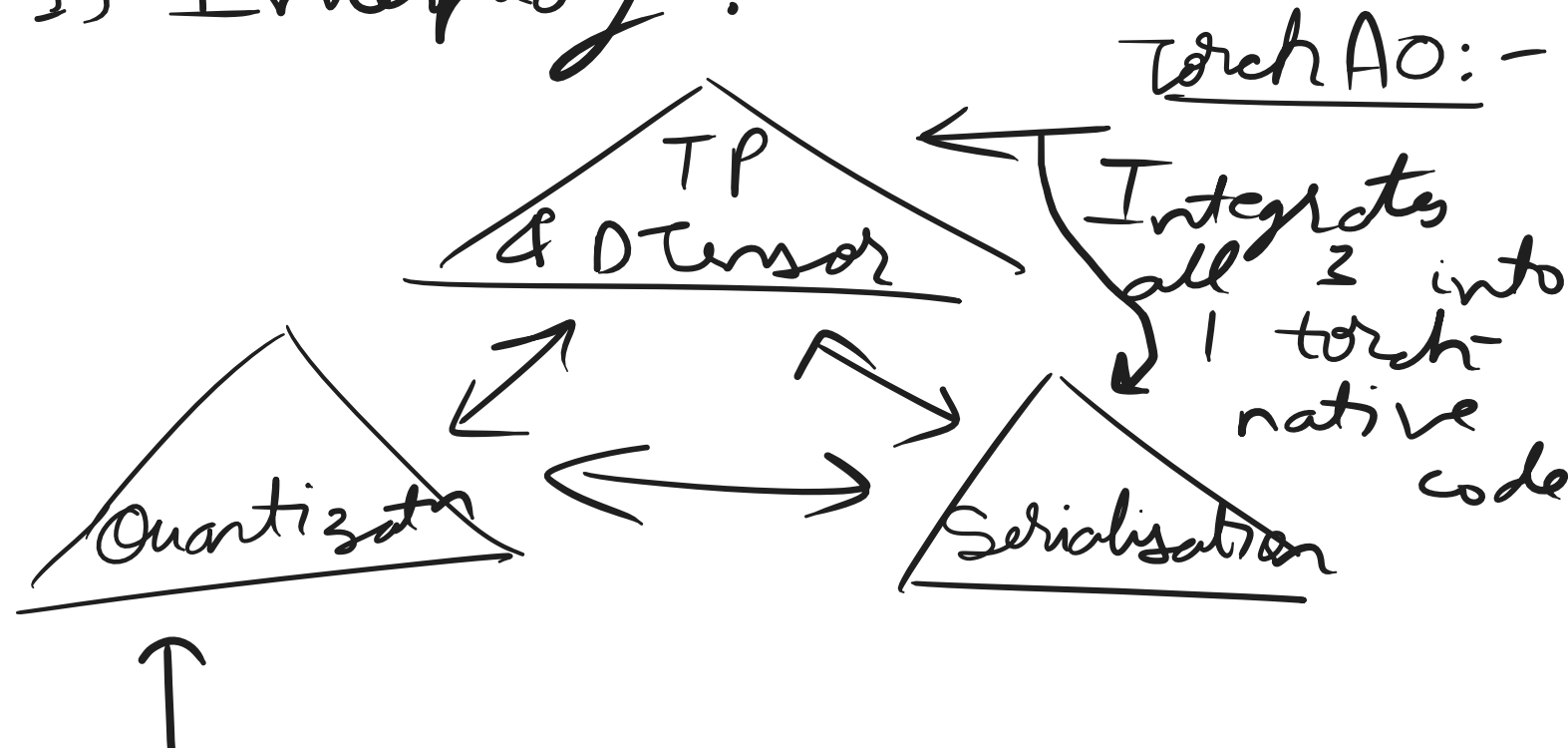


Goals:- 1) Large batch size (memory)
2) Better speed (compute)

Problems:- 1) Perf drops w/t large batch size
2) Types of quantization are limited
 ↳ ex:- Asymmetric Quant only kernels

3) Interplay :-



GemLite
provides new types of quantization

2) Perf for large batch size

#Experiment :- Llama 3.1 - 8B
w/t diff batch sizes & TP sizes.

Compute :- 8X H100 for decode

Baseline → BF16
torch.compile model

Batch size : 1, 32

TP size : 1, 4

Quant → 1) int4 weight only w/t group size 64
2) fp8 per row dynamic quant

(batch sz, tp) (1.21x)
(1,1) → 131 tkr/s → 166 tkr/s

(32,4) → 5575 tkr/s (1.10x) → 6154 tkr/s

(32,1) → 2799 tkr/s (1.28x) → 3586 tkr/s

Supports both H100 & A1000

* GemLite supports both packed & non-packed bit formats

(uint1, uint2 ...)

these are packed together for storage.

GemLite kernels Adapt across →

1) Small Batch size } GEMV for batch 1
 → Memory Bound }
2) Large Batch size } GEMM
 → Compute Bound }
 2-64 batch ⇒ GEMM-SPLITK

Single Sample Inference :-

* GEMV kernels → GEMV-SplitK for packed data

Batched → GEMM - SplitK

* 32-64 → Compute Bound

Optimal split hyperparam (K)

→ Range: 1-16

→ If K=1 ⇒ Normal GEMM kernel

↓
suitable for compute bound setting of 32 to 64 batch size

* Quantize first, shard later.

Follow-up :-

(1) GEMM-SPLITK

★ paper cause it forms basis of GEMMLITE

(2) Binary & Ternary ANNs: efficient MatMul paper.