

PHOBIC: Perfect Hashing with Optimized Bucket Sizes and Interleaved Coding

Stefan Hermann  

Karlsruhe Institute of Technology, Germany

Hans-Peter Lehmann  

Karlsruhe Institute of Technology, Germany

Giulio Ermanno Pibiri  

Ca' Foscari University of Venice, Italy
ISTI-CNR, Italy

Peter Sanders  

Karlsruhe Institute of Technology, Germany

Stefan Walzer  

Karlsruhe Institute of Technology, Germany

Abstract

A minimal perfect hash function (or MPHf) maps a set of n keys to $[n] := \{1, \dots, n\}$ without collisions. Such functions find widespread application e.g. in bioinformatics and databases. In this paper we revisit PTHash – a construction technique particularly designed for fast queries. PTHash distributes the input keys into small buckets and, for each bucket, it searches for a hash function seed that places its keys in the output domain without collisions. The collection of all seeds is then stored in a compressed manner. Since the first buckets are easier to place, buckets are considered in non-increasing order of size. Additionally, PTHash heuristically produces an imbalanced distribution of bucket sizes by distributing 60% of the keys into 30% of the buckets.

Our main contribution is to characterize, up to lower order terms, an *optimal* distribution of expected bucket sizes. We arrive at a simple, closed form solution which improves construction throughput for space efficient configurations in practice. Our second contribution is a novel encoding scheme for the seeds. We split the keys into partitions. Within each partition, we run the bucket distribution and search step. We then store the seeds in an *interleaved* manner by consecutively placing the seeds for the i -th buckets from all partitions. The seeds for the i -th bucket of each partition follow the same statistical distribution. This allows us to tune a compressor for each bucket. Hence, we call our technique PHOBIC – Perfect Hashing with Optimized Bucket sizes and Interleaved Coding.

Compared to PTHash, PHOBIC is 0.17 bits/key more space efficient for same query time and construction throughput. We also contribute a GPU implementation to further accelerate MPHf construction. For a configuration with fast queries, PHOBIC-GPU can construct a perfect hash function at 2.17 bits/key in 28 ns per key, which can be queried in 37 ns per query on the CPU.

2012 ACM Subject Classification Theory of computation → Data compression; Information systems → Point lookups

Keywords and phrases Compressed Data Structures, Minimal Perfect Hashing, GPU

Supplementary Material *Software (CPU implementation)*: <https://github.com/jermp/pt hash/tree/phobic>

Software (GPU implementation): <https://github.com/stefanfred/PHOBIC-GPU>

Software (Internal comparison): <https://github.com/stefanfred/PHOBIC-Scripts>

Software (Comparison with competitors): <https://github.com/ByteHamster/MPHF-Experiments>

Funding This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No. 882500). The project also received funding from the European Union's Horizon Europe research and innovation



programme (EFRA project, Grant Agreement Number 101093026). This work was also partially supported by DAIS – Ca’ Foscari University of Venice within the IRIDE program.

Acknowledgements This paper is based on and has text overlaps with the Master’s thesis [17] of Stefan Hermann. We refer readers to that thesis for a detailed evaluation and description of the GPU implementation.

1 Introduction

A *hash function* maps a set S of n keys to a range of integers $[m] := \{1, \dots, m\}$, regardless of whether multiple keys collide on the same output. A *perfect* hash function (PHF) on S is a mapping *without collisions*. This requires $m \geq n$. The function does not necessarily have to store the keys explicitly. It only has to store enough information to prevent collisions, which are more likely when m is close to n . In the extreme case of $m = n$, the mapping is called a *minimal* perfect hash function (MPHF). In this paper, we consider the minimal case only. PHFs find widespread practical application e.g. in compressed full-text indexes [4], computer networks [23], databases [9], prefix-search data structures [2], language models [30], bioinformatics [12, 27], and Bloom filters [8]. The three main performance attributes of an MPHF are low space consumption, fast construction, and fast queries. Concerning space, the lower bound is $\log_2(e) \approx 1.44$ bits/key [25]. Practically viable approaches can get within a few percent of the lower bound, but do so with some sacrifices in running time [19, 21]. This paper is concerned with techniques that are focused on achieving fast query times. For example, this is very important when using perfect hashing to implement a static hash table that is both space-efficient and allows fast search.

Perfect Hashing Through Bucket Placement. Perfect hashing through bucket placement takes the n keys and maps them to small buckets. For each bucket, it uses brute-force search to find a seed of a hash function such that all keys of the bucket do not collide with each other or previously placed keys. The first buckets are easier to place because the output domain is less full. Therefore, the methods insert the buckets in order of non-increasing size. While CHD [3] uses buckets of constant expected size, FCH [14] and PTHash [28, 29] set aside 30% of “heavy” buckets that receive 60% of the keys in expectation, while 70% of “light” buckets receive only 40% of the keys in expectation. This imbalance in expected bucket sizes improves construction speed by further decreasing the size of the last, hardest to place, buckets. The resulting list of seed values are stored with various compression techniques, resulting in a variety of trade-offs between space consumption and query speed.

Partitioning. Any PHF construction algorithm can be trivially parallelized by splitting the input keys into disjoint subsets. We refer to those subsets as *partitions*. The various PHFs are then logically “concatenated” into a single PHF taking the prefix-sum of the partition sizes. The respective offsets have to be looked up when querying a key, imposing some query time overhead. Each partition can be constructed independently in parallel. Partitioning is the usual approach for parallelization, and is applied to PTHash by PTHash-HEM [29].

Contribution. This paper aims at improving the space efficiency and construction speed of PTHash, while maintaining its fast query speed. There are three ingredients. Our main contribution (in Section 3) is to characterize, up to lower order terms, an optimal distribution of expected bucket sizes, effectively taking the imbalance-trick used in FCH and PTHash to its logical conclusion. The distribution is easy-to-implement and greatly improves construction time and space efficiency in practice. Our second contribution (in Section 4.1) is to improve the compression of seed values when using partitioning. Seeds are searched independently

for each partition, but compressed together. We exploit that the seeds of the i -th bucket of each partition follow the same statistical distribution. This allows for tuning a compressor for each such index i . We store the seeds in an *interleaved* manner by consecutively placing the seeds for the i -th buckets from all partitions. Finally, we contribute (in Section 4.2) an implementation for *Graphics Processing Units* (GPUs) to speed up construction.

2 Related Work

Perfect hashing is an active area of research. We provide an overview of state-of-the-art approaches. For more details, refer to Section 2 of [29].

Fingerprinting. Perfect hashing through fingerprinting [10, 26] is a technique aimed at fast construction and queries at the cost of reduced space efficiency. The idea is to map the n keys to γn positions using a hash function, where γ is a tuning parameter. A bit vector of length γn indicates positions that received exactly one key. Keys that take place in collisions are handled recursively on another layer of the same data structure. A query operation descends through the recursive layers until it finds a 1-bit, meaning that it was the only key mapping to that position. A rank operation on the bit vector for that position then gives the final MPH value. FMPH [5] and BBHash [22] are publicly available implementations of the approach. FMPHGO [5] extends on this idea using a small number of brute-force re-tries to reduce the number of colliding keys. Fingerprinting based approaches are fast to construct but are outperformed in terms of space consumption and query time by PTHash.

Brute Force. RecSplit [13] first partitions the input into sets of equal expected size. It then recursively splits the key set of each partition until sets of small constant size (usually ≤ 16) are left. Within these sets, it finds a perfect hash function by brute force. RecSplit achieves space usage of about 1.56 bits/key. The resulting splitting tree has to be traversed during querying which implies considerably higher query costs compared to PTHash. The brute force search was later improved in SIMDRecSplit [6], which also parallelizes the construction on the GPU. To the best of our knowledge, RecSplit is the only other PHF construction technique that has a GPU implementation.

Perfect Hashing Through Retrieval. In perfect hashing through retrieval, every key has a number of candidate positions, determined by different hash functions. A retrieval data structure then stores which of the choices should be used for each key. Note that this implies some query overhead compared to PTHash. Early implementations include BPZ [7] and GOV [15]. SicHash [20] reduces space consumption using a mix of different retrieval data structures and some retries. ShockHash-RS [19, 21] combines 1-bit retrieval with the brute-force approach of RecSplit and currently is the most space-efficient approach to MPHFs with as little as 1.49 bits/key [19].

3 Optimizing Bucket Sizes

Consider perfect hashing through bucket placement with n keys, for $m = n$ and B buckets, i.e. an average bucket size of $\lambda = n/B$. Previous literature overlooked the simple insight that large λ already guarantees a space consumption close to the lower bound of $\log_2 e$ bits per key, *without any assumptions* on the bucket sizes or their distribution.

► **Proposition 1.** *Any specialization of perfect hashing through bucket placement requires between $\log_2 e$ bits per key and $\log_2 e + \mathcal{O}(\frac{\log \lambda}{\lambda})$ bits per key in expectation.*

Our goal in this section need therefore only be to minimize construction time. Here we are faced with a lower bound for our family of approaches.

► **Proposition 2.** *Any specialization of perfect hashing through bucket placement has an expected construction time of $\Omega(e^\lambda/\lambda)$ per bucket.*

Propositions 1 and 2 are restated more formally as Proposition 18 and proved in Appendix A.4. It is intuitively clear (and proved in Proposition 16 in Appendix A.3) that buckets should be processed in order from largest to smallest. The only remaining degree of freedom is to choose the expected sizes of the buckets. We characterize asymptotically optimal ways of doing so, formalized by what we call bucket assignment functions and achieving a construction time of $e^{\lambda(1+\varepsilon)}$ per bucket. Proofs are found in the appendix.

3.1 Bucket Assignment Functions

Let w_1, \dots, w_B be the probability that a key hashes to bucket i for $i \in [B]$. We may assume without loss of generality that these probabilities are given in decreasing order. An equivalent view considers the prefix sums $\sigma_i := w_1 + \dots + w_i$. A key with (normalized) hash value $x \in (0, 1]$ is then assigned to bucket i if $x \in (\sigma_{i-1}, \sigma_i]$.

We can conveniently represent this information using a *bucket assignment function* $\gamma : [0, 1] \rightarrow [0, 1]$ that: interpolates the points $\{(\sigma_i, i/B) \mid 0 \leq i \leq B\}$, is increasing and smooth on $(0, 1)$, and has non-decreasing derivative. The bucket assigned to hash value $x \in (0, 1]$ is then $\lceil \gamma(x) \cdot B \rceil$. It is a non-trivial insight of this section that a single bucket assignment function (not depending on B and n) can result in good construction times for many values of B and n simultaneously.

From now on, let $\lambda := n/B$. We summarize some useful intuitions about bucket assignment functions. These intuitions are valid for large n and B (when γ , γ^{-1} , and γ' are approximately constant on intervals of length $\frac{1}{n}$ and $\frac{1}{B}$). For now, we neglect edge cases related to γ or γ^{-1} not being smooth at 0 or not being smooth at 1 (but just smooth on $(0, 1)$).

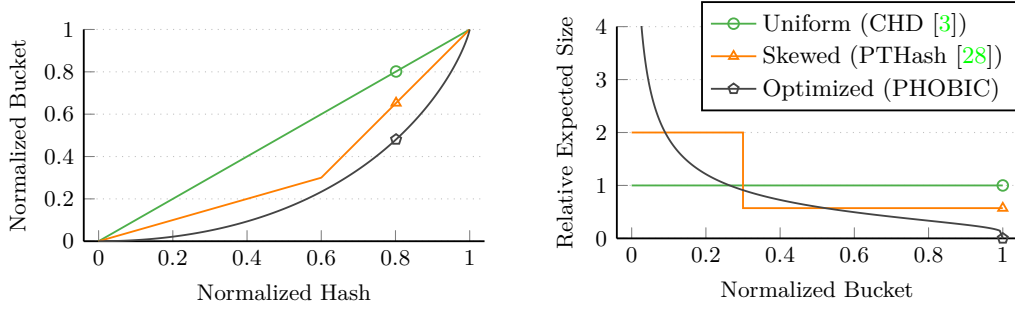
► **Intuition 3.** Let $x \in (0, 1]$ be a normalized hash and $b = \gamma(x)$ the normalized bucket index of the bucket assigned to x . Then

- (i) The expected size of the bucket assigned to x is $\lambda/\gamma'(x)$.
(Reason: In the vicinity of x and for infinitesimal δ , a δ -fraction of the hash range (used by δn keys in expectation) is shared by a $(\gamma'(x) \cdot \delta)$ -fraction of the B buckets. The quotient is $\delta n / (\gamma'(x) \delta B) = \lambda/\gamma'(x)$.)
- (ii) The expected size of the bucket with normalized index b is $\lambda/\gamma'(\gamma^{-1}(b)) = \lambda(\gamma^{-1})'(b)$.
(Follows from (i) and the inverse function rule.)
- (iii) The expected size of a bucket is decreasing in its normalized index.
(Follows from (ii) and monotonicity of γ' and γ^{-1} .)
- (iv) The expected fraction of keys with normalized hash in $(0, x]$ is x .
- (v) If $\mu > 0$ and $x_\mu \in (0, 1)$ is such that $\lambda/\gamma'(x_\mu) = \mu$ then the expected fraction of keys in buckets of size at least μ is x_μ . (Follows from (i), (iii) and (iv).)

3.2 An Optimal Bucket Assignment Function

Intuitively, we identify the following bucket assignment function to be optimal, although our precise result stated below is more subtle.

$$\beta_*(x) = x + (1-x) \cdot \ln(1-x) \text{ with derivative } \beta'_*(x) = -\ln(1-x).$$



(a) The bucket assignment functions map a normalized hash value x to a normalized bucket index $\gamma(x)$. (b) The expected bucket sizes relative to the average size λ are $(\gamma^{-1})'(b)$ for normalized bucket index b .

■ **Figure 1** Comparison of bucket assignment functions $\gamma(x)$ of related work and PHOBIC ($\gamma = \beta_*$).

For comparison, Figure 1a shows β_* as well as the bucket assignment functions used by CHD and PHash. In Figure 1b we see the distribution of expected bucket sizes, which is uniform for CHD, imbalanced for PHash, and even more imbalanced for β_* .

Recall that $\lambda = n/B$ is the average bucket size. By Proposition 2 a lower bound for the expected work is $\Omega(n \cdot e^\lambda/\lambda)$. We prove, firstly, that any bucket assignment function γ that differs from β_* leads to expected work exceeding $n \cdot e^{(1+\varepsilon)\lambda}$ for some $\varepsilon = \varepsilon(\gamma)$, provided that λ is large enough. Conversely, we show that a slight perturbation β_ε of β_* leads to an expected work of essentially at most $n \cdot e^{(1+\varepsilon)\lambda}$ for any $\varepsilon > 0$, provided that $\lambda \geq \lambda_0(\varepsilon)$ is large enough. Essentially, we can get arbitrarily close to a cost of e^λ per key and only functions close to β_* can achieve this.

Our results bound the work $w_{n,\lambda}(\gamma)$ associated with γ and involve a ‘‘coupon collector term’’ w_{coupon} , which is the work required to place buckets of size 1. We will define these more precisely below. Proofs are found in Appendices A.1 and A.2. We have reason to believe that our results generalize for the non-minimal case of $m > n$, as explained in Appendix B.

► **Theorem 4.** *Let $\gamma : [0, 1] \rightarrow [0, 1]$ be a continuous bucket assignment function that is smooth on $(0, 1)$ with non-decreasing derivative. If $\beta_* \neq \gamma$ then*

$$\exists \varepsilon > 0 : \forall \lambda \geq \lambda_0(\varepsilon) : \forall n \geq n_0(\lambda, \varepsilon) : w_{n,\lambda}(\gamma) \geq n \cdot e^{\lambda(1+\varepsilon)} + w_{\text{coupon}} \text{ whp.}$$

While this leaves the relationship between γ and $\varepsilon(\gamma)$ open, our proof suggests that any $\varepsilon < \sup_{x \in (0,1)} \frac{\beta'_*(x)}{\gamma'(x)} - 1$ is a possible choice.

► **Theorem 5.** *Let $\beta_\varepsilon(x) := \varepsilon x + (1 - \varepsilon)\beta_*(x)$ for some $\varepsilon > 0$. Then*

$$\forall \varepsilon > 0 : \forall \lambda \geq \lambda_0(\varepsilon) : \forall n \geq n_0(\lambda, \varepsilon) : w_{n,\lambda}(\beta_\varepsilon) \leq n \cdot e^{\lambda(1+\mathcal{O}(\varepsilon))} + w_{\text{coupon}} \text{ whp.}$$

By *with high probability* (whp) we mean probability $1 - \mathcal{O}(n^{-c})$ for some $c > 0$. Note that both theorems are phrased such that we may assume that n is much larger than λ and λ is much larger than $1/\varepsilon$. We give implementation details concerning the use of β_ε in Appendix C.1.

What we do not prove. Note that our analysis leaves undecided whether β_* is itself a ‘‘good’’ bucket assignment function, i.e. whether $w_{n,\lambda}(\beta_*)$ approaches e^λ in any meaningful sense. We suspect that it does. However, the perturbation simplifies the analysis and improves running times in practice. Our analysis also does not imply that the particular perturbation we choose is the best choice: there may be an alternative to β_ε such that the overall work approaches e^λ more quickly.

The work associated with a bucket assignment function. To place a bucket of size $s \in \mathbb{N}$ into a hash table of size n that already has load factor $\alpha \in [0, 1 - \frac{s}{n}]$ we repeatedly try seeds for a hash function mapping keys to $[n]$, until all keys hash to free positions. The expected cost $c_n(s, \alpha)$ associated with this task under the simple uniform hashing assumption is described precisely in Appendix A.3. We have to take into account self-collisions, i.e. while checking the keys one after the other, the load factor gradually increases and is $\alpha' = \alpha + \frac{s-1}{n}$ for the last key. For our purposes, the following bounds on $c_n(s, \alpha)$ suffice

$$(1 - \alpha)^{-s} \leq c_n(s, \alpha) \leq s \cdot (1 - \alpha')^{-s}. \quad (1)$$

This uses that $(1 - \alpha)^{-s}$ and $(1 - \alpha')^{-s}$ are lower and upper bounds on the number of seeds that have to be tried and that, to test a seed, at least 1 and at most s keys have to be considered. Now assume we are given a bucket assignment function γ as well as $n \in \mathbb{N}$, $\lambda \in \mathbb{R}_+$ and $B = n/\lambda$. By assigning keys to buckets according to γ and hash values in $(0, 1]$ we obtain buckets. Let $s_1 \geq \dots \geq s_B$ be their sizes in decreasing order. By Proposition 16 in Appendix A.3 it is advantageous to process the buckets in this order. Defining $\alpha_i := \frac{1}{n} \sum_{j=1}^{i-1} s_j$, the total cost is then $w_{n,\lambda}(\gamma) = \sum_{i=1}^B c_n(s_i, \alpha_i)$.

Note that while this describes the *expected* cost when *given* $(s_i)_{i \in [B]}$, overall $w_{n,\lambda}(\gamma)$ is still a random variable because the numbers $(s_i)_{i \in [B]}$ are random. Assume now that there are exactly k buckets of size 1 that are placed last (we may ignore buckets of size 0). For these buckets, the upper and lower bounds in Equation (1) coincide so they incur a cost of

$$w_{\text{coupon}} := \sum_{i=1}^k c(1, \frac{n-i}{n}) = \sum_{i=1}^k \frac{n}{i} = n \cdot H_k.$$

Here H_k is the k th harmonic number, which satisfies $H_k = \Theta(\log k)$. If n is sufficiently large compared to λ then we have $k \geq n^d$ whp for some constant $d > 0$, giving a cost of $\Theta(n \cdot H_{n^d}) = \Theta(n \log n)$. This dominates overall construction time if n is sufficiently large compared to λ . Our theorems list this work for buckets of size 1 as a separate term because there are techniques to mitigate the problem: The hash function may permit to directly compute for a given key x and table position i a seed for which x is mapped to i . This is the case if the seed includes an additive displacement term, as is the case in our implementation and in FCH [14].

Intuition: What makes β_* uniquely promising. For $\mu > 0$ let $x_\mu \in (0, 1)$ be such that $\lambda/\beta'_*(x_\mu) = \mu$. By Intuition 3 (v), roughly an expected x_μ -fraction of the keys (those with hashes in $[0, x_\mu]$) land in buckets of expected size at least μ . Assume for now that a bucket of expected size μ has actual size μ (ignoring the issue that μ may not be integer). Then, since we process buckets in order of increasing size, we would process a bucket of size μ when the load factor is x_μ . The expected cost for this is, by Equation (1), around $(1 - x_\mu)^{-\mu}$. Using that $\beta'_*(x) = -\log(1-x)$ gives $\mu = -\lambda/\log(1-x)$ and hence $(1 - x_\mu)^{-\mu} = (1 - x_\mu)^{\lambda/\log(1-x_\mu)} = e^\lambda$, i.e. a cost independent of μ . The idea behind Theorem 4 is that any bucket assignment function $\gamma \neq \beta_*$ fails to balance bucket sizes in this way, leading to significantly higher costs.

The simplification we made seems innocent for large μ since a bucket of large expected size typically has actual size close to its expectation. But consider $\mu = 1.5$ and cease to ignore rounding issues. If at load factor x_μ we would process buckets of size 1 half the time and buckets of size 2 half the time, the resulting costs are around $e^{\frac{2}{3}A}$ and $e^{\frac{4}{3}A}$, respectively, which does not average out to e^λ . Luckily, things are more complicated. It turns out that for small $s \in \mathbb{N}$ and assuming large λ , the *expected number* of buckets of size s is meaningfully greater than the number of buckets of *expected size* in the range $[s - 1, s + 1]$. This means

that we get more small buckets than we seem to have called for, decreasing costs at high load factors. It seems clear that the flipside of this beneficial effect must be a detrimental effect for larger bucket sizes that a proof of Theorem 5 must quantify. When it comes to *very* large bucket sizes, we bail ourselves out by using β_ε instead of β_* : since β'_ε is lower bounded by ε , the expected bucket sizes are capped at λ/ε . It is buckets of intermediate sizes that have to pay the price.

4 Fine-Grained Partitioning

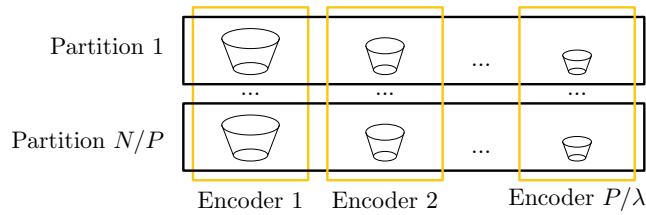
Any PHF construction can trivially be parallelized by hashing the keys into subsets of expected equal size and building a PHF for each subset in parallel. We refer to those subsets as *partitions*. The various PHFs are then logically “concatenated” into a single PHF taking the prefix sum of the partition sizes. The respective offsets have to be looked up when querying a key, imposing some query time overhead. Partitioning is widely used for a variety of construction techniques. It was also used by PTHash in the PTHash-HEM variant [29]. In this paper, we use partitions that are several magnitudes smaller than the ones used in PTHash-HEM. In itself, reducing the partition size results in only marginal construction time improvements. However, small partitions enable a new, more efficient encoding scheme which we introduce in Section 4.1. Additionally, they enable a fast GPU parallelization which we describe in Section 4.2.

Encoding the offsets of the partitions. The offsets and sizes of the many partitions require a non-negligible amount of space. We mitigate this without incurring too much query overhead by storing sizes only implicitly as the difference of two subsequent offsets. Offsets are stored as the difference to their expectation.

Hash Function. The original PTHash hash function works by XOR-ing the hash value of the key with the hash value of the seed. This reduces hash function evaluation to a simple XOR operation, which improves performance in practice. Although the technique works well on large-enough partition sizes, it might fail for small sizes because of correlations in the hash values. Given the small partition sizes we use here, we have to rely on a different technique. We use the seed value p to store two numbers, namely $p = s \cdot m + d$, where m is the actual partition size and $d \in [m]$ is an additive displacement. The position of a key x is $(h(x, s) + d) \bmod m$. While searching for a seed, we calculate $h(x, s)$ for all keys of the bucket and then only have to increment the values to obtain the positions for the next seeds. If no position is found within $d \in [m]$ we continue by incrementing s , re-calculating $h(x, s)$ and setting $d = 0$. At query time, we have to calculate the position of a key x using seed p . Note that we have $d = p \bmod m$, so we can compute the position of the key x as $(h(x, \lfloor p/m \rfloor) + p) \bmod m$.

4.1 Interleaved Coding of Seeds

Once the search has finished, the seeds found for each bucket have to be stored in some compressed manner. Ideally, the seeds should be encoded such that they require little space and are quickly accessible during querying. We mainly use *Compact* and *Golomb-Rice* encoding as building blocks for our new technique. Compact encoding is also used in the original PTHash implementation. In Compact encoding, all values are stored consecutively by concatenating their binary representation. All values use the same bit length, allowing for quick access. The bit length is chosen such that the highest seed can be accommodated. Golomb-Rice [16, 31] encoding stores the b least significant bits of each seed using compact



■ **Figure 2** Interleaved coding. Encoder i stores the seed of bucket i from all partitions.

encoding. The most significant bits are stored in unary representation. A selection structure enables access to the unary part of the seeds in constant time. We apply the formula by Kiely [18] to select b . A straightforward approach would be to encode all seeds using a single encoder. However, the seeds do not follow the same statistical distributions across different buckets, hence using the same encoder for all buckets is suboptimal. It is instead beneficial to group seeds which follow the same distribution and encode them using the same encoder. PTHash does this only partially by using two encoders – one for each expected bucket size (the so-called “front-back” compression [28]).

We now introduce our new technique. For each partition we hash to the same number of buckets $B = P/\lambda$, based on the average partition size P and average bucket size λ . The i -th bucket of a partition has the same expected size and the corresponding seed follows the same statistical distribution as the i -th bucket of any other partition. Although the idea of our optimized bucket assignment function is to give all buckets the same seed distribution, this is not achievable in practice. At least one reason for this are the discrete bucket sizes. This results in discrete jumps in the probability that a seed is found when processing one bucket after another. In *interleaved coding* we therefore employ B encoders and the i -th encoder stores the seeds of the i -th buckets of all partitions. Each encoder can thus use tuning parameters for its specific distribution (e.g., different Golomb-Rice parameters). Figure 2 gives an illustration of interleaved coding.

It is also possible to mix different encoding techniques, similar to what PTHash does. Larger buckets are accessed more often than smaller buckets because they contain more keys. Hence, it is beneficial to use an encoding technique which is optimized for fast lookup time (e.g., Compact) for the larger buckets. Conversely, the encoding for the seeds belonging to smaller buckets should be tuned for space efficiency (e.g., Golomb-Rice). To conclude this section, we point out that each of the B encoders introduces some metadata overhead (e.g., for storing its parameters). Using rather small partition sizes P decreases the number B of encoders and therefore the constant overhead.

4.2 GPU Parallelization

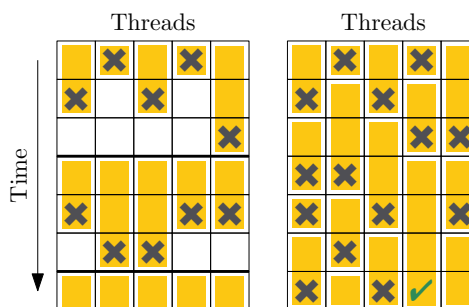
We provide a GPU implementation for even faster construction. On a GPU, each *workgroup* executes independently, typically with its own subset of data. Within each workgroup, individual *threads* execute concurrently. Threads within the same workgroup can share data and synchronize with each other through mechanisms like barriers and shared memory. Thread level parallelism allows for fine-grained parallel execution of instructions within a workgroup. However, only threads which follow the same control path and thus execute the same instruction at the same time can be executed in parallel. It is therefore crucial to avoid divergent control paths. As a first step, our parallel implementation transfers the keys to the GPU, before partitioning them. Afterwards, we sort the buckets and start the search.

■ **Algorithm 1** Seed search for one bucket.

```

shared sFound  $\leftarrow \infty$ 
shared sNext  $\leftarrow$  threadCount
seed  $\leftarrow$  threadId, keyIndex  $\leftarrow$  0
while sFound =  $\infty$  do
  isCollision  $\leftarrow$  COLL(seed, keyIndex)
  keyIndex  $\leftarrow$  keyIndex + 1
  if isCollision then
    keyIndex  $\leftarrow$  0
    seed  $\leftarrow$  ATOMADD(sNext, 1)
  else if keyIndex = bucketSize then
    sFound  $\leftarrow$  ATOMMIN(sFound, seed)

```



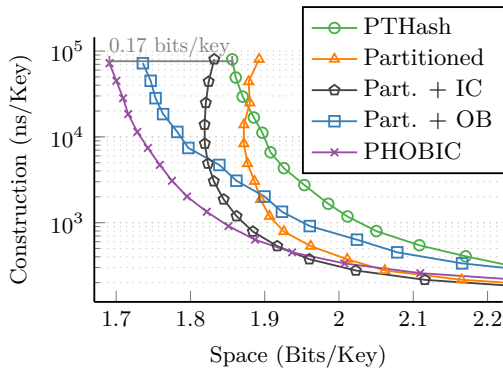
■ **Figure 3** Each box represents one seed tested by one thread. Left: Synchronized nested loop. Right: Algorithm 1 where we continue with the first key and the next seed after a collision.

Search. Our fine-grained partitioning naturally maps to the architecture of a GPU. Each partition is processed by one workgroup. The small partition sizes enable performing the search entirely using fast but small shared memory. Bucket by bucket, all threads of the workgroup cooperate to quickly find the smallest working seed. A CPU implementation would usually do so using a nested loop. The outer loop would iterate over seed values and the inner loop over keys. If a collision occurs, it would immediately leave the inner loop and continue with the first key and the next seed. However, on a GPU, leaving the inner loop would result in divergence because threads might encounter a collision after a different number of keys. This is illustrated in the left of Figure 3.

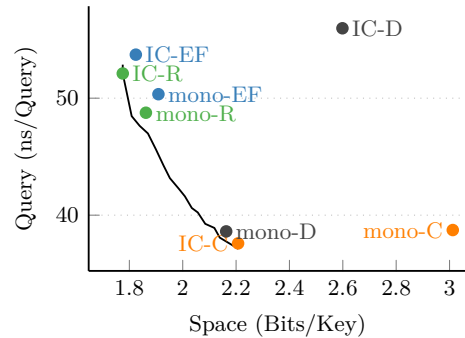
Instead, we use the technique described in [32] to emulate the behavior of the nested loop using a single loop to reduce divergence. Hence, our GPU implementation parallelizes over partitions, seeds and keys. The inner loop is emulated by incrementing the key index in each iteration. If a collision occurs, we reset the key index and emulate the behavior of the outer loop by atomically incrementing a seed counter which is shared among all threads. If the last key did not collide, we found a working seed. Multiple threads can find a seed in the same iteration. To reduce entropy, we use an atomic minimum to identify the smallest of those seeds. Note that this finds the smallest working seed overall because all threads finding a working seed must have processed each key. Therefore, if there was a smaller working seed, it would have been found in an earlier loop iteration. We give pseudocode in Algorithm 1 and illustrate the behavior on the right in Figure 3. During search, we only access shared memory and perform fast arithmetic operations. We remark that specific optimizations for our additive displacement hash function are not shown here, i.e. after calculating the initial positions we can apply new displacements using only additions.

5 Experiments

In Section 5.1, we gradually integrate our improvements to show the individual effects. Then we compare our GPU and CPU implementation with the state of the art in Section 5.2. We use a machine with an Intel Core i7-11700 CPU with 64 GiB of DDR4 RAM running Ubuntu 22.04.1. Each core has 48 KiB L1 and 512 KiB L2 data cache. As a GPU, we use an Nvidia RTX 3090 and use Vulkan 1.3.236 to interface with it. We compile using GCC 11.4.0 and compiler options `-march=native` and `-O3`. All benchmarks use random strings of random length between 10 and 50 characters as input which is adopted from previous



■ **Figure 4** Comparing original PTHash (EF, $\alpha = 0.99$) to fine-grained partitioning with additive displacement hash (mono-R). We then add interleaved coding (IC-R) and optimized bucketing (OB) individually. Putting it all together we arrive at PHOBIC. All single-threaded. There are small differences in query time.



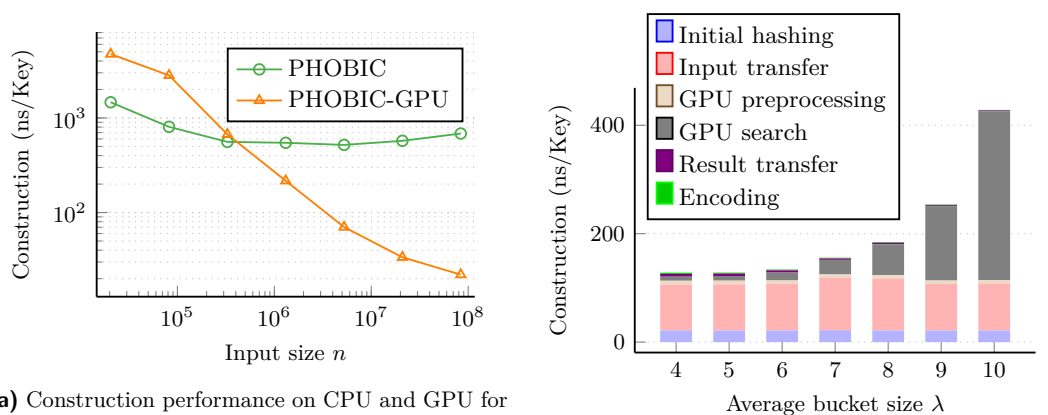
■ **Figure 5** Query time and space consumption of Elias-Fano, Dictionary, Compact, and Rice encoders. Points prefixed “mono” place all seeds into a single encoder and those prefixed “IC” use interleaved coding. The curve shows different mixtures (see Section 4) of Compact and Rice encoders in the interleaved coding.

work [6, 19, 20, 21]. Note that almost all competitors first generate master hash codes of the input. This makes the construction largely independent of the input distribution. We measure the query time by querying each key once in random order. All experiments use 100 million keys, $\lambda = 8$ and an average partition size of 2 500 if not stated otherwise. Our source code is public under the General Public License. You can find it through the links on the title page of this paper.

5.1 From PTHash to PHOBIC

We now gradually introduce our improvements to PTHash. As basic improvements, we replace the initial hash function with xxHash [1] and implement faster parallel partitioning. In all experiments, PTHash contains these changes as well to focus on our algorithmic improvements. Figure 4 gives measurements for the different improvement steps and shows them in different combinations.

Interleaved Coding. Partitioning of PTHash is already used in PTHash-HEM for parallelization [29]. PTHash-HEM uses partitions of size $\approx 10^6$. Smaller partitions only lead to minor improvements, as we show in Appendix D. However, smaller partition sizes shine when used with our newly introduced interleaved coding (Section 4.1). Interleaved coding uses P/λ encoders, where P is the expected partition size. Reducing the partition size can significantly reduce constant space overheads, as we also show in Appendix D. Figure 4 compares the technique of placing all seeds into a single Rice encoder (orange curve) to placing the keys using interleaved Rice coding (black). Interleaved coding consistently improves space efficiency by 0.06 bits/key. Figure 5 compares different combinations of encoders, which were partially used in the original implementation. Interleaved coding allows for mixing of different encoding techniques. If we use this to encode the keys using different numbers of Compact and Rice encoders, we can cover the entire query time to space trade-off in our configuration. Note that the construction performance is similar for all approaches because the encoding is fast compared to the remaining construction.



(a) Construction performance on CPU and GPU for $\lambda = 8.0$ and 8 threads, comparing different n . The GPU requires large n to fully utilize its computing resources.

(b) Different construction steps by values of average bucket size λ .

■ **Figure 6** GPU performance for different input sizes (left) and λ (right).

Optimized Bucket Assignment. In Figure 4 we also show how using the optimized bucket function affects construction speed and space. Our optimization of the bucket assignment function is particularly helpful to construct very space efficient configurations. When compared for same construction time, the optimized function is up to 0.14 bits/key more space efficient relative to the original PTHash bucket assignment.

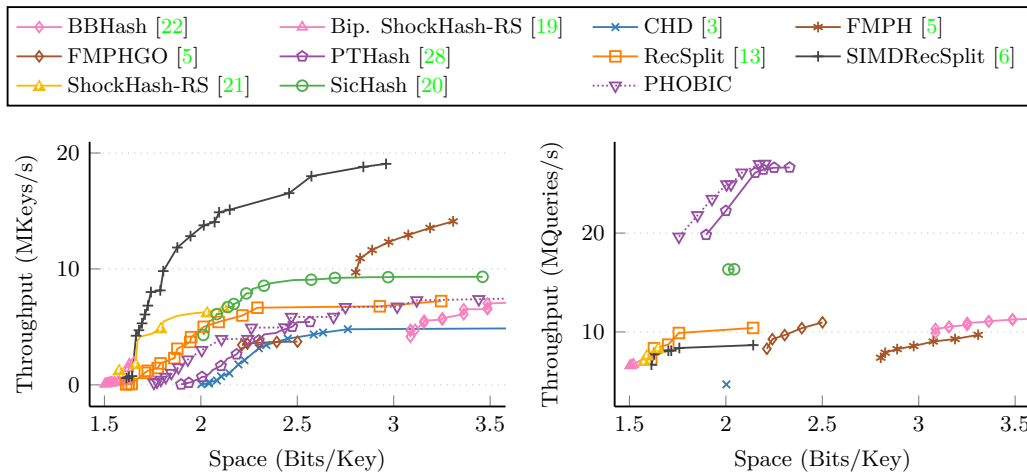
Further Remarks. With interleaved coding, another improvement originates from the secondary bucket ordering. Primarily the buckets are sorted in non-increasing size. Secondly sorting in increasing expected size reduces the space consumption by 0.044 bits/key compared to decreasing expected size. The reason for this behavior remains an open problem.

Original PTHash observed significant performance improvements by first calculating a non-minimal PHF and repairing the “gaps” afterwards. Refer to [28] for details. This trick does not result in an improvement when using PHOBIC.

GPU Parallelization. Our final contribution is a GPU implementation to speed up construction. Our implementation parallelizes over partitions, seeds and keys. The GPU implementation is mainly useful for large average bucket sizes λ . This is well illustrated in Figure 6b: For smaller values of λ , the construction time is dominated by the time to transfer the input data to the GPU. We also compare CPU and GPU construction speed for different input sizes in Figure 6a. The GPU requires a large number of input keys and thus a large number of partitions before its computing resources are fully utilized. Overall, the GPU outperforms the CPU for a sufficiently large λ and n . We use the GPU only to accelerate construction, while measuring all queries on the CPU.

5.2 Comparison to Other Methods

We compare our new approach to several other methods from the literature. First and foremost, we compare against the original PTHash [28, 29] implementation. The comparison also includes the fingerprinting approaches BBHash [22], FMPH [5], and FMPHGO [5]. We also compare against RecSplit [13] and approaches based on it, such as SIMDRecSplit [6], ShockHash-RS [21], and bipartite ShockHash-RS [19]. Finally, we also compare against CHD [3] and SicHash [20].



■ **Figure 7** Construction throughput (left) and query throughput (right) for various methods on 100 million keys and using a single processing thread.

Each method has a wide range of configurations that provide a trade-off between space, construction time, and query time. To give an initial overview, we show a Pareto front for each method in Figure 7. A configuration is on the Pareto front if no other configuration of the same method is simultaneously faster and more space efficient. For this plot we use a single thread (a multithreaded measurement would mainly show what method implemented the partitioning step most efficiently instead of focusing on the algorithmic aspects). The figure shows that PTHash and PHOBIC are clear winners in terms of query performance. Even though BBHash [22] and FMPH [5] are also focused on fast queries, they are significantly slower than PTHash and PHOBIC. Figure 4 shows that PHOBIC consistently saves about 0.17 bits/key for a large range of different construction times while maintaining the good query speed. We remark that this is a significant reduction in space considering the proximity to the space lower bound. The competitors achieving even lower space consumption (i.e. RecSplit [13], SIMDRecSplit [6], ShockHash-RS [21], and bipartite ShockHash-RS [19]) all have a rather slow query performance. However, somewhat surprisingly, SIMDRecSplit has the fastest construction even for less space efficient configurations. SicHash [20] takes a middle ground with faster construction than PHOBIC and query performance between PHOBIC and the RecSplit variants.

Table 1 gives a selection of configuration parameters for direct comparison, mostly taken from the corresponding papers. Appendix D.2 gives the same table measured on a large machine with 64 threads. Comparing configurations with the same space consumption, PHOBIC is significantly faster to construct than the original PTHash implementation. Comparing configurations that both need 1.86 bits/key and have a similar query time, PHOBIC can be constructed 83 times faster than PTHash.

On the GPU, we compare against the only available GPU construction, RecSplit-GPU [6]. Figure 9 in Appendix D.2 illustrates the comparison. Basically, we achieve the same peak construction throughput as RecSplit-GPU for the less space efficient configurations. The queries of both approaches are done on the CPU, so the fact that PHOBIC offers much faster queries applies here as well (see Figure 7). Comparing the multithreaded CPU implementation and the GPU implementation of PHOBIC, we get a construction speedup of 62 for $\lambda = 9$ with interleaved Rice coding. Note that with $\lambda = 9$, the GPU still spends a lot of its construction time on transferring the input data (see Figure 6b), but much larger values of λ are not

■ **Table 1** Performance of various methods on 100 million keys.

Method	Space (bits/key)	Query (ns/query)	Construction (ns/key)		
			1 Thread	8 Threads	Speedup
Bip. SH-RS, $n=64, b=2000$	1.52	160	5 756	1 218	4.7
CHD, $\lambda=3$	2.27	222	352	-	-
CHD, $\lambda=5$	2.07	207	2 206	-	-
FMPH, $\gamma=2.0$	3.40	100	69	17	4.0
FMPH, $\gamma=1.0$	2.80	134	99	24	4.0
SIMDRecSplit, $n=8, b=100$	1.81	124	109	20	5.2
SIMDRecSplit, $n=14, b=2000$	1.58	143	11 062	2 360	4.7
SicHash, $\alpha=0.9, p_1=21, p_2=78$	2.41	72	129	25	5.0
SicHash, $\alpha=0.97, p_1=45, p_2=31$	2.08	64	179	32	5.6
PTHash, $\lambda=4.0, \alpha=0.99, C-C$	3.19	35	314	143	2.2
PTHash, $\lambda=5.0, \alpha=0.99, EF$	2.11	54	525	252	2.1
PTHash, $\lambda=10.5, \alpha=0.99, EF$	1.86	49	82 721	35 048	2.4
PTHash-HEM, $\lambda=4.0, \alpha=0.99, C-C$	3.19	39	299	45	6.6
PTHash-HEM, $\lambda=5.0, \alpha=0.99, EF$	2.11	58	582	86	6.7
PHOBIC, $\lambda=3.9, \alpha=1.0, IC-C$	3.18	40	197	32	6.2
PHOBIC, $\lambda=4.5, \alpha=1.0, IC-R$	2.11	57	254	40	6.2
PHOBIC, $\lambda=6.5, \alpha=1.0, IC-R$	1.85	52	992	176	5.6
PHOBIC, $\lambda=9.0, \alpha=1.0, IC-R$	1.74	50	9 171	1 781	5.1
GPU + 8 CPU Threads					
PHOBIC-GPU, $\lambda=9.0, IC-C$	2.17	37		28	
PHOBIC-GPU, $\lambda=9.0, IC-R$	1.76	52		27	
PHOBIC-GPU, $\lambda=13.0, IC-R$	1.68	50		560	
PHOBIC-GPU, $\lambda=14.0, IC-R$	1.67	49		1 470	
RecSplit-GPU, $\ell=8, b=100$	1.81	126		24	
RecSplit-GPU, $\ell=14, b=2000$	1.58	147		80	
RecSplit-GPU, $\ell=18, b=2000$	1.55	135		1 732	

feasible on the CPU.

Directly comparing the performance of CPU and GPU is always difficult because of the different hardware architectures. Given that the power consumption is a major cost factor in production environments, we measure it using a Voltcraft Energy Check 3000 wattmeter. For CPU-only measurements, we dismount the GPU. The machine requires about 405 W during the search step of the GPU version and 195 W for the multithreaded CPU implementation. Thus, the above speedup of 62 translates to roughly 30 times lower energy consumption for constructing an MPHf on the GPU. Single-threaded CPU construction requires 74 W which is less energy efficient compared to multithreading.

6 Conclusion and Future Work

PHOBIC introduces optimized bucket sizes and interleaved encoding to PTHash. Our improvements result in 0.17 bits/key better space efficiency when compared to PTHash for similar construction and query speed. When compared for the same space consumption,

PHOBIC can be constructed up to 83 times faster than PTHash, while still having the same query time. Finally, our GPU implementation improves the construction by a factor of up to 62 compared to the multithreaded CPU implementation.

Future work may explore combinations of the most time efficient approaches to perfect hashing and the most space efficient approaches. Concretely, we are hopeful that a hybrid between PHOBIC and ShockHash [21] puts further trade-offs between space and time into reach.

References

- 1 xxhash github. URL: <https://github.com/Cyan4973/xxHash>.
- 2 Djamel Belazzougui, Paolo Boldi, Rasmus Pagh, and Sebastiano Vigna. Fast prefix search in little space, with applications. In *ESA (1)*, volume 6346 of *Lecture Notes in Computer Science*, pages 427–438. Springer, 2010. doi:10.1007/978-3-642-15775-2_37.
- 3 Djamel Belazzougui, Fabiano C. Botelho, and Martin Dietzfelbinger. Hash, displace, and compress. In *ESA*, volume 5757 of *Lecture Notes in Computer Science*, pages 682–693. Springer, 2009. doi:10.1007/978-3-642-04128-0_61.
- 4 Djamel Belazzougui and Gonzalo Navarro. Alphabet-independent compressed text indexing. *ACM Trans. Algorithms*, 10(4):23:1–23:19, 2014. doi:10.1145/2635816.
- 5 Piotr Beling. Fingerprinting-based minimal perfect hashing revisited. *ACM J. Exp. Algorithmics*, 28:1.4:1–1.4:16, 2023. doi:10.1145/3596453.
- 6 Dominik Bez, Florian Kurpicz, Hans-Peter Lehmann, and Peter Sanders. High performance construction of recsplit based minimal perfect hash functions. In *ESA*, volume 274 of *LIPICs*, pages 19:1–19:16. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2023. doi:10.4230/LIPICs.ESA.2023.19.
- 7 Fabiano C. Botelho, Rasmus Pagh, and Nivio Ziviani. Simple and space-efficient minimal perfect hash functions. In *WADS*, volume 4619 of *Lecture Notes in Computer Science*, pages 139–150. Springer, 2007. doi:10.1007/978-3-540-73951-7_13.
- 8 Andrei Z. Broder and Michael Mitzenmacher. Survey: Network applications of Bloom filters: A survey. *Internet Math.*, 1(4):485–509, 2003. doi:10.1080/15427951.2004.10129096.
- 9 Chin-Chen Chang and Chih-Yang Lin. Perfect hashing schemes for mining association rules. *Comput. J.*, 48(2):168–179, 2005. doi:10.1093/COMJNL/BXH074.
- 10 Jarrod A. Chapman, Isaac Ho, Sirisha Sunkara, Shujun Luo, Gary P. Schroth, and Daniel S. Rokhsar. Meraculous: De novo genome assembly with short paired-end reads. *PLOS ONE*, 6(8):1–13, 08 2011. doi:10.1371/journal.pone.0023501.
- 11 K. P. Choi. On the medians of gamma distributions and an equation of Ramanujan. *Proceedings of the American Mathematical Society*, 121:245–251, 1994.
- 12 Victoria G. Crawford, Alan Kuhnle, Christina Boucher, Rayan Chikhi, and Travis Gagie. Practical dynamic de bruijn graphs. *Bioinform.*, 34(24):4189–4195, 2018. doi:10.1093/BIOINFORMATICS/BTY500.
- 13 Emmanuel Esposito, Thomas Mueller Graf, and Sebastiano Vigna. RecSplit: Minimal perfect hashing via recursive splitting. In *ALENEX*, pages 175–185. SIAM, 2020. doi:10.1137/1.9781611976007.14.
- 14 Edward A. Fox, Qi Fan Chen, and Lenwood S. Heath. A faster algorithm for constructing minimal perfect hash functions. In *SIGIR*, pages 266–273. ACM, 1992. doi:10.1145/133160.133209.
- 15 Marco Genuzio, Giuseppe Ottaviano, and Sebastiano Vigna. Fast scalable construction of (minimal perfect hash) functions. In *SEA*, volume 9685 of *Lecture Notes in Computer Science*, pages 339–352. Springer, 2016. doi:10.1007/978-3-319-38851-9_23.
- 16 Solomon W. Golomb. Run-length encodings (corresp.). *IEEE Trans. Inf. Theory*, 12(3):399–401, 1966. doi:10.1109/TIT.1966.1053907.

- 17 Stefan Hermann. Accelerating minimal perfect hash function construction using gpu parallelization. Master's thesis, Karlsruhe Institute for Technology (KIT), 2023. doi: [10.5445/IR/1000164413](https://doi.org/10.5445/IR/1000164413).
- 18 Aaron Kiely. Selecting the Golomb parameter in Rice coding. *IPN progress report*, 42:159, 2004.
- 19 Hans-Peter Lehmann, Peter Sanders, and Stefan Walzer. Bipartite ShockHash: Pruning ShockHash search for efficient perfect hashing. *CoRR*, abs/2310.14959, 2023. doi: [10.48550/ARXIV.2310.14959](https://doi.org/10.48550/ARXIV.2310.14959).
- 20 Hans-Peter Lehmann, Peter Sanders, and Stefan Walzer. SicHash – small irregular cuckoo tables for perfect hashing. In *ALLENEX*, pages 176–189. SIAM, 2023. doi: [10.1137/1.9781611977561.CH15](https://doi.org/10.1137/1.9781611977561.CH15).
- 21 Hans-Peter Lehmann, Peter Sanders, and Stefan Walzer. Shockhash: Towards optimal-space minimal perfect hashing beyond brute-force. In *ALLENEX*. SIAM, 2024. doi: [10.1137/1.9781611977929.15](https://doi.org/10.1137/1.9781611977929.15).
- 22 Antoine Limasset, Guillaume Rizk, Rayan Chikhi, and Pierre Peterlongo. Fast and scalable minimal perfect hashing for massive key sets. In *SEA*, volume 75 of *LIPICs*, pages 25:1–25:16. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2017. doi: [10.4230/LIPICs.SEA.2017.25](https://doi.org/10.4230/LIPICs.SEA.2017.25).
- 23 Yi Lu, Balaji Prabhakar, and Flavio Bonomi. Perfect hashing for network applications. In *ISIT*, pages 2774–2778. IEEE, 2006. doi: [10.1109/ISIT.2006.261567](https://doi.org/10.1109/ISIT.2006.261567).
- 24 Colin McDiarmid. *On the method of bounded differences*, page 148–188. London Mathematical Society Lecture Note Series. Cambridge University Press, 1989. doi: [10.1017/CB09781107359949.008](https://doi.org/10.1017/CB09781107359949.008).
- 25 Kurt Mehlhorn. On the program size of perfect and universal hash functions. In *FOCS*, pages 170–175. IEEE Computer Society, 1982. doi: [10.1109/SFCS.1982.80](https://doi.org/10.1109/SFCS.1982.80).
- 26 Ingo Müller, Peter Sanders, Robert Schulze, and Wei Zhou. Retrieval and perfect hashing using fingerprinting. In *SEA*, volume 8504 of *Lecture Notes in Computer Science*, pages 138–149. Springer, 2014. doi: [10.1007/978-3-319-07959-2_12](https://doi.org/10.1007/978-3-319-07959-2_12).
- 27 Giulio Ermanno Pibiri. Sparse and skew hashing of k-mers. *Bioinformatics*, 38(Supplement_1):i185–i194, 2022.
- 28 Giulio Ermanno Pibiri and Roberto Trani. Pthash: Revisiting FCH minimal perfect hashing. In *SIGIR*, pages 1339–1348. ACM, 2021. doi: [10.1145/3404835.3462849](https://doi.org/10.1145/3404835.3462849).
- 29 Giulio Ermanno Pibiri and Roberto Trani. Parallel and external-memory construction of minimal perfect hash functions with pthash. *IEEE Trans. Knowl. Data Eng.*, 36(3):1249–1259, 2024. doi: [10.1109/TKDE.2023.3303341](https://doi.org/10.1109/TKDE.2023.3303341).
- 30 Giulio Ermanno Pibiri and Rossano Venturini. Efficient data structures for massive N -gram datasets. In *SIGIR*, pages 615–624. ACM, 2017. doi: [10.1145/3077136.3080798](https://doi.org/10.1145/3077136.3080798).
- 31 Robert F Rice. Some practical universal noiseless coding techniques. Technical report, 1979.
- 32 Peter Sanders. Emulating MIMD behaviour on SIMD-machines. In *EUROSIM*, pages 313–320. Elsevier, 1994.

A Full proofs

We begin with a more formal statement of what we need from Intuition 3. Let $\mu_i := n \cdot w_i$ be the expected size of bucket i (i.e. the expected number of keys assigned to bucket i) and again $\lambda := n/B$.

► **Observation 6.** Let $\gamma : [0, 1] \rightarrow [0, 1]$ be a continuous bucket assignment function that is smooth on $(0, 1)$ with non-decreasing derivative, $x_0 \in (0, 1)$ a hash and $i = \lceil B \cdot \gamma(x_0) \rceil$ the bucket assigned to x . Let $\mu = \lambda/\gamma'(x_0)$. Then the expected bucket sizes satisfy

- $\mu_j \geq \mu$ for all $j < i$,
- $\mu_j \leq \mu$ for all $j > i$,

Proof. Consider $j < i$ and the range $(\sigma_{j-1}, \sigma_j]$ assigned to bucket j . Since x_0 is assigned to bucket i we have $\sigma_j < x_0$. By monotonicity $\gamma'(x) < \gamma'(x_0)$ for all $x \in [\sigma_{j-1}, \sigma_j]$. Since $\gamma(\sigma_j) - \gamma(\sigma_{j-1}) = 1/B$ by construction we have

$$\frac{1}{B} = \int_{\sigma_{j-1}}^{\sigma_j} \gamma'(x) dx \leq \int_{\sigma_{j-1}}^{\sigma_j} \gamma'(x_0) dx = (\sigma_j - \sigma_{j-1}) \cdot \gamma'(x_0) = w_j \cdot \frac{\lambda}{\mu} = \frac{\mu_j \lambda}{n \mu} = \frac{\mu_j}{\mu} \frac{1}{B}.$$

Rearranging gives $\mu_j \geq \mu$. The second claim is obtained analogously. ◀

A.1 Proof of Theorem 4

We assume the context of Theorem 4. In particular, γ is a monotonic function, smooth on $(0, 1)$ and satisfies w.l.o.g. $\gamma(0) = \beta_*(0) = 0$ and $\gamma(1) = \beta_*(1) = 1$. We have also assumed that $\gamma(x) > 0$ for $x > 0$, but never defended this assumption. Let us deal with this rather silly case now, i.e. assume $\gamma(x) = 0$ for some $x > 0$. Then by monotonicity $\gamma(x') = 0$ for all $x' \in [0, x]$. The first bucket then receives at least nx keys in expectation and $\Omega(nx)$ keys with high probability. A simple argument shows that if $n \geq 1/x^3$ then the expected number of seeds that need to be tried for the first bucket is $\exp(\Omega(n^{1/3}))$. This vastly exceeds the lower bound we wish to prove even for $\varepsilon = 1$.

It was therefore w.l.o.g. that we assumed $\gamma(x) > 0$ for all $x \in (0, 1]$. Since γ' is non-decreasing this implies $\gamma'(x) > 0$ for all $x \in (0, 1)$.

We now use the assumption that $\gamma \neq \beta_*$. Let $x \in (0, 1)$ with $\gamma(x) \neq \beta_*(x)$. If $\gamma(x) > \beta_*(x)$ then there is some $y \in [x, 1)$ with $\gamma'(y) < \beta'_*(y)$ and if $\gamma(x) < \beta_*(x)$ then there is some $y \in (0, x]$ with $\gamma'(y) < \beta'_*(y)$. In both cases we have $1/\gamma'(y) > 1/\beta'_*(y)$. Because γ' and β'_* are continuous on $(0, 1)$ we can find $\varepsilon > 0$ and $z \in (0, 1 - 3\varepsilon)$ such that

$$1/\gamma'(z + 3\varepsilon) > (1 + 3\varepsilon)/\beta'_*(z).$$

Let i be the bucket that γ assigns to hash $z + 3\varepsilon$. Then all hashes in $[0, z + 3\varepsilon]$ are assigned to buckets with index at most i . Using a concentration bound and assuming large n , the number of keys assigned to buckets with index less than i is at least $n \cdot (z + 2\varepsilon)$. By Observation 6 all these buckets have expected size at least $\lambda/\gamma'(z + 3\varepsilon)$, hence expected size at least $\lambda(1 + 3\varepsilon)/\beta'_*(z)$. For large λ , most will have actual size at least $s := \lambda(1 + 2\varepsilon)/\beta'_*(z)$. Again by a concentration bound, the number of keys in buckets of actual size at least s is at least $n \cdot (z + \varepsilon)$. In particular, at least εn keys are in buckets of size at least s that are processed when the load factor is already at least z . To get a lower bound on the work this causes, we may assume that these εn keys are in buckets of size exactly s and processed at load factor exactly z . By Equation (1) the expected cost of each such bucket is then lower bounded by

$$(1 - z)^{-s} = (1 - z)^{-\lambda(1+2\varepsilon)/\beta'_*(z)} = (1 - z)^{\lambda(1+2\varepsilon)/\ln(1-z)} = e^{\lambda(1+2\varepsilon)}.$$

The last step uses that $c^{1/\ln(c)} = e$ for all $c > 0$. Multiplying this by the number $\varepsilon n/s = \mathcal{O}(\varepsilon n/\lambda)$ of such buckets yields a contribution of $n \cdot \mathcal{O}(\varepsilon/\lambda) \cdot e^{\lambda(1+2\varepsilon)}$. This is at least $n \cdot e^{\lambda(1+\varepsilon)}$ if $\lambda \geq \lambda_0(\varepsilon)$ is large enough. Since the buckets of size 1 contribute a cost of w_{coupon} and are distinct from the buckets that contributed to $n \cdot e^{\lambda(1+\varepsilon)}$, we obtain the lower bound of $n \cdot e^{\lambda(1+\varepsilon)} + w_{\text{coupon}}$ as claimed.

A.2 Proof of Theorem 5

In this section we prove Theorem 5 and assume the corresponding context. In particular, we are given $\varepsilon > 0$, which defines a bucket assignment function β_ε , we assume $\lambda \geq \lambda_0(\varepsilon)$ is large enough and $n \geq n_0(\lambda, \varepsilon)$ is large enough, which defines $B = n/\lambda$. We begin with a few definitions and give corresponding intuition.

► **Definition 7.** Let $i \in [B]$ and $s \in \mathbb{N}$. We define the following.

- $\lambda_i := n \cdot (\beta_\varepsilon^{-1}(i/B) - \beta_\varepsilon^{-1}((i-1)/B))$ is the expected size of the i th bucket. The stated formula involves the length of the range of hashes that β_ε maps to bucket i .
- s_i is the number of keys assigned to bucket i . This random variable has distribution $s_i \sim \text{Bin}(n, \lambda_i)$.
- $\alpha_s := \frac{1}{n} \sum_{i=1}^B s_i \cdot \mathbb{1}_{s_i \geq s}$ is the random fraction of keys within buckets of size at least s and therefore the load factor after all buckets of size at most s have been processed.
- $d_s := 1 - e^{-\lambda/s}$ is the deadline for bucket size s . It is the load factor up to which Equation (1) guarantees that processing a bucket of size s incurs an expected cost of at most $s \cdot e^\lambda$.

The key lemma of this section proves a suitably weakened variant of the claim “ $\forall s : \alpha_s \leq d_s$ ”, i.e. that buckets of size s or larger are handled before their deadline d_s .

► **Lemma 8 (Deadline Lemma).**

- (i) For $2 \leq s \leq \sqrt{\lambda}$: $(\frac{1-\alpha_s}{1-d_s})^s = e^{-\mathcal{O}(\varepsilon A)}$ whp.
- (ii) For $\sqrt{\lambda} < s \leq 2e^2 A/\varepsilon$: $(\frac{1-\alpha_s}{1-d_s})^s = e^{-\mathcal{O}(\varepsilon A)}$ whp.
- (iii) For $2e^2 A/\varepsilon < s < 2 \frac{\log n}{\log \log n}$: $\alpha_s \leq d_s$ whp.
- (iv) For $s = 2 \frac{\log n}{\log \log n}$: $\alpha_s = 0$ whp.

Each case is handled in a separate paragraph in Appendix A.2.1. We now check that the Deadline Lemma implies Theorem 5.

Proof of Theorem 5. We apply Lemma 8 for each $2 \leq s \leq 2 \frac{\log n}{\log \log n}$. Since each corresponding event holds whp (i.e. with probability $1 - \mathcal{O}(n^{-c})$ for some $c > 0$) and since there are $\mathcal{O}(\frac{\log n}{\log \log n})$ events, they *jointly* hold whp. We may assume that this is the case.

We have to bound the expected cost for handling the buckets by $n \cdot e^{\lambda(1+\mathcal{O}(\varepsilon))} + w_{\text{coupon}}$. Consider therefore any bucket b with actual size s . By Lemma 8, we know $s \leq 2 \frac{\log n}{\log \log n}$. If $s = 1$ then the work for b is accounted for by w_{coupon} . So assume $s \in \{2, \dots, 2 \frac{\log n}{\log \log n}\}$. After all buckets of size at least s are handled, the load factor is α_s . In particular, after b is handled, the load factor is at most α_s . From Equation (1) we get that the expected cost for handling b is at most $s \cdot (1 - \alpha_s)^{-s}$. If s falls into case (iii) of Lemma 8 then we can bound this as follows:

$$s(1 - \alpha_s)^{-s} \leq s(1 - d_s)^{-s} = s(e^{-\lambda/s})^{-s} = se^\lambda.$$

If s falls into case (i) or (ii) then we get

$$s(1 - \alpha_s)^{-s} = s(1 - d_s)^{-s} \left(\frac{1 - \alpha_s}{1 - d_s} \right)^{-s} = se^\lambda \cdot e^{\mathcal{O}(\varepsilon A)} = se^{(1+\mathcal{O}(\varepsilon))\lambda}.$$

Summing this over all buckets of size at least 2, the sizes of which sum to at most n , we get an overall expected cost of at most $n \cdot e^{(1+\mathcal{O}(\varepsilon))\lambda}$, which is the main term of our bound. ◀

A.2.1 Proof of the Deadline Lemma

We begin with observations that relate to several of the cases (i),(ii), (iii) and (iv). To sharpen up presentation we will occasionally use \approx , \gtrsim and \lesssim when a relation only holds whp and/or we suppress an error term that is clearly negligible in a given context when n is chosen large enough. This notation is not meant to indicate a gap in the proof.

We have previously encountered the random variable α_s , the fraction of keys in buckets of size at least s . This is not to be confused with the following non-random quantities.

► **Definition 9.** Let $\mu > 0$ and $s \in \mathbb{N}$. We define

- $x_\lambda := \frac{1}{n} \sum_{i=1}^B \lambda_i \cdot \mathbb{1}_{\lambda_i \geq \mu} \in [0, 1]$, the range of hashes assigned to buckets of expected size at least μ . Equivalently, the probability that a key is placed in a bucket of expected size least μ .
- $\hat{x}_\lambda = 1 - \exp(-\frac{\lambda/\mu - \varepsilon}{1 - \varepsilon})$, the unique number satisfying $\lambda/\beta'_\varepsilon(\hat{x}_\lambda) = \mu$ if $\mu \leq \lambda/\varepsilon$ and $\hat{x}_\lambda = 0$ otherwise. Equivalently this is the length of the range $[0, \hat{x}_\lambda]$ where $\lambda/\beta'_\varepsilon(x)$ attains values of at least μ .
- $\hat{\alpha}_s := \mathbb{E}[\alpha_s]$, the expected fraction of keys in buckets of size at least s .

► **Observation 10.** $\max_{i \in [B]} \lambda_i \leq \lambda/\varepsilon$.

Proof. This follows because $\beta'_\varepsilon(0) \geq \varepsilon$. See also Observation 6. ◀

► **Observation 11.** For any $\mu > 0$ we have $|\hat{x}_\lambda - x_\lambda| \leq \frac{\lambda}{\varepsilon n}$, hence in most contexts $\hat{x}_\lambda \approx x_\lambda$.

Proof. Let i_λ be the bucket assigned to hash \hat{x}_λ . By Observation 6 the buckets of expected size at least μ are precisely the buckets preceding bucket i_λ and, possibly, bucket i_λ itself. Let μ' be the expected size of i_λ . If $\mu' \geq \mu$ then the entire range $[0, \hat{x}_\lambda]$ of hashes is assigned to buckets of size at least μ , hence $x_\lambda \geq \hat{x}_\lambda$. If $\mu' \leq \mu$ then only a subset of the range $[0, \hat{x}_\lambda]$ of hashes is assigned to buckets of size at least μ , hence $x_\lambda \leq \hat{x}_\lambda$. The difference between the two cases is the μ'/n , the length of the range assigned to bucket i_λ . Since $\mu' \leq \lambda/\varepsilon$ by Observation 10, the claim follows. ◀

The quantities α_s and $\hat{\alpha}_s$ are probabilistically related as follows.

► **Lemma 12.** For any $s \in \mathbb{N}$: $\Pr[n \cdot \alpha_s - n \cdot \hat{\alpha}_s \geq \delta] \leq \exp(\frac{-2\delta^2}{ns^2})$.

Proof. We can directly apply the method of bounded differences [24]. The n hash values of keys are independent variables that determine $n \cdot \alpha_s$, the number of keys in buckets of size at least s . The expectation of $n \cdot \alpha_s$ is $n \cdot \hat{\alpha}_s$ by definition. Changing a single hash can affect $n \cdot \alpha_s$ by at most $\pm s$ with “+ s ” corresponding to moving a key from a bucket of size less than s to a bucket of size $s - 1$ and “- s ” to the reverse change. ◀

By choosing $\delta = n^{2/3}$ we get

► **Corollary 13.** For $n \geq n_0(s)$ large enough we have $|\hat{\alpha}_s - \alpha_s| \leq n^{-1/3}$ whp, hence in most contexts $\hat{\alpha}_s \approx \alpha_s$.

Proof. Apply Lemma 12 with $\delta = n^{2/3}$ and assume that $n \geq s^{12}$. ◀

Two further claims do not relate to our setting in particular.

▷ **Claim 14.** There exists a constant $c_0 > 0$ such that for any $\mu \geq 3$ and $X \sim \text{Po}(\mu)$ we have $\Pr[X \leq \mu - 3] \geq c_0$.

Proof. The median $\text{med}(\mu)$ of $\text{Po}(\mu)$ satisfies $\text{med}(\mu) \in (\mu - 1, \mu + 1)$ [11]. For large μ , no single outcome has high probability. This implies for $X \sim \text{Po}(\mu)$:

$$\begin{aligned} \Pr[X \leq \mu - 3] &\geq \Pr[X \leq \text{med}(\mu) \wedge |\text{med}(\mu) - X| \geq 4] \\ &\geq \Pr[X \leq \text{med}(\mu)] - \Pr[|\text{med}(\mu) - X| \leq 4] = 1/2 - o_{\mu \rightarrow \infty}(1). \end{aligned}$$

In other words, $\Pr[X \leq \mu - 3]$ is bounded away from 0 for large μ . For small $\mu \geq 3$ clearly $\Pr[X \leq \mu - 3] > 0$. Hence, the desired $c_0 > 0$ exists. ◀

▷ **Claim 15.** Let $\varepsilon \in (0, 1)$, $\mu \geq \frac{1}{2\pi\varepsilon^2}$, $X \sim \text{Po}(\mu)$ and $s = \mu/(1 - \varepsilon)$. Then

$$\Pr[X \geq s] \leq (e^\varepsilon(1 - \varepsilon))^s.$$

Note that $e^\varepsilon(1 - \varepsilon) < 1$, which follows from $1 - x < e^{-x}$ for $x \in \mathbb{R} \setminus \{0\}$.

Proof. Consider the sum $\Pr[X \geq s] = e^{-\mu} \sum_{i \geq s} \frac{\mu^i}{i!}$. The ratio between subsequent terms is $\mu/i \leq \mu/s = 1 - \varepsilon$. Hence, we can upper bound the sum by its first term and a geometric sum as follows:

$$\Pr[X \geq s] \leq e^{-\mu} \frac{\mu^s}{s!} \cdot \sum_{i \geq 0} (1 - \varepsilon)^i = e^{-\mu} \frac{\mu^s}{s!} \cdot \frac{1}{\varepsilon}.$$

Using Stirling's approximation for $s!$ and using $s \geq \mu \geq \frac{1}{2\pi\varepsilon^2}$ gives:

$$\Pr[X \geq s] \leq e^{-\mu} \frac{\mu^s e^s}{s^s \sqrt{2\pi s}} \cdot \frac{1}{\varepsilon} \leq e^{-\mu} \frac{\mu^s e^s}{s^s} = e^{-(1-\varepsilon)s} \frac{((1-\varepsilon)s)^s e^s}{s^s} = (e^\varepsilon(1 - \varepsilon))^s. \quad \blacktriangleleft$$

We will now consider the cases of Appendix A.2.1 one after the other.

(i) Buckets of size $2 \leq s \leq \lambda^{1/3}$. We have

$$1 - x_s \approx 1 - \hat{x}_s = \exp\left(-\frac{\lambda/s - \varepsilon}{1 - \varepsilon}\right) \quad \text{and} \quad 1 - x_{s+1} \approx 1 - \hat{x}_{s+1} = \exp\left(-\frac{\lambda/(s+1) - \varepsilon}{1 - \varepsilon}\right)$$

Using $s \leq \lambda^{1/3}$ we find $\lambda/s - \lambda/(s+1) = \frac{\lambda}{s(s+1)} = \Omega(\lambda^{1/3})$. This implies $\frac{\lambda/s - \varepsilon}{1 - \varepsilon} - \frac{\lambda/(s+1) - \varepsilon}{1 - \varepsilon} = \Omega(\lambda^{1/3})$ and for λ large enough we have $\exp\left(-\frac{\lambda/(s+1) - \varepsilon}{1 - \varepsilon}\right) \geq (1 + \frac{1}{c_0}) \exp\left(-\frac{\lambda/s - \varepsilon}{1 - \varepsilon}\right)$ where c_0 is the constant from Claim 14. Hence

$$\begin{aligned} x_s - x_{s+1} &\approx \hat{x}_s - \hat{x}_{s+1} = \exp\left(-\frac{\lambda/(s+1) - \varepsilon}{1 - \varepsilon}\right) - \exp\left(-\frac{\lambda/s - \varepsilon}{1 - \varepsilon}\right) \\ &\geq \frac{1}{c_0} \exp\left(-\frac{\lambda/s - \varepsilon}{1 - \varepsilon}\right) \geq \frac{1}{c_0} \exp\left(-\frac{\lambda/s}{1 - \varepsilon}\right). \end{aligned}$$

We now bound the probability $1 - \hat{\alpha}_s$ that a key k ends up in a bucket of size at most $s - 1$. Let $x \sim \mathcal{U}([0, 1])$ be its hash, $\mu \in (0, \lambda/\varepsilon)$ the *random variable* denoting the *expected* size of the bucket that k ends up in, and s' the number of other keys sharing the bucket with k . Conditioned on μ , we have $s' \sim \text{Bin}(n - 1, \mu/n)$, a distribution that is well known to be approximately $\text{Po}(\mu)$. In the following computation we use that $\Pr[s' \leq s - 2 \mid \mu]$ is decreasing in μ , and we use that $\Pr[\mu \in [s, s + 1]] = \Pr[x \in (x_{s+1}, x_s)] = x_s - x_{s+1}$.

$$\begin{aligned} 1 - \hat{\alpha}_s &= \Pr[s' \leq s - 2] \geq \Pr[s' \leq s - 2 \wedge \mu \in [s, s + 1]] \\ &= \Pr[s' \leq s - 2 \mid \mu \in [s, s + 1]] \cdot \Pr[\mu \in [s, s + 1]] \\ &\geq \Pr_{s'' \sim \text{Bin}(n-1, \frac{s+1}{n})}[s'' \leq s - 2] \cdot (x_s - x_{s+1}) \\ &\approx \Pr_{s'' \sim \text{Po}(s+1)}[s'' \leq s - 2] \cdot (x_s - x_{s+1}) \stackrel{\text{Obs. 14}}{\geq} c_0(x_s - x_{s+1}) \gtrsim \exp\left(-\frac{\lambda/s}{1 - \varepsilon}\right). \end{aligned}$$

This gives us the bound we desired, closing this case:

$$\left(\frac{1 - \alpha_s}{1 - d_s}\right)^s \approx \left(\frac{1 - \hat{\alpha}_s}{1 - d_s}\right)^s \geq \frac{\exp(-\frac{\lambda/s}{1-\varepsilon})^s}{\exp(-\lambda/s)^s} = \frac{\exp(-\frac{\lambda}{1-\varepsilon})}{\exp(-\lambda)} = e^{-\lambda/(1-\varepsilon)+\lambda} = e^{-\mathcal{O}(\varepsilon A)}.$$

(ii) Buckets of size $\lambda^{1/3} \leq s \leq 2e^2 A/\varepsilon$. We proceed in a similar way as in case **(i)** to bound the probability $\hat{\alpha}_s$ that a key k ends up in a bucket b of size at least s . Let again $\mu \in (0, \lambda/\varepsilon)$ be a random variable denoting the *expected* size of b and s' the number of other keys in b . We consider two overlapping events that cover all cases where b has size at least s . On the one hand, b might have large expected size. On the other hand, b might have small expected size and still have *actual* size at least s .

$$\begin{aligned} \hat{\alpha}_s &= \Pr[s' \geq s - 1] \leq \Pr[\mu \geq (1 - \varepsilon)s \vee (\mu < (1 - \varepsilon)s \wedge s' \geq s - 1)] \\ &\leq \Pr[\mu \geq (1 - \varepsilon)s] + \Pr[\mu < (1 - \varepsilon)s] \cdot \Pr[s' \geq s - 1 \mid \mu < (1 - \varepsilon)s] \\ &\leq x_{(1-\varepsilon)s} + (1 - x_{(1-\varepsilon)s}) \cdot \Pr_{s'' \sim \text{Bin}(n-1, \frac{(1-\varepsilon)s}{n})}[s'' \geq s - 1] \\ &\gtrsim x_{(1-\varepsilon)s} + (1 - x_{(1-\varepsilon)s}) \cdot \Pr_{s'' \sim \text{Po}((1-\varepsilon)s)}[s'' \geq s - 1] \\ &\leq x_{(1-\varepsilon)s} + (1 - x_{(1-\varepsilon)s}) \cdot 2((1 - \varepsilon)e^\varepsilon)^s \end{aligned}$$

The last step uses Claim 15, very conservatively accounting for a “−1” discrepancy with a factor of 2.

We now make another minor case distinction, first assuming that $(1 - \varepsilon)s \leq \lambda/\varepsilon$. In that case $1 - x_{(1-\varepsilon)s} = \exp(-\frac{\lambda/(s \cdot (1-\varepsilon)) - \varepsilon}{1-\varepsilon}) \geq \exp(-\frac{\lambda/s}{(1-\varepsilon)^2})$. We can now turn to $1 - \hat{\alpha}_s$ and quotient with $1 - d_s$. We assume ε is small enough such that $1/(1 - \varepsilon)^2 \leq 1 + 3\varepsilon$ and that $\lambda \geq \lambda_0(\varepsilon)$ is large enough such that $2s((1 - \varepsilon)e^\varepsilon)^s \leq \varepsilon$ for all $s \geq \lambda^{1/3}$.

$$\begin{aligned} 1 - \hat{\alpha}_s &\geq 1 - x_{(1-\varepsilon)s} - (1 - x_{(1-\varepsilon)s}) \cdot 2((1 - \varepsilon)e^\varepsilon)^s \\ &= (1 - x_{(1-\varepsilon)s})(1 - 2((1 - \varepsilon)e^\varepsilon)^s) \\ \Rightarrow \frac{1 - \hat{\alpha}_s}{1 - d_s} &\geq \frac{(1 - x_{(1-\varepsilon)s})}{1 - d_s} (1 - 2((1 - \varepsilon)e^\varepsilon)^s) \geq \frac{\exp(-\frac{\lambda/s}{(1-\varepsilon)^2})}{\exp(-\lambda/s)} (1 - 2((1 - \varepsilon)e^\varepsilon)^s) \\ &\geq \exp(-3\varepsilon A/s) (1 - 2((1 - \varepsilon)e^\varepsilon)^s) \\ \Rightarrow \left(\frac{1 - \hat{\alpha}_s}{1 - d_s}\right)^s &\geq \exp(-3\varepsilon A) (1 - 2((1 - \varepsilon)e^\varepsilon)^s)^s \\ &\geq \exp(-3\varepsilon A) (1 - 2s((1 - \varepsilon)e^\varepsilon)^s) \geq \exp(-3\varepsilon A) \cdot (1 - \varepsilon) \geq \exp(-4\varepsilon A). \end{aligned}$$

The last step again assumes that λ is large enough.

We still have to consider the case where $(1 - \varepsilon)s > \lambda/\varepsilon$. By Observation 10 there are no buckets of expected size $(1 - \varepsilon)s$ or larger, hence $x_{(1-\varepsilon)s} = 0$. This gives $1 - \hat{\alpha}_s \geq 1 - 2((1 - \varepsilon)e^\varepsilon)^s$, meaning the above derivation only involves the less critical term that comes out as $1 - \varepsilon \geq \exp(-\varepsilon A)$.

(iii) Buckets of size $2e^2 A/\varepsilon < s \leq 2\frac{\log n}{\log \log n}$. Our use of β_ε rather than β_* renders this case quite easy. However, now that the bound on s grows with n , shortcuts involving “for n large enough” and the associated notation \approx would now be suspect and have to be replaced with careful arguments.

We wish to bound $\hat{\alpha}_s$, the probability that at least $s - 1$ further keys join a given key in its bucket. Any bucket has expected size at most λ/ε by Observation 10. We can hence

argue

$$\hat{\alpha}_s \leq \binom{n-1}{s-1} \left(\frac{\lambda/\varepsilon}{n}\right)^{s-1} \leq \left(\frac{ne}{s-1}\right)^{s-1} \left(\frac{\lambda/\varepsilon}{n}\right)^{s-1} = \left(\frac{eA/\varepsilon}{s-1}\right)^{s-1} \leq e^{-s+1}$$

where the last step used $s \geq e^2 A/\varepsilon + 1$. Applying Lemma 12 with $\delta = n \cdot e^{-s}$ gives $\Pr[n \cdot \alpha_s - n \cdot \hat{\alpha}_s \geq n \cdot e^{-s}] \leq \exp(-ne^{-2s}/(2s^2))$. For $s \leq 2 \log n / \log \log n$ this probability is negligible. In particular, we have whp

$$\alpha_s = (\alpha_s - \hat{\alpha}_s) + \hat{\alpha}_s \stackrel{\text{whp}}{\leq} e^{-s} + \hat{\alpha}_s \leq e^{-s} + e^{-s+1} \leq e^{-s+2}.$$

In particular, α_s is *much* smaller than the deadline d_s , we can argue for instance by using that $1 - e^{-x} \geq x/2$ for $x \in [0, \frac{1}{2}]$ to get:

$$d_s = 1 - e^{-\lambda/s} \geq \lambda/(2s) \geq 1/s \geq e^{-s+2} \stackrel{\text{whp}}{\geq} \alpha_s.$$

(iv) Buckets of size at least $s = 2 \frac{\log n}{\log \log n}$. Again, the largest expected bucket size is at most λ/ε by Observation 6 and the probability that a key hashes to a specific bucket hence at most $\frac{\lambda/\varepsilon}{n}$. The probability $p_{\geq s}$ that the bucket has size at least s is therefore bounded by

$$p_{\geq s} \leq \binom{n}{s} \left(\frac{\lambda/\varepsilon}{n}\right)^s \leq \left(\frac{ne}{s}\right)^s \cdot \left(\frac{eA/\varepsilon}{n}\right)^s = \left(\frac{eA/\varepsilon}{s}\right)^s = 2^{s \cdot (\log(eA/\varepsilon) - \log s)}.$$

Plugging in $s = 2 \frac{\log n}{\log \log n}$ gives $p_{\geq s} = n^{-2+o(1)}$. By a union bound, the probability that at least one bucket has size at least s is at most $B \cdot p_{\geq s} \leq n^{-1+o(1)}$. Hence, whp no such bucket exists and we have $\alpha_s = 0$ as claimed.

A.3 Why buckets should be processed in order of decreasing size

In this section we prove Proposition 16, restated here for convenience.

► **Proposition 16.** *To minimize expected construction time in a PTHash context, buckets should be processed in order from largest to smallest.*

Even though the claim is very intuitive and CHD [3] and PTHash [28] implicitly assume its truth, the argument is surprisingly subtle and involves the following random process.

Let $k \in \mathbb{N}$, $p_1, \dots, p_k \in (0, 1)$ and let $\{0, \dots, k\}$ be a set of *states*. When taking a *step* in state $0 \leq i < k$, the successor state is state $i + 1$ with probability p_{i+1} and state 0 with probability $1 - p_{i+1}$. Let $w(p_1, \dots, p_k)$ be the expected number steps needed to reach state k from state 0.

► **Lemma 17.** *Assume $1 > p_1 > p_2 > \dots > p_k > 0$ and $1 \leq i < \frac{k}{2}$. Then*

$$w(p_1, \dots, p_{k-i}) + w(p_{k-i+1}, \dots, p_k) < w(p_1, \dots, p_i) + w(p_{i+1}, \dots, p_k).$$

Before proving Lemma 17, let us check that it implies Proposition 16.

Proof of Proposition 16. Assume that the sequence of bucket sizes in processing order is b_1, \dots, b_B and that $b_i < b_{i+1}$ for some $1 \leq i < B$. We will show that by switching the order of these two buckets, the expected work (number of hash function evaluations) decreases. For this, let $\ell = \sum_{j=1}^{i-1} b_j$ be the number of keys that are handled before the i th bucket. The work for buckets i and $i + 1$ is then

$$w\left(\frac{n-\ell}{n}, \frac{n-\ell-1}{n}, \dots, \frac{n-\ell-b_i+1}{n}\right) + w\left(\frac{n-\ell-b_i}{n}, \frac{n-\ell-b_i-1}{n}, \dots, \frac{n-\ell-b_i-b_{i+1}+1}{n}\right).$$

When processing the two buckets in swapped order the expected work becomes

$$w\left(\frac{n-\ell}{n}, \frac{n-\ell-1}{n}, \dots, \frac{n-\ell-b_{i+1}+1}{n}\right) + w\left(\frac{n-\ell-b_{i+1}}{n}, \frac{n-\ell-b_{i+1}-1}{n}, \dots, \frac{n-\ell-b_i-b_{i+1}+1}{n}\right).$$

By Lemma 17 this is less. Note that the expected work needed for other buckets is unchanged because their sizes and the table load when they are processed is unchanged. ◀

Intuition for Lemma 17. Assume a juggler knows, for some $a, b \in \mathbb{N}$, a routine S of a simple throws, a routine D of a difficult throws, and a routine M of b throws of intermediate difficulty. Performing a routine means attempting all throws in sequence, starting over after the first unsuccessful throw, and repeating until all throws succeed. Let $t(R)$ denote the expected number of throws when performing a routine R . It is intuitively plausible that

$$t(S \circ M) + t(D) < t(S) + t(D \circ M) < t(S) + t(M \circ D)$$

where “ \circ ” denotes concatenation of routines. The first inequality means that appending M to a difficult routine rather than the simple routine makes things more costly overall. The second inequality means that delaying difficult throws until the end of a routine increases its cost because failures tend to happen *after* the easier throws have already taken place.

Proof of Lemma 17 formalizes these insights.

Simple Observations. We make four observations about the function $w(p_1, \dots, p_k)$, each time justifying them briefly after stating them. The first serves as an alternative definition of $w(p_1, \dots, p_k)$.

$$w(p_1, \dots, p_k) = \sum_{i=1}^k \frac{1}{p_i \cdot \dots \cdot p_k}. \quad (2)$$

This holds because the number of visits to state i for $1 \leq i < k$ has a geometric distribution with parameter $p_i \cdot \dots \cdot p_k$ and hence has expectation $\frac{1}{p_i \cdot \dots \cdot p_k}$.

$$w(p_1, \dots, p_k) \text{ is monotonically decreasing in all of its parameters.} \quad (3)$$

This follows from Equation (2) because each summand is non-increasing in all parameters and the first summand is decreasing in all parameters.

$$\text{if } p_j < p_{j+1} \text{ for } 1 \leq j < k \text{ then } w(p_1, \dots, p_k) < w(p_1, \dots, p_{j-1}, p_{j+1}, p_j, p_{j+2}, \dots, p_k). \quad (4)$$

In other words, when swapping two adjacent parameters (at indices j and $j+1$), then having the larger parameter further to the left yields a larger value of w . This follows from Equation (2) because the summand with $i = j+1$ is then larger, and all other summands are the same.

$$w(p_1, \dots, p_k) \text{ is maximized if the parameters appear in non-increasing order,} \quad (5)$$

and any other ordering of the same set of parameters yields a smaller value.

This follows by iterating Equation (4) until the parameters are sorted in non-increasing order.

$$\text{for any } 1 \leq j < k: w(p_1, \dots, p_k) = \frac{1}{p_{j+1} \cdot \dots \cdot p_k} w(p_1, \dots, p_j) + w(p_{j+1}, \dots, p_k). \quad (6)$$

This follows from Equation (2) by separating the sum into the summands for $i > j$ and those for $i \leq j$. We are now ready to proof Lemma 17.

Lemma 17. Let $1 > p_1 > p_2 > \dots > p_k > 0$ and $1 \leq i < \frac{k}{2}$. Then

$$\begin{aligned}
& w(p_1, \dots, p_{k-i}) + w(p_{k-i+1}, \dots, p_k) \\
& \stackrel{(6)}{=} \frac{1}{p_{i+1} \dots p_{k-i}} w(p_1, \dots, p_i) + w(p_{i+1}, \dots, p_{k-i}) + w(p_{k-i+1}, \dots, p_k) \\
& = w(p_1, \dots, p_i) + \left(\frac{1}{p_{i+1} \dots p_{k-i}} - 1 \right) w(p_1, \dots, p_i) \\
& \quad + w(p_{i+1}, \dots, p_{k-i}) + w(p_{k-i+1}, \dots, p_k) \\
& \stackrel{(3)}{<} w(p_1, \dots, p_i) + \left(\frac{1}{p_{i+1} \dots p_{k-i}} - 1 \right) w(p_{k-i+1}, \dots, p_k) \\
& \quad + w(p_{i+1}, \dots, p_{k-i}) + w(p_{k-i+1}, \dots, p_k) \\
& = w(p_1, \dots, p_i) + \frac{1}{p_{i+1} \dots p_{k-i}} w(p_{k-i+1}, \dots, p_k) + w(p_{i+1}, \dots, p_{k-i}) \\
& \stackrel{(6)}{=} w(p_1, \dots, p_i) + w(p_{k-i+1}, \dots, p_k, p_{i+1}, \dots, p_{k-i}) \\
& \stackrel{(5)}{<} w(p_1, \dots, p_i) + w(p_{i+1}, \dots, p_k). \quad \blacktriangleleft
\end{aligned}$$

A.4 General bounds for Perfect Hashing with Bucket Placement

In this section we consider perfect hashing through bucket placement in general and without any partitioning (see introduction).

► **Proposition 18.** Consider perfect hashing through bucket placement with n keys and $B := n/\lambda$ buckets for some $2 \leq \lambda = o(n/\log n)$. Let T be the sum of resulting seeds values and S the number of bits when encoding the seeds using Elias δ -encoding. Then

- (i) $\mathbb{E}[T] = \Omega(ne^\lambda/\lambda)$
- (ii) $\mathbb{E}[S] = n(\log_2 e + \mathcal{O}(\frac{1}{\lambda} \log \lambda))$.

Note: (i) implies that the expected construction time per key of $e^{\lambda(1+\varepsilon)}$ achieved using β_ε in Section 3 (when ignoring buckets of size 1) is almost optimal. (ii) means that the space consumption of perfect hashing through bucket placement approaches the lower bound $n \log_2(e)$ for large λ , in the sense that it is $n(\log_2 e + o_{\lambda \rightarrow \infty}(1))$, regardless of the bucket assignment function that is used.

Proof. The first step of perfect hashing through bucket placement is to map the n keys into B buckets. The sizes s_1, \dots, s_B of these buckets are random variables. However, this arguments makes no use of their distributions and works for arbitrary bucket sizes. We focus on the second step where the buckets are placed one after the other using brute force. Assume the buckets are indexed in placement order. The i th bucket is then associated with the probability p_i that randomly hashing a set of s_i keys into a table does not cause any collision assuming $s_1 + \dots + s_{i-1}$ out of n positions are already occupied.

A useful observation is that the product of all p_i is the probability that a random function of n keys to n positions is a bijection, hence

$$\prod_{i=1}^B p_i = \frac{n!}{n^n} \text{ or equivalently } \prod_{i=1}^B 1/p_i = \frac{n^n}{n!}.$$

If the i th seed value is σ_i then its distribution is $\sigma_i \sim \text{Geom}(p_i)$ implying $\mathbb{E}[\sigma_i] = 1/p_i$. Hence

$$\mathbb{E}[T] = \mathbb{E}\left[\sum_{i=1}^B \sigma_i\right] = \sum_{i=1}^B \mathbb{E}[\sigma_i] = \sum_{i=1}^B 1/p_i \geq B \cdot \left(\frac{n^n}{n!}\right)^{1/B}.$$

The last inequality uses that the sum of the side lengths of a B -dimensional cuboid with volume $V = \frac{n^n}{n!}$ is minimized by the B -dimensional cube where all B side lengths are equal to $V^{1/B}$. Continuing with Stirling's approximation yields

$$\mathbb{E}[T] \geq B \cdot \left(\frac{n^n}{n!}\right)^{1/B} = B \left(\frac{e^n}{\mathcal{O}(\sqrt{n})}\right)^{1/B} = \frac{n}{\lambda} e^{\lambda} n^{-\mathcal{O}(\lambda/n)} = \Omega(ne^\lambda/\lambda),$$

where the last step used $\lambda = o(n/\log n)$. This proves **(i)**.

To encode a number $x \in \mathbb{N}$, Elias δ -coding requires $\lceil \log_2(x) \rceil + 2\lceil \log(\lceil \log_2(x) \rceil + 1) \rceil + 1 \leq \log_2(x) + 2\log(2 + \log x) + 1$ bits. Using that $\mathbb{E}[\log_2 X] \leq \log_2 \mathbb{E}[X]$ for any non-negative random variable X by Jensen's inequality allows us to bound

$$\begin{aligned} \mathbb{E}[S] &= \mathbb{E}\left[\sum_{i=1}^B (\log_2(\sigma_i) + 2\log_2(2 + \log_2 \sigma_i) + 1)\right] \\ &\leq \sum_{i=1}^B \log_2 \mathbb{E}[\sigma_i] + 2 \sum_{i=1}^B \log_2(2 + \log_2 \mathbb{E}[\sigma_i]) + B \\ &= \sum_{i=1}^B \log_2(1/p_i) + 2 \sum_{i=1}^B \log_2(2 + \log_2 1/p_i) + B. \end{aligned}$$

This first sum can be bounded by $\sum_{i=1}^B \log_2(1/p_i) = \log_2 \frac{n^n}{n!} \leq n \log_2 e$ using Stirling's approximation. The only thing left to prove is that the second sum is subsumed by the error term $\mathcal{O}(n \frac{\log_2 \lambda}{\lambda})$ in our stated bound. We again using Jensen's inequality and the fact that $\sum_{i=1}^B \log_2 1/p_i = \log_2 \frac{n^n}{n!} \leq n \log_2 e$.

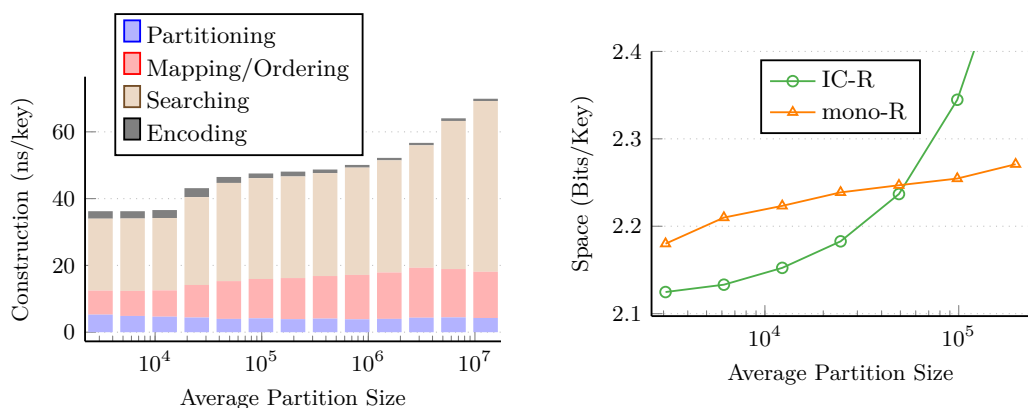
$$\begin{aligned} \sum_{i=1}^B \log_2(2 + \log_2 1/p_i) &\leq B \cdot \log_2 \left(2 + \frac{1}{B} \log_2 \frac{n^n}{n!}\right) \\ &\leq B \cdot \log_2 \left(2 + \frac{n}{B} \log_2 e\right) = \frac{n}{\lambda} \cdot \log_2 \left(2 + \lambda \log_2 e\right) = \mathcal{O}\left(n \frac{\log_2 \lambda}{\lambda}\right). \end{aligned}$$

The last step uses $\lambda \geq 2$. ◀

B Bucket assignment for non-minimal PHF

Assume we wish to construct a *non-minimal* perfect hash function using the “perfect hashing through bucket placement” framework. Hence, assume the range of the hash function is $[m]$ and $n = \alpha m$ for $\alpha \in (0, 1)$.

In Section 3 we only discussed the case $\alpha = 1$. There is, however, a natural way to obtain from a bucket assignment function γ for $\alpha = 1$ a bucket assignment function γ_α for the non-minimal case. Namely, we can define $\gamma_\alpha(x) = \frac{1}{\gamma(\alpha)} \gamma(\alpha x)$, which rescales whatever γ does on its subdomain $[0, \alpha]$. The idea is that, in order to fill a table to load factor 1, you first need to fill it to load factor α , and if γ is optimal overall then it should implicitly contain an optimal strategy for this first phase. Since γ defines on $[0, \alpha]$ how an α -fraction of the keys is assigned to buckets of largest expected size (presumably processed first), this is where we should look. This argument remains heuristic because the buckets of largest *expected* size may not turn out to be the buckets of largest *actual* size.



(a) Time for different construction steps when varying the partition size.

(b) Space consumption for placing all seeds into a single encoder (mono) and using interleaved coding when varying the expected partition size.

■ **Figure 8** Varying partition count using the original implementation and $\lambda = 4$ and $n = 100$ million keys. We use the original PTHash bucket assignment function and their hash function.

C Implementation Details

C.1 Details on the Bucket Assignment in practice

Our implementation of bucket assignment functions differs from our theoretical results in two minor ways.

Perturbation. Recall that $\beta_* : [0, 1] \rightarrow [0, 1]$ was, modulo certain qualifications, identified as the optimized bucket assignment function in Section 3. One of the qualifications was that we actually analyze a slightly perturbed function $\beta_\varepsilon(x) := \varepsilon x + (1 - \varepsilon)\beta_*(x)$. This limits expected bucket sizes to λ/ε , which helped with bounding construction times of large buckets. Limiting the bucket sizes is also useful in practice to reduce self-collisions inside the small partitions. Concretely we choose $\varepsilon = \frac{\lambda}{5\sqrt{P}}$, for an expected partition size of P . Without this perturbation ($\varepsilon = 0$), running times are noticeably worse.

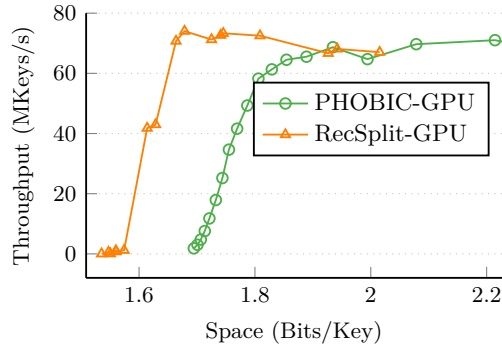
Tabulating values. The functions β_* and γ_P involve expensive arithmetic operations such as a logarithm. We achieve a significant speedup by tabulating γ_P for 2048 discrete values of x and interpolating linearly.

D Additional Experimental Data

In this section, we give additional experiments that further illustrate the data in Section 5.

D.1 From PTHash to PHOBIC

In Figure 8, we give additional details on the effect of fine-grained partitioning. As stated in the main body of the paper, using very small partitions in itself has only a marginal effect on reducing construction time. The search step does not get much faster when using very small partitions instead of medium sized ones, while there is a small additional overhead in the partitioning step. Figure 8a shows this effect. However, using small partitions enables the using of interleaved coding (see Section 4.1), whose effect is shown in Figure 8b. For small partitions, we can achieve significant space improvements; for large partitions instead, we need too many encoders and this leads to some constant overhead.



■ **Figure 9** Comparison between the GPU implementations of PTHash and RecSplit, for 100 million keys. The time is measured for the entire construction, including data transfers to/from the GPU.

D.2 Comparison to Other Methods

In Section 5.2, we have compared PHOBIC to state-of-the-art MPHFs using multithreading and the GPU. In particular, Table 1 illustrates the overall comparison. Table 2 illustrates the result of the same experiment but using a machine with more cores. This machine is equipped with an AMD EPYC Rome 7702P processor with 64 cores (with hyper-threading) and a clock frequency of 2.0 GHz. While this machine does not have a GPU, it shows how well PHOBIC scales with a large number of threads. We also note that the query times are higher compared to those measured with the machine used for Table 1 because of the different clock frequency.

Lastly, Figure 9 shows the construction throughput of PHOBIC and RecSplit on the GPU, by varying the number of bits/key. As it is clear, PHOBIC construction is as fast as RecSplit – the fastest construction method up to date. However, we remark that PHOBIC is considerably faster to query than RecSplit.

■ **Table 2** Performance of various methods on 100 million keys on a machine with 64 cores.

Method	Space (bits/key)	Query (ns/query)	Construction (ns/key)		
			1 Thread	64 Threads	Speedup
Bip. SH-RS, $n=64$, $b=2000$	1.52	425	9 365	356	26.3
CHD, $\lambda=3$	2.27	357	756	-	-
CHD, $\lambda=5$	2.07	320	5 483	-	-
FMPH, $\gamma=2.0$	3.40	216	126	9	12.9
FMPH, $\gamma=1.0$	2.80	269	190	15	12.1
SIMDRecSplit, $n=8$, $b=100$	1.81	315	303	13	22.1
SIMDRecSplit, $n=14$, $b=2000$	1.59	368	40 168	763	52.6
SicHash, $\alpha=0.9$, $p_1=21$, $p_2=78$	2.41	159	191	7	24.1
SicHash, $\alpha=0.97$, $p_1=45$, $p_2=31$	2.08	143	260	17	15.1
PTHash, $\lambda=4.0$, $\alpha=0.99$, C-C	3.19	79	504	304	1.7
PTHash, $\lambda=5.0$, $\alpha=0.99$, EF	2.11	160	1 016	307	3.3
PTHash, $\lambda=10.5$, $\alpha=0.99$, EF	1.86	150	139 184	10 312	13.5
PTHash-HEM, $\lambda=4.0$, $\alpha=0.99$, C-C	3.19	88	489	9	49.1
PTHash-HEM, $\lambda=5.0$, $\alpha=0.99$, EF	2.11	166	930	13	66.9
PHOBIC, $\lambda=3.9$, $\alpha=1.0$, IC-C	3.18	94	299	9	31.3
PHOBIC, $\lambda=4.5$, $\alpha=1.0$, IC-R	2.11	240	369	11	33.5
PHOBIC, $\lambda=6.5$, $\alpha=1.0$, IC-R	1.85	217	1 341	29	46.0
PHOBIC, $\lambda=9.0$, $\alpha=1.0$, IC-R	1.74	205	12 138	251	48.3