



## Pragmatic router FIB caching

Kaustubh Gadkari\*, M. Lawrence Weikum, Dan Massey, Christos Papadopoulos

Department of Computer Science, Colorado State University, Fort Collins, CO, United States

### ARTICLE INFO

#### Article history:

Received 19 September 2015

Revised 7 February 2016

Accepted 13 February 2016

Available online xxx

#### Keywords:

Computer networks  
Network architecture

### ABSTRACT

Several recent studies have shown that router FIB caching offers excellent hit rates with cache sizes that are an order of magnitude smaller than the original forwarding table. However, hit rate alone is not sufficient – other performance metrics such as memory accesses, robustness to cache attacks, queuing delays from cache misses, etc., should be considered before declaring FIB caching viable.

In this paper we tackle several pragmatic questions about FIB caching. We characterize cache performance in terms of memory accesses and delay due to cache misses. We study cache robustness to pollution attacks and show that in order to evict the most popular prefixes an attacker must sustain packet rates higher than the link capacity. We show that caching was robust even during a recent flare of NTP attacks. We carry out a longitudinal study of cache hit rates over four years and show the hit rate is unchanged over that duration. We characterize cache misses to determine which services are impacted by FIB caching. We conclude that FIB caching is viable by several metrics, not just impressive hit rates.

© 2016 Published by Elsevier B.V.

### 1. Introduction

The growth of forwarding table (FIB) sizes, fueled by factors such as multi-homing, traffic engineering, deaggregation and the adoption of IPv6 has led to a renewed interest in FIB caching methods. Past work (including ours) has shown repeatedly that there is significant traffic locality in the Internet that makes FIB caching beneficial. However, FIB caching has not been implemented to practice. Part of the reason might be that past work has focused on demonstrating that FIB caching is beneficial, leaving several practical and engineering questions unanswered. These include questions like how should the cache be implemented in a modern line card? Who suffers most from cache misses and how? How long does it take for a miss to be serviced? What are the memory bandwidth requirements for cache updates? How easily can one attack the cache? In this paper we address such practical questions and show that FIB caching is not only highly beneficial, but also very practical. We hope that our work answers important engineering questions and leads to renewed interest in building routers with caches.

Is FIB caching still relevant? Can't Cisco already support a million entries in their line cards? Opinions range from "it's not an issue" to "the sky is falling". We do not attempt to take a position in

this debate but seek only to inform. There are recent trends, however, that make the matter worth revisiting. One is the slow but steady growth of IPv6, which steadily adds more prefixes in the FIB. Another is the quest to build a Tb/s forwarding chip, which, for packaging reasons, will have limited on-chip memory, a perfect candidate for a cache.

In summary, we make the following contributions:

- We classify packets that cause cache misses according to their type and protocol.
- We evaluate the effect of cache misses on delay and buffer utilization.
- We evaluate the effect of caching on memory bandwidth requirements.
- We examine the behavior of the system under a cache pollution attack by someone who wants to replace popular with unpopular prefixes.

To achieve fast cache updates we propose an architecture that includes a *cacheable FIB* i.e. a FIB that does not suffer from the cache hiding problem (explained later). The cacheable FIB is derived from the standard FIB.

Briefly, our results are as follows: first, we confirm past observations that FIB caching is effective: with a cache size of 10K entries hit rates are in excess of 99.5%. Second, our classification of cache misses shows that NTP and DNS queries are the main culprits of cache misses. Not surprisingly, TCP control packets (SYNs, SYNACKs, FINs, FINACKs and RSTs) account for 70% of the TCP misses. Data packets cause approximately only a third of the

\* Corresponding author. +1 9702611699.

E-mail addresses: [kaustubh@cs.colostate.edu](mailto:kaustubh@cs.colostate.edu) (K. Gadkari), [lawrencemq@gmail.com](mailto:lawrencemq@gmail.com) (M.L. Weikum), [massey@cs.colostate.edu](mailto:massey@cs.colostate.edu) (D. Massey), [christos@cs.colostate.edu](mailto:christos@cs.colostate.edu) (C. Papadopoulos).

<http://dx.doi.org/10.1016/j.comcom.2016.02.006>

0140-3664/© 2016 Published by Elsevier B.V.

misses. Third, we show that very few packets need to be queued due to cache misses and they suffer insignificant queuing delays. Finally, we show that a cache recovers quickly when subjected to cache pollution attacks aiming to replace cache entries with unpopular prefixes. In our datasets, an attacker must send 1.2 billion packets/s over a 10G link to effectively disrupt the cache.

Our data comes from a regional ISP and thus our observations are mainly from the edge of the network. However, our study can easily be repeated for the network core by someone with access to the appropriate packet traces.

The rest of the paper is organized as follows. Section 2 introduces previous work that has looked into FIB scaling methods. In Section 3 we introduce our FIB caching solution. In Section 4 we introduce the cache hiding problem. Next, we introduce our hole filling algorithm and evaluate it in Section 4.1. In Section 5 we describe the datasets used in our evaluation. We evaluate LRU caches in Section 6 and the effect of cache misses in Section 7. In Section 8 we evaluate the robustness of our caching system when it is being attacked. Finally, we conclude in Section 9.

## 2. Related work

Approaches to reduce the FIB size have been studied for more than a decade [3,6,8,10,13,17,18,23,25,29]. These approaches fall into two broad categories – caching-based and aggregation-based approaches. Our work is independent of FIB aggregation techniques and hence we do not discuss these approaches in this work. There are other approaches that achieve FIB size reduction by reducing the RIB size [7,20,28]. We believe that our work is complementary to this work and the reduction of the RIB size will result in a more dramatic decrease in the cache size.

The idea of using route caching was first proposed by Feldmeier in 1988 [8], which was further extended by Kim et al. in [13]. Kim et al. introduce the cache hiding problem (discussed further in Section 4). To solve the cache hiding problem, the approach in [13] splits the IPv4 address space into the constituent /24 prefixes. Other approaches to handling the cache hiding problem include treating related prefixes as an atomic block so that all cache operations (insertion, lookup and deletion) are carried out on the block as a whole, using a complicated data structure with on-the-fly computation to remove interdependencies between prefixes [15] or using genetic algorithms to allow the cache policy to evolve while tracking the heavy hitter flows [27]. In this paper, we propose a *hole filling* algorithm to address the cache hiding problem, similar to the method proposed in [16]. While the algorithm presented in [16] generates the needed most-prefix on a per-packet basis, thus incurring a per packet cost, our algorithm pre-computes the prefix to add. By adding these extra prefixes, we ensure that there are no overlapping prefixes in the FIB. Consequently, a packet hitting a prefix not in the cache will cause the correct route to be fetched from the full FIB, instead of an incorrect longest prefix match with a covering prefix already in the cache. We discuss this algorithm in Section 4.1. We show that the number of extra prefixes added to the FIB is only a small fraction of the total number of FIB entries. Further, by using this *cacheable FIB*, we show that we can forward 99% or more packets along the fast path with a cache size of the order of 10,000 entries.

In [24], the authors propose a traffic offloading strategy to leverage the Zipf-like properties of Internet traffic. The proposed system, Traffic-aware Flow Offloading (TFO) is a heavy hitter selection strategy that leverages the Zipf-like properties and the stability of the heavy hitters across various time scales. TFO maintains the set of current heavy hitter flows in fast memory and sends packets from flows not in the heavy hitter set to a slower path. Results show that TFO sends a few thousand packets/second to the slow path.

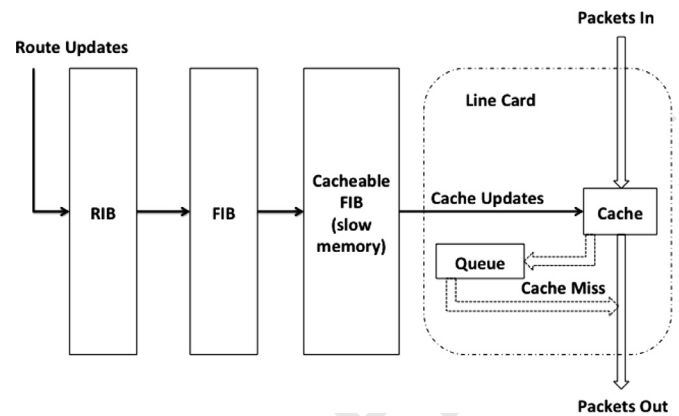


Fig. 1. Proposed caching system architecture.

Previous work [9,19] has investigated the impact of various cache poisoning attacks on caches and has proposed methods of mitigating such attacks. However, the systems investigated were software systems (web proxies) that have a different set of constraints than the hardware-based route caching system proposed in this work. Further, unlike previous work, we consider the impact of cache misses on the system performance not only in terms of hit rates achieved but also the types of requests resulting in cache misses and their implications for operators.

## 3. Cache system design

Fig. 1 shows our proposed FIB caching architecture. The RIB and FIB on the left are part of the standard hardware architecture. We add (a) a *cacheable FIB*, which resides in main memory, (b) a cache, which can take the place previously occupied by the FIB, (c) a queue to hold packets experiencing a cache miss, and (d) the appropriate logic to handle cache updates and misses. The router processor handles incoming route information as usual and computes the local RIB (routing table) and FIB as before. In a traditional architecture the router would load the entire FIB in the line card; in our architecture, the router derives the cacheable FIB from the standard FIB and pushes just the cache in the line card, which is one or two orders of magnitude smaller than the original FIB.

Our caching architecture may seem to keep two copies of the FIB, the original and cacheable FIB, which may seem a waste of memory. However, the two FIBs are very similar as we will show later and an implementor can easily use appropriate data structures to avoid duplication.

Each incoming packet incurs a lookup in the cache using the standard longest prefix match algorithm. If there is a hit the packet is immediately forwarded along the cached route. If there is a miss the packet is queued in the line card until the cache is updated with the appropriate prefix from the cacheable FIB.

Next, we elaborate on the need and derivation of a cacheable FIB.

## 4. The need for a cacheable FIB

Entries in the original FIB cannot be cached due to the *cache hiding problem*, which occurs when a less specific prefix in the cache hides a more specific prefix in the FIB. This is a result of the fact that a route-caching system looks for the longest-matching prefix only in the cache, thus missing possible longer matches in the full FIB. This can lead to incorrect forwarding, loops and packet losses.

To further understand the problem, consider an empty cache and a full routing table that only has two prefixes: 10.13.0.0/16

**Table 1**

Original non-cacheable FIB, after one iteration of hole filling algorithm and final cacheable FIB.

Prefix	OFF
(a) Non-cacheable FIB	
10.13.0.0/16	1
10.13.14.0/24	2
(b) Non-cacheable FIB after one iteration of hole filling algorithm	
10.13.0.0/17	1
10.13.128.0/17	1
10.13.14.0/24	2
(c) Cacheable FIB	
10.13.128.0/17	1
10.13.64.0/18	1
10.13.32.0/19	1
10.13.16.0/20	1
10.13.0.0/21	1
10.13.8.0/22	1
10.13.12.0/23	1
10.13.14.0/24	2
10.13.15.0/24	1

158 associated with an output interface O1, and 10.13.14.0/24 associ-  
 159 ated with an output interface O2. Suppose a packet arrives with a  
 160 destination IP address of 10.13.2.3. The router will find the longest-  
 161 matching prefix (LMP) for the destination IP address, which is  
 162 10.13.0.0/16, and will install the route [10.13.0.0/16 → O1] in the  
 163 cache. Assume that the next packet that arrives at the router has  
 164 a destination IP address of 10.13.14.5. The router will find the pre-  
 165 viously installed route [10.13.0.0/16 → O1] in the cache and will  
 166 forward the packet along O1. However, this is incorrect, since the  
 167 correct LMP for 10.13.14.5 is 10.13.14.0/24 and the packet should  
 168 have been forwarded on interface O2.

169 Past work addressed the cache hiding problem by either  
 170 caching only /24 prefixes [13], by using a complex data structure  
 171 with on-the-fly computation to eliminate inter-dependencies be-  
 172 tween prefixes [15], by adding on-the-fly to the FIB the appropriate  
 173 leaf node corresponding to the packet's destination address [16] or  
 174 by treating a set of related prefixes as an atomic block and per-  
 175 forming cache actions (insertion, lookup and delete) on the atomic  
 176 block instead of an individual prefix. Our approach is similar to the  
 177 approach presented in [16], except that we pre-calculate the entire  
 178 cacheable FIB table.

#### 179 4.1. Generating a cacheable FIB

180 We have previously seen that simply caching existing FIB en-  
 181 tries would lead to incorrect forwarding due to the cache hiding  
 182 problem. In this section we describe an algorithm to generate a  
 183 cacheable FIB, i.e. a FIB that is free from the cache hiding problem.

184 We call it the *hole filling* algorithm because it starts with the  
 185 original FIB and fills in holes between related entries to produce a  
 186 new FIB whose entries are cacheable. While the new FIB is larger  
 187 because the algorithm adds more entries, the increase is small, as  
 188 we will show later, and caching makes it irrelevant.

189 The intuition behind the algorithm is as follows – if a prefix  
 190 covers other prefixes in the table, we delete that prefix and add  
 191 its children to the FIB. We repeat this process until there are no  
 192 covering prefixes, at which point we are left with a cacheable FIB.

193 To better understand the algorithm, consider the FIB shown in  
 194 Table 1a. This FIB is non-cacheable, since the 10.13.0.0/16 prefix,  
 195 if present in the cache, will hide the 10.13.14.0/24 prefix. This FIB  
 196 needs to be transformed into a cacheable FIB.

197 We first choose the prefix from the FIB with the smallest mask  
 198 length, which is 10.13.0.0/16 in this case. Then we check if this  
 199 prefix has any children. If the selected prefix has a child, we split

**Table 2**

FIB size increase due to hole filling algorithm. The in-  
 crease in FIB size due to hole filling is minimal.

Table	Original	Cacheable	Increase
Regional ISP	441,778	468,095	5.96%
Hurricane	444,977	470,911	8.51%
Telstra	440,156	466,416	8.43%
Level-3	436,670	462,966	8.41%
AOL	438,719	464,988	8.40%
NTT-A	439,078	465,536	8.38%
ATT	438,265	464,662	8.37%
SAVIS	438,582	465,045	8.37%
Sprint	438,634	465,079	8.37%
VZWBIZX	436,821	463,220	8.35%
SWISS-COM	169,861	173,636	8.27%
KPNE	439,288	465,790	6.03%
Tiscali	438,993	465,343	6.00%
IJ	440,196	466,505	5.98%

**Table 3**

Trace statistics. We use trace T8 in our caching perfor-  
 mance analysis and trace T9 in our cache robustness  
 analysis.

No.	Date	Time	Link	No. packets
T1	3/31/09	27H	1G	7,620,972,889
T2	8/17/09	24H	1G	3,821,714,756
T3	8/3/10	24H	1G	2,084,398,007
T4	8/3/10	24H	1G	2,050,990,835
T5	12/14/11	12H	1G	625,547,727
T6	04/13/12	12H	1G	3,729,282,487
T7	2/14/13	12H	10G	22,622,946,036
T8	3/20/13	12H	10G	21,998,786,996
T9	2/21/14	12H	10G	18,435,172,172

the prefix into its two children, otherwise we add the prefix to  
 the cacheable FIB. In this example, since 10.13.0.0/16 has a child,  
 10.13.14.0/24, we split the /16 into its two constituent /17s and re-  
 move the original 10.13.0.0/16, as shown in Table 1b.

The process continues, until the non-cacheable FIB is empty,  
 and the cacheable FIB contains the leaf nodes necessary to forward  
 all packets correctly. Table 1c shows the final, cacheable FIB.

#### 4.1.1. FIB inflation due to hole filling

The hole filling algorithm produces a cacheable FIB that is  
 larger than the original. For example, the original FIB in Table 1a  
 has only two entries, while the cacheable FIB shown in Table 1c  
 has nine entries, an increase of 350%. We next investigate the ef-  
 fect of the algorithm on real FIBs.

We measured inflation on a FIB from our regional ISP, which  
 contains next hop information, and on FIBs from RouteViews [22]  
 that do not contain next hop information, which we approximate  
 using the next hop AS. Table 2 shows at worst inflation is under  
 9%. Since the cacheable FIB resides in main memory the impact is  
 very small.

## 5. Data sets and trace statistics

We used nine 12H and 24H packet traces taken at links be-  
 tween our regional ISP and one of its tier-1 providers. We captured  
 traffic using specialized hardware [1] that ensured no packets were  
 dropped during capture and packet timestamps were accurate to  
 a nanosecond resolution. We then used a trie-based longest pre-  
 fix matching algorithm to determine the matching prefix for each  
 packet. The cacheable FIB was derived from the actual FIB obtained  
 from the router where the packet trace was captured. Table 3  
 shows the trace statistics.

In an interesting twist, there was a sustained NTP reflection  
 attack in trace T9. In this attack, several comprised machines

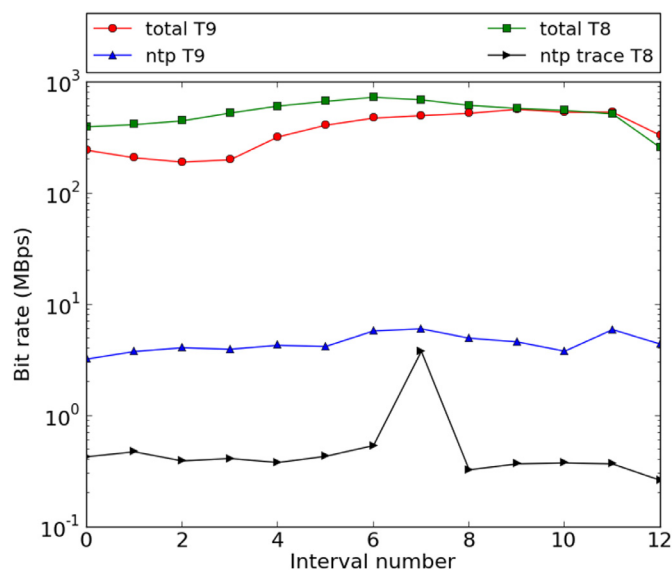


Fig. 2. Comparison of average overall bit rate and ntp bit rates per hour between a “normal” trace and a trace with attack traffic. Trace T9 was captured during an NTP attack.

belonging to our ISP’s customers were used as amplifiers to attack other hosts on the Internet using a NTPD vulnerability [2]. Fig. 2 shows the overall average bit rate per hour as well as the average NTP bit rate per hour. NTP traffic accounted for 1.3% of all traffic during the duration of the trace. In comparison, trace T8 has approximately the same byte rates as trace T9. However, NTP accounted for only 0.1% of all traffic in trace T8.

For brevity, we show statistics from trace T8 only. Results from the other traces (except T9) are similar. We use trace T9 in our cache robustness analysis in Section 8.3.

## 6. Results

In this section we present results using FIB caching emulation, using real traces to drive the emulator.

### 6.1. Cache system performance

We begin by evaluating cache performance using two standard cache replacement policies – least recently used (LRU) and least frequently used (LFU). We observed that LFU performs worse than

LRU. Qualitatively our results confirm previous observations, but it is important to show them to establish a baseline. Fig. 3a and b shows the performance of an LRU cache with varying cache sizes for trace T8 from Table 3. We plot average cache hit rates at every 5-min interval.

From Fig. 3a and b we see that both LRU and LFU consistently achieves very high hit rates. With a cache size of only 1K entries the hit rate is 98%. The hit rate increases with the cache size, reaching almost 99.9% with a cache of 10K entries. LRU achieves maximum hit rates of 99.26%, 99.74%, 99.91% and 99.96% with cache sizes of 1K, 2.5K, 5K and 10K respectively. Note that even with a cold cache the hit rate penalty is very small. For the rest of the paper, unless otherwise noted, we will use a cache size of 10K. The reason is that even with 10K entries, this is a still a very small cache compared with the original FIB size (currently around 500K). In the rest of the paper, we present results only for a LRU cache unless otherwise noted, since it clearly offers higher hit rates than LFU.

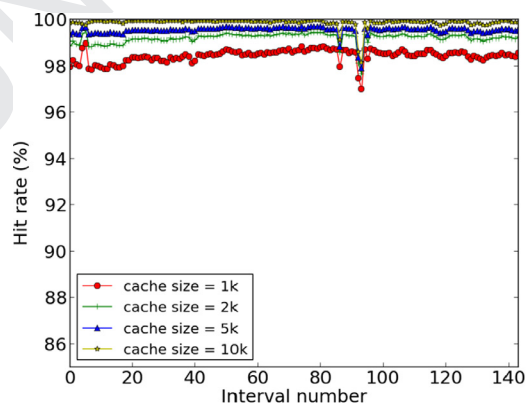
### 6.2. Caching with delayed updates

In this section, we examine the effects of an optimization of the caching strategy. Namely, instead of inserting a prefix  $P$  immediately after encountering a cache miss, the system waits until a certain pre-set threshold is crossed i.e. prefix  $P$  receives at least  $N$  packets in a given time interval  $T$ . This optimization targets the “bottom” of the cache i.e. transient prefixes that only receive a few hits – these mostly unpopular prefixes will be evicted from the cache thus resulting in cache churn. By not inserting these prefixes into the cache in the first place, we potentially save memory lookups and cache updates. Fig. 4 shows the hit rates achieved with the delayed cache updates.

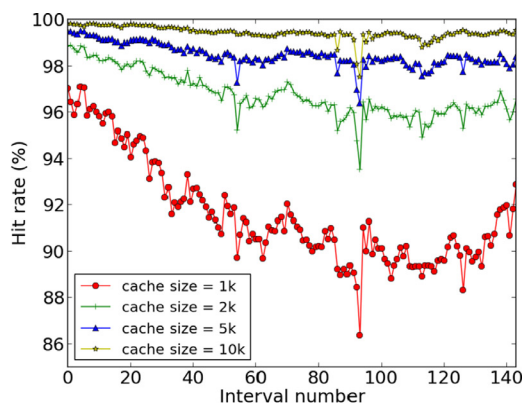
Fig. 5 a and b shows the number of packets that are forwarded along the slow path and the fraction of total packets delivered along the slow path, respectively, during each of the intervals. As expected from Fig. 4, the number of packets sent along the slow path varies inversely with the timeout value. Figs. 4 and 5 show that the cache need not be updated every time a cache miss occurs. Instead, the cache can be updated in intervals with the interim misses being forwarded along the slow path without adversely affecting the cache performance.

### 6.3. Impact of route updates

Routers periodically receive route updates from their peers. These updates may add new prefixes (announcements) to the table



(a) LRU Hit Rates



(b) LFU Hit Rates

Fig. 3. LRU and LFU hit rates for cache sizes varying from 1K to 10K.

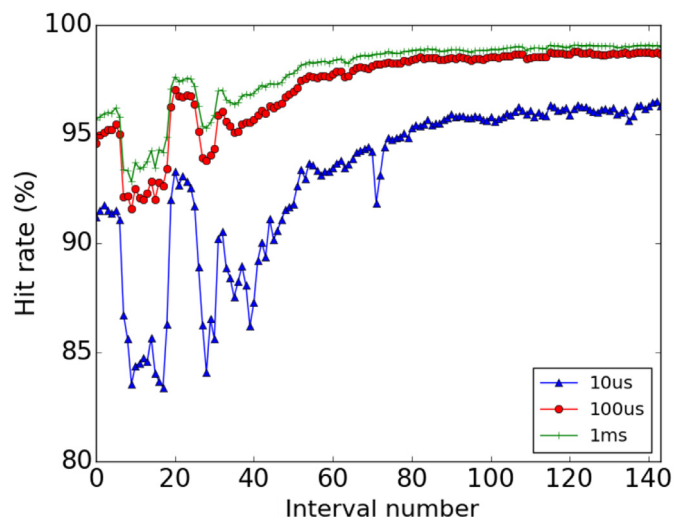


Fig. 4. Hit rates for LRU cache with delayed updates. The cache size was 10K entries.

Table 4 Update statistics.

Update	Entries	Announce	Withdraw
U1	55,718	43,720	6,263
U2	65,120	45,551	6,375
U3	73,599	46,976	4,606
Total	194,437	136,247	17,244

or withdraw existing prefixes from the table (withdrawals). Each such update may cause one or more entries in the cache to be invalidated.

To evaluate the effect of route updates we took a full RIB table from a RouteViews [22] peer and generated a cacheable FIB. We again approximate next hop information using the next hop AS from the ASPATH attribute. Then, we applied updates for the same peer to the cacheable FIB and generated a new cacheable FIB. Finally, we measured the difference between the original cacheable FIB and the newly generated cacheable FIB.

We show results for 3 sets of updates. Table 4 shows some statistics about the updates applied to the cache.

The size of the cacheable FIB generated from the original FIB went up from 458,494 to 464,512, an increase of only 1.29%. 82.5% of the prefixes in the cacheable FIB were the same as those in the

Table 5 Expansion statistics.

Data	Updated	Cacheable	Growth	Same
FIB	458,494	465,512	1.29	82.5
U1	458,725	458,763	0.008	99.97
U2	458,715	458,755	0.009	99.97
U3	458,890	458,974	0.018	99.97

original FIB (Table 5). After subsequent updates were applied, the size increase as well as the number of prefixes that changed was negligible – the average change in the size of the cacheable FIB was only 0.0018%.

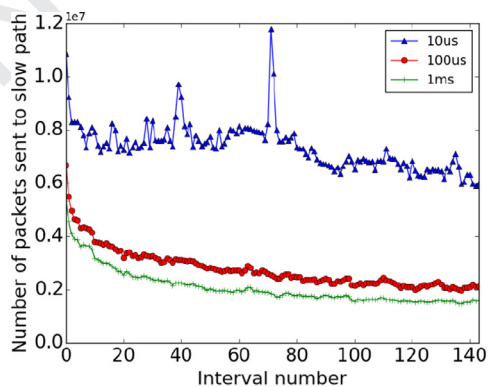
Next, we count the number of prefixes that had next hop changes, since only these prefixes will have to be invalidated if present in the cache. Our results show that, on average, only 144 prefixes in each 15 min set of updates had a next hop change. The insignificant change in the size of the cacheable FIB after applying updates, coupled with the fact that only a few hundred prefixes will have to be invalidated from the cache due to change in forwarding state, suggest that routing updates will have very little impact on the cache state and forwarding performance.

It should be noted that while this study is limited in scope to tables and updates from only one RouteViews peer, we believe that a more comprehensive study will yield similar results. Research shows that routing tables from different RouteViews peers have a very few differences [26] and hence we believe that our results are applicable to other peers as well.

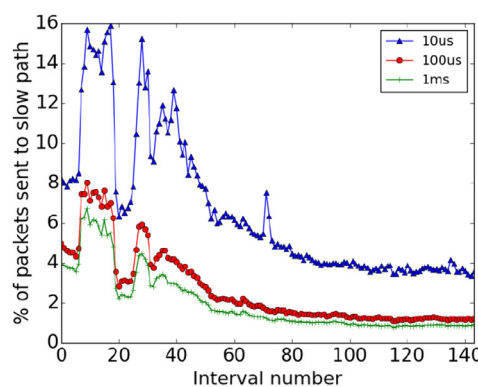
6.4. Trends in cache size

The global routing table has been growing steadily for the past several years. Measurements show that the routing table grew from 150,000 entries in 2003 to 517,802 entries at the time of this writing, representing an increase of 245% in the last decade alone [11]. Routers thus need increasing amounts of memory and processing power [12]. The conventional wisdom that FIB memory will scale at rates surpassing the rate of growth in the routing information handled by core router hardware is not true for the low-volume, customized hardware used in core routing [21].

Given the fast rate of increase of the global routing table over time, the natural question to ask is does a FIB cache face a similar rate of growth? To get some insight into the answer we measure the hit rates achieved with a cache size of 10K entries on packet traces gathered from 2009 to 2014. These traces, except trace T9, were collected at the same point in our ISP, so they truly capture



(a) Number of packets forwarded along slow path



(b) Percentage of packets forwarded along slow path

Fig. 5. Analysis of packets sent to the slow path.

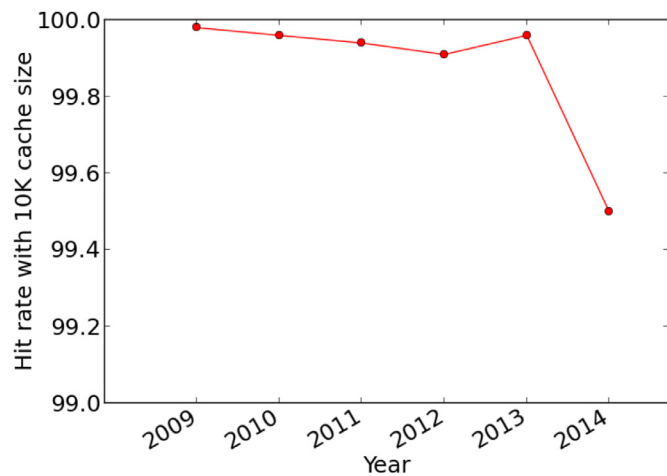


Fig. 6. Hit rates achieved by a 10K entry cache remain almost constant from 2009 to 2014. The 2014 trace contains the NTP attack traffic.

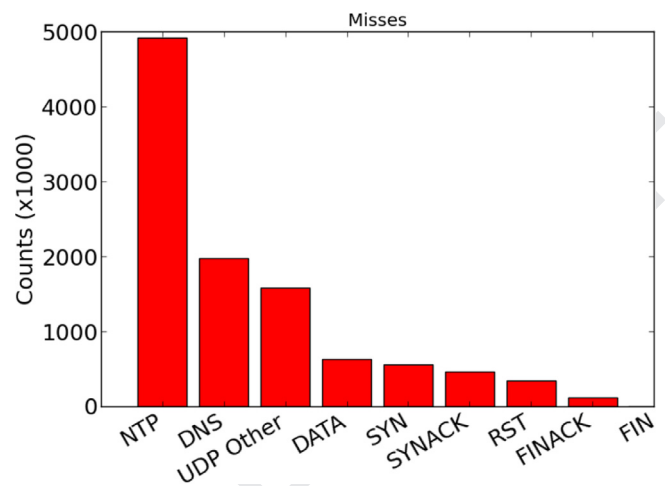


Fig. 7. Classification of packets causing cache misses.

the cache evolution over the last four years. Trace T9 was captured at another monitoring location in our ISP, but results show that the hit rates for trace T9 are similar to the hit rates for the other traces.

Fig. 6 shows the results. While there is some minor variation, the hit rates achieved with a cache size of 10K are consistently higher than 99.95%, meaning that over the past four-year period the cache size did not need to change to maintain the same cache hit rate. Thus, while the global routing table has certainly changed, traffic locality has not changed sufficiently to affect the cache size.

## 7. Analysis of cache misses

In our FIB caching system when a packet arrives and the appropriate prefix is not in the cache, the packet is queued while a cache update takes place. This introduces some delay before the packet is forwarded. Delay may affect packets in different ways. For example, delaying video packets may result in dropped frames and choppy video. On the other hand, queuing DNS packets may result in only a slightly longer resolution delay. In this section we characterize the type of packets queued due to misses and the delay they experience.

First, we classify the types of packets that cause cache misses to determine which applications are affected. Second, we determine buffer requirements and queuing delays due to cache misses. Finally, we analyze the impact of cache misses on memory bandwidth requirements in routers. We show results only for trace T8 from Table 3. Results for other traces are similar and have been omitted due to space constraints. Since trace T9 was captured during an ongoing NTP reflection attack, the cache misses in this trace were heavily influenced by the attack. We present the analysis for trace T9 later in Section 8.3.

### 7.1. Classification of cache misses

We classify packets causing cache misses as follows. First, we separate TCP and UDP packets. For TCP, we classified them as SYNs, SYNACKs, FINs, FINACKs, RSTs and DATA packets. Any TCP packet that was not a SYN, SYNACK, FIN, FINACK or RST is classified as a DATA packet. Preliminary analysis of UDP packets shows that NTP and DNS packets suffered most misses. We therefore plot NTP and DNS packets separately, with the remaining packets classified as “UDP Other”. Fig. 7 shows the classification of packets causing cache misses with an LRU cache.

We also see that the largest number of cache misses is caused by NTP packets. The reason is that while NTP requests are regular, they are also infrequent. Thus, the prefixes required to forward the NTP packets are regularly evicted from the cache, only to be re-inserted when the next NTP request occurs. In our dataset, NTP packets were destined to 69,097 unique prefixes.

After NTP, DNS packets are responsible for the largest number of cache misses. This occurs due to a DNS query which needs to be forwarded to the authoritative nameserver of the zone that owns the prefix. If the prefix is not in the cache, the DNS query will result in a cache miss. In our dataset, DNS packets hit 58,435 unique prefixes.

### 7.2. Effects of cache misses on queuing

There are several ways of dealing with packets causing cache misses. One possible strategy is to forward the packets along a slower path until the appropriate route can be inserted into the FIB/cache [4]. If this strategy is employed, there is no need for queues/buffers at the routers to store packets. However, the packets causing a cache miss have to travel a longer path, thus experiencing longer delay and possible reordering. Another strategy is to queue the packets at the router until the route is inserted into the cache. While the packets do not incur any longer paths (and hence stretch), line cards must now have enough buffer space and the appropriate logic to queue packets. In our analysis we assume this strategy for several reasons: (a) it prevents packet reordering, which router vendors try hard to avoid, (b) to the best of our knowledge it has not been evaluated before, and (c) as our results show, the buffer requirements are modest and the queues can easily fit in existing buffers.

#### 7.2.1. Packet queuing emulator

We built an emulator in order to measure queuing at the router during cache misses. Fig. 8 shows the block diagram of our emulator. We assume that it takes about 60 ns to update the cache (the lookup time for DRAM) and that a packet needs to be queued while the update is taking place.

When a packet arrives we do a longest prefix match against the cache. If there is a match the packet is forwarded and does not incur any additional delays. If there is no match against the cache the packet is queued. When it gets to the head of the queue, we do a prefix fetch from the cacheable FIB, update the cache and forward the packet.

We measure the maximum queue utilization during a given 5 min interval. We also look at the maximum queuing delay in

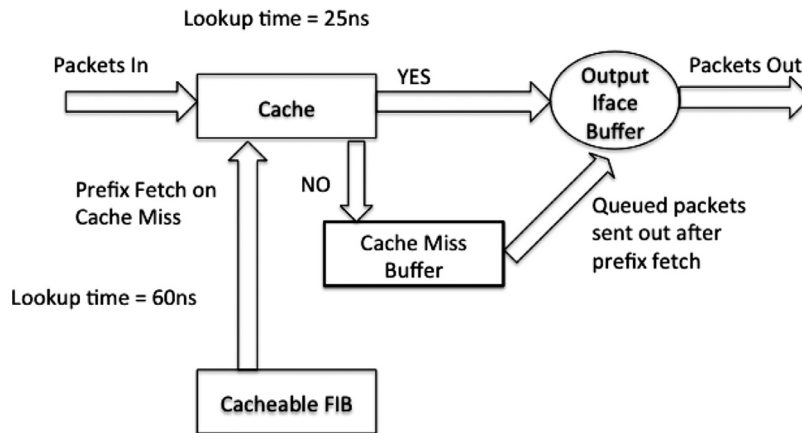


Fig. 8. Queue utilization simulator.

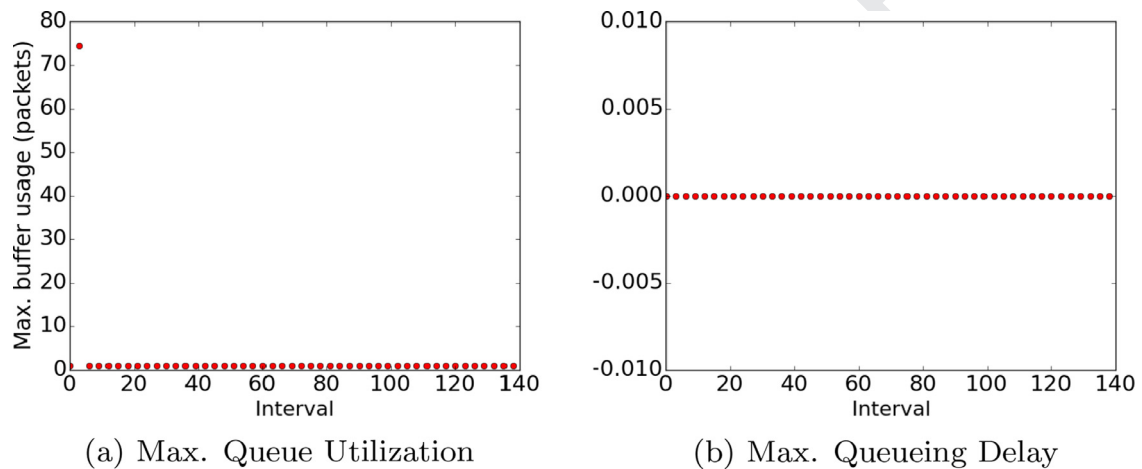


Fig. 9. Max. queue utilization and queuing delays.

our emulator i.e. the time elapsed between when a packet entered the queue and when it exited the queue. Let us assume that the time required to perform a cache lookup is  $T_C$ , the time required to perform a slow-memory lookup is  $T_S$  and the current buffer size is  $B$ . If a packet arrives at the queue at time  $T_A$  the time the packet departs the queue  $T_D$  is given by:

$$T_D = T_A + T_C + (B * T_S)$$

The queuing delay  $D$  is then

$$D = T_D - T_A$$

### 7.2.2. Evaluation

Fig. 9a shows the maximum queue utilization during 48 5-min intervals. In the first interval we see a startup effect due to the cold cache. Since there are no entries in the cache, arriving packets cause many cache misses and have to be queued, resulting in a maximum queue size of 73 packets. After this initial startup period the queue drains and queue utilization drops. The maximum queue utilization is only 1 packet, which means that assuming 1500 byte packets, we need only 1500 bytes of buffer space to buffer LRU cache misses. Even when the queue utilization peaks during the cache warm-up period, only 73 packets are queued. To store these packets, we will need only 110KB of buffer space.

Fig. 9b shows the maximum queuing delays for a LRU cache during the same intervals as above. Packets queued after a LRU cache miss suffer virtually no queuing delay. This is due to the small average queue size (1 packet per 5 min interval) and the

relatively fast cache updates (in the order of ns) which keep the backlog small.

While these numbers are specific to our dataset, we believe that they generalize well in similar environments. Moreover, the buffer requirements are in line with what is found in routers today, where the rule of thumb is that buffers should hold an RTT's worth of data.

### 7.3. Memory bandwidth requirements

Another important benefit of FIB caching is a reduction in the number of lookups that need to be performed on the full FIB, which is often kept in slow (DRAM) memory. Without caching one lookup typically needs to be performed per packet. Moreover, the next hop information is often only part of the lookup information, which may also include filtering, QoS, and other state, increasing the demand on memory bandwidth. With increasing FIB sizes and line speeds these lookups are nearing the bandwidth limit of router memory.

Caching has the potential to drastically reduce the number of lookups to external memory. For example, one may envision a system where the cache is kept on-chip and the full RIB remains on external slow DRAM. The latter needs to be accessed only when a cache miss occurs. The question then is, what are the savings in terms of memory bandwidth if one uses an on-chip cache?

Fig. 10 shows the average number of memory lookups required per second in each 5 min interval in the trace. Without caching, the number of memory lookups is equal to the packet rate, since

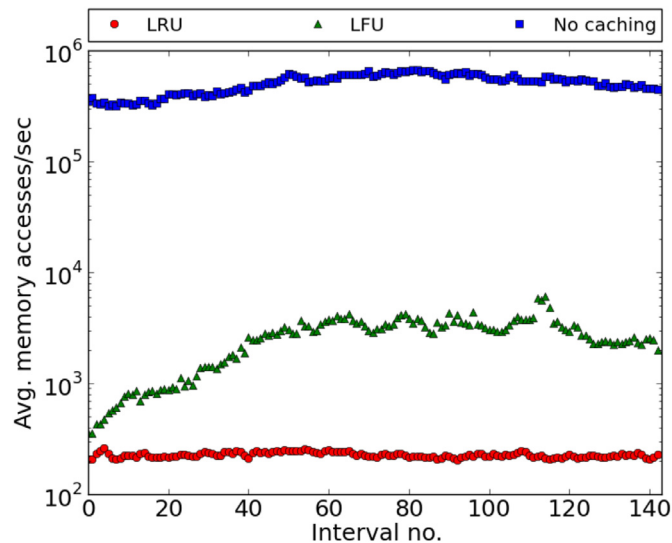


Fig. 10. Memory bandwidth requirements reduce by orders of magnitude when caching is employed.

each packet requires a memory lookup. With caching, a memory lookup is required only in case of a cache miss. We see that the number of memory accesses required with caching is several orders of magnitude lower than the number of accesses required when no caching is deployed. If the cache uses a LRU replacement policy the memory accesses are on the order of  $10^2$  accesses per second and the rate stays almost constant. This is to be expected, since we see that the hit rates achieved with LRU stay constant throughout the duration of the trace.

Caching offers an order of magnitude improvement in memory bandwidth when compared to no caching. Coupled with the fact that the required cache can easily fit on a chip and that the cache size appears to remain constant over time, caching virtually eliminates any memory bandwidth issues for current and potentially future routing table sizes. This is significant for scaling future routers.

## 8. Cache robustness

The use of FIB caching, especially with LRU, exposes routers to a cache pollution attack, where an attacker attempts to disrupt packet caching and increase the cache miss rate. An attacker can pollute the cache by continuously sending packets to numerous unpopular prefixes that would not be otherwise cached, in an attempt to evict legitimate prefixes. Depending on the attack rate, the attacker can either cause the cache to thrash, thus increasing the miss rate, or reduce the effectiveness of the cache by ensuring that at least part of it is always polluted with bogus entries. This attack will adversely affect the packet forwarding performance as well, by substantially increasing the number of packets that need to be queued while the new prefix is fetched, potentially leading to packet loss.

In this section we investigate cache pollution attacks and their feasibility. We describe a generalized threat model where the attacker tries to replace a subset of cache entries, and then estimate the attack rate required to adversely affect the cache. Next, we evaluate cache performance when subjected to a real NTP attack that happened to be present in trace T9.

### 8.1. Attacks against the cache

In this section, we describe an idealized attack against a LRU cache, in which an attacker inserts consecutive packets to at least

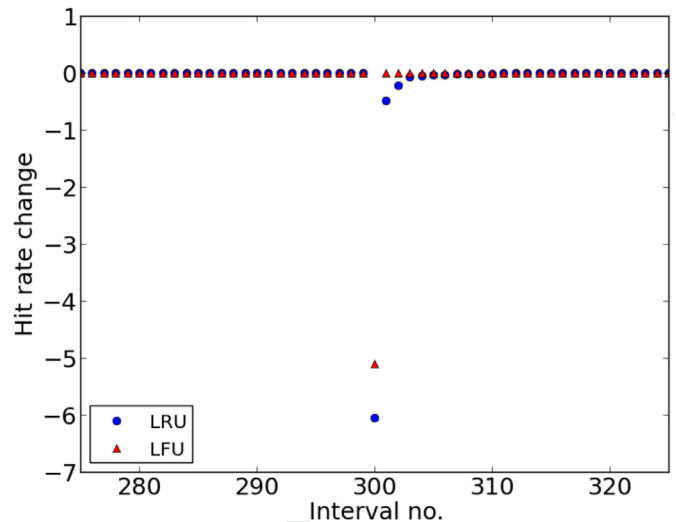


Fig. 11. Difference in cache hit rates with and without attack packets. Attack is in interval number 300.

$N$  unpopular prefixes. We assume a cache size of 10K entries and evaluate the performance of both a LRU cache.

We begin our evaluation by taking a 20 min slice from the beginning of trace T8. Then we insert a stream of 20,000 attack packets to 20,000 idle prefixes starting at 300 s into the trace. We used this attack trace as input for our caching system and measured the hit rate at 1 s intervals.

Fig. 11 shows the difference between the hit rates observed with and without the attack packets. A negative difference implies that the hit rate observed during the attack was lower than the hit rate in the same interval, but without the attack packets present. For clarity, we present data for 25 intervals before and after the attack interval.

As Fig. 11 shows, when the attack hits (at interval 300), the LRU cache hit rate reduces by almost 6%. Further, it takes 2 after the attack has stopped for the cache to recover. However, despite the attack, the cache hit rate never drops below 93%.

Next, we measure the queue size and delay under the attack. Fig. 12a shows the mean queue size with a LRU cache under attack. Normally, the queue size is close to zero. However, when the attack hits (at interval 300), the queue peaks at 9121 packets.

Fig. 12b shows the corresponding queuing delays suffered by packets. Normally, the packets suffer no queuing delay, since very few packets are queued under normal traffic conditions. However, under attack, we see that packets suffer queuing delays of about  $650\mu\text{s}$ .

It should be noted that this attack is an idealized attack because it assumes all attack packets arrive back-to-back. In reality it is hard for an attacker to inject such bursts into the normal packet stream unless traffic is sufficiently low.

### 8.2. Generalized threat model

In this section we describe a generalized threat model to attack the cache, and determine the rate at which an attacker must send packets to disrupt packet delivery. We assume that the attacker knows or can determine in advance the set of popular and unpopular prefixes at a given vantage point. We also assume that the attacker has the capability to send packets at high rate, up to the capacity of the link, either from a single host or a botnet. We also assume that the goal of the attacker is to replace legitimate cache entries with unpopular entries for the sole purpose of disrupting cache performance. In other words, we assume a deliberate



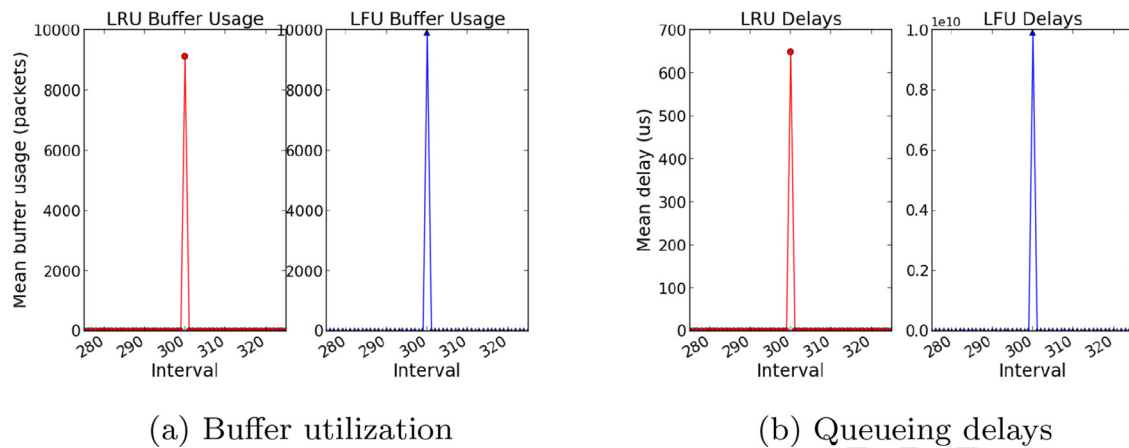


Fig. 12. Comparison of buffer utilization and queuing delays under normal and attack conditions. Attack is in interval 300.

551 attack on the cache and not a general attack on the infrastructure  
552 (e.g., DDoS). The latter will trigger other network defenses.

553 To achieve the above goal the attacker still needs to send pack-  
554 ets at a high enough rate to evict legitimate prefixes quickly from  
555 the cache. Thus, the attack cannot be stealthy since its packet rate  
556 must compete head-to-head with legitimate traffic. Next, we esti-  
557 mate the packet rate that would make such attack viable.

#### 558 8.2.1. Intensity of an effective attack

559 To determine the rate at which the attacker must send pack-  
560 ets to affect cache performance we first rank prefixes according  
561 with their popularity during a low- and high-traffic interval. Recall  
562 that our measurement interval is 5 min. We chose to investigate  
563 the per-interval behavior rather than average traffic over the en-  
564 tire 24H period in order to determine both high and low required  
565 attack rates.

566 We estimate the required attack intensity. If an attacker wants  
567 to hit the  $N$ th entry in the cache (and all the prefixes below it) it  
568 must send packets to enough prefixes to evict the  $i$ th entry and all  
569 the other entries below it. So for example, if the attacker wants to  
570 blow the entire cache, then the attacker must attack at  $P_{attack}$  rate,  
571 which must be greater than the cache size  $N$  multiplied by  $P_{top}$   
572 which is the packet rate of the most popular prefix. To generalize,  
573 the attack rate to evict the  $i$  bottom prefixes from the cache must  
574 be:

$$P_{attack} \geq P_i^* i$$

575 In the low traffic interval, the most popular prefix received  
576 33,049,971 packets in five minutes for an average packet rate of  
577 110,166 packets per second, whereas in the high traffic interval the  
578 most popular prefix received 37,079,737 packets for an average of  
579 123,599 packets per second. Thus, to replace just the most popular  
580 prefix from the cache, the attacker needs to send packets at a rate  
581 between 1,101,660,000 packets/s and 1,235,990,000 packets/s to an  
582 idle prefix. For a 10 Gb/s link and 64 byte packets the maximum  
583 packet rate possible is 19,531,250 packets/s, thus the required at-  
584 tack rate is not feasible as it far exceeds the capacity of the link.

585 Note that such an attack can be easily detected by looking for  
586 spikes in the cache miss rate. Once detected, one can imagine sev-  
587 eral defenses against this type of attack, such as pinning down at  
588 least part the historical prefix working set until the attack subsides.  
589 This poses little danger of false positives, since when prefixes sud-  
590 denly become popular (as in the case of a flash crowd), they are  
591 unlikely to do it in numbers in the order of the cache size. Thus,  
592 we believe that practical attacks of this kind can only affect part of  
593 the cache and are easily weakened by reasonably over-engineering  
594 the cache.

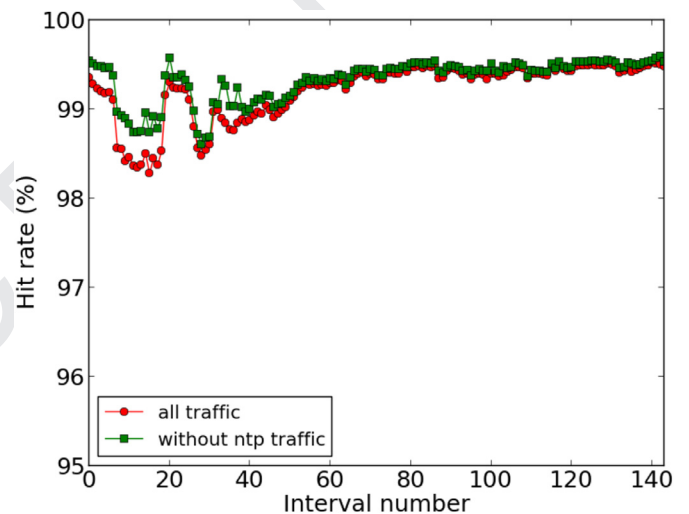


Fig. 13. Comparison of LRU hit rates for all traffic and non-NTP traffic. The cache size was set to 10K entries.

#### 595 8.3. Cache performance under a real DDoS attack

596 As described in Section 5, trace T9 was captured during an on-  
597 going NTP reflection attack. Fig. 2 shows that overall, NTP traffic  
598 accounted for approximately 1.3% of all traffic in the trace. 5659  
599 unique addresses from our regional ISP were the target of this at-  
600 tack [5]. Even though this was not an attack on the cache, we take  
601 advantage of this unfortunate incident to investigate the perfor-  
602 mance of a FIB cache under attack conditions and compare perfor-  
603 mance with a version of the same trace that has NTP traffic  
604 removed.

605 Fig. 13 shows the hit rates achieved by a 10K entry LRU cache  
606 for all traffic in trace T9 as well as the hit rates for all non-NTP  
607 traffic.

608 As Fig. 13 shows, the cache performs well under this incident.  
609 The difference in the hit rates achieved by the cache do not differ  
610 by more than 0.5% during the trace. Thus, caching remains robust.  
611 However, comparing Figs. 3a and 13, it is clear that the hit rates  
612 for trace T9 are lower than of trace T8, although the difference is  
613 less than 1%.

614 Next, we look at the breakdown of packets resulting in cache  
615 misses as shown in Fig. 14a.

616 As expected, the NTP traffic accounts for a majority of the  
617 packet misses. Out of a total of  $114 \cdot 10^6$  cache misses observed

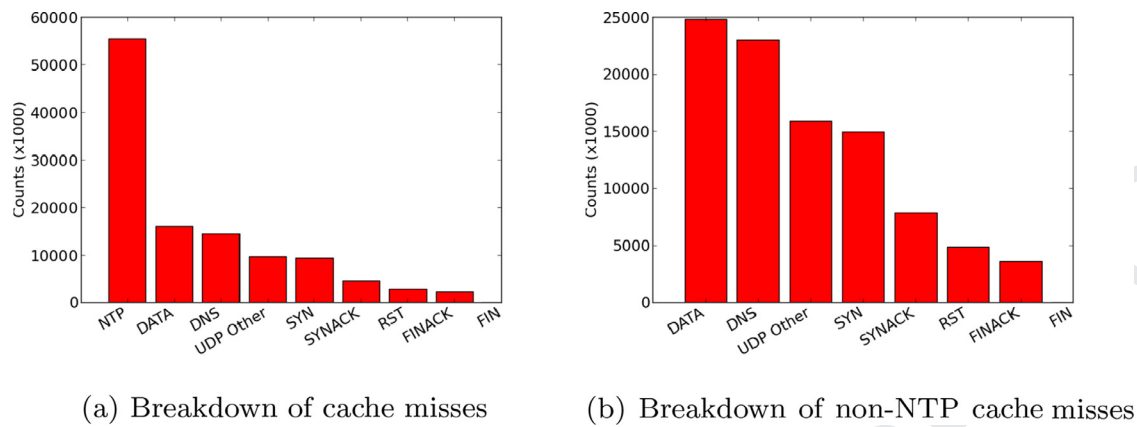


Fig. 14. Comparison of cache misses for all traffic and non-NTP traffic. The cache size was set to 10K entries.

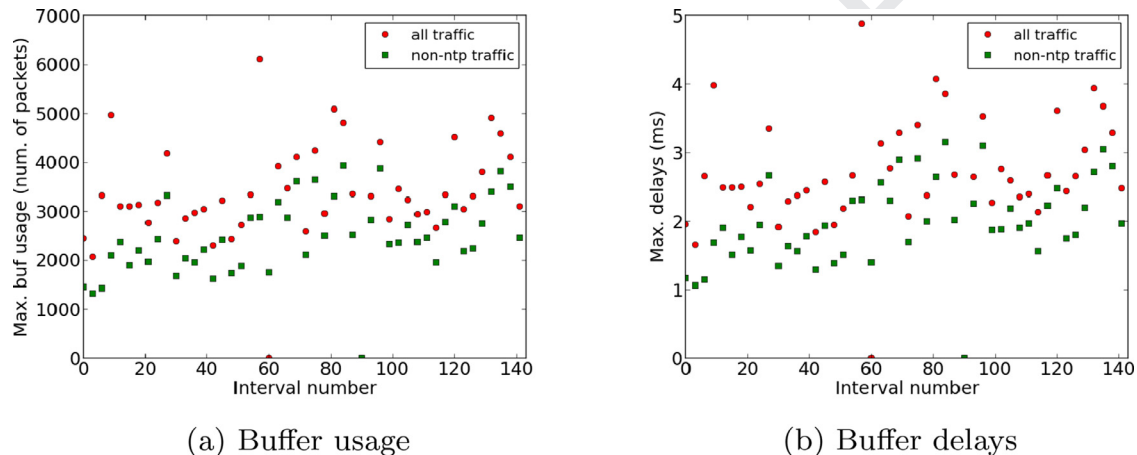


Fig. 15. Comparison of buffer usage and delays for all traffic and all non-NTP traffic. The cache size was set to 10K entries.

618 during the trace, NTP packets caused  $55 \cdot 10^6$  misses or 48%. The  
 619 next highest number of misses were caused by TCP DATA packets,  
 620 with  $25 \cdot 10^6$  misses (22%). NTP therefore caused 2.2 times more  
 621 misses than the TCP DATA packets. Fig. 14b shows a zoomed-in  
 622 version of Fig. 14a, with the NTP traffic filtered out of the trace.  
 623 We see that TCP DATA and DNS packets account for the bulk of  
 624 the misses, similar to what we see in Fig. 7.

625 To further quantify the performance of the cache during the  
 626 NTP incident, we measured the maximum buffer usage and queuing  
 627 delays incurred by packets during the ongoing NTP attacks. We  
 628 compare these with the buffer usage and queuing delays with the  
 629 NTP traffic filtered out. Fig. 15a and b shows the results.

630 As Fig. 15a and b shows, the buffer usage and queuing delays  
 631 with the attack traffic are not drastically different to those with  
 632 the attack traffic filtered out. The maximum buffer usage with the  
 633 attack traffic included is 6104 packets, compared to 3944 pack-  
 634 ets without the attack traffic. The corresponding maximum queuing  
 635 delays are 4.8 ms and 3.2 ms respectively. Thus we conclude that  
 636 the NTP incident did not put undue strain on the caching system.  
 637

638 The difference in hit rates observed for traces T8 and T9 also re-  
 639 flects in the buffer usage and buffering delays, as seen from Figs. 9  
 640 and 15. Trace T9 and T8 were captured on different links and trace  
 641 T9 had more diversity in terms of the number of prefixes carrying  
 642 packets than T8. In trace T8, the average number of unique prefixes  
 643 carrying packets in each 5 min interval we investigated was 15,310,  
 644 compared to 24,965 prefixes for trace T9. In future work, we plan  
 645 to investigate traces T8 and T9 further to quantify the differences  
 646 observed in the results.

## 9. Conclusions

647

648 In this paper we take a look at some practical considerations  
 649 in the implementation of FIB caching. We extend previous work  
 650 in significant ways by looking at practical issues, such as charac-  
 651 terization of cache misses, queuing issues (buffer occupancy and  
 652 delay), memory bandwidth requirements and robustness against  
 653 cache pollution attacks. We used traces collected at links from our  
 654 regional ISP to a Tier-1 ISP for our analysis. Our design uses a  
 655 cacheable FIB to avoid the cache hiding problem, and we presented  
 656 an algorithm to convert a regular FIB to a cacheable FIB.

657 Our work has some limitations. First, we only look at packet  
 658 traces from a single regional ISP. We therefore cannot evaluate  
 659 cache performance at core routers, where traffic may be more di-  
 660 verse causing hit rates to drop. While we do not have access to  
 661 data from core routers to answer this question (we need a packet  
 662 trace and a simultaneous snapshot of the FIB at a core router), the  
 663 tools and methodology we developed are applicable to a core en-  
 664 vironment and we plan to repeat the study once we have an ap-  
 665 propriate dataset.

666 Are there trends that may invalidate the benefits of FIB  
 667 caching? On the contrary, recent trends such as traffic concentra-  
 668 tion to major social networks, search engines that reside in data-  
 669 centers [14] and CDNs make caching even more likely to provide  
 670 benefits. In these environments, traffic becomes more focused and  
 671 likely to hit a smaller set of prefixes, resulting into a more stable  
 672 working set.

673 Other potential benefits of employing FIB caching include man-  
 674 ufacturing of cheaper router hardware and routers with longer

675 service cycles. Future terabit/s forwarding chips may employ FIB  
 676 caching for less on-chip memory or to allow the non-cached por-  
 677 tion of the FIB to be more aggressively compressed. Finally, cur-  
 678 rent architectures with the entire FIB on DRAM may also benefit  
 679 by backing away from the looming memory bandwidth wall.

## 680 References

- Q6 681 [1] Endace, <http://www.endace.com>.  
 682 [2] SANS ISC, NTP Reflection Attack – Internet Security, SANS ISC. <https://isc.sans.edu/forums/diary/NTP+reflection+attack/17300>.  
 683 [3] How to choose the best router switching path for your network, 2005. <http://www.cisco.com/warp/public/105/20.pdf>.  
 684 [4] H. Ballani, P. Francis, T. Cao, J. Wang, Making routers last longer with viaggre, in: Proceedings of the 6th USENIX Symposium on Networked Systems Design and Implementation, NSDI'09, USENIX Association, Berkeley, CA, USA, 2009, pp. 453–466.  
 685 [5] J. Czyz, M. Kallitsis, M. Gharaibeh, C. Papadopoulos, M. Bailey, M. Karir, Taming the 800 pound gorilla: the rise and decline of NTP DDOS attacks, in: Proceedings of the 2014 Internet Measurement Conference, IMC'14, ACM, New York, NY, USA, 2014, pp. 435–448, doi:10.1145/2663716.2663717.  
 686 [6] R. Draves, C. King, S. Venkatachary, B. Zill, Constructing optimal IP routing tables, in: Proceedings of the Eighteenth Annual Joint Conference of the IEEE Computer and Communications Societies, INFOCOM'99, vol. 1, 1999, pp. 88–97, doi:10.1109/INFCOM.1999.749256.  
 687 [7] D. Farinacci, V. Fuller, D. Meyer, D. Lewis, Locator/ID Separation Protocol (LISP), draft-ietf-lisp-08.txt, 2010.  
 688 [8] D. Feldmeier, Improving gateway performance with a routing-table cache, in: Proceedings of the INFOCOM'88, Networks: Evolution or Revolution? Proceedings of the Seventh Annual Joint Conference of the IEEE Computer and Communications Societies, 1988, pp. 298–307, doi:10.1109/INFCOM.1988.12930.  
 689 [9] Y. Gao, L. Deng, A. Kuzmanovic, Y. Chen, Internet cache pollution attacks and countermeasures, in: Proceedings of the 2006 IEEE International Conference on Network Protocols, ICNP '06, IEEE Computer Society, Washington, DC, USA, 2006, pp. 54–64, doi:10.1109/ICNP.2006.320198.  
 690 [10] W. Herrin, Opportunistic topological aggregation in the RIB-FIB calculation?, <http://www.ops.ietf.org/lists/rng/2008/msg01880.html>, 2008.  
 691 [11] G. Huston, BGP analysis reports. <http://bgp.potaroo.net/index-bgp.html>.  
 692 [12] G. Huston, G. Armitage, Projecting future IPV4 router requirements from trends in dynamic BGP behaviour, in: Proceedings of the Australian Telecommunication Networks and Applications Conference (ATNAC), 2006.  
 693 [13] C. Kim, M. Caesar, A. Gerber, J. Rexford, Revisiting route caching: the world should be flat, in: Proceedings of the 10th International Conference on Passive and Active Network Measurement, PAM'09, 2009.
- [14] C. Labovitz, S. Iekel-Johnson, D. McPherson, J. Oberheide, F. Jahanian, Internet inter-domain traffic, in: Proceedings of the ACM SIGCOMM 2010 Conference, 2010.  
 [15] H. Liu, Routing prefix caching in network processor design, in: Proceedings of the Tenth International Conference on Computer Communications and Networks, 2001, pp. 18–23, doi:10.1109/ICCCN.2001.956214.  
 [16] Y. Liu, S.O. Amin, L. Wang, Efficient FIB caching using minimal non-overlapping prefixes, SIGCOMM Comput. Commun. Rev. 43 (1) (2012) 14–21, doi:10.1145/2427036.2427039.  
 [17] Y. Liu, L. Wang, B. Zhang, FIFA: fast incremental FIB aggregations, in: Proceedings of the 32nd IEEE International Conference on Computer Communications, INFOCOM, 2013.  
 [18] Y. Liu, X. Zhao, K. Nam, L. Wang, B. Zhang, Incremental forwarding table aggregation, in: Proceedings of the GlobeCom Conference, 2010.  
 [19] V. Manivel, M. Ahamad, H. Venkateswaran, Attack resistant cache replacement for survivable services, in: Proceedings of the 2003 ACM Workshop on Survivable and Self-regenerative Systems in association with 10th ACM Conference on Computer and Communications Security, SSRS'03, ACM, New York, NY, USA, 2003, pp. 64–71, doi:10.1145/1036921.1036928.  
 [20] D. Massey, L. Wang, B. Zhang, L. Zhang, Proposal for scalable internet routing and addressing, draft-wang-ietf-efit-00.txt, 2007.  
 [21] D. Meyer, L. Zhang, K. Fall, Report From the IAB Workshop on Routing and Addressing, 2007.  
 [22] University of Oregon, The Route Views Project, Advanced Network Topology Center and University of Oregon. <http://www.routeviews.org/>.  
 [23] S. Richardson, Vertical aggregation: a strategy for FIB reduction, 1996.  
 [24] N. Sarrar, T.U.B. T-labs, S. Uhlig, R. Sherwood, Leveraging Zipf's law for traffic offloading 42 (1) (2012) 17–22.  
 [25] Z.A. Uzmi, M. Nebel, A. Tariq, S. Jawad, R. Chen, A. Shaikh, J. Wang, P. Francis, SMALTA: practical and near-optimal FIB aggregation, in: Proceedings of the Seventh Conference on emerging Networking Experiments and Technologies, CoNEXT'11, 2011, pp. 29:1–29:12, doi:10.1145/2079296.2079325.  
 [26] H. Yan, B. Say, B. Sheridan, D. Oko, C. Papadopoulos, D. Pei, D. Massey, IP reachability differences: myths and realities, in: Proceedings of the 2011 IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS), 2011, pp. 834–839, doi:10.1109/INFCOMW.2011.5928928.  
 [27] M. Zadnik, M. Canini, Evolution of cache replacement policies to track heavy-hitter flows, in: Proceedings of the 12th International Conference on Passive and Active Network Measurement, PAM'11, 2011, pp. 21–30.  
 [28] X. Zhang, P. Francis, J. Wang, K. Yoshida, Scaling IP routing with the core router-integrated overlay, in: Proceedings of the 2006 IEEE International Conference on Network Protocols, ICNP'06, 2006, pp. 147–156. <http://dx.doi.org/10.1109/ICNP.2006.320208>.  
 [29] X. Zhao, Y. Liu, L. Wang, B. Zhang, On the aggregatability of router forwarding tables, in: Proceedings of the IEEE INFOCOM 2010 Conference, 2010.