

# DPO, ORPO, GRPO, and PPO - Theory (Part-1)

Part-1 covers theoretical background. Refer part-2 for comparison of techniques and part-3 for the real world scenario.

|   |           |
|---|-----------|
| <b>1 Theoretical Background</b>               | <b>1</b>  |
| Core Concept                                  | 1         |
| 1.1 PPO (Proximal Policy Optimization)        | 2         |
| 1.2 DPO (Direct Preference Optimization)      | 2         |
| 1.3 ORPO (Odds Ratio Preference Optimization) | 2         |
| 1.4 GRPO (Group Relative Policy Optimization) | 2         |
| <b>2 Comparison Tables</b>                    | <b>3</b>  |
| 2.1 Comprehensive Method Comparison           | 3         |
| 2.2 Performance Comparison by Task            | 4         |
| 2.3 Resource Requirements Comparison          | 5         |
| 2.4 Training Characteristics Comparison       | 5         |
| 2.5 Pipeline Components by Method             | 6         |
| <b>3 At a Glance</b>                          | <b>7</b>  |
| 3.1 PPO (Proximal Policy Optimization)        | 8         |
| 3.2 DPO (Direct Preference Optimization)      | 8         |
| 3.3 ORPO (Odds Ratio Preference Optimization) | 9         |
| 3.4 GRPO (Group Relative Policy Optimization) | 10        |
| 3.5 Method Comparison Matrix                  | 11        |
| 3.6 Mathematical Relationships                | 11        |
| 3.7 Optimization Targets                      | 11        |
| <b>4 Emerging Trends (2025-2026)</b>          | <b>12</b> |

## 1 Theoretical Background

### Core Concept

Preference optimization methods align language models with human preferences by training models to prefer better responses over worse ones. These techniques emerged to simplify and improve upon traditional Reinforcement Learning from Human Feedback (RLHF).

## 1.1 PPO (Proximal Policy Optimization)

PPO is the **original RL algorithm used for RLHF** and the foundation of LLM alignment. It works by:

- Using a **reward model** trained on human preference data to score responses
- Optimizing the policy through **constrained policy updates** using a clipped surrogate objective
- Employing a **value function (critic network)** to estimate advantages and reduce variance
- Using **KL-divergence penalties** relative to the reference policy to prevent drastic changes

The key innovation: PPO constrains policy updates to a "trust region" through clipping, ensuring training stability while maintaining sample efficiency.

## 1.2 DPO (Direct Preference Optimization)

DPO reformulates the RLHF objective to eliminate the need for explicit reward models and complex RL training. It works by:

- Reparameterizing the RLHF objective to create an **implicit reward function** derived from the policy itself and a reference model
- Framing preference learning as a **classification problem** between chosen and rejected responses
- Using a static preference dataset with paired responses (chosen vs rejected)

The key insight: DPO indirectly solves the RLHF objective by optimizing this implicit reward, avoiding the instability of RL training.

## 1.3 ORPO (Odds Ratio Preference Optimization)

ORPO combines SFT and preference optimization into a **single monolithic training stage**. Its innovations include:

- Eliminating the reference model entirely, reducing memory by 20-30%
- Using the **odds ratio** to contrast preferred vs dispreferred responses during training
- Penalizing disfavored responses during SFT rather than treating all examples equally

ORPO measures the **absolute odds ratio within the current policy**, whereas DPO measures preference change relative to a frozen reference model.

## 1.4 GRPO (Group Relative Policy Optimization)

GRPO represents a paradigm shift introduced by DeepSeek in January 2025. It eliminates both reward models and critic networks by using:

- **Group-relative scoring:** For each prompt, generate 8-16 responses and score them against verifiable rewards (e.g., code correctness, math accuracy)
- **Advantage estimation:** Each response's advantage is computed relative to the group mean
- **Policy updates:** Increase probability of above-average responses, decrease probability of below-average ones

Recent research shows GRPO and DPO optimize fundamentally similar contrastive objectives—the difference lies in online vs offline training regimes.

## 2 Comparison Tables

### 2.1 Comprehensive Method Comparison

| Criterion                  | PPO   | DPO                  | ORPO               | GRPO                         |
|----------------------------|---|----------------------|--------------------|------------------------------|
| <b>Reference Model</b>     | Required  | Required             | Not required       | Not required                 |
| <b>Reward Model</b>        | Required (separate training)                    | Not required         | Not required       | Optional                     |
| <b>Critic Network</b>      | Required (value function)                       | Not required         | Not required       | Not required                 |
| <b>Training Stages</b>     | 3 stages (SFT → RM → PPO)                       | 2 stages (SFT → DPO) | 1 stage            | 2 stages (SFT → GRPO)        |
| <b>Data Requirements</b>   | Preference pairs for RM, then online generation | Paired preferences   | Paired preferences | Prompts + verifiable rewards |
| <b>Compute Cost</b>        | Highest (4 models in memory)                    | Medium               | Low                | High                         |
| <b>Memory Requirements</b> | 3-4× base model                                 | 2× base model        | 1× base model      | 1-2× base model              |

|                                     |                                    |                              |                       |                                    |
|-------------------------------------|------------------------------------|------------------------------|-----------------------|------------------------------------|
| <b>Learning Type</b>                | Online RL                          | Offline contrastive          | Offline hybrid        | Online RL                          |
| <b>Implementation Complexity</b>    | Very high (37+ details)            | Low                          | Low                   | Medium                             |
| <b>Training Stability</b>           | Moderate (prone to instability)    | Good                         | Very good             | Excellent                          |
| <b>Hyperparameter Sensitivity</b>   | Very high                          | Low                          | Low                   | Medium                             |
| <b>Sample Efficiency</b>            | Low (on-policy)                    | High (offline)               | High                  | Moderate                           |
| <b>Best Performance</b>             | Complex reasoning, code generation | Safety, style, dialogue      | Resource-constrained  | Math, code with verifiable outputs |
| <b>Catastrophic Forgetting Risk</b> | High                               | Low                          | Low                   | Low                                |
| <b>Fine-grained Control</b>         | Excellent                          | Limited                      | Limited               | Good                               |
| <b>Theoretical Foundation</b>       | Trust region optimization          | Implicit reward maximization | Odds ratio during SFT | Group-relative advantage           |

## 2.2 Performance Comparison by Task

| Task Type              | PPO                | DPO                | ORPO      | GRPO               |
|------------------------|--------------------|--------------------|-----------|--------------------|
| Code Generation        | ★★★★★<br>Excellent | ★★★★ Good          | ★★★ Fair  | ★★★★★<br>Very Good |
| Mathematical Reasoning | ★★★★★<br>Excellent | ★★★★ Good          | ★★★ Fair  | ★★★★★<br>Excellent |
| Dialogue/Chat          | ★★★★★<br>Very Good | ★★★★★<br>Excellent | ★★★★ Good | ★★★★ Good          |
| Safety Alignment       | ★★★★★<br>Very Good | ★★★★★<br>Excellent | ★★★★ Good | ★★★★ Good          |

|                    |                     |                     |                     |                     |
|--------------------|---------------------|---------------------|---------------------|---------------------|
| Summarization      | ★★★★★<br>Very Good  | ★★★★★★<br>Excellent | ★★★★★<br>Very Good  | ★★★★★ Good          |
| Small Models (<7B) | ★★★★★ Good          | ★★★★★<br>Very Good  | ★★★★★★<br>Excellent | ★★★ Fair            |
| Large Models (7B+) | ★★★★★★<br>Excellent | ★★★★★<br>Very Good  | ★★★★★ Good          | ★★★★★★<br>Excellent |

### 2.3 Resource Requirements Comparison

| Resource                     | PPO                            | DPO                 | ORPO             | GRPO                         |
|------------------------------|--------------------------------|---------------------|------------------|------------------------------|
| <b>GPU Memory</b>            | 3-4× base model                | 2× base model       | 1× base model    | 1-2× base model              |
| <b>Training Time</b>         | Longest                        | Medium              | Shortest         | Long                         |
| <b>Implementation Effort</b> | Very High (37+ details)        | Low                 | Very Low         | Medium                       |
| <b>Infrastructure Needs</b>  | Complex (4 models)             | Moderate (2 models) | Simple (1 model) | Moderate (vLLM server)       |
| <b>Hyperparameter Tuning</b> | Extensive                      | Minimal             | Minimal          | Moderate                     |
| <b>Data Requirements</b>     | Preference pairs + RM training | Preference pairs    | Preference pairs | Prompts + verifiable rewards |
| <b>Sample Efficiency</b>     | Low (on-policy)                | High (offline)      | High (offline)   | Moderate (online)            |

### 2.4 Training Characteristics Comparison

| Characteristic            | PPO                   | DPO    | ORPO         | GRPO         |
|---------------------------|-----------------------|--------|--------------|--------------|
| <b>Training Stability</b> | ⚠ Moderate (unstable) | ✓ Good | ✓✓ Very Good | ✓✓ Excellent |

|                                |                    |                        |                    |                           |
|--------------------------------|--------------------|------------------------|--------------------|---------------------------|
| <b>Convergence Speed</b>       | Slow               | Fast                   | Fastest            | Moderate                  |
| <b>Catastrophic Forgetting</b> | ⚠ High Risk        | ✓ Low Risk             | ✓ Low Risk         | ✓ Low Risk                |
| <b>Reward Signal Quality</b>   | Flexible (learned) | Fixed (implicit)       | Fixed (odds ratio) | Verifiable (ground truth) |
| <b>Online Learning</b>         | ✓ Yes              | ✗ No                   | ✗ No               | ✓ Yes                     |
| <b>Exploration Capability</b>  | ✓ High             | ✗ None                 | ✗ None             | ✓ Moderate                |
| <b>KL Divergence Control</b>   | Explicit penalty   | Implicit via reference | None needed        | Implicit                  |

## 2.5 Pipeline Components by Method

| Component          | PPO         | DPO | ORPO | GRPO         |
|--------------------|-------------|-----|------|--------------|
| Policy Model       | ✓           | ✓   | ✓    | ✓            |
| Reference Model    | ✓           | ✓   | ✗    | ✗            |
| Reward Model       | ✓ (trained) | ✗   | ✗    | ✓ (optional) |
| Value/Critic Model | ✓           | ✗   | ✗    | ✗            |

# 3 At a Glance

| PPO (Proximal Policy Optimization)  | DPO (Direct Preference Optimization)  | ORPO (Odds Ratio Preference Optimization)   | GRPO (Group Relative Policy Optimization)   |
|---|---|---|---|
| <p><b>SETUP</b></p>   | <p><b>SETUP</b></p>   | <p><b>SETUP</b></p>   | <p><b>SETUP</b></p>   |
| <p><b>OBJECTIVE</b></p> <p>Maximize reward while keeping policy close to reference policy <math>\pi_{ref}</math></p> $\max_{\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}} [r_{\theta}(x, y)]$ <p>s.t. <math>KL(\pi_{\theta} \parallel \pi_{ref}) \leq \epsilon</math></p>  | <p><b>OBJECTIVE</b></p> <p>Maximize log-prob of chosen over rejected relative to reference policy <math>\pi_{ref}</math></p> $\max_{\theta} \mathbb{E}_{(x, y^+, y^-) \sim \mathcal{D}} [\log \sigma(\beta \Delta_{\theta}(x, y^+, y^-))]$  | <p><b>OBJECTIVE</b></p> <p>Optimize odds ratio with explicit normalization (improves stability)</p> $\max_{\theta} \mathbb{E}_{(x, y^+, y^-)} \left[ \log \frac{\pi_{\theta}(y^+   x)}{\pi_{\theta}(y^-   x)} - \log \frac{\pi_{ref}(y^+   x)}{\pi_{ref}(y^-   x)} \right]$ | <p><b>OBJECTIVE</b></p> <p>Maximize group-relative advantage (no value function required)</p> $\max_{\theta} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y_{1:G} \sim \pi_{\theta}} \left[ \frac{1}{G} \sum_{i=1}^G \hat{A}_i \log \pi_{\theta}(y_i   x) \right]$ |
| <p><b>PPO CLIPPED OBJECTIVE</b></p> $L_{PPO}(\theta) = \mathbb{E} \left[ \min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) - \beta KL(\pi_{\theta} \parallel \pi_{ref}) \right]$ <p>where <math>r_t(\theta) = \frac{\pi_{\theta}(y_t   x_t)}{\pi_{old}(y_t   x_t)}</math></p>   | <p><b>DPO LOSS</b></p> $L_{DPO}(\theta) = -\mathbb{E}_{(x, y^+, y^-)} [\log \sigma(\beta \Delta_{\theta}(x, y^+, y^-))]$ <p>where</p> $\Delta_{\theta}(x, y^+, y^-) = \log \frac{\pi_{\theta}(y^+   x)}{\pi_{ref}(y^+   x)} - \log \frac{\pi_{\theta}(y^-   x)}{\pi_{ref}(y^-   x)}$ $\sigma(z) = \frac{1}{1 + e^{-z}}$ | <p><b>ORPO LOSS</b></p> $L_{ORPO}(\theta) = -\mathbb{E}_{(x, y^+, y^-)} [\log \sigma(\beta \Gamma_{\theta}(x, y^+, y^-))]$ <p>where</p> $\Gamma_{\theta}(x, y^+, y^-) = \log \frac{\pi_{\theta}(y^+   x) \pi_{ref}(y^-   x)}{\pi_{\theta}(y^-   x) \pi_{ref}(y^+   x)}$     | <p><b>GROUP ADVANTAGE</b></p> $A_i = r_i - \frac{1}{G} \sum_{j=1}^G r_j$ $\hat{A}_i = \frac{A_i - \mu_A}{\sigma_A + \epsilon}$ $\mu_A = \frac{1}{G} \sum_{j=1}^G A_j, \quad \sigma_A^2 = \frac{1}{G} \sum_{j=1}^G (A_j - \mu_A)^2$                                |
| <p><b>ADVANTAGE (GAE)</b></p> $\hat{A}_t = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}$ <p>e.g., <math>\delta_t = r_t + \gamma V(x_{t+1}) - V(x_t)</math></p>   | <p><b>OPTIONAL REGULARIZATION</b></p> $L = L_{DPO} + \lambda KL(\pi_{\theta} \parallel \pi_{ref})$  | <p><b>OPTIONAL REGULARIZATION</b></p> $L = L_{ORPO} + \lambda KL(\pi_{\theta} \parallel \pi_{ref})$   | <p><b>GRPO LOSS</b></p> $L_{GRPO}(\theta) = -\mathbb{E}_x \mathbb{E}_{y_{1:G}} \left[ \frac{1}{G} \sum_{i=1}^G \hat{A}_i \log \pi_{\theta}(y_i   x) \right] - \beta KL(\pi_{\theta} \parallel \pi_{ref})$   |
| <p><b>VALUE LOSS</b></p> $L_V(\theta) = \mathbb{E} [(V_{\theta}(x_t) - R_t)^2]$   |   |   |   |
| <p><b>TOTAL LOSS</b></p> $L_{PPO total} = -L_{PPO} + c_1 L_V - c_2 H(\pi_{\theta})$   |   |   |   |
| <p><b>NOTES / KEY POINTS</b></p> <ul style="list-style-type: none"> <li>On-policy</li> <li>Uses reward model for scalar reward</li> <li>Stable with clipping + KL penalty</li> <li>Requires value function (critic)</li> </ul>  | <p><b>NOTES / KEY POINTS</b></p> <ul style="list-style-type: none"> <li>Supervised on preference pairs</li> <li>No reward model or critic needed</li> <li>Equivalent to Bradley-Terry model</li> <li>Simple, efficient, stable</li> </ul>   | <p><b>NOTES / KEY POINTS</b></p> <ul style="list-style-type: none"> <li>Odds ratio formulation (symmetric)</li> <li>Better calibrated gradients than DPO</li> <li>No reward model or critic</li> <li>Simple and efficient</li> </ul>  | <p><b>NOTES / KEY POINTS</b></p> <ul style="list-style-type: none"> <li>Sample G responses per prompt</li> <li>Uses group-relative advantage</li> <li>No critic / value model</li> <li>Efficient and scalable</li> </ul>  |
| <p><b>AT A GLANCE</b></p> <p>On-policy RL<br/>Reward model + Critic<br/>Clip + KL for stability</p>   | <p>Offline preference learning<br/>Pairwise (chosen vs rejected)<br/>Simple logistic objective</p>  | <p>Odds ratio between policies<br/>Improved gradient calibration<br/>No critic, no reward model</p>   | <p>Group sampling + relative advantage<br/>No critic, no reward model<br/>Scalable and efficient</p>  |
| <p><b>Symbols:</b> <math>\pi_{\theta}</math> : current policy   <math>\pi_{ref}</math> : reference policy   <math>r_{\theta}</math> : reward model   <math>x</math> : prompt   <math>y</math> : response   <math>y^+, y^-</math> : chosen / rejected   <math>r</math> : reward   <math>KL</math> : Kullback-Leibler divergence</p> <p><math>\sigma(z) = \frac{1}{1 + e^{-z}}</math> : sigmoid   <math>\beta</math> : inverse temperature / regularization strength   <math>\epsilon</math> : clip range   <math>G</math> : group size   <math>\mathbb{E}</math> : expectation</p> |   |   |   |

### 3.1 PPO (Proximal Policy Optimization)

**Pipeline:** Base Model → SFT → Reward Model Training → PPO Optimization

**Key Equation:**

$$L^{PPO}(\theta) = \mathbb{E}_t \left[ \min \left( r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right] - \beta D_{KL}(\pi_\theta || \pi_{ref})$$

Where:  $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{old}(a_t|s_t)}$

**Components:**

- Policy:  $\pi_\theta$
- Reference Policy:  $\pi_{ref}$
- Reward Model:  $r_\phi(x, y)$
- Value Function:  $V_\psi(x)$
- Advantage:  $\hat{A}_t = r(x, y) + V(x') - V(x)$

**Objective:**

$$\max_{\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(y|x)} [r_\phi(x, y)] - \beta D_{KL}(\pi_\theta || \pi_{ref})$$

### 3.2 DPO (Direct Preference Optimization)

**Pipeline:** Base Model → SFT → DPO

**Key Equation:**

$$L_{DPO}(\theta) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w|x)}{\pi_{ref}(y_w|x)} - \beta \log \frac{\pi_\theta(y_l|x)}{\pi_{ref}(y_l|x)} \right) \right]$$

**Implicit Reward:**

$$r(x, y) = \beta \log \frac{\pi_\theta(y|x)}{\pi_{ref}(y|x)}$$

**Simplified Form:**

$$L_{DPO} = -\log \sigma \left( \beta \left[ \log \frac{\pi_\theta(y_w|x)}{\pi_{ref}(y_w|x)} - \log \frac{\pi_\theta(y_l|x)}{\pi_{ref}(y_l|x)} \right] \right)$$

**Components:**

- Policy:  $\pi_\theta$
- Reference Policy:  $\pi_{ref}$  (frozen)
- Preference pairs:  $(x, y_w, y_l)$  where  $y_w \succ y_l$

**Probability Interpretation:**

$$P(y_w \succ y_l) = \sigma(r(x, y_w) - r(x, y_l))$$

### 3.3 ORPO (Odds Ratio Preference Optimization)

**Pipeline:** Base Model → ORPO (single stage)

**Key Equation:**

$$L_{ORPO} = L_{SFT} + \lambda \cdot L_{OR}$$

**SFT Loss:**

$$L_{SFT} = -\mathbb{E}_{(x,y)}[\log \pi_{\theta}(y|x)]$$

**Odds Ratio Loss:**

$$L_{OR} = -\mathbb{E}_{(x,y_w,y_l)} \left[ \log \sigma \left( \log \frac{\text{odds}_{\theta}(y_w|x)}{\text{odds}_{\theta}(y_l|x)} \right) \right]$$

**Odds Definition:**

$$\text{odds}_{\theta}(y|x) = \frac{\pi_{\theta}(y|x)}{1 - \pi_{\theta}(y|x)}$$

**Combined Objective:**

$$L_{ORPO} = -\log \pi_{\theta}(y_w|x) - \lambda \log \sigma \left( \log \frac{\text{odds}_{\theta}(y_w|x)}{\text{odds}_{\theta}(y_l|x)} \right)$$

**Components:**

- Policy:  $\pi_{\theta}$  (no reference model needed)
- Preference pairs:  $(x, y_w, y_l)$

### 3.4 GRPO (Group Relative Policy Optimization)

**Pipeline:** Base Model → SFT → GRPO

**Key Equation:**

$$L_{GRPO}(\theta) = \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y_i \sim \pi_\theta(y|x)} [A_i \log \pi_\theta(y_i|x)] - \beta D_{KL}(\pi_\theta || \pi_{ref})$$

**Group Advantage:**

$$A_i = r(x, y_i) - \frac{1}{G} \sum_{j=1}^G r(x, y_j)$$

**Per-Sample Update:**

$$\nabla_\theta L = \sum_{i=1}^G (r_i - \bar{r}) \nabla_\theta \log \pi_\theta(y_i|x)$$

**Components:**

- Policy:  $\pi_\theta$
- Group size:  $G$  (typically 8-16 samples per prompt)
- Reward function:  $r(x, y)$  (verifiable, e.g., code correctness)
- Group mean:  $\bar{r} = \frac{1}{G} \sum_{i=1}^G r_i$

**Contrastive Form:**

$$L_{GRPO} \approx -\mathbb{E} [\log \sigma (r(x, y_+) - r(x, y_-))]$$

Where  $y_+$  has above-average reward,  $y_-$  has below-average reward

### 3.5 Method Comparison Matrix

| Method      | Models Required                       | Loss Type            | Data Format       | KL Constraint          |
|-------------|---------------------------------------|----------------------|-------------------|------------------------|
| <b>PPO</b>  | 4 ( $\pi$ , $\pi_{ref}$ , $r$ , $V$ ) | Clipped surrogate    | Online generation | Explicit penalty       |
| <b>DPO</b>  | 2 ( $\pi$ , $\pi_{ref}$ )             | Binary cross-entropy | Preference pairs  | Implicit via reference |
| <b>ORPO</b> | 1 ( $\pi$ )                           | SFT + Odds ratio     | Preference pairs  | None                   |
| <b>GRPO</b> | 1-2 ( $\pi$ , optional $\pi_{ref}$ )  | Policy gradient      | Prompts + rewards | Optional penalty       |

### 3.6 Mathematical Relationships

DPO  $\approx$  GRPO (offline case):

$$\lim_{G \rightarrow \infty} L_{GRPO} \approx L_{DPO} \text{ when rewards are preference-based}$$

ORPO vs DPO:

$$L_{DPO} \text{ uses } \log \frac{\pi_{\theta}}{\pi_{ref}}, \quad L_{ORPO} \text{ uses } \log \frac{\text{odds}_{\theta}(y_w)}{\text{odds}_{\theta}(y_l)}$$

PPO  $\rightarrow$  DPO simplification:

DPO eliminates:  $r_{\phi}$ ,  $V_{\psi}$  by deriving closed-form optimal policy

### 3.7 Optimization Targets

| Method      | What It Maximizes  |
|-------------|--|
| <b>PPO</b>  | $\mathbb{E}[r_{\phi}(x, y)] - \beta D_{KL}(\pi_{\theta}    \pi_{ref})$ |
| <b>DPO</b>  | $\mathbb{E}[\log P(y_w \succ y_l)]$ via implicit reward                |
| <b>ORPO</b> | $\text{Likelihood}(y_w) + \lambda \cdot \text{OddsRatio}(y_w, y_l)$    |
| <b>GRPO</b> | $\mathbb{E}_{group}[(r_i - \bar{r}) \log \pi(y_i)]$                    |

## 4 Emerging Trends (2025-2026)

| Trend                               | Method                | Application                         | Impact                                     |
|-------------------------------------|-----------------------|-------------------------------------|--|
| <b>Reasoning models boom</b>        | GRPO, RLVR            | Math, coding, formal logic          | GRPO democratizes reasoning model training |
| <b>Process reward models</b>        | PPO variants          | Step-by-step feedback               | OpenAI o1/o3 style training                |
| <b>Hybrid pipelines</b>             | DPO → PPO             | Initial safety → performance tuning | Best of both worlds                        |
| <b>RLAIF (RL from AI Feedback)</b>  | DPO + synthetic data  | Scalable preference generation      | Reduces human labeling cost                |
| <b>Multi-objective optimization</b> | PPO with multiple RMs | Safety + capability + style         | Gemini, Claude approach                    |