# Parallel and External-Memory Construction of Minimal Perfect Hash Functions with PTHash

Giulio Ermanno Pibiri – Roberto Trani

**Abstract**—A *minimal perfect hash function* $f$ for a set $S$ of $n$ keys is a bijective function of the form $f : S \rightarrow \{0, \ldots, n-1\}$. These functions are important for many practical applications in computing, such as search engines, computer networks, and databases. Several algorithms have been proposed to build minimal perfect hash functions that: scale well to large sets, retain fast evaluation time, and take very little space, e.g., $2-3$ bits/key. PTHash is one such algorithm, achieving very fast evaluation in compressed space, typically several times faster than other techniques. In this work, we propose a new construction algorithm for PTHash enabling: (1) *multi-threading*, to either build functions more quickly or more space-efficiently, and (2) *external-memory processing* to scale to inputs much larger than the available internal memory. Only few other algorithms in the literature share these features, despite of their big practical impact. We conduct an extensive experimental assessment on large real-world string collections and show that, with respect to other techniques, PTHash is competitive in construction time and space consumption, but retains $2-6\times$ better lookup time.

**Index Terms**—Minimal Perfect Hashing; PTHash; Multi-Threading; External-Memory

◆

## 1 INTRODUCTION

**M**INIMAL perfect hashing (MPH) is a well-studied and fundamental problem in Computer Science. It asks to build a data structure that assigns the numbers in $[n] = \{0, \ldots, n-1\}$ to the $n$ distinct elements of a static set $S$. In other words, the resulting data structure realizes a "one-to-one" correspondence between $S$ and the integers in $[n]$, with $n = |S|$. Such bijective function $f : S \rightarrow \{0, \ldots, n-1\}$ is called a *minimal perfect hash function* (MPHF henceforth).

MPHFs are useful in all those practical situations where space-efficient storage and fast retrieval from static sets is deemed. In fact, they are employed in compressed full-text indexes [1], computer networks [2], databases [3], prefix-search data structures [4], language models [5, 6], Bloom filters and their variants [7, 8, 9], programming languages, compilers, and operating systems, just to mention some important applications.

The most interesting aspect of the MPH problem is that it ignores the behavior of $f$ on keys that do *not* belong to $S$, i.e., $f(x)$ can be any value in $[n]$ if $x \notin S$. Therefore, the MPHF data structure does not need to store the keys. As a result, pioneer work on the problem has proved a space lower bound of $n \log_2 e$ bits for the size of any MPHF [10, 11], which is approximately just 1.44 bits/key. While it is difficult to come close to this space usage, several practical algorithms exist that take little space, i.e., $2-3$ bits/key, retain very fast lookup time, and scale well to large sets, e.g., billions of keys.

Our initial investigation on the problem focused on achieving especially fast lookup [12]. We were motivated by the observation that MPHFs are usually built once and evaluated many times, thus making lookup time the most critical efficiency aspect for the MPH problem, provided that both construction time and space usage are reasonably low. PTHash [12] was proposed to accomplish this goal and shown to be significantly faster at lookup time than other techniques *while* taking a similar memory footprint. In our previous work, however, we only proposed a sequential and internal-memory construction.

In this work, we extend the treatment of PTHash and consider two important algorithmic aspects: *multi-threading* and *external-memory scaling*. While multi-threading can be used to either quicken the construction or build more space-efficient functions, temporary disk storage can be used to scale to inputs much larger than the available internal memory. Only few other algorithms in the literature support these two features, despite of their big practical impact.

A simple and elegant solution to harness both aspects is to partition the input and build, in internal memory, an independent MPHF on each partition [13]. This approach also offers good scalability as the independent partitions that fit into internal memory may be processed in parallel by multiple threads. Very importantly, this solution is valid for *any* construction algorithm. However, this approach imposes an indirection at lookup time to identify the partition, i.e., the proper MPHF to query, and additional space to store the offset of each partition. Indeed, partitioned MPHFs are usually $30-50\%$ slower and $20\%$ larger than their non-partitioned counterpart [13, 14].

**Our Contribution.** Since the construction can either run in internal or external memory, can either use one or multiple threads, we consider the four possible construction settings:

- internal-memory/sequential,

- internal-memory/parallel,
- external-memory/sequential,
- external-memory/parallel.

We present a new construction algorithm for PTHash that overcomes the overhead of the folklore partitioning approach and, yet, easily adapts to all the above four settings.

We conduct an extensive experimental assessment over real-world string collections, ranging in size from tens of millions to several billions of strings. We show that PTHash is competitive at building MPHFs with the best existing techniques that also support parallel execution and external-memory scaling. However and very importantly, PTHash retains the best lookup time, i.e., $2 - 6\times$ smaller than other techniques.

**Source Code.** Our C++ MPHF library is publicly available at https://github.com/jermp/pthash.

**Organization.** The article is structured as follows. Section 2 reviews all practical approaches to MPH and indicates which techniques support multi-threading and/or use external memory. Section 3 presents a new general PTHash construction algorithm, with Section 4 highlighting how its design seamlessly adapts to the different settings we consider. In particular, Section 4.2 and 4.3 detail how the general construction is supported in internal and external memory, respectively, with multi-threading support (Section 4.1). Section 5 presents the experimental results. PTHash is compared against all the approaches reviewed in Section 2. We conclude in Section 6.

## 2  RELATED WORK

Minimum perfect hashing has a long history of development. Our focus is on practical approaches, that is, those that have been implemented and shown to perform well on very large key sets. In fact, we point out that some theoretical constructions, like that by Hagerup and Tholey [15], can be proved to reach the space lower bound of $n \log_2 e$ bits, but only work in the asymptotic sense, i.e., for $n$ too large to be of any practical interest.

Up to date, four "classes" of different, practical, approaches have been devised to solve the problem, which we describe below in chronological order of proposal. The upper part of Table 1 summarizes the notation used to describe the algorithms.

**Hash and Displace.** The *hash and displace* technique was originally introduced by Fox, Chen, and Heath [16] in a work that was named FCH after them. That work also inspired the development of PTHash [12]. The main idea of the method is as follows.

Keys are first hashed and mapped into *non-uniform* buckets; then, the buckets are sorted and processed by falling size: for each bucket, a displacement value $d_i \in [n]$ is determined so that all keys in the bucket can be placed without collisions to positions $(h(x) + d_i) \bmod n$, for a proper universal hash function $h$. Lastly, the sequence of displacements $d_i$ is stored in compact form using $\lceil \log_2 n \rceil$ bits per value. While the theoretical analysis suggests that by decreasing the number of buckets it is possible to lower the space usage (at the cost of a larger construction time), in practice it is unfeasible to go below 2.5 bits/key for large values of $n$.

TABLE 1
General notation (upper) and PTHash-specific notation (bottom).

| | |
|---|---|
| $S$ | a set of keys |
| $n$ | the number of keys in $S$, $n = |S|$ |
| $f$ | MPHF built from $S$ |
| $f(x)$ | a value in $[n] = \{0, \ldots, n-1\}$, the result of computing $f$ on the key $x$ |
| $h$ | hash function chosen at random from a universal familiy |
| $c$ | a value that trades search efficiency for space effectiveness |
| $\alpha$ | a value in $(0, 1]$, used to define the search space $[n' = n/\alpha]$ |
| $m$ | number of buckets used for the search, $m = \lceil cn/\log_2 n \rceil$ |
| $L$ | largest bucket size |
| $P$ | pilots table |
| $x \oplus y$ | bitwise XOR between hash codes $x$ and $y$ |

In the *compressed hash and displace* (CHD) variant by Belazzougui et al. [17], keys are first *uniformly* distributed to buckets, with expected size $\lambda > 0$. For each bucket, a pair of displacements $\langle d_0, d_1 \rangle$ is determined so that all keys in the bucket can be placed without collisions to positions $(h_1(x) + d_0 h_2(x) + d_1) \bmod n$, for a given pair of hash functions $h_1, h_2$. Instead of explicitly storing a pair $\langle d_0, d_1 \rangle$ for each bucket, the index of such pair in the sequence

$$\langle 0, 0 \rangle, \ldots, \langle 0, n-1 \rangle, \ldots, \langle n-1, 0 \rangle, \ldots, \langle n-1, n-1 \rangle$$

is stored. Lastly, the sequence of indexes is stored in compressed form using the entropy coding mechanism introduced by Fredriksson and Nikitin [18], retaining $O(1)$ access time.

**Linear Systems.** In the late 90s, Majewski et al. [19] introduced an algorithm to build a MPHF exploiting a connection between linear systems and hypergraphs (almost ten years later, Chazelle et al. [20] proposed an analogous construction in an independent manner). The MPHF $f$ is found by generating a system of $n$ random equations in $t$ variables of the form

$$w_{h_1(x)} + w_{h_2(x)} + \cdots + w_{h_r(x)} = f(x) \bmod n, x \in S,$$

where $h_i : S \to [t]$ is a random hash function, and $\{w_i\}$ are $t$ variables whose values are in $[n]$. Due to bounds on the acyclicity of random graphs, if the ratio between $t$ and $n$ is above a certain threshold $\gamma_r$, the system can be almost always triangulated and solved in linear time by peeling the corresponding hypergraph. The constant $\gamma_r$ depends on the degree $r$ of the graph, and attains its minimum for $r = 3$, i.e., $\gamma_3 \approx 1.23$.

Belazzougui et al. [14] proposed a cache-oblivious implementation of the previous algorithm suitable for external memory construction and named EMPHF. The algorithm uses a compact representation of the incidence lists of the hypergraph. The compact representation is based on the observation that, for the wanted operations on the lists, it is not necessary to store actual edges. Instead, all vertices in the same position can be XORed together, hence a constant amount of memory per node is retained.

Genuzio et al. [21, 22] (GOV) demonstrated practical improvements to the Gaussian elimination technique, which is used to solve the linear system. The improvements are based on broadword programming techniques. The authors

released a very efficient implementation of the algorithm that scales well using external memory and multi-threading.

**Fingerprinting.** Müller et al. [23] introduced a technique based on fingerprinting. The general idea is as follows. All keys are first hashed in $[n]$ using a random hash function and collisions are recorded using a bitmap $B_0$ of size $n_0 = n$. In particular, keys that do not collide have their position in the bitmap marked with a `1`; all positions involved in a collision are marked with a `0` instead. If $n_1 > 0$ collisions are produced, then the same process is repeated recursively for the $n_1$ colliding keys using a bitmap $B_1$ of size $n_1$. All bitmaps, called "levels", are then concatenated together in a bitmap $B$. The lookup algorithm keeps hashing the key level by level until a `1` is hit, say in position $p$ of $B$. A constant-time ranking data structure [24, 25] is used to count the number of `1`s in $B[0..p]$ to ensure that the returned value is in $[n]$. On average, only 1.56 levels are accessed in the most succinct setting [23], which takes 3 bits/key.

Limasset et al. [26] provided an implementation of this approach, named BBHash, that is very fast in construction and scales to very large key sets using multiple threads. Furthermore, no auxiliary data structures are needed during construction except the levels themselves, meaning that the space consumed in the process is essentially that of the final MPHF. Thus, the approach is also suitable for construction in external memory.

A parameter $\gamma \geq 1$ is introduced to speedup construction and lookup time, so that bitmap $B_i$ on level $i$ is $\gamma n_i$ bits large. This clearly reduces collisions and, thus, the average number of levels accessed per lookup. However, the larger $\gamma$, the higher the space consumption.

**Recursive Splitting.** Very recently, Esposito et al. [27] proposed a new technique, named RecSplit, for building very succinct functions. RecSplit is the most compact data structure up to date, achieving 1.80 bits/key on large key sets. Also, the construction runs in expected linear time and provides fast lookup time in expected constant time.

The authors first observed that for very small sets it is possible to find a MPHF simply by brute-force searching for a bijection with suitable codomain. Then, the same approach is applied in a divide-and-conquer manner to solve the problem on larger inputs.

Given two parameters $b$ and $\ell$, the keys are divided into buckets of average size $b$ using a random hash function. Each bucket is recursively split until a block of size $\ell$ is formed and for which a MPHF can be found using brute-force, hence forming a rooted tree of splittings. The parameters $b$ and $\ell$ provide different space/time trade-offs. While providing a very compact representation, the evaluation performs one memory access for each level of the tree that penalizes the lookup time.

The authors indicate that the approach is amenable to good parallelization as the buckets may be processed by multiple threads in parallel. However, there is no public parallel implementation of this algorithm yet.

## 3 A GENERAL CONSTRUCTION

In this section, we present a general construction algorithm for PTHash, based on few abstract steps, that is suitable for

```
1  build(S, c, α) :
2  │  blocks = map(S, c)
3  │  buckets = merge(blocks)
4  │  P = search(buckets, α)
5  │  encode(P)
```

Alg. 1. Construction for an input set $S$, with parameters $c$ and $\alpha$.

both parallel and external-memory settings. While the focus of this section is on the algorithm rather than on the specific implementation, Section 4 will detail how the proposed construction can be easily implemented to support multi-threading, in either internal or external memory. Refer to the bottom part of Table 1 for the notation used to describe PTHash.

In short, PTHash maps the keys into non-uniform buckets and processes the buckets by falling size. For each bucket $i$, it searches for $k_i$ – the *pilot* of bucket $i$ – an integer able of placing all the keys in the bucket to positions

$$f(x) = (h(x) \oplus h(k_i)) \bmod n$$

without collisions. Lastly, the collection of such pilots, one for each bucket, is materialized in a *pilots table* (PT) and indicated with $P$ in Table 1. This is what the MPHF data structure actually stores. In particular, the table is stored in *compressed* format using a compressor supporting random access to $P[i]$ in $O(1)$ worst-case time.

Let us now consider how the pilots table is computed. Besides the input set $S$, the construction uses two parameters, $c > \log_2 e$ and $0 < \alpha \leq 1$, respectively affecting the number of buckets and the size of the search space (see also the bottom part of Table 1). The general construction algorithm is composed of three steps, namely *map*, *merge*, and *search*, followed by an encoding step that compresses the table as anticipated before. These steps, explained in the following, are summarized in Alg. 1.

### 3.1 Map

The first step aims at mapping each key of $S$ to a fixed number of bytes with the help of an universal hash function $h$. During this step, keys are "logically" associated to one of the $m = \lceil cn / \log_2 n \rceil$ buckets that are indicated with the integers in $[0, m)$. Specifically, the bucket identifier for a key $x$ is obtained via the function

$$\text{bucket}(x) = \begin{cases} h'(x) \bmod p_2 & x \in S_1 \\ p_2 + (h'(x) \bmod (m - p_2)) & \text{otherwise} \end{cases},$$

where $S_1 = \{x | (h'(x) \bmod n) < p_1\}$ and $h'$ is, like $h$, a random hash function. The thresholds $p_1$ and $p_2$ are set to, respectively, $0.6n$ and $0.3m$ so that the mapping of keys into buckets is *skewed*: roughly 60% of the keys are mapped into 30% of the buckets. We point out that the skew distribution of keys into buckets is very important for the *efficiency* of the search step (and the final compression effectiveness), but we defer the reader to our previous paper [12] for a discussion about this fact.

For each key $x \in S$ we generate a pair of the form

$$\langle id, hash \rangle = \langle \text{bucket}(x), h(x) \rangle$$

```
1   search(buckets, α) :
2   │   n' = n/α                                ▷ search space
3   │   m = |buckets|                           ▷ number of buckets
4   │   L = |buckets[0].hashes|                 ▷ largest bucket size
5   │   P = ∅
6   │   positions = ∅
7   │   allocate the bitmap taken[0..n' − 1] with all 0s
8   │   for i = 0; i < m; i = i + 1 :
9   │   │   bucket = buckets[i]
10  │   │   k_i = 0                             ▷ pilot for i-th bucket
11  │   │   while true :
12  │   │   │   clear positions
13  │   │   │   j = 0
14  │   │   │   for ; j < |bucket.hashes|; j = j + 1 :
15  │   │   │   │   p = (bucket.hashes[j] ⊕ h(k_i)) mod n'
16  │   │   │   │   if taken[p] = 1 :
17  │   │   │   │   │   k_i = k_i + 1
18  │   │   │   │   └   break                   ▷ try next pilot
19  │   │   │   └   add p to positions
20  │   │   │   if j = |bucket.hashes| :
21  │   │   │   │   if positions contains duplicates :
22  │   │   │   │   │   k_i = k_i + 1
23  │   │   │   │   └   continue                ▷ try next pilot
24  │   │   │   │   add ⟨bucket.id, k_i⟩ to P   ▷ save pilot
25  │   │   │   │   for all p in positions : taken[p] = 1
26  │   │   └   └   break
27  │   sort(P)
28  └   return P
```

Alg. 2. Search procedure for an input collection of *buckets*, with parameter $\alpha$.

and divide the $n$ pairs into $K$ *blocks*, each of (roughly) equal size. The pairs in each block are sorted by, first, the *id* component, then by *hash*. Blocks are formed either because: keys are evenly distributed among $K$ threads and processed in parallel in internal memory, or a block of pairs is flushed to disk when internal memory is exhausted. We will better explain these two scenarios in Section 4.2 and 4.3.

The map step takes $O(n \log(n/K))$ expected time using quick-sort and assuming the computation of $h(x)$ and $h'(x)$ to take constant time. If we use $u$-bit hash codes, the amount of temporary space taken by the map step is $n(\lceil \log_2 m \rceil + u)$ bits.

### 3.2 Merge

The $K$ blocks output by the previous step are then merged to group together all the hash codes for a given bucket. Specifically, the pairs are merged by their *id* component only. In this way we are also able to check that all keys were correctly hashed to distinct hash codes[1]. Moreover, since pairs having the same *id* are also sorted by *hash*, checking for duplicates is efficient (i.e., scanning of the pairs).

1. In case of a duplicate hash code, any MPHF construction must fail. However, the failing probability is very low and depends only on $n$ and $u$. Indeed, we never saw a single collision during our experiments using $u = 64$ bits.

During this step we also accomplish to sort the buckets by non-increasing size, as follows[2]. We allocate $L$ buffers, where $L$ is the largest bucket size. Suppose the merge for a bucket is complete and the bucket contains $k$ hash codes. We append its identifier and the $k$ hash codes to the $k$-th buffer. In this way we have all the buckets of the same size written contiguously.

Each bucket in the $k$-th buffer is a contiguous span of $k + 1$ integers, for $0 \leq k < L$. Therefore, the total space taken by the $L$ buffers is $m\lceil \log_2 m \rceil + nu + O(L) \simeq m\lceil \log_2 m \rceil + nu$ bits since $L$ is very small compared to $n$ (e.g., $L < 40$ in our experiments with billions of keys). Once the merge is complete for all buckets, the input blocks of pairs are destroyed.

The overall complexity is that of the preliminary merge phase which is carried out using a min-heap data structure of size $O(K)$ memory words, hence $O(n \log K)$ time.

The output of this step is the set of such $L$ buffers, so that the next step, the search, has to logically process the buffers from index $L - 1$ down to index 0. However, we indicate the $L$ buffers with the abstraction *buckets* in Alg. 1. This abstract collection *buckets* must only support *sequential iteration* in the wanted bucket order (from largest to smallest bucket size). Sequential iteration is an access pattern of crucial importance as to avoid expensive memory reads for the buckets (especially in the external-memory setting). We are going to assume that each element of the collection is a *bucket* – an object made by an array of *hashes*, plus a unique identifier $0 \leq id < m$.

### 3.3 Search

The search procedure is illustrated in Alg. 2 and it is the core of the whole construction. The procedure keeps track of occupied positions in the search space $[n' = n/\alpha] = \{0, \ldots, n' - 1\}$ using a bitmap, $taken[0..n' - 1]$. For each $bucket = buckets[i]$ we search for an integer $k_i$ such that the position assigned to $h(x) \in bucket.hashes$ is $p = (h(x) \oplus h(k_i)) \mod n'$ and $taken[p] = 0$. As anticipated, the integer $k_i$ is called the pilot for *bucket*. If the search for $k_i$ is successful, i.e., $k_i$ places all hashes of *bucket* into unused positions of *taken* (and there are no hashes in *bucket* that are mapped to the same slot of *taken*, i.e., the **if** in line 21 fails), then $k_i$ is saved in the pilots table $P$ (line 24) and the positions are marked as occupied via $taken[p] = 1$ (line 25); otherwise, the next integer $k_i + 1$ is tried.

At this stage of the description, $P$ is modelled as an abstract collection of pairs $\langle id, pilot \rangle$. The collection is then sorted by *id* to materialize the final pilots table (line 27). In this way, the code of the search can adapt to either work in internal or external memory.

Both positions $p$ (line 15) and pilots $k_i$ are 64-bit integers. Therefore, the search consumes $n'$ bits for the bitmap *taken*, plus $64L$ bits for the array *positions* because at most $L$ positions can be added to it, plus $m(\lceil \log_2 m \rceil + 2 \times 64)$ bits for $P$.

2. The steps *map* and *merge* have the same role of the steps called *mapping* and *ordering* in the original MOS framework introduced by FCH [16]. However, in our construction, the ordering of the buckets is carried out during both steps: sorting of the pairs during *map* and displacing of the pairs into the $L$ buffers during *merge*.

```
1  P[i] :
2  |  block = ⌊i/b⌋
3  |  offset = i mod b
4  |  w = W[block + 1] − W[block]
5  |  position = W[block] × b + offset × w
6  |  x = B[position, position + w]
7  |  return x
```

Alg. 3. Algorithm for retrieving the $i$-th element of the pilots table when represented in partitioned compact (PC) form. The notation $x = B[a, b]$, $a < b$, means that $b − a$ bits are read from the bitmap $B$ starting at position $a$ and written to an integer value $x$.

Lastly, we give the following theorem which relates the performance of the search, in terms of expected CPU time, to the parameters $c$ and $\alpha$.

**Theorem 1.** The expected time of the search, for $n$ keys and parameters $c > \log_2 e$ and $0 < \alpha \le 1$, is $O(n^{1+\Theta(\alpha/c)})$.

*Proof* – Given in the Appendix.

The search space is of size $n' = n/\alpha$ which is actually larger than $n$ if $\alpha < 1$. This makes the search for pilots faster by lowering the probability of hitting already taken positions in the bitmap. However, the output of $f$ must be guaranteed to be minimal, i.e., a value in $[n]$, not in $[n']$. We therefore need a mechanism to assign the slots left unused by $f$ in its codomain $[n]$. Suppose $F$ is the list of unused positions *up to*, and including, position $n − 1$. Then there are $|F|$ keys that are mapped to positions $p_i \ge n$, that can fill such unused positions. We materialize an array $free[0..n' − n − 1]$, where $free[p_i − n] = F[i]$, for each $i = 0, \ldots, |F| − 1$. To compute the (uncompressed) *free* array we need $\Theta(n)$ time and $64(n' − n)$ bits.

### 3.4 Encode

The PTHash data structure stores the two compressed tables, $P$ and *free*. For the *free* array, we always use Elias-Fano [28, 29], noting that it only takes $(n' − n)(\lceil \log_2 \frac{n}{n'−n} \rceil + 2 + o(1))$ bits. The pilots table $P$, instead, can be compressed using *any* compressor for integer sequences that supports constant-time random access to its $i$-th integer $P[i]$. (It is also desirable that compressing $P$ runs in linear time, i.e., $\Theta(m)$ time. We only consider compressors with this complexity in the article.)

Therefore, we have three degrees of freedom for the tuning of PTHash, namely the choice of

- the compressor for $P$,
- the size of the search space, tuned with $\alpha$, and
- the number of buckets, tuned with $c$,

that allow to obtain different space/time trade-offs. We point the interested reader to our previous paper [12] for an overview and discussion of the achievable trade-offs. For example, a good balance between space effectiveness and lookup efficiency is obtained using the front-back dictionary-based encoding (D-D) [12], with $\alpha = 0.94$ and $c = 7.0$. More compact representations can be obtained, instead, by using Elias-Fano (EF) on the prefix-sums of $P$ at the price of a slowdown at query time. We are going to use both configurations in the experiments as reference points.

Here we propose another compressed representation that supports constant-time random access, and named *partitioned compact* (PC) encoding hereafter. The pilots table is divided into $\lceil m/b \rceil$ blocks of size $b$ each (except possibly the last one; in our implementation we use $b = 256$). For each block, we compute the minimum bit-width $w_i$ necessary to represent its maximum element, i.e., $w_i = \lceil \log_2(max + 1) \rceil$ (if $max = 0$, then $w_i$ is set to 1), so that every integer in the block can be represented using $w_i$ bits. The whole pilots table is represented by concatenating the representations of all blocks in a bitmap $B$ that takes $O(1) + b \sum_{i=0}^{\lceil m/b \rceil − 1} w_i$ bits. To support random access we also materialize in a separate array the sequence

$$W = [0, W(0), W(1), W(2), \ldots, W(\lceil m/b \rceil − 2)]$$

where $W(j) = \sum_{i=0}^{j} w_i$. The algorithm for retrieving $P[i]$ is illustrated in Alg. 3: it requires two memory accesses (one to $W$, the other to $B$) and no branches, thus indicating that fast evaluation is retained. The space effectiveness of the encoding is expected to be good for the same reasons why front-back compression works well [12, Section 4.2]: the entropy of $P$ is lower at the front and higher at the back, with a smooth transition between the two regions. This means that the bit-widths $\{w_i\}$ tend to generally increase when moving from $i = 0$ to $i = \lceil m/b \rceil − 1$. Indeed, the PC encoding can be considered as a generalization of front-back compression.

## 4 IMPLEMENTATION DETAILS

In this section, we detail concrete implementations of the PTHash construction from Section 3. These implementations are meant to deliver good practical performance, hence, exploit multi-threading and external-memory scaling. To this end, we first present in Section 4.1 a fundamental building block – a *parallel search* procedure. In Section 4.2 we assume that the whole input set $S$ can be processed in internal memory. In Section 4.3, instead, we aim at building the PTHash data structure for very large sets that cannot be processed in internal memory, thus it is necessary to use temporary disk storage.

For all constructions, the pairs $\langle id, hash \rangle$, defined during the map step in Section 3.1, consume an integral number of bytes, $q$, e.g., $q = 12$ if $id$ is a 4-byte identifier (allowing up to $2^{32}$ buckets) and $hash$ is a 8-byte long hash code.

### 4.1 Parallel Search

As already noted in Section 3, the search step is the core of the whole PTHash construction and the most time-consuming step. Therefore we would like to exploit all the target machine cores to search pilots more efficiently. But parallelizing the algorithm in Alg. 2 presents some serious limitations in that we cannot compute pilots in parallel, say, $K$ pilots for $K$ different buckets, because the displacement of keys in the $i$-th bucket depends on the displacement of *all* previous buckets.

However, *while a thread may not finalize the computation of a pilot, it can try some pilots anyway given the current bit-configuration of the bitmap and, thus, potentially discard many pilots that surely cannot work because of the occupied slots.* Fig. 1

buckets

already processed

computation

| | | |
|---|---|---|
| **thread 2** | 5690 | |
| thread 3 | 5691 | |
| thread 4 | 5692 | |
| thread 5 | 5693 | |
| thread 6 | 5694 | |
| thread 7 | 5695 | |
| thread 0 | 5696 | |
| thread 1 | 5697 | |

to be processed

Fig. 1. A graphical representation of the parallel search procedure with 8 threads, when processing buckets from index 5690 to index 5697. The bars to the right of the buckets represent the amount of computation that threads have to do before computing the first successful pilot. The *shaded part* of the bars corresponds to the work done for *unsuccessful* pilots (those that create collisions in the bitmap). In this example, thread 2 (in bold font) has to commit its work, therefore it can execute until success (the shaded part covers its entire bar). The crucial point is that: *the shaded part can be computed in parallel*, hence saving time compared to the sequential implementation.

illustrates a pictorial representation of this idea. Based on this idea, we proceed as follows.

We spawn $K$ parallel threads, working on $K$ consecutive buckets, one thread per bucket. Each thread advances in the search of a pilot independently, pausing the search when a pilot that works for the *current state* of the bitmap is found. At that point, either the thread is processing the first bucket to commit, or it is waiting for some other thread to update the bitmap. In the former case, the thread (i) concludes the pilot search, (ii) updates the bitmap, and (iii) starts processing the next bucket. Only one thread at a time can thus update the bitmap. In the latter case, the waiting thread wakes up after a bitmap update and continues the pilot search according to the new state of the bitmap, pausing as soon as it finds a new working pilot.

The crucial point, now, is to manage the turns between the threads in an efficient way. To do so, we label the threads with unique identifiers, from 0 to $K - 1$, and maintain the following invariant: thread $i$ processes the bucket of index $0 \le j < m$ where $i = j \bmod K$ (see Fig. 1 for an example with $K = 8$). *It follows that the threads must finalize their respective computation and commit changes to the bitmap following the identifier order to guarantee correctness.* We guarantee this ordering using a *shared* identifier, periodically indicating the thread that is allowed to commit. Therefore, a thread advances the search as far as it can and checks (inside a loop, sometimes called "busy waiting") the shared identifier.

If the shared identifier is equal to its own identifier, then the thread can commit his work. Note that we do *not* require a lock/unlock mechanism to synchronize the threads that would sensibly erode the benefit of parallelism, but rely on the *value* assumed by the shared identifier. This is a rather important point to obtain good practical performance.

At the beginning of the algorithm, the shared identifier is equal to 0, hence only the first thread is allowed to conclude the search and update the bitmap. Then the shared identifier is set to 1 by the first thread. The second thread can now commit, and so on. In general, the turn of the thread $i$ comes when all threads $0 \le j < i$ have committed (when the thread $K - 1$ commits, the shared identifier is set back to 0).

### 4.2 Internal Memory

Let us assume that the construction algorithm has enough internal memory available. Then we have to specify only few details given the flexibility of the algorithm in Section 3.

- The blocks output by the map step are materializd as in-memory vectors of pairs. Let $K$ be the number of parallel threads to use. The map step spawns $K$ threads. Each thread allocates a vector of $n/K$ pairs and fills it by forming pairs from a partition of $S$ consisting in $n/K$ keys. Once the vector is filled up, it is sorted. Since the threads work on different partitions of $S$, the map with $K$ threads achieves nearly ideal speed-up.

- The collection of buckets output by the merge step consists in $L$ in-memory vectors of hash codes. There is no meaningful opportunity for parallelization during the merge step.

- The auxiliary arrays used by the search step, i.e., *taken* and *positions* are materialized as in-memory arrays of 64-bit integers. The pilots table $P$ in the pseudocode of Alg. 2 is also implemented with a 64-bit integer array. Since $P$ fits in internal memory, the saving of a pilot in line 24 simplifies to $P[bucket.id] = k_i$ so that $P$ is already sorted and the sort($P$) step in line 27 is not performed at all. The search can be executed in parallel using $K$ threads as we explained in Section 4.1.

### 4.3 External Memory

When not enough internal memory is available for the construction, e.g., because the size of $S$ is very large, it is mandatory to resort to external memory. The general algorithm described in Section 3 easily adapts to this scenario because the steps of the algorithm do *not* change, rather the output of map and merge is written to (resp. read from) disk. Let $M$ indicate the maximum amount of internal memory (in bytes) that the construction is allowed to use, and $q$ be the number of bytes taken by a $\langle id, hash \rangle$ pair.

- During the map step, we allocate an in-memory vector of size $M/q$ pairs. Whenever the vector fills up, it is sorted (possibly in parallel) and flushed to a new file on disk, as to guarantee that at most $M$ bytes of internal memory are used during the process.

- Let $K$ be the number of files created during the map step, i.e., $K = \lceil (qn)/M \rceil$. The pairs from these $K$ files are merged together to create the collection of $L$ buffers representing the buckets. Again, the $L$ buffers

are formed in internal memory without violating the limit of $M$ bytes. Whenever the limit is reached, the content of the buffers is accumulated to $L$ files on disk.

- The buckets are read from the $L$ files and processed by the search procedure, possibly by its parallel implementation described in Section 4.1.

We point out that the bitmap *taken* is *always* kept in internal memory (and the array *positions* as well), because its access pattern is *random* and the search would be slowed down to unacceptable rates if the bitmap were resident on disk.

Therefore, taking into account that the seach needs $n' = n/\alpha$ bits of internal memory, we are left with $M' = M - n'/8$ bytes available to store the pilots. The pilots are accumulated in an in-memory vector of $M'/q$ pairs $\langle id, pilot \rangle$ (also these pairs consume $q$ bytes each as those used during the map step). Whenever the $M$-byte limit is reached, the vector is flushed to disk and emptied. Lastly, the different files containing the pilots are merged together to obtain the final $P$ table – the sort$(P)$ step in line 27 of the algorithm in Alg. 2.

## 5 EXPERIMENTAL RESULTS

In this section of the article, we report on the results of an extensive experimental analysis that, as mentioned in Section 1, considers the four possible settings to run the construction:

- internal-memory/sequential,
- internal-memory/parallel,
- external-memory/sequential,
- external-memory/parallel.

**Hardware.** We use a server machine equipped with 8 Intel i9-9900K cores (@3.60 GHz), 64 GB of RAM DDR3 (@2.66 GHz), and running Linux 5 (64 bits).

For all experiments where parallelism is enabled *we use 8 parallel threads*, one thread per core.

Each core has two private levels of cache memory: 32 KiB L1 cache (one for instructions and one for data); 256 KiB for L2 cache. A shared L3 cache spans 16 MiB. All cache levels have a line size of 64 bytes. All datasets are read from a Western Digital Red mechanical disk with 4 TB of storage and a rotation rate of 5400 rpm (SATA 3.1, 6 GB/s).

**Software.** The implementation of PTHash is written in C++ and available at https://github.com/jermp/pthash. We use native C++ threads to support parallel execution, i.e., without relying on other frameworks such as Intel's TBBs or OpenMP that would otherwise limit the portability of our software. For the experiments reported in the article, the code was compiled with gcc 9.2.1 using the flags: `-std=c++17 -pthread -O3 -march=native`.

**Methodology.** Construction time is reported in total seconds, taking the average between 3 runs. Lookup time is measured by looking up every single key in the input using a single core of the processor, and taking the average time between 5 runs. For inputs residing on disk, we load batches of 8 GB of keys in internal memory and measure lookup time on each batch. The final reported time is the average among the measurements on all batches. Lookup time is

TABLE 2
String collections used in the experiments.

| Collection | num. strings | avg. bytes/string |
|---|---|---|
| UK2005 URLs | 39 459 925 | 71.37 |
| ClueWeb09-B URLs | 49 937 704 | 54.72 |
| GoogleBooks 2-gr | 665 752 080 | 17.21 |
| TweetsKB IDs | 1 983 291 944 | 18.75 |
| ClueWeb09-Full URLs | 4 780 950 911 | 67.28 |
| GoogleBooks 3-gr | 7 384 478 110 | 23.30 |

TABLE 3
Methods and implementations.

| Method | Implementation | |
|---|---|---|
| | multi-threading | external memory |
| PTHash | ✓ | ✓ |
| PTHash-HEM | ✓ | ✓ |
| FCH [16] | | |
| CHD [17] | | |
| EMPHF [14] | | ✓ |
| EMPHF-HEM [14] | | ✓ |
| GOV [21, 22] | ✓ | ✓ |
| BBHash [26] | ✓ | ✓ |
| RecSplit [27] | | |

reported in nanoseconds per key (ns/key). The space of the MPHF data structures is reported in bits per key (bits/key).

A testing detail of particular importance is that, before running a construction algorithm, *the disk cache is cleared* to ensure that the whole input dataset is read from the disk.

**Datasets.** We use some real-world string collections as input datasets. Table 2 reports the basic statistics for these collections. All datasets are publicly available for download by following the corresponding link in the References.

We use natural-language $q$-grams as they are in widespread use in IR, NLP, and Machine-Learning applications. Specifically, we use the 2-3-grams from the English GoogleBooks corpus, version 2 [30], as also used in prior work on the problem [14, 26]. URLs are interesting as they represent a sort of "worst-case" input given their very long average length. We use those of the Web pages in the ClueWeb09 dataset [31] (category B and full dataset), and those collected in 2005 from the UbiCrawler [32] relative to the .uk domain [33]. TweetsKB [34] is a collection of unique tweets identifiers (IDs), corresponding to tweets collected from February 2013 to December 2020.

As TweetsKB IDs, ClueWeb09-Full URLs, and GoogleBooks 3-gr do not fit in the internal memory of our test machine (64 GB), we are going to use these three datasets to benchmark the construction algorithms in external memory. The other three datasets can be processed in internal memory.

**Methods and Implementations.** PTHash is compared against the state-of-the-art methods reviewed in Section 2. Table 3 indicates what methods support parallel execution and external-memory scaling. The "HEM" suffix used in the tables stands for *heuristic external memory* [13] and refers to the approach of partitioning the input and building an independent MPHF on each partition.

Whenever possible, we used the implementations avail-

TABLE 4
*Internal-memory* construction time, space, and lookup time, for a range of methods.
Numbers in parentheses refer to the parallel construction using 8 threads.
All PTHash configurations use $\alpha = 0.94$ and $c = 7.0$.

| Method | UK2005 URLs | | | ClueWeb09-B URLs | | | GoogleBooks 2-gr | | |
|---|---|---|---|---|---|---|---|---|---|
| | construction (seconds) | space (bits/key) | lookup (ns/key) | construction (seconds) | space (bits/key) | lookup (ns/key) | construction (seconds) | space (bits/key) | lookup (ns/key) |
| PTHash (D-D) | 9 (5) | 3.11 | 46 | 12 (6) | 3.07 | 49 | 410 (156) | 2.97 | 63 |
| PTHash (PC) | 9 (5) | 2.82 | 49 | 12 (6) | 2.80 | 52 | 407 (153) | 2.67 | 69 |
| PTHash (EF) | 9 (5) | 2.50 | 59 | 12 (6) | 2.49 | 63 | 408 (154) | 2.38 | 121 |
| PTHash-HEM (D-D) | 8 (2) | 3.11 | 52 | 10 (2) | 3.08 | 57 | 167 (34) | 2.97 | 72 |
| PTHash-HEM (PC) | 8 (2) | 2.82 | 54 | 10 (2) | 2.80 | 59 | 164 (34) | 2.67 | 80 |
| PTHash-HEM (EF) | 8 (2) | 2.50 | 66 | 10 (2) | 2.49 | 70 | 164 (34) | 2.38 | 159 |
| FCH ($c = 3.0$) | 1138 | 3.00 | 52 | 1438 | 3.00 | 55 | – | – | – |
| FCH ($c = 4.0$) | 266 | 4.00 | 60 | 267 | 4.00 | 64 | 9110 | 4.00 | 54 |
| FCH ($c = 5.0$) | 119 | 5.00 | 68 | 107 | 5.00 | 71 | 3225 | 5.00 | 55 |
| CHD ($\lambda = 4$) | 32 | 2.17 | 170 | 43 | 2.17 | 185 | 1251 | 2.17 | 410 |
| CHD ($\lambda = 5$) | 84 | 2.07 | 173 | 115 | 2.07 | 182 | 3923 | 2.07 | 410 |
| CHD ($\lambda = 6$) | 306 | 2.01 | 173 | 429 | 2.01 | 178 | 15583 | 2.01 | 406 |
| EMPHF | 10 | 2.61 | 126 | 13 | 2.61 | 135 | 276 | 2.61 | 211 |
| EMPHF-HEM | 9 | 3.49 | 152 | 11 | 3.30 | 158 | 148 | 3.44 | 287 |
| GOV | 39 (14) | 2.23 | 87 | 47 (15) | 2.23 | 94 | 613 (170) | 2.23 | 170 |
| BBHash ($\gamma = 1.0$) | 50 (12) | 3.10 | 154 | 66 (14) | 3.06 | 174 | 1248 (189) | 3.08 | 273 |
| BBHash ($\gamma = 2.0$) | 31 (7) | 3.71 | 133 | 40 (9) | 3.71 | 150 | 666 (102) | 3.71 | 204 |
| RecSplit ($\ell = 5, b = 5$) | 6 | 2.95 | 153 | 8 | 2.95 | 148 | 132 | 2.95 | 244 |
| RecSplit ($\ell = 8, b = 100$) | 37 | 1.79 | 113 | 47 | 1.79 | 118 | 724 | 1.80 | 202 |
| RecSplit ($\ell = 12, b = 9$) | 225 | 2.23 | 93 | 284 | 2.23 | 107 | 3789 | 2.23 | 225 |

able from the original authors. We include a link to the source code in the References section. All implementations are in C/C++, except for the construction of GOV which is only available in Java (but lookup time is measured using a C program). Moreover, we also tested the algorithms *using the same parameters as suggested by their respective authors* to offer different trade-offs between construction time and space effectiveness. Below we report some details.

FCH is the only algorithm that we re-implemented (in C++) faithfully to the original paper. For the CHD method we were unable to use $\lambda = 7$ for more than a few thousand keys. The EMPHF library also includes the corresponding HEM implementation of the algorithm, EMPHF-HEM. The authors of BBHash also considers $\gamma = 5.0$ in their own work but we obtained a space larger than 6.8 bits/key, so we excluded this value of $\gamma$ from the analysis.

### 5.1 Internal Memory

Table 4 shows the performance of the methods on the datasets that can be processed entirely in internal memory. The numbers in parentheses refer to the parallel construction using 8 threads; moreover, all PTHash configurations use $\alpha = 0.94$ and $c = 7.0$.

Let us first consider the sequential construction and formulate some general observations. PTHash is the fastest MPHF data structure at evaluation time. FCH offers similar lookup performance but PTHash is more space-efficient and much faster at construction time (FCH with $c = 3.0$ takes too much time to run over the GoogleBooks 2-gr dataset). Moreover, PTHash offers good space effectiveness, albeit not the best: around 3.2 bits/key using the D-D encoding but less using PC, i.e., around 2.7 – 2.8 bits/key, and even

less using Elias-Fano (EF), i.e., around 2.4 – 2.5 bits/key. For methods achieving a similar space, such as BBHash, PTHash is 2× faster at lookup time, and even better at construction time. Compared to more space-efficient methods, such as RecSplit or CHD, PTHash is 2 – 6× faster at lookup time with better construction time.

For the HEM implementation of PTHash we fix the average partition size to $b = 5 \times 10^6$ and let the number of partitions be approximately $n/b$. Therefore, we use 8, 16, and 128 partitions for UK2005 URLs, ClueWeb09-B URLs, and GoogleBooks 2-gr respectively. It is worth noting that PTHash-HEM is faster to build than PTHash, and the gap increases by increasing $n$. For example, building 128 partitions takes 170 seconds on the GoogleBooks 2-gr dataset instead of 410 seconds, hence 2.45× faster. This is a direct consequence of the complexity of the search that is *not* linear in $n$ as per Theorem 1. Thus, building $n/b$ functions each for $\approx b$ keys is faster than building a *single* function for $n$ keys.

Another meaningful point to observe is that PTHash-HEM imposes a penalty at lookup time of 10 – 16% with respect to PTHash with no space overhead. As a comparison, EMPHF-HEM imposes a penalty of 17 – 36% at lookup time and of 26 – 33% of space usage compared to EMPHF. It is desirable that the space of the HEM representation with $r$ partitions is very similar to that of the un-partitioned data structure built from the same input with $m = \lceil cn/\log_2 n \rceil$ buckets. To this end, we search for each MPHF partition using exactly $\lfloor m/r \rfloor$ buckets so to guarantee that the total number of buckets used by the HEM data structure is $m$. In fact, note that using $\lfloor m/r \rfloor$ buckets for the $i$-th partition is different than letting the search over $n_i$ keys to use $\lceil cn_i/\log_2 n_i \rceil$ buckets because of the non-

Fig. 2. Internal-memory construction time of PTHash and PTHash-HEM on the GoogleBooks 2-gr dataset, by varying $c$, and for two values of $\alpha$. The parallel construction uses 8 threads. The horizontal line marks the lowest construction time achieved by sequential PTHash, that is for $c = 7.0$. Note the log-scale of the y-axis.



Fig. 3. Space/time trade-offs for PTHash and PTHash-HEM by varying $c$ and compressor (D-D, PC, EF), and for two values of $\alpha$. Performance was measured on the GoogleBooks 2-gr dataset.

linearity of the $\log_2$ function. In particular, it will be that $\lfloor m/r \rfloor < \lceil cn_i / \log_2 n_i \rceil$ assuming an average partition size $n_i$ close to $n/r$.

We now discuss the results for parallel constructions. Since a parallel construction has the clear advantage of reducing construction time, it can be used to either:

- build functions more quickly, for a *fixed space budget*;
- build more compact functions by selecting tighter parameters $c$ and $\alpha$ for the search, for a *fixed construction-time budget*.

The numbers in parentheses in Table 4 show that construction time improves significantly. The parallel algorithm reduces construction time by $2 - 2.5\times$, thus allowing PTHash to build faster than its direct competitors that also exploit multi-threading, like GOV and BBHash. Better speedup

can be achieved, not surprisingly, by PTHash-HEM because partitions can be built independently in parallel.

Fig. 2 highlights, instead, that parallelism can also be used to build more space-efficient functions. In the figure, the horizontal line marks the lowest construction time achieved by sequential PTHash (for $c = 7.0$), so that it is easy to see the *same* time is spent by the parallel PTHash for a *smaller c*. For example, $c$ can be reduced from 7.0 to 5.0 for $\alpha = 0.94$, and from 7.0 to 5.2 for $\alpha = 0.99$, to build a smaller MPHF within the same time-budget using 8 parallel threads. The companion Fig. 3 shows the space consumption and lookup time corresponding to such values of $c$. Continuing the same example, for the D-D encoder, space improves from 2.97 to 2.63 bits/key for $\alpha = 0.94$, and from 2.89 to 2.52 bits/key for $\alpha = 0.99$.

TABLE 5
*External-memory* construction time, space, and lookup time, for a range of methods.
Numbers in parentheses refer to the parallel construction using 8 threads.
All PTHash configurations use $\alpha = 0.94$ and $c = 7.0$.

| Method | TweetsKB IDs | | | ClueWeb09-Full URLs | | | GoogleBooks 3-gr | | |
|---|---|---|---|---|---|---|---|---|---|
| | construction (seconds) | space (bits/key) | lookup (ns/key) | construction (seconds) | space (bits/key) | lookup (ns/key) | construction (seconds) | space (bits/key) | lookup (ns/key) |
| PTHash (D-D) | 2105 (1058) | 3.07 | 80 | 7234 (4869) | 2.96 | 120 | 9770 (5865) | 2.91 | 91 |
| PTHash (PC) | 2098 (1051) | 2.61 | 101 | 7161 (4859) | 2.58 | 175 | 9756 (5736) | 2.56 | 143 |
| PTHash (EF) | 2098 (1052) | 2.36 | 188 | 7225 (4788) | 2.32 | 214 | 9649 (5849) | 2.31 | 208 |
| PTHash-HEM (D-D) | 989 (581) | 2.85 | 116 | 4651 (3632) | 2.75 | 152 | 5215 (3510) | 2.71 | 135 |
| PTHash-HEM (PC) | 982 (580) | 2.61 | 128 | 4522 (3541) | 2.58 | 192 | 5015 (3366) | 2.57 | 190 |
| PTHash-HEM (EF) | 983 (580) | 2.36 | 201 | 4627 (3631) | 2.32 | 235 | 5179 (3512) | 2.31 | 230 |
| EMPHF | 4021 | 2.61 | 207 | 24862 | 2.61 | 231 | 37731 | 2.61 | 220 |
| EMPHF-HEM | 909 | 3.03 | 279 | 3980 | 3.31 | 304 | 5606 | 3.06 | 304 |
| GOV | 2678 (1463) | 2.23 | 192 | 8228 (5400) | 2.23 | 232 | 10782 (6461) | 2.23 | 242 |
| BBHash ($\gamma = 1.0$) | 4750 (931) | 3.07 | 294 | 19360 (18391) | 3.07 | 320 | 20178 (9554) | 3.07 | 305 |
| BBHash ($\gamma = 2.0$) | 2365 (590) | 3.71 | 217 | 11074 (10348) | 3.71 | 236 | 10254 (5404) | 3.71 | 235 |

Fig. 3 also shows that the partitioned compact encoding (PC) introduced in Section 3 is always more effective than D-D with only a minor lookup penalty. PC may be preferable to EF for its better lookup performance despite the latter being still more space-efficient. We conclude with two more illustrative examples.

- Compared to RecSplit ($\ell = 12, b = 9$), parallel PTHash with $\alpha = 0.99$ and $c = 4.0$ builds faster (3789 vs. 2710 seconds) and, using PC, it retains even better space, i.e., 2.12 bits/key, with $3.75\times$ faster lookup time.
- Elias-Fano (EF) allows PTHash to break the 2 bits/key "barrier", indeed consuming only 1.98 bits/key for $\alpha = 0.99$ and $c = 4.0$: we point out that, up to date, only RecSplit was able to do so. Using that configuration of parameters, PTHash-HEM is only slightly larger than RecSplit ($\ell = 8, b = 100$), that is, 1.98 vs. 1.80 bits/key, but builds nearly $2\times$ faster with $1.5\times$ better lookup time.

## 5.2 External Memory

Table 5 shows the performance of the methods that support external-memory scaling, on datasets comprising several billions of strings. For PTHash and PTHash-HEM we fixed the amount of internal memory to 16 GB (out of the 64 GB available) as we did not observe any appreciable difference by using less or more memory. The number of partitions used by PTHash-HEM is set using the same methodology of Section 5.1, i.e., we fix the average partition size to $b = 5 \times 10^6$ and let the number of partitions be approximately $n/b$. Therefore, we use 400, 1000, and 1500 partitions for TweetsKB IDs, ClueWeb09-Full URLs, and GoogleBooks 3-gr respectively.

Results for external memory are *consistent* with those discussed for internal memory: the construction time of PTHash is competitive or better than that of other methods, for similar space effectiveness but better lookup time.

We point out, again, that the parallel search algorithm described in Section 4.1 is a key ingredient to achieve efficient parallel construction, even in external memory. PTHash-HEM offers further reductions in construction time

at the cost of a lookup penalty. However, note that the parallel speedup is partially eroded in external memory compared to internal memory, due to I/O operations that are merely sequential.

Considering Table 5, we now illustrate some comparisons using the largest dataset GoogleBooks 3-gr ($\approx 7.4$ billion strings) as example, noting that very similar considerations are valid for the other datasets:

- PTHash with D-D is $2\times$ faster than BBHash with $\gamma = 1.0$ at construction, retains even better space usage, and $3\times$ faster lookups. Compared to GOV, it achieves similar or better construction time with $2.7\times$ faster lookups, albeit being less space-efficient (23% larger).
- PTHash-HEM with EF is practically as efficient as GOV regarding space and lookup time, but with almost $2\times$ better construction time.
- PTHash with PC is nearly $4\times$ faster at construction than EMPHF (using the sequential algorithm), retains the same space and $1.5\times$ better lookup. (The high construction time of EMPHF is due to excessive I/O operations.) The HEM implementation of PTHash outperforms the HEM implementation of EMPHF under every aspect.

## 6 CONCLUSIONS

In this work we described a construction algorithm for PTHash that enables (the joint use of) multi-threading and external-memory scaling. These two features are very important to scale to large datasets in reasonable time. Our C++ implementation is publicly available.

We presented an extensive experimental analysis and comparisons, using large real-world string collections. There are three efficiency aspects for minimum perfect hash data structures: construction time, space consumption, and lookup time. *PTHash can be tuned to match the performance of another method on one of the three aspects and, then, outperform the same method on the other two aspects.* We remark that lookup time may be the most relevant aspect in concrete applications of minimal perfect hashing. In this regard, PTHash offers very fast evaluation, e.g., $2 - 6 \times$ faster than other techniques.

## APPENDIX

*Proof of Theorem 1* – Each pilot $k_i$ is a random variable taking value $v \in \{0, 1, 2, \dots\}$ with probability depending on the current load factor of the bitmap *taken* (fraction of $1$s). It follows that $k_i$ is *geometrically distributed* with success probability $p_i$ being the probability of placing all keys in $B_i$ without collisions. Let $\alpha(i)$ be the load factor of the bitmap after buckets $B_0, \dots, B_{i-1}$ have been processed, that is $\alpha(i) = \frac{\alpha}{n} \sum_{j=0}^{i-1} |B_j|$, for $i = 1, \dots, m-1$ and $\alpha(0) = 0$ for convenience (empty bitmap). Then the probability $p_i$ can be modeled as $p_i = (1 - \alpha(i))^{|B_i|}$.

Since $k_i$ is geometrically distributed, the probability that $k_i = v$, i.e., the probability of having success after $v$ failures ($v + 1$ total trials), is $p_i(1 - p_i)^v$, with expected value $\mathbb{E}[k_i] = \frac{1}{p_i} - 1 = \left(\frac{1}{1-\alpha(i)}\right)^{|B_i|} - 1$.

By linearity of expectation and taking into account that each trial takes at most $|B_i|$ time, the expected running time is at most

$$\sum_{i=0}^{m-1} |B_i| \cdot \left(\frac{1}{1 - \alpha(i)}\right)^{|B_i|}. \tag{1}$$

To bound (1) from above, we first introduce some helper definitions. Let $L$ be the largest bucket size and $C_\ell$ the set of all buckets of size $L - \ell$. Also, define $\beta(\ell)$ as the load factor of the bitmap after all buckets in the sets $C_0, \dots, C_{\ell-1}$ have been processed, with $\beta(0) = 0$. Using these definitions, the previous sum (1) can be rewritten as[3]

$$\sum_{\ell=0}^{L-1} \sum_{i=0}^{|C_\ell|-1} (L - \ell) \cdot \left(\frac{1}{1 - (\beta(\ell) + \alpha i (L - \ell)/n)}\right)^{L-\ell}. \tag{2}$$

We estimate the innermost sum with an integral and bound it from above using a Riemann-Sum, taking into account that the argument of the sum is an increasing function of $i$. (In fact, given an increasing function $f(x)$, $\int_a^b f(x)\,dx \le \sum_{i=a}^b f(i) \le \int_a^{b+1} f(x)\,dx$ holds true.) We obtain

$$(2) \le \sum_{\ell=0}^{L-1} \int_0^{|C_\ell|} (L - \ell) \cdot \left(\frac{1}{1 - (\beta(\ell) + \alpha y (L - \ell)/n)}\right)^{L-\ell} dy$$

$$= \sum_{\ell=0}^{L-1} \int_{1-\beta(\ell)}^{1-\beta(\ell+1)} (L - \ell) \cdot \left(\frac{1}{x}\right)^{L-\ell} \left(-\frac{n}{\alpha(L - \ell)}\right) dx$$

$$= \frac{n}{\alpha} \sum_{\ell=0}^{L-1} \int_{1-\beta(\ell+1)}^{1-\beta(\ell)} \left(\frac{1}{x}\right)^{L-\ell} dx. \tag{3}$$

Since keys are not uniformly distributed into the buckets, $\beta(\ell) \le \Delta\ell$ with $\Delta = \alpha/L$ being the average increase in the load factor caused by each set $C_\ell$. Therefore we have that

$$(3) \le \frac{n}{\alpha} \sum_{\ell=0}^{L-2} \int_{1-\Delta(\ell+1)}^{1-\Delta\ell} \left(\frac{1}{x}\right)^{L-\ell} dx + Z, \tag{4}$$

where the last term

$$Z = \frac{n}{\alpha} \int_{1-\alpha\cdot(1-1/n)}^{1-\alpha\cdot(1-1/n)+\Delta} \frac{1}{x}\,dx$$

$$= \frac{n}{\alpha} \log_e \left(\frac{1 - \alpha \cdot (1 - 1/n) + \Delta}{1 - \alpha \cdot (1 - 1/n)}\right) \tag{5}$$

comes from the buckets of size 1 (class $C_{L-1}$)[4].

Now we overestimate each integral in the sum with $\Delta/(1 - \Delta(\ell+1))^{L-\ell}$ (the area of a rectangle whose base is $\Delta$ and height is $(1/x)^{L-\ell}$ in the left extreme $x = 1 - \Delta(\ell+1)$) and obtain that

$$(4) \le \frac{n\Delta}{\alpha} \sum_{\ell=0}^{L-2} \frac{1}{(1 - \Delta(\ell+1))^{L-\ell}} + Z$$

$$= \frac{n}{L} \sum_{\ell=2}^{L} \left(\frac{1}{1 - \alpha + \frac{\alpha}{L}(\ell-1)}\right)^\ell + Z$$

by replacing $\Delta$ with $\alpha/L$ and some manipulation. Since $(1/(1 - \alpha + \frac{\alpha}{L}(\ell-1)))^\ell$ is maximum around $\ell = L/(3\alpha)$, then

$$\sum_{\ell=2}^{L} \left(\frac{1}{1 - \alpha + \frac{\alpha}{L}(\ell-1)}\right)^\ell < (L-1)\left(\frac{1}{4/3 - \alpha(1 + 1/L)}\right)^{\frac{L}{3\alpha}}$$

$$< (L-1)2^{\alpha L}, \text{ for all } L > L',$$

where $L'$ is a small constant. The expected running time is therefore $O(n2^{\alpha L})$, which is $O(n^{1+\Theta(\alpha/c)})$ because $L = \Theta(\log n / c)$.

∎

## REFERENCES

[1] Djamal Belazzougui and Gonzalo Navarro. Alphabet-independent compressed text indexing. *ACM Transactions on Algorithms*, 10(4):1–19, 2014.

[2] Yi Lu, Balaji Prabhakar, and Flavio Bonomi. Perfect hashing for network applications. In *2006 IEEE International Symposium on Information Theory*, pages 2774–2778. IEEE, 2006.

[3] Chin-Chen Chang and Chih-Yang Lin. Perfect hashing schemes for mining association rules. *The Computer Journal*, 48(2):168–179, 2005.

[4] Djamal Belazzougui, Paolo Boldi, Rasmus Pagh, and Sebastiano Vigna. Fast prefix search in little space, with applications. In *European Symposium on Algorithms*, pages 427–438. Springer, 2010.

[5] Giulio Ermanno Pibiri and Rossano Venturini. Efficient data structures for massive n-gram datasets. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 615–624, 2017.

[6] Giulio Ermanno Pibiri and Rossano Venturini. Handling massive *N*-gram datasets efficiently. *ACM Transactions on Information Systems*, 37(2):25:1–25:41, 2019.

[7] Andrei Broder and Michael Mitzenmacher. Network applications of bloom filters: A survey. *Internet mathematics*, 1(4):485–509, 2004.

[8] Bin Fan, Dave G Andersen, Michael Kaminsky, and Michael D Mitzenmacher. Cuckoo filter: Practically better than bloom. In *Proceedings of the 10th ACM International on Conference on emerging Networking Experiments and Technologies*, pages 75–88, 2014.

---

3. In (1), the maximum value assumed by the function $\alpha(\cdot)$ is $\alpha \cdot (1 - 1/n)$ because, when processing the last key, the bitmap has a load factor of $\alpha \cdot (1 - 1/n)$. The same holds true for the function $\beta(\cdot)$ in (2).

4. Interestingly, our problem with $L = 1$ and $\alpha = 1$ reduces to the *coupon collector's* problem that has a complexity of $\Theta(n \log n)$ [35]. In this specific case, $\Delta$ would be $\alpha \cdot (1 - 1/n)$, which is the increase in the load factor by adding all the buckets in the class $C_{L-1}$. In fact, using this value of $\Delta$ and $\alpha = 1$ in (5) we would obtain a complexity of $\Theta(n \log n)$.

[9] Thomas Mueller Graf and Daniel Lemire. Xor filters: Faster and smaller than bloom and cuckoo filters. *Journal of Experimental Algorithmics (JEA)*, 25:1–16, 2020.

[10] Michael L Fredman, János Komlós, and Endre Szemerédi. Storing a sparse table with o(1) worst case access time. *Journal of the ACM (JACM)*, 31(3):538–544, 1984.

[11] Kurt Mehlhorn. On the program size of perfect and universal hash functions. In *23rd Annual Symposium on Foundations of Computer Science*, pages 170–175. IEEE, 1982.

[12] Giulio Ermanno Pibiri and Roberto Trani. PTHash: Revisiting FCH Minimal Perfect Hashing. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021. To appear.

[13] Fabiano C Botelho, Rasmus Pagh, and Nivio Ziviani. Practical perfect hashing in nearly optimal space. *Information Systems*, 38(1):108–131, 2013.

[14] Djamal Belazzougui, Paolo Boldi, Giuseppe Ottaviano, Rossano Venturini, and Sebastiano Vigna. Cache-oblivious peeling of random hypergraphs. In *2014 Data Compression Conference*, pages 352–361. IEEE, 2014. URL https://github.com/ot/emphf.

[15] Torben Hagerup and Torsten Tholey. Efficient minimal perfect hashing in nearly minimal space. In *Annual Symposium on Theoretical Aspects of Computer Science*, pages 317–326. Springer, 2001.

[16] Edward A Fox, Qi Fan Chen, and Lenwood S Heath. A faster algorithm for constructing minimal perfect hash functions. In *Proceedings of the 15th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 266–273, 1992.

[17] Djamal Belazzougui, Fabiano C Botelho, and Martin Dietzfelbinger. Hash, displace, and compress. In *European Symposium on Algorithms*, pages 682–693. Springer, 2009. URL https://github.com/bonitao/cmph.

[18] Kimmo Fredriksson and Fedor Nikitin. Simple compression code supporting random access and fast string matching. In *International Workshop on Experimental and Efficient Algorithms*, pages 203–216. Springer, 2007.

[19] Bohdan S Majewski, Nicholas C Wormald, George Havas, and Zbigniew J Czech. A family of perfect hashing methods. *The Computer Journal*, 39(6):547–554, 1996.

[20] Bernard Chazelle, Joe Kilian, Ronitt Rubinfeld, and Ayellet Tal. The bloomier filter: an efficient data structure for static support lookup tables. In *Proceedings of the fifteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 30–39. Citeseer, 2004.

[21] Marco Genuzio, Giuseppe Ottaviano, and Sebastiano Vigna. Fast scalable construction of (minimal perfect hash) functions. In *International Symposium on Experimental Algorithms*, pages 339–352. Springer, 2016. URL https://github.com/vigna/Sux4J.

[22] Marco Genuzio, Giuseppe Ottaviano, and Sebastiano Vigna. Fast scalable construction of ([compressed] static— minimal perfect hash) functions. *Information and Computation*, page 104517, 2020.

[23] Ingo Müller, Peter Sanders, Robert Schulze, and Wei Zhou. Retrieval and perfect hashing using fingerprinting. In *International Symposium on Experimental Algorithms*, pages 138–149. Springer, 2014.

[24] Guy Jacobson. Space-efficient static trees and graphs. In *30th annual symposium on foundations of computer science*, pages 549–554. IEEE Computer Society, 1989.

[25] Gonzalo Navarro. *Compact data structures: A practical approach*. Cambridge University Press, 2016.

[26] Antoine Limasset, Guillaume Rizk, Rayan Chikhi, and Pierre Peterlongo. Fast and scalable minimal perfect hashing for massive key sets. In *16th International Symposium on Experimental Algorithms*, volume 11, pages 1–11, 2017. URL https://github.com/rizkg/BBHash.

[27] Emmanuel Esposito, Thomas Mueller Graf, and Sebastiano Vigna. Recsplit: Minimal perfect hashing via recursive splitting. In *2020 Proceedings of the Twenty-Second Workshop on Algorithm Engineering and Experiments (ALENEX)*, pages 175–185. SIAM, 2020. URL https://github.com/vigna/sux.

[28] Robert Mario Fano. On the number of bits required to implement an associative memory. *Memorandum 61, Computer Structures Group, MIT*, 1971.

[29] Peter Elias. Efficient storage and retrieval by content and address of static files. *Journal of the ACM*, 21(2):246–260, 1974.

[30] GoogleBooks Corpus, version 2.0, http://storage.googleapis.com/books/ngrams/books/datasetsv2.html. 2012.

[31] The ClueWeb09 Dataset, https://lemurproject.org/clueweb09/webGraph.php. 2009.

[32] Paolo Boldi, Bruno Codenotti, Massimo Santini, and Sebastiano Vigna. Ubicrawler: A scalable fully distributed web crawler. *Software: Practice and Experience*, 34(8):711–726, 2004.

[33] UK-2005 URLs, http://data.law.di.unimi.it/webdata/uk-2005. 2005.

[34] Pavlos Fafalios, Vasileios Iosifidis, Eirini Ntoutsi, and Stefan Dietze. TweetsKB: A public and large-scale rdf corpus of annotated tweets. In *European Semantic Web Conference*, pages 177–190. Springer, 2018. URL https://data.gesis.org/tweetskb.

[35] Sheldon Ross. *A first course in probability*. Pearson, 9th edition, 2014.

**Giulio Ermanno Pibiri** is a postdoctoral researcher at ISTI-CNR, High Performance Computing Lab. He received a Ph.D. in Computer Science from the University of Pisa in 2019. His research interests involve data compression algorithms for indexing massive datasets, data structures, and information retrieval with focus on efficiency. For more information: http://pages.di.unipi.it/pibiri.

**Roberto Trani** received a Ph.D. in Computer Science from the University of Pisa, in 2020. He is a postdoctoral researcher at the National Research Council of Italy (ISTI-CNR), within the High Performance Computing Laboratory, and his research interests are machine learning, algorithms, information retrieval, and high performance computing. For more information: http://hpc.isti.cnr.it/roberto-trani.