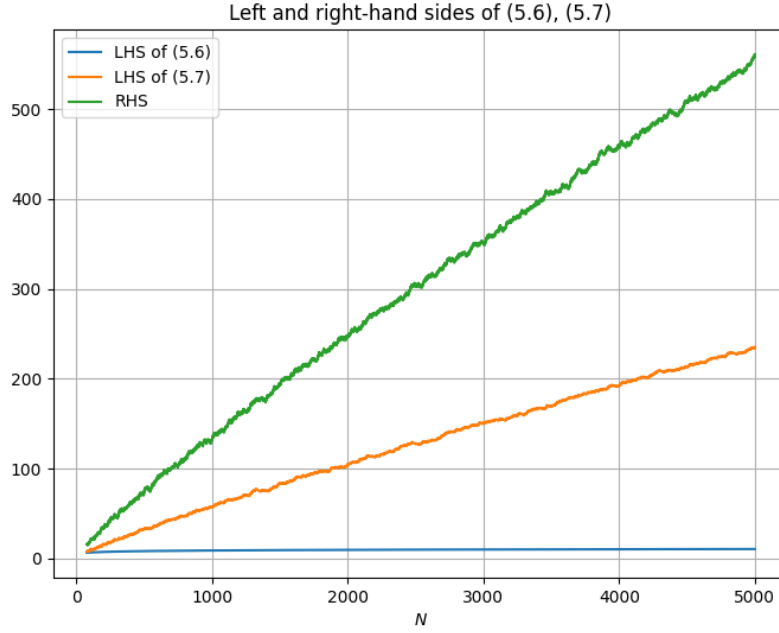


FIGURE 11. A plot of the left and right-hand sides of (5.6), (5.7) for  $80 \leq N < 5000$ .

## 6. REARRANGING THE STANDARD FACTORIZATION

In this section we describe an approach to establishing lower bounds on  $t(N)$  by starting with the standard factorization  $\{1, \dots, N\}$ , dividing out some small prime factors from some of the terms, and then redistributing them to other terms. This approach was introduced in [10] to give lower bounds of the shape  $\frac{t(N)}{N} \geq \frac{3}{16} + o(1)$  (by redistributing powers of two only) and  $\frac{t(N)}{N} \geq \frac{1}{4} + o(1)$  (by redistributing powers of two and three). With computer assistance, we are also able to show that  $\frac{t(N)}{N} \geq \frac{1}{3} + o(1)$  for sufficiently large  $N$ , in a simpler fashion than the method used to prove Theorem 1.3(iv) in the next section.

We need some notation. Define a *downset* to be a finite set  $\mathcal{D}$  of natural numbers, with the property that if  $d \in \mathcal{D}$  then all factors of  $d$  lie in  $\mathcal{D}$ . If  $\mathcal{D}$  is a downset and  $d \in \mathcal{D}$ , we define the density  $\sigma_{d,\mathcal{D}}$  to be the product

$$\sigma_{d,\mathcal{D}} \prod_{p < P_+(d) \text{ or } pd \in \mathcal{D}} \left(1 - \frac{1}{p}\right)$$

where  $P_+(d)$  is the largest prime factor of  $d$  (with the convention  $P_+(1) = 1$ ). We will later establish the identity

$$\sum_{d \in \mathcal{D}} \frac{\sigma_{d,\mathcal{D}}}{d} = 1. \quad (6.1)$$

**Example 6.1.** The set  $\mathcal{D} = \{1, 2, 4\}$  is a downset, with  $\sigma_{1,\mathcal{D}} = \sigma_{2,\mathcal{D}} = \frac{1}{2}$  and  $\sigma_{4,\mathcal{D}} = 1$ . The set  $\mathcal{D}' = \{1, 2, 3, 4\}$  is also a downset with  $\sigma_{1,\mathcal{D}'} = \frac{1}{3}$ ,  $\sigma_{2,\mathcal{D}'} = \frac{1}{2}$ ,  $\sigma_{3,\mathcal{D}'} = \frac{1}{2}$ ,  $\sigma_{4,\mathcal{D}'} = 1$ . The identity

(6.1) becomes

$$\frac{1/2}{1} + \frac{1/2}{2} + \frac{1}{4} = 1$$

for the former downset and

$$\frac{1/3}{1} + \frac{1/2}{2} + \frac{1/2}{3} + \frac{1}{4} = 1$$

for the latter downset.

**Proposition 6.2** (Criterion for asymptotic lower bound). *Let  $\mathcal{D}$  be a downset that contains at least one prime  $p_0$ , let  $0 < \alpha < 1$ , and suppose we have non-negative reals  $a_\ell$  for all natural numbers  $\ell$  obeying the following conditions:*

(i) *For all primes  $p$ , one has*

$$\sum_{\ell=1}^{\infty} v_p(\ell) a_\ell \leq \sum_{d \in \mathcal{D}} v_p(d) \frac{\sigma_{d, \mathcal{D}}}{d}. \quad (6.2)$$

(ii) *For every natural number  $\ell$ , we have*

$$\sum_{\ell' > \ell} a_{\ell'} > \sum_{d \in \mathcal{D}} \sigma_{d, \mathcal{D}} \min\left(\frac{1}{d}, \frac{\alpha}{\ell}\right). \quad (6.3)$$

*Then  $t(N) \geq \alpha N$  for all sufficiently large  $N$ .*

Informally,  $\mathcal{D}$  represents the factors that one removes (as greedily as possible) from the standard factorization  $\{1, \dots, N\}$  of  $N!$  to free up some prime factors, and  $a_\ell$  is the proportion of elements in the resulting subfactorization that are to be multiplied by  $\ell$  (again, in a greedy fashion) to (hopefully) bring the subfactorization into  $t$ -admissibility. The condition (6.2) asserts that enough primes are freed up by the first step to “afford” the second step, while (6.3) is the assertion that the  $a_\ell$  have enough mass at large  $\ell$  to bring even the smallest elements of the subfactorization  $t$ -admissible.

*Proof.* We first make a small technical modification to the sequence  $a_\ell$ . If  $p \in \mathcal{D}$ , then the right-hand side of (6.2) is positive. If equality holds here, then one of the  $a_\ell$  with  $\ell$  divisible by  $p$  is positive; but one can reduce this quantity slightly without violating (6.3), since  $a_\ell$  only impacts finitely many cases of these strict inequalities. Thus, we may assume without loss of generality that the inequality (6.2) is strict for all  $p \in \mathcal{D}$ .

From (6.2) and (2.1) we see that

$$\sum_{\ell=1}^{\infty} a_\ell \log \ell \leq \sum_{d \in \mathcal{D}} v_p(d) \frac{\sigma(d)}{d}.$$

In particular, from the dominated convergence theorem we have

$$r \sum_{\ell \geq p_0^r} a_\ell \rightarrow 0 \quad (6.4)$$

as  $r \rightarrow \infty$ .

Let  $N$  be sufficiently large. It will suffice to obtain an  $\alpha N$ -admissible subfactorization  $\mathcal{B}'$  of  $N!$  of cardinality  $N$ .

Every natural number can be uniquely factored as  $n = p_1 \dots p_r$  for some primes  $p_1 \leq \dots \leq p_r$ . If we let  $d := p_1 \dots p_k$  be the largest initial segment of this factorization that lies in  $\mathcal{D}$ , then we have  $n = dm$  where  $m$  lies in the set  $A_{d,\mathcal{D}}$  of natural numbers not divisible by any prime  $p$  with  $p < P_+(d)$  or  $pd \in \mathcal{D}$ . Conversely, every pair  $d, m$  with  $d \in \mathcal{D}$  and  $m \in A_{d,\mathcal{D}}$  arises exactly once in this manner. We thus have the partition

$$\mathbb{N} = \bigsqcup_{d \in \mathcal{D}} d \cdot A_{d,\mathcal{D}}.$$

From the Chinese remainder theorem, the  $A_{d,\mathcal{D}}$  have density  $\sigma_{d,\mathcal{D}}$ , in the sense that

$$|A_{d,\mathcal{D}} \cap I| = \sigma_{d,\mathcal{D}}|I| + O(1)$$

for all intervals  $I$ , where we allow the implied constants to depend on the downset  $\mathcal{D}$ . Summing the densities, we conclude (6.1).

From this partition, we also have

$$\{1, \dots, N\} = \bigsqcup_{d \in \mathcal{D}} d \cdot (A_{d,\mathcal{D}} \cap [1, N/d])$$

and on multiplying we obtain a factorization

$$N! = \prod_{d \in \mathcal{D}} d^{|A_{d,\mathcal{D}} \cap [1, N/d]|} \times \prod \mathcal{B}$$

where  $\mathcal{B}$  is the multiset

$$\mathcal{B} := \bigsqcup_{d \in \mathcal{D}} \left( A_{d,\mathcal{D}} \cap \left[ 1, \frac{N}{d} \right] \right).$$

Thus  $\mathcal{B}$  has cardinality  $|\mathcal{B}| = N$ , and is a subfactorization of  $N!$  with surplus

$$\begin{aligned} v_p \left( \frac{N!}{\prod \mathcal{B}} \right) &= \sum_{d \in \mathcal{D}} v_p(d) \left| A_{d,\mathcal{D}} \cap \left[ 1, \frac{N}{d} \right] \right| \\ &= N \sum_{d \in \mathcal{D}} v_p(d) \frac{\sigma_{d,\mathcal{D}}}{d} + O(1). \end{aligned} \tag{6.5}$$

In particular this multiset is in balance for all primes  $p \notin \mathcal{D}$ .

The multiset  $\mathcal{B}$  will contain elements that are smaller than  $\alpha N$ ; but we can compute the number of such elements fairly precisely. Indeed, for any natural number  $\ell$ , the number of elements of  $\mathcal{B}$  that are less than  $\alpha N/\ell$  is

$$\sum_{d \in \mathcal{D}} \left| A_{d,\mathcal{D}} \cap \left[ 1, \frac{N}{d} \right] \cap \left[ 1, \frac{\alpha N}{\ell} \right] \right| = N \sum_{d \in \mathcal{D}} \sigma_{d,\mathcal{D}} \min \left( \frac{1}{d}, \frac{\alpha}{\ell} \right) + O(1). \tag{6.6}$$

We now form a modification  $\mathcal{B}'$  of the multiset  $\mathcal{B}$  by multiplying each element of  $\mathcal{B}$  by an appropriate natural number to make it at least  $\alpha N$ . More precisely, we perform the following algorithm.

- Initialize  $\mathcal{B}'$  to be empty.

- Choose a large natural number  $r$ .
- For each  $k \geq 1$ , and each element  $m$  of  $\mathcal{B}$  that lies in the interval  $[\alpha N/p_0^{r+k}, \alpha N/p_0^{r+k-1})$ , add  $p_0^{r+k}m$  to  $\mathcal{B}'$ .
- For the  $\lfloor (\sum_{\ell \geq p_0^r} a_\ell)N \rfloor$  smallest elements  $m$  of  $\mathcal{B}$  that are at least  $\alpha N/p_0^r$ , add  $p_0^r m$  to  $\mathcal{B}'$ .
- For the  $\lfloor a_{p_0^{r-1}} N \rfloor$  next smallest elements  $m$  of  $\mathcal{B}$  that are at least  $\alpha N/p_0^r$ , add  $(p_0^r - 1)m$  to  $\mathcal{B}'$ .
- Repeat the previous step with  $p_0^r - 1$  replaced in turn by  $a_{p_0^{r-2}}, \dots, a_2$ , halting if no further elements of  $\mathcal{B}$  remain.
- For any remaining elements  $m$  of  $\mathcal{B}$  that have not been involved in any previous step, add  $m$  to  $\mathcal{B}'$ .

It is clear that  $\mathcal{B}'$  has the same cardinality  $N$  as  $\mathcal{B}$ . If  $p \notin \mathcal{D}$ , the hypothesis (ii) forces  $a_\ell = 0$  for all  $\ell$  that are divisible by  $p$ ; because of this,  $\mathcal{B}'$  remains in balance at those primes. For the primes  $p$  in  $\mathcal{D}$ , we see from construction that the  $p$ -surplus of  $\mathcal{B}'$  has decreased from  $\mathcal{B}$  by at most

$$\begin{aligned} v_p \left( \frac{N!}{\prod \mathcal{B}} \right) - v_p \left( \frac{N!}{\prod \mathcal{B}'} \right) &\leq \sum_{\ell < p_0^r} v_p(\ell) a_\ell N \\ &\quad + \sum_{\ell \geq p_0^r} r a_\ell N \\ &\quad + \sum_{k=1}^{\infty} (r+k) |\mathcal{B} \cap [\alpha N/p_0^{r+k}, \alpha N/p_0^{r+k-1})|. \end{aligned}$$

The final summand is empty unless  $k = O(\log N)$ , and then by (6.6) we can crudely bound

$$\begin{aligned} \sum_{k=1}^{\infty} (r+k) |\mathcal{B} \cap [\alpha N/p_0^{r+k}, \alpha N/p_0^{r+k-1})| &\ll \sum_{k=1}^{\infty} (r+k) \frac{N}{p_0^{r+k}} + \log N \\ &\ll \frac{r}{p_0^r} N + \log N. \end{aligned}$$

If  $r$  is sufficiently large, and then  $N$  sufficiently large depending on  $r$ , this quantity is smaller than any given multiple of  $N$ ; the same is true for  $\sum_{\ell \geq p_0^r} r a_\ell N$  thanks to (6.4). If we now compare with (6.5) and the strictness of (6.2), we conclude for such choices of  $r, N$  that

$$v_p \left( \frac{N!}{\prod \mathcal{B}'} \right) \geq 0$$

for all  $p \in \mathcal{D}$ , thus  $\mathcal{B}'$  remains a subfactorization of  $N!$ .

It remains to verify that  $\mathcal{B}'$  is  $\alpha N$ -admissible. From an inspection of the algorithm, we see that the only way this can fail to be the case is if, for some  $1 \leq \ell < p_0^r$ , the number of elements

of  $\mathcal{B}$  in  $[\alpha N/p'_0, \alpha N/\ell)$  exceeds the quantity

$$\begin{aligned} & \lfloor a_{\ell+1}N \rfloor + \cdots + \lfloor a_{p'_0-1}N \rfloor + \left\lfloor \left( \sum_{\ell \geq p'_0} a_\ell \right) N \right\rfloor, \\ &= N \sum_{\ell' > \ell} a_{\ell'} + O(p'_0) \end{aligned} \tag{6.7}$$

since this quantity is the number of smallest elements of  $\mathcal{B}$  that are at least  $\alpha N/p'_0$  by factors greater than  $\ell$ . By (6.6), the number of elements of  $\mathcal{B}$  in  $[\alpha N/p'_0, \alpha N/\ell)$  is at most

$$N \sum_{d \in \mathcal{D}} \sigma_{d,D} \max \left( \frac{1}{d}, \frac{\alpha}{\ell} \right) + O(1),$$

so from (6.3) we see that this quantity cannot exceed (6.7) if  $N$  is sufficiently large depending on  $r$ . Thus  $\mathcal{B}'$  is  $\alpha N$ -admissible, and the claim follows.  $\square$

**Example 6.3.** Let  $0 < \alpha < 3/16$  and  $\mathcal{D} = \{1, 2, 4\}$ . If we set  $a_{2^r} = \frac{3}{2^{r+3}}$  for  $r \geq 1$ , and  $a_\ell = 0$  for all other  $\ell$ , then one can calculate that

$$\sum_{\ell=1}^{\infty} a_\ell v_2(\ell) = \frac{3}{4} = \sum_{d \in \mathcal{D}} v_2(d) \frac{\sigma(d)}{d}$$

and

$$\sum_{\ell > 2^r} a_{2^r} = \frac{3}{2^{r+3}} > \frac{2\alpha}{2^r} = \sum_{d \in \mathcal{D}} \sigma_{d,D} \min \left( \frac{1}{d}, \frac{\alpha}{2^r} \right)$$

for any  $r \geq 1$ . From this one can readily check that the hypotheses of Proposition 6.2 are satisfied, and so we recover the bound  $\frac{t(N)}{N} \geq \frac{3}{16} - o(1)$  from [10].

If one makes the ansatz  $a_{2^r} = c/2^r$  for some  $c > 0$  and all  $r \geq r_0$ , with  $a_\ell = 0$  for all other  $\ell > 2^{r_0}$ , then the task of locating weights  $a_r$  obeying the hypotheses of Proposition 6.2 becomes a linear programming problem. Numerically<sup>7</sup>, we were able to locate such weights for  $\alpha = 1/3$ ,  $\mathcal{D} = \{d : 1 \leq d \leq 2^{11}\}$ , and  $r_0 = 11$ , thus establishing that  $t(N) \geq N/3$  for sufficiently large  $N$ .

We can also recover a weak version of Theorem 1.3(iv) with this method:

**Proposition 6.4** (Asymptotic lower bound). *If  $0 < \alpha < 1/e$ , then one has  $t(N) \geq \alpha N$  for all sufficiently large  $N$ .*

*Proof.* We select the following parameters:

- A sufficiently small real  $c_0 > 0$  (which can depend on  $\alpha$ );
- A sufficiently large natural number  $M$  (which can depend  $c_0, \alpha$ );
- A sufficiently large natural number  $C_0$  (which can depend on  $M, c_0, \alpha$ );
- A sufficiently large prime  $p_-$  (which can depend on  $C_0, M, c_0, \alpha$ ); and
- A prime  $p_+$  with  $\log p_+ \asymp \log^2 \log p_-$ ; and

<sup>7</sup><https://github.com/teorth/erdos-guy-selfridge/tree/main/src/dnup>

Let  $\mathcal{D}$  be the set of all numbers  $d$  which are either of the form  $2^m$  for  $0 \leq m \leq M$ , or  $2^m p$  for  $0 \leq m \leq M$  and  $p_- \leq p \leq p_+$ . This is a downset, and the densities  $\sigma_{d,D}$  can be computed explicitly as

$$\sigma_{2^m,D} = \frac{1}{2} \mu_+; \sigma_{2^m p,D} = \frac{1}{2} \mu_p$$

for  $0 \leq m < M$  and

$$\sigma_{2^M,D} = \mu_+; \sigma_{2^M p,D} = \mu_p,$$

where

$$\mu_p := \prod_{p_- \leq p' < p} \left(1 - \frac{1}{p'}\right)$$

and

$$\mu_+ := \prod_{p_- \leq p' \leq p_+} \left(1 - \frac{1}{p'}\right).$$

The identity (6.1) is then equivalent to the telescoping identity

$$\sum_{p_- \leq p \leq p_+} \frac{\mu_p}{p} + \mu_+ = 1. \quad (6.8)$$

We can now define the weights  $a_\ell$  as follows:

- (i) If  $\ell = 2^m$  for some  $m \geq 1$ , we set  $\alpha_\ell := C_0 \frac{m}{2^m} \mu_+$ .
- (ii) If  $\ell = 2^{C_0} p$  for some  $p_- \leq p \leq 2^{C_0} p_-$ , we set  $\alpha_\ell := \mu_p/p$ .
- (iii) If  $\ell = p$  for some  $2^{C_0} p_- < p < p_+/2$ , we set  $\alpha_\ell := c_0 \mu_p/p$ .
- (iv) If  $\ell = 2p$  for some  $2^{C_0} p_- < p < p_+/2$ , we set  $\alpha_\ell := (1 - c_0) \mu_p/p$ .
- (v) If  $\ell = 2^{C_0+m} p$  for some  $p_+/2 \leq p \leq p_+$  and  $m \geq 1$ , we set  $\alpha_\ell = \frac{m}{2^m} \mu_p/p$ .
- (vi) In all other cases, we set  $\alpha_\ell = 0$ .

By Proposition 6.2, it suffices to verify the conditions (6.2), (6.3). We begin with (6.2). If  $p$  is not equal to 2 or in the range  $[p_-, p_+]$ , then both sides of (6.2) vanish. If  $p$  is in  $[p_-, p_+]$ , then a routine computation shows that both sides are equal to  $\mu_p$ . For  $p = 2$ , the right-hand side can be simplified using (6.8) to

$$\sum_{0 \leq m \leq M}^* \frac{1}{2} \frac{m}{2^m} = 1 - O\left(\frac{M}{2^M}\right)$$



where the  $*$  indicates that the final term  $m = M$  is doubled (thus  $\sum_{0 \leq m \leq M}^* f(m) = \sum_{0 \leq m < M} f(m) + 2f(M)$ ). Meanwhile, the left-hand side can be computed to equal

$$\begin{aligned} & \sum_{m=1}^{\infty} \frac{C_0 m^2}{2^m} \mu_+ \\ & + \sum_{p_- \leq p \leq 2^{C_0} p_-} C_0 \frac{\mu_p}{p} \\ & + \sum_{2^{C_0} p_- < p < p_+/2} (1 - c_0) \frac{\mu_p}{p} \\ & + \sum_{p_+/2 \leq p \leq p_+} \sum_{m=1}^{\infty} \frac{(C_0 + m)m}{2^m} \frac{\mu_p}{p} \end{aligned}$$

which simplifies using (6.8) and summation in  $m$  to

$$1 - c_0 + O_{C_0} \left( \mu_+ + \sum_{p_- \leq p \leq 2^{C_0} p_-} \frac{\mu_p}{p} + \sum_{p_+/2 \leq p \leq p_+} \frac{\mu_p}{p} \right).$$

From Mertens' theorem (or Lemma 2.2) one has

$$\mu_p = \frac{\log p_-}{\log p} \left( 1 + O \left( \frac{1}{\log^{10} p} \right) \right) \quad (6.9)$$

and similarly

$$\mu_+ = \frac{\log p_-}{\log p_+} \left( 1 + O \left( \frac{1}{\log^{10} p_+} \right) \right) \quad (6.10)$$

and (6.2) for  $p = 2$  then follows from the choice of parameters after a brief calculation.

It remains to verify (6.3). We first consider the case  $\ell < 2^{C_0} p_-$ . In this case, the left-hand side simplifies using (6.8) to

$$\sum_{m: 2^m > \ell} \frac{C_0 m}{2^m} \mu_+ + (1 - \mu_+)$$

and the right-hand side similarly simplifies to

$$\sum_{0 \leq m \leq M}^* \frac{1}{2} \min \left( \frac{1}{2^m}, \frac{\alpha}{\ell} \right) \mu_+ + (1 - \mu_+).$$

Thus it suffices to show that

$$\sum_{0 \leq m \leq M}^* \frac{1}{2} \min \left( \frac{1}{2^m}, \frac{\alpha}{\ell} \right) < \sum_{m: 2^m > \ell} \frac{C_0 m}{2^m}. \quad (6.11)$$

But the left-hand side is  $\ll \frac{\log(2+\ell)}{\ell}$  and the right-hand side is  $\gg C_0 \frac{\log(2+\ell)}{\ell}$ , giving the claim.

Next, we consider the case  $\ell \geq p_+$ . The left-hand side of (6.3) is at least

$$\sum_{m: 2^m > \ell} \frac{C_0 m}{2^m} \mu_+ + \sum_{p_+/2 \leq p \leq p_+} \sum_{m: 2^{m+C_0} p > \ell} \frac{m}{2^m} \frac{\mu_p}{p}$$

and the right-hand side is at most

$$\sum_{0 \leq m \leq M}^* \frac{1}{2} \min\left(\frac{1}{2^m}, \frac{\alpha}{\ell}\right) \mu_+ + \sum_{0 \leq m \leq M}^* \frac{1}{2} \sum_{p_- \leq p \leq p_+} \min\left(\frac{1}{2^m p}, \frac{\alpha}{\ell}\right) \mu_p. \quad (6.12)$$

By (6.11) it suffices to show that

$$\sum_{p_+/2 \leq p \leq p_+} \sum_{m: 2^{m+C_0} p > \ell} \frac{m}{2^m} \frac{\mu_p}{p} > \sum_{0 \leq m \leq M}^* \frac{1}{2} \sum_{p_- \leq p \leq p_+} \min\left(\frac{1}{2^m p}, \frac{\alpha}{\ell}\right) \mu_p.$$

The left-hand side can be computed using (6.9) and the prime number theorem to be

$$\gg 2^{C_0} \frac{\log p_-}{\log^2 p_+} \frac{\log(2 + \ell/p_+)}{\ell}.$$

By (6.9) and Lemma 2.2, the right-hand side may be bounded by

$$\ll (\log p_-) \sum_{m=0}^{\infty} \int_{p_-}^{p_+} \min\left(\frac{1}{2^m t}, \frac{1}{\ell}\right) \frac{dt}{\log^2 t}.$$

We can perform the  $m$  summation and bound this by

$$\ll \frac{\log p_-}{\ell} \int_{p_-}^{p_+} \log\left(2 + \frac{\ell}{t}\right) \frac{dt}{\log^2 t},$$

and the claim then follows from a routine computation.

Finally, we consider the case  $2^{C_0} p_- < \ell < p_+$ . Note that if we redefined  $a_{\ell'}$  to make rules (iii), (iv) apply for all  $p_- \leq p \leq p_+$ , and delete rules (ii) and (v), then this amounts to transferring the mass of  $a_{\ell'}$  from larger  $\ell'$  to smaller  $\ell'$ , so that the sum  $\sum_{\ell' > \ell} a_{\ell'}$  does not increase. From this observation, we see that we can lower bound the left-hand side of (6.3) by

$$\sum_{m: 2^m > \ell} \frac{C_0 m}{2^m} \mu_+ + \sum_{p_- \leq p \leq p_+} \frac{\mu_p}{p} (c_0 1_{p > \ell} + (1 - c_0) 1_{2p > \ell}),$$

while the right-hand side is at most (6.12). By (6.11), it suffices to show that

$$\sum_{p_- \leq p \leq p_+} \frac{\mu_p}{p} (c_0 1_{p > \ell} + (1 - c_0) 1_{2p > \ell}) \geq \sum_{0 \leq m \leq M}^* \frac{1}{2} \sum_{p_- \leq p \leq p_+} \min\left(\frac{1}{2^m p}, \frac{\alpha}{\ell}\right) \mu_p. \quad (6.13)$$

Applying (6.9), Lemma 2.2, we can bound the right-hand side of (6.13) by

$$\left(1 + O\left(\frac{1}{\log^{10} p_-}\right)\right) \frac{\log p_-}{2} \sum_{0 \leq m \leq M}^* \int_{p_-}^{p_+} \min\left(\frac{1}{2^m t}, \frac{\alpha}{\ell}\right) \frac{dt}{\log^2 t}.$$

The minimum  $\min(\frac{1}{2^m t}, \frac{\alpha}{\ell})$  is equal to  $\frac{1}{2^m t}$  for  $t \geq \ell/(2^m \alpha)$  and  $\frac{\alpha}{\ell}$  for  $t < \ell/(2^m \alpha)$ . Routine estimation then gives

$$\begin{aligned} & \int_{p_-}^{p_+} \min\left(\frac{1}{2^m t}, \frac{\alpha}{\ell}\right) \frac{dt}{\log^2 t} \\ & \leq \frac{1}{2^m \log \frac{\ell}{\alpha 2^m}} - \frac{1}{2^m \log p_+} + \frac{1}{2^m \log^2 \frac{\ell}{\alpha 2^m}} + O_M\left(\frac{1}{\log^3 \ell}\right) \\ & = \frac{1}{2^m \log \ell} - \frac{1}{2^m \log p_+} + \frac{m \log 2 - \log \frac{1}{e\alpha}}{2^m \log^2 \ell} + O_M\left(\frac{1}{\log^3 \ell}\right). \end{aligned}$$

Performing the  $m$  summation, we conclude that the right-hand side of (6.13) is at most

$$(\log p_-) \left( \frac{1}{\log \ell} - \frac{1}{\log p_+} + \frac{\log 2 - \log \frac{1}{e\alpha}}{\log^2 \ell} + O_M\left(\frac{1}{\log^3 \ell}\right) \right).$$

Meanwhile, by (6.9), the left-hand side of (6.13) can be computed to be

$$\left(1 + O\left(\frac{1}{\log^{10} p_-}\right)\right) (\log p_-) \left( c_0 \sum_{\ell < p \leq p_+} \frac{1}{p \log p} + (1 - c_0) \sum_{\ell/2 < p \leq p_+} \frac{1}{p \log p} \right).$$

Applying Lemma 2.2 and evaluating the integrals, we can bound this by

$$\left(1 + O\left(\frac{1}{\log^{10} p_-}\right)\right) (\log p_-) \left( c_0 \left( \frac{1}{\log \ell} - \frac{1}{\log p_+} \right) + (1 - c_0) \left( \frac{1}{\log(\ell/2)} - \frac{1}{\log p_+} \right) \right)$$

which after routine Taylor expansion simplifies to

$$(\log p_-) \left( \frac{1}{\log \ell} - \frac{1}{\log p_+} + \frac{\log 2 - c_0 \log 2}{\log^2 \ell} + O\left(\frac{1}{\log^3 \ell}\right) \right).$$

Since  $\alpha < 1/e$ , we can ensure that  $c_0 \log 2 < \log \frac{1}{e\alpha}$  by taking  $c_0$  small enough. The claim (6.13) then follows.  $\square$

## 7. THE ACCOUNTING EQUATION

Given a  $t$ -admissible multiset  $\mathcal{B}$  (which we view as an approximate factorization of  $N!$ ), we can apply the fundamental theorem of arithmetic (2.1) to the rational number  $N! / \prod \mathcal{B}$  and rearrange to obtain the *accounting equation*

$$\mathcal{E}_t(\mathcal{B}) + \sum_p v_p \left( \frac{N!}{\prod \mathcal{B}} \right) \log p = \log N! - |\mathcal{B}| \log t \quad (7.1)$$

where we define the  $t$ -excess  $\mathcal{E}_t(\mathcal{B})$  of the multiset  $\mathcal{B}$  by the formula

$$\mathcal{E}_t(\mathcal{B}) := \sum_{a \in \mathcal{B}} \log \frac{a}{t}. \quad (7.2)$$

**Example 7.1.** Suppose one wishes to factorize  $5! = 2^3 \times 3 \times 5$ . The attempted 3-admissible factorization  $\mathcal{B} := \{3, 4, 5, 5\}$  has a 2-surplus of  $v_2(5! / \prod \mathcal{B}) = 1$ , is in balance at 3, and has a 5-deficit of  $v_5(\prod \mathcal{B} / 5!) = 1$ , so it is not a factorization or subfactorization of  $5!$ . The 3-excess of this multiset is

$$\mathcal{E}_3(\mathcal{B}) = \log \frac{3}{3} + \log \frac{4}{3} + \log \frac{5}{3} + \log \frac{5}{3} = 1.3093 \dots$$

and the accounting equation (7.1) becomes

$$1.3093 \dots + \log 2 - \log 5 = 0.3930 \dots = \log 5! - 4 \log 3.$$

If one replaces one of the copies of 5 in  $\mathcal{B}$  with a 2, this erases both the 2-surplus and the 5-deficit, and creates a factorization  $\mathcal{B}' = \{2, 3, 4, 5\}$  of  $5!$ ; the 3-excess now drops to

$$\mathcal{E}_3(\mathcal{B}') = \log \frac{2}{3} + \log \frac{3}{3} + \log \frac{4}{3} + \log \frac{5}{3} = 0.3930 \dots,$$

bringing the accounting equation back into balance.

In view of Remark 1.2, one can now equivalently describe  $t(N)$  as follows:

**Lemma 7.2** (Equivalent description of  $t(N)$ ).  *$t(N)$  is the largest quantity  $t$  for which there exists a  $t$ -admissible subfactorization of  $N!$  with*

$$\mathcal{E}_t(\mathcal{B}) + \sum_p v_p \left( \frac{N!}{\prod \mathcal{B}} \right) \log p \leq \log N! - N \log t.$$

One can view  $\log N! - N \log t$  as an available “budget” that one can “spend” on some combination of  $t$ -excess and  $p$ -surpluses. For  $t$  of the form  $t = N/e^{1+\delta}$  for some  $\delta > 0$ , the budget can be computed using the Stirling approximation (2.4) to be  $\delta N + O(\log N)$ . The non-negativity of the  $t$ -excess and  $p$ -surpluses recovers the trivial upper bound (1.2); but note that any prime  $p > \frac{t}{\lfloor \sqrt{t} \rfloor}$  must inevitably contribute at least  $\log \frac{\lfloor t/p \rfloor}{t/p}$  to the  $t$ -excess if it is to appear in the multiset  $\mathcal{B}$ . By pursuing this line of reasoning, one can obtain an alternate proof of Lemma 5.1; see [16, Lemma 2.1].

## 8. MODIFIED APPROXIMATE FACTORIZATIONS

In this section we present and then analyze an algorithm that starts with an *approximate* factorization  $\mathcal{B}^{(0)}$  of  $N!$ , which is  $t$ -admissible but omits all tiny primes, and is approximately in balance in small and medium primes, and attempts to “repair” this factorization to establish a lower bound of the form  $t(N) \geq t$ .

To describe the criterion for the algorithm to succeed, it will be convenient to introduce the following notation. For  $a_+, a_- \in [0, +\infty]$ , we define the asymmetric norm  $|x|_{a_+, a_-}$  of a real number  $x$  by the formula

$$|x|_{a_+, a_-} := \begin{cases} a_+ |x| & x \geq 0 \\ a_- |x| & x \leq 0, \end{cases}$$

with the usual convention  $+\infty \times 0 = 0$ . If  $a_+, a_-$  are finite, this function is Lipschitz with constant  $\max(a_+, a_-)$ . One can think of  $a_+$  as the “cost” of making  $x$  positive, and  $a_-$  as the “cost” of making  $x$  negative.

The analysis of the algorithm is now captured by the following proposition.

**Proposition 8.1** (Repairing an approximate factorization). *Let  $N, K$  be natural numbers, and let  $1 \leq t \leq N$  be an additional parameter obeying the conditions*

$$\frac{t}{K} \geq \sqrt{N}; \quad \frac{t}{K^2} \geq K \geq 5. \quad (8.1)$$

*We also assume that there are additional parameters  $\kappa_* > 0$  and  $0 \leq \gamma_2, \gamma_3 < 1$ , such that there exist 3-smooth numbers*

$$t \leq 2^{n_2} 3^{m_2}, 2^{n_3} 3^{m_3} \leq e^{\kappa_*} t \quad (8.2)$$

*such that*

$$2m_2 \leq \gamma_2 n_2; \quad n_3 \leq 2\gamma_3 m_3. \quad (8.3)$$

*We define the “norm” of a pair  $n, m$  of real numbers by the formula*

$$\|(n, m)\|_\gamma := \max \left( \frac{n - 2\gamma_2 m}{1 - \gamma_2}, \frac{2m - \gamma_3 n}{1 - \gamma_3} \right).$$

*Let  $\mathcal{B}^{(0)}$  be a  $t$ -admissible multiset of natural numbers, with all elements of  $\mathcal{B}^{(0)}$  at most  $(t/K)^2$ , and suppose that one has the inequalities*

$$\sum_{i=1}^8 \delta_i \leq \delta \quad (8.4)$$

*and*

$$\sum_{i=1}^7 \alpha_i \leq 1 \quad (8.5)$$

where

$$\delta_1 := \frac{1}{N} \mathcal{E}_t(\mathcal{B}^{(1)}) \quad (8.6)$$

$$\delta_2 := \frac{1}{N} \sum_{t/K < p \leq N} f_{N/t}(p/N) \quad (8.7)$$

$$\delta_3 := \frac{\kappa_{4.5}}{N} \sum_{3 < p_1 \leq t/K} \left| v_{p_1} \left( \frac{N!}{\prod \mathcal{B}^{(0)}} \right) \right| \quad (8.8)$$

$$\delta_4 := \kappa_{4.5} \sum_{K < p_1 \leq t/K} A_{p_1} \quad (8.9)$$

$$\delta_5 := \kappa_{4.5} \sum_{3 < p_1 \leq K} |A_{p_1} - B_{p_1}|_{\frac{\log p_1}{\log(t/K^2)}, 1} \quad (8.10)$$

$$\delta_6 := \frac{\kappa_{4.5}}{N} \quad (8.11)$$

$$\delta_7 := \frac{\kappa_*}{\log t} \left( \log \sqrt{12} - B_2 \log 2 - B_3 \log 3 \right) \quad (8.12)$$

$$\delta_8 := \frac{2(\log t + \kappa_*)}{N} \quad (8.13)$$

$$\delta := \frac{1}{N} \log N! - \log t \quad (8.14)$$

$$\alpha_1 := \frac{1}{N} \left\| \left( v_2 \left( \prod \mathcal{B}^{(0)} \right), v_3 \left( \prod \mathcal{B}^{(0)} \right) \right) \right\|_\gamma \quad (8.15)$$

$$\alpha_2 := \|(B_2, B_3)\|_\gamma \quad (8.16)$$

$$\alpha_3 := \frac{\log \frac{t}{K} + \kappa_{**}}{N \log \sqrt{12}} \sum_{3 < p_1 \leq t/K} \left| v_{p_1} \left( \frac{N!}{\prod \mathcal{B}^{(0)}} \right) \right| \quad (8.17)$$

$$\alpha_4 := \frac{1}{\log \sqrt{12}} \sum_{K < p_1 \leq t/K} \left( \log \frac{t}{p_1} + \kappa_{**} \right) A_{p_1} \quad (8.18)$$

$$\alpha_5 := \frac{1}{\log \sqrt{12}} \sum_{3 < p_1 \leq K} |A_{p_1} - B_{p_1}|_{\frac{\log p_1}{\log(t/K^2)}, (\log K^2 + \kappa_{**}), \log p_1 + \kappa_{**}} \quad (8.19)$$

$$\alpha_6 := \frac{\log t + \kappa_{**}}{N \log \sqrt{12}} \quad (8.20)$$

$$\alpha_7 := \max \left( \frac{\log(2N)}{(1 - \gamma_2)N \log 2}, \frac{\log(3N)}{(1 - \gamma_3)N \log \sqrt{3}} \right) \quad (8.21)$$

$$\kappa_{**} := \max(\kappa_{4.5, \gamma_2}^{(2)}, \kappa_{4.5, \gamma_3}^{(3)}) \quad (8.22)$$

$$A_{p_1} := \frac{1}{N} \sum_m v_{p_1}(m) |\{a \in \mathcal{B}^{(0)} : a = mp \text{ for a prime } p > t/K\}| \quad (8.23)$$

$$B_{p_1} := \frac{1}{N} \sum_{m \leq K} v_{p_1}(m) \sum_{\frac{t}{m} \leq p < \frac{t}{m-1}} \left\lfloor \frac{N}{p} \right\rfloor, \quad (8.24)$$

with the convention that the upper bound  $p < \frac{t}{m-1}$  in (8.24) is vacuous when  $m = 1$ . Then  $t(N) \geq t$ .

In practice, the parameter  $K$  will be quite small compared to  $N$ , and the quantities  $\gamma_2, \gamma_3, \kappa_*$  will also be somewhat smaller than 1.

**Remark 8.2.** In the notation of this proposition, Lemma 5.1 can essentially be interpreted as a necessary condition  $\delta_2 \leq \delta$  for  $t(N) \leq t$  to be provable; to use the above proposition effectively, it is thus desirable to have all the other  $\delta_i, i \neq 2$  terms be as small as possible. The criterion in Lemma 7.2 can similarly be rewritten as  $\delta_1 + \delta_9 \leq \delta$ , where

$$\delta_9 := \frac{1}{N} \sum_p \left| v_{p_1} \left( \frac{N!}{\prod \mathcal{B}^{(0)}} \right) \right|_{\log p, \infty}.$$

In practice,  $\delta_9$  is too large (or infinite) for this criterion to be directly useful; the algorithm below is intended to replace this large quantity with something much smaller, and in particular to utilize tiny primes to gain factors such as  $\kappa_L$  for various  $L$  in the bounds of the main  $\delta_i$  terms besides the “non-negotiable”  $\delta_2$ . The secondary condition (8.5) can be interpreted as a requirement that “enough” tiny primes are available in the factorization of  $N!$  to perform such adjustments.

The rest of this section will be devoted to the proof of this proposition. It will be convenient to divide the primes into four classes:

- *Tiny primes*  $p = 2, 3$ .
- *Small primes*  $3 < p \leq K$ .
- *Medium primes*  $K < p \leq t/K$ .
- *Large primes*  $p > t/K$ .

Initially, the multiset  $\mathcal{B}^{(0)}$  may have the “wrong” number of factors at large primes. We fix this by applying the following modifications to  $\mathcal{B}^{(0)}$ :

- (a) Remove all elements of  $\mathcal{B}^{(0)}$  that are divisible by a large prime  $p > t/K$  from the multiset.
- (b) For each large prime  $p > t/K$ , add  $v_p(N!)$  copies of  $p \lceil t/p \rceil$  to the multiset.

We let  $\mathcal{B}^{(1)}$  be the multiset formed after completing both Step (a) and Step (b). We make two simple observations:

- (A) Since the elements of  $\mathcal{B}^{(0)}$  are at most  $(t/K)^2$ , all the elements removed in Step (a) are of the form  $mp$  where  $m \leq t/K$ .
- (B) For each large prime  $p$  considered in Step (b), one has  $v_p(N!) = \lfloor N/p \rfloor$  by (2.3) and (8.1), while  $\lceil t/p \rceil \leq K \leq t/K$  (again by (8.1)).

From this, we see that  $\mathcal{B}^{(1)}$  is automatically  $t$ -admissible, and in balance at any large prime  $p > t/K$ :

$$v_p \left( \frac{N!}{\prod \mathcal{B}^{(1)}} \right) = 0.$$

For medium primes  $K < p_1 \leq t/K$ , one can have some increase in the  $p_1$ -surplus coming from Step (a), which is described by (8.23):

$$v_{p_1} \left( \frac{N!}{\prod \mathcal{B}^{(1)}} \right) = v_{p_1} \left( \frac{N!}{\prod \mathcal{B}^{(0)}} \right) + N A_{p_1}.$$

For small or tiny primes  $p \leq K$ , one also has some possible decrease in the  $p_1$ -surplus coming from Step (b), which is described by (8.24):

$$v_{p_1} \left( \frac{N!}{\prod \mathcal{B}^{(1)}} \right) = v_{p_1} \left( \frac{N!}{\prod \mathcal{B}^{(0)}} \right) + N(A_{p_1} - B_{p_1}).$$

In particular, we have from (8.15), (8.16) and the triangle inequality that

$$\frac{1}{N} \left\| \left( v_2 \left( \prod \mathcal{B}^{(1)} \right), v_3 \left( \prod \mathcal{B}^{(1)} \right) \right) \right\|_\gamma \leq \alpha_1 + \alpha_2. \quad (8.25)$$

Each element removed in Step (a) reduces the  $t$ -excess, while each element  $p \lceil t/p \rceil$  added in Step (b) increases the  $t$ -excess by  $\log \frac{\lceil t/p \rceil}{t/p}$ , so each large prime  $t/K < p \leq N$  contributes a net of  $\lfloor \frac{N}{p} \rfloor \log \frac{\lceil t/p \rceil}{t/p} = f_{N/t}(p/N)$  to the  $t$ -excess. Thus by (8.6), (8.7) we have

$$\frac{1}{N} \mathcal{E}_t(\mathcal{B}^{(1)}) \leq \delta_1 + \delta_2. \quad (8.26)$$

Now we bring the multiset  $\mathcal{B}^{(1)}$  into balance at small and medium primes  $3 < p \leq t/K$ . We make the following observations:

- (C) If an element in  $\mathcal{B}^{(1)}$  is divisible by some small or medium prime  $3 < p \leq t/K$ , and one replaces  $p$  by  $\lceil p \rceil_{4.5}^{(2,3)}$  in the factorization of that element, then the  $p$ -deficit decreases by one, while (by Lemma 2.1) the  $t$ -excess increases by at most  $\kappa_{4.5}$ , and the quantity  $\|(v_2(\prod \mathcal{B}^{(1)}), v_3(\prod \mathcal{B}^{(1)}))\|_\gamma$  increases by at most  $\frac{\log p + \kappa_{**}}{\log \sqrt{12}}$ . All other  $p_1$ -surpluses or  $p_1$ -deficits for  $p_1 \neq 2, 3, p$  remain unaffected.
- (D) If one adds an element of the form  $m \lceil t/m \rceil_{4.5}^{(2,3)}$  to  $\mathcal{B}^{(1)}$  for some  $m \leq t/K$  that is the product of small or medium primes  $3 < p \leq t/K$ , then the  $p$ -surpluses at small or medium primes  $p$  decrease by  $v_p(m)$ , while (by Lemma 2.1) the  $t$ -excess increases by at most  $\kappa_{4.5}$ , and the quantity  $\|(v_2(\frac{N!}{\prod \mathcal{B}^{(1)}}), v_3(\frac{N!}{\prod \mathcal{B}^{(1)}}))\|_\gamma$  increases by at most  $\frac{\log(t/m) + \kappa_{**}}{\log \sqrt{12}}$ . The  $p$ -surpluses or  $p$ -deficits at medium or large primes remain unaffected.

With these observations in mind, we perform the following modifications to the multiset  $\mathcal{B}^{(1)}$ .

- (c) If there is a  $p_1$ -deficit  $v_{p_1}(\prod \mathcal{B}^{(1)}/N!) > 0$  at some small or medium prime  $3 < p_1 \leq t/K$ , then we perform the replacement of  $p_1$  in one of the elements of  $\mathcal{B}^{(1)}$  with  $\lceil p_1 \rceil_{4.5}^{(2,3)}$



as per observation (C), repeated  $v_{p_1}(\prod B^{(1)}/N!)$  times, in order to eliminate all such deficits.

- (d) If there is a  $p$ -surplus  $v_p(\prod N!/B^{(1)}) > 0$  at some medium prime  $K < p \leq t/K$ , we add the element  $p \lceil t/p \rceil_{4.5}^{(2,3)}$  to  $B^{(1)}$  as per observation (D),  $v_p(\prod N!/B^{(1)})$  times, in order to eliminate all such surpluses at medium primes.
- (d') If there are  $p$ -surpluses  $v_p(\prod N!/B^{(1)}) > 0$  at some small primes  $3 < p \leq K$ , we multiply all these primes together, then apply the greedy algorithm to factor them into products  $m$  in the range  $t/K^2 < m \leq t/K$ , plus at most one exceptional product in the range  $1 < m \leq t/K$ . For each of these  $m$ , add  $m \lceil t/m \rceil_{4.5}^{(2,3)}$  to  $B^{(1)}$  as per observation (D), to eliminate all such surpluses at small primes.

Call the multiset formed from  $B^{(1)}$  formed as the outcome of applying Steps (c), (d), (d') as  $B^{(2)}$ . The product of all the primes arising in Step (d') has logarithm equal to

$$\sum_{3 < p_1 \leq K} \left| v_{p_1} \left( \frac{N!}{\prod B^{(1)}} \right) \right|_{\log p_1, 0} = \sum_{3 < p_1 \leq K} \left| v_p \left( \frac{N!}{\prod B^{(0)}} \right) \right|_{\log p_1, 0}$$

and hence the number of non-exceptional  $m$  arising in (d') is at most

$$\sum_{3 < p_1 \leq K} \left| v_p \left( \frac{N!}{\prod B^{(0)}} \right) \right|_{\frac{\log p_1}{\log(t/K^2)}, 0}.$$

The total excess of  $B^{(2)}$  is increased in Step (c) by at most

$$\kappa_{4.5} \sum_{3 < p_1 \leq t/K} \left| v_{p_1} \left( \frac{N!}{\prod B^{(1)}} \right) \right|_{0,1} = \kappa_{4.5} \sum_{3 < p_1 \leq t/K} \left| v_{p_1} \left( \frac{N!}{\prod B^{(0)}} \right) + N(A_{p_1} - B_{p_1}) \right|_{0,1},$$

in Step (d) by at most

$$\kappa_{4.5} \sum_{K < p_1 \leq t/K} \left| v_{p_1} \left( \frac{N!}{\prod B^{(1)}} \right) \right|_{1,0} = \kappa_{4.5} \sum_{K < p_1 \leq t/K} \left| v_{p_1} \left( \frac{N!}{\prod B^{(0)}} \right) + N A_{p_1} \right|_{1,0},$$

and in Step (c) by at most

$$\kappa_{4.5} \left( 1 + \sum_{3 < p_1 \leq K} \left| v_{p_1} \left( \frac{N!}{\prod B^{(0)}} \right) \right|_{\frac{\log p_1}{\log(t/K^2)}, 0} \right).$$

From the triangle inequality and (8.26), (8.8), (8.9), (8.10), (8.11), we then have

$$\frac{1}{N} \mathcal{E}_t(B^{(2)}) \leq \sum_{i=1}^6 \delta_i. \quad (8.27)$$

Similarly, the quantity  $\frac{1}{N} \|(v_2(\prod B^{(1)}), v_3(\prod B^{(1)}))\|_\gamma$  is increased in Step (c) by at most

$$\frac{1}{N \log \sqrt{12}} \sum_{3 < p_1 \leq t/K} \left| v_{p_1} \left( \frac{N!}{\prod B^{(0)}} \right) + N(A_{p_1} - B_{p_1}) \right|_{0, \log p_1 + \kappa_{**}},$$

in Step (d) by at most

$$\frac{1}{N \log \sqrt{12}} \sum_{K < p_1 \leq t/K} \left| v_{p_1} \left( \frac{N!}{\prod B^{(0)}} \right) + N A_{p_1} \right|_{\log(t/p_1) + \kappa_{**}, 0},$$

and in Step (d') by at most the sum of

$$\frac{1}{N \log \sqrt{12}} \sum_{3 < p_1 \leq K} \left| \nu_{p_1} \left( \frac{N!}{\prod \mathcal{B}^{(0)}} \right) + N(A_{p_1} - B_{p_1}) \right|_{\log(K^2) + \kappa_{**}, 0}$$

and

$$\frac{1}{N \log \sqrt{12}} (\log t + \kappa_{**})$$

so by (8.25), (8.17), (8.18), (8.19), (8.20), and the triangle inequality we have

$$\frac{1}{N} \|(\nu_2(\prod \mathcal{B}^{(2)}), \nu_3(\prod \mathcal{B}^{(3)}))\|_\gamma \leq \sum_{i=1}^6 \alpha_i. \quad (8.28)$$

By construction, the multiset  $\mathcal{B}^{(2)}$  is  $t$ -admissible, and in balance at all small, medium, and large primes  $p > 3$ ; thus  $N! / \prod \mathcal{B}^{(2)} = 2^n 3^m$  for some integers  $n, m$ . From (8.28), (8.5), (2.3), (8.21) we have

$$\begin{aligned} n - 2\gamma_2 m &= \nu_2(N!) - 2\gamma_2 \nu_3(N!) - \left( \nu_2 \left( \prod \mathcal{B}^{(2)} \right) - 2\gamma_2 \nu_3 \left( \prod \mathcal{B}^{(2)} \right) \right) \\ &\geq \nu_2(N!) - 2\gamma_2 \nu_3(N!) - N(1 - \gamma_2) \sum_{i=1}^6 \alpha_i \\ &> N - \frac{\log N}{\log 2} - 1 - \gamma_2 N - N(1 - \gamma_2)(1 - \alpha_7) \\ &= N(1 - \gamma_2)\alpha_7 - \frac{\log(2N)}{\log 2} \\ &\geq 0 \end{aligned}$$

and similarly

$$\begin{aligned} 2m - \gamma_3 n &= 2\nu_3(N!) - \gamma_3 \nu_2(N!) - \left( 2\nu_3 \left( \prod \mathcal{B}^{(2)} \right) - \gamma_3 \nu_2 \left( \prod \mathcal{B}^{(2)} \right) \right) \\ &\geq 2\nu_3(N!) - \gamma_3 \nu_2(N!) - N(1 - \gamma_3) \sum_{i=1}^6 \alpha_i \\ &> N - \frac{\log N}{\log \sqrt{3}} - 2 - \gamma_3 N - N(1 - \gamma_3)(1 - \alpha_7) \\ &= N(1 - \gamma_3)\alpha_7 - \frac{\log(3N)}{\log \sqrt{3}} \\ &\geq 0. \end{aligned}$$

From (8.3) and Cramer's rule we conclude that that  $(n, 2m)$  lies in the non-negative linear span of  $(n_2, 2m_2)$ ,  $(n_3, 2m_3)$ , thus

$$(n, 2m) = \beta_2(n_2, 2m_2) + \beta_3(n_3, 2m_3) \quad (8.29)$$

for some reals  $\beta_2, \beta_3 \geq 0$ . We now create the multiset  $\mathcal{B}^{(3)}$  by adding  $\lfloor \beta_2 \rfloor$  copies of  $2^{n_2} 3^{m_2}$  and  $\lfloor \beta_3 \rfloor$  copies of  $2^{n_3} 3^{m_3}$  to  $\mathcal{B}^{(2)}$ . By (8.2), this multiset remains  $t$ -admissible, and each element

added increases the  $t$ -excess by at most  $\kappa_*$ . The number of such elements can be upper bounded using (8.29), (2.3) as

$$\begin{aligned}
 \lfloor \beta_2 \rfloor + \lfloor \beta_3 \rfloor &\leq \beta_2 + \beta_3 \\
 &\leq \frac{1}{\log t} (\beta_2(n_2 \log 2 + m_2 \log 3) + \beta_3(n_3 \log 2 + m_3 \log 3)) \\
 &= \frac{1}{\log t} (n \log 2 + m \log 3) \\
 &\leq \frac{1}{\log t} ((v_2(N!) - N B_2) \log 2 + (v_3(N!) - N B_3) \log 3) \\
 &\leq \frac{1}{\log t} \left( N \log 2 + \frac{N}{2} \log 3 - N B_2 \log 2 - N B_3 \log 3 \right) \\
 &= \frac{N \log \sqrt{12}}{\log t} - \frac{N(B_2 \log 2 + B_3 \log 3)}{\log t}.
 \end{aligned}$$

By (8.27), (8.12), we thus have

$$\frac{1}{N} \mathcal{E}_t(\mathcal{B}^{(3)}) \leq \sum_{i=1}^7 \delta_i. \quad (8.30)$$

Meanwhile by construction we see that  $\mathcal{B}^{(3)}$  is a subfactorization of  $N!$  that is in balance at all non-tiny primes, with tiny prime surpluses bounded by

$$v_2 \left( \frac{N!}{\prod \mathcal{B}^{(3)}} \right) \leq n_2 + n_3; \quad v_3 \left( \frac{N!}{\prod \mathcal{B}^{(3)}} \right) \leq m_2 + m_3.$$

and thus by (8.2), (8.13), we thus have

$$\frac{1}{N} \sum_p v_p \left( \frac{N!}{\prod \mathcal{B}^{(3)}} \right) \log p \leq \frac{\log 2^{n_2} 3^{m_2} + \log 2^{n_3} 3^{m_3}}{N} \leq \delta_8$$

and thus by (8.30), (8.4) we have

$$\mathcal{E}_t(\mathcal{B}^{(3)}) + \sum_p v_p \left( \frac{N!}{\prod \mathcal{B}^{(3)}} \right) \log p \leq \log N! - N \log t.$$

Applying Lemma 7.2, we conclude that  $t(N) \geq t$  as claimed.

## 9. ESTIMATING TERMS

In order to use Proposition 8.1 for a given choice of  $N, t$ , we need to find a  $t$ -admissible multiset  $\mathcal{B}^{(0)}$  and parameters  $K, \kappa_*, \gamma_2, \gamma_3$  obeying (8.1) as well as good upper bounds on the quantities  $\delta_i, i = 1, \dots, 8$  and  $\alpha_i, i = 1, \dots, 7$ , that can either be evaluated asymptotically or numerically. Many of these terms will be straightforward to estimate; we discuss only the more difficult ones here.

We introduce a further natural number parameter  $A$  and define

$$\sigma := \frac{3N}{At}. \quad (9.1)$$

We let  $\mathcal{B}^{(0)}$  be the multiset of 3-rough elements of the interval  $(t, t(1 + \sigma)]$ , with each element repeated precisely  $A$  times. This is clearly  $t$ -admissible. It has no presence at tiny primes, so

$$\alpha_1 = 0. \quad (9.2)$$

We will also introduce an auxiliary parameter  $L$  to assist us with the estimates. The influence of the parameters  $A, K, L$  on the other parameters  $\delta_i, \alpha_i$  (and  $\gamma_2, \gamma_3, \kappa_{**}$ ) can be roughly summarized as follows:

- $\gamma_2, \gamma_3 \asymp \log L / \log N$ ; assuming this quantity is small enough, we have  $\kappa_{**} \asymp 1$ .
- $\delta_1 \asymp 1/A$ .
- $\delta_2 \asymp 1/\log N$ .
- $\delta_3 \asymp A/K \log N$  and  $\alpha_3 \asymp A/K$ .
- $\delta_5 \asymp \log^{O(1)} K / \log^2 N$  and  $\alpha_2, \alpha_5 \asymp \log^{O(1)} K / \log N$ .
- $\delta_7 \asymp \kappa_L / \log N$ .
- $\delta_4, \delta_6, \delta_8, \alpha_4, \alpha_6, \alpha_7$  will be lower order terms.

We will quantify these relationships more precisely below, but they already suggest that one should take  $A$  to only be moderately large (e.g., of logarithmic size), that  $K$  should only be slightly larger than  $A$ , and that  $L$  should be significantly smaller than  $N$ .

We use the notation  $\sum^*$  to denote summation restricted to 3-rough numbers, thus for instance  $\sum_{a < k \leq b}^* 1$  denotes the number of 3-rough numbers in  $(a, b]$ . We have a simple estimate for such counts:

**Lemma 9.1.** *For any interval  $(a, b]$  with  $0 \leq a \leq b$  one has  $\sum_{a < k \leq b}^* 1 = \frac{b-a}{3} + O_{\leq}(4/3)$ .*

*Proof.* By the triangle inequality, it suffices to show that  $\sum_{0 < k \leq x}^* 1 - \frac{x}{3} = O_{\leq}(2/3)$  for all  $x \geq 0$ . This is easily verified for  $0 \leq x \leq 6$ , and the left-hand side is 6-periodic in  $x$ , giving the claim; see Figure 12.  $\square$

This lets us estimate  $\delta_1$ :

**Lemma 9.2.** *We have*

$$\delta_1 \leq \frac{3N}{2tA} + \frac{4}{N}.$$

*Proof.* By definition, we have

$$\mathcal{E}_t(\mathcal{B}^{(1)}) = A \sum_{t < n \leq t(1+\sigma)}^* \log \frac{n}{t}.$$

By the fundamental theorem of calculus, this is

$$A \int_0^{t\sigma} \sum_{t < n \leq t+h}^* 1 \frac{dh}{t+h}.$$


 FIGURE 12. The function  $\sum_{k \leq x}^* 1 - \frac{x}{3}$ .

Bounding  $\frac{1}{t+h}$  by  $\frac{1}{t}$  and applying Lemma 9.1, (9.1), we conclude that

$$\mathcal{E}_t(\mathcal{B}^{(1)}) \leq A \int_0^{3N/A} \left( \frac{h}{3} + \frac{4}{3} \right) \frac{dh}{t} = \frac{3N^2}{2tA} + 4.$$

and the claim follows.  $\square$

To construct  $\gamma_2, \gamma_3, \kappa_*, n_2, m_2, n_3, m_3$ , we introduce another parameter  $L \geq 1$  and assume that

$$t > 3L. \quad (9.3)$$

We define  $n_2, n_3, m_2, m_3$  by setting

$$2^{n_2} 3^{m_2} := 2^{n_0} \lceil t/2^{n_0} \rceil^{\langle 2,3 \rangle}; \quad 2^{n_3} 3^{m_3} := 3^{m_0} \lceil t/3^{m_0} \rceil^{\langle 2,3 \rangle}$$

where  $2^{n_0}, 3^{m_0}$  are the largest powers of 2, 3 respectively that are at most  $t/L$ . By construction and (2.6), (8.2) holds with

$$\kappa_* = \kappa_L. \quad (9.4)$$

We have

$$2m_2 \leq \frac{\log \lceil t/2^{n_0} \rceil^{\langle 2,3 \rangle}}{\log \sqrt{3}} \leq \frac{\log(2L) + \kappa_L}{\log \sqrt{3}}$$

and

$$n_2 \geq n_0 \geq \frac{\log t - \log(2L)}{\log 2};$$

similarly

$$n_3 \leq \frac{\log(3L) + \kappa_L}{\log 2}$$

and

$$2m_3 \geq \frac{\log t - \log(3L)}{\log \sqrt{3}}.$$

We conclude that (8.3) holds with

$$\begin{aligned} \gamma_2 &:= \frac{\log 2}{\log \sqrt{3}} \frac{\log(2L) + \kappa_L}{\log t - \log(2L)} \\ \gamma_3 &:= \frac{\log \sqrt{3}}{\log 2} \frac{\log(3L) + \kappa_L}{\log t - \log(3L)}; \end{aligned} \tag{9.5}$$

one can of course also take larger values of  $\gamma_2, \gamma_3$  if desired. This lets us compute the quantity  $\kappa_{**}$  defined in (8.22).

To estimate  $\delta_3, \alpha_3$  we use

**Lemma 9.3.** *For every  $3 < p \leq t/K$ , one has*

$$v_p\left(\frac{N!}{\prod B^{(1)}}\right) = O_{\leq}\left(\frac{4A+3}{3} \left\lceil \frac{\log N}{\log p} \right\rceil\right). \tag{9.6}$$

*Proof.* One has

$$\begin{aligned} v_p(\prod B^{(1)}) &= A \sum_{t < n \leq t(1+\sigma)}^* v_p(n) \\ &= A \sum_{1 \leq j \leq \frac{\log N}{\log p}} \sum_{t/p^j < n \leq t(1+\sigma)/p^j}^* 1 \\ &= A \sum_{1 \leq j \leq \frac{\log N}{\log p}} \left( \frac{N}{p^j A} + O_{\leq}(4/3) \right) \\ &= \frac{N}{p-1} - O_{\leq}^+\left(\frac{1}{p-1}\right) + O_{\leq}\left(\frac{4A}{3} \left\lceil \frac{\log N}{\log p} \right\rceil\right) \\ &= \frac{N}{p-1} - O_{\leq}^+\left(\left\lceil \frac{\log N}{\log p} \right\rceil\right) + O_{\leq}\left(\frac{4A}{3} \left\lceil \frac{\log N}{\log p} \right\rceil\right). \end{aligned}$$

Meanwhile, from (2.3) one has

$$v_p(N!) = \frac{N}{p-1} - O_{\leq}^+\left(\left\lceil \frac{\log N}{\log p} \right\rceil\right)$$

and the claim follows. □

**Corollary 9.4.** *One has*

$$\delta_3 \leq \frac{(4A+3)\kappa_{4.5}}{3N} \left( \pi\left(\frac{t}{K}\right) + \frac{\log N}{\log 5} \pi\left(\sqrt{N}\right) \right)$$

and

$$\alpha_3 \leq \frac{(4A+3) \left( \log \frac{t}{K} + \kappa_{**} \right)}{3N \log \sqrt{12}} \left( \pi\left(\frac{t}{K}\right) + \frac{\log N}{\log 5} \pi\left(\sqrt{N}\right) \right).$$

*Proof.* This is immediate from Lemma 9.3 and (8.17), (8.8) after noting that  $\lfloor \frac{\log N}{\log p} \rfloor \leq 1 + \frac{\log N}{\log 5} 1_{p \leq \sqrt{N}}$  for  $3 < p \leq t/K$ .  $\square$

The main quantities left to estimate are the quantities  $\delta_4, \delta_5, \alpha_4, \alpha_5$  that involve  $A_{p_1}$ . By construction of  $\mathcal{B}^{(0)}$ , we have

$$A_{p_1} = \frac{1}{N} \sum_m^* v_{p_1}(m) \sum_{\substack{\frac{t}{K}, \frac{t}{m} < p \leq \frac{t(1+\sigma)}{m}}} A.$$

In particular, for  $p > K(1 + \sigma)$  the quantity  $A_{p_1}$  vanishes entirely:

$$A_{p_1} = 0. \quad (9.7)$$

For the remaining primes  $3 < p \leq K(1 + \sigma)$  one has

$$A_{p_1} = \frac{A}{N} \sum_{m \leq K(1+\sigma)}^* v_{p_1}(m) \left( \pi \left( \frac{t(1+\sigma)}{m} \right) - \pi \left( \frac{t}{\min(m, K)} \right) \right). \quad (9.8)$$

In practice, these expressions can be adequately controlled by Lemma 2.2, as can the quantities  $B_{p_1}$ .

## 10. THE ASYMPTOTIC REGIME

With the above estimates, we can now establish the lower bound in Theorem 1.3(iv). Thus we aim to show that  $t(N) \geq t$  for sufficiently large  $N$ , where

$$t := \frac{N}{e} - \frac{c_0 N}{\log N} + \frac{N}{\log^{1+c_1} N} \asymp N \quad (10.1)$$

and  $0 < c_1 < 1$  is a small absolute constant. We use the construction of the previous section with the parameters

$$A := \lfloor \log^2 N \rfloor \quad (10.2)$$

$$K := \lfloor \log^3 N \rfloor \quad (10.3)$$

$$L := N^{0.1}, \quad (10.4)$$

so from (9.1) one has

$$\sigma = \frac{3N}{tA} \asymp \frac{1}{A} \asymp \frac{1}{\log^2 N}. \quad (10.5)$$

The conditions (8.1), (9.3) are easily verified for  $N$  large enough.

By (9.4), (10.4), and Lemma 2.1(ii) we have

$$\kappa_* \ll \log^{-c} N$$

for some absolute constant  $c > 0$ . From (9.5), (10.1), (10.4) we have

$$\gamma_2 = \frac{1}{10} \frac{\log 2}{\log \sqrt{3}} + O \left( \frac{1}{\log N} \right), \quad \gamma_3 = \frac{1}{10} \frac{\log \sqrt{3}}{\log 2} + O \left( \frac{1}{\log N} \right)$$

and hence by (8.22), (2.10), (2.11) we have for sufficiently large  $N$  that

$$\kappa_{**} \ll 1.$$

By Proposition 8.1, it thus suffices to establish the inequalities (8.4), (8.5). Several of the quantities  $\delta, \delta_i, \alpha_i$  can now be immediately estimated using (9.2), (9.2), Corollary 9.4, (2.4), and the prime number theorem:

$$\begin{aligned} \delta_1 &\ll \frac{1}{A} \asymp \frac{1}{\log^2 N} \\ \delta_3 &\ll \frac{A}{K \log N} \asymp \frac{1}{\log^2 N} \\ \delta_6 &\ll \frac{1}{N} \\ \delta_7 &\ll \frac{\kappa_*}{\log N} \ll \frac{1}{\log^{1+c} N} \\ \delta_8 &\ll \frac{\log N}{N} \\ \delta &= \frac{ec_0}{\log N} + \frac{e}{\log^{1+c_1} N} + O\left(\frac{1}{\log^2 N}\right) \end{aligned}$$

$$\begin{aligned} \alpha_1 &= 0 \\ \alpha_3 &\ll \frac{A}{K} \asymp \frac{1}{\log N} \\ \alpha_6, \alpha_7 &\ll \frac{\log N}{N} \end{aligned}$$

On the interval  $(t/NK, 1]$ , the function  $f_{N/t}$  is piecewise monotone with  $O(K)$  pieces, and bounded by 1, so its augmented total variation norm is  $O(K)$ . Applying (8.7) and Lemma 2.2 (with classical error term), we have

$$\begin{aligned} \delta_2 &\leq \frac{1}{\log(t/K)} \int_{t/NK}^1 f_{N/t}(x) dx + O\left(\frac{1}{\log^2 N}\right) \\ &\leq \frac{1}{\log N} \int_{1/eK}^{N/et} f_{N/t}(etx/N) dx + O\left(\frac{1}{\log^2 N}\right) \end{aligned}$$

where we have used (1.9) to manage error terms. Similarly to the proof of Proposition 5.2, the function  $f_{N/t}(etx/N)$  differs from  $f_e(x)$  outside of an exceptional set of measure  $O(1/\log N)$ , and hence by (1.6) (and (1.9)) we have

$$\delta_2 \leq \frac{ec_0}{\log N} + O\left(\frac{1}{\log^2 N}\right).$$

To finish the verification of the conditions (8.4), (8.5), it will suffice to show that

$$\delta_4, \delta_5 \ll \frac{(\log \log N)^{O(1)}}{\log^2 N} \tag{10.6}$$



and

$$\alpha_2, \alpha_4, \alpha_5 \ll \frac{(\log \log N)^{O(1)}}{\log N}. \quad (10.7)$$

By Mertens' theorem (or Lemma 2.2) and (8.9), (8.10), (8.16), (8.18), (8.19), (10.5), it suffices to show that

$$A_{p_1}, B_{p_1} \ll \frac{(\log \log N)^{O(1)}}{p_1 \log N} \quad (10.8)$$

for all  $p_1 \leq K(1 + \sigma)$  (recalling from (9.7) that  $A_{p_1}$  vanishes for any larger  $p_1$ ), as well as the variant

$$|A_{p_1} - B_{p_1}|_{0,1} \ll \frac{(\log \log N)^{O(1)}}{p_1 \log^2 N} \quad (10.9)$$

for  $3 < p_1 \leq K$ .

For (10.8) we use (9.8), (8.24), and the crude bound

$$\nu_{p_1}(m) \ll 1_{p_1|m} \log \log N \quad (10.10)$$

for  $m \leq K(1 + \sigma)$ , and reduce to showing that

$$\frac{A}{N} \sum_{m \leq K(1+\sigma)} 1_{p_1|m} \left( \pi \left( \frac{t(1+\sigma)}{m} \right) - \pi \left( \frac{t}{\min(m, K)} \right) \right) \ll \frac{(\log \log N)^{O(1)}}{p_1 \log N}$$

and

$$\frac{1}{N} \sum_{m \leq K} 1_{p_1|m} \sum_{\substack{t \\ \frac{t}{m} \leq p < \frac{t}{m-1}}} \left\lfloor \frac{N}{p} \right\rfloor \ll \frac{(\log \log N)^{O(1)}}{p_1 \log N}.$$

But from the Brun–Titchmarsh inequality (or Lemma 2.2) and (10.5) one has

$$\pi \left( \frac{t(1+\sigma)}{m} \right) - \pi \left( \frac{t}{\min(m, K)} \right) \ll \frac{t\sigma}{m \log N} \ll \frac{N}{Am \log N}$$

and

$$\sum_{\substack{t \\ \frac{t}{m} \leq p < \frac{t}{m-1}}} \left\lfloor \frac{N}{p} \right\rfloor \ll \frac{tm}{m^2 \log N} \ll \frac{N}{m \log N}$$

and the claim then follows from summing the harmonic series.

It remains to show (10.9). For  $3 < p_1 \leq K$ , we see from (9.8), (10.5), (10.10) and Lemma 2.2 (with classical error term) that

$$\begin{aligned} A_{p_1} &\geq \frac{1}{N} \sum_{m \leq K(1+\sigma)}^* \nu_{p_1}(m) \left( \frac{At\sigma}{m \log N} + O \left( \frac{(\log \log N)^{O(1)} At\sigma}{m \log^2 N} \right) \right) \\ &= \frac{1}{\log N} \sum_{m \leq K(1+\sigma)}^* \nu_{p_1}(m) \frac{3}{m} + O \left( \frac{(\log \log N)^{O(1)}}{\log^2 N} \right) \\ &= \frac{1}{\log N} \sum_{m \leq K}^* \nu_{p_1}(m) \frac{3}{m} + O \left( \frac{(\log \log N)^{O(1)}}{\log^2 N} \right) \end{aligned}$$

and similarly from (8.24), (10.10), and Lemma 2.2 (again with classical error term)

$$\begin{aligned}
B_{p_1} &\leq \frac{1}{N} \sum_{m \leq K} v_{p_1}(m) \sum_{\frac{t}{m} \leq p < \frac{t}{m-1}} \frac{N}{p} \\
&\leq \frac{1}{N} \sum_{m \leq K} v_{p_1}(m) \left( \frac{N}{\log(t/m)} \int_{t/m}^{t/(m-1)} \frac{dx}{x} + O\left(\frac{N}{\log^{10} N}\right) \right) \\
&\leq \frac{1}{\log N} \sum_{m \leq K} v_{p_1}(m) \log \frac{m}{m-1} + O\left(\frac{(\log \log N)^{O(1)}}{\log^2 N}\right)
\end{aligned}$$

so it will suffice to establish the inequality

$$\sum_{m \leq K} v_{p_1}(m) \log \frac{m}{m-1} \leq \sum_{m \leq K}^* v_{p_1}(m) \frac{3}{m} \quad (10.11)$$

for all  $p_1 > 3$ .

Writing  $v_{p_1}(m) = \sum_{j \geq 1} 1_{p_1^j | m}$ , it suffices to show that

$$\sum_{m \leq K; p_1^j | m} \frac{3}{m} 1_{(m,6)=1} - \log \frac{m}{m-1} \geq 0.$$

Making the change of variables  $m = p_1^j n$ , it suffices to show that

$$\sum_{n \leq K'} \frac{3}{n} 1_{(n,6)=1} - p_1^j \log \frac{p_1^j n}{p_1^j n - 1} \geq 0$$

for any  $K' > 0$ . Using the bound

$$\log \frac{p_1^j n}{p_1^j n - 1} = \int_{p_1^j n - 1}^{p_1^j n} \frac{dx}{x} \leq \frac{1}{p_1^j n - 1}$$

and  $p_1^j \geq 5$ , we have

$$p_1^j \log \frac{p_1^j n}{p_1^j n - 1} \leq \frac{1}{n - 0.2}$$

and so it suffices to show that

$$\sum_{n \leq K'}^* \frac{3}{n} 1_{(n,6)=1} - \frac{1}{n - 0.2} \geq 0. \quad (10.12)$$

Since

$$\sum_{n=1}^{\infty} \frac{1}{n - 0.2} - \frac{1}{n} = \psi(0.8) - \psi(1) = 0.353473 \dots,$$

where  $\psi$  here denotes the digamma function rather than the von Mangoldt summatory function, it will suffice to show that

$$\sum_{n \leq K'} \frac{3}{n} 1_{(n,6)=1} - \frac{1}{n} \geq 0.4. \quad (10.13)$$



FIGURE 13. A plot of (10.12), (10.13).

This can be numerically verified for  $K' \leq 100$ , with substantial room to spare for  $K'$  large; see Figure 13. On a block  $6a - 1 \leq n \leq 6a + 4$  with  $a > 1$ , the sum is positive:

$$\begin{aligned} \sum_{6a-1 \leq n \leq 6a+4}^* \frac{3}{n} - \frac{1}{n} &= \left( \frac{1}{6a-1} - \frac{1}{6a} \right) + \left( \frac{1}{6a-1} - \frac{1}{6a+2} \right) \\ &\quad + \left( \frac{1}{6a+1} - \frac{1}{6a+3} \right) + \left( \frac{1}{6a+1} - \frac{1}{6a+4} \right) \\ &> 0. \end{aligned}$$

The inequality for  $K' > 100$  is then easily verified from the  $K' \leq 100$  data and the triangle inequality.

## 11. GUY–SELFIDGE CONJECTURE

We now establish the Guy–Selfridge conjecture  $t(N) \geq N/3$  in the range

$$N \geq N_0 := 10^{11}.$$

We will apply Proposition 8.1 with the construction in Section 9 and the choice of parameters

$$t := N/3$$

$$A := 189$$

$$K := 293$$

$$L := 4.5;$$

the choice of  $A$  and  $K$  was obtained after some numerical experimentation. In particular, by (9.1) we have

$$\sigma = \frac{3N}{At} = \frac{9}{189} = 0.047619 \dots$$

One can readily check the required conditions (8.1), (9.3) for  $N \geq N_0$ , so it remains to verify the hypotheses (8.4), (8.5) of Proposition 8.1 in this range. Some of the quantities in these hypotheses involve sums over large ranges, such as  $(t/K, N]$ ; but one can use Lemma 2.2 to obtain adequate upper or lower bounds on such quantities, leaving one with sums over short

ranges such as  $p \leq K$  or  $p \leq K(1 + \sigma)$ . As such, all of the bounds needed can be quickly computed even for very large  $N$  with simple computer code<sup>8</sup>.

Many of the bounds we will use will be monotone decreasing in  $N$ , so that they only need to be tested at the left endpoint  $N = N_0$ . However, this is not the case for all of the bounds, as some involve subtracting one monotone quantity from another. For those estimates, we will initially only establish bounds in two extreme cases,  $N = N_0$  and  $N \geq 10^{70}$ , and discuss how to cover the intervening ranges  $N_0 < N < 10^{70}$  at the end of the section.

We now bound some of the terms appearing in Proposition 8.1. From Lemma 2.1 we have

$$\kappa_{4.5} = \log \frac{4}{3} = 0.28768 \dots$$

From (9.5) one can take

$$\gamma_2 := \frac{\log 2}{\log \sqrt{3}} \frac{\log(2L) + \kappa_L}{\log(N_0/3) - \log(2L)} = 0.1423165 \dots$$

and

$$\gamma_3 := \frac{\log \sqrt{3}}{\log 2} \frac{\log(3L) + \kappa_L}{\log(N_0/3) - \log(3L)} = 0.1059116 \dots$$

for all  $N \geq N_0$ ; by (8.22) and some calculation we then have

$$\kappa_{**} \leq 6.830101 \dots$$

From (2.4) one has

$$\delta \geq \log N - \log t = \log \frac{3}{e} = 0.0986122 \dots$$

for all  $N \geq N_0$ .

We will use this lower bound as our unit of reference for all other  $\delta_i$  quantities, bounding them by suitable multiples of  $\delta$ . For instance, from (9.2) one has

$$\delta_1 \leq \frac{9}{2A} + \frac{4}{N_0} \leq 0.241447\delta$$

for all  $N \geq N_0$ .

From (8.7) and Lemma 2.2, and the monotonicity of  $E(N)/N$ , one has

$$\begin{aligned} \delta_2 &\leq \frac{\int_{1/3K}^1 f_3(x) dx}{\log(t/K)} + \frac{\|f_3\|_{\text{TV}((1/3K, 1])} E(N)}{\log(t/K) N} \\ &\leq \frac{0.919785}{\log(N_0/3K)} + \frac{1159.795}{\log(N_0/3K)} \frac{E(N_0)}{N_0} \\ &\leq 0.504735\delta \end{aligned}$$

<sup>8</sup>[https://github.com/teorth/erdos-guy-selfridge/blob/main/src/python/interval\\_computations.py](https://github.com/teorth/erdos-guy-selfridge/blob/main/src/python/interval_computations.py)

for all<sup>9</sup>  $N \geq N_0$ . For  $N \geq 10^{70}$  we may replace  $N_0$  by  $10^{70}$  and obtain the significantly better bound

$$\delta_2 \leq 0.060410\delta.$$

From Corollary 9.4 and (2.13) one has

$$\begin{aligned} \delta_3 &\leq \frac{(4A+3)\kappa_{4.5}}{3} \left( \frac{1}{3K \log \frac{N}{3K}} + \frac{1.2762}{3K \log^2 \frac{N}{3K}} + \frac{\log N}{\sqrt{N} \log \sqrt{N} \log 5} + \frac{1.2762 \log N}{\sqrt{N} \log^2 \sqrt{N} \log 5} \right) \\ &\leq \frac{(4A+3)\kappa_{4.5}}{3} \left( \frac{1}{3K \log \frac{N_0}{3K}} + \frac{1.2762}{3K \log^2 \frac{N_0}{3K}} + \frac{\log N_0}{\sqrt{N_0} \log \sqrt{N_0} \log 5} + \frac{1.2762 \log N_0}{\sqrt{N_0} \log^2 \sqrt{N_0} \log 5} \right) \\ &\leq 0.051574\delta \end{aligned}$$

for all  $N \geq N_0$ .

We skip  $\delta_4, \delta_5, \delta_7$  for now. From (8.11) we have

$$\delta_6 \leq \frac{\kappa_{4.5}}{N_0} \leq 3 \times 10^{-11} \delta$$

for all  $N \geq N_0$ , and from (8.13) we have

$$\delta_8 \leq \frac{2(\log(N_0/3) + \kappa_{4.5})}{N_0} \leq 6 \times 10^{-10} \delta$$

for all  $N \geq N_0$ , so these two terms are negligible in the analysis.

From (9.2) we have

$$\alpha_1 = 0.$$

We skip  $\alpha_2, \alpha_4, \alpha_5$  for now. From Corollary 9.4 and (2.13) one has

$$\begin{aligned} \alpha_3 &\leq \frac{4A+3}{3 \log \sqrt{12}} \left( \log \frac{N}{3K} + \kappa_{**} \right) \\ &\quad \times \left( \frac{1}{3K \log \frac{N}{3K}} + \frac{1.2762}{3K \log^2 \frac{N}{3K}} + \frac{\log N}{\sqrt{N} \log \sqrt{N} \log 5} + \frac{1.2762 \log N}{\sqrt{N} \log^2 \sqrt{N} \log 5} \right). \end{aligned}$$

Expanding out the product, one can check that all terms are non-increasing in  $N$ ; so we may substitute  $N_0$  for  $N$  in the right-hand side, which after some calculation gives

$$\alpha_3 \leq 0.361121$$

for all  $N \geq N_0$ . From (8.20) we have

$$\begin{aligned} \alpha_6 &\leq \frac{\log(N_0/3) + \kappa_{**}}{N_0 \log \sqrt{12}} \\ &\leq 3 \times 10^{-10} \end{aligned}$$

---

<sup>9</sup>Despite the seemingly large numerator, the second term is in fact negligible in the regime  $N \geq N_0$ , due to the square root type decay in  $E(N)/N$ .

for all  $N \geq N_0$ , and similarly from (8.21) we have

$$\alpha_7 \leq \max \left( \frac{\log(2N_0)}{(1-\gamma_2)N_0 \log 2}, \frac{\log(3N_0)}{(1-\gamma_3)N_0 \log \sqrt{3}} \right) \leq 6 \times 10^{-10}$$

for all  $N \geq N_0$ . so the contribution of these two terms are negligible.

Conveniently<sup>10</sup>, the choice of parameters  $A, K$  ensure that there are no primes in the range

$$293 = K < p \leq K(1 + \sigma) = K(1 + \sigma) = 306.952 \dots$$

and thus

$$\delta_4 = \alpha_4 = 0$$

for all  $N \geq N_0$ .

The remaining terms  $\delta_5, \delta_7, \alpha_2, \alpha_5$  to estimate involve the quantities  $A_{p_1}, B_{p_1}$  defined in (8.23), (8.24), and require a bit more care. For  $B_{p_1}$ , we can split the expression as

$$B_{p_1} = \sum_{m \leq K} v_{p_1}(m) \sum_{k: a_{k,m} < b_{k,m}} \frac{k}{N} (\pi(Nb_{k,m}) - \pi(Na_{k,m}))$$

where

$$a_{k,m} := \max \left( \frac{1}{3m} -, \frac{1}{k} \right); \quad b_{k,m} := \max \left( \frac{1}{3(m-1)} -, \frac{1}{k-1} \right)$$

where the  $-$  denotes the subtraction of an infinitesimal quantity to reflect the restriction to the range  $\frac{t}{m} \leq p < \frac{t}{m-1}$  rather than  $\frac{t}{m} < p \leq \frac{t}{m-1}$ . Using Lemma 2.2 (and a limiting argument), we can upper bound this quantity by

$$B_{p_1} \leq \sum_{m \leq K} v_{p_1}(m) \sum_{k: a_{k,m} < b_{k,m}} \left( \frac{k}{2 \log(Na_{k,m})} + \frac{k}{2 \log(Nb_{k,m})} \right) (a_{k,m} - b_{k,m}) + 2 \frac{E(Nb_{k,m})}{N \log(Na_{k,m})}$$

and lower bound it by

$$B_{p_1} \geq \sum_{m \leq K} v_{p_1}(m) \sum_{k: a_{k,m} < b_{k,m}} \frac{k \left( 1 - \frac{2}{\sqrt{a_{k,m}N}} \right)}{\log(N(a_{k,m} + b_{k,m})/2)} (a_{k,m} - b_{k,m}) - 2 \frac{E(Nb_{k,m})}{N \log(Na_{k,m})}$$

We caution here that while the upper bound for  $B_{p_1}$  is monotone decreasing in  $N$ , the lower bound does not have a favorable monotonicity property, particularly as it will be used when *subtracting* copies of  $B_{p_1}$  rather than *adding* then.

From the monotonicity of the upper bound, one can use (8.16) to calculate that

$$\alpha_2 \leq 0.269878$$

for all  $N \geq N_0$ . For (8.12), subtraction is involved, and one must proceed with more caution. For  $N = N_0$ , one has

$$\delta_7 \leq 0.11359\delta.$$

<sup>10</sup>Even if this were not the case, the quantities  $\delta_4, \alpha_4$  should be viewed as lower order terms, and are far smaller than several of the other  $\delta_i$  or  $\alpha_i$  for typical choices of parameters.

For  $N \geq 10^{70}$ , we simply discard the negative terms here and obtain the bound

$$\delta_7 \leq \frac{\kappa_{4.5} \log \sqrt{12}}{\log(10^{70}/3)} \leq 0.02212\delta$$

As for the  $A_{p_1}$ , we know from (9.7) that this vanishes unless  $3 < p_1 \leq K(1 + \sigma)$ . From (9.8) and Lemma 2.2 one has the upper bound

$$\begin{aligned} A_{p_1} \leq & \sum_{m \leq K(1+\sigma)}^* \left( \frac{Av_{p_1}(m)}{2 \log(N/3 \min(m, K))} + \frac{Av_{p_1}(m)}{2 \log(N(1 + \sigma)/3m)} \right) \left( \frac{1 + \sigma}{3m} - \frac{1}{3 \min(m, K)} \right) \\ & + \frac{Av_{p_1}(m)}{\log(N/3 \min(m, K))} \frac{2E(N(1 + \sigma)/3m)}{N} \end{aligned}$$

and the lower bound

$$\begin{aligned} A_{p_1} \geq & \sum_{m \leq K(1+\sigma)}^* \frac{Av_{p_1}(m) \left( 1 - \frac{2}{\sqrt{N_0/3 \min(m, K)}} \right)}{\log((N/3 \min(m, K) + N(1 + \sigma)/3m)/2)} \left( \frac{1 + \sigma}{3m} - \frac{1}{3 \min(m, K)} \right) \\ & - \frac{Av_{p_1}(m)}{\log(N/3 \min(m, K))} \frac{2E(N(1 + \sigma)/3m)}{N}. \end{aligned}$$

Again, the upper bound is monotone decreasing in  $N$ , but the lower bound does not have a favorable monotonicity. At  $N = N_0$ , one can calculate using these bounds and (8.10), (8.19) to obtain

$$\delta_5 \leq 0.06203\delta$$

$$\alpha_5 \leq 0.31418$$

which, when combined with the previous bounds, gives

$$\sum_{i=1}^8 \delta_i \leq 0.9740\delta$$

and

$$\sum_{i=1}^7 \alpha_i \leq 0.9452$$

at  $N = N_0$ , thus verifying (8.4), (8.5) in those cases.

For  $N \geq 10^{70}$ , we use the triangle inequality to crudely upper bound

$$\delta_5 \leq \kappa_{4.5} \sum_{3 < p_1 \leq K} \frac{\log p_1}{\log(t/K^2)} A_{p_1} + B_{p_1}$$

and

$$\alpha_5 \leq \frac{1}{\log \sqrt{12}} \sum_{3 < p_1 \leq K} \frac{(\log K^2 + \kappa_{**}) \log p_1}{\log(t/K^2)} A_{p_1} + (\log p_1 + \kappa_{**}) B_{p_1}.$$

The bounds available for the right-hand side are now monotone in  $N$ , and one can calculate that

$$\begin{aligned}\delta_5 &\leq 0.077301\delta \\ \alpha_5 &\leq 0.184975\end{aligned}$$

for  $N \geq 10^{70}$ . This is better than the previous bound for  $\alpha_5$ . For  $\delta_5$ , the bound is slightly worse, but this is more than compensated for by the improved bounds on  $\delta_2$ ,  $\delta_7$ , and (8.4), (8.5) can be verified here with significant room to spare.

This completes the proof of Theorem 1.3(iii) (and hence Theorem 1.3(ii)) in the cases  $N = N_0$  and  $N \geq 10^{70}$ . It remains to cover the intermediate range  $N_0 < N \leq 10^{70}$ . Here we adopt the perspective of interval arithmetic. If  $N$  is constrained to a given interval, such as  $[10^{11}, 5 \times 10^{11}]$ , we can use the worst-case upper and lower bounds for  $A_{p_1}, B_{p_1}$  to obtain conservative upper bounds on the most delicate quantities  $\delta_5, \delta_7, \alpha_2, \alpha_5$ , thus potentially verifying the conditions (8.4), (8.5) simultaneously for all  $N$  in such an interval. As it turns out, there is enough room to spare in these estimates, particularly for large  $N$ , that this strategy works using only a small number of intervals; specifically, by considering  $N$  in the intervals

$$[10^{11}, 5 \times 10^{11}]; \quad [5 \times 10^{11}, 10^{14}]; \quad [10^{14}, 10^{20}]; \quad [10^{20}, 10^{70}]$$

one can check that such bounds are sufficient to verify (8.4), (8.5) in these cases. This now verifies Theorem 1.3(ii), (iii) for all  $N \geq 10^{11}$ . (In fact, with more effort, this verification can be pushed down to  $N \geq 6 \times 10^{10}$  using the same choice of parameters  $A, K, L$ .)

#### APPENDIX A. DISTANCE TO THE NEXT 3-SMOOTH NUMBER

We now establish the various claims in Lemma 2.1. We begin with part (iii). The claim (2.7) is immediate from (2.6), (2.5). Now prove (2.8), (2.9). If we write  $\lceil x/12^a \rceil^{(2,3)} = 2^b 3^c$ , then by (2.6) we have

$$b \log 2 + c \log 3 \leq \log x - a \log 12 + \kappa_L,$$

while from definition of  $a$  we have

$$\log x - a \log 12 \leq \log(12L). \tag{A.1}$$

We now compute

$$\begin{aligned}\frac{\nu_2(\lceil x \rceil_L^{(2,3)}) - 2\gamma \nu_3(\lceil x \rceil_L^{(2,3)})}{1 - \gamma} &= \frac{2a + b - 2\gamma(a + c)}{1 - \gamma} \\ &\leq 2a + \frac{\log x - a \log 12 + \kappa_L}{(1 - \gamma) \log 2} \\ &= \frac{\log x}{\log \sqrt{12}} + \left( \frac{1}{(1 - \gamma) \log 2} - \frac{1}{\log \sqrt{12}} \right) (\log x - a \log 12) \\ &\quad + \frac{\kappa_L}{(1 - \gamma) \log 2}\end{aligned}$$



giving (2.8) from (A.1); similarly, we have

$$\begin{aligned} \frac{2v_3(\lceil x \rceil_L^{(2,3)}) - \gamma v_2(\lceil x \rceil_L^{(2,3)})}{1 - \gamma} &= \frac{2(a + c) - \gamma(2a + b)}{1 - \gamma} \\ &\leq 2a + \frac{2(\log x - a \log 12 + \kappa_L)}{(1 - \gamma) \log 3} \\ &= \frac{\log x}{\log \sqrt{12}} + \left( \frac{2}{(1 - \gamma) \log 3} - \frac{1}{\log \sqrt{12}} \right) (\log x - a \log 12) \\ &\quad + \frac{\kappa_L}{(1 - \gamma) \log \sqrt{3}} \end{aligned}$$

giving (2.9) from (A.1).

To prove parts (i) and (ii) of Lemma 2.1, we establish the following lemma to upper bound  $\kappa_L$ .

**Lemma A.1.** *If  $n_1, n_2, m_1, m_2$  are natural numbers such that  $n_1 + n_2, m_1 + m_2 \geq 1$  and*

$$\frac{3^{m_1}}{2^{n_1}}, \frac{2^{n_2}}{3^{m_2}} \geq 1$$

*then*

$$\kappa_{\min(2^{n_1+n_2}, 3^{m_1+m_2})/6} \leq \log \max \left( \frac{3^{m_1}}{2^{n_1}}, \frac{2^{n_2}}{3^{m_2}} \right).$$

*Proof.* If  $\min(2^{n_1+n_2}, 3^{m_1+m_2})/6 \leq t \leq 2^{n_2-1}3^{m_1-1}$ , then we have

$$t \leq 2^{n_2-1}3^{m_1-1} \leq \max \left( \frac{3^{m_1}}{2^{n_1}}, \frac{2^{n_2}}{3^{m_2}} \right) t, \quad (\text{A.2})$$

so we are done in this case. Now suppose that  $t > 2^{n_2-1}3^{m_1-1}$ . If we write  $\lceil t \rceil^{(2,3)} = 2^n 3^m$  be the smallest 3-smooth number that is at least  $t$ , then we must have  $n \geq n_2$  or  $m \geq m_1$  (or both). Thus at least one of  $\frac{2^{n_1}}{3^{m_1}} 2^n 3^m$  and  $\frac{3^{m_2}}{2^{n_2}} 2^n 3^m$  is an integer, and is thus at most  $t$  by construction. This gives (A.2), and the claim follows.  $\square$

Some efficient choices of parameters for this lemma are given in Table 4. For instance,  $\kappa_{4.5} \leq \log \frac{4}{3} = 0.28768 \dots$  and  $\kappa_{40.5} \leq \log \frac{32}{27} = 0.16989 \dots$ . In fact, since  $\lceil 4.5 + \varepsilon \rceil^{(2,3)} = 6$  and  $\lceil 40.5 + \varepsilon \rceil^{(2,3)} = 48$  for all sufficiently small  $\varepsilon > 0$ , we see that these bounds are sharp (And similarly for the other entries in Table 4); this establishes part (i).

**Remark A.2.** It should be unsurprising that the continued fraction convergents  $1/1, 2/1, 3/2, 8/5, 19/12, \dots$  to

$$\frac{\log 3}{\log 2} = 1.5849 \dots = [1; 1, 1, 2, 2, 3, 1, \dots]$$

are often excellent choices for  $n_1/m_1$  or  $n_2/m_2$ , although other approximants such as  $5/3$  or  $11/7$  are also usable.

$n_1$	$m_1$	$n_2$	$m_2$	$\min(2^{n_1+n_2}, 3^{m_1+m_2})/6$	$\log \max(3^{m_1}/2^{n_1}, 2^{n_2}/3^{m_2})$
1	1	<b>1</b>	<b>0</b>	$1/2 = 0.5$	$\log 2 = 0.69314 \dots$
<b>1</b>	<b>1</b>	2	1	$2^2/3 = 1.33 \dots$	$\log(3/2) = 0.40546 \dots$
3	2	<b>2</b>	<b>1</b>	$3^2/2 = 4.5$	$\log(2^2/3) = 0.28768 \dots$
3	2	<b>5</b>	<b>3</b>	$3^4/2 = 40.5$	$\log(2^5/3^3) = 0.16989 \dots$
<b>3</b>	<b>2</b>	8	5	$2^{10}/3 = 341.33 \dots$	$\log(3^2/2^3) = 0.11778 \dots$
<b>11</b>	<b>7</b>	8	5	$2^{18}/3 = 87381.33 \dots$	$\log(3^7/2^{11}) = 0.06566 \dots$
19	12	<b>8</b>	<b>5</b>	$3^{17}/2 \approx 6.4 \times 10^7$	$\log(2^8/3^5) = 0.05211 \dots$
19	12	<b>27</b>	<b>17</b>	$3^{29}/2 \approx 3.4 \times 10^{13}$	$\log(2^{27}/3^{17}) = 0.03856 \dots$
19	12	<b>46</b>	<b>29</b>	$3^{41}/2 \approx 1.8 \times 10^{19}$	$\log(2^{46}/3^{29}) = 0.02501 \dots$

TABLE 4. Efficient parameter choices for Lemma A.1. The parameters used to attain the minimum or maximum are indicated in **boldface**. Note how the number of rows in each group matches the terms 1, 1, 2, 2, 3, ... in the continued fraction expansion.

Finally, we establish (ii). From the classical theory of continued fractions, we can find rational approximants

$$\frac{p_{2j}}{q_{2j}} \leq \frac{\log 3}{\log 2} \leq \frac{p_{2j+1}}{q_{2j+1}} \quad (\text{A.3})$$

to the irrational number  $\log 3 / \log 2$ , where the convergents  $p_j/q_j$  obey the recursions

$$p_j = b_j p_{j-1} + p_{j-2}, \quad q_j = b_j q_{j-1} + q_{j-2}$$

with  $p_{-1} = 1, q_{-1} = -1 = 0, p_0 = b_0, q_0 = 1$ , and

$$[b_0; b_1, b_2, \dots] = [1; 1, 1, 2, 2, 3, 1 \dots]$$

is the continued fraction expansion of  $\frac{\log 3}{\log 2}$ . Furthermore,  $p_{2j+1}q_{2j} - p_{2j}q_{2j+1} = 1$ , and hence

$$\frac{\log 3}{\log 2} - \frac{p_{2j}}{q_{2j}} = \frac{1}{q_{2j}q_{2j+1}}. \quad (\text{A.4})$$

By Baker's theorem (see, e.g., [3]),  $\frac{\log 3}{\log 2}$  is a Diophantine number, giving a bound of the form

$$q_{2j+1} \ll q_{2j}^{O(1)} \quad (\text{A.5})$$

and a similar argument (using  $p_{2j+2}q_{2j+1} - p_{2j+1}q_{2j+2} = -1$ ) gives

$$q_{2j+2} \ll q_{2j+1}^{O(1)}. \quad (\text{A.6})$$

We can rewrite (A.3) as

$$\frac{3^{q_{2j}}}{2^{p_{2j}}}, \frac{2^{p_{2j+1}}}{3^{q_{2j+1}}} \geq 1$$

and routine Taylor expansion using (A.4) gives the upper bounds

$$\frac{3^{q_{2j}}}{2^{p_{2j}}}, \frac{2^{p_{2j+1}}}{3^{q_{2j+1}}} \leq \exp\left(O\left(\frac{1}{q_{2j}}\right)\right).$$

From Lemma A.1 we obtain

$$K_{\min(2^{p_{2j}+p_{2j+1}}, 3^{q_{2j}+q_{2j+1}})/6} \ll \frac{1}{q_{2j}}.$$

The claim then follows from (A.5), (A.6) (and the obvious fact that  $\kappa$  is monotone non-increasing after optimizing in  $j$ ).

**Remark A.3.** It seems reasonable to conjecture that  $c$  can be taken to be arbitrarily close to 1, but this is essentially equivalent to the open problem of determining that the irrationality measure of  $\log 3 / \log 2$  is equal to 2.

## APPENDIX B. ESTIMATING SUMS OVER PRIMES

In this appendix we establish Lemma 2.2. The key tool is

**Lemma B.1** (Integration by parts). *Let  $(y, x]$  be a half-open interval in  $(0, +\infty)$ . Suppose that one has a function  $a : \mathbb{N} \rightarrow \mathbb{R}$  and a continuous function  $f : (y, x] \rightarrow \mathbb{R}$  such that*

$$\sum_{y < n \leq z} a_n = \int_z^y f(t) dt + C + O_{\leq}(A)$$

*for all  $y \leq z \leq x$ , and some  $C \in \mathbb{R}$ ,  $A > 0$ . Then, for any function  $b : (y, x] \rightarrow \mathbb{R}$  of bounded total variation, one has*

$$\sum_{y < n \leq x} b(n)a_n = \int_x^y b(t)f(t) dt + O_{\leq}(A\|b\|_{\text{TV}^*(y,x]}). \quad (\text{B.1})$$

*Proof.* If, for every natural number  $y < n \leq x$ , one modifies  $b$  to be equal to the constant  $b(n)$  in a small neighborhood of  $n$ , then one does not affect the left-hand side of (B.1) or increase the total variation of  $b$ , while only modifying the integral in (B.1) by an arbitrarily small amount. Hence, by the usual limiting argument, we may assume without loss of generality that  $b$  is locally constant at each such  $n$ . If we define the function  $g : (y, x] \rightarrow \mathbb{R}$  by

$$g(z) := \sum_{y < n \leq z} a_n - \int_z^y f(u) du - C$$

then  $g$  has jump discontinuities at the natural numbers, but is otherwise continuously differentiable, and is also bounded uniformly in magnitude by  $A$ . We can then compute the Riemann–Stieltjes integral

$$\int_{(y,x]} b dg = \sum_{y < n \leq x} b(n)a_n - \int_y^x f(t)b(t) dt.$$

Since the discontinuities of  $g$  and  $b$  do not coincide, we may integrate by parts to obtain

$$\int_{(y,x]} b dg = b(x)g(x) - b(y^+)g(y^+) - \int_{(y,x]} g db.$$

The left-hand side is  $O_{\leq}(A\|b\|_{\text{TV}^*(y,x]})$ , and the claim follows.  $\square$

We now prove (2.14). In fact we prove the sharper estimate

$$\sum_{y < p \leq x} b(p) \log p = \int_y^x b(t) \left(1 - \frac{2}{\sqrt{t}}\right) dt + O_{\leq}(\|b\|_{\text{TV}^*((y,x])} \tilde{E}(x)) \quad (\text{B.2})$$

where

$$\tilde{E}(x) := 0.95\sqrt{x} + \min(\max(\varepsilon_0, \varepsilon_1(x)), \varepsilon_2(x), \varepsilon_3(x))1_{x \geq 10^{19}} \quad (\text{B.3})$$

and

$$\begin{aligned} \varepsilon_0(x) &:= \frac{\sqrt{x}}{8\pi} \log x (\log x - 3) \\ \varepsilon_1(x) &:= 1.12494 \times 10^{-10} \\ \varepsilon_2(x) &:= 9.39(\log^{1.515} x) \exp(-0.8274\sqrt{\log x}) \\ \varepsilon_3(x) &:= 0.026(\log^{1.801} x) \exp(-0.1853(\log^{3/5} x)(\log \log x)^{-1/5}) \end{aligned}$$

From using the  $\varepsilon_2$  term, it is clear that

$$\tilde{E}(x) \ll x \exp(-c\sqrt{\log x})$$

for some absolute constant  $c > 0$ ; and by using the  $\varepsilon_0, \varepsilon_1$  term and routine calculations one can show that

$$\tilde{E}(x) \leq E(x)$$

for all  $x \geq 1423$ .

Observe that  $\tilde{E}$  is monotone non-decreasing. Thus by Lemma B.1, to show (B.2) will suffice to show that

$$\sum_{p \leq x} \log p = x - \sqrt{x} + O_{\leq}(\tilde{E}(x)) = \int_0^x \left(1 - \frac{2}{\sqrt{t}}\right) dt + O_{\leq}(\tilde{E}(x))$$

for all  $x \geq 1423$ .

For  $1423 \leq x \leq 10^{19}$ , this claim follows from [5, Theorem 2]. For  $x > 10^{19}$ , we apply [4, (6.10), (6.11)] to conclude that

$$\sum_{p \leq x} \log p = \psi(x) - \psi(\sqrt{x}) + O_{\leq}(1.03883(x^{1/3} + x^{1/5} + 2(\log x)x^{1/13})),$$

where  $\psi(x) := \sum_{n \leq x} \Lambda(n)$  is the usual von Mangoldt summatory function. From [15, Theorems 10, 12] we have

$$\psi(\sqrt{x}) = \sqrt{x} + O_{\leq}(0.18\sqrt{x}).$$

Since

$$0.18\sqrt{x} + 1.03883(x^{1/3} + x^{1/5} + 2(\log x)x^{1/13}) \leq 0.95\sqrt{x}$$

in this range of  $x$ , it suffices to show that

$$\psi(x) = x + O_{\leq}(\min(\max(\varepsilon_0(x), \varepsilon_1(x)), \varepsilon_2(x), \varepsilon_3(x)))$$

for  $x > 10^{19}$ . The claims for  $i = 2, 3$  follow from [12, Theorems 1.1, 1.4]. In [4, Theorem 2, (7.3)], the bound

$$\psi(x) = x + O_{\leq}(\varepsilon_0(x))$$

is established whenever  $x \geq 5000$  and  $4.92 \frac{x}{\sqrt{\log x}} \leq T$ , where  $T$  is a height up to which the Riemann hypothesis has been established. Using the value  $T = 3 \times 10^{12}$  from [13], we can therefore cover the range  $10^{19} < x < e^{55}$  (in fact we could go up to  $e^{58.33} \approx 2.1 \times 10^{25}$ ). For  $x \geq e^{55}$ , we can use [4, Table 2] (the value  $T = 2.445 \times 10^{12}$  used there following from [13]).

**Remark B.2.** Assuming the Riemann hypothesis, the  $\varepsilon_1, \varepsilon_2, \varepsilon_3$  terms in the definition of  $\tilde{E}(x)$  may be deleted, since [4, (7.3)] then holds for all  $x \geq 5000$ .

The claim (2.16) now follows from (2.14) by setting  $b(t) := \frac{1}{\log t}$ .

### APPENDIX C. COMPUTATION OF $c_0$ AND RELATED QUANTITIES

In this appendix we give some details regarding the numerical estimation of the constants  $c_0, c'_1, c''_1, c_1$  defined in (1.6), (5.2), (5.3), (5.4).

We begin with  $c_0$ . As one might imagine from an inspection of Figure 3, direct application of numerical quadrature converges quite slowly due to the oscillatory singularity. To resolve the singularity, we can perform a change of variables  $x = 1/y$  to express  $c_0$  as an improper integral:

$$c_0 = \frac{1}{e} \int_1^\infty \lfloor y \rfloor \log \frac{\lceil y/e \rceil}{y/e} \frac{dy}{y^2}. \quad (\text{C.1})$$

Next, observe<sup>11</sup> that

$$\begin{aligned} \frac{1}{e} \int_e^\infty y \log \frac{\lceil y/e \rceil}{y/e} \frac{dy}{y^2} &= \sum_{k=1}^\infty \int_{ke}^{(k+1)e} y \log \frac{k+1}{y/e} \frac{dy}{y^2} \\ &= \frac{1}{e} \sum_{k=1}^\infty \int_k^{k+1} (\log(k+1) - \log y) \frac{dy}{y} \\ &= \frac{1}{2e} \sum_{k=1}^\infty \log^2 \left( 1 + \frac{1}{k} \right) \\ &= 0.1797439053 \dots; \end{aligned}$$

The value here was computed in interval arithmetic by subtracting off the asymptotically similar sum  $\frac{1}{2e} \sum_{k=1}^\infty \frac{1}{k^2} = \frac{1}{2e} \frac{\pi^2}{6}$ , summing the resulting partial sum up to  $k = 10^5$ , bounding the tail of the sum rigorously.

$$\frac{1}{e} \int_1^e \lfloor y \rfloor \log \frac{e}{y} \frac{dy}{y^2} = \frac{2}{e^2} - \frac{\log 2}{2e} = 0.143173268 \dots$$

and hence

$$c_0 = \frac{1}{2e} \sum_{k=1}^\infty \log^2 \left( 1 + \frac{1}{k} \right) + \frac{2}{e^2} - \frac{\log 2}{2e} - \frac{1}{e} \int_e^\infty \{y\} \log \frac{\lceil y/e \rceil}{y/e} \frac{dy}{y^2}$$

where  $\{x\} := x - \lfloor x \rfloor$ . The integrand here lies between 0 and  $1/y^3$ , so the integral for  $y \geq T$  lies between 0 and  $1/2T^2$ . Truncating to say  $T = 10^5$  and performing the integral exactly, one can evaluate

$$\frac{1}{e} \int_e^\infty \{y\} \log \frac{\lceil y/e \rceil}{y/e} \frac{dy}{y^2} = 0.018498162 \dots$$

so that

$$c_0 = 0.30441901 \dots$$

<sup>11</sup>We thank an anonymous commenter on the blog of one of the authors for this suggestion.

A similar calculation (which we omit) reveals that

$$\begin{aligned} c'_1 &= \sum_{k=1}^{\infty} \frac{1 + \log(k+1)}{2e} \log^2 \left(1 + \frac{1}{k}\right) - \frac{1}{3e} \log^3 \left(1 + \frac{1}{k}\right) \\ &\quad + \frac{6}{e^2} - \frac{\log^2 2 + \log 2 + 3}{2e} \\ &\quad - \frac{1}{e} \int_e^{\infty} \{y\} (\log y) \log \frac{[y/e]}{y/e} \frac{dy}{y^2} \\ &\approx 0.3702051 \dots \end{aligned}$$

Computing the sum  $c''_1$  to reasonable accuracy requires some further analysis. From the crude bound

$$0 \leq \frac{1}{k} \log \left( \frac{e}{k} \left\lceil \frac{k}{e} \right\rceil \right) \leq \frac{e}{k^2}$$

and the integral test, one has the simple tail bound

$$0 \leq \sum_{k=K+1}^{\infty} \frac{1}{k} \log \left( \frac{e}{k} \left\lceil \frac{k}{e} \right\rceil \right) \leq \frac{e}{K}$$

but the convergence rate here is slow. To accelerate the convergence, we write  $\lceil \frac{k}{e} \rceil = \frac{k}{e} + \{-\frac{k}{e}\}$  and use the more precise Taylor approximation

$$\frac{e \left\{ -\frac{k}{e} \right\}}{k^2} - \frac{e^2 \left\{ -\frac{k}{e} \right\}^2}{2k^3} \leq \frac{1}{k} \log \left( \frac{e}{k} \left\lceil \frac{k}{e} \right\rceil \right) \leq \frac{e \left\{ -\frac{k}{e} \right\}}{k^2}.$$

Bounding  $\{-k/e\}$  by one, we have the tail bound

$$0 \leq \sum_{k=K+1}^{\infty} \frac{e^2 \left\{ -\frac{k}{e} \right\}^2}{2k^3} \leq \frac{e^2}{4K^2}$$

so the main task is then to control the simplified tail

$$\sum_{k=K+1}^{\infty} \frac{e \left\{ -\frac{k}{e} \right\}}{k^2}.$$

From the integral test one has

$$\frac{e}{2(K+1)} \leq \sum_{k=K+1}^{\infty} \frac{\frac{e}{2}}{k^2} \leq \frac{e}{2K}$$

so one can instead look at the normalized tail

$$\sum_{k=K+1}^{\infty} e \frac{\left\{ -\frac{k}{e} \right\} - \frac{1}{2}}{k^2}.$$

The Erdős–Turán inequality states that, for any absolutely convergent non-negative weights  $c_k$ , any interval  $I \subset [0, 1]$  of length  $|I|$ , and any real numbers  $\xi_k$ , and any  $N \geq 1$ , one has

$$\left| \sum_k c_k (1_I(\xi_k \bmod 1) - |I|) \right| \leq \frac{1}{N+1} \sum_k c_k + \sum_{n=1}^N \left( \frac{2}{\pi n} + \frac{2}{N+1} \right) \left| \sum_k c_k e^{2\pi i n \xi_k} \right|;$$

see the inequality<sup>12</sup> after [18, Theorem 20]. Applying this for  $I = [0, h]$  and then averaging in  $h$  from 0 to 1, we conclude that

$$\left| \sum_k c_k \left( \{\xi_k\} - \frac{1}{2} \right) \right| \leq \frac{1}{N+1} \sum_k c_k + \sum_{n=1}^N \left( \frac{2}{\pi n} + \frac{2}{N+1} \right) \left| \sum_k c_k e^{2\pi i n \xi_k} \right|.$$

In particular, we have

$$\left| \sum_{k=K+1}^{\infty} e^{\left\{ -\frac{k}{e} \right\} - \frac{1}{2}} \right| \leq \frac{1}{N+1} \sum_{k=K+1}^{\infty} \frac{e}{k^2} + \sum_{n=1}^N \left( \frac{2e}{\pi n} + \frac{2e}{N+1} \right) \left| \sum_{k=K+1}^{\infty} \frac{e^{-2\pi i n k/e}}{k^2} \right|.$$

To estimate the exponential sum

$$S_{n,K} := \sum_{k=K+1}^{\infty} \frac{e^{-2\pi i n k/e}}{k^2}$$

observe from shifting  $k$  by one that

$$S_{n,K} = e^{-2\pi i n/e} \sum_{k=K}^{\infty} \frac{e^{-2\pi i n k/e}}{(k+1)^2} = e^{-2\pi i n/e} S_{n,K} + O_{\leq} \left( \frac{1}{(K+1)^2} + \sum_{k=K+1}^2 \frac{1}{k^2} - \frac{1}{(k+1)^2} \right)$$

and hence on summing the telescoping series

$$|S_{n,K}| \leq \frac{2}{|e^{-2\pi i n/e} - 1|(K+1)^2} = \frac{1}{(K+1)^2 \sin(\pi n/e)}.$$

Because the irrationality measure of  $e$  is 2, this will give error terms of the shape  $O(\log K/K^2)$  if one sets  $N \approx K/\log K$ . Setting for instance  $K = 10^6$ ,  $N = 10^5$ , an interval arithmetic computation then gives

$$c_1'' = 1.679578996 \dots$$

and thus by (5.4)

$$c_1 = 0.7554808 \dots$$

## REFERENCES

- [1] K. Alladi, C. Grinstead, *On the decomposition of  $n!$  into prime powers*, J. Number Theory **9** (1977) 452–458.
- [2] S. F. Assmann, D. S. Johnson, D. J. Kleitman, J. Y.-T. Leung, *On a dual version of the one-dimensional bin packing problem*, J. Algorithms **5** (1984) 502–525.
- [3] A. Baker, G. Wüstholz, *Logarithmic forms and Diophantine geometry*, New Math. Monogr., 9 Cambridge University Press, Cambridge, 2007.
- [4] J. Büthe, *Estimating  $\pi(x)$  and related functions under partial RH assumptions*, Math. Comp., 85(301), 2483–2498, Jan. 2016.
- [5] J. Büthe, *An analytic method for bounding  $\psi(x)$* . Math. Comp., **87** (312), 1991–2009.
- [6] P. Dusart, *Explicit estimates of some functions over primes*, Ramanujan J. **45** (2018) 227–251.
- [7] P. Erdős, *Some problems in number theory*, in Computers in Number Theory, Academic Press, London New York, 1971, pp. 405–414.
- [8] P. Erdős, *Some problems I presented or planned to present in my short talk*, Analytic number theory, Vol. 1 (Allerton Park, IL, 1995) (1996), 333–335.

<sup>12</sup>In the cited reference, only the special case in which  $c_k$  is a uniform probability distribution function on  $\{1, \dots, M\}$  is discussed, but it is easy to see that the argument in fact works for arbitrary absolutely convergent non-negative weights  $c_k$ .

- [9] P. Erdős, R. Graham, *Old and new problems and results in combinatorial number theory*, Monographies de L'Enseignement Mathématique 1980.
- [10] R. K. Guy, *Unsolved Problems in Number Theory*, 3rd Edition, Springer, 2004.
- [11] R. K. Guy, J. L. Selfridge, *Factoring factorial  $n$* , Amer. Math. Monthly **105** (1998) 766–767.
- [12] D. Johnston, A. Yang, *Some explicit estimates for the error term in the prime number theorem*, J. Math. Anal. Appl., **527** (2) (2023), Paper No. 127460.
- [13] D. Platt, T. Trudgian, *The Riemann hypothesis is true up to  $3 \cdot 10^{12}$* , Bull. Lond. Math. Soc. **53** (2021), no. 3, 792–797.
- [14] H. Robbins, *A Remark on Stirling's Formula*, Amer. Math. Monthly **62** (1955) 26–29.
- [15] J. Rosser, L. Schoenfeld, *Approximate formulas for some functions of prime numbers*, Illinois J. Math. **6** (1962), 64–94.
- [16] T. Tao, *Decomposing factorials into bounded factors*, preprint, 2025. <https://arxiv.org/abs/2503.20170v2>
- [17] T. Tao, *Verifying the Guy–Selfridge conjecture*, Github repository, 2025. <https://github.com/teorth/erdos-guy-selfridge>.
- [18] J. D. Vaaler, *Some extremal functions in Fourier analysis*, Bulletin (New Series) of the American Mathematical Society, Bull. Amer. Math. Soc. (N.S.) **12**(2), 183–216, (April 1985).

UNAFFILIATED, ATHENS, GA 30605

*Email address:* boris.alexeev@gmail.com

UVA DEPARTMENT OF MATHEMATICS, CHARLOTTESVILLE, VA 22903

*Email address:* auj4kq@virginia.edu

MIT DEPARTMENT OF MATHEMATICS, CAMBRIDGE, MA 02139.

*Email address:* drew@math.mit.edu

UCLA DEPARTMENT OF MATHEMATICS, LOS ANGELES, CA 90095-1555.

*Email address:* tao@math.ucla.edu

???

*Email address:* ???

GOOGLE, MOUNTAIN VIEW, CA

*Email address:* kevinventullo@google.com