Frictional Agent Alignment Framework: Slow Down and Don't Break Things

Anonymous ACL submission

Abstract

AI support of collaborative interactions entails 001 002 mediating potential misalignment between interlocutor beliefs. Common preference alignment methods like DPO excel in static settings, but struggle in dynamic collaborative tasks where the explicit signals of interlocutor beliefs are sparse and skewed. We propose the Frictional Agent Alignment Framework (FAAF), to generate precise, context-aware "friction" that prompts for deliberation and re-examination of existing evidence. FAAF's two-player objective 012 decouples from data skew: a frictive-state policy identifies belief misalignments, while an intervention policy crafts collaborator-preferred responses. We derive an analytical solution to this objective, enabling training a single policy 017 via a simple supervised loss. Experiments on three benchmarks show FAAF outperforms competitors in producing concise, interpretable friction and in OOD generalization. By aligning LLMs to act as adaptive "thought partners"not passive responders-FAAF advances scalable, dynamic human-AI collaboration.

1 Introduction

037

041

When collaborating to solve problems, humans continually interrogate each other's intentions and assumptions (Stalnaker, 2002; Asher and Gillies, 2003; Klein et al., 2005). With the rapid integration of generative AI, exemplified by large language models (LLMs), into personal, educational, business, and even governmental workflows, AI systems will increasingly be called upon to act as collaborators with humans; to adequately fill this role, AIs must be able to recapitulate the reflection and deliberation that makes human-human collaboration successful, but also causes temporary slowdowns in dialogue while interlocutors construct a *common ground* on which to collectively reason we will call this phenomenon **friction**.

"Friction" in this sense is something that LLMs struggle with. To prompt an interlocutor to reflect



Figure 1: FAAF conditions responses on both the dialogue context x and representation of the "frictive" (belief) state ϕ , to generate outputs that prompt for reflection, deliberation, and verification of evidence.

upon their assumptions requires that one have an approximate understanding of what those assumptions are and entail (Lewis and Sarkadi, 2024). This is predicated upon a *theory of mind* (ToM; Premack and Woodruff (1978)), which is likewise a challenge for LLMs (Sap et al., 2022; Ullman, 2023).

To address this, we present the **Frictional Agent Alignment Framework (FAAF)**, a novel approach to aligning LLMs to be adept *collaborators* in dialogue-driven tasks. Unlike common preference alignment approaches which focus predominantly on reward differences between textual surface forms to generate the best possible completions as a sequence of actions, FAAF takes a statedriven approach based on the notion of a *frictive state*—a dynamic natural language representation that integrates task context and the beliefs of participants as they change over time (Fig. 1). We use this state-wise representation to train "friction agent" models aligned to prompt collaborators to-

061

062ward reflection and deliberation in shared tasks, to063help them resolve conflicting beliefs and assump-064tions that result in frictive states. Our results on two065challenging collaborative task datasets and variants066show that FAAF's belief state conditioning consis-067tently produces output that is more relevant, impact-068ful on the dialogue, and thought-provoking than069competing methods. Our key contributions are:

- a novel offline LLM alignment framework for collaborative agents based not on pairwise reward differences between responses, but rather on advantage of an intervention over a dialogue state representation;
- development of an LLM "agent" that inserts interventions into collaborative dialogues to prompt participants toward reflection, deliberation, and common ground;
- evaluations on three challenging collaborative task datasets that show the advantages of FAAF over competing alignment methods.

2 Related Work

074

075

081

090

091

100

101

103

104

105

106

107

108

109

110

111

RLHF-inspired preference alignment in LLMs has become a cornerstone of developing generative AI systems that cater to user preferences (Stiennon et al., 2020). Both "offline" approaches like Direct Preference Optimization (DPO; Rafailov et al. (2024b)), Identity Preference Optimization (IPO; Azar et al. (2024)) and other supervised methods (Meng et al., 2024; Hong et al., 2024; Fisch et al., 2024; Pal et al., 2024; Nath et al., 2024b) and "online" methods (Schulman et al., 2017; Pang et al., 2024) focus predominantly on preference samples often sourced from datasets like Reddit TL;DR (Völske et al., 2017) or Ultrafeedback (Cui et al., 2024) for algorithm development.

These methods excel in generating summaries or completions that reflect human preferences including on single-turn human-AI interaction datasets like SGD (Rastogi et al., 2020) or MultiWOZ (Zang et al., 2020; Ye et al., 2022), but are often ill-equipped to handle the complexities of real-world multiparty interactions, where communication occurs across diverse modalities, including sparse and ambiguous spoken dialogues between multiple collaborators (Karadzhov et al., 2023; Khebour et al., 2024b).

A key challenge in these multiparty shared task settings is the scarcity of annotated data, particularly where interventions emerge contextually but sparsely (Karadzhov et al., 2023; Khebour et al., 2024b). While preference data generated with AIfeedback is a viable option (Li et al., 2023b; Yuan et al., 2024), DPO-trained models depend crucially on the sampling or data-generating distribution due to its Bradley-Terry (BT) model of "implicit rewards," limiting their applications to dialoguedriven settings where preferences may be intransitive (Tversky, 1969) or change over time. This data-dependence is true even for more sophisticated methods that optimize on human utility (Ethayarajh et al., 2024), discard the BT assumption (Azar et al., 2024), or use iterative online approaches (Rosset et al., 2024; Pang et al., 2024). Game-theoretic approaches to reduce this dependence focus on optimizing a "general preference model" (Munos et al., 2023; Calandriello et al., 2024) that does not suffer from this data-bias. But these have limited practical application due to their compute-intensive nature, often requiring the storage and computations with intermediate-stage policies during training (Choi et al., 2024). In contrast, FAAF avoids this datadependence by explicitly conditioning policies on belief-misalignment in a simple "one-step" supervised manner without requiring computations of complicated mixture policies during training.

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

3 Definitions

Let us first define key terms we rely on.

Frictive state Entailed by Clark (1996)'s common ground, or the set of beliefs shared by interlocutors, a frictive state arises during a collaborative task when different interlocutors have contradictory beliefs about a task-relevant proposition (i.e., one believes p and another sees evidence against p). This can be realized as a formal model of agent beliefs in an evidence-based dynamic epistemic logic (van Benthem et al., 2014; Pacuit, 2017), or a natural language description thereof, as we use. Different evidence leads to different predictions of future trajectories (Craik, 1943). Thus frictive states, though sparse in dialogues, can critically delay or preclude success in a collaboration due to unresolved misunderstandings. The occurrence of a frictive state may not guarantee task failure, as p may be trivial to actual task completion. Therefore, in a functionally frictive state, the lack of common ground prohibits progress on the task, or presents a significant risk of failure unless it is resolved.

Friction intervention Friction can indicate an impasse (the frictive state), but can also be used to resolve it, through a *friction intervention* that

inserts into the dialogue indirect prompting to the 163 participants to reevaluate their beliefs and incor-164 rect assumptions or positions in light of available 165 evidence (Oinas-Kukkonen and Harjumaa, 2009), 166 rather than accepting possibly erroneous presuppositions inherent in the dialogue. Importantly, a 168 frictive intervention may be non-contradictory to 169 the individual beliefs on display (i.e., neither assert-170 ing p nor $\neg p$), but rather slows down the dialogue for reflection and deliberation, such as the probing 172 utterances in Karadzhov et al. (2023) and Nath et al. 173 (2024c). In the context of LLMs and FAAF, the fric-174 tion agent constitutes a language model aligned 175 toward the capacity to make frictive interventions.¹ 176 An ideal friction agent does not intervene arbitrar-177 ily, which would cause distraction in collaborative 178 tasks, but is conditioned to resolve the lack of common ground between human collaborators.

4 Task Formulation and Background

181

182

184

186

190

191

192

193

194

198

199

201

210

Let f be a frictive intervention (utterance) that is not required to contradict any particular belief encapsulated in a frictive state ϕ , and let the human preference probability $\mathcal{P}(f \succ \phi)$ be the probability that an expert annotator would prefer f over maintaining ϕ , given prior dialogue history, x. An RLHF-based approach to LLM alignment toward an optimal policy π_f^* would assume a partition function $Z^*(\phi, x)$ that normalizes the probabilities of all possible responses (see Appendix B for more details). While the optimal policy formulation is closed form, the dependence on Z^* makes it practically intractable to estimate it for LLMs since Z^* is a summation over the set of all possible sequences of tokens in the tokenizer, often requiring methods like importance sampling (Korbak et al., 2022) or ensembling models (Go et al., 2023) for an unbiased estimate. This problem remains even if the set of friction interventions \mathcal{F} were a restricted subset of the space of all possible actions \mathcal{Y} . To overcome this, prior RLHF and Preference-based RL (Wirth et al., 2017) literature suggests supervised learning algorithms for obtaining an optimal policy induced under the expectation over a preference dataset. These offline methods, such as DPO (Rafailov et al., 2024b), IPO (Azar et al., 2024), or Kahneman-Tversky Optimization (KTO; (Ethayarajh et al., 2024)), either rely on the BT model of preferences (Bradley and Terry, 1952) where the optimal pol-

¹We use π_f to denote the friction agent which generates high-quality interventions, but refer to it as the "optimal policy" for consistency with RLHF literature.

icy can be induced from a static preference dataset using implicitly-defined pointwise rewards, or assume that alignment is conducted with access to a non-biased data-generation or "sampling" distribution μ from which π_f^* can be learned using pairwise preferences without adopting a strictly BT assumption (Azar et al., 2024).² These approaches would give us the following formulation for π_f^* :

$$\pi_f^* = \frac{\pi_{\text{ref}} \exp\left(\beta^{-1} \mathbb{E}_{f \sim \mu(\cdot \mid x)} \Psi(\mathcal{P}(f \succ \phi \mid x))\right)}{\frac{\phi \sim \mu(\cdot \mid x)}{Z^*(\phi, x)}}, \quad (1)$$

where $\Psi(p)$ is the identity mapping for IPO, and $\log\left(\frac{\mathcal{P}}{1-\mathcal{P}}\right)$ (inverse sigmoid) for DPO and KTO.

While the practicality of these supervised algorithms is a clear advantage, their dependence on preference data selected via a sampling-is a limitation in reconstructing the human preference probability \mathcal{P} . This is particularly true for collaborative dialogue tasks where common ground changes over time, meaning that the occurrence of frictive states is dynamic, and where participants may not intervene due to variables obscure to a language model, such as not realizing the existence of a frictive state or judging the frictive state to be non-functional (Sec. 3). Operationally, even if the true underlying preferences (\mathcal{P}) of collaborators are transitive and consistent, constructing a preference dataset for use with existing offline training methods is not straightforward as dialogues may be skewed or sparse (Khebour et al., 2024b). When using generative AI to create denser training data, even high-capacity LLMs like GPT-4 are prone to various forms of biases such as toward length (Lambert et al., 2024) or certain linguistic registers. Therefore, the core motivation of FAAF is as followshow do we train a high-quality friction agent that can leverage the inherent scalability of offline alignment methods and reconstruct the true underlying preference distribution while still being robust to the data skew that may arise when sampling a preference dataset, whether using generative AI or from real-life collaborative dialogues?

4.1 FAAF Objective

We define a novel two-player adversarial optimization objective J_{FAAF}^* (Eq. 2). Specifically, given a reference model π_{ref} and a regularization parameter $\beta \in \mathbb{R}_+$, our goal is to learn two interdependent 213 214

215 216

211

212

217

218

219 220

221

223

224

225

226

228

229

230

231

232

233

234

235

236

237

238

240

241

242

243

244

245

246

247

248

249

250

251

253

254

255

²By "sampling," we mean those actions that make it to the preference annotation phase after being sampled with the data-generator μ .

256"collaborative" policies: (i) a frictive state policy257 π_{ϕ}^* that generates the most semantically rich fric-258tive states ϕ , capturing tensions or uncertainties259(in the form of first-order beliefs of dialog partici-260pants) in dialogue, and (ii) a friction intervention261policy π_f^* that generates constructive interventions262f, conditioned on the frictive state, to improve dis-263course clarity and converge onto a common ground264between participants. Mathematically,

$$J_{\mathsf{FAAF}}^{*} = \min_{\pi_{\phi}} \max_{\pi_{f}} \mathbb{E} \sum_{\substack{x \sim \rho \\ \phi \sim \pi_{\phi}(\cdot \mid x) \\ f \sim \pi_{f}(\cdot \mid \phi, x)}} \begin{bmatrix} \mathcal{P}(f \succ \phi \mid x) \\ -\beta D_{\mathsf{KL}}(\pi_{f} \parallel \pi_{\mathsf{ref}} \mid \phi, x) \\ +\beta D_{\mathsf{KL}}(\pi_{\phi} \parallel \pi_{\mathsf{ref}} \mid x) \end{bmatrix}.$$
(2)

267

270

271

275

276

279

287

291

292

Notice how the optimal intervention policy π_f^* , by definition of the inner max operator, generates interventions that are, on average, most preferred by collaborators, while the first KL-divergence term, defined as $D_{\text{KL}}(\cdot \mid \phi, x)$, stabilizes learning in π_f^* by keeping it closer to a reference model. Compared to a standard RLHF objective, the additional KL term $D_{\mathrm{KL}}(\pi_{\phi} \parallel \pi_{\mathrm{ref}} \mid x)$ forces the frictive state policy π^*_{ϕ} to be adversarially robust, in that it must ensure that sampled frictive states ϕ $\sim \pi_{\phi}^*$ cannot be exploited by π_f^* to generate subpar interventions that remain too close to the reference model. Thus, FAAF serves as an agent policy that adapts to dialogues over time: the frictive state policy searches for the most immediate tension points or exposes the lack of common ground between task participants, while the intervention policy generates outputs that remain grounded in the particulars of the relevant frictive state (e.g., regarding the correct task items or propositions), and naturally and intuitively prompts for reflection and deliberation on these points. The key takeaway is that optimal friction interventions should not be arbitrary interventions in the dialogue, but should surface the presuppositions that gave rise to the most logically necessary frictive state, making interventions precise and interpretable.

4.2 Dataset Annotation and Generation

In Sec. 2, we discuss why common preference optimization datases such as Ultrafeedback, Reddit TL;DR, SGD, or MultiWOZ are not appropriate for FAAF's collaborative task use case. Therefore, we consider two collaborative task datasets to evaluate FAAF—DeliData (Karadzhov et al., 2023) and the Weights Task Dataset (WTD; Khebour et al. (2024a)). These datasets also exemplify the data sparsity problem with deliberation and friction in collaboration (Sec. 4). DeliData contains dialogues from 500 groups of 5 attempting the Wason Card task (Wason, 1968), which involves reasoning about if a card with a specific characteristic (e.g., even number on one side) must have a different characteristic (e.g., a vowel on the other). Karadzhov et al. (2023) annotated DeliData with "probing" interventions, or naturally-occurring friction that prompts for reasoning and deliberation without introducing new information. However, these amount to an average of only 3.46 probing interventions per group, out of 17,110 total utterances. WTD is an audiovisual dataset of 10 triads collaborating to deduce the weights of differentlycolored blocks and infer the pattern describing them, and is similarly sparse. We annotated WTD for naturally-occurring friction given a definition following Oinas-Kukkonen and Harjumaa (2009).³ Two annotators annotated half the groups each while a third annotated all 10. They then collectively adjudicated each annotation following the definition. Cohen's κ between initial and final annotations was 0.632, indicating substantial agreement. An average of 4 naturally-occurring friction interventions per group were found in the WTD.

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

346

347

Training Dataset Construction This extreme sparsity does not capture anything close to the possible frictive states available in the combinatorics of the problem space, and so motivated the need for data augmentation to construct sufficiently diverse preference datasets for training and evaluation. We used GPT-40 as a high-capacity LLM for our sampling distribution μ . We used a *self-rewarding* approach (Yuan et al., 2024) to simultaneously generate candidate interventions (and their rationales) and assign them rewards, which naturally induced an implicit preference ranking. We provided GPT-40 with sequences of h utterances from each dialogue in the two datasets, and prompted it to label frictive states and generate friction interventions following colloquial renderings of the definitions in Sec. 3.⁴ Finally, we conducted contrastive pairing of "winning" and "losing" interventions f_w and f_l

³Frictive interventions in this setting act as indirect persuasion (Oinas-Kukkonen and Harjumaa, 2009) where participants are passively prompted to reevaluate their beliefs and assumptions or propositions, in light of incoming goal-specific evidence. See Appendix A for the complete definition.

 $^{^{4}}h$ was set to 15 for DeliData and 10 for WTD.

429

430

431

432

433

434

435

436

437

438

439

440

397

399

400

401

402

403

357

348

361 363

373 374

375

389

391

with the corresponding dialogue history x to construct the final preference datasets for each task, comprising tuples of x, frictive state ϕ , f_w , and f_l .

For each dataset, we conducted additional taskappropriate augmentation. For DeliData, we constructed alternative tuples where the specific cards mentioned in the original data were replaced with other cards of the same classes that preserved the relevant rule (e.g., replacing even numbers with other even numbers, consonants with other consonants, etc.). This resulted in 68,618 preference samples for training, with average μ -assigned reward for preferred samples of 8.03 and for dispreferred samples of 3.96 (out of 10). We held out 50 randomly-sampled dialogues for testing.

Since the original WTD contains only 10 dialogues, holding one or two out for evaluation would adversely impact the data distribition. Therefore we used Shani et al. (2024)'s method to generate novel simulated collaborative conversations about the Weights Task, providing a task descriptions and ground-truth values for the weights. GPT-40 was prompted to role-play personality-facet combinations from the Big 5 personality types (Goldberg, 2013), and for each labeled frictive state ϕ we generated and scored 6 friction interventions. This resulted in two distinct versions of the WTD preference dataset. The Simulated WTD friction dataset consisted of 56,698 training preference samples, with mean scores of 8.48 (preferred interventions) and 6.01 (dispreferred). 54 dialogues were held out for testing. The Original WTD friction dataset (see above) contained 4,299 preference samples (preferred mean score 8.36, dispreferred 6.35). These were all retained for an OOD evaluation of FAAF trained on the Simulated WTD data. See Appendix D for more details on data generation.

Human Validation We conducted a human evaluation to assess the quality of the GPT-generated friction intervention on a random representative subset of 50 pairwise samples each from both the DeliData and WTD generated test datasets.⁵ For each sample, 2 annotators were asked to choose which of the two candidate interventions was more appropriate for provoking participants' reflection to help them advance in their task without being given the solution. Average Cohen's κ on WTD samples was 0.58 and on DeliData samples was 0.92, indicated substantial to near complete agreement on

which was the better intervention, and indicates that the preferred/dispreferred friction distinction sourced from GPT-40 as μ aligns with human judgments. See Appendix D.4 for more.

4.3 Deriving the Empirical FAAF Loss

While the data is constructed using a standard pairwise preference format, the FAAF optimization conditions upon the dialogue context x and textual rendering of the frictive state ϕ . To derive an empirical offline (supervised) preference learning loss from the two-player objective (Eq. 2), we use a divide-and-conquer approach. Deriving the inner maximization loop of Eq. 2 results in an analytical expression of the optimal frictive intervention policy, π_f^* (see Appendix B.1, Eq. 8). However, we observe that π_f^* in its analytical form (Eqs. 1 and 8) is not fully expressive since it does not contain the optimal frictive-state policy π_{ϕ}^{*} term. Therefore, we derive π_{ϕ}^* using a Lagrangian formulation (see Appendix C for details) that expresses the preference for any intervention f_1 over f_2 analytically in terms of **both** the optimal friction intervention policy $(\pi_f^*(\cdot \mid \phi, x))$ and the optimal frictive-state policy $(\pi_{\phi}^{*}(\cdot|x))$. This allows us to use a straightforward supervised (ℓ_2) objective—similar in spirit to IPO (Azar et al., 2024)—that empirically regresses the predicted preference expression derived from $\pi_f^*(\cdot \mid \phi, x)$ and $\pi_\phi^*(\cdot \mid x)$ to the observed relative preferences $p(f_1 \succ f_2 \mid x)$ (relative to ϕ), assuming access to a large enough preference-annotated dataset of friction interventions. Notably, this objective is optimized by a *single* parametrized policy that leverages the inherent expressivity of LLMs and induces a unique global minimum in the space of policies (see Theorem 2 in Appendix B.1). Algorithm 1 shows the full training algorithm.⁶

5 Experimental Setup

Training Setup and Baselines We use Meta-Llama-3-8B-Instruct (AI@Meta, 2024)⁷ for all experiments including baselines. For an in-depth evaluation of FAAF, we include the Supervised-Finetuned (SFT) model as well as the base instruct model generations in our experiments. For "offline" contrastive approaches, we choose

⁵WTD samples include both Original and Simulated interventions.

⁶For compactness reasons here we represent all policies π as parameterized by weights θ . Similarly to approaches such as Choi et al. (2024), because we formulate two distinct policies with the preference equation, we can empirically enforce it using ℓ_2 loss and learn it with a single expressive policy parameterized by θ .

⁷https://huggingface.co/meta-llama/ Meta-Llama-3-8B-Instruct

- **Require:** Training data \mathcal{D}_{μ} containing tuples (x, ϕ, f_w, f_l) , where x: prompt, ϕ : frictive state, f_w : preferred response, f_l : non-preferred response.
- 1: Define likelihood ratios: 2: $\Delta R = \log \left(\frac{\pi_{\theta}(f_{w}|\phi, x)}{\pi_{\theta}(f_{w}|-1, x)}\right) - \log \left(\frac{\pi_{\theta}(f_{x}|\phi, x)}{\pi_{\theta}(f_{w}|-1, x)}\right)$

$$2: \Delta P' = \log \left(\frac{\pi_{\text{ref}}(f_w | \phi, x)}{\pi_{\text{ref}}(f_w | x)} \right) = \log \left(\frac{\pi_{\theta}(f_u | \phi, x)}{\pi_{\text{ref}}(f_u | x)} \right)$$

3: $\Delta R' = \log\left(\frac{\pi_{\theta(J|W)}(x)}{\pi_{\text{ref}}(f_w|x)}\right) - \log\left(\frac{\pi_{\theta(J|W)}}{\pi_{\text{ref}}(f_l|x)}\right)$ 4: Loss function: $\mathcal{L} = \mathbb{E}_{\mathcal{D}\mu}\left[(1 - \beta(R + R'))^2\right]$

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

- 5: Gradient update: $\nabla_{\theta} \mathcal{L} = \mathbb{E}_{\mathcal{D}_{\mu}} [-2\beta \delta \nabla_{\theta} \log(R \cdot R')],$ where $\delta = 1 - \beta (\log R + \log R')$
- 6: Update policy parameters θ using gradient descent

DPO (Rafailov et al., 2024b) and IPO (Azar et al., 2024) and for "online" approaches, we include Proximal Policy Optimization (PPO; Schulman et al. (2017)) baseline. For SFT, we employ rejection sampling (Xu et al., 2023) to maximize the likelihood of interventions that receive high rewards under μ . For SFT, DPO, and IPO, the respective losses are computed only on the output tokens and frictive states ϕ , excluding dialogue context tokens. This training approach ensures that the models learn to generate effective interventions while maintaining contextual understanding. For PPO, we train an OPT 1.3B (Zhang et al., 2022) reward model on each dataset using a standard Bradley-Terry loss (Stiennon et al., 2020) over preference pairs. For *ablations*, we consider variants of FAAF that ablate the different likelihood ratios—FAAF $_{\Delta R}$ keeps only the ϕ -conditioned implicit rewards in the FAAF objective (line 2 in Algorithm 1), and FAAF $_{\Delta R'}$ removes ϕ -conditioning (keeping only line 3). See Appendix D.5 for more details on training and hyperparameters.

Evaluation Strategies As LLM generation is 463 open-ended, we employ an LLM-as-a-judge (us-464 ing GPT-40) "win-rate" evaluation method where 465 a high-capacity model is prompted to select its 466 preference, given two completions. First, we sam-467 pled friction interventions from all competing mod-468 els on 500 randomly sampled prompts from the 469 DeliData, Simulated WTD and Original WTD 470 test sets. Next, we conducted two evaluations 471 using said completions, one with a preference-472 model (Munos et al., 2023) and another with a 473 474 reward-model (Hong et al., 2024). Since GPT-40 also served as the data generation distribution μ , 475 preference-model evaluation compares the two pre-476 sented choices and nothing else in the data, mitigat-477 ing lingering bias toward μ (Munos et al., 2023). 478

Within preference-based evaluation settings, we adopt the framework proposed by Cui et al. (2024) to retrieve utility scores across seven friction dimensions, building on insights from Chen and Schmidt (2024). Specifically, we assess relevance and align*ment with rationale and golden samples*⁸ to determine how well a friction intervention aligns with surface-level semantics. Meanwhile, actionability, specificity, thought-provoking, and impact measure its expected long-term influence on behavior, reasoning, and decision-making. The LLM-judge assigns Likert-type scores across these dimensions, providing a fine-grained evaluation of task-specific preference desiderata. These scores are collected in a pairwise fashion where π_{θ} -generated interventions f_i from a baseline are compared with π_{ref} -generated counterparts, f_j . We positionally swap these interventions in the evaluation prompt (Fig. 7) for each API call and average the scores for each of the seven dimensions over two runs to mitigate positional bias (Lambert et al., 2024) in computing the final win rates. Specifically, for any pair of interventions (f_i, f_j) , let $s(x, f_*)$ denote the score estimate⁹ for intervention f_* given context x. The win-rate percentage for a run is computed as $100 \times \frac{1}{N} \sum_{m=1}^{N} \mathbf{1}\{s(x^{(m)}, f_i^{(m)}) >$ $s(x^{(m)}, f_i^{(m)})$, where N is the total number of samples, and $x^{(m)}$ represents the context of the m^{th} sample. These results are reported in Table 1.

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

502

503

504

505

507

508

509

510

511

512

513

514

515

516

517

518

519

520

522

The above evaluation tests for preference alignment advantage of the aligned model over π_{ref} . For a more robust evaluation, we compare FAAF's generations "head-to-head" against all baselines. Here, instead of the preference model, we utilize the trained OPT 1.3B Reward Model (RM) as described in our PPO training setup. These pointwise estimates of rewards provide a more accurate assessment of the advantage provided by FAAF proposed approach when directly pitted against other alignment baselines. Specifically, we compare FAAF_{\Delta R}, FAAF_{\Delta R'} as well as our full objective baseline (FAAF_{\Delta(R+R')}) against all chosen baselines. We compute the reward accuracy (or winrates) similarly and report our results in Table 2.

⁸A subset of these golden friction interventions was used for human evaluations (see Appendix D.4).

⁹In this evaluation, "overall" (first column in Table 1) is computed based on the judge's choice of winner *after* rating all other dimensions. As such, $s(x, f_*)$ represents scores over these fine-grained friction preference desiderata, and "overall" does not necessarily represent an average or aggregate of the other dimensions but rather a binary judgment based on them.

616

617

618

619

620

621

574

6 Results and Analysis

523

525

526

527

530

532

535

537

538

539

540

541

543

544

545

546

547

548

551

552

Table 1 shows that in the eyes of the the LLM-judge, FAAF models have a consistently greater advantage over the SFT model π_{ref} than other baselines across the 7 preference dimensions and overall. For instance, in "overall" preference on the DeliData test samples, FAAF achieves a 75.7% win-rate over π_{ref} , surpassing PPO (68.9%), DPO (70.8%), and IPO (70.1%). On the WTD datasets, win rates for all models are higher, reflecting π_{ref} 's weakness with the underspecified nature of WTD dialogues; alignment on this data has a greater net effect on win-rates than the generally less ambiguous Deli-Data. On WTD FAAF is a clear all-around winner, at 90.9% (vs. DPO's 89.0% and 82.0%) and 91.5% (vs. DPO's 82.9% and IPO's 83.0%) on the Original and Simulated WTD datasets, respectively.

We find that FAAF's win-rates on dimensions such as *actionability* and *gold-alignment* are somewhat lower compared to other dimensions possibly reflecting that multiple kinds of interventions may be appropriate in context. However, across dimensions like *thought-provoking* and *rationale-fit* we find that FAAF improves 5-6%, or even up to 12% over equivalent PPO, IPO, and DPO win-rates. PPO's win-rates consistently lag across all datasets (this is particularly pronounced on the WTD data), indicating the challenge that the dimensions of friction pose for a standard approach. DPO is typically FAAF's closest competitor against π_{ref} , with the narrowest average gap in win-rates.

Robustness to OOD Generalization FAAF maintains superior performance on the Original WTD 555 556 dataset (Overall: +1.9% over DPO, +8.9% over IPO, and +14.9% over PPO). No model was ex-557 plicitly aligned to this data, and so this result 558 shows FAAF's robustness to OOD settings compared to other approaches. This is particularly noteworthy given that the Original WTD dataset 561 comprises word-for-word transcriptions of actual 562 human dialogues-with disfluencies, sentence frag-563 ments, etc.-which differs markedly from the grammatical, structured text typically found in LLM 565 training data or the preference pair samples in the **Simulated WTD** data. That FAAF generalizes well to organic human data provides a strong basis 569 of confidence that a FAAF-aligned agent, jointlyconditioned on the dialogue transcript and frictive state rendering ϕ , could effectively intervene in 571 and mediate real collaborations, where dialogues are often sparse, informal, and structurally distinct 573

from LLM-generated text (Martins et al., 2020).

DPO, although also optimized against ϕ as part of the context (Sec. 5), suffers from the Longestcommon-subsequence problem (Pal et al., 2024)¹⁰ due to the Bradley-Terry preference model assumption where dependence on the context via DPO's log-partition term is effectively canceled in gradient estimates. In contrast, FAAF's combined ΔR and $\Delta R'$ regularization (Algorithm 1) avoids missing such signals in its learning, thereby allowing it to capture more nuanced human preferences.

Does ϕ -conditioning help FAAF learn more accurate preferences? Table 2 shows results from the trained OPT-1.3B RM's evaluation of the full $\mathsf{FAAF}_{\Delta(R+R')}$ objective and its ablated variants— ϕ conditioned FAAF $_{\Delta R}$ and unconditioned FAAF $_{\Delta R'}$ — "head-to-head" against all baselines, including the base Meta-Llama-3-8B-Instruct model. Across the three datasets, FAAF win-rates computed with pointwise reward estimates on sampled interventions exceed 80%, on average, against the base and SFT models, consistent with prior work (Hong et al., 2024). We also find that while explicit conditioning on ϕ provides clear advantages (e.g., +6.6%) vs. Base on Simulated WTD, +14% vs. PPO), and even the unconditioned version consistently wins over baselines, neither term alone achieves the robust performance of $FAAF_{\Delta(R+R')}$.

Both IPO and FAAF use a squared ℓ_2 loss. IPO's performance against FAAF's ablations suggests that this structural similarity makes it more competitive with FAAF (FAAF ablations beat IPO 53.6% and 54.2% on DeliData and 58.0% and 62.0% on Original WTD, compared to anywhere from a 68-85% win rate against the Base model). In general, these ablations demonstrate that neither variant alone is sufficient. $\mathsf{FAAF}_{\Delta(R+R')}$ (the full objective) shows consistently stronger performance against IPO (79.6% on DeliData, 73.7% on WTD Sim., 74.0% on WTD Orig.) while maintaining high win rates across other baselines (~81% vs Base/SFT, \sim 74% vs. DPO). These results, in light of the trends observed previously in OOD evaluation, suggest that while FAAF $_{\Delta R}$ learns rich ϕ -conditioned preferences, the additional regularization term $\Delta R'$ enables better reward space exploration and generalized preference learning. The combination is crucial for robust performance.

¹⁰The LCS issue in DPO, where gradient signals from tokens shared by winning and losing responses are ignored, is well-studied (Pal et al., 2024; Zhang et al., 2024; Rafailov et al., 2024a).

Policy	Overall	Ac	Ga	Im	Rf	Re	Sp	Th
				DELIDATA	A			
PPO	$68.9_{\pm 1.5}$	$59.9_{\pm 1.5}$	$65.4_{\pm 1.5}$	$68.6_{\pm 1.5}$	$64.9_{\pm 1.5}$	$65.1_{\pm 1.5}$	$71.1_{\pm 1.4}$	$64.0_{\pm 1.5}$
IPO	$70.1_{\pm 1.4}$	$61.2_{\pm 1.5}$	$65.7_{\pm 1.5}$	$69.3_{\pm 1.5}$	$65.3_{\pm 1.5}$	$65.5_{\pm 1.5}$	$72.1_{\pm 1.4}$	$64.1_{\pm 1.5}$
DPO	$70.8_{\pm 1.4}$	$61.0_{\pm 1.5}$	$66.8_{\pm 1.5}$	$69.6_{\pm 1.5}$	$66.1_{\pm 1.5}$	$67.5_{\pm 1.5}$	$72.2_{\pm 1.4}$	66.2 ± 1.5
FAAF	$75.7_{\pm 1.4}$	$65.6_{\pm 1.5}$	$69.5_{\pm 1.5}$	$75.0_{\pm 1.4}$	$72.0_{\pm 1.4}$	$71.1_{\pm 1.4}$	$75.3_{\pm 1.4}$	$70.4_{\pm 1.4}$
			V	VTD Origii	NAL			
PPO	$76.0_{\pm 4.3}$	$74.0_{\pm 4.4}$	$75.0_{\pm 4.3}$	$75.0_{\pm 4.3}$	$67.0_{\pm 4.7}$	$70.0_{\pm 4.6}$	$73.0_{\pm 4.4}$	$74.0_{\pm 4.4}$
IPO	$82.0_{\pm 3.8}$	$87.0_{\pm 3.4}$	$75.0_{\pm 4.3}$	$84.0_{\pm 3.7}$	$75.0_{\pm 4.3}$	$80.0_{\pm 4.0}$	$88.0_{\pm 3.2}$	$78.0_{\pm 4.1}$
DPO	$89.0_{\pm 3.1}$	$92.0_{\pm 2.7}$	$82.0_{\pm 3.8}$	$89.0_{\pm 3.1}$	$84.0_{\pm 3.7}$	$87.0_{\pm 3.4}$	$89.0_{\pm 3.1}$	$79.0_{\pm 4.1}$
FAAF	$90.9_{\pm 2.9}$	$81.8_{\pm 3.9}$	$84.8_{\pm 3.6}$	$90.9_{\pm 2.9}$	$86.9_{\pm 3.4}$	$89.9_{\pm 3.0}$	$88.9_{\pm 3.1}$	$90.9_{\pm 2.9}$
			W	TD SIMULA	ATED			
PPO	$73.6_{\pm 1.5}$	$69.7_{\pm 1.5}$	$64.9_{\pm 1.6}$	$74.2_{\pm 1.5}$	$67.6_{\pm 1.6}$	$71.9_{\pm 1.5}$	$78.1_{\pm 1.4}$	$78.3_{\pm 1.4}$
IPO	$83.0_{\pm 1.3}$	$74.8_{\pm 1.4}$	$78.4_{\pm 1.4}$	$82.9_{\pm 1.3}$	$76.9_{\pm 1.4}$	$81.4_{\pm 1.3}$	$82.5_{\pm 1.3}$	$83.2_{\pm 1.2}$
DPO	82.9 ± 1.3	80.4 ± 1.3	75.8 ± 1.4	$81.3_{\pm 1.3}$	72.9 ± 1.5	$76.3_{\pm 1.4}$	80.2 ± 1.3	$79.2_{\pm 1.4}$
FAAF	$91.5_{\pm 0.9}$	$87.5_{\pm 1.1}$	$87.1_{\pm 1.1}$	$90.1_{\pm 1.0}$	$82.0_{\pm 1.3}$	$85.1_{\pm 1.2}$	$90.3_{\pm 1.0}$	$90.1_{\pm 1.0}$

Table 1: Win-rates (%) against the SFT model (π_{ref}) for all alignment methods on sampled interventions (temperature of 0.7, top-*p* of 0.9) from 500 randomly-sampled prompts from DeliData and WTD evaluation sets, according to GPT-40. Metrics: Ac (*Actionability*), Ga (*Gold-alignment*), Im (*Impact*), Rf (*Rationale-fit*), Re (*Relevance*), Sp (*Specificity*), and Th (*Thought-provoking*). The LLM-as-a-judge evaluation follows Cui et al. (2024). Average win rates are reported over two runs, with positional swapping to mitigate position bias.

Dataset	Policy	Win-rate vs. Base	Win-rate vs. SFT	Win-rate vs. DPO	Win-rate vs. IPO	Win-rate vs. PPO
DeliData	$FAAF_{\Delta R'}$	$82.2_{\pm 1.7}$	$78.8_{\pm 1.8}$	$74.0_{\pm 1.9}$	$53.6_{\pm 2.2}$	$79.2_{\pm 1.8}$
	$FAAF_{\Delta R}$	$85.8_{\pm 1.5}$	$81.4_{\pm 1.7}$	$73.2_{\pm 1.9}$	$54.2_{\pm 2.2}$	$73.4_{\pm 1.9}$
	$FAAF_{\Delta(R+R')}$	$86.2_{\pm 1.5}$	$84.0_{\pm 1.6}$	$75.6_{\pm 1.9}$	$79.6_{\pm 1.8}$	$76.0_{\pm 1.9}$
WTD Orig.	$FAAF_{\Delta R'}$	$78.0_{\pm 5.8}$	$78.0_{\pm 5.8}$	$76.0_{\pm 6.0}$	$58.0_{\pm 6.9}$	$58.0_{\pm 6.9}$
	$FAAF_{\Delta R}$	$68.0_{\pm 6.5}$	$74.0_{\pm 6.2}$	$72.0_{\pm 6.3}$	$62.0_{\pm 6.8}$	$70.0_{\pm 6.4}$
	$FAAF_{\Delta(R+R')}$	$84.0_{\pm 5.1}$	$76.0_{\pm 6.0}$	$74.0_{\pm 6.2}$	$74.0_{\pm 6.2}$	$82.0_{\pm 5.4}$
WTD Sim.	$FAAF_{\Delta R'}$	$79.1_{\pm 1.9}$	$80.2_{\pm 1.8}$	$70.4_{\pm 2.1}$	$68.6_{\pm 2.1}$	$60.8_{\pm 2.3}$
	$FAAF_{\Delta R}$	$85.7_{\pm 1.6}$	$80.8_{\pm 1.8}$	$70.8_{\pm 2.1}$	$72.2_{\pm 2.1}$	$74.8_{\pm 2.0}$
	$FAAF_{\Delta(R+R')}$	$88.0_{\pm 1.5}$	$83.7_{\pm 1.7}$	$72.8_{\pm 2.0}$	$73.7_{\pm 2.0}$	$75.1_{\pm 2.0}$

Table 2: Win rates of of FAAF variants—FAAF $_{\Delta R'}$ (not ϕ -conditioned), FAAF $_{\Delta R}$ (ϕ -conditioned), and FAAF $_{\Delta(R+R')}$ (full objective)—against competing methods in pairwise comparisons (temperature of 0.7, top-p of 0.9). All alignment baselines are SFT-initialized and Meta-Llama-3-8B-Instruct is used as Base.

7 Conclusion

622

623

626

630

632

635

636

637

641

FAAF introduces a novel perspective on LLM alignment, focusing on the problem of generating outputs that elicit reasoning and reexamination of assumptions and evidence in a collaborative context. This critical capacity can help avert collaboration failure due to groups or individuals proceeding hastily according to their own preconceptions (Koschmann, 2016), such that a fragile common ground collapses. We proposed a novel twoplayer objective with an analytical form that can be optimized using a single policy (Sec. 4). Through evaluations on three datasets representing two different collaborative tasks, and with detailed ablations (Sec. 6), we showed that FAAF bests other common preference alignment methods in performance against a reference model, and that FAAF's simultaneous conditioning on both the frictive state ϕ and surface context x is critical to its success.

initions of "friction" in human-AI collaboration (Sec. 3). Friction creates opportunity for negotiation of intents toward a common goal, and space for accountability and collaborative reasoning. These moments may result in a net slower interaction, but are critical to eventual task success. The study of friction has broad applicability to fields like discourse studies, team science, and education (Sønneland, 2019; Collins et al., 2024; Sutton and Rao, 2024), and is something we believe the NLP community would do well to invest effort in. Counter to AI being sold as a speed and efficiency multiplier, our formulation of alignment to specialize in friction shifts LLMs from mere responders to being "thought partners," and sets a new standard for dynamic, dialogue-centric environments.

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

658

659

660

Limitations

FAAF addresses only the question of aligning language models to generate friction conditioned upon

In the process, we also put forth operational def-

a task state where the terms of the task (though 661 not the solution) are known, rather than toward a general response generation problem such as instruction following or summarization. Our goal is to train an LLM aligned toward the generation of interventions that prompt reflection and deliberation, and not a general dialogue agent/chatbot. In 667 our results we have shown that common alignment methods of the kind used in dialogue or chatbot alignment are inferior to FAAF in ability to generate 670 these kinds of utterances. This does not necessarily mean that FAAF is superior to other methods in 672 aligning for human preference in other tasks, and as discussed in Sec. 4.2, it is not clear that this 674 would be a meaningful comparison because of the 675 domain difference.

> Of course, real humans are infamous for flummoxing the most theoretically-rigorous AI systems and so the performance of FAAF (or any other alignment method) in a real multiparty collaborative setting remains an open question. FAAF provides a theoretically-grounded and empirically-validated basis of confidence for success, and we have focused on the alignment technique in this paper (and thus framed this paper as a preference alignment paper) and demonstrated feasibility on challenging collaborative task datasets, but real-time user studies remain the topic of future work.

679

684

701

702

704

708

710

712

Although we motivate FAAF based in part on theory of mind (Sec. 1), we do not claim that it imbues an LLM with ToM and acknowledge that FAAF aligned models could still inherit potential biases (say, from pretrains) in generating interventions as well as risks of overly confident or misaligned suggestions that could derail group dynamics. Instead, we use an "agentic" framework that trains a model to perform interventions for a desired effect (Russell and Norvig, 2016; Krishnaswamy et al., 2022). This is not to be confused with senses of LLM-agents such as "tool using" agents (Liu et al., 2024). Within this framework, we render the frictive state ϕ in plain English text to make it amenable to LLM input, but as briefly mentioned in Sec. 3, frictive states have a formal defintion based on evidence-based dynamic epistemic logic: a mental model $\mathcal{M} = \langle A, W, E, V \rangle$ consists of agents A, worlds W, evidence relation E defining accessibility between worlds, and valuation function V. This allows the agent to assess alternatives and predict future developments from past events (Craik, 1943). Thus, other formal structures to encode the frictive state could be explored

but were out of scope for this paper.

Finally, in terms of computational limitations, while we constructed FAAF in a way that addresses data skewness and evaluated in a manner that sought to mitigate biases in the data generation distribution μ , we cannot guarantee for certain that our results are bias-free. And, FAAF still requires a reference model to be kept in memory, which leads to some additional compute requirements.

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

References

AI@Meta. 2024. Llama 3 model card.

- Nicholas Asher and Anthony Gillies. 2003. Common ground, corrections, and coordination. *Argumentation*, 17:481–512.
- Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. 2024. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pages 4447–4455. PMLR.
- R. A. Bradley and M. E. Terry. 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345.
- Daniele Calandriello, Daniel Guo, Remi Munos, Mark Rowland, Yunhao Tang, Bernardo Avila Pires, Pierre Harvey Richemond, Charline Le Lan, Michal Valko, Tianqi Liu, et al. 2024. Human alignment of large language models through online preference optimisation. *arXiv preprint arXiv:2403.08635*.
- Zeya Chen and Ruth Schmidt. 2024. Exploring a behavioral model of "positive friction" in human-ai interaction. *Preprint*, arXiv:2402.09683.
- Eugene Choi, Arash Ahmadian, Olivier Pietquin, Matthieu Geist, and Mohammad Gheshlaghi Azar. 2024. Robust chain of thoughts preference optimization. In *Seventeenth European Workshop on Reinforcement Learning*.
- Herbert H Clark. 1996. Using language. Cambridge university press.
- Katherine M Collins, Valerie Chen, Ilia Sucholutsky, Hannah Rose Kirk, Malak Sadek, Holli Sargeant, Ameet Talwalkar, Adrian Weller, and Umang Bhatt. 2024. Modulating language model experiences through frictions. *CoRR*.
- Kenneth James Williams Craik. 1943. *The nature of explanation*, volume 445. CUP Archive.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, et al. 2024. Ultrafeedback: Boosting language models with scaled ai feedback.

764 765	In Forty-first International Conference on Machine
105	Learning.
766	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and
767	Luke Zettlemoyer. 2024. Qlora: Efficient finetuning
768	of quantized llms. Advances in Neural Information
769	Processing Systems, 36.
770	Kawin Ethayaraih Winnie Xu, Niklas Muennighoff
774	Den Jurafaku and Douwa Kiala 2024 Model align
770	mont as prospect theoretic optimization. In Forth first
773	International Conference on Machine Learning.
774	Adam Fisch, Jacob Eisenstein, Vicky Zayats, Alekh
775	Agarwal, Ahmad Beirami, Chirag Nagpal, Pete Shaw,
776	and Jonathan Berant. 2024. Robust preference opti-
777	mization through reward model distillation. Preprint,
778	arXiv:2405.19316.
779	Dongyoung Go, Tomasz Korbak, Germán Kruszewski
780	Jos Rozen, Nahveon Rvii and Marc Dymetman
781	2023 Aligning language models with preferences
782	through f-divergence minimization arXiv preprint
792	arYiv:2302.08215
103	urxiv.2502.08215.
784	Lewis R Goldberg. 2013. An alternative "description
785	of personality": The big-five factor structure. In
786	Personality and Personality Disorders, pages 34–47.
787	Routledge.
700	Harbort Dayl Crice 1075 Logic and conversation Sur
788 789	tax and semantics, 3:43–58.
790	Barbara J Grosz and Candace L Sidner. 1986. Attention,
791	intentions, and the structure of discourse. Computa-
792	tional linguistics, 12(3):175–204.
702	Ari Holtzman, Ian Buye, Li Du, Maxwell Forbes, and
793	Vaiin Chai 2010 The aurious area of neural text
794	degeneration arViv preprint arViv 1004 00751
/90	degeneration. <i>urxiv preprint urxiv.1904.09751</i> .
796	Jiwoo Hong, Noah Lee, and James Thorne. 2024.
797	ORPO: Monolithic preference optimization without
798	reference model. In Proceedings of the 2024 Con-
799	ference on Empirical Methods in Natural Language
800	Processing, pages 11170–11189.
204	Ahmad Hussain, Mahamad Madhat Gabar, Eyad Elyan
001	and Chrising Joyne, 2017. Imitation learning: A sure
002	and Christina Jayne. 2017. Initiation learning: A Sur-
803	vey of learning methods. ACM Computing Surveys
804	(CSUR), 50(2):1-35.
805	Georgi Karadzhov, Tom Stafford, and Andreas Vlachos.
806	2023. Delidata: A dataset for deliberation in multi-
807	party problem solving. <i>Proceedings of the ACM on</i>
808	Human-Computer Interaction, 7(CSCW2):1–25.
	Theshim Whatever D' to a De with the De De the
809	Ibrahim Knebour, Richard Brutti, Indrani Dey, Rachel
510	Dickier, Keisey Sikes, Kenneth Lai, Mariah Bradford,
811	Brittany Cates, Paige Hansen, Changsoo Jung, et al.
812	2024a. When text and speech are not enough: A
813	multimodal dataset of collaboration in a situated task.
814	Journal of Open Humanities Data, 10(1).

Ibrahim Khalil Khebour, Kenneth Lai, Mariah Bradford, Yifan Zhu, Richard A. Brutti, Christopher Tam, Jingxuan Tu, Benjamin A. Ibarra, Nathaniel Blanchard, Nikhil Krishnaswamy, and James Pustejovsky. 2024b. Common ground tracking in multimodal dialogue. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 3587–3602, Torino, Italia. ELRA and ICCL.

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

- Gary Klein, Paul J Feltovich, Jeffrey M Bradshaw, and David D Woods. 2005. Common ground and coordination in joint activity. *Organizational simulation*, 53:139–184.
- Tomasz Korbak, Hady Elsahar, Germán Kruszewski, and Marc Dymetman. 2022. On reinforcement learning and distribution matching for fine-tuning language models with no catastrophic forgetting. *Advances in Neural Information Processing Systems*, 35:16203–16220.
- Matthew A Koschmann. 2016. The communicative accomplishment of collaboration failure. *Journal of Communication*, 66(3):409–432.
- Nikhil Krishnaswamy, William Pickard, Brittany Cates, Nathaniel Blanchard, and James Pustejovsky. 2022. The voxworld platform for multimodal embodied agents. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1529– 1541.
- Nathan Lambert, Valentina Pyatkin, Jacob Daniel Morrison, Lester James Validad Miranda, Bill Yuchen Lin, Khyathi Raghavi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hanna Hajishirzi. 2024. RewardBench: Evaluating reward models for language modeling. *ArXiv*, abs/2403.13787.
- Peter R Lewis and Ştefan Sarkadi. 2024. Reflective artificial intelligence. *Minds and Machines*, 34(2):1–30.
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023a. Camel: Communicative agents for" mind" exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008.
- Zekun Li, Wenhu Chen, Shiyang Li, Hong Wang, Jing Qian, and Xifeng Yan. 2023b. Controllable dialogue simulation with in-context learning. *Preprint*, arXiv:2210.04185.
- Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su, Jun Zhu, et al. 2024. Llava-plus: Learning to use tools for creating multimodal agents. In *European Conference on Computer Vision*, pages 126– 142. Springer.
- Ilya Loshchilov, Frank Hutter, et al. 2017. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101*, 5.

- 871 872
- 8
- 8
- 87

- 88
- 8

88

890 891 892

89

8 8 8

> 900 901 902

903 904 905

906 907 908

909

> 919 920

926

927

Shengyu Mao, Xiaohan Wang, Mengru Wang, Yong Jiang, Pengjun Xie, Fei Huang, and Ningyu Zhang. 2024. Editing personality for large language models. In CCF International Conference on Natural Language Processing and Chinese Computing, pages 241–254. Springer.

- Pedro Henrique Martins, Zita Marinho, and André F. T. Martins. 2020. Sparse text generation. *Preprint*, arXiv:2004.02644.
- Yu Meng, Mengzhou Xia, and Danqi Chen. 2024. Simpo: Simple preference optimization with a reference-free reward. *Preprint*, arXiv:2405.14734.
- Rémi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland, Zhaohan Daniel Guo, Yunhao Tang, Matthieu Geist, Thomas Mesnard, Andrea Michi, et al. 2023. Nash learning from human feedback. *arXiv preprint arXiv:2312.00886*.
- Abhijnan Nath, Shadi Manafi Avari, Avyakta Chelle, and Nikhil Krishnaswamy. 2024a. Okay, let's do this! modeling event coreference with generated rationales and knowledge distillation. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 3931–3946.
- Abhijnan Nath, Changsoo Jung, Ethan Seefried, and Nikhil Krishnaswamy. 2024b. Simultaneous reward distillation and preference learning: Get you a language model who can do both. *arXiv preprint arXiv:2410.08458*.
- Abhijnan Nath, Videep Venkatesha, Mariah Bradford, Avyakta Chelle, Austin Youngren, Carlos Mabrey, Nathaniel Blanchard, and Nikhil Krishnaswamy. 2024c. "any other thoughts, hedgehog?" linking deliberation chains in collaborative dialogues. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 5297–5314.
- Harri Oinas-Kukkonen and Marja Harjumaa. 2009. Persuasive systems design: Key issues, process model, and system features. *Communications of the association for Information Systems*, 24(1):28.

OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Conn, Andrew Tulloch, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka

Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, 928 Barret Zoph, Behrooz Ghorbani, Ben Leimberger, 929 Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin 930 Zweig, Beth Hoover, Blake Samic, Bob McGrew, 931 Bobby Spero, Bogo Giertler, Bowen Cheng, Brad 932 Lightcap, Brandon Walkin, Brendan Quinn, Brian 933 Guarraci, Brian Hsu, Bright Kellogg, Brydon East-934 man, Camillo Lugaresi, Carroll Wainwright, Cary 935 Bassin, Cary Hudson, Casey Chu, Chad Nelson, 936 Chak Li, Chan Jun Shern, Channing Conger, Char-937 lotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, 938 Chong Zhang, Chris Beaumont, Chris Hallacy, Chris 939 Koch, Christian Gibson, Christina Kim, Christine 940 Choi, Christine McLeavey, Christopher Hesse, Clau-941 dia Fischer, Clemens Winter, Coley Czarnecki, Colin 942 Jarvis, Colin Wei, Constantin Koumouzelis, Dane 943 Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, 944 David Carr, David Farhi, David Mely, David Robin-945 son, David Sasaki, Denny Jin, Dev Valladares, Dim-946 itris Tsipras, Doug Li, Duc Phong Nguyen, Duncan 947 Findlay, Edede Oiwoh, Edmund Wong, Ehsan As-948 dar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, 949 Eric Kramer, Eric Peterson, Eric Sigler, Eric Wal-950 lace, Eugene Brevdo, Evan Mays, Farzad Khorasani, 951 Felipe Petroski Such, Filippo Raso, Francis Zhang, 952 Fred von Lohmann, Freddie Sulit, Gabriel Goh, 953 Gene Oden, Geoff Salmon, Giulio Starace, Greg 954 Brockman, Hadi Salman, Haiming Bao, Haitang 955 Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, 956 Heather Whitney, Heewoo Jun, Hendrik Kirchner, 957 Henrique Ponde de Oliveira Pinto, Hongyu Ren, 958 Huiwen Chang, Hyung Won Chung, Ian Kivlichan, 959 Ian O'Connell, Ian O'Connell, Ian Osband, Ian Sil-960 ber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya 961 Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, 962 Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub 963 Pachocki, James Aung, James Betker, James Crooks, 964 James Lennon, Jamie Kiros, Jan Leike, Jane Park, 965 Jason Kwon, Jason Phang, Jason Teplitz, Jason 966 Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Var-967 avva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui 968 Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, 969 Joaquin Quinonero Candela, Joe Beutler, Joe Lan-970 ders, Joel Parish, Johannes Heidecke, John Schul-971 man, Jonathan Lachman, Jonathan McKay, Jonathan 972 Uesato, Jonathan Ward, Jong Wook Kim, Joost 973 Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, 974 Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, 975 Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai 976 Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin 977 Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, 978 Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, 979 Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle 980 Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lau-981 ren Workman, Leher Pathak, Leo Chen, Li Jing, Lia 982 Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lil-983 ian Weng, Lindsay McCallum, Lindsey Held, Long 984 Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kon-985 draciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, 986 Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine 987 Boyd, Madeleine Thompson, Marat Dukhan, Mark 988 Chen, Mark Gray, Mark Hudnall, Marvin Zhang, 989 Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, 990 Max Johnson, Maya Shetty, Mayank Gupta, Meghan 991

Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunninghman, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. 2024. Gpt-4o system card. Preprint, arXiv:2410.21276.

992

993

995

1001

1002

1003

1004

1007

1010

1012

1013

1014

1015

1017

1019

1020

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

- Alizée Pace, Jonathan Mallinson, Eric Malmi, Sebastian Krause, and Aliaksei Severyn. 2024. West-of-n: Synthetic preferences for self-improving reward models. *Preprint*, arXiv:2401.12086.
- Eric Pacuit. 2017. Neighborhood semantics for modal logic. Springer.
 - Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White. 2024. Smaug: Fixing failure modes of preference optimisation with dpo-positive. *Preprint*, arXiv:2402.13228.
- Keyu Pan and Yawen Zeng. 2023. Do llms possess a personality? making the mbti test an amazing eval-

uation for large language models. *arXiv preprint arXiv:2307.16180*.

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1067

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

- Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar, and Jason Weston. 2024. Iterative reasoning preference optimization. *arXiv preprint arXiv:2404.19733*.
- Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. 2019. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *Preprint*, arXiv:1910.00177.
- David Premack and Guy Woodruff. 1978. Does the chimpanzee have a theory of mind? *Behavioral and brain sciences*, 1(4):515–526.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR.
- Rafael Rafailov, Joey Hejna, Ryan Park, and Chelsea Finn. 2024a. From r to q* : Your language model is secretly a q-function. *arXiv preprint arXiv:2404.12358*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024b. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.
- Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 8689–8696.
- Corby Rosset, Ching-An Cheng, Arindam Mitra, Michael Santacroce, Ahmed Awadallah, and Tengyang Xie. 2024. Direct nash optimization: Teaching language models to self-improve with general preferences. *ArXiv*, abs/2404.03715.
- Stuart J Russell and Peter Norvig. 2016. Artificial intelligence: a modern approach. Pearson.
- Maarten Sap, Ronan Le Bras, Daniel Fried, and Yejin Choi. 2022. Neural theory-of-mind? on the limits of social intelligence in large lms. In *Proceedings of the* 2022 Conference on Empirical Methods in Natural Language Processing, pages 3762–3780.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Lior Shani, Aviv Rosenberg, Asaf Cassel, Oran Lang, Daniele Calandriello, Avital Zipori, Hila Noga, Orgad Keller, Bilal Piot, Idan Szpektor, Avinatan Hassidim, Yossi Matias, and Rémi Munos. 2024. Multiturn reinforcement learning from preference human feedback. *Preprint*, arXiv:2405.14655.

- 1107 1108 1109 1110 1111 1112 1113 1114 1115 1116 1117 1118 1119 1120 1121 1122 1123 1124 1125 1126 1127 1128 1129 1130 1131 1132 1133 1134 1135 1136 1137 1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152
- 1153 1154
- 1155 1156
- 1157 1158
- 1159
- 1160
- 1161

- Margrethe Sønneland. 2019. Friction in fiction: A study of the importance of open problems for literary conversations. L1-Educational Studies in Language and *Literature*, 19:1–28.
- Robert Stalnaker. 2002. Common ground. Linguistics and philosophy, 25(5/6):701-721.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. Advances in Neural Information Processing Systems, 33:3008-3021.
- Robert I Sutton and Huggy Rao. 2024. The friction project: How smart leaders make the right things easier and the wrong things harder. Random House.
- Corbyn Terpstra, Ibrahim Khebour, Mariah Bradford, Brett Wisniewski, Nikhil Krishnaswamy, and Nathaniel Blanchard. 2023. How good is automatic segmentation as a multimodal discourse annotation aid? In Proceedings of the 19th Joint ACL-ISO Workshop on Interoperable Semantics (ISA-19), pages 75-81.
- Jingxuan Tu, Kyeongmin Rim, Bingyang Ye, Kenneth Lai, and James Pustejovsky. 2024. Dense paraphrasing for multimodal dialogue interpretation. Frontiers in artificial intelligence, 7:1479905.
- Amos Tversky. 1969. Intransitivity of preferences. Psychological review, 76(1):31.
- Tomer Ullman. 2023. Large language models fail on trivial alterations to theory-of-mind tasks. arXiv preprint arXiv:2302.08399.
- Johan van Benthem, David Fernández-Duque, and Eric Pacuit. 2014. Evidence and plausibility in neighborhood structures. Annals of Pure and Applied Logic, 165(1):106-133.
- Michael Völske, Martin Potthast, Shahbaz Syed, and Benno Stein. 2017. Tl; dr: Mining reddit to learn automatic summarization. In Proceedings of the Workshop on New Frontiers in Summarization, pages 59-63.
- Peter C Wason. 1968. Reasoning about a rule. Quarterly journal of experimental psychology, 20(3):273-281.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. Preprint, arXiv:2201.11903.
- Sarah Wiegreffe, Jack Hessel, Swabha Swayamdipta, Mark Riedl, and Yejin Choi. 2022. Reframing human-ai collaboration for generating free-text explanations. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 632-658.

Sarah Wiegreffe, Ana Marasović, and Noah A Smith. 2021. Measuring association between labels and free-text rationales. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language *Processing*, pages 10266–10284.

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

- Christian Wirth, Riad Akrour, Gerhard Neumann, and Johannes Fürnkranz. 2017. A survey of preferencebased reinforcement learning methods. The Journal of Machine Learning Research, 18(1):4945–4990.
- Jing Xu, Andrew Lee, Sainbayar Sukhbaatar, and Jason Weston. 2023. Some things are more cringe than others: Preference optimization with the pairwise cringe loss. arXiv preprint arXiv:2312.16682.
- Rui Yang, Ruomeng Ding, Yong Lin, Huan Zhang, and Tong Zhang. 2024. Regularizing hidden states enables learning generalizable reward model for llms. Preprint, arXiv:2406.10216.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models, 2023. URL https://arxiv. org/pdf/2305.10601. pdf.
- Fanghua Ye, Jarana Manotumruksa, and Emine Yilmaz. 2022. MultiWOZ 2.4: A multi-domain task-oriented dialogue dataset with essential annotation corrections to improve state tracking evaluation. In Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue, pages 351–360, Edinburgh, UK. Association for Computational Linguistics.
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2024. Self-rewarding language models. arXiv preprint arXiv:2401.10020.
- Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara, Raghav Gupta, Jianguo Zhang, and Jindong Chen. 2020. MultiWOZ 2.2 : A dialogue dataset with additional annotation corrections and state tracking baselines. In Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI, pages 109-117, Online. Association for Computational Linguistics.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. arXiv preprint arXiv:2205.01068.
- Xuan Zhang, Chao Du, Tianyu Pang, Qian Liu, Wei Gao, and Min Lin. 2024. Chain of preference optimization: Improving chain-of-thought reasoning in llms. arXiv preprint arXiv:2406.09136.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging LLM-as-a-judge with MT-Bench and Chatbot Arena. In NeurIPS Datasets and Benchmarks Track.

1232

1233

1234

1236 1237

1238

1239

1240

1241

1219

Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. 2008. Maximum entropy inverse reinforcement learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA.

A Functional Definition and Samples of Naturally-Occurring Friction

The functional operative definition of friction in collaborative contexts that we used is given below. This definition was used when annotating the WTD for naturally-occurring frictive utterances, and used to construct the prompt for friction intervention generation, following work by Oinas-Kukkonen and Harjumaa (2009) and Karadzhov et al. (2023).

FUNCTIONAL DEFINITION OF FRICTION IN COLLABORATIVE TASKS

Frictive interventions in this setting acts as indirect persuasion (Oinas-Kukkonen and Harjumaa, 2009) where participants are passively prompted to reevaluate their beliefstates and incorrect assumptions or propositions, in light of incoming goal-specific evidence. We define productive or positive friction as interventions that act as indirect persuasion: agentic interventions that prompt participants to reevaluate their beliefs and assumptions about the task state, primarily but not solely, in light of incoming evidence (say, occurrences in the physical environment or a correct "declaration" previously occurring in the dialogue that any participant missed) that negates their preconceived notions about the state of the task. We call this indirect persuasion since we do not want our friction agent to directly offer hints about the task and thereby biasing task performance or negatively affecting the deliberation process that is proven to beneficial for successful task completion in reasoningbased, collaborative tasks (Karadzhov et al., 2023).

Table 3 shows a sample friction annotation and training sample from the Weights Task Dataset, consisting of the dialogue history x, GPT-40-identified frictive state ϕ , rationale, and preferred and dispreferred friction interventions f_w and f_l . Because the WTD is a multimodal dataset, the transcriptions we use are enriched using *dense paraphrasing* (Tu et al., 2024), a textual enrichment technique that uses the multimodal channels to de-

contextualize referents and in this case transform contextually-dependent phrasings such as demonstratives to explicit denotations of the content. For example, under dense paraphrasing, "seems like *these* might be about the same" while the speaker in the video is pointing to the red and blue blocks becomes "seems like *red block, blue block* might be about the same." The dense paraphrased utterances are included as part of the publicly-available WTD (Khebour et al., 2024a,b). 1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1283

B Frictive-state conditioning and RLHF

In its simplest formulation within Chain-of-Thought (CoT) settings (Wei et al., 2023), the friction agent is modeled as a policy distribution π_f that sequentially generates frictive states, sampling $\phi_i \sim \pi_f(\cdot \mid x, \phi_1, \dots, \phi_{i-1})$, and ultimately producing the final friction intervention $f \sim \pi_f(\cdot \mid$ $x, \phi_1, \ldots, \phi_n$). Here, x represents the dialogue history, f denotes the intervention, and ϕ consists of sequentially sampled frictive state tokens, analogous to "thoughts" in standard CoT-based reasoning frameworks (Yao et al., 2023). Unlike standard CoT-based alignment, which relies on selfrewarding strategies, we frame friction agent alignment within preference-based RL (PbRL; Wirth et al. (2017)). Prior work (Zhang et al., 2024) shows CoT frameworks benefit significantly from contrastive signals in preference learning.

In this setting, we define the human preference probability $\mathcal{P}(f \succ \phi)$ as the probability that an expert annotator would prefer f over maintaining the frictive state ϕ , given prior dialogue history, x. The key insight is that to retrieve the optimal policy π_f^* , we can leverage established methods from RLHF and PbRL by formulating the problem as a KL-divergence constrained minimum relative entropy optimization (Ziebart et al., 2008), a wellknown approach with a closed-form solution (Peng et al., 2019).

$$J^*_{\rm RLHF}(\pi_f) =$$
 1281

$$\max_{\pi_f} \mathbb{E}_{f \sim \pi_f} \left[\mathcal{P}(f \succ \phi \mid x) \right] -$$
 12

$$\beta D_{\mathrm{KL}}(\pi_f \parallel \pi_{\mathrm{ref}}). \tag{3}$$

This formulation (Eq. 3)—where J^*_{RLHF} enforces1284a KL-based "soft"-constraint on the parametric1285form of π^*_f wrt the reference policy π_{ref} —provides1286crucial tradeoffs between training stability and bal-1287ancing exploration vs exploitation. Specifically,1288

Field	Content
Dialogue History (x)	P1: i guess if red block red one's ten grams P2: we got red ten P1: seems like red block, blue block might be about the same P3: i would agree yeah so blue block one's ten P3: Alright let's see if we can find a twenty P3: Too heavy so P2: Way too heavy P2: this is a sensitive scale P2: Looks like about twenty P1: that's looking pretty even P3: Alright let's see if we can find a thirty P1: so yellow block one is noticeably heavier than P2: probably yellow block big sucker P1: the purple is no ne P1: making sure that purple block didn't have the weight at the bottom P2: it's just stuff written at the bottom that's a so red block, green block's a ten and a twenty right now right that's looking P2: Well P2: red block, blue block, green block's a look, purple block check that purple block's not also a twenty P1: yeah it looks a little P2: cause it um just purple block one etner P3: is blue block one a twenty P2: ok so purple block, green block, yellow block, purple block're only in increments of ten purple: if red block, blue block, green block, yellow block seems like the thirty takes it past but P2: it's so sensitive P2: if red block, blue block, green block, yellow block, purple block're only in increments of ten purple block has to be
Frictive state (ϕ)	P2 initially identified the red and blue blocks as both 10 grams and has speculated about the green at 20 grams, but is uncertain about the actual weights of the yellow and purple blocks.
Rationale	P2 suggests that the red, green, and yellow blocks are all in increments of ten. Encourage a double-check.
Preferred Friction (f_w)	Since the purple block seems heavier and we're unsure about its exact weight, should we double-check the increments of ten assumption? Maybe the purple block doesn't fit this pattern.
Dispreferred Friction (f_l)	You know, the purple block being heavier might actually mean the blocks aren't increasing consistently at all. What if the increments are random, like 10, 15, or 25 grams, and we're forcing a pattern that isn't there?

Table 3: A transcribed, sparse collaborative dialogue from the Weights Task Dataset (Khebour et al., 2024a) with frictive states and friction interventions. Preferred and dispreferred friction interventions are shown at the bottom. Positive friction interventions prompt participants to reevaluate their assumptions with frictive states (evolution of beliefs and rationales) providing indirect hints and directions. Here, P2's uncertainty about the green block and assumption of weight increment by 10g is addressed by the positive friction. In contrast, the dispreferred intervention introduces randomness, instigating the group to abandon structured reasoning.

 J_{RLHF}^* ensures that π_f^* retrieves the best possible 1289 preference probabilities for its generated interven-1290 tions f, as assigned by $\mathcal{P}(f \succ \phi)$, whether over 1291 the distribution of preferences encoded in an of-1292 fline dataset (Rafailov et al., 2024b) or from online sampling $\sim \pi_f$ during training (Schulman et al., 2017) while being distributionally close to an 1295 "already-good" imitator, the Supervised-finetuned 1296 (SFT) reference model (Hussein et al., 2017). No-1297 tice that unlike standard RLHF, we formulate $J^*_{\mathbf{RLHF}}$ such that π^*_f takes the form $\pi^*_f(\cdot \mid \phi, x)$ 1299 is explicitly conditioned on the frictive state ϕ , 1300 apart from x. This is intentional since we hy-1301 pothesize that an ideal friction agent does not in-1302 tervene arbitrarily, causing distraction in collaborative tasks and is conditioned to resolve the lack 1304 of common ground thereof between human collaborators, by definition—as observed in ϕ . While prior work (Choi et al., 2024; Zhang et al., 2024) explores preference alignment in LLMs in such 1308 CoT-conditioned scenarios, we provide a more principled approach to proving the existence and 1310 the uniqueness of π_f^* that J_{RLHF}^* seeks to retrieve. Mathematically, 1312

$$\pi_f^* = \frac{\pi_{\text{ref}} \exp(\beta^{-1} \mathcal{P}(f \succ \phi | x))}{Z^*(\phi, x)}, \qquad (4)$$

where $Z^* = \sum_{f'} \pi_{\rm ref} \exp(\beta^{-1} \mathcal{P}(f' \succ \phi | x))$ is 1314 the partition function which is fixed and does not 1315 depend on f and can be safely ignored in the op-1316 timization of J_{RLHF}^* (Rafailov et al., 2024b). See 1317

1313

Appendix B.1 and Equation (8) for the full-proof and optimal policy form respectively.

1318

1319

1320

1321

1326

1327

1331

1332

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

B.1 Existence and uniqueness of the optimal friction intervention policy

In order to derive an empirical offline (supervised) 1322 preference learning loss from the complicated two-1323 staged FAAF-alignment objective defined in Equa-1324 tion (2), we use a divide and conquer approach-1325 our core insight here is to express the preference of interventions conditioned on the frictive states in terms of two mutually supportive "twin" policies. 1328 As such, we first derive the inner maximization 1329 loop of Eq. 2 to get an analytical expression of the 1330 optimal frictive intervention policy, π_f^* as shown in the proof for Eq. 8. However, we observe that π_f^* in its analytical form is not fully expressive since it 1333 does *not* contain the optimal frictive-state policy π^*_{ϕ} 1334 term. Therefore, we propose a novel method to de-1335 rive π_{ϕ}^* using a Lagrangian formulation. We show 1336 the detailed derivation for this part in Appendix C including supporting results from Lemma 3 and Lemma 6.

This above result is one of our main contributions since it lets us express the preference for any intervention f_1 over f_2 analytically in terms of **both** the optimal friction intervention policy $(\pi_f^*(\cdot \mid \phi, x))$ and the optimal frictivestate policy $(\pi_{\phi}^*(\cdot|x))$. Finally, this core result is used to propose a straightforward supervised (ℓ_2) objective-similar in spirit to IPO (Azar et al., 2024)—that empirically regresses the predicted preference expression derived from $\pi_f^*(\cdot \mid \phi, x)$ and $\pi_{\phi}^{*}(\cdot|x)$ to the observed relative preferences $p(f_1 \succ f_2 \mid x)$ (relative to ϕ), assuming access to a large-enough preference-annotated dataset of frictive interventions. Notably, this objective is optimized by a *single* parametrized policy that leverages the inherent expressivity of LLMs with billions of parameters.

1348

1349

1350

1351

1353

1354

1355

1357

1358

1359

1361

1362

1365

1366

1368

1370

1371 1372

1374

1376

1378

1379

1380

1385

1386

In particular, this FAAF objective formulation avoids some of the policy degeneracy issues that popular supervised "offline" alignment algorithms like Direct Preference Optimization (DPO) (Rafailov et al., 2024b) face due to its unbounded rewards. Additionally, unlike Fisch et al. (2024), our regression objective works directly on preference labels and does not require an external reward model in avoiding such degeneracies. Finally, we also prove that FAAF-trained policies are unique solutions in the policy space in Theorem 2.

For completeness, we first prove the existence of the optimal friction/frictive intervention policy that solves the inner maximization of our two-part minimax objective. The structural solution to this objective is well-studied in the RL/control-theory literature including popular frameworks in preference alignment in LLMs (Ziebart et al., 2008; Peng et al., 2019; Rafailov et al., 2024b; Azar et al., 2024) as well as Chain-of-Thought (CoT)-based preference alignment frameworks (Choi et al., 2024). We show how it specifically applies to our unique parametrization. Our proof follows similar logic as Azar et al. (2024). Let us recall two-part minimax objective (Eq. 2) for clarity here:

1382
$$J_{\mathsf{FAAF}}^{*} = \min_{\pi_{\phi}} \max_{\pi_{f}} \mathbb{E} \underset{\substack{x \sim \rho \\ \phi \sim \pi_{\phi}(\cdot | x) \\ f \sim \pi_{f}(\cdot | \phi, x)}} x_{\phi \sim \pi_{\phi}(\cdot | x)} \left[\mathcal{P}(f \succ \phi \mid x) -\beta D_{\mathsf{KL}}(\pi_{f} \parallel \pi_{\mathsf{ref}} \mid \phi, x) +\beta D_{\mathsf{KL}}(\pi_{\phi} \parallel \pi_{\mathsf{ref}} \mid x) \right].$$
(5)

For fixed π_{ϕ} , the inner maximization reduces to our regularized objective:

1387
$$\mathcal{L}_{\beta}(\pi_f) =$$

1387
$$\mathcal{L}_{\beta}(\pi_{f}) = \mathbb{E}_{f \sim \pi_{f}}[p(f \succ \phi | x)] -$$
1388
$$\beta D_{\mathrm{KL}}(\pi_{f} \parallel \pi_{\mathrm{ref}} | \phi, x),$$
1389
$$= \sum_{f} \pi_{f}(f | \phi, x) p(f \succ \phi | x) -$$
1390
$$\beta D_{\mathrm{KL}}(\pi_{f} \parallel \pi_{\mathrm{ref}} | \phi, x), \quad (6)$$

where $f \in \mathcal{F}$ is from a finite friction token alpha-1391 bet \mathcal{F} , $p(f \succ \phi | x)$ maps elements of \mathcal{F} to the 1392 utility of generating a frictive intervention f de-1393 fined as the preference of f over the frictive-state 1394 ϕ , given context $x, \beta \in \mathbb{R}^*_+$ is a strictly positive 1395 real number, and π_f, π_{ref} are conditional proba-1396 bility distributions. In particular, notice that the 1397 conditional probability distribution $\pi_f(f|\phi, x)$ can 1398 be identified as a positive real function satisfying: 1399

$$\sum_{f} \pi_{f}(f|\phi, x) = 1.$$
 (7) 140

1401

1402

Now, if we define the optimal friction intervention policy π_f^* as:

$$\pi_f^*(f|\phi, x) = \frac{\pi_{\text{ref}}(f|\phi, x) \exp(\beta^{-1}p(f \succ \phi|x))}{Z^*(\phi, x)},$$
(8) 1403

1, where $Z^*(\phi, x)$ recalling Eq. =1404 $\sum_{f'} \pi_{\text{ref}}(f'|\phi, x) \exp(\beta^{-1}p(f' \succ \phi|x)), \text{ then,}$ 1405 under the previous definitions, we have: 1406

$$\pi_f^* = \operatorname*{arg\,max}_{\pi_f} \mathcal{L}_\beta(\pi_f) \tag{9}$$
 1407

$$\begin{aligned} & \frac{\mathcal{L}_{\beta}(\pi_{f})}{\beta} = \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \frac{p(f \succ \phi|x)}{\beta} - D_{\mathrm{KL}}(\pi_{f} \parallel \pi_{\mathrm{ref}}|\phi, x), \\ &= \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \frac{p(f \succ \phi|x)}{\beta} - \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \log\left(\frac{\pi_{f}(f|\phi, x)}{\pi_{\mathrm{ref}}(f|\phi, x)}\right), \\ &= \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \left(\frac{p(f \succ \phi|x)}{\beta} - \log\left(\frac{\pi_{f}(f|\phi, x)}{\pi_{\mathrm{ref}}(f|\phi, x)}\right)\right), \\ &= \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \left(\log(\exp(\beta^{-1}p(f \succ \phi|x))) - \log\left(\frac{\pi_{f}(f|\phi, x)}{\pi_{\mathrm{ref}}(f|\phi, x)}\right)\right), \\ &= \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \log\left(\exp(\beta^{-1}p(f \succ \phi|x))\frac{\pi_{\mathrm{ref}}(f|\phi, x)}{\pi_{f}(f|\phi, x)}\right), \\ &= \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \log\left(\frac{\pi_{\mathrm{ref}}(f|\phi, x) \exp(\beta^{-1}p(f \succ \phi|x))}{\pi_{f}(f|\phi, x)}\right), \\ &= \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \log\left(\frac{\pi_{\mathrm{ref}}(f|\phi, x) \exp(\beta^{-1}p(f \succ \phi|x))}{\pi_{f}(f|\phi, x)}\right), \\ &= \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \log\left(\frac{\pi_{\mathrm{ref}}(f|\phi, x) \exp(\beta^{-1}p(f \succ \phi|x))}{\pi_{f}(f|\phi, x)}\right), \\ &= \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \log\left(\frac{\pi_{\mathrm{ref}}(f|\phi, x) \exp(\beta^{-1}p(f \succ \phi|x))}{Z^{*}(\phi, x)}\right), \\ &= \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \log\left(\frac{\pi_{\mathrm{ref}}(f|\phi, x)}{\pi_{f}(f|\phi, x)}\right) + \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \log Z^{*}(\phi, x), \\ &= \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \log\left(\frac{\pi_{\mathrm{ref}}^{*}(f|\phi, x)}{\pi_{f}(f|\phi, x)}\right) + \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) \log Z^{*}(\phi, x), \\ &= -D_{\mathrm{KL}}(\pi_{f} \parallel \pi_{f}^{*}) + \log Z^{*}(\phi, x), \quad (\text{using normalization } \sum_{f \in \mathcal{F}} \pi_{f}(f|\phi, x) = 1) \end{aligned}$$

By definition of the KL divergence, we know that $\pi_f^* = \arg \max_{\pi_f} \left[-D_{\text{KL}}(\pi_f \parallel \pi_f^*) \right]$ and as:

$$-D_{\mathrm{KL}}(\pi_f \parallel \pi_f^*) = \frac{\mathcal{L}_{\beta}(\pi_f)}{\beta} - \log Z^*(\phi, x)$$

where $\log Z^*(\phi, x)$ is the partition function (Peng et al., 2019; Rafailov et al., 2024b) and has no dependency on π_f and $\beta \in \mathbb{R}^*_+$ is a strictly positive real number. Therefore, the argmax of $-D_{\mathrm{KL}}(\pi_f \parallel \pi_f^*)$ coincides with that of $\mathcal{L}_{\beta}(\pi_f)$, concluding the proof. \Box

Lemma 1 (Value of Inner Maximization). When 1408 substituting the optimal friction intervention policy 1409 π_f^* , as derived in Eq. 8, into Eq. 5, the objective in 1410 Eq. 5 reduces to: 1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

$$J_{\mathsf{FAAF}}^* = \min_{\pi_{\phi}} \mathbb{E}_{x \sim \rho, \phi \sim \pi_{\phi}(\cdot|x)} [\beta \log(Z^*(\phi, x)) + \beta D_{KL}(\pi_{\phi} || \pi_{ref} |x)]$$
(10)

Proof. Substituting π_f^* into the KL divergence term:

$$D_{KL}(\pi_f^* || \pi_{ref} | \phi, x)$$

$$= \mathbb{E}_{f \sim \pi_f^*} \bigg[\log(\pi_f^*(f | \phi, x)) - \log(\pi_{ref}(f | \phi, x)) \bigg]$$

$$= \mathbb{E}_{f \sim \pi_f^*} \bigg[\frac{p(f \succ \phi | x)}{\beta} - \log(Z^*(\phi, x)) \bigg]$$
(11)

The original objective becomes:

$$p(f \succ \phi | x) - \beta \mathbb{E}_{f \sim \pi_f^*} \left[\frac{p(f \succ \phi | x)}{\beta} - \log(Z^*(\phi, x)) \right]$$
$$= \beta \log(Z^*(\phi, x))$$
(12)

The result follows by substituting this value back into the full objective.

С **Derivation of Optimal Frictive State Policy**

We begin with the reduced objective function after solving the inner maximization as shown in Lemma 1.

$$J_{\mathsf{FAAF}}^* = \min_{\pi_{\phi}} \mathbb{E}_{x \sim \rho, \phi \sim \pi_{\phi}(\cdot|x)} \left[\beta \log(Z^*(\phi, x)) + \beta D_{KL}(\pi_{\phi} || \pi_{\mathrm{ref}} |x) \right]$$
(13)

The Kullback-Leibler divergence term expands as follows:

1424
$$D_{KL}(\pi_{\phi}||\pi_{\mathrm{ref}}|x) = \mathbb{E}_{\phi \sim \pi_{\phi}} \left[\log \frac{\pi_{\phi}(\phi|x)}{\pi_{\mathrm{ref}}(\phi|x)} \right]$$
(14)

Substituting this back into our objective:

$$J_{\mathsf{FAAF}}^* = \min_{\pi_{\phi}} \mathbb{E}_{\phi \sim \pi_{\phi}} \Big[\beta \log(Z^*(\phi, x)) + 1426 \Big] \Big]$$

$$\beta \log(\pi_{\phi}(\phi|x)) - \beta \log(\pi_{\text{ref}}(\phi|x)) \Big]$$
(15)

1425

Since π_{ϕ} must be a valid probability distribu-1428 tion satisfying $\sum_{\phi} \pi_{\phi}(\phi|x) = 1$, we introduce a 1429 Lagrange multiplier λ and define the correspond-1430 ing Lagrangian function to derive the optimality 1431 conditions:

$$L(\pi_{\phi}) = \mathbb{E}_{\phi \sim \pi_{\phi}} \Big[\beta \log(Z^*(\phi, x)) +$$
 1433

$$\beta \log(\pi_{\phi}(\phi|x)) - \beta \log(\pi_{\text{ref}}(\phi|x))] +$$
 1434

$$\lambda\left(1-\sum_{\phi}\pi_{\phi}(\phi|x)\right)$$
1433

$$= \sum_{\phi} \pi_{\phi}(\phi|x) \bigg[\beta \log(Z^*(\phi, x)) +$$
 143

$$\beta \log(\pi_{\phi}(\phi|x)) - \beta \log(\pi_{ref}(\phi|x)) \right]$$
143

$$+\lambda\left(1-\sum_{\phi}\pi_{\phi}(\phi|x)\right).$$
 (16) 143

Now, to find the optimal policy $\pi^*_{\phi}(\phi|x)$, we 1439 take the derivative of the Lagrangian with respect 1440 to $\pi_{\phi}(\phi|x)$ and equate it to zero: 1441

$$\frac{\delta L}{\delta \pi_{\phi}(\phi|x)} = \beta \log(Z^*(\phi, x)) +$$
 1442

$$\beta \frac{\delta}{\delta \pi_{\phi}} \left[\pi_{\phi}(\phi|x) \log(\pi_{\phi}(\phi|x)) \right]$$
1443

$$-\beta \log(\pi_{ref}(\phi|x)) - \lambda = 0.$$
 (17) 1444

1450

From the standard functional derivative of 1449
entropy
$$\frac{\delta}{\delta \pi_{\phi}} \left[\pi_{\phi}(\phi|x) \log(\pi_{\phi}(\phi|x)) \right] = 1 + 1440$$

 $\log(\pi_{\phi}(\phi|x))$, we obtain: 1449

$$\beta \log(Z^*(\phi, x)) + \beta(1 + \log(\pi_{\phi}(\phi|x))) - 1444 \\ \beta \log(\pi_{ref}(\phi|x)) + \lambda = 0.$$
(18) 1444

Rearranging the terms:

$$\log(\pi_{\phi}(\phi|x)) = \log(\pi_{ref}(\phi|x)) -$$

$$\log(Z^*(\phi, x)) - \frac{\lambda}{\beta} - 1. \quad (19)$$
1452

(

1454

1455 1456

1457

1458

1464

1480

1481

1482

1483

1484

1485

1486

1487

Notice that without losing any generality, we can 1460 parametrize the above optimal frictive-state policy 1461 with any outcome f consistent with the structure 1462 in Eq. 22 as follows: 1463

$$\pi_{\phi}^{*}(f|x) = \frac{\pi_{\text{ref}}(f|x)}{Z(x)} e^{-\beta \log(Z^{*}(f,x))}.$$
 (23)

Taking the exponential on both sides:

tion, we define the normalization constant:

 $Z(x) = \sum_{\phi} \pi_{\text{ref}}(\phi|x) e^{-\beta \log Z^*(\phi,x)}.$

Thus, the optimal frictive-state policy is:

 $\pi_{\phi}^{*}(\phi|x) = \frac{\pi_{\text{ref}}(\phi|x)e^{-\beta \log Z^{*}(\phi,x)}}{Z(x)}.$

 $\pi_{\phi}(\phi|x) = e^{-1 - \frac{\lambda}{\beta}} \pi_{\text{ref}}(\phi|x) e^{-\log Z^{*}(\phi,x)}.$ (20)

To ensure $\pi_{\phi}(\phi|x)$ is a valid probability distribu-

(21)

(22)

1465 Note that although this formulation of the optimal frictive-state policy $(\pi_{\phi}^*(\phi|x))$ is an analyt-1466 ical solution to J^* from Eq. 13, we still need to 1467 represent $\pi_{\phi}^{*}(\phi|x)$ in terms of the optimal friction 1468 intervention policy, $\pi_f^*(\cdot \mid \phi, x)$ proposed in Eq. 8 1469 and the preference probabilities $p(f \succ \phi | x)$, the 1470 preference probability of the friction f over the 1471 frictive-state ϕ , given context x. This is crucial 1472 to derive the empirical FAAF optimization objec-1473 tive that can be used for standard offline learning. 1474 Therefore, to represent the $p(f \succ \phi | x)$ in terms of 1475 $\pi_f^*(\cdot \mid \phi, x)$, we take the logarithm of Eq. 8 on both 1476 sides and some algebra, we obtain: 1477

1478
$$\log(\pi_{f}^{*}(f|\phi, x)) = \frac{p(f \succ \phi|x)}{\beta} + \log(\pi_{\text{ref}}(f|\phi, x)) - \log(Z^{*}(\phi, x))$$

Multiplying both sides by β and rearranging terms, we obtain:

$$p(f \succ \phi | x) = \beta [\log(\pi_f^*(f | \phi, x)) - \log(\pi_{\text{ref}}(f | \phi, x)) + \log(Z^*(\phi, x))].$$
(24)

Similar to Munos et al. (2023), Azar et al. (2024), and Choi et al. (2024), we can apply the identity that $p(\phi \succ \phi | x) = \frac{1}{2}$ and substitute $f = \phi$ into the previous equation and derive:

$$\frac{1}{2} = \beta [\log(\pi_f^*(\phi | \phi, x)) - 1488]$$

 $\log(\pi_{\text{ref}}(\phi|\phi, x)) + \log(Z^*(\phi, x))].$ (25)1489

Solving for $\log(Z^*(\phi, x))$ gives:

$$\log(Z^*(\phi, x)) = \frac{1}{2\beta} -$$
1491

$$[\log(\pi_f^*(\phi|\phi, x)) - \log(\pi_{\text{ref}}(\phi|\phi, x))].$$
 (26) 1492

Substituting this back into Eq. 24 results in:

$$p(f\succ\phi|x)=\beta[\log(\pi_f^*(f|\phi,x))-1494$$

$$\log(\pi_{\mathrm{ref}}(f|\phi, x))$$
 1495

1490

1493

$$+ \frac{1}{2\beta} - (\log(\pi_f^*(\phi|\phi, x)) - 1496)$$

$$\log(\pi_{\rm ref}(\phi|\phi, x)))]$$
149

$$= \beta \log \left(\frac{\pi_f^*(f|\phi, x)}{\pi_{\text{ref}}(f|\phi, x)} \right) + \frac{1}{2}$$
149

$$-\beta \log\left(\frac{\pi_f^*(\phi|\phi, x)}{\pi_{\text{ref}}(\phi|\phi, x)}\right).$$
 (27) 149

The $\log\left(\frac{\pi_f^*(\phi|\phi,x)}{\pi_{\rm ref}(\phi|\phi,x)}\right)$ term in the above step is a 1500 self-referential term signifying the friction interven-1501 tion policy's $(\pi_f^*(\cdot \mid \phi, x))$ estimate of the frictive 1502 state given ϕ . However, this term does *not* provide much information on the regularized preference in 1504 terms of the frictive state policy. Recall that our outer minimization objective operates over $\pi_{\phi}(\cdot|x)$. Fortunately, we can use our results from Lemma 3 1507 and Lemma 6 to express Eq. 27 in terms of the op-1508 timal frictive state policy $\pi_{\phi}^{*}(\cdot|x).$ Therefore, from 1509 Lemma 6 we can express π_f^* and π_{ref} as follows: 1510 1511

For the optimal policy π_f^* :

$$\log(\pi_{f}^{*}(\phi|\phi, x)) = \log(\pi_{\phi}^{*}(\phi|x)) -$$

$$\log(\pi_{\phi}^{*}(f|x))$$
(28)
1513

For the reference policy π_{ref} :

$$\log(\pi_{\text{ref}}(\phi|\phi, x)) = \log(\pi_{\text{ref}}(\phi|x)) -$$
1515

$$\log(\pi_{\rm ref}(f|x))$$
 (29) 1516

1514

Now, substituting these expressions into Equa-1517

tion
$$(27)$$
, we get:

153

1536

1537

1519
$$p(f \succ \phi | x) = \beta \Big[\log(\pi_f^*(f | \phi, x)) - \log(\pi_{ref}(f | \phi, x)) + \log(\pi_{ref}(f | \phi, x)) \Big]$$

$$\log(\pi_{\mathrm{ref}}(f|\phi, x)) +$$

521
$$\frac{1}{2\beta} - \log(\pi_{\phi}^*(\phi|x)) +$$

1522
$$\log(\pi_{\phi}^*(f|x)) - \log(\pi_{\text{ref}}(\phi|x)) +$$

1523
$$\log(\pi_{\mathrm{ref}}(f|x))$$

$$= \beta \Big[\log(\pi_f^*(f|\phi, x)) -$$

1525
$$\log(\pi_{ref}(f|\phi, x)) + \frac{1}{126} - \left(\log(\pi_{\phi}^{*}(\phi|x)) - \log(\pi_{\phi}^{*}(\phi|x))\right)$$

1527
$$2\beta$$
 $(\log(n_{\phi}(\varphi|x)))$
 $\log(\pi_{ref}(\phi|x)) -$

$$\left(\log(\pi_{\phi}^*(f|x)) - \log(\pi_{\text{ref}}(f|x)))\right)\right]$$

1529
$$= \beta \left[\log \left(\frac{\pi_f^*(f|\phi, x)}{\pi_{\text{ref}}(f|\phi, x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f|x)}{\pi_{\text{ref}}(f|x)} \right) - \log \left(\frac{\pi_\phi^*(f|x)}{\pi_{\text{ref}}(f|x)} \right) - \log \left(\frac{\pi_\phi^*(f|x)}{\pi_{\text{ref}}(f|x)} \right) \right]$$

1531
$$\log\left(\frac{\pi_{\phi}^{*}(\phi|x)}{\pi_{\mathrm{ref}}(\phi|x)}\right) \right].$$
(30)

Now, replacing
$$f$$
 by f_1 in $p(f \succ \phi | x)$:

$$p(f_1 \succ \phi | x) = \beta \left[\log \left(\frac{\pi_f^*(f_1 | \phi, x)}{\pi_{\text{ref}}(f_1 | \phi, x)} \right) + \frac{1}{2} \left(\frac{\pi_f^*(f_1 | \phi, x)}{\pi_f^*(f_1 | x)} \right) \right]$$

1534
$$\frac{1}{2\beta} + \log\left(\frac{\pi_{\phi}^{*}(f_{1}|x)}{\pi_{\text{ref}}(f_{1}|x)}\right) - \log\left(\frac{\pi_{\phi}^{*}(\phi|x)}{\pi_{\text{ref}}(f_{1}|x)}\right) \right]$$
(3)

$$\log\left(\frac{\pi_{\phi}^*(\phi|x)}{\pi_{\rm ref}(\phi|x)}\right)$$
 (31)

Similarly, expressing f_2 in $p(f \succ \phi | x)$, we obtain:

1538
$$p(f_2 \succ \phi | x) = \beta \left[\log \left(\frac{\pi_f^*(f_2 | \phi, x)}{\pi_{\text{ref}}(f_2 | \phi, x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) - \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) - \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \frac{1}{2\beta} + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left(\frac{\pi_\phi^*(f_2 | x)}{\pi_{\text{ref}}(f_2 | x)} \right) + \log \left($$

1539
$$\frac{1}{2\beta} + \log\left(\frac{\pi_{\phi}^*}{\pi_{\text{ref}}}\right)$$

1540
$$\log\left(\frac{\pi_{\phi}^{*}(\phi|x)}{\pi_{\mathrm{ref}}(\phi|x)}\right)$$
(32)

Now, expressing $p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x)$, 1541 the relative preference probability of f_1 over f_2 1542 given ϕ and x, we observe that $\log\left(\frac{\pi_{\phi}^*(\phi|x)}{\pi_{\text{ref}}(\phi|x)}\right)$ 1543 terms cancel out and we derive: 1544

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$
1545

$$\beta \left[\log \left(\frac{\pi_f^*(f_1 | \phi, x)}{\pi_{\text{ref}}(f_1 | \phi, x)} \right) + \frac{1}{2\beta} + \frac{1}{2\beta} \right]$$
 1546

$$\log\left(\frac{\pi_{\phi}^{*}(f_{1}|x)}{\pi_{\mathrm{ref}}(f_{1}|x)}\right) - \log\left(\frac{\pi_{\phi}^{*}(\phi|x)}{\pi_{\mathrm{ref}}(\phi|x)}\right) \right]$$
1547

$$-\beta \left[\log \left(\frac{\pi_f^*(f_2|\phi, x)}{\pi_{\text{ref}}(f_2|\phi, x)} \right) + \frac{1}{2\beta} + \frac{1}{2\beta} + \frac{1}{2\beta} \right]$$
1548

$$\log\left(\frac{\pi_{\phi}^*(f_2|x)}{\pi_{\rm ref}(f_2|x)}\right) - \log\left(\frac{\pi_{\phi}^*(\phi|x)}{\pi_{\rm ref}(\phi|x)}\right)$$
1549

$$= \beta \left[\log \left(\frac{\pi_{f}^{*}(f_{1}|\phi, x)}{\pi_{\text{ref}}(f_{1}|\phi, x)} \right) - \log \left(\frac{\pi_{f}^{*}(f_{2}|\phi, x)}{\pi_{\text{ref}}(f_{2}|\phi, x)} \right) \right]$$
1550
$$\left(\pi_{f}^{*}(f_{1}|x) \right) = \left(\pi_{f}^{*}(f_{2}|x) \right)$$

$$+\log\left(\frac{\pi_{\phi}^{*}(f_{1}|x)}{\pi_{\mathrm{ref}}(f_{1}|x)}\right) - \log\left(\frac{\pi_{\phi}^{*}(f_{2}|x)}{\pi_{\mathrm{ref}}(f_{2}|x)}\right) \right]$$
(33)

This above result is one of our core contributions 1552 since it lets us express the relative preference of 1553 any friction intervention f_1 over f_2 given a frictive 1554 state (ϕ) analytically in terms of **both** the optimal 1555 friction intervention policy (($\pi_f^*(\cdot \mid \phi, x)$)) and the 1556 optimal frictive state policy $(\pi_{\phi}^*(\cdot|x))$: 1557

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \succ \phi | x) =$$

$$p(f_1 \succ \phi | x) - p(f_2 \vdash \phi | x) =$$

$$p(f_1 \vdash \phi | x) - p(f_2 \vdash \phi | x) =$$

$$p(f_1 \vdash \phi | x) - p(f_2 \vdash \phi | x) =$$

$$p(f_1 \vdash \phi | x) - p(f_2 \vdash \phi | x) =$$

$$p(f_1 \vdash \phi | x) - p(f_2 \vdash \phi | x) =$$

$$p(f_1 \vdash \phi | x) - p(f_2 \vdash \phi | x) =$$

$$p(f_1 \vdash \phi | x) - p(f_2 \vdash \phi | x) =$$

$$\beta \left[\log \left(\frac{\pi_f(J_1|\phi, x)}{\pi_{\text{ref}}(f_1|\phi, x)} \right) - \log \left(\frac{\pi_f(J_2|\phi, x)}{\pi_{\text{ref}}(f_2|\phi, x)} \right) + \log \left(\frac{\pi_\phi^*(f_1|x)}{\pi_{\text{ref}}(f_2|x)} \right) - \log \left(\frac{\pi_\phi^*(f_2|x)}{\pi_{\text{ref}}(f_2|x)} \right) \right]$$

$$(1559)$$

$$+ \log\left(\frac{\pi_{\phi}^{*}(f_{1}|x)}{\pi_{\mathrm{ref}}(f_{1}|x)}\right) - \log\left(\frac{\pi_{\phi}^{*}(f_{2}|x)}{\pi_{\mathrm{ref}}(f_{2}|x)}\right) \right]$$
(34)

Following a standard approach for empirical es-1561 timation of the LHS (Azar et al., 2024) in the above 1562 equation, one can learn both the optimal friction 1563 intervention policy π_f^* and the frictive-state pol-1564 icy π_{ϕ}^* using a trainable policy π_{θ} , parametrized 1565 with θ . The core insight here is to exploit the expressive nature of LLMs' hidden representations 1567 with billions of parameters to learn a single opti-1568 mal policy. A reasonable choice here is to train 1569 π_{θ} through an ℓ_2 loss (Fisch et al., 2024) that en-1570 forces the relative preference ordering between any 1571 pair of friction interventions (f_1, f_2) with implicit 1572 reward estimates from the RHS of Eq. 34. How-1573 ever, unlike (Fisch et al., 2024), our approach in 1574 enforcing this constraint does not require access to 1575

an external reward model or an "oracle" for point-1576 wise reward estimates, assuming we have access to 1577 labeled preference feedback in samples. Addition-1578 ally, the ℓ_2 formulation avoids placing a unbounded 1579 logit or a inverse sigmoid function over the preference since this has been shown to cause non-trivial 1581 policy degeneracy issues in learning algorithms 1582 like DPO (Azar et al., 2024). Applying this ℓ_2 loss, 1583 we derive: 1584

$$\mathcal{L}_{\pi_{\theta}} = \mathbb{E} \underset{\substack{\phi \sim \pi_{\theta}(\cdot|x) \\ f_{1}, f_{2} \sim \pi_{\theta}(\cdot|\phi, x)}}{x_{ref}(f_{1}|\phi, x)} \left(p(f_{1} \succ \phi|x) - p(f_{2} \succ \phi|x) - \beta \left[\log \left(\frac{\pi_{\theta}(f_{1}|\phi, x)}{\pi_{ref}(f_{1}|\phi, x)} \right) - \log \left(\frac{\pi_{\theta}(f_{2}|\phi, x)}{\pi_{ref}(f_{2}|\phi, x)} \right) + \log \left(\frac{\pi_{\theta}(f_{1}|x)}{\pi_{ref}(f_{1}|x)} \right) - \log \left(\frac{\pi_{\theta}(f_{2}|x)}{\pi_{ref}(f_{2}|x)} \right) \right] \right)^{2}$$
(35)

Since the friction dataset \mathcal{D}_{μ} sampled from μ contains preference-annotated pairs (f_w, f_l) given ϕ and x, the preference probabilities can be expressed using indicator functions as $p(f_w \succ f_l|x) = \mathbb{E}[\mathbf{1}(f_w \succ f_l|x)] = 1$ and $p(f_l \succ f_w|x) = \mathbb{E}[\mathbf{1}(f_l \succ f_w|x)] = 0$. Furthermore, the difference $p(f_w \succ f_l|x) - p(f_l \succ f_w|x) =$ 1 - 0 = 1 aligns with the formulation $p(f_1 \succ \phi|x) - p(f_2 \succ \phi|x)$ when $f_1 = f_w$ and $f_2 = f_l$. Therefore, we can write our final FAAF-alignment empirical objective function $(\hat{\mathcal{L}})$ as follows:

 $\hat{\mathcal{L}}(\pi_{\theta}) = \mathbb{E}_{(x,\phi,f_w,f_l) \sim \mathcal{D}_{\mu}}$

 $\log\left(\frac{\pi_{\theta}(f_l|\phi,x)}{\pi_{\mathrm{ref}}(f_l|\phi,x)}\right) +$

 $\left(1-\beta \Bigg[\log \left(\frac{\pi_{\theta}(f_w | \phi, x)}{\pi_{\rm ref}(f_w | \phi, x)}\right) - \right.$

1600

1601

1589

1590

1591

1592

1593

1594

1595

1597

1598

1599

160

1603
$$\log\left(\frac{\pi_{\text{ref}}(f_w|x)}{\pi_{\text{ref}}(f_l|x)}\right)^{-1}$$
1604
$$\log\left(\frac{\pi_{\theta}(f_l|x)}{\pi_{\text{ref}}(f_l|x)}\right)^{-1}$$
(36)

 $\int \pi_{\theta}(f_w|x)$

1605where (f_w, f_l) represent the winning (preferred)1606and losing (less preferred) friction interventions1607respectively in each annotated pair.

$$\hat{\mathcal{L}}(\pi_{\theta}) = \mathbb{E}_{(x,\phi,f_w,f_l) \sim \mathcal{D}_{\mu}}$$
1608

$$\left(1 - \beta \left[\log\left(\frac{\pi_{\theta}(f_w|\phi, x)\pi_{\text{ref}}(f_l|\phi, x)}{\pi_{\theta}(f_l|\phi, x)\pi_{\text{ref}}(f_w|\phi, x)}\right) + 1609\right]\right)$$

$$\log\left(\frac{\pi_{\theta}(f_w|x)\pi_{\text{ref}}(f_l|x)}{\pi_{\theta}(f_l|x)\pi_{\text{ref}}(f_w|x)}\right)\right)^2 \tag{37}$$

$$\hat{\mathcal{L}}(\pi_{ heta}) = \mathbb{E}_{(x,\phi,f_w,f_l) \sim \mathcal{D}_{\mu}}$$
 1611

$$\left(1 - \left[\underbrace{\beta \log\left(\frac{\pi_{\theta}(f_w | \phi, x) \pi_{\text{ref}}(f_l | \phi, x)}{\pi_{\theta}(f_l | \phi, x) \pi_{\text{ref}}(f_w | \phi, x)}\right)}_{\Delta R} + 1612\right]$$

$$\underbrace{\beta \log \left(\frac{\pi_{\theta}(f_w | x) \pi_{\text{ref}}(f_l | x)}{\pi_{\theta}(f_l | x) \pi_{\text{ref}}(f_w | x)} \right)}_{\Delta R'} \right] \right)^2 \tag{38}$$

where ΔR and $\Delta R'$ represent implicit reward dif-1614 ferences (Rafailov et al., 2024b; Azar et al., 2024), 1615 the former being explicitly conditioned on the fric-1616 tive state ϕ , with no such conditioning on the latter. 1617 Theorem 2 (Uniqueness of FAAF Empirical Loss). 1618 We prove this by contradiction. Let μ be the 1619 sampling distribution that samples friction interventions for the preference dataset, and assume 1621 $\text{Supp}(\mu) = \text{Supp}(\pi_{\text{ref}})$. Then the FAAF loss $\mathcal{L}(\pi)$ 1622 has a unique solution in policy space $\in \Pi$. 1623 *Proof.* Assume by contradiction that there exist two distinct optimal policies $\pi_A, \pi_B \in \Pi$. By their definition, $\hat{\mathcal{L}}(\pi_A) = \hat{\mathcal{L}}(\pi_A) = 0$ as π_A and π_B are global minima. Consider (s_{ϕ}^A, s^A) and (s_{ϕ}^B, s^B) as their respective logit parameterizations where:

$$\pi_k(f|\phi) = \frac{\exp(s_{\phi}^k(f))}{\sum_{f'} \exp(s_{\phi}^k(f'))}$$
$$\pi_k(f) = \frac{\exp(s^k(f))}{\sum_{f'} \exp(s^k(f'))} \quad \text{for } k \in \{A, B\}$$

where $\pi_k(f|\phi)$ and $\pi_k(f)$ are the ϕ -conditioned and ϕ -unconditioned policies. By the structure of our FAAF loss from Equation (38):

$$\hat{\mathcal{L}}(\pi) = \mathbb{E}_{f, f' \sim \mu} \Big[\big(1 - \beta (\Delta s_{\phi} + \Delta s) \big)^2 \Big] \ge 0$$

Notice that adding a constant c to all logits of s_{ϕ} or logits of s (directionally denoted as the $(c, \ldots, c) \in \mathbb{R}$) does not affect either policy probabilities due to softmax normalization. For $\hat{\mathcal{L}}(\pi)$, this is the *only* direction where the loss function might not be strictly convex. Outside of these directions, any change in the logits would increase $\mathcal{L}(\pi)$ with strict convexity as a consequence for $\alpha \in (0, 1)$, implying:

$$\hat{\mathcal{L}}(\alpha \pi_1 + (1 - \alpha)\pi_2) < \alpha \hat{\mathcal{L}}(\pi_1) + (1 - \alpha)\hat{\mathcal{L}}(\pi_2) \\ = \alpha(0) + (1 - \alpha)(0) = 0$$

where the equality follows from π_1, π_2 being global minima, by definition. This contradicts the non-negativity of $\hat{\mathcal{L}}$, which proves the uniqueness of the FAAF objective.

 $\hat{\mathcal{L}}(\pi_{\theta})$ has no dependence on log-partition terms involving $Z^*(\phi, x)$ and $Z^*(x)$ Our final FAAF empirical objective loss in Eq. 38 has no dependence on either partition function terms. This makes it convenient for practical applications. In fact, similar to DPO's derivation (Rafailov et al., 2024b), these log-partition terms effectively cancel out in formulating the frictive state-conditioned and unconditioned implicit rewards, scaled by the KL-strength parameter β . In its essence, $\mathcal{L}(\pi_{\theta})$ regresses the DPO-based implicit rewards ($\Delta R'$ term) with an additional ϕ -conditioned reward term $(\Delta R \text{ term})$ onto the empirically observed preference probabilities, labeled with preference labels from \mathcal{D}_{μ} . Notice that without the ΔR term, $\mathcal{L}(\pi_{\theta})$ reduces to a structurally similar form as IPO (Azar et al., 2024), differing a constant scaling term β . This suggests that under this condition, both $\mathcal{L}(\pi_{\theta})$ and IPO objective likely have similar qualitative loss landscapes though convergence rates and optimal solutions would differ—while both lead π_{θ} toward a reward-consistent preference alignment.

This also explains the somewhat similar performance of the IPO baseline and $FAAF_{\Delta R'}$ in both DeliData and WTD OPT 1.3B reward model-based win-rate evaluations, where $FAAF_{\Delta R'}$ achieves comparatively middling win-rates (Table 2).

Lemma 3 (Sequential Choice Decomposition in Friction Agent Optimization). Consider the minimax optimization between frictive-state policy π_{ϕ} and friction intervention policy π_f where we seek to generate optimal friction interventions f from frictive states ϕ :

$$J^* = \min_{\pi_{\phi}} \max_{\pi_f} \mathbb{E} \max_{\substack{x \sim \rho \\ \phi \sim \pi_{\phi}(\cdot \mid x)}} \left[p(f \succ \phi \mid x) - f_{\sigma = \pi_f(\cdot \mid \phi, x)} \right]$$
 165

$$\beta D_{\mathrm{KL}}(\pi_f \parallel \pi_{\mathrm{ref}} \mid \phi, x) + \beta D_{\mathrm{KL}}(\pi_\phi \parallel \pi_{\mathrm{ref}} \mid x) \right]$$
(39)

For any policy π (either optimal friction policy1659 π_f^* or reference policy π_{ref}), the sequential choice1660probability decomposes as:1661

$$\pi(\phi|\phi, x) = \frac{\pi(\phi|x)}{\pi(f|x)}$$
(40) 166

1646

1647

1648

1649

1650

1651

1652

1653

1654

1656

1645

Proof. The key insight in deriving this decomposition lies in understanding how optimal friction interventions are generated sequentially from frictive states. For the optimal friction policy π_f^* , consider its probability space $P_{\pi_f^*}$. By definition of conditional probability, we have $\pi_f^*(\phi|\phi, x) = \frac{P_{\pi_f^*}(\phi,\phi|x)}{P_{\pi_f^*}(\phi|x)}$. This term is crucial as it captures the policy's propensity to maintain a frictive state rather than generate a friction intervention. Under choice independence^{*a*} within this policy space assuming a Markovian nature of friction intervention generation, we have $P_{\pi_f^*}(\phi,\phi|x) = P_{\pi_f^*}(\phi|x)P_{\pi_f^*}(\phi|x)$. With policy-specific preference probability symmetry (Munos et al., 2023; Fisch et al., 2024), the probability $P_{\pi_f^*}(\phi|x) + P_{\pi_f^*}(f|x) = 1$, reflecting the binary choice between maintaining a frictive state or generating a friction intervention, we obtain $\pi_f^*(\phi|\phi, x) = \frac{\pi_f^*(\phi|x)}{\pi_f^*(f|x)}$, where the optimality of $\pi_f^*(f|x)$ ensures $\pi_f^*(\phi|x) \le \pi_f^*(f|x)$. A similar argument can be made in the case of π_{ref} , the reference policy, where π_{ref} 's initialization with the supervised-finetuned (SFT) model on friction interventions ensures $\pi_{ref}(f|x) \ge \pi_{ref}(\phi|x)$. This decomposition is fundamental to the minimax objective, J^* , as it enables expressing the KL-regularized preference probability in terms of base policy probabilities while preserving the structure necessary for optimal friction intervention generation from frictive states.

^{*a*}Assuming a single-step bandit setting (Rafailov et al., 2024b,a), choice independence holds since each frictive-state intervention is independent of past episodes. Using conditional probability, we express the joint probability under any policy π as $P_{\pi}(\phi, \phi \mid x) = P_{\pi}(\phi \mid \phi, x)P_{\pi}(\phi \mid x)$. By choice independence, the probability of selecting ϕ at the second step does not depend on the first selection given x, i.e., $P_{\pi}(\phi \mid \phi, x) = P_{\pi}(\phi \mid x)$. Substituting this, we obtain $P_{\pi}(\phi, \phi \mid x) = P_{\pi}(\phi \mid x)P_{\pi}(\phi \mid x)$.

The sequential choice decomposition provides 1663 crucial insight into determining optimal timing for 1664 friction interventions. In other words, this decom-1665 position has an interesting implication in deciding 1666 when is a friction intervention most desirable or cost-effective. Specifically, our derived identity 1668 $\pi(\phi|\phi,x)=\frac{\pi(\phi|x)}{\pi(f|x)}$ establishes a natural threshold mechanism through the ratio $\tau(x) = \frac{\pi_f^*(\phi|x)}{\pi_f^*(f|x)}$. 1670 When $\tau(x) \approx 1$, the policy maintains the current frictive state ϕ , while $\tau(x) \ll 1$ triggers a 1672 friction intervention f. This mechanism emerges 1673 naturally from the preference probability $p(f \succ$ 1674 $\phi|x) = \beta[\log(\pi_f^*(f|\phi, x)) - \log(\pi_{ref}(f|\phi, x)) +$ $\frac{1}{2\beta} - (\log(\pi_f^*(\phi|\phi, x)) - \log(\pi_{\text{ref}}(\phi|\phi, x)))] \text{ in our minimax objective } J^*, \text{ where } \pi_f^* \text{ optimally gener-}$ 1676 1677 ates interventions when the likelihood ratio indi-1678 cates low confidence in the current frictive state ϕ . However, exploring this sequential decomposition 1680 and determining optimal timing in interventions is 1681 outside the scope of this paper. As such, we leave that for future work. Lemma 4 (Uniqueness of Intervention Thresholds). 1684

The threshold $\tau(x) = \frac{\pi_f^*(\phi|x)}{\pi_f^*(f|x)}$ uniquely determines optimal intervention policy π_f^* . *Proof.* We prove uniqueness by contradiction. Consider two potentially optimal policies π_f^1 and π_f^2 with corresponding thresholds $\tau_1(x)$ and $\tau_2(x)$. Assume $\tau_1(x) \neq \tau_2(x)$ but both policies are optimal. By optimality, their contributions to the objective J^* must be equal for any observation tuple x, f and ϕ :

$$\beta [\log(\pi_f^1(f|\phi, x)) - \log(\pi_{\text{ref}}(f|\phi, x)) - (\log(\tau_1(x)) - \log(\pi_{\text{ref}}(\phi|\phi, x)))] = \beta [\log(\pi_f^2(f|\phi, x)) - \log(\pi_{\text{ref}}(f|\phi, x)) - (\log(\tau_2(x)) - \log(\pi_{\text{ref}}(\phi|\phi, x)))]$$
(41)

Simplifying and rearranging terms:

$$\log(\pi_f^1(f|\phi, x)) - \log(\tau_1(x)) = \log(\pi_f^2(f|\phi, x)) - \log(\tau_2(x))$$
(42)

However, by the strict convexity of KL divergence and Jensen's inequality:

$$D_{\mathrm{KL}}(\pi_{f}^{1} \parallel \pi_{\mathrm{ref}} \mid \phi, x) + D_{\mathrm{KL}}(\pi_{f}^{2} \parallel \pi_{\mathrm{ref}} \mid \phi, x) > 2D_{\mathrm{KL}}(\frac{\pi_{f}^{1} + \pi_{f}^{2}}{2} \parallel \pi_{\mathrm{ref}} \mid \phi, x)$$
(43)

This inequality implies that a mixed policy $\pi_f^{\text{avg}} = \frac{\pi_f^1 + \pi_f^2}{2}$ would achieve a lower KL divergence cost due to strict convexity and equal expected reward (regularized preference probabilities) from the equality of optimal policies. Therefore, π_f^{avg} would achieve strictly better objective value than both π_f^1 and π_f^2 , contradicting their assumed optimality. This proves threshold uniqueness. The contradiction arises because:

$$J^*(\pi_f^{\text{avg}}) > \frac{1}{2} [J^*(\pi_f^1) + J^*(\pi_f^2)]$$
(44)

which is impossible if both π_f^1 and π_f^2 were truly optimal.

Corollary 5 (Uniqueness of Optimal Policy Under Threshold Identity). If two optimal intervention policies π_f^1 and π_f^2 satisfy the same threshold condition $\tau_1(x) = \tau_2(x)$ for all x, then $\pi_f^1 = \pi_f^2$.

Proof. Assume for contradiction that two distinct optimal policies π_f^1 and π_f^2 satisfy the threshold condition $\frac{\pi_f^1(\phi|x)}{\pi_f^1(f|x)} = \frac{\pi_f^2(\phi|x)}{\pi_f^2(f|x)} = \tau(x)$. Define the mixed policy $\pi_f^{avg} = \frac{1}{2}(\pi_f^1 + \pi_f^2)$, which preserves the threshold as $\tau_{avg}(x) = \tau(x)$ due to linearity, implying π_f^{avg} is also optimal. Now, applying Jensen's inequality to the KL divergence term in the objective, we obtain $D_{KL}(\pi_f^{avg} \parallel \pi_{ref} \mid \phi, x) \leq \frac{1}{2}D_{KL}(\pi_f^1 \parallel \pi_{ref} \mid \phi, x) + \frac{1}{2}D_{KL}(\pi_f^2 \parallel \pi_{ref} \mid \phi, x)$. Strict convexity ensures a strict inequality whenever $\pi_f^1 \neq \pi_f^2$ on a set of positive measure where $\sup(\pi_{ref}) > 0$, implying $J^*(\pi_f^{avg}) < \frac{1}{2}[J^*(\pi_f^1) + J^*(\pi_f^2)]$. This contradicts the assumed optimality of π_f^1 and π_f^2 , proving that they must be identical. □

Lemma 6 (Policy Ratio Equivalence). For the optimal friction policy π_f^* and the optimal frictive-state policy π_{ϕ}^* , the following expectation-based ratio holds:

$$\mathbb{E}_{\substack{x \sim \rho \\ \phi \sim \pi_{\phi}^{*}(\cdot|x) \\ f \sim \pi_{f}^{*}(\cdot|\phi,x)}} \left[\frac{\pi_{f}^{*}(\phi|x)}{\pi_{f}^{*}(f|x)} \right] = \mathbb{E}_{\substack{x \sim \rho \\ \phi \sim \pi_{\phi}^{*}(\cdot|x) \\ f \sim \pi_{f}^{*}(\cdot|\phi,x)}} \left[\frac{\pi_{\phi}^{*}(\phi|x)}{\pi_{\phi}^{*}(f|x)} \right]$$
(45)

1696

1691

1693

1694

1695

Proof. We show that both the policy ratios simplify to the same value under the expectation. We begin by taking the expectation over the preference probability formulation^a:

$$\mathbb{E}\left[p(f \succ \phi \mid x)\right] = \mathbb{E}\left[\beta\left(\log(\pi_f^*(f \mid \phi, x)) - \log(\pi_{\text{ref}}(f \mid \phi, x)) + \log Z^*(\phi, x)\right)\right].$$
 (46)

We first represent the ratios of the optimal frictive intervention policies (LHS of this lemma) for any tuple (x, ϕ, f) in terms of their parametric representations from Eq. 8 as follows:

$$\frac{\pi_f^*(\phi|x)}{\pi_f^*(f|x)} = \frac{\pi_{\text{ref}}(\phi|x)}{\pi_{\text{ref}}(f|x)} e^{\beta^{-1}(p(\phi \succ f|x) - p(f \succ \phi|x))} \quad (\log Z^*(x) \text{ cancels out})$$
(47)

Take the expectation on both sides and apply^b the identity $p(\phi \succ f \mid x) = 0$:

$$\mathbb{E}\left[\frac{\pi_f^*(\phi|x)}{\pi_f^*(f|x)}\right] = \mathbb{E}\left[\frac{\pi_{\text{ref}}(\phi|x)}{\pi_{\text{ref}}(f|x)}e^{-\beta^{-1}p(f\succ\phi|x)}\right] \quad (\text{since } p(\phi\succ f\mid x) = 0).$$
(48)

Notice that by definition in Eq. 2, the optimal friction intervention policy $\pi_f^*(\cdot|\phi, x)$ is KLconstrained wrt to the reference policy $\pi_{ref}(\cdot|\phi, x)$. So under the expectation, the following has to be true for $\pi_f^*(\cdot|\phi, x)$ to be optimal:

$$\mathbb{E}\left[\pi_f^*(f|\phi, x)\right] \approx \mathbb{E}\left[\pi_{\text{ref}}(f|\phi, x)\right].$$
(49)

Substituting the preference probability formulation $p(f \succ \phi | x)$ from Eq. 24 in Eq. 48 and applying the KL-regularization approximation in Eq. 49 we derive that:

$$\mathbb{E}\left[e^{-\left(\log(\pi_f^*(f|\phi,x)) - \log(\pi_{\mathrm{ref}}(f|\phi,x)) + \log Z^*(\phi,x)\right)}\right] \approx \mathbb{E}\left[\frac{Z^*(f,x)}{Z^*(\phi,x)}\right].$$
(50)

Using this substitution, we rewrite Eq. 48 as:

$$\mathbb{E}\left[\frac{\pi_{f}^{*}(\phi|x)}{\pi_{f}^{*}(f|x)}\right] = \mathbb{E}\left[\frac{\pi_{\mathrm{ref}}(\phi|x)}{\pi_{\mathrm{ref}}(f|x)}e^{-\left(\log(\pi_{f}^{*}(f|\phi,x)) - \log(\pi_{\mathrm{ref}}(f|\phi,x)) + \log Z^{*}(\phi,x)\right)}\right]$$
$$= \mathbb{E}\left[\frac{\pi_{\mathrm{ref}}(\phi|x)}{\pi_{\mathrm{ref}}(f|x)}\frac{Z^{*}(f,x)}{Z^{*}(\phi,x)}\right].$$
(51)

Similarly, for the optimal frictive state policy ratio we derive:

$$\mathbb{E}\left[\frac{\pi_{\phi}^{*}(\phi|x)}{\pi_{\phi}^{*}(f|x)}\right] = \mathbb{E}\left[\frac{\frac{\pi_{\mathrm{ref}}(\phi|x)}{Z_{\phi}^{*}(x)}e^{-\beta^{-1}\log Z^{*}(\phi,x)}}{\frac{\pi_{\mathrm{ref}}(f|x)}{Z_{\phi}^{*}(x)}e^{-\beta^{-1}\log Z^{*}(f,x)}}\right] = \mathbb{E}\left[\frac{\pi_{\mathrm{ref}}(\phi|x)}{\pi_{\mathrm{ref}}(f|x)}\frac{e^{-\log Z^{*}(\phi,x)}}{e^{-\log Z^{*}(f,x)}}\right] \quad (Z_{\phi}^{*}(x) \text{ cancels})$$
(52)

$$= \mathbb{E}\left[\frac{\pi_{\mathrm{ref}}(\phi|x)}{\pi_{\mathrm{ref}}(f|x)}\frac{Z^*(f,x)}{Z^*(\phi,x)}\right].$$
(53)

Thus,
$$\mathbb{E}\left[\frac{\pi_{f}^{*}(\phi|x)}{\pi_{f}^{*}(f|x)}\right] = \mathbb{E}\left[\frac{\pi_{\phi}^{*}(\phi|x)}{\pi_{\phi}^{*}(f|x)}\right].$$

^{*a*}For clarity, the expectation \mathbb{E} is taken over $x \sim \rho, \phi \sim \pi_{\phi}^*(\cdot \mid x), f \sim \pi_f^*(\cdot \mid \phi, x)$ throughout the proof, but this is omitted in the notation when the context is clear.

^bSince learning occurs in a supervised setting with preference-annotated data, the probability follows as
$$p(f \succ \phi \mid x) = \mathbb{E}[1(f \succ \phi \mid x)] = 1$$
, implying $p(\phi \succ f \mid x) = 0$.

1706

1707

1708

1709

1710

1711

1712

1713

1714

1715

1716

1718

1719

1720

1721

1723

1724

1725

1726

1727

1728

1730

1731

1732

1733

1735

1736

1737

1738

1739

1741

1742

1743

1744

1745

1746

1697

D Operationalizing μ : Frictive State and Friction Intervention Generations

In order to train and evaluate our baselines along with FAAF for friction intervention generation in collaborative tasks, we carry out a series of data augmentation procedures using GPT-40 (denoted as μ) in order to construct two diverse preference datasets. For details on our choice of datasets, please refer to Section 4.2.

In this section, we provide procedural details of our friction intervention datasets, that were generated out of the original Weights Task and DeliData dataset. For all our data-generation experiments, we use a high-capacity LLM (GPT-40) (OpenAI et al., 2024) as our sampling distribution μ , as defined in Section 4. In particular, we utilize a selfrewarding LLM approach (Yuan et al., 2024; Xu et al., 2023; Rosset et al., 2024) to simultaneously generate and assign rewards to μ -generated interventions, since previous work (Pace et al., 2024; Meng et al., 2024) provides evidence that such synthetic preference-data generation still leads to higher-quality reward models and preferencealigned policies. Prior work (Zheng et al., 2023) provides substantial evidence that this approach leads to more high-quality LLM-as-a-judge-based evaluations especially for conversational benchmarks (Lambert et al., 2024). Additionally, reward assignments for sampled intervention naturally provides an implicit preference rankingwhich we use for constructing our respective preference datasets. After these data-generation experiments, we further conduct filtering and contrastive pairing of a "winning" (f_w) or preferred interventions and "losing" (f_l) or dispreferred interventions along with their corresponding dialogue histories (x) to create our final preference datasets for each augmented dataset.

D.1 DeliData Friction Intervention Preference Dataset

In order to generate frictive state and friction interventions in the DeliData dataset, we use the prompt shown in Figure 2. In order to contextualize the extraction of frictive states, we only provide h = 15 previous utterances in each dialogue group (group_id) assuming that frictive states are likely to be present within a "attentionalstate" (Grosz and Sidner, 1986) window that describes the focused part in the discourse. This technique allows us to avoid unnecessary api-calls while also providing a more focused dialogue con-1747 text to GPT-40. Additionally, since this dataset al-1748 ready contains manual human annotations of "prob-1749 ing" interventions (which are a subset of friction 1750 interventions as per our definitions), we explicitly 1751 guide the data-generator to exclude probing inter-1752 ventions in extracting the frictive states. Note that 1753 each functionally-frictive state (denoted as ϕ), as 1754 extracted by GPT-40, resulted in two friction in-1755 terventions, f_w and f_l . In total, this generation 1756 process led to 6238 (x, ϕ, f_w, f_l) tuples after keep-1757 ing 50 randomly sampled dialogue groups separate 1758 for the evaluation set, our of which 476 (33) were 1759 probing interventions in train (test) partitions. Ad-1760 ditionally, we carry out another round of training 1761 pair augmentations since 6238 samples is very less 1762 compared to popular preference alignment datasets 1763 like Ultrafeedback (Cui et al., 2024) which con-1764 tains roughly 62k¹¹ training preference pairs. The 1765 average rewards for the preferred and dispreferred 1766 interventions assinged by μ are 8.03 and 3.96 re-1767 spectively (rated out of 10). 1768

1769

1770

1771

1772

1773

1774

1775

1776

1777

1778

1779

1780

1781

1782

1783

1784

1785

1786

1787

1788

1789

1790

As such, for each training tuple (x, ϕ, f_w, f_l) , we generate N augmented versions (x', ϕ', f'_w, f'_l) by applying a replacement mapping $R: \Sigma \to \Sigma'$ N times, where Σ represents the original set of card values (vowels¹², odd numbers, and even numbers), and Σ' represents their replacements. The replacement function R is defined as follows: Each vowel $v \in \{A, E, O, U\}$ is replaced with another vowel v' such that $v' \in \{A, E, O, U\} \setminus \{v\}$, where v' is sampled uniformly at random from the remaining vowels. Similarly, each odd number $o \in$ $\{1, 3, 5, 7, 9\}$ is replaced with another odd number o' such that $o' \in \{1, 3, 5, 7, 9\} \setminus \{o\}$, where o' is sampled uniformly at random. Likewise, each even number $e \in \{0, 2, 4, 6, 8\}$ is replaced with another even number e' such that $e' \in \{0, 2, 4, 6, 8\} \setminus \{e\}$, with e' sampled uniformly at random. For example, if an utterance contains reference to card "A" and "6", the rules of the Wason Card task still applies equivalently for, say, "E" and "8"-while keeping the reasoning consistent with the original utterance and the utterance with replacement. We

¹¹We found 14 samples where GPT-40 did not return any strings for the frictive state description. We filtered out these samples from our training set.

¹²We did not replace instances of "I" to avoid noise from mistakenly replacing first-person references in the dialogues. Additionally, since vowels constitute the majority of prompted card solutions vs. consonants, applying our replacement function R for vowels was enough to generate ~62k additional samples, comparable to Ultrafeedback (Cui et al., 2024)

	Train			Test		
	Min	Max	Mean ± Std	Min	Max	Mean ± Std
Dialogue History	16	824	288.43 ± 132.97	25	733	291.88 ± 118.03
Belief State	8	140	32.93 ± 15.92	20	140	47.99 ± 28.38
Chosen Friction	6	60	24.03 ± 4.65	9	39	22.05 ± 5.55
Chosen Rationale	8	78	22.84 ± 8.67	10	78	29.61 ± 13.33
Rejected Friction	9	45	23.95 ± 4.10	10	41	22.16 ± 5.11
Rejected Rationale	8	73	19.60 ± 6.89	10	59	26.04 ± 11.61

Table 4: Token Length Statistics for DeliData Preference Dataset using the Meta-Llama-3-8B-Instruct tokenizer.

Field	Train			Test		
	Min	Max	Mean ± Std	Min	Max	Mean ± Std
Dialogue History	4	1464	227.83 ± 189.48	4	1031	235.04 ± 180.36
Belief State	11	65	30.55 ± 6.65	17	54	30.47 ± 6.29
Chosen Friction	10	45	21.20 ± 5.12	11	42	21.08 ± 5.10
Chosen Rationale	10	35	20.38 ± 3.44	12	32	19.67 ± 3.38
Rejected Friction	6	32	15.88 ± 3.68	7	29	15.57 ± 3.75
Rejected Rationale	8	41	20.10 ± 3.51	12	30	19.88 ± 3.47

Table 5: Token Length Statistics for **WTD Simulated Friction** dataset using the Meta-Llama-3-8B-Instruct tokenizer.

Field	Train		Test			
	Min	Max	Mean ± Std	Min	Max	Mean ± Std
Dialogue History	16	555	309.88 ± 81.11	25	555	316.46 ± 79.16
Belief State	41	140	84.94 ± 15.58	41	140	84.95 ± 16.27
Chosen Friction	9	31	16.85 ± 3.47	9	27	16.87 ± 3.49
Chosen Rationale	24	78	44.19 ± 8.46	26	78	44.43 ± 8.59
Rejected Friction	9	31	17.12 ± 3.51	10	28	17.23 ± 3.41
Rejected Rationale	24	73	40.00 ± 6.62	24	59	39.89 ± 6.43

Table 6: Token Length Statistics for **WTD Original Friction** dataset using the Meta-Llama-3-8B-Instruct tokenizer.

Personality Type	Facet	Description
Extraversion	Assertiveness Sociability Activity Level Excitement Seeking Positive Emotions	Tends to take charge and speak confidently. Enjoys engaging with others and maintaining conversation. Shows high energy and enthusiasm. Looks for novel and stimulating experiences. Expresses optimism and cheerfulness.
Neuroticism	Anxiety Depression Vulnerability Self-Consciousness Anger	Shows worry and concern about potential mistakes. Tends to be pessimistic and doubtful. Easily becomes overwhelmed or stressed. Shows hesitation and uncertainty. Can become frustrated and irritated easily.
Agreeableness	Trust Altruism Compliance Modesty Sympathy	Readily trusts others and their suggestions. Shows concern for others' success and well-being. Tends to avoid conflicts and agree with others. Downplays own contributions and abilities. Shows understanding and empathy towards others.

Table 7: Descriptions of our chosen 3 personality types and facet combinations from the Big Five framework that we use for simulated friction generation on the Weights Task.

apply this replacement mapping across all fields in tuples (x', ϕ', f'_w, f'_l) in the training set. This led to training set of 68,618 preference pairs. Note that 1793 we only apply this augmentation for the training set 1794 to generate a reasonably large preference dataset for more robust training signals. Table 4 shows a 1796 detailed breakdown of the token-length statistics of the DeliData Friction preference dataset using the 1798 Meta-Llama-3-8B-Instruct tokenizer.

1791

1792

1797

1801

1802

1803

1804

1808

1809

1810

1812

1813

1814

1816

1817

1818

1819

1820

1821

1822

1824

1825

1826

1827

1828

1830

1831

1833

1835

1836

1837

1838

1839

1841

D.2 WTD Friction Intervention Preference Dataset

WTD "Original" Friction dataset Unlike the DeliData dataset, which includes pre-annotated probing interventions as natural friction points, the Weights Task dataset (Khebour et al., 2024a) consists of dense-paraphrased utterances transcribed manually (Terpstra et al., 2023) and with Whisper (Radford et al., 2023), making friction interventions sparse due to its multimodal nature. Manual inspection found only 3-4 frictive interventions per group, yielding \approx 30-40 samples—insufficient for training an effective agent without overfitting, especially for LLMs with billions of parameters. As such, we carry out two phases of dataaugmentations and preference annotations. In our first round, we generate the WTD Original Friction dataset which contains annotations of frictivestates and friction interventions. Similar to DeliData preference annotations Appendix D.1, we use a self-rewarding LLM set-up to first generate these states and interventions in an autoregressive manner and prompt μ to rate them in the same api-call, for each frictive state extraction. Since WTD dialogues can be substantially long (> 200utterances) for certain groups, we only consider a non-overlapping window of 10 previous utterances as context history h = 10 for a more robust grounding for μ ; See Fig. 2 for details on the prompt used constructing the WTD Original Friction dataset. This process led to 4299 (470) training (testing) preference pairs. Preferred interventions achieved mean scores (mean±std) of 8.36±1.12 (train) and 8.35±1.08 (test), while dispreferred interventions scored 6.35±1.13 (train) and 6.36±1.11 (test), demonstrating consistent preference margins across splits.

Note that we do not use WTD Original Friction for training any of our baselines—but use it for out-of-domain distribution (OOD) evaluation (see Sec. 4.2). This allows us to more extensively evaluate FAAF in checking test-time OOD generalization (Rafailov et al., 2024b; Choi et al., 2024) against baselines-where OOD generalization is a major limitation in supervised preference alignment algorithms that depend crucially on the sampling distribution (Yang et al., 2024; Fisch et al., 2024).

1842

1843

1844

1845

1846

1847

1848

1849

1851

1852

1853

1854

1856

1857

1858

1859

1860

1861

1862

1863

1864

1865

1866

1868

1869

1870

1871

1872

1873

1874

1876

1877

1878

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889

WTD "Simulated" Friction dataset Additionally, for a more robust training and in order to evaluate multi-turn preference alignment in interventions, we use (Shani et al., 2024)'s method to generate novel full collaborative conversations using the weight-definitions of the original WTD environment. This method is more akin to "Westof-N" sampling (Pace et al., 2024) techniques that allow synthetic data generations with high-capacity LLMs-where highest and lowest rewarded candidates naturally form preference pairs. As shown in Fig. 4, we sample a full dialogue at once using μ , while providing initial task-related guidelines and gold-truth labels of actual weights of the five blocks in the WTD dataset. For example, we explicitly prompt μ to role-play (Li et al., 2023a) the triad consisting of three participants in the weight-deduction process. Furthermore, to generate more realistic utterances, we utilize participant personality-facet combinations (Pan and Zeng, 2023; Mao et al., 2024) from Big 5^{13} personality classifications (Goldberg, 2013) as additional attributes in the prompt. In other words, each sampled full-dialogue contains a unique combination of these personality-facet combinations for each participant (total 3,375 combinations).

Similarly, for each sampled frictive state within a conversation (dialogue), we generated N = 6 friction interventions with corresponding effectiveness scores in resolving the frictive state. Since WTD data does not contain any probing intervention samples, in order to further ground these generations to the task, we also provide a one-shot example of a naturally occurring friction intervention (marked with P1(f) in Fig. 4). In total, out of the expected 3,375 personality-facet combinations (3*5 unique combinations for each participant), 3,362 were successfully generated using μ and parsed. Finally, to create the preference pairs, for each frictive state, we paired the lowest scoring response with all the higher scoring ones, akin to the West-of-N technique. This resulted in 56, 689 preference pairs

¹³See Tab. 7 for our full set of personality-type and facet combinations. Similar to (Mao et al., 2024), we choose three personality types from Big 5 framework for consistency.

1890 after excluding 54 dialogues (amounting to 800 preference pairs) for the test set. This process fi-1891 nally resulted in the WTD Simulated Friction 1892 dataset. Preferred interventions achieved mean 1893 1894 scores of 8.48±1.52 (train) and 8.51±1.50 (test), 1895 while dispreferred interventions scored 6.01±0.88 (train) and 6.08±0.87 (test), demonstrating consis-1896 tent preference margins across splits. 1897

Personality Types	P1	P2	P3
Extraversion	4740	4889	4741
Neuroticism	5928	5591	5921
Agreeableness	4573	4761	4579

Table 8: Friction Count for Participants

FRICTION GENERATION PROMPT: DELIDATA DATASET

System: You are an expert in collaborative reasoning and dialogue analysis. Your task is to detect *frictive states* and generate *friction interventions* that resolve them in group dialogue. A frictive state occurs when a participant makes a claim that contradicts another participant's belief model (i.e., their assumed understanding of the rule or task constraints), leading to misalignment in reasoning that could hinder progress. Friction interventions encourage self-reflection in participants and prompt them to reevaluate these contradicting beliefs and assumptions.

User: Analyze this dialogue about the Wason card selection task. Participants see four cards showing numbers or letters and must test this rule: "All cards with vowels on one side have an even number on the other." Remember that the correct answer is to select a vowel and an odd number. Provide [N] frictive states with their resolutions in the following JSON format. For each state, include both a preferred and less preferred intervention that could help resolve the conflict. Additionally, provide a one-sentence rationale for your intervention.

IMPORTANT: - Do not analyze utterances labeled as "probing" or statements immediately before them, as these frictive states have already been detected.

- For each frictive state detected, you should:
- * Identify the dialogue index where it occurs
- * Summarize the conflicting beliefs
- * Explain why the contradiction affects reasoning

Here is the provided dialogue: [Dialogue]

Message ID: [index_where_friction_occurs] Contradiction: [describe_the_conflicting_beliefs] Contradiction Reason: [explain_why_the_contradiction_affects_reasoning]

Preferred Intervention:

Statement: [your_friction_intervention] Rationale: [your_rationale] Score: [your_score]

Less Preferred Intervention:

Statement: [your_friction_intervention] Rationale: [your_rationale] Score: [your_score]

Figure 2: Delidata (Karadzhov et al., 2023) Friction Generation Prompt. We use GPT-40 as our sampling distribution μ and prompt it to simultaneously generate frictive states and friction interventions. For diversity, we use the default temperature of 1. This process implicitly provides us with preference rankings between intervention, via the reward scores. See Section 3 for definitions of frictive states and friction interventions. Note that we exclude already-present "probing" interventions in this generation process since are present in the original Delidata annotations.

FRICTION GENERATION PROMPT: WEIGHTS TASK DATASET (WTD ORIGINAL)

System: You are an expert in collaborative reasoning and dialogue analysis. Your task is to detect *frictive states* and generate *friction interventions* that resolve them in group dialogue. A frictive state occurs when a participant makes a claim that contradicts another participant's belief model (i.e., their assumed understanding of the rule or task constraints), leading to misalignment in reasoning that could hinder progress. Friction interventions encourage self-reflection in participants and prompt them to reevaluate these contradicting beliefs and assumptions.

User: Analyze this dialogue about the Weights Task dataset. Three participants (P1, P2, and P3) are collaborating to determine the weights of colored blocks using a scale.

Block Weights (in grams):

- Red block: 10g
- Blue block: 10g
- Green block: 20g
- Purple block: 30g
- Yellow block: 50g

Game Rules:

- 1. Participants can only weigh two blocks at a time
- 2. They are told the red block's weight at the start
- 3. All other block weights are initially unknown
- 4. Scale slider is not needed (blocks are in 10g increments)

Provide [N] frictive states with their resolutions in the following JSON format. For each state, include both a preferred and less preferred intervention that could help resolve the conflict. Additionally, provide a one-sentence rationale for your intervention. Here is the provided dialogue: [Dialogue]

Message ID: [index_where_friction_occurs] Contradiction: [describe_the_conflicting_beliefs] Contradiction Reason: [explain_why_the_contradiction_affects_reasoning]

Preferred Intervention:

Statement: [your_friction_intervention] Rationale: [your_rationale] Score: [your_score]

Less Preferred Intervention:

Statement: [your_friction_intervention] Rationale: [your_rationale] Score: [your_score]

Figure 3: Weights Task dataset (Khebour et al., 2024b) Friction Generation Prompt. We use GPT-40 as our sampling distribution μ and prompt it to simultaneously generate frictive states and friction interventions. For diversity, we use the default temperature of 1.

FRICTION GENERATION PROMPT: WEIGHTS TASK DATASET (WTD SIMULATED)

System: You are an expert in collaborative reasoning and dialogue analysis. Your task is to *generate a complete dialogue* where participants (P1, P2, P3) discuss which block to measure next and how to measure them. The three participants have distinct personality types that influence their behavior and dialog must reflect these personality traits in their communication style and behavior. The dialog is considered complete when all block weights are measured and agreed upon. Additionally, identify frictive states within the dialogue and provide N friction interventions at these points.

[Definition: Frictive State] [Definition: Friction Intervention]

User: Three participants (P1, P2, P3) work together in the Weights Task to determine the weights of colored blocks (red=10g, blue=10g, green=20g, purple=30g, yellow=50g). They can only weigh two blocks at a time, start knowing only the red block's weight, and use a scale with 10g increments (no slider needed).

Your tasks:

Generate a full dialogue until all weights are correctly identified and agreed upon.

Identify frictive states where reasoning misalignment occurs.

Provide N friction interventions with their corresponding rationales at these points. Rank them by effectiveness in resolving the conflict. Assign each a quality score from 1 to 10.

P1 has {personality_type} personality type with high {personality_facet}.

Here is an example dialogue where friction statements are labeled as (f). Actions of participants are provided within "[]" blocks.

P2: [pointing towards the purple block first and then towards the blue block] I think this one is purple and this one is blue.

P3: [reading from the laptop screen] Ok so blue is ten and purple is

P3: [looking at the blocks and asking a rhetorical question] Thirty

P1 (f): [putting green and red blocks on the left side of the scale and purple block on the right side] Yes verify real quick but I think it is

P2: [observing the balanced scale] Yes thirty

P1: [removing green, red, and purple blocks from the scale] Yeah we got them yeah

Generated Dialogue:

[Full_generated_dialogue_until_completion]

Message ID: [index_where_friction_occurs] Contradiction: [describe_the_conflicting_beliefs] Contradiction Reason: [explain_why_the_contradiction_affects_reasoning]

Friction Interventions:

Statement: [your_friction_intervention] Rationale: [your_rationale] Score: [your_score]

Figure 4: "Simulated" Weights Task dataset (WTD Simulated) Friction Generation Prompt. To ground these friction interventions with personality-traits of the participants, we use (Mao et al., 2024)'s prompting framework with personality-facet combinations. We use GPT-40 as our sampling distribution μ and prompt it to simultaneously generate frictive states and friction interventions. For diversity, we use the default temperature of 1.

1925

1926

1927

1928

1930

1931

1932

1933

1934

1935

1936

1937

1938

1940

1942

1943

1944

D.3 Tie Counts: GPT-40 Evaluation

Fig. 5 shows the tie-count distribution (baselines vs. 1899 SFT model completions) over our 7 preference di-1900 mensions on DeliData (top), Simulated WTD (mid-1901 dle) and Original WTD (bottom) datasets, when 1902 evaluated for win-rate computations using scores 1903 assigned by the LLM-judge (GPT-40). To avoid 1904 1905 positional bias in the placement of the sampled completions (friction interventions), we swap the 1906 positions of the two candidate samples in each run 1907 and then report the mean tie-count across each preference dimension. On average, Fig. 5 reveals that 1909 the LLM-judge have lower raw-agreement on di-1910 mensions such as consistency of the friction inter-1911 vention with its rationale (rationale_fit), relevance 1912 1913 and thought_proving on all three datasets compared to aspects like gold-alignment, specificity and im-1914 pact. This is expected since surface-level alignment 1915 with the golden samples are easier to assign a clear 1916 preference compared to metrics like rationale con-1917 sistency especially when interventions from both 1918 the candidate and the opponent are well-justified. 1919 Consistent with our results from Table 1, we find that FAAF model tends to tie less than other baselines on average. This trends is more pronounced 1922 in the WTD datasets consistent with FAAF's overall 1923 performance as shown in our main results.

D.4 Human Validation of Generated Friction Interventions

Following previous work that evaluates LLMgenerated annotations and outputs (Wiegreffe et al., 2021, 2022; Nath et al., 2024a,c), in addition to choosing the winning intervention, we asked the human annotators¹⁴ to evaluate the candidates in each sample across dimensions of reasoning, specificity, and thought provoking. Annotators were asked to rate both candidate interventions on a 5point Likert-type scale. For analysis, we bucketed the ratings together by valence—1 & 2: negative valence (-1), 3: neutral valence (0), and 4 & 5: positive valence (1), and calculated average valences and Krippendorff's α and Cohen's κ . We find that the average valence ratings of the various dimensions is low, very close to neutral, as are the α and κ values ($\alpha = 0.276, \kappa = 0.205$ on DeliData samples, $\alpha = -0.265$, $\kappa = 0.004$ on WTD). There is little agreement on the qualities of the friction

statement which suggests that although the anno-1945 tators usually have strong agreement that there is 1946 a clear winner for each pair (see Sec. 4.2), there 1947 is a lot of subjectivity on the qualities of these ut-1948 terances. While the winning utterance was judged 1949 to be better at prompting reflection or redirecting 1950 the dialogue, it may nor be entirely clear to the 1951 annotators why. In addition, these qualities are 1952 loosely-derived from other human-LLM validation 1953 frameworks, which usually align somewhat with 1954 how LLMs themselves score things, which is often 1955 based on specific detail and level of informativity. 1956 These might not actually be the best qualities to 1957 emphasize in a collaborative dialogue, because they tend to violate Gricean principles (Grice, 1975) in 1959 a collaborative context, due to informativity and specific detail leading to redundancy, violating the 1961 maxim of quantity, etc. 1962

¹⁴Our two human annotators have the following demographic breakdown: both male, college undergraduates, one Caucasian, one African, both fluent English speakers.



Figure 5: Comparison of average tie counts of baselines against SFT model over two runs across our 7 distinct dimensions (metrics) when evaluated using our GPT-4o-based LLM-as-a-judge evaluation in a "preference"-based setting (see Fig. 7)—on DeliData (top), Simulated WTD (middle) and Original WTD (bottom). Note that there were no ties in GPT's "overall" preference between a baseline vs SFT model.

D.5 Training Settings and Hyperparameters

1963

1964

1965

1966

1967

As motivated in Sec. 4, FAAF-aligned π_{θ} learns to distinguish signals that determine why a particular intervention is more preferred by explicitly conditioning its implicit reward estimation on

the frictive-state ϕ . This allows the model to esti-1968 mate the true preference distribution \mathcal{P} by balanc-1969 ing its load, from learning both with and without ϕ -conditioning, given a context. This is empirically seen in Fig. 6 (top), where π_{θ} displays a bal-1972

anced¹⁵ learning of "preference-strengths" between the winning and losing response (via the winning and losing response rewards as well as margins conditioned on ϕ), subject to the KL-regularization strength parameter β . We use the TRL Library's trainer classes for efficient multi-GPU training.

1973

1974

1975

1976

1978

1979

1980

1981

1983

1984

1985

1986

1987

1988

1989

1990

1992

1993

1994

1995

1996

1997

1998

2001

2004

2009

2010

2012

2013

2014

2015

2016

Hyperparameters for baselines All our preference alignment baselines:DPO (Rafailov et al., 2024b), IPO (Azar et al., 2024) and PPO (Schulman et al., 2017) are initialized with the Supervisedfinetuned (SFT) models that were trained on the winning responses (f_w) of DeliData and Simulated WTD training sets, following prior work to ensure the SFT model has reasonable support over the winning responses generated from μ .

For SFT models, we initialize them from the base meta-llama/Meta-Llama-3-8B-Instruct model in order to leverage its instruction following and general conversational abilities (AI@Meta, 2024). Due to compute constraints, we conducted all our training experiments with LoRA (Low-Rank Adaptation of Large Language Models), where LoRA $\alpha = 16$, LoRA dropout = 0.05 and a LoRA R of 8 was used in training with the PEFT¹⁶ and SFTT¹⁷ trainers from the TRL library. We use the bitsandbytes¹⁸ library to load our models in 4-bit quantization for more cost-efficient training.

Additionally, as mentioned in Sec. 5, we only compute the loss on completions (includes both frictive states ϕ and interventions f_w) using a ConstantLengthDataset format for more efficient training. We use a learning-rate (LR) of 1e-4with AdamW (Loshchilov et al., 2017; Dettmers et al., 2024) optimization with a cosine LR scheduler with a weight-decay of 0.05 and 100 warm-up steps. We train the SFT models for 6000 steps (\approx 1.5 epochs with approximately 58k samples) with an effective batch-size of 16 (gradient accumulation of 4) that reasonably achieves convergence on a 5% validation set randomly sampled from the training sets of both datasets. For context-length, we use a maximum length of 4096 tokens.

Offline baselines For DPO and IPO, we use similar LoRA settings with a max_length (including

¹⁶https://huggingface.co/docs/peft/index

both prompts and responses) for 4096 tokens with 2017 a max_prompt_length of 1024 tokens that only 2018 minimally filters our preference pairs that exceed 2019 this length, and helps avoid out-of-memory (OOM) issues during training. We train for 2000 steps 2021 with an effective batch size of 32 and an LR of 5e - 6, following default settings. Note that for 2023 IPO, we normalize the log-probabilities of the pre-2024 ferred and the dispreferred responses using their token-lengths. 2026

2027

2031

2033

2034

2036

2040

2041

2042

2043

2048

2050

2051

2053

2054

2058

2060

2061

2062

2063

PPO baseline For PPO, we additionally training an OPT 1.3B reward model (RM) following prior work (Hong et al., 2024) using a standard Bradley-Terry loss formulation using the TRL reward modeling library.¹⁹ Due to PPO's excessive compute requirements, for policy training, we use an effective batch size of 8 with a mini-batch size of 4 and gradient accumulation per 2 steps and train for 4,000 batches for two epochs. We constrain response tokens to be between 180 and 256 tokens using a LengthSampler while the queries are truncated to 1,024 tokens, with LR of 3e-6 for DeliData and 1.41e - 6 for Simulated WTD. For sampling response tokens, we use a top-p of 1.0 for diversity. We found that subtracting the baseline reward for the golden friction interventions (f_w) from the RM-assigned rewards stabilizes training. Therefore, we report results using this method in Table 1 and Table 2.

FAAF Training Settings For training FAAF, we use a batch size of 8 with the same PEFT/LoRA settings mentioned above and train for 2000 steps with a slightly smaller LR of 5e - 7, due to the smaller batch-sizes. For efficiency, we compute both the ϕ -conditioned ($\pi_{\theta}(f|\phi, x)$) and unconditioned $(\pi_{\theta}(f|x))$ policy logits in parallel within each forward pass. The winning (f_w) and losing (f_l) intervention pairs for each conditioning type are batched together, requiring only two forward passes total per batch. We implement this using a modified version of the DPO Trainer²⁰ from TRL, adapting it to handle the dual policy outputs. For data preprocessing, we filter pairs exceeding max_length of 2,500 and 3,000 tokens in DeliData and Simulated WTD respectively, with max_prompt_length set to 1024 tokens. Following standard practice, we compute token-length normalized log-probabilities for more

¹⁵By balance, we mean that *both* ϕ -conditioned and ϕ unconditioned implicit rewards capture preference strengths from the data.

¹⁷https://huggingface.co/docs/trl/en/sft_ trainer

¹⁸https://huggingface.co/docs/transformers/ main/en/quantization/bitsandbytes

¹⁹https://github.com/huggingface/trl/blob/main/ trl/trainer/reward_trainer.py

²⁰https://huggingface.co/docs/trl/main/en/dpo_ trainer

stable training. For the KL-regularization hyper-2064 parameter β , we conducted an ablation study over 2065 $\beta \in \{10, 5, 1, 0.01\}$. As shown in Fig. 6, $\beta = 10$ 2066 achieves optimal performance across multiple metrics: (1) higher implicit reward accuracy in both 2068 2069 ϕ -conditioned and unconditioned policies, (2) better reward margins between winning and losing 2070 interventions, and (3) more stable convergence of 2071 the FAAF loss, while NLL loss or cross-entropy loss is relatively lower than lower β values. Notably, 2073 while smaller β values (e.g., $\beta = 0.01$) fail to dis-2074 tinguish preference margins effectively, $\beta = 10$ 2075 provides sufficient reward margins. We therefore 2076 use $\beta = 10$ for all FAAF experiments reported in 2077 our results. 2078

Training Hardware We train all our models that 2079 require a reference model in memory on two Nvidia 2080 A100 GPUs, while the OPT 1.3B reward model 2081 (full-parameter training) and the SFT model were trained on a single A100 GPU. Training a single baseline for 2000 steps roughly took 12 hours of GPU compute, but PPO models that were trained for 4000 minibatches of size 8 took roughly 24 hours to train until convergence. 2087

2084



Figure 6: Ablation study of FAAF's β hyperparameter ($\beta \in \{10, 5, 1, 0.01\}$) during training on the Simulated WTD data (top-half) and DeliData datasets (bottom-half) across 2k and 1k training steps respectively. Higher β values (e.g., $\beta = 10$) show better implicit reward estimation as shown in Reward Accuracy plots and estimated preference-strengths (Reward Margins), while very small values ($\beta = 0.01$) fail to distinguish preferences effectively. $\beta = 10$ also minimizes NLL and FAAF losses suggesting model stability and better convergence. As such, we report results with FAAF models trained with $\beta = 10$ in Table 1 and Table 2

E Friction Intervention Evaluation Prompts and Sampled Representative Interventions

2088

2089

2090

2104

2105

Fig. 7 shows prompt used for friction intervention 2091 assessments in an LLM-as-a-judge format. We use 2092 a standard format (Cui et al., 2024) but adapt fric-2094 tion preference dimensions to collaborative taskspecific settings. This prompt systematically scores 2095 friction interventions on 7 target dimensions of 2096 friction intervention quality such as correct reason-2097 ing, consistency with the agent's justification for 2098 friction, alignment with golden friction samples, 2099 clarity etc. For sampling from GPT-40, we use 2100 standard settings with a nucleus sampling parame-2101 ter (top-p) (Holtzman et al., 2019) and temperature 2102 of 1. 2103

> Tables 9–12 show some representative interventions from each baseline and FAAF.

PAIRWISE LLM-AS-A-JUDGE EVALUATION PROMPT: FRICTION INTERVENTIONS

System: You are an expert evaluating the quality of friction interventions in collaborative problem-solving.

Game-definition: Participants (P1, P2, P3) are solving a block-weighing puzzle. They can only weigh two blocks at a time and know the red block is 10g. They must determine weights of all blocks (blue=10g, green=20g, purple=30g, yellow=50g) but don't know these values initially. A friction intervention is an indirect persuasion statement that prompts self-reflection and reevaluation of assumptions, like asking "Are we sure?" or suggesting to revisit steps. You must rate each intervention (between 1 to 5) along these **dimensions** given the json format below.

[Dialogue] [Gold intervention] [Intervention A] [Rationale A] [Intervention B] [Rationale B]

You must a choice between which of two interventions is more preferable and provide one sentence explanation at the end.

1. Relevance: How well does the intervention address key issues or assumptions in the reasoning process?

2. Gold Alignment: How well does the friction intervention align with the golden friction sample?3. Actionability: Does the friction intervention provide actionable guidance or suggest concrete steps for participants to improve their reasoning?

4. Rationale Fit: How well does the provided rationale align with the preference for the friction intervention?

5. Thought-Provoking: Encourages self-reflection

6. Specificity: Does the intervention pinpoint specific flaws, assumptions, or gaps?

7. Impact: To what extent does the friction intervention have the potential to change the course of the participants' reasoning?

Format your response as follows:

A: relevance: [1-5], gold_alignment: [1-5], actionability: [1-5], rationale_fit: [1-5], thought_provoking: [1-5], specificity: [1-5], impact: [1-5]B: similar format Winner: ['A' or 'B'] Rationale: [One sentence explanation]

Figure 7: Evaluation prompt used for friction intervention assessments in an LLM-as-a-judge format.

Prompt	You are an expert in collaborative task analysis and reasoning. Your task is to analyze the dialogue history involving three participants (P1, P2, and P3) trying to deduce the weights of certain blocks. For each dialogue:
Gold (GPT-40) Friction	Let's not assume anything about the yellow block's weight until we weigh it against the purple one!
Model Outputs FAAF	Belief State: P1 believes the green block weighs 20 grams, P2 is unsure about the exact weight of the green block, and P3 is cautious about making assumptions about the purple block's weight.Friction: Wait, are we jumping to conclusions about the green block's weight? We've only confirmed it's less than 30 grams, but we haven't ruled out the possibility that it's actually 15 or 25 grams. Let's make sure we're not overestimating its weight before moving on to the next block.
DPO	 Belief State: P1 believes that the blue block weighs 10 grams since it balances with the red block. P2 is considering the possibility that the green block might weigh 20 grams since it feels heavier, but is uncertain and open to further measurement. P3 is considering the weight of the purple block, is curious, and wants to measure it to confirm its weight, but is not assuming a specific value yet. Friction: Hold on, maybe we're jumping to conclusions about the green block's weight. What if it's not 20 grams? Can we really confirm that by just feeling it? Let's consider other possibilities and make sure we're not relying on unverified assumptions.
PPO	Belief State:P1 believes the green block is probably 20 grams based on its perceived "heft" (assumed connection between weight and feeling). P2 is cautious and thinks they shouldn't jump to conclusions. P3 is open-minded and willing to consider multiple options.Friction:You know, we're getting ahead of ourselves by assuming the green block weighs 20 grams based on its size or feeling. Can we take a step back and really think about this? We've already found one block with the same weight as the red block. What if we weigh the green block against something else, like the yellow block? That would give us a better sense of where its weight really falls.
SFT	Belief State: P1 believes the yellow block is likely 50g. P2 is unsure about the yellow block's weight but thinks it might be 50g. P3 is unsure about the yellow block's weight and thinks it might be 50g or heavier. Friction: What if we measure the yellow block against the green block? That would give us some clarity about its weight.

Table 9: Comparison of model-generated friction interventions on the evaluation prompts of the Simulated WTD dataset.

Prompt	You are an expert in collaborative task analysis and reasoning. Your task is to analyze the dialogue history involving three participants (P1, P2, and P3) trying to deduce the weights of certain blocks. For each dialogue: <belief_state> Identify reasoning flaws or misunderstandings. </belief_state> <rationale> Justify the need for intervention and its impact. </rationale> <friction> Generate an intervention to prompt reflection and alignment. </friction> User Dialogue: P1: this is nine rb so P3: oh no they don't have any hints here P3: Eighty gram sounds seems like a ninety gram to me but P3: it's cool we have to i guess eighty alright P2: Eighty let's let's go for eighty Yeah P3: I would go for ninety but you guys are going for eighty Alright Cool Okay P1: I would have gone for one thirty P2: Eighty there you go Yeah P4: For the final part of this test read the scenario and estimate the weight of the missing second mystery block P2: what's this P4: You will have two attempts to estimate and explain your answer you can now continue the survey P1: can you read it out loud or something P3: well the factory that creates the wooden blocks did not send us the second mystery block P2: Okay P3: you have two opportunities to estimate the weight of the second block explain in one two sentences how you determined your estimate indicate the weight of the second mystery block P1: fibonacci like let's say that's the first one according to our translate thirteen one thirty P2: eight mystery block 's the first one so eight plus five thirteen yeah there is one thirty eight plus five P3: is mystery block the first mystery block P2: Yeah P3: and we already did that now this is yeah oh okay okay okay okay right P2: yeah the second one means according to that statement is the second one mystery block two P2: So eight plus five one thirty P1: makes sense to me as much sense as i can make of it #### Assistant:
Gold (GPT-40) Friction	Let's not assume anything about the yellow block's weight until we weigh it against the purple one!
Model Outputs FAAF	 Belief State: P1 believes that the second mystery block is 130 grams, based on the Fibonacci sequence and the weight of the first mystery block. P2 believes that the second mystery block is also 130 grams, following P1's reasoning. P3 is unsure and questioning the assumption that the second mystery block is the same as the first mystery block. Friction: Hey guys, I think we're making a mistake here. We already estimated the weight of the first mystery block, and now we're assuming the second mystery block is the same? That doesn't seem right. Can we take a step back and think about this? Maybe there's another way to approach this?
DPO	 Belief State: P1 believes the pattern is linear and straightforward, assuming each block weighs 5 grams more than the previous one (e.g., 13 = 8 + 5, 18 = 13 + 5). P2 also assumes a linear pattern but thinks the mystery block is the second one. P3 is skeptical about the pattern, especially the value 80. Friction: Wait a minute, are we too quick to assume the pattern is linear? What if there's another factor at play? For instance, what if the weights are related to the block's surface area or material composition? Could our assumptions be leading us down the wrong path?
PPO	Belief State:P1 believes that the first mystery block has a weight of 130g, as they mentioned "fibonaccilike, let's say that's the first one, according to our translate, thirteen, one thirty". P2 believes that the secondmystery block has a weight of 135g, as they calculated "eight plus five, thirteen, yeah, there is one thirty eightplus five". P3 is unsure and hesitant, questioning the accuracy of the previous estimates.Friction:Hey, let's take a step back and re-examine our assumptions.Before we jump to conclusions, canwe re-read the instructions and see if there's any additional information we missed? Maybe we can approachthis problem from a different angle, rather than relying on our initial guesses?
SFT	Belief State: P1 believes the second mystery block is likely 130 grams based on the Fibonacci sequence. P2believes the second mystery block is 80 grams based on the initial estimates. P3 believes the second mysteryblock could be either 80 grams or 90 grams, expressing confusion about the weight.Friction: I think we might have misinterpreted the blocks again.

Table 10: Comparison of model-generated friction interventions on the evaluation prompts of the Original WTD dataset.

Prompt	You are an expert in collaborative task analysis and reasoning. Participants must test the rule: "All cards with vowels have an even number on the other side." A common mistake is verifying only one direction—ignoring the need to check whether odd-numbered cards might have vowels. This incomplete reasoning risks false validation. For each dialogue: <belief_state> Identify reasoning flaws or misunderstandings. </belief_state> <rationale> Justify the need for intervention and its impact. </rationale> <friction> Generate an intervention to prompt reflection and alignment. </friction> User Dialogue: Alpaca: The task we just submitted Narwhal: "the instructions aren't clear" Falcon: Which card did you choose? I chose U and 4 Narwhal: So which cards to turn? Alpaca: So what did everyone submit? I did U and 5. Duck: "I agree but we're supposed to flip the U and 4, right?" Narwhal: U and B Duck: Why would you flip B? Narwhal: Because they are both letters. Falcon: We should choose cards with vowels or an even number, I guess just U and 4. Narwhal: U Narwhal: U Alpaca: Yes, we all said the letter U is one of the cards to flip. ### Assistant:
	THE ASSISTANT.
Gold (GPT-40) Friction	How do you see consonants related to the even number requirement? Is it necessary in our current tests?
Model Outputs FAAF	 Belief State: Alpaca and Narwhal seem to be focusing on the "vowel" part of the rule, but not considering the "even number" part. Duck is trying to correct Narwhal, but Narwhal is not listening to the correction. Falcon is partially correct in suggesting that they should choose cards with vowels or even numbers, but doesn't specify which ones. Friction: Hey, wait a minute! We're only considering one side of the rule. What about the other side? What if we flip a card with an odd number and it has a vowel? Wouldn't that mean we need to flip it too?
DPO	Belief State: None Friction: Can we clarify the reasoning behind our choices? Let's try to think about both directions of the rule. What would happen if we choose a card with an odd number? Would it necessarily mean it doesn't have a vowel, or could it still have a vowel? How about the other direction? If we choose a card with a vowel, would it necessarily have an even number?
PPO	Belief State: Narwhal states that 'B' is not a vowel, which implies a misunderstanding of the rule regarding non-vowel cards. This affects reasoning as it leads to a potential oversight of the necessity to test odd-numbered cards, which could lead to incorrect conclusions about the rule. Friction: What does that mean for the other cards we should check?
SFT	Belief State: Duck states 'the other one is 5,' which implies they believe the 4 is already accounted for, while it is unclear whether 4 is being considered as part of the rule. The lack of clarity on whether the 4 is being considered could lead to confusion about the relevance of the odd-numbered card in the context of the rule. Friction: Could we clarify if the 4 is being considered as part of the rule or not?

Table 11: Comparison of model-generated friction interventions on the evaluation prompts of the DeliData dataset.

Prompt	You are an expert in collaborative task analysis and reasoning. Participants must test the rule: "All cards with vowels have an even number on the other side." A common mistake is verifying only one direction—ignoring the need to check whether odd-numbered cards might have vowels. This incomplete reasoning risks false validation. For each dialogue: <belief_state> Identify reasoning flaws or misunderstandings. </belief_state> <frictionale> Justify the need for intervention and its impact. <friction> Generate an intervention to prompt reflection and alignment. </friction> User Dialogue: Dolphin: "thinking about it, I guess it wouldn't matter what the C has on it" Llama: would be E be one to turn? Dolphin: "wouldn't you want to also check the 5, to test the theory? because if there is a vowel on the other side, that would definitely need to be turned Kiwi: Oh true, the 5 might need to be turned as well Dolphin: E would definitely need to be turned Kiwi: Oh true, the 5 might need to be turned as well Dolphin: on, maybe E & 5 is the answer? Llama: I agree with the 5. unsure with the 6, but i dont think C should be turned Llama: E and 5 seems like the best answer to me as well Kiwi: Yes E and 5 Dolphin: ok, that is what i am going to chose, e & 5 Kiwi: 5 Kiwi: 5 Llama: same Llama: E,5 Dolphin: it says we have to participate in the chat for at least 5 minutes before making our final selection. have we been chatting for 5 minutes? ### Assistant:</frictionale>
Gold (GPT-40 Friction	What makes you think 'just E' is sufficient? Can you clarify your reasoning?
Model Outputs FAAF	Belief State: The participants seem to be engaging in a discussion about the rule "All cards with vowels have an even number on the other side." However, there appears to be a misunderstanding about the directionality of the rule. They are focusing on testing the rule by checking only one direction (vowels to even numbers) without considering the other direction (even numbers to vowels). Friction: Can we flip both sides of the coin to test the rule? Let's not just focus on one direction. What if we assume the rule holds true for odd numbers as well? How would that affect our understanding of the situation?
DPO	Belief State: It appears that the participants are making progress in their understanding of the rule, but there is still a potential contradiction in their reasoning. Friction: Hey team, let's reflect on our reasoning. When we say 'all cards with vowels have even numbers,' are we only considering one direction or both? Are we making assumptions about certain cards being irrelevant? Can you think of any scenarios where an odd-numbered card might have a vowel on the other side? Let's take a moment to consider all possible outcomes before making our final decision.
PPO	Belief State: Kiwi suggests turning the 5 and E, but there is a misunderstanding about the importance of the 5 in relation to the rule. The implication is that Kiwi believes the 5 is crucial to test the rule, which could lead to incorrect conclusions if the group doesn't fully understand why the 5 is necessary. Friction: What do you think is the significance of the 5 in relation to the rule?
SFT	Belief State: Kiwi's statement 'E and 5' suggests a final conclusion without addressing potential errors, while it doesn't acknowledge the need for a more thorough discussion. This creates a misalignment in the reasoning process, as Kiwi may not fully consider the implications of their conclusion, potentially leading to incorrect submissions.Friction: Could you explain why you're confident in choosing 'E and 5' without further discussion?

Table 12: Comparison of model-generated friction interventions on the evaluation prompts of the DeliData dataset.