

# LLM Training & Reinforcement Learning.

Pierre Lepagnol

2025-02-14

# Sommaire

Qu'est-ce le Reinforcement Learning ?

Comment c'est appliqué au LLMs ?

Qu'est-ce le Reinforcement Learning ?

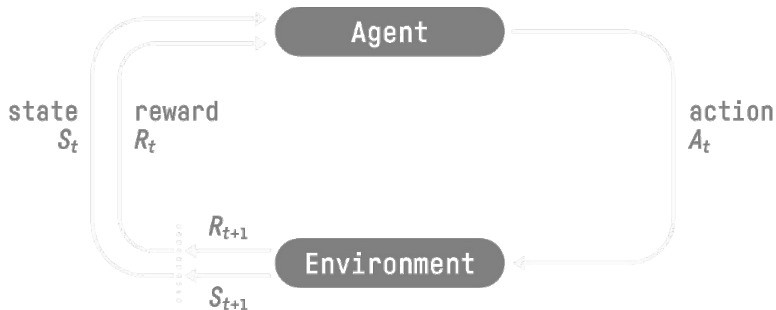
# Définitions & Basics



Un **agent** qui interagit avec l'**environnement** en prenant des **actions** de manière à maximiser la fonction de **valeur** (sommées des récompenses).

L'objectif du RL: apprendre la **politique optimale** qui maximise les récompenses futures par **essais et erreurs**.

# Composants essentiels du RL: Actions



- ▶ **Discrètes** : Exemple - LLM, où chaque token est une action est une action spécifique.
- ▶ **Continues** : Exemple - contrôler un robot, où bouger le bras mécanique est une action continue.

# Politique ( $\pi$ ) & Fonction de valeur ( $V$ )

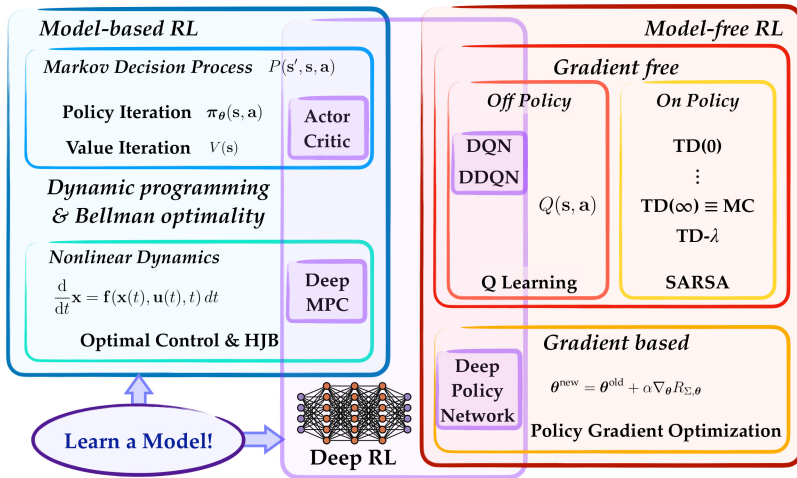
## Politique

- ▶ **Stratégie** (ensemble de règles) que l'agent suit pour décider quelle action entreprendre dans un état donné: Exemple : Réseau de neurones, une liste de If-then-Else, etc
- ▶ Peut être **déterministe** ou **probabiliste**.
- ▶ Formalisme :  $\pi(a|s)$ , ce qui signifie la probabilité de prendre l'action **a** dans l'état **s**.

## Fonction de valeur ( $V$ )

- ▶ Souvent la somme des **Somme des récompenses futures**
  - ▶ Mesure à quel point il est **bon** d'être dans un état particulier.
- Rappel Objectif: apprendre la **politique optimale** qui maximise les récompenses futures par **essais et erreurs**.

# Different algorithms dde RL



► Policy-Based methods : TRPO, PPO, GRPO

# Policy gradient method

- ▶ TRPO
- ▶ PPO
- ▶ GRPO



Comment c'est appliqué au LLMs ?

# Exemple de Deepseek-R1

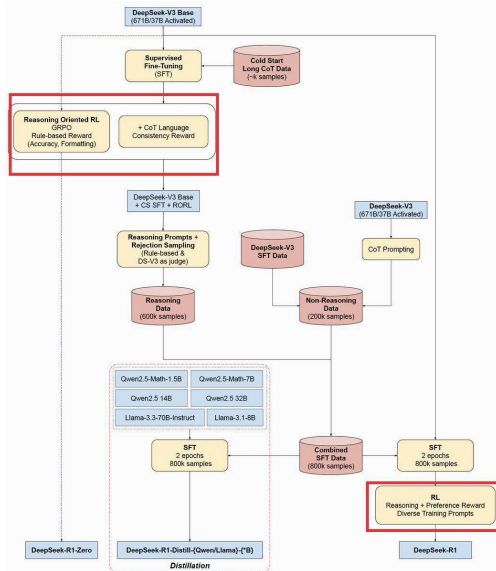


Figure 1: TrainingPipeline de DeepSeek R1

# RL during Post-Training

Alternative au RL (DPO)

# RL during Pre-Training

- GRPO is an improved version of Proximal Policy Optimization (PPO).

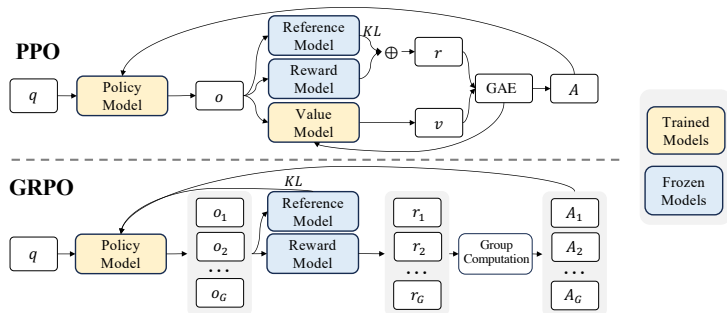


Figure 2: Comparaison PPO & GRPO issue de DeepSeekMath