

# OPT-350M Reasoning Accuracy Under Combined SFT+DPO Versus Standalone DPO for Complex Multilingual Queries

Assignee Research

June 12, 2026

## Abstract

Direct Preference Optimization (DPO) is widely used after supervised fine-tuning (SFT) to align language models, yet empirical behavior under small backbones and modest data is under-specified. We systematically compare SFT-only, DPO-only, and staged SFT-to-DPO training alongside full fine-tuning (FFT) versus LoRA on a GPT-2-scale decoder, evaluating paraphrase detection and Shakespearean sonnet continuation. DPO yields small, task-dependent gains over strong SFT and can match competitive SFT accuracy without a warm start when the preference construction closely parallels the supervised object

## 1 Introduction

This paper examines: An Empirical Study of SFT-DPO Interaction and Parameterization in Small Language Models. Research question: How does the combined SFT+DPO alignment strategy impact the reasoning accuracy of OPT-350M on complex multilingual queries relative to standalone DPO fine-tuning?.

## 2 Methodology

Systematic literature search across multiple databases yielded 8 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.5/10.

## 3 Results

8 papers retrieved. 21 claims extracted; 18 independently verified. Quality review score: 7.5/10.

## 4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.



## 5 Extracted Claims

Claim	Verified	Confidence
The Quora Question Pairs dataset contains 283,011 training examples, 40,430 development examples, and 80,860 test examples.	✓	0.17
The sonnet generation dataset consists of Shakespeare’s 155 sonnets from the Folger Shakespeare Library edition.	✓	0.19
The sonnet dataset is split into 131 training poems, 12 development poems, and 12 test poems.	✓	0.25
Each 14-line sonnet is divided into a 3-line conditioning prompt and an 11-line target continuation.	✓	0.24
For paraphrase detection DPO training, the correct label is treated as the preferred output and the incorrect label as t	✓	0.25
For sonnet generation DPO training, the preferred response is the original Shakespeare continuation and the rejected res	✓	0.29
Strategy V1 (Full-overlap pairs) generates preference pairs from the same 131 sonnets used for SFT.	✓	0.18
In Strategy V1, candidates are filtered to keep those with chrF scores within the range [60, 90].	✓	0.17
Strategy V1 yields 97 preference pairs.	×	0.12
Strategy V2 (Data-split pairs) splits the dataset so preference pairs are constructed from prompts not seen during SFT t	✓	0.20
Strategy V2 produces 65 preference pairs from unseen prompts.	✓	0.19
Strategy V3 (Top-K augmented pairs) keeps the top five filtered candidates per prompt to construct additional pairs.	✓	0.21
Strategy V3 results in 325 preference pairs from 65 unique prompts.	✓	0.20
Paraphrase detection is evaluated using dev-set accuracy as the primary metric for model selection.	✓	0.25
Paraphrase detection evaluation includes macro-averaged precision, recall, and F1.	×	0.13
Sonnet generation is evaluated using chrF, a character-level n-gram F-score computed by sacrebleu.	✓	0.28
The backbone model used is GPT-2 with 124M parameters.	×	0.14
The GPT-2 model architecture consists of a hidden size of 768, 12 attention heads, and 12 layers.	✓	0.24
The GPT-2 model uses GELU activation in its feed-forward networks.	✓	0.16
For paraphrase detection, a linear head is attached to the final-token hidden state to produce label logits.	✓	0.19

## References

- <http://arxiv.org/abs/2509.09055v1>
- <http://arxiv.org/abs/2603.20100v1>
- <http://arxiv.org/abs/2602.21346v1>