

PaCHash: Packed and Compressed Hash Tables

Florian Kurpicz ✉ 

Karlsruhe Institute of Technology, Germany

Hans-Peter Lehmann ✉ 

Karlsruhe Institute of Technology, Germany

Peter Sanders ✉ 

Karlsruhe Institute of Technology, Germany

Abstract

We introduce PaCHash, a hash table that stores its objects contiguously in an array without intervening space, even if the objects have variable size. In particular, each object can be compressed using standard compression techniques. A small search data structure allows locating the objects in constant expected time. PaCHash is most naturally described as a static external hash table where it needs a constant number of bits of internal memory per block of external memory. However, PaCHash can be dynamized and is also useful for internal memory, having lower space consumption than all previous approaches even when considering only objects of identical size. For example, in some sense it beats a lower bound on the space consumption of k -perfect hashing. An implementation for fast SSDs needs about 5 bits of internal memory per block of external memory, requires only one disk access (of variable length) per search operation and has internal search overhead small compared to the disk access cost.

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

Keywords and phrases compressed data structure, external hash table, perfect hashing

Supplementary Material All implementations presented in this paper and scripts to reproduce our experimental evaluation are available under the GPLv3 license.

PaCHash and Separator implementation: <https://github.com/ByteHamster/PaCHash>

Scripts for reproduction of results: <https://github.com/ByteHamster/PaCHash-Experiments>

1 Introduction

Hash tables support constant time key-based retrieval of objects and are one of the most widely used data structures. *Compressed* data structures store data in a space efficient way, preferably approaching the information theoretical limit, and support various kinds of operations without the need to decompress the entire data structure first [26, 1, 22, 51]. There has been intensive previous work on both subjects but, surprisingly, the intersection leaves big gaps. There is a lot of work on hash tables which need little more space than just the stored objects themselves [8, 3, 23, 30, 44]. However, all these approaches are only space efficient for objects of identical size which makes it impossible to compress the objects with variable bit-length codes. Currently, most hash tables for objects of variable size store references from table entries to the data which entails a space overhead of at least $\log N$ bits per object, where N is the total size of all objects in the table. Throughout this paper, $\log x$ stands for $\log_2 x$. See Section 2 for an introduction of basic techniques and Appendix A for a summary of the notation.

PaCHash eliminates fragmentation by *packing* the objects contiguously in memory without leaving free space. This makes it impossible to use the approach of most previous hash tables to directly use the hash function value to (approximately) locate the objects. Instead, PaCHash uses a highly space efficient search data structure that translates hash function values to memory locations. More precisely, objects are first hashed to *bins*. The bins are

stored contiguously in m blocks of size B . PaCHash essentially stores one bin index per block using a searchable compressed representation which enables finding the block(s) where a bin is stored. In Section 4, we describe the data structure in more detail and in Section 5 we analyze it. Basically, for a tuning parameter a , the expected number of block reads to retrieve an object x of size $|x|$ is about $1 + 1/a + |x|/B$ while the internal memory data structure needs $2 + \log(a)$ bits per block. We also discuss even smaller representations.

There is little previous work for objects of variable size (see Section 3). For objects of identical size s , the most space efficient previous solutions are based on minimal perfect hashing (MPH) [18, 7] and require a constant number of bits per object. PaCHash approximates this when choosing $B = s$, also needing a (slightly larger) constant number of bits per object but lower construction time. The picture changes when we look at larger block sizes $B = ks$ and the corresponding approach of *minimal k -perfect hashing (MkPH)* [7]. Now, PaCHash still needs only a constant number of bits per block, while there is a *lower bound* of $\Omega(\log k)$ bits per block using MkPH (see Appendix C).

Section 6 describes different implementation variants of PaCHash including fully internal and fully external versions. Section 7 describes experiments for an external implementation. Section 8 summarizes the results and discusses possible directions for further research.

Our Contribution. In this paper, we design the new hash table PaCHash. The data structure supports objects of variable size with space overhead close to competitors that only support objects of identical size. We analyze it thoroughly in a variant of the external memory model. Finally, we compare our implementation with competitors from the literature. As close contenders, we also implement *Separator Hashing* [27, 33] and *Cuckoo Hashing* [5, 43] with adaptations that partially allow variable size objects.

2 Preliminaries

Monotonic Sequences and Bit Vectors. Our index data structure mainly consists of a compressed representation of a monotonically increasing sequence $p = \langle p_1, \dots, p_k \rangle$ of integers in the range $1..U$. Searching boils down to predecessor queries in p , i.e., given a query integer i , the largest sequence element $\leq i$ is returned.

A well-known practical solution is *Elias-Fano coding* [17, 21] which splits each p_i . The $\log(U/k)$ least significant bits are directly stored in an array L requiring $k \log(U/k)$ bits of space. The $\log(k)$ most significant bits form a monotonic sequence of integers $H = \langle u_1, \dots, u_k \rangle$ in the range $0..k$. H is stored in a bit vector of size $2k + 1$ where u_i is represented as a 1-bit in position $i + u_i$. The total space usage therefore is $k(2 + \log(U/k)) + 1$ bits. A predecessor query in p executes a $select_0$ query in H (finding the i -th 0-bit in H) which locates a cluster of entries in L that must contain the sought element. Using additional space $o(k)$, $select_0$ queries can be answered in constant time [11]. In contrast to the general case, we will show that searching the cluster takes expected constant time in our application.

One can also interpret p as the positions of 1-bits in a sparse bit vector which enables even more compact representations. For example, using Succincter [45], about $k(1.44 + \log(U/k)) + 1$ bits are achievable which is almost information theoretically optimal. In Section 4.2 we give an even more compact format exploiting additional structure in the bit vector.

Model of Computation. We describe our results in a variant of the external memory model [50] adapted to a situation where objects are compressed to variable length sequences of bits. We have a *fast memory* of size M bits. Accesses to a large *external memory* are I/Os

■ **Table 1** Space efficient object stores from the literature. To unify the notation, we convert all values so that they refer to objects of size $s = 256$ bytes stored in blocks of $B = 4096$ bytes. The load factor α is given in percent and the internal space I_b in bits per *block* of $B/s = 16$ objects.

Method	I_b	α	I/Os	Method	I_b	α	I/Os
Extendible Hashing [20]	* ³	90	1	SILT LogStore [35]	832	100	1
Larson et al. [34]	96	<96	1	SkimpyStash [15]	32	≤ 98	8
SILT SortedStore [35]	51	100	1	PaCHash , $a = 1$	2	99.95	2.06 ⁴
Linear Separator [32]	8	85	1	PaCHash , $a = 8$	5	99.95	1.19 ⁴
Separator [27, 33]	6	98	1	(b) Dedicated variable size object stores.			
Robin Hood [9]	3	99	1.3	(a) Stores for objects of identical size. Can be used for objects of variable size by using indirection or for some methods by accepting significantly lower load factors. For details about how PaCHash can utilize identical size objects, see Section 5.			
Ramakrishna et al. [47]	4	80	1				
Jensen, Pagh [29]	0	80	1.25				
Cuckoo [5, 43]	0	<100	2				
PaCHash , $a = 1$	2	100	2 ⁴				
PaCHash , $a = 8$	5	100	1.13 ⁴				

to blocks of B consecutive bits. In contrast to the original model, we analyze both I/Os and internal work. $scan(N)$ denotes the cost (I/Os *and* internal work) of scanning N bits of data.¹ $sort(N)$ denotes the cost of sorting N bits.² In particular, we are interested in a high load factor, which is N divided by the total external space usage.

3 Related Work

The following section introduces related data structures from the literature. Table 1 provides an overview over the most important parameters. There are close contenders in the form of *object stores* from the database literature. BerkeleyDB [41] uses a B⁺-Tree [13] of order d , where each node branches between d and $2d$ times. LevelDB [28] and RocksDB [19] use a Log-Structured Merge tree [42], which stores multiple levels of a static data structure with increasing size. Insertions go into the first level and when a level gets too full, it is merged into the next level. SILT’s *LogStore* [35], Facebook *Haystack* [6] and *FAWN* [2] simply store a pointer of size $\Omega(\log N)$ to each object. Real world instances often store very small objects [40], so the pointers add a considerable amount of overhead.

Sorted Objects. *LevelDB*’s static part [28] stores objects in key order, enabling range searches and common-prefix-compression. *SortedStore* in SILT [35] sorts the objects by their hashed key and uses entropy coded tries as an index. Pagh [43] proposes to sort the n objects by a hash function with range $\geq n^3$. The internal memory stores the first hash function value mapped to each block. This data structure can be queried using a predecessor data structure in time $O(\log \log n)$.

¹The internal work may depend on the encoding of the data. For example, we may need $\Theta(N)$ machine instructions, or, a faster encoding may enable bit-parallel processing in $O(N/\log n)$.

²This entails $(N/B)(1 + \lceil \log_{M/B}(N/M) \rceil)$ I/Os. In this paper algorithms with linear internal work are possible exploiting random integer keys. The cost also includes (de)coding overhead as in *scan* operations.

³The space usage per block is logarithmic in the number of blocks.

⁴PaCHash performs one I/O of variable size which is faster than the competitors’ multiple I/Os.

External Hash Tables. In external hash tables, each table cell corresponds to a fixed size block. A common technique to support variable size objects is using indirection by internally storing a pointer to the object contents, possibly inlining parts of the objects [35, Section 4]. NVMKV [39] and KallaxDB [10] use an SSD as one large hash table and rely on SSD internals to handle empty blocks in a space efficient way. Overflowing blocks due to hash collisions can be handled with perfect hashing [34, 47] or using one of the following techniques.

With *Hashing with Chaining*, objects of overflowing blocks are stored in linked lists. *SkimpyStash* [15] chains objects using an external successor pointer for each object. This trades internal memory space for latency because of multiple dependent I/Os. Jensen and Pagh’s [29] data structure reserves parts of the external memory as a buffer to reduce the need for chaining. *Extendible Hashing* [20] keeps a balanced tree of blocks. Overflowing blocks are split into two children indexing one more bit of the hashed key.

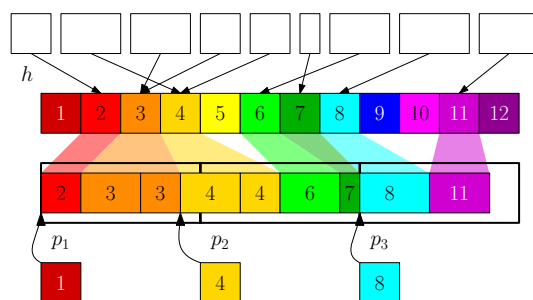
Another method for resolving collisions is *open addressing*, where each object could be located in multiple blocks. *Cuckoo Hashing* [44, 16] locates each object in one of two (or more [23]) independently hashed blocks. Queries can load both blocks in parallel to reduce latency. With *Separator Hashing* [27, 33], each object has a sequence of blocks it could be stored in and a corresponding sequence of signatures. When a block overflows, the objects with the highest signature values are pushed out to the next block in their respective sequences. The internal memory stores the highest signature value of the objects placed in each block. A query follows the object’s sequence of blocks and stops when it finds a separator that is larger than the corresponding signature. *Linear hashing with separators* [32] is a dynamic variant with a linear probe sequence. *External Robin Hood Hashing* [9] is similar to linear separator hashing, but it instead pushes out objects that are closest to their respective home address. For each block, the internal memory stores the smallest distance of its objects to their respective home address.

4 The PaCHash Data Structure

We now present PaCHash in detail – a hash table which considerably improves on the data structures from the literature. It natively supports variable size objects without the need for indirection or empty cells. It needs only a few bits of internal memory per *block* and still needs only one single I/O operation (of variable length) per query. PaCHash consists of an *external part* subdivided into m blocks of B bits each that store the actual objects and an *internal part* that allows finding the blocks storing an object. Figure 1 gives an example for the external and internal memory data structures.

4.1 External Object Representation

PaCHash stores the objects sorted by a hash function h with a rather small domain, namely $h : K \rightarrow 1..am$, where K is the set of possible keys, m is the number of blocks and a is a tuning parameter that we assume to be a power of two. The hashes can collide and therefore group the objects into am bins. The objects are now basically stored contiguously. “Basically” means that blocks may also contain information needed to find the first object or bin stored in them. Refer to Section 6 for a discussion of alternative encodings. Our default assumption is as follows: Each external block stores an offset of size $d = \log B$ bits indicating the bit where the first bin in the block starts. The remaining space stores the objects contiguously where an object may have an arbitrary size in bits. No space is left between subsequent objects. In particular, object representations may overlap block boundaries. We assume that objects are encoded in a self-delimiting way, i.e., when we know where an object starts, we can also find



■ **Figure 1** Example of PaCHash with $n = 9$ objects and $m = 3$ blocks. Using the hash function h , the objects are mapped to 12 bins shown as colors, i.e., $a = 4$. The bin content is then contiguously written to the external memory blocks. The internal memory index p stores the first bin intersecting with each block. Note that locating bin 8 will return the range 2..3, i.e., block 2 is loaded superfluously because there is no preceding empty bin that can encode whether it overlaps into the previous block. All other bins are located optimally.

its end. For example, we could have a prefix-free code for the objects. Construction first sorts the objects by their hash function value. Then it scans the sorted objects, constructing both the external and the internal data structure along the way. Refer to Section 5 for more details. If the internal data structure gets lost, for example due to a power outage, it can be re-generated using a single scan over the external memory data.

4.2 Internal Memory Data Structure

Given a bin b , the internal memory data structure p can be used to determine a (near-)minimal range $i..j$ of block indices such that b is stored in that range. When performing a query, that block range can then be loaded from external memory and scanned for the sought key. In practice, the resulting latency is often close to that of loading a single block since it includes only one disk seek. Conceptually, p stores a sequence $\langle p_1, \dots, p_m \rangle$ where p_i specifies the first bin whose data is at least partially contained in block i .⁵ We can use a predecessor query on p to determine i . When the predecessor is b itself, we also need to load the previous block. Another predecessor query or scanning then determines j , as illustrated by the pseudocode in Algorithm 1. To get the most out of this specification, we take empty bins into account: When a bin starts exactly at a block boundary and has an empty predecessor, we store that predecessor. This implies that if (and only if) a bin b starts at a block boundary and the previous bin $b - 1$ is nonempty, retrieving bin b will load one block too much. Note that p is a monotonically increasing sequence of integers which can be represented with different methods and trade-offs.

Elias-Fano Coding. A standard technique for storing monotonic sequences is Elias-Fano coding. A way to interpret the vector H of upper bits of an Elias-Fano coded sequence is that it stores the number of items having each possible combination of most significant bits in unary coding. To locate the predecessor of item $b = au + \ell$ in the sequence, we calculate $select_0(u - 1)$ on the upper bits H , which gives us the start of a cluster of entries that all have most significant bits u . The corresponding index in L can be calculated by subtracting

⁵An alternative would be to store the first bin that *starts* in each block. This introduces a special case when a block is fully overlapped by a bin and needs slightly more work when performing queries.

■ **Algorithm 1** A query for an object x calls $\text{locate}(x)$, loads the returned block range, and scans the blocks to find the object content. Determining the range boils down to predecessor queries on p .

```

Function locate( $x$ )
   $b := h(x)$ 
  find  $i$  such that  $p_{i-1} < b$  and  $p_i \geq b$  // Predecessor query
  if  $p_i = b$  then  $i := i - 1$  // Parts of  $b$  could be in previous block
  find the first  $j$  such that  $p_j > b$  // Predecessor query or scan  $p$  starting at  $i$ 
  return  $i..(j - 1)$ 

```

$(u - 1)$. We scan the cluster to find the largest index i with $p_i \leq b$. In our case, this takes constant expected time (see Lemma 5). The internal memory usage is $m(2 + \log(a) + o(1))$ bits (see Lemma 2).

Bit Vector with Succincter. It is also possible to store p as a bit vector with rank and select support. An item p_i at position i is then represented as a 1-bit in position $i + p_i$. The position of the predecessor of a bin b can be found in constant time by calculating $\text{select}_0(b) - b$. The actual value can be calculated using a select_1 query. Because the bit vector is sparse, we can use Succincter [45] to compress it and its rank and select structures down to about $m(1.44 + \log(a + 1) + o(1))$ bits (see Lemma 3).

Entropy Coding. We observed that in practice, the bit vector is considerably more regular than a truly random one and thus allows additional compression. This can be made fast by splitting it into ranges that are compressed individually, e.g., using dictionary compression.

5 Analysis

We now formalize the properties of PaCHash in Theorem 1. Then we prove the Theorem, discussing some variants and implications on the way. Section 5.1 considers construction time and final space consumption, while Section 5.2 looks at I/Os and internal work of queries.

► **Theorem 1.** *Consider n objects of total size N bits which are stored in m blocks of size B . Let $d \in 0.. \log B$ be an encoding-dependent number of bits needed to specify where the first bin or object of a block starts and $\bar{B} = B - d$ be the payload size per block, i.e., $m = N/\bar{B}$. For a parameter a , let a random hash function map the objects to am bins.*

Then, PaCHash with Elias-Fano coding needs $m(2 + \log a + o(1))$ bits of internal memory and $N(1 + d/\bar{B})$ bits of external memory. The construction cost is the same as that of sorting the objects using am random integer keys. The expected time for retrieving an object of size $|x|$ bits is constant plus the time for scanning $1 + |x|/\bar{B} + 1/a$ blocks. The unsuccessful search time is the same except that $|x|$ is replaced by 0.

5.1 Construction

Assuming that the set of input objects is stored in compressed form on external memory, we mainly need to sort the objects by their hash function value. In our model, this has complexity $\text{sort}(N)$. In most practically relevant situations, this can even be done in $O(\text{scan}(N))$ using integer sorting. Refer to Appendix B for details.

The sorted representation is then scanned and basically copied to the output, only adding d bits of information within each block, which allow a query to initialize the scanning

operation. What d depends on the concrete encoding of the data, ranging from $d = 0$ for objects of identical size or for 0-terminated strings to $d = \log(\mathbb{B})$ bits when we explicitly encode the starting position of an object or bin. Refer to Section 6 for examples.

► **Lemma 2.** *When using Elias-Fano coding to store p , the index needs $2 + \log a + o(1)$ bits of internal memory per block and can be constructed in time $O(m)$.*

Proof. p consists of $k = m$ integers $\leq am = U$. Inserting this into the space usage of Elias-Fano coded sequences (see Section 2) gives us $\text{space}(p) = k(2 + \log(U/k)) + 1 = m(2 + \log(am/m)) + 1 = m(2 + \log a) + 1$. The select_0 data structure on the upper bits H can be stored in $o(m)$ bits [11]. Each of the m insertions into the sequence can be done in constant time while generating the external object representation. The construction of the select_0 data structure takes time $O(m)$. ◀

As we show in Appendix C, minimum k -perfect hash functions need $\Omega(\log k)$ bits per block for identical size objects, while we show above that PaCHash needs a constant number. In a way, PaCHash therefore breaks the theoretical lower space bounds of MkPHFs while keeping the same $O(1)$ query time. Choosing parameter a large can bring the number of I/O operations arbitrarily close to optimal, independently of k .

► **Lemma 3.** *When using Succincter [45] to store p , the index needs $1.4427 + \log(a + 1) + o(1)$ bits of internal memory per block.*

Proof (Sketch, for full proof see Appendix D). Using Succincter, i.e., [45, Theorem 2] with a length- $(a + 1)m$ bit vector containing m ones, we can represent the internal memory index using only $\log \binom{(a+1)m}{m} + o(m) \leq m(1.4427 + \log(a + 1)) + o(m)$ bits, which results in the space mentioned above per external memory block. ◀

5.2 Query

► **Lemma 4.** *Retrieving an object x of size $|x|$ from a PaCHash data structure loads $\leq 1 + |x|/\bar{\mathbb{B}} + 1/a$ consecutive blocks from the external memory in expectation (setting $|x| = 0$ if x is not in the table).⁶*

Proof. We first derive the expected number of blocks overlapped by the bin $b_x = h(x)$ that x is stored in. We then analyze the edge case that PaCHash sometimes loads one additional block unnecessarily even though it is not overlapped. The expected size $\mathbb{E}(|b_x|)$ of b_x is the sum of $|x|$ and all other objects that are mapped to it:

$$\begin{aligned} \mathbb{E}(|b_x|) &= |x| + \sum_{y \in S, y \neq x} |y| \mathbb{P}(y \in b_x) \\ &\leq |x| + \sum_{y \in S} |y| \mathbb{P}(y \in b_x) = |x| + \sum_{y \in S} |y| \cdot \frac{1}{am} = |x| + \bar{\mathbb{B}}m \cdot \frac{1}{am} = |x| + \frac{\bar{\mathbb{B}}}{a} \end{aligned}$$

⁶Using fewer estimates in the proof one can derive a bound of $1 + \frac{|x| - c + 1 - e^{-\beta}}{\bar{\mathbb{B}}} + \frac{1}{a}$ where $\beta = \frac{n\bar{\mathbb{B}}}{Na}$ is the average number of objects per bin and c is the greatest common divisor of $\bar{\mathbb{B}}$ and all object sizes. In particular, for objects of identical size dividing \mathbb{B} , the bound is close to $1 + 1/a$.

Let X denote the number of blocks overlapped by bin b_x . Assuming that the block boundaries and bin boundaries are statistically independent,⁷ and using the linearity of the expected value, we get $\mathbb{E}(X) = 1 + (\mathbb{E}(|b_x|) - 1)/\bar{B} = 1 + |x|/\bar{B} + 1/a - 1/\bar{B}$.

At a position i , the sequence p stores the first bin b_i that intersects with block i . Most of the time, this also means that b_i extends into block $i - 1$, which is why queries load that block as well. When a bin starts *exactly* at a block boundary, though, the previous block is not actually needed. Because bin boundaries are statistically independent of block boundaries, the probability of that happening is $1/\bar{B}$.⁸

We get the result by putting together the expected blocks overlapped by a bin and the probability for loading one single block too much. For negative queries, we are interested in the size of the bin that x would be hashed to, so we can simply set $|x| = 0$. ◀

► **Lemma 5.** *When using Elias-Fano coding for the index data structure of PaCHash, the range of blocks containing the bin of an object x can be found in expected constant time.*

Proof. A query for an object x consists of four steps. First, we hash x to get the corresponding bin $b_x = au + \ell$. We then execute a constant time [11] $select_0$ query on the upper bits H . That gives us the start of a cluster of entries in the sequence that all have the same $\log(m)$ most significant bits u . We need to iterate over the cluster entries which are $< b_x$ until we find the predecessor. Each cluster entry corresponds to a stored bin index. Let us bound the expected size $\mathbb{E}(Y_u)$ of all bins that have most significant bits u and are $< b_x$.

$$\begin{aligned} \mathbb{E}(Y_u) &= \sum_{y \in S} |y| \cdot \mathbb{P}(h(y) \text{ has MSB} = u \text{ and } h(y) < h(x)) \\ &\leq \sum_{y \in S} |y| \cdot \mathbb{P}(h(y) \text{ has MSB} = u) = \frac{1}{m} \sum_{y \in S} |y| = \frac{m\bar{B}}{m} = \bar{B} \end{aligned}$$

The expected number of cluster entries we need to scan is therefore $\mathbb{E}(Y_u)/\bar{B} = 1$. The practical implementation then further scans the cluster to find the last block overlapping b_x . This takes non-constant time $O(1 + |x|/\bar{B})$, which is not a problem since a proportional number of blocks are loaded anyway. However, we strengthen the lemma by observing that we can also use another $select_0$ query followed by a backward scan of the cluster. ◀

6 Variants and Refinements

Up until now, PaCHash was described as a static, external hash table for objects of variable size. The following section describes variants of the scheme.

Object Encoding. Instead of storing objects contiguously with a self-delimiting encoding, PaCHash allows for a wide range of other options, as shown in Table 2. In general, we have a tradeoff between the space needed for setting up object decoding in a block and the strength of assumptions made on object representation. We single out the particularly important

⁷We can guarantee the independence by cyclically shifting the data structure, i.e., we set the offset of the first block to a random number in $0..(\bar{B} - 1)$ and let the last bins wrap around into the first block.

⁸When the preceding bin b_{-1} is empty, PaCHash stores that empty bin in p , as described in Section 4. This means that the probability of unnecessary block loads actually is smaller, namely $\frac{1}{\bar{B}}(1 - \mathbb{P}(|b_{-1}| > 0))$, where $\mathbb{P}(|b_{-1}| > 0) = \left(1 - \frac{1}{am}\right)^n \approx e^{-\frac{n}{am}}$ is the probability of b_{-1} being empty.

■ **Table 2** External space overhead of d bits per block in order to facilitate scanning that block. The term $+1$ when $d \neq 0$ is needed for the case that no object starts in a block.

d	Case Description
0	Identical object sizes, zero terminated strings and analogous cases
$\lceil \log(w + 1) \rceil$	Objects that use variable bit-length encoding with $\leq w \leq B$ bits
$\lceil \log(W/w + 1) \rceil$	Objects of size divisible by w with $W = \min(B, \max \text{ object size})$
$\lceil \log(B) \rceil$	Explicit storage of a starting position of a bin

case of objects of identical size where we can calculate the block offset at query time and therefore need no external space overhead. When the object size divides the block size, it can be shown that the expected number of I/O operations is close to $1 + 1/a$.

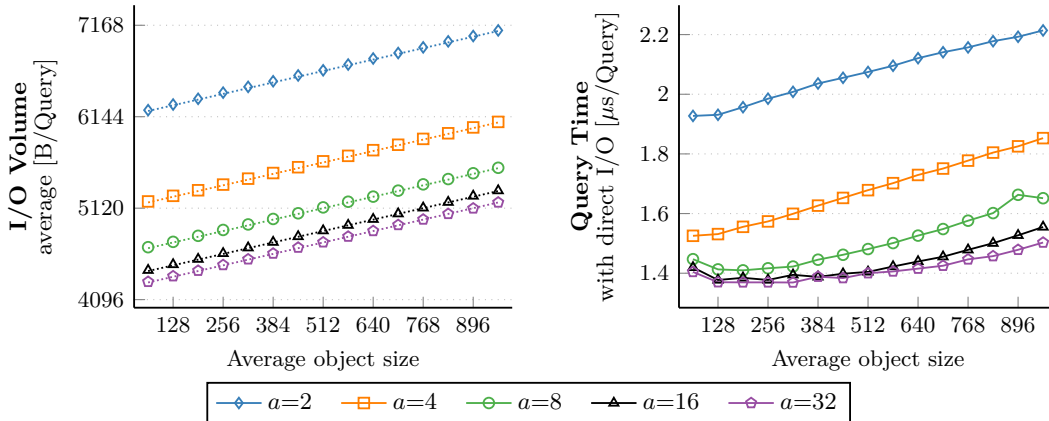
Memory Locations. PaCHash can be stored fully externally. By doing so, the number of I/Os for a query is increased by three (two I/Os to query the rank and select data structure on the bit vector of the Elias-Fano coding and one I/O to get the remaining bits). The number of I/Os can be reduced by interleaving the arrays of the Elias-Fano coding. PaCHash is also interesting as a purely internal data structure since it allows for configurations that need less space than any previous approach, even for objects of identical size. A variant that simplifies the external memory representation is to store the d bits of offsets per block in an internal memory data structure, possibly interleaved with the Elias-Fano representation. A variant enabling faster scanning of blocks separates keys and values [37], for example by storing $\log B$ bits of offset for each object.

Functional Enhancements. Because PaCHash sorts objects by their hashed key, *range queries* with respect to the original keys are not immediately possible. Litwin and Lomet [36] implement range queries for hash tables by partitioning the key space into smaller pieces. An index tree then leads to a number of small (PaCHash) tables that are fully scanned. Order-preserving hash functions [25] are another alternative.

PaCHash can be made *dynamic* using standard techniques like a Log-Structured Merge Tree [42, 38]. Merging multiple PaCHash data structures is possible efficiently. The idea is to construct the hash function h by first hashing to a larger range and then mapping it linearly to the range am . When updating h to the new total number of blocks, the objects of both input data structures are already sorted and can be merged with a linear sweep.

7 Experiments

Experimental Setup. We run our experiments on an Intel i7 11700 processor with 8 cores and a base clock speed of 2.5 GHz. We use a Samsung 980 Pro NVMe SSD with a capacity of 1 TB. The machine runs Ubuntu 21.10 with Linux 5.13.0. We use the GNU C++ compiler version 11.2.0 with optimization flags `-O3 -march=native`. Externally, each block of size $B = 2^{15}$ bits (4096 bytes) stores a table of 8 byte keys and 2 byte object offsets. During construction, we sort pointers to the objects using IPS²Ra [4]. Unless otherwise specified, the index is an Elias-Fano coded sequence based on sds!s [26] arrays of flexible bit width and the select data structures by Kurpicz [31]. For the I/O operations, we use `io_uring`. Query operations keep a queue of 128 asynchronous requests in flight.

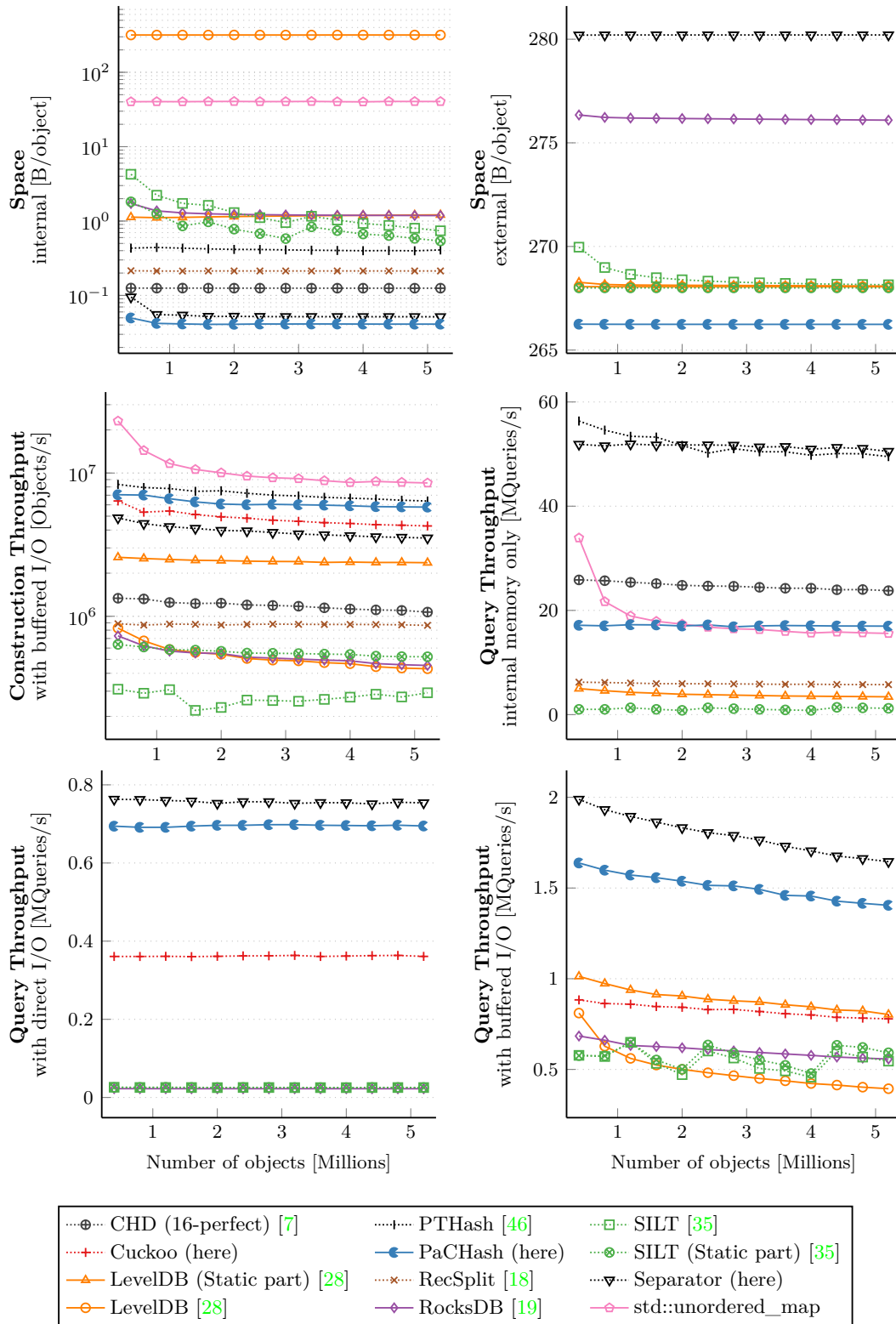


■ **Figure 2** Dependence of I/O volume and query time on the average object size s . Sizes are normal distributed with variance $s/5$, rounded to the next positive integer. Dotted lines show theoretic I/O volumes, while marks show measurements. Note that the measurements closely match the analysis. Using other distributions and plotting over the returned objects’ sizes gives equivalent results.

Competitors. To our knowledge, there is no existing implementation of a hash table for variable size objects that is simultaneously aimed at low internal memory usage and few I/O operations. As the main competitors, we choose LevelDB [28], RocksDB [19] and SILT [35]. To abstract from the different implementations of I/O operations, we also extract the internal memory index (address calculation) from some competitors and compare it to `std::unordered_map`, as well as the perfect hash functions RecSplit [18], CHD [7, 14], and PTHash [46]. We also implement Separator Hashing [27, 33] and Cuckoo Hashing [5, 43]. In contrast to the original papers, our implementations can be used with objects of variable size $\leq B$ when setting the load factor low enough. Note that decreasing the load factor increases the number of blocks and therefore the space needed for indexing. The construction of PaCHash always succeeds, while it can fail for Separator and Cuckoo Hashing depending on the preselected load factor or tuning parameter. Refer to Appendix E.2 for details.

PaCHash Configurations. Figure 2 plots the bytes read per query, depending on the average object size and parameter a . It confirms the results of our theoretical analysis in practice. The parameter a provides a trade-off between internal space usage and query performance. The throughput of the Elias-Fano representation increases when parameter a gets larger because the SSD needs to load fewer blocks. We also see that (at least for larger a) query times grow more slowly with object size than the I/O volume. We choose $a = 8$ for the comparison with competitors because it achieves a good balance between space usage (≈ 5 bits/block) and throughput ($\approx 700k$ Queries/second). In Appendix E.3, we compare the space usage and performance of different index data structures for PaCHash using real world size distributions from Twitter, Wikipedia and a protein database. The entropy coded bit vector saves up to one bit of internal memory per block for small a . While it comes with a performance penalty caused by decompression, it is fast enough that it can be useful for some applications. Succincter provides space usage lower than Elias-Fano but has no implementation.

Comparison with Competitors. Figure 3 compares PaCHash to other hash table data structures – see Appendix E.1 for the exact configurations used. These plots use identical size objects in order to allow for a large set of competitors. Perhaps the closest contender



■ **Figure 3** Comparison with competitor object stores using objects of identical size 256 bytes. Keys are 8 byte random strings. Dotted lines indicate that the methods only support objects of identical size.

to PaCHash is the Separator method where our implementation allows variable object size ($\leq B$). It needs comparable internal space and has faster queries (always a single block access). However, Separator not only has slower construction, but it also cannot achieve a load factor close to 100% except for objects with identical size when the block size is divisible by the object size. Appendix E.2 gives details showing load factors between 85% and 95% in typical cases. The perfect hashing methods CHD and RecSplit have similar problems with respect to variable size objects and are more expensive with respect to internal space and construction costs. While PTHash offers fast construction and queries, it does not support variable size objects and needs more internal space. Cuckoo hashing needs no internal space but has more expensive queries and the same problem with high load factors as Separator or perfect hashing. The object stores LevelDB, RocksDB, and SILT have much larger internal space requirements *and* some external overhead. In part this comparison is unfair since they have additional functionality like dynamic operation. For SILT and LevelDB we have been able to extract the static part but still get considerably more space and lower performance than PaCHash. Comparing query throughput is complicated because of different file access modes, internal caching, and history dependent performance for the actual SSD accesses (the controller uses caching and rearranges data outside the control of the user). We have therefore looked at two different access methods and also at only the index data structure. However, overall, we get a consistent picture with Separator being the fastest method followed by PaCHash. A comparison with the vanilla internal hash table `std::unordered_map` is also instructive. We naturally get faster construction and high internal space consumption. However, surprisingly, access to the internal data structure is only faster than PaCHash for very small inputs that fit into cache. In Appendix E.4, we show that the same observations apply to objects of variable size.

8 Conclusion and Future Work

With PaCHash, we present a static hash table that can store variable size (possibly compressed) objects in a very space efficient way. The objects are stored contiguously without the usual need for empty space to equalize the nonuniformity in assignment by a hash function. This is facilitated by an index data structure that needs only a constant number of internal memory bits per external memory block. In constant expected time, it yields a near-optimal range of blocks that contain the sought object. Our implementation of PaCHash considerably outperforms previous object stores for variable size objects and even matches or outperforms systems that are purely internal memory or only handle objects of identical size.

Future work might include the integration of PaCHash into dynamic external memory object stores as well as the engineering of fast and space efficient purely internal memory variants. On the theoretical side, we would like to better understand the space requirements of bit vectors with entropy coding as well as lower bounds. This includes relations to different variants of perfect hashing. Although our current analysis assumes random hash functions, PaCHash may also be provably efficient for more realistic simple hash functions.

Acknowledgements. The authors would like to thank Peter Dillinger and Stefan Walzer for early discussions leading to this paper. 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).



References

- 1 Rachit Agarwal, Anurag Khandelwal, and Ion Stoica. Succinct: Enabling queries on compressed data. In *NSDI*, pages 337–350. USENIX Association, 2015.
- 2 David G. Andersen, Jason Franklin, Michael Kaminsky, Amar Phanishayee, Lawrence Tan, and Vijay Vasudevan. FAWN: a fast array of wimpy nodes. In *SOSP*, pages 1–14. ACM, 2009. doi:10.1145/1629575.1629577.
- 3 Yuriy Arbitman, Moni Naor, and Gil Segev. Backyard cuckoo hashing: Constant worst-case operations with a succinct representation. In *FOCS*, pages 787–796. IEEE Computer Society, 2010. doi:10.1109/FOCS.2010.80.
- 4 Michael Axtmann, Sascha Witt, Daniel Ferizovic, and Peter Sanders. Engineering in-place (shared-memory) sorting algorithms. *ACM Trans. Parallel Comput.*, 9(1):2:1–2:62, 2022. doi:10.1145/3505286.
- 5 Yossi Azar, Andrei Z. Broder, Anna R. Karlin, and Eli Upfal. Balanced allocations (extended abstract). In *STOC*, pages 593–602. ACM, 1994. doi:10.1145/195058.195412.
- 6 Doug Beaver, Sanjeev Kumar, Harry C. Li, Jason Sobel, and Peter Vajgel. Finding a needle in haystack: Facebook’s photo storage. In *OSDI*, pages 47–60. USENIX Association, 2010.
- 7 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.
- 8 Michael A. Bender, Alex Conway, Martin Farach-Colton, William Kuszmaul, and Guido Tagliavini. All-purpose hashing. *CoRR*, abs/2109.04548, 2021.
- 9 Pedro Celia. External robin hood hashing. Technical report, Computer Science Department, Indiana University. TR246, 1988.
- 10 Xubin Chen, Ning Zheng, Shukun Xu, Yifan Qiao, Yang Liu, Jiangpeng Li, and Tong Zhang. Kallaxdb: A table-less hash-based key-value store on storage hardware with built-in transparent compression. In *DaMoN*, pages 3:1–3:10. ACM, 2021. doi:10.1145/3465998.3466004.
- 11 David Clark. *Compact PAT trees*. PhD thesis, University of Waterloo, 1997. URL: <http://hdl.handle.net/10012/64>.
- 12 Yann Collet. LZ4: Extremely fast compression algorithm. <https://github.com/lz4/lz4>.
- 13 Douglas Comer. The ubiquitous B-tree. *ACM Comput. Surv.*, 11(2):121–137, 1979. doi:10.1145/356770.356776.
- 14 Davi de Castro Reis, Djamel Belazzougui, Fabiano Cupertino Botelho, and Nivio Ziviani. CMPH - C minimal perfect hashing library. <http://cmph.sourceforge.net/>, 2012.
- 15 Biplob K. Debnath, Sudipta Sengupta, and Jin Li. Skimpystash: RAM space skimpy key-value store on flash-based storage. In *SIGMOD Conference*, pages 25–36. ACM, 2011. doi:10.1145/1989323.1989327.
- 16 Martin Dietzfelbinger and Christoph Weidling. Balanced allocation and dictionaries with tightly packed constant size bins. *Theor. Comput. Sci.*, 380(1-2):47–68, 2007. doi:10.1016/j.tcs.2007.02.054.
- 17 Peter Elias. Efficient storage and retrieval by content and address of static files. *J. ACM*, 21(2):246–260, 1974. doi:10.1145/321812.321820.
- 18 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.
- 19 Facebook. RocksDB. a persistent key-value store for fast storage environments. <https://rocksdb.org>, 2021.
- 20 Ronald Fagin, Jürg Nievergelt, Nicholas Pippenger, and H. Raymond Strong. Extendible hashing - A fast access method for dynamic files. *ACM Trans. Database Syst.*, 4(3):315–344, 1979. doi:10.1145/320083.320092.
- 21 Robert Mario Fano. On the number of bits required to implement an associative memory. Technical report, MIT, Computer Structures Group, 1971. Project MAC, Memorandum 61".

- 22 Paolo Ferragina and Giovanni Manzini. Indexing compressed text. *J. ACM*, 52(4):552–581, 2005. doi:10.1145/1082036.1082039.
- 23 Dimitris Fotakis, Rasmus Pagh, Peter Sanders, and Paul G. Spirakis. Space efficient hash tables with worst case constant access time. *Theory Comput. Syst.*, 38(2):229–248, 2005. doi:10.1007/s00224-004-1195-x.
- 24 Matteo Frigo, Charles E. Leiserson, Harald Prokop, and Sridhar Ramachandran. Cache-oblivious algorithms. In *FOCS*, pages 285–298. IEEE Computer Society, 1999. doi:10.1109/SFFCS.1999.814600.
- 25 Anil K. Garg and C. C. Gotlieb. Order-preserving key transformations. *ACM Trans. Database Syst.*, 11(2):213–234, 1986. doi:10.1145/5922.5923.
- 26 Simon Gog, Timo Beller, Alistair Moffat, and Matthias Petri. From theory to practice: Plug and play with succinct data structures. In *SEA*, volume 8504 of *Lecture Notes in Computer Science*, pages 326–337. Springer, 2014. doi:10.1007/978-3-319-07959-2_28.
- 27 Gaston H. Gonnet and Per-Åke Larson. External hashing with limited internal storage. *J. ACM*, 35(1):161–184, 1988. doi:10.1145/42267.42274.
- 28 Google. LevelDB is a fast key-value storage library written at google. <https://github.com/google/leveldb>, 2021.
- 29 Morten Skaarup Jensen and Rasmus Pagh. Optimality in external memory hashing. *Algorithmica*, 52(3):403–411, 2008. doi:10.1007/s00453-007-9155-x.
- 30 Donald E. Knuth. *The Art of Computer Programming, Volume III: Sorting and Searching*. Addison-Wesley, 1973.
- 31 Florian Kurpicz. Engineering compact data structures for rank and select queries on bit vectors. *CoRR*, abs/2206.01149, 2022.
- 32 Per-Åke Larson. Linear hashing with separators - A dynamic hashing scheme achieving one-access retrieval. *ACM Trans. Database Syst.*, 13(3):366–388, 1988. doi:10.1145/44498.44500.
- 33 Per-Åke Larson and Ajay Kajla. File organization: Implementation of a method guaranteeing retrieval in one access. *Commun. ACM*, 27(7):670–677, 1984. doi:10.1145/358105.358193.
- 34 Per-Åke Larson and M. V. Ramakrishna. External perfect hashing. In *SIGMOD Conference*, pages 190–200. ACM Press, 1985. doi:10.1145/318898.318916.
- 35 Hyeontaek Lim, Bin Fan, David G. Andersen, and Michael Kaminsky. SILT: a memory-efficient, high-performance key-value store. In *SOSP*, pages 1–13. ACM, 2011. doi:10.1145/2043556.2043558.
- 36 Witold Litwin and David B. Lomet. The bounded disorder access method. In *ICDE*, pages 38–48. IEEE Computer Society, 1986. doi:10.1109/ICDE.1986.7266204.
- 37 Lanyue Lu, Thanumalayan Sankaranarayanan Pillai, Hariharan Gopalakrishnan, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. Wisckey: Separating keys from values in ssd-conscious storage. *ACM Trans. Storage*, 13(1):5:1–5:28, 2017. doi:10.1145/3033273.
- 38 Chen Luo and Michael J. Carey. LSM-based storage techniques: a survey. *VLDB J.*, 29(1):393–418, 2020. doi:10.1007/s00778-019-00555-y.
- 39 Leonardo Mármlol, Swaminathan Sundararaman, Nisha Talagala, and Raju Rangaswami. NVMKV: A scalable, lightweight, ftl-aware key-value store. In *USENIX Annual Technical Conference*, pages 207–219. USENIX Association, 2015.
- 40 Rajesh Nishtala, Hans Fugal, Steven Grimm, Marc Kwiatkowski, Herman Lee, Harry C. Li, Ryan McElroy, Mike Paleczny, Daniel Peek, Paul Saab, David Stafford, Tony Tung, and Venkateshwaran Venkataramani. Scaling memcache at facebook. In *NSDI*, pages 385–398. USENIX Association, 2013.
- 41 Michael A. Olson, Keith Bostic, and Margo I. Seltzer. Berkeley DB. In *USENIX Annual Technical Conference, FREENIX Track*, pages 183–191. USENIX, 1999.
- 42 Patrick E. O’Neil, Edward Cheng, Dieter Gawlick, and Elizabeth J. O’Neil. The log-structured merge-tree (LSM-tree). *Acta Informatica*, 33(4):351–385, 1996. doi:10.1007/s002360050048.

- 43 Rasmus Pagh. Basic external memory data structures. In *Algorithms for Memory Hierarchies*, volume 2625 of *Lecture Notes in Computer Science*, pages 14–35. Springer, 2003. doi:[10.1007/3-540-36574-5_2](https://doi.org/10.1007/3-540-36574-5_2).
- 44 Rasmus Pagh and Flemming Friche Rodler. Cuckoo hashing. *J. Algorithms*, 51(2):122–144, 2004. doi:[10.1016/j.jalgor.2003.12.002](https://doi.org/10.1016/j.jalgor.2003.12.002).
- 45 Mihai Patrascu. Succincter. In *FOCS*, pages 305–313. IEEE Computer Society, 2008. doi:[10.1109/FOCS.2008.83](https://doi.org/10.1109/FOCS.2008.83).
- 46 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).
- 47 M. V. Ramakrishna and Walid R. Tout. Dynamic external hashing with guaranteed single access retrieval. In *FODO*, volume 367 of *Lecture Notes in Computer Science*, pages 187–201. Springer, 1989. doi:[10.1007/3-540-51295-0_127](https://doi.org/10.1007/3-540-51295-0_127).
- 48 Peter Sanders, Kurt Mehlhorn, Martin Dietzfelbinger, and Roman Dementiev. *Sequential and Parallel Algorithms and Data Structures - The Basic Toolbox*. Springer, 2019. doi:[10.1007/978-3-030-25209-0](https://doi.org/10.1007/978-3-030-25209-0).
- 49 Baris E. Suzek, Hongzhan Huang, Peter B. McGarvey, Raja Mazumder, and Cathy H. Wu. Uniref: comprehensive and non-redundant uniprot reference clusters. *Bioinform.*, 23(10):1282–1288, 2007. doi:[10.1093/bioinformatics/btm098](https://doi.org/10.1093/bioinformatics/btm098).
- 50 Jeffrey Scott Vitter and Elizabeth A. M. Shriver. Algorithms for parallel memory I: two-level memories. *Algorithmica*, 12(2/3):110–147, 1994. doi:[10.1007/BF01185207](https://doi.org/10.1007/BF01185207).
- 51 Feng Zhang, Jidong Zhai, Xipeng Shen, Onur Mutlu, and Wenguang Chen. Efficient document analytics on compressed data: Method, challenges, algorithms, insights. *Proc. VLDB Endow.*, 11(11):1522–1535, 2018. doi:[10.14778/3236187.3236203](https://doi.org/10.14778/3236187.3236203).

A Symbols

Table 3 summarizes the most important symbols used in this paper.

■ **Table 3** Symbols used in this paper

S	Set of objects	$m = N/\bar{B}$	Number of blocks
n	Number of objects	B	Block size (bits)
N	Total size of objects (bits)	$\bar{B} = B - d$	Payload data per block
p	Internal index data structure	$d \in 0.. \log B$	Encoding-dependent number of bits
a	Tuning parameter: Bins per block		to store position of first bin of block

B External Sorting

We now show that the external sorting needed during construction of a PaCHash data structure can be done in scanning complexity using very modest additional assumptions. First note that the problem of sorting objects during construction is easy when the average object size exceeds the block size, i.e., $N/n > B$ and thus $n < N/B$. In that case, a variant of bucket sort that maps the keys to $O(n)$ buckets runs with linear internal expected work and $O(n + N/B) = O(N/B)$ I/Os [48, Theorem 5.9].

On the other hand, the average object size N/n must be at least $\log n$ since we are looking at objects with unique keys. For the remaining case $\log n \leq N/n \leq B$, we additionally make a *tall cache assumption* quite usual for external memory [24] where $M > B^2$. Since the index data structure has at least N/B bits, we also know that $M \geq N/B$. A single scan of the input can partition it into pieces of size about

$$\frac{N}{M/B} \leq \frac{N}{(N/B)/B} = B^2 \leq M$$

which fit into internal memory. Moreover, since the average object size is $\geq \log n$, we can afford to replace the objects in an internally sorted fragment of the input by key-pointer pairs which once more allows us to use bucket sort – this time running in internal memory.

C Lower Space Bounds of Perfect Hashing

The lower bound for the space usage of a minimum k -perfect hash function for objects of identical size approaches $n \cdot (\log(e) + \log(k!/k^k)/k)$ [7]. Using Stirling's approximation, we derive a new lower space bound that is easier to interpret.

$$\begin{aligned} n \cdot (\log(e) + \log(k!/k^k)/k) &\approx n \cdot \left(\log(e) + \log \left(\frac{\sqrt{2\pi k} (k/e)^k}{k^k} \right) / k \right) \\ &= n \cdot \left(\log(e) + \log(\sqrt{2\pi k} (1/e^k)) / k \right) = n \cdot \left(\log(e) + \frac{\log(\sqrt{2\pi k})}{k} - \frac{\log(e^k)}{k} \right) \\ &= n \cdot \left(\log(e) + \frac{\log((2\pi k)^{1/2})}{k} - \log(e) \right) = \frac{n}{k} \cdot \frac{1}{2} \log(2\pi k) \end{aligned}$$

The value n/k is the number of blocks, so Mk PHFs need $\Omega(\log k)$ bits of space per block.

D Space Usage of Succincter

Now, we show in more detail how we can achieve the memory requirements of the internal memory index of PaCHash using the Succincter rank and select data structure [45].

Full Proof of Lemma 3. Remember that the internal memory data structure p of PaCHash stores m integers in the range $1..am$ and must support predecessor queries. We represent all integers in a bit vector of length $(a+1)m$, using the same idea used for the most significant bits in Elias-Fano coding. That is, each of the m integers p_i is represented as a 1-bit in position $i + p_i$. Answering predecessor queries (which we do not consider here) becomes harder to analyze, as we have no information about the distribution of 1-bits in the bit vector.

Using Succincter, we can store a size- u bit vector that contains n ones and supports rank and select queries using only $\log \binom{u}{n} + \frac{u}{\log u} + \tilde{O}\left(u^{\frac{3}{4}}\right)$ bits. Since we have a length- $(a+1)m$ bit vector that contains m ones, we require $\log \binom{(a+1)m}{m} + \frac{(a+1)m}{\log((a+1)m)} + \tilde{O}\left(\left((a+1)m\right)^{\frac{3}{4}}\right)$ bits of space. We now show the upper bound for required memory using Lemma 7 and $\tilde{O}\left(\left((a+1)m\right)^{\frac{3}{4}}\right) = o(m)$.

$$\begin{aligned}
 \log \binom{(a+1)m}{m} + o(m) &< \log \left(\sqrt{\frac{(a+1)}{2\pi am}} \left(\frac{(a+1)^{a+1}}{a^a} \right)^m e^{-\frac{1}{12m+1}} \right) + o(m) \\
 &= \underbrace{\log \sqrt{\frac{(a+1)}{2\pi am}}}_{\leq 0} + \log \left(\left(\frac{(a+1)^{a+1}}{a^a} \right)^m \right) + \underbrace{\log e^{-\frac{1}{12m+1}}}_{\leq 0} + o(m) \\
 &\leq \log \left(\left(\frac{(a+1)^{a+1}}{a^a} \right)^m \right) + o(m) \\
 &= m \left((a+1) \log(a+1) - a \log a \right) + o(m) \\
 &= m \left(a \log \left(\frac{a+1}{a} \right) + \log(a+1) \right) + o(m) \\
 &\leq m (1.4427 + \log(a+1)) + o(m)
 \end{aligned}$$

The last inequality is due to the fact that $a \log \left(\frac{a+1}{a} \right)$ converges to $1.4427 \approx \frac{1}{\ln 2}$ from below. Overall, we require less than $1.4427 + \log(a+1) + o(1)$ bits for each external memory block. ◀

► **Lemma 6.** *Using Succincter for representing monotonic sequences is almost space optimal.*

Proof. In Lemma 3 we have already seen that Succincter needs close to $m(\log(e) + \log(a+1))$ bits of space. $\binom{am}{m}$ is the number of *strictly* monotonic sequences of m numbers in the range $1..am$ and thus a lower bound for the number of monotonic sequences. Using Lemma 7 once more, we get

$$\log \binom{am}{m} \approx m \left((a-1) \log \left(\frac{a}{a-1} \right) + \log a \right)$$

bits as a lower bound. Looking at the difference divided by m (i.e. bits per block), we get

$$\begin{aligned} a \log \frac{a+1}{a} + \log(a+1) - (a-1) \log \frac{a}{a-1} - \log a &= a \log \frac{a^2-1}{a^2} + \log \frac{a+1}{a-1} \\ &= \frac{\log e}{a} + O\left(\frac{1}{a^3}\right). \end{aligned}$$

This difference (obtained using Taylor series development) is much smaller than the $\log e + \log(a+1)$ bits per block needed by the Succincter data structure – at least for sufficiently large a . ◀

► **Lemma 7.** For any $c > 1, n > 0$, let $f(n, c) := \sqrt{\frac{c}{(c-1)2\pi n}} \left(\frac{c^c}{(c-1)^{c-1}}\right)^n$, then

$$f(n, c) \left(1 - \frac{c^2 - c + 1}{12c(c-1)n}\right) < \binom{cn}{n} < f(n, c) e^{-\frac{1}{12n+1}} = f(n, c) \left(1 - \frac{1}{12n} + O\left(\frac{1}{n^2}\right)\right).$$

Proof. We use the identity $\binom{cn}{n} = \frac{(cn)!}{n!(cn-n)!}$ as well as Stirling's approximation

$$\sqrt{2\pi m} \left(\frac{m}{e}\right)^m e^{\frac{1}{12m+1}} < m! < \sqrt{2\pi m} \left(\frac{m}{e}\right)^m e^{\frac{1}{12m}}.$$

For the upper bound we get

$$\begin{aligned} \binom{cn}{n} &= \frac{(cn)!}{n!(cn-n)!} < \frac{\sqrt{2\pi cn} \left(\frac{cn}{e}\right)^{cn} e^{\frac{1}{12cn}}}{\sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{1}{12n+1}} \cdot \sqrt{2\pi(c-1)n} \left(\frac{(c-1)n}{e}\right)^{(c-1)n} e^{\frac{1}{12(c-1)n+1}}} \\ &= \sqrt{\frac{c}{(c-1)2\pi n}} \cdot \left(\frac{c^c}{(c-1)^{c-1}}\right)^n \cdot e^{\frac{1}{12cn} - \frac{1}{12n+1} - \underbrace{\frac{1}{12(c-1)n+1}}_{\leq 12cn}}. \end{aligned}$$

The claim follows by observing that the **leftmost** and **rightmost** term in the exponent of e cancel out in the estimation. The asymptotic expansion of the upper bound can be obtained using Taylor series expansion.

Similarly, for the lower bound we get

$$\begin{aligned} \binom{cn}{n} &= \frac{(cn)!}{n!(cn-n)!} > \frac{\sqrt{2\pi cn} \left(\frac{cn}{e}\right)^{cn} e^{\frac{1}{12cn+1}}}{\sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{1}{12n}} \cdot \sqrt{2\pi(c-1)n} \left(\frac{(c-1)n}{e}\right)^{(c-1)n} e^{\frac{1}{12(c-1)n}}} \\ &= \sqrt{\frac{c}{(c-1)2\pi n}} \cdot \left(\frac{c^c}{(c-1)^{c-1}}\right)^n \cdot e^{\frac{1}{12cn+1} - \frac{1}{12n} - \frac{1}{12(c-1)n}} \\ &> \sqrt{\frac{c}{(c-1)2\pi n}} \cdot \left(\frac{c^c}{(c-1)^{c-1}}\right)^n \cdot \left(1 - \frac{c^2 - c + 1}{12c(c-1)n}\right). \end{aligned}$$

E Additional Experimental Data

E.1 Competitor Configuration

In Table 4, we give an overview over the most important parameters selected for the competitors. Our benchmark code is available in the supplementary materials.

Table 4 Configurations of competitors

Competitor	Configuration parameters
CHD [7]	Load factor 0.98. $k = 16$ collisions. Bin size 12.
Cuckoo (here, based on [5, 43])	2 alternative positions for each object, loaded in parallel to reduce latency. Streamed queries with <i>await any</i> . Load factor 0.95. Random walk insertion.
LevelDB [28]	No compression. Construction using a single, large write batch. No Bloom filters.
PaCHash (here)	$a = 8$. External blocks store a table of keys and offsets. Streamed queries with <i>await any</i> .
PTHash [46]	“Optimizing the general trade-off” [46] with $\alpha = 0.94, c = 7$, D-D Encoding.
RecSplit [18]	Leaf size $\ell = 8$. Bucket size $b = 2000$.
RocksDB [19]	Block cache disabled. No memory mapping or WAL. Queries use batches of size 64. No Bloom filters.
Separator (here, based on [27, 33])	6 bit separators. Load factor 0.96. Streamed queries with <i>await any</i> .
SILT [35]	<code>testCombi.xml</code> configuration from original repository.
std::unordered_map	8 byte keys. 64 bit pointers to object contents.

E.2 Maximum Load Factors

Both Separator Hashing [27, 33] and Cuckoo Hashing [5, 43] provide very high maximum load factors for identical size objects. Our new implementations of both variants provide a limited support for variable size objects. However, as we can see in Figure 4, the maximum achievable load factor highly depends on the average object size and the distribution of object sizes. PaCHash, in contrast, natively supports arbitrary object sizes and is independent of the distribution.

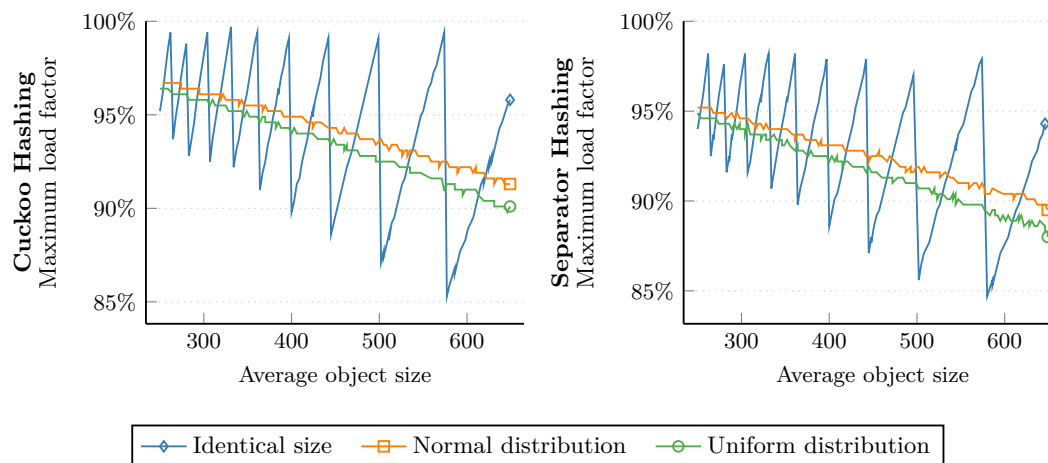
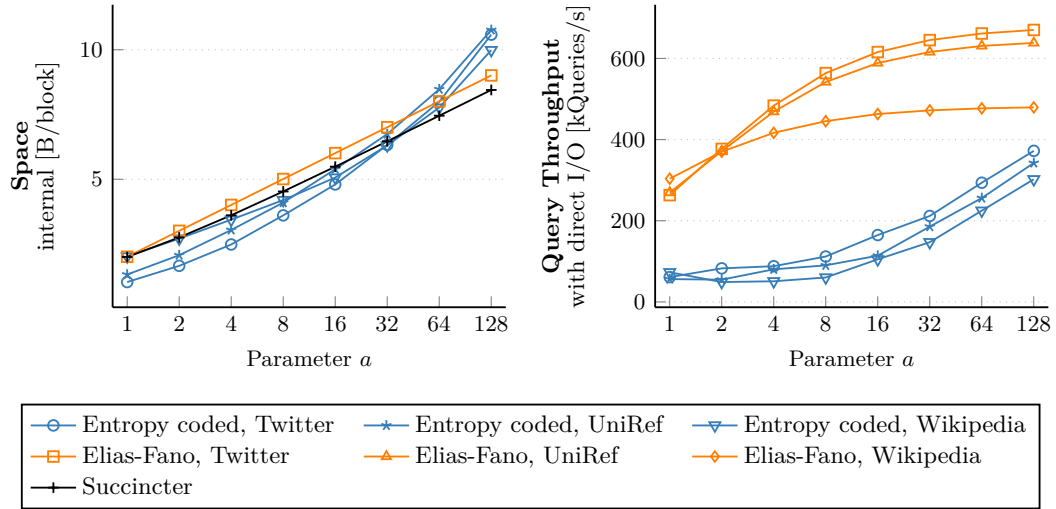


Figure 4 Maximum achievable load factor with different distributions of object sizes of our implementations of Separator Hashing and Cuckoo Hashing. For an average object size s , the normal distribution has a variance of $s/5$ and the uniform random sizes are drawn from $[0.25s, 1.75s]$

E.3 PaCHash Indexes with Real World Data Sets

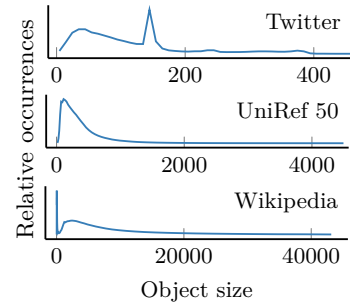
Figure 5 compares throughput and space usage of PaCHash using real world size distributions and different index data structures. The Twitter data set contains tweets from August 1st – August 5th 2021 and has only small objects. The UniRef 50 protein database [49] contains some objects larger than the block size and the LZ4 compressed [12] English Wikipedia from November 2021 contains significantly larger objects. See Table 5 for details.



■ **Figure 5** PaCHash with real world data sets using different index data structures. There is no practical implementation of Succincter [45], so we only give calculated values and no throughput. The space usage of Elias-Fano and Succincter is independent of the object size distribution.

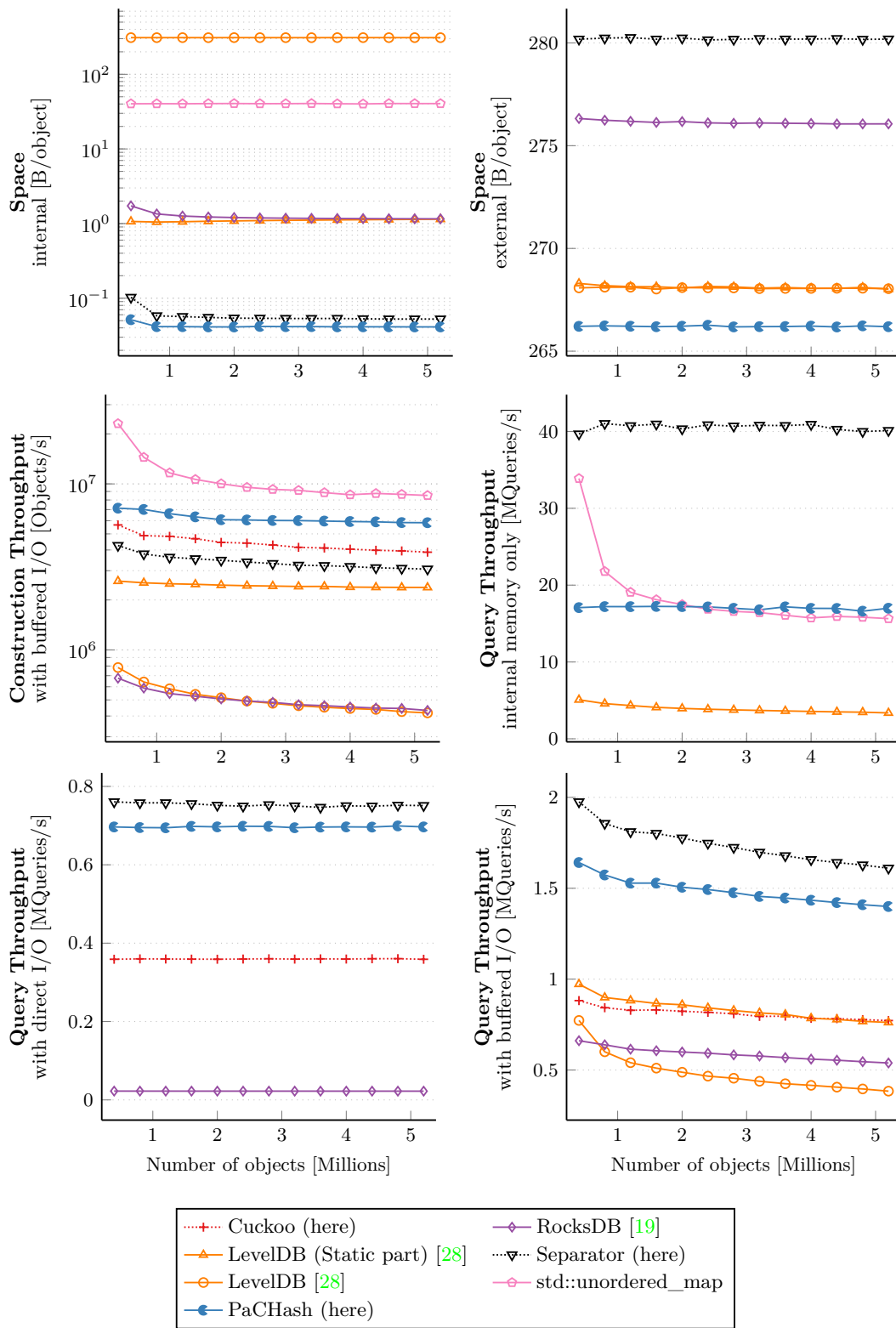
■ **Table 5** Twitter, UniRef, and Wikipedia real world data sets we use for benchmarks. The peak of the Wikipedia data set at ≈ 70 bytes is caused by redirects.

	Twitter	UniRef 50	Wikipedia
Objects n	20 238 968	48 531 431	16 181 427
Average size	115 B	281 B	1731 B
Median size	94 B	194 B	77 B
Maximum size	560 B	45 KB	272 KB
Total size N	2.4 GB	13.2 GB	26.3 GB
Objects $> B$	0%	0.08%	12%



E.4 Comparison with Competitors using Variable Size Objects

In Figure 6, we repeat the measurement from Figure 3 using objects of variable size. This rules out some competitors but overall leads to the same results as for identical size objects.



■ **Figure 6** Comparison with competitor object stores using objects of uniform random size $\in [128, 384]$ bytes. Keys are 8 byte random strings.