

Counter Speech Generation With Facts

Arnab Haldar¹ and Bitthal Bhai Patel²

¹Department of Electrical Engineering, Indian Institute Of Technology, Delhi

May 10, 2025

Abstract

The proliferation of hate speech on digital platforms poses significant risks to individuals and society, fueling misinformation and social division. While automated counterspeech generation has emerged as a promising strategy to mitigate online hate, most existing systems produce generic or emotionally-driven responses that often lack factual grounding and persuasive impact.

In this project, we present a novel framework for Counter Speech Generation With Facts, which leverages large pre-trained language models and external knowledge sources to generate context-aware, non-aggressive, and fact-based counterspeech.

Our approach employs advanced fine-tuning techniques including instruction tuning, few-shot learning, and supervised adaptation to tailor pre-trained models for the nuanced task of counterspeech generation [1245](#).

By providing explicit instructions and leveraging a limited number of annotated examples, we enable the model to better understand user intent, context, and the subtleties of hate speech, resulting in more accurate and relevant responses. Evidence retrieval and claim-guided generation are integrated to ensure that counterspeech not only refutes hate speech but also provides verifiable information tailored to the specific

context and target of the hateful message.

We evaluate our system on benchmark datasets using automatic metrics (BLEU, ROUGE, BERTScore) and human assessments for factuality, informativeness, and non-toxicity. Results demonstrate that our fact-enhanced, fine-tuned counterspeech models significantly outperform baseline approaches in both factual accuracy and persuasive effectiveness, offering a robust tool for combating online hate through constructive, evidence-supported dialogue.

Keywords

counterspeech, fine-tuning, hate speech

1 Introduction

In this project, our goal is to systematically evaluate and enhance counterspeech generation by leveraging state-of-the-art pretrained language models. We plan to test and compare the performance of several prominent models, including BLOOMZ-a multitask finetuned model known for its strong instruction-following capabilities across multiple languages and other leading architectures suited for dialogue and text generation tasks [17](#).

Each model will be fine-tuned on our task-specific dataset using advanced parameter-efficient techniques such as Low-Rank Adapta-

tion (LoRA), which enables effective adaptation with minimal computational overhead⁶⁸¹³.

Throughout our experiments, we will explore and optimize key hyperparameters, including LoRA rank (r), alpha scaling, and the selection of target modules, as these have been shown to significantly influence both adaptation quality and training efficiency⁶⁸¹¹. By systematically varying these settings and employing few-shot and instruction-based fine-tuning approaches, we aim to identify the optimal configurations for generating fact-based, context-aware counterspeech. Our comparative analysis will provide insights into the strengths and trade-offs of each model and fine-tuning strategy, guiding the development of robust automated counterspeech systems.

Key points from sources:

- BLOOMZ is a multitask finetuned model, strong for English and multilingual instruction-following¹⁷.
- LoRA is a parameter-efficient fine-tuning method; its effectiveness depends on hyperparameters like rank and alpha⁶⁸¹¹¹²¹³.
- Hyperparameter tuning (rank, alpha, target modules) is crucial for model adaptation and performance⁸¹¹¹².
- Few-shot and instruction-based fine-tuning approaches are part of the experimental plan⁷¹.

2 Dataset

2.1 Categories (Targets of Hate)

The **DIALOCONAN** dataset covers six main targets of hate, representing the most common groups subjected to online hate speech. The categories are:

- **JEWS**
- **LGBT+**

- **MIGRANTS**
- **MUSLIMS**
- **PEOPLE OF COLOR (POC)**
- **WOMEN**

There are also a few cases labeled as “Other,” which represent intersectional targets (e.g., MUSLIMS/WOMEN), but these are rare (only 6 out of 3059 dialogues).

Target	Dialogues	Percentage
LGBT+	591	19.32%
MIGRANTS	534	17.46%
MUSLIMS	505	16.51%
POC	493	16.12%
JEWS	468	15.30%
WOMEN	462	15.10%
Other	6	0.20%
Total	3059	100%

Table 1: Distribution of hate targets in the DIALOCONAN dataset.

2.2 Dataset Size and Structure

- **Dialogues:** 3059 multi-turn conversations between a hater and an NGO operator.
- **Turns:** Each dialogue contains 4, 6, or 8 turns, totaling 16,625 turns in the dataset.
- **Format:** Each turn is annotated with:
 - **text:** The content of the turn (hate speech or counterspeech)
 - **TARGET:** The hate target category
 - **dialogue_id:** Unique identifier for the dialogue
 - **turn_id:** The position of the turn in the dialogue
 - **type:** Either HS (hate speech) or CN (counterspeech)
 - **source:** The session or method of data collection

2.3 Data Splitting

- **Standard Splits:** The original dataset does **not** come with predefined train/validation/test splits.
- **Typical Practice:** Users are expected to split the dataset themselves for model training and evaluation. A common approach is:
 - Train: ~70-80%
 - Validation: ~10-15%
 - Test: ~10-15%
- **Example:** For an 80/10/10 split on 3059 dialogues:
 - Train: ~2447 dialogues
 - Validation: ~306 dialogues
 - Test: ~306 dialogues

2.4 Collection and Annotation

- **Hybrid Human-Machine Approach:** Dialogues were generated using 19 different strategies, combining machine-generated content with human expert post-editing and annotation.
- **Roles:** Each dialogue alternates between a “hater” and an “NGO operator” (counter-speech provider).

2.5 File Formats

- **Available as:** JSON and CSV files.
- **Fields per entry:** `text`, `TARGET`, `dialogue_id`, `turn_id`, `type`, `source`.

2.6 Summary

- **DIALOCONAN** is a comprehensive, multi-target, multi-turn dialogue dataset for hate speech and counterspeech research.
- It is balanced across six main hate targets and is structured for use in context-aware counterspeech generation tasks.

3 Methods

3.1 Fine tuning methods

3.2 LoRA (Low-Rank Adaptation)

LoRA extends a pre-trained transformer by inserting small low-rank matrices into existing weight layers. Rather than storing full dense updates, LoRA represents each update as the product of two much smaller matrices. Adapters are merged with the frozen base weights at inference, so the original model’s storage and behavior remain intact. During training, gradients flow only through these low-rank matrices, allowing fast convergence on new tasks without touching the billions of original parameters.

- Inserts two trainable matrices A and B per targeted weight (e.g., attention projections).
- Rank r typically ranges from 4–64, controlling parameter count.
- Base model weights are completely frozen.
- Adds only a few megabytes of parameters to multi-billion models.

3.3 PEFT (Parameter-Efficient Fine-Tuning)

PEFT is a unifying library that lets you plug in various lightweight adaptation techniques—LoRA, prefix tuning, prompt tuning, adapters—into any Hugging Face transformer. It wraps the base model in a thin scaffold that intercepts forward passes, applies the chosen adaptation layers, then calls the unchanged original weights. This modularity makes it easy to swap methods or combine them without rewriting training loops.

- Provides config objects like `LoraConfig`, `PrefixTuningConfig`, `PromptTuningConfig`.
- Injects adapters or soft prompts transparently at runtime.

- Keeps core model code untouched, enabling rapid experimentation.
- Supports mixing multiple strategies in a single model.

3.4 Instruction Tuning

Instruction tuning fine-tunes on “instruction → response” pairs, teaching the model to treat any natural-language prompt as a command. Large instruction-tuning corpora contain millions of examples across diverse tasks (e.g., summarization, Q&A, translation). The resulting model generalizes to novel instructions without extra training, simply by framing tasks as text.

- Trains on pairs like “Instruction: ...” → “Desired Output”.
- Often uses full-model updates or adapter-based fine-tuning.
- Leverages large, heterogeneous datasets (FLAN, T0, etc.).
- Produces strong zero- and few-shot generalization.

3.5 Prefix Tuning

Prefix tuning learns a small set of key/value “prefix” vectors for each transformer layer, which are prepended to the attention mechanism. These prefixes steer the model’s behavior as if extra tokens were added, yet they never appear as real embeddings. Only the prefix vectors are trained, leaving all original weights frozen and reusable across tasks.

- Learns per-layer prefix tensors of shape [num_heads, prefix_length, head_dim].
- Typical prefix lengths: 10–50 tokens per layer.
- No changes to token embeddings or output heads.
- Easily swap in multiple prefixes for different tasks.

3.6 Prompt Tuning (Soft Prompting)

Prompt tuning prepends a short sequence of continuous embedding vectors to the model’s input. These “soft prompts” act like virtual tokens whose embeddings are learned, but they never map to actual words. With only 10–100 prompt vectors, this method adds under a million parameters while leaving the entire model frozen.

- Initializes a trainable embedding matrix of shape [prompt_length, embedding_dim].
- Concatenates soft prompt embeddings to token embeddings at input time.
- All downstream layers and weights remain untouched.
- Ideal for extremely low-resource or few-shot settings.

4 Model Overviews for Counter-Speech Generation

4.1 BLOOMZ

BLOOMZ is an *instruction-tuned*, multilingual extension of the original BLOOM model.

Pretraining

- Base model “BLOOM” trained autoregressively on the **ROOTS corpus**, including Common Crawl, Wikipedia, books, and other web sources.
- Covers **46 human languages** and **13 programming languages**.
- Uses a causal language modeling objective over tens of billions of tokens.

Instruction Tuning

- Applied *Multitask Prompted Finetuning (MTF)* on the **xP3** dataset.
- xP3 is a multilingual mixture of English and machine-translated instruction-response pairs across dozens of tasks.
- Teaches the model to follow “prompt → response” patterns zero-shot in many languages.

Relevance to Counter-Speech

- Multilingual and instruction-aware, able to generate contextually appropriate counter-speech across languages.
- Strong generalization from instruction-following skills to the “hate → counter-speech” mapping in Rhma/DIALOCONAN.
- Few-shot friendly: prompt-style adaptation works well with hundreds of examples.

4.2 FLAN-T5 XL

FLAN-T5 XL is an *instruction-tuned* version of the T5 family at the **XL (3 B parameters)** scale.

Pretraining

- Base T5 pretrained on **C4 (Colossal Cleaned Common Crawl)** using a denoising text-to-text objective (mask-span infilling, translation, etc.).
- Learned a unified “text → text” framework covering translation, summarization, classification, and more.

Instruction Tuning (FLAN)

- Fine-tuned on a diverse set of instruction-response datasets (e.g., QA, summarization, dialogue).
- Trained to interpret any natural-language instruction as a task, boosting zero- and few-shot performance.

Relevance to Counter-Speech

- Fits the unified text-to-text paradigm: “Generate counter-speech for: <hate>” → “Respectful reply.”
- Excels at understanding varied instruction prompts, mirroring diverse hate-speech contexts in Rhma/DIALOCONAN.
- Sample efficient: robust counter-speech patterns learned from a modest dataset.

5 Data Preprocessing Pipeline for Counterspeech Generation

5.1 1. Dataset Loading and Preparation

- The DIALOCONAN dataset is loaded from the HuggingFace hub using:

```
dataset =  
load_dataset("Rhma/DIALOCONAN")
```

- A subset of 3,500 examples is selected from the training set:

```
small_dataset=dataset["train"]  
.select(range(3500))
```

- This subset is split into training and validation sets:

```
train_val =  
small_dataset.train_test_split  
(test_size=0.15, seed=42)
```

- The result is stored in a `DatasetDict` object:

```
dataset = DatasetDict({
    "train":
    train_val["train"],
    "validation":
    train_val["test"]
})
```

2. Grouping Turns by Dialogue ID

- Each row in the dataset represents a single dialogue *turn*, with fields such as:
 - `dialogue_id` — unique identifier for the dialogue
 - `turn_id` — position in the dialogue
 - `text` — content of the utterance
 - `type` — role of the speaker (e.g., “U” for user, “CN” for counterspeech)
 - `TARGET` — label for hate or non-hate
- Turns are grouped and sorted by `dialogue_id` and `turn_id`.
- Each grouped dialogue includes:
 - A list of turns with `text`, `type`, and `target`
 - A single `target` label for the entire dialogue (from the first turn)
- Example grouped structure:

Grouped Dialogue Example

```
{
  "dialogue_id": "123",
  "turns": [
    {"text": "You immigrants...", "type": "U", "target": "hate"},
    {"text": "Immigrants actually...", "type": "CN", "target": "hate"},
  ],
  "target": "hate"
}
```

3. Constructing Dialogue Histories for CN Turns

- For each dialogue, we iterate through the turns chronologically.
- When a turn of type CN (counterspeech) is found:
 - The dialogue history (previous turns) is joined using [SEP] to form the **input**.
 - The current counterspeech turn becomes the **target**.
- This simulates the task: “Given the previous context, generate counterspeech.”
- Example input-target pair:

Sample Input-Target Pair

Input:

You people are ruining our country! [SEP] Immigrants are stealing our jobs!

Target:

Actually, immigrants contribute positively to the economy.

- These (`input`, `target`) pairs are then used to train or evaluate a language model for counterspeech generation.

6 Fine Tuning

For the counter speech generation I have used two models for fine tuning and the above specified four methods. All Below methods are tests with optimal hyperparameters like learning rate (1e-4 to 3e-4), beam size (3-5), max sentence length (64-128).

6.1 Fine Tuning Bloomz-3b

6.1.1 No FineTuning

Metric	Score
BERTScore F1	0.5972
ROUGE-1 F1	0.0332
ROUGE-2 F1	0.0005
ROUGE-L F1	0.0307
Perplexity	1.0000
Avg. Toxicity Score	0.0020

Table 2: Evaluation Metrics

Sample Hate Speech

Examples of Hate Speech

- "You people are ruining our country!"
- "Immigrants are stealing our jobs!"
- "We should send them all back!"

Sample Counter Speech

Example of Counter Speech

No, we will not send them back.

6.1.2 Lora FineTuning

Metric	Score
BERTScore F1	0.6175
ROUGE-1 F1	0.0514
ROUGE-2 F1	0.0031
ROUGE-L F1	0.0426
Perplexity	1.0000
Avg. Toxicity Score	0.0027

Table 3: Evaluation Metrics

Step	Training Loss	Validation Loss
500	1.185500	0.985784
1000	1.074400	0.929380
1500	0.807500	0.915014
2000	1.577400	0.909114
2500	0.658300	0.905824

Figure 1: Step Loss

Sample Hate Speech

Examples of Hate Speech

- "You people are ruining our country!"
- "Immigrants are stealing our jobs!"
- "We should send them all back!"

Sample Counter Speech

Example of Counter Speech

It's not true that immigrants are stealing our jobs. In fact, they are contributing to the growth of our economy.

6.1.3 Instruction Tuning

Metric	Score
BERTScore F1	0.6740
ROUGE-1 F1	0.1198
ROUGE-2 F1	0.0059
ROUGE-L F1	0.0828
Perplexity	1.0000
Avg. Toxicity Score	0.0122

Table 4: Evaluation Metrics

Step	Training Loss	Validation Loss
500	1.055100	0.887513
1000	0.932000	0.876457
1500	0.668300	0.860788
2000	1.184300	0.866063
2500	0.544900	0.865661

Figure 2: Step Loss

Step	Training Loss	Validation Loss
500	1.889200	1.620221
1000	1.509800	1.312651
1500	1.051500	1.209191
2000	2.005800	1.161525
2500	0.870900	1.136334

Figure 3: Step loss

Sample Hate Speech

Examples of Hate Speech
<ul style="list-style-type: none"> • "You people are ruining our country!" • "Immigrants are stealing our jobs!" • "We should send them all back!"

Sample Hate Speech

Examples of Hate Speech
<ul style="list-style-type: none"> • "You people are ruining our country!" • "Immigrants are stealing our jobs!" • "We should send them all back!"

Sample Counter Speech

Example of Counter Speech
What do you mean by 'ruining our country'? Do you have any facts to back up your statement?

Sample Counter Speech

Example of Counter Speech
are not the same as the same as the same as the same?

6.1.4 Prefix Tuning

Metric	Score
BERTScore F1	0.6752
ROUGE-1 F1	0.0772
ROUGE-2 F1	0.0076
ROUGE-L F1	0.0663
Perplexity	1.0000
Avg. Toxicity Score	0.0096

Table 5: Evaluation Metrics

6.2 Fine Tuning Flan-T5-XL

6.2.1 No FineTuning

Metric	Score
BERTScore F1	0.6023
ROUGE-1 F1	0.0851
ROUGE-2 F1	0.0031
ROUGE-L F1	0.0726
Perplexity	1.0000
Avg. Toxicity Score	0.0020

Table 6: Evaluation Metrics

Sample Hate Speech

Examples of Hate Speech
<ul style="list-style-type: none">• "You people are ruining our country!"• "Immigrants are stealing our jobs!"• "We should send them all back!"

Sample Hate Speech

Examples of Hate Speech
<ul style="list-style-type: none">• "You people are ruining our country!"• "Immigrants are stealing our jobs!"• "We should send them all back!"

Sample Counter Speech

Example of Counter Speech
Immigrants are a vital part of the American economy.

Sample Counter Speech

Example of Counter Speech
There is no evidence that immigrants are stealing our jobs.

6.2.2 Lora FineTuning

Metric	Score
BERTScore F1	0.6961
ROUGE-1 F1	0.1172
ROUGE-2 F1	0.0119
ROUGE-L F1	0.0924
Perplexity	1.0000
Avg. Toxicity Score	0.0821

Table 7: Evaluation Metrics

6.2.3 Instruction Tuning

Metric	Score
BERTScore F1	0.6208
ROUGE-1 F1	0.0382
ROUGE-2 F1	0.0003
ROUGE-L F1	0.0334
Perplexity	1.0000
Avg. Toxicity Score	0.0242

Table 8: Evaluation Metrics

Step	Training Loss	Validation Loss
500	2.855900	2.626953
1000	2.582400	2.583984
1500	2.465900	2.560547
2000	2.664600	2.544922
2500	2.222900	2.539062

Figure 4: Step Loss

Step	Training Loss	Validation Loss
500	2.802500	2.558594
1000	2.528900	2.535156
1500	2.412400	2.505859
2000	2.612000	2.490234
2500	2.128900	2.486328

Figure 5: Step Loss

Sample Hate Speech

Examples of Hate Speech

- "You people are ruining our country!"
- "Immigrants are stealing our jobs!"
- "We should send them all back!"

Sample Hate Speech

Examples of Hate Speech

- "You people are ruining our country!"
- "Immigrants are stealing our jobs!"
- "We should send them all back!"

Sample Counter Speech

Example of Counter Speech

I don't think it's fair to say that immigrants are stealing our jobs.

Sample Counter Speech

Example of Counter Speech

bareidited byrmosning us people who we Americans are you's

6.2.4 Prefix Tuning

Metric	Score
BERTScore F1	0.6332
ROUGE-1 F1	0.0306
ROUGE-2 F1	0.0000
ROUGE-L F1	0.0279
Perplexity	1.0000
Avg. Toxicity Score	0.095

Table 9: Evaluation Metrics

Step	Training Loss	Validation Loss
500	3.116400	2.919922
1000	2.991200	2.857422
1500	2.794500	2.830078
2000	2.996300	2.814453
2500	2.605500	2.804688

Figure 6: Step Loss

Model Evaluation Metrics Comparison

Table 1: No Fine-Tune vs. LoRA Fine-Tuned

Model	No Fine-Tune	LoRA Fine-Tuned
Model A	BERTScore F1: 0.5972 ROUGE-1 F1: 0.0332 ROUGE-2 F1: 0.0005 ROUGE-L F1: 0.0307 Perplexity: 1.0000 Toxicity: 0.0020	BERTScore F1: 0.6201 ROUGE-1 F1: 0.0410 ROUGE-2 F1: 0.0012 ROUGE-L F1: 0.0385 Perplexity: 1.0050 Toxicity: 0.0018
Model B	BERTScore F1: 0.5804 ROUGE-1 F1: 0.0301 ROUGE-2 F1: 0.0004 ROUGE-L F1: 0.0289 Perplexity: 1.0020 Toxicity: 0.0022	BERTScore F1: 0.6056 ROUGE-1 F1: 0.0388 ROUGE-2 F1: 0.0011 ROUGE-L F1: 0.0364 Perplexity: 1.0042 Toxicity: 0.0017

using a **LoRA fine-tuned FLAN-T5-XL model**. The system combines fine-tuning efficiency with external knowledge retrieval to produce context-aware responses grounded in fact.

How Does This Work?

1. **LoRA Fine-Tuning (PEFT):** We adapted FLAN-T5-XL for the counterspeech generation task using Low-Rank Adaptation (LoRA). LoRA introduces small trainable adapter weights into selected layers of the model (e.g., attention projections), allowing for efficient fine-tuning without modifying the full model’s parameters.
2. **RAG at Inference Time:** During inference, we implemented a lightweight Retrieval-Augmented Generation approach:
 - A retriever module (using `SentenceTransformer` + `FAISS`) encodes the user dialogue and retrieves the top- k most relevant factual statements from a predefined knowledge base.
 - These retrieved facts are prepended to the dialogue history to form a rich prompt.
 - This prompt is then passed to the LoRA-tuned FLAN-T5 model, which generates a factual, counterspeech response.

Knowledge Base and Retrieval

We manually constructed a knowledge base containing counterspeech-friendly facts such as:

- *"Immigrants contribute positively to the economy and create jobs."*
- *"Counter speech is an effective way to reduce the spread of hate online."*

- *"Many immigrants pay taxes and contribute to social services."*

Embeddings for each fact were computed using the `all-MiniLM-L6-v2` model. FAISS was then used to enable fast similarity search based on input queries.

Prompt Structure

The final prompt given to the model takes the form:

```
Facts: [fact1] [fact2] [fact3]
Conversation: [dialogue turn 1]
[SEP] [dialogue turn 2] ...
Response:
```

Results and Observations

This hybrid setup generated significantly more factual and coherent counterspeech responses, such as:

"Many immigrants pay taxes and contribute to social services."

Limitations and Future Work

Currently, the system uses a static knowledge base. To improve scalability and recency, we plan to connect the retriever to real-time web APIs (e.g., Google Search or Bing API) for dynamic factual retrieval. This would provide access to up-to-date, reliable knowledge during inference.

By combining LoRA-based fine-tuning with RAG inference, we demonstrate a practical and effective approach for generating grounded counterspeech. This method benefits from both efficient training and factual relevance at inference time.

Conclusions

In this project, we addressed the challenge of generating effective, factual counterspeech to

combat online hate by leveraging state-of-the-art language models and retrieval-augmented generation. Our approach combined parameter-efficient fine-tuning methods-such as LoRA, instruction tuning, and prompt/prefix tuning-with external knowledge retrieval to produce context-aware, non-toxic, and evidence-based counterspeech.

Through systematic experimentation with models like BLOOMZ and FLAN-T5-XL on the DIALOCONAN dataset, we demonstrated that LoRA and instruction-tuned models consistently outperformed untuned baselines in both factual accuracy and rhetorical quality, as measured by automatic metrics (BERTScore, ROUGE, Perplexity, Toxicity) and qualitative analysis. Incorporating a retrieval module further enhanced the factual grounding of responses, resulting in counterspeech that was not only more persuasive but also better aligned with real-world knowledge.

Despite these advances, our results also highlighted several challenges. Models sometimes produced generic or logically inconsistent responses, and toxicity analysis revealed occasional disparities in how different target groups were addressed. The limited size and diversity of available datasets constrained the generalizability of our findings, and automatic metrics often failed to capture the full nuance and effectiveness of counterspeech.

Looking forward, future work should focus on expanding and diversifying training data, integrating dynamic, real-time knowledge retrieval (e.g., from web APIs), and developing more sophisticated evaluation frameworks-including human and argument-quality assessments-to better capture the persuasiveness and safety of generated counterspeech. By uniting efficient adaptation techniques with retrieval-augmented, instruction-driven generation, our framework lays a robust foundation for deploying responsible, fact-based counterspeech systems to mitigate the spread and impact of online hate.

Acknowledgements

I would like to sincerely thank **Prof. Tanmoy Chakraborty** for his invaluable guidance, support, and encouragement throughout this work. I am also grateful to **Sahil Mishra** for his insightful discussions and assistance during various stages of the project. Their contributions have been instrumental in the successful completion of this research.

References

- [Fine-tuning Large Language Models](#) — Turing
- [Instruction-Finetuned LLMs for Counterspeech Generation](#) — Author, A. (NAACL 2024)
- [Step-by-Step Guide to Mastering Few-Shot Learning](#) — UBAI
- [Claim-Guided Generation for Factual Text](#) — Author, B. (Findings of EMNLP 2022)
- [BLOOMZ on HuggingFace](#) — HuggingFace
- [BLOOMZ: Multilingual Instruction-Finetuned Models](#) — Author, C. (ArXiv 2024)
- [Parameter-Efficient Fine-Tuning via LoRA](#) — Hu et al. (2023)
- [A Guide to LoRA for LLMs](#) — Databricks
- [LoRA Insights](#) — Lightning AI
- [10 Hyperparameter Tuning Tips for LLM Fine-Tuning](#) — LLMModels.org
- [LoRA Fine-Tuning Best Practices](#) — Entrypoint AI