# Reinforcement Learning for LLMs

## Part 2

**Learn how agents improve their decision making strategy**

**Bhavishya Pandit**

# INTRODUCTION



COMMAND

Sit

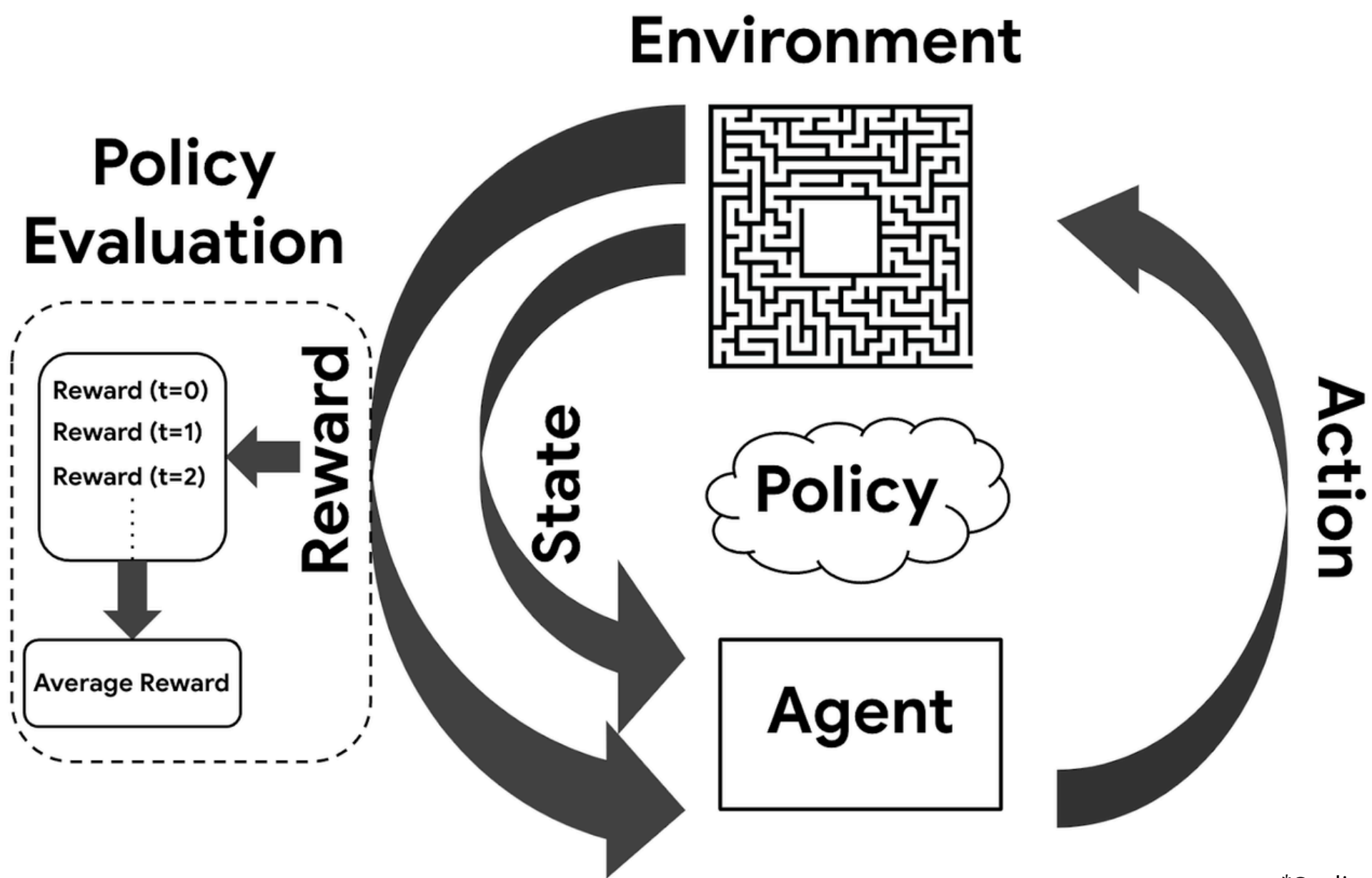Action
(stand)

COMMAND

Sit

Action
(sit)

Penalty

Reward

Reinforcement Learning (RL) is a type of machine learning where models learn through experiencemuch like how we do in real life. The model, known as an agent, interacts with an environment, takes actions, receives feedback in the form of rewards or penalties, and gradually improves its decision-making over time.

Think of it like training a dog to sit. At first, it doesn't respond correctly. You correct it. With repeated practice and rewards for the right actions, the dog eventually learns to sit on command.

That's the essence of RL learning by doing, improving through feedback.

If you're new to RL or want a quick refresher, check out our previous post linked in the comment section below!

**Bhavishya Pandit**

# WHAT IS POLICY

## Policy Evaluation

Reward (t=0)
Reward (t=1)
Reward (t=2)

Average Reward

**Reward**

**State**

## Environment

**Policy**
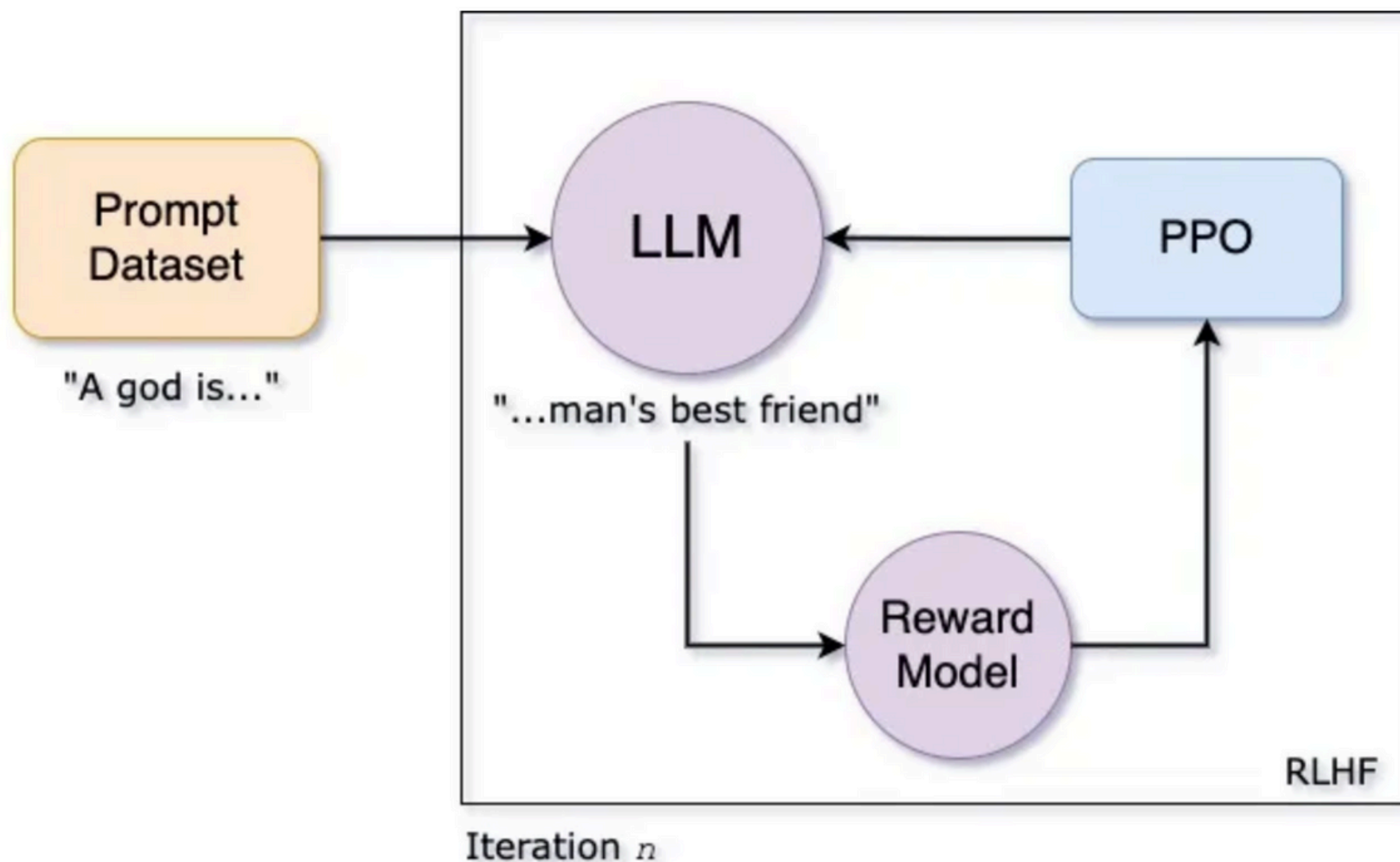
**Action**

**Agent**

*Credit:research.google

In Reinforcement Learning, a policy is like the agent's strategy or game plan. It tells the agent what action to take in a given situation. Think of it as a function that maps observations (or states) to actions. The better the policy, the better the agent performs in its environment.

Some key terms to know:
- **State**: The current situation or condition the agent is in.
- **Action**: What the agent decides to do in a state.
- **Reward**: Feedback the agent gets after taking an action positive or negative.
- **Policy evaluation**: The process of adjusting the policy to maximize total rewards over time.

**Bhavishya Pandit**
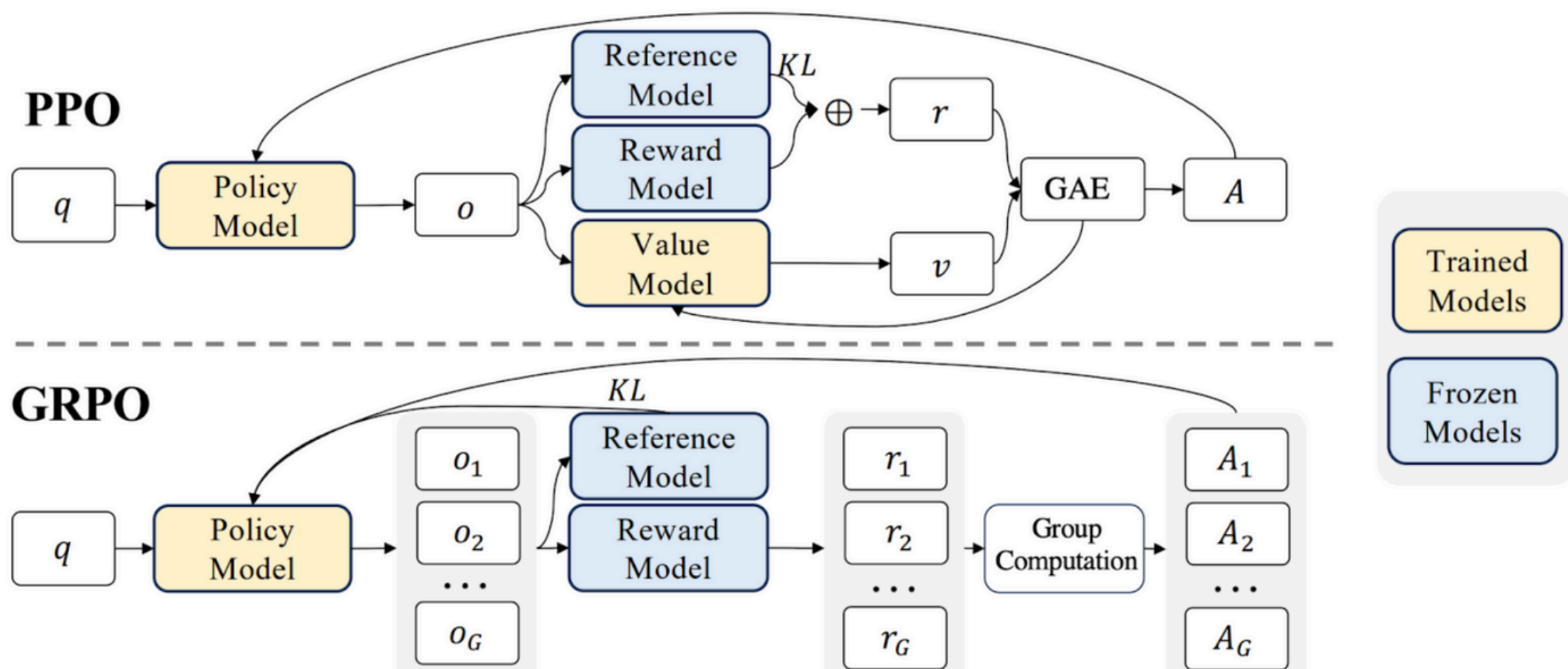
# WHAT IS PPO



Iteration $n$

PPO is a reinforcement learning technique.It improves a model (like an LLM or robot) step by step by learning from rewards.Instead of making big jumps, PPO takes small, safe updates to the model's behavior.Goal is Balance learning efficiently while not forgetting what it already knows.

Think of it like training a dog: You give it small rewards for sitting correctly and correct it gently when it's wrongover time, it learns the right behavior.

As you can see in the image, the process starts with a prompt like "A god is...", and the LLM responds (e.g., "...man's best friend"). A Reward Model evaluates the response, and PPO uses that score to update the LLM. This loop repeats, helping the model improve with each cycle.

**Bhavishya Pandit**

# WHAT IS GRPO

GRPO is a reinforcement learning method that boosts a model's reasoning by comparing a group of responses and learning from the best one. This group-based approach is more stable and efficient than scoring answers individually, as seen in PPO.

Imagine a classroom where students solve the same problem. Instead of grading each separately, the teacher highlights the best answerand everyone learns from it. Over time, the entire class improves
.
IIn PPO, Generalized Advantage Estimation (GAE) helps the agent evaluate how much better an action is compared to the average, enhancing learning efficiency.
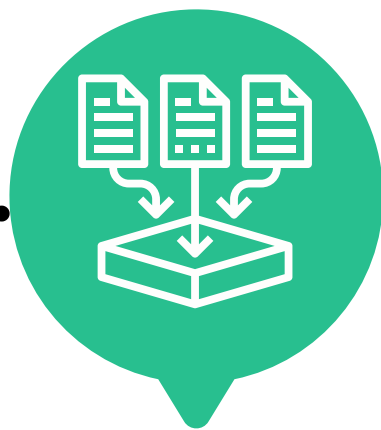
So while PPO rewards one answer at a time, GRPO learns from the strongest in a batch. This makes it especially powerful for tasks like math or logic, where comparing multiple responses helps build better reasoning.

**Bhavishya Pandit**

# HOW GRPO WORKS

Group Relative Policy Optimization (GRPO) enhances language model training by evaluating response groups, selecting the best, and refining updates with KL penalties. This boosts efficiency, reduces costs, and improves reasoning.

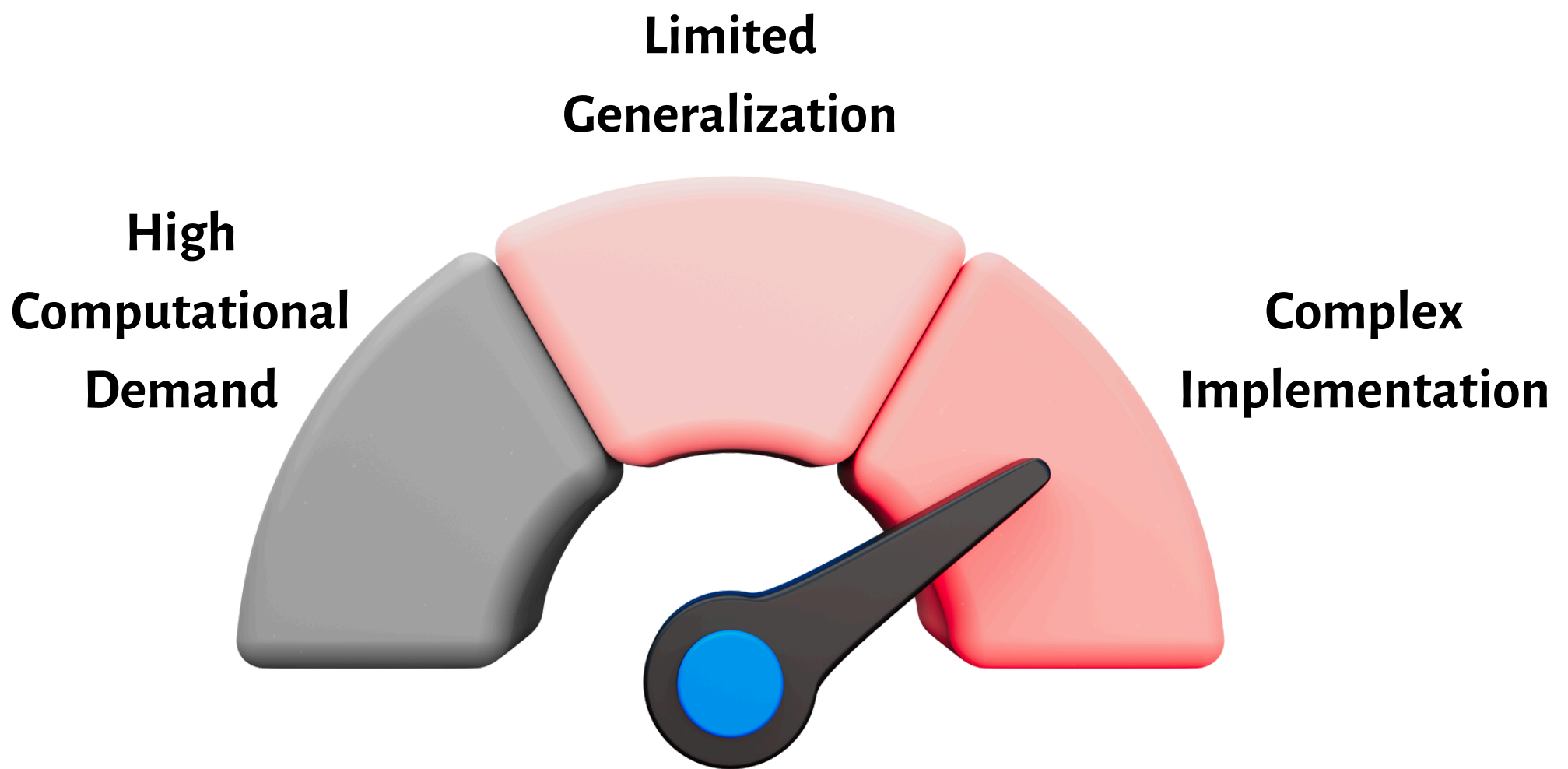| Group-Based Comparisons | Clipped Policy Gradient | KL Divergence Penalty |
|---|---|---|

**Group-Based Comparisons**: GRPO creates multiple responses for each prompt and compares them within the group learning from the best one. This approach leads to more stable and efficient training than scoring answers individually.

**Clipped Policy Gradient**: Controls how much the model's behavior can change at once. By limiting sudden shifts, it ensures learning stays stable and smooth avoiding overcorrections and keeping training on track.

**KL Divergence Penalty**: Keeps the updated policy close to the original (reference) model by penalizing large deviations. This helps maintain stability during training and ensures the model doesn't drift too far from what it already knows.

# LIMITATIONS

Despite its advantages, GRPO does have certain limitations and challenges:

**Limited Generalization**

**High Computational Demand**

**Complex Implementation**

**High Computational Demand:** Though more efficient than PPO, GRPO still demands high resources due to multiple completions per prompt.

**Limited Generalization:** Excels in structured reasoning but may require domain-specific tuning for broader RL applications.

**Complex Implementation:** Requires careful hyperparameter tuning and reward model design, making optimization challenging.

**Bhavishya Pandit**

# Stay Ahead with Our Newsletter! 🚀

👉 **Subscribe now and never miss an update!**
🔗 **https://bhavishyapandit9.substack.com/**

**Join our newsletter for:**

- Step-by-step guides to mastering complex topics
- Industry trends & innovations delivered straight to your inbox
- Actionable tips to enhance your skills and stay competitive
- Insights on cutting-edge AI & software development



💡 Whether you're a developer, researcher, or tech enthusiast, this newsletter is your shortcut to staying informed and ahead of the curve.
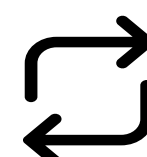
**Bhavishya Pandit**

# Follow to stay updated on Generative AI

👍

**LIKE**

💬

**COMMENT**

🔁

**REPOST**