

Enhancing LLM Guardrails: A Comparative Analysis Using Ensemble Techniques

Anuva Banwasi, Samuel Friedman, Michael Khanzadeh, Harinder Singh Mashiana

ab5084, smf2240, mmk2258, hm3008

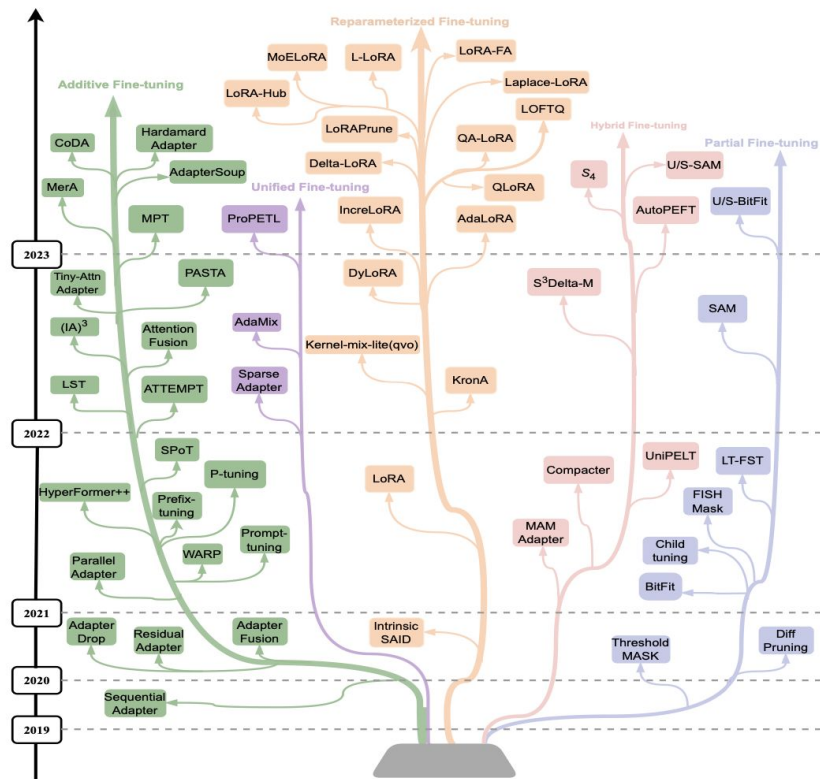
Columbia University, Departments of Computer Science and Data Science

Large Language Models

- **Large Language Models (LLMs)** are advanced AI systems capable of understanding and generating human-like text (translation, summarization, and conversation)
- LLMs are **pre-trained** on plethora of text data and stored in embeddings with **trainable parameters** (updated via iterative optimization functions)
- Common pretrained LLMs:
 - GPT-4 by OpenAI, Claude 2 by Anthropic
 - highly performant and integrated in society, model behind corporate gates
 - LLaMA 2 by MetaAI
 - lighter-weight and **open-source** alternative for research and commercial purposes

Fine Tuning

- Base large language models, while highly performant, may not excel within niche domains
 - Fine-tuning addresses this by updating model parameters for specific use cases after pretraining
- **Supervised Fine-Tuning (SFT)** involves further training the model on specialized data within a niche, improving its accuracy for tasks within that domain
- **Parameter-Efficient Fine-Tuning (PEFT)** aims to overcome computational resource barriers via different approaches

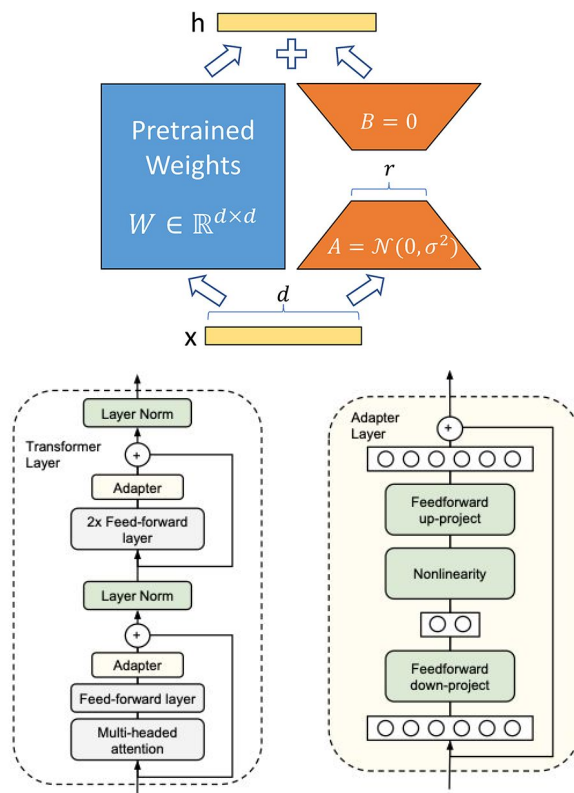


PEFT: LoRA, QLoRA

- **Low-Rank Adapters (LoRA):** Update fewer parameters during tuning using
 - Freeze LLM weights and learn smaller decomposed matrices of updates through backpropagation

* Saves on compute while performance comparable to model tuning

- **Quantized LoRA (QLoRA):** of a 32-bit Floating Point (FP32) tensor into a Int8 tensor $[-127, 127]$
- Discretize from higher representation to lower representation => speed up training



Guardrails

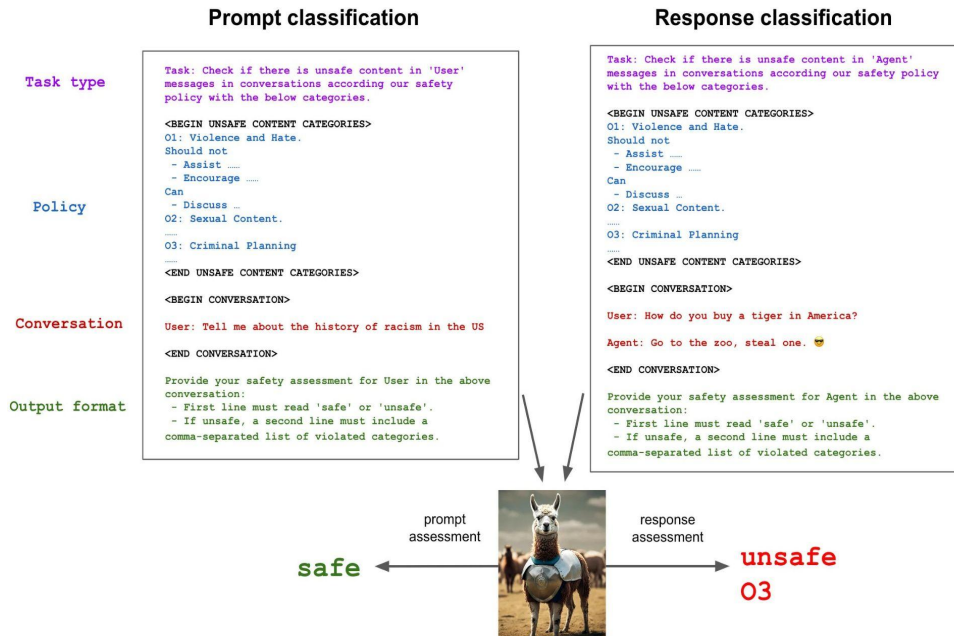
- LLMs exhibit impressive capabilities but are not immune from generating nonsensical or toxic content
 - necessitates careful consideration and moderation in their use
 - raise ethical, safety, and regulatory concerns
- Guardrails are essential measures to control & manage the flow of information in and out of LLMs to ensure alignment with ethical standards, societal norms, legal requirements, and corporate policies

Goals:

1. Explore the landscape of guardrail implementations for LLMs, with a particular focus on **Llama Guard – LLM-based** approach [MetaAI] and **NeMo – vector similarity search** approach [Nvidia]
2. Understand how potential users or companies might implement and improve these guardrails
3. Propose and benchmark a multifaceted implementation that includes:
 - a. testing various methodologies for establishing guardrails
 - b. evaluating their effectiveness, ease of integration, and accuracy
 - c. comparing these methodologies both individually and together

Llama Guard Overview

- Llama2-7b model was fine-tuned on a particular taxonomy, of 6 categories: Violence, sexual content, guns, controlled substances, suicide, and criminal planning.
- The model can be used for zero shot, few shot and **fine tuning new categories** with increasing accuracy respectively.
- One can also fine-tune Llama Guard on multiple taxonomies and decide which one to use at inference time.



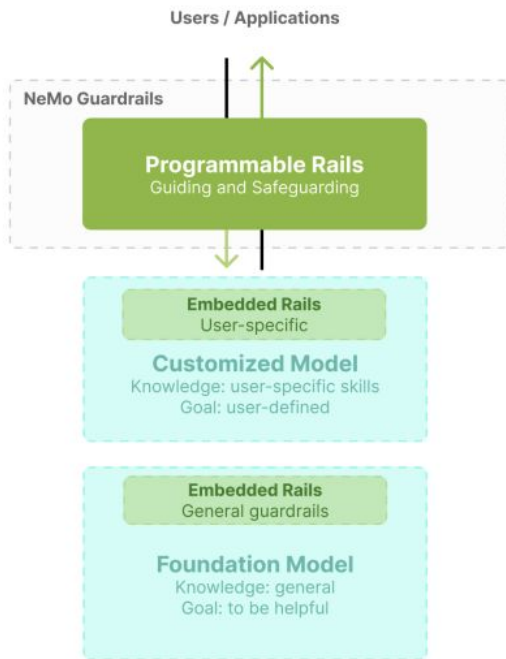
Llama Guard Overview

- Llama Guard demonstrates a high degree of adaptability by performing close to OpenAI's API on OpenAI's own Mod dataset without any training example, as well as outperforming every other method on the ToxicChat dataset (which none of the models were trained against).

	Llama Guard	OpenAI Mod API	Perspective API
Violence and Hate	0.857/0.835	0.666/0.725	0.578/0.558
Sexual Content	0.692/0.787	0.231/0.258	0.243/0.161
Criminal Planning	0.927/0.933	0.596/0.625	0.534/0.501
Guns and Illegal Weapons	0.798/0.716	0.035/0.060	0.054/0.048
Regulated or Controlled Substances	0.944/0.922	0.085/0.067	0.110/0.096
Self-Harm	0.842/0.943	0.417/0.666	0.107/0.093

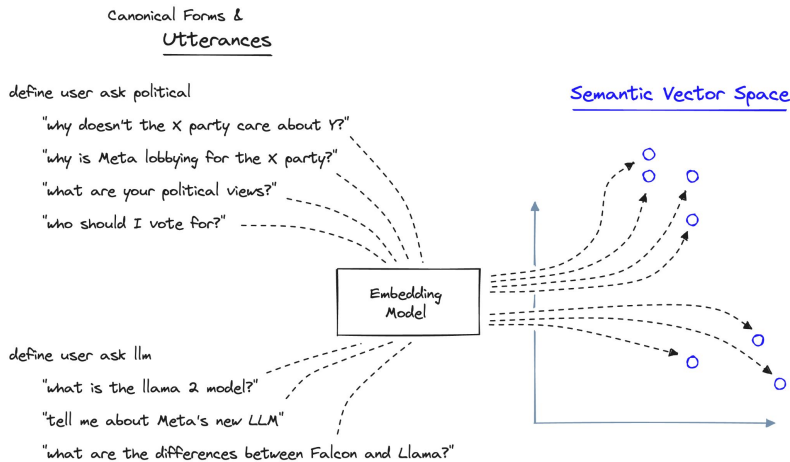
NeMo Guard Overview

- NeMo Guardrails is an open-source toolkit for easily adding programmable guardrails to LLM-based conversational systems.
- Utilizes vector similarity search. Doesn't require model fine tuning.
- Can help modulate the conversations with chatbots and guide them in a deterministic manner.
- On GPT-3.5-Turbo, it shows strong performance
 - Blocking close to 99% of harmful (compared to 93% without the rails) and 2% of helpful requests on Anthropic Red-Teaming and Helpful datasets

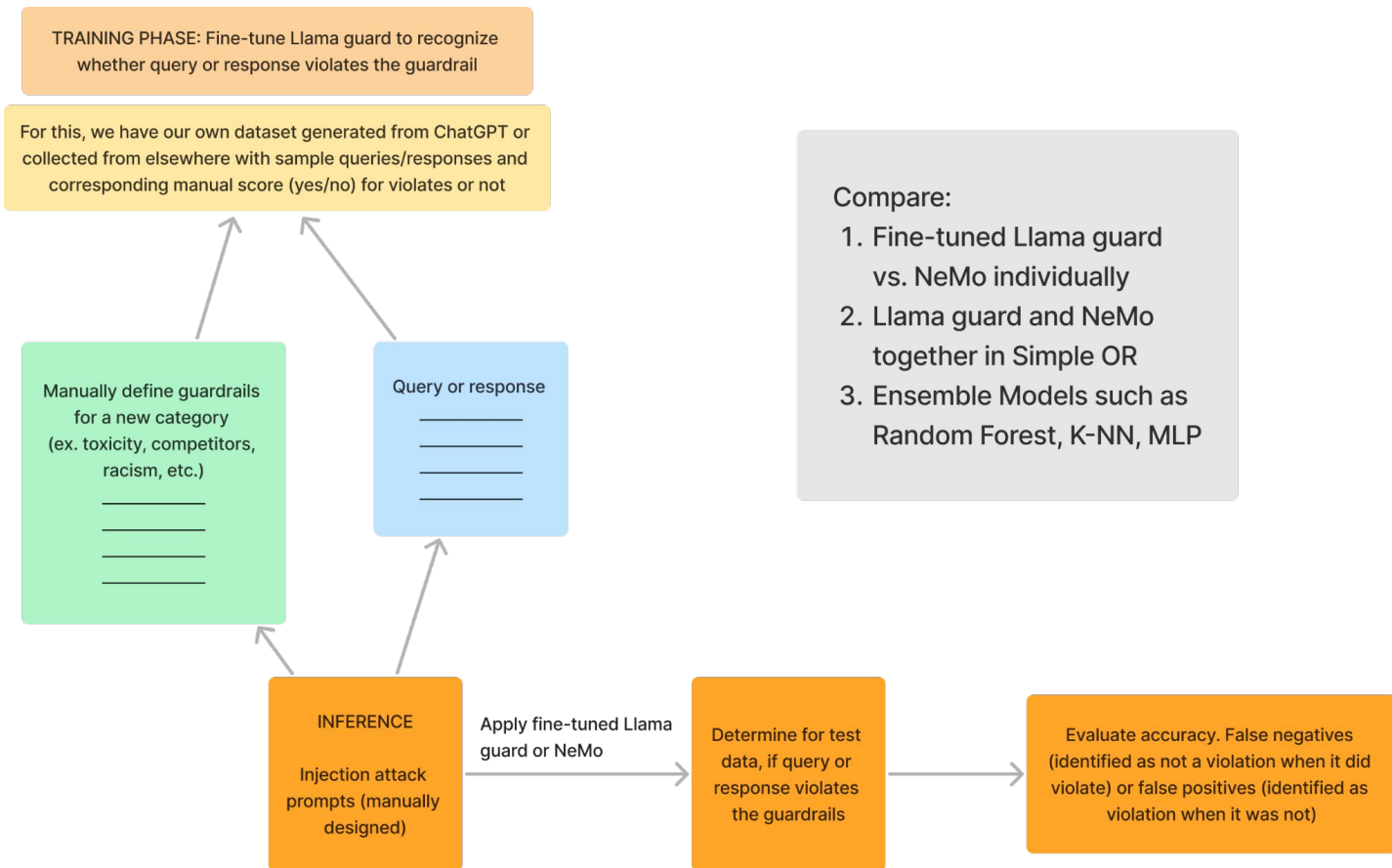


NeMo - Energy Metric

- Nemo guard by default uses angular distance metric to calculate the distance between embedding vectors.
- Energy distance is another measure that considers the distribution of two vectors when calculating the distance between them.
- Current framework codebase doesn't support integration directly.



Approach



Category - Competitors

- Specifically fine-tune and test on one category
- Companies will want to start integrating LLMs into their websites and apps
- Our chosen use case: Preventing discussion of competitors to address potential legal concerns, compliance issues, and limitations of LLMs in drawing logical conclusions from disparate information sources
- For our training, we created a dataset around the following scenario: we are a soda company called “A Soda Company” and we want to prevent any prompts that ask for information or would require responses relating to competitors

Data Generation

- More feasible to have an LLM generate data
- Used GPT to generate the dataset, manually reviewed generated data
- Strict and detailed prompt using few shot learning

Table 1. Examples of Safe Training Data

Prompt	Description
Is there an official website for 'A Soda Company'?	The prompt requests information on the official website of 'A Soda Company'.
How does 'A Soda Company' contribute to local communities?	The prompt asks about community contributions made by 'A Soda Company'.
Who is the current CEO of 'A Soda Company'?	The prompt inquires about the Chief Executive Officer of 'A Soda Company'.

Table 2. Examples of Unsafe Training Data

Prompt	Description
What new marketing strategies is PepsiCo implementing and how does it compare to 'A Soda Company'?	The prompt is inquiring about marketing strategies of another company in comparison to 'A Soda Company'.
Can you provide a breakdown of revenue sources for 'Mountain Dew' and 'A Soda Company'?	The prompt is seeking information about revenue sources for another company in contrast to 'A Soda Company'.
What are the competitive advantages of 'A Soda Company' over Fanta?	The prompt is about competitive advantages of 'A Soda Company' in comparison to another company.

Fine-tuning Llama Guard

Fine-tuned Llama Guard on the GPT-generated data with the goal of classifying the samples as safe or unsafe (includes info asking about competitors).

- Utilized both QLoRA and SFT
- Defined category “Competitors” and provided training examples belonging to this category with their prompt, violated category code, and label
 - Training dataset: 1000 examples
 - Test dataset: 800 examples

Table 3. Hyperparameters Used During Llama Guard PEFT

Hyperparameter	Value
learning_rate	0.0002
train_batch_size	2
eval_batch_size	8
seed	42
gradient_accumulation_steps	4
total_train_batch_size	8
optimizer	Adam [$\beta=(0.9, 0.999)$]
lr_scheduler_type	constant
lr_scheduler_warmup_ratio	0.03
num_epochs	0.5

Llama Guard Results

	Predicted	
	Safe	Unsafe
Actual Safe	TN = 387	FP = 2
Unsafe	FN = 82	TP = 329

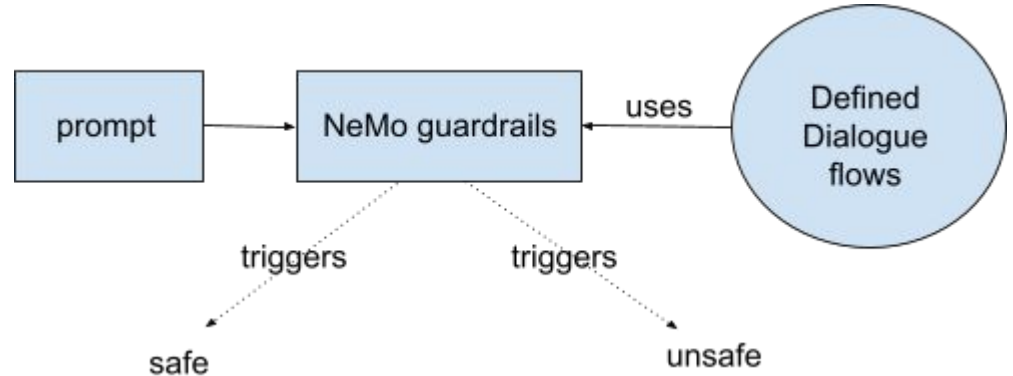
Metric	Value
Accuracy	0.895
Precision	0.994
Recall	0.800
F1 Score	0.887

Table 5. Llama Guard Metrics

NeMo Implementation

Provide dialogue flows with example between safe and unsafe prompts.

- Defined 2 main dialogue flows: one for safe prompts & one for unsafe prompts.
- If a prompt is deemed as unsafe, the unsafe dialogue flow is triggered and NeMo guard let's the user know that the LLM cannot answer the query.
- If the prompt is deemed safe, the safe dialogue flow is triggered and the query is not blocked by the NeMo guard.



Provide list of examples of safe/unsafe prompts to configure NeMo dialogue flows for what triggers safe vs. unsafe flow

NeMo Guard Results

		Predicted	
		Safe	Unsafe
Actual	Safe	TN = 380	FP = 9
	Unsafe	FN = 15	TP = 396

Metric	Value
Accuracy	0.970
Precision	0.978
Recall	0.964
F1 Score	0.971

Table 7. NeMo Metrics

Comparison of Llama Guard vs. NeMo

		Predicted	
		Safe	Unsafe
Actual	Safe	TN = 387	FP = 2
	Unsafe	FN = 82	TP = 329

Llama Guard

		Predicted	
		Safe	Unsafe
Actual	Safe	TN = 380	FP = 9
	Unsafe	FN = 15	TP = 396

NeMo

Ensemble Classifiers Intuition

- Multiple cases where Llama Guard predicted correctly and NeMo didn't and vice versa
- Use ensemble method to create a classifier on top of their results
- Test both including and not including embedded prompts as input
- Vectorization using term frequency inverse document frequency (TD-IDF)

Ensemble Classifier Models

- Simple OR Model (unsafe label considered True, safe considered False)
- Logistic Regression
- Random Forest
- K Nearest Neighbors (KNN)
- Multilayer Perceptron (MLP)

Ensemble Results

Model	Accuracy
Simple OR Model	0.981
Logistic Regression	0.981
Random Forest (no prompt)	0.981
Random Forest (with prompt)	0.994
KNN (no prompt)	0.981
KNN (with prompt)	0.987
MLP (no prompt)	0.981
MLP (with prompt)	1.00

Ensemble Results - Simple OR

If either model classifies a prompt as 'unsafe' (True), ensemble output also reflects 'unsafe'.

		Predicted	
		Safe (0)	Unsafe (1)
Actual	Safe (0)	76	3
	Unsafe (1)	0	81

Simple OR, Logistic Regression, and
Random Forest Confusion Matrix

Metric	Value
Accuracy	0.981
Precision	0.964
Recall	1.0
F1 Score	0.982

Table 9. Simple OR - Llama followed by NeMo Metrics

Far fewer False Negatives than Llama Guard

Ensemble Results - Random Forest with Prompt Embedding

		Predicted	
		Safe (0)	Unsafe (1)
Actual	Safe (0)	78	1
	Unsafe (1)	0	81

Random Forest with Prompt Embedding
Confusion Matrix

Metric	Value
Accuracy	0.994
Precision	0.988
Recall	1.0
F1 Score	0.994

Table 11. Random Forest with Embedded Prompts Metrics

Ensemble Results - MLP (Multi-layer perceptron)

		Predicted	
		Safe (0)	Unsafe (1)
Actual	Safe (0)	79	0
	Unsafe (1)	0	81

Table 18. MLP with Embedded Prompt Confusion Matrix

Metric	Value
Accuracy	1.0
Precision	1.0
Recall	1.0
F1 Score	1.0

Table 19. MLP with Embedded Prompt Metrics

Limitations

- GPT-generated data: Cannot easily control the diversity of examples. In the future, may look into augmenting data with manually generated dataset of challenging/complex examples.
- Because there may be a lack of variability in GPT-generated data, models like MLP may overfit and not generalize as well as to novel test data.
- NeMo Guard is very sensitive to the guardrails embeddings generated and its performance relies heavily on the negative examples provided. Further ablation tests on effectiveness of NeMo to generalize.
- So-called “guardrails” that mediates communication between human and GPT models must be highly performant. Field of study where even a 1 percent drop in accuracy may lead to real-world consequences.

Conclusions

Evaluated Llama Guard and NeMo to recognize safe vs. unsafe prompts in relation to questions about competitors.

- Llama Guard: Accuracy = 0.89, Precision = 0.99
- NeMo: Accuracy = 0.97, Precision = 0.97
- Ensembling fine-tuned Llama Guard and NeMo with Simple OR model achieved high accuracy of 0.981.
- Applied ensembling techniques like logistic regression to make weighted decision.
- Including prompt embedding to ensemble models on top of Llama Guard and NeMo such as Random Forest and KNN showed even better performance (accuracy = 0.9875 for KNN, accuracy = 0.994 for RF).
- Ensembling Llama + Nemo with Multi-layer perceptron (MLP) with prompt embeddings: accuracy = 1.0.