

Groq

Fast inference on GPU's

Tokens drive compute.

Fastest inference speed for models

TTFB is the most important metric for LLM's.
in a real-time setting

Groq Chip

→ First built shw

→

Chip & GroqCard ⇒ GroqNode

& GroqNode ⇒ GroqRack

Process anything seq in nature. V FAST

GroqChip:

Start from sand (i.e the chip)

AIM: H/w that's easy to program

Jonathan - founder of Google. Inventor of TPU

TPU hard to program

AIM: SW should be easy on H/w

How did we get here?

End of Dennard Law →

Single core ILP hit a wall

Multicore

Shift parallelism responsibility from H/w to SW

How to fix it?

Custom H/w for applications.

AI algos are well behaved (data flow is static, algo, etc).

Problem is memory is not deterministic

This is why AI compilers struggle

Solution: Make 100% deterministic hardware

- LPU
- Simpler than H100
 - 14 nm
 - 3 tech holder
 - No HBM
 - No silicon imposer

Mojo? Architecture

Lot of waiting for H100. Adds to latency.
Using HBM has penalty

LPU

- Fully deterministic
- No cache. Direct memory?
- How do you load Llama2 - 70B model into a device that doesn't have HBM?
- Entire system is deterministic
- All chips are synced to act like one device i.e like a megachip.
- Within 2-3 μ s, have access to

TB worth of memory

- can do the above in a deterministic way
- Soons to access any other device within range
- Assign work to idle units fast
- Low latency because everything is pre-scheduled & orchestrated by software
- Massively Scale.

→ Grog has shown 600x improvement in cyber security over GPUs.

→ Work in fin space. Fast processing & of data

→ See papers on website.

→ Also work with sparse models.

→ TSP: Tensor Streaming Processor
But CPU built for LM

Building blocks

SIMD Unit: 320 - element vector.

Lightweight instruction dispatch.

Different types of special SIMD units

MM - Matrix Vector

VMM - Vector Vector

SMM - Data reshapes

MEM - On chip SRAM

Flat memory

Significant bandwidth to compute! 80 TB/s

⇒ Almost 100x the bandwidth of HBM

Single one dimensional interconnect for inter FU communication.

Instruction set:

320 element vector ops. Compared to 10k instructions per

explicit resource selection

(instructions vs
pytorch)

Cool parts:

→ Activation functions like ReLU, TanH are part of the core instruction set.

MLIR used in compiler.

Orxog: ⇒ 35 developers

Nvidia ⇒ 50K developers

Orxog only focuses on inference for now

hnm chip next!!

For 70b LLM ⇒ didn't answer it.

Interconnect:-

Usual 3 layers :- S/W, Network & Compute
CPU S/W + N/W + Compute
S/W controlled network

- Use Slw to sync chips to act like one big chip (layer)
- Use 'DragonFly' mode of connection.
- One effective global clock
- No routing
- 3/NS \Rightarrow Access to 2 TB of SRAM.

Group can reach max bandwidth in a few ns.
This is why inference is fast.

Slw scheduled communication

Efficient use of resources is the key reason.

Scalability

- Linear increase in perf as number of LPU's increase

GPU VS LPU:-

Both are matrix processing engines.

CPU: Have a 'production' shop.
DGX box is limited of 8 packed chips

GPU: Assembly line for making tokens
No HBM
Don't go through switches to move data between devices
Linear movement.
Low latency
Lower power consumption.

Power metric

Joules per token
GPU is 10x better.

Reasons

- No storage in HBM (because no HBM)
- Don't access NIC
- Because of scaling, don't need to fight back pressure, adaptive routing.