# Fine tune a LLM: how much memory do I need?

**RTX 4090
24 GB**

**tuning a 2B model**

| | |
|---|---|
| 4GB | **Original model** |
| 4GB | **Gradient** |
| 8 GB | **Optimizer** |
| ~6GB | **Activation** |
| | Something else |

2B (FP16) = 2*16/8 = 4GB

Same size with the model

Two tensors: mean + variance

Model architecture, batch size, sequence length etc, ~1.5x model size
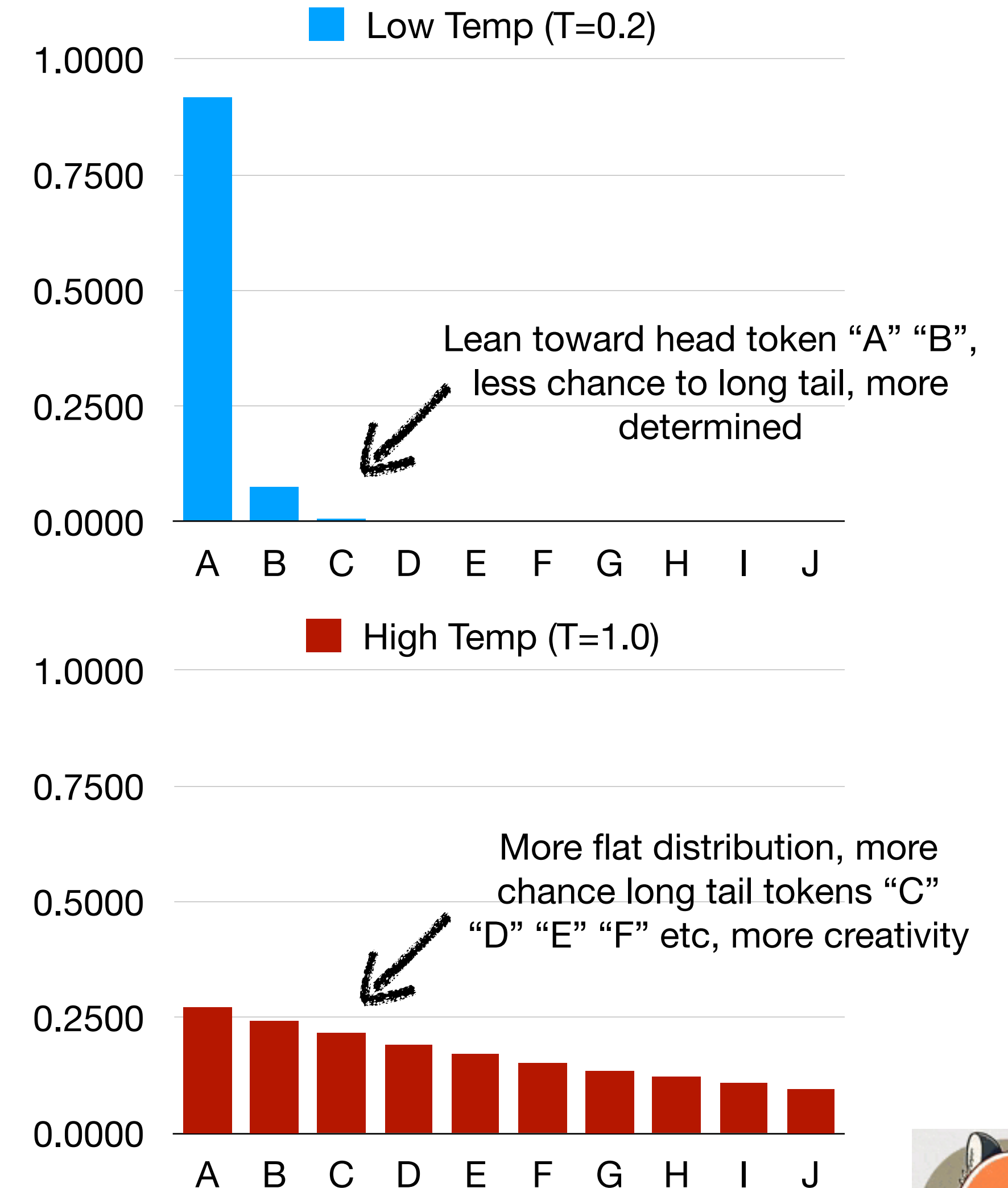
# Understand temperature in LLM

```
response = openai.ChatCompletion.create(
    model='gpt-4o',
    temperature=0.7,
    max_tokens=30,
    messages=[{
        'role': 'user',
        'content': question
    }],
)
```

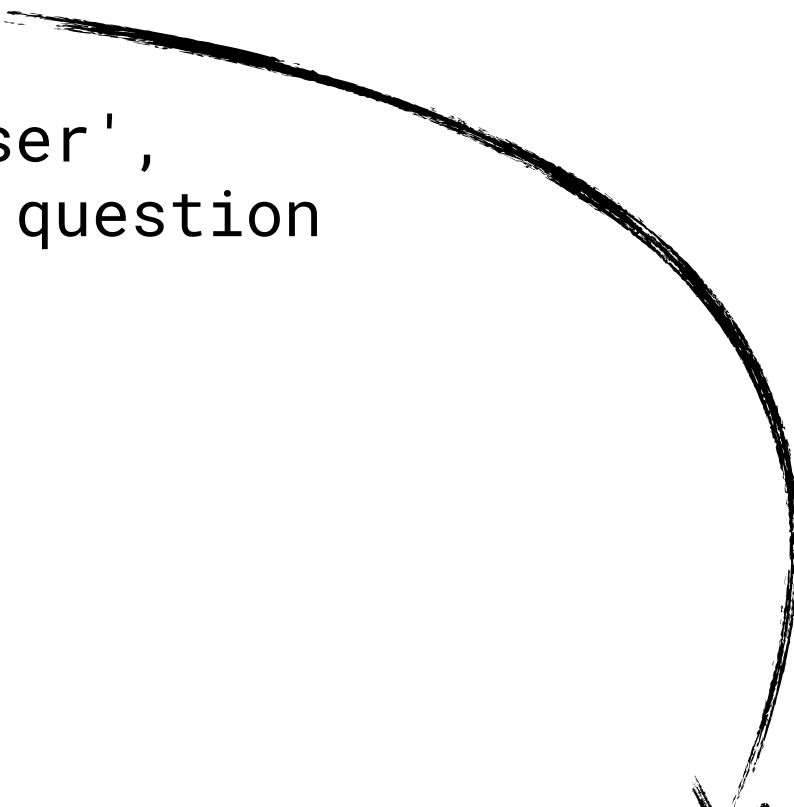$$P(x_i) = \frac{\exp(z_i/T)}{\sum_{j=1}^{N} \exp(z_j/T)}$$

Softmax with Temperature
Hinton et al "Distilling the Knowledge in a Neural Network" 2015
https://arxiv.org/abs/1503.02531

T = 0.2

Low Temp (T=0.2)

Lean toward head token "A" "B", less chance to long tail, more determined

T = 1.0

High Temp (T=1.0)

More flat distribution, more chance long tail tokens "C" "D" "E" "F" etc, more creativity
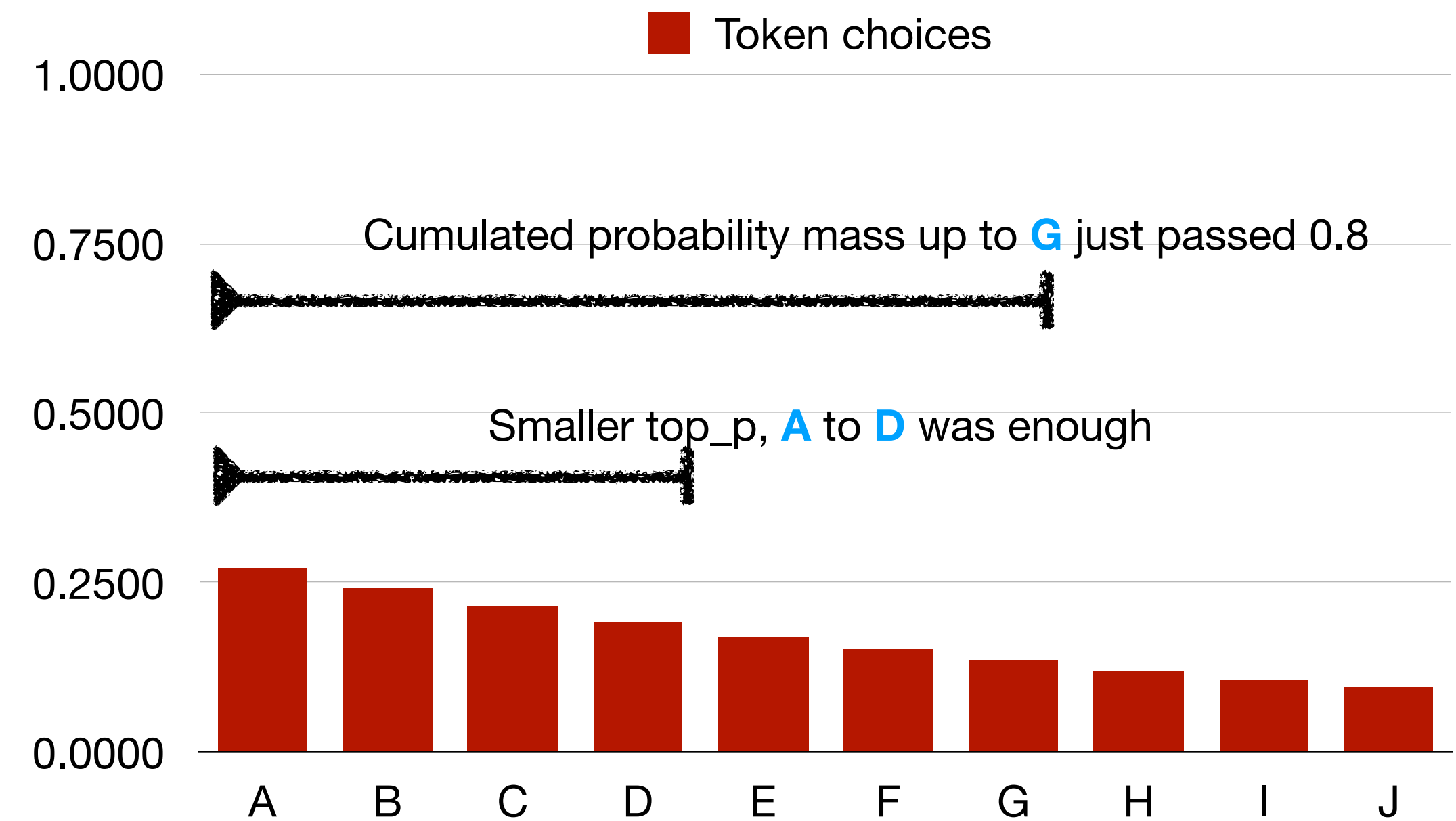
# Understand top_p in LLM

```
response = openai.ChatCompletion.create(
    model='gpt-4o',
    max_tokens=30,
    top_p=0.8,
    messages=[{
        'role': 'user',
        'content': question
    }],
)
```

$$\sum_{x \in V^{(p)}} P(x|x_{1:i-1}) \geq p$$

Top P or "Nucleus sampling"
Holtzman et al, "The Curious Case of Neural Text Degeneration" 2019
https://arxiv.org/abs/1904.09751



Cumulated probability mass up to **G** just passed 0.8

Smaller top_p, **A** to **D** was enough

More flexible than top K to the shape of the
probability distribution, across different contexts.

https://www.linkedin.com/in/liuhongliang/

# Understand Boltzmann distribution and neural networks

**What computer scientists see**

**What physicists see**

interpret neural network input as
probabilities instead of numbers

Probability of class **i** given
input **x** in a neural network

$$P(\text{class } i|\mathbf{x}) = \text{softmax}(z_i) = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

logits (pre-softmax
activations)

Probability of a
certain state

Energy of the state

$$P(\mathbf{s}) = \frac{1}{Z} \exp\left(-\frac{E(\mathbf{s};\boldsymbol{\theta})}{T}\right)$$

Partition function
(normalization)

Temperature
parameter

Optimization and generalization
can make statistical sense.

# Understand prompting
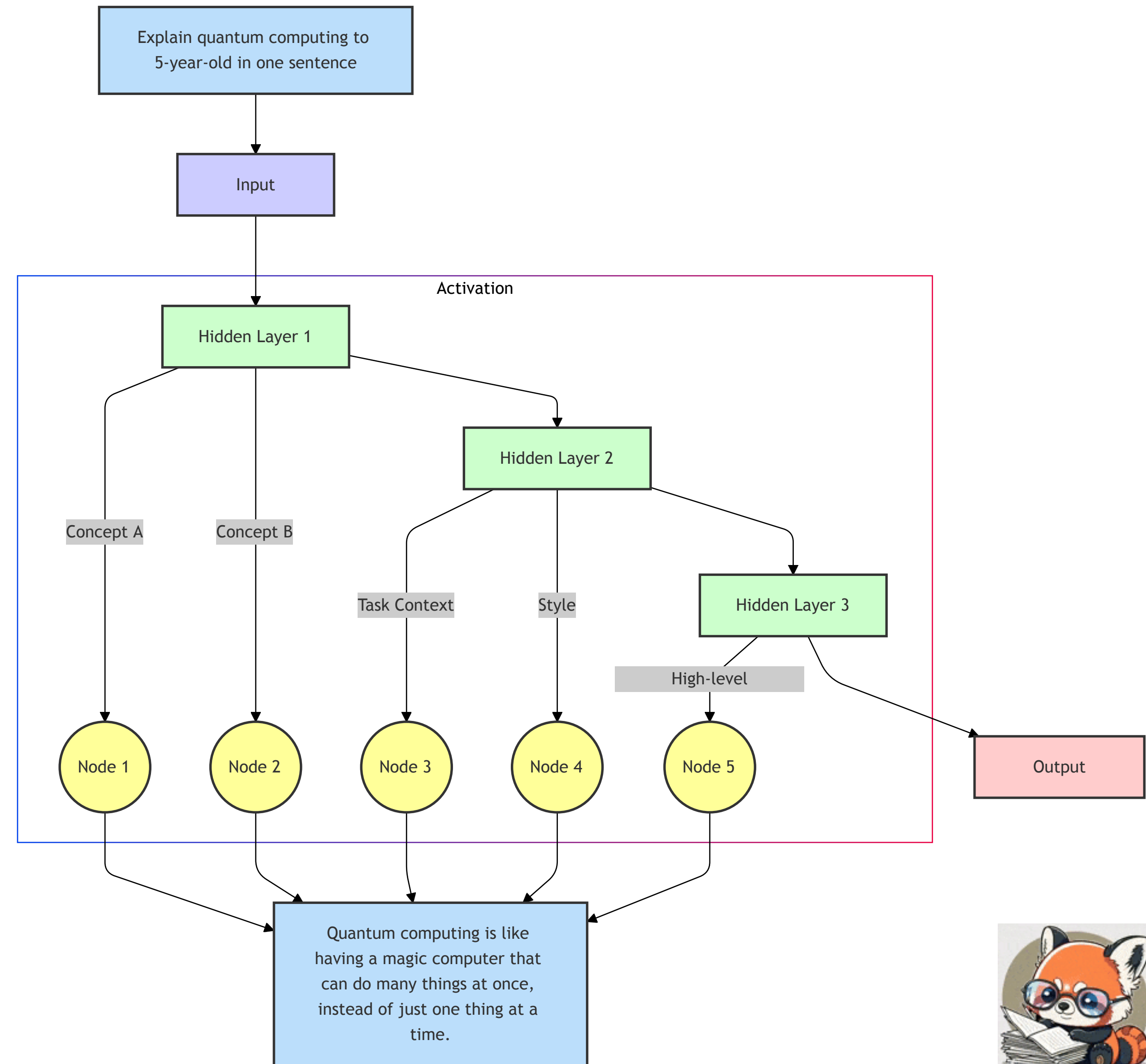
Prompting introduces a task-specific bias to the model's output distribution by its activation patterns, effectively aligning the target task with the LLM's pre-trained task manifold.

Task gradient
$$\nabla_{\text{task}} = \frac{\partial \mathcal{L}_{\text{task}}}{\partial h_L}$$

$$P(y_t|\text{prompt}, x) = \text{softmax}(f(h_L + \beta \nabla_{\text{task}})/\tau)$$

$$h_L = g_L(...g_2(g_1([E_{\text{prompt}}; E_x])))$$

Hidden states

$$\beta = \gamma \cdot \text{sim}(h_L, h_{\text{task}})$$

Alignment strength



https://www.linkedin.com/in/liuhongliang/

# Understand tool using

LLMs translate language tasks into tool actions by computing P(tool|context) through attention-based alignment between task requirements and tool capabilities.

$$P(\text{action}|\text{state}) = \int P(\text{action}|\text{intention})P(\text{intention}|\text{state})dI$$

Connect state observation to action selection

$$P(\text{tool}|\text{context}) = \text{softmax}(\frac{h_{\text{tool}}^{\text{ctx}} W_{\text{out}}}{\tau})$$

Context-aware **tool selection**

Temperature control

$$\alpha_{\text{tool}} = \text{softmax}(\frac{h_{\text{context}} W_Q (h_{\text{tool}} W_K)^T}{\sqrt{d_k}})$$

**Attention** of tools

$$h_{\text{context}} = \text{LayerNorm}([h_{\text{task}}; h_{\text{tools}}])$$

Hidden states of **context**



Input

Available Tools | Task Description

Tool Recognition Layer

Tool Embeddings | Task Embedding

Pattern Matching
P(tool_match|task)

Tool Selection Layer

Tool Scoring
score = f(relevance, confidence)

Tool Validation
constraints & requirements

Final Selection
argmax(valid_tools)

Output

Selected Tool | Tool Parameters

# Understand chain-of-thoughts (CoT)

Chain-of-thought (CoT) emerges from attention mechanism building up a working memory of key-value pairs from each reasoning step, while hidden states evolve by attending to this memory to compute next steps, creating a computational cycle where each step can query and build upon previous computations.

Attention-based Memory Access

$$c_t = \text{softmax}\left(\frac{h_t W^Q (M_t^K)^T}{\sqrt{d}}\right) M_t^V$$

Hidden State Evolution

$$h_t = f(\text{attention}(h_{t-1}, [h_1, ..., h_{t-1}]))$$

$$M_t = [K_t, V_t] = [(h_1 W^K, h_1 W^V), (h_2 W^K, h_2 W^V), ..., (h_t W^K, h_t W^V)]$$

KV Memory Collection (working memory)



Input

Think step by step, solve

Hidden State

Initial Hidden State

Hidden State $h_1$

Hidden State $h_2$

Final Hidden State

Query

Output

Step 1

Step 2

Final

Query

Context

Query

Context

Attention

K/V Memory Store

Context

Query

Context

Attention

Context

# Understand LoRA ranks

$$d_{\text{head}} = k$$

rank

$$B \in \mathbb{R}^{r \times k}$$

$$d_{\text{model}} = d$$

$$A \in \mathbb{R}^{d \times r}$$

rank

$$\Delta W = AB$$

LoRA fine tune weight update

Rank = 4

$$\Delta W_{low} = A_{d \times 4} B_{4 \times k}$$

$$|\Delta W|* = \sum_{i=1}^{4} \sigma_i \leq 4 |\Delta W| F$$

spread limit = 4

simpler, more concentrated updates
Better preservation of base model knowledge
Low risk of overfitting

Nuclear Norm (sum all singular values)

$$|\Delta W|* = \sum_{i=1}^{\min(d,k)} \sigma_i$$

$$\frac{}{|\Delta W| F = \sqrt{\sum_{i=1}^{\min(d,k)} \sigma_i^2}} \leq r$$

spread limit or "rank"
(power of LoRA fine-tune)

Frobenius Norm as
sqrt(sum of squared singular values)

$$\sum_{i=1}^{r} \sigma_i \leq \sqrt{r} \sqrt{\sum_{i=1}^{r} \sigma_i^2}$$

Cauchy-Schwarz Interpretation

Rank = 32

$$\Delta W_{high} = A_{d \times 32} B_{32 \times k}$$

spread limit = 32

$$|\Delta W|* = \sum_{i=1}^{32} \sigma_i \leq 32 |\Delta W| F$$

more complex, distributed changes
More expressive power and complex adaptations
High risk of overfitting

https://www.linkedin.com/in/liuhongliang/

# Understand LLM inference time

**Total time**

**Position embeddings**
$$T_{pos} = \mathcal{O}(s \cdot d)$$

**+**

**Self-attention**
$$T_{attn} = \mathcal{O}(s \cdot d \cdot h)$$

**Feed-forward network**
$$T_{ffn} = \mathcal{O}(s \cdot d \cdot 4d)$$

**Layer norm**
$$T_{ln} = \mathcal{O}(s \cdot d)$$

**Number of layers** ✖
32 layers in Llama 8B

**+**

**Final layer norm**
$$T_{ln} = \mathcal{O}(s \cdot d)$$

**+**

**Output projection**
$$T_{proj} = \mathcal{O}(s \cdot d \cdot v)$$

Example: Llama 8B (FP32) on T4 GPU
512 tokens input 100 tokens output

### Input Processing - 3.5ms

Input

↓

Tokenization -

↓

Token Embedding - 0.5ms

↓

Init KV Cache - 2ms

### First Forward Pass - 151.2ms

Position Embeddings - 0.2ms

↓

32 Layer Stack - 150ms

↓

Final LayerNorm - 0.4ms

↓

Output Projection -

### Token Generation Loop - 2000ms total

Token Check

↓

KV Cache Access - 0.5ms

↓

Token Embedding - 0.2ms

↓

Layer Stack - 18ms

↓

Sampling - 1ms

↓

New Token

$s$ = sequence length (e.g 512)
$d$ = hidden dimension (4096)
$h$ = number of attention heads (32)
$v$ = vocabulary size (e.g 32k)

https://www.linkedin.com/in/liuhongliang/

Total Runtime: 2.15 seconds

# Understand speculative decoding

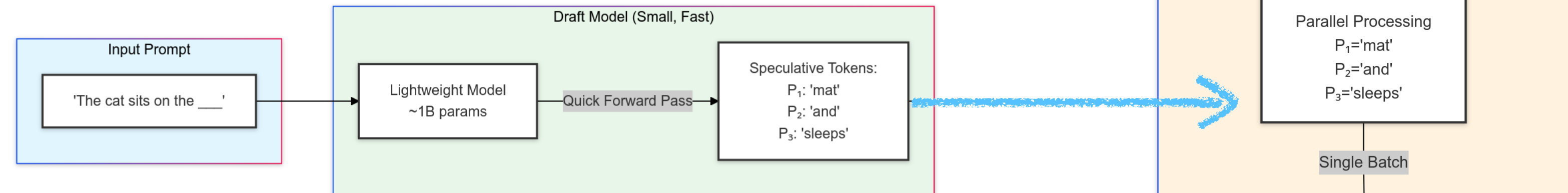Speculative decoding uses a small model to quickly guess multiple next tokens that a large model can verify in parallel, replacing sequential token generation with batch processing when predictions are correct.

## Input Prompt
'The cat sits on the ___'

## Draft Model (Small, Fast)
Lightweight Model ~1B params

Quick Forward Pass

Speculative Tokens:
$P_1$: 'mat'
$P_2$: 'and'
$P_3$: 'sleeps'

## Main Model (Large)
Parallel Processing
$P_1$='mat'
$P_2$='and'
$P_3$='sleeps'

Single Batch

Verify Probabilities

$P_1$ matches

## Token-by-Token Verification

$P_2$ matches

Accept 'mat'

$P_3$ mismatch

Accept 'and'

Reject 'sleeps'

## Final Processing
Regenerate from position 3
'The cat sits on the mat and purrs'

## Sequential Generation vs Parallel Validation

**Sequential Decoding**
Token 1 Forward Pass
Token 2 Forward Pass
Token 3 Forward Pass
Total Time

**Speculative Decoding**
Draft Model (3 tokens)
Parallel Validation
Total Time

**Compute Units**
Sequential GPU Compute
Parallel GPU Compute

0 0 0 0 0 1 1 1 1 1 2 2 2 2 2 3

https://www.linkedin.com/in/liuhongliang/