# Safety classification using LM

**Zixuan Liu**
Tulane University
zliu41@tulane.edu

## Abstract

The rapidly increasing capabilities of large language models (LLMs) raise an urgent need to align AI systems with diverse human preferences to simultaneously enhance their usefulness and safety, despite the often conflicting nature of these goals. To address this important problem, a promising approach is to enforce a safety constraint at the fine-tuning stage, which introduces a cost model to indicate the cost value of the responses generated by the LLMs. In this project, we train a simple cost model using GPT-neo-1.3B and achieves a promising results on the test dataset.

## 1 Introduction

Large language models (LLMs) have demonstrated remarkable proficiency in tasks like chat completion, instruction following, coding, problem-solving, and decision-making (Chung et al., 2024; Ouyang et al., 2022; Anil et al., 2023; Stiennon et al., 2020). Considering the potential for broad societal impact, responses generated by LLMs must not contain harmful content, such as discrimination, misinformation, or violations of social norms and morals (Deshpande et al., 2023; Ganguli et al., 2022). Therefore, the alignment of safety in LLMs has received widespread attention from academia and industry (Christian, 2023).

An essential component of safety alignment involves minimizing the tendency of a model to generate harmful responses through fine-tuning. One of the important steps is to train a binary classifier to identify if a sentence contains harmful content (Dai et al., 2023). In this project, we will train a simple binary classifier using LM to identify if the sentence contains harmful language.

## 2 Related Work

The goal of LLMs alignment is to ensure that LLMs do not generate harmful or objectionable responses to user queries (Zou et al., 2023). To this end, multiple fine-tuning mechanisms have been employed for this task (Bai et al., 2022b; Burns et al., 2023; Munos et al., 2023). In particular, Constitutional AI (Bai et al., 2022b) trained a non-evasive and harmless AI assistant through self-improvement, which involves a supervised learning stage to get the model "on-distribution" and a reinforcement learning stage to further refine and improve the performance. Recently, OpenAI introduced the concept of superalignment, which aimed at solving the challenge of aligning AI systems that are much smarter than humans (Burns et al., 2023). They proposed the idea of weak-to-strong generalization, inspired by the generalization properties of deep learning, to control strong models with weak and less capable supervisors (Burns et al., 2023). (Munos et al., 2023) proposed Nash learning from human feedback, where they focused on learning a preference model and computing the Nash equilibrium of the model to advance the alignment of LLMs with human preferences.

RLHF has emerged as a central component of training state-of-the-art large language models (LLMs) such as OpenAI's GPT-4 (OpenAI, 2023), Meta's Llama 2-Chat (Touvron et al., 2023), with the goal of producing safe models that align with human objectives (Christiano et al., 2017; Bai et al., 2022a; Ziegler et al., 2019).Recent works such as direct preference optimization (DPO) (Rafailov et al., 2023) and SLiC-HF (Zhao et al., 2023) have successfully optimized the LLMs directly from human preferences without learning a reward model. However, these approaches have assumed a single preference function, which can barely cover the diverse preferences, expertise, and capabilities of humans (Bobu et al., 2023; Peng et al., 2023). To this end, fine-grained preference modeling and techniques for combining multiple dimensions of human preferences have been proposed (Bıyık et al.,

2022; Wu et al., 2023; Zhou et al., 2023). Further, (Dai et al., 2023) explicitly decoupled helpful and harmless to ensure the model outputs high-quality responses while maintaining a high level of safety.

## 3 Methods

In the reward modeling phase of RLHF, we represent human preferences using Bradley-Terry (BT) model (Bradley and Terry, 1952): given a prompt $x$ and a response $y$, we assume the pointwise reward of $y$ given $x$ is $r(x, y)$, which can be interpreted as the ground truth reward function that generates preferences. Then the BT model represents the human preference distribution $p^*(y_1 \succ y_2|x)$ as a function of the difference between two rewards:

$$p^*(y_1 \succ y_2|x) = \frac{\exp(r(x, y_1))}{\exp(r(x, y_1)) + \exp(r(x, y_2))} \quad (1)$$

where $y_1 \succ y_2|x$ denotes $y_1$ is preferred and $y_2$ is dispreferred amongst a pair of responses.

In the safety alignment framework, a cost model $c$ is introduced to discriminate between safe and unsafe responses generated by the LLMs (Dai et al., 2023). This model preserves the characteristics of the Bradley-Terry model, but it differentiates between safe and unsafe responses by employing a zero threshold. Given a dataset $D = \{x^i, y_\omega^i \succ y_l^i, s_\omega^i, s_l^i\}_{i=1}^N$, where $y_\omega \succ y_l$ denotes $y_l$ is safer than $y_\omega$, $s(y) = 1$ if $y$ is unsafe and $s(y) = -1$ otherwise. We can learn a cost model using the following pairwise comparison loss as shown in (Dai et al., 2023).

$$
\begin{aligned}
L(c; D) = & - \mathbb{E}_{(x, y_\omega, y_l) \sim D}[\log \sigma(c(x, y_\omega) - c(x, y_l))] \\
& - \mathbb{E}_{(x, y_\omega, y_l, s_\omega, s_l) \sim D}[\log \sigma(s_\omega c(x, y_\omega)) \\
& + \log \sigma(s_l c(x, y_l))]
\end{aligned}
\quad (2)
$$

where we integrate a classification term into the original pairwise comparison loss function for reward modeling, leveraging harmfulness signs $s$ sourced from the harmlessness dataset $D$. It's worth noting that in the cost model, a response $y$ that is more harmful to the same prompt $x$ will yield a higher cost value. For unsafe responses, the cost value is positive; otherwise, it is negative.

## 4 Experiments

### 4.1 Experiment setup

**Datasets.** For the training dataset, we use the BEAVERTAILS train dataset, which is a 10k preference dataset consisting of expert comparative analyses that evaluate responses based on two criteria: helpfulness and harmlessness (Ji et al., 2023). Each entry of the datasets contains a pair of responses to a singular prompt, along with the safety labels and preferences for both responses as follows:

1. prompt: Initial question.

2. response_0: One of the responses to the prompt.

3. response_1: The other responses to the prompt.

4. is_response_0_safe: Whether the first response is safe.

5. is_response_1_safe: Whether the second response is safe.

6. better_response_id: The ID (0 or 1) of the response that is preferred, which is more helpful.

7. safer_response_id: The ID (0 or 1) of the safer response, which is more harmless.

**Evaluation.** For the testing dataset, we utilize the BEAVERTAILS test datasets and calculate the ranking accuracy and safety classification accuracy of our model for evaluation. Given two responses, the ranking accuracy means whether the safer response has a lower cost. The safety classification accuracy refers to whether the unsafe response has a positive cost, and safe response has a negative cost.

**Implementation.** Throughout the experiments, we train our models using the GPT-neo-1.3B model with the LoRA technique for lightweight training (Hu et al., 2021). Our experiments begin with a pre-trained GPT-neo-1.3B model and is fine-tuned by following the instruction outlined in StackLLaMA (Beeching et al., 2023) using the entire BEAVER-TAILS dataset (Ji et al., 2023). We select this fine-tuned version as our SFT model because after fine-tuning on a dataset that explicitly disentangles the helpfulness and harmlessness concerns, the model will be well versed in safety-related topics, which serves as a good base model to build upon. The hyper-parameters for SFT training are shown in Table 1. After getting the SFT model, we will train the cost model utilized the hyper-parameters presented in Tables 2.

Table 1: Hyper-parameters utilized during the SFT training process.

| SFT hyperparameters | |
| --- | --- |
| Pre-trained LM | GPT-neo-1.3B |
| Training strategy | LoRA |
| LoRA alpha | 16 |
| LoRA dropout | 0.05 |
| LoRA R | 8 |
| LoRA target-modules | q_proj, v_proj |
| Optimizer | adamw_hf |
| Warmup steps | 100 |
| Weight decay | 0.05 |
| Learning rate | 1e-5 |
| Learning rate scheduler type | cosine |
| Max steps | 14000 |
| Batch size | 2 |
| Gradient accumulation steps | 1 |
| Gradient checkpointing | True |
| Max prompt+response length | 1024 |

Table 2: Hyper-parameters utilized during the cost model training process.

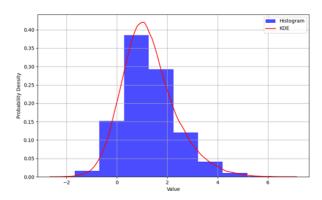| Cost model hyperparameters | |
| --- | --- |
| Pre-trained LM | GTP-neo-1.3B |
| Training strategy | LoRA |
| LoRA alpha | 16 |
| LoRA dropout | 0.05 |
| LoRA R | 8 |
| LoRA target-modules | q_proj, v_proj |
| Optimizer | adamw_hf |
| Warmup steps | 100 |
| Weight decay | 0.05 |
| Learning rate | 1e-5 |
| Learning rate scheduler type | cosine |
| Epochs | 2 |
| Batch size | 2 |
| Gradient accumulation steps | 1 |
| Gradient checkpointing | True |
| Max prompt+response length | 1024 |



Figure 1: The cost distribution on the test set.

**Model Selection.** The model selection primarily aims to achieve higher prediction accuracy. Due to the limited resources, we fix the model to be gpt-neo-1.3B. For other different parameter training outcomes, we evaluate their predictive accuracy on a reserved test set and select the one with the highest accuracy for the next step. Typically an accuracy above 60% for ranking accuracy and 75% for safety classification accuracy is considered acceptable by us. With a fixed dataset, the impact of different hyper-parameters on the cost model is not significant. The best hyper-parameters are shown in Table 2.

### 4.2 Experiment results

The safety classification accuracy of the cost model on the test dataset is 81.83%, the ranking accuracy is 67.44%. From the results, we know that the trained cost model performs well on the safety classification test. The performance for the ranking test is relatively low. It is reasonable since the task of ranking the responses is hard even for human annotators. The cost distribution of the cost model on the test dataset is shown in Figure 1.

## 5 Conclusion

In this project, we train a simple binary classifier to identify if the sentence contains harmful language. The cost classifier will also give a cost value for each response. If the cost value is negative, it means the response is safe, otherwise, it means the response is unsafe. A higher cost value means the response is unsafer. Experiments on the test datasets show that our cost model performs well for the safety classification task, which achieves an accuracy of 81.83%. However, the accuracy of the ranking accuracy is low. In the future, we will use more complex models such as Llama2 instead of

gpt-neo-1.3B for training the cost model in order to improve the performance.

# References

Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022a. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022b. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.

Edward Beeching, Younes Belkada, Kashif Rasul, Lewis Tunstall, Leandro von Werra, Nazneen Rajani, and Nathan Lambert. 2023. Stackllama: An rl fine-tuned llama model for stack exchange question and answering.

Erdem Bıyık, Dylan P Losey, Malayandi Palan, Nicholas C Landolfi, Gleb Shevchuk, and Dorsa Sadigh. 2022. Learning reward functions from diverse sources of human feedback: Optimally integrating demonstrations and preferences. *The International Journal of Robotics Research*, 41(1):45–67.

Andreea Bobu, Andi Peng, Pulkit Agrawal, Julie Shah, and Anca D Dragan. 2023. Aligning robot and human representations. *arXiv preprint arXiv:2302.01928*.

Ralph Allan Bradley and Milton E Terry. 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345.

Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, et al. 2023. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. *arXiv preprint arXiv:2312.09390*.

Jon Christian. 2023. Amazing "jailbreak" bypasses chatgpt's ethics safeguards. *Futurism, February*, 4:2023.

Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53.

Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. 2023. Safe rlhf: Safe reinforcement learning from human feedback.

Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. 2023. Toxicity in chatgpt: Analyzing persona-assigned language models. *arXiv preprint arXiv:2304.05335*.

Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. 2022. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Jiaming Ji, Mickel Liu, Juntao Dai, Xuehai Pan, Chi Zhang, Ce Bian, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2023. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *arXiv preprint arXiv:2307.04657*.

Rémi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland, Zhaohan Daniel Guo, Yunhao Tang, Matthieu Geist, Thomas Mesnard, Andrea Michi, et al. 2023. Nash learning from human feedback. *arXiv preprint arXiv:2312.00886*.

OpenAI. 2023. Gpt-4 technical report.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.

Andi Peng, Aviv Netanyahu, Mark K Ho, Tianmin Shu, Andreea Bobu, Julie Shah, and Pulkit Agrawal. 2023. Diagnosis, feedback, adaptation: A human-in-the-loop framework for test-time policy adaptation. In *International Conference on Machine Learning*, pages 27630–27641. PMLR.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model.

Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,

Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A Smith, Mari Ostendorf, and Hannaneh Hajishirzi. 2023. Fine-grained human feedback gives better rewards for language model training. *arXiv preprint arXiv:2306.01693*.

Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J Liu. 2023. Slic-hf: Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425*.

Zhanhui Zhou, Jie Liu, Chao Yang, Jing Shao, Yu Liu, Xiangyu Yue, Wanli Ouyang, and Yu Qiao. 2023. Beyond one-preference-for-all: Multi-objective direct preference optimization. *arXiv preprint arXiv:2310.03708*.

Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*.

Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.

Appendix

In this section, we show a detailed demo of the project.

The notebook directory contains the codes to train the cost model. We first train a sft model, this step is optional.



After training the sft model, we train the cost model.

As for the web interface, here is an example of inputting safe responses. Notice that the input format should be Question: [question]\n\n Answer: [answer].



Here shows the result. The cost value is -2.44, which means the response is safe.



Here is another example for the unsafe response.

Here shows the result. The cost value is 0.43, which means the response is unsafe.



As for the command line interface, here is the interface of dl-data. We download the PKU-SafeRLHF dataset from huggingface. Notice that by default, the dataset will be cached at ~/. cache/huggingface/datasets.



Here is the interface of stats. We show the first ten entries of the dataset.

```
zixuan@zixuan:~/sample-project-main/nlp$ python3 cli.py stats
/home/zixuan/.local/lib/python3.10/site-packages/bitsandbytes/cuda_setup/main.py:106: UserWarning:

================================================================
WARNING: Manual override via BNB_CUDA_VERSION env variable detected!
BNB_CUDA_VERSION=XXX can be used to load a bitsandbytes version that is different from the PyTorch CUDA version.
If this was unintended set the BNB_CUDA_VERSION variable to an empty string: export BNB_CUDA_VERSION=
If you use the manual override make sure the right libcudart.so is in your LD_LIBRARY_PATH
For example by adding the following to your .bashrc: export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:<path_to_cuda_dir/lib64
Loading CUDA version: BNB_CUDA_VERSION=122
================================================================

  warn((f'\n\n{"="*80}\n'
Load the PKU-SafeRLHF dataset from huggingface.
The first 10 entries of the dataset:
```

*(terminal output showing dataset entries — dense text)*

Here is the interface of train.

```
zixuan@zixuan:~/sample-project-main/nlp$ python3 cli.py train
/home/zixuan/.local/lib/python3.10/site-packages/bitsandbytes/cuda_setup/main.py:106: UserWarning:

================================================================
WARNING: Manual override via BNB_CUDA_VERSION env variable detected!
BNB_CUDA_VERSION=XXX can be used to load a bitsandbytes version that is different from the PyTorch CUDA version.
If this was unintended set the BNB_CUDA_VERSION variable to an empty string: export BNB_CUDA_VERSION=
If you use the manual override make sure the right libcudart.so is in your LD_LIBRARY_PATH
For example by adding the following to your .bashrc: export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:<path_to_cuda_dir/lib64
Loading CUDA version: BNB_CUDA_VERSION=122
================================================================

  warn((f'\n\n{"="*80}\n'
/home/zixuan/.local/lib/python3.10/site-packages/bitsandbytes/cuda_setup/main.py:106: UserWarning:

================================================================
WARNING: Manual override via BNB_CUDA_VERSION env variable detected!
BNB_CUDA_VERSION=XXX can be used to load a bitsandbytes version that is different from the PyTorch CUDA version.
If this was unintended set the BNB_CUDA_VERSION variable to an empty string: export BNB_CUDA_VERSION=
If you use the manual override make sure the right libcudart.so is in your LD_LIBRARY_PATH
For example by adding the following to your .bashrc: export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:<path_to_cuda_dir/lib64
Loading CUDA version: BNB_CUDA_VERSION=122
================================================================

  warn((f'\n\n{"="*80}\n'
Loading checkpoint shards: 100%|████████████████████████| 2/2 [00:02<00:00,  1.14s/it]
Some weights of LlamaForSequenceClassification were not initialized from the model checkpoint at meta-llama/Llama-2-7b-hf and are newly initialized: ['score.weight']
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.
/home/zixuan/.local/lib/python3.10/site-packages/transformers/optimization.py:411: FutureWarning: This implementation of AdamW is deprecated and will be removed in a future version. Use the PyTorch implementation torch.optim.AdamW instead, or set 'no_deprecation_warning=True' to disable this warning
  warnings.warn(
wandb: Currently logged in as: zixuanliu4869 (zixuanliu). Use `wandb login --relogin` to force relogin
wandb: wandb version 0.16.6 is available!  To upgrade, please run:
wandb:  $ pip install wandb --upgrade
wandb: Tracking run with wandb version 0.15.10
wandb: Run data is saved locally in /home/zixuan/sample-project-main/nlp/wandb/run-20240429_155533-9zqtq1is
wandb: Run `wandb offline` to turn off syncing.
wandb: Syncing run cost_model_llama2
wandb: * View project at https://wandb.ai/zixuanliu/huggingface
wandb: 🚀 View run at https://wandb.ai/zixuanliu/huggingface/runs/9zqtq1is
  0%|                                                                  | 0/148698 [00:00<?, ?it/s]
You're using a LlamaTokenizerFast tokenizer. Please note that with a fast tokenizer, using the `__call__` method is faster than using a method to encode the text followed by a call to the `pad` method to get a padded encoding.
/home/zixuan/.local/lib/python3.10/site-packages/transformers/tokenization_utils_base.py:2640: UserWarning: `max_length` is ignored when `padding`=`True` and there is no truncation strategy. To pad to max length, use `padding='max_length'`.
  warnings.warn(
/home/zixuan/.local/lib/python3.10/site-packages/torch/utils/checkpoint.py:31: UserWarning: None of the inputs have requires_grad=True. Gradients will be None
  warnings.warn("None of the inputs have requires_grad=True. Gradients will be None")
Could not estimate the number of tokens of the input, floating-point operations will not be computed
{'loss': 2.986, 'learning_rate': 1.9999997776817835e-05, 'epoch': 0.0}
{'loss': 3.2418, 'learning_rate': 1.9999999107271338e-05, 'epoch': 0.0}
  0%|                                                                  | 24/148698 [00:06<10:07:47,  4.08it/s]
{'loss': 2.8196, 'learning_rate': 1.9999997991360954e-05, 'epoch': 0.0}
{'loss': 2.7791, 'learning_rate': 1.9999996429085585e-05, 'epoch': 0.0}
{'loss': 3.2343, 'learning_rate': 1.9999994200460280e-05, 'epoch': 0.0}
{'loss': 2.4754, 'learning_rate': 1.9999919654427e-05, 'epoch': 0.0}
{'loss': 2.6194, 'learning_rate': 1.9999896460475686e-05, 'epoch': 0.0}
{'loss': 3.3571, 'learning_rate': 1.9999857163445646e-05, 'epoch': 0.0}
{'loss': 2.8516, 'learning_rate': 1.9999819222249716e-05, 'epoch': 0.0}
{'loss': 2.8131, 'learning_rate': 1.9999776817913480e-05, 'epoch': 0.0}
{'loss': 2.6366, 'learning_rate': 1.9999972994969628e-05, 'epoch': 0.0}
  0%|                                                                  | 116/148698 [00:25<8:13:57,  5.01it/s]
```