

Diplomarbeit

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

Anatol Constantin Fürst

7th January 2024

Technische Universität Dresden
Fakultät Informatik
Institut für Systemarchitektur
Professur Betriebssysteme

Betreuender Hochschullehrer: Prof. Dr.-Ing. Horst Schirmeier
Betreuender Mitarbeiter: 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: 24. November 2023
Einzureichen am: 9. Februar 2024

Prof. Dr.-Ing. Horst Schirmeier
Betreuender Hochschullehrer

Selbständigkeitserklärung

Hiermit erkläre ich, dass ich diese Arbeit selbstständig erstellt und keine anderen als die angegebenen Hilfsmittel benutzt habe.

Dresden, den 7. Januar 2024

Anatol Constantin Fürst

Abstract

...abstract ...

write ab-
stract

Contents

List of Figures	XIII
List of Tables	XV
1 Introduction	1
1.1 A Section	1
1.2 Another Section	1
1.3 Yet Another Section	1
1.4 Test commands	1
1.5 Test Special Chars	1
2 Technical Background on Intel DSA	3
2.1 Architecture	3
2.2 HW/SW Setup	7
2.3 Microbenchmarks	7
2.4 Evaluation	7
3 Design	9
3.1 Introduction VAMPIR	9
3.2 Analysis of Applicability of DSA	9
4 Implementation	11
5 Evaluation	13
6 Future Work	15
7 Conclusion And Outlook	17
Bibliography	21

Todo list

write abstract	VII
adopt title page	1
adopt disclaimer	1
write introduction	1
add content	1
Figure: Come up with a mindblowing figure.	2
consider adding projected use cases as in the architecture specification here . . .	3
provide microbenchmarks with multiple configurations and for many use cases .	7
evaluate the benchmarks and conclude with projected use cases - may use the cases from dsaspec/guide	7
write implementation	11
write evaluation	13
write future work	15
write conclusion	17

List of Figures

1.1	Short description	2
1.2	A mindblowing figure	2
2.1	Intel Data Streaming Accelerator (DSA) Internal Block Diagramm . . .	4
2.2	DSA Internal Block Diagramm	6

List of Tables

1.1	Some interesting numbers	1
-----	------------------------------------	---

1 Introduction

adopt title page

adopt disclaimer

write introduction

1.1 A Section

Referencing other chapters: 2 3 4 5 6 7

Name	Y	Z
<i>Foo</i>	20,614	23 %
<i>Bar</i>	9,914	11 %
<i>Foo + Bar</i>	30,528	34 %
<i>total</i>	88,215	100 %

Table 1.1: Various very important looking numbers and sums.

More text referencing Table 1.1.

1.2 Another Section

Citing [Bel05] other documents [Bel05; Boi06] and Figure 1.1.

Something with umlauts and a year/month date: [BD04].

And some online resources: [Gre04], [Hub89]

1.3 Yet Another Section

add content

1.4 Test commands

DROPS L⁴LinuxNOVA QEMU memcpy A sentence about BASIC. And a correctly formatted one about ECC.

1.5 Test Special Chars

Before you start writing your thesis please make sure that your build setup compiles the following special chars correctly into the PDF! If for example ß is printed as 'SS' then you should fix this! There are a few hints in the repository in `preamble/packages.txt`.

ö ä ü Ö Ä Ü ß < >



Figure 1.1: A long description of this squirrel figure. Image taken from http://commons.wikimedia.org/wiki/File:Sciurus-vulgaris_hernandeangelis_stockholm_2008-06-04.jpg

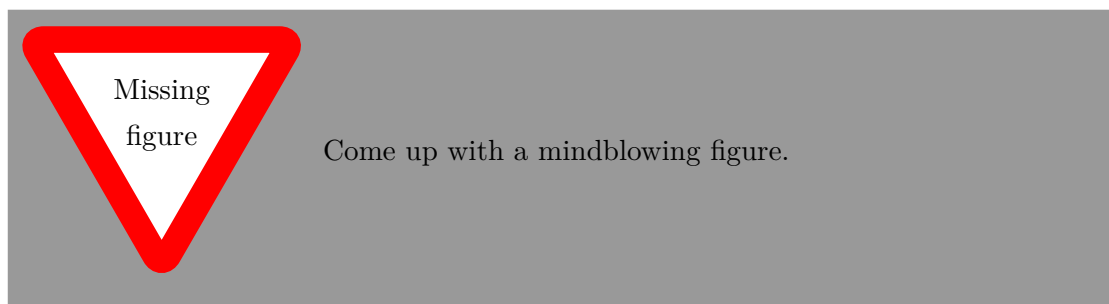


Figure 1.2: A mindblowing figure

2 Technical Background on Intel DSA

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. [Cor22a, p. 15]

Introduced with the 4th generation of Intel Xeon Scalable Processors [Cor22b], the DSA promises to alleviate the CPU from ‘common storage functions and operations such as data integrity checks and deduplication’ [Cor22b]. 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.

2.1 Architecture

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 [Cor22a].

2.1.1 Hardware Architecture

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

consider adding projected use cases as in the architecture specification here

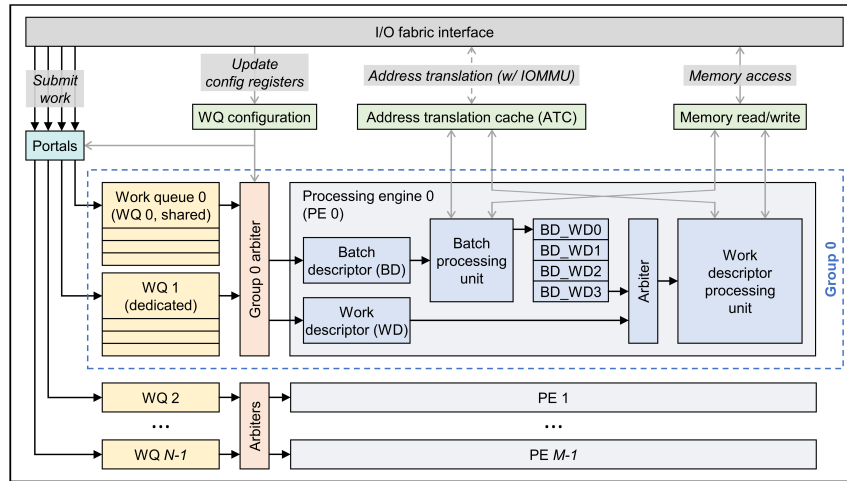


Figure 2.1: Taken from Figure 1a of [al23]

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, 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’ [Cor22a, 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’ [Cor22a, p. 27].

2.1.2 Software View

Due to efforts by intel programmers, since Linux Kernel 5.10 [Cora, Installation Instructinos], there exists a driver for the DSA [Corb] which has no counterpart in the Windows OS-Family [Cora, Installation Instructinos], meaning code developed without an alternative path will not work there. To interface with the driver and perform configuration operations, intels libaccl-conf [Corc] 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 [al23, 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 Intels Data Mover Library [Cora] 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.

- 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 [Cor22a, p. 30], `sfence` or `mfence` should be executed before pushing drain descriptor [Cor22a, p. 32]
- cache control flag in descriptor controls whether writes are directed to cache or to memory [Cor22a, p. 31] effects on copy from DRAM > HBM unknown

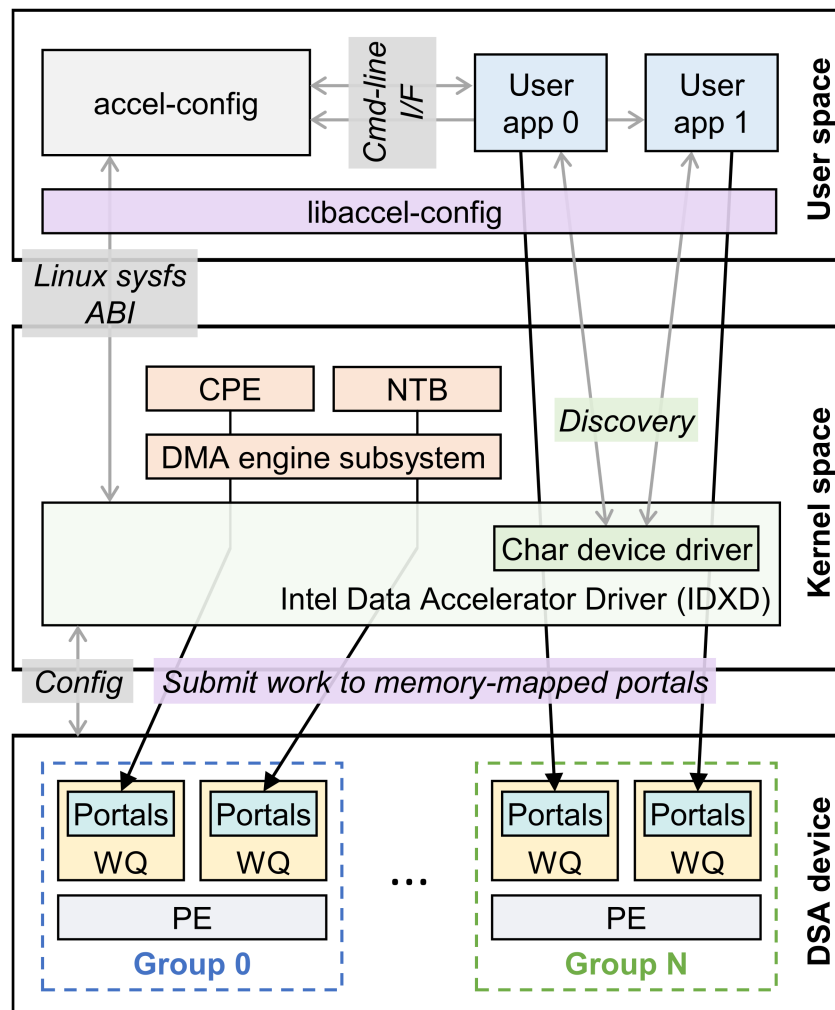


Figure 2.2: Taken from Figure 1b of [al23]

2.2 HW/SW Setup

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"

2.3 Microbenchmarks

2.4 Evaluation

provide microbenchmarks with multiple configurations and for many use cases

evaluate the benchmarks and conclude with projected use cases - may use the cases from dsaspec/guide

3 Design

3.1 Introduction VAMPIR

- Hardware Overview with CPU/RAM/HBM/NUMA-Nodes in Graph
- Overview of Software with query-pipeline

3.2 Analysis of Applicability of DSA

- Benchmark the amount of time spent on memory operations in VAMPIR
- Back-reference to the Microbenchmarks and conclusion on possible gains

4 Implementation

...implementation ...

write imple-
mentation

5 Evaluation

...evaluation ...

write evaluation

6 Future Work

...future work ...

write future
work

7 Conclusion And Outlook

...conclusion ...

write conclusion

Glossary

A

ATC

... desc ...

B

BAR

... desc ...

D

DSA

... desc ...

DWQ

... desc ...

E

Engine

... desc ...

ENQCMD

... desc ...

G

Group

... desc ...

I

IOMMU

... desc ...

M

MOVDIR64B

... desc ...

P

PASID

... desc ...

PCIe Deferrable Memory Write Request

... desc ...

S

SWQ

... desc ...

W

WQ

... desc ...

Bibliography

- [al23] Reese Kuper et al. *A Quantitative Analysis and Guideline of Data Streaming Accelerator in Intel® 4th Gen Xeon® Scalable Processors*. May 2023. URL: <https://arxiv.org/pdf/2305.02480.pdf> (visited on 7th Jan. 2024).
- [BD04] Michael Becher and Maximillian Dornseif. ‘Feuriges Hacken - Spaß mit Firewire’. In: *21C3: Proceedings of the 21st Chaos Communication Congress*. Dec. 2004.
- [Bel05] Fabrice Bellard. ‘QEMU, a fast and portable dynamic translator’. In: *Proceedings of the USENIX Annual Technical Conference, FREENIX Track*. 2005, pp. 41–46.
- [Boi06] Adam Boileau. ‘Hit by a Bus: Physical Access Attacks with Firewire’. In: *RUXCON*. 2006.
- [Cora] Intel Corporation. *Intel Data Mover Library Documentation*. URL: <https://intel.github.io/DML/index.html> (visited on 7th Jan. 2024).
- [Corb] Intel Corporation. *Intel IDXD Driver for Linux Kernel Repository*. URL: <https://github.com/intel/idxd-driver> (visited on 7th Jan. 2024).
- [Corc] Intel Corporation. *Intel Libaccel-Config Repository*. URL: <https://github.com/intel/idxd-config> (visited on 7th Jan. 2024).
- [Cor22a] Intel Corporation. *Intel® Data Streaming Accelerator Architecture Specification*. 16th Sept. 2022. URL: <https://www.intel.com/content/www/us/en/content-details/671116/intel-data-streaming-accelerator-architecture-specification.html> (visited on 15th Nov. 2023).
- [Cor22b] Intel Corporation. *New Intel® Xeon® Platform Includes Built-In Accelerators for Encryption, Compression, and Data Movement*. Dec. 2022. URL: <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).
- [Gre04] Tom Green. *1394 Kernel Debugging Tips and Tricks*. Slide presentation at the WinHEC 2004. 2004. URL: http://download.microsoft.com/download/1/8/f/18f8cee2-0b64-41f2-893d-a6f2295b40c8/DW04001_WINHEC2004.ppt (visited on 3rd June 2009).
- [Hub89] William S. Huber. ‘Operating system debugger’. 4819234 (Needham, MA). Apr. 1989. URL: <http://www.freepatentsonline.com/4819234.html>.