

Revisiting Database Storage Optimizations on Flash

ABSTRACT

The database storage hierarchy has been heavily optimized for the performance characteristics of disks. Storage managers typically employ row- or column-oriented storage layouts, or a combination, to improve the I/O performance of different query workloads with disks. The recent rise of flash memory-based solid-state drives (SSDs) significantly change the performance characteristics of storage: these drives provide an order of magnitude lower read/access latencies, significantly higher read bandwidths, and most importantly, negligible seek overheads.

In light of these differences, we analyze major storage optimizations for read-optimized databases. We examine the benefits of row and column-oriented storage layouts on flash SSDs. Our measurements span through different workload variations, including selectivity, projectivity and concurrency that affect query processing on flash. Further, we also investigate the cost and benefits of a set of database optimizations, including data compression, prefetching, and indexes on flash SSDs. Our analytical models back our experimental evaluation of the performance tradeoffs of these optimizations.

Three of our key findings are: (1) SSDs scale up linearly with concurrent execution of database queries and outperform disks by up to a factor of two, (2) the low seek cost on SSDs makes column-storage a better choice for laying out data on a variety of flash devices, (3) and that while data compression is useful to further leverage the bandwidth of flash, database prefetching has less benefit for flash storage. Finally, we present a list of design implications of our findings on future database and operating systems for effectively embracing flash storage.

1. INTRODUCTION

Tape is Dead, Disk is Tape, Flash is Disk, RAM locality is King.

– Jim Gray [15]

For decades, databases have been optimized for the performance characteristics of magnetic disks, such as their long access and seek latencies and high sequential bandwidth [14, 27]. For example, databases often prefetch large buffers to amortize the cost of I/O over more data. However, solid-state disks (SSDs), built on flash memory, have recently achieved large capacity and high performance, making them a promising replacement for disks in many workloads.

SSDs represent a major advancement for storage management in database systems. To date, most uses of the flash technology have focused on their high random read capacity: a single mid-market device may provide 35,000 random I/O reads per second, while the fastest disks achieve barely 300. Thus, SSD usage in data management has been limited to the domain of transaction processing, where small random accesses are common.

However, there has been little investigation of the use of SSDs in decision support systems for analytical data processing. These workloads benefit from higher sequential bandwidths of SSDs, their small form factors and their low-power operation. A farm of slow, expensive and power-hungry disk arrays can be replaced with large SSDs optimized for selection, projection and scan queries used for business-intelligence applications and data warehouses. These applications deploy read-optimized databases for these workloads. In particular, such systems are tailored for read-only queries and are updated by bulk-loading with large database relations periodically [30, 18]. SSDs, with order of magnitude faster access latencies and high bandwidths, are well suited for these applications when combined with a separate write-optimized staging area for periodic updates [28].

Most read-optimized databases employ several techniques for improving the performance of magnetic disks [18, 20]. In particular, database storage managers:

- use row or column-oriented storage layouts [13, 30] or their combinations [9, 16, 17] to reduce the cost of I/O to disks;
- compression to improve the effective bandwidth of disks at the cost of increased CPU overheads [6, 14, 33];
- database and file-system prefetching to amortize seek costs by reading ahead additional contiguous pages from disk; [29].
- reordering, scheduling and delaying I/O requests to minimize seeking between different datasets on disk [22].

In the light of widely different performance characteristics for SSDs, the cost and benefits of these optimizations may change as compared to disks.

In this paper, we revisit these storage optimizations on flash storage. We experiment with a high-performance database storage manager [1] and workloads based on the TPC-H specification [5] to isolate the performance impact of different database storage layouts for SSDs. Our experiments use densely packed pages that resemble closely with the characteristics of various commodity read-optimized databases [7, 10]. Our measurements span through a range of devices and workload variations. With this study, we hope to inspire the database and OS research community to reconsider these optimizations originally designed for disks as many applications migrate to flash storage. Specifically, we address the following questions to unravel different performance tradeoffs for data processing on SSDs through our experiments and analysis:

- What is the impact of different storage layouts on database query processing on flash storage? How do the performance tradeoffs for row and column stores differ for SSDs when compared with near-line and enterprise disk configurations?
- What is the impact of different query processing workloads, such as changed relation size, selectivity, or concurrency of different queries, on SSDs? Does workload affect performance differently on flash than on disk?
- What are the costs and benefits of optimizations such as data compression, storage indexing and database prefetching when used with SSDs?

The rest of the paper is structured as follows. We review the basics of flash memory and solid-state disks in Section 2, followed by a description of modern database storage hierarchy in Section 3. Section 4 describes our experimental methodology, discussing our query workloads, storage manager, measurement framework and storage devices used for experiments. Section 5 describes our findings and analytical models for performance tradeoffs of different storage optimizations. Section 6 presents a list of design implications on future database and operating systems for effectively embracing flash storage. Finally, we present related work in Section 7 and conclude in Section 8.

2. FLASH STORAGE BACKGROUND

As prices drop and write performance improves, non-volatile NAND flash memory has become a viable storage replacement for hard disks. Solid-state disks, built of multiple flash memory chips, commonly provide a drop-in replacement for hard disks to avoid the need for new device drivers. With additional mechanisms incorporated in the device firmware called the Flash Translation Layer (FTL), SSDs mask the differences between flash and disk storage technologies. However, SSDs differ from hard disks in three major ways relevant to data-analytics workloads: I/O performance, cost, and power consumption.

I/O Performance. SSD performance differs from disks both in transfer rate and seek time. Most importantly, flash media provides significantly lower random read latencies (0.1ms vs. 4-8ms for disks). In addition, a single SSD may internally contain many flash chips, allowing RAID-like increase in I/O bandwidth within a single device. Thus, sequential read performance can be much higher than disks (250 MB/s for mid-range SSDs vs. 100 MB/s for the fastest disks). However, write performance for flash may be slower than disk, because blocks must first be erased. Better flash devices maintain a pool of clean blocks to absorb writes thus reducing the need to wait for erasing a block [8].

Device	Sequential (MB/s)		Random 4K-I/O/s	
	Read	Write	Read	Write
HDD	80	70	120-300/s	
USB flash	11.7	4.3	150/s	20/s
SSD	250	170	35K/s	3.3K/s
PCI-e flash	700	600	102K/s	101K/s

Table 1: Disk and NAND flash memory performance: Hard disks exhibit a small variance in performance due to its mechanical nature. In contrast, flash memory devices present a wide range in performance due to different host interfaces and significant internal parallelism.

In contrast to disks, which present small variance in performance due to their mechanical nature (seek times and rotation speed have a narrow range), flash storage devices exhibit a wide range of performance. Table 1 shows the performance of a variety of devices. Inexpensive and low-end devices such as USB flash sticks or camera memories offer moderate read bandwidth but have poor random-write performance. Solid-state disks (SSD), with a standard SATA interface provide much better performance, up to 3 times the fastest hard disks. This is mainly attributed to the device firmware, which implements intelligent block mapping schemes, parallel I/O accesses to multiple flash chips and write buffering [23]. High-end flash drives connected with the PCI-e interconnect interface and dedicated device drivers (rather than using the existing SATA drivers) are even faster [2]. The variance in flash performance arises from two sources. First, an SSD can incorporate additional banks of flash chips, allowing more throughput through parallelism. Second, an SSD can incorporate smarter FTLs that are better able to conceal the costs of erasing flash before writing.

Cost. Until recently, flash memory was far more expensive than either disk or DRAM. The density of flash memory chips has doubled 14 times in the last 19 years, which is faster than the Moore’s law for processors. This trend is expected to continue at least until a density of 32 GB/chip in the next few years [3]. Mid-range SSDs currently cost approximately \$2.8/GB (quote of Intel X-25M SATA SSD, as of October 2009), which is 2–10 times expensive per byte than enterprise and near-line disks. This arises from the manufacturing process of SSDs, which requires expensive wafer fabs. In contrast, for workloads demanding high random I/O operations per second (IOPS), flash SSDs are about 50 times cheaper than a configuration of disks supporting the same number of IOPS. In addition, price-per-MBPS for sequential throughput is comparable to disks, as a single SSD can deliver nearly triple the bandwidth of a single disk.

Power. Unlike disks, flash does not have any mechanical or moving parts. Hence, flash devices consume significantly lower power while operating and almost zero power when idle. The typical power consumption for SSDs range between 0.15–2 W when active and as low as 0.06 W when idle [21]. In contrast, power consumption for SATA disks is between 13–18 W, or six to ten times greater than a SSD. At 10 cents per kilowatt-hour, the cost of a single disk for continuous three year activity would be about \$47, and almost \$100 when including the cost of cooling and power dis-

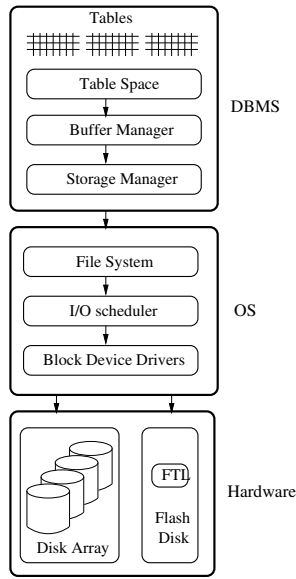


Figure 1: Database Storage Hierarchy: Buffer and storage managers employ different mechanisms to optimize for the performance of storage device and query workload. In addition, file systems and the operating system prefetch data and schedule block I/O requests submitted to the underlying device for amortizing disk seeks.

tribution, while a SSD can be powered for just \$10. This price difference increases further with an increase in the number of disks and expensive hardware controllers used for providing comparable device level parallelism through RAID arrays. Thus, SSDs tend to be price-competitive with disks when considering the complete cost of both device and power.

3. DATABASE STORAGE MANAGEMENT

Disk access can be a dominant cost for databases workloads, so database management systems carefully manage all I/O. DBMS storage managers also account for the storage device performance characteristics, such as seek time, access latency, and sequential bandwidth. Thus, databases lay out data and optimize their access patterns to minimize the cost of I/O.

Figure 1 shows a typical database storage hierarchy on Linux (sophisticated and special-purpose database systems may differ). Multiple database relations and storage indexes are clustered together in logical table spaces that are laid out as files on disks. The buffer manager maintains a pool of memory buffers to cache data in memory. Lower down the stack, the storage manager is responsible for most of the I/O to the underlying storage device, deciding which block to retrieve and when. The storage manager may directly access the device or may use the file system to perform I/O on its behalf.

Data Layout. One of the major focus of this paper is studying the impact of the storage manager’s data layout. The two major layout organizations are row-oriented, in which entire rows of a database relation are stored contiguously, or column-oriented, in which attribute values of each column in a table are stored contiguously. Row stores are more effective if the entire row is read, while column stores improve performance if only a small number

of attributes are projected from each tuple. Other hybrid storage layouts, such as PAX [9], DMG [17] and column abstractions [16], mix row- and column-oriented storage.

Compression. Compression can improve query performance by trading CPU processing for more effective use of disk and memory bandwidth. Column stores enable compression by storing all values of an attribute together, for example by replacing data values with indexes into a dictionary. Therefore, database administrators frequently use different compression schemes during the physical design phase of database schemas to optimize for both performance and storage space [6, 14, 33].

Prefetching. Database storage managers prefetch data that is not needed immediately. Prefetching for disks provides two benefits. First, reading more data at a time amortizes the high random seek latencies for disks over larger sequential requests. Second, prefetching overlaps I/O with computation, so that data is already available in memory when it is finally requested [29]. Storage managers in modern database systems, for example the SQL Server Enterprise, may prefetch up to 1024 pages (each page is 8 kilobytes).

OS and Device Optimizations. Within the OS, the I/O scheduler merges, reorders and delays requests to optimize the performance on the underlying storage device (presumably disks or disk arrays). The device provides the final layer of I/O scheduling. For disks, the controller may again reorder or buffer requests to improve performance based on the current location of the disk head. For SSDs, scheduling occurs in the FTL, which improves performance by remapping logical block addresses to physical flash addresses.

In summary, the database storage hierarchy embeds different disk-oriented optimizations at various levels. Both storage managers and file systems employ data prefetching to reduce access latencies, and optimize data layout to reduce the number of seeks. Finally, the disk scheduler and device drivers both re-order operations to minimize seeks as well. This paper revisits the cost and benefits of these disk-oriented optimizations for flash storage.

4. EXPERIMENTAL SETUP

The I/O performance of query processing in database systems is effected by the query workload, storage management and the characteristics of the storage device. This section presents our experimental methodology to investigate the impact of each of them for flash storage.

We focus our study on queries commonly used for large-scale data analysis. To model this workload, we use different select, project and scan queries based on TPC-H workload specifications [5]. To isolate the impact of various storage manager optimizations, we use a high-performance query engine [18] that uses column-oriented storage layouts and PostgreSQL 8.3, a widely popular open source database server [4] (Section 4.2). We measure performance on a variety of storage devices with different performance characteristics (Section 4.3) to quantify the impact of device performance on query performance. All our experiments are repeated multiple times, and we report the average over ten executions.

FIXME: [Do we need PostgreSQL?]

4.1 Query Workload

We focus our study on data-analysis workloads. **FIXME: [Describe what data analysis is first.]** This workload consists of selection, projection and scans over large relations, but few updates. Thus, these workloads are generally run on read-optimized databases **FIXME: [explain what this means: can ignore features needed for writing, such as concurrency control]**. Database relations are periodically updated in bulk from a separate write-optimized staging area, where new data is aggregated. This workload forms the basis of TPC-H [5].

Flash storage provides ample opportunities to optimize the performance of such queries because of its high read bandwidth and low random access latency. Hence, we study the performance of different variants of the following select, project and scan queries:

```
select T.a1, T.a2, T.a3 ... from T
where Predicate P(T.a1)
```

In this query, T represents the database relation; a_1, a_2, a_3 are different attributes and P is a sargable predicate on the first attribute. To isolate the effects of different storage layouts, we do not use storage indexes to accelerate queries. However, we use bitmap indexes to investigate the impact of SSDs for storing database indexes. We vary the number of attributes projected in each query from one to all. Similarly, we change the selectivity factor from 0.1% to 100% (low selectivity implies less qualified tuples) by modifying the predicate P. The number of columns projected and the number of rows selected in a query have a direct influence on its execution time.

Our experiments use two tables LINEITEM and ORDERS, that are based on TPC-H benchmark specification. We choose these tables to isolate the effects of width of relations and to ensure direct comparison of our results with earlier studies [18]. We use the official TPC-H toolkit [5] to populate these tables with data values for different attributes. For our experiments, LINEITEM represent a wide relation and has a tuple width of 150 bytes containing 16 attributes per tuple. ORDERS has a tuple width of 32 bytes and contains 7 attributes per tuple. To ensure a fair comparison between the two tables, we scale them to have the same number of rows: scaling LINEITEM by 10x and ORDERS by 40x ensures that both relations have 60 million tuples. LINEITEM takes over 9GB of disk space and ORDERS takes over 2GB. For most of our experiments, the size of these relations is sufficient enough to analyze the steady state I/O performance of different storage devices.

4.2 Data Manager

Query Engine. We focus on comparing the performance trade-offs of row and column stores for flash devices. In order to isolate the impact of storage layout, we use the query engine implemented by Harizopoulos et al. [18], which is available online [1]. While some commodity and research database systems implement column stores, such as C-Store [30] and MonetDB [10], they provide extensive performance optimizations for query processing such as in-memory database kernels built on virtual memory primitives, multi-threaded parallelism and vector storage for columns. These optimizations tend to blur the fundamental impact of row and column stores for flash devices which is the focus of our experiments.

The query engine used can operate on both row and column-oriented data. It has been used in previous published work [18, 20], and thus ensures a direct comparison of row and columns stores for flash de-

vices. Furthermore, the query engine uses zero-copy direct I/O, and transfers data directly from the storage device to user-space buffers without an explicit buffer pool.

Scanners. The query engine pre-compiles the queries and pipelines their execution for operating on the output blocks. Scanners reconstruct the tuples, apply predicates and extract the projected attributes and combine them for materialization later. Both the row and column stores use densely-packed pages on disk. The scanners for row stores read data pages from disks into an I/O scan buffer and then decode the columns from each page. Column store scanners, in contrast, read multiple files (chunks) from disk corresponding to the columns projected until the output tuple buffer is full. Each projected column is examined only at positions where the predicate was satisfied by the scan of the preceding column. This reduces disk I/O at the cost of additional seeking within a column.

Application Parameters. We tune the configuration parameters of the query engine for high performance. The major parameters we tune are: I/O depth (prefetching distance), I/O unit (scan buffer) size, page size, and block (materialized tuple buffer) size. We find that the most significant parameters are I/O depth and I/O unit size. We use an I/O unit of 128 KB and an I/O depth (prefetch read-ahead distance) of 6 MB (48 I/O units) unless otherwise specified. In addition to these application-level parameters, our experiments require careful configuration of operating system and storage device parameters, which we discuss in Section 4.3 and Section 4.4 respectively.

Data Compression. Compression can improve scan performance by trading CPU processing for more effective use of disk bandwidth. Flash devices have higher bandwidths and thus may benefit from compression in a different manner than disks. We use three different compression schemes for our experiments - bit-packing, dictionary and FOR-delta. Bit-packing stores each attribute using only the minimum number of bits in the maximum value of its domain. Dictionary-based compression uses an array with all distinct values of the attribute and stores each attribute as an index number to that array (similar to a hash lookup). FOR-delta (Frame-Of-Reference) uses a base value per page and stores deltas for attributes with it (see [6, 33, 14] for more details on these compression schemes). The performance differences of the compression schemes have been studied earlier [6], so we only present results for the best mechanism.

Database Indexes. Storage indexes improve the execution time for processing a query by directly seeking to the selected row. As our optimized storage manager does not support indexes, we instead use PostgreSQL 8.3 [4] with the same tables (LINEITEM and ORDERS) for these experiments. The PostgreSQL query planner selects the primary and secondary indexes for the columns used in query predicates. We investigate the impact of multiple indexes by using additional AND predicates. **FIXME: [We probably don't want to do this: Finally, we study the effects of hot and cold caches for index pages by varying the selectivity of the query and the size of available memory on the system.]**

4.3 Measurement Platform

Platform and Measurement Tools: We perform all measurements on a 2.5GHz Intel Core 2 Quad system configured with 4GB DRAM and 3MB L2 cache per core, running Ubuntu 8.0.4 (Linux kernel 2.6.24). We verify our results for the elapsed time for query execution using both performance counters and the Posix *time* utility.

Device	Sequential (MB/s)		Random 4K-I/O/s		Latency ms
	Read	Write	Read	Write	
Disk	80	68	120-300/s		4-5
SSD-Fast	250	70	35K/s	3.3K/s	0.1
SSD-Medium	69	20	7K/s	66/s	0.2
SSD-Slow	25	20	6K/s	136/s	0.6

Table 2: Performance characteristics of storage devices used: SSD-Fast, SSD-Medium and SSD-Slow represent different price points and performance. In general, SSDs substantially outperform disks for random read IOPS.

Furthermore, we instrument the Linux kernel with the Linux *blk-trace* mechanism to intercept and trace the I/O requests at the block layer in the operating system. These traces, along with the Linux *iostat* utility, enables us to monitor disk activity, such as the number of seeks performed during a query. We use ext2 file system for both disk and flash devices. While the journaling ext3 file system is more commonly used in practice, its read performance is identical to ext2 but the journal requires extra updates to ensure consistency after a crash.

4.4 Storage Device Characteristics

There is high variance in the performance of flash SSDs, with performance roughly corresponding to price. We therefore use three flash devices (solid-state disks) fabricated by three major SSD manufacturers at different price points. Among the three, two use a SATA 2.0 interface and the third uses SATA 1.0 interface. Since our intention is not to compare the performance of these competing SSDs, we refer to them as SSD-Fast, SSD-Medium and SSD-Slow from faster to slower devices. SSD-Fast is a relatively high-end device with SSD-Medium and SSD-Slow being intermediate and low-end devices. Table 2 presents the measured performance for these devices, which differs from advertised data-sheet values. Sequential read bandwidth, random read I/O operations per second and seek latency of the three devices are most relevant for our experiments. We use a Seagate Barracuda disk (SATA 2.0, 7200RPM) and simply refer to it as disk. **FIXME: [Update for new disks – probably drop results for this disk altogether.]** Column and row stores are in general affected by I/O bandwidth of the system, which RAID striping can improve. However, the general shape of our graphs for disks is similar to earlier studies with and without disk arrays [18, 20, 32], so our results can be scaled appropriately for disk arrays.

Flash devices usually require tuning to attain optimal performance. We enable Native Command Queuing (NCQ) for SSD-Fast and configure the system BIOS to treat SATA devices in native, rather than compatible mode, to boost its performance. The device I/O queue depth is configured to 32 for both SSD-Fast and disk to ensure a fair comparison. Finally, we enable on-disk prefetching for all devices unless otherwise mentioned.

5. PERFORMANCE STUDY

Database management systems use different storage layouts and other mechanisms such as data compression, prefetching and storage indexing to optimize the performance of different query workloads. We focus our measurement study on three questions surrounding these components:

- What is the impact of different storage technologies and lay-

outs on the performance of database query processing?

- What are the costs and benefits of different disk-oriented storage optimizations for flash storage?
- How does the performance of database storage vary across different query workloads?

We experiment with different disk, SSDs and device configurations and answer the first question in Section 5.1. Next, we evaluate the cost and benefits of different storage optimizations such as database prefetching, data compression and storage indexing in Section 5.2. Finally, we investigate a variety of query workloads to measure the performance of flash database storage in Section 5.3.

5.1 Database Storage Layouts

In this section, we experimentally evaluate and analytically model the performance of different database storage layouts across a range of devices.

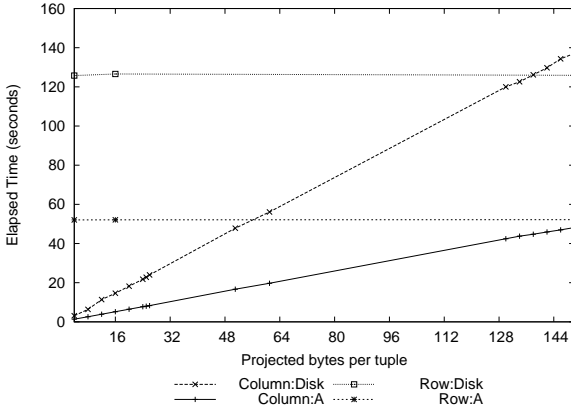
5.1.1 How does the performance tradeoffs for row and column stores differ for flash as compared to disks?

To answer this question, we compare the performance of row stores with column stores on the high-end flash SSD-Fast with disk. We measure scan performance with select queries for a selectivity factor of 10%. Column-store effectiveness increases when only a few attributes are retrieved, so fewer bytes are read from the device. Thus, we vary the number of attributes projected per tuple for our experiments. We use the LINEITEM table with a tuple width of 150 bytes and 16 attributes.

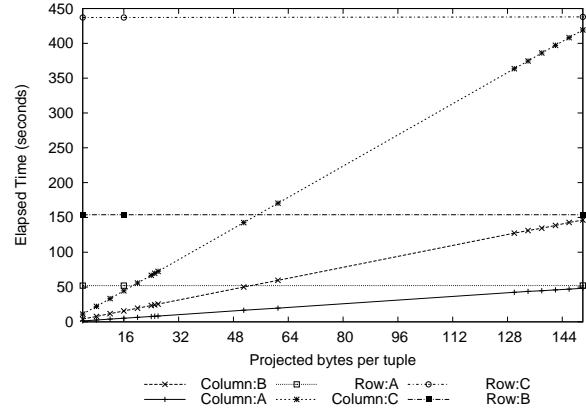
Figure 2(a) shows time to complete the select query using row and column store layouts on flash and disk with different projectivity factors. Both row and column store layouts for SSD-Fast outperform disk. Row stores for SSD-Fast are twice as fast as for disks, which reflects the difference in sequential read bandwidths of the two devices; row stores tend to saturate the I/O capacity of the device.

For disk, column store performance degrades quickly as the number of projected columns increases from 1 to 16 because of an increase in both I/O and CPU overheads for processing more columns. I/O wait time increases by factor of 8 when increasing the number of columns, and total CPU time for the benchmarks increases by a factor of **FIXME: [50 – this seems way to high]**. Thus, for disks, we observe a cross-over point when 90% of the tuple is projected at which point row stores become more efficient than column stores. **FIXME: [It is not clear what CPU and I/O overheads are. is this the time spent waiting for I/O and the CPU usage? The CPU usage increase by 50 seems way too high.]**

For a deeper understanding of why column store performance degrades with projectivity, we instrument the Linux kernel with *blk-trace* and trace each I/O request submitted to the device driver. We identify as *seeks* all requests that are at least 63 disk sectors apart from the previous request. At 100% projectivity, 867 seeks occur for column stores – almost 10 times greater than for row stores. At 10% projectivity, column stores issue only 83 seeks – roughly the same number for row stores. This suggests that as projectivity increases, column stores spend more time seeking, both across columns and within columns to skip the attributes which do not satisfy the predicate. This overhead adds to the time both row and



(a) Flash vs. Disk



(b) Variance in Flash performance

Figure 2: Performance of row and column stores on flash devices and disk. Flash devices at different price points exhibit different performance. However, in contrast to disks, column stores always outperform row stores for all flash devices.

FIXME: [These figures are way too small. They need to be two-across, not three-across. Or, they need to be scaled. Figure b in particular is impossible to read.]

column stores spend for reading the relation. Thus, it results in a crossover point where column stores perform worse with disks.

The shape of the performance curves for column stores on flash is similar to that for disk. However, column stores *always* perform better than row stores for SSD-Fast. SSD-Fast provides much lower seek latency (0.1ms vs. 4ms) than disks, which prevents a crossover between the column- and row-store curves even at high projectivities. For disks, the penalty for seeks at 100% projectivity takes at least 3.5 seconds, while they take less than 0.1 seconds for SSD-Fast.

We now present a simple analytically model of the performance of the two storage layouts to predicting the crossover point. We assume that select query workloads are I/O bound with negligible CPU overhead because there is a significant overlap between the CPU and I/O times. Let R be the size of the relation in megabytes and B the bandwidth of the storage device in megabytes/s. For a row store layout, the query completion time is given as the ratio of the two quantities:

FIXME: [We say ahead that we do do 867 seeks with column stores, but we have only 16 columns. hence, the number of columns is not actually the right thing to use when computing the number of seeks. Should we be multiplying by the 10% selectivity here as well so we get the predicted number of 867 seeks?]

$$t_r = R/B \quad (1)$$

For a column store layout, query execution time is also affected due to seeking between different columns. Let k be the number of attributes projected, C be the average size of each chunk (column file), and l be the seek latency of the device.

$$t_c = \alpha \cdot k \cdot (C/B) + \beta \cdot k \cdot l \quad (2)$$

In this equation, the elapsed time for column stores Equation 2 has two components: the time to read columns at full sequential bandwidth and the time to seek between columns. α and β adjust for imperfect I/O behavior. α reflects that not all columns are of equal width and that not all columns are read at the full sequential bandwidth. This is because ext2 and most other file systems do not lay out data in a perfectly sequential manner. Hence, α is the reduction in the sustained transfer bandwidth. Since column stores result in seeking both across and within columns to skip attribute values that do not satisfy the predicate, the term β adjusts the number of seeks; in most cases β is typically greater than one. **FIXME:** [does β depend on selectivity?]

The crossover point occurs when the performance of row stores equals that of column stores as modeled by Equation 1 and 2 respectively. We derive the following formulation for k , the number of columns projected at crossover:

$$k \approx \frac{R}{\alpha \cdot C + \beta \cdot B \cdot l} \quad (3)$$

For a device with no seek cost, the crossover point never occurs and k approximates the total number of attributes in the relation, which equals $\frac{R}{C}$. For devices with higher seek cost l , the crossover occurs when the bandwidth-delay product, $\beta \cdot B \cdot l$ becomes a substantial fraction of the column size C in the denominator.

For disk, l can be as high as 4-5 milliseconds and thus we find a crossover. In contrast, l for flash devices is typically an order of magnitude lower at one-tenth of a millisecond. Hence, the bandwidth-latency product in the denominator vanishes for flash devices. This results in k equal to the total number of attributes in the tuple, thereby explaining why there is no crossover for flash.

5.1.2 Do the performance tradeoffs differ across device models and disk configurations?

Flash devices at different price points provide widely varying performance due to different internal levels of parallelism (as in RAID

for disks) and sophistication of write-buffering algorithms [8]. Thus, the performance on a high-end SSD may not carry through to cheaper devices. We repeat our experiments on an intermediate flash SSD-Medium and a low-end flash SSD-Slow. Figure 2(b) shows row and column store performance for these two devices when compared to SSD-Fast.

We observe two important features. First, regardless of the device performance characteristics, column stores always perform better than row stores for all flash devices; there is no crossover between the two storage layouts. This again conforms with our formulation for predicting crossover for flash devices since the bandwidth-latency product for both SSD-Medium and SSD-Slow is half that of SSD-Fast and much smaller than a column.

Second, the performance of row and column stores for SSD-Medium is comparable to that of disk. This is because SSD-Medium possess lower read bandwidth than disk (69MB/s vs. 80MB/s), which is compensated by its an order of magnitude better seek latency (0.2ms vs. 4ms) to some extent. On the other hand, SSD-Slow always performs worse than disk for both row and column stores because of its considerably low sustained read bandwidth (25MB/s) and relatively higher seek latency (0.6ms).

While flash devices exhibit internal parallelism, near-line storage disks can be configured in RAID arrays to provide device level parallelism and improved performance. We construct a RAID-0 software disk array with two SATA disks, capable of delivering up to 160 MB/s bandwidth (80 MB/s per disk). However, the mechanical nature of the disk arrays limit the performance of column store layouts and it degrades with increased projectivity **FIXME: [How? While the configuration provides throughput similar to SSD-Fast, its seek time is still much higher. Thus, we still observe a crossover for the two layouts when 90% of the tuple is projected.]**

5.2 Database Storage Optimizations

In addition to customized storage layouts, database storage managers implement numerous optimizations to improve performance by hiding or reducing the cost of disk accesses. We next investigate the impact of three such optimizations on flash storage: compression, prefetching, and indexing.

5.2.1 How does database compression perform on flash storage?

As we have shown, performance of select and project queries is constrained by I/O bandwidth. Thus, compression offers an opportunity to improve I/O performance at the cost of additional CPU usage to compress and uncompress the relations. We use ORDERS for our analysis, which compresses a 32 byte tuple to 12 bytes.

Figure 3(a) shows the speedup in performance of row and column stores for different storage devices at 100% projectivity. We make two major observations about the effects of compression. First, compression greatly benefits both disks and SSDs but by different amounts. The query time for disk is reduced by up to 56%. However, compression improves the performance of both the SSD by up to 63%. The benefit of compression comes from increasing the effective I/O bandwidth while leaving seek latencies unchanged. Thus, for SSDs, where seek latency is negligible, compression offers greater benefits for reducing total I/O time. Second, row stores benefit from compression more than column stores at 100% projec-

tivity. This becomes more clear when we observe the performance of the two layouts as we vary the number of compressed attributes projected per tuple in Figure 3(b) and the elapsed time breakdown for 100% projectivity in Figure 3(c).

We plot the query execution time for flash SSD-Fast and SSD-Medium in Figure 3(b). As we increase projectivity, the CPU cost for re-assembling tuples from columns grows even higher with compression for column stores. Therefore, the benefit of increased effective bandwidth is negated by the extra CPU time spent generating output data. Furthermore, with compression the system cannot overlap CPU utilization with I/O as effectively for column stores. Figure 3(c) shows the breakdown of the elapsed time at 100% projectivity for compressed row and column stores on disk and SSD-Fast. In this figure, I/O time represents the fraction of the elapsed time that is neither user nor system time. Similar results are obtained for flash SSD-Medium and SSD-Slow. The salient feature of this figure is the time spent in usermode, which represents the time to uncompress data. Row store layouts require less processing, and hence are better able to overlap CPU usage with I/O. The column store, in contrast, spends more of its time reconstructing tuples and hence does not keep the device busy, leading to longer execution.

We now analyze the performance of row and column stores with compression by extending our original model described in Equation 3. For simplicity, we redefine R and C as the average size of the relation and a chunk (column file) after compression. As shown in Figure 3(b), row store performance is independent of the number of projected columns and is only dominated by I/O bandwidth, thus we reuse Equation 1 for its query execution time. **FIXME: [Unlike column stores, row stores are able to overlap computation with I/O, and thus the cost of uncompressing data is negligible.]**

However, columns stores are less able to overlap the extra CPU time to uncompress data with I/O because of their less-regular I/O patterns. Therefore, the CPU cost of compression is proportional to $k \cdot C$, the product of the number of columns projected and the average size of a chunk that is uncompressed. We multiply this product by γ to adjust for the differences in compression speeds and ratios for different schemes and attribute values. For flash devices with negligible seek latencies, we rewrite Equation 3 with the added CPU overheads for compression to compute k_{flash} , the number of columns projected at crossover:

$$k_{flash} \approx \frac{R}{\alpha \cdot C + \gamma \cdot B \cdot C} \quad (4)$$

The extra CPU cost of compression in the denominator explains the left shift of the crossover point for flash SSD-Fast and SSD-Medium in Figure 3(b). Furthermore, this favors row stores at higher projectivity.

In summary, compression benefits flash devices and disks by reducing the data to be fetched. However, flash devices benefit more, because bandwidth is a greater part of the I/O cost. In addition, the additional processing to reconstruct tuples can push the performance of column stores below row stores at high projectivities. This shows that although compression has high I/O benefits for flash devices, it can still tradeoff for column store performance because of its high CPU costs.

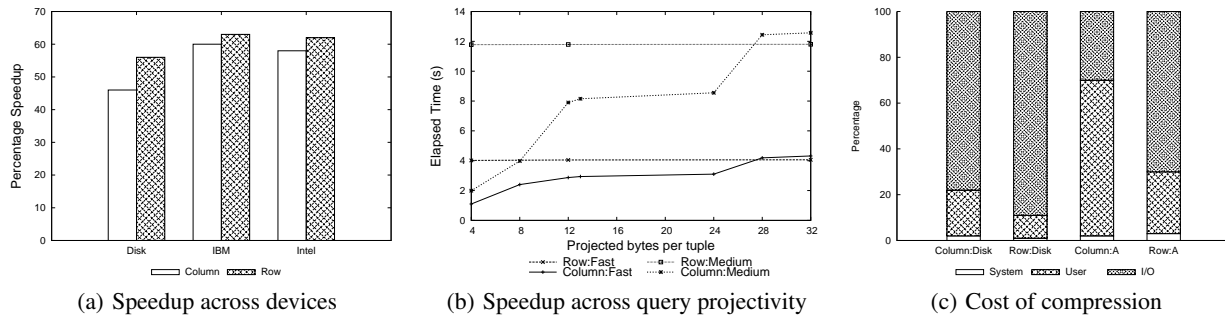


Figure 3: Performance of row and column stores on flash devices and disk with data compression. Compression improves the effective I/O bandwidth at the cost of additional CPU overheads. Row and column stores benefit differently from compression on different devices.

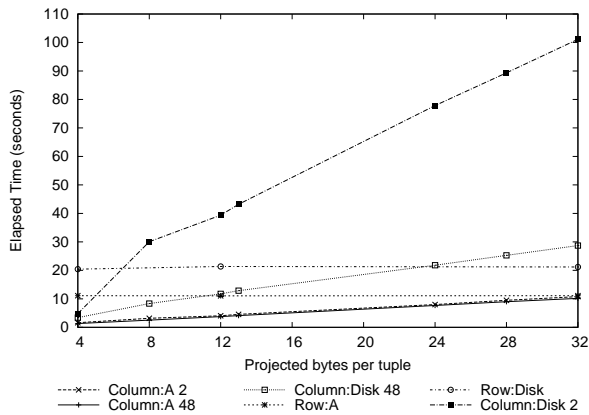


Figure 4: Database Prefetching: For disks, a high read-ahead is beneficial for amortizing seeks and improving the query execution time. In contrast, prefetching has no impact on the performance of column stores on flash devices.

5.2.2 Does prefetching benefit flash like disks?

Database storage managers prefetch data that is not needed immediately. Prefetching for disks provides two benefits. First, reading more data at a time amortizes the high random seek latencies over larger sequential requests. Second, prefetching overlaps I/O with computation, so that data is already available in memory when it is finally requested [29]. Unlike disks, flash devices possess an order of magnitude lower read access latency and low penalty for random access.

We measure the benefits of prefetching for the performance of row and column store layout by scanning ORDERS on the four devices. Figure 4 shows the query execution time with read-aheads of 256 KB (2 I/O units) and 6 MB (48 I/O units). We also measure the number of seeks for each data point. We observe that prefetching does not improve row store performance, because there are few seeks to be amortized. In addition, the number of seeks for row stores is fairly constant regardless of projectivity. Hence, we only show a single row-store curve for each of the two devices.

However, as shown earlier in Section 5.1.1, column stores incur more seeks as projectivity increases. Therefore, column stores benefit differently for the two devices with prefetching. For disk, as we increase the prefetch read-ahead from 256 KB to 6MB, there is a

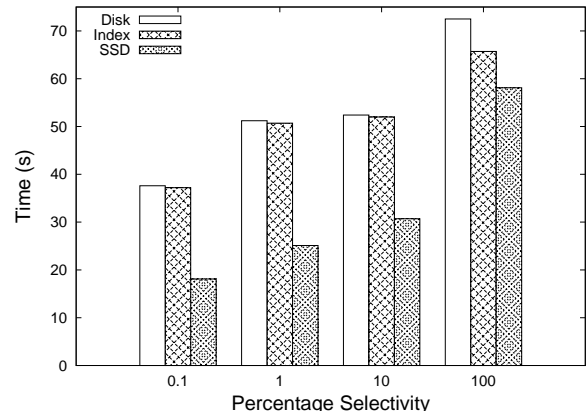


Figure 5: Impact of storage indexing on query processing for variable selectivities on PostgreSQL.

FIXME: [Change index to index-ssd]

dramatic decrease in the query execution time for column stores. This is because the number of times column stores seek decreases with an increase in the read-ahead: at 100% projectivity, the a 6 MB read-ahead causes 182 seeks, while a 256 KB read-ahead, it shoots up to 2941. Thus, prefetching compensates for disk performance by reducing the number of seeks.

However, for flash devices the cost for random access is only a small component of the access time. This is reflected in the curves for column store performance for flash SSD-Fast. There is hardly any difference in the performance of SSD-Fast as the prefetch read-ahead is reduced. Similar results are obtained for SSD-Medium and SSD-Slow.

In summary, storage optimizations that compensate for slow seeks, such as prefetching, are no longer required for flash devices, while those that improve effective bandwidth, such as compression, is still useful.

5.2.3 How do storage indexes perform on flash?

Storage indexes accelerate database query processing by directly seeking to selected rows instead of scanning all rows. As our optimized storage manager does not support indexes, we instead use PostgreSQL, which uses a row-store, for these tests. We configure PostgreSQL create bitmap indexes and use additional predicates

over different attributes **FIXME: [I don't know what this means - use additional predicates. Does it mean we change our query to reference multiple tables, as in a join? If so, give the new query].** Use of a table can change performance, as the table itself requires additional I/O, and use of a table can lead to seeking between selected tuples in the data tables.

We investigate the impact of storing indexes on disk and flash devices in three configurations. We compare storing both index and relation on (i) a single disk and (ii) a single SSD. Compared to the actual data, are accessed more frequently and are much smaller. Therefore, we also evaluate (iii) storing the index on an SSD and relation on disk, which allows a much larger dataset than the if it is all on an SSD. We use ORDERS relation and flash SSD-Fast for these experiments. We create a bitmap index on the first column with a total size of 375 megabytes and vary the selectivity, which effects the utilization of the index. Lower selectivity means the index is more useful, while at higher selectivity there is less opportunity to skip over unnecessary rows.

freffig:index1 plots the query execution time for these three configurations. We make three observations on the impact of indexing. First, we observe no difference in the benefit of index between disk and SSD: in both cases, the change in selectivity leads to a similar change in performance, indicating that index behavior affects performance on both devices similarly. **FIXME: [Explain why this is: shouldn't there be more seeks with an index at low selectivity?]**

Second, we observe that at low selectivities, storing the index on an SSD and the data on a disk has little benefit. The performance is similar because PostgreSQL aggressively prefetches index data to hide the cost of access, and thus reading the index has little impact on performance. However, we observe that at high selectivities, storing the index on an SSD improves performance by 10% compared to the disk-only case, because here the database is unable to completely conceal the cost of index access, so the faster SSD improves performance.

Third as we increase selectivity beyond 10%, indexing is less beneficial since almost all pages must be accessed from the relation. Thus, we find plain sequential scans outperform indexed scans by up to 34% for all devices at 100% selectivity. **FIXME: [we need to describe this test above]** With the addition of multiple predicates and secondary indexes on different columns, there is up to 8% increase in the CPU cost for reconstructing tuples for both disk and SSD.

5.3 Query Workloads

Database workloads, in terms of the query selectivity, width of data, and concurrency of access, can also accentuate differences between flash and disk storage.

5.3.1 Does the width of tuples affect this tradeoff?

For disks, prior work has shown that row stores perform better with narrower tables because narrow tuples can be packed tightly in a read-optimized page and thus can be scanned much faster [18, 20] The width of tuples also changes the number of bytes retrieved for each tuple. Narrow tuple size also directly affects both the bandwidth and seek components in Equation 2 due to a reduction in the size of columns.

To investigate the impact of tuple width, we repeat our experiments

with the ORDERS table, which has only 7 attributes per tuple with a total size of 32 bytes (in comparison, LINEITEM has 16 attributes in 150 bytes). Figure 6(a) plots the performance of row and column stores against projectivity on SSD-Fast and disk. We do not show the results for SSD-Medium and SSD-Slow for brevity since they were similar to SSD-Fast but scaled to their lower sequential bandwidths (as access latencies are similar).

For flash device SSD-Fast, the row and column store performance is similar to the LINEITEM table, and column stores still outperform row stores. Thus, column stores provide high performance on flash regardless of tuple width.

For disk, however, we observe the crossover point where row stores perform better occurs earlier for narrow tuples. Column stores perform worse than row stores when projecting more than 75% of the tuple. As compared to LINEITEM in Figure 2(a), the crossover point shifts left because row stores perform less I/O per tuple, so the opportunity to improve performance with column stores is lower. As each column is smaller, the seeks between columns have proportionally more impact on performance: column stores seek about 182 times and row stores only 19 times.

5.3.2 Does selectivity of the query affects this trade-off?

FIXME: [I think this would be better presented with two graphs: one of performance, and another one for CPU utilization. The reason is that we cannot very effectively show both at the same time – the stacked bars we had before when we tried concealed the change in CPU utilization. Showing a graph similar to 6a for flash/disk would help here.]

FIXME: [This result is more about row/column stores than about flash. How does this result have anything to do with flash?]

The selectivity of a query decides the number of tuples read while scanning that are discarded because they do not match the predicate. We investigate its impact by varying the selectivity of our queries from 0.1% to 100% on a log scale for both relations. To highlight the difference between row and column stores, Figure 6(b) shows the fraction of user and system times for queries that yield variable selectivities when executed on LINEITEM table stored on flash SSD-Fast. **FIXME: [Why no disk data?]** Row stores scan through the whole relation regardless of selectivity, so performance changes little. **FIXME: [This result isn't shown, as we don't show absolute I/O wait time. i think you mean I/O time, as clearly most if overlaps with computation]** Similarly, column stores scan through all the chunks corresponding to the columns being projected. Therefore, selectivity does not have a significant impact on the absolute I/O wait time for both row and column stores.

However, we observe that with increased selectivity there is an increase in the CPU component of the elapsed time only for column stores on both disk and SSD-Fast. As we vary selectivity for column stores, the total CPU utilization increases up to 59% percent at 100% selectivity for SSD-Fast. It remains constant between 5–8% for row stores regardless of selectivity. The increase in CPU time for column stores is mainly attributed to the extra work done by each scan node of the query engine which is driven by a separate value iterator.

Nevertheless, we do not find any crossover for SSD-Fast, regardless of the increase in the CPU component even at 100% selectivity

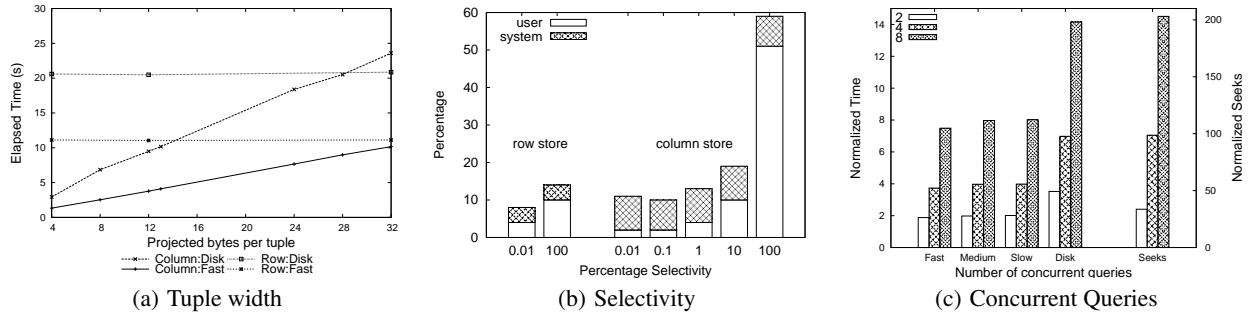


Figure 6: Performance as a function of workload for row and column stores.

FIXME: [Figure A has to say disk or SSD somewhere. Figure A is too small (the lines are hard to distinguish), and the text on figure C is too small also]

because the increased CPU time still overlaps with I/O. Similar results are obtained for the two other slower flash devices. However, with faster devices, such as the Fusion-IO ioXtreme PCI-e SSD that provides up to 600 MB/s sequential read bandwidth [2], there will be less overlap between the I/O and CPU times, and the increase in CPU utilization for column stores may tip performance to favor row stores at high selectivity. **FIXME: [Say something about if I/O bandwidth scaling exceed CPU capacity scaling?]**

We now extend our earlier models described in Equation 3 and 4 to explain the performance impact of selectivity. At low selectivity, column stores are always I/O bound and there is sufficient overlap between CPU and I/O wait time. Hence, our model described in Equation 3 is applicable without any modification.

However, as selectivity is increased, we must account for the time taken for reconstructing extra tuples. We extend Equation 4 accounting for CPU cost for compression to model the performance of queries with high selectivity. For simplicity, we redefine γ as the CPU cost for iterating through each selected row. Therefore, the total CPU overhead is the time taken for iterating through the column file corresponding to the attribute used in the predicate condition plus the selected fraction of the column files corresponding to the rest of the projected $k-1$ attributes. At a high selectivity factor of s , we rewrite Equation 4 to compute k_{flash} , the number of columns projected at crossover:

$$k_{flash} \approx \frac{R}{\alpha \cdot C + s \cdot \gamma \cdot B \cdot C} \quad (5)$$

For simplicity, we approximate the numerator and denominator in Equation 5 to match Equation 4. The CPU cost for reconstructing extra rows is given by the product $s \cdot \gamma \cdot B \cdot C$ in the denominator. As shown in Figure 6(b), increase in selectivity s for column stores also increases the percentage of this product accounting for CPU overhead. Therefore, this tends to favor row stores and explains less overlap in CPU and I/O wait time.

FIXME: [This seems unrelated and could be dropped.] Recent work by Tsirogiannis et al. [32] also shows that using PAX architecture [9], a hybrid of row and column stores, for storage layouts does not provide any improvement in query execution time for variable selectivity factors unless additional optimizations are used, such as reading only mini-pages that correspond to attributes

used in the selection predicate and using fully or partially sorted attributes. Furthermore, such optimizations only improve the performance of PAX layout for low selectivity queries over row stores. We discuss this more in Section 7.

5.3.3 How do concurrent queries scale with flash and disk?

Database systems may perform poorly with concurrent queries that cause competing disk traffic [24]. Such competing traffic can turn multiple sequential workloads into a collectively seek-bound workload that performs poorly on disks. The database, file system, and I/O scheduler of the operating system may try to minimize seeks by clustering nearby I/O requests. In some cases, concurrent scan queries to the same relation can be optimized by sharing the same scanner [19], so we analyze scans of different relations. We measure the performance with a single row or column store select query on an instance of ORDERS table, while competing against a variable number of concurrent row store select queries on an instance of LINEITEM table.

Figure 6(c) plots the query execution time of concurrent row store select queries normalized to one individual query for all three flash devices and disk, as we increase the degree of concurrency. Against the right Y-axis, we show the number of seeks. Column store performance is similar, so we do not include the results.

With a single query, performance is equal to that in **FIXME: [which figure]** and there are few seeks. However, when the number of concurrent scans increases to two, the number of seeks shoot up quickly and disk performs worse than both SSD-Fast and SSD-Medium, despite the 6 megabytes read-ahead that tries to amortize the cost of seeks. With an SSD, each individual query takes twice as long to complete, while with disk, each query takes 4 times longer. As the degree of concurrency increases, execution time increases linearly with the number of concurrent queries for all SSDs. This demonstrates that bandwidth is the most significant factor for concurrent queries with SSDs, as the seeks incurred have little impact on performance.

For disk, though, execution time increases twice as fast, and with 8 concurrent queries performs **FIXME: [14 times slower]** than for a single query. Unlike SSDs, seek times dominate performance for disk at high concurrency. Thus, performance for disk would be much better if the two queries were run sequentially rather than concurrently as the system is unable to effectively schedule the two

competing I/O streams to achieve maximum performance.

We also measure the impact of the operating system with two different I/O schedulers: CFQ and NOOP. CFQ batches up all the asynchronous requests from different processes in a number of queues with different I/O priorities. Since it tries to be fair for serving the requests from the different queues, it trades off with the seek distance between successive requests submitted to the device driver. In contrast, NOOP inserts all requests in a single FIFO queue and submits them as soon as possible.

We observe marginal differences in the performance of flash devices with the two schedulers. The performance difference between NOOP and CFQ is less than 3 percent, indicating that scheduling is less necessary for SSDs. For disks, we find a 13% improvement in NOOP performance against CFQ. This suggests that the workload submitted to the device has degenerated to a large extent to ensure fairness between different query executions and can not be improved by the scheduler **FIXME: [Does this mean that NOOP is better and CFQ makes it worse? This is pretty unclear]**. This is also visible in Figure 6(c) with a steep rise in the number of seeks. It rises by up to 33 times for merely two competing queries. For the 2 GB scan, this yields an average request length of 2.2 MB, which may be larger than the window of requests the CFQ scheduler considers for reordering. Furthermore, CPU utilization is as low as 4% for all cases which shows that all workloads are I/O bound.

These results suggest that achieving fairness for the different concurrent queries at the block layer may prove difficult with disks, and demands careful planning within the database or application. In contrast, flash devices require little scheduling to achieve high performance, and thus naturally perform well with concurrent queries.

6. FINDINGS AND IMPLICATIONS

The goal of this paper is to evaluate different components and optimizations in the database storage hierarchy for flash storage. We present a holistic view of mechanisms spanning the design of database, and OS I/O scheduling, and the characteristics of different storage devices. Our analytical models back our experimental findings on the performance tradeoffs of these mechanisms.

In this section, we present the design implications on future database and operating systems for effectively embracing flash storage.

Storage Layouts. We find that, unlike on disks, column stores outperform row stores on flash devices for a wide variety of query workloads. Similar to disks, they outperform row stores for low projectivity queries because they make better use of I/O bandwidth. Unlike disks, this tradeoff holds for high projectivity queries as well because SSDs possess negligible seek overheads. Our findings are consistent across different flash device models and disk configurations. At a high level, these findings make a strong argument database storage layouts that improve effective utilization of bandwidth best suit the performance characteristics of flash storage.

Database Compression. We find that data compression, which optimizes I/O bandwidth, has a greater benefit for SSDs than for disks, because I/O bandwidth accounts for a greater portion of performance. However, upcoming faster flash devices over new host interconnects, such as Fusion-I/O ioXtreme SSD over PCI-e bus [2] and Sun F5100 flash arrays over SAS interfaces [31], can effectively reduce the overlap between CPU and I/O wait times. Such devices may lead to CPU-bound workloads that do not benefit from

compression, or will require new compression schemes that balance CPU and I/O utilization.

Database Prefetching. We find that prefetching contiguous blocks to compensate for slow disk seeks is no longer beneficial for flash storage. Furthermore, the fast random access of SSDs provides new opportunities for the redesign of database prefetching. Rather than prefetching only sequential data, database storage managers can leverage their knowledge about the block access patterns of different query workloads. For example, they can effectively prefetch more distant pages with better temporal locality by using stride prefetching at a negligible cost of seeking on flash.

Storage Indexing. For low selectivity queries, indexes accelerate execution time by caching index pages in main memory for both disks and SSDs. However, for high selectivity queries, the effective utilization of indexes increases and storing indexes on flash offers significant performance improvement.

FIXME: [This is pretty peripheral to our paper and raises many other issues. We should probably cut it out.] Furthermore, recent technological advancements have also solved the problem of write amplification on mid and high-range SSDs to a large extent [12]. Considering such trends, future database systems can assume that frequent in-place updates to indexes are no longer slow on flash. Hence, database storage managers can benefit from the design of new hybrid systems that use flash for index tablespaces.

Concurrency. We find that the performance of flash storage scales linearly with increase in the degree of concurrency. In contrast, competing database queries can degenerate into a seeking workload and significantly degrade performance. Database mechanisms that optimize for disk performance by sharing the same scanner across different queries [19] can be significantly simplified. Flash storage offers a cleaner alternative for redesigning such database mechanisms and also providing scalable and faster performance.

Disk Scheduling. In addition to database mechanisms for managing seeking workloads, operating systems also cluster nearby I/O requests by reordering or delaying them. We find less than 3 percent difference between the performance of NOOP and CFQ I/O scheduling at the block layer in Linux kernel. Therefore, flash storage requires rethinking the design of a light-weight block layer in the operating system which keeps up with an order of magnitude low access latencies of flash than disks.

7. RELATED WORK

This paper draws on past work investigating database optimizations for flash and disk storage. We categorize this work into two broad classes: data layouts for flash and disk storage, and measurement studies on understanding the performance characteristics of flash devices.

Database Storage Layouts. Traditionally, database systems have mostly used the N-ary storage model (NSM), a page-based storage layout to store tuples contiguously. To save on the memory and disk bandwidths for queries projecting on a small fraction of tuples, Copeland et al. first proposed the decomposition storage model (DSM) [13]. Recently, more DSM-like (column store) commercial products and research prototypes have appeared, such as SybaseIQ, Vertica, C-Store [30] and MonetDB/X100 [10]. PAX (Partition Attribute Across) [9] is a hybrid approach which uses a DSM-like organization within a NSM page, thereby optimizing for

memory bandwidth. All these layouts trade I/O performance between different workloads. Harizopoulos et al. first investigated the performance tradeoffs for row and column stores [18]. Later, Holloway et al. [20] and Abadi et al. [7] answered many unresolved questions by focusing on a wider variety of scenarios. However, these studies only focus on the performance characteristics of disks, which widely differ from flash devices.

The most closely related work which focuses on flash-based database storage is by Tsirogiannis et al. [32, 28]. The authors investigate the suitability of a hybrid column-based page layout, based on PAX architecture, and compare it with NSM on a single flash device. They propose a new scan and join algorithm which leverages the column-based page layout to improve read efficiency. In contrast, we focus on the broader performance differences between disk and flash. We isolate the scenarios where the performance tradeoffs between row and column stores differ for flash devices from disks and provide analytical models for such differences. Furthermore, we analyze the tradeoffs for other disk-oriented optimizations like data compression, prefetching and I/O scheduling and highlight the additional benefits of flash devices for concurrent workloads.

Flash Measurement Benchmarks. Many studies have benchmarked the read and write performance of different flash devices to reveal their internals and provide hints for their optimal usage for different access patterns [8, 11, 12, 26]. Agrawal et al. present a taxonomy of design tradeoffs for the internal organization of SSDs [8]. They find that SSD performance is highly sensitive to workload and that FTL design choices greatly impact performance. Bouganim et al. describe uFLIP, a benchmark for measuring the response times for different flash access patterns [11]. Chen et al. present a measurement study investigating the intrinsic characteristics and system implications of solid-state disks [12]. They confirm the lack of write amplification for mid to high-range SSDs with new FTL designs. At a high level, similar to our work, all these studies acknowledge the high variance in the performance characteristics across different flash devices.

8. CONCLUSIONS

Database storage has been heavily optimized for disks over the last few decades. Compared to disks, flash devices provide an order of magnitude lower read/access latencies, much higher bandwidths and negligible seek overheads. In the light of these differences, we revisit major database storage optimizations in this paper, including data layouts, compression, database prefetching and indexes on flash. We analytically model the performance tradeoffs of these mechanisms for flash storage across different workload variations. Our study provides interesting design implications on future database and operating systems for effectively embracing flash storage.

9. REFERENCES

- [1] C-Store: A Column-Oriented Database. <http://db.csail.mit.edu/projects/cstore>.
- [2] Fusion-IO ioXtreme PCI-e SSD Datasheet. http://www.fusionio.com/ioxtreme/PDFs/ioXtremeDS_v.9.pdf.
- [3] NYTimes: Counting Down to the End of Moore's Law, May 2009. <http://tinyurl.com/o2nz2j>.
- [4] PostgreSQL Database Server. <http://www.postgresql.org>.
- [5] TPC-H Toolkit. <http://www.tpc.org/tpch>.
- [6] D. J. Abadi, S. Madden, and M. Ferreira. Integrating compression and execution in column-oriented database systems. In *SIGMOD*, 2006.
- [7] D. J. Abadi, S. R. Madden, and N. Hachem. Column-stores vs. row-stores: How different are they really? In *SIGMOD*, 2008.
- [8] N. Agrawal, V. Prabhakaran, T. Wobber, J. Davis, M. Manasse, and R. Panigrahy. Design tradeoffs for ssd performance. In *USENIX*, 2008.
- [9] A. Ailamaki, D. J. DeWitt, M. D. Hill, and M. Skounakis. Weaving relations for cache performance. In *VLDB*, 2001.
- [10] P. Boncz, M. Zukowski, and N. Nes. Monetdb/x100: Hyper-pipelining query execution. In *CIDR*, 2005.
- [11] L. Bouganim, B. por Jonsson, and P. Bonnet. uflip: Understanding flash io patterns. In *CIDR*, 2009.
- [12] F. Chen, D. A. Koufaty, and X. Zhang. Understanding intrinsic characteristics and system implications of flash memory based solid state drives. In *SIGMETRICS*, 2009.
- [13] G. P. Copeland and S. N. Khoshafian. A decomposition storage model. In *SIGMOD*, 1985.
- [14] J. Goldstein, R. Ramakrishnan, and U. Shaft. Compressing relations and indexes. In *ICDE*, 1998.
- [15] J. Gray. Tape is dead, disk is tape, flash is disk, ram locality is king, Dec. 2006. <http://tinyurl.com/d2enxp>.
- [16] A. Halverson, J. Beckmann, J. Naughton, and D. J. DeWitt. A comparison of c-store and row-store in a common framework. In *Technical Report, University of Wisconsin-Madison, TR1566*, 2006.
- [17] R. A. Hankins and J. M. Patel. Data morphing: An adaptive cache-conscious storage technique. In *VLDB*, 2003.
- [18] S. Harizopoulos, V. Liang, D. J. Abadi, and S. Madden. Performance tradeoffs in read-optimized databases. In *VLDB*, 2006.
- [19] S. Harizopoulos, V. Shkapenyuk, and A. Ailamaki. Qpipe: A simultaneously pipelined relational query engine. In *SIGMOD*, 2005.
- [20] A. L. Holloway and D. J. Dewitt. Read-optimized databases, in depth. In *VLDB*, 2008.
- [21] Intel. X-25 mainstream ssd datasheet, May 2009. <http://download.intel.com/design/flash/nand/mainstream/mainstream-sata-ssd-datasheet.pdf>.
- [22] S. Iyer and P. Druschel. Anticipatory scheduling: A disk scheduling framework to overcome deceptive idleness in synchronous IO. In *SOSP*, 2001.
- [23] H. Kim and S. Ahn. Bplru: A buffer management scheme for improving random writes in flash storage. In *USENIX FAST*, 2008.
- [24] H. T. Kung and J. T. Robinson. On optimistic methods for concurrency control. In *ACM TODS, Volume-6, Issue 2*, 1981.
- [25] Y. Li, B. Hey, Q. Luo, and K. Yi. Tree indexing on flash disks. In *ICDE*, 2009.
- [26] D. Myers. On the use of nand flash memory in high-performance relational databases. In *MIT MSc. Thesis*, 2008.
- [27] D. A. Patterson, G. A. Gibson, and R. H. Katz. A case for redundant array of inexpensive disks (raid). In *SIGMOD*, 1988.
- [28] M. A. Shah, S. Harizopoulos, J. L. Wiener, and G. Graefe. Fast scans and joins using flash drives. In *Fourth Workshop on Data Management on New Hardware (DaMoN)*,

SIGMOD, 2008.

- [29] E. Shriver, C. Small, and K. A. Smith. Why does file system prefetching work? In *USENIX*, 1999.
- [30] M. Stonebraker, D. Abadi, A. Batkin, X. Chen, M. Cherniack, M. Ferreira, E. Lau, A. Lin, S. Madden, E. O'Neil, P. O'Neil, A. Rasin, N. Tran, and S. Zdonik. C-store: A column oriented database. In *VLDB*, 2005.
- [31] Sun-Online. Sun Storage F5100 Flash Array.
<http://www.sun.com/F5100>.
- [32] D. Tsirogiannis, S. Harizopoulos, M. A. Shah, J. L. Wiener, and G. Graefe. Query processing techniques for solid state drives. In *SIGMOD*, 2009.
- [33] T. Westmann, D. Kossman, S. Helmer, and G. Moerkotte. The implementation and performance of compressed databases. In *SIGMOD Rec*, 29(3), 2000.
- [34] C.-H. Wu, T.-W. Kuo, and L. P. Chang. An efficient r-tree implementation over flash memory storage systems. In *GIS*, 2003.
- [35] D. Zeinalipour-Yazti, S. Lin, V. Kalogeraki, D. Gunopulos, and W. A. Najjar. Microhash: An efficient index structure for flash-based sensor devices. In *USENIX FAST*, 2005.