

Low-Cost Data Deduplication for Virtual Machine Backup in Cloud Storage

Wei Zhang*, Tao Yang*, Gautham Narayanasamy*, and Hong Tang[†]

* University of California at Santa Barbara, [†] Alibaba Inc.

Abstract

In a virtualized cloud cluster, frequent snapshot backup of virtual disks improves hosting reliability; however, it takes significant memory resource to detect and remove duplicated content blocks among snapshots. This paper presents a low-cost deduplication solution scalable for a large number of virtual machines. The key idea is to separate duplicate detection from the actual storage backup instead of using inline deduplication, and partition global index and detection requests among machines using fingerprint values. Then each machine conducts duplicate detection partition by partition independently with minimal memory usage. Another optimization is to allocate and control buffer space for exchanging detection requests and duplicate summaries among machines. Our evaluation shows that the proposed multi-stage scheme uses a small amount of memory while delivering a satisfactory backup throughput.

1 Introduction

Periodic archiving of virtual machine (VM) snapshots is important for long-term data retention and fault recovery. For example, daily backup of VM images is conducted automatically at Alibaba which provides the largest public cloud service in China. The cost of frequent backup of VM snapshots is high because of the huge storage demand. This issue has been addressed by storage data deduplication [8, 15] that identifies redundant content duplicates among snapshots. One architectural approach is to attach a separate backup system with deduplication support to the cloud cluster, and every machine periodically transfers snapshots to the attached backup system. Such a dedicated backup configuration can be expensive, considering that significant networking and computing resource is required to transfer raw data and conduct signature comparison.

This paper seeks for a low-cost architecture option and considers that a backup service uses the existing cloud computing resource. Performing deduplication adds significant memory cost for comparison of content fingerprints. Since each physical machine in a cluster hosts many VMs, memory contention happens frequently. Cloud providers often wish that the backup service only consumes small or modest resources with a minimal impact to the existing cloud services. Another challenge is that deletion of old snapshots compete for

computing resource as well, because data dependence created by duplicate relationship among snapshots adds processing complexity.

Among the three factors - time, cost and deduplication efficiency, one of them has to be compromised for the other two. For instance, if we were building a deduplication system that has a high rate of duplication detection and has a very fast response time, it would need a lot of memory to hold fingerprint index and cache. This leads to a compromise on cost. Our objective is to lower the cost incurred while sustaining the highest de-duplication ratio and a sufficient throughput in dealing with a large number of VM images.

The traditional approach to deduplication is an inline approach which follows a sequence of block reading, duplicate detection, and non-duplicate block write to the backup storage. Our key idea is to first perform parallel duplicate detection for VM content blocks among all machines before performing actual data backup. Each machine accumulates detection requests and then performs detection partition by partition with minimal resource usage. Fingerprint based partitioning allows highly parallel duplicate detection and also simplifies reference counting management. The tradeoff is that every machine has to read dirty segments twice and that some deduplication requests are delayed for staged parallel processing. With careful parallelism and buffer management, this multi-stage detection scheme can provide a sufficient throughput for VM backup.

2 Background and Related Work

At a cloud cluster node, each instance of a guest operating system runs on a virtual machine, accessing virtual hard disks represented as virtual disk image files in the host operating system. For VM snapshot backup, file-level semantics are normally not provided. Snapshot operations take place at the virtual device driver level, which means no fine-grained file system metadata can be used to determine the changed data.

Backup systems have been developed to use content fingerprints to identify duplicate content [8, 9]. Offline deduplication is used in [5, 2] to remove previously written duplicate blocks during idle time. Several techniques have been proposed to speedup searching of duplicate fingerprints. For example, the data domain method [15] uses an in-memory Bloom filter and

a prefetching cache for data blocks which may be accessed. An improvement to this work with parallelization is in [12, 13]. As discussed in Section 1, there is no dedicated resource for deduplication in our targeted setting and low memory usage is required so that the resource impact to other cloud services is minimized. The approximation techniques are studied in [3, 6, 14] to reduce memory requirement with a tradeoff of the reduced deduplication ratio. In comparison, this paper focuses on full deduplication without approximation.

Additional inline deduplication techniques are studied in [7, 6, 10]. All of the above approaches have focused on such inline duplicate detection in which deduplication of an individual block is on the critical write path. In our work, this constraint is relaxed and there is a waiting time for many duplicate detection requests. This relaxation is acceptable because in our context, finishing the backup of required VM images within a reasonable time window is more important than optimizing individual VM block backup requests.

3 System Design

We consider deduplication in two levels. The first level uses coarse-grain segment dirty bits for version-based detection [4, 11]. Our experiment with Alibaba’s production dataset shows that over 70 percentage of duplicates can be detected using segment dirty bits when the segment size is 2M bytes. This setting requires OS to maintain segment dirty bits and the amount of space for this purpose is negligible. In the second level of deduplication, content blocks of dirty segments are compared with the fingerprints of unique blocks from the previous snapshots. Our key strategies are explained as follows.

- **Separation of duplicate detection and data backup.** The second level detection requires a global comparison of fingerprints. Our approach is to perform duplicate detection first before actual data backup. That requires a prescanning of dirty VM segments, which does incur an extra round of VM reading. During VM prescanning, detection requests are accumulated. Aggregated deduplicate requests can be processed partition by partition. Since each partition corresponds to a small portion of global index, memory cost to process detection requests within a partition is small.
- **Buffered data redistribution in parallel duplicate detection.** Let *global index* be the meta data containing the fingerprint values of unique snapshot blocks in all VMs and the reference pointers to the location of raw data. A logical way to distribute detection requests among machines is based on fingerprint values of content blocks. Initial data blocks follows the VM distribution among machines and the detected duplicate sum-

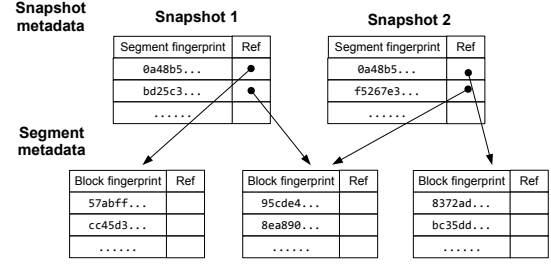


Figure 1: Metadata structure of a VM snapshot.

mary should be collected following the same distribution. Therefore, there are two all-to-all data redistribution operations involved. One is to map detection requests from VM-based distribution to fingerprint based distribution. Another one is to map duplicate summary from fingerprint-based distribution to VM based distribution. The redistributed data needs to be accumulated on the disk to reduce the use of memory. To minimize the disk seek cost, outgoing or incoming data exchange messages are buffered to bundle small messages. Given there are $p \times q$ partitions where p is the number of machines and q is the number of fingerprint-based partitions at each machine, space per each buffer is small under the memory constraint for large p or q values. This counteracts the effort of seek cost reduction. We have designed an efficient data exchange and disk data buffering scheme to address this.

We assume a flat architecture in which all p machines that host VMs in a cluster can be used in parallel for deduplication. A small amount of local disk space and memory on each machine can be used to store global index and temporary data. The real backup storage can be either a distributed file system built on this cluster or use another external storage system.

The representation of each snapshot in the backup storage has a two-level index structure in the form of a hierarchical directed acyclic graph as shown in Figure 1. A VM image is divided into a set of segments and each segment contains content blocks of variable-size, partitioned using the standard chunking technique with 4KB as the average block size. The snapshot metadata contains a list of segments and other meta data information. Segment metadata contains its content block fingerprints and reference pointers. If a segment is not changed from one snapshot to another, indicated by a dirty bit embedded in the virtual disk driver, its segment metadata contains a reference pointer to an earlier segment. For a dirty segment, if one of its blocks is duplicate to another block in the system, the block metadata contains a reference pointer to the earlier block.



Figure 2: Processing flow of Stage 1 (dirty segment scan and request accumulation), Stage 2 (fingerprint comparison and summary output), and Stage 3 (non-duplicate block backup).

The data flow of our multi-stage duplicate detection is depicted in Figure 2. In Stage 1, each machine independently reads VM images that need a backup and forms duplicate detection requests. The system divides each dirty segment into a sequence of chunk blocks, computes the meta information such as chunk fingerprints, sends a request to a proper machine, and accumulates received requests into a partition on the local temporary disk storage. The partition mapping uses a hash function applied to the content fingerprint. Assuming all machines have a homogeneous resource configuration, each machine is evenly assigned with q partitions of global index and it accumulates corresponding requests on the disk. There are two options to allocate buffers at each machine. 1) Each machine has $p \times q$ send buffers corresponding to $p \times q$ partitions in the cluster since a content block in a VM image of this machine can be sent to any of these partitions. 2) Each machine allocates p send buffers to deliver requests to p machines; it allocates p receive buffers to collect requests from other machines. Then the system copies requests from each of p receive buffers to q local request buffers, and outputs each request buffer to one of the request partitions on the disk when this request buffer becomes full. Option 2, which is depicted in Figure 2, is much more efficient than Option 1 because $2p + q$ is much smaller than $p \times q$, except for the very small values. As a result, each buffer in Option 2 has a bigger size to accumulate requests and that means less disk seek overhead.

Stage 2 is to load disk data and perform fingerprint comparison at each machine one request partition at a time. At each iteration, once in-memory comparison between an index partition and request partition is com-

pleted, duplicate summary information for segments of each VM is routed from the fingerprint-based distribution to the VM-based distribution. The summary contains the block ID and the reference pointer for each detected duplicate block. Each machine uses memory space of the request partition as a send buffer with no extra memory requirement. But it needs to allocate p receive buffers to collect duplicate summary from other machines. It also allocates v request buffers to copy duplicate summary from p receive buffers and output to the local disk when request buffers are full.

Stage 3 is to perform real backup. The system loads the duplicate summary of a VM, reads dirty segments of a VM, and outputs non-duplicate blocks to the final backup storage. Additionally, the global index on each machine is updated with the meta data of new chunk blocks. When a segment is not dirty, the system only needs to output the segment meta data such as a reference pointer. There is an option to directly read dirty blocks instead of fetching a dirty segment which can include duplicate blocks. Our experiment shows that it is faster to read dirty segments in the tested workload. Another issue is that during global index update after new block creation, the system may find some blocks with the same fingerprints have been created redundantly. For example, two different VM blocks that have the same fingerprint are not detected because the global index has not contained such a fingerprint yet. The redundancy is discovered and logged during the index update and can be repaired periodically when necessary. Our experience is that there is a redundancy during the initial snapshot backup and once that is repaired, the percentage of redundant blocks due to concurrent processing is insignif-

icant.

The above steps can be executed by each machine using one thread to minimize the use of computing resource. The disk storage usage on each machine is fairly small for storing part of global index and accumulating duplicate detection requests that contain fingerprint information. We impose a memory limit M allocated for each stage of processing at each machine. The usage of M is controlled as follows and space allocation among buffers is optimized based on the relative ratio between the cross-machine network startup cost and disk access startup cost such as seek time. Using a bigger buffer can mitigate the impact of slower startup cost.

- For Stage 1, M is divided for 1) an I/O buffer to read dirty segments; 2) $2p$ send/receive buffers and q request buffers.
- For Stage 2, M is divided for 1) space for hosting a global index partition and the corresponding request partition; 2) p receive buffers and v summary buffers.
- For Stage 3, M is divided for 1) an I/O buffer to read dirty segments of a VM and write non-duplicate blocks to the backup storage; 2) summary of duplicate blocks within dirty segments.

Snapshot deletion. Each VM will keep a limited number of automatically-saved snapshots and expired snapshots are normally deleted. We adopt the idea of mark-and-sweep [6]. A block or a segment can be deleted if its reference count is zero. To delete useless blocks or segments periodically, we read the meta data of all snapshots and compute the reference count of all blocks and segments in parallel. Similar to the multi-stage duplicate detection process, reference counting is conducted in multi-stages. Stage 1 is to read the segment and block metadata to accumulate reference count requests in different machines in the fingerprint based distribution. Stage 2 is to count references within each partition and detect those records with zero reference. The backup data repository logs deletion instructions, and will periodically perform a compaction operation when its deletion log is too big.

4 Old Performance Analysis and Comparison

We assume a flat architecture that we use all machines in a cluster to host virtual machines, and also evenly host raw data and meta data of the temporarily accumulated requests. We call global index to be the meta data of all non-duplicate chunks such as chunk fingerprints and reference pointers.

Following parameters are used to analyze the performance of our system.

- p is the number of machines in a cluster. These

machines can run in parallel for backup. The request buckets are evenly distributed among these machines.

- v is the number of virtual machines per machine. At Alibaba, $v = 25$.
- x is the number of snapshots saved for each VM.
- k is the number of iterations to complete all virtual machine backup. Each iteration performs v/k backups.
- t is the amount of temporary disk space used per physical machine for deduplication.
- m is the amount of memory used per each physical machine for deduplication. Our goal is to minimize
- s is the average size of virtual machine image. At Alibaba data we have tested, $s = 40GB$.
- d_1 is the average deduplication ratio using segment-based dirtbit. $s*d_1$ represents the amount of data items that are duplicates and can be avoided for backup. For Alabalba dataset tested, $d_1=77\%$.
- d_2 is the average deduplication ratio using content chunk fingerprints after segment-based deduplication. For Alaba dataset tested, $d_2 = 50\%$.
- b is the average disk bandwidth for reading from local storage at each machine.
- q is the number of buckets to accumulate requests at each machine. Thus the total number of buckets is $p * q$.
- c is the chunk block size in bytes. In practice $c = 4KB$.
- u is the record size of detection request per block. In practice,
- $u=40$. That includes block ID and fingerprint.
- m is the maximum memory allocated for deduplication purpose. A g fraction used for machine-machine network request buffering and $(1-g)$ fraction used for memory-disk bucket buffering.
- e is the size of a duplicate summary record for each chunk block.
- α_n is the startup cost for sending a message in a cluster. α_d is the startup cost such as seek for disk IO. β is the time cost for in-memory duplicate comparison.

The system keeps at most x copies of snapshots for each VM on average. The total size of global content fingerprints is $x * s * v / c * u * (1 - d_1) * (1 - d_2)$ where c is the average chunk size and u is the meta data size of each chunk fingerprint. In practice $c = 4K$ and u/c is about 100. $x = 10$ in the case of Alibaba cloud.

Define $r = sv(1 - d_1)/(ck)$ which is the total number of duplicate detection requests issued at each machine and at each iteration.

We first discuss the memory usage and processing time of 3 steps. For Step 1, the buffer for sending requests from one machine to another has a size of $g*m/p$,

and with such a buffering, the total number of outgoing communication messages from each machine to other machines can be

$$rup/(g * m)$$

The total amount of data communicated among machines is relatively small: rup in the cluster, distributed among p machines.

Once every machine receives detection requests and divide them into buckets, it writes the content to the disk once the buffer is full. The buffer for each bucket is $(1 - g)m/q$ and the total number of disk write requests issued after the bucket buffer is full is:

$$ruq/((1 - g) * m)$$

The total time for step 1 which reads VM images and write accumulated detection requests is:

$$r(c + u)/b + ru/m(\alpha_n p/g + \alpha_d q/(1 - g))$$

For Step 2, part of memory at each machine is to hold a bucket of global index and accumulated requests. That is

$$m_b = x * r * u * k(1 - d_2)/q + r * u/q$$

Thus the memory requirement for this portion can be made very small when setting a large q . On the other hand, as the system detects duplicates per hash bucket, we need to allocate buffer space for receiving duplicate summary for each VM. The total buffer size is $m - m_b$ which is used evenly for v VMs.

The size of the duplicate summary for each bucket is

$$S_{sum} = sv(1 - d_1)e/(kcq)$$

We can buffer the outcome of multiple buckets. The total buffer factor is

$$(m - m_b)/S_{sum}.$$

The final bucket buffer for each VM is still fairly small, and writing such a buffer to the disk may involve two I/O requests (one to fetch the old block, and one is to update). The total seek cost involved

$$2 * v * \alpha_d * q/((m - m_b)/S_{sum}) = 2vre\alpha_d/(m - m_b)$$

Thus the total time of Step 2 takes

$$(x * r * k * u * (1 - d_2) + r * u)/b_d + r * \beta + 2v * r * e\alpha_d/(m - m_b).$$

The key cost of step 3 is to read the nonduplicate parts of each VM and output the backend storage. The time of Step 3 takes:

$$2r * c * (1 - d_2)/b_d$$

That assumes that when a content chunk is not a duplicate, there is a significant number of non-duplicate chunks following that chunk.

Thus the total time to process all v virtual machines after k iterations are:

$$\begin{aligned} & k[r(c + u)/b + ru/m(\alpha_n p/g + \alpha_d q/(1 - g)) \\ & + (x * r * k * u * (1 - d_2) + r * u)/b_d + r * \beta + \\ & 2v * r * e\alpha_d/(m - m_b) + 2r * c * (1 - d_2)/b_d] \end{aligned}$$

subject to conditions that

$$m - m_b > 0$$

The total disk requirement per machine for hosting the global index and meta data of accumulated requests is:

$$x * r * k * u * (1 - d_2) + r * u.$$

That is not so big, and is acceptable as we show later.

5 New Performance Analysis and Comparison

The system keeps at most x copies of snapshots for each VM on average. The total size of global content fingerprints is $x * s * v/c * u * (1 - d_1) * (1 - d_2)$ where c is the average chunk size and u is the meta data size of each chunk fingerprint. In practice $c = 4K$ and u/c is about 100. $x = 10$ in the case of Alibaba cloud.

Define $r = sv(1 - d_1)/(ck)$ which is the total number of duplicate detection requests issued at each machine and at each iteration.

We first analyze the time cost of 5 stages assuming an evenly distributed load across the machines. Later we analyze the additional costs associated with an imbalanced load.

In Stage 1, the dirty segments are read from the virtual disk, the hash of each block is computed, and dedup requests are sent to the machine hosting the blocks' respective partitions. Since we must read r blocks from disk, send r dedup requests, and then save the requests to temporary files (one for each partition), the time for the first stage can be expressed as:

$$b_r r c + \alpha_n \frac{ur}{m_n} + b_w r u$$

In Stage 2, each partition index is read from disk, then the dedup requests for that partition are processed and the results are written back out to disk to be sent in Stage 3. The results are broken into 3 groups: duplicate blocks, new blocks, and dup-with-new blocks, which are duplicates of blocks that are new to this batch.

Let n be the total number of index entries in the system. $n = (pvx)(1 - d_1)(1 - d_2)\frac{e}{c}$, which represents the

index cost for all the deduped data in the system. Since each machine holds a constant number q partitions, and the partitions should be uniform in size as they are from the hash of the block, the total number of index entries at each machine should be n/p

The cost of Stage 2 is:

$$b_r r u + b_r \frac{ne}{p} + r\beta + b_w r e$$

In Stage 3 the new block results from Stage 2 are sent to the requesters, and in Stage 3b the new blocks are written out to the storage system. We will now mostly be dealing with the $r(1 - d_2)$ blocks that are new to the system. In 3b the dedup results must be read and the actual disk blocks for each new block must be re-read before they can be sent to the block store.

The cost of Stage 3 is:

$$b_r r(1 - d_2)e + \alpha_n \frac{er(1 - d_2)}{m_n} + b_w r(1 - d_2)e$$

and the cost of Stage 3b is:

$$b_r r(1 - d_2)(e + c) + b_b r(1 - d_2)c$$

After the new blocks have been written to the block store, and references to them have been obtained, those references must be returned to the partition index holder so that those blocks may be deduped in the future. This Stage 4 (read the new index entries, return them to the partition master, and save the received references) costs:

$$b_r r(1 - d_2)e + \alpha_n \frac{er(1 - d_2)}{m_n} + b_w r(1 - d_2)e$$

Stage 4b consists of updating the partition index with the new block references from Stage 4, and costs:

$$b_r r(1 - d_2)e + b_w r(1 - d_2)e$$

In the final stage (Stage 5), dup-with-new references are returned to the requesters, so that the snapshot recipes may be updated with references to those blocks. This process costs:

$$b_r r d_3 e + \alpha_n \frac{er d_3}{m_n} + b_w r e$$

Memory Requirements:

In every stage we need 1 disk read buffer, And then additionally we need the following:

Stage 1 network buffers, and q disk write buffers

Stage 2 n/q partition index space, and p disk write buffers (for dedup responses)

Stage 3 network buffers, v disk write buffers (for dedup responses)

Stage 3b 1 disk write buffer (to write out new blocks)

Stage 4 network buffers, q disk write buffers (to write out new refs)

Stage 4b 1 disk write buffer (to update partition index with new blocks)

Stage 5 network buffers

5.1 A Comparison with Other Approaches

The memory space requirement for the data domain approach with bloom filter is:

$$x * r k u (1 - d_2) / r$$

where r is the bloom filter with about 1:10 ratio in practice. The disk space used is

$$x * r * k * u * (1 - d_2).$$

6 Evaluation

We have implemented and evaluated a prototype of our multi-stage deduplication scheme on a cluster of dual quad-core Intel Nehalem 2.4GHz E5530 machines with 24GB memory. Our implementation is based on Alibaba's Xen cloud platform [1, 14]. Objectives of our evaluation are: 1) Analyze the deduplication throughput and effectiveness for a large number of VMs. 2) Examine the impacts of buffering during metadata exchange.

We have performed a trace-driven study using a 1323 VM dataset collected from a cloud cluster at Alibaba's Aliyun. For each VM, the system keeps 10 automatically-backed snapshots in the storage while a user may instruct extra snapshots to be saved. The backup of VM snapshots is completed within a few hours every night. Based on our study of its production data, each VM has about 40GB of storage data usage on average including OS and user data disk. Each VM image is divided into 2 MB fix-sized segments and each segment is divided into variable-sized content blocks with an average size of 4KB. The signature for variable-sized blocks is computed using their SHA-1 hash.

The seek cost of each random IO request in our test machines is about 10 milliseconds. The average I/O usage of local storage is controlled about 50MB/second for backup in the presence of other I/O jobs. Noted that a typical 1U server can host 6 to 8 hard drives and deliver over 300MB/second. Our setting uses 16.7% or less of local storage bandwidth. The final snapshots are stored in a distributed file system built on the same cluster.

The total local disk usage on each machine is about 8GB for the duplicate detection purpose, mainly for global index. Level 1 segment dirty bits identify 78% of duplicate blocks. For the remaining dirty segments, block-wise full deduplication removes about additional 74.5% of duplicates. The final content copied to the backup storage is reduced by 94.4% in total.

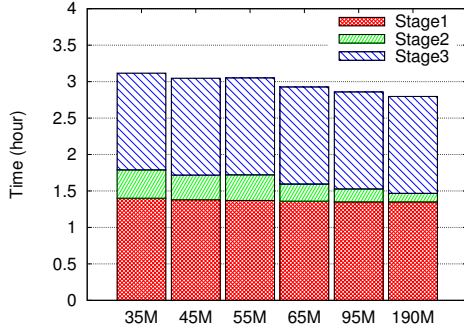


Figure 3: Parallel time when memory limit varies.

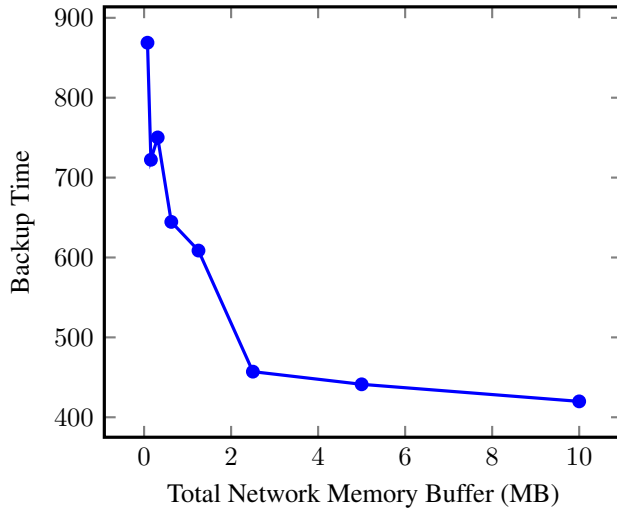


Figure 4: Backup time for varying amounts of memory allocated to network communication. Other Settings: 10 Machines, 5x40GB VMs per machine, write buffer 6.25MB, Read Buffer 128KB

Figure 3 shows the total parallel time in hours to backup 2500 VMs on a 100-node cluster when limit M imposed on each node varies. This figure also depicts the time breakdown for Stages 1, 2, and 3. The time in Stages 1 and 3 is dominated by the two scans of dirty segments, and final data copying to the backup storage is overlapped with VM scanning. During dirty segment reading, the average number of consecutive dirty segments is 2.92. The overall processing time does not have a significant reduction as M increases to 190MB. The aggregated deduplication throughput is about 8.76GB per second, which is the size of 2500 VM images divided by the parallel time. The system runs with a single thread and its CPU resource usage is 10-13% of one core. The result shows the backup with multi-stage deduplication for all VM images can be completed in about 3.1 hours with 35MB memory, 8GB disk overhead and a small CPU usage. As we vary the cluster size p ,

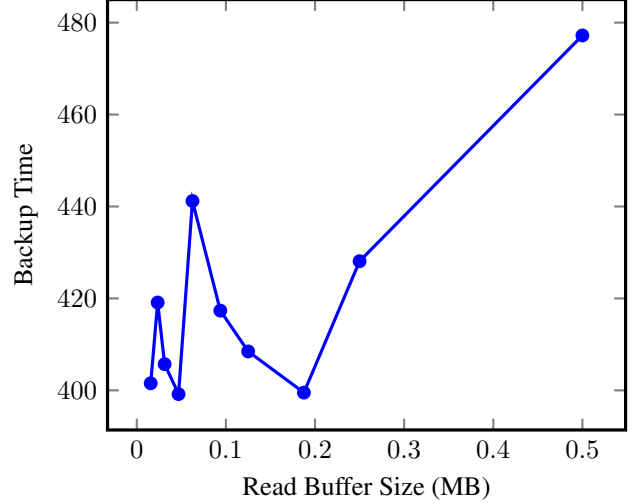


Figure 5: Backup time for varying amounts of memory allocated to disk read buffering. Other Settings: 10 Machines, 5x40GB VMs per machine, network buffers 2.5MB, write buffer 6.25MB

the parallel time does not change much, and the aggregated throughput scales up linearly since the number of VMs is $25p$.

Table 2 shows performance change when limit $M=35$ MB is imposed and the number of partitions per machine (q) varies. Row 2 is memory space required to load a partition of global index and detection requests. When $q = 100$, the required memory is 83.6 MB and this exceeds the limit $M = 35$ MB. Row 3 is the parallel time and Row 4 is the aggregated throughput of 100 nodes. Row 5 is the parallel time for using Option 1 with $p \times q$ send buffers described in Section 3. When q increases, the available space per buffer reduces and there is a big increase of seek cost. The main network usage before performing the final data write is for request accumulation and summary output. It lasts about 20 minutes and each machine exchanges about 8MB of metadata per second with others during that period, which is 6.25% of the network bandwidth.

7 Conclusion Remarks

The contribution of this work is a low-cost multi-stage parallel deduplication solution. Because of separation of duplicate detection and actual backup, we are able to evenly distribute fingerprint comparison among clustered machine nodes, and only load one partition at time at each machine for in-memory comparison.

The proposed scheme is resource-friendly to the existing cloud services. The evaluation shows that the overall deduplication time and throughput of 100 machines are satisfactory with about 8.76GB per second for 2500 VMs. During processing, each machine uses 35MB

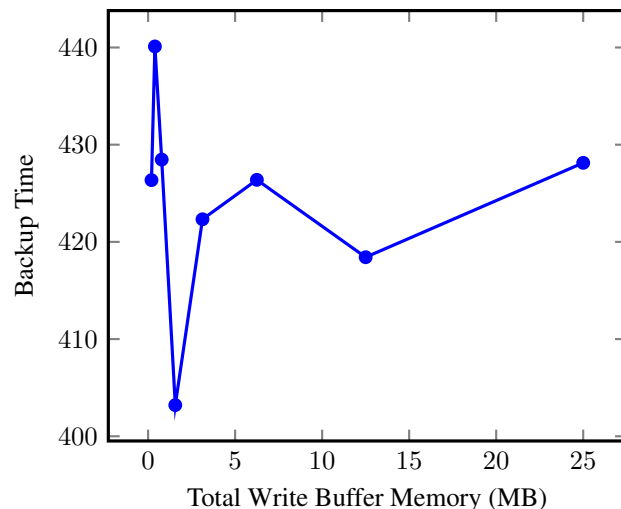


Figure 6: Backup time for varying amounts of memory allocated to network communication. Other Settings: 10 Machines, 5x40GB VMs per machine, network buffer memory 2.5MB, Read Buffer 128KB

memory, 8GB disk space, and 10-13% of one CPU core with a single thread execution. Our future work is to conduct more experiments with production workloads.

Acknowledgment. We thank Michael Agun, Renu Tewari, and the anonymous referees for their valuable comments. This work is supported in part by NSF IIS-1118106. Any opinions, findings, conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

References

- [1] Alibaba Aliyun. <http://www.aliyun.com>.
- [2] C. Alvarez. NetApp Deduplication for FAS and V-Series Deployment and Implementation Guide. NetApp. Technical Report TR-3505, 2011.
- [3] D. Bhagwat, K. Eshghi, D. D. E. Long, and M. Lillibridge. Extreme Binning: Scalable, parallel deduplication for chunk-based file backup. In *IEEE MASCOTS '09*, pages 1–9, 2009.
- [4] A. T. Clements, I. Ahmad, M. Vilayannur, and J. Li. Decentralized deduplication in san cluster file systems. In *USENIX ATC'09*, 2009.
- [5] EMC. Achieving storage efficiency through EMC Celerra data deduplication. White Paper, 2010.
- [6] F. Guo and P. Efstathopoulos. Building a high-performance deduplication system. In *USENIX ATC'11*, pages 25–25, 2011.
- [7] M. Lillibridge, K. Eshghi, D. Bhagwat, V. Deolalikar, G. Trezis, and P. Camble. Sparse Indexing: Large Scale, Inline Deduplication Using Sampling and Locality. In *FAST'09*, pages 111–123, 2009.
- [8] S. Quinlan and S. Dorward. Venti: A New Approach to Archival Storage. In *FAST '02*, pages 89–101, 2002.
- [9] S. Rhea, R. Cox, and A. Pesterev. Fast, inexpensive content-addressed storage in foundation. In *USENIX ATC'08*, pages 143–156, Berkeley, CA, USA, 2008. USENIX Association.
- [10] K. Srinivasan, T. Bisson, G. Goodson, and K. Voruganti. idedup: latency-aware, inline data deduplication for primary storage. In *FAST'12*, 2012.
- [11] M. Vrable, S. Savage, and G. M. Voelker. Cumulus: Filesystem backup to the cloud. In *FAST'09*, pages 225–238, 2009.
- [12] J. Wei, H. Jiang, K. Zhou, and D. Feng. MAD2: A scalable high-throughput exact deduplication approach for network backup services. In *IEEE MSST'10*, pages 1–14, May 2010.
- [13] T. Yang, H. Jiang, D. Feng, Z. Niu, K. Zhou, and Y. Wan. Debar: A scalable high-performance de-duplication storage system for backup and archiving. In *IEEE IPDPS*, pages 1–12, 2010.
- [14] W. Zhang, H. Tang, H. Jiang, T. Yang, X. Li, and Y. Zeng. Multi-level selective deduplication for vm snapshots in cloud storage. In *IEEE CLOUD'12*, pages 550–557.
- [15] B. Zhu, K. Li, and H. Patterson. Avoiding the disk bottleneck in the data domain deduplication file system. In *FAST'08*, pages 1–14, 2008.

p	the number of machines in the cluster
v	the number of VMs per machine. At Alibaba, $v = 25$
x	is the number of snapshots saved for each VM. At Alibaba, $x = 10$
t	the amount of temporary disk space used per machine for deduplication
m	the amount of memory used per machine for deduplication. Our goal is to minimize this
s	the average size of virtual machine image. At Alibaba, from our collected data, $s = 40GB$
d_1	the average deduplication ratio using segment-based dirty-bit. $d_1 = 77\%$
d_2	the average deduplication ratio using content chunk fingerprints after segment-based deduplication. For Alaba dataset tested, $d_2 = 50\%$
d_3	the average number of dup-with-new blocks, as a fraction of r (defined below)
b_r	the average disk bandwidth for reading from local storage at each machine
b_w	the average disk bandwidth for writing to local storage at each machine
b_b	average write bandwidth to back-end storage (block store)
q	the number of buckets to accumulate requests at each machine. (total number of buckets is $p * q$)
c	the chunk block size in bytes. In practice $c = 4KB$
u	the record size of detection request per block. In practice, $u=40$. That includes block ID and fingerprint
e	the size of a duplicate summary record for each chunk block
m_n	the memory allocated to network send & receive buffering. Total network memory is $2m_n$, with each buffer of size m_n/p
α_n	the latency for sending a message in a cluster
β	time cost for in-memory duplicate comparison

Table 1: Modeling parameters and symbols.

Table 2: Performance when $M=35MB$ and q varies.

#Partitions (q)	100	250	500	750	1000
Index+request (MB)	83.6	33.5	16.8	11.2	8.45
Total Time (Hours)	N/A	3.12	3.15	3.22	3.29
Throughput GB/s	N/A	8.76	8.67	8.48	8.30
Total time (Option 1)	N/A	7.8	11.7	14.8	26