

# Bayesian Preference Elicitation with Language Models

Anonymous ACL submission

## Abstract

001 There is increasing interest in using language  
002 models (LMs) not only for answering queries  
003 but also for *gathering information* about the  
004 preferences of human users. This preference  
005 data can be used to fine-tune LMs via re-  
006 ward modeling or completing goal-oriented  
007 tasks. However, LMs have been shown to  
008 struggle with crucial aspects of preference  
009 learning: quantifying uncertainty, modeling  
010 mental states, and posing highly informative  
011 questions. These challenges have been ad-  
012 dressed in other areas of machine learning,  
013 such as Bayesian Optimal Experimental De-  
014 sign (BOED), which focuses on designing in-  
015 formative queries within a well-defined feature  
016 space. But these methods, in turn, have his-  
017 torically been difficult to scale and apply to  
018 real-world problems (e.g. involving text and  
019 images), in which simply identifying the rel-  
020 evant features can be challenging. We intro-  
021 duce **OPEN** (Optimal Preference Elicitation  
022 with Natural language) a framework that uses  
023 BOED to guide the choice of informative ques-  
024 tions and an LM to extract features and translate  
025 abstract BOED queries into natural language  
026 questions. By combining the flexibility of LMs  
027 with the precision of BOED, OPEN can opti-  
028 mize queries for informativity while remaining  
029 adaptable to real-world domains. Conducting  
030 user studies, OPEN outperforms existing LM-  
031 and BOED-based methods for preference elicita-  
032 tion.

## 033 1 Introduction

034 Understanding users’ complex preferences and re-  
035 quirements is necessary for accurately and safely  
036 automating real-world tasks. Modeling preferences  
037 in a domain of interest requires both understand-  
038 ing *what features* of the domain are relevant to  
039 model, and *how important* these features are rela-  
040 tive to each other (Lin et al., 2022; Lindner et al.,  
041 2022; Sadigh et al., 2017; Fürnkranz and Hüller-  
042 meier, 2012). For example, when building a con-

043 tent recommendation system, relevant features may  
044 be different article topics—e.g. science, politics, or  
045 celebrity culture. A user may decide whether or  
046 not to read an article based on the strength of their  
047 preferences along these different axes.

048 With the widespread use of preference data to  
049 align language models (LMs) with human users,  
050 there has been growing interest in using LMs them-  
051 selves to elicit information about human prefer-  
052 ences. Past work focusing on prompting LMs  
053 to ask questions about user preferences (Li et al.,  
054 2023; Lin et al., 2023; Piriyaakulkij et al., 2023)  
055 has shown that LMs are capable of identifying  
056 decision-relevant features in domains spanning con-  
057 tent recommendation, software engineering, and  
058 moral reasoning. Nonetheless, prompting offers  
059 limited control over LMs’ information-gathering  
060 strategies and often fails to produce questions that  
061 are useful or informative.

062 By contrast, a long line of work on optimal ex-  
063 perimental design, such as dueling bandits (Heckel  
064 et al., 2019) and Bayesian preference learning ap-  
065 proaches (Foster et al., 2019; Gal et al., 2017;  
066 Houlshby et al., 2011), has developed methods to  
067 efficiently infer user preferences from a limited  
068 number of interactions. But these methods also  
069 have their limitations—they have primarily been  
070 applied to tasks with simple, highly-structured data  
071 where feature spaces are constrained and interac-  
072 tions are short. How can we get the best of both  
073 worlds? *Can we learn similarly sample-efficient*  
074 *models of user preferences in complex, open-ended*  
075 *tasks?*

076 In this paper, we introduce **Optimal Preference**  
077 **Elicitation with Natural language (OPEN)**—a  
078 framework that leverages the complementary as-  
079 pects of LMs and Bayesian Optimal Experiment  
080 Design (BOED) methods (see Figure 1).

081 OPEN uses an LM to select the relevant features  
082 the user likely cares about (Figure 1, *steps 1, 2*),  
083 and a Bayesian model to select optimal pairwise

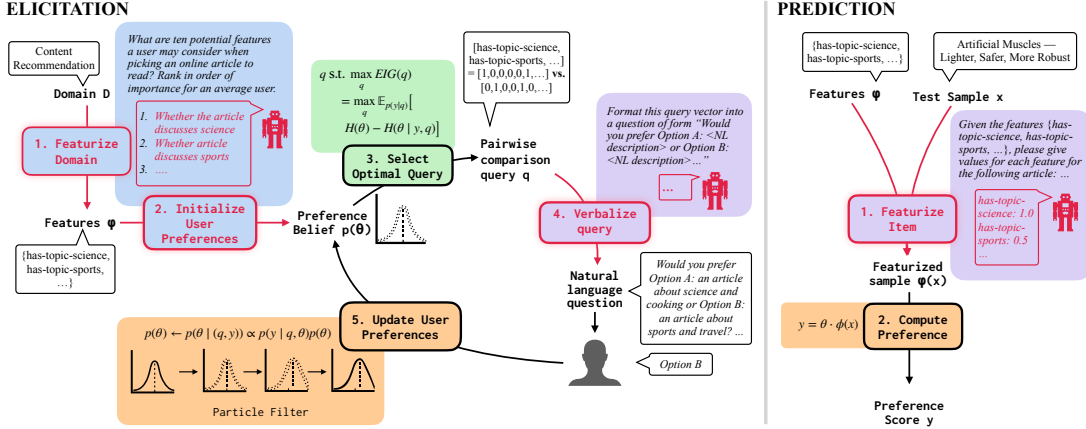


Figure 1: Overview of the OPEN framework. In red are the parts where we use a language model. During the *elicitation* stage, first, a domain  $D$  is featurized into feature  $\phi$  with a language model, which also gives us a ranking of importance over features (which is used to initialize a prior  $p(\theta)$  over user preferences). Based on the prior user preferences, the optimal pairwise comparison query  $q$  is computed, which is then verbalized using an LM into natural language. The user response is then taken to update the prior over beliefs. During the *prediction* stage, a LM is used to featurize a test sample according to the featurization  $\phi$  derived from the elicitation stage, then a preference score is computed using the elicited preferences  $\theta$ .

084 comparison queries (step 3), which are translated  
 085 into natural language by a LM (step 4). In OPEN,  
 086 we utilize a LM to provide feature coverage and a  
 087 natural conversational interface with the user (steps  
 088 1, 2, 4), while leveraging a Bayesian model to track  
 089 feature weightings and select informative questions  
 090 (steps 3, 5). Using OPEN to elicit user preferences  
 091 in a content recommendation domain, we find that  
 092 it outperforms both LM- and BOED-based prefer-  
 093 ence elicitation approaches.

## 094 2 Preliminaries & Background

095 We aim to actively model human preferences  
 096 through iterative interactions with a user. Over  
 097 the course of a conversation, we construct and up-  
 098 date a model through a series of user queries, each  
 099 designed such that they (1) **maximize information**  
 100 **gain** of the users’ preferences and (2) can reliably  
 101 be answered by users and **elicit informative re-**  
 102 **sponses**.

### 103 2.1 Bayesian Optimal Experimental Design

104 We use Bayesian Optimal Experimental Design  
 105 (BOED; Lindley, 1956, 1972; Rainforth et al.,  
 106 2023) a mathematical abstraction for selecting opti-  
 107 mal queries in an active learning setup, as the basis  
 108 for our interaction framework. BOED is a com-  
 109 mon approach to preference learning and function  
 110 estimation (Foster et al., 2019; Gal et al., 2017;  
 111 Houthby et al., 2011) and has been used across

112 many different domains and disciplines (Cavagnaro  
 113 et al., 2009; Dehideniya et al., 2018; Dushenko  
 114 et al., 2020; Vanlier et al., 2012).

115 In BOED, the goal is to select an experimental  
 116 design that maximizes the information gain for pa-  
 117 rameters of interest. We begin with a predictive  
 118 model  $p(y | q, \theta)$  which defines the relationship be-  
 119 tween an experiment  $q$ , an experimental outcome  
 120  $y$ , and the parameters of interest  $\theta$ . The prior  $p(\theta)$   
 121 describes an initial belief about the parameters.

122 We can formalize the information we gain about  
 123  $\theta$  from each experiment as:

$$124 \text{IG}(q, y) = H[\theta] - H[\theta | y, q], \quad (1)$$

125 where  $H$  is the Shannon entropy. However, be-  
 126 cause this notion of information gain relies on  $y$ ,  
 127 the outcome of the experiment, we cannot select  
 128 the optimal experiment until after the outcome has  
 129 been observed. To find the optimal design, we in-  
 130 stead take the expectation of the information gain  
 131 for the marginal distribution of  $y$  over all possible  
 132 outcomes  $p(y|q) = \mathbb{E}_{p(\theta)}[p(y|\theta, q)]$ ,<sup>1</sup> yielding:

$$133 \text{EIG}(q) = \mathbb{E}_{p(y|q)}[\text{IG}(q, y)]. \quad (2)$$

134 Thus, the optimal experiment can be defined as:

$$135 q^* = \text{argmax}_q(\text{EIG}(q)). \quad (3)$$

<sup>1</sup>This is equivalent to the mutual information between  $\theta$  and  $y$  given  $q$ .

136	In our preference learning setting, we want to	3. <b>Selecting an Optimal Question:</b> The	184
137	pick the pairwise comparison question ( $q$ ) that will	Bayesian model samples all possible pair-	185
138	allow us to learn about the user’s preferences ( $\theta$ ) as	wise comparison questions and selects the	186
139	efficiently as possible. At each interaction, $y$ is the	maximally informative question, using Equa-	187
140	user’s response to $q$ . BOED is uniquely powerful	tion (3).	188
141	in our sequential setting—enabling us to, after each		
142	interaction, incorporate users’ answers ad-hoc to	4. <b>Querying the User:</b> The LM translates $q^*$	189
143	greedily update our belief of the their preferences.	into a NL question for the user (this can also	190
144		be thought of as $\phi^{-1}$ ).	191
145		5. <b>Posterior Update:</b> Given the user’s response	192
146		to the question, we compute the Bayesian pos-	193
147		terior $p(\theta   y, q)$ .	194
148		6. <b>Prediction:</b> The Bayesian model predicts the	195
149		user’s response to each test case, using the	196
150		Bradley–Terry model of human preferences.	197
151			
152		Even for the linear preference model described in	198
153		Section 2.2, computation of the posterior $p(\theta   y, q)$	199
154		(and thus computation of $EIG(q)$ ) is intractable.	200
155		We implement a tractable approximation using	201
156		a Bayesian Particle Filter (Gordon et al., 1993;	202
157		Doucet and Johansen, 2008; Elfring et al., 2021).	203
158		As each particle can be thought of as a plausible	204
159		instantiation of user preferences $\theta$ , we refer to indi-	205
160		vidual particles as <b>personas</b> $\mathbf{p}_i$ .	206
161		Crucially, though, the OPEN does not neces-	207
162		sitate the use of particle filtering or even BOED.	208
163		Our framework fundamentally provides a domain-	209
164		agnostic approach to active and principled pref-	210
165		erence learning: enabling any uncertainty-based	211
166		preference-learning algorithm (multi-armed ban-	212
167		dots, Bayesian active learning, etc.) to interface	213
168		with a user in NL.	214
169		Below, we detail our implementation of the afor-	215
170		mentioned components of OPEN.	216
171			
172		3.1 <b>Featurization</b>	217
173		We prompt the LM with a general description of	218
174		the domain of interest (e.g. content recommenda-	219
175		tion of news articles), and ask it to produce natural	220
176		language descriptions of pertinent features in that	221
177		domain, $\mathcal{F}$ . For example, in the content recommen-	222
178		dation domain, the LM outputs “ <i>whether the article</i>	223
179		<i>discusses science</i> ” as a pertinent feature. This list	224
180		of natural language features can then be converted	225
181		to a function $\phi$ that transforms test samples to fea-	226
182		ture values by prompting another LM with $\mathcal{F}$ and	227
183		a test sample $x$ (see Section 3.6).	228
184			
185		3.2 <b>Initializing User Preferences</b>	229
186		When prompting the LM for features in Section 3.1,	230
187		we simultaneously ask it to <i>rank</i> features from <i>most</i>	231
188			
189			
190			
191			
192			
193			
194			
195			
196			
197			
198			
199			
200			
201			
202			
203			
204			
205			
206			
207			
208			
209			
210			
211			
212			
213			
214			
215			
216			
217			
218			
219			
220			
221			
222			
223			
224			
225			
226			
227			
228			
229			
230			
231			

232 to *least* important. These rankings are then used to  
 233 initialize the prior for OPEN’s belief of the user’s  
 234 preference,  $p(\theta)$ . Because we assume that  $\theta$  is  
 235 linear in the feature space, we model  $p(\theta)$  as a  
 236 Bayesian linear model with each feature weight  
 237 parameterized by a normal distribution:

$$238 \quad \begin{aligned} p(\theta_i) &\sim \mathcal{N}(0, \sigma^2) \cdot w_i, \\ w_i &= 1.2 - 0.12i \end{aligned} \quad (4)$$

239 where  $\sigma^2$  is the base variance for each feature be-  
 240 fore scaling, and  $w_i$  captures the relative impor-  
 241 tance of features from highest to lowest.

242 We construct our initial sample of  $N$  personas  
 243 by sampling each element of each persona independ-  
 244 dently from  $p(\theta_i)$ .

### 245 3.3 Selecting the Optimal Question

246 To efficiently learn user preferences, we aim to  
 247 select the pair of options that maximize the EIG  
 248 regarding the user’s preferences at each interaction  
 249 step. This process is grounded in the BOED frame-  
 250 work, as described earlier. Given a set of features,  
 251  $\mathcal{F}$ , we define the space of all possible pairwise  
 252 comparison questions as:

$$253 \quad \mathcal{Q} = \{(\mathbf{o}_a, \mathbf{o}_b) \mid \mathbf{o}_a, \mathbf{o}_b \in \{0, 1\}^{|\mathcal{F}|}, \\ \sum_{i=1}^{|\mathcal{F}|} (\mathbf{o}_a)_i = \sum_{j=1}^{|\mathcal{F}|} (\mathbf{o}_b)_j = K\},$$

254 where each option  $\mathbf{o}_a$  and  $\mathbf{o}_b$  is represented as a  
 255 binary vector indicating the presence (1) or absence  
 256 (0) of each feature in the comparison. In our exper-  
 257 iments, we set  $|\mathcal{F}| = 10$  and  $K = 2$ .<sup>2</sup> The goal is  
 258 to select the question  $(\mathbf{o}_a^*, \mathbf{o}_b^*) \in \mathcal{Q}$  that maximizes  
 259 the EIG about the user’s preferences  $\theta$ .

260 Recalling Equation 3, the EIG for the optimal  
 261 pairwise comparison  $(\mathbf{o}_a^*, \mathbf{o}_b^*)$  with which to query  
 262 the user is then given by:

$$263 \quad (\mathbf{o}_a^*, \mathbf{o}_b^*) = \operatorname{argmax}_{(\mathbf{o}_a, \mathbf{o}_b) \in \mathcal{Q}} \operatorname{EIG}(\mathbf{o}_a, \mathbf{o}_b)$$

### 264 3.4 Querying the User: Mapping Pairwise 265 Comparisons to NL

266 Given the optimal pairwise comparison, we prompt  
 267 a LM to convert the feature vectors into a NL ques-  
 268 tion for the user. This is equivalent to taken the  
 269 inverse of the featurization function  $\phi^{-1}$ . We in-  
 270 clude the set of NL features from the LM (Section

<sup>2</sup>From preliminary testing, we found that including ten total features ( $|\mathcal{F}| = 10$ ) and two features ( $K = 2$ ) in each option, best balanced users’ mental effort with the informativity of the question.

3.1) and the  $K$  present features. We prompt the  
 LM with the feature set and pairwise comparison  
 to develop a real-world example for each option ad-  
 hoc. For example, in our content recommendation  
 setting, this is an example article summary. We dis-  
 play these examples in the user interface alongside  
 each option in the pairwise comparison question.

### 278 3.5 Posterior Update

279 At each interaction step  $t$ , having observed the  
 280 user’s selection  $y_t$  between  $(\mathbf{o}_a, \mathbf{o}_b)$ , we update  
 281 our belief about the user’s preferences ( $\theta$ ) based  
 282 on their response ( $y_t$ ) to the pairwise comparison  
 283 question ( $q_t$ ) by computing the posterior:

$$284 \quad p(\theta \mid y_t, q_t, \mathcal{H}) = \frac{p(y_t, q_t \mid \theta) \cdot p(\theta \mid \mathcal{H})}{p(y_t, q_t, \mathcal{H})}, \quad (5)$$

285 where  $\mathcal{H} = \{(q, y)\}_{0 \dots t-1}$  denotes the conversa-  
 286 tion history (prior sequence of questions and an-  
 287 swers). We approximate the update in Equation (5)  
 288 using our particle filter. At each step  $t$ , we reweight  
 289 each particle  $\mathbf{p}_i$  by

$$290 \quad \begin{aligned} w_i &= p(y_t \mid \mathbf{o}_a, \mathbf{o}_b, \mathbf{p}_i^\top) \\ &= \sigma(\mathbf{p}_i^\top \phi(\mathbf{o}_a) - \mathbf{p}_i^\top \phi(\mathbf{o}_b)). \end{aligned} \quad (6)$$

291 Then, we fit a Gaussian to the distribution of  $w_i \mathbf{p}_i$ s.  
 292 Finally, we resample new personas from this distri-  
 293 bution. For a detailed overview of our algorithm,  
 294 see Appendix B.

### 295 3.6 Prediction

296 After each user interaction, we evaluate our model  
 297 via a set of test cases, which are also presented to  
 298 the user at the end of the study. For our test cases,  
 299 we use pairwise comparison questions  $(x_a, x_b)$ ,  
 300 where each  $x_a, x_b$  are taken from real-world ex-  
 301 amples from the domain (e.g. for content recom-  
 302 mendation, the headline and lede of a New York  
 303 Times news article). For each test-case question,  
 304 we select  $y \in [x_a, x_b]$  according to:

$$305 \quad \max_y \mathbb{E}_\theta [p(y \mid \theta, (x_a, x_b))]$$

306 (where the expectation is taken over our distribu-  
 307 tion of  $p(\theta)$  at that point in time). Following the  
 308 Bradley-Terry model, we define preferences as

$$309 \quad p(x_a \mid \theta, (x_a, x_b)) = \sigma(\theta \cdot \phi(x_a) - \theta \cdot \phi(x_b))$$

$$310 \quad p(x_b \mid \theta, (x_a, x_b)) = 1 - p(x_a \mid \theta, (x_a, x_b))$$

311 To compute  $\phi(x_a)$  and  $\phi(x_b)$  (recall we were given  
 312 a set of natural language descriptions of features

313	$\mathcal{F}$ ), we prompt a LM to give us concrete values	4.3 Human Experiment Details	359
314	for each features, based on the natural language	We recruit native-English-speaking human partici-	360
315	feature descriptions and the test-set sample.	pants from Prolific (Palan and Schitter, 2017). We	361
316	<b>4 Experimental Setup</b>	recruited $\sim 40$ people for each setting. <sup>3</sup> Mirroring	362
317	We use OPEN to investigate how incorporat-	Li et al. (2023), we conducted our prolific experi-	363
318	ing LMs alongside explicit information-theoretic	ments in two steps: (1) <i>elicitation</i> : participants	364
319	computations (via BOED) can improve existing	were asked to answer questions about their pref-	365
320	preference-learning methodology and provide a	erences for 5 minutes; (2) <i>prediction</i> : participants	366
321	novel NL interface with which to approach per-	were then presented with the 15 pairwise compari-	367
322	sonalized machine learning.	son questions and asked which of the two options	368
323	We outline OPEN’s experimental details, the	they preferred. We recruited new participants for	369
324	three different baselines which we compare to, the	each test across all experiments.	370
325	evaluation framework, and the human participants’	<b>4.4 Domain: Content Recommendation</b>	371
326	interface.	We consider the task of content recommendation–	372
327	<b>4.1 Hyperparameters</b>	recommending news articles for a user to read.	373
328	In the featurization step, we query the LM for	Each approach is evaluated based on its ability to	374
329	$ \mathcal{F}  = 10$ features to define the domain. When	predict which news articles the user would prefer	375
330	sampling pairwise comparison questions, each op-	on a set of 15 pairwise comparisons. Each com-	376
331	tion contains $K = 2$ features. We use GPT-4 as the	parison contains the lede and headline from New	377
332	LM for all of our experiments.	York Times articles, hand-collected by the authors.	378
333	<b>4.2 Baselines</b>	We select this task (1) because people’s preferences	379
334	We consider three baselines:	have a large amount of variance and (2) to perform	380
335	<b>LM-only Open-Ended Questions</b> Following Li	a direct comparison to Li et al. (2023)’s method-	381
336	et al. (2023)’s best-performing method in their	ology. Further information about the articles is	382
337	content-recommendation setting, we prompt a LM	available in Appendix C.	383
338	to ask informative, open-ended questions, without	<b>4.5 Evaluation</b>	384
339	any explicit optimization. The LM is provided the	We evaluate OPEN and the above baselines after	385
340	complete conversation history of the interaction be-	each user interaction on a test set of 15 pairwise	386
341	fore asking each subsequent question. The user is	comparison questions, each comparing two differ-	387
342	able to answer the question free-form.	ent real-world news articles.	388
343	<b>LM-only Pairwise Comparison Questions</b> We	<b>Methods</b> We evaluate using two different meth-	389
344	prompt a LM to ask an informative, pairwise com-	ods: in <b>OPEN prediction</b> , we use the OPEN	390
345	parison question, without any explicit optimization.	prediction method described in Section 3.6 to an-	391
346	The LM is provided the complete conversation his-	swer each test set question. In <b>LM prediction</b> ,	392
347	tory of the interaction and the feature set from Sec-	we prompt a LM to answer the test question, con-	393
348	tion 3.1 before asking each subsequent question. To	ditioned on the conversation history (up to some	394
349	answer the question, the user selects either ‘Option	turn), similar to the evaluation approach taken by	395
350	A’ or ‘Option B’ based on their preference.	Li et al. (2023).	396
351	<b>User Self-Mapping Pairwise Comparisons to NL</b>	<b>Metrics</b> After each conversation turn, we evalu-	397
352	We provide the user the feature vectors of the opti-	ate the test-set accuracy at that point in time, using	398
353	mal pairwise comparison question as determined	one of the two methods above. After the conversa-	399
354	from Section 3.3 as well as the NL descriptions of	tion has ended, we calculate the <b>Time-Integrated</b>	400
355	each feature. The user must use these descriptions	<b>Delta Accuracy (TIDA)</b> , or the integral of the	401
356	to interpret the vector comparison query, then se-	“ $\Delta$ accuracy over time curve”, where $\Delta$ accuracy is	402
357	lect either ‘Option A’ or ‘Option B’ based on their	defined as the difference between test-set accuracy	403
358	preference.	<sup>3</sup> Additional details regarding the user studies are available	
		in Appendix D	

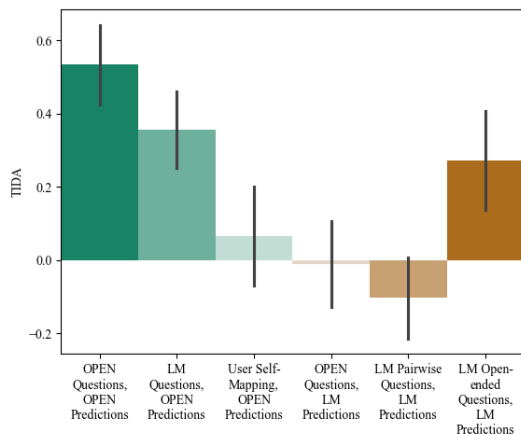


Figure 2: **Time-Integrated Delta Accuracy (TIDA) for each method.** We report the integral of the delta accuracy over time. OPEN improves over naive prompting-based approaches, all while improving transparency and reducing computational cost. Error bars are one standard error.

404 at the *current* interaction step and the initial inter-  
 405 action step. This metric rewards us for eliciting  
 406 the user preferences quickly—even if two elicitation  
 407 transcripts result in same  $\Delta$  accuracy, the one that  
 408 arrives at a higher  $\Delta$  accuracy earlier in time is  
 409 rewarded.

## 410 5 Results

411 We present the TIDA scores for OPEN against five  
 412 different combinations of question generation and  
 413 prediction in Figure 2 (refer to Sections 3 and 4.2  
 414 for question generation methods and Section 4.5  
 415 for prediction methods). Below, we discuss the  
 416 main takeaways from our results.

417 **OPEN is better at making predictions aligned**  
 418 **with human preferences when compared to LM-**  
 419 **only approaches.** Comparing the *OPEN Questions,*  
 420 *OPEN Predictions* bar to *OPEN Questions,*  
 421 *LM Predictions* bar, we find that OPEN can make  
 422 predictions that better align with human prefer-  
 423 ences than a prompted LM, from the same sequence  
 424 of questions. This implies that LMs are still subpar  
 425 at *in-context-learning* of human preferences from  
 426 demonstrations, compared to Bayesian methods  
 427 which explicitly keep track of the human’s pref-  
 428 erence function; it underscores the importance of  
 429 explicitly monitoring human preferences.

430 **OPEN is better at eliciting human preferences**  
 431 **compared to LM-only approaches.** Comparing

the *OPEN Questions, OPEN Predictions* bar to  
 the comparable *LM Pairwise Questions, OPEN*  
*Predictions* bar, we find that *OPEN Pairwise Questions*  
 outperforms *LM-only Pairwise Questions* for  
 elicitation, indicating that when asking questions,  
 there is value in tracking human preferences and  
 querying based on explicit notions of uncertainty.

Of course, one benefit of LM elicitation ap-  
 proaches is that LMs are not constrained to asking  
 only a certain type of question. Thus, as a top-  
 line, we elicit open-ended questions from LMs (*LM*  
*Open-ended Questions, LM Predictions*). Though  
 open-ended questions far outperform pairwise ques-  
 tions (both ones selected by OPEN and selected  
 by the LM) when evaluated with a LM, this im-  
 provement is overshadowed by the benefit of being  
 able to use *OPEN prediction*, which is incompat-  
 ible with LM-only elicitation methods. However,  
 this suggests that future lines of work could ex-  
 plore eliciting open-ended questions compatible  
 with OPEN predictions.

**OPEN is better at eliciting human preferences**  
**compared to BOED-only approaches.** When  
 we remove the verbalization step (*User Self-*  
*mapping, OPEN Predictions*), we find that users  
 become significantly worse at accurately specifying  
 their preferences. Thus, having a natural interface  
 for interacting with the user, which we accomplish  
 by verbalizing questions with a LM, is essential for  
 being able to accurately and efficiently elicit user  
 preferences.

## 463 6 Analysis

**The Importance of Feature Weightings** To un-  
 derstand the importance of *feature weightings* as  
 opposed to absolute rankings of features, we asked  
 users to order the set of NL features by importance.  
 We tested if a linear or exponential weighting of  
 their self-reported rankings were able to accurately  
 predict their preferences on the test cases. We com-  
 pared the user’s absolute ranking to (1) the absolute  
 rankings of OPEN by averaging OPEN’s rankings  
 across all personas at the last interaction step and  
 (2) the performance of OPEN at the last interaction  
 step with the learned feature weights.

Our findings, in Figure 4, indicate that the re-  
 lative weighting of features is critical to OPEN’s  
 prediction, and that user’s self-reported rankings  
 are not a good indicator of their actual prefer-  
 ences. This suggests that methods that rely solely  
 on human-written specifications, such as prompts,

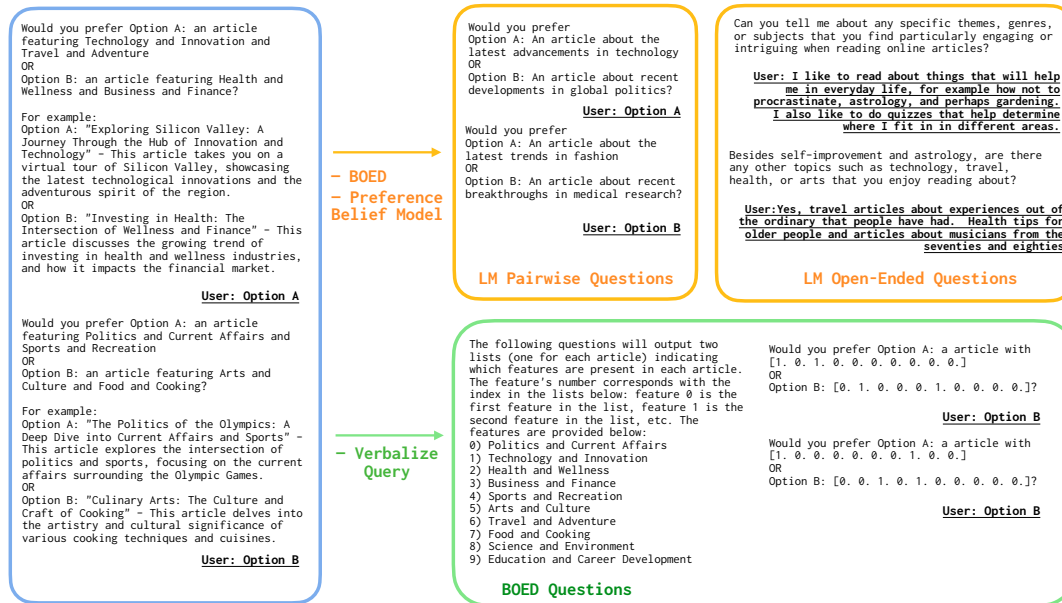


Figure 3: Sample transcripts from OPEN vs. Baselines

482 may not be as useful as methods that are able to  
 483 keep track of the precise relative weights between  
 484 features, as feature weights are typically much  
 485 harder for users to specify than feature rankings.

486 **Qualitative Analysis** At the end of the study,  
 487 we collected open-ended feedback from the users  
 488 regarding their experience interacting with the sys-  
 489 tem. Across the different elicitation methods, we  
 490 found:

491 1. **LM’s open-ended questions to be repetitive**  
 492 **and overreliant on the LM’s prior over user**  
 493 **preferences.** For example one user noted:  
 494 *"the chat bot did not seem to pick up on*  
 495 *many parts of my input, as though it had its*  
 496 *own 'agenda'.*" Another user observed that  
 497 the model repeated questions: *"the chatbot*  
 498 *continued to ask the same exact question in*  
 499 *the same way, it should have different varia-*  
 500 *tions and questions to dive deeper rather than*  
 501 *repeating itself."* The transcripts from these  
 502 users corroborate their observations.

503 2. **LM’s pairwise questions (without an ex-**  
 504 **PLICIT feature set) failed to probe for feature**  
 505 **weighting.** A user commented that: *"I’m not*  
 506 *sure what the purpose here is, but if it was to*  
 507 *pinpoint my interests than they should have*  
 508 *been comparing my previous choices together*  
 509 *as well!"*

510 3. **Users struggled to self-map features to NL**  
 511 **questions.** One user stated: *"I was thoroughly*  
 512 *confused, however I think I figured it out and*  
 513 *clicked accordingly."*

514 We also collected empirical feedback on partici-  
 515 pants’ perceived effort across the different meth-  
 516 ods. We found that they reported consistent levels  
 517 of mental demand across all methods (including the  
 518 *LM Open-ended Questions*), indicating that OPEN  
 519 questions were not more challenging for users than  
 520 other settings. Further analysis and visualization of  
 521 users’ empirical feedback is included in Appendix  
 522 D. Transcripts from each our settings can be found  
 523 in Figure 3.

524 **7 Related Work**

525 **Task Ambiguity and Underspecification** Prior  
 526 research has studied ambiguity in the context of  
 527 language-guided tasks (Lake et al., 2019; Tamkin  
 528 et al., 2023). In particular, in the context of prompt  
 529 engineering (Brown et al., 2020), writing complete  
 530 and unambiguous task specifications can be dif-  
 531 ficult (Li et al., 2023; 278, 1984). Furthermore,  
 532 recent developments have encouraged fine-tuning  
 533 approaches such as reinforcement learning with  
 534 human feedback (RLHF) (Ziegler et al., 2019;  
 535 Christiano et al., 2017) and direct preference op-  
 536 timization (DPO) (Rafailov et al., 2023) to better  
 537 align LLM behaviors with humans’ (underspeci-

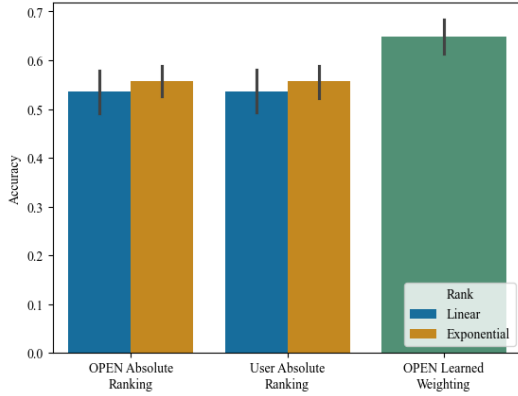


Figure 4: Accuracy for OPEN’s absolute feature rankings, users’ self reported absolute feature rankings, and OPEN’s learned weights over features as described in Section 6. Our analysis indicates removing access to precise weights significantly hurts performance. Error bars are one standard error.

538 fied) preferences. These techniques, however, all  
 539 assume access to existing human preference data.  
 540 Our goal in this work is to examine how to elicit  
 541 such data efficiently.

### 542 Classical Preference Learning Techniques

543 Learning human preferences is a rich area of work  
 544 that has spanned many fields over past few decades.  
 545 Much work has gone into *optimally* querying for  
 546 user preferences. Aside from Bayesian Optimal Ex-  
 547 perimental Design, which we use in this paper, ap-  
 548 proaches have included conjoint analysis (STERN,  
 549 1990; Arora and Huber, 2001; Kuhfeld et al., 1994),  
 550 polyhedral methods (Toubia et al., 2007), multi-  
 551 armed bandits (Lu et al., 2010), and dueling band-  
 552 its (Heckel et al., 2019), etc. (Vayanos et al., 2021)

553 Active learning is an area of machine learning  
 554 focused on interactively querying an expert to label  
 555 new data points, often selected from a pool (Settles,  
 556 2010). Uncertainty-based active learning computes  
 557 the next example to show the user based on explicit  
 558 notions of uncertainty (Lewis and Catlett, 1994).

559 However, these approaches are often only appli-  
 560 cable to simple domains with known feature spaces  
 561 (or existing pools of examples), and may not be  
 562 straightforward to derive for complex domains.

563 **LM Preference Learning** More recently, Li et al.  
 564 (2023); Piriyaikulij et al. (2023) introduced ap-  
 565 proaches to active elicit user preferences with lan-  
 566 guage models. They demonstrated LMs are able  
 567 to actively elicit user preferences across a wide  
 568 variety of domains beyond certain classical tech-

569 niques like active learning. However, our work  
 570 demonstrates that LMs are still subpar at tracking  
 571 and using feature weightings to ask informative  
 572 questions.

## 573 8 Discussion & Future Work

574 We presented OPEN, a domain-agnostic frame-  
 575 work for combining Bayesian Optimal Experimen-  
 576 tal Design with LMs—using BOED to guide the  
 577 choice of informative questions, and the LM to ex-  
 578 tract environment-relevant features and translate  
 579 abstract feature queries into real NL questions. Be-  
 580 yond the improvements demonstrated by our find-  
 581 ings, our system has two additional advantages over  
 582 LM-only methods: (1) *we can improve the trans-  
 583 parency of LMs* through extracting NL features.  
 584 In our framework the NL features from the LM  
 585 accompanied by an external uncertainty-driven  
 586 model conditioned on the features, provide greater  
 587 transparency than an LM. With OPEN, we can  
 588 understand how important each feature is to the  
 589 elicitation and prediction process. (2) *We can im-  
 590 prove privacy while reducing computational costs.*  
 591 By abstracting prediction to an external, linear sys-  
 592 tem, we can model user behavior without the need  
 593 for LM intervention. In OPEN, the Bayesian Partic-  
 594 le Filter models each user’s behavior separately  
 595 while still incorporating the high-level, contextual  
 596 information encoded in LMs. Furthermore, at test  
 597 time, conducting inference on this linear system  
 598 is orders of magnitude cheaper and yields better  
 599 results than an LM.

600 OPEN’s flexibility allows for the exploration  
 601 of other, more intensive preference-learning ap-  
 602 proaches such as variational Bayesian methods  
 603 (e.g. Foster et al., 2019) or multi-armed bandits  
 604 (e.g. Lindner et al., 2022). We are also excited by  
 605 the potential to further incorporate aspects of real-  
 606 world learning environments. For example, future  
 607 work could investigate using OPEN to develop  
 608 an adaptive feature space—as the model’s uncer-  
 609 tainty increases, the model could query for and  
 610 optimize new features that might better explain pat-  
 611 terns within the data. Additionally, recent advance-  
 612 ments in fine-tuning models from human feedback  
 613 (Ouyang et al., 2022; Rafailov et al., 2023; Etha-  
 614 yarajh et al., 2024) found significant differences in  
 615 data curation and model performance when chang-  
 616 ing the experimental design. OPEN enables future  
 617 research to explore the impact of this parameter in  
 618 an active preference elicitation environment.

619	<b>Limitations</b>		
620	Although we provide the steps to build towards		
621	better active uncertainty-driven preference learning		
622	methods, our work has several limitations. In our		
623	study, OPEN was constrained to pairwise queries		
624	in a content recommendation setting. Also, we		
625	operated under a fixed feature space and with a		
626	closed-source LM. Future work could explore ex-		
627	panding along each of these axes: investigating		
628	other preference-learning domains, incorporating		
629	open-ended questions, expanding the feature space		
630	in conversation, and testing other LMs. Further-		
631	more, future research should explore the impact		
632	of OPEN across additional real-world settings and		
633	population groups.		
634	<b>Ethical Considerations</b>		
635	Our work presents both ethical benefits and risks.		
636	Understanding user preferences in underspecified		
637	environments is crucial to avoiding real-world		
638	issues with AI systems such as spreading hate		
639	speech, generating illicit or copyrighted content,		
640	as well as perpetuating bias and stereotypes. Fur-		
641	thermore, as more systems are deployed without		
642	expert-supervision, ensuring that they can align		
643	themselves with users' values is crucial to ensuring		
644	safe and healthy interactions. However, in align-		
645	ing with user preferences, it is also possible that		
646	systems may align with unwanted or dangerous ten-		
647	dencies. Developing guidelines for models' value-		
648	systems is crucial to ensuring the long-term success		
649	of human-AI interactions. Furthermore, working		
650	with preference data inherently presents risks as it		
651	can be used maliciously to manipulate individuals.		
652	It is necessary to ensure that preference data, when		
653	used by AI systems, is collected robustly and with		
654	supervision.		
655	<b>User Study</b> We used the Prolific platform ( <a href="#">Palan</a>		
656	<a href="#">and Schitter, 2017</a> ) to conduct all of our human		
657	subject experiments. We paid workers in line with		
658	Prolific guidelines and solicited initial feedback		
659	via pilot studies to make the interactive experience		
660	more enjoyable for participants. Anecdotally, many		
661	participants expressed enjoying the study:		
662	1. <i>"This was a very interesting exercise. Thank</i>		
663	<i>you for the opportunity."</i>		
664	2. <i>"it was really interesting. good job."</i>		
665	3. <i>"The task was straightforward and I liked that</i>		
666	<i>there were diverse topics to choose from. It felt</i>		
		<i>like I had the chance to express my preferences</i>	667
		<i>on various subjects"</i>	668
	<b>Reproducibility Statement</b> We provide thor-		669
	ough experimental details including the user in-		670
	terface and LM prompts in the Appendix. We will		671
	also release our codebase and anonymized user		672
	data on GitHub.		673
	<b>References</b>		674
	1984. <a href="#">Ieee guide for software requirements specifica-</a>		675
	<a href="#">tions. IEEE Std 830-1984</a> , pages 1–26.		676
	Neeraj Arora and Joel Huber. 2001. <a href="#">Improv-</a>		677
	<a href="#">ing Parameter Estimates and Model Prediction</a>		678
	<a href="#">by Aggregate Customization in Choice Experi-</a>		679
	<a href="#">ments. Journal of Consumer Research</a> , 28(2):273–		680
	283. <a href="#">_eprint: https://academic.oup.com/jcr/article-</a>		681
	<a href="#">pdf/28/2/273/17927222/28-2-273.pdf</a> .		682
	Ralph Allan Bradley and Milton E. Terry. 1952. <a href="#">Rank</a>		683
	<a href="#">analysis of incomplete block designs: I. the method</a>		684
	<a href="#">of paired comparisons. Biometrika</a> , 39(3/4):324–		685
	345.		686
	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie		687
	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind		688
	Neelakantan, Pranav Shyam, Girish Sastry, Amanda		689
	Askell, Sandhini Agarwal, Ariel Herbert-Voss,		690
	Gretchen Krueger, Tom Henighan, Rewon Child,		691
	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,		692
	Clemens Winter, Christopher Hesse, Mark Chen, Eric		693
	Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,		694
	Jack Clark, Christopher Berner, Sam McCandlish,		695
	Alec Radford, Ilya Sutskever, and Dario Amodei.		696
	2020. <a href="#">Language models are few-shot learners</a> .		697
	Juan Carlos Candeal-Haro and Esteban Induráin-Eraso.		698
	1995. <a href="#">A note on linear utility. Economic Theory</a> ,		699
	6(3):519–522.		700
	Daniel Cavagnaro, Jay Myung, Mark Pitt, and Janne Ku-		701
	jala. 2009. <a href="#">Adaptive design optimization: A mutual</a>		702
	<a href="#">information-based approach to model discrimination</a>		703
	<a href="#">in cognitive science. Neural computation</a> , 22:887–		704
	905.		705
	Paul F Christiano, Jan Leike, Tom Brown, Miljan Mar-		706
	tic, Shane Legg, and Dario Amodei. 2017. <a href="#">Deep</a>		707
	<a href="#">Reinforcement Learning from Human Preferences</a> .		708
	In <a href="#">Advances in Neural Information Processing Sys-</a>		709
	<a href="#">tems</a> , volume 30. Curran Associates, Inc.		710
	Mahasen B. Dehideniya, Christopher C. Drovandi, and		711
	James M. McGree. 2018. <a href="#">Optimal bayesian design</a>		712
	<a href="#">for discriminating between models with intractable</a>		713
	<a href="#">likelihoods in epidemiology. Computational Statis-</a>		714
	<a href="#">tics Data Analysis</a> , 124:277–297.		715
	A. Doucet and A. M. Johansen. 2008. <a href="#">A tutorial on</a>		716
	<a href="#">particle filtering and smoothing: Fifteen years later</a> .		717

718	Sergey Dushenko, Kapildeb Ambal, and Robert D. McMichael. 2020. <a href="#">Sequential bayesian experiment design for optically detected magnetic resonance of nitrogen-vacancy centers</a> . <i>Phys. Rev. Appl.</i> , 14:054036.	D. V. Lindley. 1956. <a href="#">On a measure of the information provided by an experiment</a> . <i>The Annals of Mathematical Statistics</i> , 27(4):986–1005.	771 772 773
723	Jos Elfring, Elena Torta, and René van de Molengraft. 2021. <a href="#">Particle filters: A hands-on tutorial</a> . <i>Sensors</i> , 21(2).	D. V. Lindley. 1972. <i>Bayesian Statistics</i> . Society for Industrial and Applied Mathematics.	774 775
726	Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. 2024. <a href="#">Kto: Model alignment as prospect theoretic optimization</a> .	David Lindner, Sebastian Tschiatschek, Katja Hofmann, and Andreas Krause. 2022. <a href="#">Interactively learning preference constraints in linear bandits</a> .	776 777 778
729	Adam Foster, Martin Jankowiak, Eli Bingham, Paul Horsfall, Yee Whye Teh, Tom Rainforth, and Noah Goodman. 2019. <a href="#">Variational bayesian optimal experimental design</a> .	Tyler Lu, David Pal, and Martin Pal. 2010. <a href="#">Contextual multi-armed bandits</a> . In <i>Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics</i> , volume 9 of <i>Proceedings of Machine Learning Research</i> , pages 485–492, Chia Laguna Resort, Sardinia, Italy. PMLR.	779 780 781 782 783 784
733	Johannes Fürnkranz and Eyke Hüllermeier. 2012. <i>Preference Learning</i> , pages 2669–2672. Springer US, Boston, MA.	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. <a href="#">Training language models to follow instructions with human feedback</a> .	785 786 787 788 789 790 791 792
736	Yarin Gal, Riashat Islam, and Zoubin Ghahramani. 2017. <a href="#">Deep bayesian active learning with image data</a> .	Stefan Palan and Christian Schitter. 2017. Prolific.ac—a subject pool for online experiments. <i>Journal of Behavioral and Experimental Finance</i> , 17:22–27.	793 794 795
738	N.J. Gordon, D.J. Salmond, and A.F.M. Smith. 1993. <a href="#">Novel approach to nonlinear/non-gaussian bayesian state estimation</a> . <i>IEE Proc. F Radar Signal Process. UK</i> , 140(2):107.	Top Piriyaakulij, Volodymyr Kuleshov, and Kevin Ellis. 2023. <a href="#">Active preference inference using language models and probabilistic reasoning</a> .	796 797 798
742	Reinhard Heckel, Nihar B. Shah, Kannan Ramchandran, and Martin J. Wainwright. 2019. <a href="#">Active ranking from pairwise comparisons and when parametric assumptions do not help</a> . <i>The Annals of Statistics</i> , 47(6):3099 – 3126.	Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. <a href="#">Direct preference optimization: Your language model is secretly a reward model</a> .	799 800 801 802
747	Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. 2011. <a href="#">Bayesian active learning for classification and preference learning</a> .	Tom Rainforth, Adam Foster, Desi R Ivanova, and Freddie Bickford Smith. 2023. <a href="#">Modern bayesian experimental design</a> .	803 804 805
750	Warren F. Kuhfeld, Randall D. Tobias, and Mark Garratt. 1994. <a href="#">Efficient experimental design with marketing research applications</a> . <i>Journal of Marketing Research</i> , 31(4):545–557.	Dorsa Sadigh, Anca D. Dragan, S. Shankar Sastry, and Sanjit A. Seshia. 2017. <a href="#">Active preference-based learning of reward functions</a> . In <i>Robotics: Science and Systems</i> .	806 807 808 809
754	Brenden M. Lake, Tal Linzen, and Marco Baroni. 2019. <a href="#">Human few-shot learning of compositional instructions</a> . In <i>Annual Meeting of the Cognitive Science Society</i> .	Burr Settles. 2010. <a href="#">Active learning literature survey</a> .	810
758	David D. Lewis and Jason Catlett. 1994. <a href="#">Heterogeneous Uncertainty Sampling for Supervised Learning</a> . In William W. Cohen and Haym Hirsh, editors, <i>Machine Learning Proceedings 1994</i> , pages 148–156. Morgan Kaufmann, San Francisco (CA).	HAL STERN. 1990. <a href="#">A continuum of paired comparisons models</a> . <i>Biometrika</i> , 77(2):265–273. <a href="https://academic.oup.com/biomet/article-pdf/77/2/265/5618454/77-2-265.pdf">_eprint: https://academic.oup.com/biomet/article-pdf/77/2/265/5618454/77-2-265.pdf</a> .	811 812 813 814
763	Belinda Z. Li, Alex Tamkin, Noah Goodman, and Jacob Andreas. 2023. <a href="#">Eliciting human preferences with language models</a> .	Alex Tamkin, Kunal Handa, Avash Shrestha, and Noah Goodman. 2023. <a href="#">Task ambiguity in humans and language models</a> . In <i>The Eleventh International Conference on Learning Representations (ICLR)</i> .	815 816 817 818
766	Jessy Lin, Daniel Fried, Dan Klein, and Anca Dragan. 2022. <a href="#">Inferring rewards from language in context</a> .	Olivier Toubia, John Hauser, and Rosanna Garcia. 2007. <a href="#">Probabilistic polyhedral methods for adaptive choice-based conjoint analysis: Theory and application</a> . <i>Marketing Science</i> , 26(5):596–610.	819 820 821 822
768	Jessy Lin, Nicholas Tomlin, Jacob Andreas, and Jason Eisner. 2023. <a href="#">Decision-oriented dialogue for human-ai collaboration</a> . <i>arXiv preprint arXiv:2305.20076</i> .		



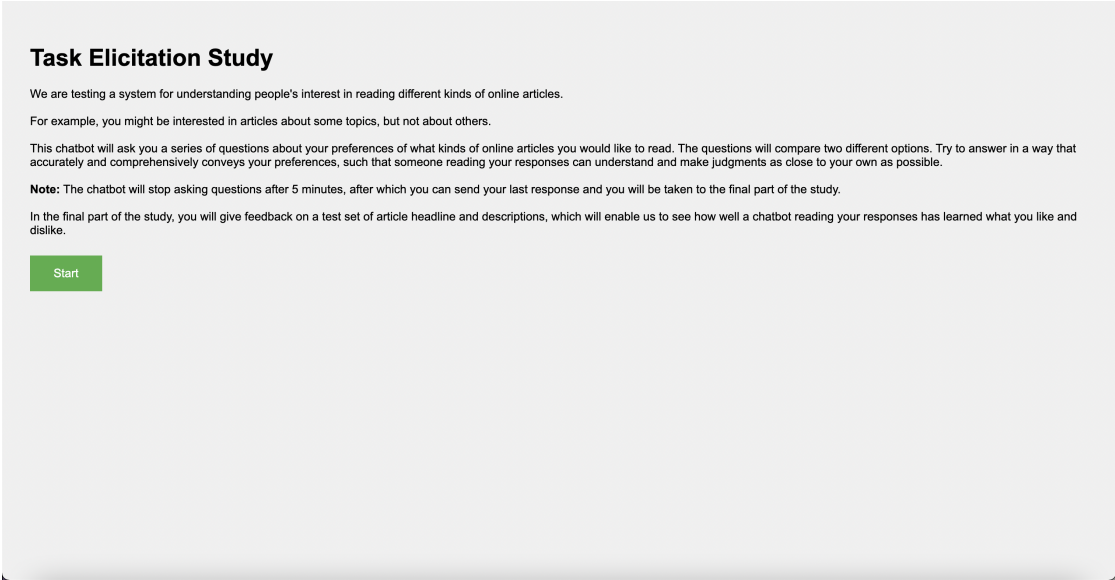


Figure 5: Start screen

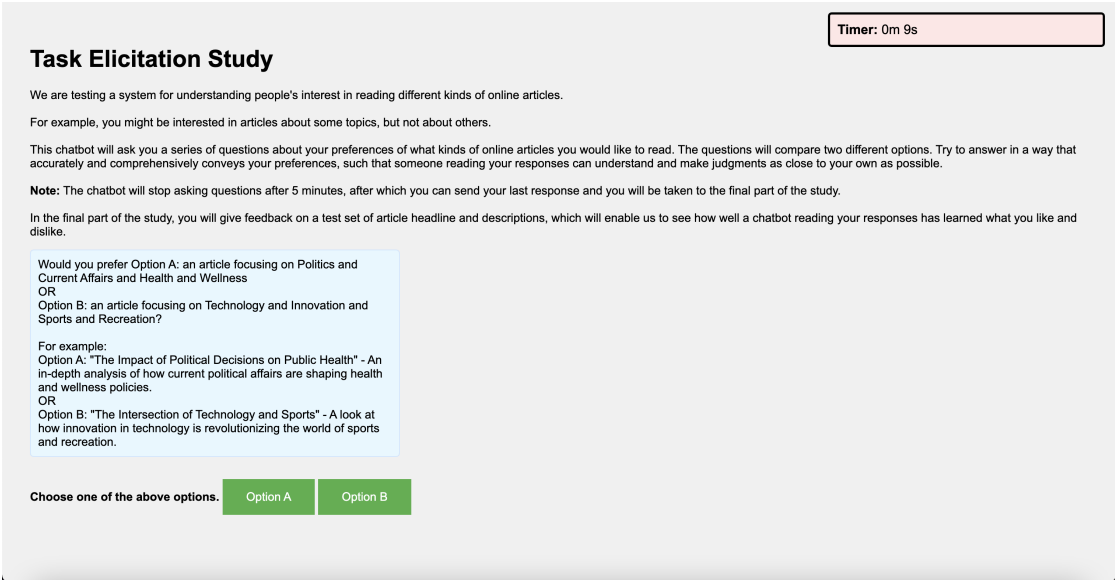


Figure 6: Question-answering process

### Task Elicitation Study

Please answer the following questions about the conversation you just had with the chatbot.

How mentally demanding was interacting with the chatbot?

Very Low Very High

1 2 3 4 5 6 7

How comprehensively do you feel characterized your preferences about the task?

Very Poorly Very Well

1 2 3 4 5 6 7

Figure 7: Collecting user feedback metrics

### Task Elicitation Study

Please indicate which of the two following articles you would prefer to read. Optionally, you may provide an explanation for your decision for each example.

"In Norway, the Electric Vehicle Future Has Already Arrived" - About 80 percent of new cars sold in Norway are battery-powered. As a result, the air is cleaner, the streets are quieter and the grid hasn't collapsed. But problems with unreliable chargers persist. (Option A)

"Two Coronations, 70 Years Apart" - The Visual Parallels Between Elizabeth II and Charles III's Coronations (Option B)

"Why China Is Tightening Its Oversight of Banking and Tech" - A series of regulatory changes approved this week reflect the increasingly centralized control of Xi Jinping, newly confirmed for a third term as China's president. (Option A)

"For a Seasonal Taste of Milan: Candied Chestnuts" - Many residents favor the marrons glacés prepared by Giovanni Galli 1911, one of the best-known confectioners in the city. (Option B)

**Your Response(s)**  
Your submitted response(s) are provided for reference, but please make decisions based on your present intuition, not strictly based on these responses.

1. Would you prefer Option A: an article focusing on Politics and Current Affairs and Health and Wellness  
OR  
Option B: an article focusing on Technology and Innovation and Sports and Recreation?

For example:  
Option A: "The Impact of Political Decisions on Public Health" - An in-depth analysis of how current political affairs are shaping health and wellness policies.  
OR  
Option B: "The Intersection of Technology and Sports" - A look at how innovation in technology is revolutionizing the world of sports and recreation.  
[Option A](#)

Figure 8: Evaluation—user answering pairwise comparison test cases

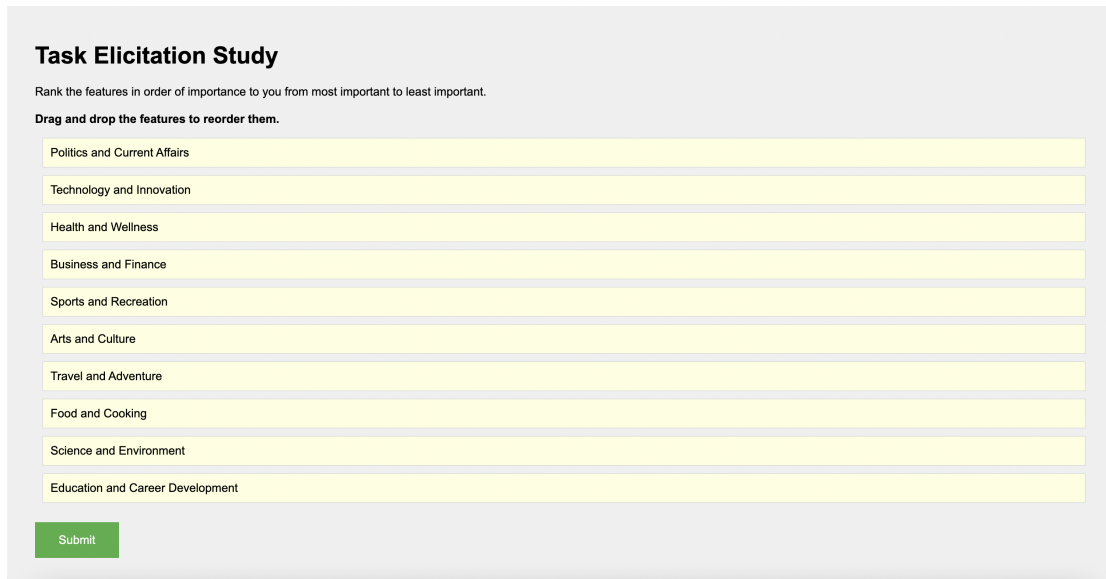


Figure 9: User reorders NL features from most to least important

907	1. feedback_challenge: <i>"How mentally demanding was interacting with the chatbot / writing your answer?"</i>	axes. Unsurprisingly, in the <i>User-Self Mapping</i>	935
908		case, users went back to look at their conversation	936
909		history much less frequently than other in other	937
910	2. feedback_new_issues_interaction: <i>"To what extent did the chatbot raise issues or aspects about your preferences that you hadn't previously considered?"</i>	methods. Interestingly, users felt that LM Pairwise	938
911		questions covered their preferences well pre-test,	939
912		but this ranking dropped significantly post-test (and	940
913		did not translate to a higher accuracy on the test	941
914	3. feedback_interaction_coverage_pretest: <i>"How comprehensively do you feel the chatbot's questions / your answer characterized your preferences about the task?"</i>	set).	942
915			
916		<b>F LM Prompts</b>	943
917		Below, we include the prompts used for the differ-	944
918	4. feedback_interaction_coverage_posttest: <i>"After seeing the examples in the *second* part of the task, how well do you feel the chatbot / the answer you wrote (in the first part of the task) covered the important issues or aspects of these examples?"</i>	ent LM-based tasks.	945
919			
920		<b>F.1 LM Open-ended Questions</b>	946
921	5. feedback_testcase_use_history: <i>"When performing the *second* part of the task, to what extent did you refer back to your conversation history / answer from the first part of the task?"</i>	Your task is to learn what topics a user is	947
922		interested in reading online article	948
923		about. People's interests are broad, so	949
924	6. feedback_lm_experience: <i>"How much experience have you had (if any) with interacting with language models (e.g. ChatGPT, GPT4, etc.)?"</i>	you should seek to understand their	950
925		interests across many topics; in other	951
926		words, go for breadth rather than depth.	952
927		Do not assume a user has given a	953
928		complete answer to any question, so make	954
929		sure to keep probing different types of	955
930		interests.\n\nPrevious questions: {	956
931		interaction_history_formatted}.\n\n	957
932		nGenerate the most informative open-	958
933	As we can see, OPEN was not more challeng-	ended question that, when answered, will	959
934	ing than the other methods across the different	reveal the most about the desired	960
		behavior beyond what has already been	961
		queried for above. Make sure your	962
		question addresses different aspects of	963
		their preferences than the questions	964
		that have already been asked. At the	965
		same time however, the question should	966
		be bite-sized, and not ask for too much	967
		at once. Phrase your question in a way	968
		that is understandable to non-expert	969

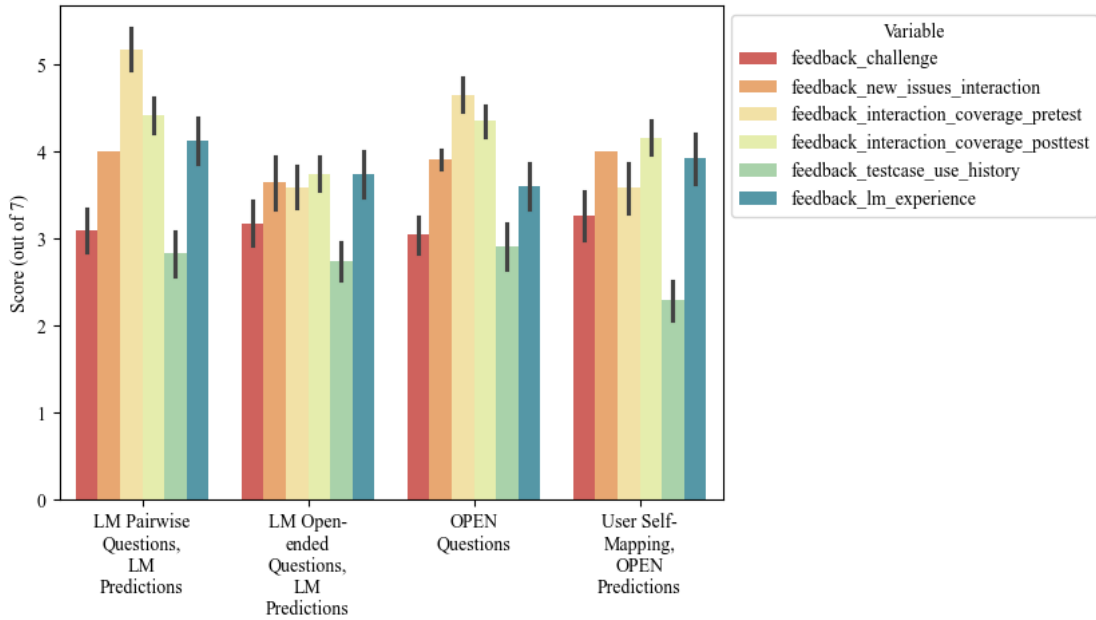


Figure 10

970 humans; do not use any jargon without explanation. Generate the question and  
 971 nothing else. Provide your output in the format: Question: <question> 1003

972 1004  
 973 1005  
 1006  
 1007  
 1008  
 1009  
 1010  
 1011  
 1012  
 1013  
 1014  
 1015  
 1016  
 1017  
 1018  
 1019  
 1020  
 1021  
 1022  
 1023  
 1024  
 1025  
 1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042

974 **F.2 LM Pairwise Questions, LM Predictions**

975 Your task is to learn what topics a user is interested in reading online article about.  
 976 People’s interests are broad, so you should seek to understand their interests across  
 977 many topics; in other words, go for breadth rather than depth. Do not assume a user has  
 978 given a complete answer to any question, so make sure to keep probing different types of  
 979 interests. Previous questions: { interaction\_history\_formatted}. Generate  
 980 the most informative pairwise comparison question that, when answered, will reveal  
 981 the most about the desired behavior beyond what has already been queried for above.  
 982 Make sure your question addresses different aspects of their preferences than the  
 983 questions that have already been asked. At the same time however, the question should  
 984 be bite-sized, and not ask for too much at once. Phrase your question in a way that is  
 985 understandable to non-expert humans; do not use any jargon without explanation. Generate  
 986 the pairwise comparison question and nothing else. Provide your output in the  
 987 format: Would you prefer Option A: <first article option> OR Option B: <second  
 988 article option>?  
 989  
 990  
 991  
 992  
 993  
 994  
 995  
 996  
 997  
 998  
 999  
 1000  
 1001

1002 **F.3 LM Pairwise Question, OPEN Predictions**

1043 one sentence description] conversation below:\n{preferences}\n 1104  
 nThe question is: Based on these 1105  
 preferences, which of the following two 1106  
 1044 **F.4 Extracting Environment Features** articles would the user prefer?\n 1107  
 nOption A: {test\_case\_0}\nOR\nOption B: 1108  
 {test\_case\_1} 1109

1045 Your task is to learn what topics a user is  
 1046 interested in reading online articles  
 1047 about. Enumerate 10 binary topics (  
 1048 features) that may impact user's  
 1049 decisions when choosing which article to  
 1050 read. People's interests are broad, so  
 1051 you should seek to understand their  
 1052 interests across many topics; in other  
 1053 words, go for breadth rather than depth  
 1054 .\n\nOnly include the most important  
 1055 features in your list. Do not include  
 1056 features for which the user's preference  
 1057 would be obvious. Order the features  
 1058 from the most to least likely of  
 1059 interest.\n\nYour output should be in  
 1060 the following format:\n1) <first feature  
 1061 >\n2) <second feature>

1062 **F.5 Mapping BOED Pairwise Comparison to**  
 1063 **NL**

1064 Create two specific, real-world examples of  
 1065 news articles someone might be  
 1066 interested in reading based on the  
 1067 following question which juxtaposes two  
 1068 different news articles:\n\n"Would you  
 1069 prefer Option A: an article with {  
 1070 pairwise\_comparison\_0}\nOR\nOption B: {  
 1071 pairwise\_comparison\_1}?".\n\nThis  
 1072 question instantiates two articles based  
 1073 on their feature values. Features lie  
 1074 on a spectrum ranging from 0.0 to 1.0,  
 1075 where 0.0 corresponds to the absence of  
 1076 that feature and 1.0 indicates an  
 1077 extremely high presence of it.\n\nMake  
 1078 sure to maintain the relative difference  
 1079 between the two articles when  
 1080 generating the descriptions.\n\nProvide  
 1081 your output in the format:\nOption A: [  
 1082 Option A title and simple, one sentence  
 1083 description]\nOR\nOption B: [Option B  
 1084 title and simple, one sentence  
 1085 description]

1086 **F.6 LM Evaluation**

1087 Provide your best guess and the probability  
 1088 that it is correct (0.0 to 1.0) for the  
 1089 following question. Give ONLY the guess  
 1090 and probability, no other words or  
 1091 explanation. If you are unsure take your  
 1092 best guess (between Option A and Option  
 1093 B). For example:\n\nGuess: <most likely  
 1094 guess--either Option A or Option B--as  
 1095 short as possible; not a complete  
 1096 sentence!\n\nProbability: <the  
 1097 probability between 0.0 and 1.0 that  
 1098 your guess is correct, without any extra  
 1099 commentary whatsoever; just the  
 1100 probability!>.\n\nA user has a  
 1101 particular set of preferences over what  
 1102 articles they would like to read. They  
 1103 have specified their preferences in a