# Dynamic Precision Scheduling (DPS) for Efficient LLM Inference

Leonid Alexeyev, Jonathan Huang, Sujin Lee, Darin Mao, Devam Shrivastava

# Table of Contents

# Motivation

- **LLMs are powerful but computationally expensive**
    - high inference latency and memory cost
- **Quantization** reduces the precision of the floating point tensors
    - theoretically, this reduces the computational cost
- Generally, quantization is performed on entire models
- Quantize some layers only? Yes, shown in LSAQ paper
    - can we automatically optimize before inference time?

We want to build a system to autonomously determine which layers to quantize, to achieve a balance of performance and output quality

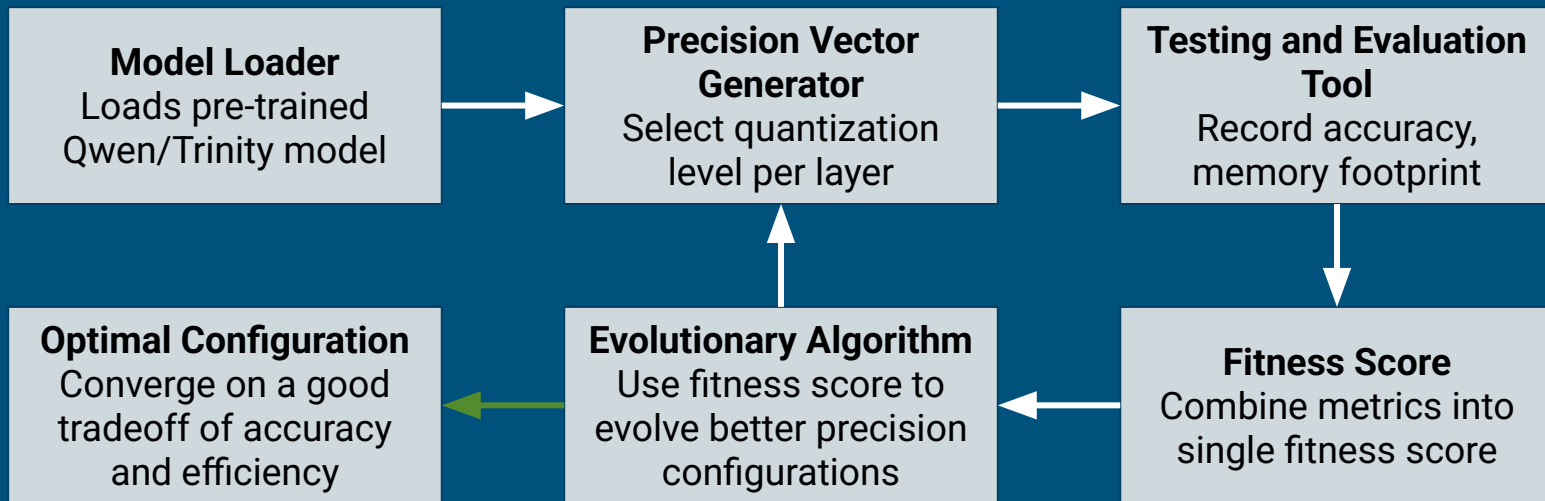- The search space is exponential in size (quantization levels ^ layers), so our system needs to be smarter

# System Architecture

# System Architecture

# System Architecture

1. **Model Management Layer:**

   - Loads the base model, applies per-layer bit-width settings, runs inference, and collects evaluation metrics.

2. **Precision Vector System**

   - Creates and manages per-layer precision vectors; for example, [8, 4, 4, 16, 8, 8, …]

3. **Evaluation Pipeline**

   - Run the quantized model on dataset prompts and computes accuracy, latency, and memory footprint.

# System Architecture

4. **Fitness Module**

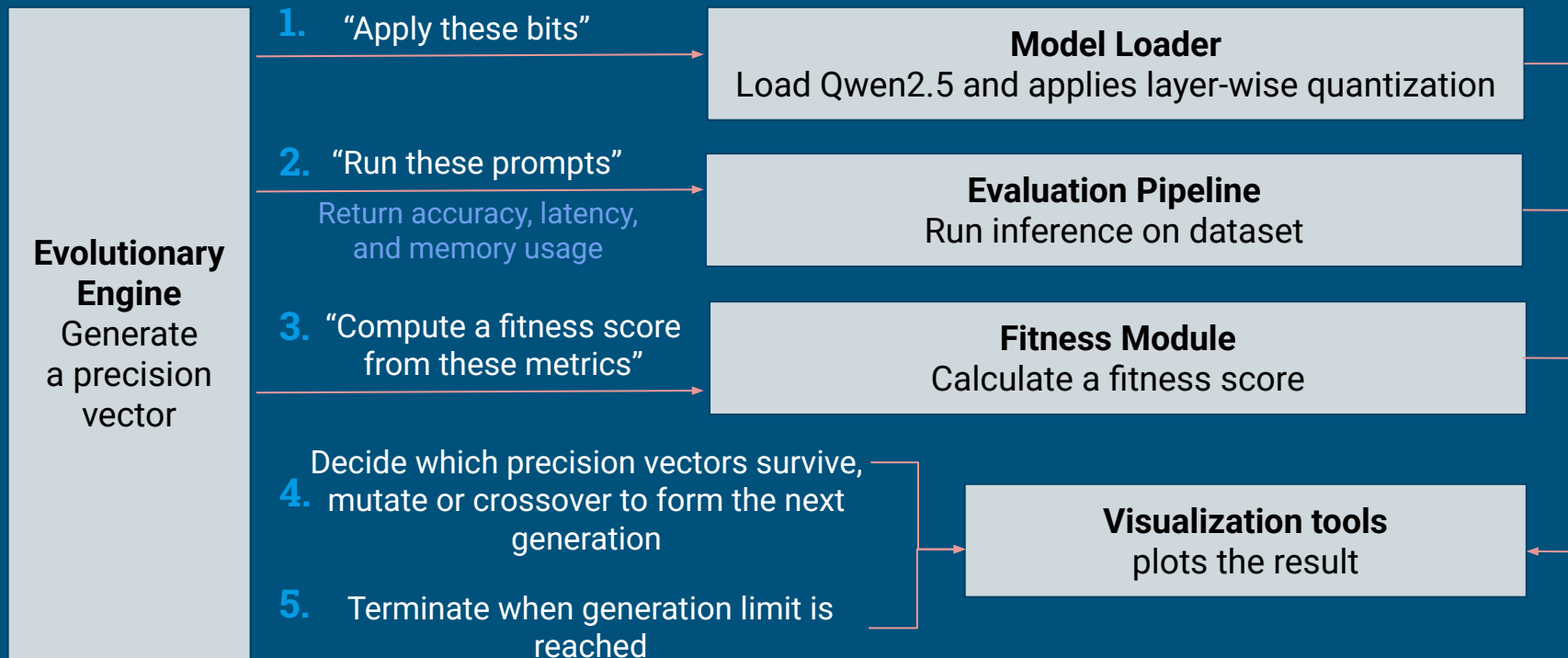   - Convert evaluation metrics into a single numerical fitness score for optimization.

5. **Evolutionary Search Engine**

   - Search for Pareto-optimal per-layer precision configurations through evolutionary optimization

6. **Visualization**

   - Produce plots for Pareto frontier (accuracy vs memory), Fitness progression, Bit-width heatmaps and 3D trade-off

# System Architecture

**Evolutionary Engine**
Generate a precision vector

1. "Apply these bits"

2. "Run these prompts"
Return accuracy, latency, and memory usage

3. "Compute a fitness score from these metrics"

4. Decide which precision vectors survive, mutate or crossover to form the next generation

5. Terminate when generation limit is reached

**Model Loader**
Load Qwen2.5 and applies layer-wise quantization

**Evaluation Pipeline**
Run inference on dataset

**Fitness Module**
Calculate a fitness score

**Visualization tools**
plots the result

# System Architecture

## Our components

- Evolutionary search algorithm
- Precision vector system
- Custom quantization wrapper around HF models
- Dataset evaluation and timing pipeline
- LLM-based autograder
- Fitness scoring module
- Pareto frontier visualization
- Precision heatmap visualization

## Open-Source Components

- HuggingFace Transformers (Model loading and inference)
- BitsAndBytes (4-bit, 8-bit quantization backend)
- Qwen2.5-7B-Instruct (base LLM)
- PyTorch (GPU inference)
- Matplotlib/Seaborn

# System Architecture

Experiments were run on 8x H200 GPUs rented from Vast.ai

- 141GB VRAM per GPU
- 428.2 TFLOPS total performance
- Python multiprocessing library used for coordinating model runs and collecting results

Evolutionary Engine

Precision Vectors (Population)

| GPU 0 | GPU ... | GPU 7 |
|---|---|---|
| LLM Layer | LLM Layer | LLM Layer |
| LLM Layer | LLM Layer | LLM Layer |
| **Run** MATH-50 **Evaluate** Acc/Lat | **Run** MATH-50 **Evaluate** Acc/Lat | **Run** MATH-50 **Evaluate** Acc/Lat |

Evolutionary Fitness Aggregation + Selection

# Dataset

# Dataset

- MATH-500: 500 problems across algebra, geometry, number theory, calculus
    - All 500 samples test, no need for train/validation since we are using this for evaluation
    - Columns are formatted in LaTeX



- "MATH-50": A subset created by our team for faster evaluation, due to resource constraints

# Dataset Example 1: test/precalculus/807.json

## Problem

Convert the point $(0, 3)$ in rectangular coordinates to polar coordinates. Enter your answer in the form $(r, \theta)$, where $r > 0$ and $0 \le \theta < 2\pi$.

## Solution

We have that $r = \sqrt{0^2 + 3^2} = 3$. Also, if we draw the line connecting the origin and $(0, 3)$, this line makes an angle of $\frac{\pi}{2}$ with the positive $x$-axis.

Therefore, the polar coordinates are $\boxed{\left(3, \frac{\pi}{2}\right)}$.

Answer:
(3, π/2)

Subject:
Precalculus

Level:
2

# Dataset Example 2: test/number_theory/515.json

## Problem

What is the smallest positive perfect cube that can be written as the sum of three consecutive integers?

## Solution

The sum of three consecutive integers takes the form $(k-1)+(k)+(k+1) = 3k$ and hence is a multiple of 3. Conversely, if a number $n$ is a multiple of 3, then $n/3 - 1$, $n/3$, and $n/3 + 1$ are three consecutive integers that sum to give $n$. Therefore, a number is a sum of three consecutive integers if and only if it is a multiple of 3. The smallest positive perfect cube that is a multiple of 3 is $3^3 = \boxed{27}$.

Answer:
27

Subject:
Number Theory

Level:
3

# Dataset Example 3: test/intermediate_algebra/1197.json

## Problem

Let $p(x)$ be a polynomial of degree 5 such that

$$p(n) = \frac{n}{n^2 - 1}$$

for $n = 2, 3, 4, \ldots, 7$. Find $p(8)$.

## Solution

Let $q(x) = (x^2 - 1)p(x) - x$. Then $q(x)$ has degree 7, and $q(n) = 0$ for $n = 2, 3, 4, \ldots, 7$, so

$$q(x) = (ax + b)(x - 2)(x - 3) \cdots (x - 7)$$

for some constants $a$ and $b$.

We know that $q(1) = (1^2 - 1)p(1) - 1 = -1$. Setting $x = 1$ in the equation above, we get

$$q(1) = 720(a + b),$$

so $a + b = -\frac{1}{720}$.

We also know that $q(-1) = ((-1)^2 - 1)p(-1) + 1 = 1$. Setting $x = -1$ in the equation above, we get

$$q(-1) = 20160(-a + b),$$

so $-a + b = \frac{1}{20160}$. Solving for $a$ and $b$, we find $a = -\frac{29}{40320}$ and $b = -\frac{3}{4480}$. Hence,

$$q(x) = \left(-\frac{29}{40320}x - \frac{3}{4480}\right)(x - 2)(x - 3) \cdots (x - 7)$$

$$= -\frac{(29x + 27)(x - 2)(x - 3) \cdots (x - 7)}{40320}.$$

In particular,

$$q(8) = -\frac{(29 \cdot 8 + 27)(6)(5) \cdots (1)}{40320} = -\frac{37}{8},$$

so

$$p(8) = \frac{q(8) + 8}{8^2 - 1} = \boxed{\frac{3}{56}}.$$
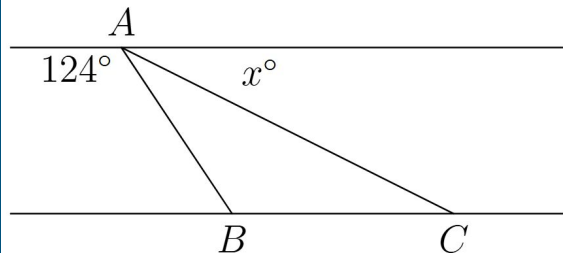
## Answer:
3 / 56

## Subject:
Intermediate Algebra

## Level:
5

# Dataset Example 4: test/geometry/434.json

## Problem

$\overline{BC}$ is parallel to the segment through $A$, and $AB = BC$. What is the number of degrees represented by $x$?
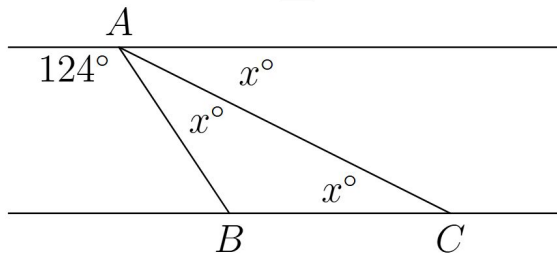


## Solution

Angle $\angle BCA$ and the angle we're trying to measure are alternate interior angles, so they are congruent. Thus, $\angle BCA = x°$.

Since $AB = BC$, we know that $\triangle ABC$ is isosceles with equal angles at $C$ and $A$. Therefore, $\angle BAC = x°$.

The sum of the three angles at $A$ is $180°$, since they form a straight angle. Therefore,

$$124 + x + x = 180,$$

which we can solve to obtain $x = \boxed{28}$.



## Answer:
28

## Subject:
Geometry

## Level:
1

# Dataset Example 5: test/counting_and_probability/134.json

## Problem

In how many ways can 8 people sit around a round table if 3 of the people – Pierre, Rosa, and Thomas – all want to sit together? (Two seatings are considered the same if one is a rotation of the other.)

## Solution

First choose three consecutive seats for Pierre, Rosa, and Thomas. It doesn't matter which three consecutive seats that we choose, since any three such seats can be rotated to any other such seats. Once the three seats are chosen, there are $3!$ ways to seat the three friends there. The other five seats are for the other five people, so there are $5!$ ways to seat them there. The answer is $3! \times 5! = \boxed{720}$.

**Answer:**
720

**Subject:**
Counting and Probability

**Level:**
4

# Experimental Results

# Experimental Setup

- Environment: Vast.ai, Pace/ICE
- GPU: 8x H200 (Vast.ai), 2x H200 (Pace/ICE)
- Model: Qwen2.5-7B-Instruct, Trinity-Mini, Mistral-7B
- Datasets: MATH-50
- Quantization Levels: 4, 8, 16 bits
- Metrics: Accuracy, inference latency, VRAM usage, loss

# Experiment Screenshots (All models)



generation progress

status of 8 GPUs

# Experimental Result #1

- test/precalculus/807.json
- answer is (3, π/2)


- fully int8 quantized model produces the correct answer

To convert the point from rectangular (Cartesian) coordinates $(x, y) = (0, 3)$ to polar coordinates $(r, \theta)$, we use the following formulas:

1. The distance from the origin (radius), $r$, is given by:

$$r = \sqrt{x^2 + y^2}$$

2. The angle $\theta$ is given by:

$$\theta = \arctan\left(\frac{y}{x}\right)$$

However, we need to adjust this value based on the quadrant of the point.

Let's calculate these values for our point $(0, 3)$.

**Step 1: Calculate $r$**

$$r = \sqrt{0^2 + 3^2} = \sqrt{9} = 3$$

**Step 2: Calculate $\theta$**

For the point $(0, 3)$, since $x = 0$ and $y > 0$, the point lies on the positive y-axis. Therefore, the angle $\theta$ is $\frac{\pi}{2}$ radians (or 90 degrees).

Thus, the polar coordinates for the point $(0, 3)$ are:

$$(r, \theta) = (3, \frac{\pi}{2})$$

Final Answer: (3, π/2)

# Experimental Result #2

- test/intermediate_algebra/1994.json
- answer is p - q

- fully int8 quantized model produces the correct answer



The inner sum is simply the sum of $\frac{1}{n^3}$ taken $n-1$ times:

$$\sum_{j=1}^{n-1} \frac{1}{n^3} = (n-1) \cdot \frac{1}{n^3} = \frac{n-1}{n^3}.$$

Therefore, the sum $S$ becomes:

$$S = \sum_{n=2}^{\infty} \frac{n-1}{n^3}.$$

We can split the fraction $\frac{n-1}{n^3}$ into two simpler fractions:

$$\frac{n-1}{n^3} = \frac{n}{n^3} - \frac{1}{n^3} = \frac{1}{n^2} - \frac{1}{n^3}.$$

So, the sum $S$ can be written as:

$$S = \sum_{n=2}^{\infty} \left( \frac{1}{n^2} - \frac{1}{n^3} \right).$$

We can separate this into two sums:

$$S = \sum_{n=2}^{\infty} \frac{1}{n^2} - \sum_{n=2}^{\infty} \frac{1}{n^3}.$$

Notice that $\sum_{n=2}^{\infty} \frac{1}{n^2}$ is just the series for $p$ starting from $n=2$, which is $p-1$. Similarly, $\sum_{n=2}^{\infty} \frac{1}{n^3}$ is just the series for $q$ starting from $n=2$, which is $q-1$.

Thus, we have:

$$S = (p-1) - (q-1) = p - q.$$

Therefore, the sum $\sum_{j=1}^{\infty} \sum_{k=1}^{\infty} \frac{1}{(j+k)^3}$ in terms of $p$ and $q$ is:

$$p - q.$$

# Experimental Result #3

- test/algebra/1072.json
- answer is 243/625


- fully int8 quantized model doesn't produce an answer
- half quantized model produces the wrong answer
- fully unquantized model produces the right answer



Now, calculate $\left(\frac{3}{5}\right)^7$:

$$\left(\frac{3}{5}\right)^7 = \frac{3^7}{5^7} = \frac{2187}{78125}$$

Substitute this back into the expression for $a_8$:

$$a_8 = \frac{125}{9} \cdot \frac{2187}{78125} = \frac{125 \cdot 2187}{9 \cdot 78125}$$

Simplify the numerator and the denominator:

$$125 \cdot 2187 = 273375$$

$$9 \cdot 78125 = 703125$$

Thus,

$$a_8 = \frac{273375}{703125}$$

Simplify the fraction by finding the greatest common divisor (GCD) of 273375 and 703125. The GCD is 375:

$$\frac{273375 \div 375}{703125 \div 375} = \frac{729}{1875}$$

Further simplify by dividing both the numerator and the denominator by their GCD, which is 27:

$$\frac{729 \div 27}{1875 \div 27} = \frac{27}{69.4444} \approx \frac{27}{69.4444} = \frac{27}{69.4444} = \frac{27}{69.4444} = \frac{27}{69.4444}$$

Finally, the simplified form is:

$$\frac{27}{69.4444} = \frac{27}{69.4444} = \frac{27}{69.4444} = \frac{27}{69.4444}$$

Therefore, the eighth term of the sequence is:

$$\frac{27}{69.4444}$$

# Model Selection

These preliminary experimental results show that the evaluated model is a good fit for our project

- The VRAM consumption is low enough to load multiple times on one GPU
    - This is necessary due to how our system produces quantized models
- Different quantization levels produce measurably different accuracy
    - Some models, though around the same size, produce extremely poor accuracy
    - These models would not show any meaningful results

# Measurement Result #1 (QWEN): Loss vs. Memory

- Total of 160 models were tested
- Steep downward slope on the graph
  - Accuracy increases with very slight increases in memory
  - Pareto frontier lies entirely on low-memory side
- Higher memory budgets do not guarantee better accuracy



Pareto Trade-off: Loss vs Memory

# Measurement Result #1 (Trinity): Loss vs. Memory

- Steep horizontal slope on the graph
  - Accuracy increases with very slight increases in memory
  - Pareto frontier lies entirely on low-loss side
- Trinity: Built for math purposes, causing conflict with dataset
  - But still evident decrease in loss, correlating to better/more accurate responses



Pareto Trade-off: Loss vs Memory

# Measurement Result #1 (Orchestrator) : Loss vs. Memory

- Due to some constraints, the evolutionary search was terminated early
- Loss is not good, regardless of quantization
- Within the explored range, loss appears largely insensitive to memory, suggesting that performance may be more constrained by model
- This makes sense, as math is not what this model is designed for



Pareto Trade-off: Loss vs Memory

# Measurement Result #1 (Mistral) : Cross Entropy Loss vs. Memory

- In this result, quantization did not seem to help with loss (better large models)

- Distinct groups of memory with different sizes

# Measurement Result #2 (QWEN): Loss vs. Latency

- Steep downward slope, but with slight curve outward
  - Loss decreases as latency increases (very slightly)
- In our experimental results, the speed increase from quantization did not have a significant effect on trade off



Pareto Trade-off: Cross Entropy Loss vs Latency

# Measurement Result #2 (Trinity): Loss vs. Latency

- Downward slope
  - Loss decreases as latency increases
- Clear view of low latency and cross entropy loss for Trinity
  - Once again reliant on pretraining of the model

# Measurement Result #2 (Orchestrator): Loss vs. Latency

- Most configurations operate within a very narrow low-latency range

- Cross entropy loss shows relatively high variability

- The single high-latency point acts as an outlier

# Measurement Result #2 (Mistral): Loss vs. Latency

- Vertical line, all models centered around about 25 ms latency
- Some outliers of higher latency for lower cross entropy than main group



Pareto Trade-off: Cross Entropy Loss vs Latency

# Measurement Result #3 (All models): Loss vs. Memory vs. Latency

- Best pretrained models lie in low-latency, low-memory corner
  - Best configuration are faster and smaller, implying stronger quantization
- High-memory models consistently have higher latency without significant accuracy benefit

# Measurement Result #4 (Qwen): Bit Allocation Heatmap

- Most models do not have uniform bit widths
  - However, certain layers consistently get assigned similar bit sizes
- Best-performing models use non-uniform precision vectors



Per-Layer Bit Allocation Heatmap

# Measurement Result #4 (Trinity): Bit Allocation Heatmap

- Most models do not have uniform bit widths
  - However, certain layers consistently get assigned similar bit sizes
- Best-performing models use non-uniform precision vectors



Per-Layer Bit Allocation Heatmap

# Measurement Result #4 (Orchestrator): Bit Allocation Heatmap

- Most models do not have uniform bit widths
  - However, certain layers consistently get assigned similar bit sizes
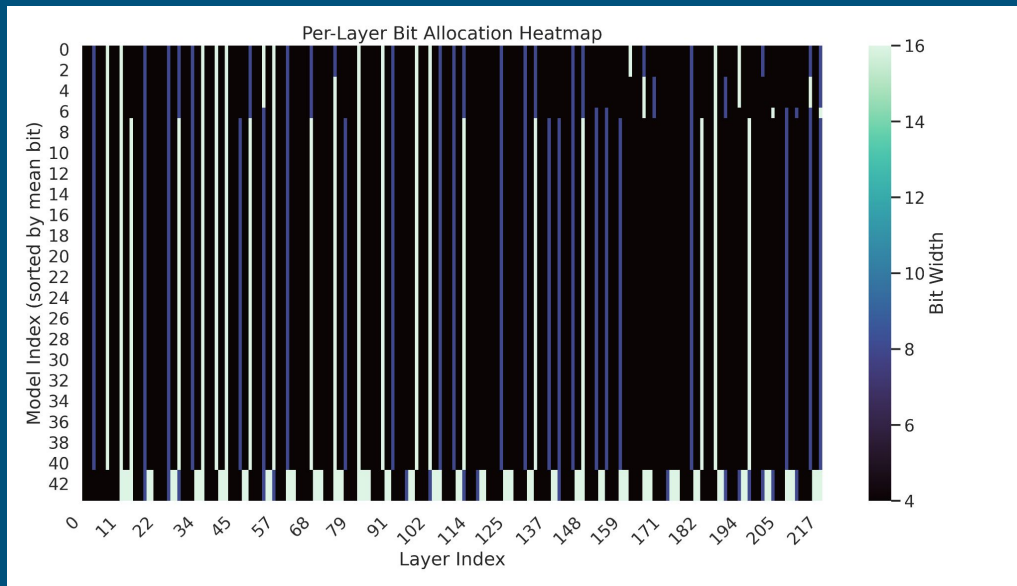- Best-performing models use non-uniform precision vectors

# Measurement Result #4 (Mistral): Bit Allocation Heatmap

- Most models do not have uniform bit widths
  - However, certain layers consistently get assigned similar bit sizes
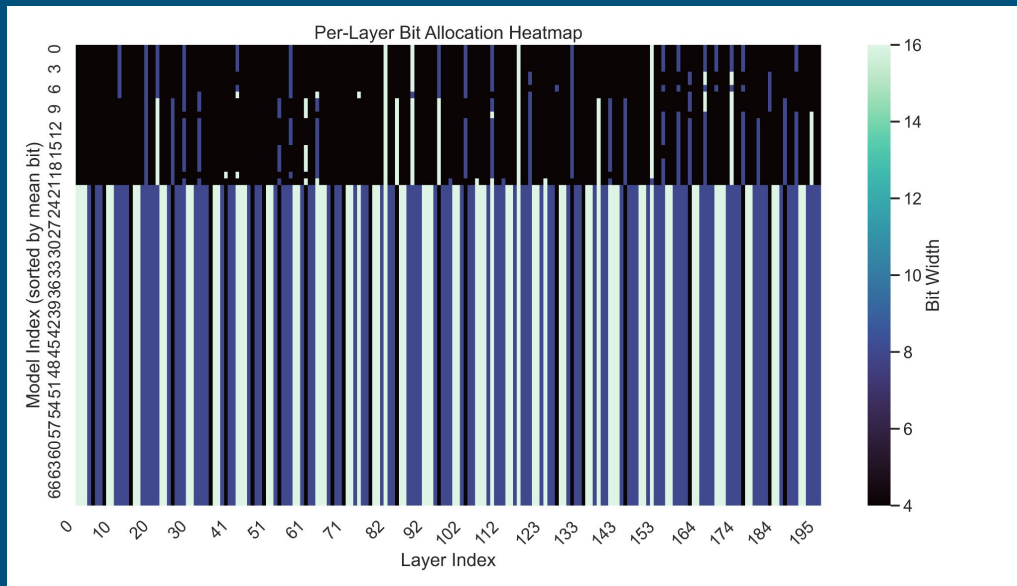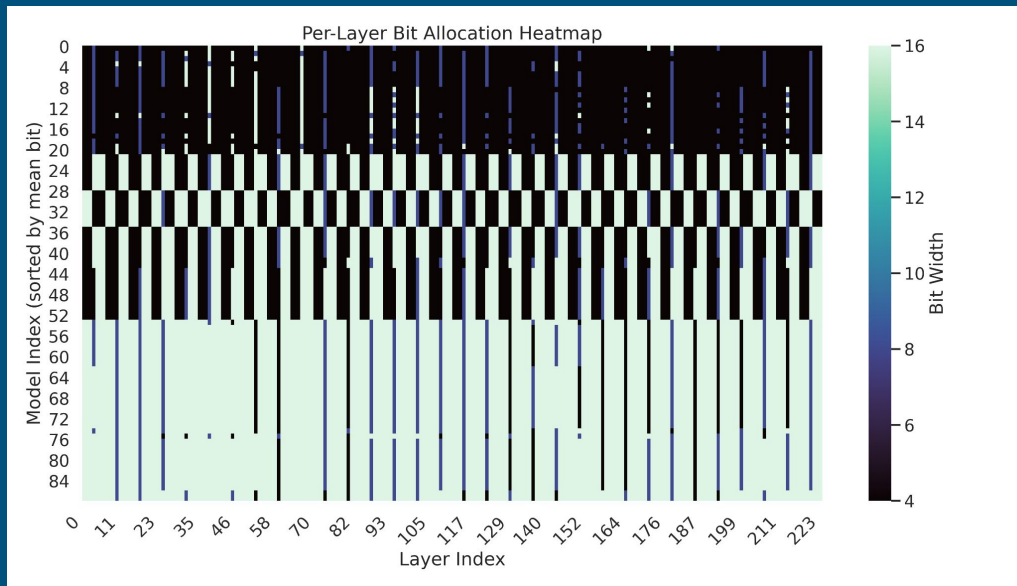- Best-performing models use non-uniform precision vectors

# Project Summary

# Key Contributions

**Successfully created an automated tool that finds pareto-optimal layer quantization configurations for LLMs under arbitrary memory and latency constraints.**

```
1    latency_lut:
2    - key:
3       - matmul
4       - 128
5       - 3584
6       - 3584
7       - 2
8       us: 22.3411201313138
9    - key:
10      - matmul
11      - 128
12      - 3584
13      - 3584
14      - 3
15      us: 22.3411201313138
16   - key:
17      - matmul
18      - 128
19      - 3584
20      - 3584
21      - 4
22      us: 22.3411201313138
```

**LUT Table For Latency Estimation**
Created a LUT table to enable fast, hardware-aware latency estimation.

**Accurate VRAM Computation**
Computed VRAM usage from parameter count x bit width

**Parallel GPU Evaluation**
Parallel Evaluation across multiple GPUs for faster processing

**Pareto Frontier Visualization**
Created and generated 2D/3D pareto frontier plots for easy referencing

# Lessons Learned

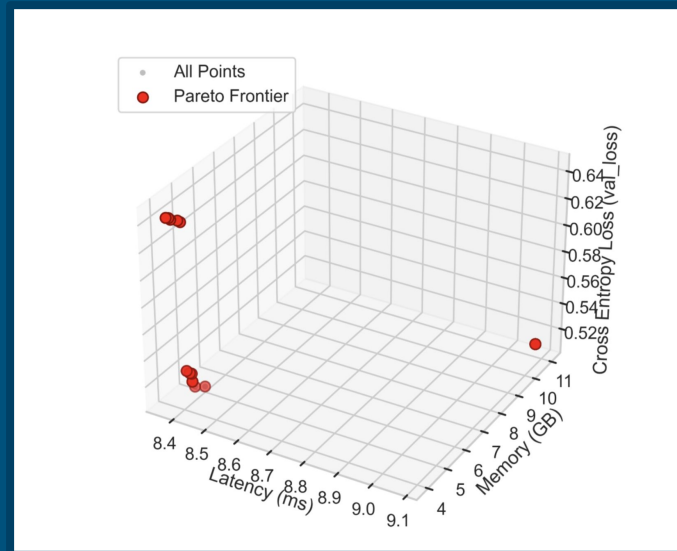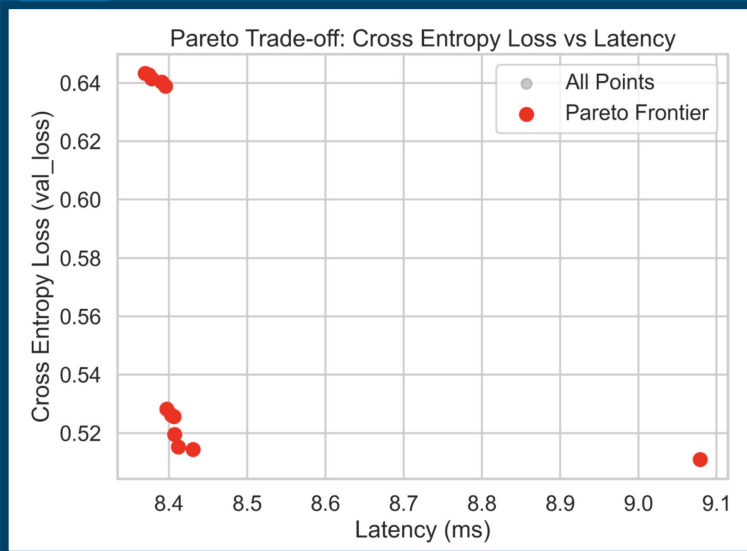**Accuracy vs Cross Entropy Evaluation**

- Cross entropy loss - measures how much the model's next-token predictions diverge from Math-500 gold solutions.
- Accuracy Metric - Directly computes correctness to see if generated answer is correct or not

**Why Cross Entropy?** - Sensitive to small perturbations in quantized weights providing smooth, stable optimization

**Why Accuracy Performed Better** - The model wasn't trained on Math-500, token level cross entropy loss produced noisy comparisons. There was really bad convergence overall. Accuracy offered a more clear ranking of "better vs worse" models.
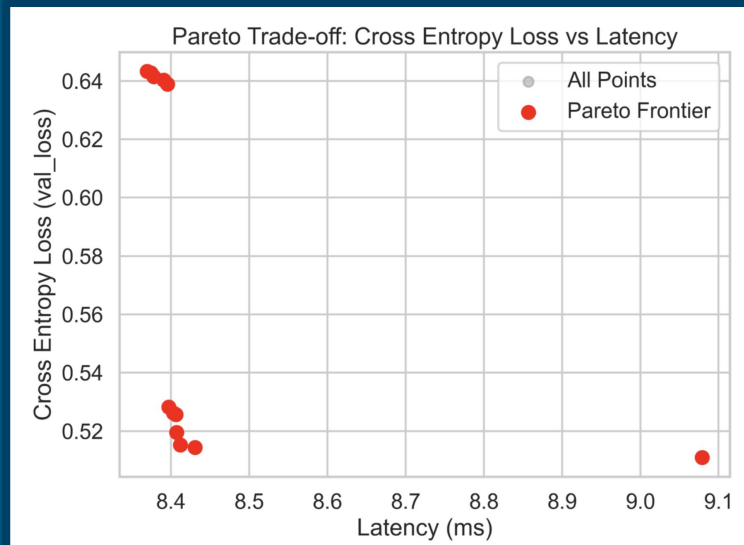
# Lessons Learned - Preliminary Results (QWEN)



Note: You can't see the non-pareto points because they cluster around the Pareto Frontier!

# Lessons Learned - Preliminary Results (QWEN)



Pareto Trade-off: Cross Entropy Loss vs Latency

"n_all": 646.0,
"n_frontier": 81.0,
"frontier_ratio": 0.12538699690402477,
"objectives_all": {
  "latency_ms": {
    "count": 69.0,
    "min": 8.36955105460832,
    "p25": 8.407168718608318,
    "median": 9.07914864820832,
    "p75": 9.07914864820832,
    "max": 9.07914864820832,
    "mean": 8.87077684122861,
    "std": 0.315163011811839
  },
  "memory_gb": {
    "count": 646.0,
    "min": 3.928752128,
    "p25": 4.421451776,
    "median": 4.882038784,
    "p75": 10.676994048,
    "max": 12.497321984,
    "mean": 7.07940405260681,
    "std": 3.006206082828667
  },
  "val_loss": {
    "count": 646.0,
    "min": 0.5109182275003857,
    "p25": 0.5309930532595988,
    "median": 0.5514976450690516,
    "p75": 0.6375011700287736,
    "max": 0.6512670234243483,
    "mean": 0.5792923453484726,
    "std": 0.05284727858420227
  },
},

"bits_all": {
  "counts": {
    "4": 70436,
    "8": 31834,
    "16": 24346
  },
  "fractions": {
    "4": 0.556296202691603,
    "8": 0.251421621280091,
    "16": 0.19228217602830605
  }
},
"bits_frontier": {
  "counts": {
    "4": 6621,
    "8": 5538,
    "16": 3717
  },
  "fractions": {
    "4": 0.4170445956160242,
    "8": 0.3488284202569917,
    "16": 0.23412698412698413
  }
},
"per_layer_all": [
  {
    "layer": 0,
    "mean": 8.910216718266254,
    "std": 5.703835080711555,
    "hist": {
      "4": 351,
      "8": 46,
      "16": 249
    }
  },

"objectives_frontier": {
  "latency_ms": {
    "count": 69.0,
    "min": 8.36955105460832,
    "p25": 8.407168718608318,
    "median": 9.07914864820832,
    "p75": 9.07914864820832,
    "max": 9.07914864820832,
    "mean": 8.87077684122861,
    "std": 0.315163011811839
  },
  "memory_gb": {
    "count": 81.0,
    "min": 3.928752128,
    "p25": 4.296671232,
    "median": 11.024728064,
    "p75": 11.024728064,
    "max": 11.024728064,
    "mean": 8.322950637037035,
    "std": 3.3449665092870706
  },
  "val_loss": {
    "count": 81.0,
    "min": 0.5109182275003857,
    "p25": 0.5109182275003857,
    "median": 0.5109182275003857,
    "p75": 0.5256021411882507,
    "max": 0.6432530486363405,
    "mean": 0.5360111221675021,
    "std": 0.04862506209113528
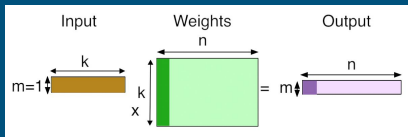  },
},

Summary stats demonstrate that quantization is indeed happening and frontier had higher precision overall.

# More Lessons Learned

## Hardware-Aware ML Model Engineering



- Predicting memory costs from GEMM kernels
- Working with memory bandwidth and quantization overhead
- GPU profiling with warm-up runs, CUDA synchronization, and caching

## Understanding Layer Quantization

```
{
  "layer": 140,
  "mean": 10.56,
  "std": 5.8030139952072695,
  "hist": {
    "4": 72,
    "8": 11,
    "16": 92
  }
},
```

```
{
  "layer": 15,
  "mean": 5.142857142857143,
  "std": 2.4511554959459674,
  "hist": {
    "4": 135,
    "8": 35,
    "16": 5
  }
},
```

- Understanding which groupings of layers are quantization sensitive through evolutionary testing

## Developing an LLM Evaluation Pipeline

```
__init__.py          frontier.py
bnb_quantize.py      group_ties.py
cache.py             logging_utils.py
config.py            lut.py
core_evolutionary.py main_search.py
data_models.py       objectives.py
evaluator.py         utils.py
```

- GPU parallelism
- Cache aware
- Stateful - not needing to load a new model every time due to high latency costs

# Future Improvements

- **Pluggable APIs -** Currently our model requires a lot of extra code
  - bitsandbytes currently requires its own file for layer quantization
    - Potentially could create an API that allows you to pass in a quantization scheme (1 bit, 2 bit, etc).
  - Additional APIs for figures and results
    - Adding an API where you input memory or latency budgets and get returned a clear quantization scheme
- Limited generalizability across diverse model architectures
  - Applying the method to new models requires substantial model-specific implementation
- Clearer results - graphs can tell a story but there must be a better way to present the different layer schemes.

# README Explanation

- Project Overview & Technical Stack: Describes AutoLLMTuner as an evolutionary algorithm framework for LLM quantization,
    - Detailed information about datasets (MATH-500), performance measurement tools (BitsAndBytes, hardware benchmarking, caching), and external resources (Qwen2.5-7B model, frameworks)
- Setup & Usage: Provides installation requirements, configuration files explanation, and instructions for running experiments
    - Hardware benchmarking, evolutionary search, HPC execution, visualization generation
- Experimental Results & Analysis: Comprehensive analysis of test runs showing that optimal solutions achieve 94% accuracy at 5.51 GB memory (vs 12.49 GB baseline),
    - Detailed performance metrics, Pareto analysis, bit allocation strategies, and visualization explanations
- Code Architecture: Documents the project structure with core components (evolutionary algorithm, quantization backend, visualization tools)
    - Explains how evaluation, frontier computation, and results generation are integrated throughout the pipeline

# Open Source References

- HuggingFace Transformers
  (Model loading and inference) - https://huggingface.co/docs/transformers/en/index

- HuggingFace MATH500

  (Dataset for testing) - https://huggingface.co/datasets/HuggingFaceH4/MATH-500

- BitsAndBytes
  (4-bit, 8-bit quantization backend) - https://huggingface.co/docs/bitsandbytes/en/index

- Qwen2.5-7B-Instruct
  (base LLM) - https://huggingface.co/Qwen/Qwen2.5-7B-Instruct

- PyTorch (GPU inference) - https://pytorch.org/

- Matplotlib - https://matplotlib.org/

- Seaborn - https://seaborn.pydata.org/

# Q&A

## Thank You!!