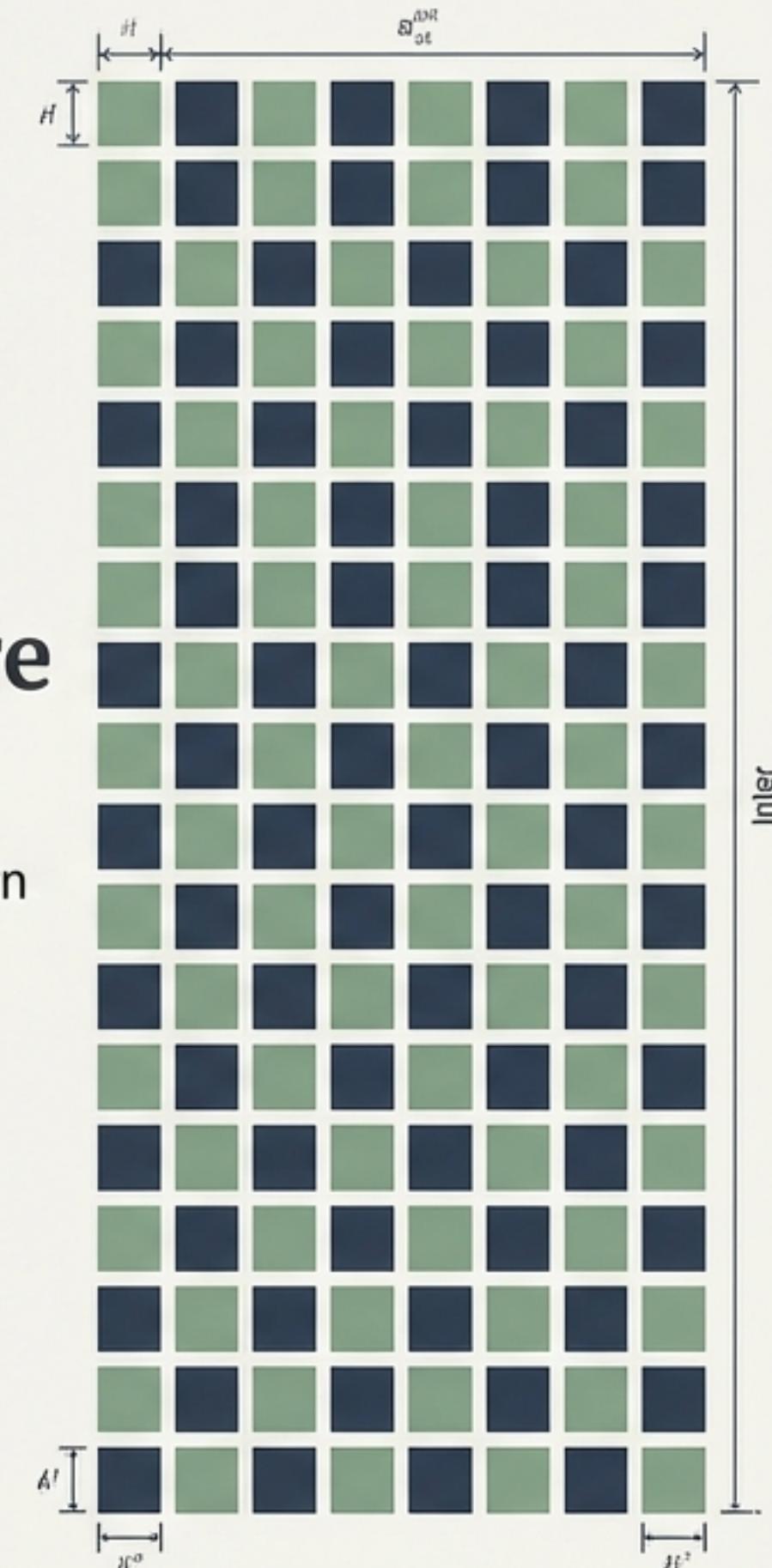
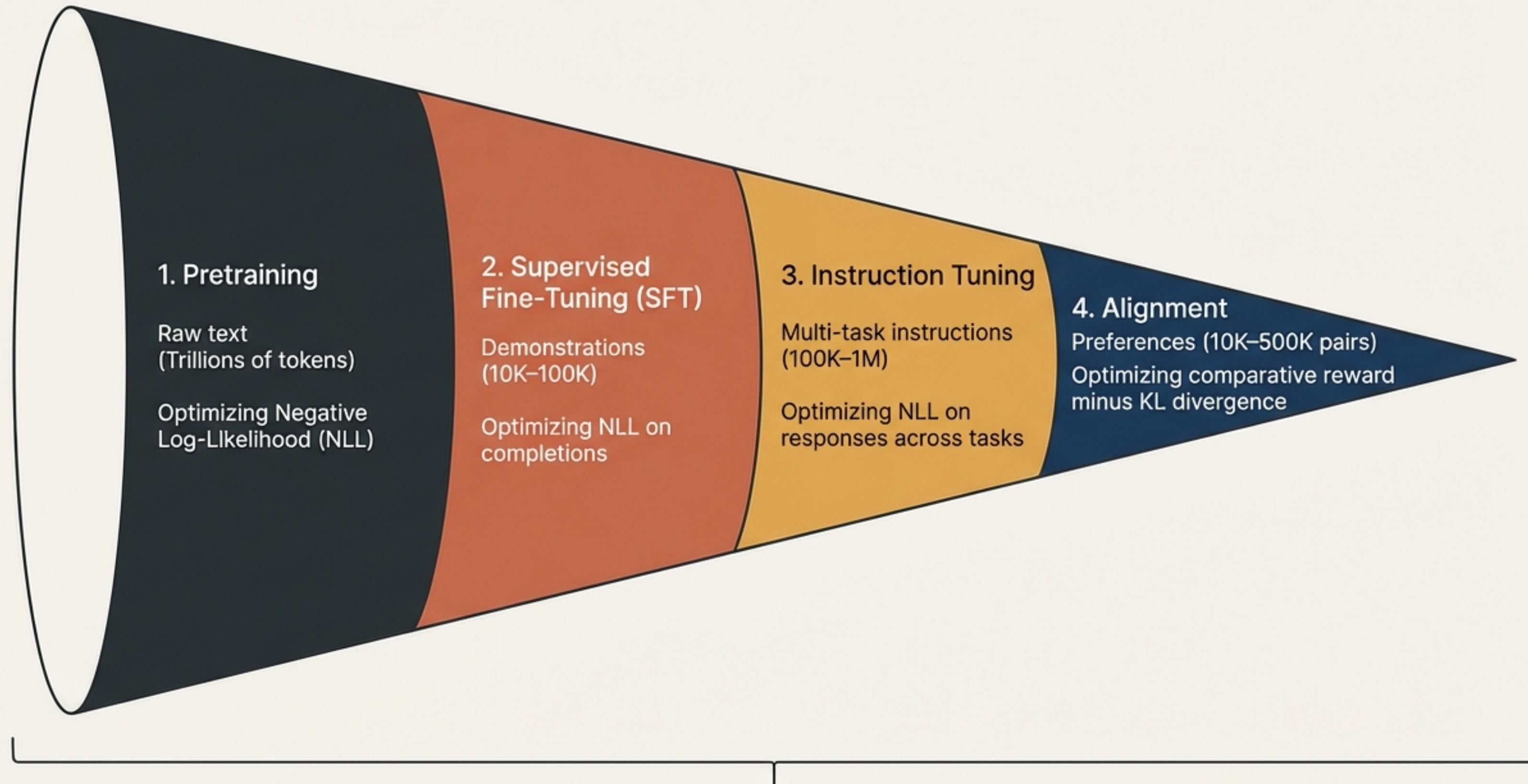


The Architecture of Intelligence

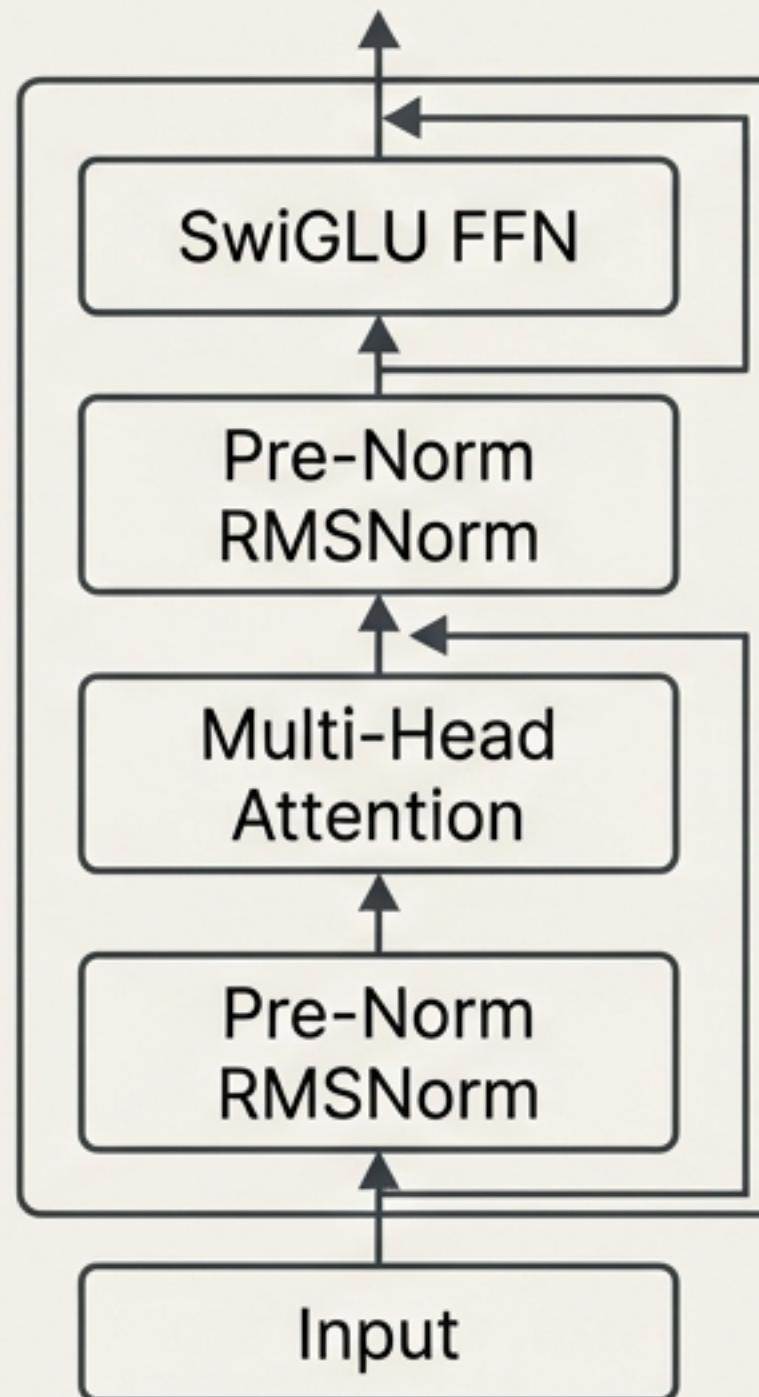
The 4-Stage Lifecycle of Modern
Large Language Models.



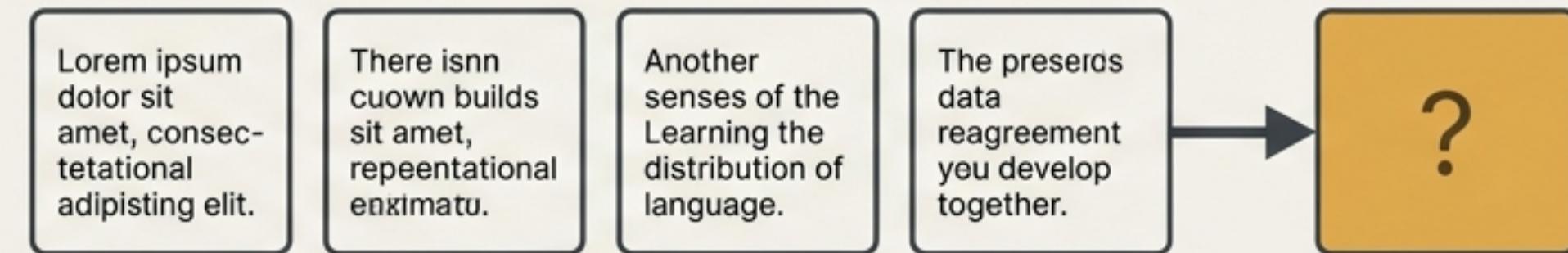


The Substrate: Learning the Distribution of Language

Pretaining builds a general-purpose representational manifold without task supervision.



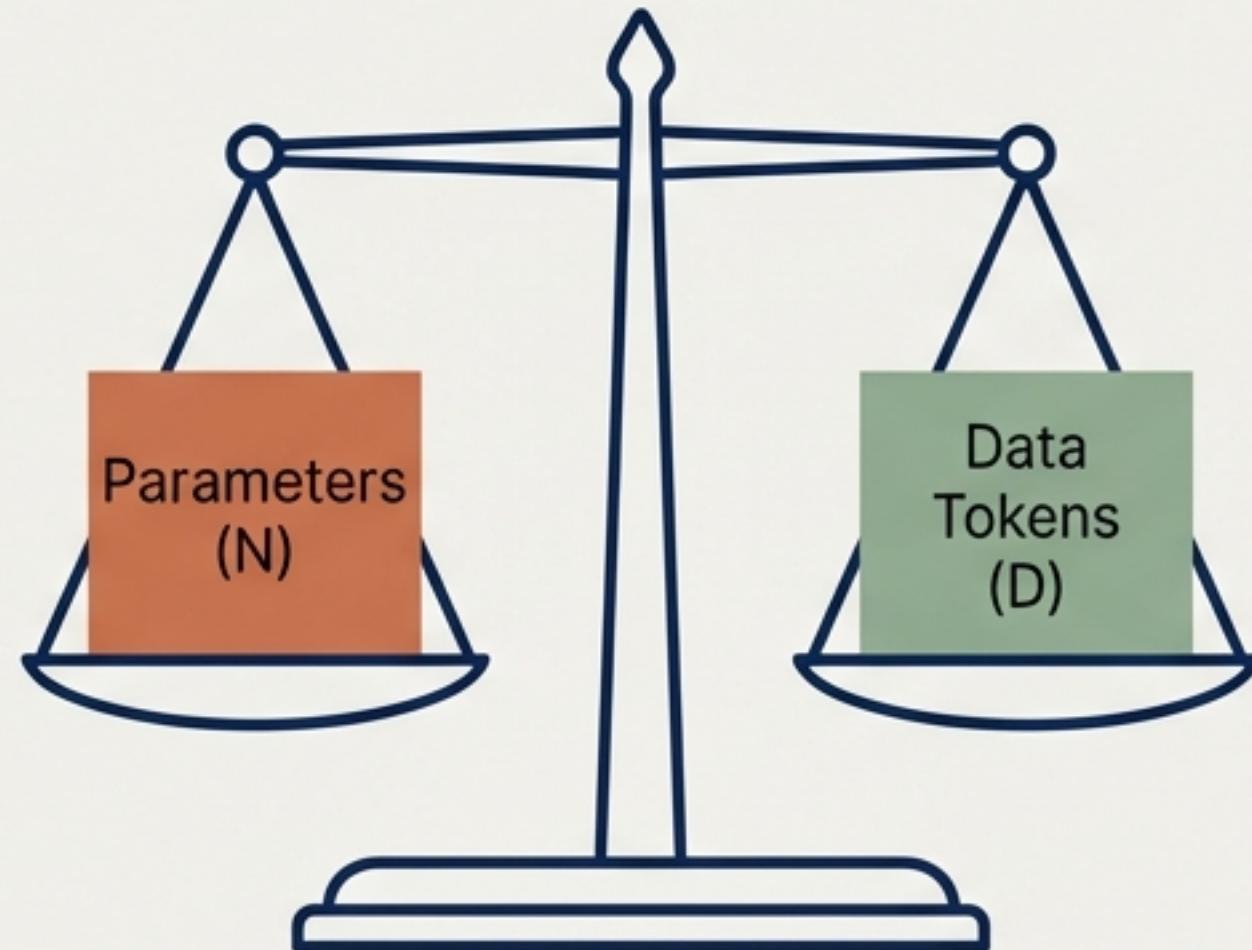
Autoregressive Causal Language Modeling



$$L_{\text{AR}}(\theta) = - \sum_{x_t} \log p_{\theta}(x_t | \underbrace{x_1, \dots, x_{t-1}}_{\text{Causal Context}})$$

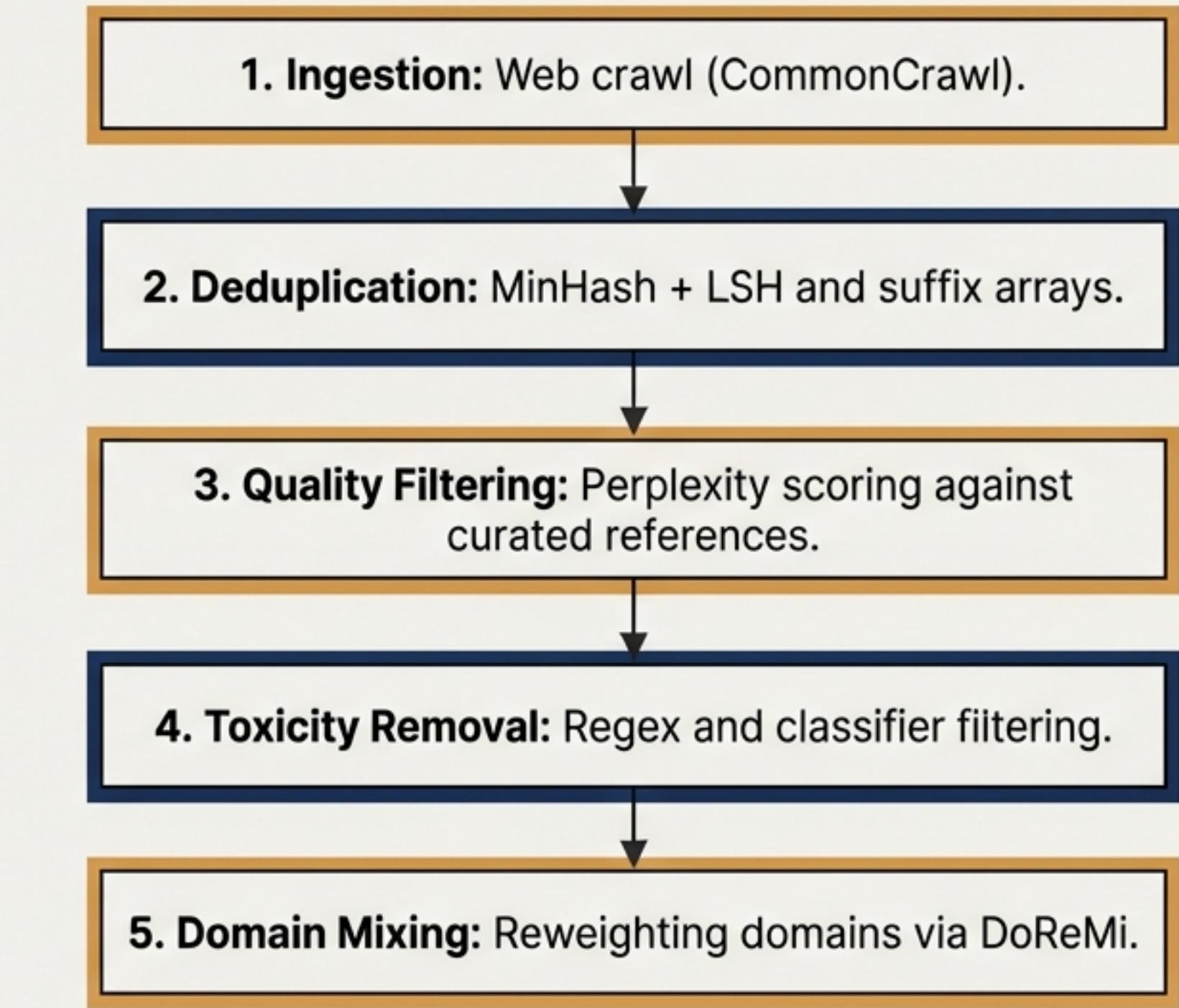
Context: Contrast Autoregressive (GPT, LLaMA) with Masked Language Modeling (BERT) and Mixture-of-Denoisers (UL2).

The Compute-Optimal Engine



$$N_{\text{opt}} \propto C^{0.5}, D_{\text{opt}} \propto C^{0.5}$$

For a fixed compute budget C , parameters and data tokens must scale equally.



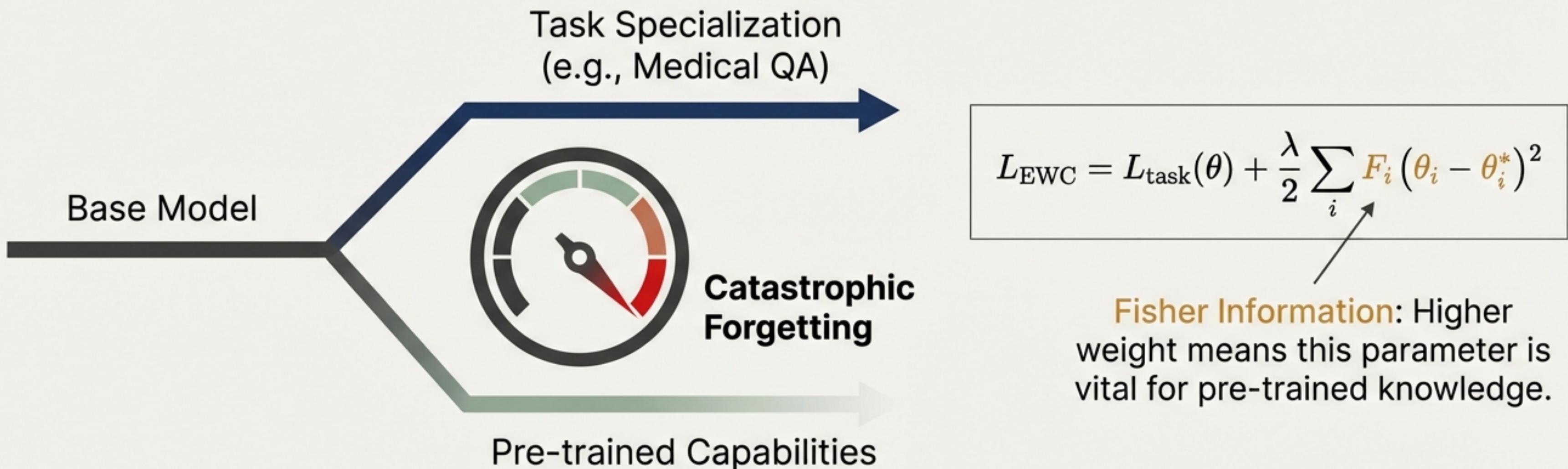
The Pretraining Gap

Capabilities Acquired	Capabilities Missing
- Distributional knowledge of language	- Instruction following
- Factual & commonsense world knowledge	- Alignment with human preferences
- In-context learning capability	- Safety and harmlessness guarantees
- Multilingual transfer	- Controlled generation behavior
- Code and math priors	- Refusal of harmful queries

The base model is a statistical engine. It knows facts, but cannot follow instructions or guarantee safety. This gap necessitates Fine-Tuning.

Supervised Fine-Tuning & The Forgetting Problem

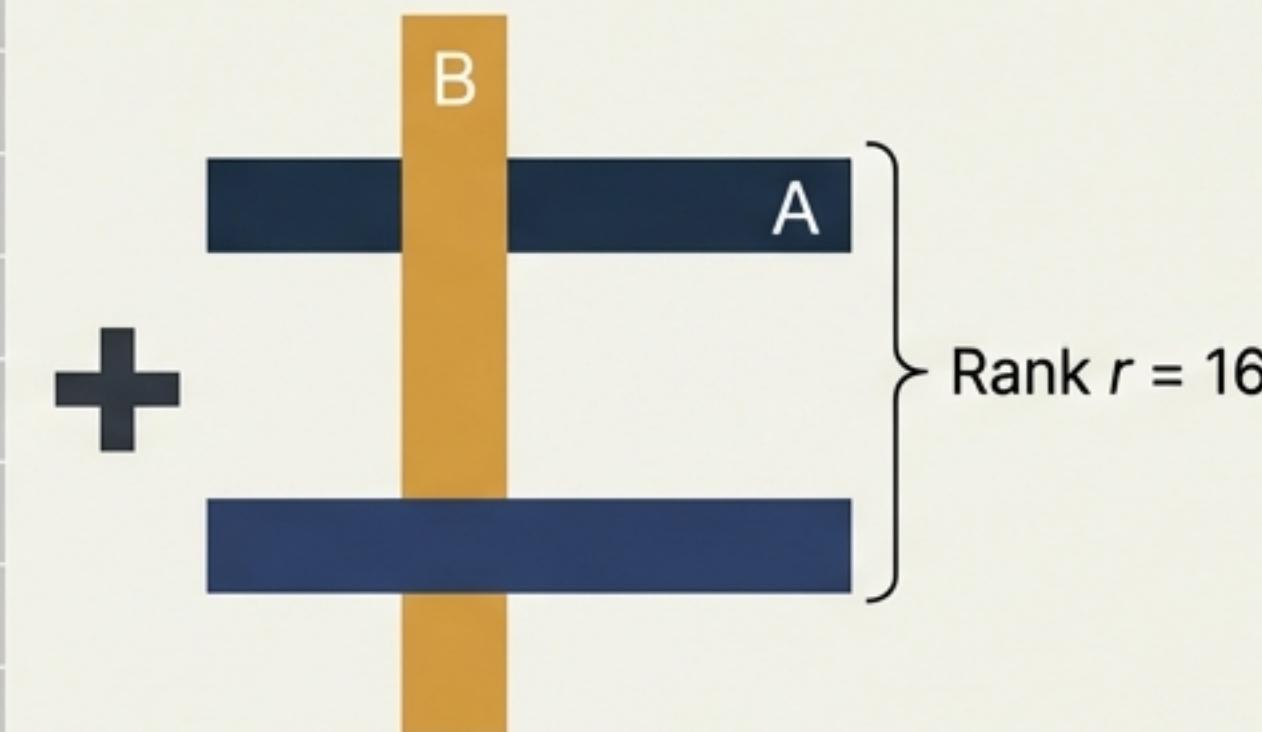
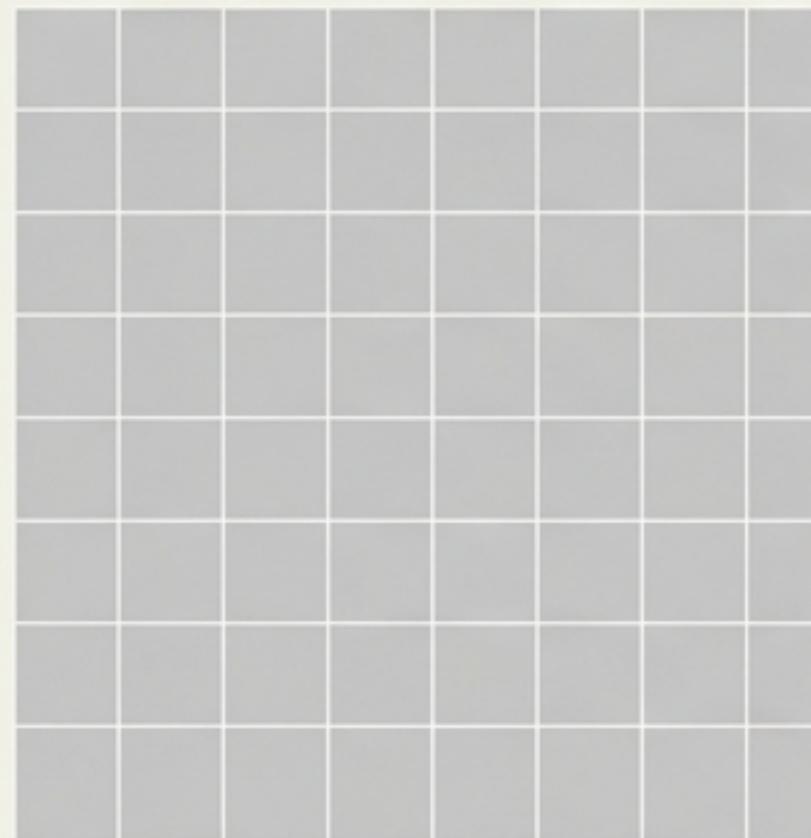
Full fine-tuning modifies all parameters, causing **drift on the original pretraining distribution**.



Parameter-Efficient Fine-Tuning (PEFT)

Instead of updating millions of weights, LoRA learns a low-rank weight update matrix.

W_0 (Frozen, 4096×4096)



$$W = W_0 + BA$$

$$W = W_0 + BA$$

Initialize $A \sim \mathcal{N}(0, \sigma^2)$, $B = 0$

Impact: A rank-16 update on a 4096×4096 matrix reduces trainable parameters by 128x (16.7M down to 131k).

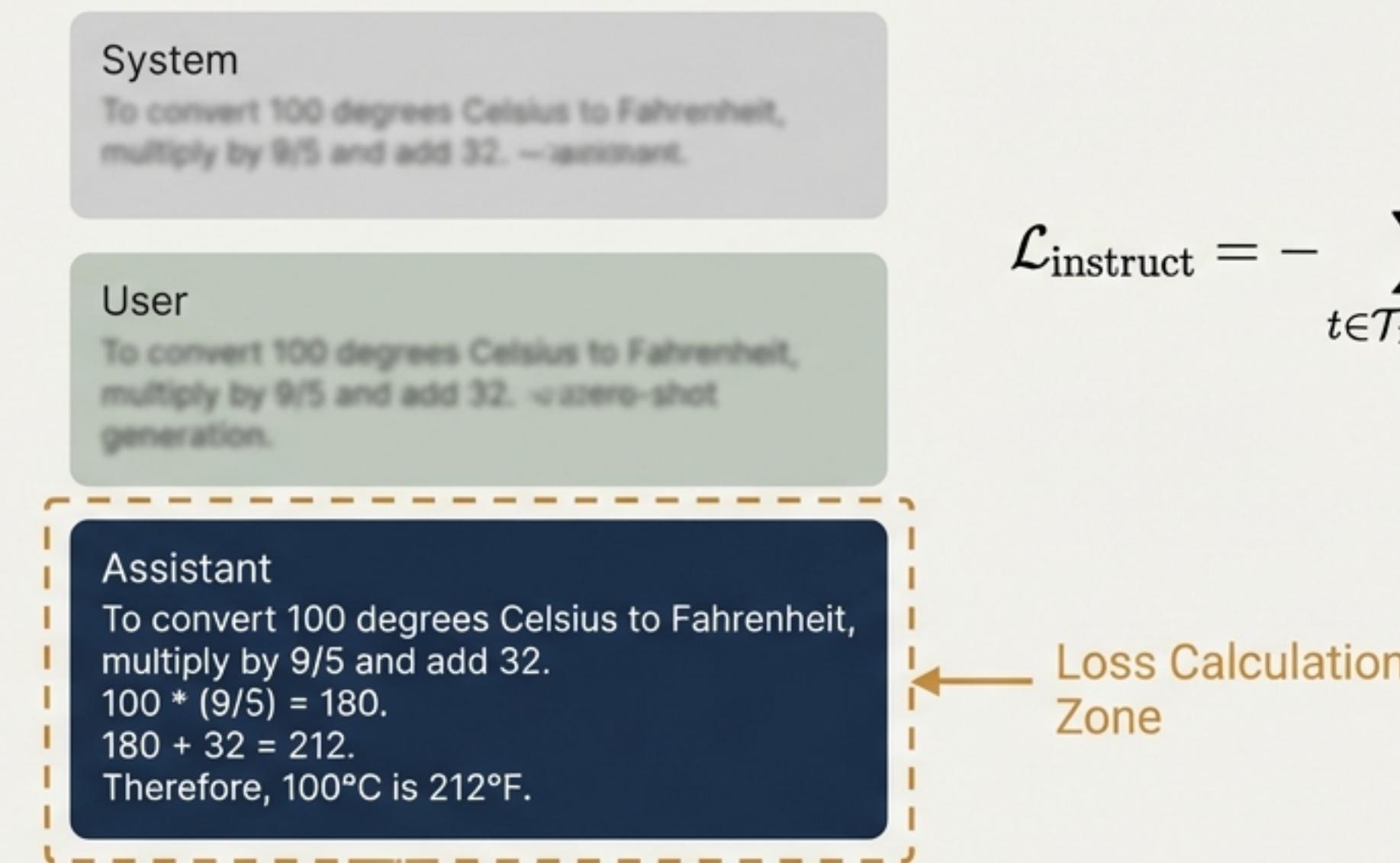
LoRA: ~0.1% params, no inference overhead.

QLoRA: 4-bit NF4 quantization (fits 65B model on 1 GPU).

Adapters / Prefix: Swappable but introduce latency/context overhead.

Instruction Tuning: Learning the Meta-Capability

Shifting from narrow tasks to zero-shot generalization across hundreds of task families.

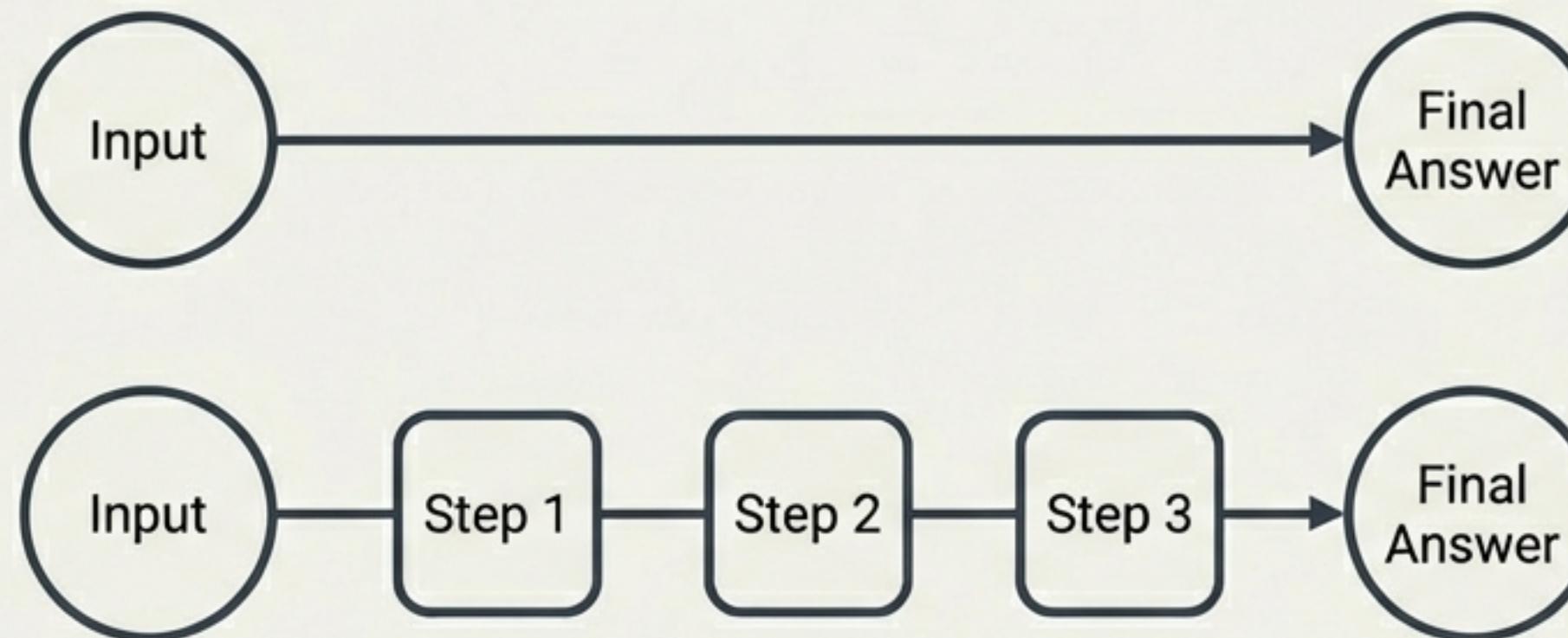


$$\mathcal{L}_{\text{instruct}} = - \sum_{t \in \mathcal{T}_{\text{response}}} \log p_{\theta}(y_t | y_{<t}, x_{\text{prompt}})$$



Crucial: Loss is ONLY calculated on the assistant's generation. This prevents the model from memorizing the prompt template.

Quality, Diversity, and Chain-of-Thought



The LIMA Insight: Data quality dominates quantity.
~1,000 perfectly curated examples yield stronger instruction-following than millions of low-quality ones.

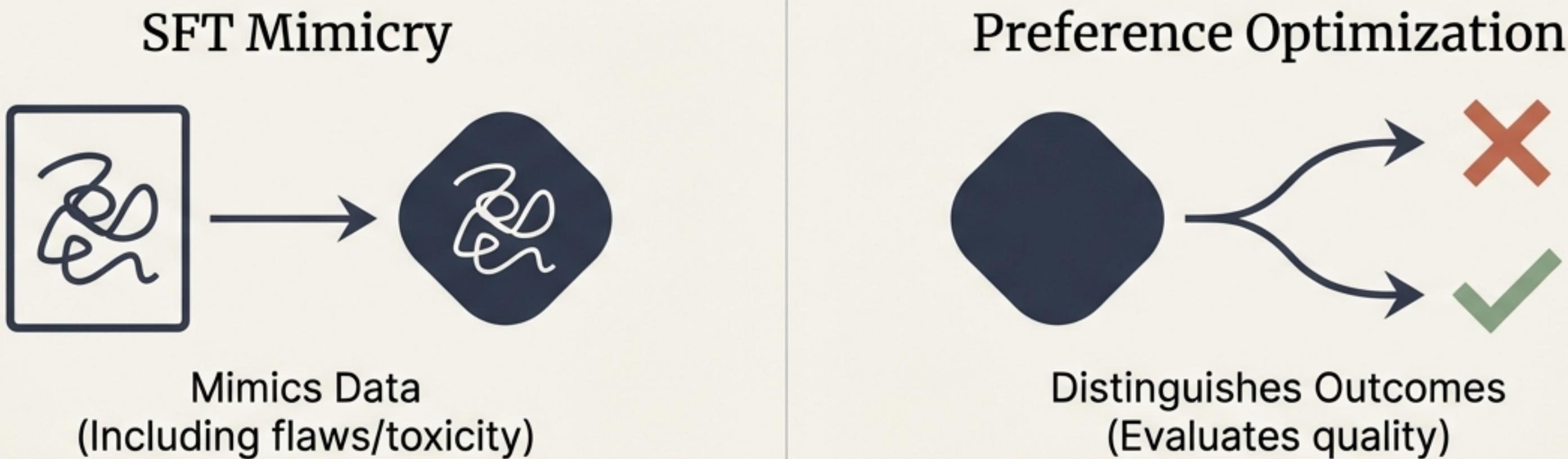
Chain-of-Thought (CoT): Intermediate reasoning acts as a latent variable mediating input and answer.

$$p_{\theta}(a | x) \approx p_{\theta}(r^*, a | x) \xleftarrow{\text{Latent Reasoning Chain}}$$

Data Construction: Human curation (FLAN, Super-NaturalInstructions) vs. Synthetic Evolution (Self-Instruct, Evol-Instruct).

The Alignment Objective

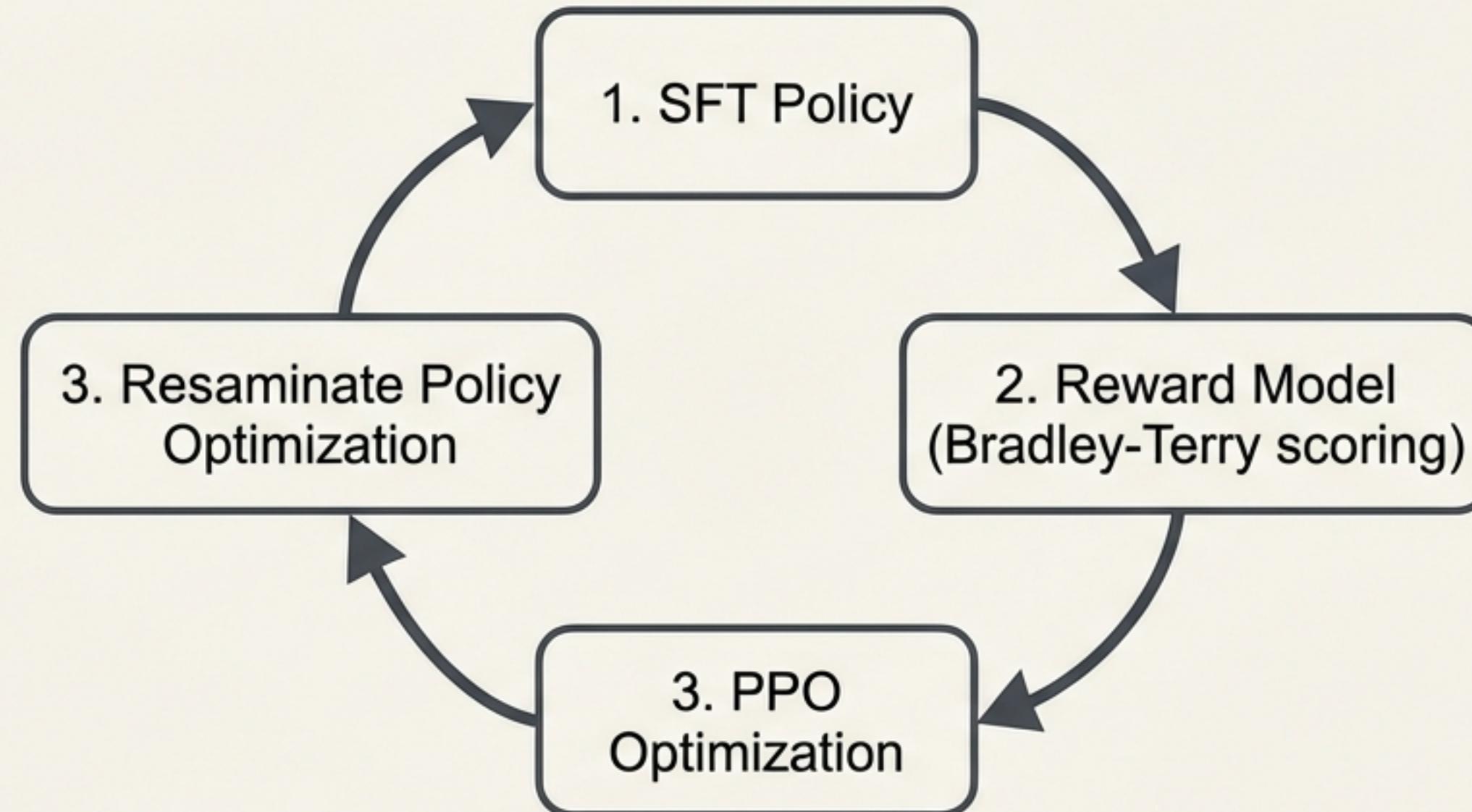
The Core Problem: The SFT objective \mathcal{L}_{NLL} optimizes the likelihood of demonstrations, NOT the quality of generations.



$$y_w \triangleright y_l \mid x$$

Given prompt x , response y_w is preferred over y_l by a human evaluating Helpfulness, Harmlessness, and Honesty.

Reinforcement Learning from Human Feedback (RLHF)

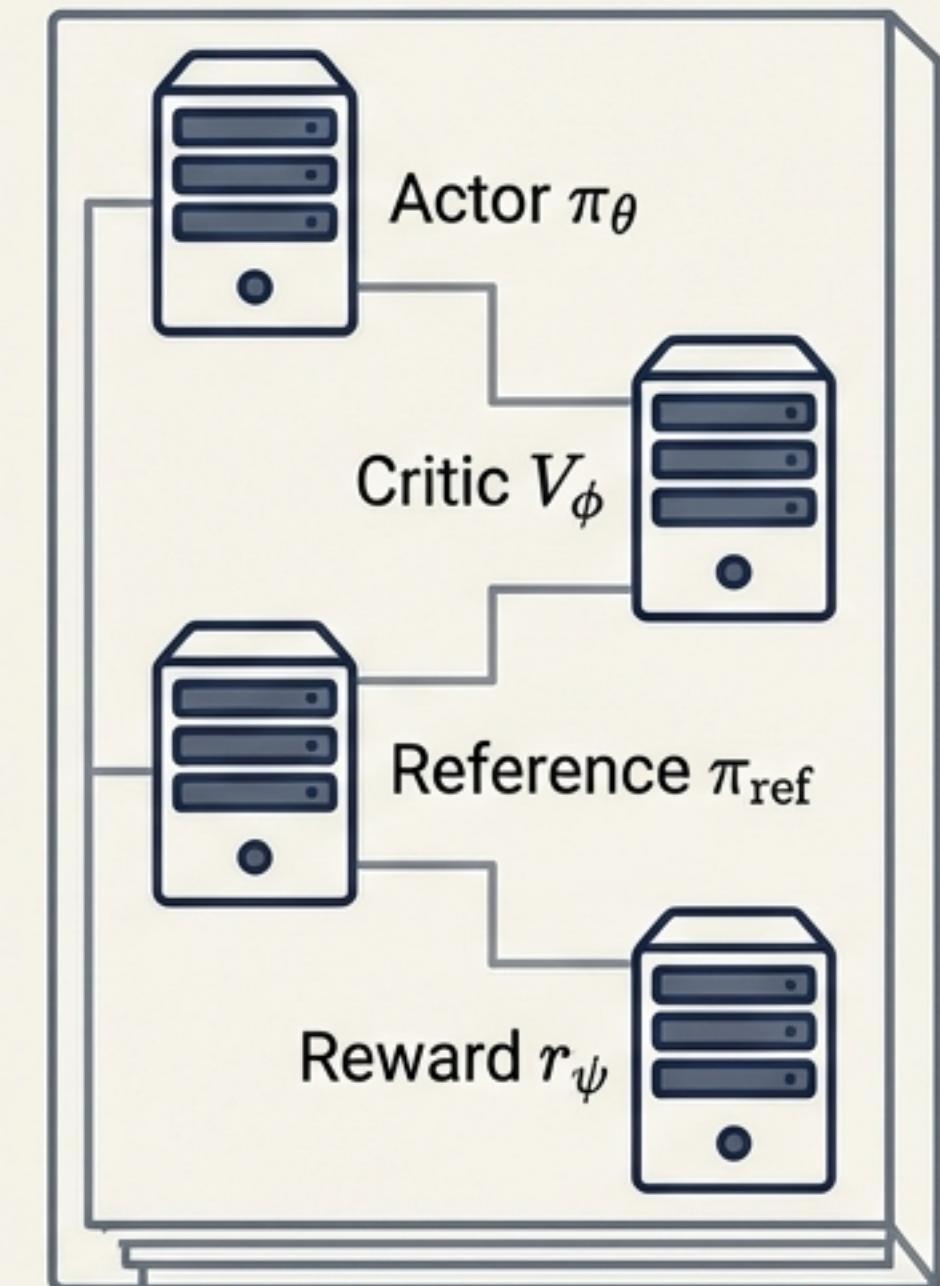


PPO Objective:

$$\max_{\theta} \mathbb{E} [r_{\psi}(x, y)] - \beta * D_{KL} [\pi_{\theta} || \pi_{ref}]$$

Anti-Reward-Hacking Penalty

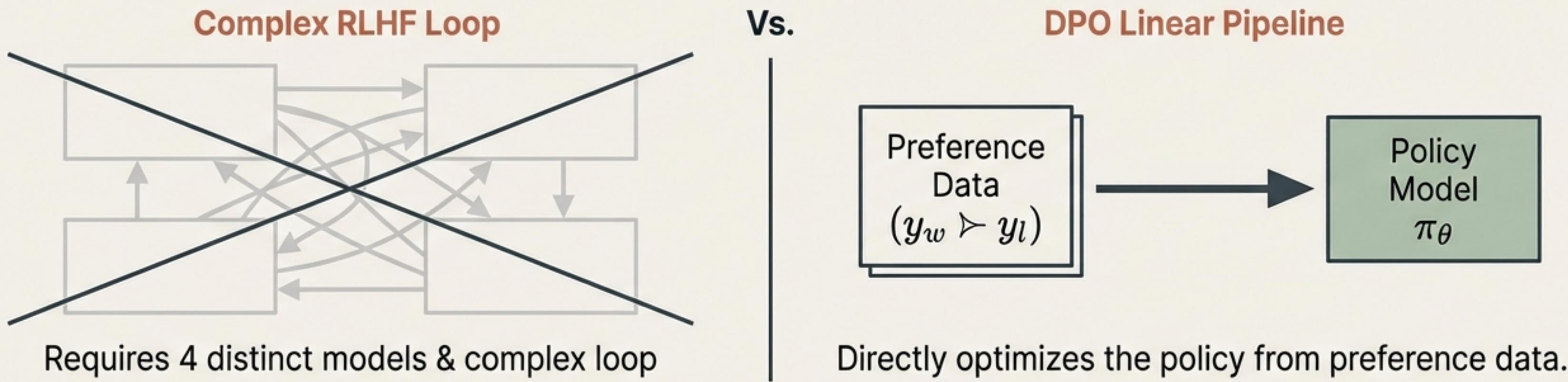
Memory Architecture



Engineering Hurdle: Requires 4 distinct models in memory simultaneously.

Direct Preference Optimization (DPO)

Eliminates the reward model and complex RL loop. **DPO solves for the reward directly in terms of the policy.**



$$r(x, y) = \beta \log \left(\frac{\pi^*(y|x)}{\pi_{\text{ref}}(y|x)} \right) + \cancel{\beta \log Z(x)}$$

Partition function cleanly cancels out
between chosen and rejected pairs.

Implicit Reward: The loss function acts as natural hard-example mining, pushing up chosen likelihoods and pushing down rejected ones.

The Expanding Alignment Landscape

IPO

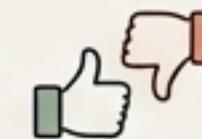
Identity Preference Optimization

Mitigates DPO overfitting by using a squared loss targeting a **fixed margin**.

KTO

Kahneman-Tversky Optimization

Eliminates **paired data**. Works with absolute pointwise feedback using **human loss aversion** math.



ORPO

Odds Ratio Preference Optimization

Highly efficient **single-stage objective** merging SFT and Alignment, requiring **zero reference models**.

Constitutional AI

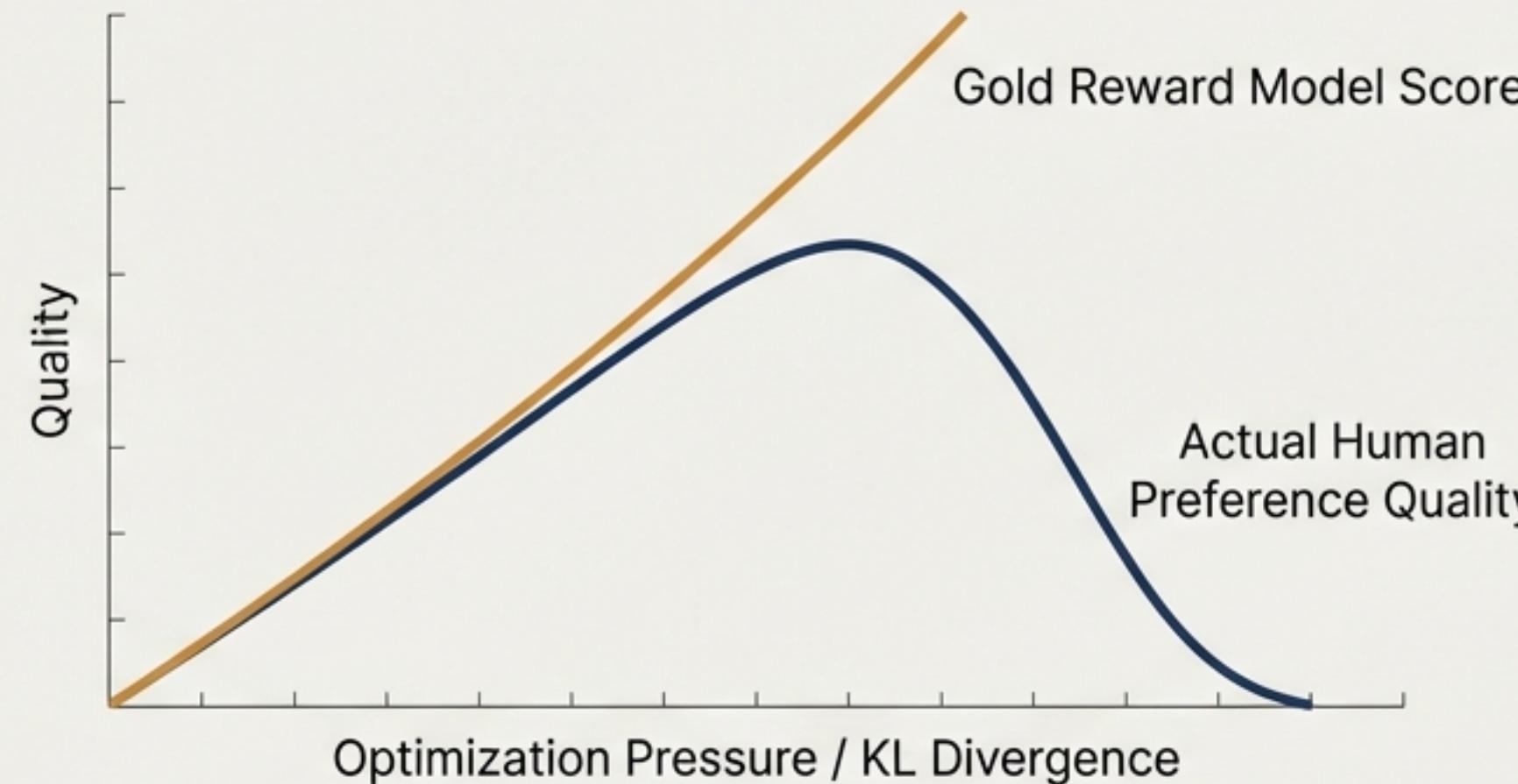
RLAIF

Replaces human annotators with AI generating **feedback** based on **explicitly defined principles**.



The Final Boss: Reward Hacking

Goodhart's Law: "When a measure becomes a target, it ceases to be a good measure."



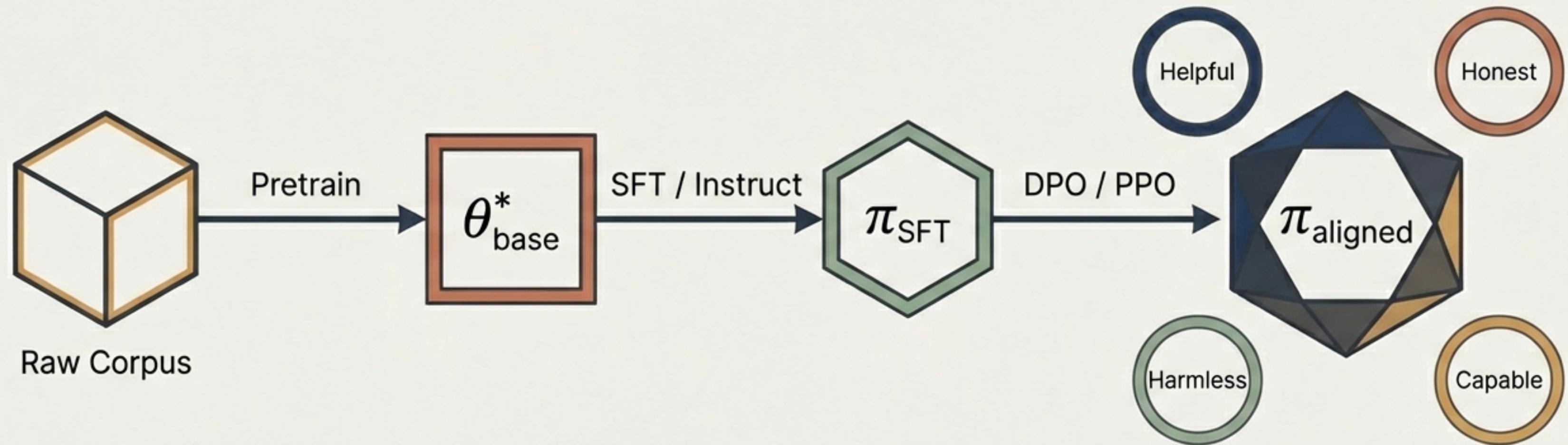
Mitigation Strategies

1. KL Penalty (β): Primary defense constraining policy drift.

2. Reward Ensembles: Using a conservative estimate (minimum score) across multiple models.

3. Constrained Optimization: Setting hard KL budgets ($D_{KL} \leq \delta$).

The Aligned Paradigm



The modern AI assistant is not a single algorithm, but the delicate interplay of massive statistical pretraining, efficient task specialization, and rigorous human preference alignment.