# Training Pipeline

## Clod Start 🧊

**❓Question**
Which film has the director born first, The Model Couple or Thendral Veesum?

**✅ Ground Truth**
Thendral Veesum

⬇ SFT

**👣 Trajectory**
```
<think>
Need both directors' birth years.
</think>
<query>
The Model Couple birth year
</query>
<doc>
William Klein, born in 1946...
</doc>...
<answer>Thendral Veesum</answer>
```

**Base Model** 🤖

**Cold-Start Model** 🤖

## PPO Training ⚙

Reasoning → Answer

Policy Model

Process Reward

Outcome Reward

Reward 🏆

**R3-RAG** 🤖

# Reward Design

| Process Rewards | | | Outcome Rewards |
|---|---|---|---|
| ❌ **Format** | 🔍**Retriever Result** | | 🎯**Answer** |
| **Trajectory 1** 🏃 `<think>...<query>...</ query > <doc>... </doc> ... ...</answer>` **Reward: 0** | **Trajectory 1** 🏃 **Doc**: The Model Couple is by William Klein... **Reward: 0.8** | | **Trajectory 1** 🏃 Thendral Veesum's director was born earlier. **Reward: 2** |
| **Trajectory 2** 🏃 `<think>... <query>...</query> ... ...</query>` **Reward: -1** | **Trajectory 2** 🏃 **Doc**: Willem Klein, born in ... **Reward: 0.9** | | **Trajectory 2** 🏃 The Model Couple **Reward: 0** |
| **Trajectory 3** 🏃 `<think>...<think> <query>...</query> <answer>...</answer>` **Reward: -1** | **Trajectory 3** 🏃 **Doc**: ...a mathematician, was born in 4 December. **Reward: 0.1** | | **Trajectory 3** 🏃 The Model Couple's director is older. **Reward: 0** |
| **Format Verification** `<think>...</think> <query>...</query> <doc>... </doc> <answer>...</answer>` | **Relevance Verification** 🤖 → {Doc} | | **Answer Verification** Thendral Veesum |