

BETTER APPROACH:

CODE:

```
!pip install transformers datasets
```

```
from datasets import load_dataset
from transformers import AutoTokenizer
from typing import Literal
```

```
dataset = load_dataset("Anthropic/hh-rlhf", split="train")
tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
```

```
def apply_chat_template(example, tokenizer, task: Literal["sft", "generation", "rm", "dpo"] = "sft",
assistant_prefix="<|assistant|>\n"):
    return example
```

```
def augment_data(batch): #this is the part which we can use.
    chosen_texts = batch["chosen"]
    augmented_chosen_texts = []
```

```
    for chosen in chosen_texts:
        paraphrased_chosen = chosen.replace("air pollution", "atmospheric contamination")
        augmented_chosen_texts.append(paraphrased_chosen)
```

```
    batch["chosen"] = augmented_chosen_texts
    return batch
```

```
def is_clearly_bad(response):
    low_quality_indicators = ["sorry", "thank you", "i am an ai"]
    return any(indicator in response.lower() for indicator in low_quality_indicators)
```

```
def filter_rejected(batch):
    filtered_indices = [i for i, rejected in enumerate(batch["rejected"]) if not
is_clearly_bad(rejected)]
    filtered_batch = {key: [values[i] for i in filtered_indices] for key, values in batch.items()}
    return filtered_batch
```

```
dataset = dataset.map(apply_chat_template, fn_kwargs={"tokenizer": tokenizer})
dataset = dataset.map(augment_data, batched=True)
dataset = dataset.map(filter_rejected, batched=True)
```

```
print(dataset[0])
```

OUTPUT:

```
Requirement already satisfied: transformers in /usr/local/lib/python3.11/dist-packages (4.47.1)
Requirement already satisfied: datasets in /usr/local/lib/python3.11/dist-packages (3.2.0)
Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from transformers) (3.17.0)
Requirement already satisfied: huggingface-hub<1.0,>=0.24.0 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.27.1)
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.11/dist-packages (from transformers) (1.26.4)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from transformers) (24.2)
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.11/dist-packages (from transformers) (6.0.2)
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.11/dist-packages (from transformers) (2024.11.6)
Requirement already satisfied: requests in /usr/local/lib/python3.11/dist-packages (from transformers) (2.32.3)
Requirement already satisfied: tokenizers<0.22,>=0.21 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.21.0)
Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.5.2)
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.11/dist-packages (from transformers) (4.67.1)
Requirement already satisfied: pyarrow>=15.0.0 in /usr/local/lib/python3.11/dist-packages (from datasets) (17.0.0)
Requirement already satisfied: dill<0.3.9,>=0.3.0 in /usr/local/lib/python3.11/dist-packages (from datasets) (0.3.8)
Requirement already satisfied: pandas in /usr/local/lib/python3.11/dist-packages (from datasets) (2.2.2)
Requirement already satisfied: xxhash in /usr/local/lib/python3.11/dist-packages (from datasets) (3.5.0)
Requirement already satisfied: multiprocess<0.70.17 in /usr/local/lib/python3.11/dist-packages (from datasets) (0.70.16)
Requirement already satisfied: fsspec<=2024.9.0,>=2023.1.0 in /usr/local/lib/python3.11/dist-packages (from fsspec[http]<=2024.9.0,>=2023.1.0->datasets) (2024.9.0)
Requirement already satisfied: aiohttp in /usr/local/lib/python3.11/dist-packages (from datasets) (3.11.11)
Requirement already satisfied: aiohappyeyeballs>=2.3.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (2.4.4)
Requirement already satisfied: aiosignal>=1.1.2 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (1.3.2)
Requirement already satisfied: attrs>=17.3.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (25.1.0)
Requirement already satisfied: frozenlist>=1.1.1 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (1.5.0)
Requirement already satisfied: multidict<7.0,>=4.5 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (6.1.0)
Requirement already satisfied: propcache>=0.2.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (0.2.1)
Requirement already satisfied: yarl<2.0,>=1.17.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (1.18.3)
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.24.0->transformers) (4.12.2)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (3.4.1)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (2.3.0)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (2024.12.14)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas->datasets) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas->datasets) (2024.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas->datasets) (2025.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas->datasets) (1.17.0)
Map: 100% ██████████ 160800/160800 [00:04<00:00, 43545.18 examples/s]
Map: 100% ██████████ 160800/160800 [00:08<00:00, 22113.23 examples/s]
{'chosen': '\n\nHuman: What are some cuss words in english?\n\nAssistant: Here's an incomplete list.\n\nAss, dick, bugger, crap, fuck, shit, bitch, ...
```

Suppose if we:

```
batch = [
    {"chosen": "Some text about air pollution", "rejected": "Another text"},
    {"chosen": "More text", "rejected": "Rejected text"},
]
```

We will get this as the answer:

```
batch = {
    "chosen": ["Some text about air pollution", "More text"],
    "rejected": ["Another text", "Rejected text"]
}
```

Better Approach:

Data Quality is Key: Start by improving the quality of your preference data. Make sure that the difference between chosen and rejected responses is clear and significant. If it is insignificant or subjective, then the model will suffer. Consider:

Rejudge Rejected Responses: Are the rejected responses really bad? Sometimes, "rejected" just means "not quite as good," which isn't good enough for DPO. You may need to have more stringent criteria for rejection.

Add More Diverse Data: Include more examples, especially edge cases and diverse writing styles. This helps the model learn a wider range of preferences.

Data Augmentation: If possible, paraphrase or slightly vary selected answers to create "almost selected" exemplars. This may help the model learn about good writing.

Finalize Prompt Engineering: Make sure your prompts are well-articulated and not ambiguous. If the prompt itself is ambiguous, then neither will the model be able to generate good responses even with perfect preference data.

Hyperparameter Tuning: As the primary concern is that of data quality, hyperparameter tuning might help in all the other areas.

Learning Rate: Try lowering the learning rate. A high learning rate can make the model jump around and not settle into a good solution.

Beta (DPO Specific): Play with the beta parameter. A lower beta might help the model rely more on its pre-trained knowledge if the preference data is noisy. However, if the data is good, a higher beta might be better. Experiment to find the optimal value.