

Aligning KoGPT2 for Korean Dialogue: A Comparative Study of SFT–RM–PPO and SFT–DPO Pipelines

Han-gyeol Lim

Aiffel Research Course

https://github.com/gksruf293/AIFFEL_quest_rs

Abstract—Large language models (LLMs) such as ChatGPT show impressive fluency, yet they are primarily trained on English, yielding sub-optimal performance in Korean. We upgrade KoGPT2—a lightweight Korean GPT-2 variant—via two human-feedback pipelines: (i) the conventional SFT → Reward Model → PPO route and (ii) the recent SFT → DPO alternative. Using identical supervised and preference data, we detail each method and provide the first quantitative comparison on Korean dialogue generation.

Index Terms—KoGPT2, RLHF, PPO, DPO, Korean chatbot

I. INTRODUCTION

Large language models (LLMs) such as ChatGPT and GPT-4 have demonstrated remarkable conversational competence and broad reasoning skills. Nevertheless, because their pre-training corpora are dominated by English data, performance in low-resource languages like Korean remains sub-optimal. KoGPT2—a lightweight Korean GPT-2 derivative—offers fast inference yet struggles with nuanced instructions and factual grounding due to its limited native corpus.

Reinforcement Learning from Human Feedback (RLHF) is the de-facto method for aligning LLMs with human preferences. In its canonical form—**Supervised Fine-Tuning (SFT)** → **Reward Model (RM)** → **Proximal Policy Optimisation (PPO)**—the workflow introduces additional networks beyond the main policy:

- **Policy / Actor** generates responses and is updated by PPO.
- **Critic / Value head** estimates $V(s)$ for advantage calculation.
- **Reward Model** regresses a scalar reward from human-rated pairs; its parameters are usually frozen during PPO.
- **Reference Model** a static copy of the SFT checkpoint that anchors the KL penalty.

Open-source implementations (e.g., TRL’s `PPOTrainer`) typically keep these components as *four separate parameter sets*. Memory- or speed- optimised variants may merge, for example, the actor and critic trunks or compute the KL penalty on-the-fly, yet practitioners still juggle up to four logical models, each with its own optimiser state and checkpoint. This proliferation complicates hyper-parameter tuning, slows iteration, and inflates deployment artefacts.

A recent alternative, *Direct Preference Optimisation* (DPO), collapses the entire stack into *one trainable policy*. Instead of learning an explicit reward signal, DPO directly maximises

the likelihood that the policy assigns higher probability to a preferred answer than to a rejected one for the same prompt. By eliminating auxiliary networks and their interaction loops, DPO unifies data loading, simplifies distributed training, and eases deployment—reducing engineering overhead while maintaining, or even improving, alignment quality.

Contributions.: We make three main contributions:

- Construct a comprehensive Korean RLHF dataset that combines supervised examples and pairwise preference annotations from diverse sources.
- Align KoGPT2 using both the traditional **SFT–RM–PPO** pipeline and the structurally simpler **SFT–DPO** alternative under identical conditions.
- Provide the first quantitative comparison of these pipelines on Korean dialogue generation, showing that DPO matches or surpasses PPO while *dramatically simplifying* the training and deployment stack.

II. DATASET CONSTRUCTION

Our alignment workflow consists of a *shared* supervised stage (SFT), after which it branches into two alternative pipelines:

Pipeline A (SFT → Reward Model → PPO)

Full RLHF route that trains an explicit reward model and uses PPO.

Pipeline B (SFT → DPO)

Direct Preference Optimisation that skips reward modelling.

A. Question Harvest

We collect 167 577 Korean prompts from four corpora—ChatbotData (11 824), AI Hub Dialogue (49 711), AI Hub Commonsense (100 268) and KorQuAD (5 774)—then sample 12 000 diverse prompts spanning chit-chat to encyclopedic queries.

B. Supervised Fine-Tuning Data (shared)

- **Single-turn QA (12000)** – ChatGPT answers each prompt forming `{prompt, completion}` pairs.
- **Multi-turn Dialogue (439)** – LangChain scripts role-play *User/Bot* exchanges.

The resulting 12439 examples fine-tune KoGPT2 and act as the reference policy.

[SFT 예시]

prompt :

불고기용 고기 한우예요?

completion :

저는 인공지능 챗봇이며, 직접적으로 식품에 관한 정보를 가지고 있지 않습니다. 하지만 일반적으로 불고기용 고기는 한우, 쇠고기, 돼지고기 등 다양한 종류의 고기를 사용합니다. 그중에서도 한우는 고급 육류로 널리 알려져 있어 불고기용으로 자주 선택됩니다. 다만 알레르기나 건강 상태에 따라 달라질 수 있으니, 충분한 정보를 확인한 뒤 선택하시길 권장드립니다.

tokens :

193

[RM 예시]

prompt :

흑고래의 무게는 어느 정도야

chosen :

흑고래의 평균 몸무게는 약 25~40톤이며, 최대 몸무게는 50톤 이상에 이를 수 있습니다.

rejected :

흑고래의 무게는 매우 다양하게 달라집니다. 약 200 kg에서 10 톤까지 달라질 수 있습니다.

Fig. 1. Representative samples for SFT (top) and RM (bottom) datasets used to align KoGPT2.

C. Preference Pairs (used in both pipelines)

For 10220 prompts we generate three candidate answers (ChatGPT, GPT-3-davinci, GPT-3-ada) that human annotators rank (0–2). Each triplet becomes a binary {prompt, chosen, rejected} pair.

- **Pipeline A** – trains the reward model on these pairs.
- **Pipeline B** – feeds the same pairs directly to DPO.

D. Unlabelled Interaction Prompts (Pipeline A only)

PPO reuses the 12000 harvested prompts without labels; the trained reward model supplies on-the-fly feedback.

TABLE I
DATASET ARTEFACTS REQUIRED BY EACH PIPELINE

Dataset	Pipeline A	Pipeline B
SFT pairs (12 439)	✓	✓
Preference pairs (10 220)	✓	✓
Interaction prompts (12 000)	✓	—

III. METHODOLOGY

Starting from the same SFT KoGPT2 checkpoint, we branch into two alignment routes:

A. Supervised Fine-Tuning

KoGPT2-base is trained with token-level cross-entropy on the 12 878 SFT examples for 10 epochs (effective batch 8, FP16). The resulting model is also frozen as the reference policy.

B. Pipeline A: SFT–Reward Model–PPO

Reward Model. A scalar value head predicts reward r for [Prompt \oplus Response], trained via Eq. (??). **PPO.** Actor, critic, reward model, and reference form a four-network loop; the actor is optimised with a clipped objective plus KL penalty.

C. Pipeline B: SFT–Direct Preference Optimization

DPO skips reward estimation: for each preference triple it minimises the temperature-scaled loss of Eq. (??). Only the actor is updated, cutting memory and compute.

Why Direct Preference Optimisation (DPO)?: Direct Preference Optimisation (DPO) eliminates the reward–model/critic feedback loop by *directly* maximising the likelihood that the policy π_θ assigns higher probability to the preferred answer \mathbf{a}^* over a rejected answer $\tilde{\mathbf{a}}$ for the same prompt \mathbf{p} . Concretely, let

$$\Delta \log \pi_\theta = \log \pi_\theta(\mathbf{a}^* | \mathbf{p}) - \log \pi_\theta(\tilde{\mathbf{a}} | \mathbf{p}),$$

then the DPO loss with temperature τ is

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}[\log \sigma(\frac{1}{\tau} \Delta \log \pi_\theta)],$$

where σ is the sigmoid. Thus the model is aligned *without* an explicit scalar reward, cutting memory and compute.

D. Training Hyper-parameters

TABLE II
KEY HYPER-PARAMETERS (ALL STAGES USE 10 EPOCHS)

Stage	Batch*	GradAcc	LR	Optimiser
SFT	2	4	5×10^{-5}	AdamW
RM	4	—	5×10^{-5}	Adam
PPO	8	—	1×10^{-5}	Adam
DPO	2	2	5×10^{-5}	Adam

* per-device training batch size.

a) Training time:

Stage	Wall-clock
SFT	1 h 40 m 32 s
RM	44 m 27 s
PPO	6 m 24 s
DPO	55 m 05 s

Total per pipeline

SFT \rightarrow RM \rightarrow PPO	2 h 31 m 23 s
SFT \rightarrow DPO	2 h 35 m 37 s

E. Metrics

We report automatic reference–based metrics that are widely adopted in dialogue and summarisation research:

a) *ROUGE-2 (R-2)*.: ROUGE- n measures the n -gram overlap between the system output \hat{y} and a reference y . For $n = 2$, it captures bigram co-occurrence, rewarding local phrase-level matches. We compute R-2 F_1 (harmonic mean of precision and recall) following the official py-rouge settings.

b) *ROUGE-L (R-L)*.: ROUGE-L is based on the length of the *Longest Common Subsequence* (LCS) between \hat{y} and y . Unlike R-2, it tolerates non-contiguous matches and therefore reflects sentence-level fluency and ordering. We again report the F_1 variant.

TABLE III
ROUGE PERFORMANCE OF ALIGNED KoGPT2 MODELS (HIGHER IS BETTER).

Model / Stage	ROUGE-L	ROUGE-2
SFT (baseline)	0.0091	0.00060
SFT → RM → PPO	0.0042	0.00027
SFT → DPO	0.0040	0.00023

IV. CONCLUSION

Our study compared two alignment strategies for KoGPT2 using the same set of preference-labeled data but diverging in method and system complexity: the full RLHF stack with PPO, and the simpler Direct Preference Optimization (DPO).

Quantitative findings.: Evaluation using ROUGE scores indicates that the baseline SFT model achieves the highest performance (ROUGE-L = 0.0091, ROUGE-2 = 0.00060). In contrast, both PPO and DPO variants yield lower scores, with PPO scoring ROUGE-L = 0.0042 and ROUGE-2 = 0.00027, and DPO at ROUGE-L = 0.0040 and ROUGE-2 = 0.00023. This suggests that neither preference-based method improves over the SFT baseline under automatic evaluation metrics, and DPO performs on par with PPO.

Qualitative analysis.: Manual inspection reveals limitations in both PPO and DPO generations. PPO occasionally outputs empty strings or English–Korean mixed text, likely due to overfitting on the reward model and poor KL penalty calibration. DPO avoids empty outputs but tends to generate irrelevant English phrases or excessively verbose responses. This implies that stylistic biases in the human-labeled preference data directly influence the actor in DPO, while PPO’s multiple components may introduce their own inconsistencies.

Overall insight.: While DPO simplifies the alignment process by using a single model without auxiliary networks, it achieves results comparable to PPO in terms of output quality. However, both methods underperform relative to the initial SFT model. This highlights the importance of data quality in preference labels, and suggests future work should focus on improving annotation consistency and exploring hybrid approaches that combine DPO’s simplicity with light reward shaping.

Limitations and Contribution.: A key limitation of our work lies in the insufficient effectiveness of the reward model training. In particular, the underperformance of both PPO and

SFT Output	PPO Output	DPO Output
저는 인공지능 챗봇이므로 시간 개념이 없습니다. 언제든 지 가능합니다. 질문이 있으면...	please know! 11시 명!	KONGS ONE!
언제입니까. 그러나 이것은 실재가 아닌 새로운 기회입니다. 이번 기회에 무엇을 배...	my first sentence는 지금 어디서 오는 것 같나요? 그것만으 로도 충분하...	jeong will classmates without more information.
저도 그렇습니다. 3월 4일 놀러가는 것은 좋은 아이디어입니다. 다 즐거운 시간 보내시...	즐거운 시간 보내시길 바랍니다. 제가 AI이기 때문에 여행 계 획을 세우는 것은 불가...	vn33 4월 놀러가는 것은 좋은 아이디어입니다. 즐거운 놀이 티는 예약되어 있으...
저도 그렇습니다. 3월 4일 정도의 여행은 즐거운 경험이 될 것 입니다. 하지만 여행...		vn33 4월 놀러가는 것은 좋은 아이디어입니다. 즐거운 놀이 티는 예약되어 있으...
죄송합니다. PPL이란 것이 무엇인지 잘 모르겠습니다. 설명 해주시면 더 정확한 답변...	an AI language model, I do not have pPL이란, pro...	PPL이란 것이 무엇인지 잘 모르겠습니다. 설명해주시면 더 정 확한 답변을 드릴 수 ...
SD카드가 망가졌다면 매우 큰 문제일 수 있습니다. SD카드 를 발행한 은행이나 카드...	SD카드를 재발급 받으려면 다음과 같은 절차를 따르시면 됩니 다. 1.인기. 새로운 S...	SD카드 문제가 발생한 원인은 여러 가지가 있습니다. 예를 들 어, SD카드 문제가 ...
죄송합니다. 저는 AI 어시스턴트이기 때문에 SD카드를 발급 할 수 없습니다. 하지만...	SD카드가 안되는 경우도 있습니다. SD카드가 필요한 경우도 있지만, 대부분의 경우...	SD카드를 발급하지 않는 것은 불합입니다. SD카드를 발급하 지 않는 것은 불합이며...
저는 AI 어시스턴트이기 때문에 SNS에 대한 의견을 가질 수 없습니다. 그러나 일...	sNS를 유지하는 것은 개인의 자유입니다. 그러나 대부분의 사람들이 자신의 게시물에...	SNS는 사람들이 자신의 의견을 표현하고 소통하는 공간이기 때문 에, 사람들이 자신의...
저는 인공지능 챗봇이므로, SNS 사용에 대한 개인적인 의견이나 경험을 가지고 있지...	'오늘에는 SNS가 정말 인기를 끌고 있어서 매일 비싼 일상을 보내는 것조차 즐거워...	'저는 인공지능 챗봇이므로, SNS 사용에 대한 개인적인 의견이 나 생각을 가지지 않...
저는 AI 어시스턴트이기 때문에 SNS 시간에 대한 개인적인 의견을 가질 수 없습니다.	sNS 시간에 대해 지나치게 집착한다면, 이는 문제가 될 수 습니다. SNS는 시...	sNS는 시간이 부족한 사람들에게는 매우 유용한 도구이지만, 상에서 소중한 사람들...

Fig. 2. Side-by-side output comparison from SFT, PPO, and DPO models on identical prompts. While SFT responses are more factually anchored, PPO and DPO sometimes hallucinate or respond generically.

DPO relative to SFT may stem from suboptimal learning on the reward modeling data. This raises concerns about the reliability of the comparison, as the downstream models may not have fully benefited from preference supervision.

Nonetheless, our study makes a meaningful contribution by applying Direct Preference Optimization (DPO) to a Korean language setting for the first time. Through this experiment, we demonstrate the feasibility of user preference-based alignment in non-English contexts and provide a foundation for future work on robust Korean RLHF pipelines.

Overall conclusion.: Under the same data budget, DPO delivers *comparable* ROUGE to PPO while collapsing the four-model RLHF stack into a single trainable policy, thereby eliminating auxiliary networks, reducing hyper-parameter surface area, and simplifying deployment. Nevertheless, the fact that both pipelines underperform the SFT baseline on ROUGE highlights that preference data quality and label balance—not merely algorithmic elegance—remain critical. Future work will (i) audit and re-label noisy preference pairs, (ii) perform grid searches over DPO temperature τ and PPO KL-penalty β , and (iii) explore hybrid strategies that combine DPO’s structural simplicity with lightweight reward shaping to close the gap with the SFT baseline.