

# GenrAltOR: Generative AI for 'Omics Research

Daniel Claborne<sup>1</sup>, Matthew Jensen<sup>1</sup>, Javier E. Flores<sup>1</sup>, Lisa Bramer<sup>1</sup>,  
and Samantha Erwin<sup>1</sup>¶

<sup>1</sup> Pacific Northwest National Laboratory ¶ Corresponding author

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

## Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Abhishek Tiwari](#)

Submitted: 29 July 2025

Published: unpublished

## License

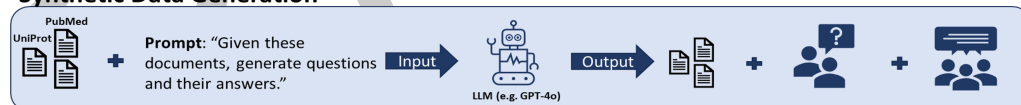
Authors of papers retain copyright  
and release the work under a  
Creative Commons Attribution 4.0  
International License ([CC BY 4.0](#)).

## Summary

We present software to fine-tune Llama 3 using Retrieval Augmented Fine-Tuning (RAFT) on proteomics literature. The package includes a command-line interface (CLI) that simplifies the fine-tuning process. The resulting model answers biomolecule-related queries using biological context.

Advances in LLMs enhance human-dependent tasks like interpreting biomolecule sets identified as important for model predictions. Domain experts can query an LLM for biological insights and draw inferences contextualized by the LLM's knowledge. Most LLMs are general-purpose, trained on broad datasets like social media. These lack the domain-specific language needed for 'omics queries. Our work uses RAFT (Zhang et al., 2024) to adapt an open-source LLM into an AI-assistant for domain experts. RAFT fine-tunes models by including irrelevant context in question-answer tasks, making them more robust than traditional retrieval-augmented generation (RAG) systems (Zhang et al., 2024). We provide a CLI to perform RAFT, from data collection to training, enabling researchers to create their own AI-assistants for biological research.

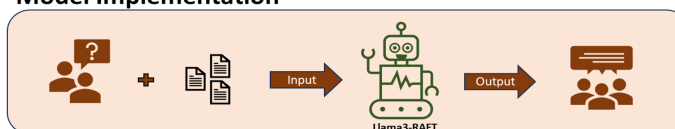
### Synthetic Data Generation



### Retrieval Augmented Fine Tuning (RAFT)



### Model Implementation



**Figure 1:** GenrAltOR Process Overview. Synthetic training data are generated using ChatGPT-4o. These question-answer-context triplicates are then used to fine-tune Llama 3 in a RAFT context. The output RAFT model is then implemented/evaluated on a hold-out set of generated triplicates.

## Statement of Need

RAFT has shown up to 30% performance gains over standard fine-tuning and RAG (Zhang et al., 2024). However, implementing RAFT presents many technical barriers for researchers.

23 Our software offers an easy-to-use package for performing RAFT in proteomics and serves as a  
 24 template for extending RAFT capabilities to other areas of molecular biology.

25 An overview of our development of a RAFT model is provided by [Figure 1](#). The general steps  
 26 are:

- 27 1. PubMed abstracts and UniProt data were retrieved using public APIs.
- 28 2. Data chunks were processed with GPT-4o to generate synthetic question-answer pairs.  
 29 Context chunks were grouped using text embeddings. An example is shown in [Figure 2](#).
- 30 3. Contexts were augmented with random 'distractor' documents to vary relevancy levels.
- 31 4. A training split of synthetic data was used to fine-tune Llama 3 for text completion.

32 We evaluated RAFT-Llama 3 against base Llama 3 using AlignScore ([Zha et al., 2023](#)).

Question: "[... CONTEXT ...] What techniques were used to establish the protein composition of chromatographic fractions?"

Answer: "[...] <ANSWER>: Two-dimensional polyacrylamide gel electrophoresis, N-terminal sequencing, endoproteinase Lys-C cleavage followed by peptide sequencing, comparison with ribosomal protein databases, and matrix-assisted laser-desorption ionization mass spectrometry (MALDI-MS)."

**Figure 2:** Example of a QA-pair generated using sampled context. The [... CONTEXT ...] chunk is a collection of documents that may or may not contain the text used to generate the QA-pair

## 33 Example Implementation

34 The process is encapsulated in a Python package with a command-line interface. First,  
 35 we retrieve context in the form of article abstracts and information about protein-  
 36 protein interactions and associated pathways, starting with a list of protein identifiers in  
 37 data/examples/uniprot.txt:

```
python -m genrator data:context \
--uniprot_ids=./data/examples/uniprot.txt \
--output_dir=./data
```

38 This will produce two files in ./data, one (uniprot\_context\_results...) with the  
 39 raw results of querying UniProt for pathway information and abstracts, and the other  
 40 (uniprot\_context\_postprocessed...) with context derived from those results and usable by  
 41 the RAFTDatasetPack class from the llama-index Python package ([Liu, 2022](#)).

42 To create synthetic question-answer pairs from this context, use the CLI with an OpenAI API  
 43 key and optionally a Hugging Face API key:

```
# set keys
export HF_TOKEN=<your-hf-token>
export OPENAI_API_KEY=<your-oai-key>

python -m genrator raft:data \
  --embed local \
  --context_path /path/to/context.txt \
  --output_path ./data/training/hf_dataset
```

44 This produces the training dataset at ./data/training/hf\_dataset, loadable via:

```
from datasets import load_from_disk

dataset = load_from_disk('/path/to/save_data_folder')
```

45 To fine-tune Llama 3 with RAFT, use the CLI target train:raft:

```
python -m genrator train:raft \
  -t /path/to/raft_data \
  -m meta-llama/Meta-Llama-3.1-8B \
  -n data/finetuned
```

46 The fine-tuned model is saved in data/finetuned and can be loaded via:

```
from transformers import (
    AutoModelForCausalLM,
    AutoTokenizer
)

tokenizer = AutoTokenizer.from_pretrained('./data/finetuned', padding_side="left")
model = AutoModelForCausalLM.from_pretrained("./data/finetuned")
```

47 Commands are configurable via flags or environment variables. Use --help for options:

```
python3 -m genrator train:raft --help
```

48 In addition to unifying data processing, model training, and result evaluation, this effort  
49 extended the Python package llama-index by:

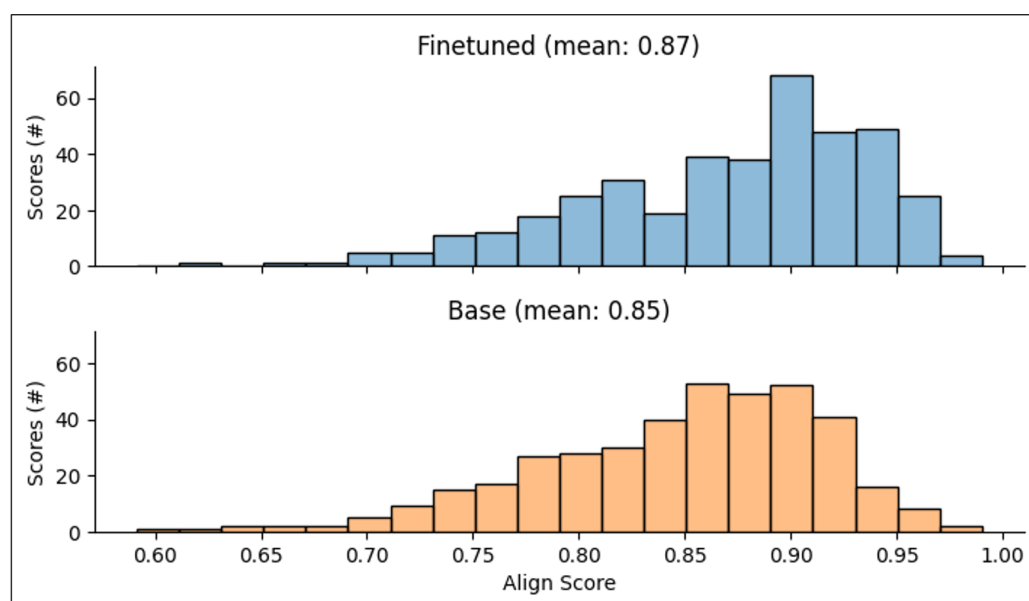
- 50 ■ Allowing configurable system prompts for generating questions.
- 51 ■ Modifying the get\_chunks method to respect the chunk\_size argument, optimizing text
- 52 length for training.

## 53 Results

54 We evaluated AlignScores on held-out question-context-answer triplets, benchmarking against  
55 base Llama-3. Figure 3 compares RAFT-Llama 3 (blue) and base Llama 3:

- 56 ■ RAFT shows a slightly higher mean AlignScore.
- 57 ■ RAFT's distribution is more left-skewed (towards lower scores).

58 These results suggest RAFT marginally improves response quality. Additionally, our software  
59 streamlines RAFT implementation, enabling researchers to develop RAFT-Llama 3 models.



**Figure 3:** Distribution of AlignScores for the RAFT-Llama 3 ("Finetuned") and the base Llama3 model.

## Discussion/Limitations

Some challenges remain:

- Resources** Large context chunks strained GPU memory, requiring limitation of context size.
- Proteomics focus.** The package focuses on proteomics; adapting to other biomolecules requires custom code, though training is domain-agnostic once data are ready.
- Dependencies.** Fine-tuning involves many dependencies, risking version mismatches (e.g., CUDA, PyTorch).

RAFT showed marginal AlignScore improvement over base Llama-3 on our QA task. This likely reflects our limited evaluation: a single metric (AlignScore) on synthetic data that could benefit from expert curation.

Our software package allows users to create context for performing RAFT in the proteomics domain given a list of proteins of interest. We hope this provides an easy-to-use base for researchers to explore the relationships between proteins identified in their experiments and expand the use of RAFT to different domain areas.

## Acknowledgements

The research described herein was funded by the Generative AI for Science, Energy, and Security Science & Technology Investment under the Laboratory Directed Research and Development Program at Pacific Northwest National Laboratory (PNNL), a multiprogram national laboratory operated by Battelle for the U.S. Department of Energy. This work was also supported by the Center for AI and Center for Cloud Computing at PNNL.

## References

- Liu, J. (2022). *LlamaIndex*. [https://github.com/jerryliu/llama\\_index](https://github.com/jerryliu/llama_index)

- 83 Zha, Y., Yang, Y., Li, R., & Hu, Z. (2023). *AlignScore: Evaluating Factual Consistency with*  
84 *a Unified Alignment Function*. arXiv. <https://doi.org/10.48550/ARXIV.2305.16739>
- 85 Zhang, T., Patil, S. G., Jain, N., Shen, S., Zaharia, M., Stoica, I., & Gonzalez, J. E. (2024).  
86 *RAFT: Adapting Language Model to Domain Specific RAG*. arXiv. [https://doi.org/10.](https://doi.org/10.48550/ARXIV.2403.10131)  
87 [48550/ARXIV.2403.10131](https://doi.org/10.48550/ARXIV.2403.10131)

DRAFT