

ine Tuning with TRL





Sergio Paniego Blanco (@sergiopaniego)

Machine Learning Engineer @ Hugging Face









From fine-tuning to TRL (fine-tuning+alignment)



"I fine-tuned the model and ..." 😅

It can't chat naturally
It can't reason or explain
It gives incoherent answers
It gives inconsistent answers
It's not aligned with human preferences
(unhelpful or unsafe)
It doesn't follow instructions



. . .

From fine-tuning to TRL (fine-tuning+alignment)

Classic fine-tuning adapts the model to **your dataset**, not necessarily to **your intentions**.

- Dataset (X→Y) → Pretrained model → Fine-tuned model
- Optimizes likelihood: predicts correctly based on the data
- Doesn't guarantee alignment with human preferences

We may need something stronger (TRL!), tools that can teach the model, not just fit it (TRL?!).

Data → Model → Fine-tuned model (Still says weird stuff)



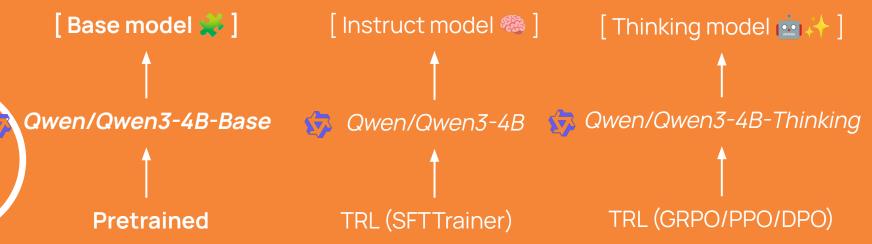
Where does TRL fits in the model lifecycle?

Pre-training: teaches general capacities such as language use, broad reasoning, and world knowledge

Mid-training: imparts domain knowledge (code, medical databases, internal docs...)

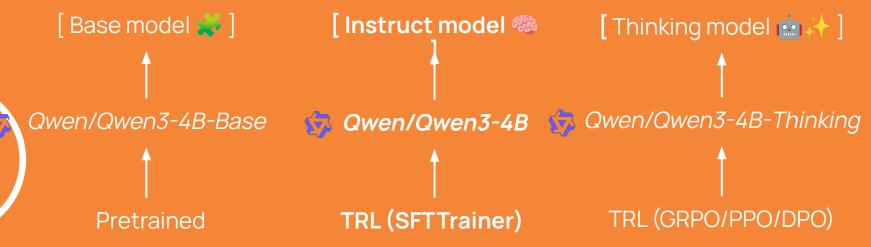
Post-training/Alignment (Text!): learns to behave as we want (instruction following, reasoning...)

TRL is not for pretraining → it's for the final stages, where we make models useful, safe, and aligned



Base: predicts next token Raw text, high quantity of data, masked language modelling, industrial compute

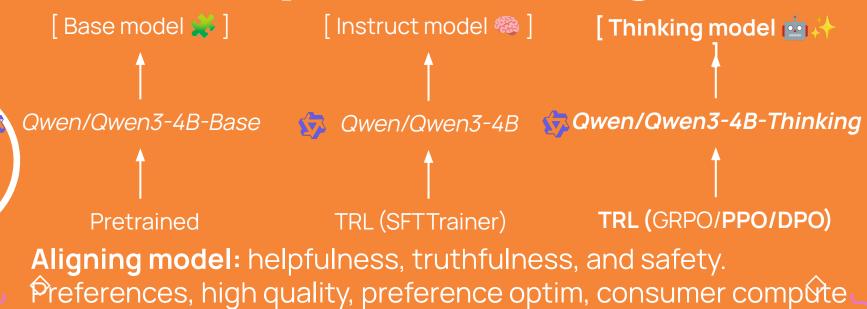
"Translate this sentence → Translate this sentence into..."



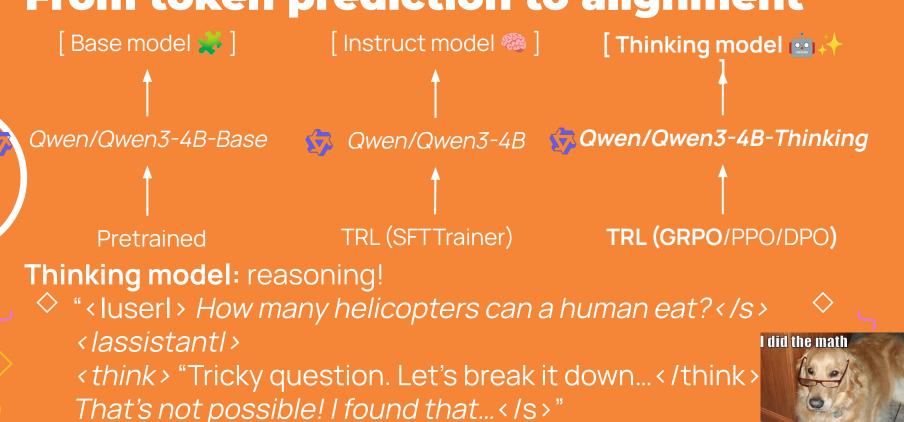
Instruct: follows instructions
Structured responses, high quality data, reduced compute

"<luserl> Translate this sentence </s>
<lassistantl> Sure! Here's the translation </s>"

Example dataset



- - "<luserI> How many helicopters can a human eat?</s>
 - (a) < lassistantl > Humans cannot eat helicopters < /s >
 - (a) < lassistantl > That's not possible! < /s > "
- Example dataset



Example dataset

TRL: Fine-tuning + alignment!

TRL: HF library that offers specialized trainers (SFT, GRPO, PPO, DPO...)

- TRL bridges the gap between instruction following and full alignment pipelines.
- Compatible with accelerate, PEFT, kernels, trackio, vLLM...
- Multimodal models support...
- Post-training: SFT, alignment (DPO), reasoning (GRPO)





Trainer VS TRL

TRL trainers are **built on top of transformers Trainer**, enhancing its functionality.

Each trainer **abstracts a full training method**, so you don't need to reimplement everything from scratch.

Gives ML engineers **ready-to-use recipes** for SFT, GRPO, PPO, DPO...

Focus: practicality, reproducibility, and reduced boilerplate

```
from trl import SFTTrainer
from datasets import load_dataset

trainer = SFTTrainer(
    model="Qwen/Qwen3-0.6B",
    train_dataset=load_dataset("trl-lib/Capybara", split="train"),
)
trainer.train()
```

SFT and GRPO

We'll focus here on two trainers. But there are more ••

Supervised Fine-Tuning (SFT)

- Teach the model to follow instructions
 Uses human-labeled data (prompt → response)
- Improves helpfulness

Group Relative Policy Optimization (GRPO)

- Train with online reinforcement learning
- Uses reward models or human preference signals
- Improves alignment and reasoning



TRL taxonomy

Online methods

(learn from feedback) Model generates reponses and learns from reward signals during training

- GRPO
- RLOO
- OnlineDPO
- NashMD 4
- XPO
- PPO

Reward modeling

Training a model to judge responses

- PRM
- Reward

Offline methods

(data-driven) Model learns from human preference datasets or supervised signals.

- SFT
- DPO
- ORPO
- BCO
- CPO
- KTO

Knowledge distillation

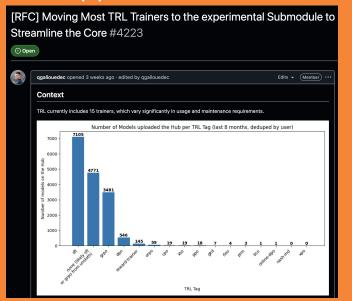
Transferring knowledge from a stronger model to a smaller one

GKD



SFT and GRPO are all you need?

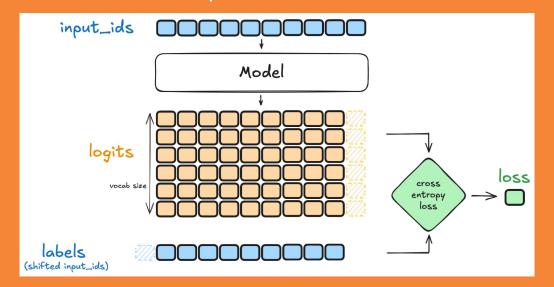
SFT and GRPO dominate the scene (apparently)
TRL includes +10 trainers → research moves fast!
Hard to maintain? There's even an open issue in the repo
New training methods appear!



Supervised Fine-Tuning

SFT is the simplest and most commonly used method to adapt a LM to a target dataset

- Goal: minimize negative log-likelihood of a sequence
- Simple, stable, and super effective.



Supervised Fine-Tuning

SFT Trainer supports both **standard** and **conversational** dataset formats.

```
# Standard
language_modeling_example = {"text": "The sky is blue."}

# Conversational
messages = [
    {"role": "user", "content": "Hello, how are you?"},
    {"role": "assistant", "content": "I'm doing great. How can I help you today?"},
    {"role": "user", "content": "I'd like to show off how chat templating works!"},
]
```

SFT is easy to instantiate and train

Only needs a model (even just the name!) and a dataset from the Hub to train

```
from trl import SFTTrainer
from datasets import load_dataset

trainer = SFTTrainer(
    model="Qwen/Qwen3-0.6B",
    train_dataset=load_dataset("trl-lib/Capybara", split="train"),
)
trainer.train()
```



SFTTrainer and SFTConfig

```
from trl import SFTTrainer, SFTConfig
from datasets import load_dataset
trainer = SFTTrainer(
   model="Qwen/Qwen3-0.6B",
   train_dataset=load_dataset("trl-lib/Capybara", split="train"),
   args = SFTConfig(
     per_device_train_batch_size = 2,
     gradient_accumulation_steps = 8,
     num train epochs = 1,
     learning_rate = 2e-4,
     optim = "paged adamw 8bit",
     logging_steps=1,
     report to="trackio",
     push_to_hub=True,
   peft_config=LoraConfig(),
trainer.train()
```

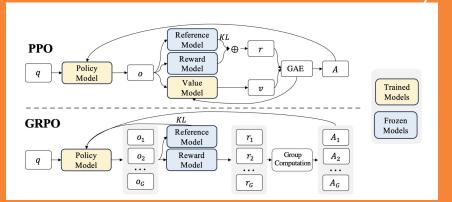


Group Relative Policy

GRPO was described plantage to paper.

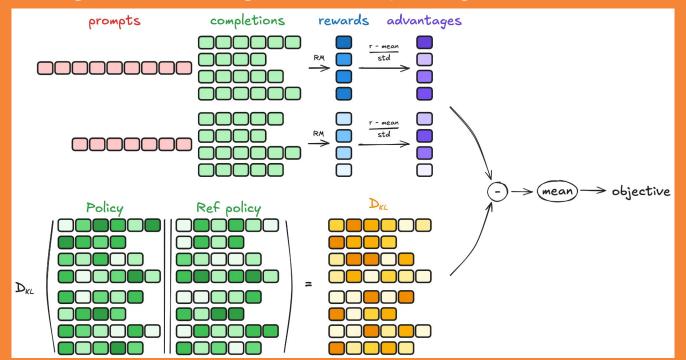
Variant of Proximal Policy Optimization (PPO), that enhances maths reasoning abilities while optimizing memory usage

- Online method
- Optimizes behavior using rewards from preference models or heuristics (reward functions)
- Successor of PPO: more stable and memory-efficient



GRPO

Steps: generating completions, computing advantage, estimating the KL divergence, computing the loss



GRPOTrainer and GRPOConfig

```
. . .
from trl import GRPOTrainer, GRPOConfig
import re
def format reward func(completions, **kwargs):
    """Reward function that checks if the completion has a specific format."""
    pattern = r"^<think>.*?</think><answer>.*?</answer>$"
    completion_contents = [completion[0]["content"] for completion in completions]
   matches = [re.match(pattern, content) for content in completion_contents]
    return [1.0 if match else 0.0 for match in matches]
trainer = GRPOTrainer(
   model="Qwen/Qwen2-0.5B-Instruct",
    reward funcs=[format reward func, ...],
   args=GRPOConfig(...),
    train dataset=load dataset("trl-lib/tldr", split="train"),
trainer.train()
```

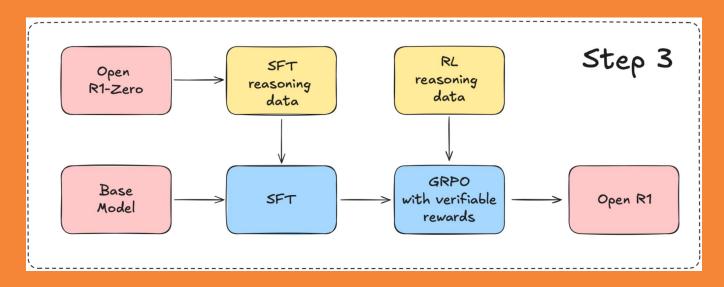
GRPOConfig detail

```
• • •
from trl import GRPOConfig
training_args = GRPOConfig(
    max_completion_length=512, # default: 256
    num_generations=8, # default: 8
    max prompt length=512, # default: 512
```



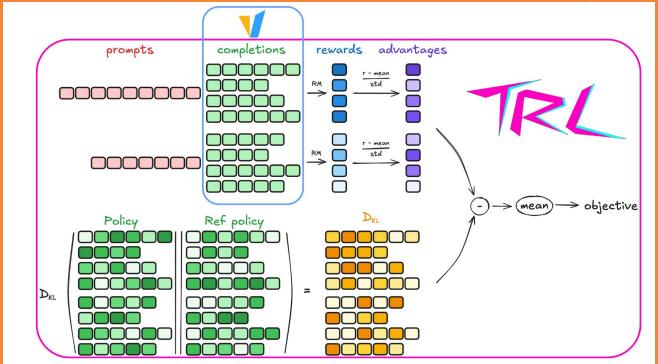
Initiative by HF to replicate and extend the techniques behind DeepSeek-R1 🐳

This model was built using SFT + GRPO!



Can we improve online training?

Some TRL trainers work online /, the model generates completions during training to compute a reward signal.











Can we improve online training?

*** Bottleneck! Generation is slow and inefficient **

*** vLLM solves this problem and TRL supports it! You can call it from the CLI (yes, TRL has a CLI)

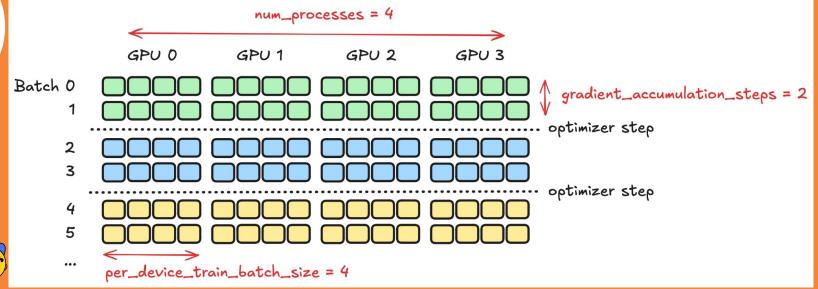


```
● ● ● ● CUDA_VISIBLE_DEVICES=0,1,2,3 trl vllm-serve --model Qwen/Qwen2.5-7B --tensor-parallel-size 2 --data-parallel-size 2
```

```
from trl import GRPOConfig

training_args = GRPOConfig(
    ...,
    use_vllm=True,
    vllm_mode="server", # default value, can be omitted. Can also be "colocate"
)
```

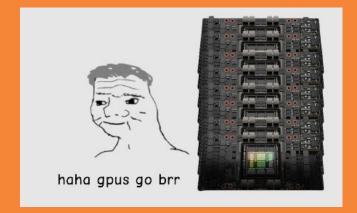
TRL also integrates seamlessly with accelerate and DeepSpeed (model sharing, zero redundancy optimized, mixed precision training, offloading...)





In the config we specify num of machines, num of processes...

accelerate launch --config_file examples/accelerate_configs/deepspeed_zero2.yaml train.py





• • •

Train long contexts via **context parallelism** (accelerate): by splitting the sequence across multiple GPUs Even **ND-parallelism** (N model replicas, N devices, tensor parallelism, context parallelism)

```
accelerate launch --config_file context_parallel_2gpu.yaml train.py
```

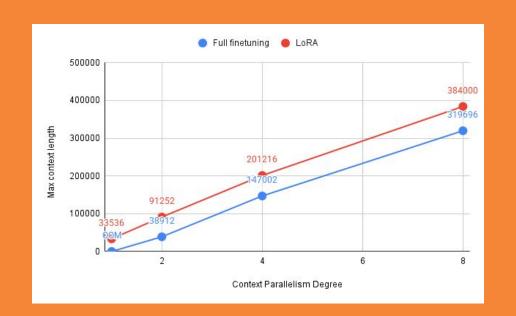




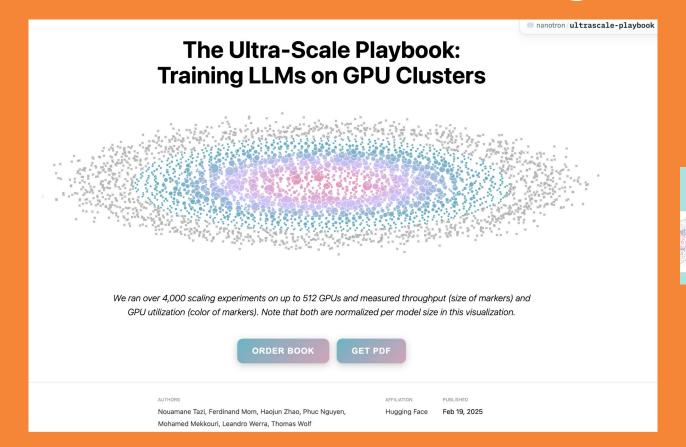




Context Parallelism (CP) with *Qwen/Qwen3-8B* scales to over **300k tokens in 8 GPUs**







The Ultra-Scale Playbook: Training LLMs on GPU Clusters

Optimizing all the things!

Transformers already provides us with some optimizing strategies like gradient checkpointing and we also have LoRA/QLoRA. In addition, TRL has:

- Truncation
- Packing
- Padding-free
- Padding sequences to a multiple
- Liger kernels W LIGER KERNEL
- Activation offloading

We developed a notebook for training a **14B** model using QLoRA in <u>free Colab</u>!

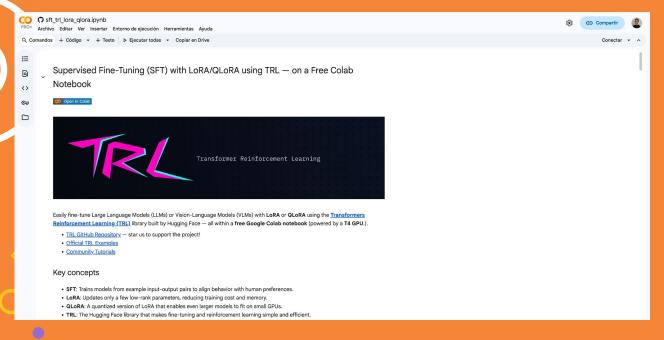


Optimizing is super easy



Optimizing for GPU poors

We developed a notebook for training a **14B** model using QLoRA in <u>free Colab</u>!







HF jobs + TRL

Models can be trained using HF's infra via HF Jobs

- trl-jobs cli package
- hf jobs (cli or python)

```
# pip install trl-jobs
trl-jobs sft --model_name Qwen/Qwen3-0.6B --dataset_name trl-lib/Capybara
```

```
hf jobs uv run \
---flavor a100-large \
---with trl \
---secrets HF_TOKEN \
train.py
```

```
from huggingface_hub import run_uv_job

run_uv_job(
    "train.py",
    dependencies=["trl"],
    flavor="a100-large",
    secrets={"HF_TOKEN": "hf_..."},
)
```











Multimodality!



Vision Language Models (VLMs) are getting stronger, but aligning them to human preferences still matters. TRL supports VLMs in many trainers, including examples and recipes.

- DPO (Direct Preference Optimization)
- MPO (Mixed Preference Optimization)
- GRPO (Group Relative Policy Optimization)
 - GSPO (Group Sequence Policy Optimization)





- RLOO (Reinforce Leave One Out)
- Online DPO

Multimodality!!



Native SFT support for VLMs

Just can still use your own collator

```
from trl import SFTConfig, SFTTrainer
from datasets import load_dataset

trainer = SFTTrainer(
    model="Qwen/Qwen2.5-VL-3B-Instruct",
    args=SFTConfig(max_length=None), # To avoid truncation that may remove image tokens during training train_dataset=load_dataset("trl-lib/llava-instruct-mix", split="train"),
)
trainer.train()
```



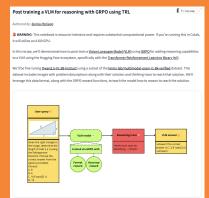
Multimodality!!!



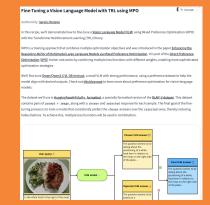
Shared scripts and recipes!

examples/scripts/dpo_vlm.py	This script shows how to use the <u>DPOTrainer</u> to fine-tune a Vision Language Model to reduce hallucinations using the <u>openbmb/RLAIF-V-Dataset</u> dataset.
examples/scripts/evals/judge_tldr.py	This script shows how to use <u>HfPairwiseJudge</u> or <u>OpenAlPairwiseJudge</u> to judge model generations.
examples/scripts/gkd.py	This script shows how to use the <u>GKDTrainer</u> to fine-tune a model.
trl/scripts/grpo.py	This script shows how to use the <u>GRPOTrainer</u> to fine-tune a model.
examples/scripts/grpo_vlm.py	This script shows how to use the <u>GRPOTrainer</u> to fine-tune a multimodal model for reasoning using the <u>Imms-lab/multimodal-open-r1-8k-verified</u> dataset.
examples/scripts/gspo.py	This script shows how to use GSPO via the <u>GRPOTrainer</u> to fine-tune model for reasoning using the <u>Al-MO/NuminaMath-TIR</u> dataset.
examples/scripts/gspo_vlm.py	This script shows how to use GSPO via the <u>GRPOTrainer</u> to fine-tune a multimodal model for reasoning using the <u>Imms-lab/multimodal-open-r1-8k-verified</u> dataset.
examples/scripts/kto.py	This script shows how to use the <u>KTOTrainer</u> to fine-tune a model.
examples/scripts/mpo_vlm.py	This script shows how to use MPO via the <u>DPOTrainer</u> to align a model based on preferences using the <u>HuggingFaceH4/rlaif-v_formatted</u> dataset and a set of loss weights with weights.
examples/scripts/nash_md.py	This script shows how to use the <u>NashMDTrainer</u> to fine-tune a model.
examples/scripts/online_dpo.py	This script shows how to use the <u>OnlineDPOTrainer</u> to fine-tune a model.
examples/scripts/online_dpo_vlm.py	This script shows how to use the <u>OnlineDPOTrainer</u> to fine-tune a a Vision Language Model.

examples/scripts/sft_vlm.py	This script shows how to use the <u>SETTrainer</u> to fine-tune a Vision Language Model in a chat setting. The script has only been tested with <u>LLaVA 1.5</u> , <u>LLaVA 1.6</u> , and <u>LLama-3.2-118-Vision-Instruct</u> models so users may see unexpected behaviour in other model architectures.
examples/scripts/sft_vlm_gemma3.py	$This script shows how to use the \underline{\textbf{SFTTrainer}}\ to fine-tune\ a\ Gemma3\ model\ on\ vision\ to\ text\\ tasks.$
examples/scripts/sft_vlm smol_vlm.py	This script shows how to use the <u>SFTTrainer</u> to fine-tune a SmolVLM model.







Fine-tuning SmolVLM using direct preference optimization (DPO) with TRL To Companies on a consumer GPU

Authored by: Sergio Panies

In this recipe, we'll guide you through fine-tuning a smol "Nision Language Model (VLM) with Direct Preference Optimization (OPO) using the Transformer Reinforcement Learning (TRL) library to demonstrate how you can tailor VLMs to suit your specific needs, even when working with coassimer grade OPUs.

We'll fine-tune <u>SmolVLM</u> using a preference dataset to help the model align with desired outputs. SmolVLM is a highly performant and memory-efficient model, making it an ideal choice for this task. If you're new to <u>Preference Optimization</u> for Januarian or visional parasurae models, the choice task this for it can include this immediation.

The dataset we'll use is <u>HuppineFaceH4r/fall*</u> formatted, which contains pairs of prompt + image along with a chosen and xejected answer for each pair. The goal of this fine-tuning process is to make the model consistently prefer the **chosen answers** from the dataset reducing ballularitations.

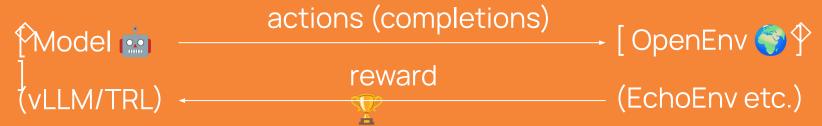
This notebook has been tested using an NVIDIA L4 GPU.







- **Solution OpenEnv** → framework for interactive TRL environments **Introducing OpenEnv**
- (Gymnasium for LLMs)
- Runs envs locally or as backend servers
- Find ready-to-use OpenEnv envs on the Hub
- Integrates with TRL → use `rollout_func` in GRPO to replace text generation with environment interaction
- Models can act, observe, and learn, not just predict text.



(rollout_func): bridges generation & env steps

Agents 🧟 via OpenEnv

```
from envs.echo_env import EchoEnv, EchoAction
from trl import GRPOConfig, GRPOTrainer
client = EchoEnv.from docker image("echo-env:latest")
def rollout_func(prompts, args, processing_class):
    response = requests.post("http://0.0.0.0:8000/generate/", json=payload)
    for msg in completions text:
        env result = client.step(EchoAction(message=msg))
        env rewards.append(env result.reward)
    result["env_reward"] = env_rewards
    return result
def reward_from_env(completions, **kwargs):
    """Extract environment rewards passed via rollout_func kwargs."""
    env_rewards = kwargs.get("env_reward", [])
    return [float(reward) for reward in env_rewards] if env_rewards else [0.0] * len(completions)
dataset = Dataset.from dict({"prompt": ["You are an AI that interacts with an *Echo* environment. Word to echo:"] * 64})
trainer = GRPOTrainer(
    model="Qwen/Qwen2.5-0.5B-Instruct",
    train dataset=dataset,
    rollout func=rollout func, # Use custom rollout
    args=GRPOConfig(
```

Resources 😵

Deep RL Course

- <u>Scripts</u>
- **Notebooks**
- Cookbooks
- **Example guides**
- **Blogs**
- **Courses**

















Thanks!



