

# Take Out the TraChe: Maximizing (Tra)nsactional Ca(che) Hit Rate

Audrey Cheng<sup>†</sup>  
Natacha Crooks<sup>†</sup>

David Chu<sup>†</sup>  
Joseph M. Hellerstein<sup>†</sup>

Terrance Li<sup>†</sup>  
Ion Stoica<sup>†</sup>

Jason Chan<sup>†</sup>  
Xiangyao Yu<sup>‡</sup>

<sup>†</sup>UC Berkeley

<sup>‡</sup>University of Wisconsin–Madison

## Abstract

Most caching policies focus on increasing object hit rate to improve overall system performance. However, these algorithms are insufficient for transactional workloads. In this work, we define a new metric, transactional hit rate, to capture when caching reduces latency for transactions. We present DeToX, a caching system that leverages transactional dependencies to make eviction and prefetching decisions. DeToX is able to significantly outperform single-object alternatives on real-world workloads and popular OLTP benchmarks, providing up to a 1.3x increase in transaction hit rate and 3.4x improvement in cache efficiency.

## 1 Introduction

To improve latency at scale, application developers often layer caching systems, such as Memcached [69] and Redis [2], over standard data stores. These systems traditionally optimize for *object hit rate*, or how often requested objects can be served from cache. Consequently, current caching policies fail to capture the transactional nature of many application workloads. On a production workload from Meta [26], we find that up to 90% of objects cached by least recently used (LRU) and least frequently used (LFU), two popular caching algorithms, do not have any impact on latency despite high object hit rates. Existing policies fail to capture the *all-or-nothing* property of transactions: all objects requested in parallel must be present in cache, or there will be little performance improvement because latency is dictated by the slowest access.

Accordingly, object hit rate is the wrong objective for transactional workloads. Instead, we propose a new metric, *transactional hit rate*, or how often objects requested in parallel can all be served from cache. This metric precisely captures when the cache reduces latency for transactions.

In this paper, we present DeToX, the first high-performance caching system that optimizes for transactional hit rate. In accordance with standard caching algorithms, DeToX assigns scores to objects and evicts those with the lowest values.

As such, its policy is easily adaptable to existing caching systems. To rank objects in the transactional context, DeToX leverages the following insight: objects accessed in parallel within the same transaction should be scored together since they must all be cached to reduce transactional latency.

While scoring keys together might seem simple, the structure of transactional workloads complicates matters. Unlike previous work on caching for parallel jobs [11] and web applications [7, 10, 18, 90, 91], transactions need to be modeled as non-trivial directed acyclic graphs (DAGs) of read and write operations [22, 94]. Crucially, some keys within a transaction are accessed in parallel, but others are not. Consequently, a transaction’s latency is determined by its *critical length*, or the number of sequential accesses on its longest, non-cached path (transactional hit rate captures the reduction of critical length). Rather than considering all keys in a transaction together, we must focus on caching the *groups* of keys that reduce critical length.

Implementing a caching policy based on grouping presents several significant challenges. (1) For an arbitrary transaction, there can be an exponential number of groups, making scoring prohibitively expensive. (2) Identifying groups requires inferring transactional DAGs through static analysis, which may not always be possible. (3) Objects that are accessed by different transactions can belong to different groups, which have varying latency benefits if cached, and we need to capture these disparities.

We address each of these issues in DeToX. (1) To reduce the overhead of an exponential number of groups, we introduce the notion of *interchangeable* keys: if two keys can replace each other in any group and still reduce critical length, then they can be represented by the same group. Interchangeable keys drastically curb the number of groups that need to be scored. (2) When transactional DAGs are not accessible, we propose a simplified policy that dynamically infers groups based on which requests are executed in parallel (termed *levels*). (3) Finally, we account for group membership when scoring keys to ensure these values precisely reflect each object’s contribution to transactional hits.

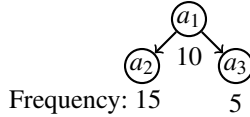


Figure 1: GetLinkedAccounts transaction.

Moreover, while our approach is primarily targeted at eviction, it also enables prefetching (Section 6). Our prefetching policy tracks dependencies within transactions to preemptively bring groups of items into the cache.

Our eviction and prefetching algorithms are implemented in DeToX, which presents a key-value API that supports drivers for Redis [2], Postgres [3], and TiKV [4]. We evaluate our system on real-world workloads from TAOBench [26], a social network benchmark that models Meta’s production workloads, as well as standard OLTP benchmarks (Epinions [37], SmallBank [87], and TPC-C [33]). Compared to single-object caching algorithms and systems, including ChronoCache [45], GDSF [27], LIFE [11], LFU, and LRU, our algorithm can achieve up to a 1.3x increase in *transactional hit rate*, leading to a 3.4x improvement in cache efficiency (defined as the least amount of cache space required to achieve a particular transactional hit rate). For a Redis-Postgres setup, this translates into 31% higher throughput and 30% lower latency.

Our transactional hit rate metric prioritizes *latency* and exposes a new trade-off in caching enabled by the cloud’s elastic resources: optimizing for latency versus reducing system load. In contrast, single-object policies focus on maximizing *object hit rate* to decrease load to the data store but do not always improve transaction request times.

In summary, we make the following contributions:

- We define a new metric, transactional hit rate, to evaluate the latency reduction of caching for transactions (Section 3).
- We provide the first formalization of transactional caching, and we prove that the problem is NP-Hard (Section 3.4).
- We present a new caching system, DeToX, that leverages transactional dependency information to optimize for transactional hit rate and significantly improve performance on popular workloads (Sections 4 – 8).

## 2 Motivation

In this section, we illustrate why single-object eviction algorithms perform poorly for transactional workloads. Specifically, we show that a well-known optimality result in caching does not hold for transactions and that popular caching algorithms achieve low transactional hit rates.

### 2.1 Object Hit Rate is Insufficient

Most existing cache eviction algorithms focus on maximizing object hit rate, or the fraction of single object requests served from cache. However, this approach fails to capture the

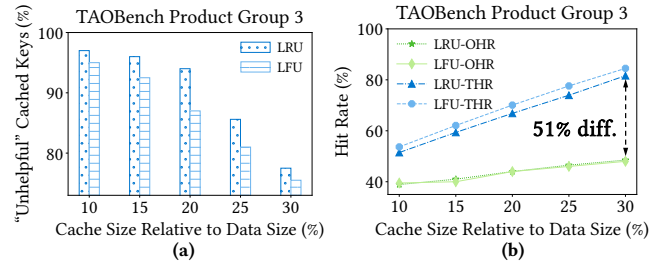


Figure 2: Single-object policy performance.

inter-object dependencies within transactions. Consider for example a simple transaction `GetLinkedAccounts` that returns secondary bank accounts  $a_2$  and  $a_3$  linked with a primary account  $a_1$  (Figure 1). This transaction must first read  $a_1$  before accessing both secondary accounts  $a_2$  and  $a_3$  in parallel. Thus,  $a_1$ ,  $a_2$  and  $a_3$  are all on the longest path of the transaction. If we cache  $a_1$ , we can reduce the end-to-end latency of the transaction. However, if we additionally cache  $a_2$ , the overall latency does not improve because we still need to access  $a_3$  from disk. In fact, caching either  $a_2$  or  $a_3$  individually does not improve performance; transaction latency remains equal to the case in which *neither* key was cached. On the other hand, caching *both*  $a_2$  and  $a_3$  does improve latency.

Transactions have an implicit *all-or-nothing* property on groups of objects that traditional caching algorithms fail to capture. This can lead popular eviction algorithms, such as LRU and LFU, to make poor caching decisions. Consider a situation in which, over all transactions,  $a_2$  is more frequently accessed than  $a_1$  and  $a_3$ . LFU and LRU would choose to evict  $a_1$  and  $a_3$  over  $a_2$ , resulting in no latency improvement for this transaction. In this case, a “hot” (frequently accessed) key  $a_2$  is requested in parallel alongside a “cold” (rarely accessed) key  $a_3$ . If all accesses of  $a_2$  are sent in parallel with requests to different cold keys, there is no benefit to caching  $a_2$  unless all these cold keys are cached. In effect, cold keys can “contaminate” (degrade the cacheability of) hot keys like  $a_2$ .

**Real-world workloads.** This observation is not limited to our simple example: we find that single-object eviction algorithms also perform poorly for complex, real-world workloads. Figure 2a illustrates that over 90% of cached keys do not have any impact on latency (“unhelpful” keys) for the Product Group 3 workload of TAOBench [26]. The root cause is simple: these algorithms optimize for *object hit rate* (OHR) rather than transactional hit rate (THR). As we see in Figure 2b, LRU and LFU achieve high object hit rates but up to 51% lower transactional hit rates. Transactions in this workload access either a combination of hot keys and warm keys, or hot keys and cold keys. Single-object algorithms, which use only individual object features to score keys, retain only hot keys but evict most warm keys and all cold keys. As a result, they achieve few transactional hits. A transactionally-aware policy would instead recognize that cold keys contaminate their associated hot keys and prioritize retaining only the hot and warm keys that are accessed together.

T	Keys accessed	Cache state	Optimal cache state
1	$a_1, a_2, a_3$	-	-
2	$a_4, a_5, a_6$	$a_1, a_2, a_3$	$a_1, a_2, a_3$
3	$a_4, a_5, a_7$	$a_1, a_4, a_5$	$a_1, a_2, a_3$
4	$a_1, a_2, a_3$	$a_1, a_4, a_5$	$a_1, a_2, a_3$

Figure 3: Non-optimality of Belady.

## 2.2 Optimality

Our observations also have theoretical implications. We find that Belady [16], the offline, optimal eviction algorithm for uniformly-sized objects does not make the best decisions for maximizing transactional hit rate. This policy evicts keys that are accessed furthest in the future but fails to take into account whether these keys generate transactional hits.

We prove that Belady is *not* optimal even for the simplest case of uniformly-sized transactions with uniformly-sized objects (Figure 3). In this example, we have four transactions with a cache size of 3.  $T_1$  and  $T_4$  access keys  $a_1, a_2, a_3$ , while  $T_2$  accesses  $a_4, a_5, a_6$  and  $T_3$  accesses  $a_4, a_5, a_7$ . Belady chooses to first cache  $a_1, a_2, a_3$  and then replaces the last two keys with  $a_4, a_5$  since these keys give object hits (but no latency reduction) for  $T_3$ . However, keeping  $a_2, a_3$  in the cache would lead to a transactional hit (and latency improvement) for  $T_4$ .

## 2.3 Towards a new approach

Our results highlight how single-object caching strategies yield low transactional hit rates by storing many unhelpful objects. Web caching algorithms suggest a way forward: they acknowledge the need to cache multiple objects together (e.g., page-level hit ratio) but only consider flat dependencies [11, 91]. In contrast, transactions can have complicated topologies with multiple levels of dependencies.

To develop a transactionally-aware caching system, we must address three challenges: (1) formalizing caching in the transactional context, including optimality analysis (Section 3), (2) identifying which groups of objects lead to transaction hits, given the potentially complex structure of transactions (Section 4.1), and (3) scoring the individual objects in these groups to determine which objects to store in the cache (Section 4.2). In our design, we are careful to emphasize compatibility with existing caching systems, such as Memcached and Redis, so that our approach can be easily implemented for greater applicability.

## 3 Transactional Caching

In this section, we formalize the transactional caching problem. We define a new metric, *transactional hit rate*, to capture the latency reduction of caching transactions.

```

1 id = SELECT cId FROM ACCOUNTS WHERE name = cName
2 s = SELECT savings FROM SAVINGS WHERE cId = id
3 c = SELECT checking FROM CHECKING WHERE cId = id
4 return s + c

```

Listing 1: Code for Figure 4. The dependencies for Lines 2 and 3 on the output of Line 1 are highlighted in red.

## 3.1 Transactions

Transactions consist of read and write requests that must be applied atomically [22]. Some of these operations are independent and can execute in parallel, while others are *dependent* on the result of preceding operations. For instance, a read operation may query a key determined by the return value of a previous operation. As a result, these operations must be run sequentially. In effect, transaction execution can be captured by a DAG of operations. More formally, we apply the notion of a logical dependency, generalizing the model from Wu et al. [94]:

**Definition 1** (*Logical dependency*). Given two operations  $tp$  and  $t$  of a transaction, an operation  $t$  is *logically dependent* on operation  $p$  if  $p$  determines the key or value accessed by  $t$ .

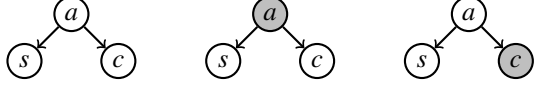
Traditionally, these dependencies are not captured by the system, which observes only sequences of reads and writes. In practice, these relationships can be captured statically through program analysis or specified at run time by the developer. Together, operations and logical dependencies define a transaction execution graph:

**Definition 2** (*Transaction execution graph*). For transaction  $T$ , a *transaction execution graph*  $G = (V, E)$  is a DAG, where each vertex in  $V$  represents a pair  $(x, X)$  of a read or write operation to key  $x$  in table  $X$ , and each edge in  $E$  represents a logical dependency between operations.

Each transaction execution graph corresponds to a transaction type:

**Definition 3** (*Transaction type*). Transactions of the same type have identical execution graphs when only considering tables.

We infer transaction execution graphs and their resulting types through static analysis, as done in prior work [36, 94]. Note that we only extract table accesses and graph structure; the individual keys accessed by transactions are known only at run time. As such, we make no assumptions about the DAG structure and support general-purpose, interactive transactions. For example, the SmallBank workload [87] contains the transaction types: Amalgamate, Balance, DepositChecking, SendPayment, TransactSavings, and WriteCheck. For the Balance transaction (Listing 1), requests to both the Savings ( $S$ ) and Checkings ( $C$ ) tables are dependent on the result of the read to the Accounts table ( $A$ ). The corresponding execution graph consists of three nodes, one for each operation, and logical dependencies  $r[A] \rightarrow r[S]$  and  $r[A] \rightarrow r[C]$ . While the reads to  $S$  and  $C$  are independent



(a) Cache state:  $\{\}$  (b) Cache state:  $\{a\}$  (c) Cache state:  $\{c\}$

Figure 4: SmallBank Balance transactions.

and can be executed in parallel, they cannot proceed until after the read to  $A$  finishes. At run time, a Balance transaction that reads the keys  $a, s, c$  from the tables  $A, S, C$  respectively can be mapped onto the same execution graph (Figure 4a).

## 3.2 Cache

The previous section presents the notion of a transaction, including the logical dependencies that constrain its execution. We now formalize how *caching* affects transactions, drawing from Abrams et al. [48] for notation.

**Definition 4 (Cache state).** A cache state is a set of keys  $C$  for which  $|C| \leq n$ , where  $n$  is the capacity of the cache.

In line with prior work [29], we assume that the cache state does not change for the duration of each transaction.

By assumption, objects are served with lower latency from the cache than from the underlying data store. We make the simplifying assumption that requests served from the cache have *zero* latency for notational simplicity (we explore the effects of varying cache latency in Section 8.6). Under this model, transaction latency is defined by the number of sequential, *non-cached* accesses. This corresponds to the longest path in the transaction’s execution graph  $G$ , excluding vertices with cached keys.

We formalize this notion as the critical length:

**Definition 5 (Critical length).** Given a transaction  $T$  with transaction execution graph  $G$ ,  $K$  number of keys, and cache state  $C$ , the critical length is the length of the longest path from any source vertex (no incoming edges) to any sink vertex (no outgoing edges), excluding vertices corresponding to keys in  $C$ . We define the function  $L: G \times 2^K \rightarrow \mathbb{N}$  for which  $2^K$  is the powerset of all keys, such that  $L(G, C)$  is the critical length.

Given a transaction  $T$  with execution graph  $G$ ,  $L(G, \{\})$  represents the length of the longest path in  $G$  when the cache is empty. For example, Figure 4a has longest paths  $\{r[a], r[c]\}$  and  $\{r[a], r[s]\}$  with critical length  $L(G, \{\}) = 2$ . Caching key  $a$  (Figure 4b) would shorten the critical length to  $L(G, \{a\}) = 1$ , as the longest paths are reduced to  $\{r[c]\}$  and  $\{r[s]\}$ . However, caching key  $c$  (Figure 4c) does not change the critical length, since  $\{r[a], r[s]\}$  remains the longest path with  $L(G, \{c\}) = 2$ . Informally, we refer to each length reduction as a *transactional hit*.

## 3.3 Transactional Hit Rate (THR)

Having defined the necessary formalisms for transaction latency and caching, we can now introduce transactional hit rate. Informally, this metric captures how much latency improves

when caching for transactions, much like how its single-object counterpart, object hit rate, does so for individual requests.

We first present THR in the context of a single transaction:

**Definition 6 (Individual transactional hit rate).** Given transaction  $T$  with execution graph  $G$  and cache state  $C$ , the individual transactional hit rate is  $\frac{L(G, \{\}) - L(G, C)}{L(G, \{\})}$ .

The difference in critical length represents the reduction in sequential, non-cached accesses after caching. We normalize this difference by dividing by the total critical length. This metric captures the impact of caching for the execution of a single transaction (note that if the transaction execution graph is a sequential list of dependent reads, then transactional hit rate is equivalent to object hit rate). We easily extend this definition to a sequence of transactions:

**Definition 7 (Transactional hit rate).** Given a sequence of transactions  $T_1, T_2, \dots, T_m$  with execution graphs  $G_1, G_2, \dots, G_m$  and the respective cache states at the time of execution  $C_1, C_2, \dots, C_m$ , the transactional hit rate is  $\frac{\sum_{i=1}^m (L(G_i, \{\}) - L(G_i, C_i))}{\sum_{i=1}^m L(G_i, \{\})}$ .

## 3.4 Optimality Analysis

Single-object caching is a well-studied problem and is known to be NP-Hard in the general case [29]. We show that the optimal transactional caching is NP-Hard through a reduction from variable-sized caching of single objects (proof in Appendix A). In summary, we reduce each variable-sized object of size  $X$  to a transaction with  $X$  unit-sized operations.

## 4 Group Identification and Scoring

Designing an optimal caching policy is impractical for transactional caching, since it would run in exponential time. Unfortunately, traditional heuristics perform poorly for transaction hit rate (Section 2) because they fail to identify the keys that must be cached as a *group* in order to yield a transactional hit. This notion of grouping is central to developing a transactionally-aware caching policy. We proceed in two steps: first, we identify which groups of keys lead to transactional hits when cached together (*group identification*). Next, we determine what scores should be assigned to each key within a group (*group scoring*).

Figure 5 gives an overview of DeToX. Our system first extracts transaction execution graphs (Section 3) from application code and identifies groups of table accesses (Section 4.1) at compile time. The number of groups that DeToX needs to consider can be reduced at compile time through the notion of *interchangeability* (Section 5.1). DeToX then scores groups based on key accesses at run time (Section 4.2). If application code is not available, DeToX constructs approximate groups at run time by using *levels* (Section 5.2).

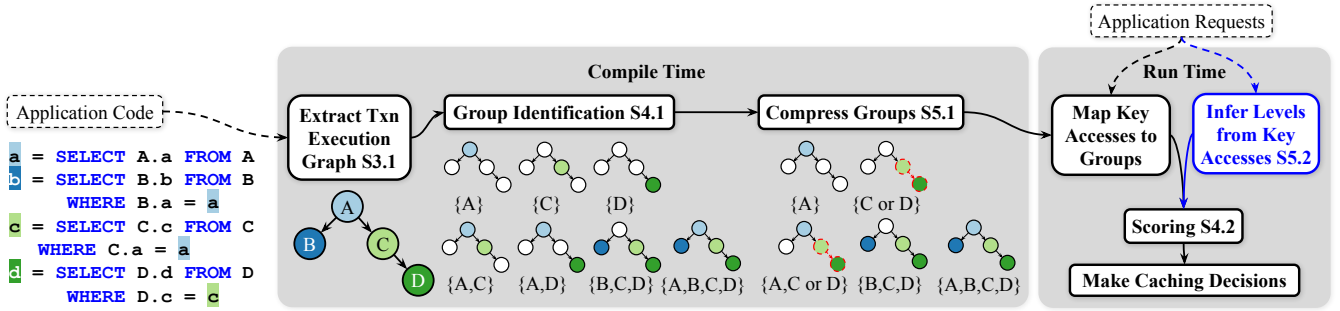


Figure 5: Overview of DeToX (gray boxes). Blue edges on the right represent the path taken if application code is not available.

## 4.1 Group Identification

Intuitively, a group is a set of keys that reduces critical length if cached together. Specifically, we define a *complete group* as one from which we cannot remove any key without increasing critical length. Completeness optimizes cache efficiency by storing the minimal subset of keys necessary to reduce latency. Formally:

**Definition 8 (Complete group).** Given a transaction  $T$  and its execution graph  $G$ , a complete group is a subset of keys  $g$  accessed in  $T$  such that  $\forall g' \subset g, L(G, g) < L(G, g')$ .

We identify complete groups of table accesses at compile time, using the transaction execution graphs  $G$  extracted via static analysis, as seen in Figure 5. A simple algorithm to identify groups is to iterate through the powerset of possible table accesses and compute their resulting reductions in critical length. These table accesses are replaced at run time with key accesses. The application passes along metadata with requests to indicate the corresponding vertex in the transaction execution graph of each key access.

Consider Figure 6a, which has a critical length of three (serial accesses of  $a, c, d$ ) and seven complete groups ( $\{a\}, \{c\}, \{d\}, \{a, c\}, \{a, d\}, \{b, c, d\}, \{a, b, c, d\}$ ). Note that  $\{c, d\}$  is not a complete group. If  $c$  and  $d$  are both cached, then the critical length is two (serial accesses of  $a, b$ ). However, only caching  $c$  already yields the same critical length (accesses to  $a, b, d$ ). Similarly,  $\{a, b\}$  is not a complete group, because it yields a critical length of two (serial accesses  $c, d$ ), which could also be achieved by just caching  $a$ .

In the worst case, the number of complete groups can be exponential in the size of the transaction, even for simple transaction topologies. Fortunately, many of these groups are, in fact, *equivalent*. We describe this notion more precisely in Section 5.1 and present an optimization that drastically reduces the number of groups that need to be considered.

## 4.2 Scoring

Caching policies typically assign scores to keys and evict keys with lower values. We adopt the same strategy by mapping complete groups to individual key scores at run time, as seen in Figure 5. This approach has two benefits: (1) we can draw from prior work on single-object caching

Parameter	Description
$\text{SCORE\_G}(\text{group})$	Score of a group
$F_{\text{group}}$	Set of all key frequencies in a group
$L_{\text{group}}$	Transactional hits of a group
$S_{\text{group}}$	Sum of key sizes in a group
$\text{SCORE\_K}(\text{key})$	Score of a key
$TS_{\text{key}}$	Sum of instance scores for a key
$F_{\text{key}}$	Frequency of a key
$A_{\text{global}}$	Global aging factor

Table 1: Scoring parameters.

algorithms, and (2) we minimize implementation changes needed for real-world caching systems.

### 4.2.1 Scoring a Group in a Single Transaction

We begin by assigning numerical scores to each group (*group scores*) with higher values representing groups that are more beneficial to cache. We draw inspiration from GDSF, a high-performing web caching algorithm [27]. GDSF considers three metrics to score keys: frequency (access count), recency, and size. Specifically, GDSF uses the following formula:  $\text{SCORE}_{\text{GDSF}}(\text{key}) = F_{\text{key}}/S_{\text{key}} + A_{\text{global}}$ , where  $F_{\text{key}}$  is frequency of the key,  $S_{\text{key}}$  is size of the key, and  $A$  is a global recency factor (described in Section 4.2.3). GDSF gives equal weight to each of these factors, and we follow this approach. We leverage frequency and size to score each group as follows (and incorporate recency into key scores in Section 4.2.3):

$$\text{SCORE\_G}(\text{group}) = \frac{\min(F_{\text{group}}) \times L_{\text{group}}}{S_{\text{group}}}$$

$F_{\text{group}}$  is a list of all key frequencies in the group.  $L_{\text{group}}$  is the number of transactional hits generated if this group is cached.  $S_{\text{group}}$  is the sum of all key sizes in the group. All scoring parameters can be found in Table 1. For the transactions in Figure 6 (which will be used as running examples), the group scores of each complete group for these transactions are shown in Figures 6b and 6d. The transaction in Figure 6a has keys  $a, b, c, d$  with frequencies of 1, 29, 99, and 50, respectively and sizes of 1. The score of group  $\{a, b, c, d\}$  is thus  $\frac{\min(1, 29, 99, 50) \times 3}{4} = 0.75$ .

**Frequency ( $F_{\text{group}}$ ).** Keys within a complete group may vary in frequency but must all be cached to yield a transactional

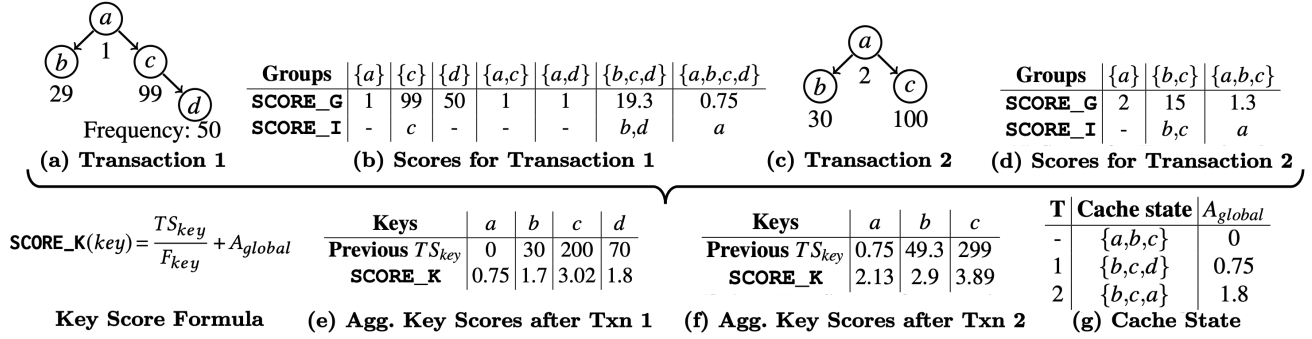


Figure 6: Example transactions and scores. Key sizes are 1, cache size is 3, and  $A_{global}$  starts at 0. The cache initially stores  $\{a,b,c\}$ .  $a$  is evicted after T1, and  $d$  is evicted after T2, with  $A_{global}$  updated on each eviction.

hit. For example, if a high-frequency key  $x$  is only associated with a group of keys  $\{y_1, \dots, y_k\}$  (each with much lower frequency than  $x$ ), then it is not beneficial to cache  $x$ . Essentially, the key with the minimum frequency determines the cacheability of the entire group. Thus, we take the minimum of all key frequencies in calculating the group score. Consider for instance the transaction in Figure 6c: key  $b$  is less frequently accessed than key  $c$  and drives down the frequency of the group  $\{b,c\}$  to  $\min(F_{group}) = \min(30, 100) = 30$ . In this example,  $b$  contaminates  $c$ .

**Critical length reduction ( $L_{group}$ ).** This parameter captures the reduction in critical length when caching a group ( $L_{group} = L(G, \{ \}) - L(G, group)$ ). Other factors being equal, groups with greater reductions are better choices to cache and should thus be assigned a higher score.

**Size ( $S_{group}$ ).**  $S_{group}$  represents the cache space needed to store the group. Since all keys in a group must be present in cache to generate a transactional hit, THR is maximized by retaining groups of smaller sizes (more groups can be cached).

Next, we describe how to go from group scores to key scores.

#### 4.2.2 Scoring Across Groups in a Single Transaction

Mapping group scores to keys is challenging: for a given transaction, a key can belong to multiple complete groups, each with a separate group score (SCORE\_G). In this section, we focus on assigning scores to keys within a single transaction; we assign each key an instance score (SCORE\_I) based on one of its group scores. We combine instance scores across transactions in Section 4.2.3.

Our algorithm leverages the insight that out of all the keys in a transaction, the highest-scoring complete group is the most beneficial set of keys to cache. Thus, our protocol first finds the complete group with the highest group score SCORE\_G and sets the instance score of all keys in that group to SCORE\_G. In Figure 6b,  $\{c\}$  has the highest group score (SCORE\_G=99), so  $c$  is assigned the instance score of 99. We then score the remaining keys of the transaction assuming that keys in the highest-scoring group will be cached.

In subsequent iterations, our algorithm finds the highest-scoring complete group that is a *superset* of all keys that

have been assigned instance scores. In Figure 6b, having scored  $c$ , the highest-scoring complete group that subsumes  $c$  is  $\{b,c,d\}$ , with a group score of 19.3. The unscored keys ( $b,d$ ) are then assigned the score of this complete group (19.3). Intuitively, this is the next set of keys that should be retained assuming that the highest-scoring complete group is already in cache. Our algorithm captures the fact that, once  $c$  is cached,  $d$  should only be cached when  $b$  is cached. The low score of  $b$  contaminates  $d$  but should not contaminate  $c$  (since  $c$  by itself can lead to a transactional hit).

The iterative process described above is repeated until all keys are scored. For our example, the next highest-scoring complete group that is a superset of  $\{b,c,d\}$  is  $\{a,b,c,d\}$ , with a group score of 0.75, which is assigned to key  $a$ , completing the scoring protocol for Figure 6b. Note that all keys will eventually be scored by this algorithm, since they are all part of the trivial complete group containing every key in the transaction.

#### 4.2.3 Scoring Across Transactions

Finally, we describe how to integrate instance key scores across multiple transactions into an *aggregate* value. This final score will be used by the system to decide which keys to evict from the cache. We adopt the following formula:

$$\text{SCORE\_K}(key) = \frac{TS_{key}}{F_{key}} + A_{global}$$

$TS_{key}$  is the sum of all instance scores from Section 4.2.2 across all transactions accessing this key.  $F_{key}$  is the frequency of this key.  $A_{global}$  is the global aging factor.

**Averaging instance scores.** To combine instance key scores into a single value for a given key, we take the running average of these scores. Each time a key is accessed, we add its instance score to the total score  $TS_{key}$  and increment  $F_{key}$  before calculating a new aggregate score. Figure 6e gives the key scores of  $a,b,c,d$  after the execution of the transaction in Figure 6a, assuming that the aging factor is initialized to 0, key size is 1, and the previous  $TS_{key}$  values are 0, 30, 200, and 70 respectively. For example,  $c$  has an instance score of 99 (Figure 6b) for the transaction in Figure 6a, a previous  $TS_{key}$  of 200, and frequency of 99, giving  $\text{SCORE\_K}(c) =$

$\frac{200+99}{99} + 0 = 3.02$  in Figure 6e. Taking the average allows us to account for contamination between different groups.

**Recency.** To account for shifts in object access distributions over time, GDSF, along with other algorithms [8], uses an aging factor to capture object recency. Since previously popular objects can remain in the cache for extended periods of time (due to their high frequencies) and prevent newly popular objects from being stored, the scores of more recent objects should be higher than those of older objects. Towards this end, GDSF applies  $A_{global}$ , a global value that is added to the score of a key upon each access to increase the scores of more recently accessed objects and age older objects out of cache. The value of  $A_{global}$  is updated each time an object is evicted and set as that object’s score. Thus, the factor increases monotonically and ensures that all accesses after this eviction will have scores higher than the last evicted key. In essence, this factor acts as a “reset” on key scores. In Figure 6,  $a$  is evicted after the transaction in Figure 6a executes, and  $A_{global}$  is set to  $a$ ’s score (0.75). This value is then added to  $SCORE\_K$  for each key accessed in the subsequent transaction (Figure 6c). For example,  $c$  has an instance score of 15 (Figure 6d), a previous  $TS_{key}$  of 299, frequency of 100, and  $A_{global}$  of 0.75, giving  $c$  an aggregate key score of  $SCORE\_K(c) = \frac{299+15}{100} + 0.75 = 3.89$  in Figure 6f.

## 5 Optimizations

While our current approach precisely captures the cacheability of each group, it can be prohibitively expensive when the number of complete groups is exponential for some transaction topologies. We address this problem in two ways. First, we observe that many complete groups capture redundant information and introduce *interchangeable groups* to avoid scoring all complete groups, reducing run time overhead. Second, we present a restricted form of grouping, *levels*, that dynamically approximates groups at run time. This technique also enables us to score keys when we do not have access to transaction code (i.e., we do not know the transaction execution graphs).

### 5.1 Interchangeability

We find that the number of complete groups can be exponential with respect to transaction size, even for simple topologies. For example, the TPC-C *Order-Status* transaction in Figure 7a has a depth of three, and the number of complete groups for this transaction is exponential with respect to its depth:  $\{c\}$ ,  $\{o\}$ ,  $\{ol_1, ol_2\}$ ,  $\{c, o\}$ ,  $\{o, ol_1, ol_2\}$ ,  $\{c, ol_1, ol_2\}$ ,  $\{c, o, ol_1, ol_2\}$  make up  $2^3 - 1 = 7$  complete groups.

We observe that transactions often contain complete groups that differ by only a single key. For instance, for every group in which  $c$  is present in Figure 7a, there exists an identical group in which  $o$  replaces  $c$  (and vice-versa). In effect, these keys can be “swapped” with each other and still produce a

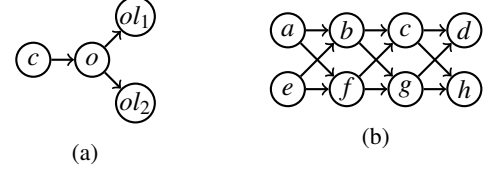


Figure 7: Transactions to demonstrate interchangeability. Figure 7a is a TPC-C *Order-Status* transaction.

complete group. This interchangeability property is powerful: if two keys can be exchanged in any complete group, then deciding to cache one key over the other is entirely dependent on the *individual* scores of these keys, as all other parameters are shared. Consequently, we do not need to calculate the scores of each their complete groups in order to score each key. Consider the groups  $\{c, ol_1, ol_2\}$  and  $\{o, ol_1, ol_2\}$  for the TPC-C *Order-Status* transaction in Figure 7a, assuming  $c$  has a higher individual score than  $o$ . Since  $c$  and  $o$  are interchangeable, we know that  $\{c, ol_1, ol_2\}$  must have a higher group score than  $\{o, ol_1, ol_2\}$ , as all other parameters (the scores of  $ol_1$  and  $ol_2$ ) are shared. Our scoring algorithm favors caching groups with higher scores, so we can avoid calculating the score of  $\{o, ol_1, ol_2\}$  at run time while determining the score for  $o$ .

We can further generalize the idea of interchangeability to *sets* of keys that can also be “swapped” with each other. Continuing the example above, the set of keys  $\{ol_1, ol_2\}$  is interchangeable with  $\{c\}$ , because any complete group that contains  $\{ol_1, ol_2\}$  will remain a complete group if  $\{ol_1, ol_2\}$  is swapped with  $\{c\}$ . We call such sets interchangeable groups:

**Definition 9 (Interchangeable groups).** Let  $s_1$  and  $s_2$  be distinct sets of keys in a transaction with execution graph  $G$ . We define  $s_1$  and  $s_2$  to be interchangeable if

- (1)  $\forall$  complete groups  $g_1$  such that  $s_1 \subseteq g_1$  and  $s_2 \cap g_1 = \emptyset$ ,  $g'_1 = g_1 \setminus s_1 \cup s_2$  is also a complete group and  $L(G, g_1) = L(G, g'_1)$ , and
- (2)  $\forall$  complete groups  $g_2$  such that  $s_2 \subseteq g_2$  and  $s_1 \cap g_2 = \emptyset$ ,  $g'_2 = g_2 \setminus s_2 \cup s_1$  is also a complete group and  $L(G, g_2) = L(G, g'_2)$ .

Like complete groups, interchangeable groups of table accesses can be identified at compile time, as seen in Figure 5. Key accesses are mapped to the vertices at run time. Computationally, interchangeability allows us to reduce the number of complete groups that need to be scored. We compress the representation of complete groups and reduce run time complexity of the scoring algorithm as follows, using Figure 7b as a running example:

• (Compile time) **Find** all interchangeable groups of vertices from the set of complete groups. The complete groups are:  $\{a, e\}$ ,  $\{b, f\}$ ,  $\{c, g\}$ ,  $\{d, h\}$ ,  $\{a, e, b, f\}$ ,  $\{c, g, b, f\}$ ,  $\{d, h, b, f\}$ ,  $\{a, e, c, g\}$ ,  $\{a, e, d, h\}$ ,  $\{c, g, d, h\}$ ,  $\{a, e, b, f, c, g\}$ ,  $\{a, e, b, f, d, h\}$ ,  $\{a, e, c, g, d, h\}$ ,  $\{c, g, b, f, d, h\}$ ,  $\{a, e, b, f, c, g, d, h\}$ . Consider replacing  $\{a, e\}$  with  $\{d, h\}$  in any complete group; the resulting group is still complete. Thus,  $\{a, e\}$  and  $\{d, h\}$  are interchangeable. Using the same

logic, we find that  $\{a,e\},\{b,f\},\{c,g\},\{d,h\}$  are all mutually interchangeable.

- (Compile time) **Compress** complete groups. Denote an access to any one of the mutually interchangeable groups— $\{a,e\},\{b,f\},\{c,g\},\{d,h\}$ —as  $[C]$ . For example,  $\{a,e,b,f,d,h\}$  becomes  $[C,C,C]$ . In this particular example, all groups of size four can be written as  $[C,C]$ , groups of size six as  $[C,C,C]$ , and groups of size eight as  $[C,C,C,C]$ . We call these representations compressed groups.

- (Run time) **Score** compressed groups by replacing vertices with individual keys in each group. Recall from Section 4.2.2 that our instance scoring algorithm scores all complete groups before greedily selecting the highest-scoring ones. With interchangeability, we no longer need to score all complete groups. Assume the minimum scores of the following interchangeable groups are:  $\{a,e\} : 1, \{b,f\} : 10, \{c,g\} : 30, \{d,h\} : 50$ . Since we know that  $\{a,e\}$  and  $\{d,h\}$  are interchangeable and that  $\{d,h\}$  has a higher score, for any complete group containing  $\{a,e\}$ , there must be another complete group containing  $\{d,h\}$  that has the same (or higher) score. Applying this intuition, the highest-scoring complete group corresponding to the compressed group  $[C,C]$  must be composed of the highest and second-highest-scoring interchangeable groups,  $\{d,h\}$  and  $\{c,g\}$  respectively.

In this example, interchangeability decreases the number of groups that need to be considered at run time from fifteen to four. Overall, interchangeability drastically reduces the number of complete groups that must be scored, lowering run time overhead.

## 5.2 Levels

For cases when we do not have access to transaction code, we design a simplified protocol to dynamically infer groups. We first define a *level* to be a set of keys in a transaction that are sent to the data store in parallel; a similar definition is used to group tasks to optimize caching for parallel job execution [11]. In practice, many applications batch parallel reads to the caching system, which often provides an explicit API to support these requests [2]. We assume that applications send requests as soon as their logical dependencies are fulfilled. For instance, the transaction in Figure 6a has levels  $\{a\}, \{b,c\}$ , and  $\{d\}$ . We have  $d$  as a standalone level since it can only be requested once the level containing both  $b$  and  $c$  has finished executing.

Levels produce identical results to our previous grouping strategies for transactions in which all keys and groups are interchangeable (e.g., Figures 7a and 7b). Many real-world workloads are comprised of such transactions (including all the ones we evaluate in Section 8). When transactions do *not* have these properties, levels can miss out on performance opportunities since they only capture a subset of all possible complete groups. For example, in Figure 6a,  $b$  and  $c$  are always scored together under levels, lowering  $c$ 's score. To

maximize transactional hits,  $b$  should instead be scored with  $d$  since both are colder keys, and  $c$  should be given a high score because caching just this key is likely to lead to a transactional hit. We measure the tradeoff between different grouping strategies in Section 8.

## 6 Prefetching

Prefetching is a popular technique to reduce the client-perceived latency of requests by caching items before they are requested [10, 24, 44, 45, 90]. We revisit this strategy in the context of transactions and design a new prefetching algorithm that uses logical dependencies to minimize latency.

Our policy leverages conditional probabilities: once key  $a$  is accessed, it may be very likely that key  $b$  will also be requested in the same transaction. Consider for example `GetLinkedAccounts` in Figure 1: the access to a primary account is almost always followed by requests to the same subsidiary accounts. Our prefetching algorithm tracks these correlated accesses and preemptively brings dependent objects into the cache ( $a_2$  and  $a_3$  are requested alongside the read to  $a_1$ ). Specifically, DeToX stores, for every request  $r$ , sets of keys in subsequent accesses that are logically dependent on  $r$ . DeToX also tracks the frequency of each set and preemptively fetches in the most popular set into cache alongside  $r$ . To bound memory overheads, we restrict the number of dependency sets that can be stored per key and set a frequency threshold below which we do not retain prefetching metadata.

## 7 Implementation

In this section, we describe our implementation of DeToX, which consists of 7K lines of Java. We adopt a standard two-tier architecture in which we layer a Redis (7.0) cache on top of a data store (Postgres (12.10) and TiKV (5.4.3) are supported). A shim layer routes requests, manages concurrency control, and enables prefetching.

### 7.1 Shim Layer

All client requests are directed to our shim layer, which mediates accesses to the cache and data store to support serializable transactions. Read requests go first to Redis. In the absence of a cache hit, the shim forwards the request to the data store and updates the cache with the result. All writes are sent directly to the data store. While our shim layer currently supports a key-value API, we can convert SQL queries to this format, as previous systems have done [34, 35, 47, 58–61, 65–67, 78, 81–84, 92, 98]. We choose to implement a stand-alone shim layer since there is limited open-source support for concurrency control between caching systems and data stores [9, 43, 45, 46, 75, 76, 85, 88]. Furthermore, our shim layer allows us to



easily plug in different systems. We will explore integrating transactional caching directly into systems in future work.

**Concurrency control.** We implement two-phase locking [22] with timeout-based deadlock detection in the shim layer to ensure serializability. The system maintains the following invariant: values in the cache will either 1) reflect the value committed in the (serializable) data store or 2) be protected by an exclusive write lock.

To achieve this, the shim acquires locks on individual objects before sending requests to either storage system. Writes are buffered at the shim layer until commit. Once values are committed in the data store, they are updated in the cache before write locks are released. To handle crashes, we rely on the data store as the source of truth, similar to previous work [46, 75, 76], and we clear the cache after failures to prevent stale reads. We view applying transaction caching to multiversioned systems as a promising avenue for future work.

**Extracting transaction types and execution graphs.** We leverage prior work [36, 94] to obtain transaction execution graphs with table accesses from application code. The widespread adoption of JDBC-style drivers presents a common interface for extracting transactions across applications.

## 7.2 Eviction

Our eviction policy scores keys as a function of their groups as well as their frequency, size, and recency. The latter three are all features that are already available in Redis, which natively supports LRU and LFU. We reuse these metrics to minimize code changes when implementing our algorithm. We make two primary modifications to Redis: we add (1) a global aging factor that is updated during eviction (as detailed in Section 4.2.3) and (2) support for scoring groups of keys. Specifically, we modify the existing method Redis provides for fetching multiple objects to delineate which keys are accessed together. We update key scores only after a transaction has completed so that we have sufficient information to calculate all group scores. Our changes involve less than 100 lines of code and suggest that DeToX can be easily integrated into any caching system. We also implement a trace-driven simulator in Python to evaluate the offline Belady and Transactional Belady algorithms.

## 8 Evaluation

In this section, we evaluate DeToX against existing caching policies on a range of different workloads. Specifically, we aim to answer the following questions:

- How does DeToX compare to single-object algorithms in terms of transactional hit rate and cache efficiency?
- What is the impact of our grouping techniques?
- What is the tradeoff between optimizing for object hit rate and transactional hit rate?

## 8.1 Experimental Setup

We run our shim layer and Postgres on separate c5a.4xlarge Amazon EC2 instances (16 CPUs, 32GB RAM) and use a memory-optimized r5.4xlarge machine (16 CPUs, 128GB RAM) for Redis. Clients run on c5a.16xlarge instances (64 CPUs, 128GB RAM). We host all machines in the same region with low network latency (0.2ms). For our experiments, we report the average of three 5-minute runs with 60 seconds of warm-up time. When an eviction is needed, we score 10 random samples and choose one to evict among these candidates. This strategy removes the overhead of maintaining a sorted list of keys without degrading performance and is popular in many caching systems [2, 79], including Redis.

**Benchmarks.** We evaluate DeToX against single-object baselines as well as the policies developed in PACMan [11] (their LFU-F is equivalent to our LFU; we evaluate their LIFE algorithm) and ChronoCache [45], a state-of-the-art prefetching system that leverages transactional dependencies. We measure performance on a range of workloads. **TAOBench** [26] is an open-source social network benchmark based on Meta’s production traces. We run the Product Group 1, 2, and 3 workloads, which represent distinct sets of (anonymized) applications at Meta that share data and use the same product infrastructure. All workloads are read-heavy and skewed, typical of most social networks. They contain point reads and writes (inserts, updates, and deletes) as well as read-only and write-only transactions. All transactions are “flat” (they contain no logical dependencies). Since transaction code is not available for this benchmark, we use levels to score groups for eviction.<sup>1</sup> We run experiments with 100M objects for a total data size of around 1 TB. **Epinions** [37] consists of nine transaction types that represent behavior observed on a consumer reviews website. We run the benchmark with 2M user and 1M items for a total data size of roughly 1 TB. **SmallBank** [87] contains six types of transactions that model a simple banking application. We configure it to run with 500M (uniformly accessed) accounts (total size of 1 TB). **TPC-C** [33], a standard OLTP benchmark, simulates the business logic of e-commerce suppliers with five types of transactions. We configure TPC-C to run with 100 warehouses (total size of 8GB). In line with prior transactional key-value stores [34, 81], we use a separate table as a secondary index on the Order table to locate a customer’s latest order in the Order-Status transaction, and on the Customer table to look up customers by their last names (for the Order-Status and Payment transactions).

## 8.2 Application Benchmark Results

We show THR over different cache sizes for all benchmark workloads in Figures 8 and 10. We omit some throughput and

<sup>1</sup>TAOBench [26] chooses to model workloads using probability distributions rather than fixed query types for adaptability.

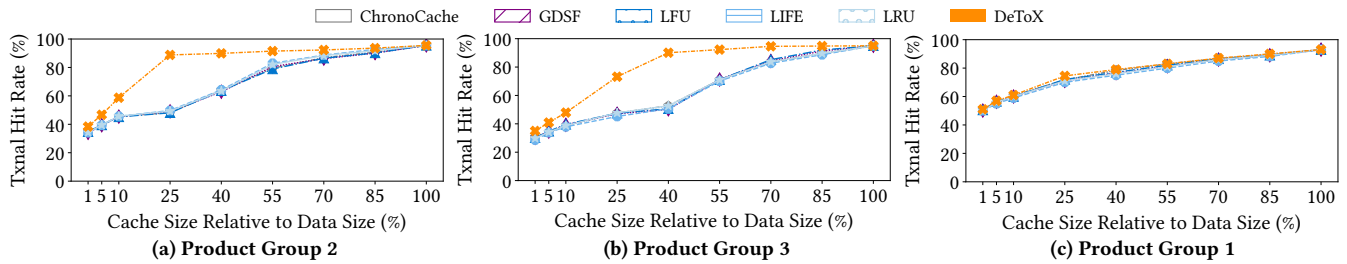


Figure 8: TAOBench THR results.

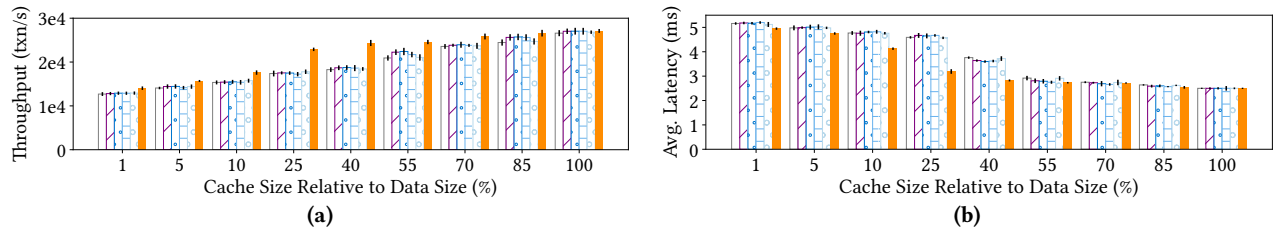


Figure 9: TAOBench PG2 results.

latency graphs for space but describe results in text. Since all transactions in these workloads have symmetric structures, there is no difference in performance between our various grouping techniques. We detail the tradeoffs between these optimizations in the next section.

**TAOBench.** DeToX obtains up to 76% higher transactional hit rates on the TAOBench PG2 and PG3 workloads compared to single-object caching algorithms (Figures 8a and 8b). DeToX achieves this with better cache efficiency: at the 25% cache size relative to data size (a common setup following the “80-20 rule”), the protocol achieves an 88% transactional hit rate while the best single-object algorithm requires 3.4x more cache space to attain the same result on PG2. Results are similar for PG3 for which the system requires a 2.2x smaller cache. Throughput increases by 31% (from 18K txns/s to 24 txns/s) for PG2 and 30% for PG3 (from 31K txns/s to 40 txns/s), while latency decreases by 30% (4.6ms to 3.2ms) for PG2 (Figure 9b) and 29% for PG3 (2.3ms to 1.6ms).

PG2 is read-dominant (>96%) with a mix of point reads, short transactions (<10 operations), and larger read transactions that span up to 40 keys. The point reads and shorter transactions make up 60% of the workload and largely access a small group of hot keys. Consequently, all algorithms achieve a THR of over 45% for small cache sizes (10% relative size). The longer read transactions follow one of two patterns: transactions access either a combination of hot and warm keys (25%), or hot and cold keys (11%). Transactions from the first category are more beneficial to cache since their keys are more frequently accessed and more likely to lead to transactional hits. There is little benefit in caching any of the keys in the second category since the cold keys contaminate all the other ones.

Under DeToX, the cache initially chooses to cache keys that belong to transactions in the first category. Thus, transactional hit rate improves as the cache size increases from 10% to

40% (Figure 8a). Past this point, the cache begins to retain more keys from transactions in the second category, but the performance benefit is limited since these requests rarely lead to transactional hits. In contrast, single-object algorithms use only individual object features to score keys, so they retain hot keys from transactions in both categories. Transactional hit rate increases slowly up to the 55% cache size at which point the cache becomes large enough to begin storing the warm keys from the first transaction category. Since the TAOBench workloads have no temporal patterns, GDSF and LFU provide slightly higher hit rates compared to LRU for all cache sizes. While LIFE uses levels, it performs poorly because it only uses the size of levels to make eviction decisions.

Similarly, in PG3, DeToX achieves better cache efficiency by not retaining contaminated keys. This workload has a smaller portion of point reads and shorter transactions (50%), so hit rates at smaller cache sizes are lower for all policies. Longer read transactions span up to 60 items and also fall into two categories. There are more transactions in the first category (33% compared to 25% in Product Group 2), so transactional hit rates grow more slowly with respect to cache size since more warm keys need to be cached.

In contrast to the other workloads, PG1 (Figure 8c) consists mainly of point reads and some short read transactions (of size four or smaller), which together make up over 97% of all requests. Our algorithm does not improve transactional hit rate over single-object policies because most hits result from standalone requests and short read transactions to a set of highly popular keys, which single-object algorithms already cache effectively. Throughput increases by 2% (from 82K txn/s to 84K txn/s), and latency decreases by 2% (from 0.61ms to 0.60ms).

ChronoCache has similar hit rates to single-object algorithms since there are no dependencies within transactions for this benchmark; the results simply reflect its eviction policy, LRU. The middleware layer, which does dependency

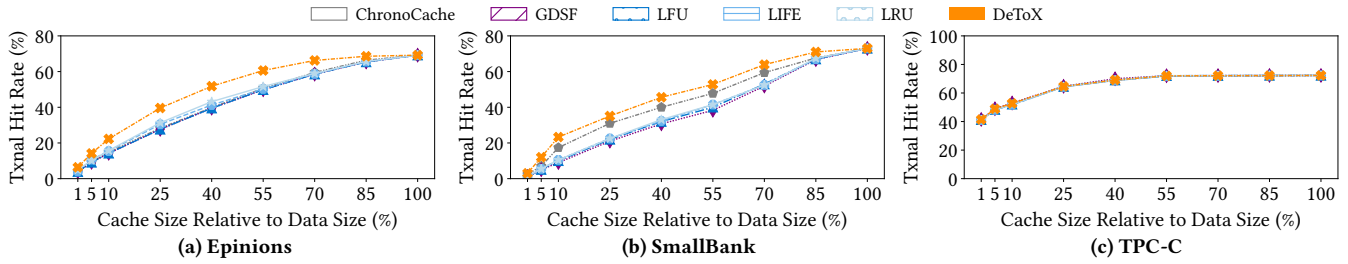


Figure 10: OLTP benchmark THR results.

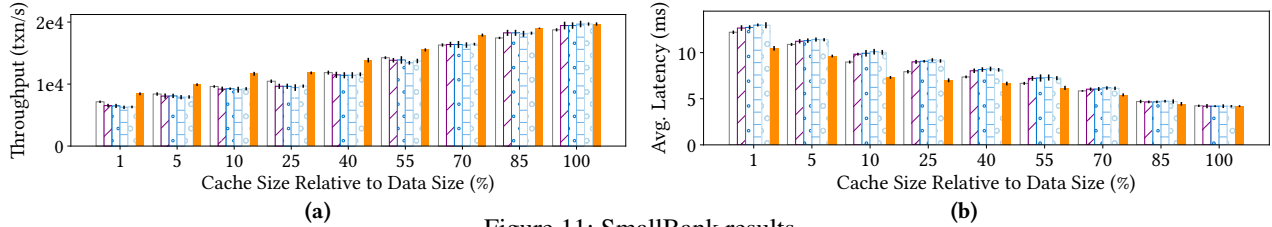


Figure 11: SmallBank results.

analysis at run time, quickly becomes the bottleneck.

**Epinions.** Epinions centers around user interactions and item reviews. It contains five read-only transactions and four update transactions. Users have both an n-to-m relationship with items (i.e., representing user reviews and ratings of items) and an n-to-m relationship with other users. There are no logical dependencies in the transactions of this workload (all operations can be parallelized).

DeToX provides up to 41% increase in transactional hit rate (Figure 10a), translating into 29% improvement in throughput (from 12K txn/s to 17K txn/s) and 25% decrease in latency (from 6.9ms to 5.5ms). At the 25% cache size, DeToX is 1.6x more efficient than the other algorithms. The transactions in Epinions request some group of objects related to a particular user or item (e.g., get all the reviews from one user), so our policy is able to successfully capture the n-to-m relationships in the data with its scoring mechanism. In contrast, the single object policies focus on caching individually popular keys without taking into account correlation between accesses. Since there are no dependencies between or within transactions for this workload, ChronoCache is unable to successfully prefetch objects.

**SmallBank.** SmallBank consists of requests to the Accounts, Checking, and Savings tables with six transaction types. Its transactions are relatively small, involving four distinct keys at most. Roughly two-thirds of operations are reads. Each customer account is materialized as three separate entries in each table and is accessed with a uniform distribution. There is high correlation between accesses to a customer’s row in the Accounts table and the customer’s rows in the other two tables.

Our algorithm provides up to a 1.3x increase in transactional hit rate (Figure 10b). The absolute hit rates remain relatively low for smaller cache sizes because of the uniform access distribution to customer accounts. Transactional hit rate increases linearly for all algorithms since more cache space directly

results in more hits. DeToX is 1.6x more efficient than the next best-performing algorithm at the 25% cache size.

We observe up to a 28% increase in throughput (from 12K txn/s to 16K txn/s) and 26% decrease in latency (from 6.8ms to 5.4ms) on this workload (Figures 11a and 11b). The long tail in access patterns and short transactions of this workload limit the benefits of our eviction algorithm over single-object alternatives, which all have similar performance.

For this workload, around two-thirds of performance improvement can be attributed to prefetching. We compare our eviction algorithm without prefetching (DeToX-E), LRU with prefetching (LRU-P), and our full policy (DeToX). DeToX-E increases throughput by 9%, LRU-P by 19%, and DeToX by 28% (graph omitted for space).

**TPC-C.** TPC-C is notably write-heavy and has transactions that can span over 50 items. Its requests tend to fall into two categories: either they access a small set of popular keys (i.e., those in the Warehouse and District tables) or a larger range of keys from a distribution with a long tail (Customer, Item, Stock). Single-object caching algorithms are designed to cache the former while the latter almost always results in transactional misses. For instance, *New-Order* accesses a key in each of the Warehouse, District, and Customer tables before requesting 10 to 15 items from the Item and Stock tables, which are chosen from a skewed distribution.

Consequently, TPC-C cannot benefit from transactional caching: most transactions access a small set of hot keys that are already in the cache (the object hit rate is >50% with a 10% cache size in Figure 10c) along with a larger set of cold keys that are unlikely to be cached and contaminate the other keys (hit rate grows slowly as cache size increases). Moreover, transactions tend to access keys in quick succession (e.g., once an order is placed, it is then processed, paid for, and delivered), so recency is especially important in this workload. All algorithms incorporate recency in some form,

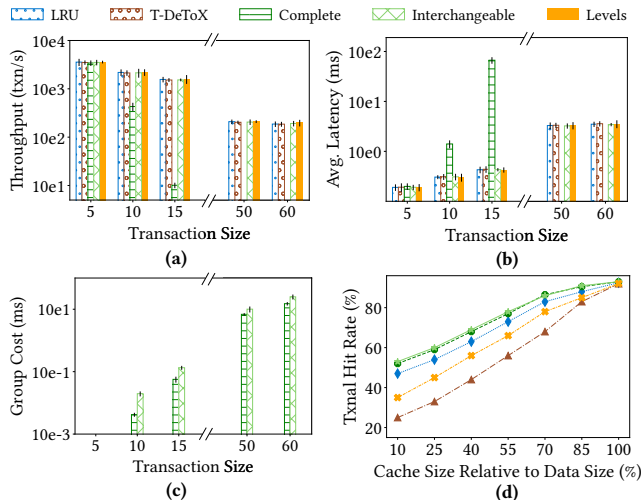


Figure 12: (a),(b),(c). Microbench. 1 (d) Microbench. 2.

so performance is similar across these policies, with up to 9K txns/s and 27ms avg. latency. DeToX performs on par with single-object policies.

### 8.3 The Need for Dependency Analysis

In this section, we investigate the relative merits of our grouping optimizations. The dependency analysis required for complete groups can impose overheads in two ways: (1) the cost of updating the scores of each key in each group and (2) metadata overhead associated with scoring. Interchangeability can reduce the number of groups that need to be scored, leading to better performance. On the other hand, levels discount unbalanced topologies while T-DeToX, a baseline that scores all keys of a transaction together, ignores dependencies. These simpler policies reduce overhead in some cases but restrict the groups that keys can belong to, leading to worse performance.

**Performance impact.** Microbenchmark 1 intentionally captures the worst-case scenario for grouping. We run a single transaction type with the topology in Figure 6a, and we extend the right branch of the graph for larger transaction sizes. Each read uniformly accesses keys at random among 10M objects. We measure throughput and latency as we increase transaction size up to 60 (equivalent to the largest transactions in the TAOBench workloads). Figures 12a and 12b show that performance for complete groups decreases dramatically as transaction size increases due to the exponential number of complete groups: for a transaction of size 15, over 16K groups have to be scored. Note that the bars for throughput and latency are omitted for complete groups for transaction sizes greater than 15 since these experiments did not finish in reasonable amount of time. In contrast, performance degradation is minimal with interchangeable groups (<5% difference compared to LRU at size 60). There are only a linear number of groups that must be scored with respect to transaction size since all keys in the right branch of this topology are interchangeable. Finally, levels offer similar performance to

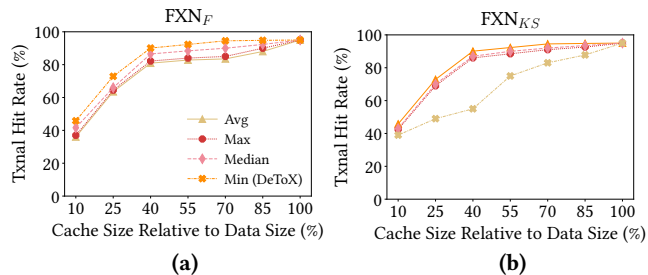


Figure 13: Scoring heuristics

LRU. Each key can only belong to one level per transaction, so larger transaction sizes do not increase overhead. The run time CPU overhead of both interchangeable groups and levels is within 5% of that of single-object algorithms for all microbenchmarks and previous benchmark workloads.

Moreover, the one-off cost of finding complete and interchangeable groups at compile time remains low: transactions of size 60 (with 100K+ groups due to worst-case topologies) require less than five minutes to process (Figure 12c). All benchmark workloads require less than 30 seconds for dependency analysis.

While dependency analysis incurs a static cost, it can lead to significant benefits compared to more basic forms of grouping (levels and T-DeToX), which ignore some or all dependency information. Microbenchmark 2 quantifies the worst-case scenarios for levels and T-DeToX. We run a single transaction type with the topology in Figure 6a in which the keys in vertices *a* and *c* are hot keys chosen from a Zipfian distribution while keys in *b* and *d* are cold keys chosen from a uniform distribution over 10M objects. Using levels causes keys in *b* and *c* to be scored together. However, keys in *b* are rarely accessed, and contaminate keys in *c*. T-DeToX makes even worse eviction decisions since it scores all keys in *a*, *b*, *c*, and *d* together. Using complete and interchangeable groups would instead cause keys in *b* and *d* to be scored together, enabling the algorithm to capture the fact that caching *c* individually reduces critical length. We find that complete and interchangeable groups significantly outperform levels (53% increase) and T-DeToX (139% increase) for THR (Figure 12d). Complete and interchangeable groups offer similar performance to LRU since these policies cache keys in *c*, which are frequently accessed.

**Memory overheads.** Metadata overhead in DeToX is low. Our algorithm stores two additional counters (total group score, individual score) per key and a global aging factor for eviction. While prefetching, DeToX stores dependency sets. On TAOBench, additional metadata takes up less than 1% of the cache space. For workloads in which prefetching is more prevalent, metadata overheads increase slightly. For example, in SmallBank, additional metadata grows to 2%. DeToX must store the dependency set associated with each transaction (1.5 keys on average).

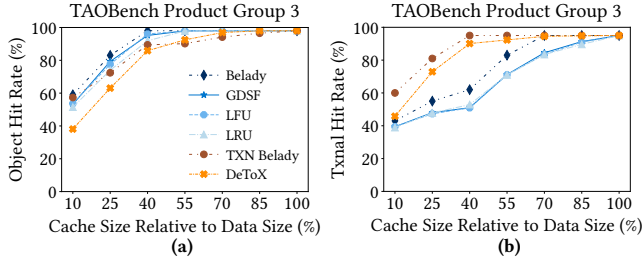


Figure 14: Object (a) and transactional (b) hit rates.

## 8.4 Scoring Heuristics

We evaluate different heuristics for calculating instance ( $F_{XNF}$ ) and aggregate scores ( $F_{XKS}$ ). DeToX uses the *minimum* frequency of keys in a group for the instance score, and *averages* instance scores to compute an aggregate score (Section 4.2). We measure transactional hit rates for simple functions (average, maximum, median, minimum) in Figure 13 for the PG3 workload (results are similar across workloads).

For assigning key instance scores, we find that, as expected, Min provides the best performance (Figure 13a). Since we only get a transactional hit if *all* keys of a group are cached, the key with the smallest frequency should have outsized impact on the group score. The other functions discount this information and thus perform worse. However, these functions still encode the all-or-nothing property of transactions to some extent since they assign the same instance scores to all keys in a particular group. As a result, we still observe higher hit rates than single-object policies.

Average and Median are the most effective functions for calculating aggregate key score (Figure 13b). Max yields a lower hit rate since it assigns each key the score of its highest-scoring group, but this may not be the most frequent group that contains this key. Min provides markedly lower performance (up to 64% lower hit rates). Each key is assigned the score of its lowest-scoring group, so most scores converge to the lowest group score (the smallest frequency of any key). As a result, most scores are low and do not differ by much.

## 8.5 OHR versus THR

There is a tradeoff between optimizing for latency and for system load. Figure 14 shows the OHR and THR of online algorithms as well as Belady and Transactional Belady (see Appendix A). As expected, Belady outperforms other algorithms for object hit rate. Conversely, DeToX and Transactional Belady give some of the lowest object hit rates. However, these two algorithms significantly outperform the other policies for transactional hit rate (and result in better throughput and latency as shown in Section 8.2). While we focus on PG3 here, we find similar results on the other workloads (omitted due to lack of space).

The difference between OHR and THR illustrates a tradeoff between reducing I/O bandwidth and optimizing for latency.

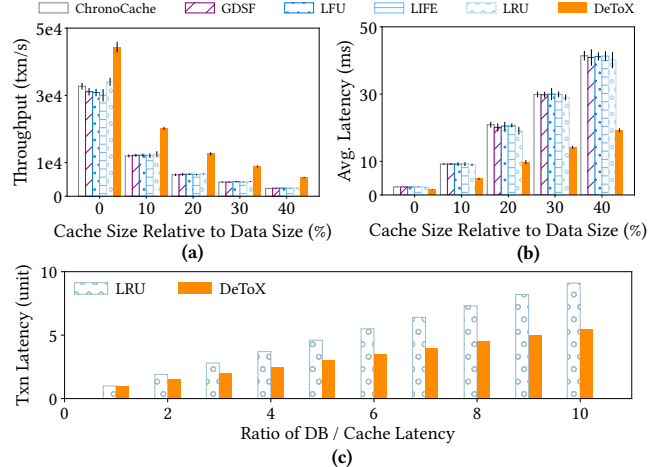


Figure 15: Network latency (a), (b) and simulation (c) results.

OHR prioritizes the absolute number of requests that can be served from cache, minimizing requests to disk. In contrast, THR focuses on the number of latency reductions for transactions, leading to lower latency and higher throughput. There are practical motivations for choosing THR as the caching objective: with increasing elasticity from cloud resources, applications often focus on latency optimization for which large wins are possible with DeToX.

## 8.6 Transactional Hit Rate

Transactional hit rate is independent of system specifics; only relative throughput and latency gains differ when cache / system latency changes. We confirm this by (1) varying this ratio (both experimentally and through simulation) and (2) evaluating DeToX with an alternative key-value store, TiKV [4].

**Network latency.** We inject latency between the shim layer and data store to simulate scenarios in which the latter is hosted in a remote cloud region. Figure 15 shows that the performance improvement with DeToX grows as network latency increases. With no additional network latency (0ms), there is a 30% increase in throughput and 29% decrease in latency between DeToX and the best single-object policy for PG3. With a WAN delay of 10ms, there is a 61% increase in throughput and 47% decrease in latency.

**Simulation results.** To illustrate the impact of cache and data store request times, we provide results for the TAOBench PG2 workload. At the 25% cache size, the THR for this workload is around 90% for DeToX and 50% for the other policies (Section 8.2). We vary request times for the cache and the data store (DB), using arbitrary units to represent latency. As we increase the ratio of DB to cache latency in Figure 15c, we find that the difference in request latency between LRU and DeToX increases from 0% to 65% as request times to the data store lengthen.

**Transactional key-value store.** We confirm that both the difference in transactional hit rate and gains in cache

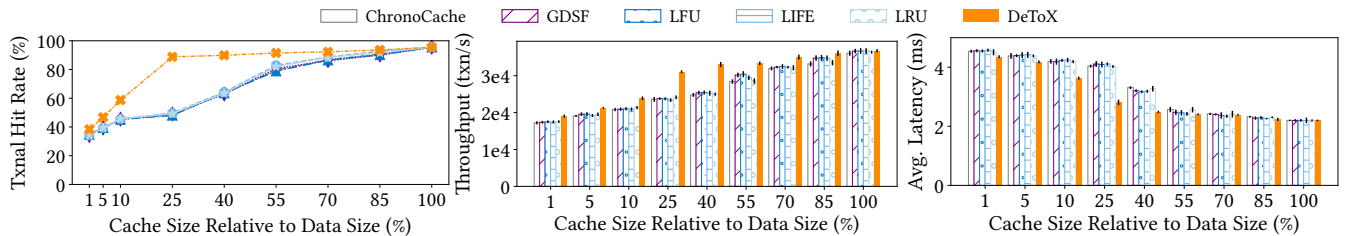


Figure 16: TAOBench PG2 results on TiKV.

efficiency (3.4x) remains identical when executing atop TiKV, demonstrating that these metrics are independent of the setup chosen (Figure 16). In contrast, as TiKV exhibits higher throughput and lower latency than Postgres, throughput and latency gains fall to 19% and 15% respectively.

## 9 Related Work

**Eviction.** There is a wide range of research on single-object caching policies that consider frequency [23, 40, 41, 53, 63], recency [32, 42, 52, 70], the number of unique keys between accesses [14, 50, 57, 64, 73], the variable sizes of objects [5, 25], and combinations of these features [6, 7, 12, 13, 15, 18, 20, 28, 49, 51, 56, 77, 80, 99]. Some specialized eviction policies optimize for flash storage [74], adapt to changing workloads [17, 19, 21, 30, 31, 39, 89], or consider network bandwidth and download time for proxy caches [93]. These previous efforts do not explicitly address how caching should be optimized for parallel accesses in transactions.

PACMan [11] presents eviction algorithms targeted towards job processing based on the all-or-nothing property: for jobs that issue tasks in parallel, latency only improves if all parallel tasks are cached. Similarly, existing literature on web caching [7, 18, 91] focuses on maximizing the *page* hit rate since latency is reduced only when all parts of a page are cached. Transactional hits in DeToX are based on a similar insight. However, DeToX addresses the issue of complex, unbalanced dependency graphs and recognizes that keys can be shared across many transactions.

**Admission algorithms.** In contrast to eviction algorithms, admission policies decide what to *allow* into the cache by enforcing a threshold based on object scores. These algorithms have often been applied alongside eviction policies [7, 21, 40, 52, 62]. While we focus on eviction and prefetching in this paper, our grouping and scoring strategies can feasibly extend to admission, which we will explore in future work.

**Prefetching.** Prefetching has been applied extensively to web caching [10, 90]. Past work focuses on web page analysis [38, 55, 68, 71, 86, 95, 96], which most stand-alone caches do not support [2, 69]. Other research [24, 44, 45, 72] centers around reducing the latency of query execution using dependency analysis. These works assume that each client issues queries sequentially, so any cache hit can

improve latency. Instead, DeToX caches in order to maximize transactional hit rate. Furthermore, none of these systems provide isolation guarantees or consider how eviction policies should be modified to handle transactions.

**Cache coherence.** Previous work combining transactions and caching focuses on maintaining isolation guarantees for cache coherence [1, 54, 76, 97]. In contrast, we focus on what objects to cache for performance. DeToX ensures serializability while optimizing for transactional hit rate.

## 10 Conclusion

In this paper, we study the problem of transactional caching. Standard caching policies fail to account for the all-or-nothing property of transactions, resulting in inefficient choices for which objects to retain in cache. In light of this issue, we provide a formal framework to quantify the latency impact of caching for transactions and introduce transactional hit rate as the key metric for this setting. We then present DeToX, a novel caching system targeting at transactional workloads. DeToX maximizes transactional hit rate by centering its caching policy around scoring groups of keys together. We consider how keys are accessed in parallel through complete groups and introduce interchangeable keys as an optimization to reduce the overhead of having to score many groups at run time. We also describe levels as a technique for cases when transaction code is not available. Our implementation is lightweight and deployable on a range of existing caching systems and data stores. DeToX improves THR by up to 1.3x and cache efficiency by up to 3.4x. This work demonstrates that many applications can benefit measurably from transactional caching.

## Acknowledgements

We thank Matt Burke, our anonymous reviewers, and our shepherd Wyatt Lloyd for their insightful feedback as well as Akshay Ravor for his engineering contributions. This work is supported by NSF CISE Expeditions Award CCF-1730628, NSF GRFP Award DGE-1752814, a Meta Next-Generation Infrastructure award, and gifts from Amazon, Astronomer, Google, IBM, Intel, Lacework, Microsoft, Nexla, Samsung SDS, and VMWare.

## References

- [1] Amazon elasticache, November 2021.
- [2] Redis, February 2021.
- [3] Postgresql, 2022.
- [4] Tikv, 2022.
- [5] Marc Abrams, Charles R. Standridge, Ghaleb Abdulla, Edward A. Fox, and Stephen Williams. Removal policies in network caches for world-wide web documents. In *Conference Proceedings on Applications, Technologies, Architectures, and Protocols for Computer Communications*, SIGCOMM '96, page 293–305, New York, NY, USA, 1996. Association for Computing Machinery.
- [6] Marc Abrams, Charles R Standridge, Ghaleb Abdulla, Stephen Williams, and Edward A Fox. Caching proxies: Limitations and potentials. 1995.
- [7] Charu Aggarwal, Joel L. Wolf, and Philip S. Yu. Caching on the world wide web. 11(1):94–107, jan 1999.
- [8] Jose Aguilar and Ernst L. Leiss. A web proxy cache coherency and replacement approach. In *Proceedings of the First Asia-Pacific Conference on Web Intelligence: Research and Development*, WI '01, page 75–84, Berlin, Heidelberg, 2001. Springer-Verlag.
- [9] Marcos K. Aguilera, Joshua B. Leners, and Michael Wal-fish. Yesquel: Scalable sql storage for web applications. In *Proceedings of the 25th Symposium on Operating Systems Principles*, SOSP '15, page 245–262, New York, NY, USA, 2015. Association for Computing Machinery.
- [10] Waleed Ali, Siti Mariyam Shamsuddin, Abdul Samad Ismail, et al. A survey of web caching and prefetching. *Int. J. Advance. Soft Comput. Appl.*, 3(1):18–44, 2011.
- [11] Ganesh Ananthanarayanan, Ali Ghodsi, Andrew Wang, Dhruba Borthakur, Srikanth Kandula, Scott Shenker, and Ion Stoica. Pacman: Coordinated memory caching for parallel jobs. In *Proceedings of the 9th USENIX Conference on Networked Systems Design and Implementation*, NSDI'12, page 20, USA, 2012. USENIX Association.
- [12] Martin Arlitt, Ludmila Cherkasova, John Dilley, Rich Friedrich, and Tai Jin. Evaluating content management techniques for web proxy caches. *SIGMETRICS Perform. Eval. Rev.*, 27(4):3–11, mar 2000.
- [13] Hyokyung Bahn, Kern Koh, S.H. Noh, and S.M. Lyul. Efficient replacement of nonuniform objects in web caches. *Computer*, 35(6):65–73, 2002.
- [14] Sorav Bansal and Dharmendra S. Modha. CAR: Clock with adaptive replacement. In *3rd USENIX Conference on File and Storage Technologies (FAST 04)*, San Francisco, CA, March 2004. USENIX Association.
- [15] Nathan Beckmann, Haoxian Chen, and Asaf Cidon. LHD: Improving cache hit rate by maximizing hit density. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*, pages 389–403, Renton, WA, April 2018. USENIX Association.
- [16] L. A. Belady. A study of replacement algorithms for a virtual-storage computer. *IBM Systems Journal*, 5(2):78–101, 1966.
- [17] Daniel S Berger. Towards lightweight and robust machine learning for cdn caching. In *Proceedings of the 17th ACM Workshop on Hot Topics in Networks*, pages 134–140, 2018.
- [18] Daniel S. Berger, Nathan Beckmann, and Mor Harchol-Balter. Practical bounds on optimal caching with variable object sizes. *Proc. ACM Meas. Anal. Comput. Syst.*, 2(2), jun 2018.
- [19] Daniel S. Berger, Benjamin Berg, Timothy Zhu, Siddhartha Sen, and Mor Harchol-Balter. RobinHood: Tail latency aware caching – dynamic reallocation from Cache-Rich to Cache-Poor. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, pages 195–212, Carlsbad, CA, October 2018. USENIX Association.
- [20] Daniel S Berger, Sebastian Henningsen, Florin Ciucu, and Jens B Schmitt. Maximizing cache hit ratios by variance reduction. *ACM SIGMETRICS Performance Evaluation Review*, 43(2):57–59, 2015.
- [21] Daniel S. Berger, Ramesh K. Sitaraman, and Mor Harchol-Balter. AdaptSize: Orchestrating the hot object memory cache in a content delivery network. In *14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17)*, pages 483–498, Boston, MA, March 2017. USENIX Association.
- [22] Philip A. Bernstein, Vassos Hadzilacos, and Nathan Goodman. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley, 1987.
- [23] Aaron Blankstein, Siddhartha Sen, and Michael J. Freedman. Hyperbolic caching: Flexible caching for web applications. In *2017 USENIX Annual Technical Conference (USENIX ATC 17)*, pages 499–511, Santa Clara, CA, July 2017. USENIX Association.
- [24] Ivan T. Bowman and Kenneth Salem. Optimization of query streams using semantic prefetching. In *Proceedings of the 2004 ACM SIGMOD international*

- conference on Management of data - SIGMOD '04. ACM Press, 2004.
- [25] Pei Cao and Sandy Irani. Cost-aware www proxy caching algorithms. In *USENIX Symposium on Internet Technologies and Systems (USITS 97)*, 1997.
- [26] Audrey Cheng, Xiao Shi, Aaron Kabcenell, Shilpa Lawande, Hamza Qadeer, Jason Chan, Harrison Tin, Ryan Zhao, Peter Bailis, Mahesh Balakrishnan, Nathan Bronson, Natacha Crooks, and Ion Stoica. Taobench: An end-to-end benchmark for social network workloads. *Proceedings of the VLDB Endowment*, 15(12):1965–1977, 2022.
- [27] Ludmila Cherkasova. *Improving WWW proxies performance with greedy-dual-size-frequency caching policy*. Hewlett-Packard Laboratories Palo Alto, CA, 1998.
- [28] Ludmila Cherkasova and Gianfranco Ciardo. Role of aging, frequency, and size in web cache replacement policies. In *International Conference on High-Performance Computing and Networking*, pages 114–123. Springer, 2001.
- [29] Marek Chrobak, Gerhard J. Woeginger, Kazuhisa Makino, and Haifeng Xu. Caching is hard—even in the fault model. *Algorithmica*, 63(4):781–794, March 2011.
- [30] Asaf Cidon, Assaf Eisenman, Mohammad Alizadeh, and Sachin Katti. Dynacache: Dynamic cloud caching. In *7th USENIX Workshop on Hot Topics in Cloud Computing (HotCloud 15)*, Santa Clara, CA, July 2015. USENIX Association.
- [31] Asaf Cidon, Assaf Eisenman, Mohammad Alizadeh, and Sachin Katti. Cliffhanger: Scaling performance cliffs in web memory caches. In *13th USENIX Symposium on Networked Systems Design and Implementation (NSDI 16)*, pages 379–392, Santa Clara, CA, March 2016. USENIX Association.
- [32] Fernando J Corbato. A paging experiment with the multics system. Technical report, MASSACHUSETTS INST OF TECH CAMBRIDGE PROJECT MAC, 1968.
- [33] The Transaction Processing Performance Council. Tpc-c, 2021.
- [34] Natacha Crooks, Matthew Burke, Ethan Cecchetti, Sitar Harel, Rachit Agarwal, and Lorenzo Alvisi. Obladi: Oblivious serializable transactions in the cloud. In *Proceedings of the 13th USENIX Conference on Operating Systems Design and Implementation, OSDI'18*, page 727–743, USA, 2018. USENIX Association.
- [35] Natacha Crooks, Youer Pu, Nancy Estrada, Trinabh Gupta, Lorenzo Alvisi, and Allen Clement. TARDiS. In *Proceedings of the 2016 International Conference on Management of Data*. ACM, June 2016.
- [36] Mohammad Dashti, Sachin Basil John, Amir Shaikhha, and Christoph Koch. Transaction repair for multi-version concurrency control. In *Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD '17*, page 235–250, New York, NY, USA, 2017. Association for Computing Machinery.
- [37] Djellel Eddine Difallah, Andrew Pavlo, Carlo Curino, and Philippe Cudre-Mauroux. Oltp-bench: An extensible testbed for benchmarking relational databases. volume 7, pages 277–288. VLDB Endowment, 2013.
- [38] Josep Domenech, Jose A. Gil, Julio Sahuquillo, and Ana Pont. Using current web page structure to improve prefetching performance. *Comput. Netw.*, 54(9):1404–1417, jun 2010.
- [39] Gil Einziger, Ohad Eytan, Roy Friedman, and Ben Manes. Adaptive software cache management. In *Proceedings of the 19th International Middleware Conference*, pages 94–106, 2018.
- [40] Gil Einziger, Roy Friedman, and Ben Manes. Tinylfu: A highly efficient cache admission policy. *ACM Transactions on Storage (ToS)*, 13(4):1–31, 2017.
- [41] Bin Fan, Hyeontaek Lim, David G. Andersen, and Michael Kaminsky. Small cache, big effect: Provable load balancing for randomly partitioned cluster services. In *Proceedings of the 2nd ACM Symposium on Cloud Computing, SOCC '11*, New York, NY, USA, 2011. Association for Computing Machinery.
- [42] Nicolas Gast and Benny Van Houdt. Transient and steady-state regime of a family of list-based cache replacement algorithms. *SIGMETRICS Perform. Eval. Rev.*, 43(1):123–136, jun 2015.
- [43] Shahram Ghandeharizadeh, Jason Yap, and Hieu Nguyen. Strong consistency in cache augmented SQL systems. In *Proceedings of the 15th International Middleware Conference on - Middleware '14*. ACM Press, 2014.
- [44] Brad Glasbergen, Michael Abebe, Khuzaima Daudjee, Scott Foggo, and Anil Pacaci. Apollo: Learning query correlations for predictive caching in geo-distributed systems, 2018.
- [45] Brad Glasbergen, Kyle Langendoen, Michael Abebe, and Khuzaima Daudjee. Chronocache: Predictive and adaptive mid-tier query result caching. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data, SIGMOD '20*, page 2391–2406, New York, NY, USA, 2020. Association for Computing Machinery.



- [46] Priya Gupta, Nickolai Zeldovich, and Samuel Madden. A trigger-based middleware cache for orms. In *Proceedings of the 12th International Middleware Conference, Middleware '11*, page 320–339, Laxenburg, AUT, 2011. International Federation for Information Processing.
- [47] Jeffrey Helt, Matthew Burke, Amit Levy, and Wyatt Lloyd. Regular sequential serializability and regular sequential consistency. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles CD-ROM*. ACM, October 2021.
- [48] Saied Hosseini-Khayat. *Investigation of Generalized Caching*. PhD thesis, USA, 1998. UMI Order No. GAX98-07761.
- [49] Xiameng Hu, Xiaolin Wang, Yechen Li, Lan Zhou, Yingwei Luo, Chen Ding, Song Jiang, and Zhenlin Wang. LAMA: Optimized locality-aware memory allocation for key-value cache. In *2015 USENIX Annual Technical Conference (USENIX ATC 15)*, pages 57–69, Santa Clara, CA, July 2015. USENIX Association.
- [50] Song Jiang, Feng Chen, and Xiaodong Zhang. Clockpro: An effective improvement of the clock replacement. In *USENIX Annual Technical Conference, General Track*, pages 323–336, 2005.
- [51] Shudong Jin and Azer Bestavros. Greedydual\* web caching algorithm: exploiting the two sources of temporal locality in web request streams. *Computer Communications*, 24(2):174–183, 2001.
- [52] Theodore Johnson and Dennis Shasha. 2q: A low overhead high performance buffer management replacement algorithm. In *Proceedings of the 20th International Conference on Very Large Data Bases, VLDB '94*, page 439–450, San Francisco, CA, USA, 1994. Morgan Kaufmann Publishers Inc.
- [53] George Karakostas and Dimitrios N Serpanos. Exploitation of different types of locality for web caches. In *Proceedings ISCC 2002 Seventh International Symposium on Computers and Communications*, pages 207–212. IEEE, 2002.
- [54] Bryan Kate, Eddie Kohler, Michael S. Kester, Neha Narula, Yandong Mao, and Robert Morris. Easy freshness with pequod cache joins. In *Proceedings of the 11th USENIX Conference on Networked Systems Design and Implementation, NSDI'14*, page 415–428, USA, 2014. USENIX Association.
- [55] Bin Lan, Stephane Bressan, Beng Chin Ooi, and Kian-Lee Tan. Rule-assisted prefetching in web-server caching. In *Proceedings of the Ninth International Conference on Information and Knowledge Management, CIKM '00*, page 504–511, New York, NY, USA, 2000. Association for Computing Machinery.
- [56] Donghee Lee, Jongmoo Choi, Jong-Hun Kim, Sam H Noh, Sang Lyul Min, Yookun Cho, and Chong Sang Kim. On the existence of a spectrum of policies that subsumes the least recently used (lru) and least frequently used (lfu) policies. In *Proceedings of the 1999 ACM SIGMETRICS international conference on Measurement and modeling of computer systems*, pages 134–143, 1999.
- [57] Cong Li. Dlirs: Improving low inter-reference recency set cache replacement policy with dynamics. In *Proceedings of the 11th ACM International Systems and Storage Conference*, pages 59–64, 2018.
- [58] Jialin Li, Ellis Michael, and Dan R. K. Ports. Eris. In *Proceedings of the 26th Symposium on Operating Systems Principles*. ACM, October 2017.
- [59] Wyatt Lloyd, Michael J. Freedman, Michael Kaminsky, and David G. Andersen. Stronger semantics for low-latency geo-replicated storage. In *10th USENIX Symposium on Networked Systems Design and Implementation (NSDI 13)*, pages 313–328, Lombard, IL, April 2013. USENIX Association.
- [60] Haonan Lu, Christopher Hodsdon, Khiem Ngo, Shuai Mu, and Wyatt Lloyd. The SNOW theorem and Latency-Optimal Read-Only transactions. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 135–150, Savannah, GA, November 2016. USENIX Association.
- [61] Haonan Lu, Siddhartha Sen, and Wyatt Lloyd. Performance-Optimal Read-Only transactions. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 333–349. USENIX Association, November 2020.
- [62] Bruce M. Maggs and Ramesh K. Sitaraman. Algorithmic nuggets in content delivery. *SIGCOMM Comput. Commun. Rev.*, 45(3):52–66, jul 2015.
- [63] Dhruv Matani, Ketan Shah, and Anirban Mitra. An o(1) algorithm for implementing the lfu cache eviction scheme. *arXiv preprint arXiv:2110.11602*, 2021.
- [64] Nimrod Megiddo and Dharmendra S. Modha. ARC: A Self-Tuning, low overhead replacement cache. In *2nd USENIX Conference on File and Storage Technologies (FAST 03)*, San Francisco, CA, March 2003. USENIX Association.
- [65] Syed Akbar Mehdi, Cody Littlely, Natacha Crooks, Lorenzo Alvisi, Nathan Bronson, and Wyatt Lloyd. I can't believe it's not causal! scalable causal consistency

- with no slowdown cascades. In *14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17)*, pages 453–468, Boston, MA, March 2017. USENIX Association.
- [66] Shuai Mu, Yang Cui, Yang Zhang, Wyatt Lloyd, and Jinyang Li. Extracting more concurrency from distributed transactions. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*, pages 479–494, Broomfield, CO, October 2014. USENIX Association.
- [67] Shuai Mu, Lamont Nelson, Wyatt Lloyd, and Jinyang Li. Consolidating concurrency control and consensus for commits under conflicts. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 517–532, Savannah, GA, November 2016. USENIX Association.
- [68] A. Nanopoulos, D. Katsaros, and Y. Manolopoulos. A data mining algorithm for generalized web prefetching. *IEEE Transactions on Knowledge and Data Engineering*, 15(5):1155–1169, 2003.
- [69] 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 *10th USENIX Symposium on Networked Systems Design and Implementation (NSDI 13)*, pages 385–398, Lombard, IL, April 2013. USENIX Association.
- [70] Elizabeth J O’neil, Patrick E O’neil, and Gerhard Weikum. The lru-k page replacement algorithm for database disk buffering. *Acm Sigmod Record*, 22(2):297–306, 1993.
- [71] Venkata N. Padmanabhan and Jeffrey C. Mogul. Using predictive prefetching to improve world wide web latency. *SIGCOMM Comput. Commun. Rev.*, 26(3):22–36, jul 1996.
- [72] Mark Palmer and Stanley B Zdonik. *Fido: A cache that learns to fetch*. Brown University, Department of Computer Science, 1991.
- [73] Sejin Park and Chanik Park. Frd: A filtering based buffer cache algorithm that considers both frequency and reuse distance. In *Proc. of the 33rd IEEE International Conference on Massive Storage Systems and Technology (MSST)*, 2017.
- [74] Seon-yeong Park, Dawoon Jung, Jeong-uk Kang, Jin-soo Kim, and Joonwon Lee. Cfrru: A replacement algorithm for flash memory. In *Proceedings of the 2006 International Conference on Compilers, Architecture and Synthesis for Embedded Systems, CASES ’06*, page 234–241, New York, NY, USA, 2006. Association for Computing Machinery.
- [75] Francisco Perez-Sorrosal, Marta Patiño Martinez, Ricardo Jimenez-Peris, and Bettina Kemme. Elastic si-cache: Consistent and scalable caching in multi-tier architectures. *The VLDB Journal*, 20(6):841–865, dec 2011.
- [76] Dan R. K. Ports, Austin T. Clements, Irene Zhang, Samuel Madden, and Barbara Liskov. Transactional consistency and automatic management in an application data cache. In *Proceedings of the 9th USENIX Conference on Operating Systems Design and Implementation, OSDI’10*, page 279–292, USA, 2010. USENIX Association.
- [77] Luigi Rizzo and Lorenzo Vicisano. Replacement policies for a proxy cache. *IEEE/ACM Transactions on networking*, 8(2):158–170, 2000.
- [78] Weihai Shen, Ansh Khanna, Sebastian Angel, Siddhartha Sen, and Shuai Mu. Rolis. In *Proceedings of the Seventeenth European Conference on Computer Systems*. ACM, March 2022.
- [79] Zhenyu Song, Daniel S. Berger, Kai Li, and Wyatt Lloyd. Learning relaxed belady for content distribution network caching. In *17th USENIX Symposium on Networked Systems Design and Implementation (NSDI 20)*, pages 529–544, Santa Clara, CA, February 2020. USENIX Association.
- [80] David Starobinski and David Tse. Probabilistic methods for web caching. *Performance evaluation*, 46(2-3):125–137, 2001.
- [81] Chunzhi Su, Natacha Crooks, Cong Ding, Lorenzo Alvisi, and Chao Xie. Bringing modular concurrency control to the next level. In *Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD ’17*, page 283–297, New York, NY, USA, 2017. Association for Computing Machinery.
- [82] Florian Suri-Payer, Matthew Burke, Zheng Wang, Yunhao Zhang, Lorenzo Alvisi, and Natacha Crooks. Basil. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles CD-ROM*. ACM, October 2021.
- [83] Adriana Szekeres, Michael Whittaker, Jialin Li, Naveen Kr. Sharma, Arvind Krishnamurthy, Dan R. K. Ports, and Irene Zhang. Meerkat. In *Proceedings of the Fifteenth European Conference on Computer Systems*. ACM, April 2020.

- [84] Cheng Tan, Changgeng Zhao, Shuai Mu, and Michael Walfish. Cobra: Making transactional Key-Value stores verifiably serializable. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 63–80. USENIX Association, November 2020.
- [85] Chuzhe Tang, Zhaoguo Wang, Xiaodong Zhang, Qianmian Yu, Binyu Zang, Haibing Guan, and Haibo Chen. Ad hoc transactions in web applications: The good, the bad, and the ugly. In *Proceedings of the 2022 International Conference on Management of Data, SIGMOD '22*, page 4–18, New York, NY, USA, 2022. Association for Computing Machinery.
- [86] Na Tang and V. Rao Vemuri. An artificial immune system approach to document clustering. In *Proceedings of the 2005 ACM Symposium on Applied Computing, SAC '05*, page 918–922, New York, NY, USA, 2005. Association for Computing Machinery.
- [87] The H-Store team. Smallbank benchmark, 2013.
- [88] Boyu Tian, Jiamin Huang, Barzan Mozafari, and Grant Schoenebeck. Contention-aware lock scheduling for transactional databases. *Proc. VLDB Endow.*, 11(5):648–662, jan 2018.
- [89] Giuseppe Vietri, Liana V Rodriguez, Wendy A Martinez, Steven Lyons, Jason Liu, Raju Rangaswami, Ming Zhao, and Giri Narasimhan. Driving cache replacement with {ML-based}{LeCaR}. In *10th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 18)*, 2018.
- [90] Jia Wang. A survey of web caching schemes for the internet. *SIGCOMM Comput. Commun. Rev.*, 29(5):36–46, oct 1999.
- [91] Justin Wang, Benjamin Berg, Daniel S. Berger, and Siddhartha Sen. Maximizing page-level cache hit ratios in largeweb services. *SIGMETRICS Perform. Eval. Rev.*, 46(2):91–92, jan 2019.
- [92] Zhaoguo Wang, Shuai Mu, Yang Cui, Han Yi, Haibo Chen, and Jinyang Li. Scaling multicore databases via constrained parallel execution. In *Proceedings of the 2016 International Conference on Management of Data, ACM*, June 2016.
- [93] Roland P. Wooster and Marc Abrams. Proxy caching that estimates page load delays. *Computer Networks and ISDN Systems*, 29(8-13):977–986, 1997.
- [94] Yingjun Wu, Chee-Yong Chan, and Kian-Lee Tan. Transaction healing: Scaling optimistic concurrency control on multicores. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD '16*, page 1689–1704, New York, NY, USA, 2016. Association for Computing Machinery.
- [95] Lifang Xu, Hongwei Mo, Kejun Wang, and Na Tang. Document clustering based on modified artificial immune network. In Guo-Ying Wang, James F. Peters, Andrzej Skowron, and Yiyu Yao, editors, *Rough Sets and Knowledge Technology*, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg.
- [96] Qiang Yang, Haining Henry Zhang, and Tianyi Li. Mining web logs for prediction models in www caching and prefetching. *KDD '01*, page 473–478, New York, NY, USA, 2001. Association for Computing Machinery.
- [97] Xiangyao Yu, Yu Xia, Andrew Pavlo, Daniel Sanchez, Larry Rudolph, and Srinivas Devadas. Sundial: Harmonizing concurrency control and caching in a distributed oltp database management system. *Proc. VLDB Endow.*, 11(10):1289–1302, jun 2018.
- [98] Irene Zhang, Naveen Kr. Sharma, Adriana Szekeres, Arvind Krishnamurthy, and Dan R. K. Ports. Building consistent transactions with inconsistent replication. volume 35, New York, NY, USA, December 2018. Association for Computing Machinery.
- [99] Yuanyuan Zhou, James Philbin, and Kai Li. The Multi-Queue replacement algorithm for second level buffer caches. In *2001 USENIX Annual Technical Conference (USENIX ATC 01)*, Boston, MA, June 2001. USENIX Association.

## A Appendix

We prove the optimal offline transactional caching problem is NP-Hard. We begin by providing intuition for how and why traditional optimal offline caching policies fail to translate to transactional caching.

### A.1 Transactional Belady

We straightforwardly adapt Belady’s optimal caching policy [16] to the transactional context by defining *Transactional Belady*, a caching policy that evicts keys that result in *transactional hits* furthest in the future. While this extension is intuitive, it does not offer optimal performance even for flat, uniformly-sized transactions that access equally-sized objects, as we prove below.

Consider the execution trace in Figure 17 with cache capacity of 5. All transactions access three keys, either all from set  $S_1: \{t, u, v, t', u'\}$  or set  $S_2: \{x, y, z, x', y'\}$ .  $T_1$  and  $T_2$  access only keys from the former group, while  $T_3$  and  $T_4$  access only keys from the latter.  $T_5$  and  $T_6$  access keys from  $S_1$  and overlap in  $v$ , while  $T_7, T_8, T_9$  overlap in  $x', y, z$  from  $S_2$ . Transactional

T	Keys accessed	Cache state	Optimal cache state
1	$t, u, v$	-	-
2	$t', u', v$	$t, u, v$	-
3	$x, y, z$	$t, u, v, t', u'$	-
4	$x', y', z$	$t, u, v, t', u'$	$x, y, z$
5	$t, u, v$	$t, u, v, t', u'$	$x, y, z, x', y'$
6	$t', u', v$	$t, u, v, t', u'$	$x, y, z, x', y'$
7	$x, y, z$	$t, u, v, t', u'$	$x, y, z, x', y'$
8	$x', y, z$	$t, u, v, t', u'$	$x, y, z, x', y'$
9	$x', y', z$	$t, u, v, t', u'$	$x, y, z, x', y'$

Figure 17: Non-optimality of Transactional Belady. Red keys indicate ones that lead to transactional hits.

Belady evicts keys that yield a transactional hit furthest in the future. After  $T_3$ 's execution, the algorithm evicts  $x, y, z$  as they would first yield a hit at  $T_7$  while the other keys would lead a hit at  $T_5$  and  $T_6$ . A similar reasoning leads the algorithm to evict  $x', y', z'$  after  $T_4$  executes. This strategy yields two transactional hits (for  $T_5$  and  $T_6$ ). Unfortunately, evicting  $x, y, z$  after  $T_3$  is the wrong decision. Keeping all keys of set  $S_2$  in the cache yields three transactional hits  $T_7, T_8, T_9$ . As a result, Transactional Belady achieves only two transactional hits, while an optimal caching policy would achieve three.

Transactional Belady does not account for shared keys across transactions. It caches set  $S_1$ , which is shared across two transactions, instead of keys in set  $S_2$ , which is shared across three transactions. Belady assumes that a cache hit closer in the future is always as valuable as a cache hit further out. This assumption holds when a single cached object provides a single object hit but breaks down when keys are shared across transactions. In these cases, an equal number of cached objects can produce varying numbers of transactional hits.

## A.2 Optimal Offline Transactional Caching is NP-Hard

We demonstrate that the optimal offline transactional caching problem (TxPolicy) is NP-Hard through a reduction from the variable-sized caching problem, CACHING(FAULT, OPTIONAL), introduced in [29].

We first provide intuition for our reduction. A page hit is only possible if the entire page is present in the cache, regardless of its size. The objective of CACHING(FAULT, OPTIONAL) is to minimize the number of page faults, or the number of pages accessed and missed. We convert each page of size  $X$  into a transaction without dependencies that accesses  $X$  operations. Therefore, there is only a transaction hit when the entire transaction is in the cache. This transforms CACHING(FAULT, OPTIONAL) into an easier version of TxPolicy with two simplifying assumptions: (1) all transactions will use unique keys, so that retaining a key in the cache from any single transaction provides no benefit to any other transaction, and (2) there are no logical dependencies. If an optimal offline transactional

caching policy exists, then through this reduction, we have the optimal policy for CACHING(FAULT, OPTIONAL).

We now formally describe CACHING(FAULT, OPTIONAL) from [29]. CACHING(FAULT, OPTIONAL) asks,

Given a set of pages  $p_1, \dots, p_k$  with sizes

$\text{SIZE}(p_1), \dots, \text{SIZE}(p_k)$ , request sequence  $r_1, \dots, r_m \in \{p_1, \dots, p_k\}$ , cache size  $C$ , and cost bound  $F$ , is there a replacement policy that serves  $r_1, \dots, r_m$  with cache size  $C$  and incurs a total fault cost at most  $F$ ?

A *fault* is incurred when  $r_i \notin C_i$ , where the FAULT parameter states that each fault has cost 1. The OPTIONAL parameter requires that  $\forall i > 1, C_i \subseteq \{C_{i-1} \cup r_i\}$ ; informally, the caching policy does not have to admit the most recent page.

We formally define the offline transactional caching problem, based on our formalisms from Section 3.

**Definition 10** (*Offline transactional caching policy*). An *offline transactional caching policy* is a function  $P$  that takes a sequence of transactions  $T_1, T_2, \dots, T_m$ , cache size  $n$ , and outputs a sequence of cache states  $C_1, C_2, \dots, C_m$ , with the following restrictions:

1.  $C_1 = \emptyset$ .
2.  $\forall i > 1, C_i \subseteq \{C_{i-1} \cup T_{i-1}\}$ .

TxPolicy asks,

Given a set of transactions  $T_1, \dots, T_m$ , cache size  $C$ , is there an offline transactional caching policy that serves  $T_1, \dots, T_m$  with cache size  $C$  and incurs at most  $F$  transactional misses? We define *transactional misses* as the number of  $i$  where  $T_i \not\subseteq C_i$ , or the number of transactions that cannot be served from cache.

**Theorem 1.** *The optimal offline transactional caching problem is NP-Hard.*

*Proof.* We reduce CACHING(FAULT, OPTIONAL) to TxPolicy through the following polynomial-time reductions. Each page  $p_i$  is reduced to a transaction  $T_i$ .  $\text{SIZE}(p_i)$  new tables are created per transaction, each with only one key. Let  $X$  be one such table. A read operation on the sole key of that table  $x \in X$  is inserted into the transaction  $T_i$ . There are no logical dependencies. Cache size  $C$  is preserved. The maximum fault cost  $F$  is converted to the maximum number of transactional misses. If there exists a policy solving the offline transactional caching problem, run it with these parameters. Its output is the output to the CACHING(FAULT, OPTIONAL) problem. CACHING(FAULT, OPTIONAL) is NP-Hard; therefore, the offline transactional caching problem is NP-Hard.  $\square$

## B Artifact Appendix

### Abstract

DeToX is a transactional caching system that leverages insights on transactional hit rate to improve caching performance for transactional workloads. DeToX is implemented as a shim layer that integrates with caching and database systems. In addition to DeToX, the artifact contains several other implementations. First, there is a modified version of ChronoCache, a middleware predictive query caching system, that measures transactional hit rate, integrates with Redis, and supports several benchmarks not available for the original system. There is also a modified version of Redis that supports several eviction algorithms, including DeToX's eviction algorithm and LIFE from the PACMan paper. Finally, there is a caching simulator that takes transaction traces as input and outputs hit rates for the offline Belady and Transactional Belady algorithms.

### Scope

The artifact enables others to run DeToX directly. All code used in the paper is made available.

### Contents

The artifact consists of a Github repository hosted at <https://github.com/audreycccheng/detox>. The repository is structured as follows:

- /chronocache - the codebase for ChronoCache
- /oltpbench-chronocache - the benchmarks for ChronoCache
- /redis - the modified version of Redis supporting transactional caching algorithms
- /simulator - the caching simulator for offline policies
- /sys - the transactional caching system
  - /benchmarks - the benchmarks for running DeToX
  - /src - the implementation of the DeToX shim layer

### Hosting

The artifact is hosted at <https://github.com/audreycccheng/detox> on the main branch at commit 604c9bd.

### Requirements

The following packages are required to run the codebase.

- mvn 3.8.5
- build-essential
- Java 17

For specific installation guides for each system, please see the Github repository.