Software-defined Hardware with Groq's Tensor Streaming Processor

A hardware-software approach



HELLO ATPESC Andrew Ling, Ph.D. Sr. Director Software



MY GOAL TODAY Introduce you to Groq's approach to SW defined hardware in the era beyond CPUs



Details...



A Software-defined Tensor Streaming Multiprocessor for Large-scale Machine Learning

Dennis Abts	Garrin Kimmell	Andrew Ling	John Kim
Groq Inc.	Groq Inc.	Groq Inc.	KAIST/Groq Inc.
Matt Boyd	Andrew Bitar	Sahil Parmar	Ibrahim Ahmee
Groq Inc.	Groq Inc.	Groq Inc.	Groq Inc.
Roberto DiCecco	David Han	John Thompson	Michael Bye
Groq Inc.	Groq Inc.	Groq Inc.	Groq Inc.
Jennifer Hwang	Jeremy Fowers	Peter Lillian	Ashwin Murth
Groq Inc.	Groq Inc.	Groq Inc.	Groq Inc.
Elyas Mehtabuddin	Chetan Tekur	Thomas Sohmers	Kris Kang
Groq Inc.	Groq Inc.	Groq Inc.	Groq Inc.
	Stanhan Marach	Jonathan Rose	

Grog Inc.

ABSTRACT

large-scale interconnection networks of tensor streaming process-

ing (TSP) elements. The system architecture includes packaging, routing, and flow control of the interconnection network of TSPs.

We describe the communication and synchronization primitives of

a bandwidth-rich substrate for global communication. This scalable

communication fabric provides the backbone for large-scale systems based on a software-defined Dragonfly topology, ultimately vielding a parallel machine learning system with elasticity to sup-

port a variety of workloads, both training and inference. We extend

the TSP's producer-consumer stream programming model to include global memory which is implemented as logically shared, but

physically distributed SRAM on-chip memory. Each TSP contributes

220 MiBytes to the global memory capacity, with the maximum

capacity limited only by the network's scale - the maximum num-

ber of endpoints in the system. The TSP acts as both a processing

element (endpoint) and network switch for moving tensors across

the communication links. We describe a novel software-controlled

networking approach that avoids the latency variation introduced

by dynamic contention for network links. We describe the topol-

Permission to make digital or hard copies of part or all of this work for personal or

classroom use is granted without fee provided that copies are not made or distributed fee profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored.

the network that serves as the fabric for a la

chine learning system with up to 10

TeraBytes of global memory accessibl

For all other uses, contact the owner/author(s ISCA '22, June 18-22, 2022, New York, NY, USA

© 2022 Copyright held by the owner/author(s)

ACM ISBN 978-1-4503-8610-4/22/06.

https://doi.org/10.1145/3470496.3527405

of end-to-end system latency.

ogy, routing and flow control to characterize the performance of

SPs and ore than 2

croseconds

CCS CONCEPTS We describe our novel commercial software-defined approach for

 Computer systems organization → Interconnection architectures.

KEYWORDS

Grog Inc.

Machine Learning, Tensor Streaming Processor, Dragonfly, Software Scheduling

1 INTRODUCTION

Historically, high-performance computing (HPC) systems were broadly categorized as capability or capacity systems. This dichotomy arises because of communication latency and bandwidth trade-offs when we apply more processing elements (PEs) to a fixed-size problem (strong scaling) with the goal of minimizing the program's execution time. Alternatively, we can deploy more PEs to increase throughput (i.e. weak scaling). This duality requires both novel chip architectures [7][1][36] in the underlying PEs and a scalable system architecture with high throughput (bisection bandwidth) and low end-to-end latency (low network diameter) for fine-grained communication necessary to efficiently handle both strong and much scaling. For example, the network demand for



vidual processing elements work in unison to collectively execute the different "layers" of a deep learning network. It is this set of sub-tasks, expressed as individual PE programs, that are distributed among the computing elements and responsible for carrying out, or executing, the specifics of the machine learning model.

The burgeoning parameter space of natural language processing (NLP) models like GPT-3 [6] use 100s-of-billions of parameters to



nnection network, and the programming model of the indi-

Groq in a Nutshell

Founded	2016
Flagship Product	Groq Tensor Streaming Processor Chip (TSP) and Software for use in AI, Machine and Deep Learning applications
Target Market(s)	Al and Machine Learning Hyperscalers, Government, HPC clusters, Autonomous vehicles, high performance edge appliances
Employees	~250

Founder/CEO Jonathan Ross



ML Compilers and Tools

DNN Library Developers

Host-IO Developers

ASIC/HW Designers

Supporting Workloads of **Today**

- What would you build given the data flow nature of Machine Learning Workloads and HPC?
 - "Nodes" in the computational graph represent **operators** and "edges" are the **operands** and **results**.



- Groq re-examined the hardware-software contract
- Created hardware that is much more predictable and streaming based
- Gave more control to software!

Have we been here before? RISC → VLIW?

With RISC, CPU architecture hides a lot of complex control logic from software

Caching, prefetching, out-of-order execution, branch prediction

Area cost of this control hardware reduces available area for raw compute

Very Long Instruction Word (VLIW) to the rescue?

The next natural step in the trend of shifting work and control from hardware to software

- Execution ordering, data prefetching, branch resolution
- Compiler responsible for exploiting instruction-level parallelism (ILP) available in a program

Itanium...

Compiler struggled to achieve good ILP due to limited view of dynamic hardware behaviour

Cache misses, dynamic mem dependencies, dynamic branches



A New Golden Age in Computing

The end of CPU hegemony

Growing demand in dataflow dominated compute

Slowing of Moore's Law and Dennard Scaling

CPU "abstraction" no longer the only foundation for developing software

Dawn of a new golden age in computing

Hennessy & Patterson: "A New Golden Age for Computer Architecture"

Lattner: "A New Golden Age of Compilers"

Karpathy: "Software 2.0" → A New Golden Age for Algorithms

In the past, hardware (CPUs) has defined our software

The hardware-software abstraction contract has been reopened!

Opportunity now to achieve "software-defined hardware"

Simplifying Compute

GroqChip™



Compute

© 2022 Grog, Inc. | Stanford EE382A

Nvidia GTX 1070



Complexity leads to more Compute Costs

GroqChip Scalable Architecture

220MB SRAM 80 TB/s of Stream Bandwidth on-chip Massive concurrency

Dense MatMul 720 TeraOP/s (1 TeraOP / s / mm²)

Vector Units



480GBps Chip-2-Chip Links

Dataflow

Instruction Control

Designing for determinism...

Design choices along the way need to accommodate the "design for determinism" design philosophy

- Hardware must <u>enable</u> the compiler and runtime interfaces to reason about program execution
 - Memory consistency model must be well understood disallowing memory references from being reordered
 - No "reactive components" like arbiters, crossbars, replay mechanisms, caches, etc
 - Software must have access to the **architectural-visible machine state** in order to intercept the data (operands) with the instruction that will execute on them.
- Compiler "knows" the exact location of every tensor on-chip
- In this way, the compiler is orchestrating the arrival of operands and the instructions which use them. The producer-consumer stream programming model allows a set of "streaming register files" to track the state of each tensor flowing through the chip.

Avoiding Complexity at the Chip Level

Conventional CPUs add features and add complexity:

- **Speculative execution** and out-of-order retirement
 - to improve instr level parallelism increases tail latency
- Implicit data flow through cache memory hierarchies introduce complexity and non-determinism
 - (e.g. DRAM → L3 → L2 → L1 → GPRs) to hide DRAM access latency & pressure - not energy or silicon efficient



Requires dynamic profiling to understand execution time and throughput characteristics of deep learning models

The TSP simplifies data flow through Stream Programming:

- A large, single-level scratchpad SRAM **fixed, deterministic latency**
- **Explicitly** allocate tensors in space and time unlocking massive memory concurrency, and compute flexibility along multiple dimensions:

[device, hemisphere, memory slice, bank, address offset]



Predictable Performance at Scale

Functionally Sliced Microarchitecture

Reorganizing the multicore mesh



IF	ID	EX	МЕМ	WB
----	----	----	-----	----

Functionally Sliced Microarchitecture

Reorganizing the multicore mesh





© 2022 Groq, Inc. | Software-defined Tensor Streaming Multiprocessor

Functionally Sliced Microarchitecture

Reorganizing the multicore mesh

Reorganize a conventional
manycore 2D meshMEM: on-chip SRAMVXM: vector unitMXM: matrix unitSXM: data reshapes

Tensor operands and results flow on **"streams"** horizontally

Instructions flow vertically executed in a SIMD manner







1\$

NET

MXM - Matrix Multiply Engines



МХМ			МХМ	
				Γ



Numeric Mode	Max Size	Supported Density	Result Tensor
int8	[1, 320] x [320, 320]	Two per MXM	int32
float16	[1, 320] x [160, 320]	One per MXM	float32

320B x 320B dot product Loads 320B x16 in 20 cycles 20 cycle execution Fully pipelined, N

Int8 & float16 Full precision expansion 32-bit accumulate Used Independently or together

VXM and Complex, Customized Functions



Dataflow begins with memory Read onto Stream Tensor Many concurrent streams are supported in programming model VXM provides a flexible and programmable fabric for Compute Compute occurs on data locality of passing Stream Tensor MEM bandwidth supports high concurrency



SXM - Switch eXecution Module

SXM		SXM	



Swiss army knife for data manipulation & Intra-vector byte operations Distributor: 4 per hemisphere perform unto mapping of input + mask to output stream within a 16 byte superlane <u>Transposer</u>: 2 per hemisphere perform intra-superlane transpose over 16 lanes for 20 superlanes <u>Permuter/Shifter</u>: arbitrary mapping of input + mask, shuffling between 320B vector elements - used for data transforms like pads/reshapes

Shift, Rotate, Distribute, Permute, Transpose, Transport to SuperLanes



Streaming Registers

Contrasts to CPU / Register file

- Memory semantics have an address **and** a direction of data flow

Streaming registers are like a **tensor assembly line** flowing eastward and westward.

Each stream holds a vector of **320 bytes** (i.e. each element in the vector is 1-byte of data)

- Multiple streams are used to represent larger multi-byte data types (eg. FP32).









An ISA That Empowers Software

Explicit time and space of instruction execution exposed by ISA to the compiler

Each Functional Unit (FU) type provides its own low-level instruction set

Number of FUs of each type and relative positions on chip are exposed to software

Compiler can choose to leverage multiple FUs for more concurrency or more pipelining

Instruction Set	Function	Instruction
	ICU	NOP N
		Ifetch
Low level		Sync
		Notify
		Config
		Repeat n,d
320-vector ops		
020 10000 000	MEM	Read a,s
		Write a,s
		Gather s, map
Explicit resource		Scatter s, map
		Countdown d
selection		Step a
Sciection		Iterations n
	VYM	unary operation
		hinary operation
Explicit		type conversions
Explicit		Relui
schedulina		TanH
Schedding		Exp
		RSqrt
	МХМ	LW
		IW
		ABC
		ACC
	SXM	Shift up/down N
		Permute map
		Distribute map
		Rotate stream
		Transpose sg16
	C2C	Deskew
		Send
		Receive

Power of Data Orchestration

Given to Groq Compiler



9roq[™]

Compilation

Traditional HPC Compiler Flow

A: The following graphic shows how cuDNN relates to other software in the stack. *Figure 11. Software stack with cuDNN.*



NVIDIA CUDNN DOCUMENTATION

https://docs.nvidia.com/deeplearning/cudnn/developer-guide/index.html



"The most dangerous phrase in the language is "We've always done it this way."

REAR ADMIRAL GRACE HOPPER

Pioneering Computer Scientist 1906–1992





GROQ[™] COMPILER

- Software-defined hardware relies on several interfaces:
 - Static-Dynamic interface compile-time versus runtime
 - Hardware-software interface exposing the architectural visible-state



Kernel-less approach to HPC Compilation

GEMM - General Matrix Multiplication

- Vector-Matrix multiplication and Matrix-Matrix multiplication are the workhorse for many ML and HPC workloads
- On-chip memory bandwidth determines how quickly we can ramp up ALUs for vector and matrix operations
- Consistent performance across a range of tensor sizes - less "hardware fitting"
- State-of-the-art (SOTA) results across a range of models and applications
 - CNNs
 - RNNs, LSTMs
 - NLP, BERT
 - Dense linear algebra V*M and M*M



GroqFamily



Scale-out organization

- Topology objectives
 - Low network diameter
 - Direct network
 - Hierarchical packaging-aware topology
- System packaging hierarchy
- Chip-to-chip (C2C) links and flow control



Top-view of node

Low-diameter Network

Minimize the number of hops in the network

- The total observed latency and variance increases with the number of hops in the network
- Dragonfly is a *hierarchical* topology minimizes the number of hops taken



Chip-to-chip (C2C) links and flow control

- Simplified communication model provides logically shared access to global SRAM which is physically distributed among the TSPs
- C2C links directly connect TSPs
- Comparing conventional RDMA with Groq C2C communication



Cholesky Factorization

- Cholesky decomposition is an important technique for solving large-scale system of linear equations.
- Compute is $N^3/3$ for a SPD matrix of size N
- Block-cyclic distributions of 320 rows across chips









(c) execution time vs problem size (input matrix size) with multiple TSPs

Multichip Compilation

Available soon in next Quarter

Inter Op Partitioning

- Model graph split into multiple graphs
- Intermediate activations become inputs to successive graphs
- Reduces weights on chip

Intra Op Partitioning

- Tile the face of the input
- Operate on a "fractional" portion of the input at any given time
- Reduces peak activation / input size on the chip



Multichip Inter Op Partitioning





Multichip Pipelining

- Invocation is composed of compute on each of the devices followed by C2C to next subgraph's device
 - Pipeline parallel execution

Input		Time		
Device 0	Inference 0, partition 0	Inference 1, partition 0	Inference 2, partition 0	Inference 3, partition 0
Device 1		Inference 0, partition 1	Inference 1, partition 1	Inference 2, partition 1
Device 2			Inference 0, partition 2	Inference 1, partition 2
Device 3				Inference 0, partition 3
Output	•	Later	ncy = 4 x (Compute + I	O time)
9009 ° © 2022 0	Groq, Inc. Software-defined Tensor Streaming Multiprocessor	Throu	ughput = Clock freq / (Compute + IO time)

Inter Op Partitioning Results Today

Linear Scaling Achieved

Scaling **BERT** encoders

Scaling to 6, 24, 48, and 96 BERT encoders on 1, 4, 8, and 16 TSPs respectively





Results



Natural Language Processing with BERT

Seq Len=128

- Spread the layers across TSPs in each node for multi-TSP model parallelism
- Mini-batch (data parallelism) across nodes for throughput
- Goal: Accelerate BERT minimizing latency and maximizing throughput
- Minimize latency variance to provide predictable throughput
 - T4, CURRENT SOTA (A100) [21] AND THIS WORK BERT-BASE LATENCY (128 SEQUENCE LENGTH).

	T4 (μs)	Current SOTA A100 (µs)	This work (µs)	Speedup (This work vs SOTA)
Average	1330	630	128.9	4.8×
95 th Percentile	1550	780	129.1	6×
99 th Percentile	1570	790	129.5	6.1×

A100 [22], T4 [23] AND TSP [13], [24] SPECIFICATIONS.

Chip	Die Area (mm ²)	Tech Process (nm)	Transistor count (B)	TDP (W)
NVIDIA T4	545	12	13.6	70
NVIDIA A100	826	7	54.2	400
Groq TSP	725	14	26.8	275

40



© 2022 Groq, Inc. | Software-defined Tensor Streaming Multiprocessor GROQ[™] COMPILER Achieving SOTA Results

LSTM-512 TOPs vs Batch Size Intel NX int8 🛛 🔵 V100 int8 GroqChip int8 T4 int8 200 150 TOPs 100 50

6

Batch Size

8



© 2022 Groq, Inc. | Software-defined Tensor Streaming Multiprocessor

0

3

Nvidia results from publicly available data on nvidia.com LSTM results from Intel sponsored paper, FPT 2020, available on intel.com

32

256

Groq Exponentially Supporting More Workloads



grog

Summary and Takeaways

- Low-latency and high-throughput are necessary...
- Delivering **predictable** and **repeatable** performance is critical for many user-facing applications
 - Batch-1 inference is important for responsiveness and delivering quality-of-service (QoS) that is impossible to do with more traditional microarchitectures using crossbars, cache hierarchies, etc.
- Determinism enables software-defined hardware and entails a design philosophy that spans both hardware and software
 - ISA is <u>not</u> about abstraction of hardware details, but about **exerting control of underlying hardware**
 - 144 independent instruction control units (ICUs) of the TSP
 - Expose the architecturally-visible state (GPRs, SRAM, instruction buffers, etc)
 - Software-based replay and exception handling
- Extending the single-chip TSP determinism to the multiprocessor using software scheduled networking to explicitly schedule tensors on the network links
- Synchronous communication model allows for large-scale machine learning





