

SCALING A RECONFIGURABLE DATAFLOW ACCELERATOR

A DISSERTATION

SUBMITTED TO THE DEPARTMENT OF ELECTRICAL ENGINEERING

AND THE COMMITTEE ON GRADUATE STUDIES

OF STANFORD UNIVERSITY

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS

FOR THE DEGREE OF

DOCTOR OF PHILOSOPHY

Yaqi Zhang

May 2020

Abstract

With the slowdown of Moore’s Law, specialized hardware accelerators are gaining tractions for delivering 100-1000x performance improvement over general-purpose processors in a variety of applications domains, such as cloud computing, biocomputing, artificial intelligence, etc. [1, 2, 3]. As the performance scaling in multicores is coming to a limit [4], a new class of accelerators—reconfigurable dataflow architectures (RDAs)—is promising in offering high-throughput and energy-efficient acceleration that keeps up with the performance demand. Instead of dynamically fetching instructions like in traditional processors, RDAs have flexible datapath that can be statically configured to spatially parallelize and pipeline the program across distributed on-chip resources. The pipelined execution model and explicitly-managed scratchpad in RDAs eliminate the performance, area, and energy overhead in dynamic scheduling and conventional memory hierarchy.

To adapt to the compute intensity in modern data-analytic workloads, particularly in the deep learning domain, RDAs are increasing to a scale that was unprecedented before. With an area footprint of 133mm² at 28nm, Plasticine is a hierarchical RDA supplying 12.3 TFLOPs of computing power [5]. Prior work has shown an up to 76x performance/watt benefit from Plasticine over a Stradix V FPGA due to an advantage in clock frequency and resource density. The increase in scale introduces new challenges in network-on-chip design to maintain the throughput and energy efficiency of RDAs. Furthermore, targeting and managing RDAs at this scale also require new strategies in mapping, memory management, and flexible control to fully utilize their compute power.

In this work, we focus on two aspects of the software-hardware co-design that impact the usability and scalability of the Plasticine accelerator. Although RDAs are flexible to support a wide range of applications, the biggest challenge that hinders the adoption of these accelerators is the required low-level knowledge in microarchitecture design and hardware constraints in order to efficiently implement a new application. To address this challenge, we introduce a compiler stack—SARA—that raises the programming abstraction of Plasticine to an imperative-style domain-specific language

(DSL) with nested control flow for general spatial architectures. Besides architecture-agnostic, this abstraction contains explicit loop constructs, enabling cross-kernel optimizations that are often not explored when programming RDAs. SARA efficiently translates imperative control constructs to a streaming dataflow graph that scale performance with distributed on-chip resources. By virtualizing resources, SARA systematically handles the physical constraints, hiding the low-level physical limitations from the programmers. To address the scalability challenge with increasing chip sizes in RDAs, we present a comprehensive study on the network-on-chip design space for RDAs [6]. We found that network performance highly correlates to network bandwidth, as supposed to latency, for RDAs with streaming dataflow execution model. Lastly, we show that a static-dynamic hybrid network design can sustain performance in a scalable fashion with high energy efficiency.

Acknowledgements

Contents

Abstract	iv
Acknowledgements	vi
1 Introduction	1
1.1 The Rising of Hardware Acceleration	1
1.2 The Need for Flexible Interconnects	4
1.3 The Gap between High-Level DSLs and Dataflow Accelerators	5
1.4 Contribution	7
1.5 Outline	8
2 Background	9
2.1 Execution Schedules of Reconfigurable Architectures	9
2.2 The Plasticine Architecture and its Programming Interface	11
2.3 The High-Level Imperative Compiler	15
3 Compiler	19
3.1 SARA Compiler Overview	19
3.2 Imperative to Dataflow Transformation	22
3.2.1 Loop Division	22
3.2.2 Virtual Context Allocation	23
3.2.3 Control Allocation	27
3.2.4 Data-Dependent Control Flow	33
3.2.5 Virtual Unit Allocation	37
3.3 Resource Allocation	37

3.3.1	Compute Partitioning	41
3.3.2	Memory Partitioning	49
3.3.3	Register Allocation	51
3.4	Optimizations	51
3.4.1	Description	51
3.4.2	Evaluation (WIP)	57
3.5	Placement and Routing	58
3.5.1	Iterative Placement	58
3.5.2	Congestion-Aware Routing	59
3.5.3	VC Allocation for Deadlock Avoidance	60
3.5.4	Runtime Analysis for Heuristic Generation	61
3.6	Evaluation (WIP)	61
4	Architecture	64
4.1	On-chip Network	64
4.1.1	Application Characteristics	66
4.1.2	Design Space for Network Architectures	69
4.1.3	Performance, Area, and Energy Modeling	72
4.1.4	Network Architecture Evaluation	76
4.2	Plasticine Specialization for RNN Serving	83
4.2.1	Mixed-Precision Support	84
4.2.2	Folded Reduction Tree	86
4.2.3	Sizing Plasticine for RNN Serving	87
4.3	Generic Banknig Support	87
5	Related Work	90
5.1	Compiler	90
5.1.1	Streaming Dataflow IRs	90
5.1.2	Spatial Compilers	90
5.2	On-Chip Network	91
5.2.1	Tiled Processor Interconnects	91
5.2.2	CGRA Interconnects	91
5.2.3	Design Space Studies	92

5.2.4	Compiler Driven NoCs (WIP)	92
6	Conclusions	93
	Bibliography	94

List of Tables

1.1	OPS comparison of different CGRAs	3
3.1	Mapping between data-structure to hardware memories	29
3.2	Interference Table	30
3.3	Resources specification of heterogeneous units and aggregation function	39
3.4	Formulation of the compute partitioning problem	42
3.5	Names and definitions used in the solver-based algorithms.	45
3.6	Solver formulation for partitioning. Expressions are presented using the Disciplined Convex Programming ruleset [7, 8]. Explanations for selected expressions can be found in the supplemental material.	46
3.7	Performance comparison of Plasticine with Tesla’s V100 GPU (Normalized throughput to transistor count in parentheses).	61
4.1	Benchmark summary	67
4.2	Network design parameter summary.	72

List of Figures

1.1	Average utilization vs. peak compute density trade-off	2
1.2	High-level performance model of a spatial architecture	4
1.3	Big picture of machine learning frameworks	6
2.1	Hierarchical pipelining and parallelization on spatial architecture	10
2.2	Plasticine chip-level architecture	11
2.3	Example PCU configuration	13
2.4	Example of Outer Product in Spatial.	16
2.5	Spatial compiler stack to target FPGAs and Plasticine	16
2.6	Spatial Example	18
3.1	(a) A naïve mapping strategy to map the control hierarchy onto Plasticine. All units are distributed across an on-chip network that can introduce unpredictable latency. The number on the edges indicates event order. Here we show a scenario where read requests from <i>C</i> do not observe the write requests from <i>B</i> that occur earlier in the program order due to network latency between <i>B</i> and <i>mem</i> . (b) Distributed controllers in SARA. Each innermost controller makes a copy of all enclosing controllers. The signals from these controllers are used to generate synchronization between distributed compute units. To address the problem in (a), the memory also needs to provide a write acknowledgment per write request for synchronization.	20
3.2	SARA Compiler Flow	21
3.3	Example of loop fission vs. loop division	23
3.4	Example of an illegal loop fission and a legal loop division	24
3.5	Context allocation and control allocation	25
3.6	Dense specialization	26

3.7	Examples for access dependency graph	31
3.8	Examples for dependency graph reduction	32
3.9	Example of dynamic loop range	33
3.10	Branching example	34
3.11	Do while example	36
3.12	Compute partitioning examples	43
3.13	Compute partitioning examples with cycle	43
3.14	Partitioning and merging algorithm comparisons	47
3.15	Memory model of different architectures	49
3.16	Memory partitioning	50
3.17	Memory strength reduction	52
3.18	Examples of route through elimination	53
3.19	Retiming	54
3.20	MLP case study	55
3.21	Reversed loop invariant hoisting	56
3.22	Optimization effectiveness	57
3.23	An example of deadlock in a streaming accelerator, showing the (a) VU data-flow graph and (b) physical placement and routes on a 2×3 network. There are input buffers at all router inputs, but only the buffer of interest is shown.	60
3.24	Scalability Evaluation	62
3.25	Performance comparison with V100 GPU	63
4.1	Area and power breakdown of Plasticine	65
4.2	Resource and bandwidth utilization	68
4.3	Application communication patterns	68
4.4	Characteristics of program graphs.	70
4.5	Switch and router power with varying duty cycle.	74
4.6	Area and per-bit energy for (a,d) switches and (b,c,f) routers. The router only has a vector granularity and can be partially clock-gated when sending scalar packets. Subplots (f) show the energy of the vector router when used for scalar values (32-bit).	75
4.7	Area breakdown for all network configurations.	76
4.8	Performance scaling with increased CGRA grid size for different networks.	77
4.9	Impact of bandwidth and flow control strategies in switches.	78

4.10	Impact of VC count and flit widths in routers.	79
4.11	Number of VCs required for dynamic and hybrid networks. (No VCs indicates that all traffic is mapped to the static network.)	80
4.12	Normalized performance for different network configurations.	80
4.13	(a,d): Normalized performance/watt. (b,e): Percentage of compute and memory PUs utilized for each network configuration. (c,f): Total data movement (hop count).	82
4.14	Geometric mean improvement for the best network configurations, relative to the worst configuration.	83
4.15	Plasticine PCU SIMD pipeline and low-precision support. Red circles are the new operations. Yellow circles are the original operations in Plasticine. In (d) the first stage is fused 1 st , 2 nd stages, and the second stage is fused 3 rd , 4 th stages of (b).	85
4.16	Variant Plasticine configuration for RNN serving	87

Chapter 1

Introduction

1.1 The Rising of Hardware Acceleration

With the end of Dennard Scaling [9], the amount of performance one can extract from a CPU is reaching a limit. To provide general-purpose flexibility, CPUs spend the majority of resources and energy on overheads, including dynamic-instruction fetching and decoding, branch prediction, and a cache hierarchy, etc., with less than 20% of the energy on the actual computation [10]. Even worse, power wall is limiting the entire multicore family to reach the doubled performance per generation scaling enabled by technology scaling in the past [4].

For this reason, hardware acceleration is emerging in various compute-intensive application domains to provide orders of magnitude acceleration, enabling algorithms that were otherwise infeasible [3, 2, 11, 12]. Examples include widely adopted General-Purpose Graphics Processing Units (GPGPUs) in computational genomics, signal processing, graph processing, and deep learning, etc. [3, 2, 1]. Moreover, many recent efforts are spent on leveraging application domain knowledge in hardware design to enable continued performance scaling while meeting the power budget [13]. As artificial intelligence receiving great success in industry and business, past years have seen a growing interest in machine learning accelerators; these accelerators contain specialized circuits for ML kernels that dramatically improve the compute efficiency [14, 15, 16, 17, 18, 19].

Nonetheless, it is non-trivial to achieve good utilization of these accelerators. While the peak FLOPS of GPU has increased by over 20x in the past 10 years, the achievable FLOPS is not increasing accordingly [20, 21]. The massive threads in GPU require embarrassingly parallel workloads to



Figure 1.1: Tradeoff between the average utilization of the peak FLOPS over a range of applications vs. the peak compute density among different architectures. The shaded area indicates the effective FLOPS achieved. With all the dynamic scheduling hardware in CPUs, such as prefetching and branch prediction, CPUs can achieve a fairly decent instruction per cycle (IPC) for data-analytic workloads that have a regular control flow. However, overheads to support flexibility, security, and programmability in CPUs result in a low FLOPS density. While GPU and fix-function accelerators have a high FLOPS density, they are prone to underutilization due to variation in application characteristics. Fine-grained reconfigurable datapath makes it easy to utilize on-chip resources of an FPGA. However, the soft logics and overheads in routing resources lead to a low resource density and clock frequency. RDAs have the right balance between flexibility and efficiency, which gives a good overall average utilization.

fully saturate the compute throughput. GPU’s bulk-synchronous nature also causes poor cache efficiency; data are spilled off-chip before getting reused [22]. On the other hand, while extremely efficient for certain models, machine learning accelerators are highly specialized for specific kernels, especially general matrix multiply (GEMM) and convolution. However, ML algorithms evolve much faster than the development and manufacture cycle of hardware, leading to inefficiency or even unsupported applications. For example, hybrid models, such as Mask R-CNN [23] and DeepLab [24], contain combinations of GEMM-compatible and incompatible operations, both computationally expensive. Fix-functional accelerators focusing on GEMM operations, such as TPU [15], have to rely on CPUs for unsupported operations or convert non-GEMM operations to GEMM operations, resulting in a severe performance gap between the peak and effective FLOPS [25]. Figure 3.25 illustrates a trade-off between average utilization of the peak FLOPS over a range of applications and the peak FLOPS available on the hardware.

Reconfigurable spatial architectures overcome this limitation by changing their datapaths based on applications’ needs. Applications are configured at the circuit-level without dynamic instruction fetching and decoding, hence improving energy-efficiency [26, 27]. In addition to instruction, data, and task-level parallelism explored by processor architectures, spatial architectures also explore instruction and task-level pipelining that further increases the compute throughput [28]. Exploring

Architectures	DySER [38]	TI [43]	REVEL [51]	Plasticine [5]	Gorgon [48]
OPS	128GOPS	64GOPS	300GOPS	12.3TFLOPS	38.4TFLOPS

Table 1.1: OPS comparison of different CGRAs. Gorgon is a Plasticine variant proposed for joint machine learning and database acceleration.

pipeline parallelism enables spatial architectures to achieve a high-throughput without massively parallelize every stage of the program. Flexible datapath also permits resource distribution proportional to the compute intensity of program stages. A key performance optimization on reconfigurable accelerators is an application-level design space exploration searching for the best resource distribution scheme that balances the compute pipeline [29].

One of the mainstream reconfigurable spatial architecture is Field Programmable Gate Arrays (FPGAs) that support fine-grain, bit-level reconfigurability with a soft logic fabric [30]. The flexible interconnect and lookup table-based logic gate can be configured to implement arbitrary datapath. FPGAs have been used to deploy services commercially [31, 32, 33] and can be rented on the AWS F1 cloud [34]. Although around for a long time, FPGAs are not broadly accepted among high-level application programmers due to their low-level programming interface and long compilation time. Fine-tuning applications on FPGAs requires expertise in digital design and takes a long development cycle, which hinder their accessibility to the general software community. As application-level accelerators, FPGAs also suffer from overhead in fine-grained reconfigurability; studies have shown that over 60% of the chip area of an FPGA is spent on routing resources [35, 26, 27].

In contrast to fine-grained reconfigurable architectures, Coarse-Grained Reconfigurable Arrays (CGRAs) are spatial architectures with coarse-grained building blocks, such as ALUs, register files, and memory controllers, distributed in a programmable, word or vector-level static interconnect [36, 37, 38, 39, 40, 41, 42]. We refer a subclass of CGRAs with dataflow-driven execution model as Reconfigurable Dataflow Architectures (RDAs)¹ [5, 43, 44, 45, 46, 47]. Lately, RDAs are emerging as a new class of spatial accelerators that retain the desired level of flexibility and energy efficiency without the area overhead and low clock frequency of bit-level reconfigurability. Our previously proposed RDA–Plasticine–has demonstrated a promising acceleration of dense, sparse, database, and streaming applications [5, 48, 49, 50]. To meet the computing demand of recent data-analytic workload, Plasticine is a large-scale RDA compared to traditional CGRAs. Table 1.1 shows the OPS comparison across a few CGRAs proposed in the prior works.

¹Dataflow overlay architectures on FPGA are technically also RDAs, which are not the primary concern in the discussion of this work.

$$\text{thrpt}_{\text{app}} = \min \left(\begin{array}{l} \frac{\text{op}}{\text{sec}} \\ \frac{\text{op}}{\text{byte}_{\text{off}}} \text{BW}_{\text{off}} \\ \frac{\text{op}}{\text{byte}_{\text{on}}} \text{BW}_{\text{on}} \\ \frac{\text{op}}{\text{byte}_{\text{net}}} \text{BW}_{\text{net}} \end{array} \right)$$

Throughput	Proportional To	
Compute	$\text{op} \sim P, D$	
Off-chip Memory	$\text{byte}_{\text{off}} \sim P, D$	
On-chip Memory	$\text{byte}_{\text{on}} \sim P, D$	
On-chip Network	$\text{byte}_{\text{net}} \sim P^2, D$	

Application-specific
Hardware-specific
P: Parallelization factor
D: Pipelining depth

Figure 1.2: High-level performance model of a spatial architecture.

The scale of Plasticine introduces challenges in network-on-chip design to sustain bandwidth requirements while staying energy-efficient. The compilation strategy also needs to explore multiple levels of concurrency in the program to saturate the compute throughput of a large-scale RDA; this strategy must introduce minimum synchronization overhead to maximize the scalability of the mapped designs.

1.2 The Need for Flexible Interconnects

Applications are mapped to RDAs by distributing computations spatially across multiple processing blocks and executing them in a pipelined, data-driven fashion. In traditional Networks-on-Chip (NoCs) for multicore systems, communication is the result of explicit message passing between parallel workers or cache misses that generate messages for coherence protocol; these traffics are bursty and relatively infrequent. On RDAs, however, applications are distributed by parallelizing and pipelining; pipelining introduces frequent and throughput-sensitive communication. Applications with different characteristics, such as compute vs. memory-bound, also exhibit very different communication patterns.

Figure 1.2 shows a high-level performance model of a spatially pipelined and parallelized reconfigurable architecture. Due to pipelined execution, the performance of a spatial architecture is mainly determined by throughput as supposed to latency. Compute, on and off-chip memory accesses, and network become an integrated pipeline; overall performance is limited by the pipeline stage with the lowest throughput. The red terms in the equation are application-specific and capture the compute, memory, or IO-bound characteristics. The blue terms are the bandwidth and FLOPS available on the hardware. While the compute and memory access in application increases linearly

with parallelization factor and pipelining depth, the required network bandwidth from the applications increases quadratically with increasing parallelism. This means as we going to larger chip size, the network bandwidth must grow super linearly to achieve a perfect performance scaling, unless exploring pipeline parallelism.

RDAs need the right amount of interconnect flexibility to achieve good resource utilization; an inflexible interconnect constrains the space of valid application mappings and hinders resource utilization. Furthermore, in the quest to increase compute density, RDA datapaths now contain increasingly coarse-grained processing blocks such as pipelined, vectorized functional units [5, 39, 52]. Plasticine, as an example, has a 512-bit vector bus, which necessitates coarser communication and higher on-chip interconnect bandwidth to avoid creating performance bottlenecks. Although many hardware accelerators with large, vectorized datapaths have fixed local networks [53], there is a need for more flexible global networks to adapt to future applications. Consequently, interconnect design for these RDAs involves achieving a balance between the often conflicting requirements of high bandwidth and high flexibility.

1.3 The Gap between High-Level DSLs and Dataflow Accelerators

Recent years have seen an explosion of hardware accelerators, primarily motivated by AI applications [18, 19, 54, 14, 55, 17, 16, 56, 57, 58, 59, 60]. Research in software infrastructures for these accelerators, however, is just emerging. Unlike CPUs, accelerators do not support a standard instruction set architecture (ISA), alleviating the overhead from layers of abstractions and the burden of backward compatibility. Nonetheless, lack of a common abstraction makes it very hard to share compiler infrastructure across accelerators.

Most RDAs comes with a co-designed software layer that is very low-level and restrictive, requiring expert knowledge in hardware architecture to efficiently target the accelerator. On the other side, machine learning frameworks [61, 62, 63, 64, 65, 66] provide high-level and succinct front-end abstraction that targets multiple hardware platforms. Yet most of theses frameworks only target mainstream accelerators, such as GPUs and FPGAs. Very few of the accelerators mentioned above can support more than a few proof-of-concept models with an end-to-end integration with these frameworks.

Part of the reason behind this is the massive engineering effort to program these accelerators.

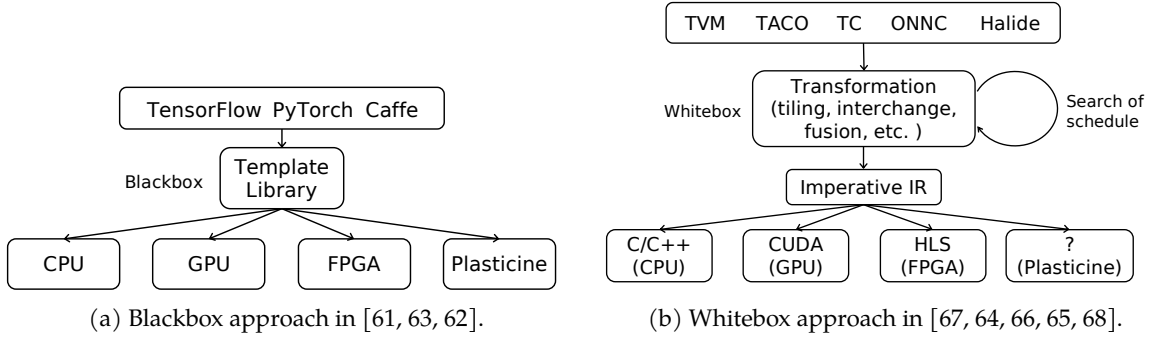


Figure 1.3: Big picture of machine learning frameworks.

Frameworks like TensorFlow [61] and PyTorch [62] use a template library approach, as shown in Figure 1.3 (a). These frameworks treat operations in the ML dataflow graph as black-box kernels and implement every operation for every backend, requiring expensive engineering hours and hardware expertise to fine-tune each accelerator. Furthermore, cross-kernel optimizations, such as fusion, often provide orders of magnitude improvement in efficiency. Recent studies in sparse ML also show promising speedup on large ML models [16]. However, sparse kernels require various data formats with very different implementation than the dense version [69]. In the black-box approach, each fused operation and data format corresponds to a new black-box kernel. As a result, the combination of required templates easily becomes untractable.

The alternative white-box approach showing in Figure 1.3 (b) uses a succinct front-end representation, such as index notations [64], to express the computation of a kernel. This representation, however, cannot be directly mapped to accelerators due to various physical constraints. The framework then uses a series of transformations and optimizations to reformulate the kernel to an imperative implementation that can be supported by the accelerator. The transformation is often target-specific and studies have demonstrated ML-based cost models to facilitate the automatic searching of backend-specific schedules [67]. Next, the compiler can easily translate the imperative intermediate representation (IR) to the imperative front-ends of different accelerators.

To enable easy integration of the white-box approach, we introduce a compiler—SARA—that raise the programming abstraction of Plasticine to an imperative, loop-based language for general reconfigurable hardware—Spatial [70]. Plasticine is a pure dataflow accelerator without a centralized scheduler that executes instructions. This design makes the architecture very scalable to achieve

FLOPS comparable to high-end GPUs and FPGAs. Nonetheless, it is not intuitive nor easy to implement a linear algebra kernel with the low-level declarative streaming configurations. By supporting imperative constructs on a data-flow architecture efficiently, SARA not only improves the programmability of Plasticine, but also enables cross kernel optimizations. These optimizations have demonstrated an average 30x speedup over a Tesla V100 GPU and 2x speedup over a Stratix 10 FPGA from a Plasticine with much less on-chip resource than the baseline architectures [71].

1.4 Contribution

This thesis work is upon the prior works on the Plasticine architecture [5] and the Spatial compiler [70]. We address the programmability and scalability challenges from a software and a hardware angle. We first introduce the Plasticine compiler—SARA—that maps Spatial applications to Plasticine. Next, we present several architectural augmentations to Plasticine to increase the flexibility in interconnect, memory access patterns, and control flow.

The key contributions of this work are summarized below:

1. We propose an imperative to dataflow transformation that implements imperative control constructs, such as nested loops and branch statements, on a pure dataflow architecture. This mapping strategy introduces minimum point-to-point synchronization that efficiently scales performance over distributed on-chip resources.
2. We introduce compiler-enforced memory consistency to preserve memory ordering in the imperative program between distributed memory accesses.
3. We illustrate a systematic approach to address physical constraints in a hierarchical RDA. We evaluate the tradeoff between a traversal-based vs. solver-based algorithms for program partitioning.
4. We use resource virtualization and optimizations to minimize fragmentation in resource allocation with heterogeneous compute tiles.
5. We demonstrate an evaluation of the scalability of SARA and effectiveness of compiler optimizations. We show a 1.9x average speedup over a Tesla V100 with a Plasticine with 8x less in on-chip resources due to the efficient use of the accelerator.
6. We analyze the key communication patterns exhibited by spatial architectures by characterizing a variety of benchmarks from different domains.

7. We demonstrate that static compiler knowledge can be used to reduce routing congestion.
8. We give a comprehensive quantitative analysis of the performance, area, and energy trade-offs involved in choosing an RDA network, using a cycle-accurate simulator and ASIC synthesis with a 28 nm industrial technology library. We explore a variety of design points, including static, dynamic, and hybrid networks, decreased flit widths and VC counts for dynamic networks and different flow-control strategies for static networks.
9. We show that bandwidth is the most critical metric to assess network performance for RDA, due to their pipelined execution nature. We found static-dynamic hybrid networks can give better energy efficiency by reducing overprovisioning in purely static networks and reducing data movement.
10. We present changes in Plasticine required to support flexible memory access patterns and flexible control flow.

1.5 Outline

The rest of this thesis is organized as follow: Chapter 2 gives a background on the execution schedule of spatial architectures, the front-end imperative language of SARA compiler, and the targeting architecture–Plasticine. Chapter 3 goes over the SARA compiler. Chapter 4 details the augmentations to the Plasticine architectures and the Network-on-chip study. Chapter 5 summarizes the related work. Chapter 6 concludes our work.

Chapter 2

Background

2.1 Execution Schedules of Reconfigurable Architectures

The key advantage of reconfigurable spatial accelerators, compared to processor-based architectures, is the ability to explore multiple levels of pipeline parallelism. In traditional Von Neumann architectures [72], like CPUs and GPUs, a computer consists of a processing unit that performs computation, a memory unit that stores the program states, and a control unit that tracks execution states and fetches the instruction to execute. This computing model inherently assumes that instructions within a program are executed in time, maximizing the flexibility to context switching between different workloads dynamically.

Reconfigurable accelerators are a direct violation of the von Neumann execution model; instructions are statically embedded in the datapath and executed in space as supposed to in time. The disadvantage of static reconfiguration is paying the resource cost for infrequently executed instructions, making it unsuitable for control-heavy workloads that traditional processors are efficient at. On the other hand, RDAs are particularly competitive in providing high-throughput, low-latency, and energy-efficiency acceleration for data-analytical workloads. Data-analytical workloads encompass wide domains of applications, including image processing, recognition, machine translation, digital signal processing, network processing, etc. These applications exhibit a rich amount of data-level parallelism with relatively static control flow.

Figure 2.1 shows an example of the hierarchical parallelism and pipelining exploit by a spatial architecture. The overall compute throughput of a parallelized and pipelined program is the product of the total parallelization factors and pipelining depth. By exploring multiple dimensions of



(a) Example program

(b) Timing of execution with hierarchical pipelining and parallelization

Figure 2.1: Hierarchical pipelining and parallelization in a spatial architecture. (a) illustrates the runtime and throughput of a hierarchically pipelined and parallelized program on a reconfigurable spatial architecture. At the inner level, instructions within each basic block are fine-grained pipelined across iterations of the innermost loop. At the outer level, the inner loops are coarse-grained pipelined across the outer loop iterations. Exploiting multiple levels of pipeline parallelism gives a total throughput of $x + y$ operations per cycle, where x and y are number of operations in the basic blocks. (b) Vectorizing the inner most loops B and C by n increases the throughput to $(x + y)n$. (c) Parallelizing the outer loop A by m further increases the throughput to $(x + y)mn$.

concurrency in the program, spatial architecture is more likely to saturate the compute throughput of the hardware for a wide range of applications. For applications that are expensive to parallelize due to irregular access patterns, spatial architectures can increase on the pipelining dimension; for application with embarrassingly parallel workloads, spatial architecture can budget most resource on increasing parallelism.

Pipelined execution is also beneficial for achieving good memory performance. Data accessed by different stages of the pipelines are stored in disjoint scratchpads instead of a shared cache, improving the address bandwidth and effective on-chip capacity. Using explicitly managed scratchpads with application knowledge also tends to improve locality and eliminate cache performance issues, such as thrashing. Across kernels, pipelined execution reduces the amount of off-chip accesses for intermediate data. SIMT architectures, like GPUs, relying on high-bandwidth DRAM technologies, such as HBM, to sustain the compute throughput of massively parallelized threads. While providing over 10x more bandwidth than traditional DDR technologies, HBM is very limited in capacity, around 16GB as supposed to on the orders of TB for DDR. As a result, the limited off-chip capacity often restricts the type of applications that GPUs can support.



Figure 2.2: Plasticine chip-level diagram

2.2 The Plasticine Architecture and its Programming Interface

Plasticine is a tile-based RDA designed for a wide range of data-intensive workloads. With an area footprint of 113mm^2 at 28-nm process, Plasticine packs 12.3 TFLOPS of compute throughput and 16 MB of on-chip memory. Without bit-level reconfiguration overhead, Plasticine runs a higher clock frequency than most FPGAs at 1GHz with a thermal design power at 49W. Previous work has shown a up to 77X performance-per-watt improvement from Plasticine over a Stratix V FPGA [5]. Figure 2.2 shows the chip-level design of Plasticine.

At high-level, Plasticine contains an array of resource tiles connected by a global interconnect. The mesh network is statically configured, providing a guaranteed in-order transmission of packet streams between any tiles. The network comes with three granularity: a 512-bit vector data bus, a 32-bit scalar data bus, and a single-bit control bus. The data networks contain a single valid bit traveled along with the data; the control network has only the valid bit without a payload, transmitting control pulses across tiles. There are three types of configurable units: the pattern compute units (PCUs) perform the most heavy lifting computation on Plasticine; the pattern memory units (PMUs) contain all distributed scratchpad on-chip; and the DRAM Address Generation Units (AGs) that generate DRAM requests going to the off-chip memory.

Like FPGAs, Plasticine can also support multiple-level of parallelization and pipelining explained in the previous section.

For the rest of this section, we will explain the native programming interfaces of the three configurable units.

Pattern Compute Unit (PCU)

As the major compute workhorse of the architecture, a PCU contains a 6-stage single instruction multiple data (SIMD) pipeline with 16 SIMD lanes. Unlike a processor core, the PCU can only statically configure six vector instructions throughout the entire execution. Additionally, the six instructions must be branch-free in order to be fully pipelined across stages. At runtime, the SIMD pipeline executes the same set of instructions over different input data. The software can configure the SIMD pipeline to depend on a set of input streams and produce a set of output streams. There are three types of streams—single-bit control streams, 32-bit word scalar streams, and 16-word vector streams—corresponding to three types of global networks. As the control stream does not have any payload, it is efficiently implemented with an up-down counter (UDC), whose value is incremented when the token arrives and decremented when the token is consumed. The UDC can produce a similar valid (not empty) and ready (not full) interface as other data FIFOs. Execution of the PCU is triggered by the arrival of its input dependencies and back pressured by the downstream buffers. PCU also contains configurable counters, which can be chained to produce the values of nested loop iterators used in the datapath.

We refer to the program graph that can be executed by the SIMD pipeline as a **compute context** or simply **context**. A context includes the branch-free instructions mapped across SIMD stages, the input and output streams, and associated counter states and control configurations. Figure 2.3 (a) shows an example of a simplified pseudo assembly code to program a PCU context. The control signals of the configurable counters, such as counter saturation or **counter done** signals, can be used to dequeue and enqueue the input and output streams, respectively. The counter bounds (i.e., min, max, and stride) can also be data-dependent using values from the scalar input streams. Compared to other dataflow architectures, the dataflow engine in Plasticine is more flexible in that it allows dynamic enqueue and dequeue window for its input and output streams. *This feature enables Plasticine to support complex control hierarchy, such as non-perfectly nested loops and branch statements, across contexts even though individual contexts can only execute instructions that are branch-free.*

Figure 2.3 (c) represents the effective execution achieved in an imperative-style program. For simplicity, we will use the imperative-style configuration in the later discussion. However, it is important to realize the instructions in the imperative-style configurations are not executed in time,

```

1  # scalar and vector input streams
2  bufferA = VecInSteram()
3  bufferB = ScalInSteram()
4  bufferC = VecInSteram()
5  outputD = VecOutSteram()
6
7  i = Counter(min=0,max=3, stride=1)
8  j = Counter(min=0,max=3, stride=1)
9  chain = CounterChain(i,j)
10
11 a = bufferA.deq(when=j.valid())
12 b = bufferB.deq(when=i.done())
13 c = bufferC.deq(when=j.done())
14
15 forward(stage=0, dst="PR0", src=c)
16 # b is broadcasted to all lanes
17 stage(0, "mul", dst="PR1", b, c)
18 # Operands are PRs from the previous stage
19 stage(1, "add", dst="PR2", oprd=["PR0", "PR1"])
20
21 # forwarding PR2 of stage 1 to stage 2
22 forward(stage=2, dst="PR2", src="PR2")
23 forward(stage=3, dst="PR2", src="PR2")
24 forward(stage=4, dst="PR2", src="PR2")
25 forward(stage=5, dst="PR2", src="PR2")
26 forward(stage=0, dst="PR3", src=j.valid())
27 forward(stage=1, dst="PR3", src="PR3")
28 forward(stage=2, dst="PR3", src="PR3")
29 forward(stage=3, dst="PR3", src="PR3")
30 forward(stage=4, dst="PR3", src="PR3")
31 forward(stage=5, dst="PR3", src="PR3")
32 outputD.enq(data="PR2", when="PR3")

```

(a) Declarative configuration

```

1  bufferA = VecSteram()
2  bufferB = ScalarStream()
3  bufferC = VecSteram()
4  outputD = VecSteram()
5
6  with Context() as ctx:
7      i = Counter(min=0,max=3, stride=1)
8      j = Counter(min=0,max=3, stride=1)
9      ctx.chain(i, j)
10     b = bufferB.deq(when=i.done())
11     c = bufferC.deq(when=j.done())
12     a = bufferA.deq(when=j.valid())
13
14     expr = a * b + c
15     outputD.enq(data=expr, when=j.valid())

```

(b) Declarative configuration with automatic register allocation

```

1  bufferA = VecSteram()
2  bufferB = ScalarStream()
3  bufferC = VecSteram()
4  outputD = VecSteram()
5
6  with Context() as ctx:
7      b = bufferB.deq()
8      for i in range(0, 3, 1)
9          c = bufferC.deq()
10         for j in range(0, 3, 1)
11             a = bufferA.deq()
12             expr = a * b + c
13             outputD.enq(expr)

```

(c) Equivalent imperative program

Figure 2.3: Pseudo PCU configuration. (a) shows the simplified declarative-style configuration for the SIMD pipeline context in a PCU. Each PCU context can produce a set of output streams as a function of input streams and counter values. Each stage contains multiple pipeline registers (PRs) that can propagate results across stages. A stage can read PRs from the previous stage and write to PRs in its current stage. Only the first stage can read streams and counter values, and only the last stage can write to streams. Other stages need to propagate the required values through PRs. (b) shows a simplified configuration where registers are implicitly allocated. The `valid` signal of the inner most counter `j` is high whenever the context is enabled. The context is implicitly triggered whenever its input streams `streamA`, `streamB`, and `streamC` are non-empty and output stream `outputD` is ready. By configuring when each stream is enqueued and dequeued using signals from the chained counters, these streams are effectively read and written within different loop bodies. (c) shows the effective execution of the declarative configuration achieved in an imperative program. We move the definitions of streams outside of the context in (b) and (c), so they can be written by other contexts. Each stream can have exactly one writer. If a stream has more than one reader, effectively the writer broadcasts the result to both input buffers of the receiver contexts.

but rather in space. There is a one-to-one translation from the declarative-style configuration to the imperative-style, but not in a reverse way. For example, in the imperative-style, instructions in the contexts must belong to a single basic block, and the loops need to be perfectly nested (only accesses to streams can occur in the outer loops).

There is no global scheduler that orchestrates the execution order among contexts—the execution is purely streaming and dataflow driven. The only way to order the execution of two independent contexts is to introduce a control **token** between two contexts acting like a dummy data-dependency. This restriction eliminates the possible long-traveling wires and communication hot spots caused by a centralized scheduler, which again improves the clock frequency and scalability of the architecture.

Pattern Memory Unit (PMU)

PMUs hold all distributed scratchpads available on-chip. Each PMU contains 16 SRAM banks with 32-word access granularity. The PMU also contains pipeline stages specialized for address computation. Unlike SIMD pipelines in PCUs, these pipeline stages are non-vectorized and can only perform integer arithmetics. The address produced by the address pipeline is broadcasted to 16 banks with a configurable offset added to each bank. In contrast to the PCU SIMD pipeline that has to be programmed atomically, the address pipeline stages within PMUs can be sliced into a write and a read context, triggered independently. In addition to compute stages, resources such as I/O ports, buffers, and counters are also shared across contexts. It is the software's responsibility to make sure the total resources consumed by all contexts do not exceed the resource limits of the PMU. All contexts within PMU have access to the scratchpads with unprotected order. To restrict a read context to access the scratchpad after the write context, the software must explicitly allocate a token from the write context to the read context. SARA automatically generates required control tokens across contexts such that the memory access order is consistent with the program order from a high-level imperative programming language.

Unlike most accelerators at this scale, Plasticine does not have any shared global on-chip memory. This design dramatically improves the memory density and scalability of the architecture by eliminating communication hotspot and bandwidth limitation of a shared memory. On the other hand, however, the burden of maintaining a consistent view of a logical memory mapped across distributed scratchpads is left to the software. The software must explicitly configure synchronizations across PCUs and PMUs, taking into account that the network can introduce unpredictable latency

on both control and data paths. One of the major contributions of SARA is to hide this synchronization burden from the programmer and still provide a programming abstraction of logical memories with configurable bandwidth and arbitrary capacity that fits on-chip.

DRAM Interface

The Plasticine architecture provides access to four DDR channels on the two sides of the tile array, as shown in Figure 2.2. Each side has a column of DRAM address generator (AG) specialized in generating off-chip requests. Like compute pipeline in PMUs, the pipeline in AG is also non-vectorized with integer arithmetics. Each AG can generate a load or a store request streams to the off-chip memory. All streams can access the entire address space of the DRAM, with in-order responses within each stream and no ordering guaranteed across streams. In streaming pipelined execution, the program can use multiple streams for DRAM accesses appeared in different locations of the program. To provide off-chip memory consistency, SARA allocates synchronizations across different streams to preserve memory order expected by the program.

2.3 The High-Level Imperative Compiler

We use Spatial [73]—a domain-specific language for reconfigurable accelerators—as the front-end of Plasticine. Spatial describes applications with imperative control constructs, such as loops and branches, augmented with parallel patterns [74]. Parallel patterns are transformation functions on memory collections that capture both access patterns and parallelization scheme of the operator. Examples of popular parallel patterns include *map* and *reduce*. Instead of using software data-structures, Spatial exposes hardware memories available on reconfigurable hardware, such as registers and SRAM, directly to programmers. These memory types allow a user to explicitly control data transfer between different levels of memory hierarchies to maximize locality. Additionally, Spatial’s language constructs include important design parameters that are essential for achieving good performance on a spatial architecture. Parameters include blocking size, loop unrolling factors, and pipelining schemes, making it easy to perform application-level design space exploration. To enable loop-level parallelization and pipelining, Spatial automatically partitions and buffers the intermediate on-chip memories. An example of outer product—element-wise multiplication of two vectors resulting in a matrix—in Spatial is shown in Figure 2.4.

```

1 // Host to accelerator register for scalar input with
2 // user annotated value for static analysis
3 val N = ArgIn[Int];
4 bound(N) = 1024
5 // 1-D DRAM size in N
6 val vecA, vecB = DRAM[T](N)
7 // 2-D DRAM size in NxN
8 val matC = DRAM[T](N, N)
9 // Loop unrolling factors
10 val op1, op2, ip: Int = ...
11 // Blocking sizes of vecA and vecB
12 val tsA, tsB: Int = ...
13
14 // Accelerator kernel
15 Accel {
16   // Parallelized by op1
17   Foreach(min=0, step=tsA, max=N, par=op1) { i =>
18     // Allocate 1-D scratchpad size in tsA
19     val tileA = SRAM[T](tsA)
20     // Load range i to i+tsA of vectorA from off- to
21     // on-chip parallelized by ip
22     tileA load vecA(i::i+tsA par ip)
23     Foreach(min=0, step=tsB, max=N, par=op2) { j =>
24       val tileB = SRAM[T](tsB)
25       tileB load vecB(j::j+tsB par ip)
26       // 2-D scratchpad
27       val tileC = SRAM[T](tsA, tsB)
28       Foreach(min=0, step=1, max=tsA) { ii =>
29         Foreach(min=0, step=1, max=tsB, par=ip) { jj =>
30           tileC(ii, jj) = tileA(ii) * tileB(jj)
31         }
32       }
33       // Store partial results to DRAM
34       matC(i::i+tsA, j::j+tsB par ip) store tileC
35     }
36   }
37 }

```

Figure 2.4: Example of Outer Product in Spatial.



Figure 2.5: Spatial compiler stack to target FPGAs and Plasticine

Spatial is a target-agnostic language for general reconfigurable architectures. The primary targets of Spatial are FPGAs. Similar to the C-based high-level synthesis languages, such as Vivado HLS [75] and SDAccel [76], Spatial provides a high-level programming interface that focuses on the algorithmic implementations of the application, hiding the low-level hardware interfaces and RTL programming from the users. To target Plasticine, we take applications described in Spatial, disabling FPGA-specific transformations, and perform architecture-specific lowering to SARA’s IR. The key transformations we take from the Spatial compiler is loop unrolling, memory buffering, and memory partitioning (explained later in Section 3.3.2). Optimizations, such as retiming and scheduling, are disabled for Plasticine, as they cannot be directly applied. Figure 2.5 summarizes the compiler flow to target FPGAs and Plasticine from Spatial.

Although SARA takes Spatial as the front-end language, SARA can be equally integrated with other imperative languages with similar control constructs, such as C-based high-level synthesis language, the back-end of the Halide IR [68], the TACO [64] compiler, etc. Nonetheless, using Spatial as our front-end has a few advantages. Rather than adapting existing languages for processor architectures like most HLS languages, Spatial is designed specifically for reconfigurable spatial architectures. Spatial’s language constructs capture the scheduling scheme that can be explored by most spatial architectures, such as coarse-grained pipelining, streaming dataflow, hierarchical parallelization, and finite state machines (FSM). These language constructs are missing from a processor-based language as they cannot be supported by a processor. Spatial IR also differs in its representation of control flow and memory constructs, which are more suitable for analyses of spatial architectures.

Most compilers, such as LLVM [77], use control flow graphs to represent the control constructs in an imperative program. The control flow graph implicitly assumes the program is executed in time, which makes it unsuitable for analysis of a reconfigurable architecture that executes the program in space. Instead, Spatial uses a control hierarchy in the IR to capture the scheduling of the program. The control hierarchy uses a tree structure to represent the nested control constructs, making it easier to analyze the relations between controllers. Figure 2.6 (b) and (c) shows an example of a control flow graph vs. a control hierarchy. The controller at each level of the hierarchy corresponds to a control construct, such as a loop, or a branch statement. A basic block is attached to each *innermost* controller including instructions and memory accesses. If the program has instructions in an outer loop, Spatial automatically inserts a *unit* controller to wrap the floating instructions. SARA takes the back-end of Spatial IR as input, which is a control hierarchy after loop unrolling shown in Figure 2.6

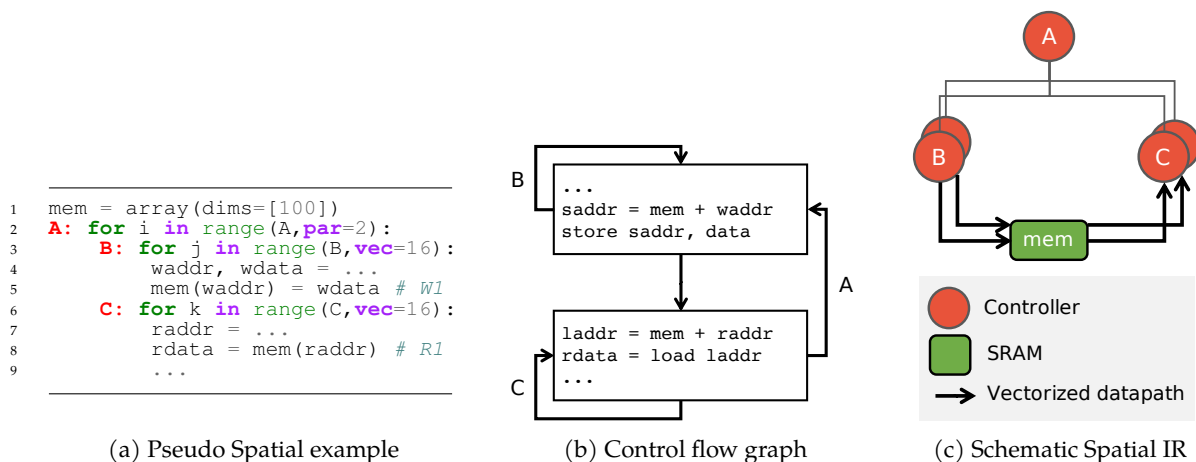


Figure 2.6: Pseudo example of SARA’s front-end language. (a) shows the an example of SARA’s front-end language. (b) shows the control-flow graph in traditional CPU-based compiler. The ‘par’ keyword indicates outer loop unrolling factor, and the ‘vec’ keyword is followed by a inner-loop vectorization factor. When an iterator is vectorized, instructions using the vectorized iterator is automatically vectorized. When unrolling the outer loop A, the enclosed loop body and next-level controllers are duplicated, as suggested in (c). Each loop in (a) corresponds to a controller in (c). The inner most controllers B and C each contain a basic block within instructions within the inner most loops.

(c).

The representation of data structures in traditional IRs often marries to the memory model of a CPU. Traditional compilers, such as LLVM [77], treat data collections as pointers to a shared global address space because of the virtual memory abstraction in CPUs. The hardware, on the back-end, implicitly manage the data movement between an on-chip cache and the off-chip memory; this scheme improves programmability of CPU at the cost of hardware complexity and unpredictable memory performance. Accelerators on the other hand often have explicitly managed on-chip scratchpads. Therefore, modeling data structures as disjoint memory spaces as supposed to pointers is more suitable for reconfigurable architectures.

Spatial is an embedded DSL in Scala. For simplicity and generality, we will use python-style pseudo code to represent the imperative input program for the rest of our discussion.

Chapter 3

Compiler

3.1 SARA Compiler Overview

In this section, we introduce the compiler framework—SARA—that targets Plasticine architecture from high-level programs described in the Spatial language. There are two challenges to map Spatial applications to Plasticine.

First, unlike an FPGA, Plasticine cannot map arbitrary RTL functionality. In the Spatial abstraction, the execution order of the program is organized by a control hierarchy, where each level of the controller schedules the execution of the next level controllers. When mapping the example in Figure 2.6 onto an FPGA, the outer controller *A* sends an enable signal to each child controller, which signals back the parent controller when completed. If the user chooses to sequentially execute the outer loop *A*, the parent controller enables the child controllers one at a time; if the user chooses to metapipeline (coarse-grain pipeline) the outer loop *A*, the outer controllers enables multiple child controllers in a pipelined fashion.

To achieve the same execution schedule on Plasticine in a naïve approach, we can map each controller in the hierarchy into a PU, sending control signals to schedule the next level controllers distributed among other PUs, as shown in Figure 3.1 (a). This strategy suffers from the expensive network round-trip delays between the parent and child controllers. To be scalable at a high clock frequency, Plasticine networks are pipelined at each switch, introducing multiple cycles of network delay across PUs on the control path. Therefore, the multi-cycle handshaking signals can lead to significant pipeline bubbles that undermine performance. Additionally, this scheme creates a communication hotspot around the parent controller *A* as loop *A* gets unrolled, which is devastating



Figure 3.1: (a) A naïve mapping strategy to map the control hierarchy onto Plasticine. All units are distributed across an on-chip network that can introduce unpredictable latency. The number on the edges indicates event order. Here we show a scenario where read requests from C do not observe the write requests from B that occur earlier in the program order due to network latency between B and mem. (b) Distributed controllers in SARA. Each innermost controller makes a copy of all enclosing controllers. The signals from these controllers are used to generate synchronization between distributed compute units. To address the problem in (a), the memory also needs to provide a write acknowledgment per write request for synchronization.

for a CGRA like Plasticine that has much less routing resource than an FPGA. Furthermore, synchronizing the compute only is insufficient to ensure memory effects are observed by the remotely distributed accessors, as shown in Figure 3.1 (a).

To address this challenge, we want to eliminate centralized outer controllers. At high-level, SARA performs **loop divisions** on each outer controller, such that all innermost controllers are perfectly nested, as shown in Figure 3.1 (b). Section 3.2.1 discusses the loop division transformation. SARA then allocates synchronization tokens across the distributed innermost controllers. All innermost controllers have their own copies of the outer controllers, which are used to generate the control tokens. These control tokens ensure the execution order of the inner controllers is the same as if they are scheduled by a centralized outer controller. Instead of synchronizing all inner controllers under an outer controller, SARA only synchronizes the ones accessing the same memory, such as B and C in Figure 2.6. This limits the synchronization among a small set of distributed nodes without impacting the result, making our design much more scalable.

The second challenge is that controllers in the Spatial hierarchy can consume an arbitrary amount of compute and memory resources, exceeding the capacity of individual PUs. For instance, a user might write an on-chip memory multiple times throughout the program with writers mapped to different PUs. The physical scratchpad, however, only has a single write port and write address

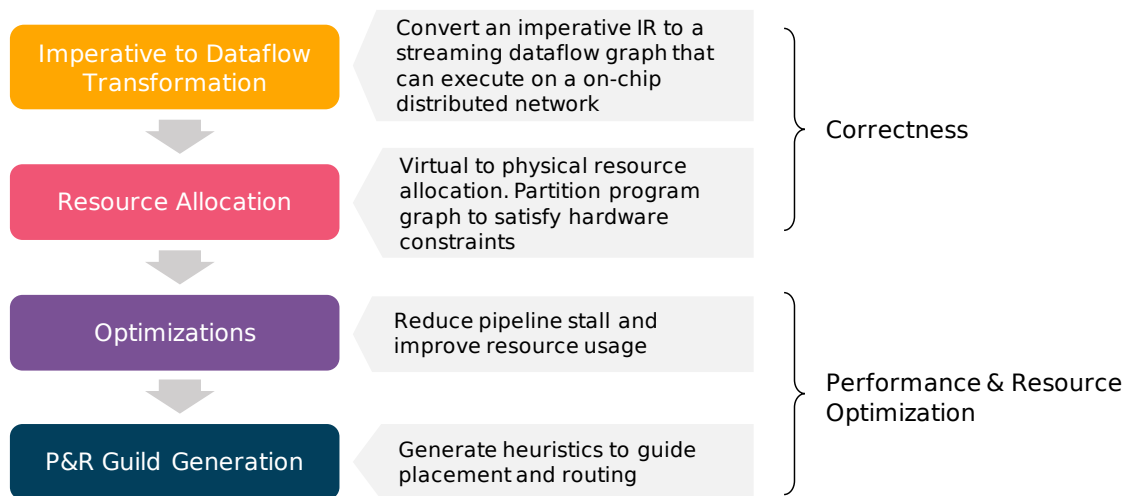


Figure 3.2: SARA Compiler Flow

pipeline. SARA maps a subset of the address computation externally in other PUs and generates control tokens among the distributed writers, time-sharing the write port, and enforcing the ordering between distributed writers. Using resource virtualization, SARA composes or time-share the physical resources when software usage exceeding the hardware limit.

In the following sections, we describe a systematic approach to compile applications described in an imperative front-end language to a purely declarative and distributed dataflow graph that can run on Plasticine. Figure 3.2 shows SARA’s compilation flow.

Section 3.2 expands on the **imperative to dataflow transformation** that addresses the first challenge. SARA allocates distributed on-chip resources to execute the program in spatially parallelized and pipelined fashion with appropriate synchronizations. A virtual unit (VU) is our intermediate representation that captures the computation mapped within the boundary of a physical unit (PU), such as a PCU and PMU. Each VU can contain multiple contexts if their aggregated resource usage can fit in a PU. The hardware can limit the maximum number of contexts a PU can support and has resources that cannot be split across contexts. Most importantly, SARA needs to ensure messages across VUs, mapped across the global network, tolerate an arbitrary amount of network latency; messages within a single VU across contexts takes only a single cycle. The transformation phase generates a virtual unit dataflow graph (VUDFG) with appropriate synchronizations, such that the parallelized and pipelined program executed over distributed on-chip resources produces the same result as a parallelized program executed in time. At the end of the allocation phase, a virtual unit can consume as many resources as the program requires.

SARA further virtualizes resource allocation and hides the underlying resource constraints from the programmers. Section 3.3 dives into the **resource allocation** phase, where SARA assigns each VU to a PU that processes the required resources. If no PU can fit a VU, SARA partitions the VU into multiple VUs to resolve the constraint violation. If there is insufficient PU or the VU cannot be partitioned, the mapping process fails with appropriate hints to the programmer for the limiting resources.

Throughout the first two phases, SARA introduces various **optimizations** that either reduce the resource cost of the VUDFG, or alleviate potential performance bottleneck in the streaming pipeline. After all VU fits in at least one type of PU, SARA performs a global optimization that merges small VUs into a larger VU to reduce resource fragmentations. Section 3.4 enumerates the optimizations SARA performs.

The output of the resource allocation phase is a VUDFG with a tagged PU type for each VU. It is up to the **placement and routing (PaR)** phase to determine where the VU will be finally placed. Right before PaR, SARA performs static analysis on the traffic pattern and generate heuristic guild for the placer to reduce routing congestion. Section 3.5 details the PaR algorithm and heuristic-guild generated by SARA.

3.2 Imperative to Dataflow Transformation

3.2.1 Loop Division

Between the front-end and the back-end abstractions of SARA, an obvious gap is that the imperative front-end language can contain arbitrarily nested control hierarchy, whereas the hardware compute engine can only execute operations that are control-free. To address this issue, we introduce a new type of transformation—loop division—for streaming reconfigurable accelerators. Similar to loop fission, loop division breaks a single loop into multiple loops. The difference is that loop fission generates a sequence of sequentially executed loops, whereas loop division generates loops executing *concurrently*. Additionally, loop fission materializes the intermediate results across fissioned loops into arrays, while loop division use queue to communicate across loops. Each loop generated from loop division can only execute if all of their input queues are not empty. Figure 3.3 gives an example of a loop fusion vs. loop division.

<pre> 1 mem = zeros(N) 2 a = rand(N) 3 b = rand(N) 4 for i in range(1, N): 5 tmp = a[i] * b[i] 6 mem[i] = tmp + i </pre>	<pre> 1 mem = zeros(N) 2 a = rand(N) 3 b = rand(N) 4 tmp = zeros(N) 5 for i in range(1, N): 6 tmp[i] = a[i] * b[i] 7 for i in range(1, N): 8 mem[i] = tmp[i] + i </pre>	<pre> 1 mem = zeros(N) 2 a = rand(N) 3 b = rand(N) 4 tmp = queue() 5 @concurrent 6 for i in range(1, N): 7 tmp.enq(a[i] * b[i]) 8 @concurrent 9 for i in range(1, N): 10 mem[i] = tmp.deq() + i </pre>
(a) Input program	(b) Loop Fission	(c) Loop Division

Figure 3.3: (b) and (c) shows the output of loop fission and loop division of the input program (a), respectively. In (b), the first loop is executed entirely before executing the second loop. The intermediate result `tmp` is materialized into an array with the same size as the loop range. In (b), the two loops can execute concurrently. The intermediate result is materialized into a queue. For each iteration, a loop can execute only if all of its queues are non-empty. The second loop can execute as soon as `tmp` receives the first element.

When executing loop division on a single-threaded CPU, the CPU must context switching between the concurrent loops and executing the one with cleared input dependencies. Like loop fission, loop division is likely worsening the performance on a processor architecture, as the worst-case memory footprint of the intermediate result `tmp` increases from $O(1)$ to $O(N)$. On RDAs, the divided loops are executing concurrently in a streaming pipelined fashion. The size of the `tmp` can be limit to a small fixed size, efficiently implemented with a hardware FIFO. Although loop transformations are generally optimizations on CPUs, loop division is a required transformation to converts an infeasible program to a feasible one for Plasticine.

Loop fission is not always safe, as it may alter the execution order of the program. Loop division, on the other hand, does not change the underlying data-dependency and is always safe. To achieve this, loop division needs to introduce additional dummy data dependencies across divided loops to enforce the correct execution order. Figure 3.4 gives an example of an invalid loop fission and a correct loop division. Section 3.2.3 gives more detail on how SARA automatically generates the dummy data-dependencies to preserve program order.

3.2.2 Virtual Context Allocation

At high-level, SARA executes the entire program in a pipelined fashion, where each basic block in the program is a pipeline stage. As a start, SARA performs loop division at each level of the control hierarchy, such that all innermost basic blocks are perfectly nested. SARA then maps each basic block to a context and each data-structure to a virtual memory, as shown in Figure 3.5 (b).

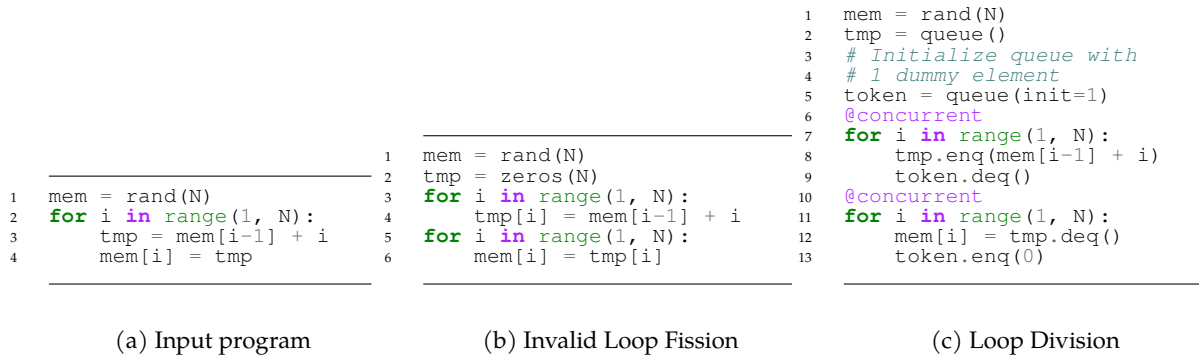


Figure 3.4: Example of an illegal loop fission and a legal loop division

A basic block maps naturally to a context, as instructions within a basic block are control-free. All ancestor controllers enclosing the basic block are duplicated into the corresponding context. With these controllers, contexts can repeatedly execute their basic blocks for expected number iterations. For each virtual memory, SARA examines all contexts reading or writing the memory, allocating synchronization tokens to enforce their access order. The token can be viewed as an access grant of the memory. By controlling how tokens are passed between the pipelined contexts, SARA can enforce that the memory is read and written with the expected ordering of an imperative program. *As long as each individual memory is accessed with a program-consistent order, the final result is identical to a sequentially executed program.* This way, SARA introduces minimum point-to-point synchronizations among small groups of contexts; contexts accessing *different* data-structures are naturally parallelized without impacting the final output. Section 3.2.3 explains how SARA allocates the synchronization tokens based on the control hierarchy of the imperative program.

Unlike traditional out-of-order execution, where both static scheduling in the software and dynamic scheduling in the hardware search for independent instructions to execute concurrently, SARA starts with executing *all* basic blocks of a program concurrently, introducing minimum synchronizations wherever necessary. Unsurprisingly, a control-heavy program with many basic blocks can easily run out of PUs. Nonetheless, data-analytic programs, which are the intended workloads for RDAs, have relatively simple control flow and abundant data-level parallelism (i.e., few basic blocks with high repetition counts). Supporting flexible control constructs, however, allows an RDA to accelerate a program with minimum host intervention, eliminating communication overheads and maximizing FLOPS utilization especially for a large-scale RDA. To map a large program that runs out of resources on Plasticine, the program graph must be sliced into multiple chunks. A single chunk is executed in-space, exploring on-chip parallelism and pipelining, while different

```

1 mem = array(dims=[100])
2 A: for i in range(A):
3     B: for j in range(B):
4         waddr, wdata = ...
5         mem(waddr) = wdata # W1
6     C: for k in range(C):
7         raddr = ...
8         rdata = mem(raddr) # R1
9         ...

```

(a) Pseudo input example

```

1 mem = VirtualMemory(dims=[100])
2
3 with Context() as ctxB:
4     A: for i in range(A):
5         B: for j in range(B):
6             waddr, wdata = ...
7             mem(waddr) = wdata # W1
8
9 with Context() as ctxC:
10    A: for i in range(A):
11        C: for k in range(C):
12            raddr = ...
13            rdata = mem(raddr) # R1
14            ...

```

(b) Initial context allocation

```

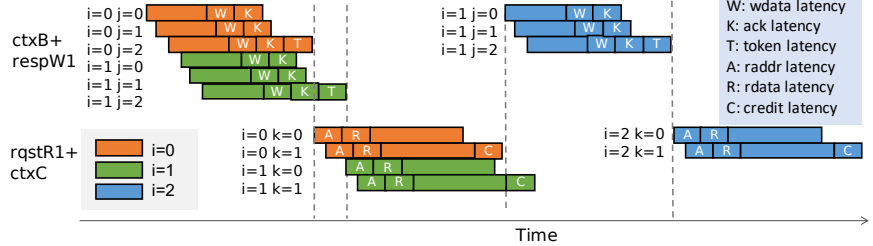
1 mem = VirtualMemory(dims=[100])
2 token = ControlStream()
3 credit = ControlStream(init=2)
4
5 with Context() as ctxB:
6     A: for i in range(A):
7         credit.deq()
8         B: for j in range(B):
9             waddr, wdata = ...
10            mem.waddr.enq(waddr)
11            mem.wdata.enq(wdata)
12 with Context() as respW1:
13     A: for i in range(A):
14         B: for j in range(B):
15             ack1 = mem.ack.deq()
16             token.enq()
17
18 with Context() as rqstR1:
19     A: for i in range(A):
20         token.deq()
21         C: for k in range(C):
22             raddr = ...
23             mem.raddr.enq(raddr)
24 with Context() as ctxC:
25     A: for i in range(A):
26         C: for k in range(C):
27             rdata = mem.rdata.deq()
28             ...
29             credit.enq()

```

(c) Request and response division



(d) Context Graph



(e) Timing on Plasticine

Figure 3.5: Context lowering and control allocation example. SARA allocates one context per basic block for B and C, shown in (b). Outer controller A is duplicated in both **ctxB** and **ctxC**. (c) SARA separates out a requesting context **rqstR1** from **ctxB** for R1 and a receiving context **respW1** from **ctxC** for W1. The resulting dataflow graph is shown in (d). To enforce the forward data-dependency between W1 and R1, SARA allocates a forward token between W1's receiving context **respW1** and R1's requesting context **rqstR1**; to enforce the loop-carried WAR dependency between R1 and W1, SARA allocates a backward token **credit** between R1's receiving context **ctxC** and W1's requesting context **ctxB**. The backward token is initialized with two elements because **mem** is double-buffered, enabling the writer for two iterations of A before back-pressured. For the forward token, the LCA controller between W1 and R1 is A. The immediate child of the LCA controller in ancestor controllers of W1 is B, therefore, the enqueue enable of the token is configured to **B.done** in **respW1**. Similarly on the receiving side, the dequeue enable of the token is **C.done** in **rqstR1**. The resulting timing of the execution is shown in (e).



Figure 3.6: Specialization for dense memory lowering if the if the memory *mem* in Figure 3.5 (a) is a dense scratchpad. (a) and (c) shows the allocated contexts. Context *rqstW1* and *rqstR1* are allocated in the same VU as the memory. The latency of the forward token and the backward credit take only one cycle, which are not shown in (d). The latency to compute read and write addresses are also not shown. For complex address calculation, SARA split the address computations to separate contexts shown in (b). The addresses can be precomputed in *rqstW1* and *rqstR1* and are back-pressured by *rqstW1'* and *rqstR1'*, respectively. This way, the latency of address calculation does not create pipeline bubble even if token and credit are frequently exchanged (e.g. when loop B and C runs a small number of iterations). *rqstW1* and *rqstR1* can be mapped outside of the memory VU if running out of address stages.

chunks are executed in-time by reconfiguring the accelerator.

3.2.3 Control Allocation

SARA uses control tokens to enforce the memory access order in an imperative program, which is the minimum synchronization necessary to produce a correct execution. In this section, we first explain *how* SARA uses tokens to enforce the memory access order between distributed contexts, achieving complex control flows on a data-flow architecture; next, we show *where* SARA inserts these tokens and how to minimize inserted tokens while enforcing the program order.

How. The first question is how do we enforce the access order of a memory from two remotely distributed contexts, given the global network and the memory itself can introduce unpredictable delays. The streaming memory interface takes a stream of read or write requests, providing a response packet for each request packet. Within a single stream, the requests are guaranteed to be served in order. Therefore, the read responses are a stream of data packets (not tagged with requested address) and the write responses are a stream of control tokens without any payload. To eliminate the round-trip latency, SARA divides each memory access in the input IR into a requesting and a receiving context, as shown in Figure 3.5 (c). For a read access, SARA maps the address generation in a separate requesting context, streaming addresses to the memory and streaming the data to the receiving context. Similarly, for a write access, SARA allocates a receiving context to accumulate the write acknowledgments for synchronizations.

To order an access A before an access B , SARA allocates a token between the *receiving* context of access A ($resp_A$) and the *requesting* context of access B ($rqst_B$). This is essential to ensure the memory effect of access A is visible to access B , as the memory might take multiple cycles to serve a request. For any two accesses enclosed by a loop in the input IR, there could be a loop-carried dependency (LCD) from the later access in the program order to the earlier access. When inserting a token for an LCD, the receiver's token buffer is always initialized with at least one token, enabling the earlier access for the first iteration of the loop. Figure 3.5 (c) shows the forward and the backward token to enforce two accesses under a loop. The initial element in backward token breaks the cyclic dependency. Borrowing the flow-control terminology[78], we refer to a backward token as a credit. If the memory is multi-buffered, the credit buffer is initialized with buffer-depth D number of credit, enabling the earlier access to execute D iterations of the enclosed loop before back pressured by the later access.

SARA wires up the enqueue/dequeue enable of a token with signals from the duplicated controllers in the writer/reader context. For dependencies of a queue, the queue should be written and read every cycle when the producer/receiver basic block is active; to achieve this, SARA configures the enqueue and dequeue signal of the token that enforces the dependency to the *valid* signals of the innermost controllers in the requesting and receiving contexts, respectively. For dependencies of a scalar variable or an array, the program expects the producer and consumer basic blocks to update and inspect the memory per iteration of their LCA controller. For a token that enforces the dependency of two accesses, SARA finds the LCA controller between the two accesses in the control hierarchy. Then within contexts containing the respective accesses, SARA finds the immediate child controllers of the LCA controller and connects its *done* signal to the enqueue and dequeue enable of the allocated token, respectively. This way, a token is produced/consumed for every iteration of the LCA controller using the local controllers within the context containing the dependent/depending access. Figure 3.5 (c) shows an example of how the enqueue and dequeue signals are configured.

By controlling *when* a token is produced and consumed, SARA is able to achieve the complex and nesting control flows in an imperative program on a dataflow accelerator. Section 3.2.4 dives into more complex control constructs with data-dependent control flows.

The requesting and receiving contexts pair are general schemes we use on general streaming memory interface with variable access latency, such as an off-chip DRAM or an on-chip SRAM when supporting sparse accesses. In common cases, we can simplify the synchronization logic for certain on-chip memory types.

Specialization for dense on-chip scratchpad Spatial can statically analyze the dense access pattern of on-chip arrays and partitions the memories to avoid bank conflicts. These memories have a guaranteed single-cycle access latency, which eliminates the need of write acknowledgments. Shown in Figure 3.6 (a), SARA allocates $rqst_A$ and $rqst_B$ contexts *in the same* VB as the accessed memory with corresponding address computation. Instead of synchronizing between $resp_A$ and $rqst_B$, SARA sets up the forward token and the backward credit between $rqst_A$ and $rqst_B$. Because memory requests are served in a single cycle, access A 's requests would be visible to access B by the time $rqst_B$ receives the token. Hence, write acknowledgment is no longer needed. Figure 3.6 shows the resulting contexts if *mem* in Figure 3.5 (a) is a dense on-chip array.

Data structure	Memory type
array (fits on-chip)	SRAM
array (does not fit on-chip)	DRAM
scalar variable	register
queue	FIFO

Table 3.1: Mapping between user declared data-structure to underlying hardware memories. Programmers explicitly specify the desired hardware type inside Spatial. In other languages, this table specifies a mapping between software data-structures and hardware memory types on Plasticine.

Specialization for non-indexable memory For most non-indexable memories like scalar variables and queues in the program, SARA directly maps them to the input buffer of the receiver context if the memory has a single writer and a single reader. Non-indexable memory with multiple readers and writers also need synchronization token across the accessing contexts to enforce a program-consistent accessing order.

Where. SARA can synchronize every pair of accesses of a memory to enforce the program order. However, that would require N^2 tokens for a memory with N accesses in the input IR ($2N^2$ with LDCs), which can be unnecessarily expensive. The second question is how can we minimize the number of inserted token while ensuring a correct execution order. To tackle this problem, SARA builds a dependency graph between all accesses of a memory in the input IR; nodes of the graph are accesses and edges represent dependencies between accesses. The dependency graph is *per data-structure* without unnecessary ordering across different memories. Next, SARA removes edges in the graph, if the removed edge does not relax the ordering across accesses. Lastly, only for each edge within the reduced dependency graph does SARA insert a token between the source and destination of the edge, using mechanism explained in the previous section.

Dependency Graph Construction For every access in the input IR, SARA checks other accesses on the same data-structure occurred earlier in the program order for a possible forward dependency, and later in the program order for a possible loop-carried dependency (LCD). SARA only inserts an edge between two accesses if they potentially interfere, which is a function of

- the type of accesses (read vs. write)
- the type of the memory (e.g. SRAM, DRAM)
- and location of the declared accesses in the control hierarchy.

Memory type	DRAM	SRAM	FIFO	Register
read-after-read (RAR)	✗	✓	✓	✗
read-after-write (RAW)	✓	✓	✓	✓
write-after-read (WAR)	✓	✓	✓	✓
write-after-write (WAW)	✓	✓	✓	✓

Table 3.2: Interference table for whether two accesses of the same memory needs to be synchronized by the software for each memory type. Two DRAM read accesses do not interfere because the DRAM interface permits multiple concurrent access streams through multiple AGs, where each AG can access the full off-chip address space. SRAMs, on the other hand, guarantees in-order response only within a single stream. Therefore, it is the software’s responsibility to ensure the read port receives requests from a single context at any point.

Table 3.1 lists hardware memories available on Plasticine and software data-structures that can be mapped onto these memories. The type of the memories matters because they have different programming interface on the hardware. Table 3.2 shows the interference relations of different memories on Plasticine. Figure 3.7 shows derived dependency graphs of a few example programs.

Dependency graph reduction SARA reduces the dependency edges in two passes, processing edges for the forward and backward dependencies separately.

For forward dependencies, SARA performs a transitive reduction (TR)[79] on the forward dependency graph. TR keeps the minimum of edges in a graph that preserves the connectivity of the original graph, enforcing the same ordering as the original dependency graph.

Next, SARA checks if each backward edge can be removed without breaking the execution order between accesses. Each backward edge represents a loop-carried dependency (LCD) and is tagged with the associated loop. A backward edge can be removed, if after removing the edge, there is a path going from the source of the removed edge to the destination, with any forward edges remained from the TR pass and *a single* backward edge tagged on the same loop as the removed edge. Figure 3.8 demonstrates the dependency reduction with examples.

The forward edges can be reduced with TR because the forward dependencies are monotonic, i.e. A depending on B and B depending on C enforces A depending on C . The backward dependencies are non-monotonic because the initial token in the destination buffer of a backward edge enables the destination of an edge to execute before the source. Therefore, having edges between A and B , and B and C do not enforce A executing before C . Hence, the backward edges cannot be concatenated to maintain dependency.



Figure 3.7: Access dependency graphs for example programs. Solid edges indicate forward dependencies and dashed edges indicate backward loop-carried dependencies (LCDs). Each access in the dependency graph will be mapped to a context executed in space. In (b) and (c), W_0 and W_1 are unrolled accesses for W in lane 0 and 1 of loop A in (a), respectively. (b) is the dependency graph if *mem* is mapped to an on-chip memory and (d) is the dependency graph if *mem* is mapped to an off-chip memory. The on-chip version does not need the cross-lane synchronization because lane 0 and 1 are implicitly ordered when SARA partitions *mem* (discussed later in Section 3.3.2). In (d), there is no forward dependency between W_1 and R_0 because the LCA controller of the two accesses is branch C, which means the two accesses cannot happen simultaneously for the same iteration of loop A . There is, however, loop-carried dependencies between R_0 and W_1 , preventing access in iteration $i + 1$ occur before iteration i of loop A between W_1 and R_0 in two distributed contexts. There is no need to enforce the LCD between an access and itself because the hardware guarantees a single request stream from the same context is served in-order. For each iteration of A , R_0 and W_1 will produce and consume a credit from each other no matter what value the condition takes. Section 3.2.4 will elaborate more on how tokens are used for branches. R_0 and R_1 do not depend on each other because they are both DRAM read accesses as explained in Table 3.2.



Figure 3.8: Dependency graphs (b,e) and reduced dependency graph (c,f) of example programs in (a,d), respectively. Solid edges indicate forward dependencies and dashed edges indicate backward loop-carried dependencies (LCDs). Colors of the LCD edges indicate the associated loop, blue for loop A and red for loop B. For forward edges (black edges), SARA uses transitive reduction (TR) to remove the redundant edges. The outputs of TR are shown in (c,f). For backward edges, an edge can be removed if there is still a path from its source to destination with all forward edges plus a single backward edge of the same loop after the edge is removed. For example from (b) to (c), edge $R2 \rightarrow R1$ is removed because there is path $R2 \rightarrow W1 \rightarrow R1$. Edge $W2 \rightarrow W1$ is removed because there is path $W2 \rightarrow R2 \rightarrow W1$ after the edge is removed in (c). Edge $R2 \rightarrow W1$ cannot be removed because there is no path from $R2$ to $W1$ with a single back edge after the edge is removed. Similarly, $R2 \rightarrow W1$ cannot be removed in (e) because there is no path from $R2$ to $W1$ after the edge is removed.



Figure 3.9: (a) shows an example program with a dynamic loop range. Expressions to generate the loop bound in line 2 of (a) belongs to a basic block that gets mapped to `ctxA` in (b). SARA treats the loop bound as a data-dependency of `ctxB`, and configures the dequeue enable of `bound` stream to counter `B.done` in `ctxB`.

3.2.4 Data-Dependent Control Flow

Using the synchronization discussed in Section 3.2.3, SARA can support control constructs that typically are not supported on a dataflow accelerator, such as branches and while loops. Most accelerators do not support control divergence. SIMT architectures, like GPUs, implement branches with predication and pay the latency penalty of both branch cases. To enable these flexible control constructs, the control path of the architecture must permit data dependencies. ?? details the required changes in Plasticine’s control path to support these features.

Dynamic Loop Range Figure 3.9 shows an example of loops with data-dependent ranges. SARA uses a context to compute the loop bounds, which are treated as input dependency to the context that maps the loop body.

Branch Condition SARA supports both predictions and branches with control divergence. Branches within the innermost loops are implemented with predication so that the innermost loop can still be vectorized across SIMD lanes. This is very similar to a SIMT architecture, where all lanes execute the *if* clause followed by the *else* clause, masking off the memory accesses in the disabled lanes. The total number of SIMD stages required is the total operations in both the *if* and the *else* clauses.

For a branch in an outer loop, the branch condition is treated as a data-dependent enable of a controller under the branch clauses. If the controller is disabled by the branch condition, its *done* signal raises to high immediately. Output tokens depending on the *done* signal will be immediately sent out. This way, the controllers under a branch does not need to execute if the outer branch condition evaluates to false.

```

1  A: for i in range(A):
2      B: even = i % 2 == 0
3      C: if even:
4          D: for j in range(D):
5              addr1, data1 = ...
6              mem(addr1) = data1 # W
7      else:
8          F: for k in range(F):
9              addr2 = ...
10             data2 = mem(addr2) # R
11

```

(a) Input program



(b) Dataflow graph

```

1  even = ScalarStream()
2  odd = ScalarStream()
3  token = ControlStream()
4  credit = ControlStream(init=2)
5  with Context() as ctxB:
6      A: for i in range(A):
7          B: tmp = i % 2 == 0
8          even.enq(tmp)
9          odd.enq(not tmp)
10 with Context() as ctxD:
11     A: for i in range(A):
12         C: if even.deq():
13             D: for j in range(D):
14                 addr1, data1 = ...
15                 mem.waddr.enq(data1)
16                 mem.wdata.enq(data1)
17                 credit.deq()
18 with Context() as rpstW:
19     A: for i in range(A):
20         C: if even.deq():
21             D: for j in range(D):
22                 mem.wack.deq()
23                 token.enq()
24 with Context() as rqstR:
25     A: for i in range(A):
26         C: if odd.deq():
27             F: for k in range(F):
28                 addr2 = ...
29                 mem.raddr.deq(addr2)
30                 token.deq()
31 with Context() as ctxF:
32     A: for i in range(A):
33         C: if odd.deq():
34             F: for k in range(F):
35                 data2 = mem.rdata.deq()
36                 ...
37                 credit.enq()

```

(c) Context configuration



(d) Timing

Figure 3.10: An example program with a outer branch statement. The LCA controller of the write access W and the read access R is branch C. Therefore, the enqueue enable of *token* and dequeue enable of *credit* is connected to *D.done*, and the dequeue enable of *token* and enqueue enable of *credit* is connected to *F.done* in their respective contexts. If *C*'s condition evaluates to *true*, *D.done* raises to high without executing iterations of loop *D*. Similarly for *F.done* if *C*'s condition evaluates to *false*. (d) shows the timing of execution with $A = 4$, $D = 3$, and $F = 2$.

Figure 3.10 shows an example with an outer branch statement. The input program in (a) writes and reads the memory *mem* on even and odd cycles of the outer loop *A*, respectively. SARA maps the three basic blocks in the program in three contexts, *ctxB*, *ctxD*, and *ctxF*, as shown in (b) and (c). For the read and write accesses, SARA allocates a context *respW* to accumulate write acknowledgments and a context *rqstR* to generate read requests. *ctxB* generates the branch conditions for both the *if* and the *else* clauses. The conditions are sent as data dependencies to contexts mapping the basic blocks under the branch conditions.

Figure 3.10 (d) shows the timing of execution. As *ctxB* has no dependencies, *ctxB* computes the conditions for all iterations of *A* in pipelined fashion and broadcasts the conditions to the receiver contexts. The *if* contexts (*ctxD+respD*) receives a `true` condition for the first iteration of *A* (A_0), executing all iterations of *D*, and passes a token to the *else* (*rqstR+ctxF*) contexts. Because *mem* is double-buffered, the credit is initialized with two elements in the receiver's input buffer, enabling the *if* contexts to execute two iterations of outer loop *A* before waiting for the *else* contexts. For the second iteration of loop *A* (A_1), the *if* contexts receives a `false` condition for branch *C*. Therefore, the enclosing loop controller *D*'s *done* signal is immediately high, sending the token to *rqstR* right away. On the other side, the A_0 of the *else* contexts is blocked by the token from the *if* contexts. As soon as the token arrives, the *else* contexts send out the credit immediately without executing loop *F* as the *else* clause is disabled for A_0 . Next, the *else* contexts check for the second token, and executes loop *F* for A_1 .

Both the *if* and the *else* contexts only execute if the enables of their enclosed branch clauses are evaluated to be true. More interestingly, Figure 3.10 (d) shows a overlapping execution of the *if* and *else* clauses across iterations of *A*. The contexts allocated for both *if* and *else* clauses are active almost all time if loop *D* and *F* runs for large number of iterations. If the latency of loop *D* and *F* are both L and the branch has 50% active for both clauses, the total runtime for N iterations of loop *A* would be on the order of $\frac{N}{2}L$ for Plasticine, assuming L and N are large. This is because the contexts mapping both clauses only executes if their clauses are enabled, which corresponds to $\frac{N}{2}$ iterations of the outer loop *A*. *The runtime is almost twice as fast as a traditional coarse-grained pipeline with hierarchical FMS schedulers, like Spatial's FPGA back-end, which runs in NL .*

Do While Loops The *do while* construct is useful to express iterative convergence algorithms or handle an external data stream with a last-bit signal that terminates the execution of the accelerator. A *do while* loop works very similarly to a loop with dynamic range, except the do-while loop has a much longer initiation interval (II) [80]. The earliest starting time for the second iteration of



Figure 3.11: An example program for a do while loop. Loop *A* in (a) is a loop running forever, with a conditional *stop* variable. Loop *A* reads and block *C* writes a value of *stop* for every iteration of loop *A*. Figure 3.11 (b) shows the dataflow graph, mapping each basic block to a context. *ctxC* computes the stop condition, which is used by both *ctxC* and *ctxB*. The break statements in (c) correspond to configuring the stop signal of the enclosing controller *A* to the output of the stop input buffer. The stop buffer is dequeued in every iteration of loop *A*. Due to LCD, the *stop* stream is initialized with `False` to enable the execution of the first iteration of *A*. The *accum* variable in (a) has two readers in block *B* and block *D*; one has a loop-carried dependency with the writer (line 9) and the other has a forward data-dependency (line 11). Therefore, the variable is mapped to two copies of *accum*, one for each reader. Because the accumulate operation is a single operation `ADD`, the accumulator gets optimized to the special accumulate pipeline register *accum1*. Without this optimization, *accum1* would be another `ScalarStream` with an initial element 0. The other *accum* becomes a scalar stream *accum2*, written by *ctxB* when loop *A* is done. *accum1*'s reset signal is connected to *A.done*. The reset signal is determined by where *accum* is declared in (a), which in this case is outside of loop *A*.

the loop is when the while condition is resolved. The while condition is a data-dependency to all contexts with basic blocks enclosed by the do-while loop. There is a loop-carried dependency between the producer of the condition within the while loop body and the do-while loop controller that consumes a value of the condition for every loop iteration. Figure 3.11 shows an example of a do-while loop.

(d) shows the timing of execution. As we can see, the (II) of loop *A* in the steady-state is bounded by the latency to evaluate *stop*, which is the number of operations in *stop* expression rounded up to a multiple of 6¹.

Procedure Calls Our front-end language Spatial currently does not capture procedure calls; all functions are inlined in the IR and do not share resources. In future work, we can extend SARA’s token mechanism to handle the procedure calls, similar to how SARA handles memory consistency. The reused hardware in a function call can be viewed as a shared resource like memories; contexts containing callers of the function must pass around tokens to serialize their accessing order based on the program order of the callers; the order only matters if the function call is stateful and have side effects.

3.2.5 Virtual Unit Allocation

After all contexts and shared memory are allocated and synchronized, SARA moves each floating context into a VU, which later maps to a PU. In the resource allocation phase, SARA partitions the big contexts into multiple VUs, and merges small contexts into a single VU. Contexts created with dense specialization mentioned earlier must stay in the same VU as their memory during partitioning. When contexts are merged into a single VU, they are still separate contexts triggered independently.

3.3 Resource Allocation

The output of the imperative to dataflow transformation discussed in Section 3.2 is a VUDFG that can execute on a Plasticine with physical units (PUs) that have infinite resources. The *Resource Allocation* phase enforces and addresses constraint violations given the specification of the Plasticine units. At the end of this phase, SARA assigns each VU in the VUDFG graph to a PU type with required resources; the placer then takes the type assignments and determines the final placement.

¹When mapping an expression with less than 6 operations to the SIMD pipeline in PCU, the latency of the expression is still 6 because the data must propagate to the end of the SIMD pipeline.

Accelerators often have heterogeneity in compute resources to improve efficiency for common special operations. In Plasticine, PMUs and AGs have specialized compute pipelines for address calculation that are less capable than the compute pipeline in PCUs. However, heterogeneity tends to reduce average utilization because different applications, and even the same application with different data sizes, can vary highly in the desired ratio among different resource [71]. A compute-bound application, for example, can heavily underutilize the AGs and PMUs. To address this problem, SARA models the virtual to physical assignment as a constraint satisfaction problem; each VU consumes a set of resources and can only be assigned to a PU if the PU processes the required resources. Table 3.3 shows the types of resources SARA models in Plasticine’s heterogeneous units. For example, special connection to off-chip memory interface is also treated as a type of resource in the AG, which forces virtual contexts accessing DRAM to map to AGs. On the other side, regular contexts with non-vectorized fixed-point operations can also be mapped to spare AGs, which improves utilization.

As shown in Algorithm 1, the *resource allocation* phase contains three steps: *constraint resolution*, *global merging*, and *virtual to physical assignment*. SARA uses a VU-PU bipartite graph (G) to keep track of potential valid assignments between the two. Initially, G is initialized to a complete bipartite graph, i.e., all VUs can be assigned to all PUs. We refer to all PUs connected to a VU v as the domain of v in G , i.e. $dom(v)$.

Constraint Resolution A list of constraint pruners, each considering a set of on-chip resources, incrementally remove the VU-PU edges that violate the resource constraints. If a VU v has an empty domain after pruning, the pruner attempts to fix the violation by decomposing the VU into multiple VUs. Not all resources are composable, and the partitioning transformation may fail. If succeeded, the partitioner generates a new set of VUs V' . SARA starts a new complete bipartite graph between V' and all resources P , and recursively prune on V' . If succeeded, the original graph G is updated with V' and their pruned resources.

Global Merging After all VUs have at least one PU in the bipartite graph, SARA triggers a global optimization that merges small VUs into a larger VU to reduce fragmentation in allocation. Each type of resource has an aggregation rule to compute how the resource cost changes if two VUs are merged together, as shown in Table 3.3. Most aggregation rules are simple, such as addition, logical or, max, or union. The in- and out-degree costs are trickier and will be detailed in Section 3.3.1.

Feature	PCU	PMU	AG	Host Unit	Aggregation Function
Vector lane width	16	16	1	1	MAX
# pipeline register (PR)	8	8	4	0	
# stages	6	10	5	0	
Scratchpad banks	0	16	0	0	SUM
Scratchpad capacity	0	256kB	0	0	
MergeBuffer	1	0	0	0	
Splitter	1	0	0	0	
Scanner	1	0	0	0	
Operation types	fix \cup float	fix	fix	\emptyset	\cup
Reduction tree	✓	✗	✗	✗	OR
Access to DRAM Interface	✗	✗	✓	✗	
Access to Host IO	✗	✗	✗	✓	
# Vector Input	6	6	4	0	G
# Scalar Inputs	6	6	4	16	
# Control Inputs	16	16	4	16	
# Vector Outputs	6	6	4	0	
# Scalar Outputs	6	6	4	16	
# Control Outputs	8	8	2	16	

Table 3.3: A list of resources SARA models in four types of configurable units in Plasticine. The host unit models the host registers I/Os. MergeBuffer, Splitter, and Scanner are new hardware units introduced in [48] and recent work to support database and sparsity in Plasticine. The aggregation function indicates how to compute the aggregated resource cost when two contexts are merged into a single virtual unit (VU). G indicates the aggregated value is the output of a graph traversal of the merged graph. How to count # I/O is discussed later in Section 3.3.1. While the aggregation function of #PR of two merged contexts is MAX, the #PR of a context is the maximum number of live variables of its dataflow graph, which is also an output of a topological traversal. This table reflects a different Plasticine configuration as the original Plasticine in [5].

```

Function alloc(V, P, pruners): /* Allocation Algorithm */
    Data: V: a set of VUs from the VUDBG
    Data: P: a set of all PUs on the hardware
    Data: pruners: a list of constraint pruners to check and fixes constraint violations
    /* Initialize a complete bipartite graph */
    G = new BipartiteGraph();
    G[V] = P;
    /* Constraint resolution */
    prune(G, pruners);
    /* Global merging */
    merge(G);
    /* Heuristic check on whether assigning all VUs in V is feasible */
    check(G);
    /* Virtual to physical assignment */
    backtracking_assign(G);

Function prune(G, pruners): /* A recursive pruning function */
    Data: G: bipartite graph between VUs and PUs
    Data: pruners: a list of constraint pruners to check and fixes constraint violations
    Result: The function update G by removing VU-PU edges that violates constraints
    guarded by pruners. The function may fail and raise an exception.
    /* All PUs on the hardware */
    P = G.values();
    for pruner in pruners:
        for v in G.keys():
            for p in G[v]:
                if pruner.cost(v) > pruner.cost(p):
                    | G[v] -= p;
            if G[v].empty():
                /* Partition VU v based on resource constraints
                registered in pruner. Not all resources can be
                partitioned and this step may fail. If succeeded, the
                function returns a new set of VUs. */
                V' = pruner.partition(v);
                G' = new BipartiteGraph();
                G'[V'] = P;
                prune(G', pruners);
                G -= v;
                G[V'] = G'[V'];

```

Algorithm 1: Resource allocation. The bipartite graph *G* contains a bi-directional many-to-many map. *G*[key] returns the set of values connecting to the key (dom(key)), and *G*[value] returns the set of keys connecting to the value. *G*[key] = value connects an edge between key and value. *G*[KeySet] = ValueSet creates all-to-all connection between KeySet and ValueSet.

```

Function check(G): /* Assignment feasibility check */
    Data: G: bipartite graph
    Result: Whether it is possible to assign all VUs in V with a different PU in P
    /* For every value set in G */
    for V in G.values().toSet():
        K = ∅;
        for v in V:
            for k in G[v]:
                if G[k] ⊂ V:
                    K += k;
            if |K| > |V|:
                return failure();
    return success();

```

Algorithm 2: Heuristic check on whether it is possible to assign all key with an value in a bipartite graph. Given there are only a few types of hardware tiles, $G.values().toSet()$ is relatively small. This algorithm roughly runs in $O(|G.keys()| \times |G.values()|)$, which is still much faster than the backtracking assignment with exponential runtime.

Virtual to Physical Assignment Next, SARA performs a quick heuristic check on the bipartite graph to see if there exists a possible assignment for all VUs with sufficient PUs (Algorithm 2), and provide feedback on the limiting resources, otherwise. Finally, SARA assigns each VU to a PU type with a backtracking search on the pruned bipartite graph.

This approach can be easily extended to handle new heterogeneous tiles in the architecture by registering the tile with existing or new types of resources with aggregation and partitioning rules. The rest of this section goes over two types of partitioning transformations—compute partitioning in Section 3.3.1 and memory partitioning in Section 3.3.2.

3.3.1 Compute Partitioning

The *compute-partitioning* phase addresses VUs using more compute resources than any PU can provide. If a VU contains multiple contexts, SARA first moves the contexts into separate VUs. If a single context exceeds the resource limit, SARA breaks down the dataflow graph in the context into multiple contexts and puts them in separate VUs. During partitioning, SARA maps each subgraph of the large dataflow graph into a new context, mirrors the control states of the original context, and streams live variables in between. We can formulate the problem of how to partition in the dataflow graph as an optimization problem, shown in Table 3.4. The partitioner “fixes” the VU v based on a single PU specification, albeit there are many potential PUs the decomposed VU can be mapped to. Currently, we use a heuristic to select a PU type from $dom(v)$ right before the compute pruning as a

Problem	Partition the dataflow graph into subgraphs such that all subgraphs satisfy the constraints of a hardware unit.
Objective	Minimize the number of partitions and connectivity across partitions.
Constraints	<p>Each partition must not exceeds the limit on the number of</p> <ul style="list-style-type: none"> • live in/out variables (I/O ports) • operations (pipeline stages), • and live variables across operations (pipeline registers), etc. <p>No <i>new</i> cycles can form across partitions other than the cycles in the original dataflow graph.</p>

Table 3.4: Formulation of the compute partitioning problem

guiding constraint for partitioning.

Because the global network is specialized to handle efficient broadcasts, the in/out-degree of a partition counts the number of unique live-in/out variables, as supposed to the number of edges across partitions. In addition, the partitioned subgraphs cannot form *new* cycles; contexts waits for all input dependencies and therefore cycles across contexts cause deadlock. Nonetheless, the original graph might contain cycles representing loop-carried dependencies, such as accumulation. For these cycles, SARA initializes the back edge of the cycle with dummy data to enable execution. Figure 3.12 shows examples of valid and invalid partitioning solutions. Figure 3.13 shows another partitioning example of a dataflow graph with cycles.

Community Detection The formulation of compute partitioning is similar to the community detection problem[81] that has a similar objective. The major difference is that the latter often takes the number of output partitions as an input to the algorithm, whereas our problem partitions until all subgraphs satisfy all constraints. Moreover, community detection algorithms do not enforce the cycle constraints. Finally, the edge connectivity in community detection counts the number of edges across partitions, as supposed to broadcast edges as in our problem.

Retiming Imbalanced data paths across partitions can cause pipeline stalls at runtime if the long-live path is not sufficiently buffered. To ensure full-throughput pipelining, SARA needs to insert retiming buffers along the imbalanced data path across partitions. Retiming introduces new VUs in addition to the partitioned VUs, which attributes to the cost in Table 3.4’s objective.



Figure 3.12: Compute partitioning examples. Solution 1 and 2 are both valid. Solution 2 is better because it has less number of broadcast edges across partitions (3 as supposed to 4 in Solution 1). Solution 3 is an illegal partitioning due to the cycle between partition 1 and 2.

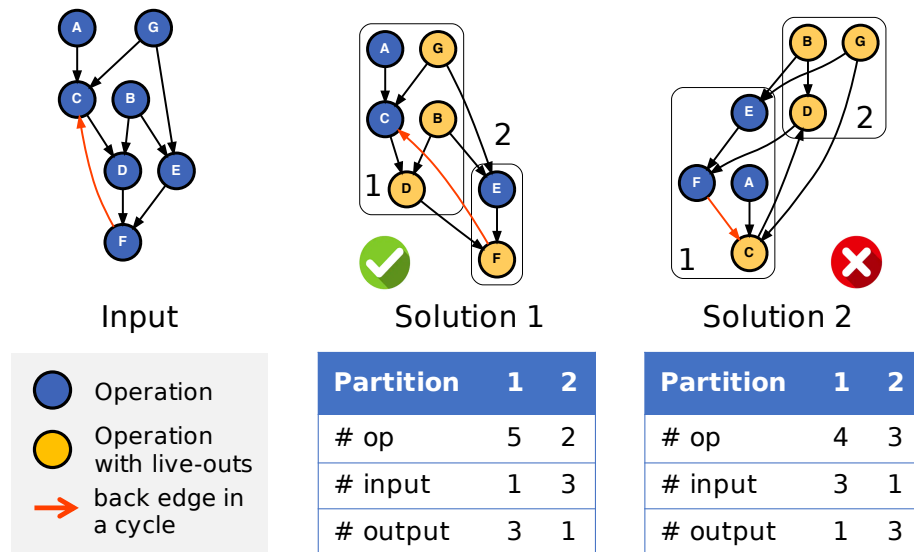


Figure 3.13: Compute partitioning examples with cycle in the dataflow graph. Solution 1 is valid because there is no cycle between partitions after removing the back edge in the original graph. Solution 2 is invalid because there is still cycle between partition 1 and 2 after removing the back edge.

In the following sections, we present two algorithms to resolve the problem described in Table 3.4: a fast traversal-based algorithm providing a decent solution, and a slow convex optimization-based algorithm providing an optimal solution.

Traversal-based Solution

To address the cycle constraint, the traversal-based algorithm performs a topological sort of the dataflow graph. The topological sort ignores the back edges of the cycles during traversal. Starting from one end of the sorted list, the algorithm iteratively adds nodes into a partition until it fits no more nodes. The algorithm then repeats the process with a new partition. This approach guarantees that no cycle is introduced with $O(V + E)$ complexity, where V and E are the numbers of vertices and edges in the dataflow graph.

The partitioning result is a function of the traversal order. We experienced with depth-first search (DFS) and breadth-first search (BFS) with forward and backward dataflow traversal orders. For DFS, we re-sort the remaining list each time starting with a new partition.

Solver-based Solution

The convex optimization solution models the problem as a node-to-partition assignment problem. Table 3.6 gives our formulation and Table 3.5 explains the notations used in Table 3.6.

At a high-level, we use a boolean matrix B to keep track of the assignment. B has dimension equals to the number of nodes in the dataflow graph by the maximum number of partitions, where $B[i, j] = 1$ indicates node i is assigned to partition j . In Table 3.6, *partition assignment* restricts each node to have a single partition assignment. The *input and output arity constraints* show the formulations that limit the number of input and output for a subgraph. These are the two most challenging constraints as we need to identify broadcast edges across partitions. To address the cycle constraint, we introduce a delay vector d_n with a size equivalent to the number of nodes. The delay vector encodes a time schedule to execute each node, whose values are selected by the solver. The *dependency constraint* enforces a node can be scheduled no earlier than its input dependencies and no later than its output dependents. Since the operations within a partition have to be triggered atomically, there is another delay vector d_p for partitions. The *delay consistency* enforces the schedule of a node equals to the schedule of its assigned partition. Finally, *constant validity* limits the range of values the delay vectors can be chosen from. In addition to enforcing the cycle constraint, these

Name	Type	Description	Definition / Default
\mathcal{N}	Constant	Enumeration of nodes to partition, numbered $\{n_i\}_i$	-
N	Constant, $\mathbb{Z}_{\geq 0}$	Number of operations to partition	$N = \mathcal{N} $
P	Constant, $\mathbb{Z}_{\geq 0}$	Number of partitions to consider	N , or from heuristic
\mathcal{E}	Constant, $\{n_i \rightarrow n_j\}$	Directed edges representing dependence	-
B	Variable, $\{0, 1\}^{N \times P}$	Boolean Partitioning Matrix	-
$\text{proj}_B(\cdot)$	$\mathbb{Z}_{\geq 0} \rightarrow \mathbb{B}$	Function to convert a positive integer into a boolean	Supplemental Materials
$\text{and}(\cdot, \cdot)$	$\{0, 1\} \times \{0, 1\} \rightarrow \{0, 1\}$	Boolean and of binary variables	Supplemental Materials
d_p	Variable, $\mathbb{Z}_{\geq 0}^P$	Vector of partition delays	-
d_n	Variable, $\mathbb{Z}_{\geq 0}^N$	Vector of node delays	-
$\text{dest}(n)$	$\mathcal{N} \rightarrow \mathcal{P}(\mathcal{N})$	The set of nodes which depend on n	$\{n' n' \in \mathcal{N} \text{ s.t. } (n \rightarrow n') \in \mathcal{E}\}$
c_o	Constant, $\mathbb{Z}_{\geq 0}$	Maximum output arity of a partition	HW Spec
c_i	Constant, $\mathbb{Z}_{\geq 0}$	Maximum input arity of a partition	HW Spec
b_d	Constant, $\mathbb{Z}_{\geq 0}$	Maximum input buffer depth	HW Spec
K	Constant, \mathbb{R}_+	Very Large Constant, used for constraint activation	$P \times N$
α_d	Hyperparameter, \mathbb{R}_+	Retime merging probability multiplier	$\frac{1}{\max\{c_o, c_i\}}$
\mathcal{C}_r	$[\mathcal{N} \rightarrow \mathbb{R}_+, \mathbb{R}_+, [\mathbb{R}_+] \rightarrow \mathbb{R}_+]$	List of per-node values, limits, and reduction functions for reducible constraints	Supplemental Materials
F	$\{0, 1\}^{N \times P}$	Feasibility matrix, whether a partition can support a node	HW Spec

Table 3.5: Names and definitions used in the solver-based algorithms.

Type	Description	Expression
Cost Function	Allocated Partitions	$\Sigma_i \text{proj}_{\mathbf{B}}(\Sigma_j B_{i,j})$
	Retiming Partitions	$\alpha_d \Sigma_{n_i \rightarrow n_j \in \mathcal{E}} \text{proj}_{\mathbf{B}}(\max\{d_n(j) - d_n(i) - b_d, 0\})$
Partition Constraint	Partition Assignment	$\forall n_i \in \mathcal{N} : \Sigma_j B_{i,j} = 1$
	Input Arity Constraint (vectorized)	$\Sigma_{n_i \in \mathcal{N}} \max\{\text{proj}_{\mathbf{B}}(\Sigma_{n_j \in \text{dest}(n_i)} B_{j,:}) - B_{i,:}, 0\} \leq c_i \times \vec{1}$
	Output Arity Constraint	$\forall p \in [0, P) :$ $\Sigma_{n_s \in \mathcal{N}} \text{and}(B_{s,p}, \text{proj}_{\mathbf{B}}(\max\{(\Sigma_{n_d \in \text{dest}(n_s)} B_{d,p}) - K \times B_{s,p}, 0\})) \leq c_o$
	Dependency Constraint	$\forall n_i \rightarrow n_j \in \mathcal{E} : d_n(i) + 1[p_i \neq p_j] \leq d_n(j)$
	Delay Consistency	$\forall n_i \in \mathcal{N} : d_n(i) \leq \min_j (d_p(j) + K - B_{i,j} \times K)$ $\forall n_i \in \mathcal{N} : d_n(i) \geq \max_j (d_p(j) + B_{i,j} \times K - K)$
	Constant Validity	$\forall n_i \in \mathcal{N} : d_n(i) \leq K$ $\forall i \in [0, P) : d_p(i) \leq K$
Merge Constraint	Feasibility Constraint	$\forall i, j \in [0, N) \times [0, P) : B_{i,j} \leq F_{i,j}$
	Reducible Constraints	$\forall j \in [0, P). \forall (c(\cdot), c_v, r(\cdot)) \in \mathcal{C} :$ $r([c(n_i) \times B_{i,j}]_{n_i \in \mathcal{N}}) \leq c_v$

Table 3.6: Solver formulation for partitioning. Expressions are presented using the Disciplined Convex Programming ruleset [7, 8]. Explanations for selected expressions can be found in the supplemental material.

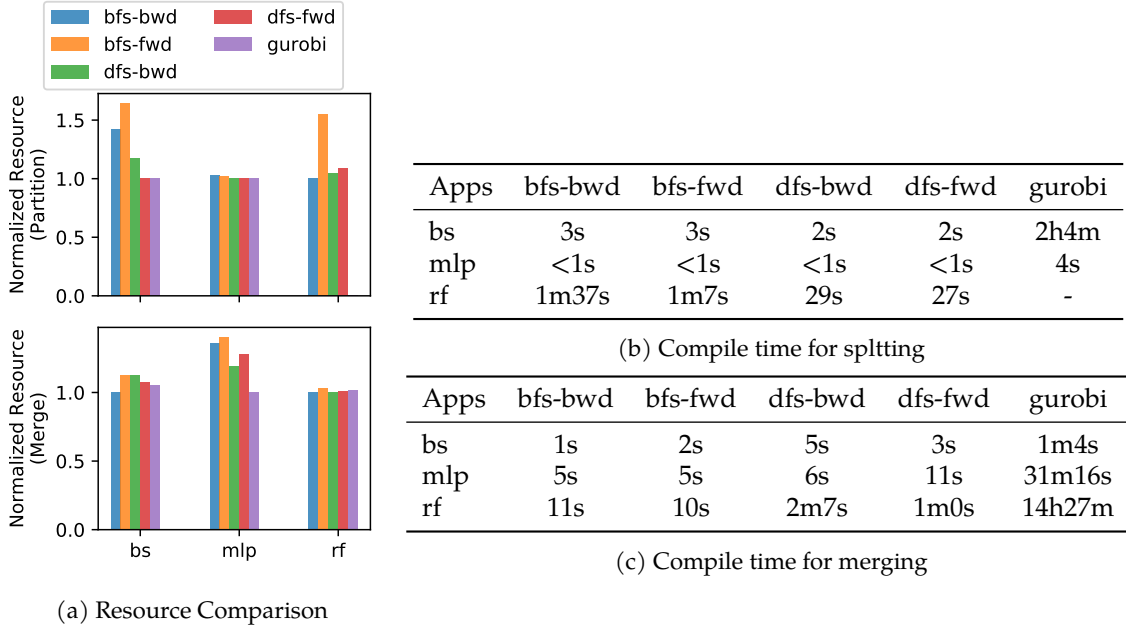


Figure 3.14: Partitioning and merging algorithm comparisons. (a) shows the normalized resource usage between different algorithms (the lower the better). (b) and (c) shows the compile time of each algorithm.

delay variables are also used to calculate where retiming is required and project the amount of introduced retiming VUs. Finally, we use the traversal-based solution to warm start the assignment matrix B and the delay vectors to reduce the solver runtime.

Comparison

TODO: Benchmark Table Figure 3.14 shows the comparison between the traversal-based and the solver-based solutions for both compute partitioning and global merging. Global merging is a global optimization merging small VUs into a large VU that can still fit in a PU. The merging algorithm is very similar to the compute partitioning algorithm, where nodes in the dataflow graph corresponds to the VUs in a VUDFG. The merging problem also has a traversal-based and solver-based solution. Section 3.4 will discuss merging in more detail.

We used a commercial solver, Gurobi [82], for the solver-based solutions. The evaluation is performed on an Intel Xeon E7-8890 CPU at 2.5GHz with 1TB DDR4 RAM. Gurobi is parallelized across ten threads for each application. To speed up convergence, we configured Gurobi with a 15%

optimality gap, i.e., the solver is allowed to early stop after the solution is more than 85% close to the optimum. The solving time increases dramatically as getting close to 100% optimum.

Figure 3.14 (a) shows the normalized resource in a number of VUs after partitioning and merging. We can see that Gurobi provides almost the best solution for all applications when it can derive an answer in a reasonable amount of time. The missing solver bar in random forest (*rf*) partitioning is due to timeout after a few days. The traversal-based algorithms can sometimes match or even outperform the solver slightly. However, because the partitioning result is a function of the traversal order, each traversal order has adversarial cases, where they can be up to 1.7x worse in resource than the best possible solution. The forward (*fwd*) traversal order schedules nodes as earlier as possible, which reduces the number of external live variables; the backward traversal minimizes the number of internal live variables across partitions. The depth-first-search (DFS) traversal order minimizes the number of live variables between partitions, albeit producing more imbalanced paths between partitions. On the other hand, breath-first-search (BFS) produces more balanced partitioning with more live variables and partitions.

There are two common graph patterns in applications that require partitioning. The first is a dataflow graph from a large basic block, which contains a small set of external live-in and live-out variables and abundant intermediate temporary variables. Such graphs typically end up with long-live variables across partitions that require retiming. **TODO: show example and discuss the other pattern.**

Figure 3.14 (b) and (c) shows the compile time for these algorithms. The single-threaded the traversal-based algorithm runs in minutes, which is significantly faster than the parallelized solver that takes hours to days. In general, the solver runtime becomes quickly unbounded with a large amount of VUs.

In summary, the solver solution provides a guaranteed close-to-optimum solution at an expensive compile time. Nonetheless, the solver-solution treats the retiming and partitioning as a joint optimization, whereas the traversal-based solution solves these two problems in two separate passes, which is less optimum. Moreover, the solver-based solution tends to produce a better result for architectures with tight I/O bound and relax stage constraint. The traversal-based solutions, on the other hand, can produce a decent solution in a short amount of time. However, the solution is a function of the traversal order; hence, the quality of the partitioning is highly sensitive to the graph structures. Certain traversal orders show a 10-50% worse in resource than the solver-based solution in Figure 3.14 (a). In practice, we can combine the two approaches and invoke the expensive solver

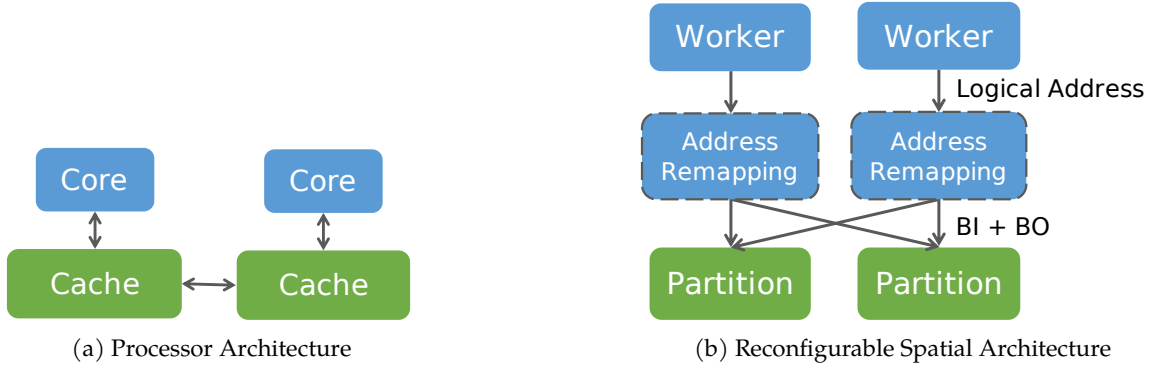


Figure 3.15: Memory model of different architectures

only when the traversal-based solution is insufficient. The quality of the traversal-based solution can be easily estimated by the resource utilization in each partition.

3.3.2 Memory Partitioning

The memory pruner addresses virtual on-chip memory exceeding the capacity and bandwidth limit of a single PMU. As we parallelize the computation, the on-chip memory must provide higher address bandwidth to sustain the compute throughput. On a processor-based architecture, this is often achieved with a separate first-level cache for each processor core, as shown in Figure 3.15 (a). The cache implements a hardware coherence protocol that synchronizes the different copies of data behind the scenes, providing the abstraction of a shared memory.

The coherence protocol is both expensive in hardware complexity and hard to scale in bandwidth for streaming pipelined execution. Instead of making redundant copies of the data, we can partition the data across different memory partitions to get additional address port on a reconfigurable accelerator, as shown in Figure 3.15 (b). Each parallel worker broadcasts the requests to all memory partitions. If the data accessed by two parallel workers live in the same bank, the bandwidth will be halved. To avoid bank conflicts, an important static analysis often used on reconfigurable accelerators is called static partitioning, or static banking [83]. For most address patterns, the compiler can derive a partitioning scheme such that every partition is accessed by a single worker at any time, which guarantees bandwidth at runtime. The output of the analysis is an expression for bank address (BA) and bank offset (BO), both are functions of the requested logical memory address. BA determines which partition each request is going to, and BO is the offset within the partition.

Figure 3.16 shows how static banking is achieved on the Plasticine architecture. Banking analysis



Figure 3.16: Example of memory partitioning across physical units. In this example, the program parallelizes the outer loop by two and vectorizes the inner loop by 4, which generates two vectorized access lanes. We need eight scratchpad banks to sustain the read bandwidth. SARA groups the 8 banks in two virtual memory partitions, P_A (bank 0-3) and P_B (bank 4-7). The address generation contexts A_0 and A_1 contains the computation for logical address and address remapping, which output BA s and BO s. SARA allocates one merge context per memory partition (M_A and M_B), merging requests from all access lanes. Inside the merge context, the shuffle operator converts the BO from access-aligned to bank-aligned. A_0 , for example, requests offsets A, B, C, and D (BO_0) from banks 1,4,5, and 3 (BA_0). The shuffle operator picks the requests belonging to partition A and outputs a BO aligned with bank 0-3. If the bank i in the partition has no request, the i th element in the output is marked as invalid. The merge context contains one shuffle operator per request lane, and uses a tree of OR operators to combine all bank-aligned BO s into the final BO_A . On the receiver side, the memory broadcasts its response to all receiver lanes. The receiver context uses an inverse shuffle operator to convert the response from bank-aligned back to access-aligned, using the BA forwarded from the address context. The response is then merged with another OR tree.

from Spatial specifies the number of banks required to sustain bandwidth for the current parallelization factors, and expressions for BAs and BOs. SARA groups the banks into virtual memories, with group size limited by the number of banks in a PMU. For each access lane, SARA allocates a context to compute the logical and remapped address, which outputs a vectorized BA and BO for each access lane. BAs and BOs are broadcasted to all memory partitions. For each memory partition, SARA allocates a merge context, merging BOs from all requesting lanes. The merge context uses a shuffle operator to transform the vectorized BO from bank order specified by BA (access-aligned) to bank order assigned to the current partition (bank-aligned). Because the static banking analysis guarantees no bank conflicts for all banks, SARA can use a OR tree to merge the bank-aligned BOs into the final BO that gets send to banks in the partition. On the receiver side, SARA uses an inverse shuffle to convert responses from partitions from bank-aligned back to access-aligned. SARA uses another OR tree to merge the access-aligned responses, which produces the final vector data requested by the access lane. With large parallelization, both the merge OR tree and the respond OR tree can be partitioned across VUs if running out of stages in a VU, and to scale in the network bandwidth.

TODO: Talk about what if banking fails.

3.3.3 Register Allocation

3.4 Optimizations

SARA performs many of the standard compiler optimizations, such as Dead Code Elimination and Constant Propagation. Some of them, however, play a much more important role for reconfigurable accelerator because they have a direct impact on resource usage. Other optimizations can be counter-intuitive, as they introduce redundant computation that reduces resources without necessarily impacting performance. In this discussion, we focus our primary objective on performance. Reducing resources is an indirect objective as resource reduction enable larger parallelization factors, which in turn improves performance. Other objectives, such as power saving, is left as future work.

3.4.1 Description

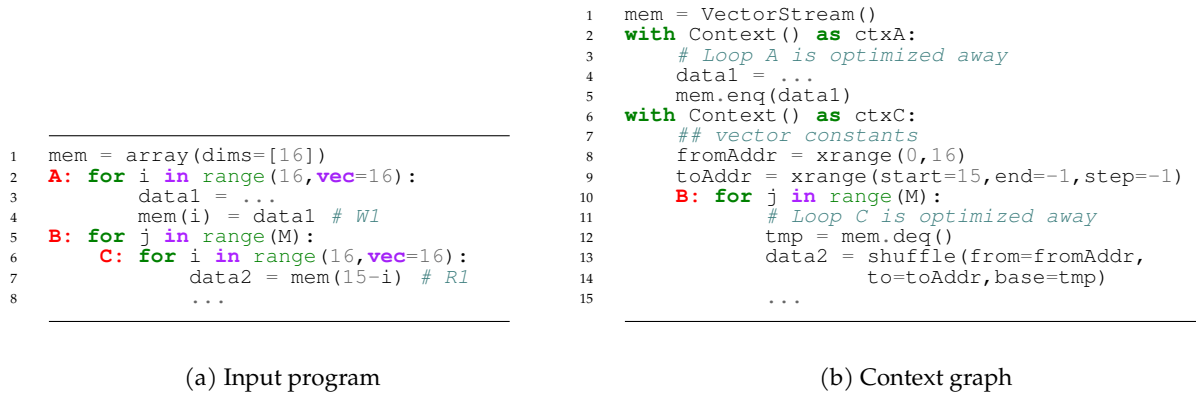


Figure 3.17: (a) shows an example where read and write address of `mem` can be constant folded because loop A and loop C are fully parallelized. Instead of mapping `mem` into a scratchpad, which heavily underutilize the capacity of the scratchpad, SARA maps the memory into a non-indexable vector stream in (b). The vector stream stores address [0-15] of the original memory, whereas the reader *R1* requests address [15-0] instead. SARA uses a `shuffle` operator to swap the vector elements based on reader's expected order. This is the same `shuffle` operator we introduced in the memory partitioning section Section 3.3.2.

Memory strength reduction (msr). Like traditional strength reduction on arithmetics, SARA replaces expensive on-chip memories with cheaper memory types whenever possible. Scratchpad, input buffer, and pipeline register are different types of on-chip memory on Plasticine from the most expensive to the cheapest. For example, SARA maps on-chip array to a scratchpad by default. However, if all of the readers and writers of the array index into memory with constant addresses, often as a result of constant propagation of full parallelized loops, SARA maps the memory to the input buffer of a context, and setup datapath to connect the reader and writer directly. Figure 3.17 shows an example of strength reducing a scratchpad into an input buffer with additional operations. Fully parallelized loops are a common outcome of an automatic design space exploration of the parallelization factors. We also use it as a way for the user to explicitly program vectorized streaming operation for Plasticine.

Route-Through Elimination (rtelm). For patterns where the content of a non-indexable memory (*M1*) is read and written to another memory (*M2*), SARA eliminates the intermediate access and feeds the output of *M1*'s writer (*W1*) directly to *M2*'s reader (*R2*) if it can prove the propagation is safe. Figure 3.18 shows examples of route-through opportunities. These patterns are often the outcome of multiple levels of compiler lowering and the *msr* optimization. Eliminating route-through patterns can simplify the control hierarchy and reduce the number of basic blocks, which

```

1 # Input program
2 accum = 0 # M1
3 out = 0 # M2
4 for i in range(N):
5     for j in range(M):
6         accum += i * j # W1
7         out = accum # R1 and W2
8     ...
9 ... = out # R2
10
11 # Optimized program
12 out = 0 # M2
13 for i in range(N):
14     for j in range(M):
15         out += i * j # W1
16     ...
17 ... = out # R2

```

(a) Example with scalar variables

```

1 # Input program
2 tmp = queue() # M1
3 out = queue() # M2
4 for i in range(N):
5     for j in range(M):
6         tmp.enq(i * j) # W1
7     for i in range(N*M):
8         out.enq(tmp.deq()) # R1 and W2
9     ...
10 ... = out.deq() # R2
11
12 # Optimized program
13 out = queue() # M2
14 for i in range(N):
15     for j in range(M):
16         out.enq(i * j) # W1
17     ...
18 ... = out.deq() # R2

```

(b) Example with queues

Figure 3.18: (a) and (b) shows route through opportunities in the program. These opportunities are often outcome of lowering and other optimizations. To proof the route through is safe, SARA needs to have access to all readers and writers of a memory, and analyze the control structure around the memory accesses. In addition to saving on the memory, the optimization eliminates the basic block in line 7 of (a) and line 8 of (b), which saves a context.

eliminates contexts.

Retiming with scratchpad (retime-mem). By default, SARA uses input buffers in PU for re-timing purposes. This option enables SARA to use scratchpad memory to retime paths requiring deep buffers. Figure 3.19 shows VUDFGs after retiming with and without scratchpad retiming. As we can see, allowing scratchpad and input buffer for retiming significantly reduces the total number of retiming units.

Crossbar datapath elimination (xbar-elm). Although in the general cases described in Section 3.3.2, crossbar data paths between accessors and memory partitions are very expensive, the BA can often be statically resolved to a constant vector for certain parallelization factors and address patterns. When BA is a constant, SARA can use this information to intelligently group banks into memory partitions that reduce the crossbar to a partial or a point-to-point connection. This is an important optimization that prevents resource scaling quadratically with increasing parallelization in the common cases.

Combined with *msr* and *rtelm* mentioned above, SARA is able to dramatically reduce the resource usage for certain parallelization factors. Figure 3.20 shows an interesting case study on an MLP presenting performance scaling and resource increase while scaling up parallelism. Figure 3.20

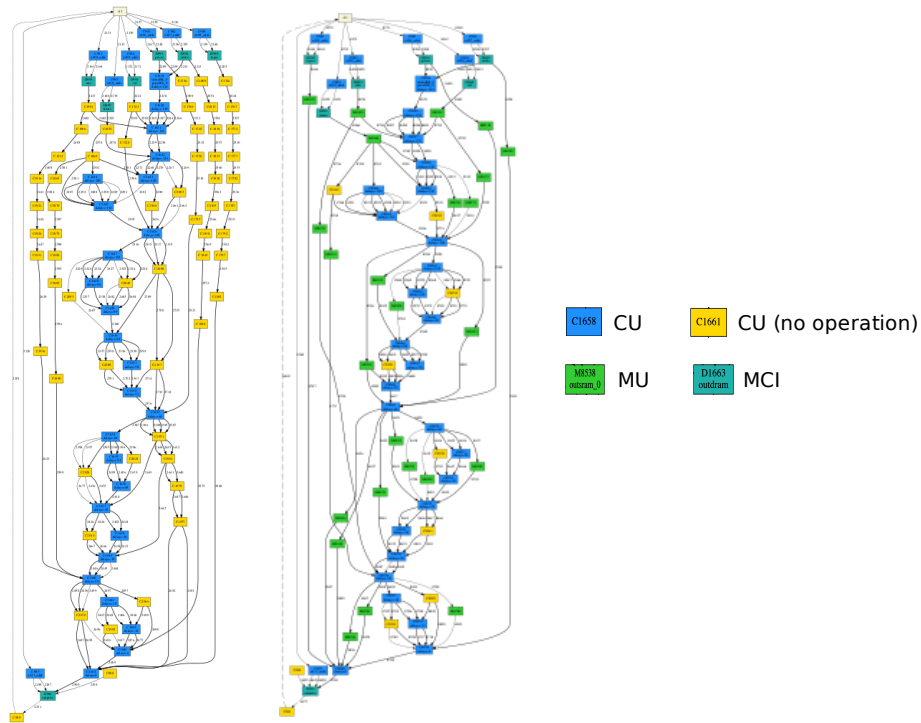


Figure 3.19: Virtual unit dataflow graph (VUDFG) after retiming. CU: virtual compute unit. MU: virtual memory unit. MCI: memory controller interface. Left: retiming with input buffer only. Right: retiming with either input buffer or scratchpad. The yellow CUs are CUs used for retiming. The original program does not have MU and the yellow MU on the right are MUs used for retiming.

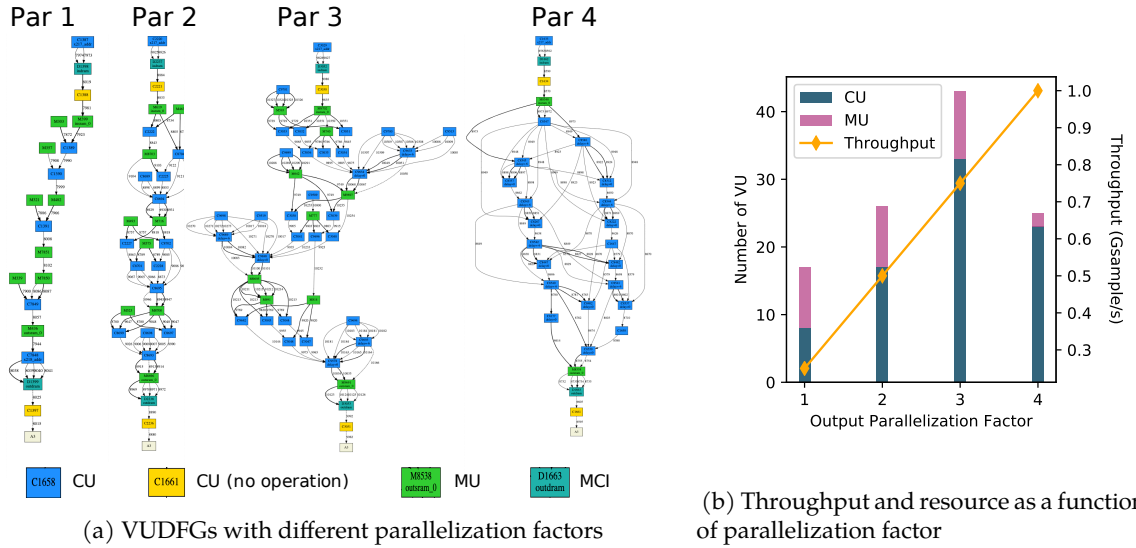


Figure 3.20: A 4x4x4 MLP case study. Light-weight models like this can be useful in ultra high-bandwidth and low-latency inference in packet processing [?]. (a) shows changes in VUDFG as we parallelize the output dimension. The input dimension of MLP is vectorized across SIMD lanes. As reaching par factor at 4, the intermediate memory in MU gets optimized away due to constant folding and rtelm on fully unrolled loops. (b) shows the throughput and resource increase as the output dimension is parallelized.

(b) demonstrates that Plasticine can have resource non-linearly increasing with parallelization factors, while having perfect performance scaling, due to compiler optimizations.

Read address duplication (dupra). Between the address context and receiver contexts in memory partitioning shown in Section 3.3.2 Figure 3.16, the imbalance datapaths between BA forwarding and request datapath can cause significant pipeline stalls, limiting the overall application performance to a small fraction of the ideal throughput. When parallelizing the memory access, the OR tree in the merge context gets partitioned into a tree of contexts, which further increases the mismatch in latency between the imbalanced datapaths. Instead of retiming BA forwarding path, this compiler flag forces SARA to duplicate the BA computation using local states within the receiver context, which could use fewer resources than retiming.

Global Merging (merge). After all VUs satisfy the hardware constraint, we perform a global optimization to compact small VUs into larger VUs. Merging has a very similar problem statement as compute partitioning (Section 3.3.1) except with more constraints. For merging, the VUDFG is the graph, the VUs are the nodes, and the merged VUs are partitions. Unlike the partitioning

```

1 mem1 = rand(N)
2 mem2 = rand(N)
3 # Block 1
4 c = a + b
5 for i in range(N):
6     # Block 2
7     mem2[i] += mem[i] * c

```

(a) Input program

```

1 mem1 = rand(N)
2 mem2 = rand(N)
3 for i in range(N):
4     # Block 2
5     c = a + b
6     mem2[i] += mem[i] * c

```

(b) Reversed loop invariant hoisting

Figure 3.21: The original program requires at least two contexts to execute Block 1 and Block 2. By moving the invariant instruction $c = a + b$ into the loop body, (b) only needs a single context instead. Because instructions within Block 2 are pipelined across loop iterations, adding instructions in the loop body introduce minimum performance impact. This transformation is beneficial until Block 2 exceeds six operations, at which point both version consume the same amount of resources.

problem, which uses a single set of cost metrics to determine if a partition is valid, the merged partition can be mapped to any PU available on-chip, each having a different set of cost metrics and a quantity limit. When adding a new VU m to a partition P in the traversal-based algorithm, we first take the union of the domains of VUs within the current partition, and intersect with the domain of m . The domain of the merged partition M is updated to all PUs within this intersection that M can fit.

$$M = P \cup \{m\} \quad (3.1)$$

$$\text{dom}(M) = \{p | p \in \{\cup \text{dom}(n) \mid \forall n \in P\} \cap \text{dom}(m), \text{cost}(M) \leq \text{cost}(p)\} \quad (3.2)$$

The caveat is that even if M has non-empty domain, the bipartite graph might not have a possible assignment after merging, as M might fit in a larger PU with insufficient quantity. Therefore, we perform the feasibility checking of the bipartite assignment at each step of merging, shown in Algorithm 2, and only merge if the assignment is feasible.

The solver-based solution combines the VU merging with VU to PU assignment as a joint problem. The output of merging gives both partition assignment as well as a PU type assignment, which eliminates the risk of infeasible assignment due to merging in the traversal-based solution.

Reversed Loop Invariant Hoisting A common loop optimization on traditional CPU is to move loop-invariant expressions outside of the loop body to reduce computation. This transformation, however, can introduce more basic blocks in the program, that translates to more contexts on Plasticine. Therefore, the reversed loop invariant hoisting, which moves loop-invariant expressions

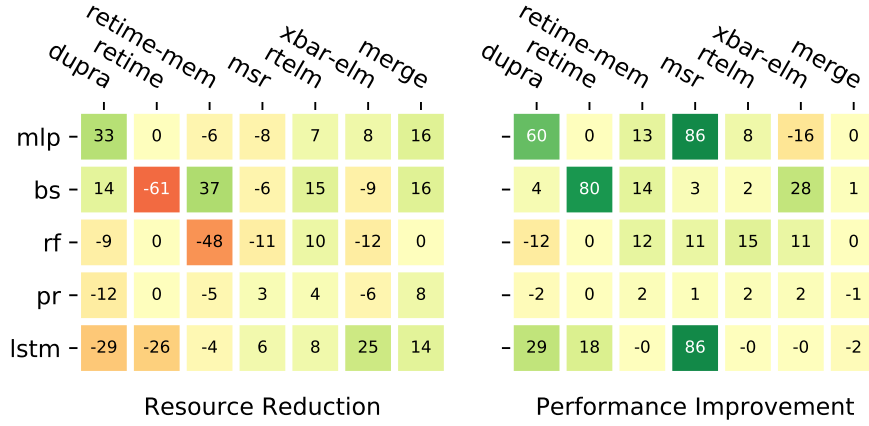


Figure 3.22: Optimization effectiveness. Percentage resource reduction or performance improvement from each optimization. The heat map shows the maximum differences when turning on/off the optimization, while fixing other optimizations at the most optimum combination for the application. Each application has a large design space with various parallelization factors in different loops, we report a design point that has the maximum impact, either positively or negatively.

into the loop body, can sometimes save resources by simplifying the control structure, as shown in Figure 3.21. Because instructions within the basic blocks are pipelined, moving instructions into the loop body increases computation without hurting throughput, which dominant the performance for spatial architecture. Currently, we rely on the user to perform this optimization manually.

3.4.2 Evaluation (WIP)

Figure 3.22 presents the evaluation of performance improvement and resource reduction from optimizations discussed in Section 3.4.1.

merge and *rtelm* are optimizations that gives a consistent resource improvement for most applications. *dupra* and *msr* helps **mlp** and **lstm**, which are heavily partitioned during the memory-partitioning stage.

dupra significantly reduces unbalancing in the data path caused due to the large merge tree generated by memory partitioning. *retime* can have a significant impact on the performance of the application at the cost of an increase in resource usage. Using scratchpad as a retiming buffer with the *retime-mem* flag can significantly reduce resource usage.

3.5 Placement and Routing

The input to the *placement and routing* (PaR) phase is a VUDFG, with each VU tagged to a type of PU. Before PaR, SARA also performs a runtime analysis of the program, annotating each edge in VUDFG with a link priority. The link priority is used to determine the routing order during PaR. PaR then iteratively places VUs and route edges in VUDFG onto the global network.

In addition to the original static network in Plasticine, we introduce a dynamic network in parallel to the static network, forming a dynamic-static hybrid network. Both purely static and purely dynamic networks are instances of the parameterized hybrid network. Section 4.1.2 will discuss the details about the hybrid network. The PaR algorithm needs to handle both network types and multiple network granularity (vector, scalar, control) at the same time.

The major difference between the static and the dynamic network is that each physical link in the static network is dedicated to a logical edge in the VUDFG for the entire duration of the execution (circuit-switching²), whereas the physical link in the dynamic network can be time-shared by multiple logical edges (packet-switching). Both static and dynamic networks are statically routed by the PaR algorithm. For a dynamic network, PaR also needs to assign virtual channels (VCs) to prevent network deadlock. The PaR algorithm configures the lookup table in each router that maps the packet header to the destination port and VC.

In the rest of this section, Section 3.5.1 discusses the placement and Section 3.5.2 reviews the routing algorithm. Section 3.5.3 explains the need for VC allocation. Lastly, Section 3.5.4 describes how SARA generates link priority.

3.5.1 Iterative Placement

At high-level, the PaR algorithm works very similarly to the FPGA PaR algorithm using simulated annealing [?]. For the initial placement, the algorithm places VUs in VUDFG in topological order. Each VU is placed to the next available PU with minimum Manhattan distance to the placed neighbors of that VU. Next, the algorithm routes all edges in the VUDFG, starting from edges with the highest link priority. If no routes are available on the static network, the route will be moved onto the dynamic network. For purely static networks, there is an “imaginary” dynamic network for this step. At the end of the PaR, if there are still routes on the fake dynamic network, PaR is considered failed for the static network. After all routes are routed, either on the static or dynamic network, the

²Technically, our static network is more restrictive than circuit-switching because circuit-switching allows deallocation after the connection is terminated. Our static network, on the other hand, cannot be reconfigured until the application terminates.

PaR evaluates the congestion cost of the current placement. Then, a genetic algorithm shuffles the VUs whose edges contribute most to congestion, and keeps the new position if it improves the route assignment. By iteratively re-placing and re-routing, the mapping process eventually converges to a good placement.

The PaR uses a heuristic cost model to rapidly evaluate placements: a penalty score is assigned as a linear function of several subscores. These include projected congestion on dynamic links, projected congestion at network injection and ejection ports, the average route length, and the length of the longest route. SARA provides a static estimate of the number of packets sent on each logical link. The PaR algorithm estimates congestion by normalizing the number of packets on each link to the program link with the highest total packet count. The most active program link sets a lower bound on the program runtime (the highest bandwidth physical link can still only send one packet per cycle), which translates to an upper bound on congestion for other links.

3.5.2 Congestion-Aware Routing

To achieve optimal performance, we use a routing algorithm that projects congestion and routes around it. Routing starts with the highest-priority routes, as determined by fanout of the broadcast edge and estimated packet count; broadcast edge with higher fanout is harder to route and hence are routed first. Using the packet count as a priority makes sure that the static network is used most efficiently. Our scheme searches a large space of routes for each link, using Dijkstra's algorithm [84] and a hop weighting function. To ensure maximum link reuse in broadcast edges, we can augment the hop weight in Dijkstra's algorithm by a multiplication factor between zero and one. When finding the shortest path between each source-destination pair, the weight of the hop is multiplied by the factor each time the hop is reused for the same broadcast edge. In other words, the reused path has a lower hop cost compared to other paths. This trick provides a balance between the shortest path and link reuse: a smaller factor encourages link reuse, even though individual source-destination pairs are not the shortest path; a factor equals to 1 ensure all source-destination pairs are routed with the shortest path. The other objective for routing broadcast links is balancing the hop counts between all destinations. This is because the shorter path will back pressure the sender before packets reaching to the longer path, causing pipeline stalls. To achieve this, we start with the source-destination pair with the longest Manhattan distance, ensuring the most far apart source-destination pair is routed on the shortest path. The consecutive source-destination pairs will try sharing the link on this path and slightly detour from their shortest path, which balances the

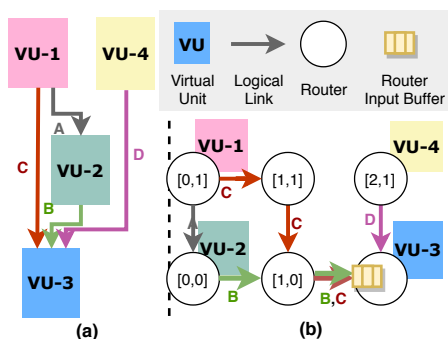


Figure 3.23: An example of deadlock in a streaming accelerator, showing the (a) VU data-flow graph and (b) physical placement and routes on a 2×3 network. There are input buffers at all router inputs, but only the buffer of interest is shown.

hop count.

3.5.3 VC Allocation for Deadlock Avoidance

Deadlock is a system pathology in dynamic routing where multiple flits form a cyclic holds on/waits for dependency on each others' buffers and prevent forward progress. Most dataflow accelerators use a streaming model, where outputs of a producer are sent over the network to one or more consumers without an explicit request; the producer is backpressure when there is insufficient buffer space. While this paradigm improves accelerator throughput by avoiding the round-trip delay of a request-response protocol, it introduces an additional source of deadlock [85].

Figure 3.23(a) shows a sample VUDFG graph, which is statically placed and routed on a 2×3 network in (b). Logical edges B and C share a physical link in the network. If C fills the buffer shared with B, VU-3 will never receive any packets from B and will not make forward progress. With streaming computation, the program graph must be considered as part of the network dependency graph, which must be cycle-free to guarantee deadlock freedom. However, this is infeasible because cycles can exist in a valid VU dataflow graph when the original program has a loop-carried dependency. Therefore, deadlock avoidance using cycle-free routing, such as dimension-order routing, does not work in our scenario. Allocating VCs to prevent multiple logical links from sharing the same physical buffer is consequently the most practical option for deadlock avoidance on streaming accelerators.

Benchmark	Throughput Unit	GPU Compiler	Latency (ms)		Throughput (Unit/s)	
			SARA	GPU	SARA	GPU
SqueezeNet (batch-1)	<i>kFrames</i>	TF+cuDNN	49.13	70.10	0.12 (1.1)	0.4
LSTM (batch-32)	<i>kSamples</i>	TF+cuDNN	3.61	6.81	8.8 (79.2)	4.7
PageRank	<i>MEdges</i>	GunRock	128.27	829.39	49	7.5
BlackScholes	<i>GOptions</i>	CUDA	0.09	0.10	88.88	80.02
Random Forest	<i>MSamples</i>	CUDA	0.10	0.32	1.04	0.32
Merge Sort	<i>GElements</i>	CUDA	0.63	2.14	6.65	1.96

Table 3.7: Performance comparison of Plasticine with Tesla’s V100 GPU (Normalized throughput to transistor count in parentheses).

3.5.4 Runtime Analysis for Heuristic Generation

In Spatial, the user can annotate runtime-variable input values to assist compiler analysis. We use these programmer annotations to compute the expected number of iterations each basic block will execute. The execution count on the basic block can future used to derive the packet count produced by these basic blocks. If the control is data-dependent, the user can annotate the data-dependent loop range, or the fraction for how frequently the branches evaluating to true. These heuristic guides can help the placer to evenly spread out the traffic. However, we do not require exact annotations for efficient placement—a rough estimate of data size is sufficient for the placer to determine the relative importance of links.

Using the estimated packet counts, the placer prioritizes highly used links for allocation to the static network and leaves infrequently used links on the dynamic network. When no annotation is provided, or loop bounds cannot be constant-propagated, the compiler estimates loop iteration counts based on the nesting depth: packets generated by the innermost loops are the more likely to be frequent. This heuristic provides a reasonable estimate of links’ priorities for routing purposes.

TODO: Elaborate on how to perform the analysis especially for streaming program

3.6 Evaluation (WIP)

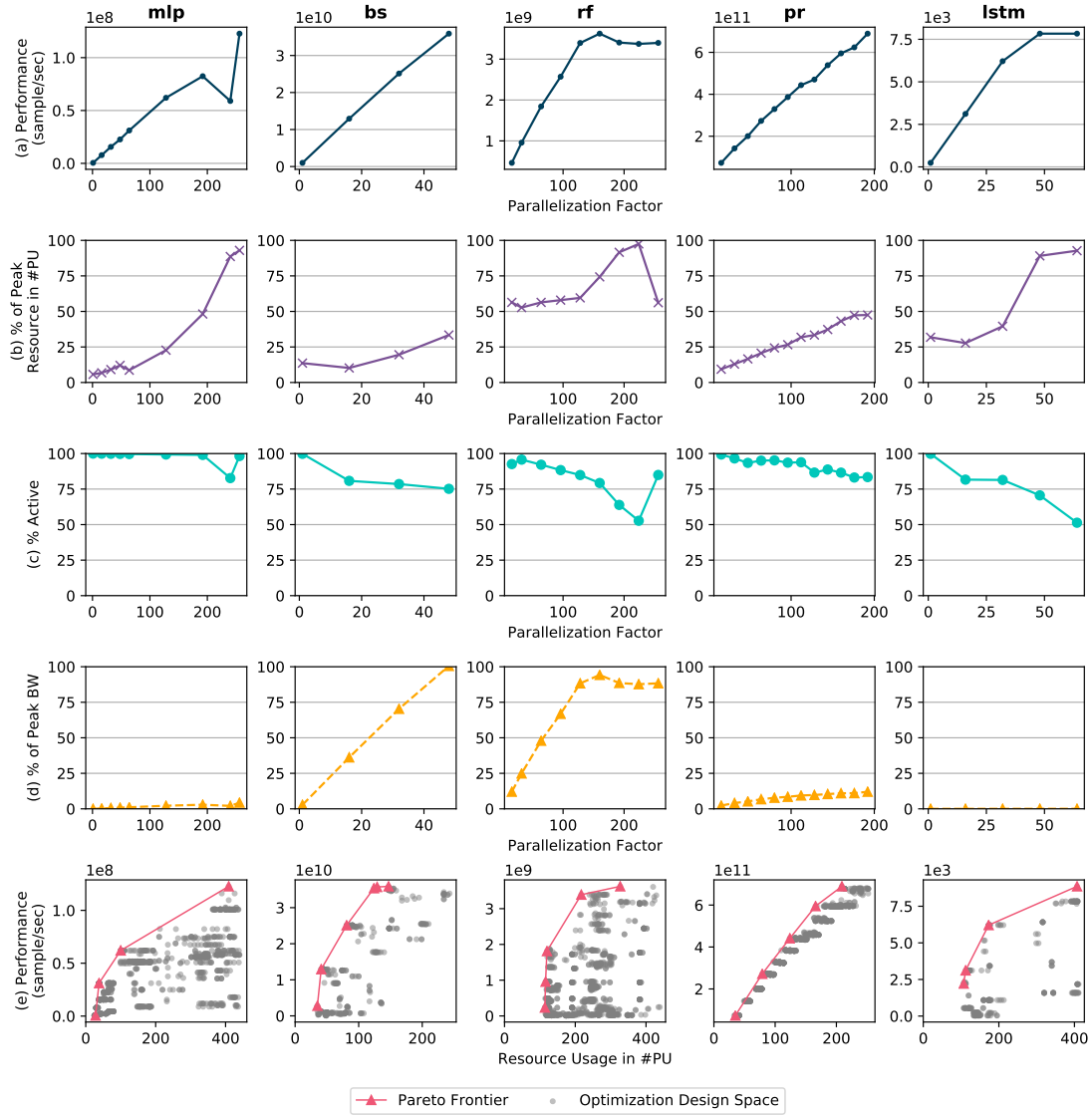


Figure 3.24: Scalability Evaluation. The first four charts show the scaling of (a) throughput, (b) resource usage, (c) runtime activation rate of PUs on the critical path of the compute pipeline, and (d) achieved HBM bandwidth, as the program gets parallelized. (e) shows the combined design space of compiler optimizations and parallelization factors on a throughput-resource curve. The pareto frontier presents the throughput shown in (a) as a function of resource increase in (b).

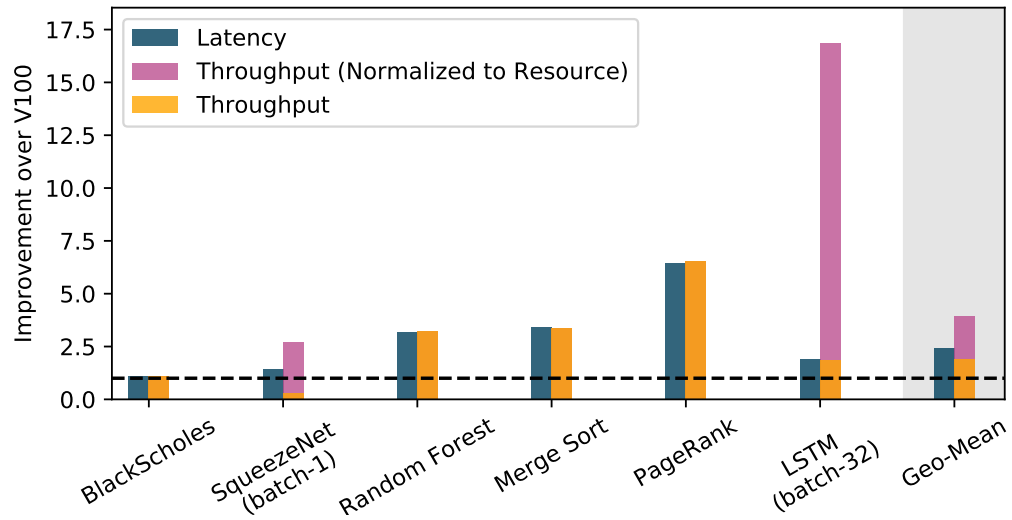


Figure 3.25: Plasticine’s latency and throughput improvement over V100 GPU. The evaluated Plasticine architecture has area footprint of $352mm^2$ at 28nm. V100 GPU has area footprint of $815mm^2$ at 12nm. Both platforms have the same off-chip bandwidth at 1TB/s with HBM technology. Yellow and blue bars show the raw measured speedup in throughput and latency, respectively. To account for the resource discrepancy, the pink bar shows the normalized throughput for compute-bound application—SqueezeNet and LSTM, which scales performance with additionally on-chip resources.

Chapter 4

Architecture

In this section, we discuss the architectural advancement on top of the original Plasticine architecture introduced in [5]. These architectural additions help increase the application coverage or improve the mapping strategies of existing applications by supporting new language constructs, data types, and improves the utilization of the hardware. Specifically, Section 4.3 layouts the datapath changes in order to support more flexible banking schemes required by general access patterns supported in Spatial; Section 4.2 discusses the hardware specialization and architectural sizing for machine learning applications; ?? provides an extensive study on on-chip network selection reconfigurable spatial architectures.

4.1 On-chip Network

Achieving scalable performance using spatial architectures while supporting diverse applications requires a flexible, high-bandwidth interconnect. Because modern CGRAs support vector units with wide datapaths, designing an interconnect that balances dynamism, communication granularity, and programmability is a challenging task.

On-chip interconnects can be classified into two broad categories: *static* and *dynamic*. Static interconnects use switches programmed at compile time to reserve high-bandwidth links between communicating units for the lifetime of the application. CGRAs traditionally employ static interconnects [86, 87]. In contrast, dynamic interconnects, or NoCs, contain routers that allow links to be shared between more than one pair of communicating units. NoC communication is typically packet-switched, and routers use allocators to fairly share links between multiple competing

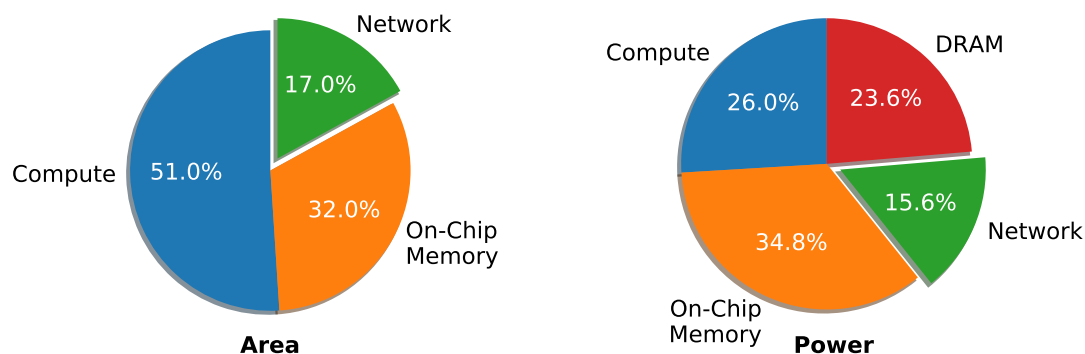


Figure 4.1: Area and power breakdown of Plasticine

packets. Although static networks are fast, they require over-provisioning bandwidth and can be underutilized when a dedicated static link is reserved for a logical link that is not 100% active. While dynamic networks allow link sharing, the area and energy cost to transmit one bit of data is higher for routers than for switches, making bandwidth scaling more expensive in dynamic networks than in static networks.

In this section, we explore the space of spatial architecture interconnect dynamism, granularity, and programmability. We start by characterizing several benchmarks' communication patterns and showing links' imbalanced bandwidth requirements, fanout, and data width in Section 4.1.1. Using these insights, we describe a hybrid network with both static and dynamic capabilities to enable both high bandwidth traffic and high resource sharing. Section 4.1.2 identifies a space of interconnection networks with static and dynamic capabilities, at multiple granularities. Next, We explain our methodology on performance, area, and power modeling in Section 4.1.3. Finally, Section 4.1.4 perform a detailed evaluation of the identified design space for a variety of benchmarks.

Because CGRAs encompass a broad range of architectures, we narrow our study on tile-based CGRAs with a streaming dataflow execution model, like Plasticine. At high-level, the architecture may contain a pool of heterogeneous compute and memory tiles (we referring as physical units (PUs)) with a global network. The network guarantees exactly-once, in-order delivery with variable latency, and communication between PUs can have varying granularities (e.g., 512-bit vector or 32-bit scalar).

To generalize the study, we introduce a variant processing engine (PE) style than Plasticine's pipeline SIMD unit. The alternative style uses time-scheduled execution, where each PU contains a vector function unit (FU) that can execute a small loop of instructions repeatedly in time. The

scheduling window is small enough that instructions are stored as part of the configuration fabric, without dynamic instruction fetch overhead. Compared to the pipelined architecture, this execution model creates more interleaved pipelining across PUs with communication that is tolerant of lower network throughput, which provides an opportunity to share links.

4.1.1 Application Characteristics

The requirements of an interconnection network are a function of the communication pattern of the application, underlying CGRA architecture, and compilation process. We identify the following key characteristics of spatially mapped applications:

Vectorized communication Recent hardware accelerators use large-granularity compute tiles (e.g., vectorized compute units and SIMD pipelines) for SIMD parallelism [5, 52], which improves compute density while minimizing control and configuration overhead. Coarser-grained computation typically increases the size of communication, but glue logic, reductions, and loops with carried dependencies (i.e., non-parallelizable loops) contribute to scalar communications. This variation in communication motivates specialization for the optimal area- and energy-efficiency: separate networks for different communication granularities.

Broadcast and incast communication A key optimization to achieve good performance on spatial reconfigurable accelerators is to explore multiple levels of parallelization and pipelining, within and across PUs. By default, pipeline parallelism across PUs introduces high-bandwidth one-to-one communication between dependent stages. To balance the pipeline throughput, we often need to parallelize the pipeline stage with the most computation, resulting in one-to-many communication when the receiver stage is parallelized, and many-to-one communication when the producer stage is parallelized. When both the producer and the consumer are parallelized, the worst case is many-to-many communication, as illustrated in Section 3.3.2. These broadcast and incast patterns introduce challenges in high-performance on-chip networks, as their bandwidth demands scale quadratically with their chip size in the worst case.

Compute to memory communication To encourage better sharing of on-chip memory capacity, many accelerators have shared scratchpads, either distributed throughout the chip or on its periphery [5, 53, 44]. Because the compute unit has no local memory to buffer temporary results, the results of all computations are sent to memory through the network. This differs from the NoCs

Benchmark	Description	Data Size
DotProduct	Inner product	1048576
OuterProduct	Outer product	1024
BlackScholes	Option pricing	1048576
TPCHQ6	TPC-H query 6	1048576
Lattice	Lattice regression [89]	1048576
GDA	Gaussian discriminant analysis	127×1024
GEMM	General matrix multiply	$256 \times 256 \times 256$
Kmeans	K-means clustering	$k=64, \text{dim}=64, n=8192, \text{iter}=2$
LogReg	Logistic regression	$8192 \times 128, \text{iter}=4$
SGD	Stochastic gradient descent for a single layer neural network	$16384 \times 64, \text{epoch}=10$
LSTM	Long short term memory recurrent neural network	1 layer, 1024 hidden units, 10 time steps
GRU	Gated recurrent unit recurrent neural network	1 layer, 1024 hidden units, 10 time steps
LeNet	Convolutional neural network for character recognition	1 image

Table 4.1: Benchmark summary

used in multi-processors, where each core has a local cache to buffer intermediate results. Studies have shown that for large-scale multiprocessor systems, network latency—not throughput—is the primary performance limiter [88]. For spatial accelerators, however, compute performance is limited by network throughput, and latency is comparatively less important.

Communication-aware compilation To study the communication patterns, we select a mix of applications from domains where hardware accelerators have shown promising performance and energy-efficiency benefits, such as linear algebra, databases, and machine learning. Table 4.1 lists the applications and their data size. Figure 4.2 shows, for each design, which resource limits performance: compute, on-chip memory, or DRAM bandwidth. DotProduct, TPCHQ6, OuterProduct, and BlackScholes are DRAM bandwidth-bound applications. These applications use few on-chip resources to achieve maximum performance, resulting in minimal communication. Lattice (a fast inference model for low-dimensional regression [89]), GDA, Kmeans, SGD, and LogReg are compute-intensive applications; for these, maximum performance requires using as much parallelization as possible. Finally, LSTM, GRU, and LeNet are applications that are limited by on-chip memory bandwidth or capacity. For compute- and memory-intensive applications, high utilization translates to a large interconnection network bandwidth requirement to sustain application throughput.

Figure 4.3(a,b) shows the communication pattern of applications characterized on the pipelined CGRA architecture, including the variation in communication granularity. Compute and on-chip

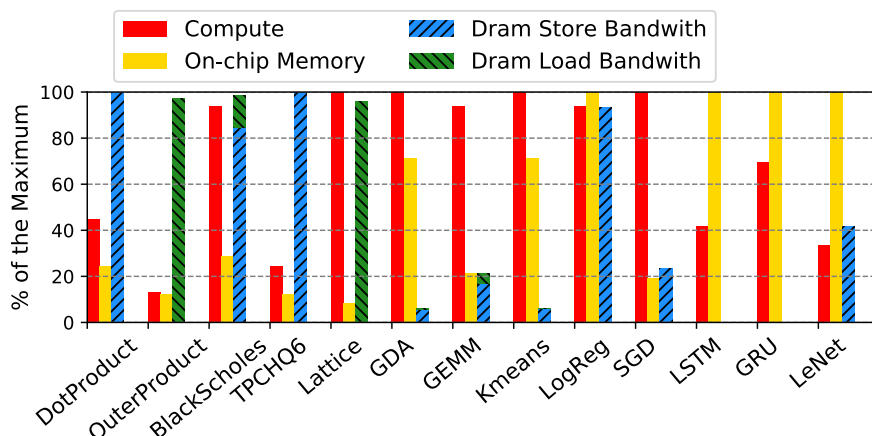


Figure 4.2: Physical resource and bandwidth utilization for various applications.

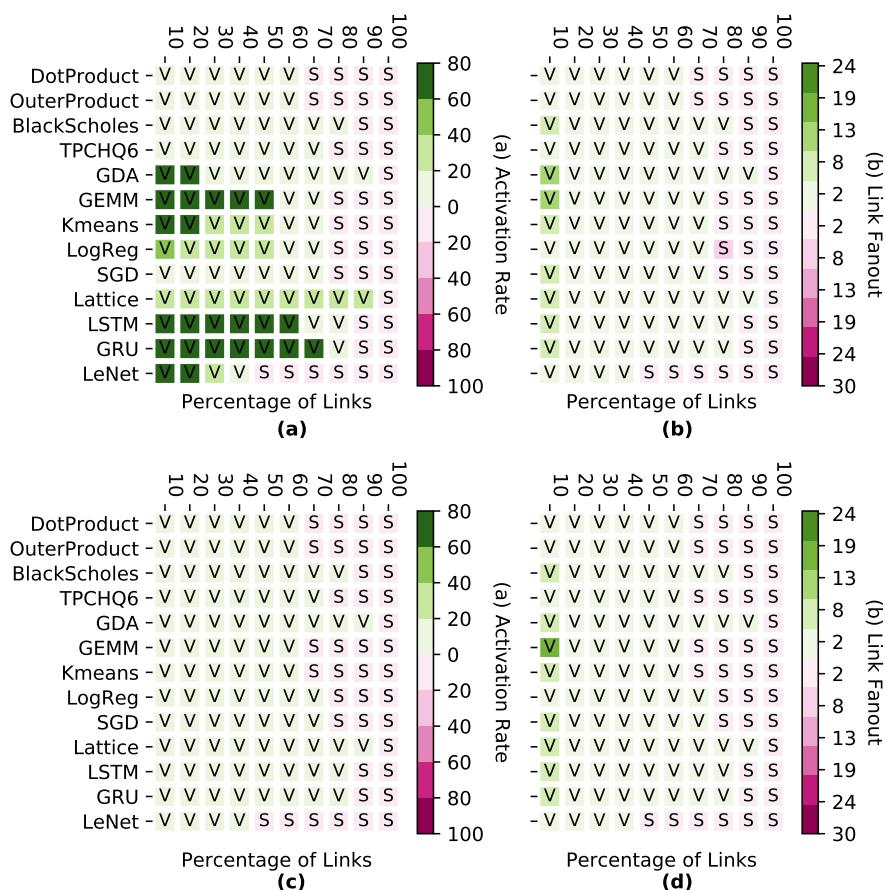


Figure 4.3: Application communication patterns on pipelined (a,b) and scheduled (c,d) CGRA architectures. (a) and (c) show the activation rate distribution of logical links at runtime. Links sorted by granularity, then rate; darker boxes indicate higher rates. The split between green and pink shows the ratio of logical vector to scalar links. (b) and (d) show the distribution of broadcast link fanouts.

memory-bound applications show a significant amount of high-bandwidth communication (links with almost 100% activity). A few of these high-bandwidth links also exhibit a high broadcast fanout. Therefore, a network architecture must provide sufficient bandwidth and efficient broadcasts to sustain program throughput. On the contrary, time-scheduled architectures, shown in Figure 4.3(c,d), exhibit lower bandwidth requirements due to the lower throughput of individual compute PUs. Even applications limited by on-chip resources have less than a 30% firing rate on the busiest logical links; this reveals an opportunity for link sharing without sacrificing performance.

Figure 4.4 shows statistics describing the VU dataflow graph (VGDFG) before and after the partitioning described in Section 3.3. The blue bars show the number of VUs, number of logical links, and maximum VU input/output degrees in the original parallelized program; the yellow and green bars show the same statistics after partitioning. Fewer VUs are partitioned for hybrid networks and dynamic networks with the time-scheduled architecture. When a VU in a pipelined CGRA consumes too many inputs or produces too many *distinct* outputs, SARA partitions it to reduce its degree and meet the input/output bandwidth constraints of a purely static network. For dynamic and hybrid networks, partitioning is not strictly necessary, but it improves performance by decreasing congestion at the network ejection port associated with a PU. We do not partition broadcasts with high output degrees because they are handled natively within the network. Finally, the output degree does not change with partitioning because most outputs with a large degree are from broadcast links.

4.1.2 Design Space for Network Architectures

We start with several statically allocated network designs, where each SIMD pipeline connects to several switches, and vary flow control strategies and network bisection bandwidth. In these designs, each switch output connects to exactly one switch input for the duration of the program. We then explore a dynamic network, which sends program data as packets through an NoC. The NoC uses a table-based routing scheme at each router to allow for arbitrary routes and tree-based broadcast routing. Finally, we explore the benefits of specialization by evaluating design points that combine several of these networks to leverage the best features of each.

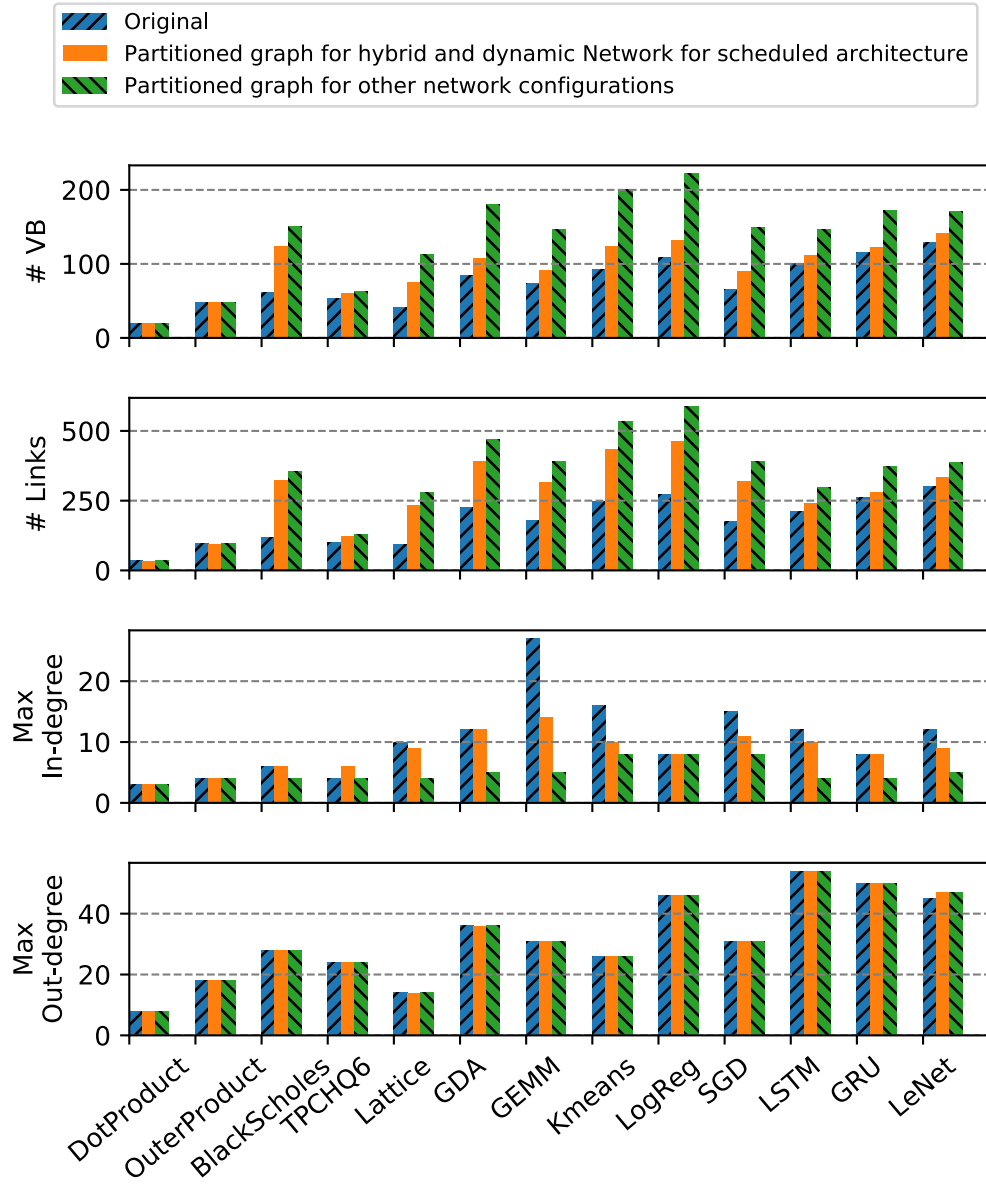


Figure 4.4: Characteristics of program graphs.

Static networks

We explore static network design points along three axes. First, we study the impact of flow-control schemes in static switches. In credit-based flow control [90], the source and destination PUs coordinate to ensure that the destination buffer does not overflow. For this design point, switches only have a single register at each input, and there is no backpressure between switches. The alternate design point uses a skid-buffered queue with two entries at each switch; using two entries enables per-hop backpressure and accounts for a one-cycle delay in stalling the upstream switch. At full throughput, the receiver will consume data as it is sent and no queue will ever fill up.

The second axis studied is the bandwidth, and therefore routability, of the static network. We vary the number of connections between switches in each direction, which trades off area and energy for bandwidth. Finally, we explore specializing static links: using a separate scalar network to improve routability at a low cost.

Dynamic networks

Our primary alternate design is a dynamic NoC using per-hop virtual channel flow control. Routing and Virtual Channel (VC) assignment are table-based: the compiler performs static routing and VC allocation, and results are loaded as a part of the routers' configurations at runtime. The router has a separable, input-first VC and switch allocator with a single iteration and speculative switch allocation [91]. Input buffers are sized just large enough (3 entries) to avoid credit stalls at full throughput. Broadcasts are handled in the network with duplication occurring at the last router possible to minimize energy and congestion. To respect the switch allocator's constraints, each router sends broadcasts to output ports sequentially and in a fixed order. This is because the switch allocator can only grant one output port per input port in every cycle, and the RTL router's allocator does not have sufficient timing slack to add additional functionality. We also explore different flit widths on the dynamic network, with a smaller bus taking multiple cycles to transmit a packet.

Because CGRA networks are streaming—each PU pushes the result to the next PU(s) without explicit request—the network cannot handle routing schemes that may drop packets; otherwise, application data would be lost. Because packet ordering corresponds directly to control flow, it is also imperative that all packets arrive in the order they were sent; this further eliminates adaptive or oblivious routing from consideration. We limit our study of dynamic networks to statically placed

Notation	Description
[S,H,D]	Static, hybrid, and dynamic network
x#	Static bandwidth on vector network (#links between switches)
f#	Flit width of a router or vector width of a switch
v#	Number of VC in router
b#	Number of buffers per VC in router
[db,cd]	Buffered vs. credit-based flow control in switch

Table 4.2: Network design parameter summary.

and routed source routing due to these architectural constraints. PUs propagate backpressure signals from their outputs to their inputs, so they must be considered as part of the network graph for deadlock purposes [85]. Furthermore, each PU has fixed-size input buffers; these are far too small to perform high-throughput, end-to-end credit-based flow control in the dynamic network for the entire program [90]. Practically, this means that no two logical paths may be allowed to conflict at *any* point in the network; to meet this guarantee, VC allocation is performed to ensure that all logical paths traversing the same physical link are placed into separate buffers.

Hybrid networks

Finally, we explore hybrids between static and dynamic networks that run each network in parallel. During static place and route, the highest-bandwidth logical links from the program graph are mapped onto the static network; once the static network is full, further links are mapped to the dynamic network. By using compiler knowledge to identify the relative importance of links—the link fanout and activation factor—hybrid networks can sustain the throughput requirement of most high-activation links while using the dynamic network for low-activation links.

4.1.3 Performance, Area, and Energy Modeling

Next section gives a quantitative analysis of the performance, area, and energy trade-offs involved in choosing a CGRA network, using benchmarks showing in Table 4.1. This section outlines our methodology on how we capture the network performance, area, and energy in our evaluations. Table 4.2 gives our notation for various network parameters.

Simulation

We use a cycle-accurate simulator to model the pipeline and scheduling delay for the two types of architectures, integrated with DRAMSim [92] to model DRAM access latency. For static networks, we model a distance-based delay for both credit-based and per-hop flow control. For dynamic networks, we integrate our simulator with Booksim [93], adding support for arbitrary source routing using look-up tables. Finally, to support efficient multi-casting in the dynamic network, we modify Booksim to duplicate broadcast packets at the router where their paths diverge. At the divergence point, the router sends the same flit to multiple output ports over multiple cycles. We assume each packet carries a unique ID that is used to look up the output port and next VC in a statically generated routing table, and that the ID is roughly the same size as an address. When the packet size is greater than the flit size, the transmission of a single packet takes multiple cycles.

Area and power

To efficiently evaluate large networks, we start by characterizing the area and power consumption of individual routers and switches used in various network configurations. The total area and energy are then aggregated over all switches and routers in a particular network. We use router RTL from the Stanford open-source NoC router [94] and our own parameterized switch implementation. We synthesize using Synopsys Design Compiler with a 28 nm technology library and clock-gating enabled, meeting timing at a 1 GHz clock frequency. Finally, we use Synopsys PrimeTime to back-annotate RTL signal activity to the post-synthesis switch and router designs to estimate gate-level power.

We found that power consumption can be broken into two types: inactive power consumed when switches and routers are at zero-load (P_{inactive} , which includes both dynamic and static power), and active power. The active power, as shown in Section 4.1.3, is proportional to the amount of data transmitted. Because power scales linearly with the amount of data movement, we model the marginal energy to transmit a single flit of data (flit energy, E_{flit}) by dividing active energy by the number flits transmitted in the testbench:

$$E_{\text{flit}} = \frac{(P - P_{\text{inactive}}) T_{\text{testbench}}}{\text{\#flit}} \quad (4.1)$$

While simulating an end-to-end application, we track the number of flits transmitted at each switch and router in the network, as well as the number of switches and routers allocated by place and

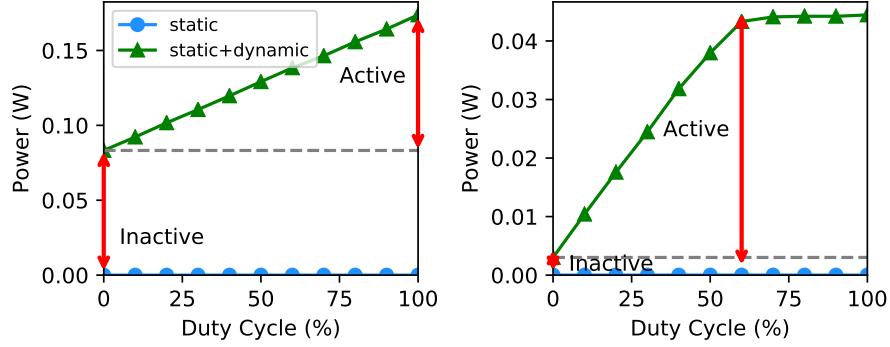


Figure 4.5: Switch and router power with varying duty cycle.

route. We assume unallocated switches and routers are perfectly power-gated and do not consume energy. The total network energy for an application on a given network (E_{net}) can be computed as:

$$E_{\text{net}} = \sum_{\text{allocated}} P_{\text{inactive}} T_{\text{sim}} + E_{\text{flit}} \# \text{flit}, \quad (4.2)$$

where P_{inactive} , E_{flit} , and $\# \text{flit}$ are tabulated separately for each network resource.

Figure 4.5 shows that switch and router power scale linearly with the rate of data transmission, but that there is non-zero power at zero-load. For simulation, the duty cycle refers to the amount of offered traffic, not accepted traffic. Because our router uses a crossbar without speedup [91], the testbench saturates the router at 60% duty cycle when providing uniform random traffic. Nonetheless, router power still scales linearly with accepted traffic.

A sweep of different switch and router parameters is shown in Figure 4.6. Subplots (d,e,f) show the energy necessary to transmit a single bit through a switch or router. Subplot (a) shows the roughly quadratic scaling of switch area with the number of links between adjacent switches. Vector switches scale worse with increasing bandwidth than scalar switches, mostly due to increased crossbar wire load. At the same granularity, a router consumes more energy a switch to transmit a single bit of data, even though the overall router consumes less power (as shown in Figure 4.5); this is because the switch has a higher throughput than the router. The vector router has lower per-bit energy relative to the scalar router because it can amortize the cost of allocation logic, whereas the vector switch has higher per-bit energy relative to the scalar switch due to increased capacitance in the large crossbar. Increasing the number of VCs or buffer depth per VC also significantly increases router area and energy, but reducing the router flit width can significantly reduce the router area.

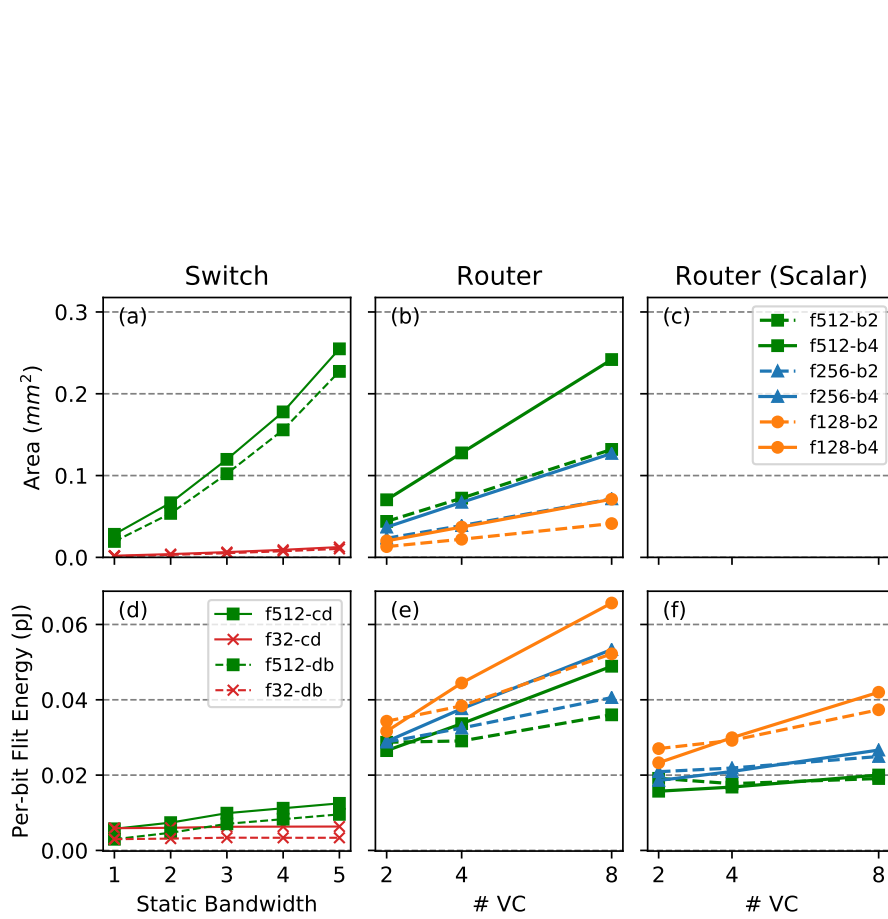


Figure 4.6: Area and per-bit energy for (a,d) switches and (b,c,f) routers. The router only has a vector granularity and can be partially clock-gated when sending scalar packets. Subplots (f) show the energy of the vector router when used for scalar values (32-bit).

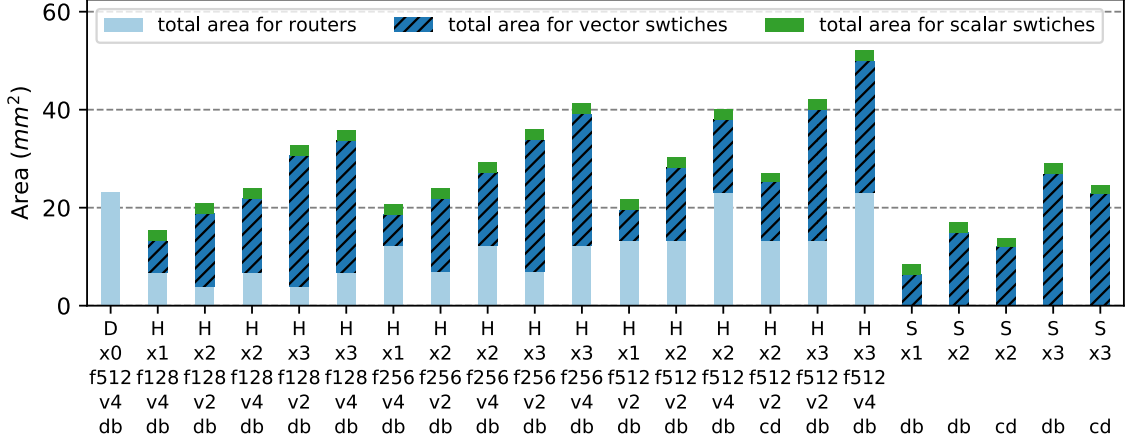


Figure 4.7: Area breakdown for all network configurations.

Overall, these results show that scaling static bandwidth is cheaper than scaling dynamic bandwidth, and a dynamic network with small routers can be used to improve link sharing for low bandwidth communication. We also see that a specialized scalar network, built with switches, adds a negligible area compared to and is more energy-efficient than the vector network. Therefore, we use a static scalar network with a bandwidth of 4 for the remainder of our evaluation, except when evaluating the pure dynamic network. The dynamic network is also optimized for the rare instances when the static scalar network is insufficient. When routers transmit scalar data, the high bits of data buffers are clock-gated, reducing energy as shown in (f). Figure 4.7 summarizes the area breakdown of all the network configurations that we evaluate.

4.1.4 Network Architecture Evaluation

We evaluate our network configurations in five dimensions: performance (perf), performance per network area (perf/area), performance per network power (perf/watt), network area efficiency (1/area), and network power efficiency (1/power). Among these metrics, performance is the most important: networks only consume a small fraction of the overall accelerator area and energy (roughly 10-20%). Because the two key advantages of hardware accelerators are high throughput and low latency, we filter out a network design point if it introduces more than 10% performance overhead. This is calculated by comparing to an ideal network with infinite bandwidth and zero latency.

For metrics that are calculated per application, such as performance, performance/watt, and power efficiency, we first normalize the metric with respect to the worst network configuration for

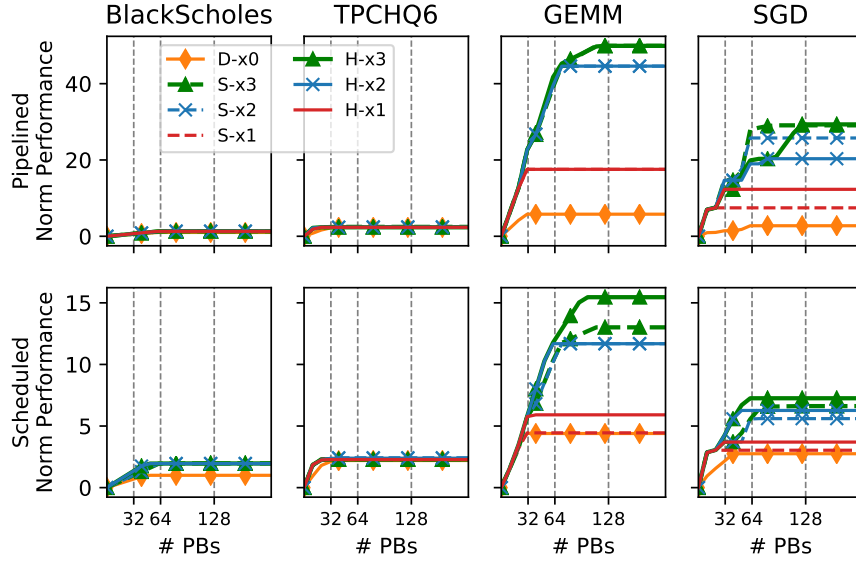


Figure 4.8: Performance scaling with increased CGRA grid size for different networks.

that application. For each network configuration, we present a geometric mean normalized across all applications. For all of our experiments, except Section 4.1.4, we use a network size of 14×14 end-point PUs. All vector networks use a vectorization factor of 16 (512 bit messages).

Bandwidth scaling with network size

Figure 4.8 shows how different networks allow several applications to scale to different numbers of PUs. For IO-bound applications (BlackScholes and TPCHQ6), performance does not scale with additional compute and on-chip memory resources. However, the performance of compute-bound applications (GEMM and SGD) improves with increased resources, but plateaus at a level that is determined by on-chip network bandwidth. This creates a trade-off in accelerator design between highly vectorized compute PUs with a small network—which would be underutilized for non-vectorized problems—and smaller compute PUs with limited performance due to network overhead. For more finely grained compute PUs, both more switches and more costly (higher-radix) switches must be employed to meet application requirements.

The scaling of time-scheduled accelerators (bottom row) is much less dramatic than that of deeply pipelined architectures (top row). Although communication between PUs in these architectures is less frequent, the scheduled architecture must use additional parallelization to match the

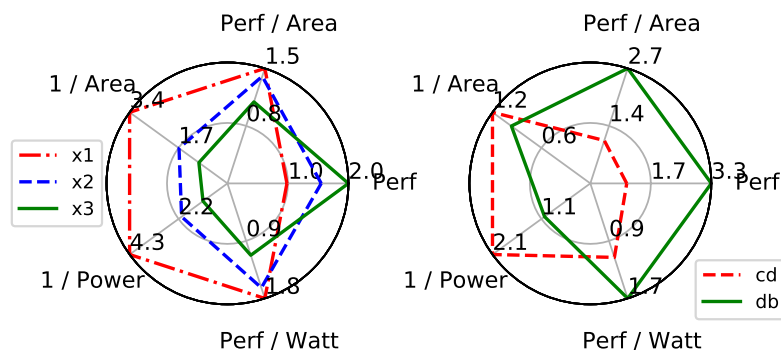


Figure 4.9: Impact of bandwidth and flow control strategies in switches.

throughput of the pipelined architecture; this translates to larger network sizes.

For pipelined architectures, both hybrid and static networks provide similar scaling with the same static bandwidth: the additional bandwidth from the dynamic network in hybrid networks does not provide additional scaling. This is mostly due to a bandwidth bottleneck between a PU and its router, which prevents the PU from requesting multiple elements per cycle. Hybrid networks tend to provide better scaling for time-scheduled architectures; multiple streams can be time-multiplexed at each ejection port without losing performance.

Bandwidth and flow control in switches

In this section, we study the impact of static network bandwidth and flow control mechanisms (per-hop vs. end-to-end credit-based). On the left side of Figure 4.9, we show that increased static bandwidth results in a linear performance increase and a superlinear increase in area and power. As shown in Section 4.1.4, any increase in accelerator size must be coupled with increased network bandwidth to effectively scale performance. This indicates that network overhead will increase with the size of an accelerator.

The right side of Figure 4.9 shows that, although credit-based flow control reduces the amount of buffering in switches and decreases network area and energy, application performance is significantly impacted. This is the result of imbalanced data-flow pipelines in the program: when there are parallel long and short paths over the network, there must be sufficient buffer space on the short path equal to the product of throughput and the difference in latency. Because performance is our most important metric, credit-based flow control is not feasible, especially because the impact of bubbles increases with communication distance, and therefore network size.

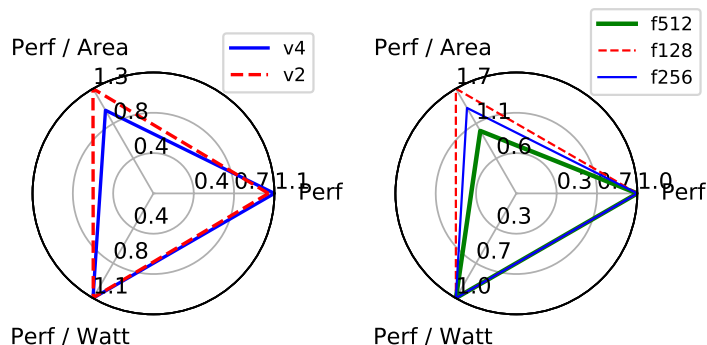


Figure 4.10: Impact of VC count and flit widths in routers.

VC count and reduced flit width in routers

In this experiment, we study the area-energy-performance trade-off between routers with different VC counts. As shown in Section 4.1.3, using many VCs increases both network area and energy. However, using too few VCs may force roundabout routing on the dynamic network or result in VC allocation failure when the network is heavily utilized. Nonetheless, the left side of Figure 4.10 shows minimal performance improvement from using more VCs.

Therefore, for each network design, we use a VC count equal to the maximum number of VCs required to map all applications to that network. Figure 4.11 shows that the best hybrid network configurations with 2x and 3x static bandwidth require at most 2 VCs, whereas the pure dynamic network requires 4 VCs to map all applications. Because dynamic network communication is infrequent, hybrid networks with fewer VCs provide both better energy and area efficiency than networks with more VCs, even though this constrains routing on the dynamic network.

We also explore the effects of reducing dynamic network bandwidth by using smaller routers; as shown in Section 4.1.3, routers with smaller flits have a much smaller area. Ideally, we could scale static network bandwidth while using a low-bandwidth router to provide an escape path and reduce overall area and energy overhead. The right side of Figure 4.10 shows that, for a hybrid network, reducing flit width improves area efficiency with minimal performance loss.

Static vs. hybrid vs. dynamic networks

Figure 4.12 shows the normalized performance for each application running on several network configurations. For some applications, the bar for S-x1 is missing; this indicates that place and route failed for all unrolling factors. For DRAM-bound applications, the performance variation between

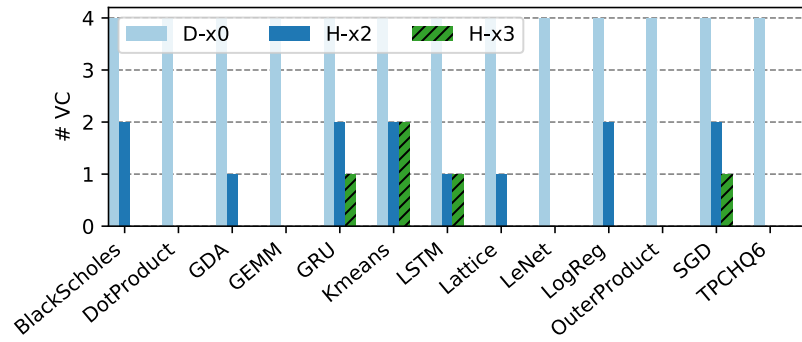


Figure 4.11: Number of VCs required for dynamic and hybrid networks. (No VCs indicates that all traffic is mapped to the static network.)

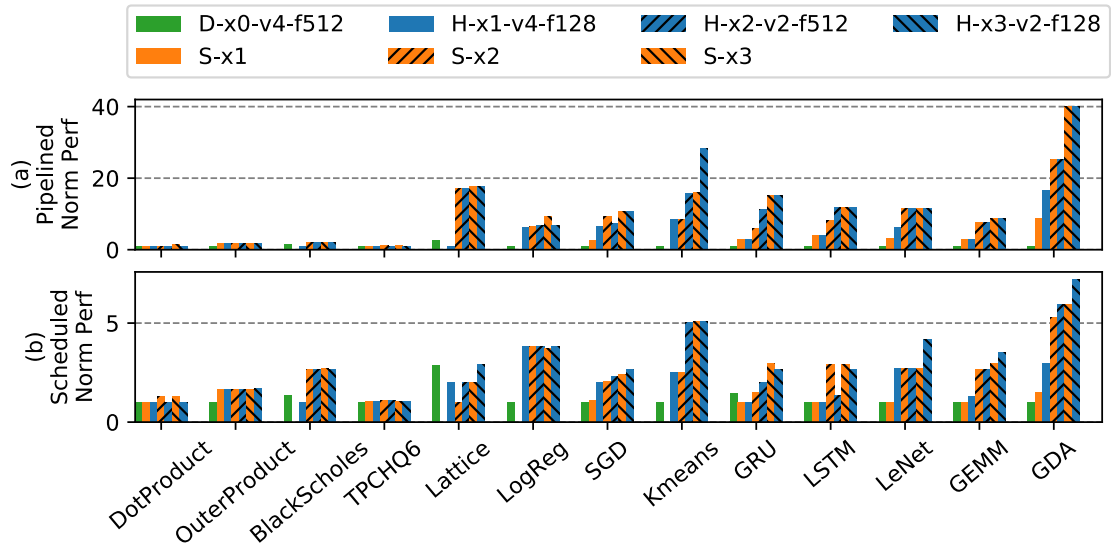


Figure 4.12: Normalized performance for different network configurations.

different networks is trivial because only a small fraction of the network is being used. In a few cases (Kmeans and GDA), hybrid networks provide better performance due to slightly increased bandwidth. For compute-bound applications, performance primarily correlates with network bandwidth because more bandwidth permits a higher parallelization factor.

The highest bandwidth static network uses the most PUs, as shown in Figures 4.13(b,e), because it permits more parallelization. It also has more data movement, as shown in (c,f), because PUs can be distributed farther apart. Due to bandwidth limitations, low-bandwidth networks perform best with small unrolling factors—they are unable to support the bisection bandwidth of larger program graphs. This is evident in Figures 4.13(b,e), where networks D-x0-v4-f512 and S-x2 have small PU utilizations.

With the same static bandwidth, most hybrid networks have better energy efficiency than the corresponding pure static networks, even though routers take more energy than switches to transmit the same amount of data. This is a result of allowing a small amount of traffic to escape onto the dynamic network: with the dynamic network as a safety net, static PaR tends to converge to better placements with less overall communication. This can be seen in Figures 4.13(c,f), where most static networks have larger hop counts than the corresponding hybrid network; hop count is the sum of all runtime link traversals, normalized per-application to the network configuration with the most hops. Subplots (e,f) show that more PUs are utilized with static networks than hybrid networks. This is because the compiler imposes less stringent IO constraints on PUs when partitioning for the hybrid network (as explained in Section ??), which results in fewer PUs, less data movement, and greater energy efficiency for hybrid networks.

In Figure 4.14, we summarize the best perf/watt and perf/area (among network configurations with <10% performance overhead) for pipelined and scheduled CGRA architectures. Pure dynamic networks are not shown because they perform poorly due to insufficient bandwidth. On the pipelined CGRA, the best hybrid network provides a 6.4x performance increase, 2.3x better energy efficiency, and a 6.9x perf/area increase over the worst network configuration. The best static network provides 7x better performance, 1.2x better energy efficiency, and 6.3x better perf/area. The hybrid network gives the best perf/area and perf/watt, with a small degradation in performance when compared to the static network. On the time-scheduled CGRA, both static and hybrid networks have an 8.6x performance improvement. The hybrid network gives a higher perf/watt improvement at 2.2x, whereas the static network gives a higher perf/area improvement at 2.6x. Overall, the hybrid networks deliver better energy efficiency with shorter routing distances by allowing an

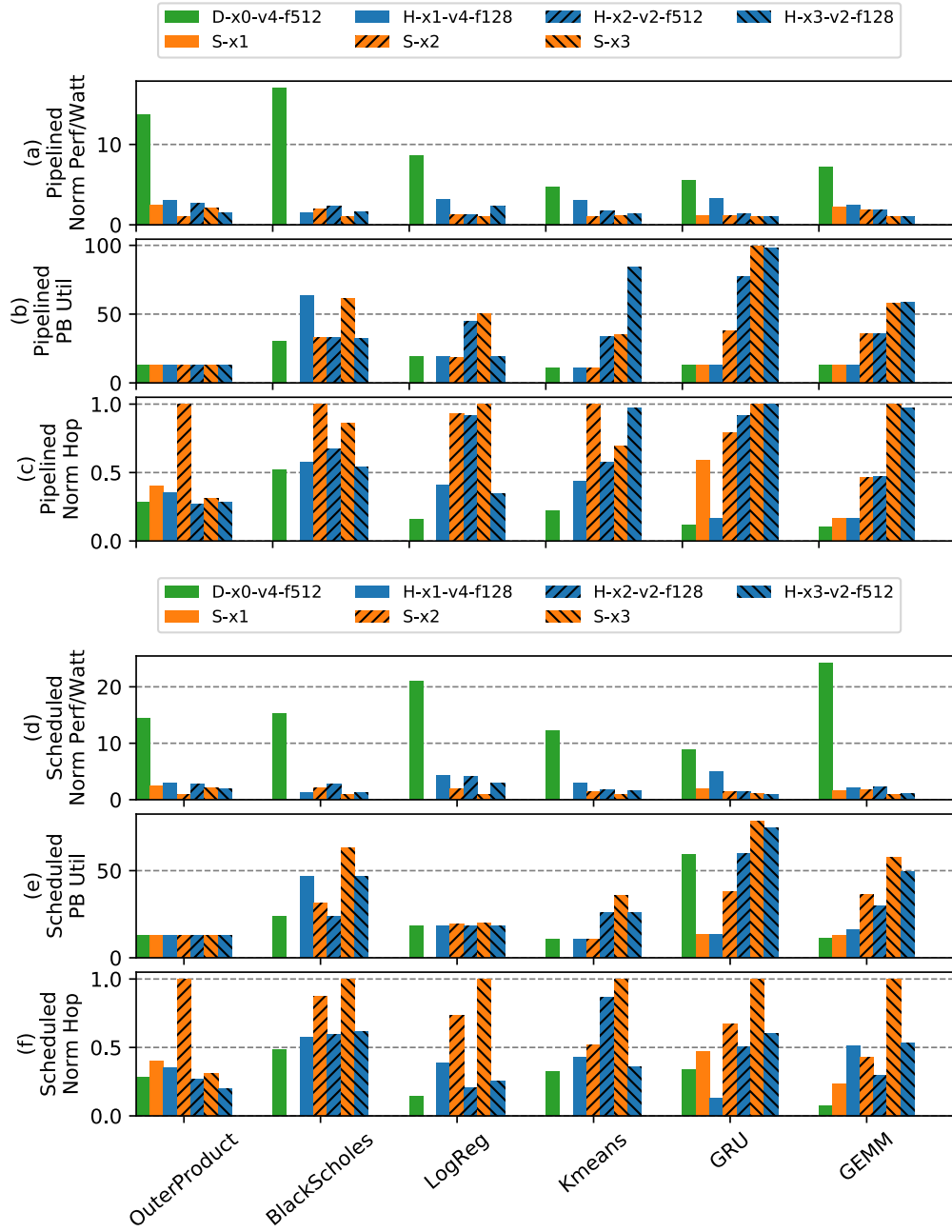


Figure 4.13: (a,d): Normalized performance/watt. (b,e): Percentage of compute and memory PUs utilized for each network configuration. (c,f): Total data movement (hop count).

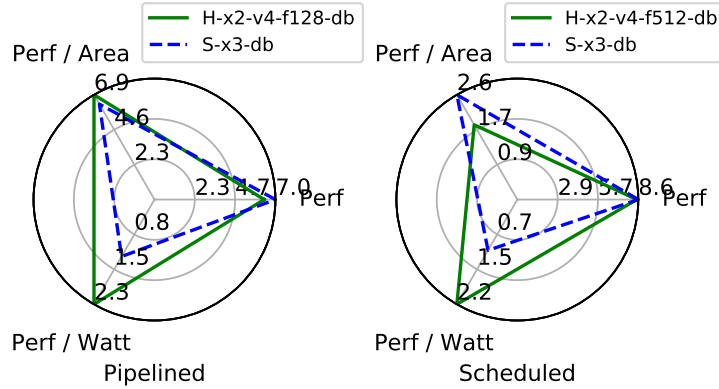


Figure 4.14: Geometric mean improvement for the best network configurations, relative to the worst configuration.

escape path on the dynamic network.

4.2 Plasticine Specialization for RNN Serving

Low-precision inference are commonly used to reduce the memory footprint of deep learning models and increase the compute density of hardware accelerators. Plasticine, however, only supports 32-bit operations and datapath. Using RNN serving as a motivating example, this section discusses the necessary architecture augmentation and specialization needed to efficiently map real-time inference on Plasticine.

Recurrent Neural Networks (RNNs) are a class of sequence models that play a key role in low-latency, AI-powered services in datacenters [95, 96]. These applications have stringent tail latency requirements, within the window of milliseconds, for real-time human-computer interactions. An example of such workloads is Google Translate; inference runs concurrently as the user types. Efficient acceleration of RNN requires flexible datapath to support global optimizations beyond BLAS kernels, which is where dataflow accelerators like Plasticine would shine. To meet this low-latency requirement, the other prerequisite is that everything has to stay on-chip. To achieve this, we introduce limited mixed-precision support with changes localized to PCUs. The enhancement supports the commonly used precision in machine learning without introducing massive overhead from fine-grained reconfigurability.

In the rest of this section, Section 4.2.1 proposes the necessary micro-architectural changes to support low-precision arithmetics on Plasticine. Section 4.2.2 introduces a folded reduction structure in

PCU's SIMD that improves the function unit (FU) utilization. Section 4.2.3 discusses architectural parameter selection for Plasticine to serve RNN applications efficiently.

4.2.1 Mixed-Precision Support

Previous works [95, 96] have shown that low-precision inference can deliver promising performance improvements without sacrificing accuracy. In the context of reconfigurable architectures such as FPGAs, low-precision inference not only increases compute density, but also reduces the required on-chip capacity for storing weights and intermediate data.

To support low-precision arithmetics without sacrificing coarse-grained reconfigurability, we introduce two low-precision struct types in Spatial: a tuple of 4 8-bit and 2 16-bit floating-point numbers, `4-float8` and `2-float16` respectively. Both types pack multiple low-precision values into single-precision storage. We support only 8 and 16-bit precisions, which are commonly seen in deep learning inference hardware. Users can only access values that are 32-bit aligned. This constraint guarantees that the microarchitectural change is only local to the PCU. PMU and DRAM access granularity remains intact from the original design.

Figure 4.15 (a) shows the original SIMD pipeline in a Plasticine PCU. Each FU supports both floating-point and fix-point operations. When mapping applications on Plasticine, the innermost loop body is vectorized across the lanes of the SIMD pipeline, and different operations of the loop body are mapped to different stages. Each pipeline stage contains a few pipeline registers (PRs) that allow propagation of live variables across stages. Special cross-lane connections as shown in red in Figure 4.15 enable reduction operations. To support 8-bit element-wise multiplication and 16-bit reduction, we add 4 opcodes to the FU, shown in Figure 4.15 (b). The 1st and 3rd stages are element-wise, low-precision operations that multiply and add 4 8-bit and 2 16-bit values, respectively. The 2nd and 4th stages rearrange low-precision values into two registers, and then pad them to higher precisions. The 5th stage reduces the two 32-bit value to a single 32-bit value using the existing add operation. From here, we can use the original reduction network shown in Figure 4.15 (a) to complete the remaining reduction and accumulates in a 32-bit connection.

With 4 lanes and 5 stages, a PCU first reads 16 8-bit values, performs 8-bit multiplication followed by rearrangement and padding, and then produce 16 16-bit values after the second stage. The intermediate values are stored in 2 PRs per lane. Next, 16 16-bit values are reduced to 8 16-bit values and then rearranged to 8 32-bit value in 2 PRs per lane. Then, the element-wise addition in a 32-bit value reduces the two registers in each line into 4 32-bit values. These values are fed through

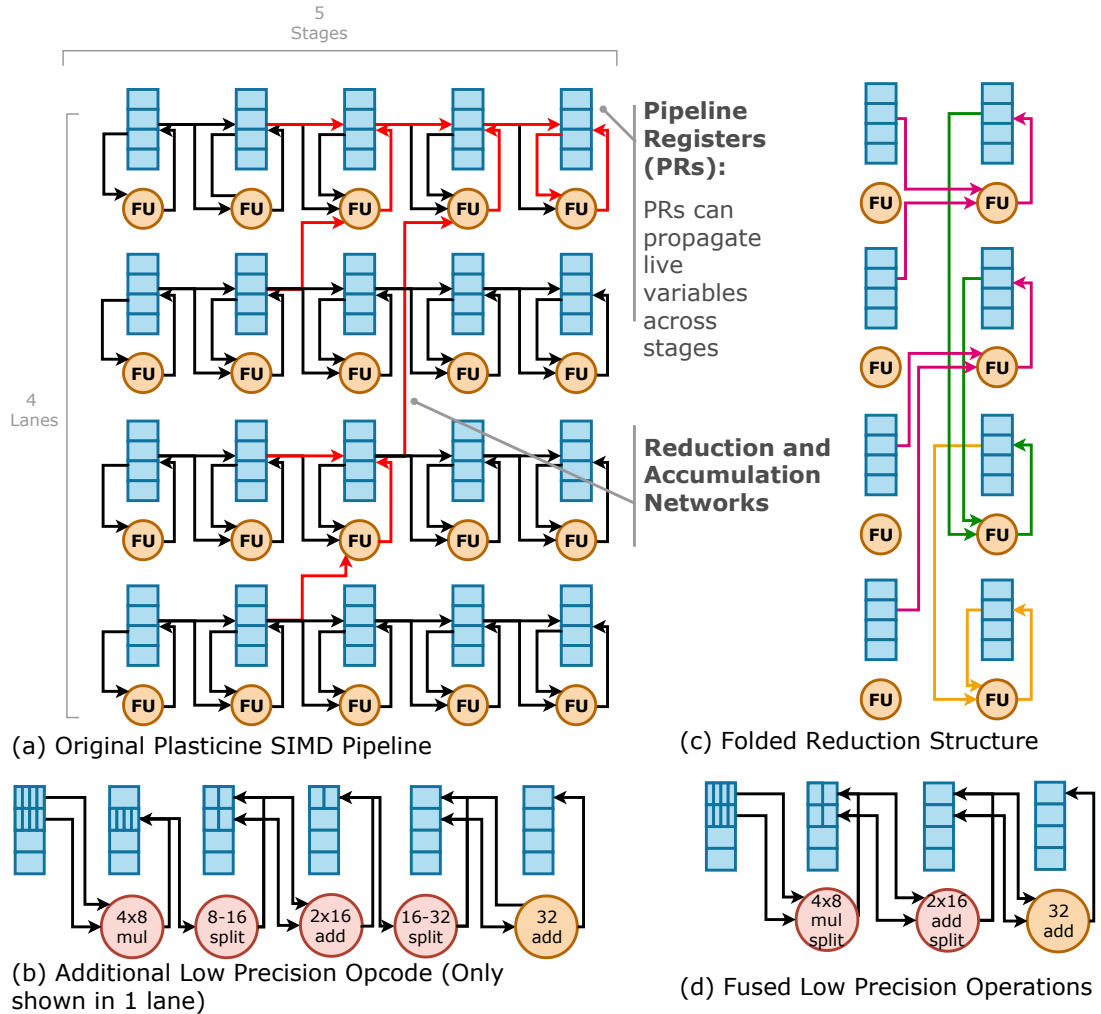


Figure 4.15: Plasticine PCU SIMD pipeline and low-precision support. Red circles are the new operations. Yellow circles are the original operations in Plasticine. In (d) the first stage is fused 1st, 2nd stages, and the second stage is fused 3rd, 4th stages of (b).

the reduction network that completes the remaining reduction and accumulation in two plus one stages.

In a more aggressive specialization, we can fuse the multiply and rearrange into the same stage. We also fuse the first low-precision reduction with the next rearrange shown in Figure 4.15 (d). In this way, we can perform the entire low-precision map-reduce in 2 stages in addition to the original full precision reduction. To maximize hardware reuse, we assume that it is possible to construct a full precision FU using low-precision FUs.

4.2.2 Folded Reduction Tree

The original reduction tree shown in Figure 4.15 (a) requires $\log_2(\#_{LANE}) + 1$ number of stages for $\#_{LANE}$ operations, leading to a low utilization of the FUs in the SIMD pipeline. The reduction tree also restricts the SIMD pipeline to have at least $\log_2(\#_{LANE}) + 1$ stages to compute a full reduction within a PCU. For a SIMD pipeline with 16 lanes, the FU utilization is only 53.33% in reduction stages.

To improve FU utilization, we introduce a folded reduction structure that performs $\#_{LANE}$ operations entirely within a single-stage pipelined over multiple cycles. Figure 4.15 (c) shows the folded reduce-accumulate structure. Instead of feeding the output of the reduction operation to the next stage, this structure folds the output back to the PR of the next unused FU in the same stage. The entire reduction plus accumulation is still fully pipelined in $\log_2(\#_{LANE}) + 1$ cycles with no structural hazard. In the original design, only a single register among the PRs have the special reduction tree connection. The downside of the folded structure is that no other live variables can be propagated after this reduction stage, unless adding multiple folded trees. In applications, we rarely see the need for multiple live variables after the reduction operation other than the accumulator itself. Therefore, it makes the most sense to put the folded reduction tree only in the last stage of a SIMD pipeline.

With the fused reduced-precision operations and the folded reduction tree, a PCU is able to perform 64 8-bit map-reduce in 4 stages, and 32 16-bit map-reduce in 3 stages. All the operations are still completed in $2 + \log_2(\#_{LANE})$ cycles, i.e. 8, 7, and 6 cycles for 8-, 16-, and 32-bit operations, respectively.

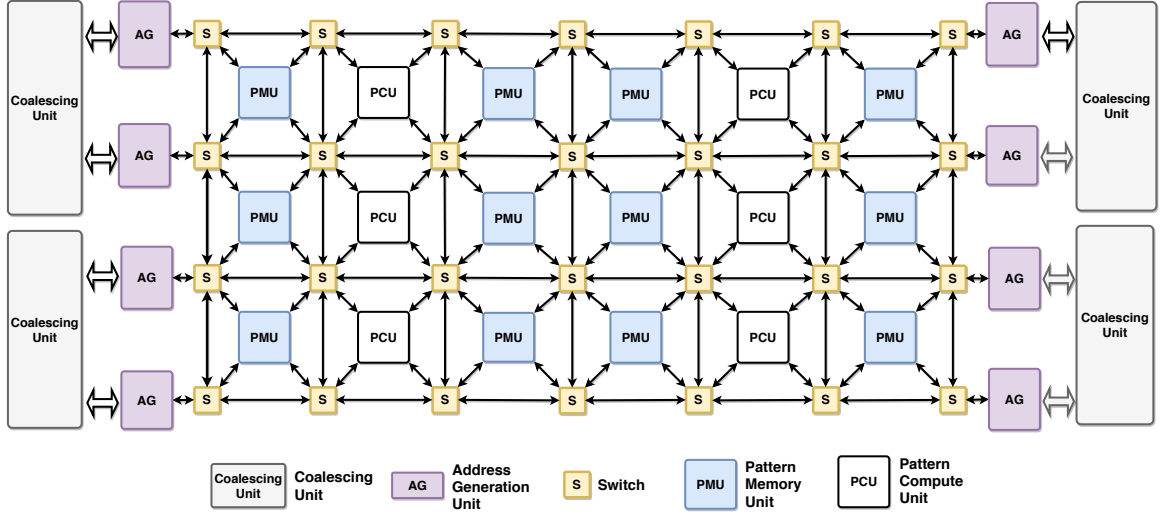


Figure 4.16: Variant Plasticine configuration for RNN serving with 2:1 ratio for PMU and PCU

4.2.3 Sizing Plasticine for RNN Serving

Evaluating an RNN cell containing N hidden units and N input features requires $2N^2$ computations and $N^2 + N$ memory reads. With a large N , the compute to memory ratio is 2:1. The original Plasticine architecture uses a checkerboard layout with a 1 to 1 ratio between PCU and PMU. A PCU has six stages and 16 lanes, and a PMU has 16 banks, which gives a 6:1 ratio between compute resource and on-chip memory read bandwidth. As a result of this layout, on-chip memory read bandwidth becomes the bottleneck for accelerating RNN serving applications. Figure 4.16 shows a specialized Plasticine configuration for RNN serving and general machine learning. Specifically, we choose a 2 to 1 PMU-PCU ratio with 4 stages in each PCU.

4.3 Generic Banknig Support

To support Spatial’s generic banking scheme for on-chip SRAM mentioned in Section 3.3.2, we need to introduce a few microarchitectural changes to Plasticine’s datapath. As a background, the static banking analysis searches an address remapping scheme, which we refer as banking scheme, that remaps the logical address space into partitions, such that addresses accessed in parallel belongs to different partitions. Each partition corresponds to a physical memory bank that comes with an additional address port. If the compiler can find a way to partition the data such that all parallel

workers are sent to different physical banks, the memory is guaranteed to feed all workers at full-throughput, which scales performance with parallelism. To search for the banking schemes, the compiler analyzes the access patterns of the memory. An access pattern refers to groups of address expressions that can be queried concurrently. For each access pattern, there can be multiple banking schemes that can sustain the access bandwidth. However, not all schemes cost the same amount of resources; some are more expensive in logic, and others are more expensive in memory.

The partitioning logic are two parametrizable equations for the bank address (BA) and the bank offset (BO), both are functions of the original user-requested logical addresses. BA is an expression that selects the partition for each request, and BO is the offset within the partition. The banking analyzer searches the parameters in the equation and injects these computations to memory accesses in the program. The first addition to Plasticine is the ability to compute independent BI and BO for each vector lanes. We introduce vectorized address pipelines in PMUs, which was originally a scalar pipeline.

The generic banking scheme also requires a full crossbar datapath between all access lanes (workers) and physical banks. To support this, we introduce the shuffle operator shown in Figure 3.16. The shuffle operator takes three vectorized operands: a from address \vec{FA} , a to address \vec{TA} , and a base \vec{B} . The output O equals to B shuffled from the order specified in \vec{FA} to the order in \vec{TA} . Specifically,

$$\forall i \in [1, |\vec{O}|], O_i = \begin{cases} B_j, & \text{if } \exists TA_i = FA_j \\ 0 & \end{cases} \quad (4.3)$$

$$|\vec{FA}| = |\vec{B}| \quad (4.4)$$

$$|\vec{TA}| = |\vec{O}| \quad (4.5)$$

Figure 3.16 gives concrete examples of inputs and outputs of the shuffle operator. On the requester side, BA is the FA , a vector constant corresponding to the banks statically assigned to the partition is the TA , and BO is the base. On the receiver side, the vector constant is the FA , the BA is the TA , and the data response D is the base. If a lane in the TA is not in the FA , the corresponding lane is marked as invalid with value 0. As a result, an invalid vector ORed with a valid vector is still the valid vector. For address going to the scratchpad, we use the first bit of the 32-bit address as the predication bit of the access; if the predicate bit equals zero, the request will be dropped by the memory. Therefore, a valid address zero would be xF0000000 to distinguish with an invalid address

0.

The shuffle operator can be expensive, requiring a 16×16 32-bit crossbar with parallel integer comparators for equality. We expect to add one or two shuffle operators per SIMD pipeline, and the operator can be pipelined across multiple stages to meet timing.

Chapter 5

Related Work

5.1 Compiler

5.1.1 Streaming Dataflow IRs

Although many works claim to emit efficient and information-rich dataflow IRs for the downstream compilers, very few of them can capture the high-level parallel patterns and implementation details that are critical to RDA mappings. For example, TensorFlow [61] emits dataflow IR composed of tensor operations. However, its IR lacks information on the parallel patterns within these operations. In contrast, most of the streaming languages [97, 98, 99] are not able to extract nested loop-level parallelism from modern data-intensive applications. For example, StreamIt [97], a language tailored for streaming computing, also adopts distributed control as in SARA. However, it lacks the necessary language features to describe deeply and irregularly nested loops that are common in modern data-intensive applications.

Kernel-specific loop optimization

- Tangram: Optimized Coarse-Grained Dataflow for Scalable NN Accelerators

5.1.2 Spatial Compilers

Most previous works [100, 28] only consider allocating resources at the same level. SARA takes a more general assumption by co-allocating resources at multiple levels of an accelerator’s hierarchy.

5.2 On-Chip Network

Multiple decades of research have resulted in a rich body of literature, both in CGRAs [86, 87] and on-chip networks [101]. We discuss relevant prior work under the following categories:

5.2.1 Tiled Processor Interconnects

Architectures such as Raw [102] and Tile [103] use scalar operand networks [104], which combine static and dynamic networks. Raw has one static and two dynamic interconnects: the static interconnect is used to route normal operand traffic, one dynamic network is used to route runtime-dependent values which could not be routed on the static network, and the second dynamic network is used for cache misses and other exceptions. Deadlock avoidance is guaranteed only in the second dynamic network, which is used to recover from deadlocks in the first dynamic network. However, as described in Section ??, wider buses and larger flit sizes create scalability issues with two dynamic networks, including higher area and power. In addition, our static VC allocation scheme ensures deadlock freedom in our single dynamic network, obviating the need for deadlock recovery. The dynamic Raw network also does not preserve operand ordering, requiring an operand reordering mechanism at every tile.

TRIPS [105] is a tiled dataflow architecture with dynamic execution. TRIPS does not have a static interconnect, but contains two dynamic networks [106]: an operand network to route operands between tiles, and an on-chip network to communicate with cache banks. Wavescalar [107] is another tiled dataflow architecture with four levels of hierarchy, connected by dynamic interconnects that vary in topology and bandwidth at each level. The Polymorphic Pipeline Array [108] is a tiled architecture built to target mobile multimedia applications. While compute resources are either statically or dynamically provisioned via hardware virtualization support, communication uses a dynamic scalar operand network.

5.2.2 CGRA Interconnects

Many previously proposed CGRAs use a word-level static interconnect, which has better compute density than bit-based routing [109]. CGRAs such as HRL [41], DySER [38], and Elastic CGRAs [110] commonly employ two static interconnects: a word-level interconnect to route data and a bit-level interconnect to route control signals. Several works have also proposed a statically scheduled interconnect [111, 112, 113] using a modulo schedule. While this approach is effective

for inner loops with predictable latencies and fixed initiation intervals, variable latency operations and hierarchical loop nests add scheduling complexity that prevents a single modulo schedule. HyCube [42] has a similar statically scheduled network, with the ability to bypass intermediate switches in the same cycle. This allows operands to travel multiple hops in a single cycle, but creates long wires and combinational paths and adversely affects the clock period and scalability.

5.2.3 Design Space Studies

Several prior studies focus on tradeoffs with various network topologies, but do not characterize or quantify the role of dynamism in interconnects. The Raw design space study [114] uses an analytical model for applications as well as architectural resources to perform a sensitivity analysis of compute and memory resources focused on area, power, and performance, without varying the interconnect. The ADRES design space study [115] focuses on area and energy tradeoffs with different network topologies with the ADRES [36] architecture, where all topologies use a fully static interconnect. KressArray Xplorer [116] similarly explores topology tradeoffs with the KressArray [37] architecture. Other studies explore topologies for mesh-based CGRAs [117] and more general CGRAs supporting resource sharing [118]. Other tools like Sunmap [119] allow end users to construct and explore various topologies.

5.2.4 Compiler Driven NoCs (WIP)

Other prior works have used compiler techniques to optimize various facets of NoCs. Some studies have explored statically allocating virtual channels [120, 121] to multiple concurrent flows to mitigate head-of-line blocking. These studies propose an approach to derive deadlock-free allocations based on the turn model [122]. While our approach also statically allocates VCs, our method to guarantee deadlock freedom differs from the aforementioned study as it does not rely on the turn model. Ozturk et al. [123] propose a scheme to increase the reliability of NoCs for chip multiprocessors by sending packets over multiple links. Their approach uses integer linear programming to balance the total number of links activated (an energy-based metric) against the amount of packet duplication (reliability). Ababei et al. [124] use a static placement algorithm and an estimate of reliability to attempt to guide placement decisions for NoCs. Kapre et al. [125] develop a workflow to map applications to CGRAs using several transformations, including efficient multicast routing and node splitting, but do not consider optimizations such as non-minimal routing.

Chapter 6

Conclusions

Reconfigurable dataflow accelerators (RDAs) are a promising class of spatial accelerators that deliver higher performance-to-resource efficiency while capturing a large application space. However, to sustain these benefits of RDAs, scalability must be taken into account in both software stack and hardware design.

On the software side, we address this challenge with a distributed asynchronous control scheme and resource virtualization. We develop a compiler, SARA, that constructs a dataflow graph from an imperative program with nested control flow. We use dataflow tokens to enforce sequential consistency across distributed memory accesses, incurring minimum communication overhead. Furthermore, SARA efficiently compose distributed resource to provide an abstraction of larger logical resources to the programmers. Our optimizations effectively eliminate pipelining stall and reduce resource fragmentation with heterogeneous resources. Lastly, our evaluation shows that SARA achieves of 2x averaged speedup and an area-normalized 4x speedup over a Tesla V100 GPU.

On the hardware side, we show that the best network design depends on both applications and the underlying accelerator architecture. Network performance correlates strongly with bandwidth for streaming accelerators, and scaling raw bandwidth is more area- and energy-efficient with a static network. We show that the application mapping can be optimized to move less data by using a dynamic network as a fallback from a high-bandwidth static network. This static-dynamic hybrid network provides a 1.8x energy-efficiency and 2.8x performance advantage over the purely static and purely dynamic networks, respectively.

Bibliography

- [1] C. Kachris and D. Soudris, “A survey on reconfigurable accelerators for cloud computing,” in *2016 26th International Conference on Field Programmable Logic and Applications (FPL)*, pp. 1–10, 2016.
- [2] S. Sarkar, T. Majumder, A. Kalyanaraman, and P. P. Pande, “Hardware accelerators for bio-computing: A survey,” in *Proceedings of 2010 IEEE International Symposium on Circuits and Systems*, pp. 3789–3792, 2010.
- [3] S. Aluru and N. Jammula, “A review of hardware acceleration for computational genomics,” *IEEE Design Test*, vol. 31, no. 1, pp. 19–30, 2014.
- [4] H. Esmailzadeh, E. Blem, R. S. Amant, K. Sankaralingam, and D. Burger, “Dark silicon and the end of multicore scaling,” in *2011 38th Annual International Symposium on Computer Architecture (ISCA)*, pp. 365–376, June 2011.
- [5] R. Prabhakar, Y. Zhang, D. Koeplinger, M. Feldman, T. Zhao, S. Hadjis, A. Pedram, C. Kozyrakis, and K. Olukotun, “Plasticine: A reconfigurable architecture for parallel patterns,” in *Proceedings of the 44th Annual International Symposium on Computer Architecture, ISCA ’17*, (New York, NY, USA), pp. 389–402, ACM, 2017.
- [6] Y. Zhang, A. Rucker, M. Vilim, R. Prabhakar, W. Hwang, and K. Olukotun, “Scalable interconnects for reconfigurable spatial architectures,” in *Proceedings of the 46th International Symposium on Computer Architecture, ISCA ’19*, (New York, NY, USA), pp. 615–628, ACM, 2019.
- [7] M. Grant, *Disciplined Convex Programming*. Phd. thesis, Stanford University, 2014.
- [8] Y. Y. Michael Grant, Steven Boyd, “Disciplined convex programming,” 2019.

- [9] R. H. Dennard, F. H. Gaensslen, V. L. Rideout, E. Bassous, and A. R. LeBlanc, "Design of ion-implanted mosfet's with very small physical dimensions," *IEEE Journal of Solid-State Circuits*, vol. 9, pp. 256–268, Oct 1974.
- [10] R. Hameed, W. Qadeer, M. Wachs, O. Azizi, A. Solomatnikov, B. C. Lee, S. Richardson, C. Kozyrakis, and M. Horowitz, "Understanding sources of inefficiency in general-purpose chips," in *Proceedings of the 37th Annual International Symposium on Computer Architecture, ISCA '10*, (New York, NY, USA), pp. 37–47, ACM, 2010.
- [11] M. A. Dias and D. A. P. Ferreira, "Deep learning in reconfigurable hardware: A survey," in *2019 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pp. 95–98, 2019.
- [12] A. Gielata, P. Russek, and K. Wiatr, "Aes hardware implementation in fpga for algorithm acceleration purpose," in *2008 International Conference on Signals and Electronic Systems*, pp. 137–140, 2008.
- [13] J. L. Hennessy and D. A. Patterson, "A new golden age for computer architecture," *Commun. ACM*, vol. 62, pp. 48–60, Jan. 2019.
- [14] Y. Chen, T. Luo, S. Liu, S. Zhang, L. He, J. Wang, L. Li, T. Chen, Z. Xu, N. Sun, and O. Temam, "Dadiannao: A machine-learning supercomputer," in *Proceedings of the 47th Annual IEEE/ACM International Symposium on Microarchitecture, MICRO-47*, (Washington, DC, USA), pp. 609–622, IEEE Computer Society, 2014.
- [15] N. P. Jouppi, C. Young, N. Patil, D. Patterson, G. Agrawal, R. Bajwa, S. Bates, S. Bhatia, N. Boden, A. Borchers, R. Boyle, P.-I. Cantin, C. Chao, C. Clark, J. Coriell, M. Daley, M. Dau, J. Dean, B. Gelb, T. V. Ghaemmamghami, R. Gottipati, W. Gulland, R. Hagmann, C. R. Ho, D. Hogberg, J. Hu, R. Hundt, D. Hurt, J. Ibarz, A. Jaffey, A. Jaworski, A. Kaplan, H. Khaitan, D. Killebrew, A. Koch, N. Kumar, S. Lacy, J. Laudon, J. Law, D. Le, C. Leary, Z. Liu, K. Lucke, A. Lundin, G. MacKean, A. Maggiore, M. Mahony, K. Miller, R. Nagarajan, R. Narayanaswami, R. Ni, K. Nix, T. Norrie, M. Omernick, N. Penukonda, A. Phelps, J. Ross, M. Ross, A. Salek, E. Samadiani, C. Severn, G. Sizikov, M. Snelham, J. Souter, D. Steinberg, A. Swing, M. Tan, G. Thorson, B. Tian, H. Toma, E. Tuttle, V. Vasudevan, R. Walter, W. Wang, E. Wilcox, and D. H. Yoon, "In-datacenter performance analysis of a tensor processing unit," *SIGARCH Comput. Archit. News*, vol. 45, pp. 1–12, June 2017.

- [16] S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. A. Horowitz, and W. J. Dally, "Eie: Efficient inference engine on compressed deep neural network," in *Proceedings of the 43rd International Symposium on Computer Architecture, ISCA '16*, (Piscataway, NJ, USA), pp. 243–254, IEEE Press, 2016.
- [17] Y.-H. Chen, T. Krishna, J. S. Emer, and V. Sze, "Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks," *IEEE Journal of Solid-State Circuits*, vol. 52, no. 1, pp. 127–138, 2017.
- [18] M. Gao, X. Yang, J. Pu, M. Horowitz, and C. Kozyrakis, "Tangram: Optimized coarse-grained dataflow for scalable nn accelerators," in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS '19*, (New York, NY, USA), p. 807–820, Association for Computing Machinery, 2019.
- [19] F. Akopyan, J. Sawada, A. Cassidy, R. Alvarez-Icaza, J. Arthur, P. Merolla, N. Imam, Y. Nakamura, P. Datta, G. Nam, B. Taba, M. Beakes, B. Brezzo, J. B. Kuang, R. Manohar, W. P. Risk, B. Jackson, and D. S. Modha, "Truenorth: Design and tool flow of a 65 mw 1 million neuron programmable neurosynaptic chip," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 34, no. 10, pp. 1537–1557, 2015.
- [20] K. Rupp, "www.karlrupp.net."
- [21] J. W. Richardson, A. D. George, and H. Lam, "Performance analysis of gpu accelerators with realizable utilization of computational density," in *2012 Symposium on Application Accelerators in High Performance Computing*, pp. 137–140, 2012.
- [22] J. Hestness, S. W. Keckler, and D. A. Wood, "Gpu computing pipeline inefficiencies and optimization opportunities in heterogeneous cpu-gpu processors," in *2015 IEEE International Symposium on Workload Characterization*, pp. 87–97, 2015.
- [23] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2980–2988, 2017.
- [24] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 834–848, 2018.

- [25] C. Guo, Y. Zhou, J. Leng, Y. Zhu, Z. Du, Q. Chen, C. Li, M. Guo, and B. Yao, "Balancing efficiency and flexibility for dnn acceleration via temporal gpu-systolic array integration," 2020.
- [26] B. H. Calhoun, J. F. Ryan, S. Khanna, M. Putic, and J. Lach, "Flexible circuits and architectures for ultralow power," *Proceedings of the IEEE*, vol. 98, pp. 267–282, Feb 2010.
- [27] K. K. W. Poon, S. J. E. Wilton, and A. Yan, "A detailed power model for field-programmable gate arrays," *ACM Trans. Des. Autom. Electron. Syst.*, vol. 10, pp. 279–302, Apr. 2005.
- [28] M. Budiu, G. Venkataramani, T. Chelcea, and S. C. Goldstein, "Spatial computation," in *Proceedings of the 11th International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS XI, (New York, NY, USA), pp. 14–26, ACM, 2004.
- [29] D. Koeplinger, R. Prabhakar, Y. Zhang, C. Delimitrou, C. Kozyrakis, and K. Olukotun, "Automatic generation of efficient accelerators for reconfigurable hardware," in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 115–127, June 2016.
- [30] I. Kuon, R. Tessier, and J. Rose, "Fpga architecture: Survey and challenges," *Foundations and Trends in Electronic Design Automation*, vol. 2, pp. 135–253, 01 2007.
- [31] A. Putnam, A. M. Caulfield, E. S. Chung, D. Chiou, K. Constantinides, J. Demme, H. Esmaeilzadeh, J. Fowers, G. P. Gopal, J. Gray, M. Haselman, S. Hauck, S. Heil, A. Hormati, J.-Y. Kim, S. Lanka, J. Larus, E. Peterson, S. Pope, A. Smith, J. Thong, P. Y. Xiao, and D. Burger, "A reconfigurable fabric for accelerating large-scale datacenter services," in *Proceeding of the 41st Annual International Symposium on Computer Architecture*, ISCA '14, (Piscataway, NJ, USA), pp. 13–24, IEEE Press, 2014.
- [32] J. Ouyang, S. Lin, W. Qi, Y. Wang, B. Yu, and S. Jiang, "Sda: Software-defined accelerator for largescale dnn systems," *Hot Chips 26*, 2014.
- [33] K. Guo, L. Sui, J. Qiu, S. Yao, S. Han, Y. Wang, and H. Yang, "From model to fpga: Software-hardware co-design for efficient neural network acceleration," in *2016 IEEE Hot Chips 28 Symposium (HCS)*, pp. 1–27, Aug 2016.
- [34] A. AWS, "Amazon ec2 f1 instances." <https://aws.amazon.com/ec2/instance-types/f1>.

- [35] I. Kuon, R. Tessier, and J. Rose, "Fpga architecture: Survey and challenges," *Found. Trends Electron. Des. Autom.*, vol. 2, pp. 135–253, Feb. 2008.
- [36] B. Mei, S. Vernalde, D. Verkest, H. De Man, and R. Lauwereins, "Adres: An architecture with tightly coupled vliw processor and coarse-grained reconfigurable matrix," in *Field Programmable Logic and Application* (P. Y. K. Cheung and G. A. Constantinides, eds.), (Berlin, Heidelberg), pp. 61–70, Springer Berlin Heidelberg, 2003.
- [37] R. Kress, "A fast reconfigurable alu for xputers," 1996.
- [38] V. Govindaraju, C.-H. Ho, and K. Sankaralingam, "Dynamically specialized datapaths for energy efficient computing," in *Proceedings of the 2011 IEEE 17th International Symposium on High Performance Computer Architecture*, HPCA '11, (Washington, DC, USA), pp. 503–514, IEEE Computer Society, 2011.
- [39] S. C. Goldstein, H. Schmit, M. Moe, M. Budiu, S. Cadambi, R. R. Taylor, and R. Laufer, "Piperench: a coprocessor for streaming multimedia acceleration," in *Proceedings of the 26th International Symposium on Computer Architecture* (Cat. No.99CB36367), pp. 28–39, May 1999.
- [40] M. Mishra, T. J. Callahan, T. Chelcea, G. Venkataramani, S. C. Goldstein, and M. Budiu, "Tartan: Evaluating spatial computation for whole program execution," in *Proceedings of the 12th International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS XII, (New York, NY, USA), pp. 163–174, ACM, 2006.
- [41] M. Gao and C. Kozyrakis, "Hrl: Efficient and flexible reconfigurable logic for near-data processing," in *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pp. 126–137, March 2016.
- [42] M. Karunaratne, A. K. Mohite, T. Mitra, and L.-S. Peh, "Hycube: A cgra with reconfigurable single-cycle multi-hop interconnect," in *Proceedings of the 54th Annual Design Automation Conference 2017*, DAC '17, (New York, NY, USA), pp. 45:1–45:6, ACM, 2017.
- [43] A. Parashar, M. Pellauer, M. Adler, B. Ahsan, N. Crago, D. Lustig, V. Pavlov, A. Zhai, M. Gambhir, A. Jaleel, R. Allmon, R. Rayess, S. Maresh, and J. Emer, "Triggered instructions: A control paradigm for spatially-programmed architectures," in *Proceedings of the 40th Annual International Symposium on Computer Architecture*, ISCA '13, (New York, NY, USA), pp. 142–153, ACM, 2013.

- [44] T. Nowatzki, V. Gangadhar, N. Ardalani, and K. Sankaralingam, "Stream-dataflow acceleration," in *2017 ACM/IEEE 44th Annual International Symposium on Computer Architecture (ISCA)*, pp. 416–429, June 2017.
- [45] C. Farabet, B. Martini, B. Corda, P. Akselrod, E. Culurciello, and Y. LeCun, "Neuflow: A runtime reconfigurable dataflow processor for vision," in *CVPR 2011 WORKSHOPS*, pp. 109–116, 2011.
- [46] M. Oh, C. Lee, S. Lee, Y. Seo, S. Kim, J. Wang, and C. Sungchung Park, "Convolutional neural network accelerator with reconfigurable dataflow," in *2018 International SoC Design Conference (ISOCC)*, pp. 42–43, 2018.
- [47] A. Niedermeier, J. Kuper, and G. Smit, "Dataflow-based reconfigurable architecture for streaming applications," in *2012 International Symposium on System on Chip (SoC)*, pp. 1–4, 2012.
- [48] M. Vilim, A. Rucker, Y. Zhang, S. Liu, and K. Olukotun, "Gorgon: Accelerating machine learning from relational data," in *Proceedings of the 47th International Symposium on Computer Architecture, ISCA '20*, 2020.
- [49] R. Singhal, Y. Zhang, J. Ullman, R. Prabhakar, and K. Olukotun, "Efficient multiway hash join on reconfigurable hardware," in *Technology Conference on Performance Evaluation and Benchmarking*, 05 2019.
- [50] R. Prabhakar, O. A. Olukotun, C. Kozyrakis, and C. Re. PhD thesis, 2018.
- [51] J. Weng, S. Liu, Z. Wang, V. Dadu, and T. Nowatzki, "A hybrid systolic-dataflow architecture for inductive matrix algorithms," in *2020 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pp. 703–716, 2020.
- [52] J. Noguera, C. Dick, V. Kathail, G. Singh, K. Vissers, and R. Wittig, "Xilinx project everest: 'hw/sw programmable engine'," *Hot Chips* 30, 2018.
- [53] J. Fowers, K. Ovtcharov, M. Papamichael, T. Massengill, M. Liu, D. Lo, S. Alkalay, M. Haselman, L. Adams, M. Ghandi, S. Heil, P. Patel, A. Sapek, G. Weisz, L. Woods, S. Lanka, S. K. Reinhardt, A. M. Caulfield, E. S. Chung, and D. Burger, "A configurable cloud-scale DNN processor for real-time AI," in *45th ACM/IEEE Annual International Symposium on Computer Architecture, ISCA 2018, Los Angeles, CA, USA, June 1-6, 2018*, pp. 1–14, 2018.

- [54] M. N. Bojnordi and E. Ipek, "Memristive boltzmann machine: A hardware accelerator for combinatorial optimization and deep learning," in *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pp. 1–13, 2016.
- [55] D. Liu, T. Chen, S. Liu, J. Zhou, S. Zhou, O. Teman, X. Feng, X. Zhou, and Y. Chen, "Pudian-nao: A polyvalent machine learning accelerator," in *Proceedings of the Twentieth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS '15*, (New York, NY, USA), p. 369–381, Association for Computing Machinery, 2015.
- [56] X. Liu, M. Mao, B. Liu, H. Li, Y. Chen, B. Li, Yu Wang, Hao Jiang, M. Barnell, Qing Wu, and Jianhua Yang, "Reno: A high-efficient reconfigurable neuromorphic computing accelerator design," in *2015 52nd ACM/EDAC/IEEE Design Automation Conference (DAC)*, pp. 1–6, 2015.
- [57] W. Qadeer, R. Hameed, O. Shacham, P. Venkatesan, C. Kozyrakis, and M. A. Horowitz, "Convolution engine: Balancing efficiency & flexibility in specialized computing," in *Proceedings of the 40th Annual International Symposium on Computer Architecture, ISCA '13*, (New York, NY, USA), p. 24–35, Association for Computing Machinery, 2013.
- [58] B. Reagen, P. Whatmough, R. Adolf, S. Rama, H. Lee, S. K. Lee, J. M. Hernández-Lobato, G. Wei, and D. Brooks, "Minerva: Enabling low-power, highly-accurate deep neural network accelerators," in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 267–278, 2016.
- [59] D. Shin, J. Lee, J. Lee, and H. Yoo, "14.2 dnpu: An 8.1tops/w reconfigurable cnn-rnn processor for general-purpose deep neural networks," in *2017 IEEE International Solid-State Circuits Conference (ISSCC)*, pp. 240–241, 2017.
- [60] L. Song, Y. Wang, Y. Han, X. Zhao, B. Liu, and X. Li, "C-brain: A deep learning accelerator that tames the diversity of cnns through adaptive data-level parallelization," in *2016 53rd ACM/EDAC/IEEE Design Automation Conference (DAC)*, pp. 1–6, 2016.
- [61] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zheng, "Tensorflow: A system for large-scale machine learning," in *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation, OSDI'16*, (Berkeley, CA, USA), pp. 265–283, USENIX Association, 2016.

- [62] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, “Caffe: Convolutional architecture for fast feature embedding,” *arXiv preprint arXiv:1408.5093*, 2014.
- [63] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems* 32 (H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, eds.), pp. 8024–8035, Curran Associates, Inc., 2019.
- [64] F. Kjolstad, S. Kamil, S. Chou, D. Lugato, and S. Amarasinghe, “The tensor algebra compiler,” *Proc. ACM Program. Lang.*, vol. 1, pp. 77:1–77:29, Oct. 2017.
- [65] W. Lin, D. Tsai, L. Tang, C. Hsieh, C. Chou, P. Chang, and L. Hsu, “Onnc: A compilation framework connecting onnx to proprietary deep learning accelerators,” in *2019 IEEE International Conference on Artificial Intelligence Circuits and Systems (AICAS)*, pp. 214–218, 2019.
- [66] N. Vasilache, O. Zinenko, T. Theodoridis, P. Goyal, Z. DeVito, W. S. Moses, S. Verdoolaege, A. Adams, and A. Cohen, “Tensor comprehensions: Framework-agnostic high-performance machine learning abstractions,” *CoRR*, vol. abs/1802.04730, 2018.
- [67] T. Chen, T. Moreau, Z. Jiang, L. Zheng, E. Yan, M. Cowan, H. Shen, L. Wang, Y. Hu, L. Ceze, C. Guestrin, and A. Krishnamurthy, “Tvm: An automated end-to-end optimizing compiler for deep learning,” in *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation, OSDI’18, (USA)*, p. 579–594, USENIX Association, 2018.
- [68] J. Ragan-Kelley, A. Adams, D. Sharlet, C. Barnes, S. Paris, M. Levoy, S. Amarasinghe, and F. Durand, “Halide: Decoupling algorithms from schedules for high-performance image processing,” *Commun. ACM*, vol. 61, p. 106–115, Dec. 2017.
- [69] S. Chou, F. Kjolstad, and S. Amarasinghe, “Format abstraction for sparse tensor algebra compilers,” *Proc. ACM Program. Lang.*, vol. 2, pp. 123:1–123:30, Oct. 2018.
- [70] D. Koeplinger, M. Feldman, R. Prabhakar, Y. Zhang, S. Hadjis, R. Fiszels, T. Zhao, L. Nardi, A. Pedram, C. Kozyrakis, and K. Olukotun, “Spatial: A language and compiler for application accelerators,” in *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2018, (New York, NY, USA)*, pp. 296–311, ACM, 2018.

- [71] T. Zhao, Y. Zhang, and K. Olukotun, "Serving recurrent neural networks efficiently with a spatial accelerator," in *Proceedings of the 2nd SysML Conference*, 2019.
- [72] J. von Neumann, "First draft of a report on the edvac," *IEEE Annals of the History of Computing*, vol. 15, no. 4, pp. 27–75, 1993.
- [73] D. Koeplinger, M. Feldman, R. Prabhakar, Y. Zhang, S. Hadjis, R. Fiszels, T. Zhao, L. Nardi, A. Pedram, C. Kozyrakis, and K. Olukotun, "Spatial: A language and compiler for application accelerators," in *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2018*, (New York, NY, USA), pp. 296–311, ACM, 2018.
- [74] K. J. Brown, A. K. Sujeeth, H. J. Lee, T. Rompf, H. Chafi, M. Odersky, and K. Olukotun, "A heterogeneous parallel framework for domain-specific languages," in *Proceedings of the 2011 International Conference on Parallel Architectures and Compilation Techniques, PACT '11*, (Washington, DC, USA), pp. 89–100, IEEE Computer Society, 2011.
- [75] "Vivado high-level synthesis." <http://www.xilinx.com/products/design-tools/vivado/integration/esl-design.html>, 2016.
- [76] Xilinx, "The xilinx sdaccel development environment." https://www.xilinx.com/publications/prod_mktg/sdx/sdaccel-backgrounder.pdf, 2014.
- [77] C. Lattner and V. Adve, "Llvm: A compilation framework for lifelong program analysis & transformation," in *Proceedings of the International Symposium on Code Generation and Optimization: Feedback-Directed and Runtime Optimization, CGO '04*, (USA), p. 75, IEEE Computer Society, 2004.
- [78] W. Dally and B. Towles, *Principles and Practices of Interconnection Networks*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2003.
- [79] A. V. Aho, M. R. Garey, and J. D. Ullman, "The transitive reduction of a directed graph," *SIAM Journal on Computing*, vol. 1, no. 2, pp. 131–137, 1972.
- [80] B. R. Rau and C. D. Glaeser, "Some scheduling techniques and an easily schedulable horizontal architecture for high performance scientific computing," *SIGMICRO Newsl.*, vol. 12, p. 183–198, Dec. 1981.

- [81] M. A. Javed, M. S. Younis, S. Latif, J. Qadir, and A. Baig, "Community detection in networks: A multidisciplinary review," *Journal of Network and Computer Applications*, vol. 108, pp. 87 – 111, 2018.
- [82] L. Gurobi Optimization, "Gurobi optimizer reference manual," 2019.
- [83] Y. Wang, P. Li, and J. Cong, "Theory and algorithm for generalized memory partitioning in high-level synthesis," in *Proceedings of the 2014 ACM/SIGDA International Symposium on Field-programmable Gate Arrays*, FPGA '14, (New York, NY, USA), pp. 199–208, ACM, 2014.
- [84] E. W. Dijkstra, "A note on two problems in connexion with graphs," *Numerische mathematik*, vol. 1, no. 1, pp. 269–271, 1959.
- [85] A. Hansson, K. Goossens, and A. Rădulescu, "Avoiding message-dependent deadlock in network-based systems on chip," *VLSI design*, vol. 2007, 2007.
- [86] R. Tessier, K. Pocek, and A. DeHon, "Reconfigurable computing architectures," *Proceedings of the IEEE*, vol. 103, pp. 332–354, March 2015.
- [87] K. Choi, "Coarse-grained reconfigurable array: Architecture and application mapping," *IPSJ Transactions on System LSI Design Methodology*, vol. 4, pp. 31–46, 2011.
- [88] D. Sanchez, G. Michelogiannakis, and C. Kozyrakis, "An analysis of on-chip interconnection networks for large-scale chip multiprocessors," *ACM Trans. Archit. Code Optim.*, vol. 7, pp. 4:1–4:28, May 2010.
- [89] E. Garcia and M. Gupta, "Lattice regression," in *Advances in Neural Information Processing Systems*, pp. 594–602, 2009.
- [90] X. Wang, P. Liu, M. Yang, and Y. Jiang, "Avoiding request–request type message-dependent deadlocks in networks-on-chips," *Parallel Computing*, vol. 39, no. 9, pp. 408–423, 2013.
- [91] W. J. Dally and B. P. Towles, *Principles and Practices of Interconnection Networks*. Elsevier, 2004.
- [92] D. Wang, B. Ganesh, N. Tuaycharoen, K. Baynes, A. Jaleel, and B. Jacob, "Dramsim: A memory system simulator," *SIGARCH Comput. Archit. News*, vol. 33, pp. 100–107, Nov. 2005.
- [93] N. Jiang, J. Balfour, D. U. Becker, B. Towles, W. J. Dally, G. Michelogiannakis, and J. Kim, "A detailed and flexible cycle-accurate network-on-chip simulator," in *Performance Analysis of Systems and Software (ISPASS), 2013 IEEE International Symposium on*, pp. 86–96, IEEE, 2013.

- [94] D. U. Becker, *Efficient Microarchitecture for Network-on-Chip Routers*. PhD thesis, Stanford University, Palo Alto, 2012.
- [95] J. Fowers, K. Ovtcharov, M. Papamichael, T. Massengill, M. Liu, D. Lo, S. Alkalay, M. Haselman, L. Adams, M. Ghandi, *et al.*, “A configurable cloud-scale dnn processor for real-time ai,” in *Proceedings of the 45th Annual International Symposium on Computer Architecture*, pp. 1–14, IEEE Press, 2018.
- [96] N. P. Jouppi, C. Young, N. Patil, D. Patterson, G. Agrawal, R. Bajwa, S. Bates, S. Bhatia, N. Boden, A. Borchers, *et al.*, “In-datacenter performance analysis of a tensor processing unit,” in *Computer Architecture (ISCA), 2017 ACM/IEEE 44th Annual International Symposium on*, pp. 1–12, IEEE, 2017.
- [97] M. I. Gordon, W. Thies, M. Karczmarek, J. Lin, A. S. Meli, A. A. Lamb, C. Leger, J. Wong, H. Hoffmann, D. Maze, and S. Amarasinghe, “A stream compiler for communication-exposed architectures,” in *Proceedings of the 10th International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS X*, (New York, NY, USA), pp. 291–303, ACM, 2002.
- [98] P. M. Phothisilimthana, T. Jelvis, R. Shah, N. Totla, S. Chasins, and R. Bodik, “Chlorophyll: Synthesis-aided compiler for low-power spatial architectures,” in *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '14*, (New York, NY, USA), pp. 396–407, ACM, 2014.
- [99] N. Trifunovic, H. Palikareva, T. Becker, and G. Gaydadjiev, “Cloud deployment and management of dataflow engines,” in *Proceedings of the 1st International Workshop on Next Generation of Cloud Architectures, CloudNG:17*, (New York, NY, USA), pp. 5:1–5:6, ACM, 2017.
- [100] T. Nowatzki, M. Sartin-Tarm, L. De Carli, K. Sankaralingam, C. Estan, and B. Robatmili, “A general constraint-centric scheduling framework for spatial architectures,” in *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '13*, (New York, NY, USA), pp. 495–506, ACM, 2013.
- [101] N. E. Jerger, T. Krishna, and L.-S. Peh, “On-chip networks, second edition,” *Synthesis Lectures on Computer Architecture*, vol. 12, no. 3, pp. 1–210, 2017.

- [102] M. B. Taylor, J. Kim, J. Miller, D. Wentzlaff, F. Ghodrat, B. Greenwald, H. Hoffman, P. Johnson, J.-W. Lee, W. Lee, A. Ma, A. Saraf, M. Seneski, N. Shnidman, V. Strumpen, M. Frank, S. Amarasinghe, and A. Agarwal, "The raw microprocessor: A computational fabric for software circuits and general-purpose programs," *IEEE Micro*, vol. 22, pp. 25–35, Mar. 2002.
- [103] D. Wentzlaff, P. Griffin, H. Hoffmann, L. Bao, B. Edwards, C. Ramey, M. Mattina, C.-C. Miao, J. F. Brown III, and A. Agarwal, "On-chip interconnection architecture of the tile processor," *IEEE Micro*, vol. 27, pp. 15–31, Sept. 2007.
- [104] M. B. Taylor, W. Lee, S. Amarasinghe, and A. Agarwal, "Scalar operand networks: On-chip interconnect for ilp in partitioned architectures," in *Proceedings of the 9th International Symposium on High-Performance Computer Architecture, HPCA '03*, (Washington, DC, USA), pp. 341–353, IEEE Computer Society, 2003.
- [105] K. Sankaralingam, R. Nagarajan, H. Liu, C. Kim, J. Huh, D. Burger, S. W. Keckler, and C. R. Moore, "Exploiting ilp, tlp, and dlp with the polymorphous trips architecture," in *30th Annual International Symposium on Computer Architecture, 2003. Proceedings.*, pp. 422–433, June 2003.
- [106] P. Gratz, C. Kim, K. Sankaralingam, H. Hanson, P. Shivakumar, S. W. Keckler, and D. Burger, "On-chip interconnection networks of the trips chip," *IEEE Micro*, vol. 27, pp. 41–50, Sept. 2007.
- [107] S. Swanson, A. Schwerin, M. Mercaldi, A. Petersen, A. Putnam, K. Michelson, M. Oskin, and S. J. Eggers, "The wavescalar architecture," *ACM Trans. Comput. Syst.*, vol. 25, pp. 4:1–4:54, May 2007.
- [108] H. Park, Y. Park, and S. Mahlke, "Polymorphic pipeline array: A flexible multicore accelerator with virtualized execution for mobile multimedia applications," in *2009 42nd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pp. 370–380, Dec 2009.
- [109] A. Ye and J. Rose, "Using bus-based connections to improve field-programmable gate-array density for implementing datapath circuits," *IEEE Trans. Very Large Scale Integr. Syst.*, vol. 14, pp. 462–473, May 2006.
- [110] Y. Huang, P. Ienne, O. Temam, Y. Chen, and C. Wu, "Elastic cgras," in *Proceedings of the ACM/SIGDA International Symposium on Field Programmable Gate Arrays, FPGA '13*, (New York, NY, USA), pp. 171–180, ACM, 2013.

- [111] B. Van Essen, A. Wood, A. Carroll, S. Friedman, R. Panda, B. Ylvisaker, C. Ebeling, and S. Hauck, "Static versus scheduled interconnect in coarse-grained reconfigurable arrays," in *Field Programmable Logic and Applications, 2009. FPL 2009. International Conference on*, pp. 268–275, IEEE, 2009.
- [112] G. Dimitroulakos, M. D. Galanis, and C. E. Goutis, "Exploring the design space of an optimized compiler approach for mesh-like coarse-grained reconfigurable architectures," in *Parallel and Distributed Processing Symposium, 2006. IPDPS 2006. 20th International*, pp. 10–pp, IEEE, 2006.
- [113] C. Nicol, "A coarse grain reconfigurable array (cgra) for statically scheduled data flow computing."
- [114] C. A. Moritz, D. Yeung, and A. Agarwal, "Exploring optimal cost-performance designs for raw microprocessors," in *Proceedings. IEEE Symposium on FPGAs for Custom Computing Machines (Cat. No.98TB100251)*, pp. 12–27, April 1998.
- [115] F. Bouwens, M. Berekovic, A. Kanstein, and G. Gaydadjiev, "Architectural exploration of the adres coarse-grained reconfigurable array," in *Reconfigurable Computing: Architectures, Tools and Applications* (P. C. Diniz, E. Marques, K. Bertels, M. M. Fernandes, and J. M. P. Cardoso, eds.), (Berlin, Heidelberg), pp. 1–13, Springer Berlin Heidelberg, 2007.
- [116] R. Hartenstein, M. Herz, T. Hoffmann, and U. Nageldinger, "Kressarray xplorer: A new cad environment to optimize reconfigurable datapath array architectures," in *Proceedings 2000. Design Automation Conference (DAC)*, pp. 163–168, Jan 2000.
- [117] N. Bansal, S. Gupta, N. Dutt, A. Nicolau, and R. Gupta, "Network topology exploration of mesh-based coarse-grain reconfigurable architectures," in *Proceedings Design, Automation and Test in Europe Conference and Exhibition*, vol. 1, pp. 474–479, Feb 2004.
- [118] Y. Kim, R. N. Mahapatra, and K. Choi, "Design space exploration for efficient resource utilization in coarse-grained reconfigurable architecture," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 18, pp. 1471–1482, Oct 2010.
- [119] S. Murali and G. D. Micheli, "Sunmap: a tool for automatic topology selection and generation for nocs," in *Proceedings. 41st Design Automation Conference, 2004.*, pp. 914–919, July 2004.

- [120] M. A. Kinsy, M. H. Cho, T. Wen, E. Suh, M. van Dijk, and S. Devadas, "Application-aware deadlock-free oblivious routing," in *Proceedings of the 36th Annual International Symposium on Computer Architecture, ISCA '09*, (New York, NY, USA), pp. 208–219, ACM, 2009.
- [121] K. S. Shim, M. H. Cho, M. Kinsy, T. Wen, M. Lis, G. E. Suh, and S. Devadas, "Static virtual channel allocation in oblivious routing," in *2009 3rd ACM/IEEE International Symposium on Networks-on-Chip*, pp. 38–43, May 2009.
- [122] C. J. Glass and L. M. Ni, "The turn model for adaptive routing," in [1992] *Proceedings the 19th Annual International Symposium on Computer Architecture*, pp. 278–287, May 1992.
- [123] O. Ozturk, M. Kandemir, M. J. Irwin, and S. H. Narayanan, "Compiler directed network-on-chip reliability enhancement for chip multiprocessors," *ACM Sigplan Notices*, vol. 45, no. 4, pp. 85–94, 2010.
- [124] C. Ababei, H. S. Kia, O. P. Yadav, and J. Hu, "Energy and reliability oriented mapping for regular networks-on-chip," in *Proceedings of the Fifth ACM/IEEE International Symposium on Networks-on-Chip*, pp. 121–128, ACM, 2011.
- [125] N. Kapre and A. Dehon, "An NoC traffic compiler for efficient FPGA implementation of sparse graph-oriented workloads," *International Journal of Reconfigurable Computing*, vol. 2011, 2011.