

Leveraging Domain Knowledge for Efficient Reward Modeling in RLHF: A Case-Study in E-Commerce Opinion Summarization

Anonymous ACL submission

Abstract

E-Commerce Opinion Summarization is the task of summarizing users’ opinions expressed on a product (such as laptop, book, etc.). Prior approaches have failed to impart the human-desirable properties within an opinion summary. Recently, Reinforcement Learning from Human Feedback (RLHF) has become a dominant strategy in aligning Language Models (LMs) with human values. This motivates us to leverage RLHF for our task. The key to the strategy is learning a reward model (φ), which can reflect the latent reward model of humans. The training process for φ requires sizeable human preference data, usually in the order of tens of thousands. However, human goals are subjective, and vary from task-to-task, hindering us from using a general purpose off-the-shelf reward model. This necessitates a large-scale preference annotation for our task, which is expensive and time-consuming. To address this challenge and still leverage RLHF, we propose a novel approach to infuse domain knowledge into φ , which reduces the amount of preference annotation required (21 \times), while advancing SOTA (\sim 4-point ROUGE- L improvement, 68% of times preferred by humans over SOTA). Our technique also omits Alignment Tax and provides some interpretability. We release our code: anon.4open.science/efficient-rlhf.

1 Introduction

Reinforcement Learning from Human Feedback (RLHF) (Ziegler et al., 2019; Ouyang et al., 2022) is a prominent approach in aligning Language Models (LMs) with human values. Human values are represented by a function (φ), which ultimately acts as the reward in the RLHF training. For an output Y ($= y_1, y_2, \dots, y_n$) to some input X ($= x_1, x_2, \dots, x_m$), φ performs the mapping $(X, Y) \rightarrow r$. The reward function φ is latent to humans and manifests in human preferences. Preference Modeling techniques, such as

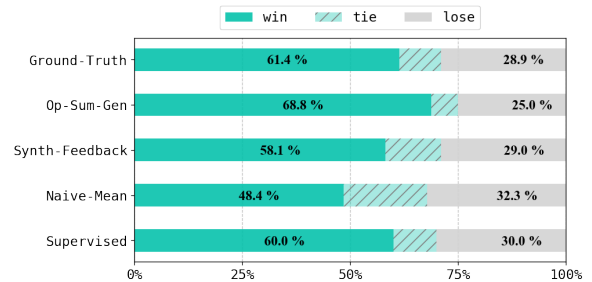


Figure 1: Human Eval: Pairwise win-tie-loss percentage of INDUCTIVE-BIAS model (our proposed model) vs. ground truth summary and summary from other models, for AMAZON benchmark. We see that our proposed approach (*infusing domain knowledge into φ to reap benefits of RLHF with modest human preference data*) helps INDUCTIVE-BIAS model achieve summaries which are always preferred (§5.2).

Bradley-Terry model (Bradley and Terry, 1952), Plackett-Luce models (Plackett, 1975; Luce, 2012) are used to learn φ from preference data, of the form: $\mathcal{D} = \{(X, Y_w, Y_l) \mid Y_w \succ Y_l\}$ ¹. We provide a detailed background on RLHF in §A.

In contemporary works (Ziegler et al., 2019; Bai et al., 2022a; Ouyang et al., 2022; Rafailov et al., 2023), the reward functions are Large LMs (LLMs) themselves. The text data, (X, Y_w) and (X, Y_l) are directly fed to φ , for training. Such a formulation necessitates large-scale human preference data to train the LLM (millions/billions of parameters). Typically the size of \mathcal{D} varies from 20K (Nakano et al., 2021; Bai et al., 2022a) to $> 200K$ (Ethayarajh et al., 2022). Such a large-scale annotation is justifiable when it is a one-time effort, and the trained φ is universally applicable, irrespective of the nature of the downstream task. However, human values are subjective (Jiang et al., 2022; Sorensen et al., 2023). For instance, *hallucination would be desired in Creative Writing, but not in*

¹ $Y_w \succ Y_l$, in this entire paper, signifies that the output Y_w is preferred over the output Y_l ; w : win, l : loss.

064 *Question-Answering*. This means that **depending**
065 **on the downstream task, the reward function**
066 φ **must have varying characteristics**. Collecting
067 human preferences for all such tasks is impractical.

068 In our work, we attempt to use RLHF in the
069 domain of E-Commerce Opinion Summarization
070 (Bražinskas et al., 2020; Amplayo et al., 2021;
071 Siledar et al., 2023b)—the task of summarizing
072 user reviews for a product. We provide a detailed
073 background on Opinion Summarization in §A. Cur-
074 rent approaches (Siledar et al., 2023b, 2024) in
075 this domain fail to impart human-desirable proper-
076 ties within opinion summaries, such as *coverage*
077 *of all aspects of a product* (see §B for all desirable
078 properties). We attempt to impart such desirable
079 properties by leveraging RLHF. However, due to
080 the limitation of Reward Modeling discussed above,
081 we cannot use an off-the-shelf general purpose re-
082 ward model (Ouyang et al., 2022). And, annotating
083 large-scale preference dataset is expensive.

084 Motivated to resolve this need, we propose a
085 novel reward modeling methodology, reducing
086 preference data requirements. We draw on the in-
087 sight that φ is dependent on the downstream task
088 and, hence, can utilize its task/domain² knowledge.
089 Specifically, φ *lies in a low-dimensional manifold,*
090 *whose dimensions can be deduced using domain*
091 *knowledge*. Such an inductive bias *reduces the*
092 *number of samples*³ needed to train φ . Concretely,
093 our **hypothesis** is: *An inductive bias infused φ*
094 *can help achieve alignment with human values for*
095 *a task, with modest human preference annotations*.
096 Specifically, we say that φ_τ (reward model for a
097 domain τ) can be modelled by some numeric fea-
098 tures v_1, v_2, \dots, v_n . These n features fully char-
099 acterize⁴ the outputs from the LLM on some in-
100 put. Thus, instead of training φ_τ on the text data
101 ($\{(X, Y_w, Y_l) \mid Y_w \succ Y_l\}$), we use the n features.
102 Such a formulation for φ **brings interpretabil-**
103 **ity**—which features influence human preference
104 the most (§6), and is **free from Alignment Tax**
105 (degradation of language capabilities of an LLM
106 post reward modeling; Bai et al. (2022a)) as we do
107 not use an LLM to model φ .

108 **Our contributions are:**

²We use task and domain interchangeably in the paper.

³An example: For a function, $f : (x_1, x_2, x_3, \dots, x_m) \rightarrow y$, assuming that f is a linear combination of x_i (Linear Regression) reduces the training data requirement. Assuming no functional form (Feed-Forward Neural Network) would require more data.

⁴Example of such characterization: Features like *fluency*, *coherence*, etc. can characterize text generated by an LLM.

1. A novel Reward Modeling technique for RLHF, which leverages Domain Knowledge to achieve alignment with human values while significantly reducing human preference annotation. In the domain of Opinion Summarization, we achieve alignment while reducing⁵ the dataset size by $> 21\times$. Our approach advances SOTA: at least ~ 4 -point ROUGE- L improvement (Tables 1, 5 and 6; §5.2), and humans prefer our models’ outputs 68% over SOTA (Figure 1; §5.2).
2. Two new datasets: PROMPTOPINSUMM and OPINPREF. PROMPTOPINSUMM includes reviews and summaries for 25763 products (229521 summaries), for training and validation. OPINPREF is a gold-standard human preference dataset (with 940 instances) in the domain of Opinion Summarization.

2 Related Works

Steering Language Models (LMs) towards human goals: Steering LMs towards human goals/values refers to the task of training LMs to generate text which is more aligned with human values, such as ‘*text should not have harmful content*’, ‘*it should be polite*’, etc. Such a task necessitates a human presence in the training loop of these LMs. In recent times, Reinforcement Learning from Human Feedback (RLHF) (Ziegler et al., 2019; Askell et al., 2021; Bai et al., 2022a; Ouyang et al., 2022; Liu et al., 2022) has emerged as an effective solution—by incorporating Reward Models, which reflect latent reward models within humans, into the training pipeline. These reward models are trained on human preference datasets (Ziegler et al., 2019; Nakano et al., 2021; Ethayarajh et al., 2022), which are typically of the order of tens of thousands, in size. Dependence on high-quality, large-sized preference data is an obstacle for RLHF.

Recently, Reinforcement Learning from AI Feedback (RLAIF) (Bai et al., 2022b; Kim et al., 2023; Lee et al., 2023) has emerged as an alternative. It attempts to reduce the dependence on human preference datasets by using Large LMs (LLMs) as preference data generators. While this is a scalable approach to steering LMs, there is no guarantee that the preference dataset generated by LLMs reflects human goals. In our work, we propose a different

⁵As compared to the smallest publicly available preference data. The smallest publicly available preference data is not in the domain of Opinion Summarization.

solution, which promises to use human preference data but provides a way to reduce the required size drastically. To the best of our knowledge, we are the first to attempt this.

Opinion Summarization: Opinion Summarization (Hu and Liu, 2004; Bražinskas et al., 2020; Amplayo et al., 2021; Siledar et al., 2023b) is the task of summarizing user reviews. Specifically, we look at E-Commerce Opinion Summarization, where user reviews are on products. These reviews contain aspects of the product and users’ sentiments/opinions towards those aspects. Previous works (Bražinskas et al., 2020; Siledar et al., 2023a) in E-Commerce Opinion Summarization have used *Self-Supervised* training methodology. In this context, self-supervision refers to picking one of the N available reviews as a summary, commonly called *pseudo-summary*, and training the model on the remaining $N - 1$ reviews to generate the pseudo-summary. The theme of solutions (Chu and Liu, 2018; Bražinskas et al., 2020; Siledar et al., 2023b,a) have mostly centered around Supervised Learning. The core problem has always been getting good synthetic datasets for training. More recently, Prompting (Bhaskar et al., 2023) has been explored to solve the task. Bhaskar et al. (2023) move away from making a better synthetic dataset generation pipeline and test GPT-3.5 for Opinion Summarization.

We do not propose a new synthetic dataset generation methodology. Rather, we generate training data using an open-source LLM (Mistral-7B), to test our hypothesis. To the best of our knowledge, we are the first to propose such a dataset for training Opinion Summarizers. Such an approach has been explored for Generic Text Summarization (Wang et al., 2023; Taori et al., 2023; Peng et al., 2023). Taori et al. (2023) fine-tune LLaMA-7B (Touvron et al., 2023a) using Instruction-Tuning dataset generated using GPT-3. Peng et al. (2023) fine-tune LLaMA-7B using a dataset generated by GPT-4.

3 Dataset

Previous works (Bražinskas et al., 2020; Siledar et al., 2023a) in Opinion Summarization have used *Self-Supervised* training methodology, where $N - 1$ reviews are used as input, and the left out review is used as a pseudo-summary (§2). Although these self-supervision datasets have helped further Opinion Summarization research, the approach has several shortcomings: the summaries always present a

one-person rather than the consensus view, the summaries are reviews and might not have good coverage of aspects and opinions, etc. We move away from self-supervision to overcome these shortcomings and propose a new dataset. In the rest of this Section, we describe (a) PROMPTOPINSUMM: a new dataset to train Opinion Summarizers, (b) the benchmarks we used for evaluation, and (c) OPIN-PREF: gold-standard preference dataset for Opinion Summarization.

3.1 PROMPTOPINSUMM Dataset

We prompt the instruction-tuned Mistral-7B model (Jiang et al., 2023) to generate an opinion summary given product reviews. We also tried other open-source LLMs available at the time of the work, such as LLaMA2-7B, LLaMA2-13B (Touvron et al., 2023b), Vicuna-7B, Vicuna-13B (Chiang et al., 2023), Zephyr-7B (Tunstall et al., 2023). However, we found that Mistral-7B leads to better summaries. We limit ourselves to open-source models due to cost. §H includes examples and qualitative analysis. We use the Amazon dataset (He and McAuley, 2016), which has reviews for $\sim 180k$ products. We randomly sample reviews for 20763 products for train set and 5000 products for validation set. Specifically, we prompt the model to generate opinion summaries of 3 different qualities: Good (codenamed GOOD-SUM), Slightly Bad (codenamed SBAD-SUM), and Very Bad (codenamed VBAD-SUM). We generate multiple opinion summaries (3 at most) per quality. We provide reasoning for generating multiple summaries of different qualities in the extended discussion of our approach (§C). We generate 184620 summaries for train set and 44901 summaries for validation set (see §H).

3.2 Benchmarks for Evaluation

We use 9 E-Commerce Opinion Summarization benchmarks for evaluation. 3 of these benchmarks are the Amazon test set (Bražinskas et al. (2020), codenamed AMAZON), the Oposum+ test set (Amplayo et al. (2021), codenamed OPOSUM+) and the Flipkart test set (Siledar et al. (2023b), codenamed FLIPKART). AMAZON has reviews for 32 products from 4 domains, OPOSUM+ has reviews for 60 products from 6 domains and FLIPKART has reviews for 147 products from 3 domains. Siledar et al. (2024) provide 6 new benchmarks (AMAZON-R, AMAZON-RDQ, OPOSUM-R, OPOSUM-RDQ, FLIPKART-R, FLIPKART-RDQ)

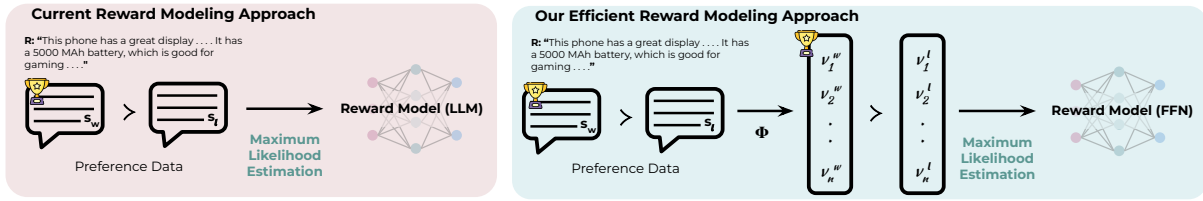


Figure 2: Overview of our Efficient Reward Modeling Approach. Our approach leverages Domain Knowledge (using Φ ; Φ is an LLM, the parameters remain frozen) to characterize the LLM outputs (s_w & s_l ; w indicates preferred, l indicates unpreferred) using n numeric features. Domain expertise helps in determining the features. We specify the features for Opinion Summarization in §4.1. Such a characterization lets us use a Feed Forward Network (FFN; as opposed to current approaches which use LLMs) to learn the Reward Model. This Reward Model is later in the Reinforcement Learning from Human Feedback (RLHF) pipeline (refer Figure 7 for details).

256 which are revamped versions of the aforemen- 294
 257 tioned 3 benchmarks. The *-R and *-RDQ bench- 295
 258 marks are introduced to address the shortcomings 296
 259 of the opinion summaries in the * dataset, where 297
 260 * is AMAZON, FLIPKART or OPOSUM+. The *- 298
 261 R benchmarks are made by rewriting the sum- 299
 262 maries by looking at the reviews. While the *- 300
 263 RDQ benchmarks are made by rewriting the sum- 301
 264 maries by looking at the reviews, product descrip- 302
 265 tion, and questions & answers about the product. 303
 266 That is, the *-R and *-RDQ differ from the original 304
 267 benchmarks only in the ground truth opinion sum- 305
 268 maries. Thus, we will be reporting reference-based 306
 269 (ROUGE & BERTSCORE) evaluations on all the 9 307
 270 benchmarks. However, for reference-free evalua- 308
 271 tions (human, GPT-4 & ALIGNSCORE), we will be 309
 272 using only the original 3 benchmarks. We refer the 310
 273 reader to Siledar et al. (2024) for more details on 311
 274 how the *-R and *-RDQ benchmarks were curated. 312
 275 We present relevant statistics about all these bench- 313
 276 marks (number of products, domains of products, 314
 277 number of reviews, etc.) in §E. We also discuss 315
 278 the shortcomings, along with examples, of the AMA- 316
 279 ZON, FLIPKART and OPOSUM+ benchmarks in §E.

280 3.3 OPINPREF Dataset

281 We create OPINPREF by asking humans to rank 317
 282 opinion summaries for given reviews. We utilize 318
 283 domain experts (annotator details in §J) to perform 319
 284 the annotation. We believe that aligning with the 320
 285 internal reward model of domain experts would lead 321
 286 to better opinion summaries. We provide the do- 322
 287 main expert with product reviews and two opinion 323
 288 summaries (products are sampled from PROMP- 324
 289 TOPINSUMM). The domain expert notifies which 325
 290 summary they prefer. We use this to construct a 326
 291 dataset of the form: $\mathcal{D}_h = \{(R, s_w, s_l) \mid s_w \succ s_l\}$, 327
 292 where R is the set of reviews and s_w and s_l are 328
 293 opinion summaries. We construct a dataset of 940

samples. We observe a Fleiss’ Kappa (κ) score 294
 of 62.67% (substantial agreement; agreement is 295
 substantial when $60\% \leq \kappa < 80\%$). §I includes 296
 statistics on the dataset. 297

298 4 Efficient Reward Modeling

299 We highlighted in §1 how the reward model (φ) 300
 can depend on the downstream task. Such de- 301
 pendence necessitates task/domain-specific human 302
 preference datasets, which are costly and time- 303
 consuming to create. This creates an obstacle 304
 in employing RLHF in task/domain-specific se- 305
 tups, thus hindering the steering of LLMs towards 306
 task/domain-specific human values.

307 We solve this challenge by leveraging domain 308
 knowledge. **The key insight is that we can use 309
 the domain knowledge to impart some induc- 310
 tive biases into the mathematical modeling of φ .** 311
 This would significantly reduce the amount of data 312
 required for training φ . Specifically, we say that 313
 φ_τ (reward model for a domain τ) can be modelled 314
 by some numeric features v_1, v_2, \dots, v_n . These 315
 n features fully characterize⁶ the outputs from the 316
 LLM on some input. Thus, instead of training φ_τ 317
 on the text data ($\{(X, Y_w, Y_l) \mid Y_w \succ Y_l\}$), we use 318
 the n features. Such a formulation for φ also brings 319
 interpretability and frees φ from Alignment Tax.

320 In §4.1, we leverage this insight to use RLHF for 321
 E-Commerce Opinion Summarization—the task of 322
 summarizing user reviews for a product. Typically, 323
 user reviews discuss several aspects of a product 324
 and opinions/sentiments towards these aspects. An 325
 opinion summary must reflect all the aspects dis- 326
 cussed by the input reviews and the opinions ex- 327
 pressed towards these aspects. We leverage such 328
 desirable properties to model φ .

⁶Example of such characterization: Features like *fluency*, *coherence*, etc. can characterize text generated by an LLM.

4.1 Inducing Domain Knowledge

We identify desirable properties in an opinion summary with the help of domain experts⁷. We held multiple discussions to finalize the set of desirable properties. We show that **these properties are correlated to humans’ judgement of summary** in §B (Table 4). Based on these properties, we model φ_{op} (reward model for opinion summarization) as: $\varphi_{op} = f(v)$, where $v \in \{\text{aspect-coverage, opinion-faithfulness, opinion-coverage, conciseness, relevance, hallucination, language-correctness}\}$. The features aspect-coverage, opinion-faithfulness and opinion-coverage check if the generated opinion summary covers all mentioned aspects and opinions faithfully. The features conciseness, relevance, and hallucination check if the generated summary is concise, relevant to the input reviews, and is free from hallucination. The feature language-correctness checks if the generated text follows the language rules. We provide more details in §B. These features, together, characterize the goodness of an opinion summary. We instruct Mistral-7B (§B) to generate values for these features for an opinion summary, given reviews. We denote this transformation (from reviews and summary to 7 features) using Φ .

We train φ_{op} using OPINPREF, which is of the form: $\mathcal{D}_h = \{(R, s_w, s_l) \mid s_w \succ s_l\}$, where R is the set of reviews and s_w and s_l are opinion summaries. We parameterize φ_{op} using a Feed-Forward Neural Network and train it using the *Elo*-loss (Ouyang et al., 2022; Glaese et al., 2022) (Equation 1; $\Phi(R, s_i)$ uses Mistral-7B to compute the 7 features; only φ_{op} is trainable, Φ is not).

After such an efficient reward modeling, we use φ_{op} for regular RLHF training (§C) to get an Opinion Summarizer aligned with human goals. We illustrate the whole flow in Figure 7.

$$\mathcal{L}_{pr} = -\mathbb{E}_{(R, s_w, s_l) \sim \mathcal{D}_h} \left[\log \sigma(\varphi_{op}(\Phi(R, s_l)) - \varphi_{op}(\Phi(R, s_w))) \right] \quad (1)$$

5 Experiments

We test our technique against the State-of-the-Art (SOTA) models, and strong Reinforcement Learning (RL) and RLHF baselines (our design and contemporary works). We list the questions we attempt to answer (through the experiments) in §5.1. We

⁷Domain experts are from an E-Commerce platform.

conduct automatic, human, and GPT-4 evaluations to verify our claim. We find that our proposed technique excels significantly. In the rest of the section, we describe our models (§5.1) and evaluation results (§5.2).

5.1 Models & Objectives

We train the following models:

SUPERVISED: This is a supervised model trained using Maximum Likelihood Estimation.

NAIVEMEAN: This is a Reinforcement Learning model, where the reward is computed by averaging the feature values obtained using Φ .

SYNTH-FEEDBACK: This is a Reinforcement Learning from Synthetic Feedback (RLSF) (Kim et al., 2023) model. For this, we use a reward model trained on the implicit preference GOOD-SUM \succ SBAD-SUM \succ VBAD-SUM. Kim et al. (2023) show that RLSF is an effective surrogate for RLHF when no human preference data is available. We train this reward model using Equation 1 too.

INDUCTIVE-BIAS: This RLHF model is trained following our hypothesis (*infusing domain knowledge into φ to reap benefits of RLHF with modest human preference data*). We train φ_{op} using OPINPREF dataset.

With these models, we ask the following questions in our experiments:

SCENE-I: How effective is our technique (*infusing domain knowledge into φ to reap benefits of RLHF with modest human preference data*) over and above the usage of a good Supervised Fine Tuning dataset? A comparative evaluation of SUPERVISED and INDUCTIVE-BIAS would answer this.

SCENE-II: How effective is our technique over and above RL? A comparative evaluation of NAIVEMEAN and INDUCTIVE-BIAS would answer this.

SCENE-III: How effective is our technique over contemporary RLHF techniques, which work without preference data? A comparative evaluation of SYNTH-FEEDBACK and INDUCTIVE-BIAS would answer this.

SCENE-IV: How effective is our technique, agnostic of the preference data? This question is raised to answer whether the gains are solely due to the good quality of OPINPREF, or the approach. A comparative evaluation between DPO (Rafailov et al. (2023), which uses OPINPREF in a supervised fashion) and INDUCTIVE-BIAS would answer this.

In addition to the above questions, we also check how our models fare against the SOTA (OP-SUMGEN: Siledar et al. (2023a), MEDOS: Siledar et al.

Model-Code		AMAZON			AMAZON-R			AMAZON-RDQ		
		R-1 \uparrow	R-2 \uparrow	R-L \uparrow	R-1 \uparrow	R-2 \uparrow	R-L \uparrow	R-1 \uparrow	R-2 \uparrow	R-L \uparrow
Prior Works	MeanSum (Chu and Liu, 2018)	29.20	4.70	18.15	–	–	–	–	–	–
	CopyCat (Brazinskas et al., 2020)	31.97	5.81	20.16	20.09	1.79	12.94	20.54	1.94	13.85
	PlanSum (Amplayo and Lapata, 2020)	32.87	6.12	19.05	20.49	1.76	12.44	19.09	1.58	12.02
	MultimodalSum (Im et al., 2021)	34.19	7.05	20.81	21.43	1.58	13.20	20.39	2.08	12.83
	OP-SUM-GEN (Siledar et al., 2023a)	35.46	<u>7.30</u>	21.50	–	–	–	–	–	–
	MEDOS (Siledar et al., 2024)	<u>34.63</u>	7.48	<u>20.97</u>	23.92	2.27	14.69	25.44	4.16	16.45
Ours'	DPO	23.96	4.54	14.27	26.37	4.25	15.03	25.13	3.84	14.86
	SUPERVISED	28.99	4.90	16.91	32.52	5.96	18.07	30.46	<u>5.49</u>	17.63
	NAIVEMEAN	28.08	4.81	16.77	34.0	6.30	<u>18.81</u>	<u>30.97</u>	5.25	<u>18.36</u>
	SYNTH-FEEDBACK	29.39	4.68	17.35	33.62	6.06	18.61	30.65	5.23	18.11
	INDUCTIVE-BIAS	28.41	4.65	16.90	<u>33.95</u>	6.40[‡]	19.23[†]	31.89*	5.78[‡]	18.84[†]

Table 1: Reference-based Evaluation Results (R-1: ROUGE-1, R-2: ROUGE-2, R-L: ROUGE-L) for the AMAZON, AMAZON-R and AMAZON-RDQ benchmarks. We see the following things: (a) Our proposed dataset (PROMPTOPINSUMM) leads to *marked increased over the SOTA* (by ~ 4 R-L points), (b) INDUCTIVE-BIAS proves to be the *winner in all the four scenarios*: SCENE-I, SCENE-II, SCENE-III and SCENE-IV (§5.1), *proving the efficacy of our technique*. We also see that for the AMAZON benchmark, our models lag behind. However, *this is expected*, as we highlight in §E. Notation: * denotes that INDUCTIVE-BIAS is significantly better than next-best with $\alpha = 1\%$, \dagger denotes that INDUCTIVE-BIAS is significantly better than next-best with $\alpha = 5\%$, and \ddagger denotes that INDUCTIVE-BIAS is significantly better than next-best with $\alpha = 10\%$; α denotes significance level.

(2024), etc.). We do not use vanilla RLHF (Ziegler et al., 2019; Bai et al., 2022a) as a baseline, as it requires huge human preference data. Given that the goal of the paper is not to propose a new RLHF technique, but rather to propose a way to use RLHF with modest human preference annotations, omitting vanilla RLHF as a baseline does not affect our conclusions in any way.

We use BART-Large (Lewis et al., 2020) for all of our models (implementation details in §F). The choice of the model is governed by two factors: (a) It provides a similar environment (model size) for comparison with SOTA, (b) We find that LLMs (Mistral-7B, LLaMA2-7B, Zephyr-7B, etc.) are already quite good at opinion summarization; thus any performance benefits (over SOTA) cannot be reliably attributed to our approach.

5.2 Evaluation Results

We test our approach on 9 benchmarks (§3.2). In the main manuscript we report the following things: (i) ROUGE evaluations on Amazon-based benchmarks (Table 1), (ii) human evaluations on AMAZON benchmark & GPT-4 evaluations (Liu et al., 2023) on AMAZON, FLIPKART and OPOSUM+ benchmarks, and (iii) ALIGNSCORE evaluations (Zha et al., 2023) on AMAZON, FLIPKART and OPOSUM+ benchmarks to check the faithfulness of the models (Table 2). We include the remaining automatic evaluations (ROUGE & BERTSCORE) in §D (Tables 5, 6 and 7). Our human/GPT-4 evaluations are reference-

free, hence, as discussed in §3.2, we report the results on only the AMAZON, FLIPKART and OPOSUM+ benchmarks. The shortcomings of these 3 benchmarks (§E) leads to unreliable conclusions for our models in reference-based evaluations on these benchmarks. Our human/GPT-4 evaluations thus complement these reference-based evaluations in proving the efficacy of our technique.

Automatic Evaluation. From Table 1, we see that our proposed models are always better than the SOTA for AMAZON-R and AMAZON-RDQ. Supervised Fine Tuning (SFT) on PROMPTOPINSUMM (SUPERVISED model) helps achieve significantly better ROUGE scores. This highlights the efficacy of our proposed PROMPTOPINSUMM dataset. From the automatic evaluations on AMAZON-R and AMAZON-RDQ, we see the following things:

Answer to SCENE-I: We see that INDUCTIVE-BIAS achieves gains over SUPERVISED. This answers the question in SCENE-I: Our technique is effective over and above using a good SFT dataset.

Answer to SCENE-II: We see that INDUCTIVE-BIAS achieves gains over NAIVEMEAN. This answers the question in SCENE-II: Our technique is effective over RL.

Answer to SCENE-III: We see that INDUCTIVE-BIAS achieves gains over SYNTH-FEEDBACK. This answers the question in SCENE-III: Our technique is effective over the SOTA RLHF technique, which works without human preference data.

Answer to SCENE-IV: We see that INDUCTIVE-

BIAS achieves gains over DPO. This verifies that gains of INDUCTIVE-BIAS can be safely attributed to the approach (not just the quality of OPINPREF).

Additionally, in the reference-free automatic evaluation (ALIGNSCORE; Table 2), we again see that our approach (INDUCTIVE-BIAS) leads to significantly more faithful summaries too.

Model	AMAZON	FLIPKART	OPOSUM+
OP-SUM-GEN	62.83	—	70.36
Ground-Truth	67.34	36.0	69.97
DPO	58.10	60.16	54.32
SUPERVISED	80.98	80.29	76.87
NAIVEMEAN	<u>82.10</u>	<u>82.91</u>	<u>78.80</u>
SYNTH-FEEDBACK	81.36	81.97	77.97
INDUCTIVE-BIAS	84.32[†]	84.10[†]	80.76[†]

Table 2: ALIGNSCORE evaluations (Zha et al., 2023) for AMAZON, FLIPKART and OPOSUM+. ALIGNSCORE (higher is better) evaluates how faithful the summaries are, given the reviews. As highlighted in §3.2, the *-R & *-RDQ differ only in the opinion summaries from the original benchmarks. Thus, we run this reference-free evaluation for just these benchmarks. We also check the faithfulness the ground truth summaries in these benchmarks. Notation: [†] indicates that INDUCTIVE-BIAS is significantly better than the next-best with a significance level of 5%.

Human/GPT-4 Evaluation. We conduct human evaluation (Figure 1) for the AMAZON benchmark, using 3 domain experts (details in §J). We observe a Fleiss’ Kappa (κ) score of 56.25% (moderate agreement; agreement is moderate when $40\% \leq \kappa < 60\%$). We ask the experts to rank the summaries (anonymized and shuffled) given the reviews. Given the rankings, we compute the fraction of pairwise wins, ties, and losses among all the models. We compare summaries from SUPERVISED, NAIVEMEAN, SYNTH-FEEDBACK, INDUCTIVE-BIAS, OP-SUM-GEN (SOTA) models and ground truth summaries. We include ground truth summaries in the evaluation to verify our claims about the quality of the benchmarks. From Figure 1, we see that INDUCTIVE-BIAS wins significantly over the competitors, further proving the efficacy of our technique.

We run GPT-4 evaluations for AMAZON, FLIPKART and OPOSUM+ benchmarks (Figures 3, 4, 5). We run GPT-4 evaluations for AMAZON, as the agreement in human evaluation was moderate. We arrive at the same conclusions as human evaluation. We prompt GPT-4 to rank the summaries (anonymized and shuffled) given the re-

views. As before, we compute the fraction of wins, ties, and losses. Again, we see that INDUCTIVE-BIAS remains a clear winner.

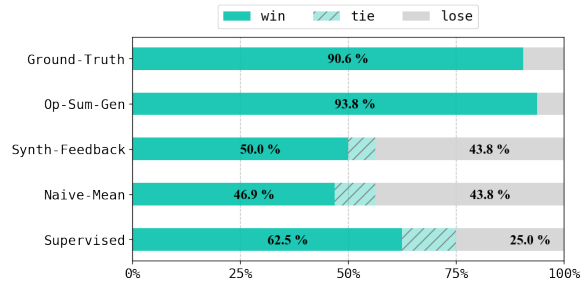


Figure 3: GPT-4 Eval: Pairwise win-tie-loss percentage of INDUCTIVE-BIAS model vs. competitors, for AMAZON benchmark.

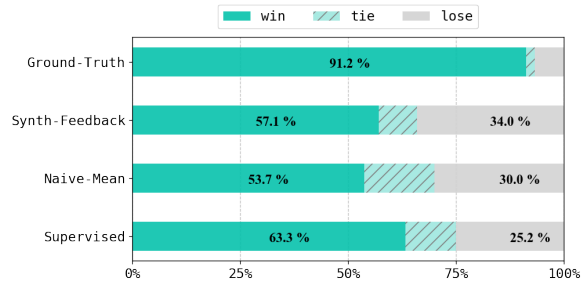


Figure 4: GPT-4 Eval: Pairwise win-tie-loss percentage of INDUCTIVE-BIAS model vs. competitors, for FLIPKART benchmark. For the FLIPKART benchmark, we do not have results from OP-SUM-GEN, as Siledar et al. (2023a) only provide aspect-specific summaries.

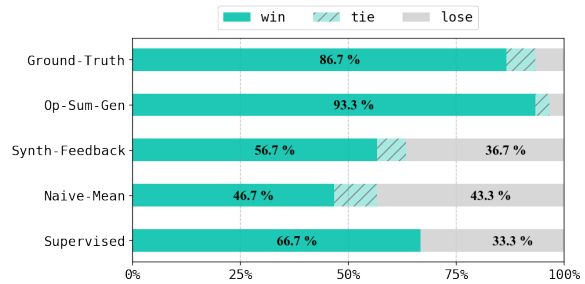


Figure 5: GPT-4 Eval: Pairwise win-tie-loss percentage of INDUCTIVE-BIAS model vs. competitors, for OPOSUM+ benchmark.

6 Analysis

We perform a two-fold analysis: (a) First, we see the domain knowledge features influence for φ_{op} , (b) Second, we see how the ground truth summary and summary from trained models fare on the domain knowledge features. This two-fold analysis

helps us understand: (a) which features influence the latent reward model within humans⁸ the most, and (b) how the ground truth summary and summary from trained models fare on these influential features. Performing well on influential features would mean the summary aligns well with the latent reward model within humans.

6.1 Analysis of φ_{op}

φ_{op} model has been trained on a set of features specified by domain experts. We analyze the relative influence of each feature on the final score assigned by φ_{op} . Doing this helps us understand an approximate importance⁹ of each of these features. We do this by varying each feature by δ ($= 0.1$) while keeping the other features constant, over multiple possible values of all features (Equation 2).

$$\Delta_i = \frac{1}{2\delta} \sum_{\mathbf{x}} (f(x_1, \dots, x_i + \delta, \dots, x_n) - f(x_1, \dots, x_i - \delta, \dots, x_n)) \quad (2)$$

Figure 6 highlights the features’ relative influence. We see that hallucination is most influential. This aligns with what our human preference annotators report—hallucination in summary is the primary cause of rejection. We see that hallucinations are mostly within the opinions in the summary. This is also reflected in Figure 6: opinion-faithfulness has significant influence. We also see that annotators prefer summaries with more specifics, i.e. they include more aspects: aspect-coverage has significant influence.

6.2 Analysis of Summaries

We analyze the top-3 performing models (in human and GPT-4 evaluations) for the following features: opinion-coverage, opinion-faithfulness, hallucination and relevance. We show the analysis only for the AMAZON benchmark in the main manuscript, we include the rest in §K. Table 3 shows the performance on these features. We see that INDUCTIVE-BIAS model fares much better than the competitors on hallucination (the most influential metric). For relevance, aspect-coverage and opinion-faithfulness, our model is fairly better than the other models.

This shows that our technique helps INDUCTIVE-BIAS model perform well on features that influence the latent reward model within humans for

⁸Note that trained φ_{op} represents the latent human reward.

⁹We call this approximate importance as the influence of a feature on the output is not necessarily its importance.

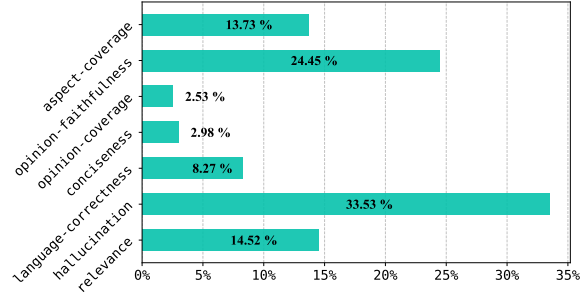


Figure 6: Relative Influence of all features in φ_{op} . All the influences sum to 100%.

Models	AC \uparrow	OPF \uparrow	RE \uparrow	HL \uparrow
IB	3.60	3.93	4.06	4.07
SF	3.43	3.73	4.04	3.94
NM	3.57	3.91	4.04	3.09

Table 3: Scores on domain knowledge-based features (AC: aspect-coverage, RE: relevance, OPF: opinion-faithfulness, HL: hallucination) on the AMAZON benchmark for top-3 models (IB: INDUCTIVE-BIAS, NM: NAIVEMEAN, SF: SYNTH-FEEDBACK). Note that for hallucination, Φ gives a higher score for less hallucination in the text.

opinion summarization. This means that our technique helps INDUCTIVE-BIAS achieve a significant alignment with the latent reward model. This conclusion **verifies our hypothesis** for opinion summarization: *A domain-knowledge infused reward model (φ_{op}) can help achieve alignment with latent reward model of humans for a task, with modest human preference annotations.*

7 Summary, Conclusion and Future Work

In this work, we provide a technique to employ RLHF for E-Commerce Opinion Summarization, using modest amount of human preferences. We achieve State-of-the-Art, while significantly reducing the size of preference data (just 940 samples). In addition to advancing SOTA and reducing preference annotations, our technique provides another two-fold benefits: (i) No Alignment Tax and (ii) Interpretability. Due to the interpretable nature, we find that our model does achieve alignment with human goals for Opinion Summarization through analysis. We propose this technique as a way of achieving alignment for Opinion Summarization. However, our technique can be applied to other tasks, with task-specific features. In the future, we will verify this approach for several Natural Language Generation tasks.

8 Ethical Considerations

We contribute two datasets in our work: PROMPTOPINSUMM, OPINPREF. These datasets are generated using an open-source model Mistral-7B (Jiang et al., 2023). We would release the datasets to further research in Opinion Summarization. For the OPINPREF, to the best of our knowledge, we have seen that it does not contain any harmful content, such as social biases, stereotypes, etc. However, we have seen that it contains products of explicit nature (sexual products). For the PROMPTOPINSUMM dataset, to the best of our knowledge, there is no presence of harmful content, such as social biases, stereotypes etc. We urge the research community to use the datasets with caution and check for potential harmfulness, based on their use-cases.

9 Limitations

A limitation of our work is it relies on a pretrained LLM (we denote it using Φ) to generate the values for the 7 features. While testing this approach on other domains, it might not be the case that a pretrained LLM can directly generate such values. The Φ model might need some fine-tuning. Additionally, the features that we propose in this paper work only for Opinion Summarization. While testing this approach on another domain, domain-experts would have to be consulted again to formulate features for that domain.

References

Reinald Kim Amplayo, Stefanos Angelidis, and Mirella Lapata. 2021. [Aspect-controllable opinion summarization](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6578–6593, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Reinald Kim Amplayo and Mirella Lapata. 2020. [Unsupervised opinion summarization with noising and denoising](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1934–1945, Online. Association for Computational Linguistics.

Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, and Jared Kaplan. 2021. [A](#)

[general language assistant as a laboratory for alignment](#). 652
653

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, T. J. Henighan, Nicholas Joseph, Saurav Kadavath, John Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Christopher Olah, Benjamin Mann, and Jared Kaplan. 2022a. [Training a helpful and harmless assistant with reinforcement learning from human feedback](#). *ArXiv*, abs/2204.05862. 654
655
656
657
658
659
660
661
662
663
664
665
666

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 2022b. [Constitutional ai: Harmlessness from ai feedback](#). 667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684

Adithya Bhaskar, Alex Fabbri, and Greg Durrett. 2023. [Prompted opinion summarization with GPT-3.5](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 9282–9300, Toronto, Canada. Association for Computational Linguistics. 685
686
687
688
689

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. 690
691
692
693

Arthur Bražinskas, Mirella Lapata, and Ivan Titov. 2020. [Unsupervised opinion summarization as copycat-review generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5151–5169, Online. Association for Computational Linguistics. 694
695
696
697
698
699

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality](#). 700
701
702
703
704
705

Eric Chu and Peter J. Liu. 2018. [Meansum: A neural model for unsupervised multi-document abstractive summarization](#). In *International Conference on Machine Learning*. 706
707
708
709

710	Kawin Ethayarajh, Yejin Choi, and Swabha Swayamdipta. 2022. Understanding dataset difficulty with \mathcal{V}-usable information . In <i>Proceedings of the 39th International Conference on Machine Learning</i> , volume 162 of <i>Proceedings of Machine Learning Research</i> , pages 5988–6008. PMLR.	<i>Processing</i> , pages 13677–13700, Singapore. Association for Computational Linguistics.	768 769
716	Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham, Jonathan Uesato, Po-Sen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth Dathathri, Rory Greig, Charlie Chen, Doug Fritz, Jaume Sanchez Elias, Richard Green, Soňa Mokrá, Nicholas Fernando, Boxi Wu, Rachel Foley, Susannah Young, Iason Gabriel, William Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu, Lisa Anne Hendricks, and Geoffrey Irving. 2022. Improving alignment of dialogue agents via targeted human judgements .	Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. 2023. Rlaif: Scaling reinforcement learning from human feedback with ai feedback .	770 771 772 773 774
729	Ruining He and Julian McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering . In <i>Proceedings of the 25th International Conference on World Wide Web</i> , WWW '16, page 507–517, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880, Online. Association for Computational Linguistics.	775 776 777 778 779 780 781 782 783
736	Minqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews . In <i>Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining</i> , KDD '04, page 168–177, New York, NY, USA. Association for Computing Machinery.	Ruibo Liu, Ge Zhang, Xinyu Feng, and Soroush Vosoughi. 2022. Aligning generative language models with human values . In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 241–252, Seattle, United States. Association for Computational Linguistics.	784 785 786 787 788 789
742	Jinbae Im, Moonki Kim, Hoyeop Lee, Hyunsouk Cho, and Sehee Chung. 2021. Self-supervised multimodal opinion summarization . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 388–403, Online. Association for Computational Linguistics.	Yang Liu, Dan Iter, Yichong Xu, Shuhang Wang, Ruo Chen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics.	790 791 792 793 794 795 796
750	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L��lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2023. Mistral 7b .	R.D. Luce. 2012. <i>Individual Choice Behavior: A Theoretical Analysis</i> . Dover Books on Mathematics. Dover Publications.	797 798 799
757	Liwei Jiang, Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Jenny Liang, Jesse Dodge, Keisuke Sakaguchi, Maxwell Forbes, Jon Borchardt, Saadia Gabriel, Yulia Tsvetkov, Oren Etzioni, Maarten Sap, Regina Rini, and Yejin Choi. 2022. Can machines learn morality? the delphi experiment .	Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Ouyang Long, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. 2021. Webgpt: Browser-assisted question-answering with human feedback . <i>ArXiv</i> , abs/2112.09332.	800 801 802 803 804 805 806 807
763	Sungdong Kim, Sanghwan Bae, Jamin Shin, Soyoung Kang, Donghyun Kwak, Kang Yoo, and Minjoon Seo. 2023. Aligning large language models through synthetic feedback . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language</i>	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 27730–27744. Curran Associates, Inc.	808 809 810 811 812 813 814 815 816 817
767		Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4 . <i>arXiv preprint arXiv:2304.03277</i> .	818 819 820
		R. L. Plackett. 1975. The analysis of permutations . <i>Journal of the Royal Statistical Society. Series C (Applied Statistics)</i> , 24(2):193–202.	821 822 823

824	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.	880
825		881
826		882
827		883
828	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms.	884
829		885
830		886
831	Tejpal Singh Siledar, Suman Banerjee, Amey Patil, Sudhanshu Singh, Muthusamy Chelliah, Nikesh Garera, and Pushpak Bhattacharyya. 2023a. Synthesize, if you do not have: Effective synthetic dataset creation strategies for self-supervised opinion summarization in E-commerce. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 13480–13491, Singapore. Association for Computational Linguistics.	887
832		888
833		889
834		890
835		891
836		892
837		893
838		894
839		895
840	Tejpal Singh Siledar, Jigar Makwana, and Pushpak Bhattacharyya. 2023b. Aspect-sentiment-based opinion summarization using multiple information sources. In <i>Proceedings of the 6th Joint International Conference on Data Science & Management of Data (10th ACM IKDD CODS and 28th COMAD), Mumbai, India, January 4-7, 2023</i> , pages 55–61. ACM.	896
841		897
842		898
843		899
844		900
845		901
846		902
847	Tejpal Singh Siledar, Rupasai Rangaraju, Sankara Sri Raghava Ravindra Muddu, Suman Banerjee, Amey Patil, Sudhanshu Shekhar Singh, Muthusamy Chelliah, Nikesh Garera, Swaprava Nath, and Pushpak Bhattacharyya. 2024. Product description and qa assisted self-supervised opinion summarization.	903
848		904
849		905
850		906
851		907
852		908
853	Taylor Sorensen, Liwei Jiang, Jena Hwang, Sydney Levine, Valentina Pyatkin, Peter West, Nouha Dziri, Ximing Lu, Kavel Rao, Chandra Bhagavatula, Maarten Sap, John Tasioulas, and Yejin Choi. 2023. Value kaleidoscope: Engaging ai with pluralistic human values, rights, and duties.	909
854		910
855		911
856		912
857		913
858		914
859	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca .	915
860		916
861		917
862		918
863		
864	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models.	
865		
866		
867		
868		
869		
870	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura,	
871		
872		
873		
874		
875		
876		
877		
878		
879		
	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models.	
	Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. 2023. Zephyr: Direct distillation of lm alignment.	
	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.	
	Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu. 2023. AlignScore: Evaluating factual consistency with a unified alignment function. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11328–11348, Toronto, Canada. Association for Computational Linguistics.	
	Daniel M. Ziegler, Nisan Stiennon, Jeff Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. <i>ArXiv</i> , abs/1909.08593.	

A Background, Terminologies & Notations

In this section, we will provide details on two important facets of this paper: **Reinforcement Learning from Human Feedback** (§A.1), and **Opinion Summarization** (§A.2).

A.1 Background: Reinforcement Learning from Human Feedback (RLHF)

Before diving into the details of RLHF, we provide the definition of Alignment, in the context of Natural Language Generation, below.

Aligning Language Models (LMs) with Humans.

Aligning LMs with humans refers to the desire/goal of making these LMs generate text which is calibrated to human values, such as ‘*text should not have harmful content*’, ‘*text should be polite*’, ‘*text should not contain illegal content*’, etc.

RLHF has emerged as a dominant paradigm in achieving alignment within Large Language Models (LLMs) (Ziegler et al., 2019; Ouyang et al., 2022). The ubiquitous recipe for RLHF follows a 3-step process: (i) Supervised Fine Tuning (SFT) stage, (ii) Reward Modeling stage, (iii) Reward Maximization stage. We describe the three processes below:

Supervised Fine Tuning (SFT). This stage involves training the LLM using human-annotated outputs (Y) on inputs (X). A dataset of the form: $\mathcal{D}_{SFT} = \{(X, Y)\}$ is used to train the LLM. Maximum Likelihood Estimation, with the objective of maximizing the likelihood of Y given X , is used.

Reward Modeling. This stage involves training a Reward Model (r), which attempts to mimic the latent reward model within humans (r^*). r^* is what humans use to judge how good a generated text is. Thus, if r is trained to be approximately equivalent to r^* , it can be used to quickly judge the goodness of a text from a human perspective, without actually using the human. Learning r is done using Preference Modeling (Bradley and Terry, 1952; Plackett, 1975). Preference Modeling involves the usage of a Preference Dataset, of the form: $\mathcal{D}_{PM} = \{(X, Y_w, Y_l) | Y_w \succ Y_l\}$, where Y_w and Y_l are two outputs sampled from the LLM on an input X , and Y_w is preferred over Y_l by the human. Step 2 in Figure 7 depicts how such a preference data is collected. Such a dataset is then used to train r using the Elo loss (Equation 3).

$$\mathcal{L}_{PM} = -\mathbb{E}_{(X, Y_w, Y_l) \sim \mathcal{D}_{PM}} \left[\log \sigma(r(X, Y_w) - r(X, Y_l)) \right] \quad (3)$$

Reward Maximization. In this stage, the trained reward model (r) is used to further train the LLM. The training is typically done using the popular Reinforcement Learning algorithm: Proximal Policy Optimization (PPO; Schulman et al. (2017)). In this stage, the LLM is asked to generate several outputs (Y_i) on a given input X . Each of these Y_i is rewarded using the trained r . This reward is then used by PPO to update the LLM¹⁰.

We refer the reader to Ouyang et al. (2022) for a more detailed discussion on RLHF.

A.2 Background: Opinion Summarization

Opinion Summarization refers to the task of summarizing users’ opinions, expressed in the form of reviews. These reviews can be about a product (books, laptops, cellphones, etc.), or a movie, concert, circus, etc. In this paper, we work with opinions expressed about a product on an E-Commerce platform (such as [amazon.com](https://www.amazon.com)). The input and output for the task of E-Commerce Opinion Summarization are specified below:

Input: *Reviews* on the product, written by users.

Output: A *summary* of these reviews.

These reviews can discuss several aspects of the product (such as *battery*, *memory*, etc. for a *cellphone*), providing varying opinions (positive, negative or mixed). An opinion summary must concisely collate all such opinions on these aspects. We provide an example below:

Input Reviews:

Review 1: *The new smartphone model has a fantastic battery life and an impressive camera. However, the screen size is too large for my liking.*

Review 2: *I love the sleek design and the battery life is outstanding. But, the phone is a bit too expensive for the features it offers.*

Review 3: *Great camera quality and the battery lasts all day. The large screen makes it difficult to use with one hand though.*

¹⁰Note that in the context of RLHF, the LLM is the RL agent.

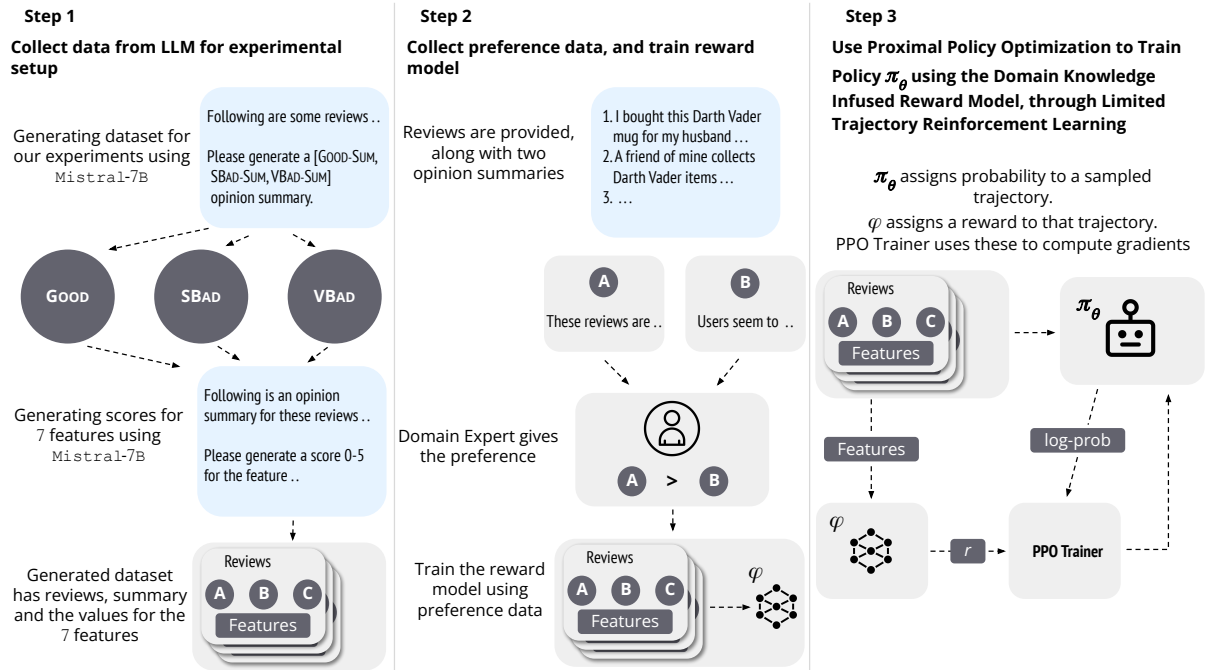


Figure 7: Overview of our approach. Step-1: We generate a new dataset for training Opinion Summarizers: PROMPTOPINSUMM, by prompting Mistral-7B model. Again, we use Mistral-7B to compute values for the 7 features discussed in §4.1. Step-2: We ask humans (domain experts) for their preference, given reviews and two opinion summaries (A, B). We use the preference data and the features to train the reward model, φ_{op} . Step-3: We sample instances from PROMPTOPINSUMM dataset; φ_{op} assigns a score to the sampled summaries, the policy, π_θ , assigns *log probabilities* to these summaries. Proximal Policy Optimization uses these to update π_θ .

Opinion Summary: *The new smartphone is praised for its excellent battery life and camera quality. However, users find the large screen size inconvenient and consider the phone to be overpriced.*

In the example, three user reviews are provided. These reviews talk about several aspects (underlined in the example) of a smartphone: *battery, design, camera*, etc. The opinion summary summarizes the users' opinions towards all such aspects, and presents it in a crisp and concise piece of text.

B Features for Reward Modeling

We use 7 domain specific features for the reward model φ_{op} . We identify these features after extensive discussions with the domain experts. For each feature we prompt Mistral-7B to generate a score within 0 and 5. We give elongated instructions, including rough rubrics, and rules to generate the scores. This is a reason why we use an instruction-tuned model. For each feature, 0 means the model is doing bad on the feature, and 5 means the model is doing good on the feature. We define all the features below:

aspect-coverage: This feature considers the aspect coverage within an opinion summary. The

feature assumes a value 5 if all the aspects of the product, mentioned in the reviews, are mentioned in the summary. If none of the aspects are picked, the feature assumes a value 0.

opinion-faithfulness: This feature considers whether the mentioned opinions/sentiments in the summary are correct, that is, they are picked correctly from the reviews. For example, if a user mentions that they are *happy* with the battery of a phone, and the summary mentions that users are *unhappy* with the battery, the summary will not be considered faithful to opinion in the review. The feature assumes a value 5 if all the opinions are faithfully reflected. If no opinion is faithfully reflected, the value would be 0.

opinion-coverage: This feature considers whether all the opinions in the input reviews are picked by the opinion summary. The feature assumes a value 5 if all the opinions are picked up. If none of the opinions are picked up, the feature assumes a value 0.

relevance: This feature checks if the summary is relevant to the input reviews (that is the product). The feature assumes a value 5 if summary is completely relevant. If it is completely irrelevant, the

feature assumes a value 0.

conciseness: This feature considers the conciseness and completeness of the opinion summary. The feature assumes a value 5 if the summary is concise and complete—not one phrase/sentence can be dropped off. It assumes a value 0 if the summary is totally incomplete, or very verbose.

hallucination: This feature considers the factuality of the opinion summary. The feature assumes a value 5 if the summary is totally factual, with respect to the input reviews. If there are a lot of hallucinations, the feature assumes a value 0.

language-correctness: This feature checks the correctness of language/text in the opinion summary. The feature assumes a value 5 if the summary is grammatically fully correct. It assumes a value 0 if the summary is very poor linguistically.

For conciseness, we do not include the prompts in the paper, we would release them as separate artifacts, with the datasets, in the camera ready version.

We also analyze how these features correlate with humans’ judgement of goodness of opinion summaries. We do this by looking at the scores for these features for preferred and dis-preferred summaries in the OPINPREF dataset. In Table 4, we see that the preferred summaries clearly have a higher score on all the features, than the dis-preferred ones. This shows that the scores correlate well with humans’ judgement of goodness.

Feature	Pref.	Dis-pref.
aspect-coverage (↑)	3.69	2.84
opinion-faithfulness (↑)	4.02	3.05
opinion-coverage (↑)	3.92	3.22
conciseness (↑)	4.05	3.44
relevance (↑)	4.10	3.10
hallucination (↑)	3.99	2.79
language-correctness (↑)	4.50	3.32

Table 4: Scores for the domain knowledge based features. We see that for all the features, the human preferred (Pref.) summaries have higher scores than the ones rejected by humans (Dis-pref.). This shows that these features correlate well with humans’ judgement of goodness of an opinion summary.

C RLHF Training Pipeline

Using the trained reward model, we follow a similar training pipeline as Bai et al. (2022a); Ouyang

et al. (2022), with a modification: *Limited Trajectory Reinforcement Learning*. Computing the transformation Φ for each generation online (during training) is expensive, especially with limited compute resources. To circumvent this, we limit the trajectories that are explored by our policy, π_θ . Specifically, we limit it to the GOOD-SUM, SBAD-SUM and VBAD-SUM trajectories in the PROMPTOPINSUMM dataset. Having varying levels of quality in PROMPTOPINSUMM is of use here—it lets the model still explore trajectories of several quality. Thus, we have an offline experience buffer, with Φ precomputed, for π_θ learn from.

We use Proximal Policy Optimization (PPO) (Schulman et al., 2017) to train our model (Equation 4). For each training step, we sample $(R, s, \Phi(R, s))$ tuples from PROMPTOPINSUMM. We use the trained φ_{op} to compute the reward for s ($= \varphi_{op}(\Phi(R, s))$). PPO uses this to update the log probability assigned by π_θ . We parameterize π_θ using a Transformer model, which takes reviews as input, and generates an opinion summary.

$$\mathcal{L}_{PPO} = -\mathbb{E}_{(R,s,\Phi(s))} \left[\varphi_{op}(\Phi(R, s)) - \beta \log \left(\frac{\pi_\theta^{RL}(s|R)}{\pi^{SFT}(s|R)} \right) \right] \quad (4)$$

D Additional Automatic Evaluation Results

In addition to the Amazon-based benchmarks (Table 1), we also report results for Flipkart and Oposum+ based benchmarks (Tables 5 and 6). As before, we see that INDUCTIVE-BIAS is almost always the winner. As before, we draw similar conclusions for SCENE-I, SCENE-II and SCENE-III: INDUCTIVE-BIAS wins, further strengthening the conclusion that our methodology is effective. We also see that, inspite of the shortcomings of the FLIPKART benchmark, our models perform similar to the SOTA.

We also include BERTSCORE evaluations for all the 9 benchmarks in Table 7. We see similar trends as ROUGE Evaluation: our models are significantly better than the SOTA in majority of the benchmarks.

For a qualitative understanding, we include generations from several models on a randomly picked sample from the AMAZON benchmark in Table 15.

Model-Code		FLIPKART			FLIPKART-R			FLIPKART-RDQ		
		R-1 ↑	R-2 ↑	R-L ↑	R-1 ↑	R-2 ↑	R-L ↑	R-1 ↑	R-2 ↑	R-L ↑
MEDOS (Siledar et al., 2024)		25.97	5.29	<u>16.05</u>	26.29	4.03	16.59	22.92	4.30	16.35
Ours'	DPO	28.85	4.10	15.55	34.23	7.86	18.62	29.96	5.25	17.28
	SUPERVISED	27.38	4.09	15.37	39.32	10.52	22.56	32.25	6.88	19.04
	NAIVEMEAN	28.34	<u>4.38</u>	16.20	40.56	10.68	22.74	<u>32.57</u>	6.67	<u>19.39</u>
	SYNTH-FEEDBACK	26.37	4.18	15.48	38.77	<u>10.99</u>	<u>22.97</u>	31.04	<u>6.98</u>	18.59
	INDUCTIVE-BIAS	<u>27.42</u>	4.21	15.71	<u>39.10</u>	11.03	23.30 [†]	33.08 [†]	7.30 [†]	19.46 [‡]

Table 5: Reference-based Evaluation Results (R-1: ROUGE-1, R-2: ROUGE-2, R-L: ROUGE-L) for the FLIPKART, FLIPKART-R and FLIPKART-RDQ benchmarks. We see the following things: (a) Our proposed dataset (PROMPTOPINSUMM) leads to *marked increased over the SOTA* (MEDOS; by ~ 6 R-L points), (b) INDUCTIVE-BIAS proves to be the *winner in all the four scenarios*: SCENE-I, SCENE-II, SCENE-III and SCENE-IV (§5.1), *proving the efficacy of our technique*. We also see that for FLIPKART benchmark, despite the shortcomings, our models perform similar to the SOTA. **Notation:** [†] denotes that INDUCTIVE-BIAS is significantly better than next-best with $\alpha = 5\%$, and [‡] denotes that INDUCTIVE-BIAS is significantly better than next-best with $\alpha = 10\%$; α denotes significance level.

Model-Code		OPOSUM+			OPOSUM-R			OPOSUM-RDQ		
		R-1 ↑	R-2 ↑	R-L ↑	R-1 ↑	R-2 ↑	R-L ↑	R-1 ↑	R-2 ↑	R-L ↑
Prior Works	MeanSum (Chu and Liu, 2018)	26.25	4.62	16.49	—	—	—	—	—	—
	CopyCat (Bražinskas et al., 2020)	27.98	5.79	17.07	22.41	2.30	13.94	22.38	2.03	14.06
	PlanSum (Amplayo and Lapata, 2020)	30.26	5.29	17.48	22.37	2.05	13.32	22.64	2.25	13.71
	MultimodalSum (Im et al., 2021)	33.08	7.46	19.75	23.35	2.98	14.53	23.73	2.80	14.70
	OP-SUM-GEN (Siledar et al., 2023a)	36.44	8.50	22.03	25.65	3.56	15.83	24.66	3.25	15.54
	MEDOS (Siledar et al., 2024)	36.57	<u>8.79</u>	21.35	26.82	3.67	15.92	26.32	3.34	16.10
Ours'	DPO	27.64	7.34	16.50	33.69	6.62	18.55	30.95	5.89	17.60
	SUPERVISED	30.57	8.02	16.90	38.32	9.10	20.35	35.69	8.17	19.28
	NAIVEMEAN	31.47	8.0	16.99	40.16	9.84	21.74	35.90	8.33	20.13
	SYNTH-FEEDBACK	<u>31.66</u>	8.86	<u>17.91</u>	<u>41.32</u>	10.40	22.23	37.85	<u>8.94</u>	<u>20.71</u>
	INDUCTIVE-BIAS	31.15	8.15	17.46	41.58 [†]	<u>10.32</u>	<u>22.02</u>	<u>37.56</u>	9.21 [†]	20.88 [‡]

Table 6: Reference-based Evaluation Results (R-1: ROUGE-1, R-2: ROUGE-2, R-L: ROUGE-L) for the OPOSUM+, OPOSUM-R and OPOSUM-RDQ benchmarks. We see the following things: (a) Our proposed dataset (PROMPTOPINSUMM) leads to *marked increased over the SOTA* (MEDOS; by ~ 6 R-L points), (b) INDUCTIVE-BIAS proves to be the *winner in almost all of the four scenarios*: SCENE-I, SCENE-II, SCENE-III and SCENE-IV (§5.1), *proving the efficacy of our technique*. We also see that for OPOSUM+ benchmark, our models lag behind. However, *this is expected*, as we highlight in §E. **Notation:** [†] denotes that INDUCTIVE-BIAS is significantly better than next-best with $\alpha = 5\%$, and [‡] denotes that INDUCTIVE-BIAS is significantly better than next-best with $\alpha = 10\%$; α denotes significance level.

E Details on the Benchmark Datasets

In this section we discuss details about the benchmarks, such as the domain of the products, summary statistics and finally highlight some shortcomings in the AMAZON, OPOSUM+ and FLIPKART datasets. AMAZON has reviews for 32 products from 4 domains: “electronics”, “home & kitchen”, “personal care”, and “clothing, shoes & jewellery”. OPOSUM+ has reviews for 60 products from 6 domains: “laptop bags”, “bluetooth headsets”, “boots”, “keyboards”, “television”, and “vacuums”. FLIPKART has reviews for 147 products from 3 domains: “laptops”, “mobiles”, and “tablets”. Table 8 includes summary statistics for the benchmarks.

Finally, now we highlight the shortcomings of the benchmark datasets in the rest of the discussion.

AMAZON: Bražinskas et al. (2020) designed the test-set in such a way that the summary has to read like a review, for instance, summary would contain ‘*I think the quality has come down over the years.*’, instead of ‘*Users think that quality has come down over years.*’. Due to this writing style, the summaries read like reviews and are often in first person—high overlap would not necessarily mean a better summary, it would rather mean a better review.

FLIPKART: Siledar et al. (2023b) generate this dataset by listing out the aspect-wise pros and cons presented within the reviews. We form an opin-

Model Code	AMAZON	AMAZON-R	AMAZON-RDQ	OPOSUM+	OPOSUM-R	OPOSUM-RDQ	FLIPKART	FLIPKART-R	FLIPKART-RDQ
OP-SUM-GEN (Siledar et al., 2023a)	88.78	86.94	86.76	86.63	86.96	86.95	–	–	–
DPO	86.45	86.60	86.37	84.39	87.35*	86.90	83.75	86.61	85.40
SUPERVISED	87.79	88.23*	87.76*	85.13	88.59*	88.02*	84.21	88.11	86.40
NAIVEMEAN	87.95	<u>88.29*</u>	<u>87.81*</u>	85.25	88.96*	88.39*	<u>84.32</u>	<u>88.29</u>	<u>86.52</u>
SYNTH-FEEDBACK	87.81	88.28*	87.74*	85.22	<u>89.08*</u>	<u>88.45*</u>	84.27	88.28	86.49
INDUCTIVE-BIAS	<u>87.98</u>	88.41*	88.16*	<u>85.33</u>	89.09*	88.46*	84.33	88.34	86.61

Table 7: BERTSCORE evaluation results on the 9 benchmark datasets. We observe a similar trend as ROUGE evaluations: SOTA is better than our models for the AMAZON and OPOSUM+ benchmarks, which is expected (§3.2). For the rest of the datasets, we see that our models are significantly better. We do not include SOTA results for Flipkart-based benchmarks, as OP-SUM-GEN only provide aspect-specific summaries for the same. * denotes gain is statistically significant compared to SOTA with significance level 1%.

Characteristic	OPOSUM+	AMAZON	FLIPKART
# domains	6	4	3
# products	60	32	147
# reviews per product	10	8	10
# summaries per product	3	3	1

Table 8: Statistics of the benchmark datasets. OPOSUM+ represents the statistics of all OPOSUM+ based benchmarks (OPOSUM+, OPOSUM-R and OPOSUM-RDQ). Similar is the case for AMAZON and FLIPKART.

ion summary by concatenating these pros and cons. Due to this, the summaries have frequent incoherent sentences.

OPOSUM+: Amplayo et al. (2021) create this benchmark by extracting sentences from the input reviews. Hence, this dataset has similar drawbacks as the AMAZON benchmark.

AMAZON

Nice boots but run a bit narrow. They look great but I think the quality has come down over the years. Still comfortable but I wish they broke in easier. I recommend these for any lady who is patient and looking for comfort.

OPOSUM+

great product for the cost . very easy to use and compatible with all of my phones ! it holds a charge great , is light enough and fits perfectly in my ear . the sound quality is great , the style is very cool and the unit feels top quality . it would drop and reconnect every 10 seconds nobody could hear me i could n't get it to unpair from the phone , there 's apparently no noise-cancellation in these . the

battery life is ... bizarre . cheap , plastic-y , and poor sound quality .

FLIPKART

Summary

Pros

Design: The full-metal Infinix INBook X1 Core i3 has a top notch and premium design.

35.56 cm (14 inch) 1920 x 1080 Pixel Full HD IPS Display: 100% sRGB with 300nits brightness ensures an excellent display.

Battery: Long-lasting battery. Gives around 8 hours of backup on normal usage.

Performance: The combination of Intel Core processor chip, high RAM size and sufficient storage capacity gives this laptop a high-speed performance experience.

Price: "Totally worth it in this price range.

Cons

Charging: Some current leakage during charging. Sometimes the laptop won't charge.

Trackpad: Not upto the mark.

Verdict: *This laptop comes with a i3 10th gen dual core processor which is suitable for normal tasks like web browsing, online classes and watching movies. Not recommended as a gaming laptop.*

Additional Information: *Can handle video editing and expandable SSD.*

F Implementation Details

We use BART-Large (Lewis et al., 2020) as our policy (π_θ) in all of the models. We do this to have a fair comparison with the state-of-the-art in Opinion Summarization. We use AdamW Optimizer

to train the models, with a weight decay of 0.05. We use a cosine learning rate scheduler. We run a hyperparameter sweep on batch size, learning rate, and learning rate warmup. We include the possible values for the sweep in Table 9. We train all of our models using $2 \times$ A100 GPUs (80GB)

Hyperparameter	Values
batch size	[64, 128, 256]
learning rate	$\sim \mathcal{U}(5e^{-6}, 5e^{-5})$
learning rate warmup	$\sim \mathcal{U}(0.2, 0.4)$

Table 9: Possible Values for Hyperparameters. For learning rate warmup, we sample the fraction of total steps the learning should be warmed up. For example, if the learning rate warmup is 0.2, it means that the learning rate will have a linear warmup for 20% of the total training steps.

For the reward model, φ_{op} , we use a Feed Forward Network for the Policy Model. We use AdamW Optimizer to train the models, with a weight decay of 0.05. As before, we run a hyperparameter sweep on batch size, learning rate, and learning rate warmup. Table 10 includes details on the hyperparameters.

Hyperparameter	Values
batch size	[32, 64, 128]
learning rate	$\sim \mathcal{U}(5e^{-3}, 1e^{-1})$

Table 10: Possible Values for Hyperparameters for the Reward Model. For learning rate warmup, we sample the fraction of total steps the learning should be warmed up. For example, if the learning rate warmup is 0.2, it means that the learning rate will have a linear warmup for 20% of the total training steps.

G Generated Summary Lengths

We analyze the generation lengths of the models, and the ground truth summary. Table 11 lists the summary lengths. We see that the DPO model generates very verbose summary. Additionally, we also see that the INDUCTIVE-BIAS model generates very concise summaries.

H Details on PROMPTOPINSUMM

Here we provide more details on the generated PROMPTOPINSUMM dataset. Table 14 includes summary statistics of the generated dataset.

Model	AMAZON	OPOSUM+	FLIPKART
Ground-Truth	60.65	85.86	129.91
NAIVEMEAN	91.09	114.67	75.48
SYNTH-FEEDBACK	80.31	115.37	71.11
OP-SUM-GEN	55.84	62.93	-
INDUCTIVE-BIAS	81.62	88.63	73.57
SUPERVISED	81.31	117.03	74.56
DPO	138.50	141.50	131.40

Table 11: Generation Length Statistics: number of words in summaries. We use NLTK to tokenize the text.

We include an example from PROMPTOPINSUMM dataset (Table 12). We show one sample from GOOD-SUM, SBAD-SUM and VBAD-SUM quality each. We do not include the reviews for conciseness. However, we include the salient aspects of the reviews. The reviews talk about the following things:

1. Great price, Nice looking / Good color, Good utility / Good protection, Quick shipping, Nice fitting, Good accessibility of the laptop while the cover is on, Good finish quality.
2. Not a good fit for older macbooks, Broken / Unusable for the original (\sim 2009 - 2010) white macbooks, cannot return return.

GOOD-SUM manages to discuss all of these things—it was able to detect that the cover does not fit specific models (highlighted in blue). It also detects that a few users like the “personal touch addition” factor of the cover.

VBAD-SUM is also a good representation of a bad summary—it totally ignores the “positive” aspects of the product and presents the “negative” aspects only. It fails at Aspect Coverage, Opinion Faithfulness and Opinion Coverage.

SBAD-SUM maintains almost a similar quality as the Good one. However, it fails to draw out certain aspects, such as “pizzaz”, “personal touch addition”, etc.

I Statistics of the OPINPREF dataset

We look at the summary statistics for the OPINPREF dataset. Table 13. We see that, interestingly, annotators prefer longer summaries—this is because these summaries contain more specifics and details from the reviews.

J Annotator Details

We include two disjoint sets of annotators in our work—first for creation of OPINPREF (3 annota-

Good: The users have shared their positive experiences with the protective covers for Macbooks. They appreciate the ease of application, vibrant colors, quick shipping, and great quality. Some have mentioned the covers fit their Macbooks perfectly, while others have shared their disappointment when they realized it did not fit their specific model. A few users have noted the covers add a personal touch and pizzazz to their devices. However, some users have reported issues with the covers not fitting their Macbooks or falling off, leaving them feeling disappointed and frustrated.

Slightly bad: Some users have reported positive experiences with the Macbook cover, praising its protective qualities, easy application, and vibrant colors. Others, however, have encountered issues with it not fitting properly on their devices or being returned due to size mismatches. Overall, the product has been described as cute, great quality, and worth the investment, though some buyers have experienced disappointment with its size compatibility and easy detachment.

Very Bad: This is a terrible product for people with old Macbooks, it doesn't fit at all! And don't even bother looking at the description carefully before buying.

Table 12: Example summaries from PROMPTOPINSUMM.

Characteristic	Value
# words in reviews	641.21
# reviews	13.08
# words in summaries	73.16
# words in preferred summaries	85.41
# words in unpreferred summaries	66.91

Table 13: Statistics of the OPINPREF dataset. We use NLTK to tokenize the text.

Split	Characteristic	μ	σ
train	# reviews	13.24	10.07
	# summaries	8.90	0.34
	# words in review	49.0	10.78
	# words in summary	78.28	34.45
validation	# reviews	10.53	6.80
	# summaries	8.98	0.16
	# words in review	48.65	10.63
	# words in summary	74.26	34.27

Table 14: Statistics of PROMPTOPINSUMM dataset. We use NLTK to tokenize the text.

tors), second for human evaluation (3 annotators). For both annotations, we use domain experts. The domain experts are NLP researchers (age group: 24 – 30) who have worked in Opinion Summarization for a long time, with publication experience (in A/A* conferences). The domain experts for human evaluation also have a similar profile. The annotators have been paid generously, based on the standard annotation rates in the geographical location.

K All Evaluation Results

We include all of the evaluation results in this section. In Tables 16, 17, 18 and 18 we include pairwise comparison results, in a win/tie/loss format. We also include results on evaluation on how the models perform on the domain features in Tables 20, 21 and 22.

1300

1301

1302

1303

1304

1305

1306

Reviews: i really like these boots. they can be a beast to get on, like any boot fit to last; once on, they are incredibly comfortable. I have had them for 4 years and they still look great - worn in but not worn out...

What great boots! They do take some breaking in - like all real leather does. Run maybe a tad small in the general spectrum, but hey, remember the break in period. Look fabulous, will last me for the next twenty years.... They bridge the fashion vs. function quandary beautifully.

Every woman needs a pair of these Frye boots! Fit is true to normal size. I wear a standard 9. These fit in the shaft for any ladies who have bigger calves as I do. You have to fight to get them on first few times and they wear stiff until fully broken in. Worth the money and break in time.

Boots excellent quality, beautiful skin, a modern model. But! they are very narrow! This is the second pair of shoes is a brand and my 8.5 size fits me perfectly, even he little more than usual. But these are very narrow! go to a cobbler, let him try to stretch. Be careful!

bought these for the wife to have something to kick around in and keep up with the kids, but she doesn't like them. Her heel lifts with every step, and they're stiff and heavy. Not an everyday casual shoe.

Don't get me wrong, they're really good boots, but they used to make them better. My last pair from a few years ago were much higher quality. Thicker sole, better full grain leather, now they only make the black in 'grease', which stains all my socks black. Keeping them, but still disappointed in Frye.

i bought it just for the heavy snow cuz my UGG got wet through.even though it is a little expensive, but the boots are nice and it was shipped to me for only two days!!!!!! i love it

They were comfortable the moment I put them on. The sizes are exact. The color the way they are made are what I expect from Frye which I've been wearing since the 70's. With a 20year break from not wearing Frye it was a pleasant surprise the quality has stood the test of time.

SUPERVISED: These Frye boots are a popular choice for women due to their **comfort** and **durability**. However, they can be a bit challenging to put on and **require a break-in period**. Some users have reported that they run small in size and are narrow, requiring a **cobbler's assistance for a proper fit**. The **leather quality is generally good**, but some users have noted that the **leather may not last as long as they had hoped**.

NAIVEMEAN: Users generally appreciate the **comfort** and **durability** of the Frye boots, with some mentioning the need for a **break-in period**. **Sizes run true** to normal, but some users find them to be a **bit narrow**. The boots are praised for their **stylish appearance** and ability to **bridge the fashion and function quandary**. However, some users have noted that they take **some time to break in** and are **not suitable for everyday casual wear**.

SYNTH-FEEDBACK: These Frye boots are praised for their **comfort** and **durability**, with some users noting the need for a **break-in period**. However, some users have noted that they run small in size and require stretching, while others find them **narrow and stiff**. The **leather is described as beautiful and of good quality**, with a **modern model**. Some users have **had them for several years** and have not worn them out.

INDUCTIVE-BIAS: These Frye boots are praised for their **comfort** and **durability**, with users expressing satisfaction with their **fit** and **appearance**. However, some users have noted that they require a **break-in period** and can be a **bit narrow**. The **leather quality is praised** for its durability and **modern model**, with some users expressing **disappointment with the lack of improved quality in recent years**.

Table 15: Example generation (randomly sampled) for some input reviews from all the models. **Olive** implies faithful/correct generation, while **red** indicates hallucinated text, or repetition. We see that only INDUCTIVE-BIAS is free from **red** text. The model closest in performance to INDUCTIVE-BIAS, the NAIVEMEAN model, misses out on two aspects: leather-quality and quality-degradation. INDUCTIVE-BIAS covers both, while being concise. We do not include DPO model in this comparison, as it was too verbose.

	SUPERVISED	NAIVEMEAN	SYNTH-FEEDBACK	INDUCTIVE-BIAS	OP-SUM-GEN
NAIVEMEAN	0.50/0.06/0.38				
SYNTH-FEEDBACK	0.44/0.12/0.44	0.40/0.09/0.5			
INDUCTIVE-BIAS	0.56 /0.09/0.28	0.46 /0.18/0.31	0.56 /0.12/0.28		
OP-SUM-GEN	0.31/0.28/0.38	0.25/0.12/0.56	0.25/0.21/0.5	0.25/0.06/ 0.68	
Ground-Truth	0.46/0.06/0.48	0.31/0.18/0.44	0.40/0.15/0.40	0.28/0.09/ 0.59	0.5/0.09/0.38

Table 16: Pairwise Win/Tie/Loss Results for all models in Human Evaluation for AMAZON benchmark. We format the data as: win/tie/loss, win specifies how many time the *row* won over the *column*.

	SUPERVISED	NAIVEMEAN	SYNTH-FEEDBACK	INDUCTIVE-BIAS	OP-SUM-GEN
NAIVEMEAN	0.63/0.12/0.25				
SYNTH-FEEDBACK	0.59/0.12/0.28	0.5/0.06/0.44			
INDUCTIVE-BIAS	0.62 /0.12/0.25	0.46 /0.09/0.44	0.5 /0.06/0.44		
OP-SUM-GEN	0.06/0.03/0.9	0.09/0.0/0.90	0.12/0.09/0.78	0.06/0.0/ 0.93	
ground-truth	0.12/0.06/0.81	0.09/0.06/0.84	0.16/0.06/0.78	0.09/0.0/ 0.90	0.68/0.09/0.22

Table 17: Pairwise Win/Tie/Loss Results for all models in GPT-4 Evaluation for AMAZON benchmark. We format the data as: win/tie/loss, win specifies how many time the *row* won over the *column*.

	SUPERVISED	NAIVEMEAN	SYNTH-FEEDBACK	INDUCTIVE-BIAS
NAIVEMEAN	0.57/0.12/0.30			
SYNTH-FEEDBACK	0.57/0.06/0.36	0.52/0.12/0.36		
INDUCTIVE-BIAS	0.63 /0.12/0.25	0.54 /0.16/0.30	0.57 /0.08/0.34	
Ground-Truth	0.10/0.06/0.84	0.06/0.01/0.92	0.07/0.01/0.91	0.06/0.02/ 0.91

Table 18: Pairwise Win/Tie/Loss Results for all models in GPT-4 Evaluation for FLIPKART benchmark. We format the data as: win/tie/loss, win specifies how many time the *row* won over the *column*.

	SUPERVISED	NAIVEMEAN	SYNTH-FEEDBACK	INDUCTIVE-BIAS	OP-SUM-GEN
NAIVEMEAN	0.56/0.03/0.4				
SYNTH-FEEDBACK	0.5/0.16/0.34	0.46/0.1/0.44			
INDUCTIVE-BIAS	0.66 /0.0/0.33	0.46 /0.1/0.44	0.56 /0.06/0.36		
OP-SUM-GEN	0.1/0.06/0.83	0.06/0.03/0.9	0.03/0.03/0.93	0.03/0.03/ 0.93	
Ground-Truth	0.13/0.13/0.73	0.1/0.033/0.8666	0.06/0.06/0.86	0.06/0.06/ 0.86	0.7/0.1/0.2

Table 19: Pairwise Win/Tie/Loss Results for all models in GPT-4 Evaluation for OPOSUM+ benchmark. We format the data as: win/tie/loss.

	AC	OPF	OPC	CC	RL	HL	LC
SUPERVISED	3.43 ± 0.20	3.71 ± 0.37	3.67 ± 0.26	3.79 ± 0.31	4.04 ± 0.37	3.89 ± 0.39	4.55 ± 0.35
NAIVEMEAN	3.56 ± 0.22	3.91 ± 0.50	3.76 ± 0.38	3.89 ± 0.36	4.04 ± 0.48	3.99 ± 0.48	4.60 ± 0.27
SYNTH-FEEDBACK	3.55 ± 0.40	3.87 ± 0.71	3.71 ± 0.43	3.94 ± 0.50	4.04 ± 0.61	3.94 ± 0.68	4.38 ± 0.92
INDUCTIVE-BIAS	3.60 ± 0.17	3.95 ± 0.40	3.85 ± 0.25	3.99 ± 0.35	4.06 ± 0.34	4.07 ± 0.43	4.65 ± 0.32
OP-SUM-GEN	3.34 ± 0.68	3.92 ± 0.79	3.70 ± 0.54	4.0 ± 0.50	4.08 ± 0.72	3.87 ± 1.08	4.05 ± 1.31
Ground-Truth	3.55 ± 0.50	3.93 ± 0.46	3.56 ± 0.31	4.08 ± 0.32	4.04 ± 0.46	3.81 ± 0.86	4.40 ± 0.45

Table 20: Intrinsic Evaluation results on the AMAZON benchmark for all the models. Legend: AC: aspect-coverage, OPF: opinion-faithfulness, OPC: opinion-coverage, CC: conciseness, RE: relevance, HL: hallucination, LC: language-correctness.

	AC	OPF	OPC	CC	RL	HL	LC
SUPERVISED	3.61 ± 0.22	4.10 ± 0.39	3.84 ± 0.33	4.04 ± 0.28	4.21 ± 0.31	4.19 ± 0.42	4.53 ± 0.27
NAIVEMEAN	3.56 ± 0.21	4.13 ± 0.41	3.84 ± 0.34	4.0 ± 0.32	4.31 ± 0.36	4.26 ± 0.34	4.54 ± 0.39
SYNTH-FEEDBACK	3.56 ± 0.25	4.09 ± 0.40	3.79 ± 0.32	4.02 ± 0.30	4.19 ± 0.34	4.19 ± 0.36	4.53 ± 0.29
INDUCTIVE-BIAS	3.63 ± 0.20	4.22 ± 0.39	3.85 ± 0.30	4.01 ± 0.28	4.26 ± 0.29	4.33 ± 0.45	4.61 ± 0.29
Ground-Truth	3.59 ± 0.15	3.88 ± 0.53	3.68 ± 0.27	4.02 ± 0.28	3.87 ± 0.59	3.67 ± 0.78	4.35 ± 0.44

Table 21: Intrinsic Evaluation results on the FLIPKART benchmark for all the models. Legend: AC: aspect-coverage, OPF: opinion-faithfulness, OPC: opinion-coverage, CC: conciseness, RE: relevance, HL: hallucination, LC: language-correctness.

.	AC	OPF	OPC	CC	RL	HL	LC
SUPERVISED	3.47 ± 0.14	3.38 ± 0.26	3.49 ± 0.06	3.64 ± 0.19	3.81 ± 0.26	3.22 ± 0.56	3.96 ± 0.32
NAIVEMEAN	3.49 ± 0.05	3.48 ± 0.06	3.5 ± 0.0	3.56 ± 0.13	3.66 ± 0.22	3.52 ± 0.33	4.1 ± 0.33
SYNTH-FEEDBACK	3.50 ± 0.03	3.41 ± 0.26	3.5 ± 0.0	3.63 ± 0.24	3.62 ± 0.20	3.32 ± 0.63	4.03 ± 0.38
INDUCTIVE-BIAS	3.54 ± 0.22	3.50 ± 0.06	3.57 ± 0.06	3.62 ± 0.19	3.65 ± 0.23	3.68 ± 0.36	4.0 ± 0.29
OP-SUM-GEN	3.39 ± 0.3	3.46 ± 0.45	3.49 ± 0.28	3.61 ± 0.40	3.58 ± 0.82	3.43 ± 0.92	3.79 ± 1.18
Ground-Truth	3.42 ± 0.22	3.475 ± 0.28	3.5 ± 0.0	3.57 ± 0.16	3.49 ± 0.28	3.21 ± 0.48	3.56 ± 0.23

Table 22: Intrinsic Evaluation results on the OPOSUM+ benchmark for all the models. Legend: AC: aspect-coverage, OPF: opinion-faithfulness, OPC: opinion-coverage, CC: conciseness, RE: relevance, HL: hallucination, LC: language-correctness.