

**Bachelors Thesis**

**Data Movement in Heterogeneous  
Memories with Intel Data Streaming  
Accelerator**

Anatol Constantin Fürst

16th January 2024

Technische Universität Dresden  
Faculty of Computer Science  
Institute of Systems Architecture  
Chair of Operating Systems

Academic Supervisors:

Prof. Dr.-Ing. Horst Schirmeier

Prof. Dr.-Ing. habil. Dirk Habich

M.Sc. André Berthold





## **Aufgabenstellung für die Anfertigung einer Bachelor-Arbeit**

Studiengang: Bachelor  
Studienrichtung: Informatik (2009)  
Name: **Constantin Fürst**  
Matrikelnummer: 4929314  
Titel: **Data Movement in Heterogeneous Memories with Intel Data Streaming Accelerator**

Developments in main memory technologies like Non-Volatile RAM (NVRAM), High Bandwidth Memory (HBM), NUMA, or Remote Memory, lead to heterogeneous memory systems that, instead of providing one monolithic main memory, deploy multiple memory devices with different non-functional memory properties. To reach optimal performance on such systems, it becomes increasingly important to move data, ahead of time, to the memory device with non-functional properties tailored for the intended workload, making data movement operations increasingly important for data intensive applications. Unfortunately, while copying, the CPU is mostly busy with waiting for the main memory, and cannot work on other computations. To tackle this problem Intel implements the Intel Data Streaming Accelerator (Intel DSA), an engine to explicitly offload data movement operations from the CPU, in their newly released Intel Xeon CPU Max processors.

The goal of this bachelor thesis is to analyze and characterize the architecture of the Intel DSA and the vendor-provided APIs. The student should benchmark the performance of Intel DSA and compare it to the CPU's performance, concentrating on data transfers between DDR5-DRAM and HBM and between different NUMA nodes. Additionally, the student should find out in what way and to what extent parallel processes copying data interfere with each other. Analyzing the performance information, the thesis should outline a gainful utilization of the Intel DSA and demonstrate its potential by extending the Query-driven Prefetching concept, which aims to speed up database query execution in heterogeneous memory systems.

Gutachter: Prof. Dr.-Ing. Dirk Habich  
Betreuer: André Berthold, M.Sc.  
Ausgehändigt am: 4. Dezember 2023  
Einzureichen am: 19. Februar 2024

Prof. Dr.-Ing. Horst Schirmeier  
Betreuender Hochschullehrer



## **Statement of Authorship**

I hereby declare that I am the sole author of this master thesis and that I have not used any sources other than those listed in the bibliography and identified as references. I further declare that I have not submitted this thesis at any other institution in order to obtain a degree.

Dresden, 16th January 2024

Anatol Constantin Fürst



## Abstract

...abstract ...

write the  
abstract





# Contents

<b>List of Figures</b>	<b>XI</b>
<b>List of Tables</b>	<b>XIII</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Technical Background</b>	<b>3</b>
2.1 High Bandwidth Memory . . . . .	3
2.2 Query Driven Prefetching . . . . .	3
2.3 Intel Data Streaming Accelerator . . . . .	3
2.4 System Setup and Configuration . . . . .	7
<b>3 Performance Microbenchmarks</b>	<b>9</b>
3.1 Benchmarking Methodology . . . . .	9
3.2 Submission Method . . . . .	9
3.3 Multithreaded Submission . . . . .	9
3.4 Multiple Engines in a Group . . . . .	10
3.5 Data Movement from DDR to HBM . . . . .	10
3.6 Analysis . . . . .	10
<b>4 Design</b>	<b>11</b>
4.1 Detailed Task Description . . . . .	11
4.2 Cache Design . . . . .	11
<b>5 Implementation</b>	<b>15</b>
5.1 Locking and Usage of Atomics . . . . .	15
5.2 Accelerator Usage . . . . .	16
<b>6 Evaluation</b>	<b>17</b>
<b>7 Conclusion And Outlook</b>	<b>19</b>
7.1 Conclusions . . . . .	19
7.2 Future Work . . . . .	19
<b>Glossary</b>	<b>21</b>
<b>Bibliography</b>	<b>23</b>



# List of Figures

2.1	Internal Archtiecture Block Diagramm Taken from Figure 1a of [3] . . .	4
2.2	Software View Block Diagramm Taken from Figure 1a of [3] . . . . .	6



# List of Tables



# 1 Introduction

---

write this  
chapter





## 2 Technical Background

---

### 2.1 High Bandwidth Memory

write introductory paragraph

---

### 2.2 Query Driven Prefetching

write this section

---

### 2.3 Intel Data Streaming Accelerator

write this section

Intel DSA is a high-performance data copy and transformation accelerator that will be integrated in future Intel® processors, targeted for optimizing streaming data movement and transformation operations common with applications for high-performance storage, networking, persistent memory, and various data processing applications. [1, p. 15]

Introduced with the 4th generation of Intel Xeon Scalable Processors [2], the DSA promises to alleviate the CPU from ‘common storage functions and operations such as data integrity checks and deduplication’ [2]. This chapter will give an overview of the architecture, software and the interaction of these two components. The reader will be familiarized with the setup and equipped with the knowledge to configure the system for a specific use case.

To be able to optimally utilize the Hardware, knowledge of its workings is required to make educated decisions. Therefore, this section describes both the workings of the DSA engine itself and the view that is presented through software interfaces. All statements are based on Chapter 3 of the Architecture Specification by Intel [1].

consider adding projected use cases as in the architecture specification here

### 2.3.1 Hardware Architecture

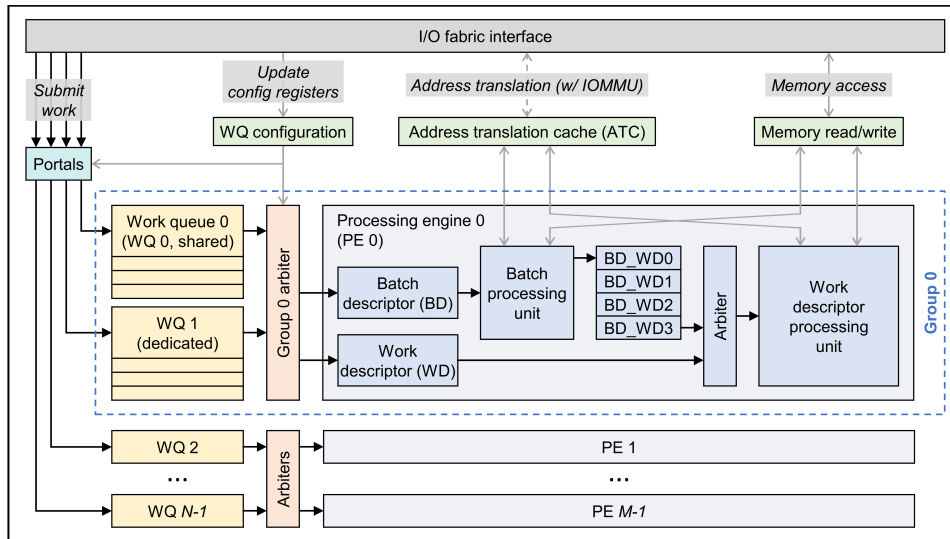


Figure 2.1:  
Internal Architecture Block Diagram  
Taken from Figure 1a of [3]

The accelerator is directly integrated into the Processor and attaches via the I/O fabric interface over which all communication is conducted. Over this interface, it is accessible as a PCIe device. Configuration therefore is done through memory-mapped registers set in the devices Base Address Register (BAR). Through these, the devices layout is defined and memory pages for work submission are set. In a system with multiple processing nodes, there may also be one DSA per node.

To satisfy different use cases, as already mentioned, the layout of the DSA may be software-defined. The structure is made up of three components, namely Work Queue (WQ)s, Engines and Groups. WQs provide the means to submit tasks to the device and will be described in more detail shortly. An Engine is the processing-block that connects to memory and performs the described task. Using Groups, Engines and WQs are tied together. This means, that tasks from one WQ may be processed from multiple Engines and that vice-versa, depending on the configuration. This flexibility is achieved through the Group Arbiter which connects the two components and acts according to the setup.

A WQ is accessible through so-called portals, which are mapped memory regions. Submission of work is done by writing a descriptor to one of these portals. A descriptor is 64 Byte in size and may contain one specific task (task descriptor) or the location of a task array in memory (batch descriptor). Through these portals, the submitted descriptor reaches a queue of which there are two types with different submission methods and use cases. The Shared Work Queue (SWQ) is intended to provide synchronized access to multiple processes and each group may only have one attached. A PCIe Deferrable Memory Write Request (DMR), which guarantees implicit synchronization, is generated

via x86 Instruction ENQCMD and communicates with the device before writing. This results in higher submission cost, compared to the Dedicated Work Queue (DWQ) to which a descriptor is submitted via x86 Instruction MOVDIR64B. The DWQ is therefore more performant but may require access control mechanisms and may only be accessed by one process at a time.

To handle the different descriptors, each Engine has two internal execution paths. One for a task and the other for a batch descriptor. Processing a task descriptor is straightforward, as all information required to complete the operation are contained within. For a batch, the DSA first reads the batch descriptor, then fetches all task descriptors for the batch from memory and processes them. An Engine can also trigger a page fault when trying to access an unloaded page and wait on its completion, if configured to do so. Otherwise, an error will be generated in this scenario.

Ordering of operations is only guaranteed for a configuration with one WQ and one Engine in a Group when submitting exclusively batch or task descriptors but no mixture. Even then, only write-ordering is guaranteed, meaning that ‘reads by a subsequent descriptor can pass writes from a previous descriptor’ [1, p. 30]. A different issue arises, should an operation fail: the DSA will continue to process the following descriptors. Care must therefore be taken with read-after-write scenarios, either by waiting for a successful completion before submitting the dependant, inserting a drain descriptor for tasks or setting the fence flag for a batch. The latter two methods tell the processing engine that all writes must be committed and, in case of the fence in a batch, abort on previous error.

An important aspect of modern computer systems is the separation of address spaces through virtual memory. The DSA must therefore handle address translation, as a process submitting a task will not know the physical location in memory which causes the descriptor to contain virtual values. For this, the Engine communicates with the Input/Output Memory Management Unit (IOMMU) and Address Translation Cache (ATC) to perform this operation. For this, knowledge about the submitting processes is required, and therefore each task descriptor has a field for the Process Address Space ID (PASID) which is filled by the ENQCMD instruction for a SWQ or set statically after a process is attached to a DWQ.

The completion of a descriptor may be signaled through a completion record and interrupt, if configured so. For this, the DSA ‘provides two types of interrupt message storage: (1) an MSI-X table, enumerated through the MSI-X capability; and (2) a device-specific Interrupt Message Storage (IMS) table’ [1, p. 27].

### 2.3.2 Software View

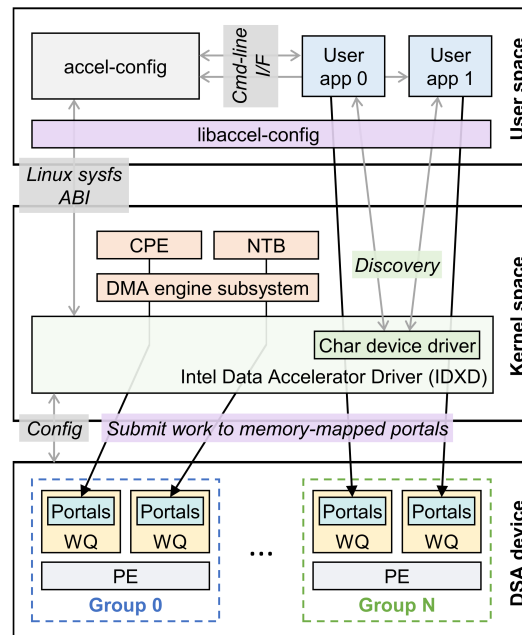


Figure 2.2:  
Software View Block Diagram  
Taken from Figure 1a of [3]

Due to efforts by intel programmers, since Linux Kernel 5.10 [4, Installation Instructions], there exists a driver for the DSA [5] which has no counterpart in the Windows OS-Family [4, Installation Instructions], meaning code developed without an alternative path will not work there. To interface with the driver and perform configuration operations, intel's `libaccel-conf` [6] user space toolset may be used which provides a command-line interface and can read configuration files to set up the device as described previously. After successful configuration, each WQ is exposed as a character device by `mmap` of the associated portal [3, p. 3].

Given the file permissions, it would now be possible for a process to submit work to the DSA via either `MOVDIR64B` or `ENQCMD` instructions, providing the descriptors by manually configuring them. This, however, is quite cumbersome, which is why Intel's Data Mover Library [4] exists. With some limitations (like lacking support for DWQs) this library presents a high-level interface that takes care of creation and submission of descriptors, some error handling and reporting. Thanks to the high-level-view the code may choose a different execution path at runtime which allows the memory operations to either be executed in hardware (on a DSA) or in software (using equivalent instructions provided by the library) which makes code based upon it automatically compatible with systems that do not provide hardware or software support.

finish this section

- drain descriptor / drain command signals completion of preceding descriptors for fencing in non-batch submissions, in batches the “fence flag” can be used to ensure ordering, failures before a fence will lead to the following descriptors being aborted [1, p. 30], `sfence` or `mfence` should be executed before pushing drain descriptor [1, p. 32]
- cache control flag in descriptor controls whether writes are directed to cache or to memory [1, p. 31] effects on copy from DRAM > HBM unknown

### 2.3.3 Programming Interface

write this section

- choice is intel data mover library
- two concepts, state-based for c-api and operation-based c++
- just explain the basics (no code) and refer to dml documentation

## 2.4 System Setup and Configuration

write this section

Give the reader the tools to replicate the setup. Also explain why the BIOS-configs are required.

Setup Requirements:

- VT-d enabled
- limit CPUPA to 46 Bits disabled
- IOMMU enabled
- kernel with iommu and idxd driver support
- kernel option "intel\_iommu=on,sm\_on"
- numa nodes for hbm access in bios



# 3 Performance Microbenchmarks

write introductory paragraph

## 3.1 Benchmarking Methodology

write this section

## 3.2 Submission Method

write this section

- submit cost analysis: best method and for a subset the point at which submit cost < time savings
- display the full opt-submitmethod graph
- maybe remeasure with higher amount of small copies? results look somewhat weird for 1k and 4k
- display the stacked bar of submit and complete time for single@1k, single@4k, single@1mib for HW-path and SW-path
- display the stacked bar of submit and complete time for batch50@1k, batch50@4k, batch50@1mib for HW-path and SW-path
- show batch because we care about the minimum task set size for a single producer (multi submit would be used for different task sets)
- conclude at which point using the DSA makes sense

## 3.3 Multithreaded Submission

write this section

- effect of mt-submit, low because SWQ implicitly synchronized, bandwidth is shared
- show results for all available core counts
- only display the 1engine tests
- show combined total throughput
- conclude that due to the implicit synchronization the sync-cost also affects 1t and therefore it makes no difference, bandwidth is shared, no guarantees on fairness

write this section

### 3.4 Multiple Engines in a Group

write this section

- assumed from arch spec that multiple engines lead to greater Performance
- reason is that page faults and access latency will be overlapped with preparing the next operation
- in the given scenario we observe the opposite, slight performance decrease
- show multsubmit 50 for both 1e and 4e
- maybe remeasure with each submission accessing different memory region?
- conclusion?

### 3.5 Data Movement from DDR to HBM

write this section

- present two copy methods: smart and brute force
- show graph for ddr->hbm intranode, ddr->hbm intrasocket, ddr->hbm intersocket
- conclude which option makes more sense (smart)
- because 4x or 2x utilization for only 1.5x or 1.25x speedup respectively
- maybe benchmark smart-copy intersocket in parallel with two smart-copies intrasocket VS. the same task with brute force

### 3.6 Analysis

write this section

- summarize the conclusions and define the point at which dsa makes sense
- minimum transfer size for batch/nonbatch operation
- effect of msubmit -> no fairness guarantees
- usage of multiple engines -> no effect
- smart copy method as the middle-ground between peak throughput and utilization
- lower utilization of dsa is good when it will be shared between threads/processes



# 4 Design

write introductory paragraph

## 4.1 Detailed Task Description

write this section

- give slightly more detailed task Description
- perspective of "what problems have to be solved"
- not "what is query driven prefetching"

## 4.2 Cache Design

The task of prefetching is somewhat aligned with that of a cache. As a cache is more generic and allows use beyond Query Driven Prefetching, the choice was made to solve the prefetching offload by implementing an offloading **Cache**. When referring to the provided implementation, **Cache** will be used from now on. The interface with **Cache** must provide three basic functions: requesting a memory block to be cached, accessing a cached memory block and synchronizing cache with the source memory. The latter operation comes in to play when the data that is cached may also be modified, requiring the entry to be updated with the source or the other way around. Due to the many possible setups and use cases, the user should also be responsible for choosing cache placement and the copy method. As re-caching is resource intensive, data should remain in the cache for as long as possible while being removed when system memory pressure due to restrictive memory size drives the **Cache** to flush unused entries.

### 4.2.1 Interface

To allow rapid integration and ease developer workload, a simple interface was chosen. As this work primarily focuses on caching static data, the choice was made only to provide cache invalidation and not synchronization. Given a memory address, **Cache::Invalidate** will remove all entries for it. The other two operations are provided in one single function, which we shall call **Cache::Access** henceforth, receiving a data pointer and size it takes care of either submitting a caching operation if the pointer received is not yet cached or returning the cache entry if it is. The cache placement and assignment of the task to accelerators are controlled by the user. In addition to the two basic operations outlined before, the user also is given the option to flush the cache using **Cache::Flush** of unused elements manually or to clear it completely with **Cache::Clear**.

As caching is performed asynchronously, the user may wish to wait on the operation. This would be beneficial if there are other threads making progress in parallel while the current thread waits on its data becoming available in the faster cache, speeding up local computation. To achieve this, the `Cache::Access` will return an instance of an object which from hereinafter will be referred to as `CacheData`. Through `CacheData::GetDataLocation` a pointer to the cached data will be retrieved, while also providing `CacheData::WaitOnCompletion` which must only return when the caching operation has completed and during which the current thread is put to sleep, allowing other threads to progress.

### 4.2.2 Cache Entry Reuse

When multiple consumers wish to access the same memory block through the `Cache`, we could either provide each with their own entry, or share one entry for all consumers. The first option may cause high load on the accelerator due to multiple copy operations being submitted and also increases the memory footprint of the system. The latter option requires synchronization and more complex design. As the cache size is restrictive, the latter was chosen. The already existing `CacheData` will be extended in scope to handle this by allowing copies of it to be created which must synchronize with each other for `CacheData::WaitOnCompletion` and `CacheData::GetDataLocation`.

### 4.2.3 Cache Entry Lifetime

By allowing multiple references to the same entry, memory management becomes a concern. Freeing the allocated block must only take place when all copies of a `CacheData` instance are destroyed, therefore tying cache entry lifetime to the lifetime of the longest living copy of the original instance. This makes access to the entry legal during the lifetime of any `CacheData` instance, while also guaranteeing that `Cache::Clear` will not have any unforeseen side effects, as deallocation only takes place when the last consumer has `CacheData` go out of scope or manually deletes it.

### 4.2.4 Usage Restrictions

As cache invalidation applies mainly to non-static data which this work does not focus on, two restrictions are placed on the invalidation operation. This permits drastically simpler cache design, as a fully coherent cache would require developing a thread safe coherence scheme which is outside our scope.

Firstly, overlapping areas in the cache will cause undefined behaviour during invalidation of any one of them. Only the entries with the equivalent source data pointer will be invalidated, while other entries with differing source pointers which, due to their size, still cover the now invalidated region, will not be invalidated and therefore the cache may and may continue to contain invalid elements at this point.

Secondly, invalidation is to be performed manually, requiring the programmer to remember which points of data are at any given point in time cached and invalidating them upon modification. No ordering guarantees will be given for this situation, possibly

leading to threads still having a pointer to now-outdated entries and continuing their progress with this.

Due to its reliance on `libnuma` for numa awareness, `Cache` will only work on systems where this library is present, excluding, most notably, Windows from the compatibility list.

#### 4.2.5 Thread Safety Guarantees

After initialization, all available operations for `Cache` and `CacheData` are fully threadsafe but may use locks internally to achieve this. In 5 we will go into more detail on how these guarantees are provided and how to optimize the cache for specific use cases that may warrant less restrictive locking.

#### 4.2.6 Accelerator Usage

Compared with the challenges of ensuring correct entry lifetime and thread safety, the application of DSA for the task of duplicating data is simple, thanks partly to Intel Data Mover Library (Intel DML) [4]. Upon a call to `Cache::Access` and determining that the given memory pointer is not present in cache, work will be submitted to the Accelerator. Before, however, the desired location must be determined which the user-defined cache placement policy function handles. With the desired placement obtained, the copy policy function then determines, which nodes should take part in the copy operation which is equivalent to selecting the Accelerators following 2.3.1. This causes the work to be split upon the available accelerators to which the work descriptors are submitted at this time. The handlers that Intel DML [4] provides will then be moved to the `CacheData` instance to permit the callee to wait upon caching completion. As the choice of cache placement and copy policy is user-defined, one possibility will be discussed in 5.



# 5 Implementation

write introductory paragraph

## 5.1 Locking and Usage of Atomics

As the usage of locking and atomics may have a significant impact on performance, their application will be discussed in detail within this section.

extend introductory paragraph

### 5.1.1 Cache State Lock

To keep track of the current cache state, a map is used internally which associates a memory address to a `CacheData` instance. In 4.2.2 we decided to reuse one cache entry for multiple consumers, requiring thread safety when accessing and extending the cache state in `Cache::Access`, `Cache::Flush` and `Cache::Clear`. The latter two both require a unique lock, preventing other calls to `Cache` from making progress while the operation is being processed. For `Cache::Access` the use of locking depends upon the caches state. At first only a shared lock is acquired for checking whether the given address already resides in cache, allowing other `Cache::Access`-operations to also perform this check. If no entry for the region is present, a unique lock is required as well when adding the newly created entry to cache, which however is a rather short operation.

In scenarios where the `Cache` is frequently tasked with flushing and re-caching by multiple threads accessing large amounts of data, leading to high memory pressure, lock contention around this lock will negatively impact performance by delaying cache access. Due to passive waiting, this impact might be less noticeable when other threads on the system are able to make progress during the wait.

### 5.1.2 CacheData Reference Counting

write this section

### 5.1.3 CacheData WaitOnCompletion

write this section

### 5.1.4 Performance Guideline

The performance impact of lock contention and atomic synchronization is not to be taken lightly, as `Cache` may be used in performance critical systems. Reducing their impact is therefore desirable which can be achieved in multiple ways. The easiest is to have one instance of `Cache` per NUMA-Node (Node) which reduces both lock contention by just serving less threads and atomic synchronization as the atomics are shared between physically close cpu cores. This requires no code modification but does not inherently reduce the amount of synchronization taking place. To achieve this reduction, restrictions

find a reference for this and use above too

find a reference for this and use above too

find a reference that shows that physical distance affects sync cost

must be placed upon the thread safety or access guarantees, which is not sensible for this generic implementation.

## 5.2 Accelerator Usage

After 4.2.6 the implementation of `Cache` provided leaves it up to the user to choose a caching and copy method policy which is accomplished through submitting function pointers at initialization of the `Cache`. In 2.4 we configured our system to have separate Nodes for accessing High Bandwidth Memory (HBM) which are assigned a Node-ID by adding eight to the Nodes ID of the Node that physically contains the HBM. Therefore, given Node 3 accesses some datum, the most efficient placement for the copy would be on Node  $3 + 8 == 11$ . As the `Cache` is intended for multithreaded usage, conserving accelerator resources is important, so that concurrent cache requests complete quickly. To get high per-copy performance while maintaining low usage, the smart-copy method is selected as described in 3.5 for larger copies, while small copies under 64 MiB will be handled exclusively by the current node. This size is quite high but due to the overhead of assigning the current thread to the selected nodes, using only the current one is more efficient. This assignment is required due to Intel DML not being Non Uniform Memory Architecture (NUMA) aware and therefore assigning submissions only to the DSA engine present on the node that the calling thread is assigned to [4].

Actually test which size makes sense, due to numa-setaffinity and scheduling overhead this will probably be much higher

# 6 Evaluation

...evaluation ...

write this  
chapter





# 7 Conclusion And Outlook

---

## 7.1 Conclusions

write introductory paragraph

---

## 7.2 Future Work

write this section

- evaluate impact of lock contention and atomics on performance
- provide optimized use case specific versions with less locking
- extend the cache implementation use cases where data is not static

write this section



# Glossary

## A

### ATC

... desc ...

## B

### BAR

... desc ...

## D

### DMR

... desc ...

### DSA

... desc ...

### DWQ

... desc ...

## E

### Engine

... desc ...

### ENQCMD

... desc ...

## G

### Group

... desc ...

## H

### HBM

... desc ...

**I****Intel DML**

... desc ...

**IOMMU**

... desc ...

**M****MOVDIR64B**

... desc ...

**N****Node**

... desc ...

**NUMA**

... desc ...

**P****PASID**

... desc ...

**S****SWQ**

... desc ...

**W****WQ**

... desc ...

# Bibliography

- [1] Intel. ‘Intel® Data Streaming Accelerator Architecture Specification’. (16th Sep. 2022), [Online]. Available: <https://www.intel.com/content/www/us/en/content-details/671116/intel-data-streaming-accelerator-architecture-specification.html> (visited on 15th Nov. 2023).
- [2] Intel. ‘New Intel® Xeon® Platform Includes Built-In Accelerators for Encryption, Compression, and Data Movement’. (Dec. 2022), [Online]. Available: <https://www.intel.com/content/dam/www/central-libraries/us/en/documents/2022-12/storage-engines-4th-gen-xeon-brief.pdf> (visited on 15th Nov. 2023).
- [3] R. K. et al. ‘A Quantitative Analysis and Guideline of Data Streaming Accelerator in Intel® 4th Gen Xeon® Scalable Processors’. (May 2023), [Online]. Available: <https://arxiv.org/pdf/2305.02480.pdf> (visited on 7th Jan. 2024).
- [4] Intel, *Intel Data Mover Library Documentation*, <https://intel.github.io/DML/index.html>. (visited on 7th Jan. 2024).
- [5] Intel, *Intel IDX Driver for Linux Kernel*, <https://github.com/intel/idx-driver>. (visited on 7th Jan. 2024).
- [6] Intel, *Intel IDX User Space Application*, <https://github.com/intel/idx-config>. (visited on 7th Jan. 2024).