

Beyond Demographics: Aligning Role-playing LLM-based Agents Using Human Belief Networks

Yun-Shiuan Chuang Krirk Nirunwiroj[†] Zach Studdiford[†] Agam Goyal
Vincent V. Frigo Sijia Yang Dhavan Shah Junjie Hu Timothy T. Rogers
University of Wisconsin-Madison

{yunshiuan.chuang, nirunwiroj, studdiford, agoyal25}@wisc.edu
{vfrigo, syang84, dshah, junjie.hu, ttrogers}@wisc.edu

Abstract

Creating human-like large language model (LLM) agents is crucial for faithful social simulation. Having LLMs role-play based on demographic information sometimes improves human likeness but often does not. This study assessed whether LLM alignment with human behavior can be improved by integrating information from empirically-derived human belief networks. Using data from a human survey, we estimated a belief network encompassing 64 topics loading on nine non-overlapping latent factors. We then seeded LLM-based agents with an opinion on one topic, and assessed the alignment of its expressed opinions on remaining test topics with corresponding human data. Role-playing based on demographic information alone did not align LLM and human opinions, but seeding the agent with a single belief greatly improved alignment for topics related in the belief network, and not for topics outside the network. These results suggest a novel path for human-LLM belief alignment in work seeking to simulate and understand patterns of belief distributions in society.

1 Introduction

With rapid advances in large language models (LLMs), there has grown increasing interest in using LLMs to simulate and understand dynamics of human communication and persuasion (Park et al., 2023, 2022; Chuang et al., 2024a; Taubenfeld et al., 2024). Current LLMs can be prompted to role-play as individuals with particular demographic traits, sometimes then producing patterns of behavior that seem remarkably human-like. For instance, when asked to report the US unemployment rate when President Obama left office, ChatGPT will provide the exact answer; but if first instructed to role-play as a typical Democrat or Republican and asked the same question, the model produces incorrect, inflated estimates that mirror patterns of partisan bias

[†]Joint second authors.

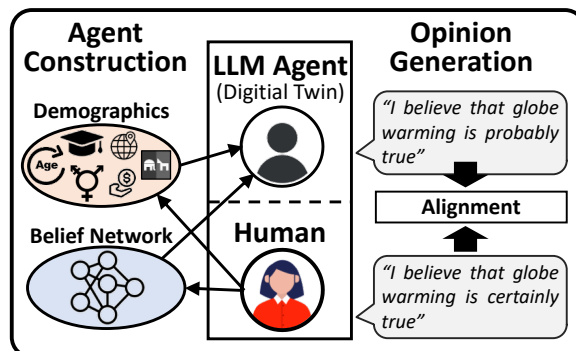


Figure 1: An LLM agent i' is constructed as the “digital twin” of a human respondent i , based on their demographic information and belief network estimated from a belief survey. We then evaluate the alignment between the opinions generated by the agent ($o_{i'}$) and those expressed by the corresponding human respondent (o_i).

in analogous human studies (Chuang et al., 2024b). Such results raise the possibility that, with strategic prompting, LLMs may serve as useful proxies for capturing the beliefs and attitudes of various socio-demographic groups.

Other recent work suggests, however, that the alignment between beliefs expressed by role-playing LLMs and matched human participants is unreliable at best. For instance, Santurkar et al. (2023) found that LLMs tuned via human feedback generally reflect opinions from liberal and well-educated demographics and that having LLMs role-play as humans with different socio-demographic traits does not remediate this tendency. Similarly, Sun et al. (2024) had LLMs offer opinions on controversial issues while role-playing as humans with varying demographic characteristics, and found that the model only reflected corresponding human opinions on one of the ten total topics. Chuang et al. (2024a) additionally found that, even when seeded with prompts specifying an initial belief that runs contrary to social consensus (e.g., “global warming is a hoax”), LLMs quickly revert to the

accepted ground-truth attitude after repeated interactions with other agents. Overall, this work suggests that LLMs fine-tuned with human feedback tend to adopt a consistent stance regardless of the demographic background they role-play—a behavior that may aid LLM fairness and value alignment, but limits their utility as models of human communicative dynamics.

This paper considers an alternative approach to aligning the attitudes expressed by role-playing LLMs and the human groups they are intended to emulate. The central idea relies on behavioral studies of human *belief networks*: the empirical observation that beliefs on different topics are not distributed at random across the population, but tend to cohere together in patterns of high-order covariation (Boutyline and Vaisey, 2017; Vlasceanu et al., 2024; Keating, 2023; Turner-Zwinkels and Brandt, 2022). For instance, people who believe that government should support social welfare programs are also more likely to believe in higher taxes on the wealthy, strong union protections, and universal health care. Thus, knowing a person’s opinion on one topic can carry rich information about their likely views on many others. Because LLMs learn from vast amounts of human-generated language data, the weights they acquire and hence patterns of behaviors they exhibit may implicitly capture the tendency for various beliefs to co-occur in human populations, providing novel leverage for alignment. Specifically, human-LLM alignment may be guided, not just by socio-demographic role-playing, but also by instructing LLMs to hold a specific opinion on a representative topic.

To test this idea, we considered a simple belief network constructed in prior work by applying factor analysis to a dataset measuring human beliefs across a diverse array of topics (Frigo, 2022). Factor analysis decomposes patterns of covariation among expressed beliefs, identifying relationships between the beliefs themselves and a set of underlying latent factors. From this analysis we identified nine *orthogonal* factors, each receiving high factor loadings (having strong associations) from several controversial beliefs. Each latent factor is associated with a distinct set of beliefs, with very little overlap between the beliefs linked to different factors. In other words, beliefs form distinct clusters with clear separation between clusters. Two example factors included a *ghost factor* grouping beliefs in various supernatural phenomena (e.g., talking to the dead) and a *partisan factor* grouping beliefs that

are typically politically polarizing in the US (e.g., effectiveness of gun control). We then considered how well the opinions of contemporary LLMs align with human participants when they role-play (a) without demographic information, (b) with demographic information only, or (c) with demographic information plus a belief on a specific topic that strongly aligns with either the same latent factor or a different latent factor. When seeding each model with such a belief, we additionally compared the effects of in-context learning (i.e. prompting) versus supervised fine-tuning. The results suggest that attention to empirically-derived human belief networks provides a useful strategy for human-LLM alignment, moreso than demographic role-playing.

2 Preliminaries: LLM Agents as Human Digital Twins

As depicted in Figure 1, we aim to construct an LLM agent i' as the i -th human’s “digital twin”, such that their opinions o on various topics x are aligned. We first use information about human i (e.g., their demographic information d) to create the corresponding LLM agent i' , and then query the agent’s opinion ($o_{i'}$) on a wide range of topics. We then evaluate the human-LLM alignment by measuring the discrepancy between the actual human opinion o_i and the LLM agent’s opinion $o_{i'}$. Note that we use the term LLM-based “agent” to refer to the digital twin because the instructed LLM is intended to produce behaviors that emulate the human individual they role-play (Park et al., 2023; Shao et al., 2023; Zhou et al., 2023).

3 Methods

3.1 Controversial Beliefs Survey

The specific opinions we assessed were taken from the *Controversial Beliefs Survey* developed in Frigo (2022). The survey measures the direction and strength of belief across 64 topics spanning broad aspects of human knowledge, including history, science, health, religion, the supernatural, economics, politics, and conspiracy theories (see Table 4 in §A for the full list of topics). Topics were selected to elicit diverse opinions about their truthfulness (hence “controversial beliefs”). Each belief was stated as a factual proposition (e.g., “States with stricter gun control laws have fewer gun deaths per capita”), and participants rated their views about the truth of the statement on a six-point Likert scale ranging from “Certainly false” to “Certainly

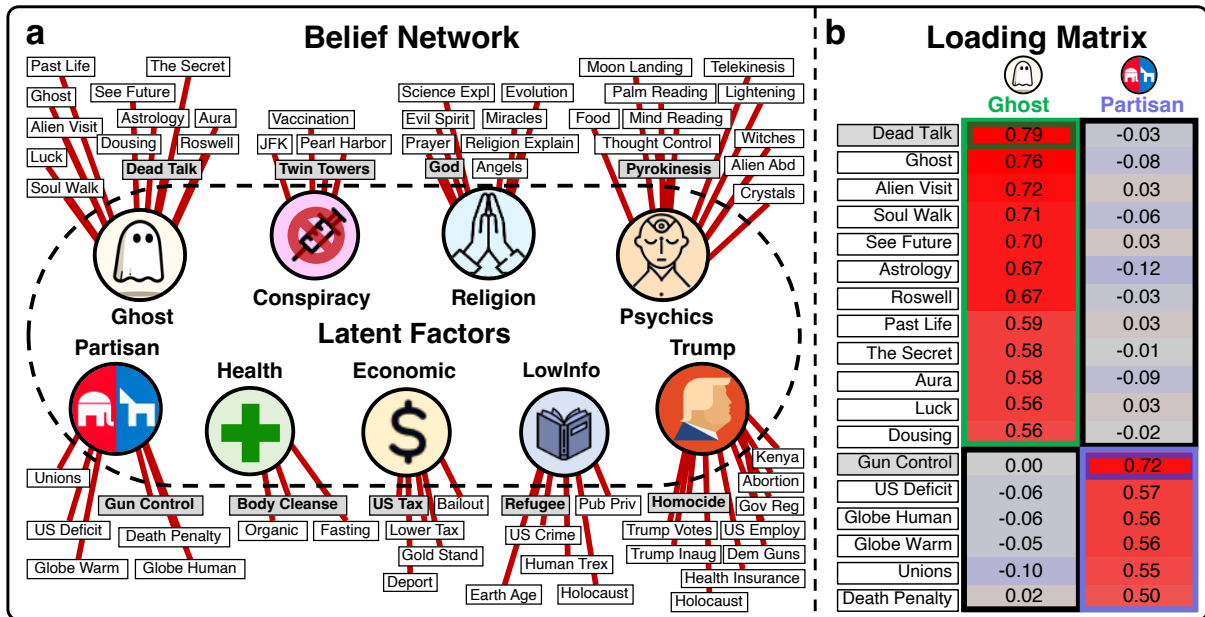


Figure 2: (a) The belief networks estimated by factor analysis from human respondents’ responses on the Controversial Beliefs Survey. The nine central nodes are the orthogonal latent factors, and the leaves (rectangles) are the 64 individual topics x . The training topics x_{train} are highlighted with grey backgrounds. (b) Factor loading matrix between two latent factors and their topics. Figure 5 shows the full factor loading matrix and Table 4 the full statement of the each topic.

true.” Responses with high numbers indicate agreement with the rational/consensus ground truth. The dataset also has extensive demographic data from respondents, including age, gender, education level, household income, urban versus rural living environment, state of residence, and political leaning.

The dataset includes ratings for $N = 564$ individuals living in the US, collected from Amazon Mechanical Turk in 2018.¹ Formally, we denote the set of 64 topics as $\mathcal{X} = \{x_j\}_{j=1}^M$ ($M = 64$). The survey dataset $\mathcal{D} = \{(d_i, x, o_i) | x \in \mathcal{X}\}_{i=1}^N$ consists of the opinion responses from N individuals, where the i -th individual having the demographic information d_i expresses an opinion o_i to the topic x . The respondents provide their opinions ($-3 \leq o_i \leq 3, o_i \neq 0$) for each statement on a 6-point Likert scale with the values -3 : Certainly false, -2 : Probably false, -1 : Lean false, $+1$: Lean true, $+2$: Probably true, $+3$: Certainly true. No neutral value was provided so participants must minimally lean in one direction or the other. The demographic and opinion data together were used to construct and evaluate the LLM agents (§3.3). The survey dataset can be obtained by contacting its authors (Frigo, 2022).

¹<https://mturk.com/>

3.2 Constructing a Belief Network using Factor Analysis

Our objective was to find independent “belief networks”—that is, groups of topics where expressed beliefs covaried across participants within each group but were independent between groups. To this end, we relied on a previous factor analysis (Frigo, 2022) that first computed a matrix of correlations in the ratings produced across participants for each pair of topics, then decomposed the resulting matrix into a set of orthogonal latent factors using principal component analysis (PCA) with Varimax rotation Kaiser (1958). The PCA yielded a factor loading matrix that encodes the loading (i.e., the association) between each topic and each latent factor. Nine latent factors were extracted based on the factor scree plot (Cattell, 1966, see §D), which together accounted for 72% of the variance in the correlation matrix. The belief network surrounding these nine factors are shown in Figure 2. For example, the *ghost* factor receives high loadings from 12 topics, all pertaining to supernatural or otherworldly beliefs; the *partisan* factor receives high loadings from 6 topics on highly polarized political issues. We referred to these topics as either belonging to the *ghost topic category* or *partisan topic category*, respectively. Hence, the

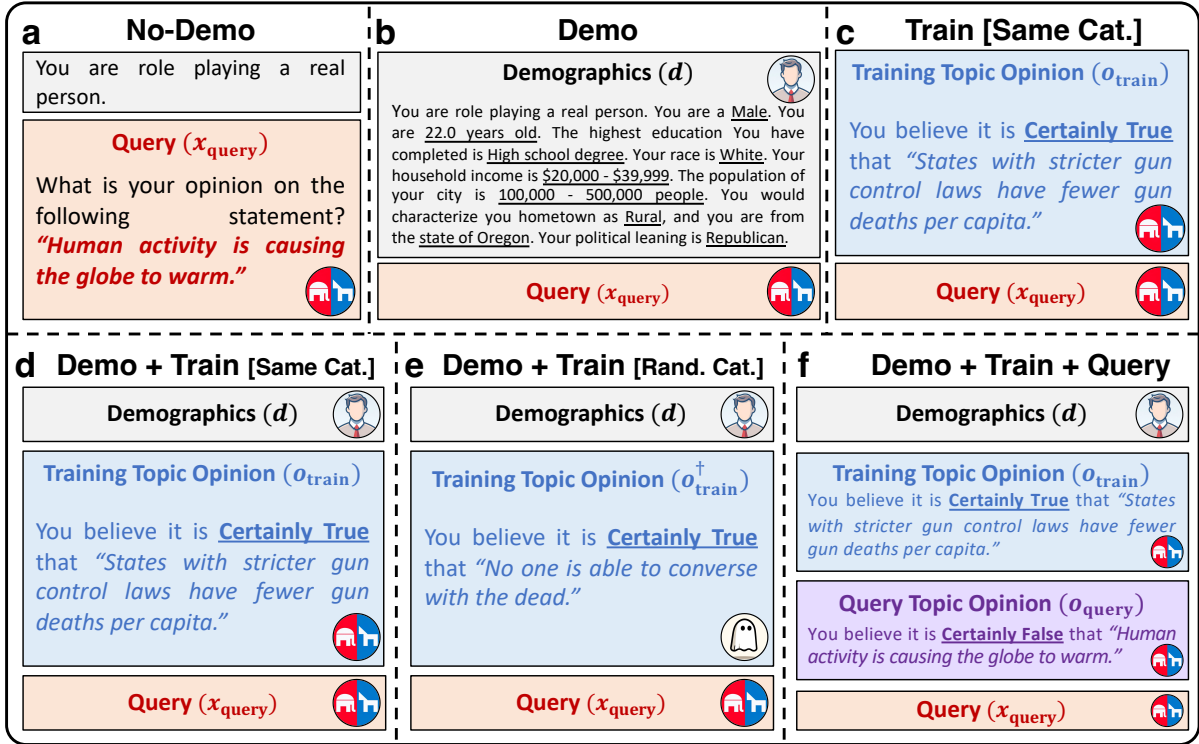


Figure 3: LLM agent construction conditions with different levels of respondent’s information. (a) “No-Demon” baseline condition where the LLM role-plays without demographic information and we directly query the LLM about its opinion on the **query topic** (x_{query}). (b) “Demo” baseline condition with demographic information (d). (c) “Train [same category]” baseline condition **training topic opinion** (o_{train} on x_{train}) from the same topic category as the query topic (in this example, they both belong to the “Partisan” category). (d) “Demo+Train [same category]” condition with demographic information plus **training topic opinion** (o_{train} on x_{train}) from the same topic category as the query topic. (e) “Demo+Train [random category]” baseline condition with demographic information, along with training topic opinion from a randomly selected topic category other than the query topic ($o_{\text{train}}^{\dagger}$ on $x_{\text{train}}^{\dagger}$) (in this example, the training topic is from the “Ghost” category). (f) “Demo+Train+Query” as an upper bound condition with both training topic opinion (from the same category) and the **query topic opinion** (o_{query} on x_{query}).

nine orthogonal latent factors resulted in nine distinct *topic categories*. We took these 64 topics and the corresponding nine latent factors as the targets for our analysis of LLM alignment. The full factor analysis results, including the full factor loading matrix of the nine factors, are reported in §F.

3.3 LLM Agent Construction

For each of the nine topic categories, we designated the topic possessing the highest loading as the model *training topic* (x_{train}). For each digital twin (role-playing LLM agent), the corresponding human opinion on the training topic (o_{train}) was used to customize the LLM agent (either through in-context learning or supervised fine-tuning, see below). Human opinions on the remaining 55 testing topics x_{test} were not provided to the LLM agent; instead, the agent’s expressed opinions o_{test} on these topics were used to evaluate their alignment with the human respondents. We hypothesized that

specifying the agent’s opinion on the training topic might elicit a shared representation that generalizes to testing topics close within the belief network (i.e., sharing the same latent factor), but not those from the other belief network.

For each human respondent i , we constructed an LLM agent i' as their “digital twin,” using a set of strategies described below. For each twin created under a given strategy, we queried the LLM agent for its opinions on the training and test topics (x_{query}), and measured how ratings generated by the digital twins correlate with the true opinions expressed by corresponding human respondents. We then assessed how this measure of human-LLM belief alignment varied with different strategies for constructing the digital twin.

In-context Learning (ICL). As shown in Figure 3, these strategies involve initializing agents via in-context learning (ICL), with different infor-

mation included in their *system message* (see §4.1 and Appendix §B for the prompts).

- a. **Baseline: No-Demo.** An LLM agent is role-playing a generic person without specific information about the human respondent (system message = “You are role playing a real person.”). This provides a performance floor since there is no way for the LLM to align with a corresponding human participant.
- b. **Baseline: Demo.** An LLM agent is constructed to role-play the i -th respondent by adding only the demographic information (d_i) in the prompt.
- c. **Baseline: Train [same category].** An LLM agent is constructed to role-play the i -th respondent by only adding the respondent’s Likert-scale opinion on the training topic ($x_{\text{train}}, o_{\text{train}}$) and is assessed on other topics from the same topic category (x_{query}) within the belief network.
- d. **Demo+Train [same category].** In addition to demographic information, the LLM receives a respondent’s Likert-scale opinion on the training topic ($x_{\text{train}}, o_{\text{train}}$) and is assessed on other topics from the same topic category (x_{query}) within the belief network. This is the critical condition of interest.
- e. **Baseline: Demo+Train [random category].** This baseline condition is similar to Demo+Train [same category], but the training topic opinion ($x_{\text{train}}^\dagger, o_{\text{train}}^\dagger$) belongs to a randomly selected topic category that is different from the query topic. This baseline allows us to determine whether adding respondent’s Likert-scale opinion is only helpful when it belongs to the same belief network as the query topic (x_{query}).
- f. **Upper Bound: Demo+Train+Query.** This condition provides the human opinion rating on both the training topic ($x_{\text{train}}, o_{\text{train}}$) and the query topic ($x_{\text{query}}, o_{\text{query}}$) during the agent construction, providing an upper bound on generalization behavior.

Supervised Fine-tuning (SFT). We also investigated whether seeding initial beliefs via supervised fine-tuning (SFT) can increase human-LLM alignment. Specifically, the correspondence between the demographic information d and the corresponding opinion o (on topic x) was used to fine-tune model weights via supervised learning,

following analogous strategies to the in-context learning approaches described above. For example, for **Demo+Train [same category]**, we first construct the dataset $\mathcal{D}_{\text{SFT}} = \{(d_i, x_{\text{train},i}), o_{\text{train},i}\}_{i=1}^N$ for each topic category. We then fine-tuned the LLM with input context providing the demographic information along with the training topic statement (d, x_{train}), and using the corresponding human Likert-scale response o_{train} as the ground-truth output. After fine-tuning, we assessed the LLM agent’s opinion on query topics x_{query} belonging to the same topic category x_{train} ². Likewise, for **Baseline: Demo+Train [random category]**, it is similar to Demo+Train [same category] condition, but the training topic opinion ($x_{\text{train}}^\dagger, o_{\text{train}}^\dagger$) is from a different topic category as the query topic x_{query} . Details of the fine-tuning procedure and the corresponding prompts are in §C and §E.

4 Experimental Settings

4.1 Configuration for LLM Agents

We evaluated the LLM agents using the following models: ChatGPT (gpt-3.5-turbo-0125; OpenAI, 2022), GPT-4o mini (gpt-4o-mini-2024-07-18), Mistral (Mistral-7B-Instruct-v0.2; Jiang et al., 2023), and LLaMA 3.1 (Llama-3.1-8B-Instruct; Touvron et al., 2023), all with temperature of 0.7. In sensitivity analyses, we consider other temperature values $T \in \{0, 1\}$. During initialization, the demographic background was incorporated into the model’s “*system messages*.” The opinion queries (x_{query}) were fed to the agent through the model’s “*user messages*.” When using in-context learning (§3.3), the training/query topic opinions were also included in the model’s “*system messages*.” The LLM agents were constructed through LangChain (Chase, 2022). For our compute resources, see §G.

4.2 Evaluation Metrics

To evaluate the “human-likeness” of the LLM agents’ opinions, for each topic category in the survey, we computed the mean absolute error (MAE_{test}) between the human opinion (o_i) and that generated by the twinned LLM agent ($o_{i'}$) across the testing topics (x_{test}). Formally, $\text{MAE}_{\text{test}} = \frac{1}{|\mathcal{X}_{\text{test}}|} \sum_{x \sim \mathcal{X}_{\text{test}}} |o_{i,x} - o_{i',x}|$, which is the mean discrepancy between the opinions of human respon-

²For example, we fine-tuned an LLM on the respondents’ opinions on the training topic for the Ghost topic category, then queried its opinion on the test topics in the Ghost topic category.

Model	Condition	Topic Categories									Average
		Ghost	Psychics	Religion	Trump	Partisan	Economic	LowIndo	Health	Conspiracy	
ChatGPT	<i>Baselines</i>										
	No-Demo	2.33	2.26	1.81	1.17	1.43	1.42	1.29	1.62	1.80	1.68
	Demo	2.58	2.28	1.87	1.23	1.41	1.51	1.21	1.66	1.51	1.70
	Train [Same Cat.]	1.48	1.46	1.80	1.18	1.36	1.48	1.23	1.60	1.76	1.48
	Demo + Train [Rand. Cat.]	2.26	1.86	1.93	1.29	1.49	1.63	1.26	1.80	1.53	1.67
	Demo + Train [Same Cat.]	1.26	1.27	1.72	1.14	1.34	1.23	1.15	1.53	1.40	1.34
	<i>Upper Bound</i>										
	Demo + Same Train + Query	0.41	0.48	0.30	0.63	0.28	0.09	0.82	0.30	0.46	0.42
	Relative Gain (%) \uparrow	60.83	56.11	9.55	15.00	6.19	19.72	15.38	9.56	10.48	22.54
	GPT-4o mini	<i>Baselines</i>									
No-Demo		1.49	1.33	1.90	1.21	1.19	1.30	1.31	2.03	1.40	1.46
Demo		1.46	1.21	1.68	1.17	1.19	1.24	1.23	1.41	1.42	1.33
Train [Same Cat.]		1.05	0.96	1.36	1.06	1.18	1.19	1.21	1.42	1.32	1.19
Demo + Train [Rand. Cat.]		1.44	1.23	1.53	1.28	1.24	1.22	1.19	1.58	1.41	1.35
Demo + Train [Same Cat.]		1.00	0.96	1.31	1.06	1.15	1.19	1.16	1.37	1.28	1.16
<i>Upper Bound</i>											
Demo + Same Train + Query		0.04	0.05	0.03	0.64	0.01	0.02	0.14	0.04	0.14	0.12
Relative Gain (%) \uparrow		32.39	21.55	22.42	20.75	3.39	4.10	6.42	2.92	10.94	13.88
Mistral		<i>Baselines</i>									
	No-Demo	1.75	1.63	1.64	1.33	1.20	1.07	1.49	1.30	1.44	1.43
	Demo	1.82	1.93	1.68	1.49	1.27	1.16	1.49	1.39	1.38	1.51
	Train [Same Cat.]	1.46	1.02	1.46	1.46	1.25	1.12	1.44	1.44	1.28	1.33
	Demo + Train [Rand. Cat.]	1.93	1.79	1.60	1.56	1.35	1.22	1.70	1.36	1.45	1.55
	Demo + Train [Same Cat.]	1.36	1.71	1.41	1.05	1.25	1.12	1.12	1.32	1.27	1.29
	<i>Upper Bound</i>										
	Demo + Same Train + Query	0.71	0.39	0.86	0.77	0.59	0.55	0.65	1.04	0.55	0.68
	Relative Gain (%) \uparrow	41.44	14.29	32.93	61.11	2.94	6.56	44.05	20.00	13.25	26.29
	LLaMA 3.1	<i>Baselines</i>									
No-Demo		2.55	2.40	1.88	1.86	2.04	2.54	1.52	1.54	2.11	2.05
Demo		2.36	2.42	1.85	1.50	1.45	2.33	1.47	1.50	2.35	1.91
Train [Same Cat.]		2.21	2.28	1.82	1.44	1.63	1.86	1.48	1.63	2.77	1.90
Demo + Train [Rand. Cat.]		2.70	2.64	2.03	1.69	1.87	2.48	1.80	1.97	2.28	2.16
Demo + Train [Same Cat.]		2.07	1.88	1.81	1.19	1.32	1.69	1.35	1.07	2.00	1.60
<i>Upper Bound</i>											
Demo + Same Train + Query		1.76	1.04	1.42	0.96	0.56	1.47	0.72	0.96	0.65	1.06
Relative Gain (%) \uparrow		48.33	39.13	9.30	57.41	14.61	74.42	16.00	79.63	21.05	39.99

Table 1: Mean absolute error (MAE_{test}) between human respondents and the corresponding LLM agents for each topic category across various LLM agent construction conditions through in-context learning (ICL). The bottom row presents the relative gain (%) as the percentage improvement from the Demo Baseline to the Upper Bound condition for the Demo + Train [Same Cat.] condition. The lower the MAE_{test} and higher the relative gain, the higher the human-LLM alignment. The condition of our main interest (i.e., Demo + Train [Same Cat.] condition) is boldfaced, which also has the best alignment.

dents and LLM agents across all test topics ($\mathcal{X}_{\text{test}}$) within the topic category. The metric MAE_{test} ranges from 0 to 4, where 0 indicates perfect agreement and 4 is the maximum possible disagreement. In addition, because we are interested in the additional value of belief network beyond demographic information, we calculate Relative Gain (%) as the percentage improvement from the Demo Baseline to the Upper Bound condition for the Demo + Train [same category] condition, i.e., Relative Gain (%) = $(MAE_{\text{test}}$ of “Baseline: Demo.” – MAE_{test} of “Demo+Train [same category]”) / $(MAE_{\text{test}}$ of “Baseline: Demo.” – MAE_{test} of “Upper Bound: Demo+Train+Query”) $\times 100$ (%). The Relative Gain is 0% if belief network provides no additional benefit, and 100% if the inclusion of belief network boosts the alignment to the supervised upper bound.

4.3 Supervised Fine-tuning (SFT)

For LLM agents constructed through supervised fine-tuning (§3.3), we used the ChatGPT model `gpt-3.5-turbo-0125`’s fine-tuning API. Critically, because the label (i.e., opinion response o) is usually not balanced in a given topic (e.g., more people believing that ghosts are real than those who don’t), we upsampled the o to ensure equal numbers of responses across the six Likert scale values. Pilot work found that, without upsampling, the fine-tuned LLM agent predominantly produced the most frequent opinion response o_{majority} in \mathcal{D}_{SFT} . Given that the primary aim of the SFT setting is to demonstrate the generalizability of our methods beyond the ICL framework, and recognizing that SFT is inherently more computationally demanding, we concentrate our investigation on two latent factors: the Ghost factor and the Partisan factor§E lists the

hyperparameters for fine-tuning.

5 Results

Demographic information alone does not align the LLM agent’s opinion. As shown in Table 1, incorporating solely the demographic information (the Demo condition) fails to align LLM agents with human respondents. The MAE_{test} of the Demo condition is similar to the No-Demo condition, indicating that the demographic information alone does not help LLM agents align with the human respondents they role-play.

Specifying the agent’s opinion on a training topic aligns other beliefs in the same network.

When the LLM is instructed to adopt the twinned human’s opinion on the training topic ($x_{\text{train}}, o_{\text{train}}$), its expressed opinions on other topics in the same belief network correlate significantly (i.e., become aligned) with the corresponding human opinions (Demo+Train [same category] condition; indicated by lower MAE_{test}). For example, when an LLM agent is initialized to believe that “some people can communicate with the dead” (the training topic x_{train}), then the LLM agent becomes more likely to also believe that “people can project their soul out of their body” (the query topic x_{query}). Concretely, when averaged across nine topic categories, the inclusion of the training topic opinion reduces MAE_{test} from 1.70 (the Demo condition; ChatGPT) to 1.34 (the Demo+Train [same category] condition), representing a 22.54% relative gain. Critically, this effect is limited to topics within the same belief network. If the training topic is from a different topic category (e.g., about the effectiveness of gun control law; Demo+Train [random category] baseline condition), the opinion of the LLM agent on the query topic remains unaligned with the corresponding human ($MAE_{\text{test}} = 1.67$). This supports our hypothesis – opinions on one topic encourage the LLM agents to align their opinions only when the topics are adjacent in the belief network.

Combining demographic information and training topic opinion reaches the best alignment.

While demographic information does not improve alignment on its own (the Demo condition), does it offer any benefit? The contrast between the Demo + Train [same category] condition and the Train [same category] baseline condition answers this question. When removing demographic information from the Demo + Train [same category] con-

Model	Demo + Train [Same Cat.]	
	[Original]	[Balanced]
ChatGPT		
Average MAE_{test}	1.34	1.41
Average Relative Gain (%) \uparrow	22.54	22.19
GPT-4-o-mini		
Average MAE_{test}	1.16	1.21
Average Relative Gain (%) \uparrow	13.88	9.91
Mistral		
Average MAE_{test}	1.29	1.31
Average Relative Gain (%) \uparrow	26.29	24.67
LLaMA 3.1		
Average MAE_{test}	1.60	1.71
Average Relative Gain (%) \uparrow	39.99	23.93

Table 2: Average MAE_{test} and average relative gain of the Demo+Train [Same Cat.] condition for the original condition (“[Original]”) and the variant where we balance the label distribution (“[Balanced]”). Note that balancing the label distribution does not change the superiority of Demo+Train [same category] condition when compared with the Demo condition.

dition, the MAE_{test} increases from 1.34 to 1.48 (ChatGPT), and the relative gain decreases from 22.54 % to 17.19 %. This shows that to reach the best alignment, both the training topic opinion and the demographic should be included.

Alignment does not reflect superficial repetition.

Does increased alignment following the Demo+Train [same category] condition arise from a model tendency to simply repeat the opinion provided for the training topic? Such a pattern might appear to lead to increased alignment simply because the training topic opinion, by definition, correlates with opinions on other topics in the same belief network. To address this concern, we conducted an additional experiment in which we balanced the label distribution in the prompting contexts by constructing reversed framing statements that entail the same semantic meaning. We then included both the original and reversed framing statements in the context. For example, for the original statement “You believe it is *certainly true* that ‘States with stricter gun control laws have *fewer* gun deaths per capita’”, the reversed frame stated “You believe it is *certainly false* that ‘States with stricter gun control laws have *more* gun deaths per capita’”. Both statements were included in the context in random order so the LLM cannot show increased alignment by merely repeating the training topic opinion. Table 2 shows that the LLMs continue to show significant alignment with human opinions (low MAE_{test}) in this case, an effect that must reflect the meaning of the joint informa-

Condition	Topic Category	
	Ghost	Partisan
<i>Baselines.</i>		
Demo	2.58	1.41
Demo + Train [Rand. Cat.]	2.31	1.35
Demo + Train [Same Cat.]	1.29	1.25
<i>Upper bound</i>		
Demo + Same Train + Truth	0.41	0.28
Relative Gain (%)	59.45	14.16

Table 3: Mean absolute error (MAE_{test}) between human respondents and the corresponding LLM agents for each topic category across various LLM agent construction conditions through supervised fine-tuning (SFT). The condition of our main interest (i.e., Demo + Train [Same Cat.]) is boldfaced, which also has the best alignment.

tion $(x_{\text{train}}, o_{\text{train}})$ rather than the opinion label o_{train} alone.

Sensitivity Analyses We evaluated the sensitivity of our result to randomness due to different temperature values when using temperature sampling. Across $T \in \{0, 0.7, 1\}$ using ChatGPT, the results showed consistent trends (Table 7).

Supervised fine-tuning yields similar results.

As shown in Table 3, when the agents are fine-tuned with a training topic x_{train} , they also express more human-like opinions on query topics belonging to the same belief network (i.e., lower MAE_{test} ; the Demo+Train [same category] condition), but not on those belonging to a different network (Demo+Train [random category] condition)—a pattern of results qualitatively similar to in-context learning.

6 Related Work

Aligning human and LLM opinions. Recent studies highlight both the potential and the limitations of using LLMs to emulate human opinions (Argyle et al., 2023; Santurkar et al., 2023; Sun et al., 2024; Feng et al., 2023; Chuang et al., 2024a,b). Argyle et al. (2023) showed that LLMs conditioned on demographic backstories can emulate human voting preferences and language use, but did not investigate topic-specific opinions. Santurkar et al. (2023) found that different models have different inherent opinions that often align with liberal, high-income, well-educated demographics, and that these opinions could not be shifted by

providing demographic role-playing information. The current paper replicates this finding, but additionally suggests that alignment may be shifted via belief networks. To the best of our knowledge no prior work has studied such effects.

Belief networks. A great deal of prior work has studied human belief networks (Boutyline and Vaisey, 2017; Vlasceanu et al., 2024; Keating, 2023; Turner-Zwinkels and Brandt, 2022; Powell et al., 2023; Devine, 2015; Jewitt and Goren, 2016; Baldassarri and Goldberg, 2014; Brandt and Sleegers, 2021) and has developed a range of approaches beyond factor analysis for characterizing these including partial correlation networks (Turner-Zwinkels and Brandt, 2022) or Bayesian networks (Powell et al., 2023). Such networks have been shown to predict “spillover effects” of attitude changes across related topics (Turner-Zwinkels and Brandt, 2022; Powell et al., 2023) in human participants, where a change in a given topic can ripple through the belief network and influence related topics. In the present study, we investigated whether we can leverage the belief network derived from human data to construct LLM agents that more accurately reflect human opinions.

7 Conclusion

We investigated the use of empirically-derived belief networks for promoting alignment of expressed beliefs between Large Language Model (LLM) agents and twinned human participants. We showed that demographic role-playing alone does not produce significant alignment, but that initializing an agent with a human opinion on one topic then aligns opinions on nearby topics within the belief network. The effect does not extend to distant topics within the network. We found similar effects for in-context learning and supervised fine-tuning, for both a proprietary and an open-source LLM. This work highlights a novel and potentially powerful means of enhancing LLM agents’ alignment with human opinions.

Limitations

The scope of topics We considered just 18 topics derived from two orthogonal latent factors identified in prior work. While the Partisan topics are of public interest and the Ghost topics explore an orthogonal dimension, future research could greatly the scope of topics.

The structure of the belief network. We considered belief networks based on two highly distinct clusters to facilitate evaluation. Other studies have used more sophisticated models, such as Bayesian networks (Powell et al., 2023), which allow for precise predictions about topic interrelations. Future work could apply such methods to better characterize belief networks.

The actions of the LLM agents. Our LLM agents expressed their opinions through Likert-scale ratings. This facilitated direct comparison with human responses but may not fully capture the expression of opinions in real-world settings like social media communication. Future studies could explore more complex actions (e.g., writing social media posts) to assess their human-likeness in realistic applications.

Ethics Statement

We aim to develop LLM agents capable of simulating realistic human communicative dynamics, including the expression of potentially harmful beliefs such as misconception about the reality of global warming. Our objective is to facilitate a deeper understanding of social phenomena like misinformation spread in order to identify strategies that mitigate these challenges effectively. Note that under the current setting, the LLM agents only produce Likert-scale ratings from a fixed set of options. Therefore, they are not able to produce unexpected harmful responses. We will release our code base solely for research purposes, and adhere to the terms of use by OpenAI's API ³ and their MIT license ⁴, as well as Mistral AI's non-production license (MNPL) ⁵.

Acknowledgements

We thank the reviewers, the area chair for their feedback. This work was funded by the Multi University Research Initiative grant from the Department of Defense, W911NF2110317 (with Rogers as Co-I), Cohesive and Robust Human-Bot Cybersecurity Teams, the John S. and James L. Knight Foundation (Award Number: MSN231314), and the National Science Foundation through the Con-

vergence Accelerator Track F: Course Correct: Precision Guidance Against Misinformation (Agency Tracking Number: 2230692; Award Number: MSN 266268).

³<https://openai.com/policies/terms-of-use>

⁴<https://github.com/openai/openai-openapi/blob/master/LICENSE>

⁵<https://mistral.ai/licenses/MNPL-0.1.md>

References

- Lisa P Argyle, Ethan C Busby, Nancy Fulda, Joshua R Gubler, Christopher Rytting, and David Wingate. 2023. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3):337–351.
- Delia Baldassarri and Amir Goldberg. 2014. Neither ideologues nor agnostics: Alternative voters’ belief system in an age of partisan politics. *American Journal of Sociology*, 120(1):45–95.
- Andrei Boutyline and Stephen Vaisey. 2017. Belief network analysis: A relational approach to understanding the structure of attitudes. *American journal of sociology*, 122(5):1371–1447.
- Mark J Brandt and Willem WA Sleegers. 2021. Evaluating belief system networks as a theory of political belief system dynamics. *Personality and Social Psychology Review*, 25(2):159–185.
- Raymond B Cattell. 1966. The scree test for the number of factors. *Multivariate behavioral research*, 1(2):245–276.
- Harrison Chase. 2022. [Langchain](#).
- Yun-Shiuan Chuang, Agam Goyal, Nikunj Harlalka, Siddharth Suresh, Robert Hawkins, Sijia Yang, Dhavan Shah, Junjie Hu, and Timothy Rogers. 2024a. Simulating opinion dynamics with networks of llm-based agents. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3326–3346.
- Yun-Shiuan Chuang, Nikunj Harlalka, Siddharth Suresh, Agam Goyal, Robert Hawkins, Sijia Yang, Dhavan Shah, Junjie Hu, and Timothy T Rogers. 2024b. The wisdom of partisan crowds: Comparing collective intelligence in humans and llm-based agents. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 46.
- Christopher J Devine. 2015. Ideological social identity: Psychological attachment to ideological in-groups as a political phenomenon and a behavioral influence. *Political Behavior*, 37:509–535.
- Shangbin Feng, Chan Young Park, Yuhan Liu, and Yulia Tsvetkov. 2023. From pretraining data to language models to downstream tasks: Tracking the trails of political biases leading to unfair nlp models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11737–11762.
- Vincent V Frigo. 2022. *An Examination of Non-Normative Belief Updating Behavior in Humans (Why Is It so Hard to Change Minds?)*. The University of Wisconsin-Madison.
- Caitlin E Jewitt and Paul Goren. 2016. Ideological structure and consistency in the age of polarization. *American Politics Research*, 44(1):81–105.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Henry F Kaiser. 1958. The varimax criterion for analytic rotation in factor analysis. *Psychometrika*, 23(3):187–200.
- David M Keating. 2023. Persuasive message effects via activated and modified belief clusters: toward a general theory. *Human Communication Research*, page hqad035.
- OpenAI. 2022. Introducing ChatGPT. <https://openai.com/blog/chatgpt>. [Accessed 13-10-2023].
- Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. *arXiv preprint arXiv:2304.03442*.
- Joon Sung Park, Lindsay Popowski, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2022. Social simulacra: Creating populated prototypes for social computing systems. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*, pages 1–18.
- Derek Powell, Kara Weisman, and Ellen M Markman. 2023. Modeling and leveraging intuitive theories to improve vaccine attitudes. *Journal of Experimental Psychology: General*, 152(5):1379.
- Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cino Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose opinions do language models reflect? In *International Conference on Machine Learning*, pages 29971–30004. PMLR.
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-llm: A trainable agent for role-playing. *arXiv preprint arXiv:2310.10158*.
- Seungjong Sun, Eungu Lee, Dongyan Nan, Xiangying Zhao, Wonbyung Lee, Bernard J Jansen, and Jang Hyun Kim. 2024. Random silicon sampling: Simulating human sub-population opinion using a large language model based on group-level demographic information. *arXiv preprint arXiv:2402.18144*.
- Amir Taubenfeld, Yaniv Dover, Roi Reichart, and Ariel Goldstein. 2024. Systematic biases in llm simulations of debates. *arXiv preprint arXiv:2402.04049*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Felicity M Turner-Zwinkels and Mark J Brandt. 2022. Belief system networks can be used to predict where to expect dynamic constraint. *Journal of Experimental Social Psychology*, 100:104279.

Madalina Vlasceanu, Ari M Dyckovsky, and Alin Coman. 2024. A network approach to investigate the dynamics of individual and collective beliefs: Advances and applications of the bending model. *Perspectives on Psychological Science*, 19(2):444–453.

Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, et al. 2023. Sotopia: Interactive evaluation for social intelligence in language agents. *arXiv preprint arXiv:2310.11667*.

A List of the 64 Topics in the Belief Survey

Table 4 shows the full statements of the 64 topics in the Belief Survey, including the topic category to which they belong according to the factor analysis result, along with whether they belong to the training or the test partition.

Topic Category	Topic Name	Topic Statement	
Ghost	Dead Talk	No one is able to converse with the dead.	
	Ghost	After someone has died it is not possible to see his or her ghost.	
	Alien Visit	Intelligent beings from outer space have not visited the Earth via spaceships.	
	Soul Walk	It is not possible for anyone to project their soul out of their body.	
	See Future	No one is capable of having visions that accurately predict future events.	
	Astrology	The position of the planets at the time of your birth has no influence on your personality.	
	Roswell	No alien spacecraft has ever crashed near Roswell, New Mexico.	
	Past Life	Nobody can accurately remember living a past life.	
	The Secret	Strongly visualizing your fondest wish does not make it more likely to become a reality.	
	Aura	Health cannot be improved by manipulating a person's aura or electrical field.	
Psychics	Luck	"Lucky streaks" where random events are more likely to favor a person are not real.	
	Dousing	Nobody can sense water using only a forked stick.	
	Pyrokinesis	Nobody can start fires just by thinking about it.	
	Thought Control	Nobody can control another's actions with their mind.	
	Food	Food dropped on the ground for less than five seconds can become contaminated.	
	Palm Reading	It is not possible to predict future life events from markings on a person's palm.	
	Telekinesis	No one is capable of moving objects with his or her mind.	
	Witches	Witches cannot influence events by using magic.	
	Mind Reading	No one is capable of reading another person's thoughts.	
	Moon Landing	US astronauts have landed on the moon.	
Religion	Crystals	Crystals do not have unexplained powers.	
	Lightening	Lightning can strike twice in the same place.	
	Alien Abd	Human beings have not been abducted by aliens from outer space.	
	God	God does not exist.	
	Prayer	Prayer cannot cure illness.	
	Angels	Angels are not real.	
	Religion Explain	Religion does not provide the most accurate explanation for how the universe came into existence.	
	Evil Spirit	It is not possible for a person's actions to be controlled by an evil spirit.	
	Science Expl	Everything that happens can eventually be explained by science.	
	Miracles	Miracles that defy the laws of nature cannot happen.	
Trump	Evolution	Species living on the Earth today have not always existed in their present form.	
	Homicide	In the US, about 80% of white homicide victims are killed by white people.	
	Trump Inaug	More people attended the inauguration of Barack Obama than the inauguration of Donald Trump.	
	Kenya	Barack Obama was born in Hawaii.	
	US Employment	The US unemployment rate in 2016 was lower than 40%.	
	Gov Reg	Government regulations do not always stifle economic growth.	
	Holocaust	The Nazi government in Germany murdered approximately 6 million Jewish people during the second world war.	
	Trump Votes	Hilary Clinton received the most overall votes in the 2016 Presidential election.	
	Abortion	Strongly Republican states have higher rates of abortion than strongly Democratic states.	
	Dem Guns	The official platform of the Democratic Party does not seek to repeal the 2nd Amendment.	
Partisan	Health Insurance	Since the Affordable Care Act (Obamacare) passed, more Americans have health insurance.	
	Gun Control	States with stricter gun control laws have fewer gun deaths per capita.	
	US Deficit	The US deficit decreased after President Obama was elected.	
	Globe Human	Human activity is causing the globe to warm.	
	Globe Warm	The global climate is rapidly growing warmer.	
	Unions	States with strong union protections have lower unemployment than states without such protections.	
	Death Penalty	States that have the death penalty have higher rates of violent crime on average.	
	Economic	US Tax	The United States doesn't have the highest federal income tax rate of any Western country.
		Deport	President G. W. Bush deported fewer undocumented immigrants than President Obama.
		Lower Tax	Lowering taxes does not always lead to economic growth.
Bailout		The rescue of big banks by the federal government aided recovery from the 2008 recession.	

	Gold Stand	Returning to the Gold Standard would make the US more vulnerable to a recession.
LowInfo	Refugee	In 2016 fewer than 100,000 refugees from the Middle East were granted permission to live in the United States.
	US Crime	The violent crime rate in the US has declined over the past 10 years.
	Earth Age	The Earth is not around 6,000 years old.
	Human Trex	The Tyrannosaurus Rex and humans did not live on the Earth at the same time.
	Pub Priv	For a given level of education, private-sector workers typically earn more than government workers.
Health	Body Cleanse	A “body cleanse” in which you consume only particular kinds of nutrients over 1-3 days does not help your body to eliminate toxins.
	Organic	Organic foods are not healthier to eat than non-organic foods.
	Fasting	Regular fasting will not improve your health.
Conspiracy	Twin Towers	The twin towers were not brought down from the inside by explosives during the 9/11 attack.
	JFK	Only one gunman was involved in the assassination of John F. Kennedy.
	Pearl Harbor	President Roosevelt did not know about the attack on Pearl Harbor ahead of time.
	Vaccination	Vaccinations cannot cause Autism.

Table 4: The statements of the 64 topics in the Belief Survey, including the topic category to which they belong according to the factor analysis result.

B The Prompts for LLM Agent Construction Through In-context Learning (ICL)

Table 5 shows the prompts we use to construct and query the LLM agents in the in-context learning setting (§3.3). Different LLM agent construction conditions include various sets of the prompt types. The parts enclosed in curly brackets “{ }” are the placeholders (e.g., {demo_age}, {query_topic_statement}), where they are filled with actual information from either the respondents or the belief survey. As shown in Figure 3 and §3.3, in the **Baseline: No-Demo** condition, only the “Query” prompt is included. In the **Baseline: Demo** condition, both the prompt types “Demographics” and “Query” are included. In the **Demo + Train** conditions (both [same category] and [random category]), the prompt types include “Demographics”, “Training Topic Opinion”, and “Query”. In the **Upper Bound: Demo + Train + Query** condition, the prompt types include “Demographics”, “Training Topic Opinion”, “Query Topic Opinion”, and “Query”.

C The Prompts for LLM Agent Construction Through Supervised Fine-tuning (SFT)

Table 6 shows the prompts we use to construct and query the LLM agents in the supervised fine-tuning setting (§3.3). The demographic information is included in the system message in the same prompt template as in §B. For the topic-specific opinions, however, instead of including them in the prompt, we formulate them as (prompt, response) pairs for supervised fine-tuning, where prompt is the input and response is the output. The prompt templates and examples are shown in Table 6.

D The Choice of Number of Factors in Factor Analysis

To determine the number of factors to retain in our factor analysis (FA), we visualize the scree plot in Figure 4. We see that the explained variance plateaus after including 9 factors (the “elbow point”). Therefore, we decide to retain 9 factors.

E Supervised Fine-tuning Details

In this section, we elaborate the different strategies used for constructing LLM agents through supervised fine-tuning.

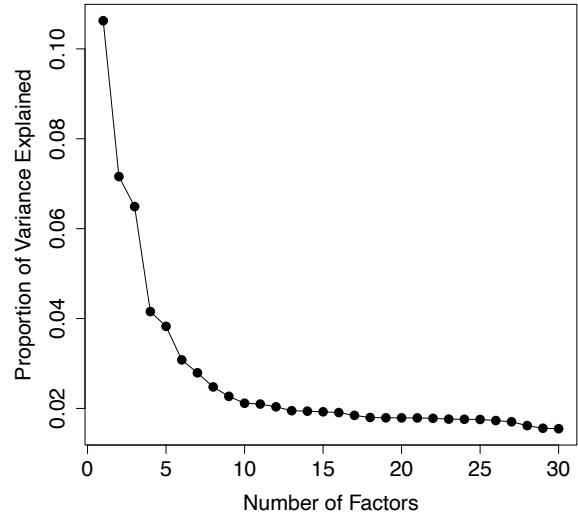


Figure 4: The scree plot of the factor analysis solution.

- Baseline: No-Demo.** Baseline without fine-tuning, (identical to same condition in ICL).
- Baseline: Demo.** Baseline without fine-tuning, identical to same condition in ICL.
- Demo+Train [same category]:** For each topic category we constructed the dataset $\mathcal{D}_{\text{SFT}} = \{(d_i, x_{\text{train},i}), o_{\text{train},i}\}_{i=1}^N$. We then fine-tuned the LLM with input context providing the demographic information along with the training topic statement (d, x_{train}) , and using the corresponding human Likert-scale response o_{train} as the target. After fine-tuning, we assessed the LLM agent’s opinion on query topics x_{query} belonging to the same topic category x_{train} ⁶. This is the critical condition of interest that tests cross-topic generalization. The verbatim prompts are in §C.
- Baseline: Demo+Train [random category]:** Similar to Demo+Train [same category] condition, but the training topic opinion $(x_{\text{train}}^\dagger, o_{\text{train}}^\dagger)$ is from a different topic category as the query topic x_{query} , allowing us to assess whether generalization is restricted to topics in the same belief category.
- Upper Bound: Demo+Train+Query.** Upper bound without fine-tuning, identical to same condition in ICL.

ChatGPT (gpt-3.5-turbo-0125) is fine-tuned through OpenAI’s fine-tuning API⁷. These

⁶For example, we fine-tuned an LLM on the respondents’ opinions on the training topic for the Ghost category, then queried its opinion on the test topics in the Ghost category.

⁷<https://platform.openai.com/docs/guides/fine-tuning>

Prompt Type	Message Type (LangChain)	Prompt Template	Example
Demographics	<i>System Message</i>	You are role playing a real person. You are a {demo_gender}. You are {demo_age} years old. The highest education You have completed is {demo_education}. Your race is {demo_race}. Your household income is {demo_income}. The population of your city is {demo_city_pop}. You would characterize your hometown as {demo_urban_rural}, and you are from the state of {demo_state}. Your political leaning is {demo_party}.	You are role playing a real person. You are a {Male}. You are {41} years old. The highest education You have completed is {Some college but no degree}. Your race is {White}. Your household income is {40,000–59,999}. The population of your city is {100,000 - 500,000}. You would characterize your hometown as {Urban (City)}, and you are from the state of {Florida}. Your political leaning is {Democrat}.
Training Topic Opinion	<i>System Message</i>	You believe that {training_topic_statement (x_{train})} is {opinion_response (o_{train})}.	You believe that {States with stricter gun control laws have fewer gun deaths per capita.} is {Probably True}.
Query Topic Opinion	<i>System Message</i>	You believe that that {query_topic_statement (x_{query})} is {opinion_response (o_{query})}.	You believe that {The global climate is rapidly growing warmer.} is {Certainly True}.
Query	<i>User Message</i>	Now, what is your opinion on the following statement using the following scale of responses? {query_topic_statement (x_{query})} is Certainly False, {query_topic_statement (x_{query})} is Probably False, {query_topic_statement (x_{query})} is Lean False, {query_topic_statement (x_{query})} is Lean True, {query_topic_statement (x_{query})} is Probably True, {query_topic_statement (x_{query})} is Certainly True. Statement: {query_topic_statement (x_{query})} Your opinion on the scale of responses:	Now, what is your opinion on the following statement using the following scale of responses? {The global climate is rapidly growing warmer.} is Certainly False, {The global climate is rapidly growing warmer.} is Probably False, {The global climate is rapidly growing warmer.} is Lean False, {The global climate is rapidly growing warmer.} is Lean True, {The global climate is rapidly growing warmer.} is Probably True, {The global climate is rapidly growing warmer.} is Lean True, {The global climate is rapidly growing warmer.} is Certainly True Statement: {The global climate is rapidly growing warmer.} Your opinion on the scale of responses:

Table 5: The prompts used for the LLM agent construction and querying in the in-context learning setting.

Prompt Template	Example Prompt	Response Template	Example Response
What is your opinion on the following statement using the following scale of responses? Certainly False that {query_topic_statement (x_{query})}, Probably False that {query_topic_statement (x_{query})}, Maybe False that {query_topic_statement (x_{query})}, Maybe True that {query_topic_statement (x_{query})}, Probably True that {query_topic_statement (x_{query})}, Certainly True that {query_topic_statement (x_{query})}. Statement: {query_topic_statement (x_{query})}. Please choose your response from the following list of options: Certainly False, Probably False, Maybe False, Maybe True, Probably True, Certainly True.	What is your opinion on the following statement using the following scale of responses? Certainly False that {States with stricter gun control laws have fewer gun deaths per capita}, Probably False that {States with stricter gun control laws have fewer gun deaths per capita}, Maybe False that {States with stricter gun control laws have fewer gun deaths per capita}, Maybe True that {States with stricter gun control laws have fewer gun deaths per capita}, Probably True that {States with stricter gun control laws have fewer gun deaths per capita}, Certainly True that {States with stricter gun control laws have fewer gun deaths per capita}. Statement: {States with stricter gun control laws have fewer gun deaths per capita} Please choose your response from the following list of options: Certainly False, Probably False, Maybe False, Maybe True, Probably True, Certainly True.	My Response: {opinion_response}	My Response: {Certainly True}

Table 6: The prompts used for the LLM agent construction and querying in the supervised fine-tuning setting.

were the hyper-parameters used in fine-tuning:

- Number of Epochs: 3
- Batch Size: 1
- Learning Rate Multiplier: 2

F The Full Factor Analysis Results

In Figure 2b in the main text, we only show the factor loading matrix of the Ghost and the Partisan factors, and the corresponding topics. In this section, we discuss the full factor analysis result.

The factor analysis reveals nine latent factors underlying the 64 topics. Figure 5 shows the full factor loading matrix. The red blocks highlight strong correlations among opinions within each factor, indicating that endorsing one conception in a cluster often predicts opinion in other conceptions within the same cluster. We assign the name of each factor based on its constituent topics: Ghost, Psychics, Religion, Trump, Partisan, Economic, LowInfo, Health, and Conspiracy. The 64 topics are categorized by which factor they have the highest loadings on. For instance, the topic about communication with the dead belongs to the

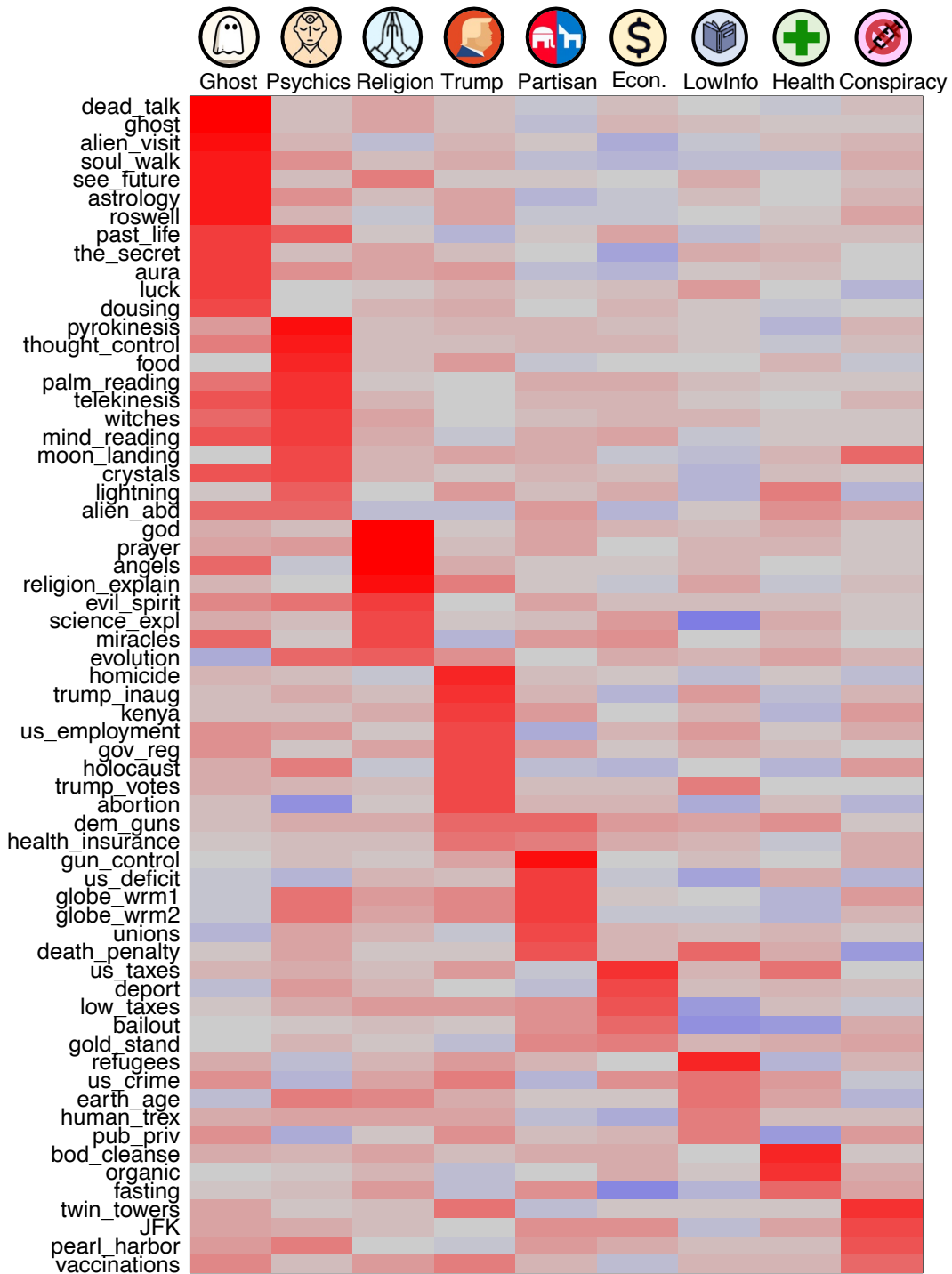


Figure 5: The factor loading matrix of the Controversial Belief Survey. The column indicates the nine factor, and the rows are the 64 topics. Red indicates topics that load highly on a factor, gray indicates near 0 loading, and blue indicates loading in the negative direction. We focus on the Ghost category and Partisan categories, highlighted by the green box and the violet box respectively. The topics in the Ghost category has minimal loading on the Partisan factor and vice versa (highlighted by the black boxes). The full statement of each topic is in Table 4 (§A).

Ghost category because it has the highest loading on the Ghost factor (Table 4 shows the full list of topics and categories).

G Compute Resources

We ran all experiments with Mistral and LLaMA 3.1 on a GPU machine equipped with 1x NVIDIA A100. The experiments with ChatGPT and GPT-4-o-mini cost about 400 USD.

Condition	Temperature		
	0	0.7	1
<i>Baselines</i>			
No-Demo	1.80	1.68	1.66
Demo	1.70	1.70	1.71
Demo + Train [Rand. Cat.]	1.66	1.67	1.68
Train [Same Cat.]	1.43	1.48	1.49
Demo + Train [Same Cat.]	1.37	1.34	1.44
<i>Upper bound</i>			
Demo + Same Train + Truth	0.42	0.42	0.53
Average Relative Gain (%)	18.57	22.54	14.64

Table 7: Average MAE_{test} and average relative gain of each LLM agent (powering by ChatGPT) construction condition across three temperature values.

H Sensitivity Analysis

Sensitivity Analyses We evaluate the sensitivity of our result to randomness due to different temperature values when using temperature sampling. Across $T \in \{0, 0.7, 1\}$ using ChatGPT, the results show consistent trends (Table 7).