

A Morsel-Driven Query Execution Engine for Heterogeneous Multi-Cores

Kayhan Dursun
Brown University
kayhan@cs.brown.edu

Carsten Binnig
TU Darmstadt
carsten.binnig@cs.tu-darmstadt.de

Ugur Cetintemel
Brown University
ugur@cs.brown.edu

Garret Swart
Oracle Corporation
garret.swart@oracle.com

Weiwei Gong
Oracle Corporation
weiwei.gong@oracle.com

ABSTRACT

Currently, we face the next major shift in processor designs that arose from the physical limitations known as the "dark silicon effect". Due to thermal limitations and shrinking transistor sizes, multi-core scaling is coming to an end. A major new direction that hardware vendors are currently investigating involves specialized and energy-efficient hardware accelerators (e.g., ASICs) placed on the same die as the normal CPU cores.

In this paper, we present a novel query processing engine called *SiliconDB* that targets such heterogeneous processor environments. We leverage the Sparc M7 platform to develop and test our ideas. Based on the SSB benchmarks, as well as other micro benchmarks, we compare the efficiency of *SiliconDB* with existing execution strategies that make use of co-processors (e.g., FPGAs, GPUs) and demonstrate speed-up improvements of up to $2\times$.

PVLDB Reference Format:

Kayhan Dursun, Carsten Binnig, Ugur Cetintemel, Garret Swart, Weiwei Gong. A Morsel-Driven Query Execution Engine for Heterogeneous Multi-Cores. *PVLDB*, 12(12): 2218-2229, 2019. DOI: <https://doi.org/10.14778/3352063.3352137>

1. INTRODUCTION

Motivation: Within the last decade, databases have undergone a major shift in designs as a result of two important hardware trends: (1) Increases in main-memory capacities made it possible to hold even large databases in RAM space, thus eliminating the I/O bottleneck when accessing data on secondary storage (such as hard disks). However, this also required databases to reconsider some fundamental decisions, such as how to layout data for efficient access by leveraging the upper levels (i.e., caches) of the memory hierarchy. (2) Also, processor designers started to exploit Moore's Law from a different perspective by increasing the

number of cores, rather than focusing on single-core performance by improving clock-frequencies.

As a result, the advent of multi-core and multi-socket processor machines with large memories has led to a multitude of parallelization strategies in databases [1, 2], including new query scheduling paradigms such as the morsel-driven execution [13] model and also other optimizations to best leverage non-uniform memory access (NUMA) architectures [1, 2, 19].

Currently, we face the next major shift in processor designs originating from physical limitations known as the *dark silicon* effect [8]. Due to thermal limitations and shrinking transistor sizes, the multi-core scaling is coming to an end, since not all the cores in a processor can be powered up at the same time. To tackle this problem, a major direction that hardware vendors are investigating is to place specialized hardware accelerators (e.g., implemented as ASICs), that require less power consumption, on the same die together as with normal CPU cores and thus leading to more power efficient designs where all processing units can be utilized at the same time. One example of such a heterogeneous multi-core environment is the Sparc M7 processor [12], which combines normal CPU cores with an ASIC design that implements a Data Analytics Accelerator (named as DAX engines) for typical main memory database operations. Another example is the Intel HARP platform [17], which combines FPGAs with normal CPU cores. We believe that in the near future there will be many more of these specialized system-on-a-chip (SoC) designs that follow the same paradigm of combining heterogeneous cores.

However, developing efficient data management systems for these emerging heterogeneous multi-cores demands a critical rethinking of the architectural design and processing assumptions. In this paper, we take a first step and investigate how parallel query execution strategies should be designed to target these environments that combine normal cores with specialized ASIC-based accelerators.

To this end, one major challenge comes from the fact that state-of-the-art parallel query execution strategies have been developed for multi-cores with homogeneous compute units, thus do not hold optimal in heterogeneous settings. Here, one important difference is that the specialized hardware accelerators usually support only a *limited set of functions* at runtime. For example, the DAX engine in the Sparc M7 processor implements only two database operations (scan

This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 12, No. 12

ISSN 2150-8097.

DOI: <https://doi.org/10.14778/3352063.3352137>

and semi-join). A second major difference in these environments is that the specialized accelerators typically expose themselves as *passive units* and cannot actively issue work request for processing. As a result, existing query execution strategies for homogeneous multi-cores are not directly applicable to the emerging heterogeneous multi-core systems.

Moreover, query processing models designed to leverage co-processors [3, 5] are not optimally suited as well, since these techniques assume a discrete architecture where the accelerators are located at the end of a slow connection bus (i.e., PCI). For such environments, the typical assumption is that the major bottleneck of query processing is the data transfers between the CPUs and co-processors. Therefore, existing solutions instead implement coarse-grained query processing strategies that offload major parts of the execution to the co-processors so that the cost of data transfers are minimized due to the reduced communication. Unlike these discrete co-processor environments, heterogeneous multi-cores commonly allow the accelerators to access main memory via the same memory bus and sometimes even share the same lowest level caches (i.e., as in Sparc M7).

Contributions: In this paper, we present a query processing engine called *SiliconDB* that implements a novel parallel query execution scheme for the emerging heterogeneous multi-core environments. Specifically, we make the following contributions: (1) We develop a new query execution scheme based on the notion of morsels [13], but adapt it for heterogeneous multi-cores to address the specialization of hardware accelerators where they support only a *limited set of functions*. (2) In order to integrate the passive DAX engines and serve them with work elements effectively in the fine-grained scheduling scheme, we implement an adaptive push-based scheduling model that leverages a cost model based on queueing-theory. Here, we define the main challenge as to identify the optimal number of work elements that needs to be pushed to the accelerators in an adaptive manner so that all compute units are effectively utilized at all times. (3) We additionally propose some novel query optimization techniques that modify the output of a classical query optimizer towards query plans that are initially estimated to be more expensive, but show that they could yield to better resource utilization, thus could actually provide lower execution times.

In this work, we use the Sparc M7 processor as our development and testing platform. However, we believe that the architecture of *SiliconDB* and its query processing model are adaptive to other instantiations of heterogeneous multi-core processors. While we provide a discussion in Section 8 about this, explicitly demonstrating it is a line of future work and out of the scope of this paper.

Outline: The rest of this paper is structured as follows: In Section 2, we discuss the overall architecture of *SiliconDB*, and then present the details of its query execution model in Section 3. In Section 4, we provide the details of our adaptive pushed-based scheduling model and then in Section 5 we present some resource-aware query optimization ideas. In Section 6, we provide the results of our experimental evaluation based on the SSB benchmark and some other micro-benchmarks. Afterwards, we discuss the related work, give an outlook on the adaptivity of *SiliconDB* to other heterogeneous environments and conclude the paper in Sections 7, 8 and 9 respectively.

2. SILICONDB OVERVIEW

In this section, we provide an overview of *SiliconDB* a query processing engine we designed to address the challenges of emerging heterogeneous multi-core environments. While we use the Sparc M7 as our target platform and evaluate our ideas on a columnar data layout, our solutions are readily adaptive to similar environments, such as Intel HARP, and different storage layouts.

2.1 System Architecture

As mentioned before, emerging heterogeneous multi-core platforms commonly host general-purpose CPU-cores and specialized accelerators on the same die, and typically provide shared access to the same memory-regions (and even to the some levels of the cache as in Sparc M7) across all these compute units. To this end, in order to attribute to NUMA-awareness and to leverage the effects of cache locality, a key idea behind *SiliconDB*'s architecture is the logical grouping of cores and accelerators, that have access to the same NUMA regions, into the so-called processing units (PUs), where each processing unit implements a morsel-based parallel execution scheme (as depicted in Fig. 1).

While in a traditional morsel-based parallel processing model, cores that reside in the same processing unit would pull and process work elements from a single shared work queue, *SiliconDB* adapts this model to deal with two major challenges: (1) Due to their hardware specialization for specific operations, the accelerators of a processing unit can process only a limited set of work elements. For example, in Sparc M7, the DAX (Data Analytics Accelerators) engines are specialized hardware units that can only support the processing of scans and semi-joins over in-memory column buffers. (2) Also the hardware accelerators are typically implemented as passive units and cannot actively pull work-elements from shared queues, thus require a CPU-core to handle their processing.

To address the first challenge, *SiliconDB*'s query compiler first decomposes a pipelined query plan (such as the one shown in Figure 2) into multiple sub-pipelines and annotates them to specify if they could be executed by DAX engines. Furthermore, each processing unit in *SiliconDB* implements a core-only work queue, and a set of "function-specific work queues" to host the work elements that are supported by the hardware accelerators.

Following query compilation, *SiliconDB* starts query processing after initializing these queues with work-elements, where a work element refers to be a small block of data that needs to be processed (i.e., a fixed number of values of a column in our case). While the work elements in function-specific queues can be processed either by a general-purpose core or an accelerator (e.g., when we have a scan or a join queue in case of Sparc M7), the work elements in a core-only work queue can be processed only by general purpose cores. Here the main challenge is to schedule these work elements residing in different queues such a way to ensure that all available compute resources are optimally utilized.

To address the second challenge and serve the passive processing units, *SiliconDB* uses cores as handlers to actively push work elements to the accelerators. In the case of a naive model, these cores would schedule individual work elements at-a-time, referring to a direct implementation of a morsel-based execution incorporating accelerators. However, this type of a scheme would not only cause non-negligible over-

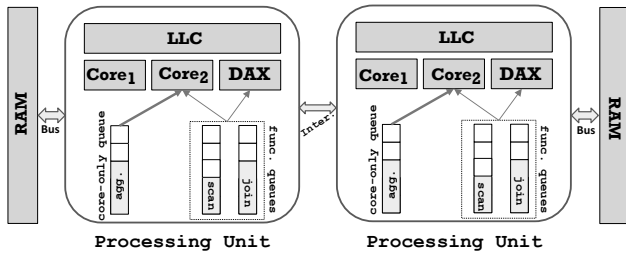


Figure 1: System Architecture

heads for the scheduling of individual work elements but would additionally cause the accelerators to become underutilized in between work-submission rounds. Since after completing the processing of an element, the accelerators would typically become idle and wait for cores to schedule a new work-element. Therefore we develop an adaptive scheduling model based on a cost-model that leverages queuing theory, mainly to decide how many work elements should be submitted to an accelerator at each round. Here another naive model would be pushing a high number of work-elements to the accelerators, but this would risk cores to turn idle (due to the lack of work-elements), especially towards the end of the query execution.

2.2 Query Processing Example

In order to intuitively show how *SiliconDB* works, we discuss the execution using a simple example. Figure 2 depicts the lifecycle of a SQL query consisting of a scan and aggregation operator. As we described before, the resource-aware query compiler first splits the query plan into two separate sub-pipelines. In our example, the first sub-pipeline of the plan represents a scan on a column `year`, which produces a bit-vector for the selected rows. Afterwards, the subsequent sub-pipeline uses these bit-vector results to identify the selected values of the column `price` table and executes the aggregation operation on them.

In this example, we assume that the query is executed in one processing unit of *SiliconDB* composed of only one core and one accelerator (as shown on the right hand side of Figure 2). For morsel-driven fine-grained parallel query processing, the processing unit uses one function-specific queue (the scan queue) to keep the work elements of the scan operator and one core-only queue to hold the aggregation work-elements that can only be processed by cores. For execution, *SiliconDB* first populates the scan-queue. Following our fine-grained operator processing model, *SiliconDB* initially splits the scan operator into multiple work-elements and places them into the scan queue to start the processing.

Afterwards, *SiliconDB* triggers several execution threads that actively start pulling work from the scan queue: one worker thread for the normal cores and one thread that is a passive-unit handler to serve the accelerator. The worker thread and the passive-unit handler are initially pinned to physical CPU-cores while they can switch their handling roles at any point of time. For execution, the core-handler as well as the passive-unit handler thread pull work elements from the scan queue and executes these work elements. After finishing one work element, all threads (core or passive unit handler) create a new aggregation work element that consumes the result of the scan. Once the first elements

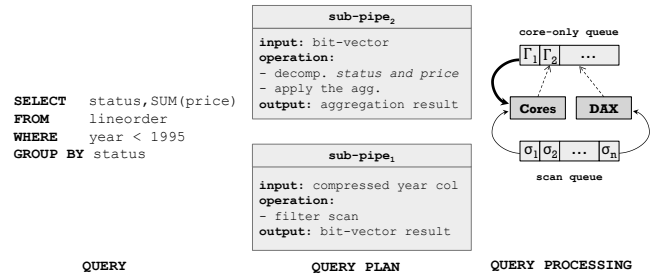


Figure 2: Query Processing

become available in the core-only queue, the core handler thread starts pulling these new elements from the core-only queue instead of working on scan work elements any further. This not only helps us to better utilize all resources but also pipeline the results between two sub-pipelines and thus leverage data locality if intermediate results are still in the cache.

3. QUERY PROCESSING

In this section, we discuss the details of *SiliconDB*'s fine-grained query processing engine and explain all the steps from query compilation to query execution including different scheduling strategies. In the end, we describe some optimization techniques that we implemented to further improve the query processing performance of *SiliconDB*.

3.1 Query Compilation

In the original morsel-driven query processing scheme [13], an incoming query is first compiled into several execution pipelines [16], where each pipeline processes small chunks of the input at-a-time, so named morsels. This morsel-based execution scheme lets the system to parallelize execution across multiple cores in an efficient way.

While *SiliconDB* compiles query plans into execution pipelines in a similar fashion, it additionally decomposes them into several *sub-pipelines* in order to address the different functionality characteristics of compute units (cores and accelerators) co-residing in the emerging multi-core environments. For instance, the specialized accelerators of the Sparc M7 processor (DAX engines) are implemented as ASIC units to support only a limited set of operators, thus cannot be used to execute full pipelines. Before we describe how we split pipelines into sub-pipelines in *SiliconDB*, we first describe the functions that a DAX engine can support for query processing. In the version we used for this paper, DAX engines are designed to support the following three database-specific functions:

(1) *Scan function*: This function scans and applies a filter on the elements of an in-memory array (e.g., an input column of a table) to detect the items that satisfy a given predicate. It returns a bit vector result representing the elements of the input column that satisfies the given query predicate.

(2) *Select function*: The select function also takes an in-memory array as input and outputs the elements represented by a given bit vector defining the selected entries.

(3) *Translate function*: In the perspective of query processing, this function can be used to execute a semi-join over

two in-memory arrays. We will give more details about this function later in the paper.

In regards to these functionality characteristics, *SiliconDB* splits the functions of a given pipeline into *DAX-executable* and *core-only* ones and the compiler tags the sub-pipelines they represent accordingly. For instance, for the example query in Figure 2 the compiler generates two sub-pipelines: Here sub-pipe_1 represents a DAX-executable operation that scans the *year* column of the input table and applies a filter, whereas sub-pipe_2 depicts the core-only executable operation of the query which applies an aggregation function on the filtered tuples of the base column.

Here it is important to note that *SiliconDB* is not restricted to execute DAX-executable sub-pipelines only on DAX engines, but also has the option to use typical cores for this purpose if necessary. The decision of what type of compute resource (cores or DAX engines) to use for a DAX-executable operation is given by our query execution scheme and it depends on the state of query processing at a given point of time as we discuss in the following subsection. We will also describe how processing DAX-executable sub-pipelines in this way instead of greedily pushing them only to DAX engines helps *SiliconDB* to provide better load balancing of query execution.

3.2 Query Execution

The query execution engine of *SiliconDB* is comprised of the so called *processing units* as shown in Figure 1. Each processing unit defines a logical component of *SiliconDB*'s processing model, where they each host a number of compute resources (CPU cores and DAX engines) that typically have access to the same memory region. Accordingly, *SiliconDB* partitions tables of a database across processing units and stores these partitions in their local memory regions to enable the benefits of NUMA-awareness.

Furthermore, the execution engine defines multiple work queues per processing unit different than a classical morsel-driven scheme that uses a single shared work queue for all cores per NUMA region. To this end, *SiliconDB* incorporates each processing unit with one *core-only queue* to hold work elements for *core-only sub-pipelines*. It also defines a set of *function-specific queues* to host the work elements of *DAX-executable sub-pipelines* that could be executed by either a DAX-engine or a regular CPU-core. In *SiliconDB*, we provide one function-specific queue for each DAX-executable function (e.g., a scan, a select, and a join queue). It is important to note that the cost of having multiple work queues (instead of a single one) for processing is negligible as the number of different sub-pipelines that could exist in a query plan is pretty small.

In order to process the work elements residing in these work-queues, *SiliconDB* pins one worker thread to each available strand of a core and also deploys a DAX handler thread inside each processing unit. While all these execution threads can actively pull work elements from corresponding work queues, only the worker threads are defined to process them directly, whereas the DAX Handler threads send them over to the passive DAX engines and monitor the state of their processing by these engines.

In the following section, we describe the details of *SiliconDB*'s work scheduling model using these worker and DAX handler threads.

3.3 Work Scheduling

The main goal of *SiliconDB*'s query processing framework is to leverage all compute resources (cores and DAX engines) effectively, so that they do not stay underutilized due to poor scheduling decisions. Now we will explain the details of the scheduling policies that the worker and DAX-handler threads follow towards satisfying this goal.

Worker Threads. As we pointed out in previous section, *SiliconDB* defines worker threads to actively pull and process work elements from core-only queues, as in a classical morsel-driven model, however with the possibility of accessing function-specific queues as well that contain DAX-executable work elements.

In regards to optimal resource utilization, here the main challenge is to decide on the order of work queues that the worker threads should be pulling the work elements from. To this end, each processing unit defines a priority map that drives its worker threads to pull work-elements from specific work queues at a specific state of query processing. More specifically, the priority mapping defines that a worker thread should first try to pull work elements from the core-only queue. In case no work elements are available in the core-only queue, the worker threads then start pulling work elements from the other function queues and work cooperatively with DAX engines on the same sub-pipelines (instead of staying idle). Algorithm 1 summarizes the overall process that the worker threads follow to implement this scheduling scheme.

Algorithm 1: Query Processing in *SiliconDB*

Input: List of Work Queues *WQ* ordered by priority

```

1 Algorithm processWorkElements(WQ):
2   while allProcessed(WQ) == false do
3     cur.queue ← getNextQueue(WQ) ;
4     sub.pipe ← getSubPipeline(cur.queue) ;
5     while cur.queue.hasNextElement() do
6       work_element ← cur.queue.nextElement() ;
7       processWorkElement(work_element) ;
8       scheduleFollowupWork(sub.pipe, work_element) ;
9     end
10  end
11 end

12 def scheduleFollowupWork(sub.pipe, work_element):
13   if sub.pipe.hasFollowing() then
14     next.pipe ← sub.pipe.getFollowing() ;
15     fw_queue ← GetWorkQueue(next.pipe) ;
16     fw_queue.addElement(work_element.getFollowing()) ;
17   end
18 end

```

The input to the query processing algorithm is an ordered-list of work queues that was generated using the priority-mapping. At the start of each iteration, the processing algorithm retrieves the first non-empty queue based on the priority-mapping that contains some work elements (line 3) and also the corresponding sub-pipeline (line 4). Afterwards, the algorithm executes the work-elements of the retrieved sub-pipeline (line 5-9) and upon their completion, it calls a procedure to schedule follow-up work elements as encoded by the query plan (line 8). The idea of the handling procedure is that it triggers the creation of follow-up work elements for a subsequent sub-pipeline in the query plan if such a sub-pipeline exists (14-16). For example, in Figure 2, the execution of a work element from sub-pipe_1 triggers the

scheduling of another work-element from sub-pipe_2 upon its completion.

DAX Handler Threads. As we described in previous section, DAX handler threads work in a different way than worker threads in order to address the passive processing nature of DAX engines. Since they cannot actively request new work-elements, DAX handler threads are responsible to optimally schedule work elements on these accelerators. Additionally they need to observe the state of the execution of the previously scheduled elements and carefully decide when to schedule new ones.

In order to comply with the main goal of *SiliconDB*, one important challenge for DAX handlers is to keep DAX engines optimally utilized during query execution. To this end, *SiliconDB* implements DAX handlers to follow an adaptive push-based scheduling strategy, which we describe in Section 4 in greater detail.

3.4 Optimizations of Execution

Pipelining for Sub-Pipelines. As we discussed previously, work elements in *SiliconDB* materialize the output of a sub-pipeline into memory and then create a follow-up work element (that would consume the materialized output) and place it on a corresponding work queue. However, pushing these follow-up work elements into a queue might cause unnecessary additional overhead and prevent cache-locality. Therefore, we extend the worker threads to provide a fusion mode between work elements of different query sub-pipelines if applicable. For instance, if a DAX-executable work element is executed by a worker thread and triggers a follow-up work element (line 16 in Algorithm 1) that is core-only executable, the worker thread immediately processes this work element instead of placing it into the core-only queue.

Furthermore, we also adapted the work element handling of DAX-handlers to provide better cache locality since DAX engines and normal cores share the same last level caches (LLC). The idea is that a DAX-handler can push work elements to the front of the work queues to maximize the chance that the output is still in the LLC. For instance, when the DAX-engine completes the processing of a scan element in Figure 2, the DAX-handler pushes the follow-up aggregation item to the front of the core-only work queue. This way, the likelihood that the input for the aggregation work element (output of the scan) is still present in the LLC is maximized when the work element is pulled by one of the worker threads.

Work Stealing. In addition to pulling work elements from the local work queues, we allow all worker threads and the DAX-handler of a processing unit to additionally steal work elements from the queues of a remote processing unit if all work queues of the local processing unit are empty. This way, we can mitigate the chance that individual processing units remain idle while others are overloaded (e.g., due to data skew).

Support for Concurrent Queries. Since *SiliconDB* defines query processing as the scheduling and execution of fine-grained work elements, the model is directly adaptive to support workloads that would consist of multiple queries. Here, we require only extra bookkeeping information tagged

along with work elements (to represent the actual queries they belong to) so that their outputs are properly directed to the corresponding query output. Even though there exist additional optimization opportunities to support concurrent queries (e.g., work sharing), these techniques are out of the scope of this paper. Therefore we implement just a baseline model where the query compiler places initial work elements into work queues with no preferential ordering and processing proceeds in a FIFO fashion as in the single query case.

4. ASYNCHRONOUS SCHEDULING

In this section, we provide the details of the scheduling model of *SiliconDB* for DAX handler threads and describe how we incorporate it with the rest of the framework in order to comply with the main goal of *SiliconDB*, optimal utilization of compute resources. To this end, we look into the problem from two perspectives:

(1) Implementing an optimal handling strategy so that system resources are not sacrificed just for observing the DAX engines (cf. Section 4.1).

(2) Providing an adaptive scheduling model that effectively utilizes the accelerators during query processing (cf. Section 4.2 and cf. Section 4.3).

4.1 Handling Model

In *SiliconDB*, we implement two different strategies to handle the scheduling of work-elements on DAX engines:

(1) *Separate Handler:* In each processing unit, one core is reserved to pull work elements from function-specific queues that hold DAX-executable work elements and assign them to a DAX engine. Furthermore, they monitor the status of ongoing work elements and potentially create follow-up items upon the completion of previously scheduled work-elements.

(2) *Piggybacked:* In this model, there is no dedicated handler thread to assign work to the accelerators. Instead, all the worker threads in a processing unit share the responsibility of scheduling work for the DAX engines. The idea is that each worker thread observes whether or not a DAX engine has finished the processing a work element, and decides if a new one should be pushed over to a DAX engine before retrieving a work element for itself.

4.2 Scheduling Policy

DAX handlers can submit a single or a group of work elements to a DAX engine at a time, which internally places them into a hardware queue and then assign these elements to the actual hardware execution pipelines for processing. It is important to note that the Sparc M7 processor implements each DAX engine to have four of these pipelines letting them to process multiple work elements simultaneously. Therefore, the actual processing of work elements on the execution pipelines is handled by an internal scheduling mechanism that the DAX handlers do not have a direct control over.

For the overall query processing scheme, the main challenge is (1) to ensure an internal hardware state where all the execution pipelines are effectively utilized and (2) the number of work elements waiting for service in the internal hardware queues are at minimum levels. While each one of these goals can be satisfied independently by applying simple heuristics, an optimal solution requires more so-

phisticated schemes to find a sweet-spot between these two important metrics.

For instance, the system could ensure that the utilization is always high by submitting a large number of work elements in each round but incurring the risk of having them to pile-up in aforementioned internal hardware queues. This not only would increase the average execution time of work elements in accelerators, but would also hurt the overall query processing performance since DAX engines might become stragglers and dominate the overall query runtime.

Similarly, the size of the internal queues can be reduced by having the system to initially submit as many work elements as the number of execution pipelines and then to supply new ones, one at a time for each completed item. However, with this approach it is clear that there is a risk for DAX engines to become underutilized since DAX handlers would first need to detect that a work element is completed before they can submit a new one. Moreover, there is some internal overhead associated with scheduling of work elements on DAX engines that further increases the latency for each element.

Therefore, the main challenge of query scheduling over DAX engines is to find an optimal number of work elements that the DAX handlers should submit to these accelerators at every iteration. For the ease of representation, we describe our ideas around the model where *SiliconDB* uses a separate thread for DAX Handlers. However, the ideas are directly applicable for the *piggybacked* case as well.

The main idea behind *SiliconDB*'s scheduling policy for the passive hardware accelerators is to define the value of an internal scheduling parameter called the *q-size*, which refers to the maximum number of work elements that can reside in the internal hardware queues of these accelerators. Accordingly, in every iteration the DAX handler thread monitors how many work elements are completed and uses the value of *q-size* to decide the number of new work elements that should be submitted to the corresponding DAX engines. Next, we describe the procedure that *SiliconDB* follows in order to assign the value of this important *q-size* parameter.

4.3 Optimal Queue Size

Before providing any further details, it is important to note that *SiliconDB* needs to adjust the value of *q-size* periodically, since the DAX engines will be processing work elements with different characteristics during query execution (i.e., different job-types or data characteristics). Moreover, the performance of DAX engines are sensitive to other run-time effects, such as the contention on the memory bus that is shared between multiple execution threads. Therefore, *SiliconDB* implements the DAX handler threads to follow an adaptive scheduling model where they continuously observe the DAX engines to collect statistics at query runtime. Then they use these observations to adjust scheduling decisions accordingly (i.e., updating the value of *q-size*) by leveraging a defined cost model.

Runtime Statistics: Since the internal state of DAX engines are not exposed to the operating system, we define and leverage an analytical model to get estimations, specifically on the two important metrics that we described in Section 4.2. In the rest of this section and in our cost models, we refer to these parameters as *utilization* and *items-waiting* to represent the utilization of the internal execution pipelines

and the number of work elements that are waiting for service in the internal hardware queues. To be able to estimate these metrics for each DAX engine, we leverage a queueing model with limited capacity defined with the $(M/M/s/k)$ mathematical notation.

The parameters of a typical queueing model consists of the system's arrival (M) and service (M) rates, the number of units that can serve the system simultaneously (s), and the maximal queue length (k) that is allowed, all respective to their orders in the notation. Then with these parameters in place, a queueing model can provide estimations about the internal state of a corresponding system. As any queue-based system can leverage queueing theory to get estimates about its performance and then uses them to improve service, in *SiliconDB* we leverage a queueing model to estimate the aforementioned *utilization* and *items-waiting* metrics and use them in a cost model to improve the scheduling scheme that the DAX Handlers follow.

In this regard, *SiliconDB* actively collects statistics about query processing such as how often DAX handlers submit work elements (to derive the arrival-rate) and how long it takes DAX engines to process these elements on average (to derive the service-rate). Then it provides these statistics as the input of the queueing-model along with two parameters representing the number of execution pipelines in each DAX engine (server-size) and the current *q-size* defined by the DAX handler.

Now in the following, we describe the details of our cost-model that the DAX handler threads use to find optimal values for the *q-size* that is based on the *utilization* and *items-waiting* estimations as provided by the queueing model and also on some throughput statistics representing query performance.

Cost Model: As we mentioned before, *SiliconDB* leverages a cost model that provides a runtime estimation value for executing a sub-pipeline (i.e., a group of work-elements). The asynchronous scheduling model of DAX handlers leverages this cost model periodically to find and adjust a *q-size* value that provides the minimum runtime estimate for the execution of a group of work-elements residing in regular time windows.

The following equations show how these runtime values are estimated based on the estimations of *utilization* and *items-waiting* for the DAX engines. For our cost model, we assume that cores and DAX engines can cooperatively work on work elements of a sub-pipeline and their outputs should be materialized before the processing of a next sub-pipeline could start.

$$tput_{\text{overall}} = tput_{\text{cores}} + tput_{\text{dax}} \cdot \frac{\text{utilization}^{W+1}}{\text{utilization}^W}$$

$$\text{cost}_{\text{runtime}} = \frac{N}{tput_{\text{overall}}} + \frac{tput_{\text{dax}}}{\text{items} - \text{waiting}}$$

As depicted with the equations above, the main components of the cost model are the $tput_{\text{cores}}$ and $tput_{\text{dax}}$ parameters, which represent the throughput that cores and DAX engines produce for an observation window, W . While the *utilization* and *items-waiting* parameters are provided by the queueing model as we described before, the scheduling

model monitors the throughput parameters at query runtime and then calculates an overall runtime estimation for the remaining N work elements.

The main idea of the cost model is that the q -size affects the *utilization* of DAX engines and thus their throughput, while the throughput of the regular cores stay stable. The change (increase or decrease) in throughput is given by the ratio between $utilization^{W+1}$ and $utilization^W$, which defines the throughput of the utilization for window $W + 1$, after changing q -size over the estimated utilization of the current window W .

Therefore, the cost model allows us to estimate the overall throughput $tput_{overall}$ for the next window $W + 1$ that results from choosing a new q -size. The throughput can then be used to compute an estimate for the total runtime for processing the remaining N work elements for a given sub-pipeline. The runtime estimate is thus given by $(N/tput_{overall})$ plus the additional time it requires to process the already scheduled work elements (*items-waiting*) that now reside in the internal hardware queues of the DAX engines (as provided by the queuing model). Algorithm 2 shows the overall procedure that *SiliconDB*'s adaptive scheduling model follows to find the optimal q -size by leveraging the defined cost model.

Algorithm 2: Finding the optimal queue size

Input : Set of work elements processed in the previous window, W
Input : Submission and Completion TimeStamps of work elements processed by DAX engines, TS
Output: The new optimal queue size, q -size

```

1 Algorithm adjustQueueSize( $W, TS$ ):
2    $a\_rate \leftarrow calculateArrivalRate(TS)$ ;
3    $s\_rate \leftarrow calculateServiceRate(TS)$ ;
4    $tput \leftarrow calculateThroughputRates(W)$ ;
5   foreach  $newQSize \in possibleQSizes$  do
6      $(utilization, items-waiting) \leftarrow$ 
        $QueueingModel(a\_rate, s\_rate, newQSize)$ ;
7      $curRuntime \leftarrow cost(tput_{overall}, utilization, items-waiting)$ ;
8     if  $curRuntime \leq minRuntime$  then
9        $minRuntime = curRuntime$ ;
10       $q\_size = curQSize$ ;
11   end
12 end
13 return  $q\_size$ ;

```

In order to find the most optimal q -size for the current state, the scheduler applies a linear search over the possible values of it, uses the estimation models as described and picks the q -size which results in the minimal estimated runtime. It is important to note that linear search is possible since the search space is sufficiently small. We simply apply a neighbor search starting with the current q -size and increment/decrement its value linearly.

In Algorithm 2, we summarize the steps the scheduler follows towards finding the most optimal q -size:

(1) For each observation window it first calculates the average arrival and service rates (line 2-3) for the work elements based on some runtime statistics and also the estimated throughput rates of cores and DAX engines (line 4).

(2) Then for each possible q -size, it first uses the queuing-model to estimate the *utilization* and *items-waiting* parameters using the observed arrival and service rates (line 6),

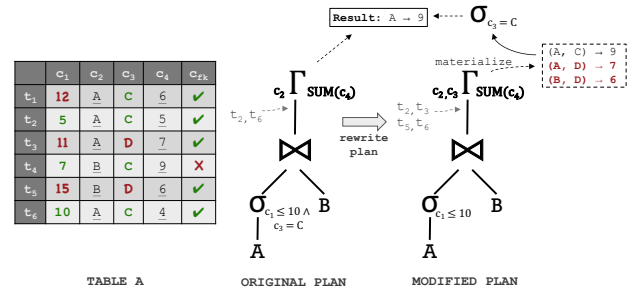


Figure 3: Additional Materialization

(3) and then using the estimated *utilization* and *items-waiting* values along with the observed throughput of compute units, it calculates a runtime value for the execution of work elements using the cost model we described before (line 7).

In the end, it returns the q -size value that produced the minimal estimated runtime.

5. QUERY OPTIMIZER

In this section, we explore query processing from a typical query optimization perspective and describe why existing query optimizers should be reconsidered for heterogeneous multi-cores. We also propose new query optimization techniques and demonstrate their potential in these new environments.

5.1 Overview

Since existing query optimizers generate query plans with the assumption of a homogeneous processing environment, there are a variety of reasons why these plans may not be applicable or optimal on heterogeneous multi-cores. For instance, a query plan might suggest the use of the operators that are not supported by the accelerators, which would render them useless for processing this plan since they can not be utilized at all (i.e., DAX-engines for the processing of a full hash-join).

Therefore, we suggest adapting the optimizer to take into account which compute resources the system can leverage to process a corresponding query plan. While this is a non-trivial problem on its own and would require a major adaptation of the optimizer, we look into this problem around the characteristics of the Sparc M7 processor with its DAX engines. Specifically, we propose simple heuristics that would rewrite a query plan produced by a classical optimizer into a potential one that would better utilize the Sparc M7 processor. With such a rewrite, we show that we can reduce the overall runtime as we demonstrate using some micro-benchmarks.

5.2 Heuristic 1: Additional Materialization

First, we propose to insert extra scan operators at the end of a query plan in order to better utilize the DAX engines. The query optimizer would normally regard this approach as an inefficient way of executing the query.

We motivate our heuristics using an example query that has an aggregation on top of a hash-join between two relations. Fig. 3 represents the original query plan that was produced by the query optimizer, which suggests applying

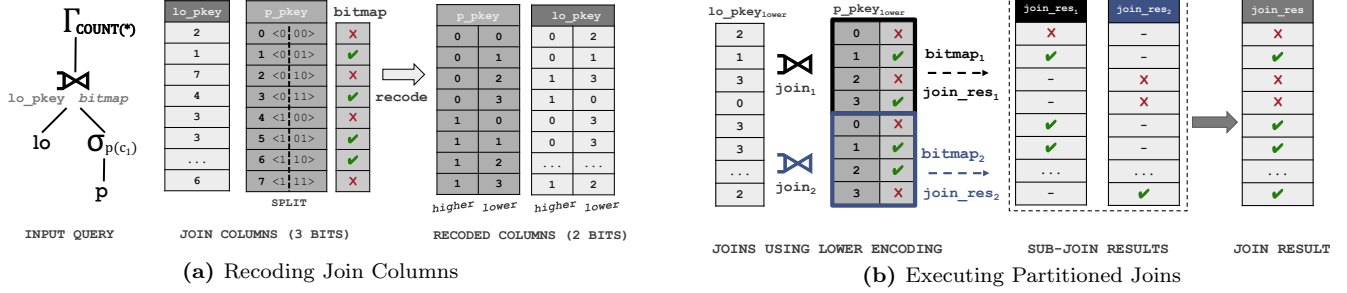


Figure 4: Operator Re-Coding

all the filters first to the relation that probes into the hash-table. In order to execute this plan on the Sparc M7 environment, the query processor would initially be able to leverage the DAX engines, but then they would stay idle during the latter parts of execution after all the scans at the bottom of the query plan are completed.

To address this problem, i.e., to allow the system to leverage the DAX engines in the later stages of query processing, we propose to modify the query plan by inserting an extra scan operator at the end to filter some temporary results produced by the aggregation.

In the example, we illustrate this idea as pulling up the second scan ($c_3 = C$) defined on the first relation and having the aggregation to first materialize its output results into a temporary buffer. Following this modification, we have the system leverage the DAX engines in order to filter out of these results (i.e., the second and third tuples in the temporary output) to produce the final query result.

To implement this type of an optimization, it is clear that the system should carefully examine if the benefits of applying it would outweigh the costs associated with it, and then provide a rewrite strategy accordingly. Here, we show the potential of such a query optimization technique using our target platform. We also provide the results of a micro-benchmark we implemented in this regard in Section 6.3.

5.3 Heuristic 2: Join Re-Coding

While DAX engines can support the functionality of a semi-join operator, they are designed to process inputs of at most 16-bit encodings. For FK-PK type of joins, this restriction requires the smaller PK relation to have at most 2^{16} keys, otherwise accelerators cannot be utilized during the execution of this operator.

We propose a technique that re-codes the inputs of the semi-join operator so that they are represented with smaller bit-encodings and adapt the query plan accordingly to have the system support queries that could not be executed by the DAX engines otherwise.

In Figure, 4 we depict an example to show how we implement this technique in *SiliconDB*. Here, the input query includes a semi-join between the *lineorder* and *part* tables of the *SSB-Benchmark*. The inputs to the operator include a bitmap result representing the selected tuples of the *part* relation and the bit-compressed *lo_pkey* foreign-key column of the *lineorder* table. Then, as to follow the implementation of a semi-join, the operator uses the values of *lo_pkey* to index into the bitmap vector and produces an output bit-result representing the tuples of the *lineorder* that satisfies the join condition. For the ease of representation, we assume that the DAX engines can support bitmap inputs

represented with at most 2-bits, but the *p_pkey* column has 8 distinct values and thus requires 3-bit encoding.

Now our goal in this example is to represent the values of the *p_pkey* column with an encoding of 2-bits and process the join accordingly so that the system can leverage the DAX engines. As we show in Fig. 4a, the idea is to create a split point over the original bit-encodings and represent the original columns as with two sub-columns that contain the high- and low-end bits of the actual encoding in respect to the applied split point.

After the *p_pkey* and *lo_pkey* columns are re-coded according to this process, we adapt the query plan to process the join in two steps represented as *join₁* and *join₂* in Fig. 4b. Both of these sub-joins share the use of *lo_pkey_{lower}* sub-column as one side of their inputs. However, they are required to use different portions of the re-coded *p_pkey_{lower}* as their bitmap inputs (*bitmap₁* and *bitmap₂*) in order to reflect it as a PK-column as required by the join operator (note the repeating pattern of the values).

After their execution, both *join₁* and *join₂* generate bit-vector results (*join_res₁* and *join_res₂*), but with some false-positives due to the fact that they were required to index into different portions of the original bitmap column. In order to eliminate these false-positives and correctly generate the final join result (*join_res*), here we use the re-coded *lo_pkey_{higher}* sub-column in the final step.

6. EXPERIMENTAL EVALUATION

In this section, we report the results of our experimental evaluation of the techniques presented in this paper. The main goal of this evaluation is to: (1) compare the proposed fine-grained execution model based on morsels over two alternative techniques that are typically used for heterogeneous environments, (2) show the significance of the adaptive scheduling model, (3) demonstrate the promise of new query optimization ideas and (4) highlight the benefits of some other optimization techniques we discussed throughout the paper.

Setup: The prototype of *SiliconDB* is implemented in C++ and compiled using GCC 4.8.2. All experiments have been executed on a machine at Oracle with 128GB of RAM and one SPARC M7 processor (32 cores, 8 DAX engines) running Solaris 11.3 as the operating system. In the SPARC M7 processor, 4 cores and 1 DAX engine share one 8 MB L3 cache. Moreover, each of the cores supports 8 strands (which are similar to hyper-threads) resulting in a total of 256 hardware threads for one server with 32 cores. In most of our experiments, we were limited to use only 8 cores and 2 DAX engines since we only had access to parts of a remote machine at Oracle. However, we were able to execute some

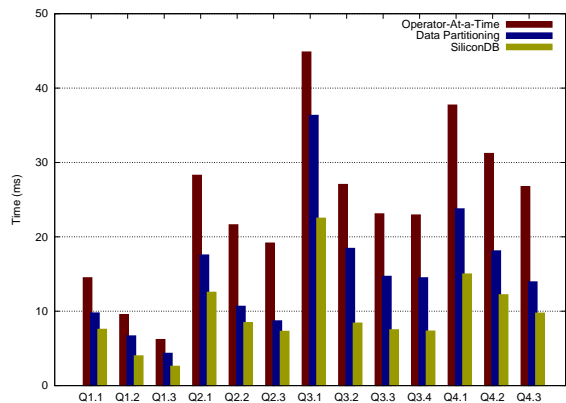


Figure 5: Star-Schema Benchmark

experiments on the full machine and we indicate the setup we used for each experiment.

6.1 Exp1: Star Schema Benchmark

In this experiment, we present the benefits of *SiliconDB*'s query processing model and compare our execution scheme against two different execution models used in the past for heterogeneous environments based on CPUs and GPUs (sitting at the end of the PCI bus).

(1) *Operator At-a-Time* [4, 6]: In this scheme, the complete execution of each operator is pushed to either cores or accelerators, depending on which resource has better execution performance on the specific operator. This model represents a coarse-grained processing scheme, commonly seen in heterogeneous co-processor environments. It is also important to note that the processing model that the execution engine of an Oracle database follows in the Sparc M7 environment aligns closest with this scheme as it pushes all scans and semi-joins to the DAX Engines.

(2) *Data Partitioning* [10]: This model partitions the input of operators across cores and accelerators before the processing of each operator starts. This model lets both cores and accelerators process some operators simultaneously, but the amount of work is statically assigned to cores/accelerators. Moreover, the processing carries in a blocking manner that avoids pipelining between operators completely.

In order to provide fair comparisons, we implemented both of these schemes in *SiliconDB* as a different mode of processing. Both for (1) and (2), we use our compilation approach to generate plans with sub-pipelines and ensured that all the work elements of a specific operator are consumed before the execution of a parent operator is started. Regarding the scheduling of work elements for the *operator-at-a-time* technique, we pushed all the work elements of the scan and semi-join operators into the accelerators, while the rest of the query plan is executed on normal cores. For the second technique, we pre-distributed the work elements of the scan and semi-join operators between cores and accelerators depending on their expected idealized performance ratio for processing these operators.

In order to show the effectiveness of *SiliconDB*'s query processing model over these two techniques, we ran queries from the complete SSB benchmark with a scale factor of 10 and provide the total runtime results as in Figure 5.

As expected, the *operator-at-a-time* technique shows the worst performance due to its coarse grained processing model, which causes cores to stay idle while accelerators are pro-

cessing work elements or vice versa. The *data-partitioning* technique improves over the *operator-at-a-time* model because it can support the co-processing of operators between cores and accelerators, thus providing improved utilization. *SiliconDB* outperforms both of these techniques by up-to 3.2x improvements over the *operator-at-a-time* model and a 2.3x speed-up over the *data-partitioning* technique. This is mainly due to the *SiliconDB*'s dynamic query scheduling model, which can adapt to run-time conditions and provide better utilization of all compute-resources. In order to better point out this fact, we present the details of resource utilization during the execution of Query 1.1 in Figure 6. This clearly shows that our scheduling strategy better utilizes all available compute resources and thus minimizes the overall runtime.

6.2 Exp. 2: Adaptive Scheduling

In this experiment, we show the significance of *SiliconDB*'s adaptive scheduling model. Our main goal is to show that finding an optimal *q-size* is necessary for an optimal query processing scheme and the value of the *q-size* should be adjusted dynamically during query execution.

We first present the result of an experiment in Figure 7 showing that the optimal *q-size* depends on different characteristics of the query workload. Here, the y-axis presents the total runtime of a scan operation that applies a filter on a compressed column with a specific bit-encoding; i.e. each line represents a different case, 8-bit and 16-bit encoding respectively. We report the results for each possible *q-size* (from 1 to 16), and the results depict how these values effect the total runtime. In both cases, we see the effect of *utilization* and *items-waiting* metrics as we discussed in Section 4.3 and that the optimal *q-size* value occurs when the system finds a sweet-spot between them. Also even more importantly, our cost model (shown in Figure 7 as well) is able to pick the optimal *q-size* value in both cases. The reason why different *q-size* are important is due to the different types of scan items processed by the system.

6.3 Exp. 3: Query Rewrites

We now show the effects of our rewrite heuristics for query optimization techniques using 8 cores and 2 DAX engines.

Exp. 3a - Join Re-Coding

In the first part, we executed a query which has a PK-FK join between the *lineorder* and *part* tables of the SSB benchmark. Implementing this operator as a semi-join requires to encode the smaller *part* relation's join-column with 18-bits of encoding as the bitmap of the semi-join operator. However, as we mentioned before, the DAX engines can support only bitmaps of up-to 16-bits, so we first execute the query with default settings which requires the joins to be handled only by CPU-cores.

Then, we applied our operator splitting approach, where we create four different small joins out of the single join operator and process the query in this manner. One disadvantage here is that the system now needs to process more joins instead of a single one, but the benefit is that we will be able to utilize more resources. The results of this rewrite are shown in Figure 8a in terms of their resource utilization and total run-time. We see that even if the re-coding approach does more work in overall, it is able to reduce the total run-time by utilizing more resources.

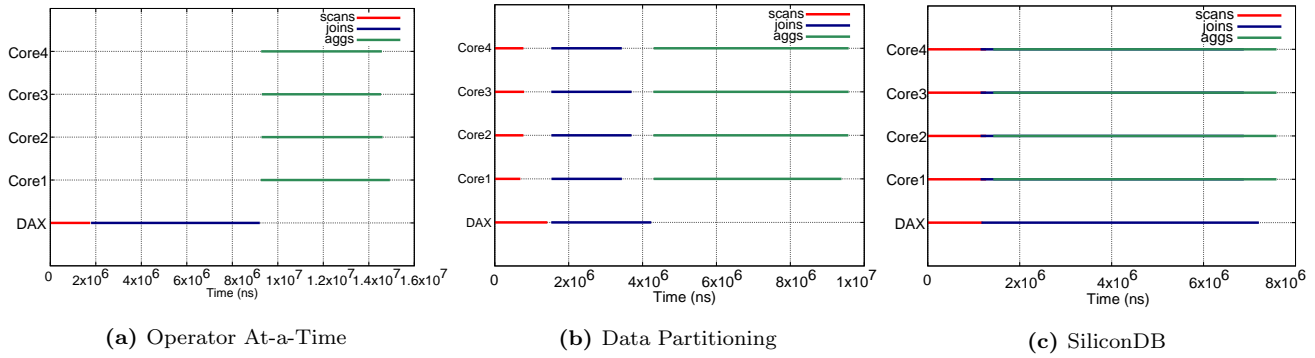


Figure 6: Resource Utilizations During SSB Query 1.1

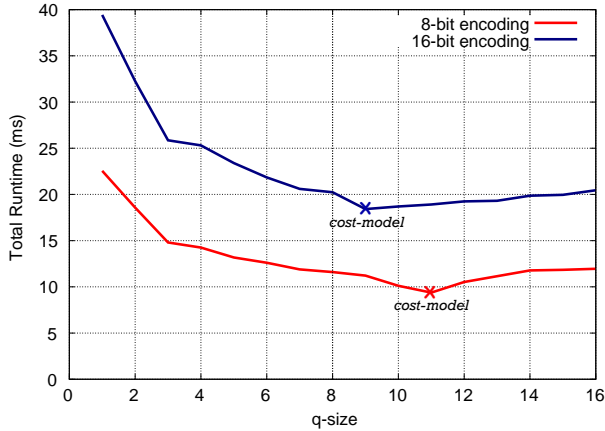


Figure 7: Optimal Queue Sizes

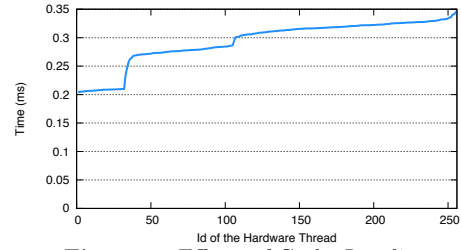


Figure 9: Effects of Cache-Locality

In this experiment, we use a fixed ratio of 0.5 to define how much of the output data the aggregation should materialize. Even if this might not be the most optimal setting, our main goal here is to show the promise and potential of query plan adaptation that leads to better resource utilization. The results of this experiment are shown in Figure 8b. Again, we can see that the rewrite reduces the overall runtime and increases the utilization of the DAX engines in the later phases of the query.

6.4 Exp. 4: Micro Benchmarks

Exp. 4a - Effects of Cache-Locality

In order to reveal the effects of cache locality, we executed the query presented in Figure 2. For this micro-benchmark we had access to the full processor and leveraged the environment using 256 different configurations.

While we used the same DAX engine for the scan sub-pipeline, for the aggregation we pinned the worker thread to a different strand at each execution. Since the Sparc M7 processor provides 32 cores each having 8 strands, thus in this microbenchmark we use 256 different executions. The output of the DAX scan is configured in a way that it fits in the L3 cache assigned to the DAX engine.

The results of the experiment are shown in Figure 9. The x -axis shows the strands we pinned the work handlers and the y -axis represents the average runtime the corresponding aggregation pipelines. We observe that the runtime of the aggregation when the software thread is pinned to one of the first 32 strands is significantly lower. The reason is that these strands are executed by the first 4 cores that share the same cache with the DAX engine. In our results, we can also see an additional increase in runtime when using the strands 128 to 256. The reason for this are NUMA effects that result from different latencies to access the two NUMA regions using an on-chip network on the SPARC M7 processor.

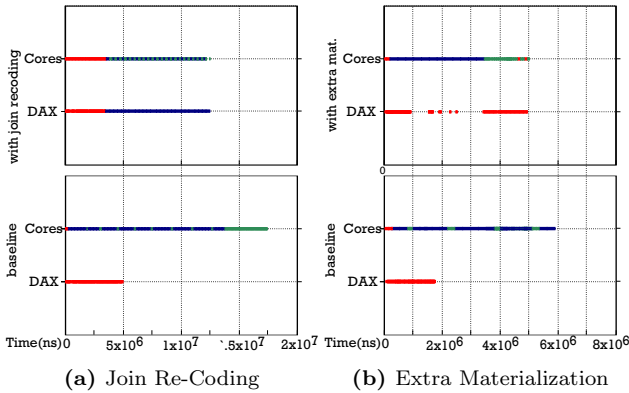


Figure 8: Rewriting Query Plans

Exp. 3b - Inserting Materialization Operators

For the second experiment, we applied the rewrite heuristic which adds additional materialization operators to better leverage DAX engines in the later phases of executing a query plan. To this end, we used the query from Figure 2 which has an aggregation on top of a scan, so the system would be able to use the DAX engines only at the start of the query. To show the effects of our suggested operator-insertion technique, we modified the query plan to have an additional having statement after the aggregation operator, so the system would first materialize some of the aggregation results and then need to apply another filtering to produce the final results.

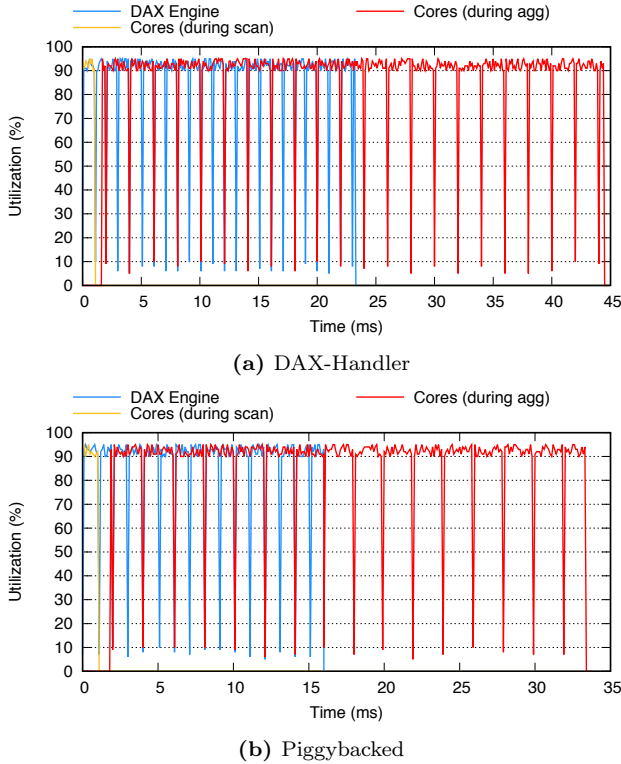


Figure 10: Efficiency of Handling Strategies

Exp. 4b - Efficiency of the Handling Model

In this experiment, we show the benefits of our execution strategies discussed in Section 4.1 when using separate and piggybacked DAX handlers. In order to show the utilization of cores/accelerator we used 4 cores and 1 DAX engine. As can be observed from Figure 10, while the utilization of both schemes is similar, the overall runtime of the piggybacked version is about 25% lower since all four cores can be used to process meaningful work instead of just handling the scheduling of DAX engine.

7. RELATED WORK

Query Processing on Co-Processors Environments:

There is a large body of work on how to leverage specialized co-processors for analytical database workloads (e.g., FPGAs [14, 15], GPUs [3, 7], etc.). The main limitation of existing co-processors has been the PCI bottleneck, therefore the main goal of solutions that integrate co-processors in an end-to-end manner into a DBMS has been to reduce the costs of overall communication [3, 6]. When using these schemes, the high speed-up rates reported for individual database operators often become negligible when considering the overall runtime that includes the communication costs between CPUs and the co-processors [5].

Some recent work has addressed coupled environments where the cores and the co-processor units are integrated on a single chip [11, 18, 20, 21]. For instance, in [18], to show the effects of using an integrated GPU, the authors propose specialized scan and aggregation operations. They were able to achieve a 3× performance boost compared to an architecture where the GPU sits at the end of the PCI bus, even though the integrated GPU has 4× lower computational power than a discrete model. Different from our paper, these approaches

focus on the integration of GPUs which can still execute arbitrary functions whereas *SiliconDB* can also efficiently integrate accelerators which provide a limited set of functions that brings new challenges as described in this paper. We believe that *SiliconDB* can also be efficiently used for these coupled CPU-GPU platforms as well. Demonstrating this claim with implementation and experimentation is beyond the scope of this paper, but we provide a thorough discussion in the following section.

In [10], authors provide a co-processing scheme for hash-joins in a similar environment. Their solution splits the execution of a hash-join into four steps and have the CPU and the GPU units to co-execute each one of these steps in a fine-grained manner. Before starting each step, they use a cost-model to find a ratio in order to split the execution between the CPU and GPU units. While this work shares the same high-level goal (i.e., utilizing all compute-resources effectively), our approach differs in various aspects. First, our focus is not to provide solutions only for specific operators but for the execution of a whole query pipeline that would optimally leverage all compute resources. Also, all the related work described above depend on static cost-model decisions that use low-level hardware parameters in order to distribute the workload between cores and co-processors. On the other hand, our scheduling strategies are designed to be adaptive at query runtime and avoid static decisions before the actual execution pipeline starts.

Fine-Grained Query Processing Models: As we mentioned throughout the paper, in [13], the authors propose a novel query processing model that splits the input of a query pipeline into equal sized partitions, called morsels, and then schedules these fine-grained elements on worker-threads that execute them in parallel. The main idea is to provide a scheduling model that is fully elastic at query run-time, so that the system can adapt its decisions accordingly. In this way, the system is able to utilize all CPU resources effectively by applying techniques such as work-stealing. While our query processing model builds on the ideas from this paper, it differs in the way we address the new challenges arising due to the characteristics of the emerging heterogeneous multi-core environments, as we described in Section 3 in more detail.

Also, our scheduling scheme has similarities with the approach used in QPipe [9]. While QPipe uses a similar queue-based concept to schedule work for different database operators (similar to our function-specific work queues), there are some important differences: First, the focus of the scheduling strategies in QPipe is on the sharing of data and work between queries, whereas our focus is to maximize the utilization of all cores and accelerators. Second, different processing units might share the same queue to enable dynamic scheduling decisions, which is a key aspect of *SiliconDB* that allows it to adapt to different co-processor architectures.

8. DISCUSSION

Since we develop and test our ideas specifically on the Sparc M7 platform, one could argue whether or not the techniques we implemented in this paper are applicable to other platforms, such as Intel Harp that combines an FPGA with normal cores, or other designs such as APUs that combine CPUs and GPUs.

In order to apply the design presented in this paper with *SiliconDB*, we require that the platform in question pro-

vides well-defined semantics regarding the characteristics of its accelerators. Such characteristics include the functionality these platforms can support (e.g., the type of operations they can execute) and their scheduling mechanisms; i.e., if they implement active or passive models (as in the case with DAX engines). Here the former is particularly important for *SiliconDB*'s query compilation and execution models in order to properly generate query sub-pipelines and incorporate them with function-specific work-queues, while the latter is crucial for the adaptation of the query scheduler in order to handle the scheduling of accelerators. Next, we first discuss these points from the perspective of environments that provide re-configurable operations, specifically FPGAs such as Intel Harp, and finish with some thoughts on APUs.

One challenge to adapt *SiliconDB* for FPGAs would be if the operator functionality provided by the FPGA would change at query run-time. However, due to the high reconfiguration costs of FPGAs this is usually not the case for query processing systems. Thus, the capability of FPGAs are typically considered to be static at query runtime, as in an ASIC-based architecture. In case new architectures provide reconfiguration options with negligible costs, this would create an interesting venue for future work, where the system could decide to change the functionality of the underlying FPGAs depending on the state of the query workload. For instance, if the FPGAs are configured to support scans initially, but the system starts to serve more joins than there are scans, it might decide to re-configure the hardware to support joins to better leverage the hardware space.

For APUs, the situation is slightly different since the GPU can provide kernels for all database operations. However, we think that the scheduling model of *SiliconDB* would still be beneficial in order to leverage all compute resources (CPU cores and GPUs) and adapt to the different speeds.

9. CONCLUSIONS

We presented *SiliconDB* that implements novel parallel query execution strategies for heterogeneous environments that combine cores with specialized ASIC-based accelerators on the same socket. We designed, implemented and experimentally tested our proposals on the Sparc M7 processor.

Our solution respects the functional and execution-level limitations of accelerator units, and aims to maximize the collective utilization of all processing elements in the system in an attempt to improve end-to-end performance. To this end, we propose cost-based scheduling models, and study query plan modifications that improve system utilization at the expense of small increases in total execution costs. Based on the SSB benchmarks, we showed that *SiliconDB* provides a speed-up of up to $2\times$ compared to alternative state-of-the-art parallel execution strategies that have been developed for heterogeneous environments.

10. ACKNOWLEDGMENTS

This research is supported in part by NSF IIS-1526639, and in part by a gift from Oracle.

11. REFERENCES

- [1] M.-C. Albutiu, A. Kemper, and T. Neumann. Massively parallel sort-merge joins in main memory multi-core database systems. *PVLDB*, 5(10):1064–1075, 2012.
- [2] C. Balkesen, G. Alonso, J. Teubner, and M. T. Özsu. Multi-core, main-memory joins: Sort vs. hash revisited. *PVLDB*, 7(1):85–96, 2013.
- [3] S. Breß. The design and implementation of cogadb: A column-oriented gpu-accelerated DBMS. *Datenbank-Spektrum*, 2014.
- [4] S. Breß et al. Automatic selection of processing units for coprocessing in databases. In *ADBIS*, 2012.
- [5] S. Breß et al. Robust query processing in co-processor-accelerated databases. In *SIGMOD*, pages 1891–1906, 2016.
- [6] S. Breß, B. Köcher, M. Heimel, V. Markl, M. Saecker, and G. Saake. Ocelot/hype: Optimized data processing on heterogeneous hardware. *PVLDB*, 7(13):1609–1612, 2014.
- [7] H. Funke, S. Breß, S. Noll, V. Markl, and J. Teubner. Pipelined query processing in coprocessor environments. *SIGMOD*, 2018.
- [8] N. Hardavellas, M. Ferdman, B. Falsafi, and A. Ailamaki. Toward dark silicon in servers. *IEEE Micro*, 2011.
- [9] Harizopoulos et al. QPipe: A Simultaneously Pipelined Relational Query Engine. In *SIGMOD*, 2005.
- [10] J. He, M. Lu, and B. He. Revisiting co-processing for hash joins on the coupled cpu-gpu architecture. *PVLDB*, 6(10):889–900, 2013.
- [11] J. He, S. Zhang, and B. He. In-cache query co-processing on coupled cpu-gpu architectures. *PVLDB*, 8(4):329–340, 2014.
- [12] G. Konstantinidis et al. Sparc m7: A 20 nm 32-core 64 mb l3 cache processor. *IEEE Journal of Solid-State Circuits*, 51:1–13, 2015.
- [13] V. Leis et al. Morsel-driven parallelism: a numa-aware query evaluation framework for the many-core age. In *SIGMOD*, 2014.
- [14] R. Mueller, J. Teubner, and G. Alonso. Data processing on fpgas. *PVLDB*, 2(1):910–921, 2009.
- [15] R. Mueller, J. Teubner, and G. Alonso. Streams on wires: A query compiler for fpgas. *PVLDB*, 2(1):229–240, 2009.
- [16] T. Neumann. Efficiently compiling efficient query plans for modern hardware. *PVLDB*, 4(9):539–550, 2011.
- [17] N. Oliver et al. A reconfigurable computing system based on a cache-coherent fabric. In *ReConFig*, 2011.
- [18] J. Power et al. Toward gpus being mainstream in analytic processing: An initial argument using simple scan-aggregate queries. In *DaMoN*, pages 1–8, 2015.
- [19] I. Psaroudakis, T. Scheuer, N. May, A. Sellami, and A. Ailamaki. Adaptive numa-aware data placement and task scheduling for analytical workloads in main-memory column-stores. *PVLDB*, 10(2):37–48, 2016.
- [20] S. Tang et al. Elastic multi-resource fairness: balancing fairness and efficiency in coupled CPU-GPU architectures. In *SC 2016, Salt Lake City, UT, USA, November 13-18, 2016*, pages 875–886, 2016.
- [21] K. Zhang et al. DIDO: dynamic pipelines for in-memory key-value stores on coupled CPU-GPU architectures. In *ICDE*, pages 671–682, 2017.