

A Mesh Architecture for Data Management of Matrix Computations

by

Adam BURDENIUK

B. Eng (Hons) Computer Systems,
The University of Adelaide, 2001

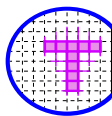
Thesis submitted for the degree of

Doctor of Philosophy

in

Electrical and Electronic Engineering
The University of Adelaide

June 2014



CHIPTec
*The Centre for High Performance Integrated
Technologies and Systems*

© 2014

Adam BURDENIUK

All Rights Reserved



**ADELAIDE
UNIVERSITY**
AUSTRALIA

Contents

| | |
|---|-----------|
| Chapter 1. Introduction | 1 |
| 1.1 Motivation | 2 |
| 1.2 The matrix processing problem | 2 |
| 1.3 Research aims and significance | 4 |
| 1.4 Introducing the Matrix Data Management Layer (MDML) . . | 5 |
| 1.5 Thesis structure | 6 |
| 1.6 Contributions and claims | 8 |
| 1.7 Publications | 10 |
| | |
| Chapter 2. A survey of parallel computer systems | 11 |
| 2.1 Multiprocessor system-on-chip architectures | 12 |
| 2.2 Specialised matrix co-processors | 13 |
| 2.2.1 Systolic processors | 14 |
| 2.2.2 The MatRISC family | 15 |
| 2.2.3 Renesas MX architecture | 19 |
| 2.2.4 The Trident architecture | 20 |
| 2.2.5 MatCore | 22 |
| 2.2.6 Discussion of matrix co-processors | 23 |
| 2.3 Communication in high-performance systems | 25 |
| 2.4 Discussion and summary | 26 |
| | |
| Chapter 3. The Matrix Data Management Layer (MDML) | 29 |
| 3.1 Background | 30 |
| 3.2 System environment | 35 |
| 3.3 Dense matrix algorithms | 38 |
| 3.4 MDML structural decomposition | 41 |

Contents

| | | |
|---|--|-----------|
| 3.4.1 | Node components | 43 |
| 3.4.2 | System-wide synchronisation | 49 |
| 3.4.3 | System modes | 49 |
| 3.5 | Summary | 50 |
| Chapter 4. Data movement within an MDML node | | 51 |
| 4.1 | Introduction | 52 |
| 4.2 | Related previous work | 53 |
| 4.3 | Common activities in matrix algorithms | 59 |
| 4.3.1 | Loop index counting | 62 |
| 4.3.2 | Code sequence selection | 64 |
| 4.3.3 | Data indexing | 66 |
| 4.3.4 | Range updating | 68 |
| 4.4 | Node architecture | 71 |
| 4.4.1 | Node pipeline | 72 |
| 4.4.2 | The matrix sequencing engine | 74 |
| 4.4.3 | Interfacing to local processing and memory resources | 83 |
| 4.4.4 | Node controller | 86 |
| 4.5 | Implementation results | 88 |
| 4.6 | Summary | 89 |
| Chapter 5. Data movement between nodes | | 91 |
| 5.1 | Introduction | 92 |
| 5.1.1 | Issues in communication system design | 92 |
| 5.1.2 | A review of synchronisation methods | 96 |
| 5.2 | Extending the MDML for multiprocessing | 99 |
| 5.3 | MDML communication architecture | 101 |
| 5.3.1 | The synchronisation manager | 103 |
| 5.3.2 | The data transport module | 109 |
| 5.3.3 | Modelling and verification | 113 |

| | | |
|--|---|------------|
| 5.3.4 | Discussion | 114 |
| 5.4 | The array controller | 115 |
| 5.4.1 | Array bootstrapping and parameter discovery | 116 |
| 5.4.2 | Edge detection | 117 |
| 5.4.3 | Array size and node location discovery | 118 |
| 5.4.4 | Network hop-size discovery | 119 |
| 5.4.5 | Modelling and verification | 119 |
| 5.5 | Summary | 121 |
| Chapter 6. Programming the MDML | | 123 |
| 6.1 | Introduction | 124 |
| 6.2 | Assumptions | 126 |
| 6.3 | Algorithms for small matrices | 127 |
| 6.3.1 | Matrix multiplication | 128 |
| 6.3.2 | LU decomposition | 134 |
| 6.3.3 | Forward substitution for a single operating point | 141 |
| 6.3.4 | Forward substitution for an aggregate of points | 146 |
| 6.3.5 | Discussion | 153 |
| 6.4 | Algorithms for large matrices | 153 |
| 6.4.1 | LU decomposition | 156 |
| 6.4.2 | Aggregate forward substitution | 167 |
| 6.5 | Impact of the MDML | 175 |
| 6.6 | Summary | 176 |
| Chapter 7. Conclusion and future work | | 177 |
| 7.1 | Conclusion | 178 |
| 7.2 | Future work | 180 |
| 7.2.1 | Development tools and operating system support | 180 |
| 7.2.2 | Generalisation of the design | 182 |
| 7.2.3 | Performance optimisations | 182 |
| 7.2.4 | Power optimisations | 183 |
| 7.2.5 | Alternate designs | 183 |

| | |
|--|------------|
| Appendix A. Format of MDML control words | 185 |
| A.1 VLIW fields for matrix acceleration | 186 |
| A.2 VLIW fields for multiprocessing | 188 |
| Appendix B. Derivations of algorithm performance | 191 |
| B.1 Non-blocked matrix multiplication | 192 |
| B.2 Non-blocked LU decomposition | 194 |
| B.3 Non-blocked forward substitution | 196 |
| B.4 Non-blocked aggregate forward substitution | 199 |
| B.5 Blocked LU decomposition | 203 |
| B.6 Blocked aggregate forward substitution | 209 |
| Appendix C. Software development toolchain | 213 |
| C.1 MDML programming model | 214 |
| C.2 Writing programs | 216 |
| C.3 Program structure | 217 |

List of Figures

Chapter 2. A survey of parallel computer systems

| | | |
|-----|---|----|
| 2.1 | The memory-memory MatRISC architecture | 15 |
| 2.2 | The load-store MatRISC architecture | 17 |
| 2.3 | The distributed MatRISC architecture | 18 |
| 2.4 | Trident system architecture | 20 |
| 2.5 | Trident matrix accelerator architectural detail | 21 |
| 2.6 | MatCore matrix accelerator architectural detail | 23 |

Chapter 3. The Matrix Data Management Layer (MDML)

| | | |
|-----|--|----|
| 3.1 | Drivers and objectives of the design | 32 |
| 3.2 | System model of the Matrix Data Management Layer | 36 |
| 3.3 | Connections within the MDML | 38 |
| 3.4 | Typical application structure | 39 |
| 3.5 | Major subsystems within a node | 42 |
| 3.6 | Flow of data through a data processing subsystem | 44 |
| 3.7 | Overview of node-level memory design | 46 |

| | | |
|-----|---|----|
| 4.1 | Form of an indexed loop structure | 54 |
| 4.2 | Kernel routine for dense matrix-matrix multiplication | 60 |

Chapter 4. Data movement within an MDML node

| | | |
|-----|---|----|
| 4.3 | Illustration of loop rebasing | 63 |
|-----|---|----|

List of Figures

| | | |
|------|---|----|
| 4.4 | Execution chains merged into a single structure | 66 |
| 4.5 | Nontrivial elements of lower triangular matrix \mathbf{C} | 70 |
| 4.6 | Pipeline detail of a single-node system | 72 |
| 4.7 | The Matrix Data Management Layer's sequencing engine | 75 |
| 4.8 | Difference engine from the matrix address sequencer module | 79 |
| 4.9 | Data-path of the loop parameter management unit | 82 |
| 4.10 | Changing line lengths (w_i) in a banded matrix | 87 |

Chapter 5. Data movement between nodes

| | | |
|------|--|-----|
| 5.1 | Schematic of a networked node | 100 |
| 5.2 | Inter-node data movement operations | 101 |
| 5.3 | A synchronisation domain and its messaging tree | 104 |
| 5.4 | Configuration example for a synchronisation domain | 105 |
| 5.5 | Interface detail of the synchronisation manager core | 108 |
| 5.6 | Interfaces of the data transport module | 110 |
| 5.7 | Internal detail of the data transport module | 112 |
| 5.8 | Transfer types supported by the data transport module | 113 |
| 5.9 | VHDL model to verify data flow in a multi-node MDML | 115 |
| 5.10 | Array control components and integration into the node | 116 |
| 5.11 | Network hop-size discovery process | 120 |
| 5.12 | Architecture of the parameter discovery module | 121 |

Chapter 6. Programming the MDML

| | | |
|-----|--|-----|
| 6.1 | Resource utilisation of small matrix-vector multiplication | 132 |
|-----|--|-----|

| | | |
|------|--|-----|
| 6.2 | Resource utilisation of LU decomposition on small matrices . . . | 140 |
| 6.3 | Behaviour of solving $L.X = A$ on small MDML arrays | 144 |
| 6.4 | Resource utilisation when solving $L.X = A$ on small matrices | 150 |
| 6.5 | Views of a partitioned matrix | 155 |
| 6.6 | An n -by- n matrix on a p -by- p array | 156 |
| 6.7 | Algorithm set for blocked LU decomposition | 160 |
| 6.8 | Blocked LU decomposition timing | 165 |
| 6.9 | Utilisation of array resources for blocked LU decomposition . | 166 |
| 6.10 | Speedup of blocked LU decomposition | 166 |
| 6.11 | Algorithm set for blocked aggregated forward substitution . . | 169 |
| 6.12 | Timing details for blocked forward substitution | 173 |
| 6.13 | Resource utilisation of blocked forward substitution | 174 |
| 6.14 | Speedup of blocked forward substitution | 174 |

Appendix A. Format of MDML control words

| | | |
|-----|---|-----|
| A.1 | VLIW fields for matrix acceleration | 186 |
| A.2 | VLIW fields for inter-node data transport | 189 |

Appendix B. Derivations of algorithm performance

| | | |
|-----|--|-----|
| B.1 | Problem definition of non-blocked matrix multiplication . . . | 192 |
| B.2 | Cycle activity for non-blocked matrix multiplication | 193 |
| B.3 | Cycle activity for non-blocked LU decomposition | 195 |
| B.4 | Cycle activity for non-blocked forward substitution | 197 |
| B.5 | Structure of matrices; aggregate forward substitution (case 1) . | 201 |

List of Figures

| | | |
|------|---|-----|
| B.6 | Cycle activity; non-blocked aggregate forward substitution (case 1) | 201 |
| B.7 | Structure of matrices; aggregate forward substitution (case 2) . | 202 |
| B.8 | Combined cycle activity; non-blocked aggregate forward substitution (case 2) | 202 |
| B.9 | Kernel routine invocations of non-blocked LU decomposition needed by block LU decomposition | 205 |
| B.10 | Kernel routine invocations to solve $(\mathbf{UL})^T = \mathbf{A}^T$ for \mathbf{L}^T as needed by block LU decomposition | 206 |
| B.11 | Kernel routine invocations to solve $\mathbf{LU} = \mathbf{A}$ for \mathbf{U} as needed by block LU decomposition | 207 |
| B.12 | Kernel routine invocations to do outer-product updates for block LU decomposition | 209 |
| B.13 | Kernel routine invocations to solve $\mathbf{LX} = \mathbf{A}$ for \mathbf{X} as needed by block forward substitution | 210 |
| B.14 | Kernel routine invocations to perform outer-product update as needed by block forward substitution | 211 |

Appendix C. Software development toolchain

| | | |
|-----|---|-----|
| C.1 | Allocation of node roles in an MDML array | 215 |
| C.2 | Styles of application mapping | 215 |
| C.3 | Programming flow for the Matrix Data Management Layer . . | 216 |

List of Tables

Chapter 3. The Matrix Data Management Layer (MDML)

| | | |
|-----|--|----|
| 3.1 | Classification of algorithms in the Basic Linear Algebra Sub-routines (BLAS) library | 39 |
|-----|--|----|

Chapter 4. Data movement within an MDML node

| | | |
|-----|--|----|
| 4.1 | Instruction breakdown of common matrix kernels by task | 62 |
| 4.2 | Transitions between loop indices depend on loop-level | 78 |
| 4.3 | Event table structure | 81 |
| 4.4 | Resources required for dependency chains of different depths | 81 |
| 4.5 | Effect of event-assisted matrix accelerator | 88 |

Chapter 5. Data movement between nodes

| | | |
|-----|---|-----|
| 5.1 | Phases of a lock-step operation | 102 |
| 5.2 | Features of techniques for sensing synchronisation loss | 107 |
| 5.3 | Operations defined for the bootstrap process | 117 |

Chapter 6. Programming the MDML

| | | |
|-----|---|-----|
| 6.1 | Scalar quantities used in this this chapter | 126 |
| 6.2 | Representing matrix, vector and scalar quantities | 126 |
| 6.3 | Data layout of a small matrix multiplication on the MDML | 130 |
| 6.4 | Behaviour of small matrix multiplication on the MDML | 132 |
| 6.5 | System activity for important matrix multiplication cases | 133 |
| 6.6 | Data layout of a small LU decomposition on the MDML | 135 |

List of Tables

| | | |
|------|--|-----|
| 6.7 | Behaviour of a small LU decomposition on the MDML | 139 |
| 6.8 | Data layout of a small forward substitution of maximum size | 141 |
| 6.9 | Behaviour of a small forward substitution on the MDML . . . | 143 |
| 6.10 | Behaviour of a small transposed forward substitution | 146 |
| 6.11 | Data layout of aggregate forward substitution | 148 |
| 6.12 | Behaviour of small aggregate forward substitution | 149 |
| 6.13 | Data layout of aggregate transposed forward substitution . . . | 151 |
| 6.14 | Behaviour of aggregate transposed forward substitution . . . | 152 |
| 6.15 | Sequencer configuration for blocked LU decomposition | 162 |
| 6.16 | Execution time for constituents of block LU decomposition . . | 163 |
| 6.17 | Repeats of constituent functions for block LU decomposition . | 163 |
| 6.18 | Sequencer configuration for sub-algorithm A1 | 170 |
| 6.19 | Sequencer configuration for subalgorithm A2 | 171 |

Appendix B. Derivations of algorithm performance

| | | |
|-----|---|-----|
| B.1 | Activities performed by non-blocked matrix multiplication . . | 193 |
| B.2 | Activities performed by non-blocked LU decomposition . . . | 195 |
| B.3 | Activities for aggregate forward substitution | 200 |

Appendix C. Software development toolchain

| | | |
|-----|---|-----|
| C.1 | API commands for matrix traversal | 219 |
|-----|---|-----|

Abstract

Using super-resolution techniques to estimate the direction that a signal arrived at a radio receiver; Tracking moving targets using particle filters; Applying advanced coding techniques to radio transmissions. Real-world applications like these rely on high-speed matrix computations. Although the computational ability of general-purpose computer architectures is growing, some numerically intensive calculations (such as those above) may benefit from specialised matrix processing hardware.

Work in this thesis builds on earlier dense matrix accelerators, which were previously implemented as coprocessors attached to general purpose processor cores. It is well known that matrix algorithms are highly parallelisable, so the coprocessor generally contains many cores to take advantage of this. However, encapsulating them within the coprocessor often prohibits their use in non-matrix computations. This limitation is addressed by the Matrix Data Management Layer (MDML), described in this thesis.

The MDML doesn't share the coprocessor-based architecture of earlier systems. Growing transistor budgets have made it feasible to embed general-purpose Commercial Off-The-Shelf (COTS) processors within the MDML. Programming is simplified as the system is homogeneous, and programmers can use familiar tools and languages. Having general-purpose processors embedded throughout the array also opens the possibility to process more than just dense matrix data.

Requirements for the MDML were synthesised by analysing common features and operations in matrix algorithms. This analysis found that data management (including sequencing, distribution and efficient reuse of data) was a significant overhead. Thus, the MDML focuses on accelerating these tasks, which are common to more than just matrix algorithms. Effort was also directed into integrating COTS components to reduce design effort, which resulted in a layered hardware design. The MDML integrates off-the-shelf memory and processing cores into a scalable on-chip mesh array.

Abstract

The MDML uses a flexible sequencing regime capable of autonomously generating memory access patterns for a variety of complex data structures. This sequencer has been tested by construction of a hardware prototype. A programming style based on sequences of instruction snippets was also developed and shown effective for loop-based algorithms.

A flexible, low-latency data transport allows the MDML to use fine-grain communication. Improved array synchronisation techniques were also proposed and implemented. These techniques reduce the performance impact of barrier synchronisation, and also reduce the frequency with which the array must synchronise.

The MDML concept was verified using both hardware modelling and theoretical analysis:

- Hardware modelling was used to verify the completeness and correctness of the architecture, and demonstrated that the MDML is practical to implement in current hardware.
- Analytic models of performance were developed for algorithms mapped to a theoretical MDML system. These were used to evaluate performance, scalability and programmability of the MDML. Algorithms that were mapped include conformal matrix multiplication, LU decomposition and forward substitution.

Thus, it is concluded that the MDML architecture provides a realistic and efficient way of achieving high-speed matrix computation.

Statement of Originality

I certify that this work contains no material that has been accepted for the award of any other degree or diploma in any university or other tertiary institution and, to the best of my knowledge and belief, contains no material previously published or written by another person, except where due reference has been made in the text. In addition, I certify that no part of this work will, in the future, be used in a submission for any other degree of diploma in any university or other tertiary institution without the prior approval of the University of Adelaide.

I give consent to this copy of the thesis, when deposited in the University Library, being available for loan, photocopying, and dissemination through the library digital thesis collection, subject to the provisions of the Copyright Act 1968.

I also give permission for the digital version of my thesis to be made available on the internet, via the University's digital research repository, the Library catalogue, the Australasian Digital Thesis Program (ADTP) and also through web search engines, unless permission has been granted by the University to restrict access for a period of time.

Signed

Date

Acknowledgments

It is with joy that I sit and write this note, and think of many great colleagues, friends and mentors who helped create this thesis.

First, I am indebted to my supervisors Assoc Professors Cheng Chew Lim and Mike Liebelt for their mentorship and guidance. They nurtured my independence as a new researcher, and patiently shaped and guided the direction of this work. The friendly assistance of school support staff, including Yadi, Rosie and Stephen, was also invaluable.

I am also grateful to John Kitchen and Gareth Parker, who encouraged me to pursue this research and gave me the time and resources in which to do it. I promise it will be done by Christmas!

I must also thank my examiners. It isn't easy to write a thesis, and it's no less easy to review one. I'm grateful for their thoughtful, constructive feedback.

Thanks also to my PhD comrades, for their friendship and encouragement over many years. I'm notably grateful to thank Kiet To for his suggestions and advice. His depth of knowledge and practical capabilities never failed to impress.

I'd like to thank my good friends: Lachie, Darryl, Bobby, Rob and Hong-Gunn... and all members of the Gang. You all supported and distracted me in equal measure.

I'm extremely thankful for my fabulous family for being pillars of emotional support and encouragement. My thesis wouldn't exist without them.

Last, but not least, thanks to my dear Jennie. Without your enthusiasm and reassurance this thesis would still be years from completion.

Definitions and acronyms

ABI Application Binary Interface

AGU Address Generator Unit

AMU Array management unit

ARQ Automatic Repeat Request

ASIC Application Specific Integrated Circuit

BAG Built-in Address Generator

BLAS Basic Linear Algebra Subroutines

BLAST Bell-Labs Layered Space-Time

CBE Cell Broadband Engine

CISC Complex Instruction Set Computer

COTS Commercial Off-The-Shelf

CPU Central Processing Unit

CSP Communicating Sequential Processes

DCT Discrete Cosine Transform

DFT Discrete Fourier Transform

DSP Digital Signal Processor

DMA Direct Memory Access

DP Data Processor

Definitions and acronyms

DRAM Dynamic RAM

DTM Data Transport Module

FEC Forward Error Correction

FIFO First-In First-Out

FLOPS Floating Point Operations

FP Floating Point

FPGA Field-Programmable Gate Array

FSM Finite State Machine

GEMM General Matrix/Matrix Multiply

GEMV General Matrix/Vector Multiply

GER General rank 1 update

GPU Graphics Processing Unit

ILP Instruction-Level Parallelism

IPC Instructions Per Clock-cycle

I/O Input-Output

JCU Jump Control Unit

KIMP Kilo-Instruction Multi-Processor

L1 Level-1

L2 Level-2

LCU Loop Control Unit

LRU Least Recently Used

MAC Multiply Accumulate

MAG Matrix Address Generator

MDML Matrix Data Management Layer

MIMD Multiple Instruction Multiple Data

MIMO Multiple Input Multiple Output

MMX MultiMedia eXtensions

MPI Message Passing Interface

MPSoC Multi-Processor System-on-Chip

MuSiC Multiple Signal Classification

NEWS A mesh-type interconnect with links in four directions. These links are named after the four main compass directions.

OFDM Orthogonal Frequency Division Multiplexing

PCI-Express A modern, high-speed interconnect technology.

PE Processing Element

PPE Power Processing Element

RAM Random Access Memory

RAW Read After Write

RISC Reduced Instruction Set Computer

RTEMS Real-Time Executive for Multiprocessor Systems

RTL Register Transfer Level

SCAP Scalable Array Processor

ScaLAPACK Scalable Linear Algebra Package

SDR Software Defined Radio

SIMD Single Instruction Multiple Data

SMP Symmetric Multi-Processor

Definitions and acronyms

SoC System-on-Chip

SPE 'synergistic' processing elements

SRAM Static RAM

SVD Singular Value Decomposition

SYMM Symmetric Matrix/Matrix Multiply

TLP Thread-Level Parallelism

TRMM Triangular Matrix/Matrix Multiply

TRSM Triangular Solve with Multiple Right-hand sides

VHDL A hardware modelling language.

VLIW Very Long Instruction Word

VLSI Very Large-scale Silicon Integration