

HPML

- * Scaling Efficiency $\Rightarrow E = \frac{t_{\text{serial}}}{t_{\text{parallel}} * p} \leq 1$
- * Strong Scaling \Rightarrow same problem size, increasing p.
- Weak Scaling \Rightarrow \uparrow problem size with \uparrow p.
- Throughput
 - * Averages $\Rightarrow t$ constant \rightarrow Arithmetic mean
 - # operations constant \rightarrow Harmonic mean
 - Speedup / avg. ratio \rightarrow Geometric mean

Benchmarks

Micro kernels, micro benchmarks, kernels, synthetic benchmark, real applications, complexity \rightarrow real workflow

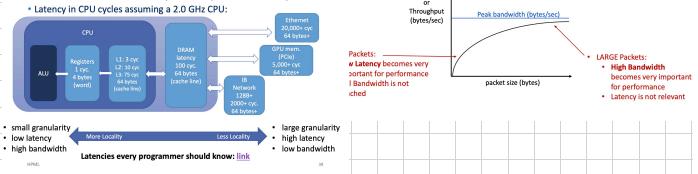
- * Compute time \rightarrow real time, CPU time, non-CPU time
real \Rightarrow sys + user

- * Profiling
 - \leftarrow Sampling \rightarrow estimate
 - counting \rightarrow actual
 - \rightarrow software counters
 - \rightarrow hardware counters

$$\text{Amdahl's law} \rightarrow S(p,s) = \frac{1}{(1-p) + p/s}$$

- * Critical path \rightarrow section of program that accounts for the most time.
- Bottleneck \rightarrow system resource that affects critical path.

Data movement Locality Principle - Latency



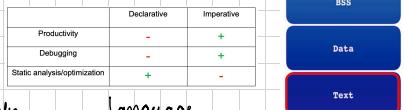
- * Disk performance \Rightarrow I/O Ops/s
For throughput multiply this by block size.

- * Computational graph \rightarrow exposes parallelism, dynamic in Pytorch, static in TF.

Declarative \rightarrow declare all, compile, compute

Impulsive \rightarrow start declaring, keep declaring and computing.

Eager vs graph \Rightarrow when/how program is executed.



- * Python optimizers - PyPy, Numba, cython

Concrete Syntax tree (CST)

\downarrow optimize

Abstract syntax tree (AST)

\downarrow CPU (Hardware)

- * Cython is general purpose, easier to distribute than Numba.

Jit optimizations \Rightarrow fusion, out of order execution and automatic lazy evaluation

\hookrightarrow work scheduling

- * Possible bottlenecks \rightarrow CPU, GPU, memory, network, I/O (disk)

$\not\rightarrow$ not deterministic

- * Profiling \rightarrow Counting, sampling, increasing overhead

\hookrightarrow high overhead

analyze

PyTorch profiler, profile, cProfile, pstats

Tool	Description
nvidia-smi	This tool provides information about GPU utilization, memory usage, and other metrics related to the NVIDIA GPU.
htop	It is a command-line tool that hierarchically displays system processes and provides insights into CPU and memory usage.
top	It is a command-line tool for monitoring CPU usage by displaying I/O statistics of processes running on your system.
gputop	It is a Python-based user-friendly command-line tool for monitoring NVIDIA GPU status. Similar to nvidia-smi, it also displays real-time GPU usage and other metrics in a user-friendly interface.
nvtop	It is a Python-based user-friendly command-line tool for monitoring NVIDIA GPU status. Similar to nvidia-smi, it also displays real-time GPU usage and other metrics in a user-friendly interface.
py-spy	It is a sampling profiler for Python that helps identify performance bottlenecks in your code.



example, an expression like $2 * v$ will be replaced with v .
Strength Reduction: Replacing more expensive operations with simpler ones. For example, $x * 2$ is replaced with $2x$.
Removing Unused Imports: If an import statement is declared but not used, it can be removed.
Optimizing Short Sequences: Replacing short sequences of instructions with more efficient ones.
Loop Invariant Code Motion: Moving loop-invariant code outside the loop.
Loop Lookup Optimization: Optimizing the lookup of variables by replacing global variable lookups with faster local variable lookups when possible.

* Peak FLOPS $\rightarrow \frac{\# \text{cores} \times \text{cycles}}{\text{s}} \times \frac{\text{FLOPs}}{\text{cycle}}$

* Arithmetic Intensity $\rightarrow \frac{\text{FLOP}}{\text{bytes}}$ \Rightarrow #Arithmetic operations / DRAM data

* There are more complex roofing models that involve cache as well as DRAM.

Cache blocking \rightarrow reduces DRAM traffic
 (\nearrow) increase arithmetic intensity

Vectorization \rightarrow increases achieved FLOPs
 (\uparrow) AI remains the same

Measure \rightarrow Analyse \rightarrow optimize

$$\theta = \theta - \alpha \nabla J(\theta) \rightarrow v_t = \gamma v_{t-1} + \alpha \nabla J(\theta - \gamma v_{t-1})$$

$$\text{Neuron} \rightarrow v_t = \gamma v_{t-1} + \alpha \nabla J(\theta - \gamma v_{t-1})$$

\hookrightarrow momentum (usually 0.9)

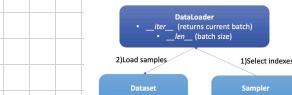
Adagrad \rightarrow separate learning rate for each parameter $\Rightarrow \theta_{t+1} = \theta_t - \frac{\alpha}{\text{diag}(a_t)} \nabla J(\theta_t)$
 \rightarrow good for sparse data but gets smaller over time

Fixed using adadelta \rightarrow only consider recent time window history.

Adam is generalisation of Adagrad.

* Python threading is concurrency, not parallelism.

* `Torch.utils.data.Datasets` \rightarrow `--len--`, `--getitem--`



- * Forking and spawning are two different start methods for new processes.
 - 1) SHWAN method: start new fresh process with minimal resource inherited from parent.
 - 2) FORK method: forks a new process and inherit all resources (files, pipes, etc.)
 - 3) FORK SERVER method: a server-process is created that forks new child processes whenever a connection arrives.
 - 4) MINIMAL method: minimal set of resources inherited
 - In general SHWAN and FORK SERVER methods are safer but slower

- * Useful Parameters:
 - batch_size: batch size
 - sampler: define which sampler to use
 - pin_memory: copy into CUDA pinned memory before returning the data
 - shuffle: shuffle data at each epoch
- * Available samplers:
 - SequentialSampler
 - RandomSampler
 - SubsetRandomSampler

* with record-function ("#forward##") for benchmarking \rightarrow time-it module

* Cuda extension \rightarrow CUDA runtime \rightarrow CUDA drivers

PTX code is considered as assembly for virtual GPU architecture.

x.cu \rightarrow x.ptx \rightarrow x.cubin \rightarrow execute

* CUDA threads execute in SIMD/SIMT. Thread may take different paths: CUDA cores, streaming multiprocessors (SM), texture processing clusters (TPC), graphics processing cluster (GPC), single-precision (FP32), double-precision (FP64).

Threads \in Blocks (\in grids) \rightarrow 1 grid/kernel

2³⁰ \rightarrow 3³⁰ generally 32

Thread synchronisation \rightarrow explicit barrier
 \uparrow
Warp as scheduling units

Streaming multi-processors (SM) \rightarrow maintains thread/block idx, manages thread execution

CUDA context for memory access, host code can be blocking/non-blocking

CUDA streams runs parallelly utilizing hardware

Shared memory is "on-chip" while local/global memory is "off-chip". Local is regular spills

Constant memory is off chip read only from kernel, aggressively cached to on chip.

* CUDA UVM - unified virtual memory \rightarrow if misaligned need 2 mem transactions

\uparrow reuse by threads \rightarrow coalesce if mem tightly packed

Tiling \rightarrow move tile from global to shared \uparrow happens if you go element n for all columns

* PyTorch CUDA ecosystem \rightarrow Extension and tools (Image and text), distributed training, mixed-precision, profiler

is a form of Parameter Server/All reduce \uparrow ring

* Parallelism approaches \rightarrow model tiling, Pipelining, Data tiling, Hybrid stage time = time without pipelining \uparrow

number of pipeline stages





P: number of processes N: total number of model parameters

- PS [centralized reduce]
 - Amount of data sent by PS by (P-1) learner processes: $N(P-1)$
 - Amount of data sent by PS to learners: $N(P-1)$
 - Amount of data sent by PS to learners: $N(P-1)$
 - Total communication cost at PS process is proportional to $2N(P-1)$

Ring All-Reduce (decentralized reduce)

- Scatter-reduce: Each process sends $\frac{N}{P}$ amount of data to (P-1) learners
 - Total amount sent (per process): $N(P-1)^2/P$
- AllGather: Each process again sends $N(P-1)^2/P$ amount of data to (P-1) learners
 - Total communication cost per process is $2N(P-1)/P$

* Allreduce got low comm overhead, and PS has to wait for staggers.

* Downpour SGD \rightarrow lots of communication, might not converge (momentum issues), no mechanism to remain in the centre.

\rightarrow Adagrad within each server shard leads to improvement.

$$\text{Decrease } \alpha \text{ by mean staleness } \alpha = \frac{\alpha_0}{\sigma} \quad \eta_j = \min \left\{ \frac{C}{\|w_j - w_{\tau(j)}\|_2^2}, \eta_{\max} \right\}$$

* Increasing power efficiency

\rightarrow Compression, quantisation, compilation

Post training quantisation Quantisation aware training

$$\begin{aligned} \text{zero point } & \Delta = \frac{2 \max(|x|)}{2^n - 1} \quad \hat{x} = \text{round}(x/\Delta) - \Delta \\ \hat{x} = \text{round}(x/\Delta) - \Delta & \Delta = \frac{\max|x| - \min|x|}{2^n - 1} \\ & \Delta = \frac{\max|x| - \min|x|}{2^n - 1} \end{aligned}$$

* Useful when distribution is not normal
Codewords are selected using clustering or training process.

Apply kmeans and fine tune centroids
Just saves storage weight.

* FP32 \rightarrow INT8 \rightarrow use scale factor for conversions

Symmetric quantization Asymmetric quantization

$WX \approx S_W(W_{\text{int}})S_X(X_{\text{int}})$ $WX \approx S_W(W_{\text{int}} - z_W)S_X(X_{\text{int}} - z_X)$

* You should simulate quantisation \rightarrow quantify and de-quantify weights.

$$s = \frac{x_{\max} - x_{\min}}{q_{\max} - q_{\min}}$$

$$q = \text{clip}\left(\text{round}\left(\frac{x}{s}\right) + \tau, q_{\min}, q_{\max}\right)$$

$$\text{argmin}_{q_{\min}, q_{\max}} \ell(X, \hat{X}(q_{\min}, q_{\max}))$$

$$\text{MSE}$$

$$q_{\min} = \min(\beta - \alpha\gamma)$$

$$q_{\max} = \max(\beta + \alpha\gamma)$$

$$z = \text{round}(q_{\min} - x_{\min})$$

$$\hat{x} = s \cdot (q - z)$$

$$\begin{aligned} \text{Cross-entropy} & \\ \text{qmin} = \min X & \\ \text{qmax} = \max X & \end{aligned}$$

* Flash attention \rightarrow minimizes I/O \rightarrow read 3 blocks at once

2. FlashAttention: main ideas

FlashAttention (Dao et al., 2022) is an exact rewrite of attention that is:

- IO-aware: minimizes reads/writes to slow global memory.
- Uses tiling so everything fits in on-chip SRAM (shared memory / registers).
- Uses a streaming softmax so you never store the full attention matrix.

* Scaling choice \rightarrow dataset size, model size
Constraints \rightarrow hardware

* Parameter efficient fine tuning \rightarrow only train a subset of weights (selected)

\rightarrow Reparameterization - LORA

\rightarrow Additive - Adapter, prompt tuning

* Constant Expressions: Evaluating constant expressions at compile time rather than at runtime. For example, an expression like $2 + 2$ will be replaced with 4.

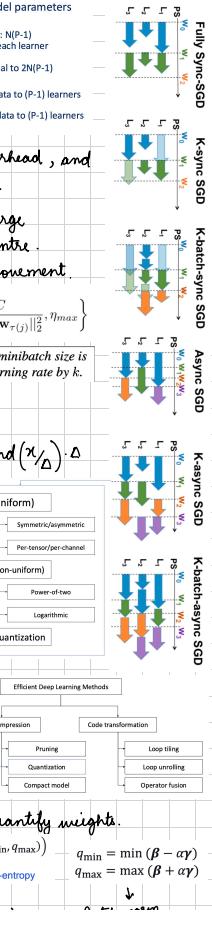
* Strength Reduction: Replacing more expensive operations with simpler ones. For example, $x * 2$ would be replaced with $2x$.

* Removing Unused Imports: If an import statement is declared but not used, it can be removed.

* Optimizing Short Sequences: Replacing short sequences of instructions with more efficient ones.

For example, replacing a sequence of LOAD_CONST instructions with a single BUILD_TUPLE instruction.

* Variable Lookup Optimization: Optimizing the lookup of variables by replacing global variable lookups with faster local variable lookups when possible.



* Pattern based pruning \rightarrow N:M, Channel pruning

Magnitude based pruning \rightarrow absolute values, L_1 norm, L_2 norm \rightarrow element wise, row wise, column wise

Scaling based pruning \rightarrow filters / output channels with small scaling would be pruned.

Percentage of zero based pruning \rightarrow because relu produces zeroes

Regression based pruning \rightarrow think of channels as coefficients

Burning ratios \rightarrow sensitivity analysis \rightarrow might not be optimal because of channel use iterative interactions.

Regularization while fine-tuning - penalize non-zero parameters, encourage smaller parameters.

Knowledge \rightarrow Relation based, feature based, response based.

$$q_i = \frac{e^{z_i/T}}{\sum_j e^{z_j/T}}$$

Soft labels give information about the wrong classes as well.

* FitNets \rightarrow train deeper and thinner students \rightarrow learn from hint layers to guided layers

* Distillation \rightarrow online, offline, self

Attention(Q, K, V) = $\text{softmax}\left(\frac{QK^T}{\sqrt{dk}}\right)V$

* Feed forward nn is easy to parallelize.

This is our sentence* 2) We embed each word* 3) Split into 8 heads. We multiply X or V with weight matrices 4) Calculate attention using the resulting Q & K matrices 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output Z layer

n all encoders other than R0, don't need embedding. I start directly with the output encoder right below this one

R0

W^O

W^Q

W^K

W^V

Q_1

K_1

V_1

Z_1

$...$

W^O

W^Q

W^K

W^V

Q_7

K_7

V_7

Z_7

W^O

Z

W^O

Z