

Online Tree Caching*

Marcin Bienkowski
Institute of Computer Science
University of Wrocław
Poland

Jan Marcinkowski
Institute of Computer Science
University of Wrocław
Poland

Maciej Pacut
Institute of Computer Science
University of Wrocław
Poland

Stefan Schmid
Department of Computer Science
Aalborg University
Denmark

Aleksandra Spyra
Institute of Computer Science
University of Wrocław
Poland

ABSTRACT

We initiate the study of a natural and practically relevant new variant of online caching where the to-be-cached items can have dependencies. We assume that the universe is a tree T and items are tree nodes; we require that if a node v is cached then the whole subtree $T(v)$ rooted at v is cached as well. This theoretical problem finds an immediate application in the context of forwarding table optimization in IP routing and software-defined networks.

We present an elegant online deterministic algorithm TC for this problem, and rigorously prove that its competitive ratio is $O(\text{HEIGHT}(T) \cdot k_{\text{ONL}} / (k_{\text{ONL}} - k_{\text{OPT}} + 1))$, where k_{ONL} and k_{OPT} denote the cache sizes of an online and the optimal offline algorithm, respectively. The result is optimal up to a factor of $O(\text{HEIGHT}(T))$.

CCS CONCEPTS

• **Theory of computation** → **Online algorithms**; *Caching and paging algorithms*; • **Networks** → *Programmable networks*; Packet-switching networks;

KEYWORDS

online algorithms, competitive analysis, caching, routers, software-defined networking, forwarding information base

ACM Reference format:

Marcin Bienkowski, Jan Marcinkowski, Maciej Pacut, Stefan Schmid, and Aleksandra Spyra. 2017. Online Tree Caching. In *Proceedings of SPAA '17, Washington DC, USA, July 24-26, 2017*, 11 pages.
<https://doi.org/10.1145/3087556.3087558>

1 INTRODUCTION

In the classic online paging problem, items of some universe are requested by a processing entity (e.g., blocks of RAM are requested by the processor). To speed up the access, computers use a faster memory, called *cache*, capable of accommodating k such items.

Upon a request to a non-cached item, the algorithm has to fetch it into the cache, paying a fixed cost, while a request to a cached item is free. If the cache is full, the algorithm has to free some space by evicting an arbitrary subset of items from the cache.

The paging problem is inherently online: the algorithm has to make decisions what to evict from the cache without the knowledge of future requests; its cost is compared to the cost of an optimal *offline* solution and the worst-case ratio of these two amounts is called *competitive ratio*. The first analysis of this basic problem in an online model was given over three decades ago by Sleator and Tarjan [30]. The problem was later considered in a variety of flavors. In particular, some papers considered a *bypassing model* [13, 17], where item fetching is optional: the requested item can be served without being in the cache, for another fixed cost (usually being at most the cost of item fetching).

In this paper, we introduce a natural extension of this fundamental problem, where items have inter-dependencies. More precisely, we assume that the universe is an arbitrary (not necessarily binary) rooted tree T and the requested items are its nodes. For any tree node v , $T(v) \subseteq T$ is a subtree rooted at v containing v and all its descendants. We require the following property: if a node v is in the cache, then all nodes of $T(v)$ are also cached. In other words, we require that *the cache is a subforest of T* , i.e., a union of disjoint subtrees of T . We call this problem *online tree caching*.

Furthermore, we assume a bypassing model and distinguish between two types of requests: a request can be either *positive* or *negative*. The positive requests correspond to “normal” requests known from caching problems: we pay 1 if the node is not cached; for a negative request, we pay 1 if the corresponding request is cached. After serving the request, we may reorganize our cache arbitrarily, but the resulting cache has to still be a subforest of T . We pay α for fetching or evicting any single node, where $\alpha \geq 1$ is an integer and a parameter of the problem. Our goal is to minimize the overall cost of maintaining the cache and serving the requests.

One interesting application for our model arises in the context of modern IP routers which need to store a rapidly increasing number of forwarding rules [1, 11]. In Section 2, we give a glimpse of this application, discussing how tree caching algorithms can be applied in existing systems to effectively reduce the memory requirements on IP routers.

1.1 Our Contributions and Paper Organization

We initiate the study of a natural new caching with bypassing problem which allows to account for tree-dependencies among

*M. Pacut and A. Spyra were supported by Polish National Science Centre grant DEC-2013/09/B/ST6/01538, M. Bienkowski by Polish National Science Centre grant 2016/22/E/ST6/00499, and S. Schmid by Aalborg University’s talent management program.

SPAA ’17, July 24-26, 2017, Washington DC, USA

© 2017 Copyright held by the owner/author(s). Publication rights licensed to Association for Computing Machinery.

This is the author’s version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Proceedings of SPAA ’17, July 24-26, 2017*, <https://doi.org/10.1145/3087556.3087558>.

items. The problem finds immediate applications, e.g., in IP routing and software-defined networking (see Section 2).

In particular, we consider the online tree caching problem within the resource augmentation paradigm: we assume that cache sizes of the online algorithm (k_{ONL}) and the optimal offline algorithm (k_{OPT}) may differ. We assume $k_{\text{ONL}} \geq k_{\text{OPT}}$ and let $R = k_{\text{ONL}}/(k_{\text{ONL}} - k_{\text{OPT}} + 1)$.

In Section 4, we present an elegant deterministic online algorithm TC for this problem. While our algorithm is simple, its analysis presented in Section 5 requires several non-trivial insights into the problem. In particular, we rigorously prove that TC is $O(h(T) \cdot R)$ -competitive, where $h(T)$ is the height of tree T . That is, we show that there exists a constant β , such that $\text{TC}(I) \leq O(h(T) \cdot R) \cdot \text{OPT}(I) + \beta$ for any input I . Note that this result is optimal up to the factor $O(h(T))$: in Appendix C, we show that the lower bound R for the paging problem [30] implies an $\Omega(R)$ lower bound for our problem for any $\alpha \geq 1$. Finally, in Section 6, we show that TC can be implemented efficiently.

1.2 Related Work on Caching

Our formal model is a novel variant of competitive paging, a classic online problem. In the framework of the competitive analysis, the paging problem was first analyzed by Sleator and Tarjan [30], who showed that algorithms LEAST-RECENTLY-USED, FIRST-IN-FIRST-OUT and FLUSH-WHEN-FULL are $k_{\text{ONL}}/(k_{\text{ONL}} - k_{\text{OPT}} + 1)$ -competitive and no deterministic algorithm can beat this ratio. In the non-augmented case when $k_{\text{ONL}} = k_{\text{OPT}} = k$, the competitive ratio is simply k .

The simple paging problem was later generalized to allow different fetching costs (weighted paging) [10, 34] and additionally different item sizes (file caching) [35], with the same competitive ratio. Asymptotically same results can be achieved when bypassing is allowed (see [13, 17] and references therein). With randomization, the competitive ratio can be reduced to $O(\log k)$ even for file caching [3]. The lower bound for randomized algorithms is $H_k = \Theta(\log k)$ [14] and is matched by known paging algorithms [2, 26].

To the best of our knowledge, the variant of caching, where fetching items to the cache is not allowed unless some other items are cached (e.g., because of tree dependencies) was not considered previously in the framework of competitive analysis. Note that there is a seemingly related problem called restricted caching [8] (there are also its variants called matroid caching [9] or companion caching [27]). Despite naming similarities, the restricted caching model is completely different from ours: there the restriction is that each item can be placed only in a restricted set of cache locations.

2 APPLICATION: MINIMIZING FORWARDING TABLES IN ROUTERS

Dependencies among to-be-cached items arise in numerous settings and are a natural refinement of many caching problems. To give a concrete example, one important application for our tree-based dependency model arises in the context of IP routers. In particular, the online tree caching problem we introduce in this paper is motivated by router memory constraints in IP-based networks. The material presented in this section serves for motivation, and is not necessary for understanding the remainder of the paper.

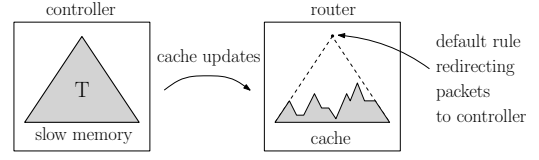


Figure 1: The router (right) caches only a subset of all rules, and rules that are not cached are answered by the controller (left) that keeps the whole tree of rules. Updates to the rules are passed by the controller to the router.

Nowadays, routers have to store an enormous number of forwarding rules: the number of rules has doubled in the last six years [1] and the superlinear growth is likely to be sustained [11]. This entails large costs for Internet Service Providers: fast router memory (usually Ternary Content Addressable Memory (TCAM)) is expensive and power-hungry [31]. Many routers currently either operate at (or beyond) the edge of their memory capacities. A solution, which could delay the need for expensive or impossible memory upgrades in routers, is to store only a subset of rules in the actual router and store all rules on a secondary device (for example a commodity server with a large but slow memory) [19–22, 29].

This solution is particularly attractive with the advent of Software-Defined Network (SDN) technology, which allows to manage the expensive memory using a software controller [19, 29]. In particular, our theoretical model can describe real-world architectures like [19, 29], that is, our model formalizes the underlying operational problems of such architectures. Our algorithm, when applied in the context of such architectures, can hence be used to prolong the lifetime of IP routers.

Setup, positive requests, fetches and evictions. The setup (see [29] for a more technical discussion) depicted in Figure 1 consists of two entities: the actual router (e.g., an OpenFlow switch) which caches only a subset of all forwarding rules, and the (SDN) controller, which keeps all rules in its less expensive and slower memory. During runtime, packets arrive at the router, and if an appropriate forwarding rule is found within the rules cached by the router, then the packet is forwarded accordingly, and the associated cost is zero. Otherwise, the packet has to be forwarded to the controller (where an appropriate forwarding rule exists); this indirection costs 1. Hence, the rules correspond to cacheable items and accesses to rules are modeled by positive requests to the corresponding items. At some chosen points in time, the caching algorithm run at the controller may decide to remove or add rules to the cache. Any such change entails a fixed cost α .¹

Tree dependencies. Note that the technical feasibility of this solution heavily depends on the rule dependencies. In the most ubiquitous scenario, the rules are prefixes of IP addresses (they are bit strings). Whenever a packet arrives, the router follows a longest matching prefix (LMP) scheme: it searches for the rule that is a prefix of the destination IP of the packet and among matching rules it

¹This cost corresponds to the transmission of a message from the controller to the router as well as the update of internal data structures of the router. Such an update of proprietary and vendor-dependent structures can be quite costly [16], but the empirical studies show it to be independent of the rule being updated [15].

chooses the longest one. In other words, if the prefixes corresponding to rules are stored in the tree², then the tree is traversed from the root downwards, and the last found rule is used. This explains why we require the cached nodes to form a subforest: leaving a less specific rule on the router while evicting a more specific one (i.e., keeping a tree node in cache while evicting its descendant) will result in a situation where packets will be forwarded according to the less specific rule, and hence potentially exit through the wrong port. The LMP scheme also ensures that the described approach is implementable: one could simply add an artificial rule at the tree root in the router (matching an empty prefix). This ensures that when no actual matching rule is found in the router (in the cache), the packet will be forwarded according to this artificial rule to the controller that stores all the rules and can handle all packets appropriately.

So far, the papers on IP rule caching avoided dependencies either assuming that rules do not overlap (a tree has a single level) [20] or by preprocessing the tree, so that the rules become non-overlapping [21, 22]. Unfortunately, this could lead to a large inflation of the routing table. A notable exception is a recent solution called CacheFlow [19]. The CacheFlow model supports dependencies even in the form of directed acyclic graphs. However, CacheFlow was evaluated only experimentally, and no worst-case guarantees were given on the overall cost of caching. Our work provides theoretical foundations for respecting tree dependencies.

Negative requests. Additionally, a rule may need to be updated. For example, due to a change communicated by a dynamic routing protocol (e.g., BGP) the action defined by a rule has to be modified. In either case, we have to update the rules at the controller: we assume that this cost is zero. (This cost is unavoidable for any algorithm, so such an assumption makes our problem only more difficult.) Furthermore, if the rule is also stored at the router, then we have to pay a fixed cost of α for updating the router (see the remark for the cost of fetches and evictions). Such penalties can be easily simulated in our model: we issue a sequence of α negative requests to the updated node. It is straightforward to show that the costs in these two models can differ by a factor of at most 2. For a formal argument, see [Appendix B](#).

Implementability. Note that the whole input (fed to a tree caching algorithm) is created at the controller: positive requests are caused by cache misses (which redirect packet to the controller) and batches of α negative requests are caused by updates sent to the dynamic routing algorithm run at the controller. Therefore, the whole tree caching algorithm can be implemented in software in the controller only. Furthermore, our algorithm is a simple counter-based scheme, which can be implemented efficiently and also fine-tuned for speed, see [Section 6](#).

Other work on forwarding table minimization. Other approaches for minimizing the number of stored rules were mostly based on *rules compression (aggregation)*, where the set of rules was replaced by another equivalent and smaller set. Optimal aggregation of a fixed routing table can be achieved by dynamic programming [12, 32], but the main challenge lies in balancing the achieved

compression and the amount of changes to the routing table in the presence of *updates* to this table. While many practical heuristics have been devised by the networking community for this problem [18, 23–25, 28, 33, 36], worst-case analyses were presented only for some restricted scenarios [6, 7]. Combining rules compression and rules caching is so far an unexplored area.

3 PRELIMINARIES

We denote the height of T by $h(T)$. For any node v , $T(v)$ denotes the subtree of T rooted at v (containing v and all its descendants). A *tree cap* rooted at v is “an upper part” of $T(v)$, i.e., it contains v and if it contains node u , then it also contains all nodes on the path from u to v . If $A \subseteq B$ are both tree caps rooted at v , then we say that A is a tree cap of B .

We assume discrete time slotted into rounds, with round $t \geq 1$ corresponding to time interval $(t - 1, t)$. In round t , the algorithm is given one (positive or negative) request to exactly one tree node and has to process it, i.e., pay associated costs (if any). Right after round t , at time t , the algorithm may arbitrarily reorganize its cache, (i) ensuring that the resulting cache is a subforest of T (i.e., if the cache contains node v , then it contains the entire $T(v)$) and (ii) preserving the cache capacity constraint. An algorithm pays α for a single node fetch or eviction. We denote the contents of the cache at round t by C_t . (As the cache changes contents only between rounds, C_t is well defined.) We assume that α is an even integer (this assumption may change costs at most by a constant factor). We assume that the algorithm starts with the empty cache.

We call a non-empty set X a *valid positive changeset* for cache C if $X \cap C = \emptyset$ and $C \cup X$ is a subforest of T , and a *valid negative changeset* if $X \subseteq C$ and $C \setminus X$ is a subforest of T . We call X a *valid changeset* if it is either valid positive or negative changeset. Note that the union of positive (negative) changesets is also a valid positive (negative) changeset. We say that the algorithm applies changeset X , if it fetches all nodes from X (for a positive changeset) and evicts all nodes from X (for a negative one). Note that not all valid changesets may be applied as the algorithm is also limited by its cache capacity (k_{ONL} for an online algorithm and k_{OPT} for the optimal offline one).

4 ALGORITHM

The algorithm TREE CACHING (TC) presented in the following is a simple scheme that follows a *rent-or-buy paradigm*: it fetches (or evicts) a changeset X if the cost associated with requests at X reaches the cost of such fetch or eviction.

More concretely, TC operates in multiple phases. The first phase starts at time 0. TC starts each phase with the empty cache and proceeds as follows. Within a phase, every node keeps a counter, which is initially zero. If at round t it pays 1 for serving the request, it increments its counter. Whenever a node is fetched or evicted from the cache, its counter is reset to zero. Note that this implies that the counter of v is equal to the number of negative (positive) requests to v since its last fetching to the cache (eviction from the cache). For a set $A \subseteq T$, we denote the sum of all counters in A at time t by $\text{cnt}_t(A)$. At time t , TC verifies whether there exists a valid changeset X , such that

- (*saturation property*) $\text{cnt}_t(X) \geq |X| \cdot \alpha$ and

²We do not have to assume that they are actually stored in a real tree; this tree is implicit in the LMP scheme.

- (maximality property) $\text{cnt}_t(Y) < |Y| \cdot \alpha$ for any valid changeset $Y \supseteq X$.

In this case, the algorithm modifies its cache applying X .

If, at time t , TC is supposed to fetch some set X , but by doing so it would exceed the cache capacity k_{ONL} , it evicts all nodes from the cache instead, and starts a new phase at time t . Such a *final eviction* might not be present in the last phase, in which case we call it *unfinished*.

In Lemma 5.1 (below), we show that at any time, all valid changesets satisfying both properties of TC are either all positive or all negative. Furthermore, right after the algorithm applies a changeset, no valid changeset satisfies saturation property.

5 ANALYSIS OF TC

Throughout the paper, we fix an input I , its partition into phases, and analyze both TC and OPT on a single fixed phase P . We denote the times at which P starts and ends by $\text{begin}(P)$ and $\text{end}(P)$, respectively, i.e., rounds in P are numbered from $\text{begin}(P) + 1$ to $\text{end}(P)$. A proof of the following technical lemma follows by induction and is presented in Appendix A.

LEMMA 5.1. *Fix any time $t > \text{begin}(P)$. For any valid changeset X for C_t , it holds that $\text{cnt}_t(X) \leq |X| \cdot \alpha$. If a changeset X is applied at time t , the following properties hold:*

- (1) X contains the node requested at round t ,
- (2) $\text{cnt}_t(X) = |X| \cdot \alpha$,
- (3) $\text{cnt}_t(Y) < |Y| \cdot \alpha$ for any valid changeset Y for C_{t+1} (note that C_{t+1} is the cache state right after application of X),
- (4) X is a tree cap of a tree from C_{t+1} if X is positive and it is a tree cap of a tree from C_t if X is negative.

In the following, we assume that no positive requests are given to nodes inside cache and no negative ones to nodes outside of it. (This does not change the behavior of TC and can only decrease the cost of OPT.)

For the sake of analysis, we assume that at time $\text{end}(P)$, TC actually performs a cache fetch (exceeding the cache size limit) and then, at the same time instant, empties the cache. This replacement only increases the cost of TC. Let k_P denote the number of nodes in the cache of TC at $\text{end}(P)$. In a finished phase, we measure it after the artificial fetch, but right before the final eviction, and thus $k_P \geq k_{\text{ONL}} + 1$; in an unfinished phase $k_P \leq k_{\text{ONL}}$.

The crucial part of our analysis that culminates in Section 5.2 is the technique of shifting requests. Namely, we modify the input sequence by shifting requests up or down the tree, so that the resulting input sequence (i) is not harder for OPT and (ii) is more structured: we may lower bound the cost of OPT on each node separately and relate it to the cost of TC.

5.1 Event Space and Fields

In our analysis, we look at a two-dimensional, discrete, spatial-temporal space, called the *event space*. The first dimension is indexed by tree nodes, whose order is an arbitrary extension of the partial order given by the tree. That is, the parent of a node v is always “above” v . The second dimension is indexed by round numbers of phase P . The space elements are called *slots*. Some slots are occupied by requests: a request at node v given at round t occupies slot (v, t) .

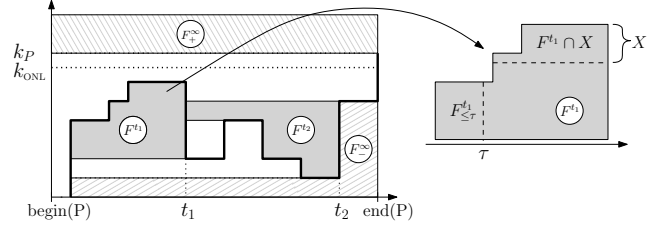


Figure 2: Partitioning of a single phase into fields for a line (a tree with no branches). The thick line represents cache contents. Possible final eviction at $\text{end}(P)$ is not depicted. F^{t_1} is a negative field and F^{t_2} is a positive one. In the particular depicted example, nodes are ordered from the leaf (bottom) to the root (top of the picture). We emphasize that for a general, branched tree, some notions (in particular fields) no longer have nice geometric interpretations.

From now on, we will identify P with a set of requests occupying some slots in the event space.

We partition slots of the whole event space into disjoint parts, called *fields*, and we show how this partition is related to the costs of TC and OPT. For any node v and time t , $\text{last}_v(t)$ denotes the last time strictly before t , when node v changed state from cached to non-cached or vice versa; $\text{last}_v(t) = \text{begin}(P)$ if v did not change its state before t in phase P . For a changeset X_t applied by TC at time t , we define the field F^t as

$$F^t = \{ (v, r) : v \in X_t \wedge \text{last}_v(t) + 1 \leq r \leq t \}.$$

That is, field F^t contains all the requests that eventually trigger the application of X_t at time t . We say that F^t ends at t . We call field F^t *positive* (*negative*) if X_t is a positive (negative) changeset. An example of a partitioning into fields is given in Figure 2. We define $\text{req}(F^t)$ as the number of requests belonging to slots of F^t and let $\text{size}(F^t)$ be the number of involved nodes (note that $\text{size}(F^t) = |X_t|$). The observation below follows immediately by Lemma 5.1.

OBSERVATION 5.2. *For any field F , $\text{req}(F) = \text{size}(F) \cdot \alpha$. All these requests are positive (negative) if F is positive (negative).*

Finally, we call the rest of the event space defined by phase P *open field* and denote it by F^∞ . The set of all fields except F^∞ is denoted by \mathcal{F} . Let $\text{size}(\mathcal{F}) = \sum_{F \in \mathcal{F}} \text{size}(F)$.

LEMMA 5.3. *For any phase P partitioned into a set of fields $\mathcal{F} \cup \{F^\infty\}$, it holds that $\text{TC}(P) \leq 2\alpha \cdot \text{size}(\mathcal{F}) + \text{req}(F^\infty) + k_P \cdot \alpha$.*

PROOF. By Observation 5.2, the cost associated with serving the requests from all fields from \mathcal{F} is $\sum_{F \in \mathcal{F}} \alpha \cdot \text{size}(F) = \alpha \cdot \text{size}(\mathcal{F})$. The cost of the cache reorganization at the fields’ ends is exactly the same. The term $\text{req}(F^\infty)$ represents the cost of serving the requests from F^∞ and $k_P \cdot \alpha$ upper-bounds the cost of the final eviction (not present in an unfinished phase). \square

5.2 Shifting Requests

The actual challenge in the proof is to relate the structure of the fields to the cost of OPT. The rationale behind our construction is based on the following thought experiment. Assume that the phase is unfinished (for example, when the cache is so large that the

whole input corresponds to a single phase). Recall that the number of requests in each field $F \in \mathcal{F}$ is equal to $\text{size}(F) \cdot \alpha$. Assume that these requests are evenly distributed among the nodes of F (each node from F receives α requests in the slots of F). Then, the history of any node v is alternating between periods spent in positive fields and periods spent in negative fields. By our even distribution assumption, each such a period contains exactly α requests. Hence, for any two consecutive periods of a single node, OPT has to pay at least α (either α for positive requests or α for negative ones, or α for changing the cached/non-cached state of v). Essentially, this shows that OPT has to pay an amount that can be easily related to $\alpha \cdot \text{size}(\mathcal{F})$.

Unfortunately, the requests may not be evenly distributed among the nodes. To alleviate this problem, we will modify the requests in phase P , so that the newly created phase P' is not harder for OPT and will “almost” have the even distribution property. In this construction, the time frame of P and its fields are fixed.

5.2.1 Legal Shifts. We say that a request placed originally (in phase P) at slot (v, t) is *legally shifted* if its new slot is $(m(v), t)$, where (i) for a positive request, $m(v)$ is either equal to v or is one of its descendants and (ii) for a negative request, $m(v)$ is either equal to v or is one of its ancestors. For any fixed sequence of fetches and evictions within phase P , the associated cost may only decrease when these actions are replayed on the modified requests.

OBSERVATION 5.4. *If P' is created from P by legally shifting the requests, then $\text{OPT}(P') \leq \text{OPT}(P)$.*

The main difficulty is however in keeping the legally shifted requests within the field they originally belonged to. For example, a negative request from F shifted at round t from node u to its parent may fall out of F as the parent may still be outside the cache at round t . In effect, a careless shifting of requests may lead to a situation where, for a single node v , requests do not create interleaved periods of positive and negative requests, and hence we cannot argue that $\text{OPT}(P')$ is sufficiently large.

In the following subsections, we show that it is possible to legally shift the requests of any field $F \in \mathcal{F}$ (i.e., shift positive requests down and negative requests up), so that they remain within F , and they will be either exactly or approximately evenly distributed among nodes of F . This will create P' with appropriately large cost for OPT.

5.2.2 Notation. We start with some general definitions and remarks. For any field F and set of nodes A , let $F \cap A = \{(v, t) \in F : v \in A\}$. Analogously, if L is a set of rounds, then let $F \cap L = \{(v, t) \in F : t \in L\}$. For any field F^t and time τ , we define

$$F_{\leq \tau}^t = F^t \cap \{t' : t' \leq \tau\}.$$

It is convenient to think that F^t evolves with time and $F_{\leq \tau}^t$ is the snapshot of F^t at time τ . Note that F^t may have some nodes not included in $F_{\leq \tau}^t$. These objects are depicted in [Figure 2](#).

We may extend the notions of req and size to arbitrary subsets of fields in a natural way. For any subset $S \subseteq F$, we call it *over-requested* if $\text{req}(S) > \text{size}(S) \cdot \alpha$.

LEMMA 5.5. *Fix any field F^t , the corresponding changeset X_t , and any time τ .*

- (1) *If F^t is negative, then for any tree cap D of X_t , the set $F_{\leq \tau}^t \cap D$ is not over-requested.*
- (2) *If F^t is positive, then for any subtree $T' \subseteq T$, the set $F_{\leq \tau}^t \cap T'$ is not over-requested.*

PROOF. As the nodes from $F_{\leq \tau}^t \cap D$ form a valid changeset at time τ , [Lemma 5.1](#) implies $\text{req}(F_{\leq \tau}^t \cap D) = \text{cnt}_\tau(F_{\leq \tau}^t \cap D) \leq |F_{\leq \tau}^t \cap D| \cdot \alpha$.

The proof of the second property is identical: As $F_{\leq \tau}^t \cap T'$ is also a valid changeset at time τ , by [Lemma 5.1](#), $\text{req}(F_{\leq \tau}^t \cap T') = \text{cnt}_\tau(F_{\leq \tau}^t \cap T') \leq |F_{\leq \tau}^t \cap T'| \cdot \alpha$. \square

By [Lemma 5.5](#) applied at $\tau = t$ and [Observation 5.2](#), we deduct the following corollary.

COROLLARY 5.6. *Fix any field F^t , the corresponding changeset X_t and any tree cap D of X_t .*

- (1) *If F^t is positive, then $\text{req}(F^t \cap D) \geq \alpha \cdot |D|$.*
- (2) *If F^t is negative, then $\text{req}(F^t \cap (X_t \setminus D)) \geq \alpha \cdot |X_t \setminus D|$.*

Informally speaking, the corollary above states that the average amount of requests in a positive field is *at least as large at the top of the field as at its bottom*. For a negative field this relation is reversed.

5.2.3 Shifting Negative Requests Up. Fix a valid negative changeset X_t applied at time t and the corresponding field F^t . We call a tree cap $Y \subseteq X_t$ *proper* if

- (1) $\text{req}(F^t \cap Y) = |Y| \cdot \alpha$ and
- (2) $F_{\leq \tau}^t \cap D$ is not over-requested for any tree cap $D \subseteq Y$ and any time $\tau \leq t$.

The first property of [Lemma 5.5](#) states that before we shift the requests of F_t , the set X_t is proper. We start with $Y = X_t$, and proceed in a bottom-up fashion, inductively using the lemma below. We take care of a single node of Y at a time and ensure that after the shift the number of requests at this node is exactly α and the remaining part of Y remains proper.

LEMMA 5.7. *Given a negative field F^t , the corresponding changeset X_t and a proper tree cap $Y \subseteq X_t$, it is possible to choose a leaf v and legally shift some requests inside Y , so that in result $\text{req}(v) = \alpha$ and $Y \setminus \{v\}$ is proper.*

PROOF. As $\text{req}(F^t \cap Y) = |Y| \cdot \alpha$, [Corollary 5.6](#) implies that any leaf of Y was requested at least α times inside F^t . We pick an arbitrary leaf v , and let $r \geq \alpha$ be the number of requests to v in F^t .

We look at all the requests to v in F^t ordered by their round. Let s be the round when $(\alpha + 1)$ -th of them arrives. We will now show that at round s , TC already has $p(v)$ in its cache. If it had not, $\{v\}$ would be a tree cap of $F_{\leq s}^t$, and by the first property of [Lemma 5.5](#), it would contain at most α requests, which is a contradiction. Hence, if we shift the chronologically last $r - \alpha$ requests from v to $p(v)$, these requests stay within F^t .

It remains to show that $Y \setminus \{v\}$ is proper after such a shift. We choose any tree cap $D \subseteq Y$ and any time $\tau \leq t$. If D does not contain $p(v)$ or $\tau < s$, then the number of requests in $F_{\leq \tau}^t \cap D$ was not changed by the shift, and hence $F_{\leq \tau}^t \cap D$ is not over-requested. Otherwise, $D \cup \{v\}$ was a tree cap in Y and by the lemma assumption, $F_{\leq \tau}^t \cap (D \cup \{v\})$ was not over-requested. As $F_{\leq \tau}^t \cap D$ has now exactly α less requests than $F_{\leq \tau}^t \cap (D \cup \{v\})$ had, it is not over-requested, either. \square

COROLLARY 5.8. *For any negative field F^t , it is possible to legally shift its requests up, so that they remain within F^t and after the modification each node is requested exactly α times.*

5.2.4 Shifting Positive Requests Down. We will now focus on the problem of shifting the positive requests down in a single positive field F^t , corresponding to a single fetch of TC at the time t . Our goal is to devise a shifting strategy, that will result in at least $\Omega(\text{size}(F^t)/h(T))$ nodes having $\alpha/2$ requests each. While this result may be suboptimal, deriving a shifting strategy for a positive field that would have the same equal distribution guarantee as the one provided by [Corollary 5.8](#) is not possible (the details are presented in the full version of the paper).

First, we prove that from any node v in the field, we can shift down a constant fraction of its requests within the field, distributing them to different nodes.

LEMMA 5.9. *Let F^t be a positive field and let X_t be the corresponding changeset fetched to the cache at time t . Fix any node $v \in X_t$ that has been requested at least $c \cdot (\alpha/2)$ times in F^t , where c is an integer. It is possible to shift down its requests to the nodes of $T(v) \cap X_t$, so that these requests remain inside F^t and $\lceil c/2 \rceil$ nodes of $T(v)$ get $\alpha/2$ requests each.*

PROOF. We order the nodes $u_1, u_2, \dots, u_{|T(v) \cap X_t|}$ of $T(v) \cap X_t$, so that $\text{last}_{u_i}(t) \leq \text{last}_{u_{i+1}}(t)$ for all i . In case of a tie, we place nodes that are closer to v first. Note that this linear ordering is an extension of the partial order defined by the tree: the parent of a node cannot be evicted later than the node itself (otherwise the cache would cease to be a subforest of T). In particular, it holds that $u_1 = v$.

We number $c \cdot (\alpha/2)$ requests to v chronologically, starting from 1. For any $j \in \{1, \dots, \lceil c/2 \rceil\}$ we look at round τ_j with the $((j-1) \cdot \alpha + 1)$ -th request to v . When this request arrives, node u_j is already present in the cache. Otherwise, we would have at least $j \cdot \alpha + 1$ requests in $F^t_{\leq \tau_j} \cap \{u_1, \dots, u_j\}$ (already in $F^t_{\leq \tau_j} \cap \{u_1\}$ alone), which would make it over-requested, and thus contradict the second property of [Lemma 5.5](#). Hence, we may take requests numbered from $(j-1) \cdot \alpha + 1$ to $(j-1) \cdot \alpha + \alpha/2$, shift them down from v to u_j , and after such modification these requests are still inside F^t . Note that for $j = 1$ requests are not really shifted, as u_1 is v itself. We perform such shift for any $j \in \{1, \dots, \lceil c/2 \rceil\}$, which yields the lemma. \square

LEMMA 5.10. *For any positive field F^t , it is possible to legally shift its requests down, so that they remain within F^t and after the modification at least $\text{size}(F^t)/(2h(T))$ nodes in F^t have at least $\alpha/2$ requests each.*

PROOF. Let X_t be the changeset corresponding to field F^t , which is fetched to the cache at time t . By [Observation 5.2](#), $\text{req}(F^t) = |X_t| \cdot \alpha$. We gather the requests at every node into groups of $\alpha/2$ consecutive requests. In every node at most $\alpha/2$ requests remain not grouped. Let $\overline{\text{req}}(X)$ denote the number of grouped requests in the set X . Clearly, $\overline{\text{req}}(F^t) \geq |X_t| \cdot \alpha/2$, i.e., there are at least $|X_t|$ groups of requests in set X_t .

Let $X_t = X_t^1 \sqcup X_t^2 \sqcup \dots \sqcup X_t^{h(T)}$ be a partition of the nodes of the tree X_t into layers according to their distance to the root. By the pigeonhole principle, there is a layer X_t^i containing at least $\lceil |X_t|/h(T) \rceil$ groups of requests (each group has $\alpha/2$ requests).

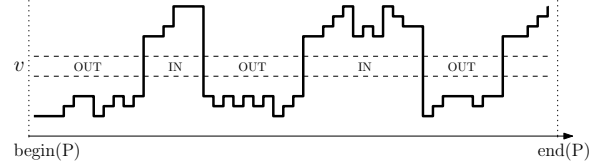


Figure 3: Partitioning of the phase into interleaving IN and OUT periods for node v . The thick line represents cache contents. The leftover OUT period (the last one) is present for node v as it has finished phase P inside TC's cache. The periods can be followed by requests contained in F^∞ .

Nodes of X_t^i are independent, i.e., for $u, v \in X_t^i$ the trees $T(u)$ and $T(v)$ are disjoint. Therefore, we may use the shifting strategy described in [Lemma 5.9](#) for each node of X_t^i separately. After such modification, at least $\lceil |X_t|/(2h(T)) \rceil \geq \text{size}(F_t)/(2h(T))$ nodes have at least $\alpha/2$ requests each. \square

5.2.5 Using Request Shifting for Bounding OPT. Finally, we may use our request shifting to relate $\text{size}(\mathcal{F}) = \sum_{F \in \mathcal{F}} \text{size}(F)$ to the cost of OPT in a single phase P . Recall that k_P denotes the size of TC's cache at the end of P . We assume that OPT may start the phase with an arbitrary state of the cache.

LEMMA 5.11. *For any phase P , $\text{OPT}(P) \geq (\text{size}(\mathcal{F})/(4h(T)) - k_P) \cdot \alpha/2$.*

PROOF. We transform P using legal shifts that are described in [Section 5.2.3](#) and [Section 5.2.4](#). That is, we create a corresponding phase P' that satisfies both [Corollary 5.8](#) and [Lemma 5.10](#). By [Observation 5.4](#), it is sufficient to show that $\text{OPT}(P') \geq (\text{size}(\mathcal{F})/(4h(T)) - k_P) \cdot \alpha/2$.

We focus on a single node v . We cut its history into interleaved periods: OUT periods, when v is outside the cache and receives positive requests, and IN periods when TC keeps v in the cache and v receives negative requests. A final (possibly empty) part corresponding to the time when v is in the F^∞ field is not accounted in OUT or IN periods, i.e., each IN or OUT period corresponds to some field $F \in \mathcal{F}$. Let p^{IN} and p^{OUT} denote the total number of IN and OUT periods (respectively) for all nodes during the phase. An example is given in [Figure 3](#).

Recall that TC starts each phase with an empty cache, and hence each node starts with an OUT period. For k_P nodes that are in TC's cache at the end of the phase (and only for them) their history ends with an OUT period not followed by an IN period. We call them *leftover periods*. Thus, $p^{\text{OUT}} = p^{\text{IN}} + k_P$. The total number of periods ($p^{\text{IN}} + p^{\text{OUT}}$) is equal to the total size of all fields, $\text{size}(\mathcal{F})$, and thus $p^{\text{OUT}} \geq \text{size}(\mathcal{F})/2$.

We call a period *full* if it has at least $\alpha/2$ requests. The shifting strategies described in the previous section ensure that all IN periods are full and at least $1/(2h(T))$ of all OUT periods are full. Thus, there are at least $p^{\text{OUT}}/(2h(T)) - k_P$ full non-leftover OUT periods; each of them together with the following IN period constitutes a *full OUT-IN pair*.

OPT has to pay at least $\alpha/2$ for the node in the course of the history described by a full OUT-IN pair: it pays α either for changing the cached/non-cached state of a node, or $\alpha/2$ for all positive requests

or $\alpha/2$ for all negative ones. Thus, $\text{OPT}(P') \geq (p^{\text{OUT}}/(2h(T)) - k_P) \cdot \alpha/2 \geq (\text{size}(\mathcal{F})/(4h(T)) - k_P) \cdot \alpha/2$. \square

5.3 Competitive Ratio

To relate the cost of OPT to TC in a single phase P , we still need to upper-bound $\text{req}(F^\infty)$ and relate $k_P \cdot \alpha$ to the cost of OPT (i.e., compare the bounds on TC and OPT provided by Lemma 5.3 and Lemma 5.11, respectively).

For the next two lemmas, we define V_{OPT} as the set of all nodes that were in OPT cache at some time of P and let $V_{\text{OPT}}^c = T \setminus V_{\text{OPT}}$. Note that V_{OPT} is a union of subforests (nodes present in OPT's cache at consecutive times), and hence a subforest itself.

LEMMA 5.12. *For any phase P , it holds that $\text{req}(F^\infty) \leq 2 \cdot k_{\text{ONL}} \cdot \alpha + 2 \cdot \text{OPT}(P)$.*

PROOF. We assume first that P is a finished phase. Then, P ends with an artificial fetch of $X_{\text{end}(P)}$ at time $\text{end}(P)$ (followed by the final eviction). We split F^∞ into two disjoint parts (see Figure 2):

$$F_-^\infty = \{(v, t) : v \in C_{\text{end}(P)}, t \geq \text{last}_v(\text{end}(P))\},$$

$$F_+^\infty = \{(v, t) : v \notin C_{\text{end}(P)} \sqcup X_{\text{end}(P)}, t \geq \text{last}_v(\text{end}(P))\}.$$

Note that F_-^∞ contains only negative requests and F_+^∞ only positive ones. As $\text{req}(F^\infty) = \text{req}(F_-^\infty) + \text{req}(F_+^\infty \cap V_{\text{OPT}}^c) + \text{req}(F_+^\infty \cap V_{\text{OPT}})$, we estimate each of these summands separately.

- Nodes from F_-^∞ are in the cache $C_{\text{end}(P)}$ and were not evicted from the cache. Thus, $\text{req}(F_-^\infty) \leq |C_{\text{end}(P)}| \cdot \alpha \leq k_{\text{ONL}} \cdot \alpha$.
- All the requests from V_{OPT}^c are paid by OPT, and hence $\text{req}(F_+^\infty \cap V_{\text{OPT}}^c) \leq \text{req}(V_{\text{OPT}}^c) \leq \text{OPT}(P)$.
- F_+^∞ is a valid changeset for cache $C_{\text{end}(P)} \sqcup X_{\text{end}(P)}$. As V_{OPT} is a subforest of T , $F_+^\infty \cap V_{\text{OPT}}$ is also a valid changeset for the cache $C_{\text{end}(P)} \sqcup X_{\text{end}(P)}$. Therefore, $\text{req}(F_+^\infty \cap V_{\text{OPT}}) \leq \text{size}(F_+^\infty \cap V_{\text{OPT}}) \cdot \alpha$, as otherwise the set fetched at time $\text{end}(P)$ would not be maximal. (TC could then fetch $X_{\text{end}(P)} \sqcup (F_+^\infty \cap V_{\text{OPT}})$ instead of $X_{\text{end}(P)}$.) Thus, $\text{req}(F_+^\infty \cap V_{\text{OPT}}) \leq |V_{\text{OPT}}| \cdot \alpha = k_{\text{OPT}} \cdot \alpha + (|V_{\text{OPT}}| - k_{\text{OPT}}) \cdot \alpha \leq k_{\text{ONL}} \cdot \alpha + \text{OPT}(P)$. The last inequality follows as — independently of the initial state — OPT needs to fetch at least $|V_{\text{OPT}}| - k_{\text{OPT}}$ nodes to the cache during P .

Hence, in total, $\text{req}(F^\infty) \leq 2 \cdot k_{\text{ONL}} \cdot \alpha + 2 \cdot \text{OPT}(P)$ for a finished phase P .

We note that if there was no cache change at $\text{end}(P)$, the analysis above would hold with $X_{\text{end}(P)} = \emptyset$ with virtually no change. Therefore, for an unfinished phase P ending with a fetch or ending without cache change at $\text{end}(P)$, the bound on $\text{req}(F^\infty)$ still holds. However, if an unfinished phase P ends with an eviction, then we look at the last eviction-free time τ of P . We now observe the evolution of field F^∞ from time τ till $\text{end}(P)$. At time τ , $\text{req}(F^\infty) \leq 2 \cdot k_{\text{ONL}} \cdot \alpha + 2 \cdot \text{OPT}(P)$. Furthermore, in subsequent times, it may only decrease: at any round F^∞ gets an additional request, but on eviction $\text{req}(F^\infty)$ decreases by α times the number of evicted nodes (i.e., at least by $\alpha \geq 1$). Hence, the value of $\text{req}(F^\infty)$ at $\text{end}(P)$ is also at most $2 \cdot k_{\text{ONL}} \cdot \alpha + 2 \cdot \text{OPT}(P)$. \square

By combining Lemma 5.3, Lemma 5.11 and Lemma 5.12, we immediately obtain the following corollary (holding for both finished and unfinished phases).

COROLLARY 5.13. *For any phase P , it holds that $\text{TC}(P) \leq O(h(T)) \cdot \text{OPT}(P) + O(h(T) \cdot (k_P + k_{\text{ONL}}) \cdot \alpha)$.*

Using the corollary above, it remains to bound the value of k_P . This is easy for an unfinished phase, as $k_P \leq k_{\text{ONL}}$ there. For a finished phase, we provide another bound.

LEMMA 5.14. *For any finished phase P , it holds that $k_P \cdot \alpha \leq \text{OPT}(P) \cdot (k_{\text{ONL}} + 1)/(k_{\text{ONL}} + 1 - k_{\text{OPT}})$.*

PROOF. First, we compute the number of positive requests in V_{OPT}^c . Let $X_{t_1}, X_{t_2}, \dots, X_{t_s}$ be all positive changesets applied by TC in P . For any t , let $X'_t = X_t \setminus V_{\text{OPT}}$. As X_t is some tree cap and V_{OPT} is a subforest of T , X'_t is a tree cap of X_t . By Corollary 5.6, the number of requests to nodes of X'_t in field F^t is at least $|X'_t| \cdot \alpha$. These requests for different changesets X_t are disjoint and they are all outside of V_{OPT} . Hence the total number of positive requests outside of V_{OPT} is at least $\sum_{i=1}^s |X'_{t_i}| \cdot \alpha$, where $\sum_{i=1}^s |X'_{t_i}| \geq |\bigcup_{i=1}^s X'_{t_i}| = |(\bigcup_{i=1}^s X_{t_i}) \setminus V_{\text{OPT}}| \geq |\bigcup_{i=1}^s X_{t_i}| - |V_{\text{OPT}}| \geq k_P - |V_{\text{OPT}}|$.

Now $\text{OPT}(P)$ can be split into the cost associated with nodes from V_{OPT} and V_{OPT}^c , respectively. For the former part, OPT has to pay at least $(|V_{\text{OPT}}| - k_{\text{OPT}}) \cdot \alpha$ for the fetches alone. For the latter part, it has to pay 1 for each of at least $(k_P - |V_{\text{OPT}}|) \cdot \alpha$ positive requests outside of V_{OPT} . Hence, $\text{OPT}(P) \geq (|V_{\text{OPT}}| - k_{\text{OPT}}) \cdot \alpha + (k_P - |V_{\text{OPT}}|) \cdot \alpha = (k_P - k_{\text{OPT}}) \cdot \alpha$. Then, $k_P \cdot \alpha \leq k_P \cdot \text{OPT}(P)/(k_P - k_{\text{OPT}})$. As the phase is finished, $k_P \geq k_{\text{ONL}} + 1$, and thus $k_P \cdot \alpha \leq (k_{\text{ONL}} + 1) \cdot \text{OPT}(P)/(k_{\text{ONL}} + 1 - k_{\text{OPT}})$. \square

THEOREM 5.15. *The algorithm TC is $O(h(T) \cdot k_{\text{ONL}}/(k_{\text{ONL}} - k_{\text{OPT}} + 1))$ -competitive.*

PROOF. Let $R = h(T) \cdot k_{\text{ONL}}/(k_{\text{ONL}} - k_{\text{OPT}} + 1)$. We split an input I into a sequence of finished phases followed by a single unfinished phase (which may not be present). For a finished phase P , we have $k_P > k_{\text{ONL}}$, and hence Corollary 5.13 and Lemma 5.14 imply that $\text{TC}(P) \leq O(R) \cdot \text{OPT}(P)$. For an unfinished phase $k_P \leq k_{\text{ONL}}$, and therefore, by Corollary 5.13, $\text{TC}(P) \leq O(h(T)) \cdot \text{OPT}(P) + O(h(T) \cdot k_{\text{ONL}} \cdot \alpha)$. Summing over all phases of I yields $\text{TC}(I) \leq O(R) \cdot \text{OPT}(I) + O(h(T) \cdot k_{\text{ONL}} \cdot \alpha)$. \square

6 IMPLEMENTATION OF TC

Recall that at each time t , TC verifies the existence of a valid changeset that satisfies saturation and maximality properties (see the definition of TC in Section 4). Here, we show that this operation can be performed efficiently. In particular, in the following two subsections, we will prove the following theorem.

THEOREM 6.1. *TC can be implemented using $O(|T|)$ additional memory, so that to make a decision at time t , it performs $O(h(T) + \max\{h(T), \deg(T)\} \cdot |X_t|)$ operations, where $\deg(T)$ is a maximum node degree in T and X_t is the changeset applied at time t ($|X_t| = 0$ if no changeset is applied).*

Let v_t be the node requested at round t . Note that we may restrict our attention to requests that entail a cost for TC, as otherwise its counters remain unchanged and certainly TC does not change

cache contents. We use [Lemma 5.1](#) to restrict possible candidates for changesets that can be applied at time t . First, we note that if a node v_t requested at round t is outside the cache, then, at time t , TC may only fetch some changeset, and otherwise it may only evict some changeset. Therefore, we may construct two separate schemes, one governing fetches and one for evictions.

In [Section 6.1](#), using [Lemma 5.1](#), we show that after processing a positive request, TC needs to verify at most $h(T)$ possible positive changesets, each in constant time, using an auxiliary data structure. The cost of updating this structure at time t is $O(h(T) + h(T) \cdot |X_t|)$.

The situation for negative changesets is more complex as even after applying [Lemma 5.1](#) there are still exponentially many valid negative changesets to consider. In [Section 6.2](#), we construct an auxiliary data structure that returns a viable candidate in time $O(h(T) + \deg(T) \cdot |X_t|)$. The update of this structure at time t can be also done in $O(h(T) + \deg(T) \cdot |X_t|)$ operations.

6.1 Positive Requests and Fetches

At any time t and for any non-cached node u , we may define $P_t(u)$ as a tree cap rooted at u containing all non-cached nodes from $T(u)$. During an execution of TC, we maintain two values for each non-cached node u : $\text{cnt}_t(P_t(u))$ and $|P_t(u)|$. When a counter at node v_t is incremented, we update $\text{cnt}_t(P_t(u))$ for each ancestor u of v (at most $h(T)$ updated values). Furthermore, if a node v changes its state from cached to non-cached (or vice versa), we update the value of $|P_t(u)|$ for any ancestor u of v (at most $h(T)$ updates per each node that changes the state). Therefore, the total cost of updating these structures at time t is at most $O(h(T) + h(T) \cdot |X_t|)$.

By [Lemma 5.1](#), a positive valid changeset fetched at time t has to contain v_t and is a single tree cap. Such a tree cap has to be equal to $P_t(u)$ for u being an ancestor of v_t . Hence, we may iterate over all ancestors u of v_t , starting from the tree root and ending at v_t , and we stop at the first node u , for which $P_t(u)$ is saturated (i.e., $\text{cnt}_t(P_t(u)) \geq |P_t(u)| \cdot \alpha$). If such a u is found, the corresponding set $P_t(u)$ satisfies also the maximality condition (cf. the definition of TC) as all valid changesets that are supersets of $P_t(u)$ were already verified to be non-saturated. Therefore, in such a case, TC fetches $P_t(u)$. Otherwise, if no saturated changeset is found, TC does nothing. Checking all ancestors of v_t can be performed in time $O(h(T))$.

6.2 Negative Requests and Evictions

Handling evictions is more complex. If the request to node v_t at round t was negative, [Lemma 5.1](#) tells us only that the negative changeset evicted by TC has to be a tree cap rooted at u , where u is the root of the cached tree containing v_t . There are exponentially many such tree caps, and hence their naïve verification is intractable. To alleviate this problem, we introduce the following helper notion. For any set of cached nodes A and any time t , let

$$\text{val}_t(A) = \text{cnt}_t(A) - |A| \cdot \alpha + \frac{|A|}{|T| + 1}.$$

Note that for any non-empty set A , $\text{val}_t(A) \neq 0$ as the first two terms are integers and $|A|/(|T| + 1) \in (0, 1)$. Furthermore, val_t is additive: for two disjoint sets A and B , $\text{val}_t(A \sqcup B) = \text{val}_t(A) + \text{val}_t(B)$. For

any time t and a cached node u , we define

$$H_t(u) = \arg \max_D \{\text{val}_t(D) : D \text{ is a non-empty tree cap rooted at } u\}.$$

Our scheme maintains the value $H_t(u)$ for any cached node u . To this end, we observe that $H_t(u)$ can be defined recursively as follows. Let $H'_t(u) = H_t(u)$ if $\text{val}_t(H_t(u)) > 0$ and $H'_t(u) = \emptyset$ otherwise. Then, for any node v and time t , by the additivity of val_t ,

$$H_t(u) = \{u\} \sqcup \bigsqcup_{w \text{ is a child of } u} H'_t(w).$$

Each cached node u keeps the value $\text{val}_t(H_t(u))$. Note that set $H_t(u)$ itself can be recovered from this information: we iterate over all children of u (at most $\deg(T)$ of them) and for each child w , if $\text{val}_t(H_t(w)) > 0$, we recursively compute set $H_t(w)$. Thus, the total time for constructing $H_t(u)$ is $O(\deg(T) \cdot |H_t(u)|)$.

During an execution of TC, we update stored values accordingly. That is, whenever a counter at a cached node v_t is incremented, we update $\text{val}_t(H_t(u))$ values for each cached ancestor u of v_t , starting from $u = v_t$ and proceeding towards the cached tree root. Any such update can be performed in constant time, and the total time is thus $O(h(T))$. For a cache change, we process nodes from the changeset iteratively, starting with nodes closest to the root in case of an eviction and furthest from the root in case of a fetch. For any such node u , we appropriately stop or start maintaining the corresponding value of $\text{val}_t(H_t(u))$. The latter requires looking up the stored values at all its children. As u does not have cached ancestors, sets H_t (and hence also the stored values) at other nodes remain unchanged. In total, the cost of updating all H_t values at time t is at most $O(h(T) + \deg(T) \cdot |X_t|)$.

Finally, we show how to use sets H_t to quickly choose a valid changeset for eviction. Recall that for a negative request v_t , the changeset to be evicted has to be a tree cap rooted at u , where u is the root of a cached subtree containing v_t . For succinctness, we use H^u to denote $H_t(u)$. We show that if $\text{val}_t(H^u) < 0$, then there is no valid negative changeset that is saturated, and hence TC does not perform any action, and if $\text{val}_t(H^u) > 0$, then H^u is both saturated and maximal, and hence TC may evict H^u .

- (1) First, assume that $\text{val}_t(H^u) < 0$. Then, for any tree cap X rooted at u , it holds that $\text{cnt}_t(X) - |X| \cdot \alpha < \text{val}_t(X) \leq \text{val}_t(H^u) < 0$, i.e., X is not saturated, and hence cannot be evicted by TC.
- (2) Second, assume that $\text{val}_t(H^u) > 0$. As $\text{cnt}_t(H^u) - |H^u| \cdot \alpha$ is an integer and $|H^u|/(|T| + 1) < 1$, it holds that $\text{cnt}_t(H^u) - |H^u| \cdot \alpha \geq 0$, i.e., H^u is saturated. Moreover, by [Lemma 5.1](#), $\text{cnt}_t(H^u) \leq |H^u| \cdot \alpha$, and therefore $\text{cnt}_t(H^u) - |H^u| \cdot \alpha = 0$, i.e., $\text{val}_t(H^u) = |H^u|/(|T| + 1)$. It remains to show that H^u is maximal, i.e., there is no valid saturated changeset $Y \supsetneq H^u$. By [Lemma 5.1](#), Y has to be a tree cap rooted at u as well. If Y was saturated, $\text{val}_t(Y) = \text{cnt}_t(Y) - |Y| \cdot \alpha + |Y|/(|T| + 1) \geq |Y|/(|T| + 1) > |H^u|/(|T| + 1) = \text{val}_t(H^u)$, which would contradict the definition of H^u .

Note that node u can be found in time $O(h(T))$, and the actual set H^u (of size $|X_t|$) can be computed in time $O(\deg(T) \cdot |X_t|)$. Therefore the total time for finding set $|X_t|$ is $O(h(T) + \deg(T) \cdot |X_t|)$.

7 CONCLUSIONS

This paper defines a novel variant of online paging which finds applications in the context of IP routing networks where forwarding rules can be cached. We presented a deterministic online algorithm that achieves a provably competitive trade-off between the benefit of caching and update costs.

It is worth noting that, in the offline setting, choosing the best static cache in the presence of only positive requests is known as a *tree sparsity* problem and can be solved in $O(|T|^2)$ time [4].

We believe that our work opens interesting directions for future research. Most importantly, it will be interesting to study the optimality of the derived result; we conjecture that the true competitive ratio does not depend on the tree height. In particular, primal-dual approaches that were successfully applied for other caching problems [3, 5, 34] may turn out to be useful also for the considered variant.

ACKNOWLEDGEMENTS

The authors would like to thank Fred Baker from Cisco, Moti Medina from the Max-Planck-Institute and Paweł Gawrychowski from University of Wrocław for useful inputs.

REFERENCES

- [1] BGP Statistics from Route-Views Data. <http://bgp.potaroo.net/bgprrpts/rva-index.html>.
- [2] Dimitris Achlioptas, Marek Chrobak, and John Noga. 2000. Competitive analysis of randomized paging algorithms. *Theoretical Computer Science* 234, 1–2 (2000), 203–218.
- [3] Anna Adamaszek, Artur Czumaj, Matthias Englert, and Harald Räcke. 2012. An $O(\log k)$ -competitive algorithm for generalized caching. In *23rd ACM-SIAM Symp. on Discrete Algorithms (SODA)*. 1681–1689.
- [4] Arturs Backurs, Piotr Indyk, and Ludwig Schmidt. 2017. Better Approximations for Tree Sparsity in Nearly-Linear Time. In *Proc. 28th ACM-SIAM Symp. on Discrete Algorithms (SODA)*. 2215–2229.
- [5] Nikhil Bansal, Niv Buchbinder, and Joseph Naor. 2012. Randomized Competitive Algorithms for Generalized Caching. *SIAM J. Comput.* 41, 2 (2012), 391–414.
- [6] Marcin Bienkowski, Nadi Sarrar, Stefan Schmid, and Steve Uhlig. 2014. Competitive FIB Aggregation without Update Churn. In *Proc. 34th IEEE Int. Conf. on Distributed Computing Systems (ICDCS)*. 607–616.
- [7] Marcin Bienkowski and Stefan Schmid. 2013. Competitive FIB Aggregation for Independent Prefixes: Online Ski Rental on the Trie. In *Proc. 20th Int. Colloq. on Structural Information and Communication Complexity (SIROCCO)*. 92–103.
- [8] Mark Brehob, Richard J. Enbody, Eric Torng, and Stephen Wagner. 2003. On-line Restricted Caching. *Journal of Scheduling* 6, 2 (2003), 149–166.
- [9] Niv Buchbinder, Shahar Chen, and Joseph Naor. 2014. Competitive Algorithms for Restricted Caching and Matroid Caching. In *Proc. 22th European Symp. on Algorithms (ESA)*. 209–221.
- [10] Marek Chrobak, Howard J. Karloff, Thomas H. Payne, and Sundar Vishwanathan. 1991. New Results on Server Problems. *SIAM Journal on Discrete Mathematics* 4, 2 (1991), 172–181.
- [11] Luca Cittadini, Wolfgang Muhlauer, Steve Uhlig, Randy Bushy, Pierre Francois, and Olaf Maennel. 2010. Evolution of internet address space deaggregation: myths and reality. *IEEE J. Sel. A. Commun.* 28, 8 (2010), 1238–1249.
- [12] Richard P. Draves, Christopher King, Srinivasan Venkatachary, and Brian D. Zill. 1999. Constructing optimal IP routing tables. In *Proc. 18th IEEE Int. Conf. on Computer Communications (INFOCOM)*. 88–97.
- [13] Leah Epstein, Csanád Imreh, Asaf Levin, and Judit Nagy-György. 2015. Online File Caching with Rejection Penalties. *Algorithmica* 71, 2 (2015), 279–306.
- [14] Amos Fiat, Richard M. Karp, Michael Luby, Lyle A. McGeoch, Daniel D. Sleator, and Neal E. Young. 1991. Competitive paging algorithms. *Journal of Algorithms* 12, 4 (1991), 685–699.
- [15] Pierre François, Clarence Filsfils, John Evans, and Olivier Bonaventure. 2005. Achieving sub-second IGP convergence in large IP networks. *ACM SIGCOMM Computer Communication Review* 35, 3, 35–44.
- [16] Danny Yuxing Huang, Ken Yocum, and Alex C. Snoeren. 2013. High-fidelity switch models for software-defined network emulation. In *Proc. 2nd ACM SIGCOMM Workshop on Hot Topics in Software Defined Networking (HotSDN)*. 43–48.
- [17] Sandy Irani. 2002. Page Replacement with Multi-Size Pages and Applications to Web Caching. *Algorithmica* 33, 3 (2002), 384–409.
- [18] Elliott Karpilovsky, Matthew Caesar, Jennifer Rexford, Aman Shaikh, and Jacobus E. van der Merwe. 2012. Practical Network-Wide Compression of IP Routing Tables. *IEEE Transactions on Network and Service Management* 9, 4 (2012), 446–458.
- [19] Naga Katta, Omid Alipourfard, Jennifer Rexford, and David Walker. 2016. CacheFlow: Dependency-Aware Rule-Caching for Software-Defined Networks. In *Proc. ACM Symposium on SDN Research (SOSR)*.
- [20] Changhoon Kim, Matthew Caesar, Alexandre Gerber, and Jennifer Rexford. 2009. Revisiting Route Caching: The World Should Be Flat. In *Proc. 10th Int. Conf. on Passive and Active Network Measurement (PAM)*. 3–12.
- [21] Huan Liu. 2001. Routing prefix caching in network processor design. In *Proc. 10th Int. Conf. on Computer Communications and Networks (ICCCN)*. 18–23.
- [22] Yaoqing Liu, Vince Lehman, and Lan Wang. 2015. Efficient FIB caching using minimal non-overlapping prefixes. *Computer Networks* 83 (2015), 85–99.
- [23] Yaoqing Liu, Beichuan Zhang, and Lan Wang. 2013. FIFA: Fast incremental FIB aggregation. In *Proc. 32nd IEEE Int. Conf. on Computer Communications (INFOCOM)*. 1213–1221.
- [24] Yaoqing Liu, Xin Zhao, Kyuhan Nam, Lan Wang, and Beichuan Zhang. 2010. Incremental Forwarding Table Aggregation. In *Proc. Global Communications Conference (GLOBECOM)*. 1–6.
- [25] Layong Luo, Gaogang Xie, Kavé Salamatian, Steve Uhlig, Laurent Mathy, and Yingke Xie. 2013. A trie merging approach with incremental updates for virtual routers. In *Proc. 32nd IEEE Int. Conf. on Computer Communications (INFOCOM)*. 1222–1230.
- [26] Lyle A. McGeoch and Daniel D. Sleator. 1991. A Strongly Competitive Randomized Paging Algorithm. *Algorithmica* 6, 6 (1991), 816–825.
- [27] Manor Mendel and Steven S. Seiden. 2004. Online companion caching. *Theoretical Computer Science* 324, 2–3 (2004), 183–200.
- [28] Gábor Rétfvári, János Tapolcai, Attila Korösi, András Majdán, and Zolán Heszberger. 2013. Compressing IP forwarding tables: towards entropy bounds and beyond. In *Proc. ACM SIGCOMM Conference*. 111–122.
- [29] Nadi Sarrar, Steve Uhlig, Anja Feldmann, Rob Sherwood, and Xin Huang. 2012. Leveraging Zipf’s law for traffic offloading. *ACM SIGCOMM Computer Communication Review* 42, 1 (2012), 16–22.
- [30] Daniel D. Sleator and Robert E. Tarjan. 1985. Amortized efficiency of list update and paging rules. *Commun. ACM* 28, 2 (1985), 202–208.
- [31] Ed Spitznagel, David E. Taylor, and Jonathan S. Turner. 2003. Packet Classification Using Extended TCAMs. In *Proc. 11th IEEE Int. Conf. on Network Protocols (ICNP)*. 120–131.
- [32] Subhash Suri, Tuomas Sandholm, and Priyank Ramesh Warkhede. 2003. Compressing Two-Dimensional Routing Tables. *Algorithmica* 35, 4 (2003), 287–300.
- [33] Zartash Afzal Uzmi, Markus E. Nebel, Ahsan Tariq, Sana Jawad, Ruichuan Chen, Aman Shaikh, Jia Wang, and Paul Francis. 2011. SMALTA: practical and near-optimal FIB aggregation. In *Proc. 7th Int. Conf. on Emerging Networking Experiments and Technologies (CoNEXT)*.
- [34] Neal E. Young. 1994. The k-Server Dual and Loose Competitiveness for Paging. *Algorithmica* 11, 6 (1994), 525–541.
- [35] Neal E. Young. 2002. On-Line File Caching. *Algorithmica* 33, 3 (2002), 371–383.
- [36] Xin Zhao, Yaoqing Liu, Lan Wang, and Beichuan Zhang. 2010. On the aggregatability of router forwarding tables. In *Proc. 29th IEEE Int. Conf. on Computer Communications (INFOCOM)*. 848–856.

A PROOF OF LEMMA 5.1

Before proving Lemma 5.1, we present the following technical claim.

CLAIM A.1. *For any phase P , the following invariants hold for any time $t > \text{begin}(P)$:*

- (1) $\text{cnt}_{t-1}(X) < |X| \cdot \alpha$ for a valid changeset X for C_t ,
- (2) $\text{cnt}_t(X) \leq |X| \cdot \alpha$ for a valid changeset X for C_t ,
- (3) any changeset X with property $\text{cnt}_t(X) = |X| \cdot \alpha$ contains the node requested at round t .

PROOF. First observe that Invariant 1 (for time t) along with the fact that round t contains only one request immediately implies that $\text{cnt}_t(X) \leq \text{cnt}_{t-1}(X) + 1 \leq (|X| \cdot \alpha - 1) + 1 = |X| \cdot \alpha$, i.e., Invariant 2 for time t . Furthermore the equality may hold only for changesets containing the node requested at round t , which implies Invariant 3 for time t .

It remains to show that **Invariant 1** holds for any step $t > \text{begin}(P)$. It is trivially true for $t = \text{begin}(P) + 1$ as $\text{cnt}_{t-1}(X) = 0$ then. Let $t + 1$ be the earliest time in phase P for which **Invariant 1** does not hold; we will then show a contradiction with the definition of TC or a contradiction with other Invariants at time t . That is, we assume that there exists a positive changeset X for C_{t+1} such that $\text{cnt}_t(X) \geq |X| \cdot \alpha$ (the proof for a negative changeset is analogous). Note that TC must have performed an action (fetch or eviction) at time t as otherwise X would be also a changeset for $C_t = C_{t+1}$ with $\text{cnt}_t(X) \geq |X| \cdot \alpha$, which means that X should have been applied by TC at time t . We consider two cases.

If TC fetches a positive changeset Y at time t , $C_{t+1} = C_t \sqcup Y$ and $\text{cnt}_t(Y) = |Y| \cdot \alpha$. Then, $Y \sqcup X$ is a changeset for C_t , and $\text{cnt}_t(Y \sqcup X) \geq |Y \sqcup X| \cdot \alpha$. This contradicts the maximality property of set Y chosen at time t by TC.

If TC evicts a negative changeset Y at time t , $C_{t+1} = C_t \setminus Y$. **Invariant 2** and the definition of TC implies $\text{cnt}_t(Y) = |Y| \cdot \alpha$, and thus, by **Invariant 3**, Y contains the node requested at round t . As $X \cap Y \subseteq C_t$, $X \cap Y$ does not have any positive requests at time t , and therefore $\text{cnt}_t(X \setminus Y) = \text{cnt}_t(X) \geq |X| \cdot \alpha \geq |X \setminus Y| \cdot \alpha$. By **Invariant 2**, $\text{cnt}_t(X \setminus Y) \leq |X \setminus Y| \cdot \alpha$, and hence $\text{cnt}_t(X \setminus Y) = |X \setminus Y| \cdot \alpha$. This contradicts **Invariant 3** as $X \setminus Y$ cannot contain the node requested at round t (because Y contains this node). \square

PROOF OF LEMMA 5.1. The inequality $\text{cnt}_t(X) \leq |X| \cdot \alpha$ is equivalent to **Invariant 2** of **Claim A.1**. Assume now that X is applied at time t . By the definition of TC, $\text{cnt}_t(X) \geq |X| \cdot \alpha$, and thus $\text{cnt}_t(X) = |X| \cdot \alpha$, i.e., **Property 2** follows. Then, **Invariant 3** of **Claim A.1** implies **Property 1**. Finally, **Invariant 1** of **Claim A.1** for time $t + 1$ is equivalent to **Property 3**.

To show **Property 4**, observe that the changeset X applied at time t cannot be a disjoint union of two (or more) valid changesets X_1 and X_2 . By **Property 2**, $|X| \cdot \alpha = \text{cnt}_t(X) = \text{cnt}_t(X_1) + \text{cnt}_t(X_2)$. If $\text{cnt}_t(X_1) < |X_1| \cdot \alpha$ or $\text{cnt}_t(X_2) < |X_2| \cdot \alpha$, then $\text{cnt}_t(X_1) + \text{cnt}_t(X_2) < (|X_1| + |X_2|) \cdot \alpha = |X| \cdot \alpha$, a contradiction. Therefore, $\text{cnt}_t(X_1) = |X_1| \cdot \alpha$ and $\text{cnt}_t(X_2) = |X_2| \cdot \alpha$. But then **Invariant 3** of **Claim A.1** would imply that both X_1 and X_2 contain a node requested at time t , which is a contradiction as they are disjoint.

Therefore, if X is a positive changeset applied at t , then X is a single tree cap of a tree from subforest C_{t+1} , and likewise if X is negative, then X is a single tree cap of a tree from subforest C_t . \square

B MINIMIZING FORWARDING TABLES USING TREE CACHING

In this section, we present a formal argument showing why we can use any q -competitive online algorithm A_T for the tree caching problem to obtain a $2q$ -competitive online algorithm A that minimizes forwarding tables.

Namely, we take any input I for the latter problem and create, in online fashion, an input I_T for the tree caching problem in a way described in **Section 2**. For any solution for I_T , we may replay its actions (fetches and evictions) on I and vice versa. However, there is one place, where these solutions may have different costs. Recall that an update of a rule stored at node v in I is mapped to a *chunk* of α negative requests to v in I_T . It is then possible that an algorithm for I_T modifies the cache *during* a chunk. An algorithm that never performs such an action is called *canonical*.

To alleviate this issue, we first note that any algorithm B for I_T can be transformed into a canonical solution B' by postponing all cache modifications that occur during some chunk to the time right after it. Such a transformation may increase the cost of a solution on a chunk at most by α and such an increase occurs only when B modifies a cache within this chunk. Hence, the additional cost of transformation can be mapped to the already existing cost of B , and thus the cost of B' is at most by a factor of 2 larger than that of B .

Furthermore, note that there is a natural cost-preserving bijection between solutions to I and canonical solutions to I_T (solutions perform same cache modifications). Hence, the algorithm A for I runs A_T on I_T , transforms it in an online manner into the canonical solution $A'_T(I_T)$, and replays its cache modification on I . Then, $A(I) = A'_T(I_T) \leq 2 \cdot A_T(I_T) \leq 2q \cdot \text{OPT}(I_T) \leq 2q \cdot \text{OPT}(I)$.

The second inequality follows immediately by the q -competitiveness of A_T . The third inequality follows by replaying cache modifications as well, but this time we take solution $\text{OPT}(I)$ and replay its actions on I_T , creating a canonical (not necessarily optimal) solution of the same cost.

C LOWER BOUND ON THE COMPETITIVE RATIO

THEOREM C.1. *For any $\alpha \geq 1$, the competitive ratio of any deterministic online algorithm for the online tree caching problem is at least $\Omega(k_{\text{ONL}}/(k_{\text{ONL}} - k_{\text{OPT}} + 1))$*

PROOF. We will assume that in the tree caching problem, evictions are free (this changes the cost by at most by a factor of two). We consider a tree whose leaves correspond to the set of all pages in the paging problem. The rest of the tree will be irrelevant.

For any input sequence I for the paging problem, we may create a sequence I_T for tree caching, where a request to a page is replaced by α requests to the corresponding leaf. Now, we claim that any solution A for I of cost c can be transformed, in online manner, into a solution A_T for I_T of cost $\Theta(\alpha \cdot c)$ and vice versa.

If upon a request r , an algorithm A fetches r to the cache and evicts some pages, then A_T bypasses α corresponding requests to leaf r , fetches r afterwards and evicts the corresponding leaves, paying $O(\alpha)$ times the cost of A . By doing it iteratively, A_T ensures that its cache is equivalent to that of A . In particular, a request free for A is also free for A_T .

Now take any algorithm A_T for I_T . It can be transformed to the algorithm A'_T that (i) keeps only leaves of the tree in the cache and (ii) performs actions only at times that are multiples of α (losing at most a constant factor in comparison to A_T). Then, fix any chunk of α requests to some leaf r' immediately followed by some fetches and evictions of A'_T leaves. Upon seeing the corresponding request r' in I , the algorithm A performs fetches and evictions on the corresponding pages. In effect, the cost of A is $O(1/\alpha)$ times the cost of A_T .

The bidirectional reduction described above preserves competitive ratios up to a constant factor. Hence, applying the adversarial strategy for the paging problem that enforces the competitive ratio $R = k_{\text{ONL}}/(k_{\text{ONL}} - k_{\text{OPT}} + 1)$ [30] immediately implies the lower bound of $\Omega(R)$ on the competitive ratio for the tree caching problem. \square

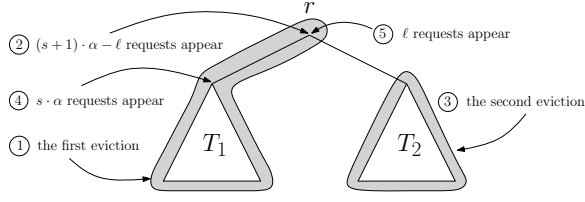


Figure 4: A troublesome example of a positive field. Numbers in circles describe the chronology of the events.

D IMPOSSIBILITY OF EXACT SHIFTING WITHIN POSITIVE FIELDS

In this section, we present an example showing that, within a positive field, we cannot shift positive requests down, obtaining α requests in every node, like we did in the case of negative requests (cf. [Corollary 5.8](#)). In our construction, the tree T consists of root r and two distinct subtrees T_1 and T_2 , each of size s and containing ℓ leaves.

Suppose that, at the beginning, TC has the entire tree T in its cache and the following ordered events happen (cf. [Figure 4](#)).

- (1) TC evicts $T_1 \cup \{r\}$ from the cache.
- (2) $(s+1) \cdot \alpha - \ell$ requests appear one by one at r . The number of requests is too small to trigger a fetch of any subtree of $T_1 \cup \{r\}$.
- (3) TC evicts T_2 from the cache.
- (4) $s \cdot \alpha$ requests appear one by one at the root of T_1 . This time, the number of requests is too small to trigger a fetch of any subtree of T .
- (5) ℓ requests appear one by one at r . After the last one appears, TC fetches the entire T to the cache.

The evictions happen because of some feasible sequence of negative requests that is irrelevant from our perspective.

Now, observe that when requests appear at the root in the second stage of our construction, T_2 is still in the cache (i.e., does not belong to the field yet). Thus, all the requests, except for the last ℓ ones can be shifted down only to nodes from T_1 . Hence, for large α and s , shifting can deliver $\Omega(\alpha)$ requests only to half of the nodes.