

14/01/24

Astronord.github.io Link:-

CPU's Architecture

Transformer Inference Optimization Toolkit

→ Streaming Processor
(SM)

→ (SM) like what is in CPU.
basic CPU block contains

Instruction Schedulers

Instruction execution Pipelines

→ [HBM], high bandwidth memory

→ special off chip memory
data is initially stored
and another block
physically bounded to CPU in
stacked layers with 1000's
of pins, provides massively
parallel data throughput
by design

[Throughput = number of items and
instructions that a CPU
can process simultaneously]

[Number of samples processed per second
by CPU]

→ SM access data and code from HBM via L2 Cache

→ acts as intermediate level between
off chip and on chip memory
and caches data that be shared
among multiple SMs

→ SM has its own

1. L1 Cache

2. Shared Memory (SRAM)

→ managed by GPU itself

→ faster than HBM, as on chip memory

→ managed by program

→ NVLink :- high bandwidth interconnect for GPUs to connect to each other.

PCIe bus :- helps GPUs to talk to outside world

[common on motherboards to transfer data]

→ Single Node :- 8 GPUs packed together

→ GPU Capability :-

① Computi Performance :- Measured by Number of Trillion float Operations per Second (TFLOPS)

② GPU memory need to store model parameters

ex: GPT : 175B parameters

Storage : 350 GB in fp16

③ Memory Bandwidth :- GB/s → Speed of bytes movement GPU to processing unit

T4 Graphics Card :- (2018)

(2018)

65 TFLOPS

40 SMs

64 KB L1 Cache

4 MB L2 Cache with

1.8 TB/s bandwidth

16 GB HBM with

300 GB/s bandwidth

A100 Graphics Card - 2020

312 TFLOPS

108 SMs - 192 KB L1 Cache

40 MB - L2 Cache

80 GB - HBM with 1.55 TB bandwidth

GPU Capabilities Grow Exponentially Fast. ↑

Some bottlenecks are :-

① Compute bound :- time spent on arithmetic operation > time spent on other operations (ex:- memory access)

② Memory Bound :- time taken by memory access > computation time
ex:- element wise operations (activation fn, dropout)
reductions (sum, softmax, normalization)

③ Overhead bound :- everything else, such as communication bound, interpreter bound, ↑

Arithmetic Intensity :- balance measured between ①, ②
ops: bytes → ratio : to know if we are in a compute or memory bound

Mathematical Example to Understand this

(Q) Find out if "Linear layer forward Pass" on AI or GPU is memory or compute bound?

Ans) Given :- $x \in \mathbb{R}^{B \times d}$ input-batch
 $w \in \mathbb{R}^{d \times d}$ weight matrix

B = batch size

d = embedding dimension

Linear layer = matrix multiplication = xw

Computation requires = $2Bd^2$ FLOPs

$$\text{Compute time} = T_{\text{compute}} = \frac{\# \text{ FLOPs}}{\text{Compute Performance}} = \frac{2Bd^2}{312 \times 10^{12}} \text{ s}$$

$2d^2$ bytes = need to read this from memory to load weight matrix w

let $B \ll d$

$$T_{\text{memory}} = \frac{\# \text{ bytes}}{\text{memory bandwidth}} = \frac{2d^2}{1.55 \times 10^{12}} \text{ s}$$

$$\text{Arithmetic Intensity} = \frac{\# \text{ flops}}{\# \text{ bytes}} \approx \frac{\text{compute performance}}{\text{memory bandwidth}}$$

$$\frac{T_{\text{compute}}}{T_{\text{memory}}} = \frac{B}{200}$$

hence, B is (batch size) smaller than 200
this is memory bounded

Increasing B to greater than 200, increases the computational time while keeping memory transfer time constant then it will be compute bound
Scenario

High Level Algorithmic Optimization

$$\text{Attention}(Q, K, V) = \text{Softmax}\left[\frac{QK^T}{\sqrt{d}}\right] \cdot V$$

$$Q \in \mathbb{R}^{L \times d}$$

$$K, V \in \mathbb{R}^{M \times d}$$

$L, M = \text{seq length.}$

$$\text{Multihead Attention}(Q, K, V) = [\text{head}_1, \dots, \text{head}_h] \cdot W^O$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad i = 1, 2, \dots, h$$

$W^Q, W^K, W^V, W^O = \text{learnable Parameters}$

If $Q = K$: Multihead Self Attention
If $Q \neq K$: Multihead Cross Attention.

Aut Generation Stages:

① Prefill

② Auto regressive Decoding

* Fact about GPU :- GPU has Computational Capacity orders of magnitude higher than memory bandwidth.

KV Cache :-

- we can Cache k and v values in the process of generation to improve efficiency, reduce redundant computational requirements [for long sequences] is called KV Cache
- In each time step, the k and v computed for the last token are added to the cache to reuse later step.
- [Issue] To store this KV Cache we need lot of HBM Capacity or either load it from CPU (DRAM) which is slower by one or 2 orders of magnitude compare to loading from HBM

Multi Query Attention / Grouped Query Attention :- (MQA)

→ Shares the cached key value pairs across multiple queries (KV cache pairs)

which reduce the memory requirements during next text generation.
leading to higher inference throughput

but this leads to cutting model parameters. and
Quality degradation

Technique to avoid quality degradation :-
Technique interpolates between MHA and MHA
that is GQA \rightarrow Grouped Query Attention.

CPU Again

CUDA Kernel Fusion

CPU \rightarrow performs thousands of operations in parallel
Called as Threads

Each thread has its own memory
Called Register

Threads running on same SM can be
grouped together \rightarrow Thread blocks

to execute at same time (max 1024)

Threads within thread blocks can load data from
Global Memory (HBM) into shared memory
(SRAM) used to write results back

Wraps :- when SM is given one or more thread blocks
to execute it partitions them to wraps.

wrap = set of 32 threads

Grid \div multiple thread blocks combined to form

Kernel : All the computations are defined in kernels

Small c++ functions executed multiple times in parallel by different threads launched as a grid

Tensor Cores : Specialized units for matrix multiplications

Kernel Fusion : Implementation of kernel which performs multiple computations at once without extra memory access

Custom CUDA kernels \rightarrow Triton
Pallas