

Fast Computation of Database Operations using Graphics Processors

Naga K. Govindaraju Brandon Lloyd Wei Wang Ming Lin Dinesh Manocha
University of North Carolina at Chapel Hill
{naga,blloyd,weiwang,lin,dm}@cs.unc.edu
<http://gamma.cs.unc.edu/DB>

ABSTRACT

We present new algorithms on commodity graphics processors to perform fast computation of several common database operations. Specifically, we consider operations such as conjunctive selections, aggregations, and semi-linear queries, which are essential computational components of typical database, data warehousing, and data mining applications. While graphics processing units (GPUs) have been designed for fast display of geometric primitives, we utilize the inherent pipelining and parallelism, single instruction and multiple data (SIMD) capabilities, and vector processing functionality of GPUs, for evaluating boolean predicate combinations and semi-linear queries on attributes and executing database operations efficiently. Our algorithms take into account some of the limitations of the programming model of current GPUs and perform no data rearrangements. Our algorithms have been implemented on a programmable GPU (e.g. NVIDIA's GeForce FX 5900) and applied to databases consisting of up to a million records. We have compared their performance with an optimized implementation of CPU-based algorithms. Our experiments indicate that the graphics processor available on commodity computer systems is an effective co-processor for performing database operations.

Keywords: graphics processor, query optimization, selection query, aggregation, selectivity analysis, semi-linear query.

1. INTRODUCTION

As database technology becomes pervasive, Database Management Systems (DBMSs) have been deployed in a wide variety of applications. The rapid growth of data volume for the past decades has intensified the need for high-speed database management systems. Many database queries and, more recently, data warehousing and data mining applications, are very data- and computation-intensive and therefore demand high processing power. Researchers have actively sought to design and develop architectures and algorithms for faster query execution. Special attention has been given

to increasing the performance of selection, aggregation, and join operations on large databases. These operations are widely used as fundamental primitives to build complex database queries and to support on-line analytic processing (OLAP) and data mining procedures. The efficiency of these operations has a significant impact on the performance of a database system.

As the current trend of database architecture moves from disk-based system towards main-memory databases, applications have become increasingly computation- and memory- bound. Recent work [2, 18] investigating the processor and memory behaviors of current DBMSs has demonstrated the significant increase of query execution time due to memory stalls (due to data miss and instruction miss), branch mispredictions, and resource stalls (due to instruction dependencies and hardware specific characteristics). Increased attention has been given to redesigning traditional database algorithms to fully utilize the available architectural features and to exploit parallel execution possibilities, minimize memory and resource stalls, and reduce branch mispredictions [1, 4, 17, 21, 28, 29, 31, 34].

1.1 Graphics Processing Units

In this paper, we exploit the computational power of graphics processing units (GPUs) for database operations. In the last decade, high-performance 3D graphics hardware has become as ubiquitous as floating-point hardware. Graphics processors are now a part of almost every personal computer, game console, or workstation. In fact, the two major computational components of a desktop computer system are its main central processing unit (CPU) and its (GPU). While CPUs are used for general purpose computation, GPUs have been primarily designed for transforming, rendering, and texturing geometric primitives, such as triangles. The driving application of GPUs has been fast rendering for visual simulation, virtual reality, and computer gaming.

GPUs are increasingly being used as co-processors to CPUs. GPUs are extremely fast and are capable of processing tens of millions of graphics primitives per second. The peak performance of GPUs has been increasing at the rate of 2.5 – 3.0 times a year, much faster than the Moore's law for CPUs. At this rate, the GPU's peak performance may move into the teraflop range by 2005 [16]. Most of this performance arises from multiple processing units and stream processing. The GPU treats the vertices and pixels constituting graphics primitives as streams. Multiple *vertex* and *pixel processing engines* on a GPU are connected via data flows. These processing engines perform simple operations in parallel.

Recently, GPUs have become programmable, allowing a user to

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 200X ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

write *fragment programs* to be executed on pixel processing engines. The pixel processing engines have direct access to texture memory and can perform vector operations with single-precision IEEE floating point arithmetic. These capabilities have been successfully exploited for many geometric and scientific applications. As graphics hardware becomes increasingly programmable and powerful, the roles of CPUs and GPUs in computing are being redefined.

1.2 Main Contributions

In this paper, we present novel algorithms for fast computation of database operations on GPUs. The operations include predicates, boolean combinations, and aggregations. We utilize the SIMD capabilities of pixel processing engines within a GPU and to perform these operations efficiently. We have used these algorithms for selection queries on one or more attributes and generic aggregation queries including selectivity analysis on large databases.

Our algorithms take into account some of the limitations of the programming model of GPUs, which prohibit data rearrangement. We present novel algorithms for performing multi-attribute comparisons, semi-linear queries, range queries, computing the *k*th largest number, and other aggregates. These algorithms have been implemented using fragment programs and have been applied to large databases composed of up to a million records. The performance of these algorithms depends on the instruction sets available for fragment programs, the number of fragment processors, and the underlying clock rate of the GPU. We also perform a preliminary comparison between GPU-based algorithms running on a NVIDIA GeForceFX 5900 Ultra graphics processor and optimized, CPU-based algorithms running on dual 2.8 GHz Intel Xeon processors.

We show that algorithms for semi-linear and selection queries map very well to GPUs and we are able to obtain significant performance improvement over CPU-based implementations. At the same time, the algorithms for aggregates obtain 2–4 times speedup over CPU-based implementations. Overall, the GPU can be used as an effective co-processor for many database operations.

1.3 Organization

The rest of the paper is organized as follows. We briefly survey related work on database operations and use of GPUs for geometric and scientific computing in Section 2. We give an overview of the graphics architectural pipeline in Section 3. We present algorithms for database operations including predicates, boolean combinations, and aggregations in Section 4. We describe their implementation in Section 5 and compare their performance with optimized CPU-based implementations. We analyze the performance in Section 6 and outline the cases where GPU-based algorithms can offer considerable gain over CPU-based algorithms.

2. RELATED WORK

In this section, we highlight the related research in main-memory database operations and general purpose computation using GPUs.

2.1 Hardware Accelerated Database Operations

Many acceleration techniques have been proposed for database operations. Alilamaki et al. analyzed the execution time of commercial DBMSs [2] and observed that almost half of the time is spent in stalls. This indicates that the performance of a DBMS can be significantly improved by reducing stalls.

Meki and Kambayashi used a vector processor to accelerate the execution of relational database operations including selection, pro-

jection, and join [21]. To utilize the efficiency of pipelining and parallelism that a vector processor provides, the implementation of each operation was redesigned to increase the vectorization rate and the vector length. The limitation of using a vector processor is that the load-store instruction can have high latency [34].

Modern CPUs have SIMD instructions that allow a single basic operation to be performed on multiple data elements in parallel. Zhu and Ross presented the SIMD implementation of many important database operations including sequential scans, aggregation, indexed searches, and joins [34]. Considerable performance gains were achieved by exploiting the inherent parallelism of SIMD instructions and reducing branch mispredictions.

Recently, Sun et al. present the use of graphics processors for spatial selections and joins [32]. They use color blending capabilities available on graphics processors to test if two polygons intersect in screen-space. Their experiments on graphics processors indicate a speedup of nearly 5 times on intersection joins and within-distance joins when compared against their software implementation. The technique focuses on pruning intersections between triangles based on their 2D overlap and is quite conservative.

2.2 General-Purpose Computing Using GPUs

In theory, GPUs are capable of performing any computation that can be mapped to the stream-computing model. This model has been exploited for ray-tracing [26], global illumination [27] and geometric computations [19].

The programming model of GPUs is somewhat limited, mainly due to the lack of random access writes. This limitation makes it more difficult to implement many data structures and common algorithms such as sorting. Purcell et al. [27] present an implementation of bitonic merge sort, where the output routing from one step to another is known in advance. The algorithm is implemented as a fragment program and each stage of the sorting algorithm is performed as one rendering pass. It can be quite slow for database operations on large databases.

GPUs have been used for performing many discretized geometric computations [19]. These include using stencil buffer hardware for interference computations [30], using depth-buffer hardware to perform distance field and proximity computations [13], and visibility queries for interactive walkthroughs and shadow generation [10].

The high throughput of the fragment processors and direct access to texture memory makes the fragment engines a powerful computation engine for certain numerical algorithms, including dense matrix-matrix multiplication [15], general purpose vector processing [33], visual simulation based on coupled-map lattices [11], linear algebra operations [14], sparse matrix solvers for conjugate gradient and multigrid [3], a multigrid solver for boundary value problems [9], etc.

3. OVERVIEW

In this section, we introduce the basic functionality available on GPUs and give an overview of the architectural pipeline. More details are given in [7].

3.1 Graphics Pipeline

A GPU is designed to rapidly transform the geometric description of a scene into the pixels on the screen that constitute a final image. Pixels are stored on the graphics card in a *frame-buffer*. The frame buffer is conceptually divided into three buffers according to the different values stored at each pixel:

- **Color Buffer:** Stores the color components of each pixel in the frame-buffer. Color is typically divided into red, green, and blue channels with an alpha channel that is used for blending effects.
- **Depth Buffer:** Stores a depth value associated with each pixel. The depth is used to determine surface visibility.
- **Stencil Buffer:** Stores a stencil value for each pixel. It is called the stencil buffer because it is typically used to enable/disable writing to portions of the frame-buffer.

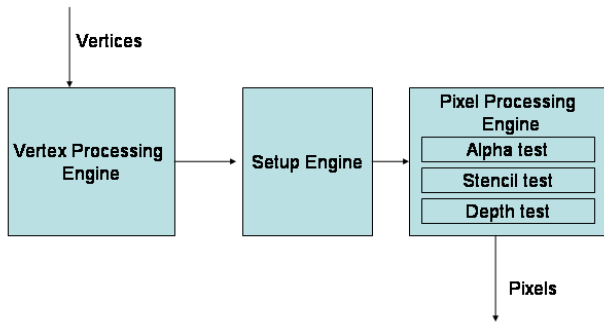


Figure 2: Graphics architectural pipeline overview: This figure shows the various units of a modern GPU. Each unit is designed to perform a specific operation efficiently.

The transformation of geometric primitives (points, lines, triangles, etc.) to pixels is performed by the graphics pipeline, consisting of several functional units, each optimized to perform a specific operation. Fig 2 shows the various stages involved in rendering a primitive.

- **Vertex Processing Engine:** This unit receives vertices as input and transforms them to points on the screen.
- **Setup Engine:** Transformed vertex data is streamed to the setup engine which generates slope and initial value information for color, depth, and other parameters associated with the primitive vertices. This information is used during rasterization to construct *fragments* at each pixel location covered by the primitive.
- **Pixel Processing Engines:** Before the fragments are written as pixels to the frame buffer, they pass through the pixel processing engines or *fragment processors*. A series of tests can be used to discard a fragment before it is written to the frame buffer. Each test performs a comparison using a user-specified relational operator and discards the fragment if the test fails.
 - **Alpha test:** Compares a fragment’s alpha value to a user-specified reference value.
 - **Stencil test:** Compares the stencil value of a fragment’s corresponding pixel with a user-specified reference value.
 - **Depth test:** Compares the depth value of a fragment to the depth value of the corresponding pixel in the frame buffer.

The relational operator can be any of the following : =, <, >, ≤, ≥, and ≠. In addition, there are two operators, *never* and *always*, that do not require a reference value.

Current generation of GPUs have a pixel processing engine that is programmable. The user can supply a custom *fragment program* to be executed on each fragment. For example, a fragment program can compute the alpha value of a fragment as a complex function of the fragment’s other color components or its depth.

3.2 Visibility and Occlusion Queries

Current GPUs can perform visibility and occlusion queries [24]. When a primitive is rasterized it is converted to fragments. Some of these fragments may or may not be written to pixels in the frame buffer depending on whether they pass the alpha, stencil and depth tests. An *occlusion query* returns the *pixel pass count*, the number of fragments that pass the different tests. We use these queries to perform aggregation computations (see Section 4).

3.3 Data Representation on the GPUs

Our goal is to utilize the inherent parallelism and vector processing capabilities of the GPUs for database operations. A key aspect is the underlying data representation.

Data is stored on the GPU as *textures*. Textures are 2D arrays of values. They are usually used to apply images to rendered surfaces. They may contain multiple channels. For example, an RGBA texture has four color channels - red, blue, green and alpha. A number of different data formats can be used for textures including 8-bit bytes, 16-bit integers, and floating point. We store data in textures in the floating-point format. This format can precisely represent integers up to 24 bits.

To perform computations on the values stored in a texture we render a single quadrilateral that covers the window. The texture is applied to the quadrilateral such that the individual elements of the texture, *texels*, line up with the pixels in the frame-buffer. Rendering the textured quadrilateral causes a fragment to be generated for every data value in the texture. Fragment programs are used to perform computations using the data value from the texture. Then the alpha, stencil, and depth tests can be used to perform comparisons.

3.4 Stencil Tests

Graphics processors use stencil tests to restrict computations to a portion of of the frame-buffer based on the value in the stencil buffer. Abstractly, we can consider the stencil buffer as a mask on the screen. Each fragment that enters the pixel processing engine corresponds to a pixel in the frame-buffer. The stencil test compares the stencil value of a fragment’s corresponding pixel against a reference value. Fragments that fail the comparison operation are rejected from the rasterization pipeline.

Stencil operations can modify on the stencil value of a fragment’s corresponding pixel. Examples of such stencil operations are

- **KEEP:** Keep the stencil value in stencil buffer. We use this operation if we do not want to modify the stencil value.
- **INCR:** Increment the stencil value by one.
- **DECR:** Decrement the stencil value by one.
- **ZERO:** Set the stencil value to zero.
- **REPLACE:** Set the stencil value to the reference value.
- **INVERT:** Bitwise invert the stencil value.

For each fragment there are three possible outcomes based on the stencil and depth tests. Based on the outcome of the tests, the corresponding stencil operation is performed:

- **Op1**: when a fragment fails the stencil test,
- **Op2**: when a fragment passes the stencil test and fails the depth test,
- **Op3**: when the fragment passes the stencil and depth tests.

We illustrate these operations with the following pseudo-code for the STENCILOP routine:

```

STENCILOP( OP1, OP2, OP3)
if (stencil test) /* perform stencil test */
/* fragment passed stencil test */
if(depth test) /* perform depth test */
/* fragment passed stencil and depth test */
perform Op3 on stencil value
else
/* fragment passed stencil test */
/* but failed depth test */
perform Op2 on stencil value
end if
else
/* fragment failed stencil test */
perform Op1 on stencil value
end if

```

4. BASIC DATABASE OPERATIONS USING GPUS

In this section, we give a brief overview of basic database operations that are performed efficiently on a GPU. Given a relational table T of m attributes (a_1, a_2, \dots, a_m) , a basic SQL query is in the form of

```

SELECT A
FROM T
WHERE C

```

where A may be a list of attributes or aggregations (SUM, COUNT, AVG, MIN, MAX) defined on individual attributes, and C is a boolean combination (using AND, OR, EXIST, NOT EXIST) of *predicates* that have the form $a_i \text{ op } a_j$ or $a_i \text{ op } \text{constant}$. The operator **op** may be any of the following: =, \neq , >, \geq , <, \leq . In essence, queries specified in this form involve three categories of basic operations: predicates, boolean combinations, and aggregations. Our goal is to design efficient algorithms to perform these operations using graphics processors.

- *Predicates*: Predicates in the form of $a_i \text{ op } \text{constant}$ can be evaluated via the depth test and stencil test. The comparison between two attributes, $a_i \text{ op } a_j$, can be transformed into a semi-linear query $a_i - a_j \text{ op } 0$, which can be executed on the GPUs.
- *Boolean combinations*: A boolean combination of predicates can always be rewritten in a conjunctive normal form (CNF). The stencil test can be used repeatedly to evaluate a series of logical operators with the intermediate results stored in the stencil buffer.
- *Aggregations*: This category includes simple operations such as COUNT, SUM, AVG, MIN, MAX, all of which can be implemented using the counting capability of the occlusion queries on GPUs.

In order for GPUs to perform these operations on a relational table, we store the attributes of each record in multiple channels of a single texel, or the same texel location in multiple textures.

4.1 Predicate Evaluation

In this section, we present novel GPU-based algorithms for performing comparisons as well as the semi-linear queries.

4.1.1 Comparison between an Attribute and a Constant

We can implement a comparison between an attribute “ tex ” and a constant “ d ” by using the depth test that utilizes graphics hardware. The stencil buffer can be configured to store the result of the depth test. This is important not only for evaluating a single comparison but also for constructing more complex boolean combinations of multiple predicates.

In order to use the depth test for performing comparisons, attribute values need to be stored in the depth buffer. We use a simple fragment program on each pixel of the screen. The fragment program copies the attribute values from the texture memory to the depth buffer.

A comparison operation against a depth value d is implemented by rendering a screen filling quadrilateral with depth d . In this operation, the rasterization hardware uses the comparison function to test each attribute value stored in the depth buffer against d . The comparison function is specified using the depth function. Routine 4.1 describes the pseudo-code for our implementation.

```

COMPARE( tex, op, d )
1 CopyToDepth( tex )
2 set depth test function to op
3 RENDERQUAD( d )

CopyToDepth( tex )
1 set up fragment program
2 RenderTexturedQuad( tex )

```

ROUTINE 4.1: COMPARE compares the attribute values stored in texture tex against d using the comparison function op . CopyToDepth called on line 1 copies the attribute values in tex into the depth buffer. CopyToDepth uses a simple fragment program on each pixel of the screen to perform the copy operation. On line 2, the depth test is configured to use the comparison operator op . The function RENDERQUAD(d) called on line 3 generates a fragment at a specified depth d for each pixel on the screen. Rasterization hardware compares the fragment depth d against the attribute values in depth buffer using the operation op .

4.1.2 Comparison between Two Attributes

The comparison between two attributes, $a_i \text{ op } a_j$, can be transformed into a special semi-linear query $(a_i - a_j \text{ op } 0)$, which can be performed very efficiently using the vector processors on GPUs. Here, we propose a fast algorithm that can perform any general semi-linear queries on GPUs.

Semi-linear Queries on GPUs

Applications encountered in Geographical Information Systems (GIS), geometric modeling, and spatial databases define geometric data objects as linear inequalities of the attributes in a relational database [25]. Such geometric data objects are called semi-linear sets. GPUs are capable of fast computation of semi-linear sets. A linear combination of m attributes is represented as:

$$\sum_{i=1}^{i=m} s_i \cdot a_i$$

where s_i are the scalar multipliers and a_i are attributes of a record in the database. The above expression can be considered as a dot product of two vectors s and a where $s = (s_1, s_2, \dots, s_m)$ and $a = (a_1, a_2, \dots, a_m)$.

```

SEMILINEAR( tex, s, op, b )
1 enable fragment program SEMILINEARFP(s, b)
2 RenderTexturedQuad( tex )

```

```

SEMILINEARFP( s, op, b )
1 a = value from tex
2 if dot( s, a ) op b
3 discard fragment

```

ROUTINE 4.2: SEMILINEAR computes the semi-linear query by performing linear combination of attribute values in *tex* and scalar constants in *s*. Using the operator *op*, it compares the the scalar value due to linear combination with *b*. To perform this operation, we render a screen filling quad and generate fragments on which the semi-linear query is executed. For each fragment, a fragment program SEMILINEARFP discards fragments that fail the query.

SEMILINEAR computes the semi-linear query:

$$(s \cdot a) \text{ op } b$$

where *op* is a comparison operator and *b* is a scalar constant. The attributes a_i are stored in separate channels in the texture *tex*. There is a limit of four channels per texture. Longer vectors can be split into multiple textures, each with four components. The fragment program SEMILINEARFP() performs the dot product of a texel from *tex* with *s* and compares the result to *b*. It discards the fragment if the comparison fails. Line 2 renders a textured quadrilateral using the fragment program. SEMILINEAR maps very well to the parallel pixel processing as well as vector processing capabilities available on the GPUs. This algorithm can also be extended to evaluate polynomial queries.

```

EVALCNF( A )
1 Clear Stencil to 1.
2 For each of  $A_i, i = 1, \dots, k$ 
3 do
4   if ( mod(i, 2) ) /* valid stencil value is 1 */
5     StencilTest to pass if stencil value is equal to 1
6     StencilOp( KEEP, KEEP, INCR )
7   else /* valid stencil value is 2 */
8     StencilTest to pass if stencil value is equal to 2
9     StencilOp( KEEP, KEEP, DECR )
10  endif
11  For each  $B_j^i, j = 1, \dots, m_i$ 
12  do
13    Perform  $B_j^i$  using COMPARE
14  end for
15  if ( mod(i, 2) ) /* valid stencil value is 2 */
16    if a stencil value on screen is 1, replace it with 0
17  else /* valid stencil value is 1 */
18    if a stencil value on screen is 2, replace it with 0
19  endif
20 end for

```

ROUTINE 4.3: EVALCNF is used to evaluate a CNF expression. Initially, the stencil is initialized to 1. This is used for performing TRUE AND A_1 . While evaluating each formula A_i , Line 4 sets the appropriate stencil test and stencil operations based on whether *i* is even or odd. If *i* is even, valid portions on screen have stencil value 2. Otherwise, valid portions have stencil value 1. Lines 11 – 14 invalidate portions on screen that satisfy $(A_1 \wedge A_2 \wedge \dots \wedge A_{i-1})$ and fail $(A_1 \wedge A_2 \wedge \dots \wedge A_i)$. Lines 15 – 19 compute the disjunction of B_j^i for each predicate A_i . At the end of line 19, valid portions on screen have stencil value 2 if *i* is odd and 0, otherwise. At the end of the line 20, records corresponding to non-zero stencil values satisfy *A*.

4.2 Boolean Combination

Complex boolean combinations are often formed by combining simple predicates with the logical operators AND, OR, NOT. In

these cases, the stencil operation is specified to store the result of a predicate. We use the function STENCILOP (as defined in Section 3.4) to initialize the appropriate stencil operation to store the result in stencil buffer.

Our algorithm evaluates a boolean expression represented as a CNF expression. Furthermore, we assume that the CNF expression has no NOT operators. If a simple predicate in this expression has a NOT operator, we can invert the comparison operation and eliminate the NOT operator. A CNF expression is represented as $A_1 \wedge A_2 \wedge \dots \wedge A_k$ where each A_i is represented as $B_1^i \vee B_2^i \vee \dots \vee B_{m_i}^i$. Each $B_j^i, j = 1, 2, \dots, m_i$ is a simple predicate. We first present an algorithm to evaluate the expression $A_1 \wedge A_2$. It can be easily extended to compute the entire CNF expression.

Suppose we compute A_1 (computed as TRUE AND A_1). We set the stencil for valid portions of A_1 to a value $x, 0 < x \leq 2$. We set the remaining stencil values to 0. We present a case for $x = 1$. The case for $x = 2$ is similar.

While evaluating A_2 , we set the stencil test to pass if the stencil value in the stencil buffer is equal to 1. We set the stencil operation to STENCILOP(KEEP, KEEP, INCR). This stencil operation ensures that for each fragment that needs to be executed by the fragment processors, its corresponding stencil value is incremented iff the fragment passed the stencil test and depth test. We do not modify the stencil value if the fragment is not valid or if it is valid but fails the depth test. The stencil operation also ensures that the maximum stencil value for any pixel on the stencil buffer is at most 2.

We perform each of the comparisons B_j^2 for $j = 1, 2, \dots, m_2$ by using the COMPARE() routine. The output of this operation generates a stencil value 2 on the screen if and only if it passes the predicate $A_1 \wedge A_2$. Next, we reset the stencil values that are equal to 1 to 0. These records passed the expression A_1 , but failed A_2 . At the end of this operation, all records with stencil value 2 passed $A_1 \wedge A_2$ and remaining records have a stencil value 0.

We proceed to the next conjunction in the CNF and evaluate it similarly with valid stencil value of 2. However, we use the stencil operation STENCILOP(KEEP, KEEP, DECR). This generates a valid stencil value of 1 at the end of evaluation of $((A_1 \wedge A_2) \wedge A_3)$. We compute the entire CNF as shown in the pseudo-code routine 4.3. Our algorithm can be easily modified to handle a boolean expression represented as DNF.

Range Queries

A range query is a common database query expressed as a boolean combination of two simple predicates. If [*low, high*] is the range for which an attribute *x* is queried, we evaluate the expression $(x \geq \text{low})$ AND $(x \leq \text{high})$ as in Section 4.2. The pseudo-code for our algorithm RANGE is given in Routine 4.4. Note that the redundant COPYTODEPTH in the second call to COMPARE need not be performed.

```

Range( tex, low, high )
1 SetupStencil()
2 COMPARE( tex,  $\geq$ , low )
3 COMPARE( tex,  $\leq$ , high )

```

ROUTINE 4.4: SETUPSTENCIL is called on line 1 to enable selection using the stencil buffer. The first pass on line 2 selects the attributes in *tex* that are at least *low*. The stencil is set to 1 for these attributes and 0 for the other. Since the stencil test has been configured to only pass where the stencil value is set to 1, the second COMPARE on line 3 only runs on the attributes previously selected. The values that are selected by this second comparison are the values that fall within the range [*low,high*].

4.3 Aggregations

Several database operations aggregate attribute values satisfying a condition. On GPUs, we can perform these operations using *occlusion queries* to return the count of records satisfying some condition.

4.3.1 COUNT

Using an occlusion query to count the number of records satisfying some condition involves three steps:

1. Initialize the occlusion query
2. Perform boolean query
3. Read back the result of the occlusion query into COUNT

Histogram A histogram partitions the data values into bins. Histograms are widely used in several database applications like selectivity estimation. We compute a histogram by partitioning the data values using a balanced binary tree. Let b be the maximum number of bits used to represent a data value. We create a balanced tree consisting of 2^b possible values. Each node in the tree represents one of the b bits as shown. We ignore portions of the tree that do not contain any of the data values.

4.3.2 MIN and MAX

The query to find the minimum or maximum value of an attribute is a degenerate case of the k th largest number. Here, we present an algorithm to generate the k th largest number.

k -th Largest Number

Computing the k -th largest number occurs frequently in several applications. We can utilize expected linear time selection algorithms such as QUICKSELECT [12] to compute the k -th largest number. Most of these algorithms require data rearrangement, which is extremely expensive on current GPUs because there is no functionality for data writes to arbitrary locations. Also, these algorithms require evaluation of conditionals and may lead to branch mispredictions on the CPU. We present a GPU-based algorithm that does not require data rearrangement. In addition, our algorithm exhibits SIMD characteristics that exploit the inherent parallelism available on the GPUs.

Our algorithm utilizes the binary data representation to compute the k -th largest value in time that is linear in the number of bits.

```

KthLargest( tex, k )
1 b_max = maximum number of bits in the values in tex
2 x = 0
3 for i = b_max-1 downto 0
4   count = Compare( tex, ≥, x + 2i )
5   if count > k - 1
6     x = x + 2i
7 return x

```

ROUTINE 4.5: KTHLARGEST computes the k -th largest attribute value in texture tex . It uses b_max passes starting from the MSB to compute the k -th largest number. During a pass i , it determines the i -th bit of the k -th largest number. At the end of b_max passes, it computes the k -th largest number in x .

The pseudocode for our algorithm KTHLARGEST is shown in routine 4.5. KTHLARGEST constructs in x the value of the k -th largest number one bit at a time starting with the most significant bit (MSB), b_max-1 . As an invariant, the value of x is maintained as less than or equal to the k -th largest value. Line 4 counts the number of values that are greater than or equal to $x + 2^i$, the tentative value

of x with the i th bit set. This count is used to decide whether to set the bit in x according to the following lemma:

Lemma 1: Let v_k be the k -th largest number in a set of values. Let $count$ be the number of values greater than or equal to a given value m .

- if $count > k - 1 : m \leq v_k$
- if $count \leq (k - 1) : m > v_k$

Proof: Trivial.

If $count > k - 1$ then the tentative value of x is smaller than the k -th largest number. In this case, we set x to the tentative value on line 6. Otherwise the tentative value is too large so we leave x unchanged. At the end of line 6, if the loop iteration is i , the first i bits from MSB of x and v_k are the same. After the last iteration of the loop, x has the value of the k -th largest number. The algorithm for the k -th smallest number is the same, except that the comparison in line 5 is inverted.

Another way of looking at this algorithm is as a variant of a binary search on an interval. Initially the interval is $[0, 2^{b_max})$, the range of numbers representable with b_max bits. We count how many values are in the upper half of the interval. If this value is greater than k then we add the value of the center of the interval to x and recurse on the upper half. Otherwise we recurse on the lower half and leave x unchanged. The recursion terminates when we arrive at an interval of width 1.

4.3.3 SUM and AVG

An accumulator is used to sum a set of data values. One way to implement an accumulator on current GPUs is to use a *mipmap* of a floating point texture. Mipmaps are multi-resolution textures consisting of multiple levels. The highest level of the mipmap contains the average of all the values in the lowest level, from which it is possible to recover the sum by multiplying by the number of values. A fragment program must be used to create a floating-point map. Computing a floating-point mipmap on current GPUs tends to be problematic for three reasons. Firstly, reading and writing floating-point textures can be slow. Secondly, if we are interested in the sum of only a subset of values, e.g. those that are greater than a given number, then introduce conditionals in the fragment program. Finally, the floating point representation may not have enough precision to give an exact sum.

Our accumulator algorithm avoids some of the problems of the mipmap method. We perform only texture reads which are more efficient than texture writes. We can calculate the precise sum to arbitrary precision and we avoid conditionals in the fragment program. One limitation of the algorithm is that it works only on integer datasets, although it can easily be extended to handle fixed-point datasets.

ACCUMULATOR sums the values stored in the texture tex utilizing the binary data representation. The sum of the values x_j in a set X can be written as:

$$\sum_{j=0}^{|X|} x_j = \sum_{j=0}^{|X|} \sum_{i=0}^k a_{ij} 2^i$$

where $a_{ij} \in \{0, 1\}$ are the binary digits of x_j and k is the maximum number of bits used to represent the values in X . There currently exists no convenient way to sum texels efficiently on current GPUs. We can, however, quickly determine the number of texels for which a particular bit i is set. If we reverse the order of the summations we get an expression that is more amenable to GPU

```

Accumulator( tex )
1 alpha test = pass with alpha ≥ 0.5
2 sum = 0
3 for i = 0 to b_max do
4   enable fragment program TestBit(i)
5   initialize occlusion query
6   RenderTexturedQuad( tex )
7   count = pixel count from occlusion query
8   sum += count * 2i
9 return sum

TestBit(i)
1 v = value from tex
2 fragment alpha = frac(v / 2(i+1))

```

ROUTINE 4.6: *Accumulator computes the sum of attribute values in texture tex. It performs b_max passes to compute the sum. Each pass computes the number of values with i-th bit set and stores it in count. This count is multiplied with 2ⁱ and added to sum. At the end of the b_max passes, the variable sum aggregates all the data values in the texture.*

computation:

$$\sum_{i=0}^k 2^i \left(\sum_{j=0}^{|X|} a_{ij} \right)$$

The inner summation is simply the number of x_j that have the i th bit set. This summation is the value of `count` calculated on lines 4-6 where we render a quad textured with `tex`.

The fragment program `TESTBIT` ensures that only fragments corresponding to texels with the i th bit set pass the alpha test. Determining whether a particular bit is set is trivial with bit-masking operations. Since current GPUs do not support bit-masking operations in fragment programs, we use an alternate approach. We observe that an integer x has its i th bit equal to 1 if and only if the fractional part of $x/2^{i+1}$ is at least 0.5. In `TestBit`, we divide each value by 2^{i+1} and put the fractional part of the result into the alpha channel. We use the alpha test to reject a fragments with alpha less than 0.5. It is possible to do the comparison and reject fragments directly in the fragment program, but it is faster in practice to use the alpha test. Pseudocode for our algorithm is shown in the routine 4.6.

`ACCUMULATOR` can sum only a subset of the records in `tex` that have been selected using the stencil buffer. Attributes that are not selected fail the stencil test and thus make no contribution to the final sum. We use the `ACCUMULATOR` algorithm to obtain `SUM`. `AVG` is obtained by computing `SUM` and `COUNT` and computed as `AVG = SUM/COUNT`.

5. IMPLEMENTATION & PERFORMANCE

We have implemented and tested our algorithms on a high end Dell Precision Workstation with dual 2.8GHz Intel Xeon Processors and an NVIDIA GeForceFX 5900 Ultra graphics processor. The graphics processor has 256MB of video memory with a memory data rate of 950MHz and can process upto 8 pixels at processor clock rate of 450 MHz. This GPU can perform IEEE single-precision floating point operations in fragment programs.

5.1 Benchmarks

For our benchmarks, we have used a database consisting of TCP/IP data for monitoring traffic patterns in local area network and wide area network and a census database [5] consisting of monthly income information. In the TCP/IP database, there are one million records in the database. In our experiments, each record has 4 at-

tributes, (*data_count, data_loss, flow_rate, retransmissions*).

Each attribute in the database is stored in as a floating-point number encoded in a 32 bit RGBA texture. The video memory available on the NVIDIA GeForce FX 5900 graphics processor can store more than 50 attributes in 1000×1000 textures, amounting to a total of 50 million values in the database. We transfer textures from the CPU to the graphics processor using an AGP 8X interface.

The census database consists of 360K records. We used four attributes for each record of this database. Our performance results on the census data are consistent with the results obtained on the TCP/IP database.

5.2 Optimized CPU Implementation

We implemented the algorithms described in section 4 and compared them with an optimized CPU implementation. We compiled the CPU implementation using Intel compiler 7.1 with full compiler optimizations¹. These optimizations include

- **Vectorization:** The compiler detects sequential data scans and generates code for SIMD execution.
- **Multi-threading:** We used the compiler switch `QParallel` to detect loops which may benefit from multi-threaded execution and generate appropriate threading calls. This option enables the CPU implementation to utilize hyper-threading technology available on Xeon processors.
- **Inter-Procedural Optimization (IPO):** The compiler performs function inlining when IPO is enabled. It reduces the function call branches, thus improving its efficiency.

For the timings, we ran each of our tests 100 times and computed the average running time for the test.

5.3 GPU Implementation

Our algorithms described in Section 4 are implemented using OpenGL API. For generating the fragment programs, we used NVIDIA's CG compiler [20]. As the code generated by the compiler is often sub-optimal, we examined the assembly code generated by the compiler and reduced the number of assembly instructions to perform the same operation.

For the counting operations, we chose to use `GL_NV_occlusion_query` for image space occlusion queries. These queries can be performed asynchronously and often, do not add any additional overhead.

5.4 Copy Operation

Various database operations, such as comparisons, selection, etc, require the data values of an attribute stored in the depth buffer. For these operations, we copy the corresponding texture into the depth buffer. A fragment program is used to perform the copy operation. Our copy fragment program implementation requires three instructions.

1. **Texture Fetch:** We fetch the texture value corresponding to a fragment.
2. **Normalization:** We normalize the texture value to the range of valid depth values $[0, 1]$.
3. **Copy To Depth:** The normalized value is copied into the fragment depth.

Figure 3 shows the time taken to copy values from textures of varying sizes into the depth buffer. The figure indicates an almost

¹http://www.intel.com/software/products/compilers/techtocics/compiler_optimization_71.pdf

linear increase in the time taken to perform the copy operation as a function of the number of records. In the future, it may be possible to copy data values from textures directly to a depth buffer and that would reduce these timings considerably. Also, the increase in clock rates of graphics processors and improved optimizations to perform depth buffer writes [23] would help in reducing these timings.

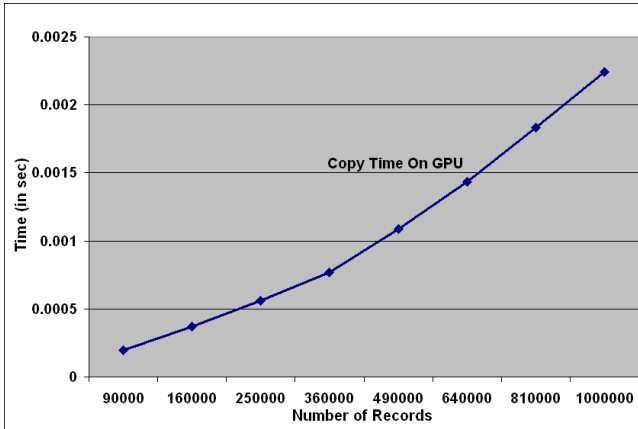


Figure 3: Plot indicating the time taken to copy data values in a texture to the depth buffer.

5.5 Predicate Evaluation

Figure 4 shows a plot of the time taken to compute a single predicate for an attribute such that the selectivity is 60%. In our experiments, we performed the operation on the first attribute of each record in the database. The plot compares a compiler generated SIMD optimized CPU code against a simple GPU implementation. The GPU timings include the computational time for evaluating the predicate, as well as the time taken to copy the attribute into the depth buffer. We observe that the GPU timings are nearly 3 times faster than the CPU timings. If we compare only the computational time on GPU, we observe that the GPU implementation is nearly 20 times faster than the SIMD optimized CPU implementation.

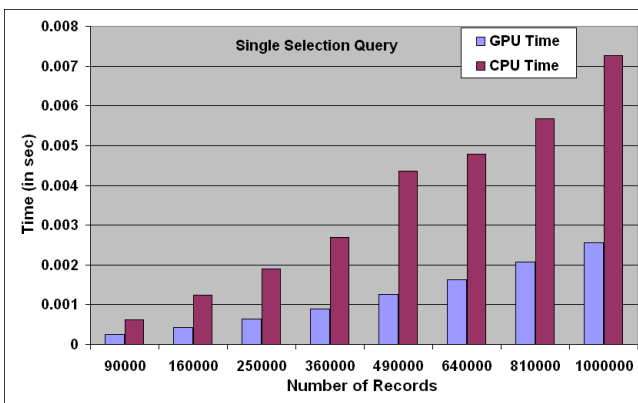


Figure 4: Execution time of a predicate evaluation with 60% selectivity by a CPU-based and a GPU-based algorithm. Timings for the GPU-based algorithm include time to copy data values into the depth buffer. Considering only computation time, the GPU is nearly 20 times faster than a compiler-optimized SIMD implementation.

5.6 Range Query

We tested the performance of RANGE by timing a range query with 60% selectivity. To ensure 60% selectivity, we set the valid range of values between the 20th percentile and 80th percentile of the data values. Again, in our tests, we used the `data_count` as our attribute. Figure 5 compares the time taken for a simple GPU implementation and a compiler-optimized SIMD implementation on CPU. In the GPU timings, we included the time taken for the copy operation. We observe that overall the GPU is nearly 4 times faster than the CPU implementation. If we consider only the computational time on GPU and CPU, we observe that the GPU is nearly 20 times faster than the compiler optimized CPU implementation.

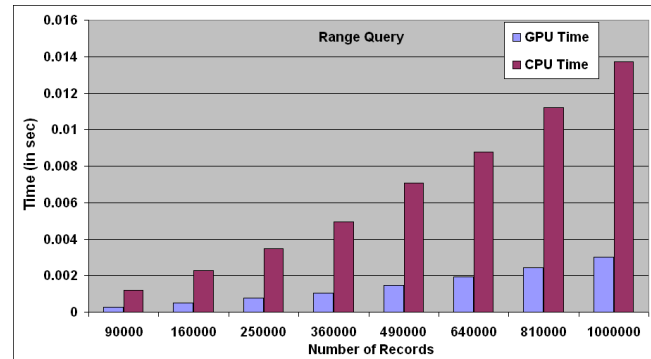


Figure 5: Execution time of a range query with 60% selectivity using a GPU-based and a CPU-based algorithm. Timings for the GPU-based algorithm include time to copy data values into the depth buffer. Considering only computation time, the GPU is nearly 20 times faster than a compiler-optimized SIMD implementation.

5.7 Multi-Attribute Query

We have tested the performance of our hardware based multi-attribute queries by varying the number of attributes and also the number of records in the database. In particular, we used queries with a selectivity of 60% for each attribute and applied the `AND` operator on the result for each attribute. In our tests, we used up to four attributes per query. For each attribute in the query, the GPU implementation copies the data values from the attribute's texture to the frame-buffer. Figure 6 shows the time taken by the GPU and CPU respectively, to perform multi-attribute queries with the varying number of attributes and records. The timings indicate that the GPU implementation is nearly 2 times faster than the CPU implementation. If we consider only the computational times on the GPU by ignoring copy times, we observe that the GPU is nearly 20 times faster than the optimized CPU implementation.

5.8 Semi-linear Query

We performed a semi-linear query on the four attributes by using a linear combination of four random floating point values and compared it against an arbitrary value. Figure 7 summarizes our timings for various number of tests on GPU and CPU. In our tests, we observe that the GPU timings are 9 times faster than an optimized CPU implementation.

5.9 K-th Largest Number

We performed three different tests to evaluate our `KTHLARGEST` algorithm on GPU. In each of these tests, we compared `KTHLARGEST` against a CPU implementation of the algorithm `QUICKSELECT` [12]. In our experiments, we used the `data_count` attribute. This

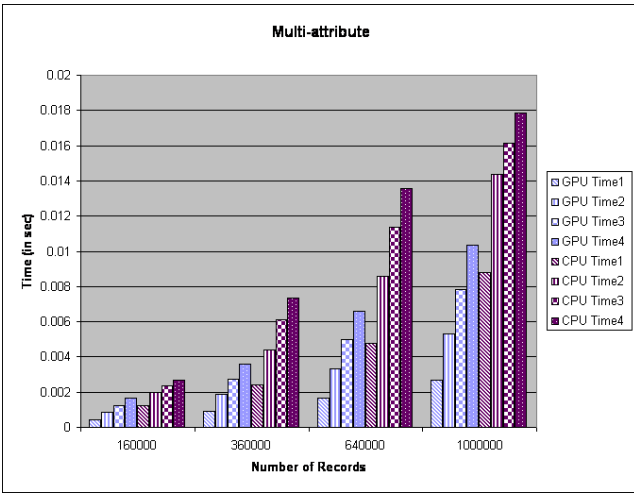


Figure 6: Execution time of a multi-attribute query with 60% selectivity for each attribute and a combination of AND operator. Time_{*i*} is the time to perform a query with *i* attributes. We show the timings for CPU and GPU-based implementations.

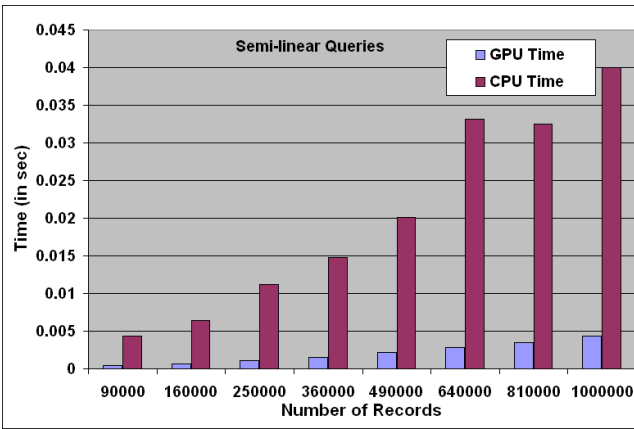


Figure 7: Execution time of a semi-linear query using four attributes of the TCP/IP database. The GPU-based implementation is almost one order of magnitude faster than the CPU-based implementation.

attribute requires 19 bits to represent the largest data value and has a high variance.

Test 1: Vary k and compute the time taken for KTHLARGEST and QUICKSELECT on the TCP/IP database. The tests were performed on 250K records in the database. Figure 8 shows the timings obtained using each of the implementations. This plot indicates that time taken by KTHLARGEST is constant irrespective of the value of k and is an interesting characteristic of our algorithm. On an average, the GPU timings for our algorithm are nearly twice as fast in comparison to the CPU implementation. It should be noted that the GPU timings include the time taken to copy values into the depth buffer. Comparing the computational times, we note that the average KTHLARGEST timings are 3 times faster than QUICKSELECT.

Test 2: In these tests, we compared the time taken by KTHLARGEST and QUICKSELECT to compute a median of a varying number of records. Figure 9 illustrates the results of our experiments. We observe that the KTHLARGEST on the GPU is nearly twice as fast as QUICKSELECT on the CPU. Considering only the computational times, we observe that KTHLARGEST is nearly 2.5 times faster

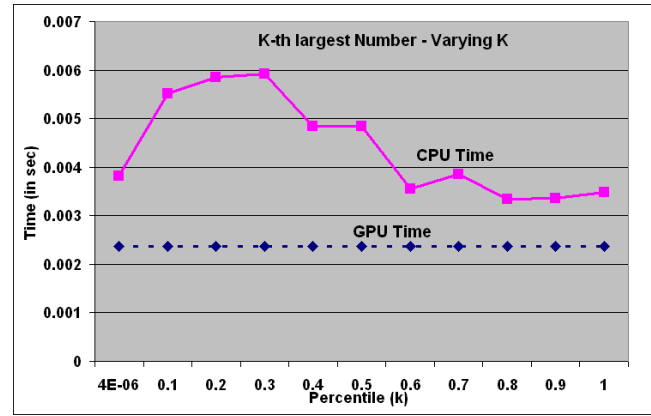


Figure 8: Time to compute k -th largest number on the data count attribute. We used a portion of the TCP/IP database with nearly 250K records.

than QUICKSELECT.

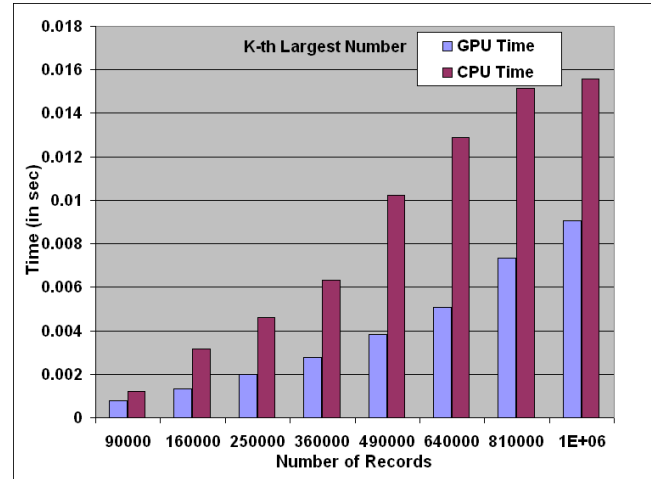


Figure 9: Time taken to compute the median using KTHLARGEST and QUICKSELECT on varying number of records.

Test 3: We also compared the time taken by KTHLARGEST and QUICKSELECT to compute a median with 80% selectivity of data values. Figure 10 indicates the time taken by KTHLARGEST and QUICKSELECT in computing the median as a function of the number of records. Our timings indicate that KTHLARGEST with 80% selectivity requires exactly the same amount of time as performing KTHLARGEST with 100% selectivity. We conclude from this observation that the use of a conditional to test for valid data has almost no effect on the running time of KTHLARGEST. For the CPU timings, we have copied the valid data into an array and passed it as a parameter to QUICKSELECT. The timings indicate that the total running time is nearly the same as that of running QUICKSELECT with 100% selectivity.

5.10 Accumulator

Figure 11 demonstrates the performance of ACCUMULATOR on the GPU and a compiler-optimized SIMD implementation of accumulator on the CPU. Our experiments indicate that our GPU algorithm is nearly 20 times slower than the CPU implementation, when including the copy time. Note that the CPUs have a much

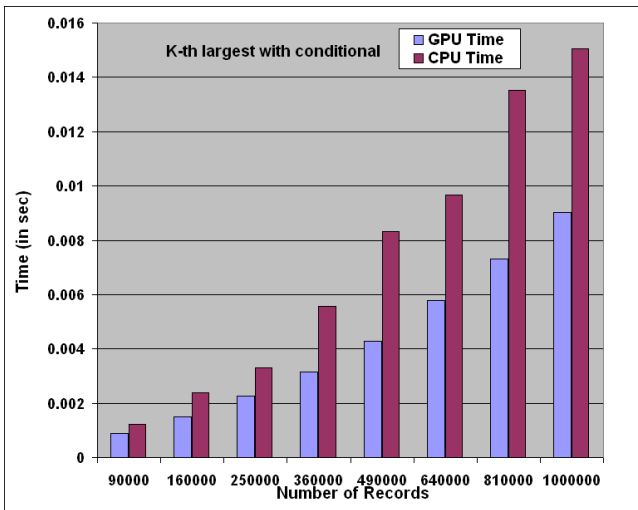


Figure 10: Time taken to compute the K-th largest number by the two implementations.

higher clock rate as compared to the GPU.

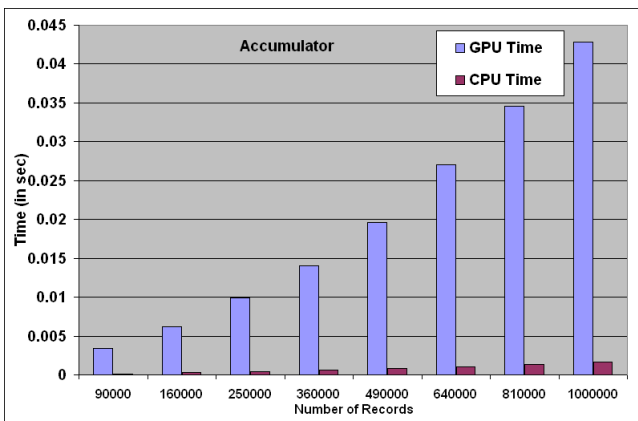


Figure 11: Time required to sum the values of an attribute on the CPU and by the GPU-based Accumulator algorithm

5.11 Selectivity Analysis

Recently, several algorithms have been designed to implement join operations efficiently using selectivity estimation [6, 8]. We compute the selectivity of a query using the COUNT algorithm (Section 4.3) We use image space occlusion queries to obtain the selectivity count. We performed the experiments described in Sections 5.5, 5.6, 5.7, 5.8. We observed that there is no additional overhead in obtaining the count of selected queries. Given selected data values scattered over a 1000×1000 frame-buffer, we can obtain the number of selected values within 0.25 ms.

6. ANALYSIS

In the previous section, we have highlighted the performance of our algorithms on different database operations and performed a preliminary comparison with some CPU-based algorithms. In this section, we analyze the performance of our algorithms and highlight the database operations for which the GPUs can offer considerable gain in performance.

6.1 Implementing Basic Operations on GPUs

There are many issues that govern the performance of the algorithms implemented on a GPU. However, some of the upcoming features in the next generation GPUs can improve the performance of these algorithms considerably.

Precision: Current GPUs have depth buffers with a maximum of 24 bits. This limited precision can be an issue. With the increasing use of GPUs in performing scientific computing, graphics hardware developers may add support for higher-precision depth buffers.

Copy Time: Several of our algorithms require data values to be copied from the texture memory to the depth buffer. Current GPUs do not offer a mechanism to perform this operation efficiently and this operation can take a significant fraction of the overall algorithm (e.g. algorithms for relational and range queries). In the future, we can expect support for this operation on GPUs which could improve the overall performance.

Integer Arithmetic Instructions: Current GPUs do not offer integer arithmetic instructions in the pixel processing engines. In addition to database operations, several image and video compression algorithms also require the use of integer arithmetic operations. Given that the fragment programs were just introduced in the last few years, the instruction sets for these programs are presently being enhanced. The instructions for integer arithmetic would reduce the timings of our ACCUMULATOR algorithm significantly.

Depth Compare Masking: Current GPUs support a boolean depth mask that enables or disables writes to a depth buffer. It is very useful to have a comparison mask specified for the depth function, similar to that specified in the stencil function. Such a mask would make it easier to test if a number has i-th bit set.

Memory Management: Current high-end GPUs have up to 512MB of video memory and this limit is increasing every year. However, due to the limited video memory, we may not be able to copy very large databases (with tens of millions of records) into GPU memory. In such situations, we would use out-of-core techniques and swap textures in and out of video memory. Another related issue is the bus bandwidth. Current PCs use AGP8x bus to transfer data from the CPU to the GPU and the PCI bus from the GPU to the CPU. With the announcement of PCI-EXPRESS bus, the bus bandwidth is going to improve significantly in the near future. Asynchronous data transfers would also improve the performance of these algorithms.

No Branching: Current GPUs implement branching by evaluating both portions of the conditional statement. We overcome this issue by using multi-pass algorithms and evaluating the branch operation using the alpha test or the depth test.

No Random Writes: The GPUs do not support random access writes, which makes it harder to develop algorithms on GPUs because they cannot use data rearrangement on GPUs.

6.2 Relative Performance Gain

We have presented algorithms for predicates, boolean combinations and aggregations. We have also performed preliminary comparison with optimized CPU-based implementations on a workstation with dual 2.8 GHz Xeon processors. Due to different clock rates and instruction sets, it is difficult to perform explicit comparisons between CPU-based and GPU-based algorithms. However, some of our algorithms perform better than others. We classify our algorithms into three categories: high performance gain, medium performance gain and low performance gain.

6.2.1 High Performance Gain

In these algorithms, we have observed an order of magnitude

speedup over CPU-based implementations. These include algorithms for semi-linear queries and selection queries. The main reason for the improved performance are:

- **Parallel Computation:** GPUs have several pixel processing engines that process multiple pixels in parallel. For example, on a GeForce FX 5900 Ultra we can process 8 pixels in parallel and reduce the computational time significantly. Also, each pixel processing engine has vector processing capabilities and can perform vector operations very efficiently. The speedup in selection queries is mainly due to the parallelism available in pixel processing engines. The semi-linear queries also exploit the vector processing capabilities.
- **Pipelining:** The GPUs are designed using a pipelined architecture. As a result, they can simultaneously process multiple primitives within the pipeline. The algorithms for handling multiple-attribute queries map well to the pipelined implementation.
- **Early Depth-Culling:** GPUs have specialized hardware that can reject fragments that will not pass the depth test early in the pipeline. Since the fragments do not have to pass through the pixel processing engines, this can lead to significant performance increases.
- **Eliminate branch mispredictions:** One of the major advantages in performing these selection queries on GPUs is that there are no branch mispredictions. Branch mispredictions can be extremely expensive on the modern CPUs. Modern CPUs use specialized schemes to predict the outcome of the branch instruction. Each branch misprediction can cost several clock cycles on current CPUs. For example, on a Pentium IV a branch misprediction can have a penalty of 17 clock cycles [22].

6.2.2 Medium Performance Gain

Several of our algorithms for database operations are only able to use a subset of the capabilities of the GPUs. In these cases, we have observed a speedup of nearly a factor of 2 to 4 times in comparison to an optimized-CPU implementation. For example, the `KTHLARGEST()` routine exhibits these characteristics. The speedup in the `KTHLARGEST()` is mainly due to the parallelism available in pixel processing engines. Given the GPU clock rate and the number of pixel processing engines, we can estimate the time taken to perform `KTHLARGEST()` under some assumptions. We assume that there is no latency in the graphics pipeline and in transmitting the pixel pass count from the GPU to the CPU. On a GeForce FX 5900 Ultra with clock rate 450MHz and processing 8 pixels per clock cycle, we can render a single 1000×1000 quad in 0.278 ms. In our experiments, we render 19 such quads to compute the *k*-th largest number. Rendering these quads should take 5.28 ms. On the other hand, the observed time for this computation is 6.6 ms. This indicates that we are utilizing nearly 80% of the parallelism in the pipeline. The observed timings are slightly higher due to the latencies in transmitting the data from the GPU to the CPU and vice-versa. A key advantage of our algorithm `KTHLARGEST()` in comparison with other parallel order statistic algorithms is that it does not require any data rearrangement. Data rearrangements can be expensive when combined with branching.

6.2.3 Low Performance Gain

In some cases, we did not observe any gain over a CPU-based implementation. Our GPU based `ACCUMULATOR` algorithm is

much slower than the CPU-based implementation. There are several reasons for this performance:

- **Lack of Integer Arithmetic:** Current GPUs do not support integer arithmetic instructions. Therefore, we used a fragment program with at least 5 instructions to test if the *i*-th bit of a texel is 1. There are several ways to implement such a feature in the hardware. A simplest mechanism is to copy the *i*-th bit of the texel into the alpha value of a fragment. This can lead to significant improvement in performance.
- **Clock Rate:** Not only are we comparing two architectures with different instruction sets, but they also have different clock rates. Our CPU implementation used top-of-the-line dual Xeon processors operating at 2.8GHz. Each Xeon processor has four SIMD processors that can perform four operations in parallel. On the other hand, the current GPU clock rate (450MHz) is much lower than the CPU clock rate. It is possible that the increasing parallelism in the GPUs and development of new instruction sets for fragment programs can bridge this gap in performance.

Our preliminary analysis indicates that it is advantageous to perform selection and semi-linear queries on GPUs. In addition, GPUs can also be used effectively to perform order statistics algorithms.

7. CONCLUSIONS AND FUTURE WORK

In this paper, we have presented novel algorithms for performing database operations on the GPUs. These include algorithms for predicates, boolean combinations, and aggregation queries. We have implemented these algorithms on a state-of-the-art GPU and highlighted its performance for following queries: relational query, range query, multi-attribute query, semi-linear query, *k*-th-largest number computation, accumulator and selectivity analysis. We have also performed preliminary comparisons with optimized implementations of CPU-based algorithms. In some cases, we observed noticeable performance gain.

We have shown that GPUs are excellent candidates for performing some of the databases operations. Some reasons for this include:

- **Commodity Hardware:** High-end GPUs are freely available on PCs, consoles and workstations. The cost of a high-end GPU is typically less than the cost of a high-end CPU (by a factor of two or more).
- **Higher Growth Rate:** Over the last decade the growth rate of GPU performance has been higher than that of CPUs. This trend is expected to continue for the next five years or so. The performance gap between the CPUs and GPUs will probably increase considerably and we can obtain improved performance for database queries on the GPUs.
- **Multiple Fragment Processors and Improved Programmability:** Current high-end GPUs already have 8 fragment processors. This number is expected to increase in the future. As the instruction sets of fragment programs improve, the running time of many of our algorithms will further decrease.
- **Useful Co-Processor:** Overall, the GPU can be used as an effective co-processor along with the CPUs. It is clear that GPU is an excellent candidate for some database operations, but not all. Therefore, it would be useful for database designers to utilize GPU capabilities alongside traditional CPU-based code.

There are many avenues for future work. It is possible to use new capabilities and optimizations to improve the performance of many of our algorithms. Furthermore, we would like to develop algorithms for other database operations and queries including sorting, join, and indexed search, and OLAP and data mining tasks such as data cube roll up and drill-down, classification, and clustering, which may benefit from multiple fragment processors and vector processing capabilities. We also plan to design GPU-based algorithms for queries on more complicated data types in the context of spatial and temporal databases and perform continuous queries over streams using GPUs.

8. REFERENCES

- [1] A. Ailamaki, D. J. DeWitt, M. D. Hill, and M. Skounakis. Weaving relations for cache performance. In Peter M. G. Apers, Paolo Atzeni, Stefano Ceri, Stefano Paraboschi, Kotagiri Ramamohanarao, and Richard T. Snodgrass, editors, *Proceedings of the Twenty-seventh International Conference on Very Large Data Bases: Roma, Italy, 11–14th September, 2001*, pages 169–180, Los Altos, CA 94022, USA, 2001. Morgan Kaufmann Publishers.
- [2] A. Ailamaki, D. J. DeWitt, M. D. Hill, and D. A. Wood. DBMSs on a modern processor: Where does time go? In *Proc. of VLDB*, pages 266–277, 1999.
- [3] J. Bolz, I. Farmer, E. Grinspun, and P. Schröder. Sparse matrix solvers on the gpu: Conjugate gradients and multigrid. *ACM Trans. on Graphics (Proc. of ACM SIGGRAPH)*, 22(3), 2003.
- [4] Peter A. Boncz, Stefan Manegold, and Martin L. Kersten. Database architecture optimized for the new bottleneck: Memory access. In *Proc. of VLDB*, pages 54–65, 1999.
- [5] Census bureau databases. <http://www.bls.census.gov/cps/>.
- [6] Zhiyuan Chen, Nick Koudas, Flip Korn, and S. Muthukrishnan. Selectively estimation for boolean queries. In *Proceedings of the nineteenth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 216–225. ACM Press, 2000.
- [7] Michael Doggett. Programmability features of graphics hardware. *ACM SIGGRAPH Course Notes # 11*, 2003.
- [8] Lise Getoor, Benjamin Taskar, and Daphne Koller. Selectivity estimation using probabilistic models. *SIGMOD Record (ACM Special Interest Group on Management of Data)*, 30(2):461–472, June 2001.
- [9] N. Goodnight, C. Woolley, G. Lewin, D. Luebke, and G. Humphreys. A multigrid solver for boundary value problems using programmable graphics hardware. *ACM SIGGRAPH/Eurographics Conference on Graphics Hardware*, 2003.
- [10] N. Govindaraju, B. Lloyd, S. Yoon, A. Sud, and D. Manocha. Interactive shadow generation in complex environments. *Proc. of ACM SIGGRAPH/ACM Trans. on Graphics*, 22(3):501–510, 2003.
- [11] M. Harris, G. Coombe, G. Scheuermann, and A. Lastra. Physically-based visual simulation on graphics hardware. *SIGGRAPH/Eurographics Workshop on Graphics Hardware*, 2002.
- [12] C. A. R. Hoare. Algorithm 65: find. *Commun. ACM*, 4(7):321–322, 1961.
- [13] K. Hoff, T. Culver, J. Keyser, M. Lin, and D. Manocha. Fast computation of generalized voronoi diagrams using graphics hardware. *Proceedings of ACM SIGGRAPH*, pages 277–286, 1999.
- [14] J. Kruger and R. Westermann. Linear algebra operators for gpu implementation of numerical algorithms. *ACM Trans. on Graphics (Proc. of ACM SIGGRAPH)*, 22(3), 2003.
- [15] E. S. Larsen and D. K. McAllister. Fast matrix multiplies using graphics hardware. *Proc. of IEEE Supercomputing*, 2001.
- [16] M. Macedonia. The gpu enters computing’s mainstream. *Computer*, October 2003.
- [17] S. Manegold, P. Boncz, and M L. Kersten. What happens during a join? Dissecting CPU and memory optimization effects. In Amr El Abbadi, Michael L. Brodie, Sharma Chakravarthy, Umeshwar Dayal, Nabil Kamel, Gunter Schlageter, and Kyu-Young Whang, editors, *VLDB 2000, Proceedings of 26th International Conference on Very Large Data Bases, September 10–14, 2000, Cairo, Egypt*, pages 339–350, Los Altos, CA 94022, USA, 2000. Morgan Kaufmann Publishers.
- [18] Stefan Manegold, Peter A. Boncz, and Martin L. Kersten. Generic database cost models for hierarchical memory systems. In Philip A. Bernstein, Yannis E. Ioannidis, Raghu Ramakrishnan, and Dimitris Papadias, editors, *Proceedings of the Twenty-eighth International Conference on Very Large Data Bases: Hong Kong, China, 20–23rd August, 2002*, pages 191–202, Los Altos, CA 94022, USA, 2002. Morgan Kaufmann Publishers.
- [19] D. Manocha. *Interactive Geometric and Scientific Computations using Graphics Hardware*. SIGGRAPH Course Notes # 11, 2003.
- [20] William R. Mark, R. Steven Glanville, Kurt Akeley, and Mark J. Kilgard. Cg: A system for programming graphics hardware in a c-like language. In *Proc. of ACM SIGGRAPH*, 2003. http://developer.nvidia.com/page/cg_main.html.
- [21] Shintaro Meki and Yahiko Kambayashi. Acceleration of relational database operations on vector processors. *Systems and Computers in Japan*, 31(8):79–88, August 2000.
- [22] <http://lava.cs.virginia.edu/bpred.html>.
- [23] Nvidia geforce fx gpus: Intellisample technology. http://www.nvidia.com/object/intellisample_tb.html.
- [24] Nv_occlusion_query specification. http://www.nvidia.com/dev_content/nvopenglspecs/GL_NV_occlusion_query.txt.
- [25] Niki Pissinou. Towards and infrastructure for temporal databases — A workshop report. *SIGMOD Record (ACM Special Interest Group on Management of Data)*, 23(1):35–51, March 1994.
- [26] T. Purcell, I. Buck, W. Mark, and P. Hanrahan. Ray tracing on programmable graphics hardware. *ACM Trans. on Graphics (Proc. of SIGGRAPH’02)*, 21(3):703–712, 2002.
- [27] T. Purcell, C. Donner, M. Cammarano, H. Jensen, and P. Hanrahan. Photon mapping on programmable graphics hardware. *ACM SIGGRAPH/Eurographics Conference on Graphics Hardware*, pages 41–50, 2003.
- [28] Jun Rao and Kenneth A. Ross. Cache conscious indexing for decision-support in main memory. In *Proc. of VLDB*, pages 78–89, 1999.
- [29] Kenneth A. Ross. Conjunctive selection conditions in main memory. In ACM, editor, *Proceedings of the Twenty-First ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems: PODS 2002: Madison, Wisconsin, June 3–5, 2002*, pages 109–120, New York, NY 10036, USA, 2002. ACM Press. ACM order number 475021.
- [30] J. Rossignac, A. Megahed, and B.D. Schneider. Interactive inspection of solids: cross-sections and interferences. In *Proceedings of ACM Siggraph*, pages 353–60, 1992.
- [31] Ambuj Shatdal, Chander Kant, and Jeffrey F. Naughton. Cache conscious algorithms for relational query processing. In Jorge B. Bocca, Matthias Jarke, and Carlo Zaniolo, editors, *20th International Conference on Very Large Data Bases, September 12–15, 1994, Santiago, Chile proceedings*, pages 510–521, Los Altos, CA 94022, USA, 1994. Morgan Kaufmann Publishers.
- [32] Chengyu Sun, Divyakant Agrawal, and Amr El Abbadi. Hardware acceleration for spatial selections and joins. In ACM, editor, *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data 2003, San Diego, California, June 09–12, 2003*, pages 455–466, New York, NY 10036, USA, 2003. ACM Press.
- [33] C. J. Thompson, S. Hahn, and M. Oskin. Using modern graphics architectures for general-purpose computing: A framework and analysis. *Proc. of IEEE/ACM International Symposium on Microarchitectures*, pages 306–317, 2002.
- [34] Jingren Zhou and Kenneth A. Ross. Implementing database operations using SIMD instructions. In Michael Franklin, Bongki Moon, and Anastassia Ailamaki, editors, *Proceedings of the ACM SIGMOD International Conference on Management of Data, June 3–6, 2002, Madison, WI, USA*, pages 145–156, New York, NY 10036, USA, 2002. ACM Press. ACM order number 475020.