

GroqChip™ Overview

SRAM Memory

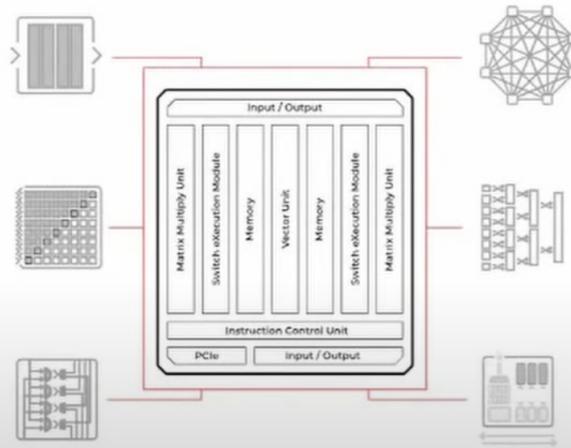
Massive concurrency
80 TB/s of BW
230MB capacity
Stride insensitive

Groq TruePoint™ Matrix

4x Engines
750 TOP/s int8
188 TFLOP/s fp16
320x320 fused dot product

Programmable Vector Units

5,120 Vector ALUs for high performance



Networking

480 GB/s bandwidth
Extensible network scalability
Multiple topologies

Data Switch

Shift, Transpose, Permuter for improved data movement and data reshapes

Instruction Control

Multiple instruction queues for instruction parallelism

groq

public 14

Architecture Empowering Software

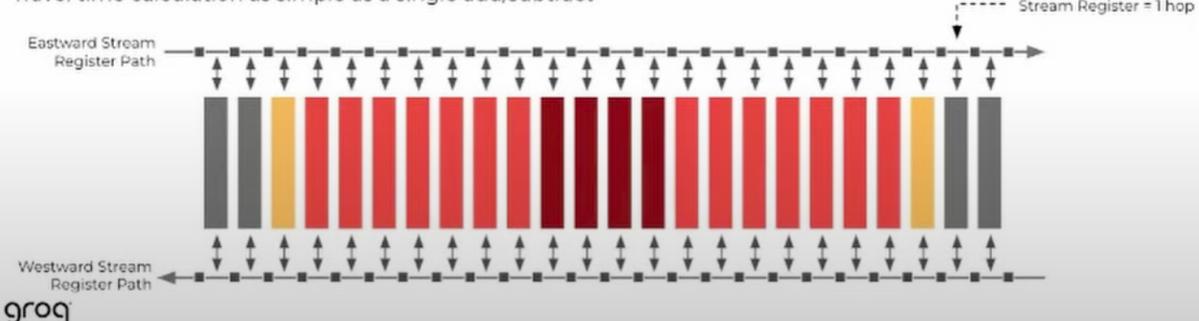
Simple, one-dimensional interconnect for inter-FU communication

Compiler can quickly reason about all data movement between FUs

- Eastward and westward paths made up of arrays of “stream registers”
- Stream register = one-cycle hop

No arbiters / queues = software can easily reason about exact data movement without simulation

Travel time calculation as simple as a single add/subtract



groq

public 23

≡ Traffic Nightmare

Stop-and-go bumper-to-bumper traffic gives **long travel times** and **poor gas mileage**.



Unpredictable nature leads to **poor utilization** of roads

groq

A Smart City for AI ≡

Non-stop point to point transport gives **shortest travel time** and **best gas mileage**.



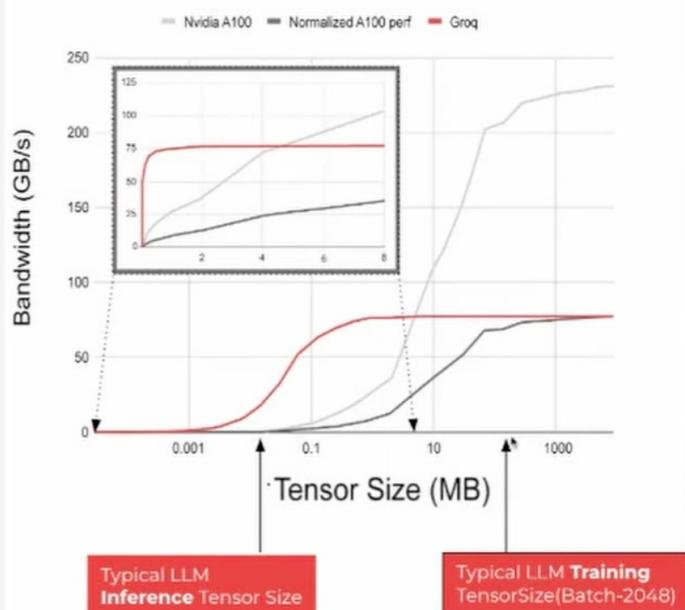
Predictability means **high utilization** and no accidents **EVER!**
Nothing is left to chance.

Public 35

AllReduce Comparison

Comparison made with an eight GPU A100 system with NCCL

- Only a handful of cycles to Read(vector) \leftrightarrow Send(vector) enables fine-grained communication across the 16 directly connected links on each TSP
- A100 system has approximately 3x higher network channel bandwidth
- When normalized, the GroqChip device matches the bandwidth at large tensor size while significantly improving bandwidth at intermediate tensor size



Results for A100 were measured on an 8 A100 GPU system with 300 GB/s of NVlink bandwidth per GPU connected through NVSwitch.
Nvidia results are from publicly available data on <https://github.com/nvidiallennlp/test/>

Public 41

Solution Diversity

Customer Problem Statement	Value Delivered by Groq
Drug discovery: Accelerate time to discovery from days to minutes	>200x speed-up when evaluating candidate COVID drugs
Cyber security: Improve accuracy and reduce false positives	>600x speed up for real-time cyber-threat anomaly detection; with superior accuracy
Fusion reactor: Enable fully predictable real-time controls systems (<1sec)	>600x speed up to make real-time plasma stabilization possible
Capital markets: Enable rapid hypothesis testing at Scale	>100x speed-up enabling rapid trading hypothesis testing
General ML: Support a diverse set of popular models	>500 common models natively compilable with performance ahead of GPUs

groq

PAGE 47

THE BATTLE OF EXPONENTIALS

Reducing TTM improves performance and energy efficiency

AI Models Growing Fast

Doubling every 3.5 months

AI workload structure and complexity changing



Running new workload on AI Accelerator architected and implemented in a process that is TTM old

AI Accelerator Architecture & Process defined TTM before Workload Execution

AI Accelerators Silicon Tech Growing Slow

Moore & Dennard Scaling Slowing Down

SoC Transistor Counts now Doubling every >40 months

Amdahl's law demands end-to-end SW/HW co-optimization

Shortening Time to Market (TTM) improves Workload Execution Efficiency

groq

PAGE 49

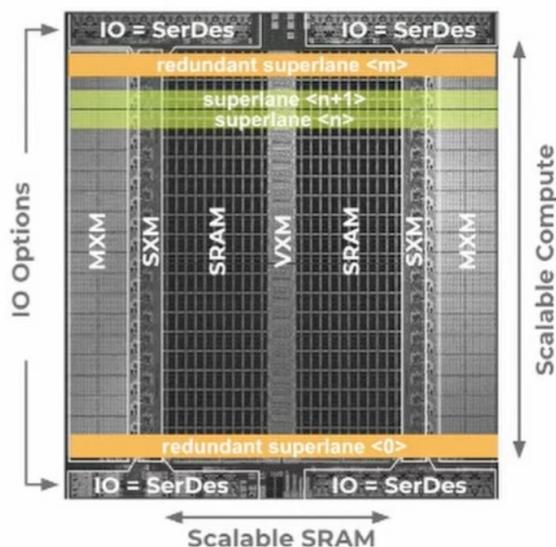
SCALABLE Silicon Tiler For Fast Time-to-market

Multiple Interconnect Options

- C2C for high-radix interconnect
- UCIe for MCM connected sidecar accelerator
- Scalable SXM for BW to/from IO and Compute

Scalable compute architecture

- SRAM scalable capacity
- VXM with scalable number of PEs
- MXM with scalable matrix sizes



groq

PUBLIC 18

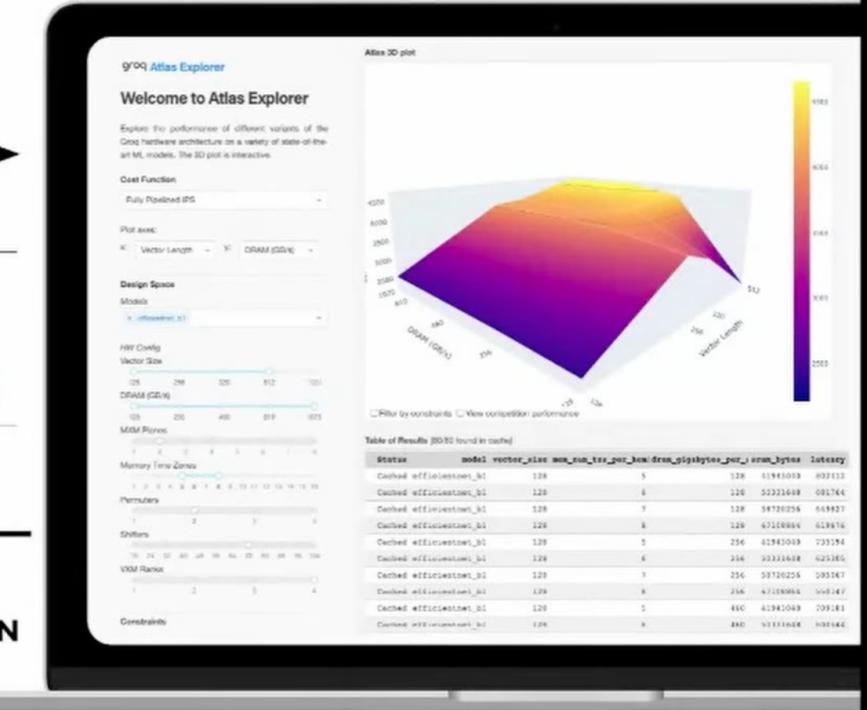
AI/HPC MODELS

DEMOS

≡ Design Exploration

**CUSTOM
ARCHITECTURE
RECOMMENDATION**

groq



SILICON TILER

Groq Silicon Tiler Ecosystem

Same
AI SOFTWARE
for the full
HW Ecosystem

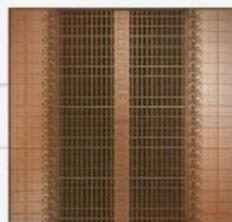
Scalable SRAM

(220–440MiB)
with 3D SRAM
extension

Scalable Compute

16 SL: 256x256
20 SL: 320x320
24 SL: 384x384

TSP core



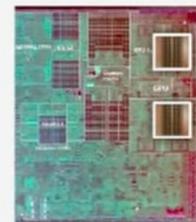
Chip



Chiplet



IP



groq

public 57

Summary

State-of-Art LLM

Performance using
Software-scheduled LPU
Compute and Networking
across 576 chips

Groq Superpowers:

Chip Determinism unlocks performance and efficiency while paving the way for more compute density with 3D chip integration

Low-diameter dragonfly network with abundant path diversity achieves better latency and peak AllReduce performance

Synchronous global communication extends multiple GroqChips into a massive multi-chip single-core cluster

Software Scheduled Deterministic Compute and Network uses **Global time** to provide efficient C2C utilization enabling massive scale



groq

public 54