

Distributed/Cluster Computing for Data Stream Mining: Draft Notes

A thesis
submitted in partial fulfilment
of the requirements for the Degree
of
Master of Science
at the
University of Waikato
by
Vladimir Petko

University of Waikato
2016

Abstract

The thesis is focused on elucidating GPU computing feasibility for clustering tasks

Acknowledgements

Contents

Abstract	i
Acknowledgements	iii
1 General Purpose GPU Computing	2
2 System Architecture	12
3 k-Nearest Neighbours	15
4 Stochastic Gradient Descent	31
4.1 Introduction	31
4.2 Approaches for SGD parallelisation	31
4.3 Hogwild!	33
4.3.1 Best Ball Optimization	33
4.3.2 Backoff scheme	33
4.4 1bit SGD	33
4.5 Stochastic Gradient Descent using OpenCL/HSA architecture	34
4.5.1 OpenCL Memory Transfer	35
5 Experimental Results	37
6 Conclusions and Future Work	41

List of Figures

1.1	CPU versus GPU hardware architecture. Reproduced from NVIDIA GPU ARCHITECTURE and CUDA PROGRAMMING ENVIRONMENT by Alan Tatourian[71]	3
1.2	GPU Memory Tiers. Reproduced from NVIDIA GPU ARCHITECTURE and CUDA PROGRAMMING ENVIRONMENT by Alan Tatourian[71]	4
1.3	GP GPU technologies tree. Reproduced from C. Nugteren, Improving the Programmability of GPU Architecture, p. 21 [55]	5
2.1	SGD classifier training activity.	13
3.1	Left; Four matrix-vector multiplication kernels designed to perform well at different shapes $m \times n$ of A . Middle; Tuning mesh. Right; Best kernel in practice. The dashed line indicates the minimum 21504 rows needed in A for full occupancy of the Nvidia Tesla C2050 card in a one-thread-per-row kernel. Note the logarithmic axes. Reproduced from High-Performance Matrix-Vector Multiplication on the GPU by Hans Henrik Brandenburg Sørensen[68]	17
3.2	Selection Algorithm Performance for $K=128$. Test configuration GPU R9 390, CPU AMD A8-7600, AMD Catalyst version 15.20.	18
3.3	k-d tree construction and NN-search pseudocode	20

3.4	Distributions with low intrinsic dimension. The purple areas in these figures indicate regions in which the density of the data is significant, while the complementary white areas indicate areas where data density is very low. The left figure depicts data concentrated near a one-dimensional manifold. The ellipses represent mean+PCA approximations to subsets of the data. Our goal is to partition data into small diameter regions so that the data in each region is well-approximated by its mean+PCA. The right figure depicts a situation where the dimension of the data is variable. Some of the data lies close to a one-dimensional manifold, some of the data spans two dimensions, and some of the data (represented by the red dot) is concentrated around a single point (a zero-dimensional manifold). Reproduced from Learning the structure of manifolds using random projections by Freund Yoav et al.[40]	21
3.5	Left: Partitioning produced by k-d tree. Right: Partitioning produced by Random Projection Tree. Reproduced from Learning the structure of manifolds using random projections by Freund Yoav et al.[40]	22
3.6	Random Projection Tree Pseudocode - Random Tree Max [34]	22
3.7	Random Projection Tree Pseudocode - Random Tree Median[34]	23
3.8	Fast Johnson-Lindenstrauss transform[22] vs. sparse matrix implementation[20] using ViennaCL. Test configuration GPU R9 390, CPU AMD A8-7600, AMD Catalyst version 15.20.	26
3.9	Z-Order Curve. Red lines highlight some region jumps. Green shows locality-preserving region.	28
3.10	k-d tree test speed	29
4.1	OpenCL SGD training process.	36

5.1	Naive KNN implementation (Radeon Mobility HD5730)	38
5.2	Stochastic Gradient Descent Training Speedup (AMD R9 390, Spectre vs MOA implementation on AMD A8-7600)	39
5.3	Stochastic Gradient Descent Training Latency (AMD R9 390 vs MOA implementation on AMD A8-7600)	40

List of Tables

1.1	Input Size and Execution Time	7
-----	---	---

Introduction

In real world applications such as industrial monitoring, sensor networks, financial data generate large unbounded streams of data which has to be processed with pre-defined response time. The processors capabilities limit the bandwidth of the stream which can be processed. Parallelizing processing algorithm will increase maximum bandwidth while maintaining the response time requirement.

Chapter 1

General Purpose GPU Computing

Introduction

The modern Graphical Processing Units (GPU) greatly outpace CPUs in arithmetic throughput and memory bandwidth for data-parallel tasks. Since 2001 the efforts were made to port data parallel algorithms to GPUs - first using shader languages such as HLSL, then with the release of Nvidia G80 in 2006 using extensions to C programming language - CUDA[12]. Presently there is a number of programming frameworks targetting specifically GPU architecture such as CUDA[53], OpenCL[11], RenderScript[16], DirectCompute[5] and more generic parallel-processing frameworks such as OpenMP[14] and AMP[3] which provide GPU backend as one of the targets. The differences in the hardware architecture between CPU and GPU is reflected in the programming model of the traditional GPU-specific languages which contain hardware architecture specific language constructs. This chapter provides an overview of GPU architecture, most known programming frameworks, lists limitations of the traditional GP GPU programming, discusses OpenCL 2.0 standard which addresses some of the limitations and Heterogenous System Architecture (HSA) an optimized platform architecture for OpenCL 2.0.

GPU Architecture

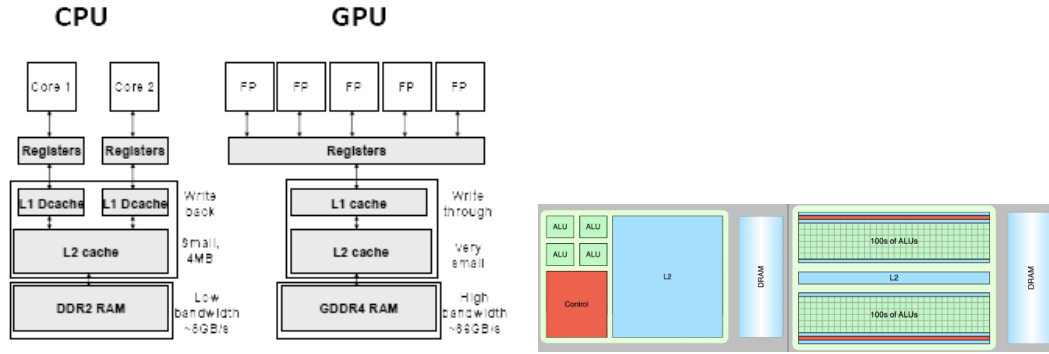


Figure 1.1: CPU versus GPU hardware architecture. Reproduced from NVIDIA GPU ARCHITECTURE and CUDA PROGRAMMING ENVIRONMENT by Alan Tatourian[71]

The main differences between modern CPU and GPU architectures are the level of parallelism and ability to directly address tiered memory. Modern CPU with 2 hex-cores support maximum of 12 threads (24 with hyperthreading) and minimal unit of execution for the NVIDIA GPU (called *wavefront*) is 32 threads. Modern GPUs implement SIMT (Single Instruction - Multiple Thread) execution model (AMD/NVIDIA desktop GPUs) first introduced by NVIDIA in G80 model[12]. The single unit of scalar instructions called *kernel* is scheduled to execute in blocks of data-parallel threads on SIMT hardware. Each instruction in a block is executed in a lock-step. The control divergence is emulated by *masking* - the device executes instructions from both branches of the conditional statement[17][53]. The CPU thread is a heavy-weight entity which is centered around execution of a specific task for an extended period of time. Whenever CPU needs to preempt the thread, the register state is stored and another thread takes over. This makes a context switch a costly operation and operating systems attempt to minimize number of context switches per second. The GPU context switch is an extremely lightweight operation and is routinely used for the latency-hiding - whenever the wavefront is waiting on data, the GPU schedules another wavefront for execution. The GPU registers are private for each thread and are not reallocated until thread execution completes.

Modern CPUs provide a flat view of the operating system memory while GPUs divide memory in tiers based on the access speed:

- *private/register* - private to the current thread
- *local* - shared within a *threadblock*
- *global* - accessible by every thread

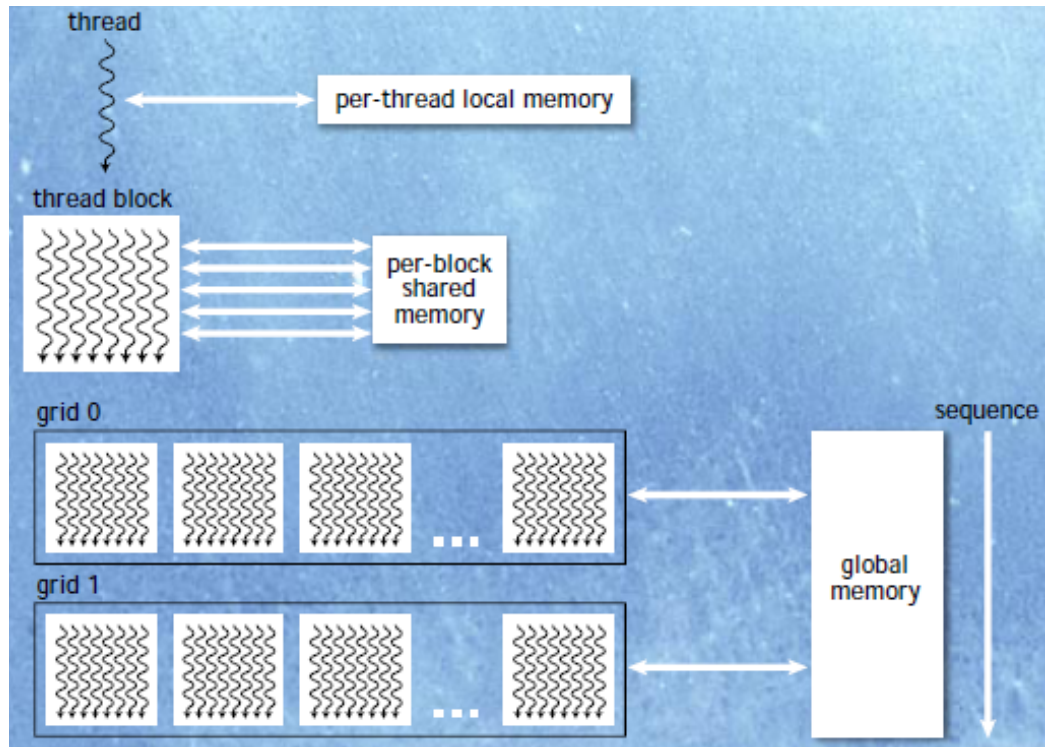


Figure 1.2: GPU Memory Tiers. Reproduced from NVIDIA GPU ARCHITECTURE and CUDA PROGRAMMING ENVIRONMENT by Alan Tatourian[71]

The GPU programming following abstractions:

- *Kernel* - a unit of execution
- *Thread* - a single unit of processed data
- *Threadblock* - a group of *threads* sharing same *kernel* and *local* memory.

The unit of scheduling is called *wavefront* in AMD terminology or *warp* in NVIDIA and typically consists of 32 threads on NVIDIA and 64 on AMD hardware. The GPU chip is equipped with a number of SIMT cores which

execute same instruction for each *warp*. Divergence of control results in under-load of the processing units and reduces performance. The branching should be reduced to wavefront granularity to avoid wasting execution cycles[67][53] It should be noted that the wavefront size is a hardware specific feature and the optimization should be performed at the run-time.

General Purpose GPU Computing Frameworks

Existing General Purpose GPU (GP GPU) computing frameworks can be classified by the level of provided hardware abstraction: high-level frameworks integrate with existing high-level programming language such as Java to provide parallel computing capabilities without exposing any hardware details[13], traditional GPU languages such as CUDA[53] expose task scheduling and memory management giving the expert user fine-tuning capabilities, low level languages provide intermediate binary format compatible with multiple hardware targets. The tree of the GP-GPU technologies is presented in the Figure 1.

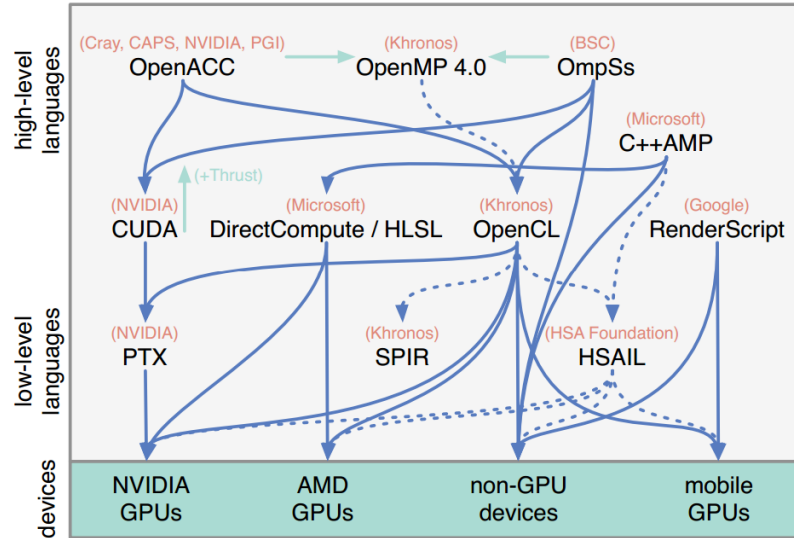


Figure 1.3: GP GPU technologies tree. Reproduced from C. Nugteren, Improving the Programmability of GPU Architecture, p. 21 [55]

High Level Languages

OpenACC and OpenMP are high level parallel programming frameworks that specify a set of annotations, environment variables and library routines for shared memory parallelism in C/C++ and Fortran programs[?][14]. Microsoft C++ AMP[3] is a C++ library which enables parallel computations for CPU and GPUs (using Microsoft DirectX Shading Language) Rootbear GPU compiler provides a transparent compilation of Java code into CUDA[59]. Aparapi provides a way to generate OpenCL kernel code from Java, theoretically allowing code which can be executed on CPU and offloaded to GPU if needed[1]. Project Sumatra is a OpenJDK project which focuses on development of the Hotspot virtual machine capable of offloading JDK 8 Stream API[10] computations to the GPU[13].

GPU-specific Languages

GPU-specific languages provide a programming model consistent with the GPU hardware implementation.

- CUDA - A programming language for NVIDIA hardware based on C language. Kernels are expressed as C-functions for one thread with parallelism defined at run-time by specifying dimensions of execution grid and thread blocks[53]
- OpenCL 1.X builds upon ideas implemented in CUDA by adding device management APIs and providing hardware-agnostic programming specification. OpenCL gives *write once-run anywhere* guarantee but does not give any performance consistency guarantees across different hardware[69].
- RenderScript - Android GPU computing component which uses OpenCL with Java binding programming model - C-style kernels and Java-based control code. RenderScript does not provide any APIs for the *work-*

group size control in the bid to provide performance portability between different devices[16].

- DirectCompute/HLSL - Microsoft Parallel Computing.

Low-Level Languages

The low level assembly representation is used to abstract compiler implementation from the actual hardware since each model or even revision may have a different instruction set. The translation is performed by *Just-In-Time* compiler before the kernel execution. Each vendor provides different low level specifications: NVIDIA CUDA uses Parallel Thread Execution and Instruction Set Architecture (PTX ISA)[15], Khronos Group specifies Standard Portable Intermediate Representation(SPIR)[18], and HSA Foundation specifies Heterogenous System Architecture Intermediate Language (HSAIL)[19].

Limitations

Input Size The massively parallel nature of GPU platforms require a certain amount of data to be passed to the kernel to achieve maximum performance. Table 1.1 shows execution time of a kernel which assigns index to each array element $X_i = i$ on AMD A8-7600. The execution time starts to increase when input size is above 1024 and remains constant for lower values. To maximum performance on AMD A8-7600 will be achieved when input size will exceed 1024 elements.

Global Size	256	512	768	1024	1280	2560	3072	3584	4608	4864
Execution Time (μ sec)	8	8	8	8	9	9	10	11	11	12

Table 1.1: Input Size and Execution Time

GPU Memory Size and Host-GPU Transfer The discrete GPU requires transfer of data from the host to the GPU memory which adds additional overhead to the computations and requires task partitioning according to the mem-

ory specification of GPU[66]. Memory transfer is a bottleneck for Aparapi and its developers allow explicit memory management[1]. This effectively reduces framework which promises CPU-GPU interoperability to the Java wrapper of the OpenCL API.

Kernel Launch There is a constant time needed to setup kernel launch which might offset any gain from parallelization if the data can be processed sequentially faster.(NB. Amdahl's law) It is impossible to schedule kernel execution from within the kernel itself requiring a mix of kernel and host code if several iterations are required.

OpenCL 2.0

OpenCL 2.0 standard[11] introduces several features which attempt to address limitations of GPU programming:

- Shared Virtual Memory - both host and kernel code share same address space thus either hiding memory transfers (discreet GPU driver stack) or if backed by the hardware architecture such as HSA eliminate its need[19]
- Dynamic Parallelism - OpenCL 2.0 allows scheduling of kernels from within a kernel without host interaction reducing host CPU bottleneck.
- Pipes - pipes feature allows passing data from kernel to kernel without processing the whole input which allows to obtain the results of computation faster.

HSA Platform

AMD introduced Heterogeneous System Architecture platform as an optimized platform architecture for OpenCL 2.0. Its specification introduces a set of requirements that allow both GPUs and CPU share same memory space, synchronize execution using signals and atomics and schedule execution both from

GPU and CPU[19]. Task execution is performed by *agents* which represent CPU or GPU nodes. The task execution is scheduled via *queues* and synchronized using *signals*. HSA memory model guarantees sequential consistency for the correctly synchronized programs.

Software Available: At the moment (Feb 2014) there is a OpenCL 2.0- $\dot{\iota}$ HSAIL compiler available[8] and a Linux-based runtime environment[7].

HSA Queues

HSA uses queues to schedule code execution. A HSA *queue* is a ringbuffer which contains *packets* with either call or synchronization parameters. The queue maintains two indexes - read index and write index. Write index is modified by the user and used to submit packets to the queue. The read index is updated by the packet processor whenever the packet is taken for execution. As soon as packet is written to the queue the ownership is taken by the HSA packet processor and it may change packet contents at any time[19]. Compared to traditional dispatch where the execution is scheduled via user-mode and kernel-mode driver layers the HSA dispatch intends to be lightweight and source-agnostic way of scheduling execution. The HSA Queues support work-stealing that is several HSA agents may be attached to the queue to share the workload.

HSA Signals

HSA uses *signals* to perform synchronization between host and kernels being executed or to signal completion of the task. A *signal* is essentially a shared memory variable modified by the HSA agent. Runtime environment provides a way to check the value of the signal or wait for the specific value.

HSA Memory Model

The sequential consistency was first defined by L. Lamport as “..the result of any execution is the same as if the operations of all the processors were

executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by its program.” Modern processors (ARM, x86, Itanium, POWER) introduce relaxed memory model to allow a range of the hardware optimizations to provide better performance by reordering load and store operations[51]. Platform specification states

The HSA memory consistency model is a relaxed model based around RCsc semantics on a set of synchronizing operations. The standard RCsc model is extended to include fences and relaxed atomic operations. In addition HSA includes concepts of memory segments and scopes.

[19] Similiar to Java Memory Model[?] it guarantees sequential consistency for the correctly synchronized programs, that is 'synchronizing operations meet the requirements for sequential consistency within each scope/segment instance'[19]. The specification introduces several memory segments: “

- Global segment, shared between all agents.
- Group segment, shared between work-items in the same work-group in a HSAIL kernel dispatch.
- Private, private to a single work-item in a HSAIL kernel dispatch.
- Kernarg, read-only memory visible to all work-items in a HSAIL kernel dispatch.
- Readonly, read-only memory visible to all agents

” Each particular memory location is always associated with one and only one segment and all operations apply to only one segment with the exception of fence operations[19]. In addition to memory segments HSA memory model introduces *scopes* : wavefront, work-group, component and system. They can be used to reduce visibility of the memory operation compared to the default supported by the segment. The global segment may use any of the specified scopes and group segment is limited to wavefront and workgroup scopes[19]. Different workgroups accessing a global variable with the workgroup scope

will work with different instances of the variable. The write serialization only applies to the operations within the segment/scope that they specify.

Implementation Notes

Sequential execution of several packets sometimes may be faster than submission of all packets and waiting for the barrier packet. For instance first scenario runs in 8 μsec per packet on AMD A8-7600 and second results in 193 μsec per packet. According to AMD support this is caused by the CPU going into power-saving mode while kernel is running. There is a constant time needed to setup kernel launch, e.g. for AMD A8-7600, it is 6 μsec using HSA.

Conclusion

The modern specifications such as OpenCL 2.0 and HSA attempt to address some of the latency issues of the GPU programming by introduction of the shared memory and lightweight dispatch and data passing mechanisms. This work will focus on the evaluation of suitability of those technologies for the latency-sensitive data stream processing.

Chapter 2

System Architecture

The implemented data stream processing library provides a set of classification algorithms for Massively Online Analysis(MOA)[27]. The library uses existing linear algebra package ViennaCL[72] which allows multiple backends such as CUDA, OpenCL, CPU and extends it with the machine learning algorithms such as nearest neighbours search and stochastic gradient descent.

MOA interface The ViennaCL library is implemented in C++ and as such requires Java Native Interface[9] to be used to interface from the Java Virtual Machine. Java Virtual Machine manages its own memory space and garbage collector may move the data at any time. JNI provides two mechanisms to access the array data from the native code. First is copying - the java pointer is locked by the critical section and array content is copied to the native array. Second skips the copying and provides direct access to the java pointer. Both involve entering and exiting a critical section and impose significant performance loss due to the copying and locking overhead. Those costs can not be avoided but can be minimized by moving them to the instance creation/-modification time - the object constructor will call the native method which allocates the native data structures and moves data from the Java storage to the native one.

The alternative solution uses *java.misc.Unsafe* class to manipulate offheap memory directly. The native code allocates GPU shared virtual memory and

passes the pointer to the Java implementation. Java code uses *java.misc.Unsafe* methods to update data in parallel to training. Figure 2.1 shows the training process of the Stochastic Gradient Descent classifier. The instances are accumulated in batches on the main thread of execution, the batches are passed to the data transfer thread that handles CPU-GPU transfer and then the training thread is triggered to run GPU kernels. Whenever the evaluation (*getVotesForInstance*) is called the threads process the last available batch and meet at the synchronization point stopping execution.

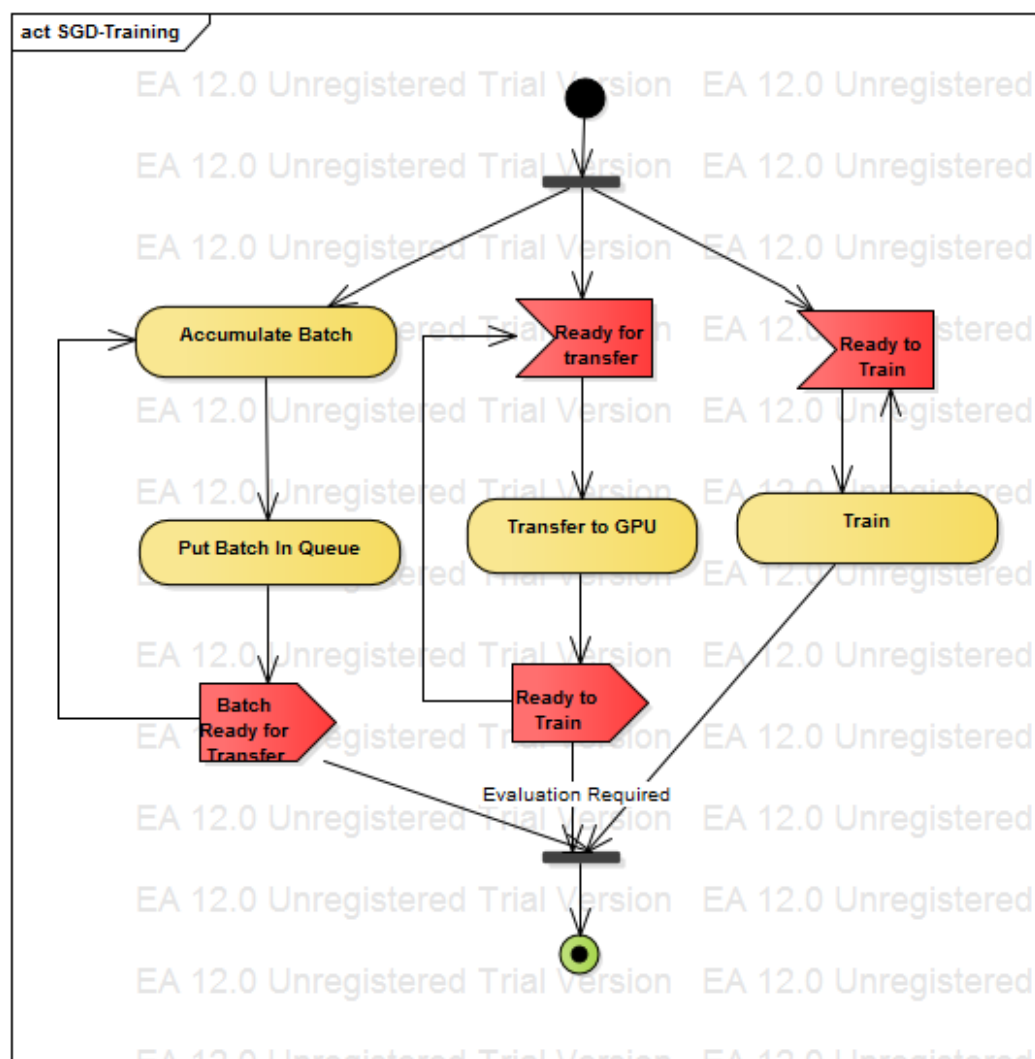


Figure 2.1: SGD classifier training activity.

GPU Memory Limits The library implementation stores the instance data as the native ViennaCL types. This implies that for GPU ViennaCL backends the data will be stored in GPU memory that may be insufficient for the larger

problems. In such case the partitioning will be used - the problem data is kept in the Java memory as collection of *weka.core.Instance* objects and offloaded to GPU-backed context on as needed basis. The *weka.core.Instance* class represents the attribute values as a vector of double precision numbers. Modern consumer GPUs provide far better floating point performance than double performance. For instance modern AMD GPUs have 8x scale that is floating point performance is 8x better than double one (R390, R290). NVIDIA GPUs have 32x scale *TODO citation*. There are several works that explore using fixed precision numbers to reduce memory requirements of the machine learning tasks[61][44]. The alternate precision implementation will impose the overhead costs of double to fixed/floating point conversion if the GPU classifier will be used for instance as a part of the meta-classifier ensemble.

HSA Backend This work adds a new HSA backend to the ViennaCL library based on the HSA Runtime[19]. This implementation is tuned for Kaveri AMD APU and uses the same set of OpenCL kernels as the OpenCL backend. In HSA backend the main system memory is transparently mapped to GPU and vice-versa, allowing to use vector or matrix element-addressing operations without first copying data to the CPU memory space.

OpenCL 2.0 features At the time of writing (driver version 1642.5) the OpenCL 2.0 features such as workgroup functions and device enqueue impose significant performance impact. Atomic locks do not result in significant performance impact. The library uses Shared Virtual Memory buffers to facilitate easy switch between OpenCL and HSA backends. The ViennaCL types are constructed from *cl_mem* representation of the SVM buffers.

Chapter 3

k-Nearest Neighbours

Introduction

k-Nearest Neighbours method[33] is a non-parametric method used for the classification and regression. It computes a given instance distance to the examples with the known label and either provides a class membership for the classification which is a class most common among nearest neighbours or an object property value which is an average of the nearest neighbours. The error rate bound by twice the Bayes error if the number of examples approaches infinity. Online implementation of the algorithm uses a sliding window of example instances updated by the data stream. Window size, instance dimensionality and allowed error bounds define the optimal approach to solving k-Nearest neighbours problem. The *Exhaustive Search* has a high computational complexity of the queries but provides a constant update time. *Exact Clustering Methods* partition the search space to achieve logarithmic query complexity at expense of the window update time. *Approximate Clustering Methods* reduce search space dimensionality to provide an approximate result with a given error bound. The GP GPU computing may be used to accelerate their runtime though success for discrete GPUs depends on developer ability to eliminate branching, optimize memory access, and avoid excessive Host-GPU transfers. The latter poses a most significant problem for online k-Nearest neighbours

implementations. This chapter reviews the *Exhaustive Search*, some of the *Exact Clustering Methods* and *Approximate Clustering Methods*, provides notes on GP GPU implementation.

Exhaustive Search

The exhaustive approach to Nearest Neighbour search is to compute distance to each instance present in the sliding window. The computational complexity of the query $O(N^d)$ where N - number of instances and d - number of attributes. The GP GPU implementation of the exhaustive search consists of distance calculation and selection phase. The distances to the query are computed as a vector-matrix multiplication or if several queries are processed at once as a matrix-matrix multiplication. GPU implementation of those routines is available as a part of libraries implementing Basic Linear Algebra Subroutines (BLAS) [6][4][72]. The selection phase finds nearest to the query out of all the computed distances. Sismanis et. al [66] provide time complexity of reduced sort algorithms, evaluate their performance on GPU and propose to interleave distance calculation and sorting phases to hide latency. The data for the distance calculation should be offloaded to GPU while it performs the sorting phase. The sliding window is partitioned if needed across multiple GPUs according to the GPU memory capabilities and does not use instances spatial information.

Distance Calculation

The optimal implementation depends on the size of the window and number of attributes present [68]. For the small instance size (≤ 100) and windows less than 10^4 elements one thread per instance implementation will provide the best solution. Best all around distance calculation should apply different strategies depending on the window size and number of attributes [68]. The alternatives are presented in the Figure 3.1.

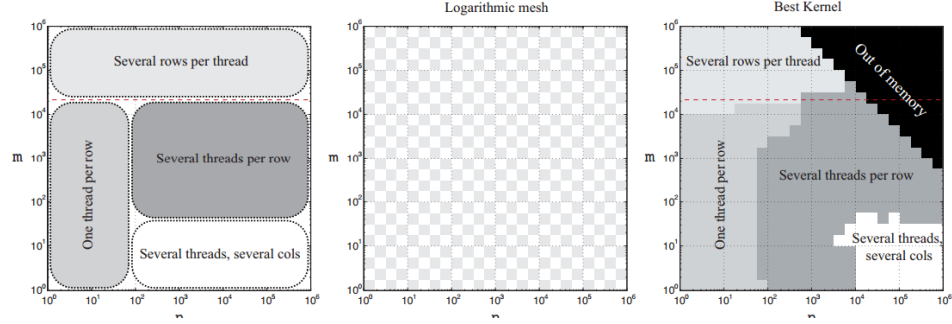


Figure 3.1: Left; Four matrix-vector multiplication kernels designed to perform well at different shapes $m \times n$ of A . Middle; Tuning mesh. Right; Best kernel in practice. The dashed line indicates the minimum 21504 rows needed in A for full occupancy of the Nvidia Tesla C2050 card in a one-thread-per-row kernel. Note the logarithmic axes. Reproduced from High-Performance Matrix-Vector Multiplication on the GPU by Hans Henrik Brandenburg Sørensen[68]

The distance calculation primitive provided as part of the master thesis uses one thread per row approach.

Selection

Alabi, et.al evaluated different selection strategies based on bucket sort algorithm and Merrill-Grimshaw implementation of radix sort[23].

The selection phase may be interleaved with the distance calculation to utilize both CPU and GPU cores and benefit from the better sort performance on low window sizes[52].

Figure 3.2 shows measured performance of different selection strategies - merge sort from AMD Bolt library[2], bitonic sort similar to reference AMD implementation and radix select based on Alabi, et. al. implementation[23]. The CPU sort and choose is a clear winner for small (65535) window sizes. The merge sort should be applied to sub 2^{25} windows and radix select (with data copy) should be used for larger window sizes. *update with clogs - i have it implemented*

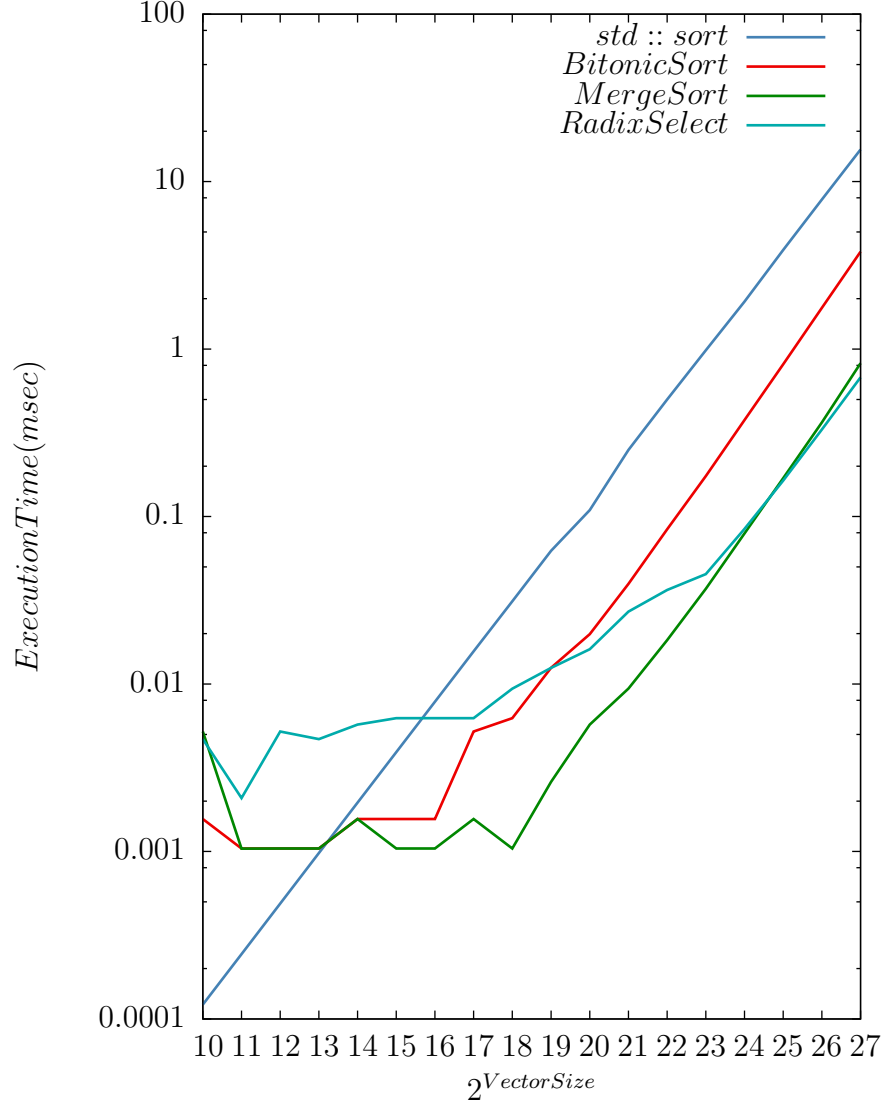


Figure 3.2: Selection Algorithm Performance for K=128. Test configuration GPU R9 390, CPU AMD A8-7600, AMD Catalyst version 15.20.

Conclusion

The *Exhaustive Search* is a basic building block of the k-nearest neighbours algorithms. It is unavoidable if we need to obtain exact solution and is used to refine approximate methods results.

Exact Clustering Methods

The clustering techniques are widely used to limit number of distance calculations needed for nearest neighbour search. Space partitioning by vector

dimensions is used by $k - d$ tree method[41], random projection tree[40] provides a data structure splitting search space along random vectors. Metric trees such as ball tree[56], cover tree[26], nearest ancestor tree[?], random ball cover[30] provide solution for finding nearest neighbours in general metric space by organizing data points in groups around some centroids.

k-d Tree

The $k - d$ tree [41] is a balanced binary tree where each node represents a set of points $P \in \langle p_1 \cdot p_n \rangle$ and its children are disjoint and almost equals sized subsets of P . The tree is constructed top-down, the initial set of points is split along the widest dimension or using other criteria until the predefined number of points in child nodes is reached. The tree can be constructed in $O(n \log n)$ time and occupies linear space. Weber *etal*[73] have shown that $k - d$ tree is outperformed by the exact calculation at moderate dimensionality ($n > 10$) and results in full processing of the data points if the number of dimensions is large enough. The $k - d$ tree requires $N \gg 2^k$ points to be more effective than exhaustive search.

The listing of the $k - d$ tree construction and nearest neighbours search pseudocode is shown in the Figure 3.3. The parallel $k - d$ tree construction on GPU utilizes breadth-first approach[74][64] - the $k - d$ tree is constructed top-down with the split criteria computed in parallel for all nodes at the specific level. The priority queue based nearest neighbours search using k-d trees as shown in Figure 3.3 does not benefit much from the GP GPU parallelism due to the branch divergence and irregular memory access patterns[43]. The $k - d$ tree search approach presented by Gieske et. al focuses on parallel execution of nearest neighbour queries in a lazy fashion. The query points are accumulated in the leaf nodes of the kd-tree until enough of them is present and then processed as a batch. This solves an issue of the GPU underutilization and low performance if leaf nodes are processed sequentially for each example[43]. Other approach would be to compute distances to the leaf nodes split planes[74]

```

1  tree_node create_tree(pointList, level)
2  {
3      int dim = select_dim(pointList); // select split dimension according to
4      // pre-defined criteria, e.g. level mod total_dimensions
5      splitVal = select_split_value( pointList, dim); // select split value
6      // according to pre-defined criteria
7      // e.g. median value of point[dim]
8
9      left = {};
10     right = {}
11     for (point : pointList )
12     {
13         if (point[dim] > splitVal)
14             right += point;
15         else
16             left += point;
17     }
18     node = {
19         .location = splitVal,
20         .dim = dim,
21         .left = create_tree(left, level + 1),
22         .right = create_tree(right, level + 1)
23     };
24     return node;
25 }
26
27 void search(Heap nearest_neighbours, tree_node root, point p)
28 {
29     if (root.is_leaf())
30         nearest_neighbours.update(root);
31     else
32     {
33         split = root.location;
34         dim = root.dim;
35         if (p[dim] < split ) // search "closest" node
36             search(nearest_neighbours, root.left, p)
37         else
38             search(nearest_neighbours, root.right, p)
39
40         distance_to_split_plane = abs(split-p[dim]);
41         distance_to_point = abs(nearest_neighbours.furthest_point()[dim]-p[dim])
42         if (distance_to_point >= distance_to_split_plane) // outer radius of NN heap
43             intersects the split plane
44             {
45                 if (p[dim] < split )
46                     search(nearest_neighbours, root.right p)
47                 else
48                     search(nearest_neighbours, root.left, p)
49             }
50     }
51 }

```

Figure 3.3: k-d tree construction and NN-search pseudocode

to provide a short-list of the leaf nodes for k -nearest neighbours search.

Random Projection Trees

The $k - d$ tree provides an effective partitioning mechanism for low data dimensionality but suffers in higher dimensions[73]. Many machine learning problems that are expressed in high dimensional space has lower intrinsic dimension as shown in Figure 3.4. Random projection tree exploits this fact

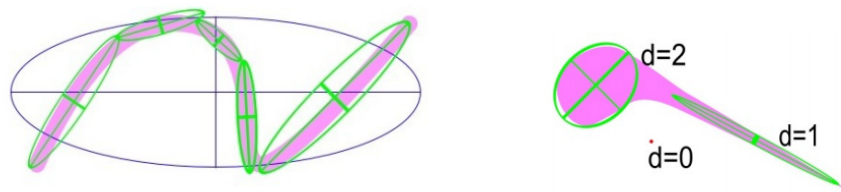


Figure 3.4: Distributions with low intrinsic dimension. The purple areas in these figures indicate regions in which the density of the data is significant, while the complementary white areas indicate areas where data density is very low. The left figure depicts data concentrated near a one-dimensional manifold. The ellipses represent mean+PCA approximations to subsets of the data. Our goal is to partition data into small diameter regions so that the data in each region is well-approximated by its mean+PCA. The right figure depicts a situation where the dimension of the data is variable. Some of the data lies close to a one-dimensional manifold, some of the data spans two dimensions, and some of the data (represented by the red dot) is concentrated around a single point (a zero-dimensional manifold). Reproduced from Learning the structure of manifolds using random projections by Freund Yoav et al.[40]

by splitting data along randomly chosen unit vectors as opposed to splitting along dimension axes in $k - d$ tree method as shown in Figure 3.5[40]. The method performs a one dimensional random projection of the data points and splits them at the median of the projections.

The random projection tree split rules are presented in Figures 3.6, 3.7.

The efficient implementation of the Random Projection Tree depends on the ability to perform random projections with low computational complexity as it is both required for tree construction and queries.

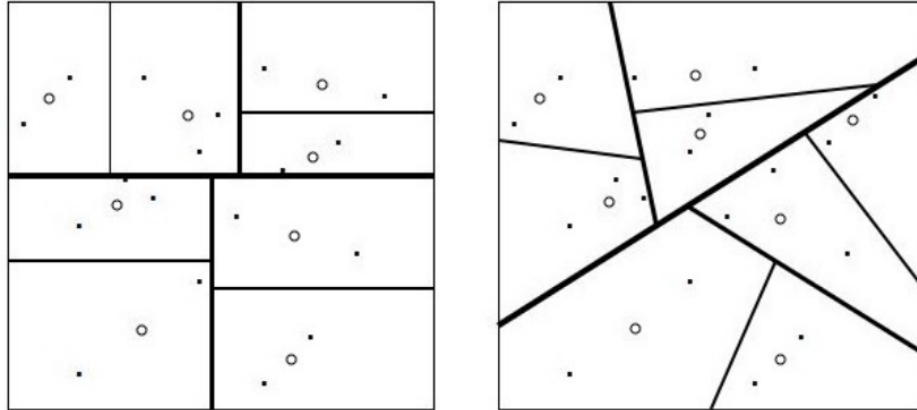


Figure 3.5: Left: Partitioning produced by k-d tree. Right: Partitioning produced by Random Projection Tree. Reproduced from Learning the structure of manifolds using random projections by Freund Yoav et al.[40]

```

1  tree_node random_tree_max(pointList, num_dimensions)
2  {
3      v = random_vector(num_dimensions);
4
5      x = pointList[ random() ];
6      y = max (distance( y in pointList, x) );
7      sigma = uniform_random(-1;1) * 6 * distance(x,y) / sqrt(num_dimensions);
8      split = median ( dot(v, x in pointList)+ sigma ) ;
9      left = {}
10     right = {}
11     for (x in pointList)
12     {
13         if (dot(v,x) <= split)
14             left += x;
15         else
16             right +=x;
17     }
18     node = {
19         .vector = v,
20         .split = split,
21         .left = create_tree(left, num_dimensions),
22         .right = create_tree(right, num_dimensions)
23     };
24     return node;
25 }

```

Figure 3.6: Random Projection Tree Pseudocode - Random Tree Max [34]

```

1
2  tree_node random_tree_mid(pointList, num_dimensions, c)
3  {
4      diameter = max( distance(x in pointList, y in pointList));
5      avg_diameter = mean(distance(x in pointList, y in pointList));
6      if ( diameter <= c* avg_diameter)
7      {
8          v = random_vector(num_dimensions);
9          split = median( dot(x in pointList, v) );
10         left = {}
11         right = {}
12         for (x in pointList)
13         {
14             if (dot(v,x) <= split)
15                 left += x;
16             else
17                 right +=x;
18         }
19         node = {
20             .rule_type = dotproduct
21             .vector = v,
22             .split = split,
23             .left = create_tree(left, num_dimensions),
24             .right = create_tree(right, num_dimensions)
25         };
26         return node;
27     }
28     else
29     {
30         meanPoint = mean(x in pointList)
31         split = median( distance(x in pointList, meanPoint));
32         left = {}
33         right = {}
34         for (x in pointList)
35         {
36             if (distance(x, meanPoint) <= split)
37                 left += x;
38             else
39                 right +=x;
40         }
41         node = {
42             .rule_type = distance
43             .mean = meanPoint,
44             .split = split,
45             .left = create_tree(left, num_dimensions),
46             .right = create_tree(right, num_dimensions)
47         };
48         return node;
49     }
50 }
51 }

```

Figure 3.7: Random Projection Tree Pseudocode - Random Tree Median[34]

Random Projection Seminal paper by Johnson and Lindenstrauss[46] established that for euclidian spaces any $x \in \mathbb{R}^n$ can be embedded into \mathbb{R}^k with $k = O(\log n / \epsilon^2)$ by projecting x in \mathbb{R}^k using projection $k \times n$ matrix Φ without distorting inter-point distances by more than $(1 \pm \epsilon)$ and $k \geq O(\log n)$. Johnson and Lindenstrauss[46] has shown that Johnson-Lindenstrauss condition holds for matrices with following properties: “S”pherical symmetry - For any orthogonal matrix $A \in O(d)$, ϕA and ϕ have the same distribution. Orthogonality - rows are orthogonal to each other Normality - the rows are unit-length vectors [21]

The lower bound of k was refined by in several papers[39][35][45][20] and with Dasgupta and Gupta[35] proving it to be $k \geq 4(\epsilon^2/2 - \epsilon^3/3)^{-1} \ln n$ for $\epsilon \in (0, 1)$. For high n this bound will still be too large to effectively employ low dimensionality search methods such as $k - d$ tree. An alternative would be to utilize very low dimensional space and then use disjunction to find desired result. This approach is essentially an iterative random projection tree search where the dataset is split along leaf nodes.

The efficient implementation of the random projection-based algorithms requires a simple approach to construct ϕ and a way to compute projection faster than naive multiplication of data point by $k \times n$ matrix. Achlioptas[20] achieved relatively sparse transformation matrix for random projection by proving that Johnson-Lindenstrauss condition holds if elements of the projection matrix are chosen independently according to the following distribution:

$$\begin{cases} +(n/3)^{-1/2}, P = 1/6 \\ 0, P = 2/3 \\ -(n/3)^{-1/2}, P = 1/6 \end{cases}$$

This method provides a 3-fold speedup over originally proposed[46] since 2/3 of the transformation matrix elements are zero. Nir Ailon and Edo Liberty[22] have developed an almost optimal random projection transformation with runtime of $O(n \log n)$ as opposed to $O(kn)$ of the naive implementation. The main

idea of the method is the application of the Heisenberg principle in its signal processing interpretation that both signal and its spectrum can not be both sharply localized. Thus applying Fourier transform to the sparse input vector will increase its support and will allow to make the transformation matrix even more sparse. To prevent the opposite - sparsification of the dense vector the input data elements signs are randomly inverted with probability $1/2$. The sparse transformation matrix used to complete the random projection[21] can be replaced by subsampled fourier transform[22]. Nir Ailon and Edo Liberty define $k = O(\delta^{-4} \log(N) \log^4 n)$ that will preserve input vector norms by a given relative error δ [22]. The reference implementation used in the work is based on Gabriel Krummenacher implementation of subsampled randomized Fourier transform <https://github.com/gabobert/fast-jlt/tree/master/fjlt>. This algorithm is well suited for GPU implementation as it consists of FFT followed by element-wise operation. The last step (select random D elements) introduces irregular memory access that can not be worked around unless multiple transformations are performed in parallel. Figure 3.8 shows comparative performance of dense matrix multiplication for random projection and Fast Johnson-Lindenstrauss transform. For the selected hardware configuration the latter starts to outperform matrix multiplication starting from $N \geq 16384$. It should be noted that FLJT has lower memory requirements than $O(kd)$ as it does not require to store dense transformation matrix and thus capable of projecting higher dimensional data on the same hardware.

Approximate Clustering Methods

The nearest neighbours search methods in high dimensional space provides little benefit over exhaustive search where an exact distance is computed to each point in the database[73][25]. The approximate methods provide means to overcome this limitation by solving the problem of finding neighbours whose distance from the query point are at most $c > 1$ times greater than distance to

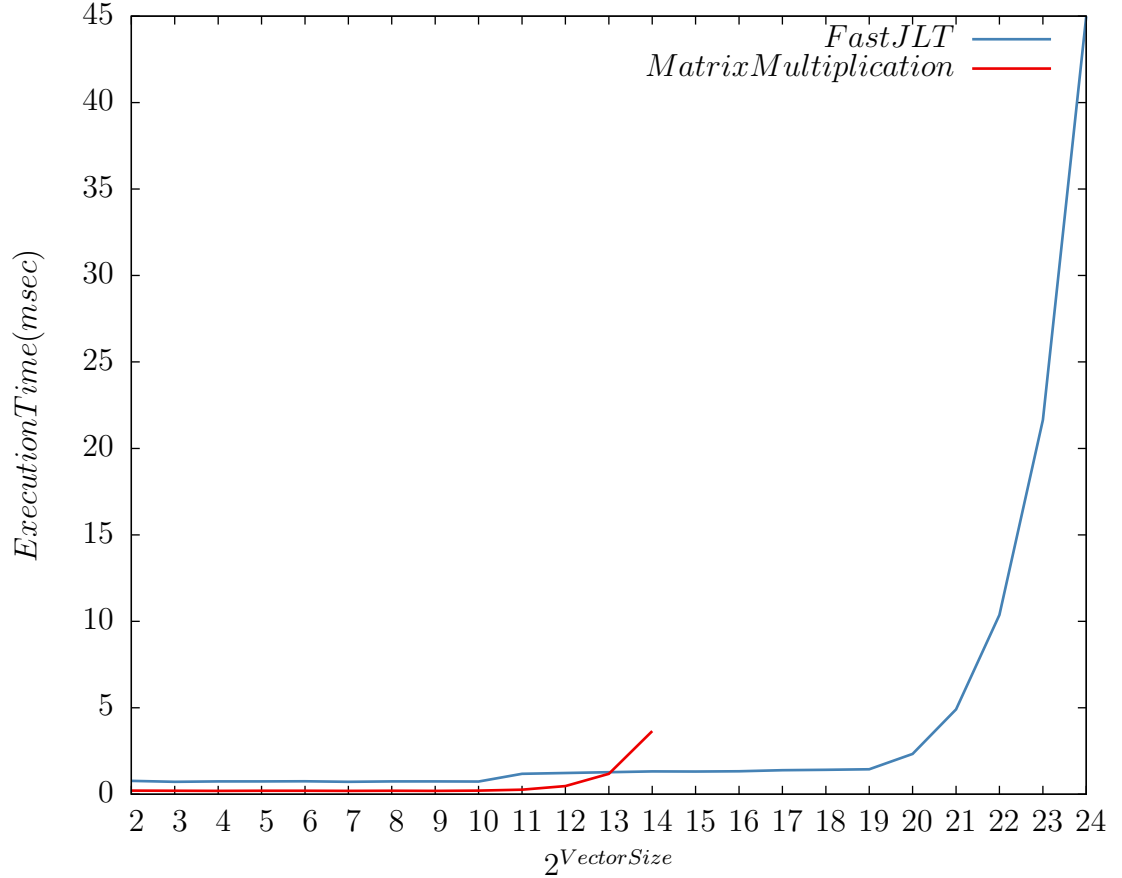


Figure 3.8: Fast Johnson-Lindenstrauss transform[22] vs. sparse matrix implementation[20] using ViennaCL. Test configuration GPU R9 390, CPU AMD A8-7600, AMD Catalyst version 15.20.

the closest neighbour. The approximate solution can be used to find exact one by computing distance to each approximate nearest neighbour and choosing closest ones. Modern approximate method use dimensionality reducing techniques such as random projection[58] and data set ordering using space filling curves[58][48][49] to improve query performance.

Locality Sensivity Hashing

Locality Sensivity Hashing[45] is a method that capitalizes on the idea that exist such hash functions $h(x), x \in \mathbb{R}^d$ that for points $p, q \in \mathbb{R}^d$, radius R and

approximation constant c

$$\begin{cases} \|p - q\| \leq R, P[h(p) = h(q)] \geq P_1 \\ \|p - q\| \geq cR, P[h(p) = h(q)] \leq P_2 \end{cases}$$

where probability $P_1 > P_2$. The LSH algorithm uses a concatenation of $M \ll d$ such functions to increase difference between P_1 and P_2 [45]. Initially it was proposed to use Hamming distance as this function satisfies required properties[45]. Later it was shown that other families of hash functions such as l_p distance[36], Jaccard coefficient [28][29], angular distance(random projection)[31] are locally sensitive. The algorithm selects L concatenations of the hash functions and uses them to transform input dataset points $v \in \mathbb{R}^d$ into lattice space \mathbb{Z}^M storing them as L hash tables. The exact query is performed by concatenating contents of the L bins corresponding to the hash codes of the query and computing exact distance. The approximation is obtained by stopping as soon as k points in cR distance from the query point is found. The method is GP GPU friendly as hash codes of the data points can be computed in massively parallel manner. The state of art GP GPU hash maps use *coockoo hashing*[24] a robust data structure allowing constant time retrieval and eliminate hash collisions that lead to the branch divergence due to the need to iterate over items in the hash bucket but require full hash map rebuild to resolve a collision.

Space Filling Curves

The space filling curves[60] are curves that traverse all points of n -dimensional space in a given region providing a mapping from n -dimensional to 1-dimensional space. The GP GPU nearest neighbours algorithms[49][48][58] utilize z -order curve introduced by G. Morton in 1966 - an ordered list of numbers composed by interleaving bits of instance attributes[?]. This curve is used due to the fact that mapping can be constructed in $O(d)$ time and is easily implemented

on GPU. The sample z-order curve is shown in Figure 3.9. The z-order curve

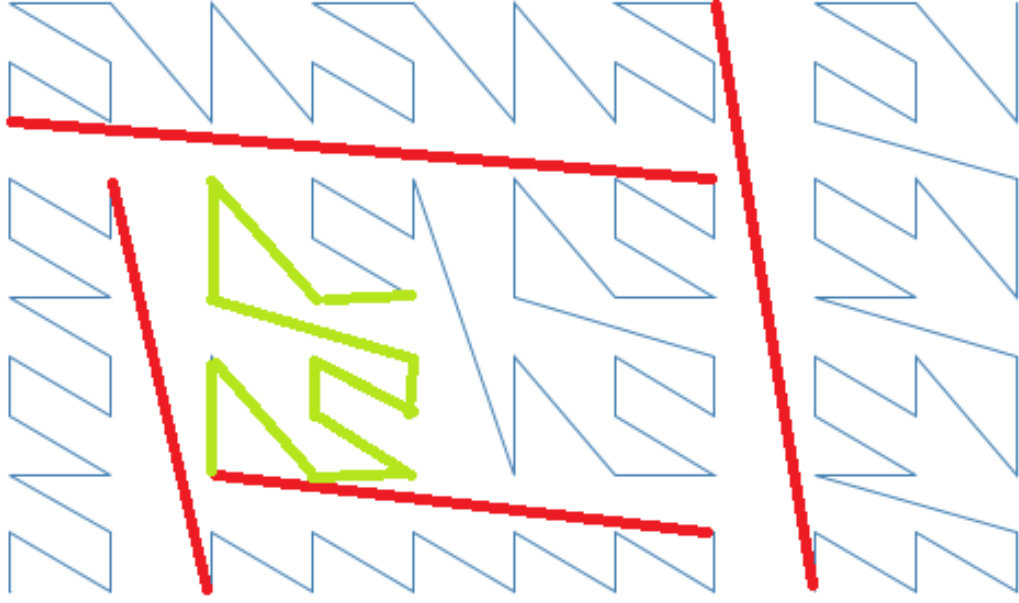


Figure 3.9: Z-Order Curve. Red lines highlight some region jumps. Green shows locality-preserving region.

mostly preserves data locality - points that are close in the n -dimensional space are also close together along the curve, but as shown in Figure 3.9 the z-order curve has jumps and to compensate for them existing GP GPU algorithms generate several curves from the input data shifted several times by some random vector[49][48] to move the points into its locality-preserving regions. Shift search method by Li et al.[48] experimentally quantifies the shift value.

Sieranoja S.[65] proposed a z-order lookup table mapping algorithm for arbitrary number of dimensions. It can be translated into GPU implementation with minimal changes.

k-Nearest Neighbours using HSA architecture

The HSA architecture opens new possibilities for GP GPU acceleration of the nearest neighbours search by embedding parallel primitives such as FJLT and exhaustive search into existing algorithms without redesigning them to conform to the GPU architecture.

k-d tree queries The k-d tree queries require all or nothing approach - the algorithm design and data structures should be implemented in GP GPU specific manner to obtain higher throughput and expense of the latency[74][43]. HSA allows to insert *exhaustive search* parallel primitive to search the tree leaves in an efficient manner. The discrete GPU would not achieve same computing node/speedup ratio due to the irregular memory access to the instance data. Figure 3.10 shows test results in synthetic test using MOA *RandomRBFGenerator* data stream with 3 numeric attributes. The test performance of the on-chip GPU of AMD A8 processor is on par with R9 390 despite latter capable of running 2560 threads as opposed to 384 of the on-chip GPU.

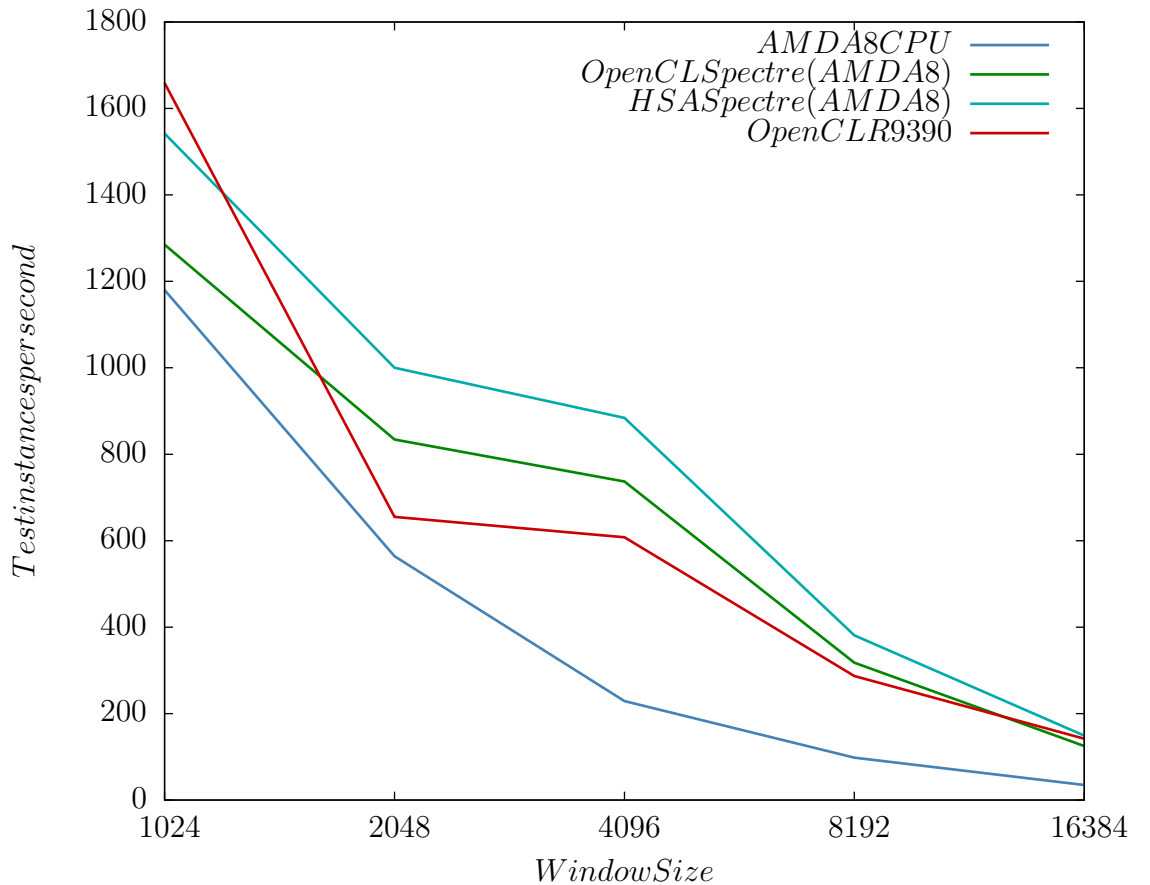


Figure 3.10: k-d tree test speed

Random Projection Tree Construction and Queries The random projection tree can be built in the breadth first manner similar to existing GPU implementations of k-d tree algorithm. This approach requires recalculation of the whole tree structure for each update of the sliding window. The online update of the random projection tree may benefit from offloading projection and thresholding to GPU similar to k-d tree example above.

TODO: Figure - synthetic test

Z-Order-based queries The significant problem in z-order query evaluation is a maintenance of the candidate lists for large k [48]. The cost of GPU-host transfer in this case is prohibitive and Li et al work around it by executing multiple queries at once so that it is more beneficial to perform sort as opposed to search[48]. The HSA implementation may keep the candidate list on CPU and perform n-ary GPU searches thus reducing query latency.

The approximate z-order based method can utilize random projections to shorten morton code length while preserving the relative distances to speed up queries and sliding window updates.

TODO: Figure - synthetic test

Conclusion The use of HSA platform may be advantageous to the discrete GP GPU computation in scenarios that require frequent host-GPU interactions or require low latency and thus make batching undesirable. This has been validated using synthetic tests that put discrete GPU in disadvantage by artificially increasing transfer/computation ratio. An example of the algorithm evaluation on the existing dataset such as MNIST[47] and real-life data streams is needed to obtain estimation of the HSA platform performance compared to discrete GP GPU computation.

Chapter 4

Stochastic Gradient Descent

4.1 Introduction

Stochastic gradient descent algorithm is an iterative optimization method used in a wide variety of machine learning tasks. It minimizes the objective function $Q(w, x)$ dependant on the set of parameters w and set of examples x by updating parameters along the gradient of the randomly chosen example x_i : $w = w - \lambda \nabla Q(w, x_i)$, where λ is the iteration step also called the learning rate. In classical learning application $Q(w, x)$ is a prediction loss function and the method aims to find the set of parameters w that provides a local minimum of an error on a training set. In online learning application $Q(w, x)$ is a regret function and the method aims to provide a best approximation of $y(x)$ where y is a classification associated with example x . The algorithm may use an average gradient of several examples - a *mini-batch* to achieve lower variance. This chapter gives an overview of the approaches for stochastic gradient descent parallelisation, describes Hogwild! and 1bit SGD algorithms used in this work, and discusses HSA algorithm implementation.

4.2 Approaches for SGD parallelisation

The stochastic gradient descent algorithm is inherently sequential. The ways to parallelise it across computing clusters are explored in a number of papers[42][54][50][75][62].

Langford et al[75] proposed a pipelined approach to SGD parallelisation. The input vector is divided into partitions that are processed in parallel by slave worker threads producing subgradients on each step. The subgradients are sent to the master worker that recomputes parameter vector and distributes it to the slave threads. The algorithm is sensitive to the delay - using parameter vector outdated by 1000 iterations increases the error from 2^{-10} to 2^{-6} on TREC Spam Track dataset[32], though smaller delays such as 10 has no effect on convergence[75]. This highlights the problem for parallel implementation of SGD - one has to minimize communication to achieve maximum computation speedup without introducing critical delay that will hurt coverage of the algorithm. The subsequent works [62][42][76] provided a way to balance the parallelism and delay either by *model parallelism* - partitioning data into independent sets and processing them in parallel [42][76], delaying update communication[62][38][57] or by removing a synchronization requirement for parameter update[54]. The state of art cluster methods use a combination of all those techniques[38][57] to perform large scale optimization.

The single node stochastic gradient implementation often deals with *data parallelism* - for a given mini-batch its gradients are computed in parallel. This operation requires a parallel calculation of the objective function given model parameters and a set of example instances - essentially a vector (parameters) by matrix (examples) multiplication to obtain current approximation. Davis and Chang[37] explore single-precision matrix-vector multiplication and conclude that due to the modern CPUs gaining last-level cache sizes and external memory bandwidth, increasing data sizes of computational problems and constrained memory size of the consumer GPU cards it creates a barrier to adoption of GP-GPU, though they suggest that on-chip GPUs may be excellent building blocks for future heterogeneous processors.

4.3 Hogwild!

The *Hogwild!* algorithm[54] uses an assumption that updates to the parameter vector w can be performed in a lock-free manner since each individual update only affect a small portion of it. The algorithm uses independent workers that have an access to the shared parameter vector w . Each worker samples uniformly at random an example x and computes an update vector $\lambda \nabla Q(w, x_i)$. The worker then proceeds to update each element of the parameter vector w in an atomic fashion. Due to multiple workers simultaneously updating the parameter vector, updates to the individual coefficients may be lost. The algorithm achieves nearly linear speedups on KDD Cup 2011[?] and Netflix prize[?] datasets[54].

4.3.1 Best Ball Optimization

The implementation of the Hogwild!algorithm (<http://i.stanford.edu/hazy/victor/Hogwild/>) contains a *best ball* autotuning method. The user picks a range of model parameters (e.g. learning rate) and the algorithm evaluates the corresponding models in parallel. At the end of the epoch the model with the lowest harmonic mean of the root mean square error is selected and its parameter vector is propagated to all other models.

4.3.2 Backoff scheme

The Hogwild! algorithm uses a constant diminishing learning rate. The algorithm uses a global synchronization point at the end of epoch of K iterations to reduce it by a constant β and continue running for the next $\beta^{-1}K$ iterations.

4.4 1bit SGD

The Hogwild! algorithm is inherently non-deterministic. The updates of the parameter vector are performed in a lock-free manner and may be lost should

several updates contain a modification of the same parameter vector element w_i . 1Bit SGD[62] approaches communication bottleneck from the different angle. It works on the same assumption as Hogwild! - the updates from each minibatch only affect a small portion of the parameter vector and adds a further constraint - only updates that are greater than a certain threshold should be communicated. In 1Bit SGD workers exchange gradient updates quantized to one bit - that is all workers share quantization constant τ and communicate gradient vector element that should be updated by this constant. The difference between computed value and quantization constant is stored locally by the worker and added to the next iteration gradient update[62].

This approach allows nearly linear speed-ups in cluster setting[70] as opposed to previous results such as[63].

4.5 Stochastic Gradient Descent using Open-CL/HSA architecture

This work implements the Hogwild! algorithm for linear models on Heterogeneous System Architecture and OpenCL platforms and compares its performance with the baseline single-threaded implementation.

The 1bit SGD thresholding is used to minimize the thread contention - only updates exceeding the quantization parameter τ are written to the shared parameter vector. This will prevent situation when near-zero update overwrites a bigger one.

The parameter τ is selected as per[62] - to minimize the square quantization error for a given column and is recomputed at the end of each algorithm iteration.

The algorithm is parametrized by number of minibatches it processes simultaneously - B . The residual quantization error is stored in matrix E with each row representing residual error for the respective mini-batch. The parallel implementation processes as follows. Sparse matrix-vector multiplication is

used to obtain vector of dot-products. For each dot product a loss function and update is computed. The bucket reduce operation is used to obtain cumulative weight update value and average it across minibatch for every minibatch. Finally a weight update kernel is launched with each thread processing separate W_{ib} - individual weight from a specific mini-batch $b \in B$. The parameter vector is updated by τ using atomic add function if the update value exceeds threshold, or added to corresponding element of residual error matrix E . Based on E and iteration number an average quantization error is calculated for each column and τ used for the next iteration is updated.

4.5.1 OpenCL Memory Transfer

The OpenCL implementation uses parallel data transfer during training as shown in Figure 4.1. The data for the next batch is transferred while GPU processes the previous one to hide latency incurred by host-GPU memory transfer.

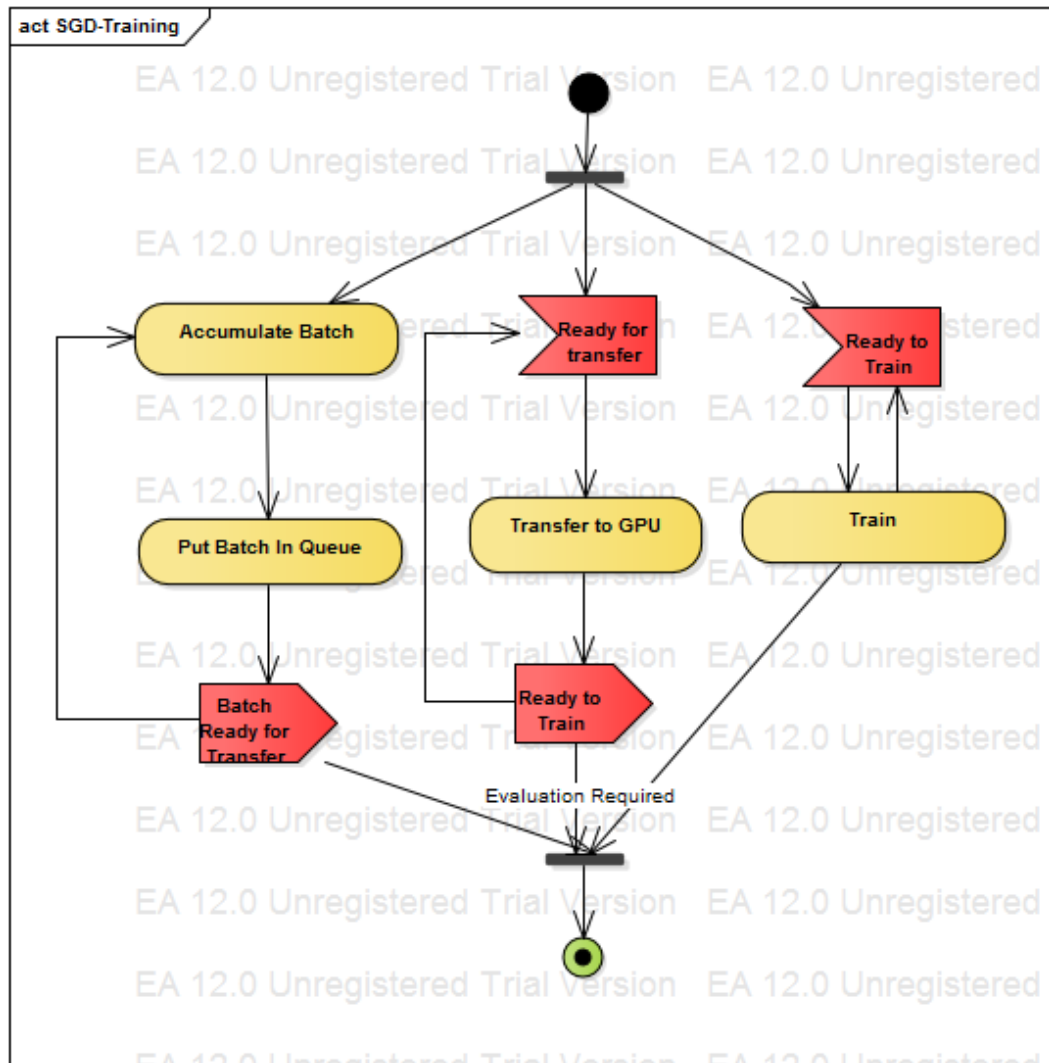


Figure 4.1: OpenCL SGD training process.

Chapter 5

Experimental Results

k-Nearest Neighbours

kNN Naive implementation

The first version of k-NN algorithm used Java implementation with JavaCL library. The results in Figure 5.1 provide overview of the achieved speedups.

Stochastic Gradient Descent

The SGD implementation training speedup for sparse instances with the double data type is presented in the Figures 5.2, 5.3. The CPU sampling profiling showed that most of the time is spent inside buffer commit function responsible for copying data from WEKA instance class into the sparse matrix. As the GPU implementation uses batching to achieve training speedup the latency is $5-10x$ worse than CPU. The profiling of the sample run on a Spectre device of OpenCL algorithm implementation with 1024 non-null attributes in an instance and 8192 instances in a batch shows that only 32% is spent training the classifier, rest is spent in the CPU-GPU data transfer.

The sparse matrix OpenCL SGD implementation is bounded by the memory copy operation as shown by the profiling data and provides a limited speedup.

TODO Provide zero-copy (HSA) speedups *TODO* Provide float speedups

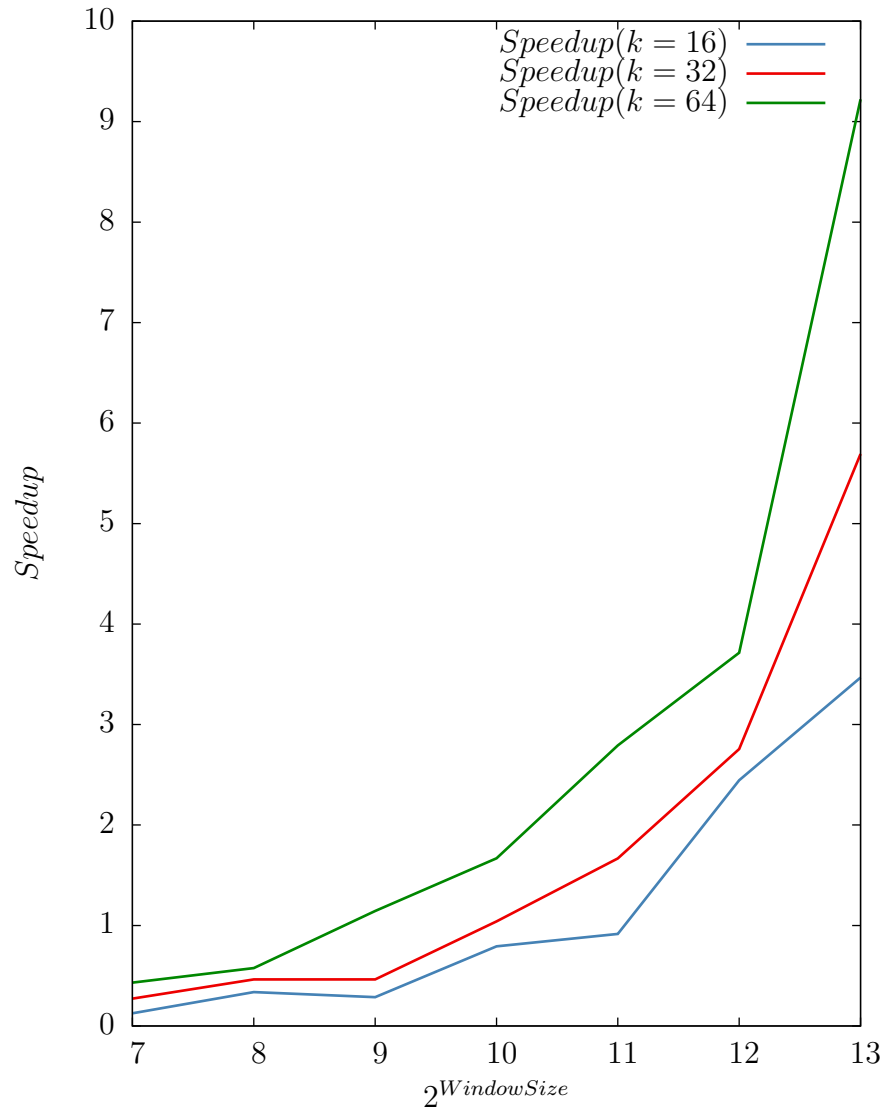


Figure 5.1: Naive KNN implementation (Radeon Mobility HD5730)

and impact on the accuracy of the algorithm for the known data sets.

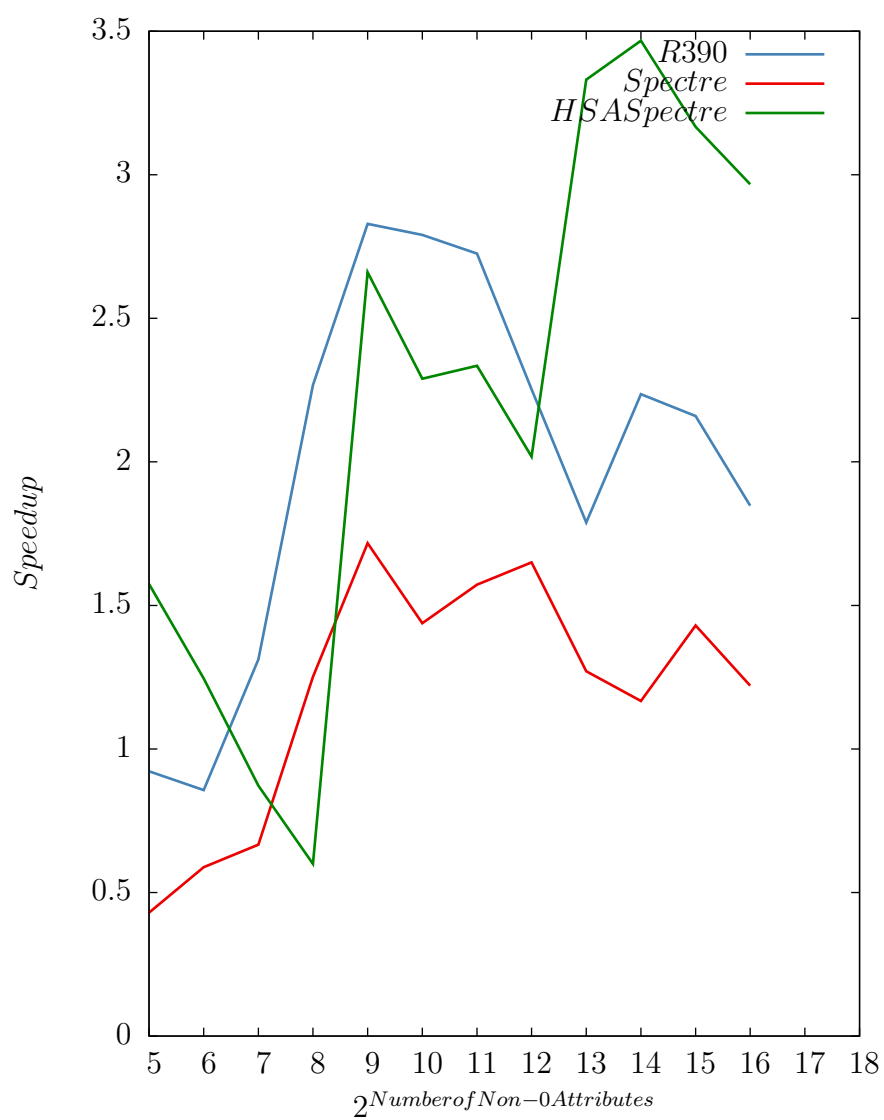


Figure 5.2: Stochastic Gradient Descent Training Speedup (AMD R9 390, Spectre vs MOA implementation on AMD A8-7600)

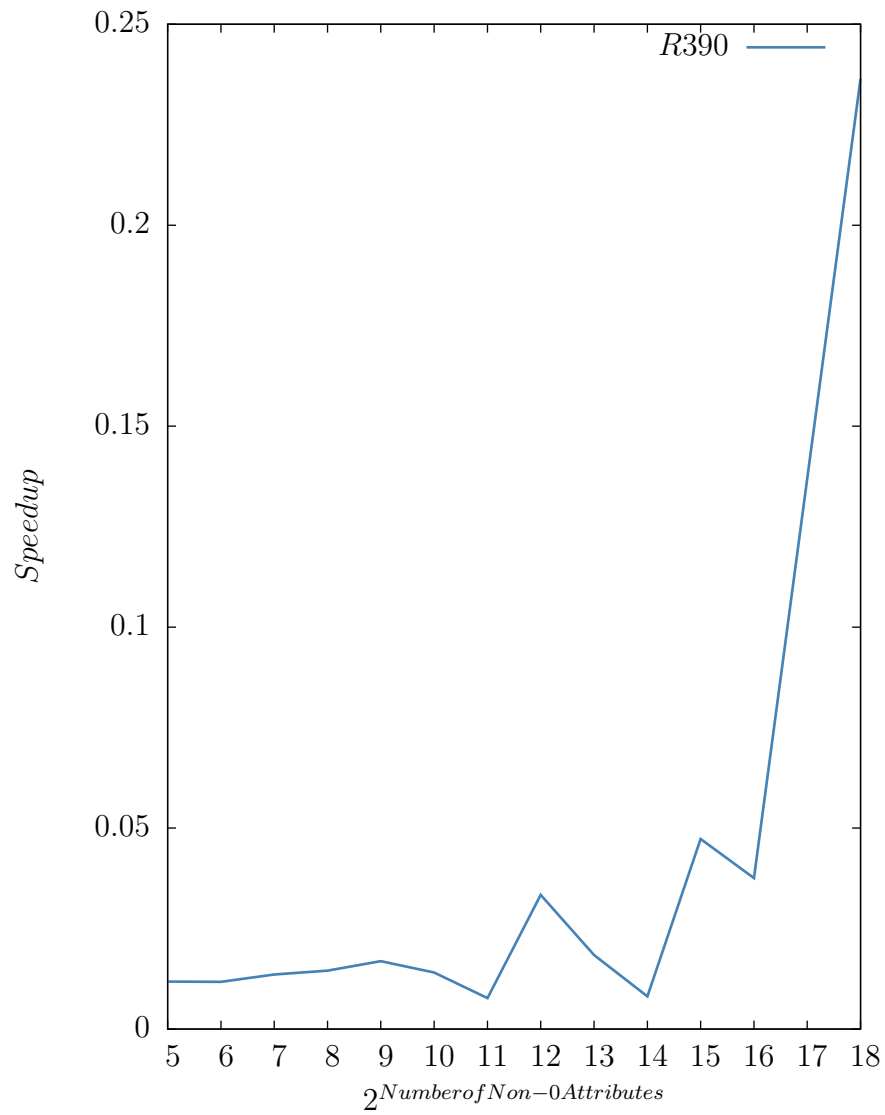


Figure 5.3: Stochastic Gradient Descent Training Latency (AMD R9 390 vs MOA implementation on AMD A8-7600)

Chapter 6

Conclusions and Future Work

TODO

References

- [1] aparapi - api for data parallel java. allows suitable code to be executed on gpu via opencl. <https://code.google.com/p/aparapi/>.
- [2] Bolt c++ template library. <http://developer.amd.com/tools-and-sdks/opencl-zone/bolt-c-template-library>.
- [3] C++ amp overview. <https://msdn.microsoft.com/en-us/library/hh265136.aspx>.
- [4] clmath - amd. <http://developer.amd.com/tools-and-sdks/opencl-zone/amd-accelerated-parallel-processing-math-libraries/>.
- [5] Compute shader overview. <https://msdn.microsoft.com/en-us/library/windows/desktop/ff476331%28v=vs.85%29.aspx>.
- [6] cublas. <https://developer.nvidia.com/cuBLAS>.
- [7] Hsafoundation. <https://github.com/HSAFoundation>.
- [8] Hsafoundation/cloc. <https://github.com/HSAFoundation/CLOC>.
- [9] Java se 7 java native interface-related apis and developer guides. <http://docs.oracle.com/javase/7/docs/technotes/guides/jni/>.
- [10] java.util.stream (java platform se 8). <http://docs.oracle.com/javase/8/docs/api/java/util/stream/package-summary.html>.
- [11] Khronos opencl registry. <https://www.khronos.org/registry/cl/>.
- [12] Nvidias next generation cuda(tm) compute architecture:fermi. http://www.nvidia.com/content/PDF/fermi_white_papers/NVIDIA_Fermi_Compute_Architecture_Whitepaper.pdf.
- [13] Openjdk: Project sumatra. <http://openjdk.java.net/projects/sumatra/>.
- [14] Openmp.org. <http://openmp.org/wp/>.

-
- [15] Ptx isa :: Cuda toolkit documentation. <http://docs.nvidia.com/cuda/parallel-thread-execution/index.html>.
 - [16] Renderscript—android developers. <http://developer.android.com/guide/topics/renderscript/compute.html>.
 - [17] Single instruction, multiple threads. http://en.wikipedia.org/wiki/Single_instruction,_multiple_threads#cite_note-spp-1.
 - [18] Spir - the first open standard intermediate language for parallel compute and graphics. <https://www.khronos.org/spir>.
 - [19] Hsa platform system architecture specification. <http://www.hsafoundation.com/?ddownload=4944>, 2014.
 - [20] Dimitris Achlioptas. Database-friendly random projections: Johnson-lindenstrauss with binary coins. *J. Comput. Syst. Sci.*, 66(4):671–687, June 2003.
 - [21] Nir Ailon and Bernard Chazelle. The fast johnson-lindenstrauss transform and approximate nearest neighbors. *SIAM Journal on Computing*, 39(1):302–322, 2009.
 - [22] Nir Ailon and Edo Liberty. An almost optimal unrestricted fast johnson-lindenstrauss transform. *ACM Transactions on Algorithms (TALG)*, 9(3):21, 2013.
 - [23] Tolu Alabi, Jeffrey D. Blanchard, Bradley Gordon, and Russel Steinbach. Fast k-selection algorithms for graphics processing units. *J. Exp. Algorithmics*, 17:4.2:4.1–4.2:4.29, October 2012.
 - [24] Dan Anthony Feliciano Alcantara. *Efficient hash tables on the gpu*. University of California at Davis, 2011.
 - [25] Alexandr Andoni and Piotr Indyk. Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions. *Commun. ACM*, 51(1):117–122, January 2008.
 - [26] Alina Beygelzimer, Sham Kakade, and John Langford. Cover trees for nearest neighbor. In *Proceedings of the 23rd international conference on Machine learning*, pages 97–104. ACM, 2006.
 - [27] Albert Bifet, Geoff Holmes, Richard Kirkby, and Bernhard Pfahringer. Moa: Massive online analysis. *J. Mach. Learn. Res.*, 11:1601–1604, August 2010.

-
- [28] Andrei Z. Broder, Steven C. Glassman, Mark S. Manasse, and Geoffrey Zweig. Syntactic clustering of the web. In *Selected Papers from the Sixth International Conference on World Wide Web*, pages 1157–1166, Essex, UK, 1997. Elsevier Science Publishers Ltd.
 - [29] Andrei Z. Broder, Steven C. Glassman, Mark S. Manasse, and Geoffrey Zweig. Syntactic clustering of the web. *Comput. Netw. ISDN Syst.*, 29(8-13):1157–1166, September 1997.
 - [30] L. Cayton. Accelerating nearest neighbor search on manycore systems. In *Parallel Distributed Processing Symposium (IPDPS), 2012 IEEE 26th International*, pages 402–413, May 2012.
 - [31] Moses S. Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings of the Thiry-fourth Annual ACM Symposium on Theory of Computing, STOC '02*, pages 380–388, New York, NY, USA, 2002. ACM.
 - [32] Gordon V Cormack and Thomas R Lynam. Trec 2005 spam track overview. In *TREC*, 2005.
 - [33] T. Cover and P. Hart. Nearest neighbor pattern classification. *IEEE Trans. Inf. Theor.*, 13(1):21–27, September 2006.
 - [34] Sanjoy Dasgupta and Yoav Freund. Random projection trees and low dimensional manifolds. In *Proceedings of the Fortieth Annual ACM Symposium on Theory of Computing, STOC '08*, pages 537–546, New York, NY, USA, 2008. ACM.
 - [35] Sanjoy Dasgupta and Anupam Gupta. An elementary proof of a theorem of johnson and lindenstrauss. *Random Struct. Algorithms*, 22(1):60–65, January 2003.
 - [36] Mayur Datar, Nicole Immorlica, Piotr Indyk, and Vahab S. Mirrokni. Locality-sensitive hashing scheme based on p-stable distributions. In *Proceedings of the Twentieth Annual Symposium on Computational Geometry, SCG '04*, pages 253–262, New York, NY, USA, 2004. ACM.
 - [37] John D Davis and Eric S Chung. Spmv: A memory-bound application on the gpu stuck between a rock and a hard place. *Microsoft Research Silicon Valley, Technical Report14 September*, 2012, 2012.
 - [38] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Andrew Senior, Paul Tucker, Ke Yang, Quoc V Le, et al.

- Large scale distributed deep networks. In *Advances in Neural Information Processing Systems*, pages 1223–1231, 2012.
- [39] Peter Frankl and Hiroshi Maehara. The johnson-lindenstrauss lemma and the sphericity of some graphs. *Journal of Combinatorial Theory, Series B*, 44(3):355–362, 1988.
- [40] Yoav Freund, Sanjoy Dasgupta, Mayank Kabra, and Nakul Verma. Learning the structure of manifolds using random projections. In *Advances in Neural Information Processing Systems*, volume 20, 2007.
- [41] Jerome H. Friedman, Jon Louis Bentley, and Raphael Ari Finkel. An algorithm for finding best matches in logarithmic expected time. *ACM Trans. Math. Softw.*, 3(3):209–226, September 1977.
- [42] Rainer Gemulla, Erik Nijkamp, Peter J. Haas, and Yannīs Sismanis. Large-scale matrix factorization with distributed stochastic gradient descent. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’11, pages 69–77, New York, NY, USA, 2011. ACM.
- [43] Fabian Gieseke, Justin Heiner mann, Cosmin Oancea, and Christian Igel. Buffer kd trees: Processing massive nearest neighbor queries on gpus. In *Proceedings of The 31st International Conference on Machine Learning*, page 172180, 2014.
- [44] Suyog Gupta, Ankur Agrawal, Kailash Gopalakrishnan, and Pritish Narayanan. Deep learning with limited numerical precision. *CoRR*, abs/1502.02551, 2015.
- [45] Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing*, STOC ’98, pages 604–613, New York, NY, USA, 1998. ACM.
- [46] William Johnson and Joram Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. In *Conference in modern analysis and probability (New Haven, Conn., 1982)*, volume 26 of *Contemporary Mathematics*, pages 189–206. American Mathematical Society, 1984.
- [47] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

- [48] Shengren Li, Lance Simons, Jagadeesh Bhaskar Pakaravoor, Fatemeh Abbasinejad, John D Owens, and Nina Amenta. kann on the gpu with shifted sorting. In *Proceedings of the Fourth ACM SIGGRAPH/Eurographics conference on High-Performance Graphics*, pages 39–47. Eurographics Association, 2012.
- [49] Michael D Lieberman, Jagan Sankaranarayanan, and Hanan Samet. A fast similarity join algorithm using graphics processing units. In *Data Engineering, 2008. ICDE 2008. IEEE 24th International Conference on*, pages 1111–1120. IEEE, 2008.
- [50] Dhruv Mahajan, S. Sathiya Keerthi, S. Sundararajan, and Léon Bottou. A parallel SGD method with strong convergence. *CoRR*, abs/1311.0636, 2013.
- [51] Luc Maranget, Susmit Sarkar, and Peter Sewell. A tutorial introduction to the arm and power relaxed memory models, 2012.
- [52] Jatin Chhugani Anthony D. Nguyen Victor W. Lee Daehyun Kim Pradeep Dubey Nadathur Satish, Changkyu Kim. Fast sort on cpus,gpus and intel mic architectures. <http://www.intel.com/content/dam/www/public/us/en/documents/technology-briefs/intel-labs-radix-sort-mic-report.pdf>.
- [53] John Nickolls, Ian Buck, Michael Garland, and Kevin Skadron. Scalable parallel programming with cuda. *Queue*, 6(2):40–53, March 2008.
- [54] Feng Niu, Benjamin Recht, Christopher R, and Stephen J. Wright. Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In *In NIPS*, 2011.
- [55] C. Nugteren. Improving the programmability of gpu architectures, 2014.
- [56] Stephen Malvern Omohundro. *Five balltree construction algorithms*. International Computer Science Institute Berkeley, 1989.
- [57] Thomas Paine, Hailin Jin, Jianchao Yang, Zhe Lin, and Thomas Huang. Gpu asynchronous stochastic gradient descent to speed up neural network training. *arXiv preprint arXiv:1312.6186*, 2013.
- [58] Jia Pan and Dinesh Manocha. Bi-level locality sensitive hashing for k-nearest neighbor computation. In *Proceedings of the 2012 IEEE 28th International Conference on Data Engineering, ICDE '12*, pages 378–389, Washington, DC, USA, 2012. IEEE Computer Society.

-
- [59] P.C. Pratt-Szeliga, J.W. Fawcett, and R.D. Welch. Rootbeer: Seamlessly using gpus from java. In *High Performance Computing and Communication 2012 IEEE 9th International Conference on Embedded Software and Systems (HPCC-ICESS), 2012 IEEE 14th International Conference on*, pages 375–380, June 2012.
- [60] Hans Sagan. Space-filling curves, 1994.
- [61] D. Sculley, Daniel Golovin, and Michael Young. Big learning with little ram. <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.298.2337&rep=rep1&type=pdf>.
- [62] Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and application to data-parallel distributed training of speech dnns. In *Interspeech 2014*, September 2014.
- [63] Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. On parallelizability of stochastic gradient descent for speech dnns. In *Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference on*, pages 235–239. IEEE, 2014.
- [64] Maxim Shevtsov, Alexei Soupikov, and Er Kapustin. Highly parallel fast kd-tree construction for interactive ray tracing of dynamic scenes.
- [65] Sami Sieranoja. High dimensional knn-graph construction using space filling curves. <http://urn.fi/urn:nbn:fi:uef-20150325>, 2015.
- [66] N. Sismanis, N. Pitsianis, and Xiaobai Sun. Parallel search of k-nearest neighbors with synchronous operations. In *High Performance Extreme Computing (HPEC), 2012 IEEE Conference on*, pages 1–6, Sept 2012.
- [67] T. Aila S.Laine, T. Karras. Megakernels considered harmful: Wavefront path tracing on gpus. https://research.nvidia.com/sites/default/files/publications/laine2013hpg_paper.pdf.
- [68] Hans Henrik Brandenborg Sørensen. High-performance matrix-vector multiplication on the gpu. In *Proceedings of the 2011 International Conference on Parallel Processing, Euro-Par’11*, pages 377–386, Berlin, Heidelberg, 2012. Springer-Verlag.
- [69] John E Stone, David Gohara, and Guochun Shi. Opencl: A parallel programming standard for heterogeneous computing systems. *Computing in science & engineering*, 12(1-3):66–73, 2010.

-
- [70] Nikko Strom. Scalable distributed dnn training using commodity gpu cloud computing. In *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.
- [71] Alan Tatourian. Nvidia gpu architecture and cuda programming environment. <http://tatourian.com/2013/09/03/nvidia-gpu-architecture-cuda-programming-environment/>.
- [72] Philippe Tillet, Karl Rupp, Siegfried Selberherr, and Chin-Teng Lin. Towards performance-portable, scalable, and convenient linear algebra. In *Presented as part of the 5th USENIX Workshop on Hot Topics in Parallelism*, Berkeley, CA, 2013. USENIX.
- [73] Roger Weber, Hans-Jörg Schek, and Stephen Blott. A quantitative analysis and performance study for similarity-search methods in high-dimensional spaces. In *Proceedings of the 24rd International Conference on Very Large Data Bases, VLDB '98*, pages 194–205, San Francisco, CA, USA, 1998. Morgan Kaufmann Publishers Inc.
- [74] Kun Zhou, Qiming Hou, Rui Wang, and Baining Guo. Real-time kd-tree construction on graphics hardware. *ACM Trans. Graph.*, 27(5):126:1–126:11, December 2008.
- [75] Martin Zinkevich, John Langford, and Alex J. Smola. Slow learners are fast. In Y. Bengio, D. Schuurmans, J.D. Lafferty, C.K.I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 2331–2339. Curran Associates, Inc., 2009.
- [76] Martin Zinkevich, Markus Weimer, Lihong Li, and Alex J Smola. Parallelized stochastic gradient descent. In *Advances in neural information processing systems*, pages 2595–2603, 2010.