

EndRE: An End-System Redundancy Elimination Service for Enterprises

Bhavish Aggarwal^{*}, Aditya Akella[†], Ashok Anand^{1†}, Athula Balachandran^{1‡},
Pushkar Chitnis^{*}, Chitra Muthukrishnan^{1†}, Ramachandran Ramjee^{*} and George Varghese^{2§}

^{*} Microsoft Research India; [†] University of Wisconsin-Madison; [‡] CMU; [§]UCSD

Abstract

In many enterprises today, WAN optimizers are being deployed in order to eliminate redundancy in network traffic and reduce WAN access costs. In this paper, we present the design and implementation of EndRE, an alternate approach where redundancy elimination (RE) is provided as *an end system service*. Unlike middleboxes, such an approach benefits both end-to-end encrypted traffic as well as traffic on last-hop wireless links to mobile devices.

EndRE needs to be fast, adaptive and parsimonious in memory usage in order to opportunistically leverage resources on end hosts. Thus, we design a new fingerprinting scheme called SampleByte that is much faster than Rabin fingerprinting while delivering similar compression gains. Unlike Rabin fingerprinting, SampleByte can also adapt its CPU usage depending on server load. Further, we introduce optimizations to reduce server memory footprint by 33-75% compared to prior approaches. Using several terabytes of network traffic traces from 11 enterprise sites, testbed experiments and a pilot deployment, we show that EndRE delivers 26% bandwidth savings on average, processes payloads at speeds of 1.5-4Gbps, reduces end-to-end latencies by up to 30%, and translates bandwidth savings into equivalent energy savings on mobile smartphones.

1 Introduction

With the advent of globalization, networked services have a global audience, both in the consumer and enterprise spaces. For example, a large corporation today may have branch offices at dozens of cities around the globe. In such a setting, the corporation’s IT admins and network planners face a dilemma. On the one hand, they could concentrate IT servers at a small number of locations. This might lower administration costs, but increase network costs and latency due to the resultant increase in WAN traffic. On the other hand, servers could be located closer to clients; however, this would increase operational costs.

This paper arises from the quest to have the best of both worlds, specifically, having the operational benefits of centralization along with the performance benefits of distribution. In recent years, protocol-independent redundancy elimination, or simply RE [20], has helped bridge the gap by making WAN communication more efficient through elimination of redundancy in traffic. Such compression is typically applied at the IP or TCP layers, for instance, using a pair of middleboxes placed at either end of a WAN link connecting a corporation’s data center and a branch office. Each box caches payloads from flows that traverse the link, irrespective of the application or protocol. When one box detects chunks of data that match entries in its cache (by computing “fingerprints” of incoming data and matching them against cached data), it encodes matches using tokens. The box at the far end reconstructs original data using its own cache and the tokens. This approach has seen increasing deployment in “WAN optimizers”.

Unfortunately, such middlebox-based solutions face two key drawbacks that impact their long-term usefulness: (1) Middleboxes do not cope well with end-to-end encrypted traffic and many leave such data uncompressed (e.g., [1]). Some middleboxes accommodate SSL/SSH traffic with techniques such as connection termination and sharing of encryption keys (e.g., [5]), but these weaken end-to-end semantics. (2) In-network middleboxes cannot improve performance over last-hop links in mobile devices.

As end-to-end encryption and mobile devices become increasingly prevalent, we believe that RE will be *forced out* of middleboxes and directly *into* end-host stacks. Motivated by this, we explore a new point in the design space of RE proposals — an *end-system redundancy elimination service* called EndRE. EndRE could supplement or supplant middlebox-based techniques while addressing their drawbacks. Our paper examines the costs and benefits that EndRE implies for clients and servers.

Effective end-host RE requires looking for small redundant chunks of the order of 32-64 bytes (because most enterprise transfers involve just a few packets each [16]). The standard Rabin fingerprinting algorithms (e.g., [20]) for identifying such fine scale redundancy are very expensive in terms of memory and processing especially on resource constrained clients such

¹A part of this work was done while the authors were interns at Microsoft Research India.

²The author was a visiting researcher at Microsoft Research India during the course of this work.

as smartphones. Hence, we adopt a novel *asymmetric* design that systematically offloads as much of processing and memory to servers as possible, requiring clients to do no more than perform basic FIFO queue management of a small amount of memory and do simple pointer lookups to decode compressed data.

While client processing and memory are paramount, servers in EndRE need to do other things as well. This means that server CPU and memory are also crucial bottlenecks in our asymmetric design. For server processing, we propose a new fingerprinting scheme called SampleByte that is much faster than Rabin fingerprinting used in traditional RE approaches while delivering similar compression. In fact, SampleByte can be up to 10X faster, delivering compression speeds of 1.5-4Gbps. SampleByte is also *tunable* in that it has a payload sampling parameter that can be adjusted to reduce server processing if the server is busy, at the cost of reduced compression gains.

For server storage, we devise a suite of highly-optimized data structures for managing meta-data and cached payloads. For example, our Max-Match variant of EndRE (§5.2.2) requires 33% lower memory compared to [20]. Our Chunk-Match variant (§5.2.1) cuts down the aggregate memory requirements at the server by 4X compared to [20], while sacrificing a small amount of redundancy.

We conduct a thorough evaluation of EndRE. We analyze several terabytes of traffic traces from 11 different enterprise sites and show that EndRE can deliver significant bandwidth savings (26% average savings) on enterprise WAN links. We also show significant latency and energy savings from using EndRE. Using a testbed over which we replay enterprise HTTP traffic, we show that latency savings of up to 30% are possible from using EndRE, since it operates above TCP, thereby reducing the number of roundtrips needed for data transfer. Similarly, on mobile smartphones, we show that the low decoding overhead on clients can help translate bandwidth savings into significant energy savings compared to no compression. We also report results from a small-scale deployment of EndRE in our lab.

The benefits of EndRE come at the cost of memory and CPU resources on end systems. We show that a median EndRE client needs only 60MB of memory and negligible amount of CPU. At the server, since EndRE is adaptive, it can opportunistically trade-off CPU/memory for compression savings.

In summary, we make the following contributions:

- (1) We present the design of EndRE, an end host based redundancy elimination service (§4).
- (2) We present new asymmetric RE algorithms and optimized data structures that limit client processing and memory requirements, and reduce server memory us-

age by 33-75% and processing by 10X compared to [20] while delivering slightly lower bandwidth savings (§5).

(3) We present an implementation of EndRE as part of Windows Server/7/Vista as well as on Windows Mobile 6 operating systems (§6).

(4) Based on extensive analysis using several terabytes of network traffic traces from 11 enterprise sites, testbed experiments and a small-scale deployment, we quantify the benefits and costs of EndRE (§7 - §9)

2 Related Work

Over the years, enterprise networks have used a variety of mechanisms to suppress duplicate data from their network transfers. We review these mechanisms below.

Classical approaches: The simplest RE approach is to compress objects end-to-end. It is also the least effective because it does not exploit redundancy due to repeated accesses of similar content. Object caches can help in this regard, but they are unable to extract cross-object redundancies [8]. Also object caches are application-specific in nature; e.g., Web caches cannot identify duplication in other protocols. Furthermore, an increasing amount of data is dynamically generated and hence not cacheable. For example, our analysis of enterprise traces shows that a majority of Web objects are not cacheable, and deploying an HTTP proxy would only yield 5% net bandwidth savings. Delta encoding can eliminate redundancy of one Web object with respect to another [14, 12]. However, like Web caches, delta encoding is application-specific and ineffective for dynamic content.

Content-based naming: The basic idea underlying EndRE is that of *content-based naming* [15, 20], where an object is divided into chunks and indexed by computing hashes over chunks. Rabin Fingerprinting [18] is typically used to identify chunk boundaries. In file systems such as LBFS [15] and Shark [9], content-based naming is used to identify similarities across different files and across versions of the same file. Only unique chunks are transmitted between file servers and clients, resulting in lower bandwidth consumption. A similar idea is used in value-based Web caching [19], albeit between a Web server and its client. Our chunk-based EndRE design is patterned after this approach, with key modifications for efficiency (§5).

Generalizing these systems, DOT [21] proposes a “transfer service” as an interface between applications and the network. Applications pass the object they want to send to DOT. Objects are split into chunks and the sender sends chunk hashes to the receiver. The receiver maintains a cache of earlier received chunks and requests only the chunks that were not found in its cache or its neighbors’ caches. Thus, DOT can leverage TBs of cache in the disks of an end host and its peers to elim-

inate redundancy. Similarly, SET [17] exploits chunk-level similarity in downloading related large files. DOT and SET use an average chunk size of 2KB or more. These approaches mainly benefit large transfers; the extra round trips that can only be amortized over the transfer lengths. In contrast, EndRE identifies redundancy across chunk sizes of 32 bytes and does not impose additional latency. It is also limited to main-memory based caches of size 1-10MB per pair of hosts (§5). Thus, EndRE and DOT complement each other.

Protocol-independent WAN optimizers. To overcome the limitations of the “classical” approaches, enterprises have moved increasingly toward protocol independent RE techniques, used in WAN optimizers. These WAN optimizers can be of two types, depending on which network layer they operate at, namely, IP layer devices [20, 3] or higher-layer devices [1, 5].

In either case, special middleboxes are deployed at either end of a WAN link to index all content exchanged across the link, and identify and remove partial redundancies on the fly. Rabin fingerprinting [18] is used to index content and compute overlap (similar to [20, 15]). Both sets of techniques are highly effective at reducing the utilization of WAN links. However, as mentioned earlier, they suffer from two key limitations, namely, lack of support for end-to-end encryption and for resource-constrained mobile devices.

3 Motivation

In exploring an end-point based RE service, one of the main issues we hope to address is whether such a service can offer bandwidth savings approaching that of WAN optimizers. To motivate the likely benefits of an end-point based RE service, we briefly review two key findings from our earlier study [8] of an IP-layer WAN optimizer [7].

First, we seek to identify the origins of redundancy. Specifically, we classify the contribution of redundant byte matches to bandwidth savings as either *intra-host* (current and matched packet in cache have identical source-destination IP addresses) or *inter-host* (current and matched packets differ in at least one of source or destination IP addresses). We were limited to a 250MB cache size given the large amount of meta-data necessary for this analysis, though we saw similar compression savings for cache sizes up to 2GB. Surprisingly, our study revealed that *over 75% of savings were from intra-host matches*. This implies that a pure end-to-end solution could potentially deliver a significant share of the savings obtained by an IP WAN optimizer, since the contribution due to inter-host matches is small. However, this finding holds good only if end systems operate with similar (large) cache sizes as middleboxes, which is impractical. This brings us to the second key finding.

Examining the temporal characteristics of redundancy, we found that the redundant matches in the WAN optimizer displayed a high degree of temporal locality with *60-80% of middlebox savings arising from matches with packets in the most recent 10% of the cache*. This implies that small caches could capture a bulk of the savings of a large cache.

Taken together, these two findings suggest that an end-point-based RE system with a small cache size can indeed deliver a significant portion of the savings of a WAN optimizer, thus motivating the design of EndRE.

Finally, note that, the focus of comparison in this section is between an IP-layer WAN optimizer with an in-memory cache (size is $O(GB)$) and an end-system solution. The first finding is not as surprising once we realize that the in-memory cache gets recycled frequently (on the order of tens of minutes) during peak hours on our enterprise traces, limiting the possibility for inter-host matches. A WAN optimizer typically also has a much larger on-disk cache (size is $O(TB)$) which may see a large fraction of inter-host matches; an end-system disk cache-based solution such as DOT [21] could capture analogous savings.

4 Design Goals

EndRE is designed to optimize data transfers in the direction from servers in a remote data center to clients in the enterprise, since this captures a majority of enterprise traffic. We now list five design goals for EndRE — the first two design goals are shared to some extent by prior RE approaches, but the latter three are unique to EndRE.

1. Transparent operation: For ease of deploy-ability, the EndRE service should require no changes to existing applications run within the data center or on clients.

2. Fine-grained operation: Prior work has shown that many enterprise network transfers involve just a few packets [16]. To improve end-to-end latencies and provide bandwidth savings for such short flows, EndRE must work at fine granularities, suppressing duplicate byte strings as small as 32-64B. This is similar to [20], but different from earlier proposals for file-systems [15] and Web caches [19] where the sizes of redundancies identified are 2-4KB.

3. Simple decoding at clients: EndRE’s target client set includes battery- and CPU-constrained devices such as smart-phones. While working on fine granularities can help identify greater amounts of redundancy, it can also impose significant computation and decoding overhead, making the system impractical for these devices. Thus, a unique goal is to design algorithms that limit client overhead by *offloading* all compute-intensive actions to *servers*.

4. Fast and adaptive encoding at servers: EndRE is

designed to opportunistically leverage CPU resources on end hosts when they are not being used by other applications. Thus, unlike commercial WAN optimizers and prior RE approaches [20], EndRE must *adapt* its use of CPU based on server load.

5. Limited memory footprint at servers and clients:

EndRE relies on data caches to perform RE. However, memory on servers and clients could be limited and may be actively used by other applications. Thus, EndRE must use as minimal memory on end-hosts as possible through the use of optimized data structures.

5 EndRE Design

In this section, we describe how EndRE’s design meets the above goals.

EndRE introduces RE modules into the network stacks of clients and remote servers. Since we wish to be transparent to applications, EndRE could be implemented either at the IP-layer or at the socket layer (above TCP). As we argue in §6, we believe that socket layer is the right place to implement EndRE. Doing so offers key performance benefits over an IP-layer approach, and more importantly, shields EndRE from network-level events (e.g., packet losses and reordering), making it simpler to implement.

There are two sets of modules in EndRE, those belonging on servers and those on clients. The server-side module is responsible for identifying redundancy in network data by comparing against a cache of prior data, and encoding the redundant data with shorter meta-data. The meta-data is essentially a set of $\langle \text{offset}, \text{length} \rangle$ tuples that are computed with respect to the client-side cache. The client-side module is trivially simple: it consists of a fixed-size circular FIFO log of packets and simple logic to decode the meta-data by “de-referencing” the offsets sent by the server. Thus, most of the complexity in EndRE is mainly on the server side and we focus on that here.

Identifying and removing redundancy is typically accomplished [20, 7] by the following two steps:

- *Fingerprinting*: Selecting a few “representative regions” for the current block of data handed down by application(s). We describe four fingerprinting algorithms in §5.1 that differ in the trade-off they impose between *computational overhead* on the server and the *effectiveness* of RE.
- *Matching and Encoding*: Once the representative regions are identified, we examine two approaches for identification of redundant content in §5.2: (1) Identifying chunks of representative regions that repeat in full across data blocks, called Chunk-Match and (2) Identifying maximal matches around the representative regions that are repeated across data blocks, called Max-Match. These two approaches differ in the trade-off be-

tween the *memory overhead* imposed on the server and the *effectiveness* of RE.

Next, we describe EndRE’s design in detail, starting with selection of representative regions, and moving on to matching and encoding.

5.1 Fingerprinting: Balancing Server Computation with Effectiveness

In this section, we outline four approaches for identifying the representative payload regions at the server that vary in the way they trade-off between computational overhead and the effectiveness of RE. In some of the approaches, computational overhead can be *adaptively tuned* based on server CPU load, and the effectiveness of RE varies accordingly. Although three of the four approaches were proposed earlier, the issue of their computational overhead has not received enough attention. Since this issue is paramount for EndRE, we consider it in great depth here. We also propose a new approach, SAMPLEBYTE, that combines the salient aspects of prior approaches.

We first introduce some notation and terminology to help explain the approaches. Restating from above, a “data block” or simply a “block” is a certain amount of data handed down by an application to the EndRE module at the socket layer. Each data block can range from a few bytes to tens of kilobytes in size.

Let w represent the size of the minimum redundant string (contiguous bytes) that we would like to identify. For a data block of size S bytes, $S \geq w$, a total of $S - w + 1$ strings of size w are potential candidates for finding a match. Typical values for w range from 12 to 64 bytes. Based on our findings of redundant match length distribution in [8], we choose a default value of $w = 32$ bytes to maximize the effectiveness of RE. Since $S \gg w$, the number of such candidate strings is on the order of the number of bytes in the data block/cache. Since it is impractical to match/store all possible candidates, a fraction $1/p$ “representative” candidates are chosen.

Let us define *markers* as the first byte of these chosen candidate strings and *chunks* as the string of bytes between two markers. Let *fingerprints* be a pseudo-random hash of fixed w -byte strings beginning at each marker and *chunk-hashes* be hashes of the variable sized chunks. Note that two fingerprints may have overlapping bytes; however, by definition, chunks are disjoint. The different algorithms, depicted in Figure 1 and discussed below, primarily vary in the manner in which they choose the markers, from which one can derive chunks, fingerprints, and chunk-hashes. As we discuss later in §5.2, the Chunk-Match approach uses chunk-hashes while Max-Match uses fingerprints.

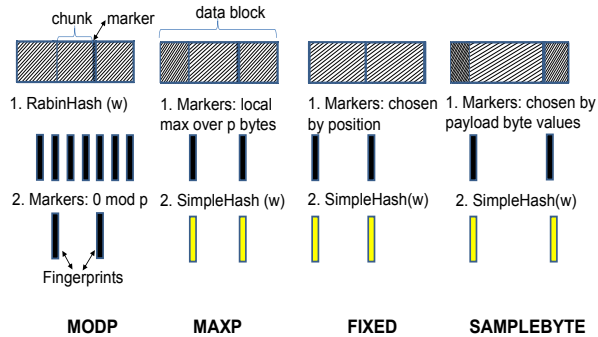


Figure 1: Fingerprinting algorithms with chunks, markers and fingerprints; chunk-hashes, not shown, can be derived from chunks

```

1 //Let  $w = 32; p = 32$ ; Assume  $len \geq w$ ;
2 //RabinHash() computes RABIN hash over a  $w$  byte window
3 MODP(data, len)
4   for( $i = 0; i < w - 1; i++$ )
5     fingerprint = RabinHash(data[i]);
6   for( $i = w - 1; i < len; i++$ )
7     fingerprint = RabinHash(data[i]);
8     if (fingerprint %  $p == 0$ ) //MOD
9       marker =  $i - w + 1$ ;
10    store marker, fingerprint in table;

```

Figure 2: MODP Fingerprinting Algorithm

5.1.1 MODP

In the “classical” RE approaches [20, 7, 15], the set of fingerprints are chosen by first computing Rabin-Karp hash [18] over sliding windows of w contiguous bytes of the data block. A fraction $1/p$ are chosen whose fingerprint value is $0 \bmod p$. Choosing fingerprints in this manner has the advantage that the set of representative fingerprints for a block remains mostly the same despite small amount of insertions/deletions/reorderings since the markers/fingerprints are chosen based on content rather than position.

Note that two distinct operations — marker identification and fingerprinting — are both handled by the same hash function here. While this *appears* elegant, it has a cost. Specifically, the per block computational cost is independent of the sampling period, p (lines 4–7 in Figure 2). Thus, this approach *cannot* adapt to server CPU load conditions (e.g., by varying p). Note that, while the authors of [20] report some impact of p on processing speed, this impact is attributed to the overhead of managing meta-data (line 10). We devise techniques in §5.2 to significantly reduce the overhead of managing meta-data, thus, making fingerprint computation the main bottleneck.

```

1 //Let  $w = 32; p = 32$ ; Assume  $len \geq w$ ;
2 //SAMPLETABLE[i] maps byte  $i$  to either 0 or 1
3 //JenkinsHash() computes hash over a  $w$  byte window
4 SAMPLEBYTE(data, len)
5   for( $i = 0; i < len - w; i++$ )
6     if (SAMPLETABLE[data[i]] == 1)
7       marker =  $i$ ;
8       fingerprint = JenkinsHash(data +  $i$ );
9       store marker, fingerprint in table;
10     $i = i + p/2$ ;

```

Figure 3: SAMPLEBYTE Fingerprinting Algorithm

5.1.2 MAXP

Apart from the conflation of marker identification and fingerprinting, another shortcoming of the MODP approach is that the fingerprints/markers are chosen based on a *global* property, i.e., fingerprints have to take certain pre-determined values to be chosen. The markers for a given block may be clustered and there may be large intervals without any markers, thus, limiting redundancy identification opportunities. To guarantee that an adequate number of fingerprints/markers are chosen uniformly from each block, markers can be chosen as bytes that are *local-maxima over each region of p bytes* of the data block [8]. Once the marker byte is chosen, an efficient hash function such as Jenkins Hash [2] can be used to compute the fingerprint. By increasing p , fewer maxima-based markers need to be identified, thereby reducing CPU overhead.

5.1.3 FIXED

While markers in both MODP and MAXP are chosen based on content of the data block, the computation of Rabin hashes and local maxima can be expensive. A simpler approach is to be content-agnostic and simply select *every p^{th} byte as a marker*. Since markers are simply chosen by position, marker identification incurs no computational cost. Once markers are chosen, S/p fingerprints are computed using Jenkins Hash as in MAXP. While this technique is very efficient, its effectiveness in RE is not clear as it is not robust to small changes in content. While prior works in file systems (e.g., [15]), where cache sizes are large ($O(TB)$), argue against this approach, it is not clear how ineffective FIXED will be in EndRE where cache sizes are small ($O(MB)$).

5.1.4 SAMPLEBYTE

MAXP and MODP are content-based and thus robust to small changes in content, while FIXED is content-agnostic but computationally efficient. We designed SAMPLEBYTE (Figure 3) to combine the robustness of a content-based approach with the computational efficiency of FIXED. It uses a 256-entry lookup table with

a few predefined positions set. As the data block is scanned byte-by-byte (line 5), a byte is chosen as a marker if the corresponding entry in the lookup table is set (line 6–7). Once a marker is chosen, a fingerprint is computed using Jenkins Hash (line 8), and $p/2$ bytes of content are skipped (line 10) before the process repeats. Thus, SAMPLEBYTE is content-based, albeit based on a single byte, while retaining the content-skipping and the computational characteristics of FIXED.

One clear concern is whether such a naive marker identification approach will do badly and cause the algorithm to either over-sample or under-sample. First, note that MODP with 32-64 byte rolling hashes was originally used in file systems [15] where chunk sizes were large (2-4KB). Given that we are interested in sampling as frequent as every 32-64 bytes, sampling chunk boundaries based on 1-byte content values is not as radical as it might first seem. Also, note that if x entries of the 256-entry lookup table are randomly set (where $256/x = p$), then the expected sampling frequency is indeed $1/p$. In addition, SAMPLEBYTE skips $p/2$ bytes after each marker selection to avoid oversampling when the content bytes of data block are not uniformly distributed (e.g., when the same content byte is repeated contiguously). Finally, while a purely random selection of $256/x$ entries does indeed perform well in our traces, we use a lookup table derived based on the heuristic described below. This heuristic outperforms the random approach and we have found it to be effective after extensive testing on traces (see §8).

Since the number of unique lookup tables is large (2^{256}), we use an offline, greedy approach to generate the lookup table. Using network traces from one of the enterprise sites we study as training data (site 11 in Table 2), we first run MAXP to identify redundant content and then sort the characters in descending order of their presence in the identified redundant content. We then add these characters one at a time, setting the corresponding entries in the lookup table to 1, and stop this process when we see diminishing gains in compression. The intuition behind this approach is that characters that are more likely to be part of redundant content should have a higher probability of being selected as markers. The characters selected from our training data were 0, 32, 48, 101, 105, 115, 116, 255. While our current approach results in a static lookup table, we are looking at online dynamic adaptation of the table as part of future work.

Since SAMPLEBYTE skips $p/2$ bytes after every marker selection, the fraction of markers chosen is $\leq 2/p$, irrespective of the number of entries set in the table. By increasing p , fewer markers/fingerprints are chosen, resulting in reduced CPU overhead.

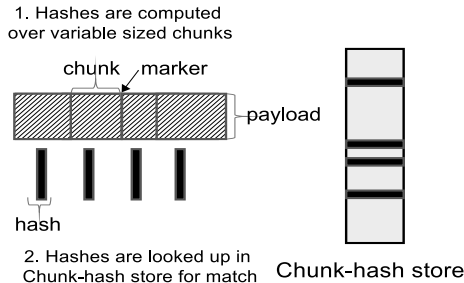


Figure 4: Chunk-Match: only chunk-hashes stored

5.2 Matching and Encoding: Optimizing Storage and Client Computation

Once markers and fingerprints are identified, identification of redundant content can be accomplished in two ways: (1) Identifying *chunks* of data that repeat in full across data blocks, called Chunk-Match, or (2) Identifying *maximal matches around fingerprints* that are repeated across data blocks, called Max-Match. Both techniques were proposed in prior work: the former in the context of file systems [15] and Web object compression [19], and the latter in the context of IP WAN optimizer [20]. However, prior proposals impose significant storage and CPU overhead.

In what follows we describe how the overhead impacts both servers and clients, and the two techniques we employ to address these overheads. The first technique is to leverage *asymmetry* between servers and clients. We propose that clients offload most of the computationally intensive operations (e.g., hash computations) and memory management tasks to the server. The second technique is to exploit the inherent *structure* within the data maintained at servers and clients to optimize memory usage.

5.2.1 Chunk-Match

This approach (Figure 4) stores hashes of the chunks in a data block in a “Chunk-hash store”. Chunk-hashes from payloads of future data blocks are looked up in the Chunk-hash store to identify if one or more chunks have been encountered earlier. Once matching chunks are identified, they are replaced by meta-data.

Although similar approaches were used in prior systems, they impose significantly higher overhead if employed directly in EndRE. For example, in LBFS [15], clients have to update their local caches with mappings between new content-chunks and corresponding content-hashes. This requires expensive SHA-1 hash computation at the client. Value-based web caching [19] avoids the cost of hash computation at the client by having the server send the hash with each chunk. However, the client still needs to store the hashes, which is a significant overhead for small chunk sizes. Also, sending

hashes over the network adds significant overhead given that the hash sizes (20 bytes) are comparable to average chunk sizes in EndRE (32-64 bytes).

EndRE optimizations: We employ two ideas to improve overhead on clients and servers.

(1) Our design carefully *offloads all storage management and computation to servers*. A client simply maintains a fixed-size circular FIFO log of data blocks. The server emulates client cache behavior on a *per-client basis*, and maintains within its Chunk-hash store a mapping of each chunk hash to the start memory addresses of the chunk in a client’s log along with the length of the chunk. For each matching chunk, the server simply encodes and sends a four-byte $\langle \text{offset}, \text{length} \rangle$ tuple of the chunk in the client’s cache. The client simply “de-references” the offsets sent by the server and reconstructs the compressed regions from local cache. This approach avoids the cost of storing and computing hashes at the client, as well as the overhead of sending hashes over the network, at the cost of slightly higher processing and storage at the server end.

(2) In traditional Chunk-Match approaches, the server maintains a log of the chunks locally. We observe that the server only needs to maintain an up-to-date chunk-hash store, but *it does not need to store the chunks themselves* as long as the chunk hash function is collision resistant. Thus, when a server computes chunks for a new data block and finds that some of the chunks are not at the client by looking up the chunk-hash store, it inserts mappings between the new chunk hashes and their expected locations in the client cache.

In our implementation, we use SHA-1 to compute 160 bit hashes, which has good collision-resistant properties. Let us now compute the storage requirements for Chunk-Match assuming a sampling period p of 64 bytes and a cache size of 16MB. The offset to the 16MB cache can be encoded in 24 bits and the length encoded in 8 bits assuming the maximum length of a chunk is limited to 256 bytes (recall that chunks are variable sized). Thus, server meta-data storage is 24 bytes per 64-byte chunk, comprising 4-bytes for the $\langle \text{offset}, \text{length} \rangle$ tuple and 20-bytes for SHA-1 hash. This implies that server memory requirement is about 38% of the client cache size.

5.2.2 Max-Match

A drawback of Chunk-Match is that it can only detect exact matches in the chunks computed for a data block. It could miss redundancies that, for instance, span contiguous portions of neighboring chunks or redundancies that only span portions of chunks. An alternate approach, called Max-Match, proposed for IP WAN optimizer [20, 7] and depicted in Figure 5, can identify such redundancies, albeit at a higher memory cost at the server.

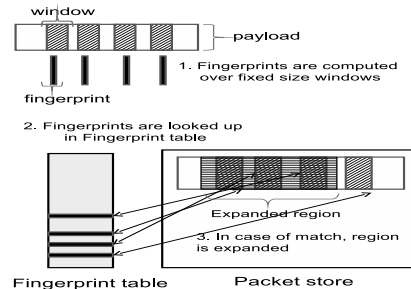


Figure 5: Max-Match: matched region is expanded

| index (implicit fingerprint, 18 bits) | fingerprint remainder (8 bits) | offset (24 bits) |
|---------------------------------------|--------------------------------|------------------|
| 0 | ... | ... |
| ... | ... | ... |
| $2^{18} - 1$ | ... | ... |

Table 1: 1MB Fingerprint store for 16MB cache

In Max-Match, fingerprints computed for a data block serve as random “hooks” into the payload around which more redundancies can be identified. The computed fingerprints for a data block are compared with a “fingerprint store” that holds fingerprints of all past data blocks. For each matching fingerprint, the corresponding matching data block is retrieved from the cache and the match region is expanded byte-by-byte in both directions to obtain the *maximal region* of redundant bytes (Figure 5). Matched regions are then encoded with $\langle \text{offset}, \text{length} \rangle$ tuples.

EndRE optimizations: We employ two simple ideas to improve the server computation and storage overhead.

First, since Max-Match relies on byte-by-byte comparison to identify matches, fingerprint collisions are not costly; any collisions will be recovered via an extra memory lookup. This allows us to significantly *limit fingerprint store maintenance overhead for all four algorithms* since fingerprint values are simply overwritten without separate bookkeeping for deletion. Further, a simple hash function that generates a few bytes of hash value as a fingerprint (e.g., Jenkins hash [2]) is sufficient.

Second, we optimize the representation of the fingerprint hash table to limit storage needs. Since the mapping is from a fingerprint to an offset value, the fingerprint itself need not be stored in the table, at least in its entirety. The index into the fingerprint table can implicitly represent part of the fingerprint and only the remaining bits, if any, of the fingerprint that are not covered by the index can be stored in the table. In the extreme case, the fingerprint table is simply a contiguous set of offsets, indexed by the fingerprint hash value.

Table 1 illustrates the fingerprint store for a cache size of 16MB and $p = 64$. In this case, the number of fingerprints to index the entire cache is simply $2^{24}/64$ or 2^{18} . Using a table size of 2^{18} implies that 18 bits of a fingerprint are implicitly stored as the index of the table. The

offset size necessary to represent the entire cache is 24 bits. Assuming we store an additional 8 bits of the fingerprint as part of the table, the entire fingerprint table can be compactly stored in a table of size $2^{18} * 4$ bytes, or 6% of the cache size. A sampling period of 32 would double this to 12% of the cache size. This leads to a significant reduction in fingerprint meta-data size compared to the 67% indexing overhead in [20] or the 50% indexing overhead in [7].

These two optimizations are not possible in the case of Chunk-Match due to the more stringent requirements on collision-resistance of chunk hashes. However, server memory requirement for Chunk-Match is only 38% of client cache size, which is still significantly lower than 106% of the cache size (cache + fingerprint store) needed for Max-Match.

6 Implementation

In this section, we discuss our implementation of EndRE. We start by discussing the benefits of implementing EndRE at the socket layer above TCP.

6.1 Performance benefits

Bandwidth: In the socket-layer approach, RE can operate at the size of socket writes which are typically larger than IP layer MTUs. While Max-Match and Chunk-Match do not benefit from these larger sized writes since they operate at a granularity of 32 bytes, the large size helps produce higher savings if a compression algorithm like GZIP is *additionally* applied, as evaluated in §9.1.

Latency: The socket-layer approach will result in fewer packets transiting between server and clients, as opposed to the IP layer approach which merely compresses packets without reducing their number. This is particularly useful in lowering completion times for short flows, as evaluated in §9.2.

6.2 End-to-end benefits

Encryption: When using socket-layer RE, payload encrypted in SSL can be compressed before encryption, providing RE benefits to protocols such as HTTPS. IP-layer RE will leave SSL traffic uncompressed.

Cache Synchronization: Recall that both Max-Match and Chunk-Match require caches to be synchronized between clients and servers. One of the advantages of implementing EndRE above TCP is that TCP ensures reliable in-order delivery, which can help with maintaining cache synchronization. However, there are still two issues that must be addressed.

First, multiple simultaneous TCP connections may be operating between a client and a server, resulting in ordering of data across connections not being preserved. To account for this, we implement a simple sequence number-based *reordering mechanism*.

Second, TCP connections may get reset in the middle of a transfer. Thus, packets written to the cache at the server end may not even reach the client, leading to cache inconsistency. One could take a *pessimistic* or *optimistic* approach to maintaining consistency in this situation. In the pessimistic approach, writes to the server cache are performed only after TCP ACKs for corresponding segments are received at the server. The server needs to monitor TCP state, detect ACKs, perform writes to its cache and notify the client to do the same. In the optimistic approach, the server writes to the cache but monitors TCP only for reset events. In case of connection reset (receipt of a TCP RST from client or a local TCP timeout), the server simply notifies the client of the last sequence number that was written to the cache for the corresponding TCP connection. It is then the client's responsibility to detect any missing packets and recover these from the server. We adopt the *optimistic approach of cache writing* for two reasons: (1) Our redundancy analysis [8] indicated that there is high temporal locality of matches; a pessimistic approach over a high bandwidth-delay product link can negatively impact compression savings; (2) The optimistic approach is easier to implement since only for reset events need to be monitored rather than every TCP ACK.

6.3 Implementation

We have implemented EndRE above TCP in Windows Server/Vista/7. Our default fingerprinting algorithm is SAMPLEBYTE with a sampling period, $p = 32$. Our packet cache is a circular buffer 1-16MB in size per pairs of IP addresses. Our fingerprint store is also allocated a bounded memory based on the values presented earlier. We use a simple resequencing buffer with a priority queue to handle re-ordering across multiple parallel TCP streams. At the client side, we maintain a fixed size circular cache and the decoding process simply involves lookups of specified data segments in the cache.

In order to enable protocol independent RE, we transparently capture application payloads on the server side and TCP payloads at the client side at the TCP stream layer, that lies between the application layer and the TCP transport layer. We achieve this by implementing a kernel level filter driver based on Windows Filtering Platform (WFP) [6]. This implementation allows seamless integration of EndRE with all application protocols that use TCP, with no modification to application binaries or protocols. We also have a management interface that can be used to restrict EndRE only to specific applications. This is achieved by predicate-based filtering in WFP, where predicates can be application IDs, source and/or destination IP addresses/ports.

Finally, we have also implemented the client-side of EndRE on mobile smartphones running the Windows

| Trace Name (Site #) | Unique Client IPs | Dates (Total Days) | Size (TB) |
|--|-------------------|--|-----------|
| Small Enterprise (Sites 1-2) | 29-39 | 07/28/08 - 08/08/08 (11) 11/07/08 - 12/10/08 (33) | 0.5 |
| Medium Enterprise (Sites 3-6) | 62-91 | 07/28/08 - 08/08/08 (11) 11/07/08 - 12/10/08 (33) | 1.5 |
| Large Enterprise (Sites 7-10) | 101-210 | 07/28/08 - 08/08/08 (11) 11/07/08 - 12/10/08 (33) | 3 |
| Large Research Lab (Site 11, training trace) | 125 | 06/23/08 - 07/03/08 (11) | 1 |

Table 2: Data trace characteristics (11 sites)

Mobile 6 OS. However, since Windows Mobile 6 does not support Windows Filtering Platform, we have implemented the functionality as a user-level proxy.

7 Evaluation approach

We use a combination of trace-based and testbed evaluation to study EndRE. In particular, we quantify bandwidth savings and evaluate scalability aspects of EndRE using enterprise network traffic traces; we use a testbed to quantify processing speed and evaluate latency and energy savings. We also report results from a small pilot deployment (15 laptops) in our lab spanning 1 week.

Traces: Our trace-based evaluation is based on full packet traces collected at the WAN access link of 11 corporate enterprise locations. The key characteristics of our traces are shown in Table 2. We classify the enterprises as small, medium or large based on the number of internal host IP addresses seen (less than 50, 50-100, and 100-250, respectively) in the entire trace at each of these sites. While this classification is somewhat arbitrary, we use this division to study if the benefits depend on the size of an enterprise. Note that the total amount of traffic in each trace is approximately correlated to the number of host IP addresses, though there is a large amount of variation from day to day. Typical incoming traffic numbers for small enterprises varied from 0.3-10GB/day, for medium enterprises from 2-12GB/day and for large enterprises from 7-50GB/day. The access link capacities at these sites varied from a few Mbps to several tens of Mbps. The total size of traffic we study (including in-bound/outbound traffic and headers) is about 6TB.

Testbed: Our testbed consists of a desktop server connected to a client through a router. In wireline experiments, the router is a dual-homed PC capable of emulating links of pre-specified bandwidth and latency. In wireless experiments, the router is a WiFi access point. The server is a desktop PC running Windows Server 2008. The client is a desktop PC running Windows Vista or Windows 7 in the wireline experiments, and a Samsung mobile smartphone running Windows Mobile 6 in the wireless experiments.

8 Costs

In this section, we quantify the CPU and memory costs of our implementation of EndRE. Though our evalua-

| Max-Match $p \rightarrow$ | Fingerprint | | InlineMatch | | Admin | |
|------------------------------|-------------|-------|-------------|-----|-------|-----|
| | 32 | 512 | 32 | 512 | 32 | 512 |
| MODP | 526.7 | 496.7 | 9.6 | 6.8 | 4.8 | 0.6 |
| MAXP | 306.3 | 118.8 | 10.1 | 7.7 | 5.2 | 0.5 |
| FIXED | 69.4 | 14.2 | 7.1 | 4.7 | 4.7 | 0.4 |
| SAMPLEBYTE(SB) | 76.8 | 20.2 | 9.5 | 6.1 | 3.0 | 0.7 |

Table 3: CPU Time(s) for different algorithms

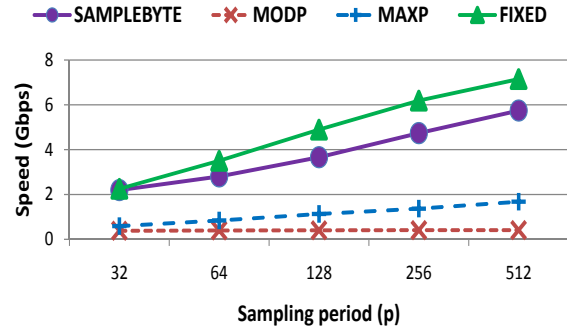


Figure 6: Max-Match processing speed

tion focus largely on Max-Match, we also provide a brief analysis of Chunk-Match.

8.1 CPU Costs

Micro-benchmarks: We first focus on micro-benchmarks for different fingerprinting algorithms using Max-Match for a cache size of 10MB between a given client-server pair (we examine cache size issues in detail in §8.2). Table 3 presents a profiler-based analysis of the costs of the three key processing steps on a single large packet trace as measured on a 2GHz 64-bit Intel Xeon processor. The fingerprinting step is responsible for identifying the markers/fingerprints; the Inline-Match function is called as fingerprints are generated; and the Admin function is used for updating the fingerprint store. Of these steps, only the fingerprinting step is distinct for the algorithms, and is also the most expensive.

One can clearly see that fingerprinting is expensive for MODP and is largely independent of p . Fingerprinting for MAXP is also expensive but we see that as p is increased, the cost of fingerprinting comes down. In the case of FIXED and SAMPLEBYTE, as expected, fingerprinting cost is low, with significant reductions as p is increased.

Finally, note that the optimizations detailed earlier for updating the fingerprint store in Max-Match result in low cost for the Admin function in all the algorithms. Since matching and fingerprinting are interleaved [20], the cost of fingerprinting and matching functions, and hence total processing speed, depend on the redundancy of a particular trace. We next compute the *average* processing speed for the different algorithms over a large set of traces.

Trace analysis: Figure 6 plots the average processing speed in Gbps at the server for Max-Match for different

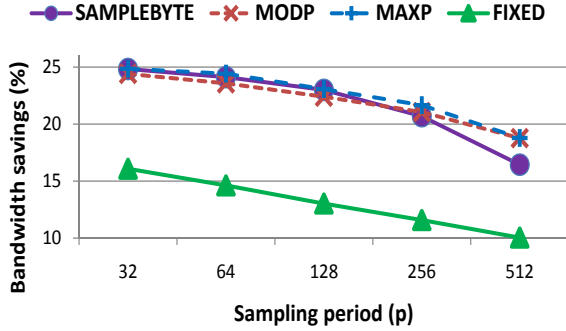


Figure 7: Max-Match bandwidth savings

fingerprinting algorithms, while Figure 7 plots the average bandwidth savings. We assume a packet cache size of 10MB. We use the 11-day traces for sites 1-10 in Table 2.

We make a number of observations from these figures. First, the processing speed of MODP is about 0.4Gbps and, as discussed in §5, is largely unaffected by p . Processing speed for MAXP ranges from 0.6 – 1.7Gbps, indicating that the CPU overhead can be decreased by increasing p . As expected, FIXED delivers the highest processing speed, ranging from 2.3 – 7.1Gbps since it incurs no cost for marker identification. Finally, SAMPLEBYTE delivers performance close to FIXED, ranging from 2.2 – 5.8Gbps, indicating that the cost of identification based on a single byte is low. Second, examining the compression savings, the curves for MODP, MAXP, and SAMPLEBYTE in Figure 7 closely overlap for the most part with SAMPLEBYTE under-performing the other two only when the sampling period is high (at $p = 512$, it appears that the choice of markers based on a single-byte may start to lose effectiveness). On the other hand, FIXED significantly under-performs the other three algorithms in terms of compression savings, though in absolute terms, the saving from FIXED are surprisingly high.

While the above results were based on a cache size of 10MB, typical for EndRE, a server is likely to have multiple simultaneous such connections in operation. Thus, in practice, it is unlikely to benefit from having beneficial CPU cache effects that the numbers above portray. We thus conducted experiments with large cache sizes (1-2GB) and found that processing speed indeed falls by about 30% for the algorithms. Taking this overhead into account, SAMPLEBYTE provides server processing speeds of 1.5 – 4Gbps. To summarize, *SAMPLEBYTE provides just enough randomization for identification of chunk markers that allows it to deliver the compression savings of MODP/MAXP while being inexpensive enough to deliver processing performance, similar to FIXED, of 1.5 – 4Gbps.*

In the case of Chunk-Match, the processing speed (not

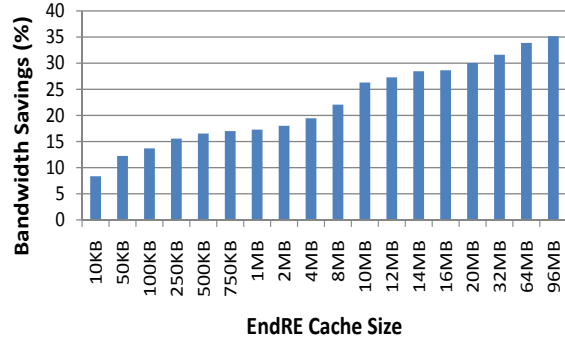


Figure 8: Cache size vs overall bandwidth savings

shown) is only 0.1-0.2Gbps. This is mainly due to SHA1 hash computation (§5.2.1) and the inability to use the fingerprint store optimizations of Max-Match (§5.2.2). We are examining if a cheaper hash function coupled with an additional mechanism to detect collision and recover payload through retransmissions will improve performance without impacting latency.

Client Decompression: The processing cost for decompression at the end host *client* is negligible since EndRE decoding is primarily a memory lookup in the client’s cache; our decompression speed is 10Gbps. We examine the impact of this in greater detail when we evaluate end-system energy savings from EndRE in §9.3.

8.2 Memory Costs

Since EndRE requires a cache per communicating client-server pair, quantifying the memory costs at both clients and servers is critical to estimating the scalability of the EndRE system. In the next two sections, we answer the following two key questions: 1) what cache size limit do we provision for the EndRE service between a single client-server pair? 2) Given the cache size limit for one pair, what is the cumulative memory requirement at clients and servers?

8.2.1 Cache Size versus Savings

To estimate the cache size requirements of EndRE, we first need to understand the trade-off between cache sizes and bandwidth savings. For the following discussion, unless otherwise stated, by cache size, we refer to the client cache size limit for EndRE service with a given server. The server cache size can be estimated from this value depending on whether Max-Match or Chunk-Match is used (§5). Further, while one could provision different cache size limits for each client-server pair, for administrative simplicity, we assume that cache size limits are identical for all EndRE nodes.

Figure 8 presents the overall bandwidth savings versus cache size for the EndRE service using the Max-Match approach (averaged across all enterprise links). Although not shown, the trends are similar for the

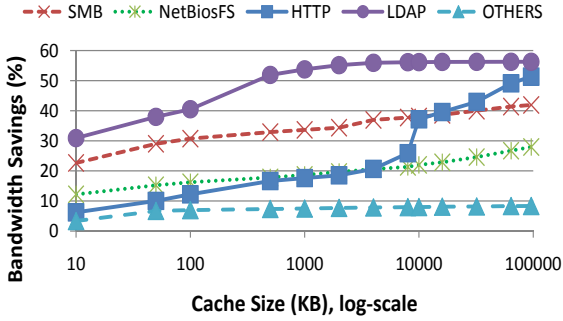


Figure 9: Cache size vs protocol bandwidth savings

Chunk-Match approach. Based on the figure, a good operating point for EndRE is at the knee of the curve corresponding to 810MB of cache, allowing for a good trade-off between memory resource constraints and bandwidth savings.

Figure 9 plots the bandwidth savings versus cache size (in log-scale for clarity) for different protocols. For this trace set, HTTP (port 80,8080) comprised 45% of all traffic, SMB (port 445) and NetBios File sharing (port 139) together comprised 26%, LDAP (port 389) was about 2.5% and a large set of protocols, labeled as OTHERS, comprised 26.5%. While different protocols see different bandwidth savings, all protocols, except OTHERS, see savings of 20+% with LDAP seeing the highest savings of 56%. Note that OTHERS include several protocols that were encrypted (HTTPS:443, Remote Desktop:3389, SIP over SSL:5061, etc.). For this analysis, since we are estimating EndRE savings from IP-level packet traces whose payload is already encrypted, EndRE sees 0% savings. An implementation of EndRE in the socket layer would likely provide higher savings for protocols in the OTHERS category than estimated here. Finally, by examining the figure, one can see the “knee-of-the-curve” at different values of cache size for different protocols (10MB for HTTP, 4MB for SMB, 500KB for LDAP, etc.). This also confirms that the 10MB knee of Figure 8 is largely due to the 10MB knee for HTTP in Figure 9.

This analysis suggests that the cache limit could be tuned depending on the protocol(s) used between a client-server pair without significantly impacting overall bandwidth savings. Thus, we use 10MB cache size only if HTTP traffic exists between a client-server pair; 4MB if SMB traffic exists, and a default 1MB cache size otherwise. Finally, while this cache size limit is derived based on static analysis of the traces, we are looking at designing dynamic cache size adaptation algorithms for each client-server pair as part of future work.

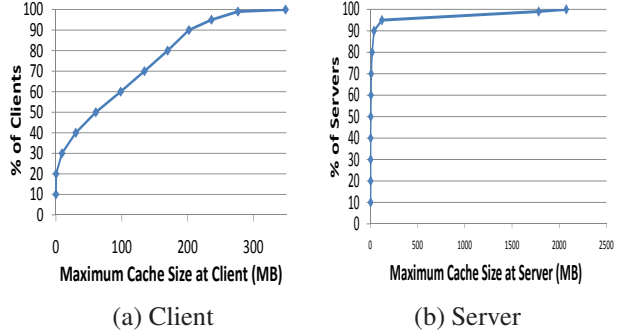


Figure 10: Cache scalability

8.2.2 Client and Server Memory Costs

Given the cache size limits derived in the previous section, we now address the critical question of EndRE scalability based on the *cumulative* cache needs at the client and server for all their connections. Using the entire set of network traces of ten enterprise sites (44 days, 5TB) in Table 2, we emulate the memory needs of EndRE with the above cache size limits for all clients and servers. We use a conservative memory page-out policy in the emulation: if there has been no traffic for over ten hours between a client-server pair, we assume that the respective EndRE caches at the nodes are paged to disk. For each node, we then compute the maximum in-memory cache needed for EndRE over the entire 44 days.

Figure 10(a) plots the CDF of the client’s maximum EndRE memory needs for all (≈ 1000) clients. We find that *the median (99 percentile) EndRE client allocates a maximum cache size of 60MB (275MB) during its operation over the entire 44-day period*. We also performed an independent study of desktop memory availability by monitoring memory availability at 1 minute intervals for 110 desktops over 1 month at one of the enterprise sites. Analyzing this data, we found that the 5, 50 and 90th percentile values of unused memory, available for use, at these enterprise desktops were 1994MB, 873MB, and 245MB, respectively. This validates our hypothesis that desktop memory resources are typically adequately provisioned in enterprises, allowing EndRE to operate on clients without significant memory installation costs.

We now examine the size of the cache needed at the server. First, we focus on Max-Match and study the net size of the cache required across all active clients at the server. Using the same enterprise trace as above, we plot the CDF of server cache size for all the servers in the trace in Figure 10(b). From the figure, we find that the maximum cache requirement is about 2GB. If it is not feasible to add extra memory to servers, say due to cost or slot limitations, the Chunk-Match approach could be adopted instead. This would reduce the maximum cache requirement by 3X (§5).

| Site | Trace Size GB | GZIP 10ms | EndRE Max-Match 10MB | | | | EndRE Max-Match+GZIP 10MB | EndRE Chunk-Match 10MB | EndRE Max-Match + DOT 10MB | IP WAN-Opt Max-Match 2GB | IP WAN-Opt Max-Match + DOT 2GB |
|----------|---------------|-----------|----------------------|------|-------|----|---------------------------|------------------------|----------------------------|--------------------------|--------------------------------|
| | | | MODP | MAXP | FIXED | SB | SB | MODP | SB | SB | SB |
| 1 | 173 | 9 | 47 | 47 | 16 | 47 | 48 | 46 | 56 | 71 | 72 |
| 2 | 8 | 14 | 24 | 25 | 19 | 24 | 28 | 19 | 33 | 33 | 33 |
| 3 | 71 | 17 | 25 | 26 | 23 | 26 | 29 | 22 | 32 | 34 | 35 |
| 4 | 58 | 17 | 23 | 24 | 20 | 24 | 31 | 21 | 30 | 45 | 47 |
| 5 | 69 | 15 | 26 | 27 | 22 | 27 | 31 | 21 | 37 | 39 | 42 |
| 6 | 80 | 12 | 21 | 21 | 18 | 22 | 26 | 17 | 28 | 34 | 36 |
| 7 | 80 | 14 | 25 | 25 | 22 | 26 | 30 | 21 | 33 | 31 | 33 |
| 8 | 142 | 14 | 22 | 23 | 18 | 22 | 28 | 19 | 30 | 34 | 40 |
| 9 | 198 | 9 | 16 | 16 | 14 | 16 | 19 | 15 | 26 | 44 | 46 |
| 10 | 117 | 13 | 20 | 21 | 17 | 21 | 25 | 17 | 30 | 27 | 30 |
| Avg/site | 100 | 13 | 25 | 26 | 19 | 26 | 30 | 22 | 34 | 39 | 41 |

Table 4: Percentage bandwidth savings on incoming links to 10 enterprise sites over 11 day trace

9 Benefits

We now characterize various benefits of EndRE. We first investigate WAN bandwidth savings. We then quantify latency savings of using EndRE, especially on short transactions typical of HTTP. Finally, we quantify energy savings on mobile smartphones, contrasting EndRE with prior work on energy-aware compression [11].

9.1 Bandwidth Savings

In this section, we focus on the bandwidth savings of different RE algorithms for each of the enterprise sites, and examine the gains of augmenting EndRE with GZIP and DOT [21]. We also present bandwidth savings of an IP WAN optimizer for reference.

Table 4 compares the bandwidth savings on incoming links to ten enterprise sites for various approaches. This analysis is based on packet-level traces and while operating at packet sizes or larger buffers make little difference to the benefits of EndRE approaches, buffer size can have a significant impact on GZIP-style compression. Thus, in order to emulate the benefits of performing GZIP at the socket layer, we aggregate consecutive packet payloads for up to 10ms and use this aggregated buffer while evaluating the benefits of GZIP. For EndRE, we use cache sizes of up to 10MB. We also emulate an IP-layer middlebox-based WAN optimizer with a 2GB cache.

We observe the following: First, performing GZIP in isolation on packets aggregated for up to 10ms provides per-site savings of 13% on average. Further, there are site specific variations; in particular, GZIP performs poorly for site 1 compared to other approaches. Second, comparing the four fingerprinting algorithms (columns 3-6 in Table 4), we see that MODP, MAXP, and SAMPLEBYTE deliver similar average savings of 25-26% while FIXED under-performs. In particular, in the case of site 1, FIXED significantly under-performs the other three approaches. This again illustrates how SAMPLEBYTE captures enough content-specific characteristics to significantly outperform FIXED. Adding GZIP com-

pression to SAMPLEBYTE improves the average savings to 30% (column 7). While the above numbers were based on Max-Match, using Chunk-Match instead reduces the savings to 22% (column 8), but this may be a reasonable alternative if server memory is a bottleneck.

We then examine savings when EndRE is augmented with DOT [21]. For this analysis, we employ a heuristic to extract object chunks from the packet traces as follows: we combine consecutive packets of the same four-tuple flow and delineate object boundaries if there is no packet within a time window(1s). In order to ensure that the DOT analysis adds redundancy not seen by EndRE, we conservatively add only inter-host redundancy obtained by DOT to the EndRE savings. We see that (third column from right) DOT improves EndRE savings by a further 6-10%, and the per-site average bandwidth savings improves to 34%. For reference, a WAN optimizer with 2GB cache provides per-site savings of 39% and if DOT is additionally applied (where redundancy of matches farther away than 2GB are only counted), the average savings goes up by only 2%. Thus, it appears that half the gap between EndRE and WAN optimizer savings comes from inter-host redundancy and the other half from the larger cache used by the WAN optimizer.

Summarizing, EndRE using the Max-Match approach with the SAMPLEBYTE algorithm provides average per-site savings of 26% and delivers two-thirds of the savings of a IP-layer WAN optimizer. When DOT is applied in conjunction, the average savings of EndRE increase to 34% and can be seen to be approaching the 41% savings of the WAN optimizer with DOT.

Pilot Deployment: We now report results from a small scale deployment. EndRE was deployed on 15 desktop/laptop clients (11 users) and one server for a period of about 1 week (09/25/09 to 10/02/09) in our lab. We also hosted a HTTP proxy at the EndRE server and users manually enabled/disabled the use of this proxy, at any given time, using a client-based software. During this period, a total of 1.7GB of HTTP traffic was delivered through the EndRE service with an average compression

| RTTs | 1 | 2 | 3 | 4 | 5 | > 5 |
|--------------|---|-----|-----|-----|-----|-----|
| Latency Gain | 0 | 20% | 23% | 36% | 20% | 22% |

Table 5: HTTP latency gain for different RTTs

of 31.2%. A total of 159K TCP connections were serviced with 72 peak active simultaneous TCP connections and peak throughput of 18.4Mbps (WAN link was the bottleneck). The CPU utilization at the server remained within 10% including proxy processing. The number of packet re-orderings was less than 1% even in the presence of multiple simultaneous TCP connections between client and server. We also saw a large number of TCP RSTs but, as reported in [10], these were mostly in lieu of TCP FINs and thus do not contribute to cache synchronization issues. Summarizing, even though this is a small deployment, the overall savings numbers match well with our analysis results and the ease of deployment validates the choice of implementing EndRE over TCP.

9.2 Latency Savings

In this section, we evaluate the latency gains from deploying EndRE. In general, latency gains are possible for a number of reasons. The obvious case is due to reduction of load on the bottleneck WAN access link of an enterprise. Latency gains may also arise from the choice of implementing EndRE at the socket layer above TCP. Performing RE above the TCP layer helps reduce the amount of data transferred and thus the number of TCP round-trips necessary for connection completion. In the case of large file transfers, since TCP would mostly be operating in the steady-state congestion avoidance phase, the percentage reduction in data transfer size translates into a commensurate reduction in file download latency. Thus, for large file transfers, say, using SMB or HTTP, one would expect latency gains similar to the average bandwidth gains seen earlier.

Latency gains in the case of short data transfers, typical of HTTP, are harder to estimate. This is because TCP would mostly be operating in slow-start phase and a given reduction in data transfer size could translate into a reduction of zero or more round trips depending on many factors including original data size and whether or not the reduction occurs uniformly over the data.

In order to quantify latency gains for short file transfers, we perform the following experiment. From the enterprise network traces, we extract HTTP traffic that we then categorize into a series of session files. Each session file consists of a set of timestamped operations starting with a connect, followed by a series of sends and receives (i.e., transactions), and finally a close.

The session files are then replayed on a testbed consisting of a client and a server connected by a PC-based router emulating a high bandwidth, long latency link, using the mechanism described in [13]. During the replay,

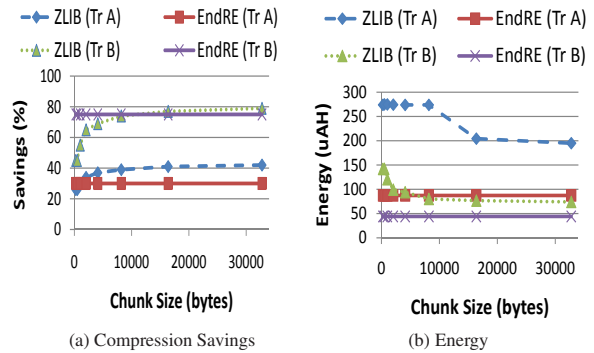


Figure 11: ZLIB vs. EndRE

strict timing is enforced at the start of each session based on the original trace; in the case of transactions, timing between the start of one transaction and the start of the next transaction is preserved as far as possible. The performance metric of interest is latency gain which is defined as the ratio of reduction in transaction time due to EndRE to transaction time without EndRE.

Table 5 shows the latency gain for HTTP for various transactions sorted by the number of round-trips in the original trace. For this trace, only 40% of HTTP transactions involved more than one round trip. For these transactions, latency gains on average ranged from 20% to 35%. These gains are comparable with the average bandwidth savings due to EndRE for this trace (~30%), demonstrating that even short HTTP transactions see latency benefits due to RE.

9.3 Energy Savings

We study the energy and bandwidth savings achieved using EndRE on Windows Mobile smartphones and compare it against both no compression as well as prior work on energy-aware compression [11]. In [11], the authors evaluate different compression algorithms and show that ZLIB performs best in terms of energy savings on resource constrained devices for decompression. We evaluate the energy and bandwidth gains using two trace files. Traces A and B are 20MB and 15MB in size, respectively, and are based on enterprise HTTP traffic, with trace B being more compressible than trace A.

We first micro-benchmark the computational cost of decompression for ZLIB and EndRE. To do this, we load pre-compressed chunks of the traces in the mobile smartphone’s memory and turn off WiFi. We then repeatedly decompress these chunks and quantify the energy cost. Figures 11(a) and (b) plot the average compression savings and energy cost of in-memory decompression for various chunk sizes, respectively. The energy measurements are obtained using a hardware-based battery power monitor [4] that is accurate to within 1mA. From these figures, we make two observations. First, as the chunk size is increased, ZLIB compression sav-

| | None | ZLIB | | | | EndRE | |
|---|------|--------|-----------|------|------|-----------|-----|
| | | Energy | % savings | | Byte | % savings | |
| | | uAh | pkt | 32KB | pkt | 32KB | pkt |
| A | 2038 | -11 | 42 | 26 | 44 | 25 | 29 |
| B | 1496 | -11 | 68 | 41 | 75 | 70 | 76 |

Table 6: Energy savings on a mobile smartphone

ings increase and the energy cost of decompression decreases. This implies that ZLIB is energy efficient when compressing large chunks/files. Second, the compression savings and energy costs of EndRE, as expected, are independent of chunk size. More importantly, EndRE delivers comparable compression savings as ZLIB while incurring an energy cost of 30-60% of ZLIB.

We now compare the performance of ZLIB and EndRE to the case of no compression by replaying the traces over WiFi to the mobile smartphone and performing in-memory decompression on the phone. In the case of ZLIB, we consider two cases: packet-by-packet compression and bulk compression where 32KB blocks of data are compressed at a time, the latter representing a bulk download case. After decompression, each packet is consumed in memory and not written to disk; this allows us to isolate the energy cost of communication. If the decompressed packet is written to disk or further computation is performed on the packet, the total energy consumed for all the scenarios will be correspondingly higher.

Table 6 shows energy and compression gains of using ZLIB and EndRE as compared to using no compression. We see that when ZLIB is applied on a packet-by-packet basis, even though it saves bandwidth, it results in increased energy consumption (negative energy savings). This is due to the computational overhead of ZLIB decompression. On the other hand, for larger chunk sizes, the higher compression savings coupled with lower computational overhead (Figure 11) result in good energy savings for ZLIB. In the case of EndRE, we find that the bandwidth savings directly translate into comparable energy savings for communication. This suggests that EndRE is a more energy-efficient solution for packet-by-packet compression while ZLIB, or EndRE coupled with ZLIB, work well for bulk compression.

10 Conclusion

Using extensive traces of enterprise network traffic and testbed experiments, we show that our end-host based redundancy elimination service, EndRE, provides average bandwidth gains of 26% and, in conjunction with DOT, the savings approach that provided by a WAN optimizer. Further, EndRE achieves speeds of 1.5-4Gbps, provides latency savings of up to 30% and translates bandwidth savings into comparable energy savings on mobile smartphones. In order to achieve these benefits,

EndRE utilizes memory and CPU resources of end systems. For enterprise clients, we show that median memory requirements for EndRE is only 60MB. At the server end, we design mechanisms for working with reduced memory and adapting to CPU load.

Thus, we have shown that the cleaner semantics of end-to-end redundancy removal can come with considerable performance benefits and low additional costs. This makes EndRE a compelling alternative to middlebox-based approaches.

Acknowledgments. We thank our shepherd Sylvia Ratnasamy and the anonymous reviewers for their comments. Aditya Akella, Ashok Anand, and Chitra Muthukrishnan were supported in part by NSF grants CNS-0626889, CNS-0746531 and CNS-0905134, and by grants from the UW-Madison Graduate School.

References

- [1] Cisco Wide Area Application Acceleration Services. http://www.cisco.com/en/US/products/ps5680/Products_Sub_Category_Home.html.
- [2] Jenkins Hash. <http://burtleburtle.net/bob/c/lookup3.c>.
- [3] Peribit Networks (Acquired by Juniper in 2005): WAN Optimization Solution. <http://www.juniper.net/>.
- [4] Power Monitor, Monsoon Solutions. <http://www.msoon.com/powermonitor/powermonitor.html>.
- [5] Riverbed Networks: WAN Optimization. <http://www.riverbed.com/solutions/optimize/>.
- [6] Windows Filtering Platform. [http://msdn.microsoft.com/en-us/library/aa366509\(V.85\).aspx](http://msdn.microsoft.com/en-us/library/aa366509(V.85).aspx).
- [7] A. Anand, A. Gupta, A. Akella, S. Seshan, and S. Shenker. Packet Caches on Routers: The Implications of Universal Redundant Traffic Elimination. In *ACM SIGCOMM*, Seattle, WA, Aug. 2008.
- [8] A. Anand, C. Muthukrishnan, A. Akella, and R. Ramjee. Redundant in Network Traffic: Findings and Implications. In *ACM SIGMETRICS*, Seattle, WA, June 2009.
- [9] S. Annapureddy, M. J. Freedman, and D. Mazires. Shark: Scaling file servers via cooperative caching. In *NSDI*, 2005.
- [10] M. Arlitt and C. Williamson. An analysis of tcp reset behavior on the internet. *ACM CCR*, 35(1), 2005.
- [11] K. C. Barr and K. Asanovic. Energy-aware lossless data compression. *IEEE Transactions on Computer Systems*, 24(3):250–291, Aug 2006.
- [12] F. Dougliis and A. Iyengar. Application-specific delta-encoding via resemblance detection. In *USENIX*, 2003.
- [13] J. Eriksson, S. Agarwal, P. Bahl, and J. Padhye. Feasibility study of mesh networks for all-wireless offices. In *MobiSys*, 2006.
- [14] J. C. Mogul, F. Dougliis, A. Feldmann, and B. Krishnamurthy. Potential benefits of delta encoding and data compression for http. In *SIGCOMM*, pages 181–194, 1997.
- [15] A. Muthitacharoen, B. Chen, and D. Mazières. A low-bandwidth network file system. *SIGOPS Oper. Syst. Rev.*, 35(5), 2001.
- [16] R. Pang, M. Allman, M. Bennett, J. Lee, V. Paxson, and B. Tierney. A first look at modern enterprise traffic. In *IMC*, 2005.
- [17] H. Pucha, D. G. Andersen, and M. Kaminsky. Exploiting similarity for multi-source downloads using file handprints. In *NSDI*, 2007.
- [18] M. Rabin. Fingerprinting by random polynomials. Technical report, Harvard University, 1981. Technical Report, TR-15-81.
- [19] S. C. Rhea, K. Liang, and E. Brewer. Value-Based Web Caching. In *12th World Wide Web Conference*, 2003.
- [20] N. T. Spring and D. Wetherall. A protocol-independent technique for eliminating redundant network traffic. In *SIGCOMM*, pages 87–95, 2000.
- [21] N. Tolia, M. Kaminsky, D. G. Andersen, and S. Patil. An architecture for Internet data transfer. In *NSDI*, 2006.