**SnT 2025**
CTBT: SCIENCE AND TECHNOLOGY CONFERENCE

**8 SEPTEMBER**
ONLINE DAY
**9 TO 12 SEPTEMBER**
AT HOFBURG PALACE, VIENNA & ONLINE

P3.5-569

# An end-to-end LLM engineering platform for fine-tuning, evaluation and registration of custom models and adapters

Evangelos Dellis, Cahya Wirawan

Preparatory Commission for the Comprehensive Nuclear-Test-Ban Treaty Organization (CTBTO),
P.O. Box 1200, 1400 Vienna, Austria

## INTRODUCTION AND MAIN RESULTS

We present an end-to-end LLM engineering platform for fine-tuning, evaluation and registration of custom models and adapters that is built on top of open-source tools. The versatility of our platform is demonstrated through various applications such as fine-tuning multimodal open-source LLMs on custom datasets for increased accuracy.

Our platform aims to accelerate the development of custom models and adapters, enabling a wide range of innovative applications across CTBTO's technologies. These adapters consist of small collections of model weights that can be dynamically loaded onto a common base LLM, enabling it to specialize itself on-the-fly for specific tasks.

# An end-to-end LLM engineering platform for fine-tuning, evaluation and registration of custom models and adapters

Evangelos Dellis, Cahya Wirawan

**P3.5-569**

## Introduction

What is the end-to-end LLM engineering platform?

o A **platform** for building, evaluating, training, monitoring and configuring LLM assistants/agents
o It is based on popular **open-source** technologies/tools
o Helps with dataset creation, **data labeling** and annotation from private CTBTO data
o Streamlines the **fine-tuning**, deployment, and management of custom models and **adapters**
o Manage **LLM API access** and enforce budgets, guardrails, logging and cost tracking

**Who is it for?** Targeted to developers and admins.

Based on popular tools that are running on NVIDIA GPUs:

✓ **vLLM:** OpenAI API compatible LLM inference engine
✓ **Qdrant:** Vector database for similarity search
✓ **Airflow:** Pipeline orchestration tool
✓ **MLFlow:** Model registry, tracking, model wrapping
✓ **Langfuse:** Observability tool for traces and metrics
✓ **Easy Dataset:** Tool for creating fine-tuning datasets
✓ **Llama Factory:** LLM fine-tuning tool
✓ **Marker:** Convert PDF to markdown
✓ **DeepEval:** LLM evaluation framework

## LLM Engineering Platform

**LLM Application Observability:**
o Inspect and debug complex logs
o Ingest traces to **Langfuse**
o Track LLM calls & retrieval/embed

**LLM Evaluations:**
o LLM-as-a-Judge
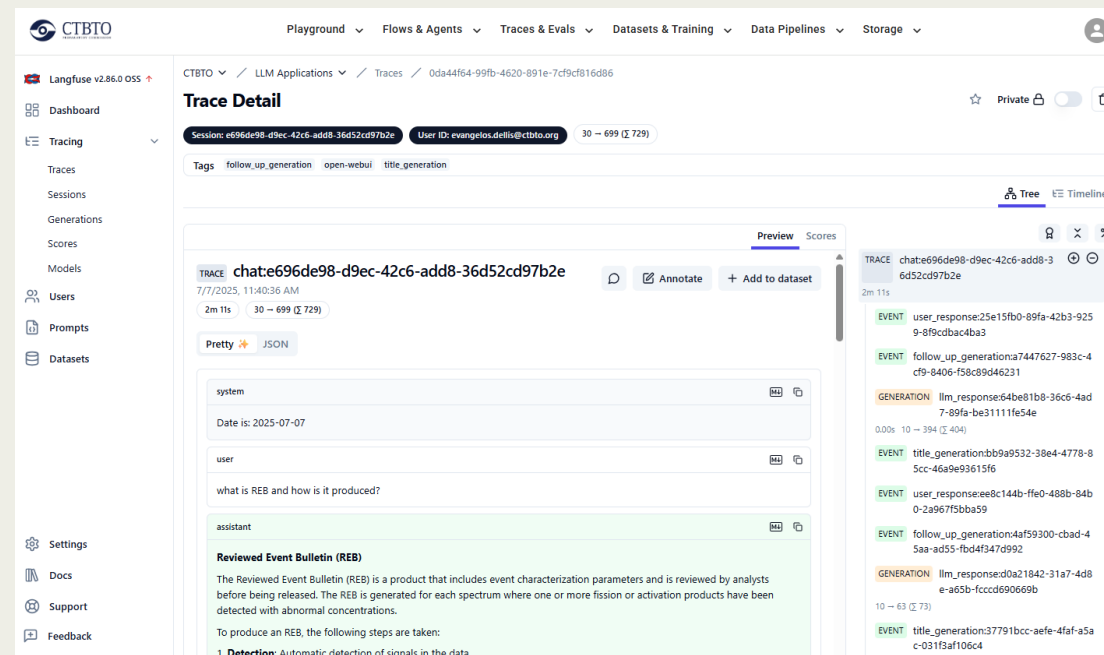o Manual annotations
o Custom evals via SDKs

**Datasets:**
o Test sets & benchmarks
o Structured experiments

**Prompt Management:**
o Manage & Version Control
o Collaborate on prompts



Evaluation is crucial for ensuring the quality and reliability of your LLM applications. **DeepEval** is an open-source framework and platform for evaluating, testing, and monitoring LLM applications:

**AI Assistant and Agent metrics:** Task Completion, Tool Correctness, Conversation Completeness
**RAG metrics: A**nswer relevancy, faithfulness, hallucination, and safety tests.
**Safety and Bias tests:** Detects Toxicity, Bias, and other vulnerabilities
**DeepEval features:** Integrates with CI/CD, offers synthetic dataset generation and many more.

CTBTO PREPARATORY COMMISSION | PUTTING AN END TO NUCLEAR EXPLOSIONS

# SnT 2025
**8 SEPTEMBER** ONLINE DAY
**9 TO 12 SEPTEMBER** AT HOFBURG PALACE, VIENNA & ONLINE
CTBT: SCIENCE AND TECHNOLOGY CONFERENCE

**An end-to-end LLM engineering platform for fine-tuning, evaluation and registration of custom models and adapters**

Evangelos Dellis, Cahya Wirawan

**P3.5-569**

## Creating fine-tuning datasets

Easy Dataset is a unified **framework** for synthesizing fine-tuning data from unstructured documents and it consists of two primary components:
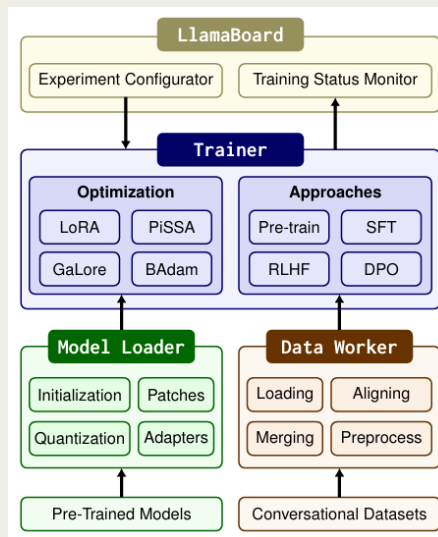
o **Adaptive document processing**: documents of different formats are processed via model-based parsing, followed by hybrid chunking to produce text chunks.

o **Data synthesis:** pairs are created for each document to guide the construction of diverse question-answer pairs. Questions are then generated to further increase diversity, and final augmented QA pairs are synthesized via knowledge-enhanced prompting to ensure factual consistency.
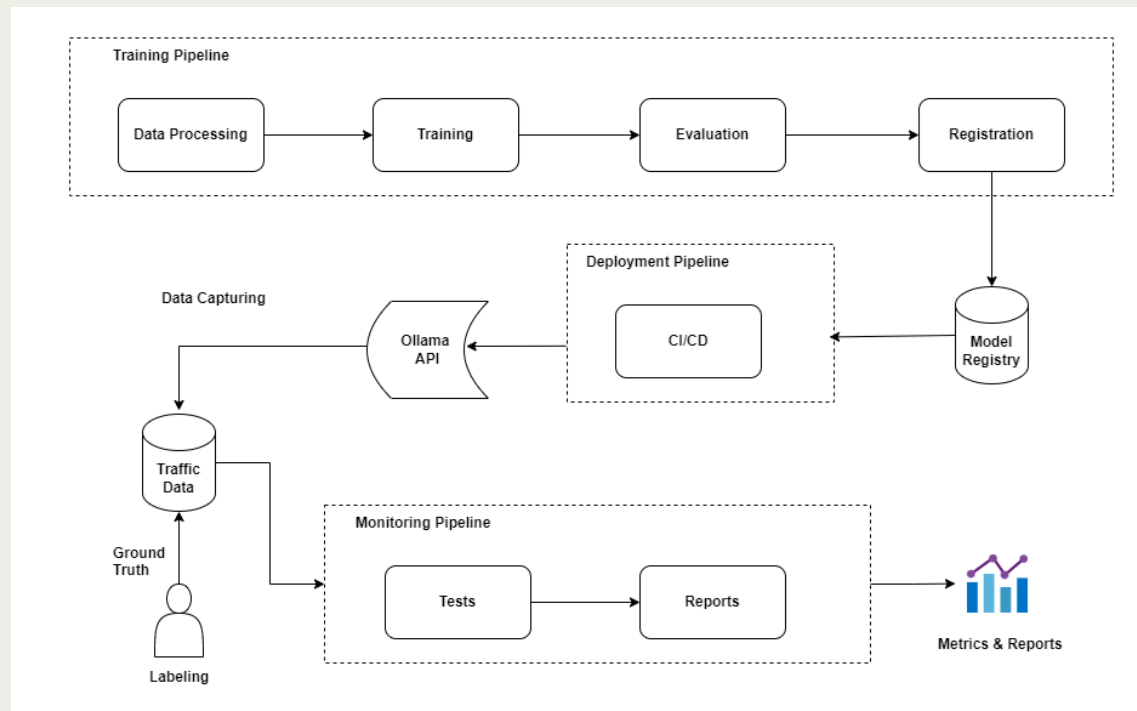


## LLM Fine-tuning using SFT





Our solution is based on **Llama Factory** which is a LLM fine-tuning tool that supports Pre-Training, Supervised Fine-Tuning, and Reward Modeling training modes.

o Various **models**: LLama, LLava, Mistral, Qwen, DeepSeek, Yi, Gemma, GPT-OSS, Phi, etc.
o Integrated **methods**: Pre-training, SFT, PPO, DPO, etc
o Scalable **resources**: LoRA and QLoRA via AQLM, AWQ, GPTQ, etc
o Advanced **algorithms**: GaLore, BAdam, APOLLO, Adam-mini, Muon, OFT, DoRA, LoRA+, LoftQ and PiSSA.
o Practical **tricks**: FlashAttention-2, Unsloth, Liger Kernel, RoPE scaling, NEFTune and rsLoRA.
o Wide **tasks**: Multi-turn dialogue, tool using, image understanding, visual grounding, video recognition, etc.
o **Experiment** monitors: LlamaBoard, TensorBoard, Wandb, MLflow, SwanLab, etc.
o Faster **inference**: OpenAI-style API, Gradio UI and CLI with vLLM worker.

CTBTO PREPARATORY COMMISSION | PUTTING AN END TO NUCLEAR EXPLOSIONS

# An end-to-end LLM engineering platform for fine-tuning, evaluation and registration of custom models and adapters

Evangelos Dellis, Cahya Wirawan

**P3.5-569**

## Our contribution towards an end-to-end LLM training system

1. **Create** Dataset:
   - use Easy Dataset
2. **Train** Model:
   - use Llama Factory
3. **Evaluate & Predict:**
   - use Llama Factory
4. **Export** merged Model:
   - use Llama Factory
5. **Deploy** Model:
   - use vLLM
6. **Import** Dataset:
   - use Langfuse
7. **Run** Experiments:
   - use Langfuse
8. **Create** Model:
   - use Ollama
9. **Push** Model to Registry:
   - use Mlflow



Building a **production LLM training system** is much more complex than training a model:

**What are the key challenges, i.e. how to:**
- Ingest large amounts of data, clean and preprocess that data and compute and serve features
- Set up a scalable training process, evaluate the model, track and version the data used
- Serve the model in a cost-effective manner, monitor the model and automate training and deployment

## Use Case: Paper Reviewer

**Source**: 9000+ NeurIPS and ICLR papers with its reviews from 2023/2024 collected from OpenReview.net

**Data processing:**
- Review normalisation using LLM API
- Cleanup the normalized reviews
- Semi synthetic dataset generation aligned by normalized review for each criteria
  - 9000- normalized reviews becomes 70.000+ reviews. A row per criteria
  - **Llama 3.1 8B** was used to generate the new dataset for 2 days using **8x V100 GPUs**

CTBTO PREPARATORY COMMISSION | PUTTING AN END TO NUCLEAR EXPLOSIONS