

Beyond Demographics: Enhancing Cultural Value Survey Simulation with Multi-Stage Personality-Driven Cognitive Reasoning

Haijiang Liu^{123*}, Qiyuan Li¹², Chao Gao¹², Yong Cao⁴,
Xiangyu Xu³, Xun Wu³, Daniel Hershcovich⁵, Jinguang Gu^{12†}

¹ School of Computer Science and Technology, Wuhan University of Science and Technology, Wuhan 430065, China

² Hubei Province Key Laboratory of Intelligent Information Processing and Real-time Industrial System,
Wuhan University of Science and Technology, Wuhan 430065, China

³ Innovation, Policy and Entrepreneurship Thrust,

The Hong Kong University of Science and Technology (Guangzhou), Guangzhou, 511453, China

⁴ Tübingen AI Center, University of Tübingen, Tübingen, 72074, Germany

⁵ Department of Computer Science, University of Copenhagen, 2100 Copenhagen, Denmark
{alecliu, simon}@ontoweb.wust.edu.cn

Abstract

Introducing **MARK**, the **M**ulti-**s**tAge Reasoning framework for cultural value survey response simulation, designed to enhance the accuracy, steerability, and interpretability of large language models in this task. The system is inspired by the type dynamics theory in the MBTI psychological framework for personality research. It effectively predicts and utilizes human demographic information for simulation: life-situational stress analysis, group-level personality prediction, and self-weighted cognitive imitation. Experiments on the World Values Survey show that MARK outperforms existing baselines by 10% accuracy and reduces the divergence between model predictions and human preferences. This highlights the potential of our framework to improve zero-shot personalization and help social scientists interpret model predictions.¹

1 Introduction

Social surveys are essential for generating insights and supporting decision-making processes in research, yet traditional human-centric experiments often incur significant costs. Therefore, recent studies leverage AI technologies with human-behaviour simulating abilities to create more accurate and interpretable simulations as a valuable alternative (Horton, 2023; Argyle et al., 2023).

As illustrated on the left side of Figure 1, solutions like demographic augmentation (Kwok et al., 2024; Zhao et al., 2024a) and opinion insertion (Hwang et al., 2023) were developed to improve model simulation responses in cultural value surveys. While these studies have primarily relied on

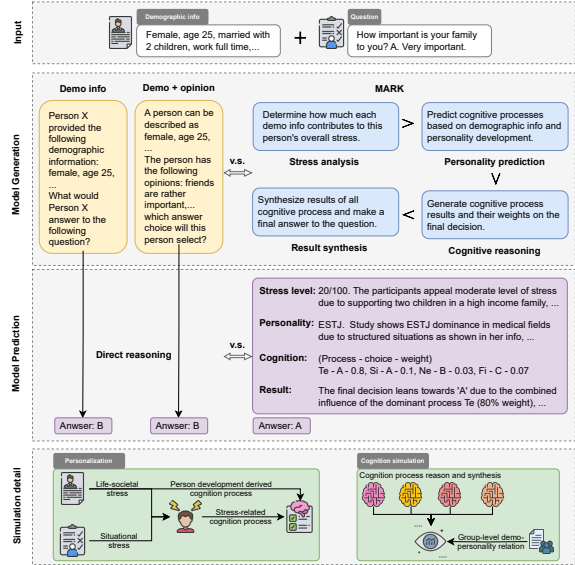


Figure 1: The MARK framework enhances cultural value survey simulation via multi-stage reasoning over MBTI personality theory with life-situational stress analysis, demographic-based group-level personality prediction, and weighted cognitive imitation. These enhancements increase simulation accuracy, decision interpretability, and response personalization.

prompting that uses personal demographic information (shown as "demo info") as input to the model, models still struggle with zero-shot personalization, and prediction interpretation is still challenging.

In this paper, we focus on **improving cultural value survey simulation by incorporating psychological theory** into the modeling of human responses. Our key insight is that human responses to survey questions are significantly influenced by their mental status (like stress) and associated cognitive processes. To capture these influences, we build a **M**ulti-**s**tAge Reasoning framework (**MARK**) that simulates how different personality types process information and make decisions. MARK draws on the MBTI personality type dynamic theory (Mitchell, 2001) to model cognitive

*Work done during his visit at The Hong Kong University of Science and Technology (Guangzhou).

†Corresponding Author: Jinguang Gu. (E-mail: simon@ontoweb.wust.edu.cn).

¹<https://github.com/alexc-l1/MARK-Multi-stAge-Reasoning-framework>

functions and their interactions with stress and personal experience, aiming to achieve accurate and psychologically grounded survey simulations than previous methods that primarily rely on demographic information.

The process begins with a **stress analysis** to evaluate societal pressures based on demographic data, which generates sociodemographic prompts that enhance the simulation. Following this, a two-step **personality prediction** is carried out to infer cognitive functions from these descriptions, utilizing MBTI personality theory. During the survey simulation, a **cognitive reasoning** stage mimics human thought processes under both life- and situational-stress, influenced by the interplay between human traits and specific questions. Finally, the process concludes with a **synthesis** stage that reviews the reasoning outcomes to produce the final answer.

To reduce personality assumption bias, we evaluate our framework on a clustered version of the World Values Survey (see Appendix D.1). Experiments show that MARK can outperform existing baselines by 10% and reduce the divergence between model predictions and human preferences. Further analysis reveals the robustness of our multi-stage reasoning framework - *MARK performs stably during simulation even with uncertain personalities*. The fidelity of our framework suggests that incorporating personality-driven cognitive processes is a promising direction for improving survey simulation accuracy while providing deeper insights into human decision-making patterns.

Our contributions are threefold:

1. We designed a multi-stage reasoning framework with personality cognition called MARK to simulate the human decision process on value surveys.
2. Experiments on simulations show that MARK outperforms other baselines by presenting higher accuracy and generalizes consistently across the U.S. and Chinese cultures.
3. Further analysis suggests that the system performs robustly with uncertain personalities, benefiting from the multi-stage design.

2 Related Work

In the field of AI for social science research, researchers have increasingly explored the use of LLMs to **simulate human responses to survey questions** (Argyle et al., 2023; Cao et al., 2025; Liu

et al., 2025). Researchers investigate LLMs' predictive ability of simulating cross-cultural variations (Kwok et al., 2024; Cao et al., 2023; AlKhamissi et al., 2024), value diversity (Santurkar et al., 2023; Zhao et al., 2024b), and modeling of complex value systems (Benkler et al., 2023; Arora et al., 2023).

Moreover, LLMs also show promise in replicating human samples across various fields, including economics (Horton, 2023), psychology (Aher et al., 2023), and other domain tasks (Luo et al., 2025). However, these studies also highlight the limitations of LLMs in capturing the intricacies of human behavior, particularly in nuanced decision-making situations (Bisbee et al., 2023).

To improve simulation accuracy, researchers have investigated various prompting strategies (Mou et al., 2024). A key approach has been the integration of demographic information into LLMs to enhance zero-shot simulations (Deshpande et al., 2023), known as **sociodemographic prompting**. However, studies have found that LLMs tend to produce stereotypes and biases under straightforward demographic prompting conditions (Cheng et al., 2023), complicating efforts to simulate human behaviors with greater precision.

To address these challenges, researchers have explored more sophisticated prompting methods. Hwang et al. (2023) proposed incorporating user opinions to better reflect specific group perspectives. Zhou et al. (2023) demonstrated that including cultural background information could enhance model transferability. Cao et al. (2025) shows that fine-tuning on first token probabilities can better align simulated and real response distributions. However, recent studies by Li et al. (2024) and Shen et al. (2024) reveal that current methods still fall short in simulating preference diversity and the intricate decision-making processes characteristic of human behavior. This limitation underscores the need for more personalized and psychologically grounded approaches to survey simulation.

3 MARK: Multi-stAge Reasoning framework

The Myers-Briggs Type Indicator (Clawson, 2008, MBTI) is a widely used personality assessment tool that categorizes individuals into 16 distinct personality types. Each personality type is characterized by the interactions between four cognitive functions that describe the mental processes in the order of how we naturally access them and their

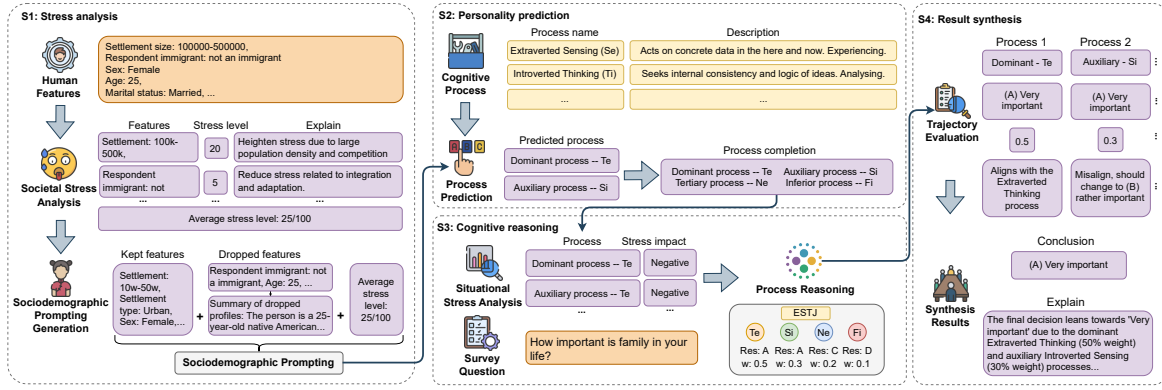


Figure 2: Overview of the proposed multi-stage reasoning framework for simulating individual cognition in survey responses. The framework comprises four stages: 1) **Stress Analysis**—evaluates life-societal pressures to generate stress levels and a summarized human profile; 2) **Personality Prediction**—predicts cognitive functions and personality traits based on the profile; 3) **Cognitive Reasoning**—assesses stress impact of specific questions, conducts reasoning using predicted functions, and assigns weights to finalize answers; and 4) **Synthesis**—reviews reasoning history, validates results, and synthesizes them into a cohesive response. Color-coded data sources: **Orange** boxes indicate survey data, **yellow** boxes represent function call outputs, and **purple** boxes denote text generated by LLMs.

development—Dominant, Auxiliary, Tertiary, and Inferior functions (Mitchell, 2001).

The Dominant Function is the most developed, conscious, and comfortable aspect of an individual’s personality, serving as their core and primary cognitive function. The Auxiliary Function supports and balances the Dominant Process, aiding in perception, judgment, and the balance between Extraversion and Introversion. The Tertiary Function is less competent and comfortable than the Dominant and Auxiliary, typically activated under stress for decision-making. The Inferior Function is the least developed, emerging under extreme stress and potentially causing individuals to feel out of control. (Details in Appendix A)

To capture how these personality dynamics influence survey responses, we develop MARK, a multi-stage reasoning framework that simulates human decision-making through the lens of personality type dynamics. Our key insight is that personality traits and their associated cognitive functions significantly impact how individuals process information and respond to survey questions, particularly under varying levels of stress.

As presented in Figure 2, MARK organizes four major reasoning stages: stress analysis, personality prediction, cognitive reasoning, and synthesis. We provide a brief summary of stage functionality in Table 1 (detailed in Appendix C).

Stage 1: Stress analysis The MBTI-type dynamic theory emphasizes the impact of stress on

triggering different thinking functions and inducing the positive or negative reasoning pattern of the process (Clark, 2002). Before the personality reasoning of each human subject, we ask the model to evaluate the features recorded in surveys and predict the stress level of each feature S_{c_i} (0 to 100). We use the average stress level of all the features as the overall stress level S_{overall} of the participant.

As personalities develop throughout life, the demographic information captured in our surveys reflects this development and provides insights into possible personality types. However, identifying which features have the greatest impact on simulation results presents a challenge. To address this, we developed an automatic sociodemographic prompting method that preserves features associated with high-stress levels and key personality interactions. Our method analyzes both relationship dependencies between features and feature complexity to maintain the integrity of preserved features while summarizing less critical ones. For any demographic information not retained, the model generates profile descriptions that maintain the overall representation of the human subject.

Stage 2: Personality prediction The MBTI type dynamics theory (Mitchell, 2001) establishes a comprehensive pattern of how cognitive functions form different personalities. The personality cognition compass (Figure 3) shows that personality types can be determined once the first two functions are selected.

Stage	Step	Input	Key Function	Output
1. Stress Analysis	Feature Evaluation Profile Filtering	Human features Stress-scored features	Assign stress levels (0-100) to each feature Retain high-stress features, summarize dropped ones	Features with stress scores Filtered feature profile
2. Personality Prediction	Process Selection	Filtered features + stress level	Select MBTI cognitive processes	D-A-T-I processes
3. Cognitive Reasoning	Stress Impact Assessment Multi-Process Reasoning	MBTI processes + survey questions Stress-impacted processes + survey	Determine if stress triggers negative aspects Generate reasoning for each process with weights	Processes with stress impact Weighted reasoning results
4. Synthesis	Result Integration	All process reasoning + weights	Evaluate consistency and synthesize final answer	Final survey response

Table 1: Overview of the four-stage multi-agent reasoning framework for human value preference simulation.

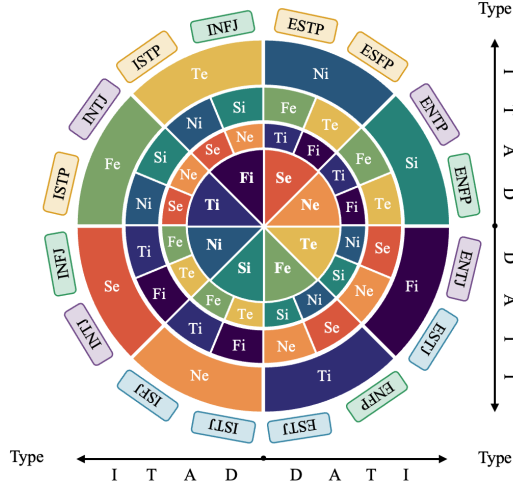


Figure 3: Personality cognition compass: how cognitive functions merge to personality types. The "DATI" represents "Dominant", "Auxiliary", "Tertiary", and "Inferior" functions. Once the first two functions are selected, the corresponding personality type is set.

While the MBTI assessment traditionally relies on self-reported behaviors, demographic factors serve as **indirect proxies** for environmental pressures that systematically influence personality development through mediating pathways (Myers, 1998; Rentfrow et al., 2013). Rather than assuming direct demographic-personality mappings, our framework **leverages demographic features to parameterize contextual priors** reflecting how occupational, educational, and geocultural environments may shape cognitive function preferences over time. This approach provides greater mechanistic sophistication than direct demographic-to-attitude prediction by modeling personality as an intermediate latent construct influenced by environmental factors.

During personality prediction, we prompt the model to predict dominant cognitive functions using demographic information as environmental context. The system selects the primary function, then applies type dynamics theory to identify auxiliary function candidates, reducing prediction errors. Once these two functions are determined, the

remaining functions are derived to form the complete personality type. The prediction algorithm is detailed in Appendix B.

Stage 3: Cognitive reasoning The MBTI cognitive process exhibits two reasoning patterns influenced by societal and situational stress (Clark, 2002): 1) Under low situational stress, life-societal stress predominates, leading to normal cognitive functioning; 2) Under high situational stress, often triggered by past trauma, individuals experience heightened stress and overuse of cognitive processes, similar to Post-traumatic Stress Disorder (Fang et al., 2020, PTSD). In this stage, we evaluate the interaction of these stress types to determine if individuals are vulnerable to cognitive overload. Detailed descriptions of normal and overused functions are provided in Appendix A.

For each thinking process, we assign stress impact levels and ask the model to reason about cognitive status and generate value preferences. To determine the confidence in each function’s result, we also assign influence weights for each process, which are used in the synthesis stage.

Stage 4: Synthesis During synthesis, we collect and document the complete reasoning trajectory. The model reviews these to ensure consistency with the explanations and process descriptions. If it finds misalignments, it rewrites the results to minimize hallucinations and verifies that appropriate weights are assigned to each process in the MBTI hierarchy.

Following this evaluation, the model synthesizes the reasoning results based on weights and selects the most appropriate survey response option. To enhance explainability, it provides a brief explanation of how weights influenced the final decision.

4 Experiment

4.1 Dataset

Our simulation samples are drawn from the 7th wave of the World Values Survey, a comprehensive global study that examines changing values and their influence on social and political life. For

Model	GLM-4-air				GPT-4o			Doubao-1.5-pro				
Global Distribution												
Avg. Metrics	1-JSD \uparrow		EMD \downarrow	κ \uparrow	1-JSD \uparrow		EMD \downarrow	κ \uparrow	1-JSD \uparrow		EMD \downarrow	κ \uparrow
Random selection [†]	0.3319		0.2848	0.0000	0.3319		0.2848	0.0000	0.3319		0.2848	0.0000
Without Demo.	0.4492		0.1398	0.0125	0.4421		0.1177	0.0138	0.3788		0.1455	0.0112
Durmus et al. (2024)	0.2664		0.1466	0.0142	0.4411		0.0981	-0.0071	0.3544		0.1280	0.0116
Liu et al. (2025)	0.4371		0.0490	0.0118	0.4620		0.0211	0.0113	0.5258		0.0218	0.0104
Demo.+Ideo.	0.4336		0.1041	0.0011	0.6404		0.0987	0.0112	0.5520		0.1011	0.0180
Demo.+Ideo.+Opinion	0.5425		0.1112	0.0180	0.6449		0.0981	0.0097	0.5695		0.0995	0.0198
Zhao et al. (2024a)	0.4267		0.1280	0.0024	0.6174		0.1113	0.0111	0.3974		0.1265	0.0123
MARK (Ours)	0.5963		0.1026	0.0127	0.6438		0.0985	0.0128	0.6111		0.1141	0.0123
Sampled Distribution												
Avg. Metrics	ACC (%)	1-JSD \uparrow	EMD \downarrow	κ \uparrow	ACC (%)	1-JSD \uparrow	EMD \downarrow	κ \uparrow	ACC (%)	1-JSD \uparrow	EMD \downarrow	κ \uparrow
Demo.+Ideo.	25.49	0.3539	0.0741	0.02	32.30	0.4075	0.0755	0.09	30.75	0.4563	0.0685	0.02
Demo.+Ideo.+Opinion	25.06	0.4313 (0.6)	0.0669	0.11	33.07	0.4069	0.0739	0.09	31.45	0.4723	0.0703	0.08
Zhao et al. (2024a)	26.98	0.3814	0.0452	0.05	36.96	0.4654	0.0610	0.12	24.23	0.3584	0.0364	0.05
MARK (Ours)	33.69	0.4348	0.0887	0.15	38.11	0.4879	0.0826	0.15	46.98	0.5195	0.0561	0.09

Table 2: MARK simulation performance on U.S. social-survey data, evaluated under both global-distribution and sampled-distribution settings. ‘Avg.’ denotes the overall mean metric values, (value) denotes the p-value of results between baseline and MARK larger than 0.05. MARK shows improvements on sampled distributions by achieving the **highest simulation accuracy** while **maintaining low divergences**, with most improvements being statistically significant. It also shows **generalization to global distributions** with unseen demographics.

evaluation purposes, we utilize survey data from the United States and China as case studies.

Setting 1: Global Distribution To evaluate the generalization of MARK on responses over the whole survey population in specific cultural backgrounds, we also compare the model-predicted responses based on sampled demographic information, presented in setting 2, with the distribution divergence and agreement with the whole survey population (*Global Distribution* setting).

Setting 2: Sampled Distribution To encompass a broader range of human features, we examine 30 features detailed in the survey, such as continent, residential area, and education level (Zhao et al., 2024a). To reduce assumption bias for individual data, representatives are selected through a group-level clustering approach (*Sampled Distribution* setting) reflecting the diversity of demographics.

We conducted a group-level demographic-based clustering analysis in the following steps: 1) Numerically transform and normalize the demographic feature attributes for human subjects from the social survey. 2) Implement the Silhouette Analysis (Oncevay et al., 2020) to determine cluster optimality and apply the K-means clustering.

Based on analysis (see Appendix D.1), we conducted the configuration of 20 clusters and sampled 20 subjects from each cluster to represent U.S. culture, utilizing a dataset of 102,800 QA pairs. To evaluate the generalizability of our method, we ap-

plied the same clustering approach to the Chinese version of the survey and selected the cluster centroids as our simulation targets. Ultimately, we obtained 20 subjects for the simulation, using a dataset of 4,340 QA pairs.

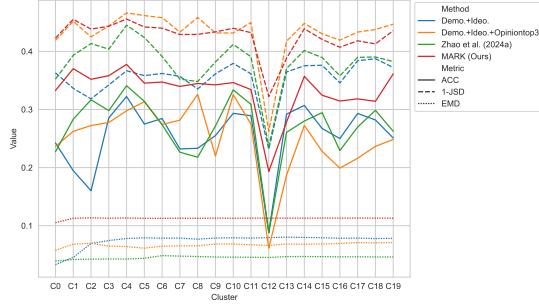
Moreover, we augmented the World Values Survey with the assumed personalities (oracle personality) of human participants to study the robustness of our framework. (See Appendix D.2)

4.2 Baselines

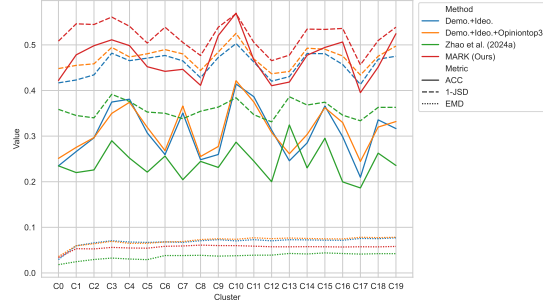
We introduce several baseline methods without personal demographic info by Durmus et al. (2024) to investigate the generalization of MARK on the *Global Distribution* setting:

Random selection[†] We present a random selection as the lower bound performance for the task. Given that random selection is independent of any model, the performance outcomes remain consistent across all models. To facilitate improved comparison, these results are reiterated in the accompanying table.

Without Demo. (Durmus et al., 2024) uses multiple-choice survey questions originally written without any indication of the personal demographic information (DP setting in their paper). **Nation Info Only** (Durmus et al., 2024; Cao et al., 2025; Liu et al., 2025) asks LLMs to predict how someone from the U.S. or China would respond to these questions. We chose two standard baseline prompt as this kind of baseline: Durmus et al. (2024) (CP



(a) Performance of GLM-4-air model.



(b) Performance of Doubao-1.5-pro model.

Figure 4: Group-level performance of survey simulation on U.S. culture in the *Sampled Distribution* setting. MARK shows significant improvement over other baseline methods on most clusters in terms of ACC and 1-JSD.

setting) and Liu et al. (2025).

We also select three baseline methods presented by Hwang et al. (2023) and Zhao et al. (2024a) using the demo. info to investigate the alignment of LLMs with individual user preferences, focusing on modeling user demographics, ideologies, and past opinions to predict future responses:

Demo.+Ideo. (Hwang et al., 2023) combines demographic information and ideological inclinations of a user to predict their responses to survey questions. **Demo.+Ideo.+OpinionTop3** (Hwang et al., 2023) combines demographic information, ideological inclinations, and the top 3 most relevant opinions to predict user responses. **Zhao et al. (2024a)** focuses on three demographic variables (i.e., the continent, resident area, and education level) to indicate a human and asks the model to predict their preferences on survey questions.

4.3 Models

All the experiment is conducted with GLM-4-air (GLM et al., 2024), GPT-4o (gpt-4o-2024-08-06) (OpenAI et al., 2024)², and Doubao-1.5-pro³ via API⁴, and the embedding model used by Demo.+Ideo.+OpinionTop3 is embedding-v3 by ZHIPU AI⁵. The experiments for the U.S. and Chinese simulations were completed with a budget of \$1000. The generation temperature is 0.9, and the max token length is 4096.

²Due to budget limitations, we only conduct simulations on cluster centers on GPT-4o model.

³<https://team.doubao.com/en/direction/llm>

⁴We also test on Llama-3.1-7B/13B, and found they are limited in conducting cognitive simulations, likely due to the complexity of instructions and human features.

⁵<http://open.bigmodel.cn/dev/howuse/model>

4.4 Metrics

We evaluate the simulation performance via the preference distribution gathered from the responses and prediction accuracy with exact text matching.

Response accuracy (ACC) ACC measures how well the model predicts the exact responses compared to the actual labels. We obtain ACC by calculating the percentage of the correct responses of all responses in the cluster and all simulation targets. Higher ACC indicates better alignment.

Distribution divergence (1-JSD, EMD) 1-JSD and EMD metrics evaluate how well the distribution of the model’s predictions aligns with the actual human distribution. (1) 1-JSD, known as 1 minus the Jensen-Shannon divergence, implemented by Durmus et al. (2024), is a symmetric measure of the similarity between two probability distributions; (2) the Earth Mover Distance (Rubner et al., 1998, EMD), known as the Wasserstein distance, quantifies the minimum amount of work required to transform one distribution into another.

Particularly for the Sampled Distribution setting, we collect responses from each cluster per question, then plot distributions for human participants and model predictions within the cluster. Then, we calculate the 1-JSD and EMD score over these distributions. We took the mean values across all clusters as average scores. Lower EMD and higher 1-JSD values indicate better alignment.

Cohen’s κ utilized by (Badshah and Sajjad, 2024) measuring the percentage of agreements on the group-level voting. We calculate the agreement of model predictions and human responses using Cohen’s κ . Higher agreement shows better performance.

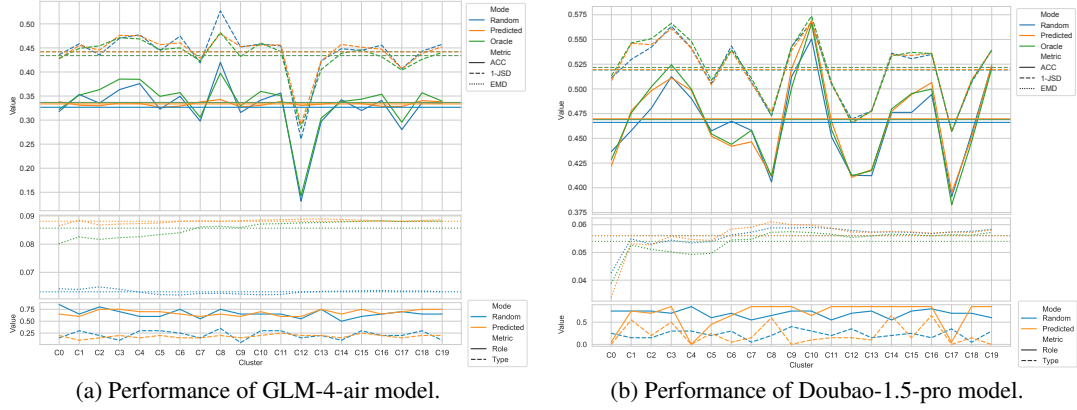


Figure 5: Group-level results of survey simulation on U.S. culture with various personality strategies in the *Sampled Distribution* setting. The horizontal line indicates the average score of each metric. MARK shows robust performance across strategies and exceeds with oracle personalities.

5 Results

This section presents the experimental results and examines the impact of incorporating personality traits for simulation. Additionally, we evaluate the system’s robustness, generalization on the Chinese culture, and significance of each component in the MARK system through ablation studies conducted on best- and worst-performing models.

5.1 Generalization to Global Distribution

Table 2 (top) reports results when each model’s aggregate predictions—conditioned solely on demographics (or nothing)—are compared to the full U.S. survey population histogram. MARK shows competitive performance: it achieves the competitive 1-JSD (0.5963 on GLM-4-air; 0.6438 on GPT-4o; 0.6111 on Doubao-1.5-pro) and the lower EMD (0.1026, 0.0985, 0.1141, respectively), indicating its simulated population closely mirrors the true preference distribution. Cohen’s κ remains at or above the level without demo. info and nation info only prompts, underscoring robust group-level agreement even under sparse conditioning.

5.2 Sampled-Distribution Performance

Table 2 (bottom) summarizes each model’s accuracy (ACC), distributional alignment (1-JSD, EMD), and category-agreement (Fleiss’ κ) when predicting the held-out human responses from our sampled participant pool. Across all three base LLMs (GLM-4-air, GPT-4o, Doubao-1.5-pro), MARK achieves the highest ACC, improving over the strongest “Demo+Ideo+Opinion” baseline by 4.3% on GLM-4-air, 8.9% on GPT-4o, and a striking 15.5% on Doubao-1.5-pro. These gains are

statistically significant ($p < 0.05$), **demonstrating its robustness in predicting exact responses.**

In tandem with higher accuracy, MARK also yields the best distributional fidelity: it attains the largest 1-JSD scores (indicating minimal Jensen-Shannon divergence) and competitive EMD values. For example, on GLM-4-air, MARK’s 1-JSD of 0.4348 and EMD of 0.0887 outperform both the pure demographic baselines (1-JSD ≈ 0.35 – 0.43 , EMD ≈ 0.045 – 0.074) and Zhao et al. (2024b). Finally, MARK’s Fleiss’ κ (0.15) more than triples the “Demo+Ideo” baseline (0.02), demonstrating that its categorical vote proportions align substantially better with the actual human label distributions.

Further analysis in Figure 4 shows that MARK presents different sensitivity over certain groups of humans with similar demographic background: For example, GLM-4-air shows a *strong performance in demographics of married working full-time jobs* like in C_0 (34.28%) and C_5 (34.59%). And Doubao-1.5-pro *excels in demographics of skilled workers for private business* as in clusters C_{10} (56.96%) and C_{19} (52.45%), where it outpaces competitors by large margins.

These superior performances of MARK can be attributed to its integration of MBTI personality cognition, which enhances its ability to simulate human-like responses in social surveys. By incorporating cognitive traits, MARK captures nuanced behavioral patterns of the demographic information and value preferences. This cognitive augmentation allows MARK to better align with real-world survey responses, leading to higher accuracy and demonstrating the importance of psychological di-

mensions in social survey simulations.

5.3 Framework Robustness

Due to limited reliable personality data, we use Gemini-2-pro to estimate the likelihood of different personalities selecting survey responses, designating the most probable as the oracle personality⁶. Three personality assignment strategies—predicted (*Predicted*), random (*Random*), and oracle (*Oracle*)—are integrated into MARK reasoning to assess robustness against personality uncertainty.

As shown in Figure 5, MARK demonstrates strong robustness across all assignment strategies on all models, with simulation accuracy (ACC) and response distribution similarity (1-JSD) remaining **stable**. For example, Doubao-1.5-pro’s ACC stays at 0.47 (± 0.04 –0.05) and 1-JSD between 0.52 and 0.53 across strategies, with similar consistency observed in GLM-4-air.

Personality similarity (role- and type-level) also has **minimal influence** on simulation performance. Although role-level MBTI similarity is moderately high, type-level alignment is much lower, especially under random assignment. Importantly, clusters with low type similarity do not show reduced ACC or increased divergence. Performance variance across clusters is nearly identical between strategies, and performance curves overlap, indicating stability even with noisy personality data.

MARK’s robustness is attributed to its design: *group-level experimental structure*, *weighted cognitive imitation*, and *stress-informed reasoning*, all of which collectively mitigate the impact of noisy personality inputs and ensure consistent performance across experimental conditions.

5.4 Generalization to Chinese Culture

We evaluate the system’s generalization to Chinese culture using both Chinese and English prompts on the *Sampled Distribution* setting.

Chinese prompts As shown in Table 3b, MARK consistently outperforms baselines in accuracy and response alignment. For Doubao-1.5-pro, MARK achieves 45.16% ACC, significantly higher than previous benchmarks, and for GLM4, it closely matches the top-performing baseline that leverages additional human opinion data. MARK also records the best distributional metrics (Avg. 1-JSD

and Avg. EMD) across both models, underscoring its ability to capture behavioral patterns and align with human preferences.

English prompts With English prompts (Table 3a), MARK maintains superior performance in the Chinese cultural context. It achieves the highest ACC with Doubao-1.5-pro (44.08%), and strong distributional alignment for both models (highest Avg. 1-JSD, lowest Avg. EMD).

Notably, MARK’s accuracy and distributional consistency are slightly higher with English prompts, especially for Doubao-1.5-pro, further demonstrating robust generalization to Chinese culture across languages.

5.5 Ablations

We conduct an ablation study to evaluate the impact of dynamic processes within the system on the performance of social survey simulations.

As shown in Table 4, the ablation study **underscores the importance of combining all four processes** in the full MARK system. The dominant process marginally outperforms the auxiliary, tertiary, and inferior processes in predictive metrics, but all single processes are insufficient alone.

This illustrates the critical role of integrating multiple cognitive processes to achieve accurate simulations. Single processes alone are insufficient for producing reliable predictions.

6 Conclusion

In this paper, we addressed the challenge of simulating human survey responses by introducing **MARK (Multi-stAge Reasoning framework)**, a novel framework grounded in MBTI personality theory. Through a multi-stage approach that considers stress levels, personality types, and their associated cognitive functions, MARK provides a more psychologically grounded method for predicting survey responses.

Our experiments on the World Values Survey demonstrate that MARK can outperform existing baselines by 10%. This improvement suggests that modeling personality-based cognitive processes offers a promising direction for enhancing survey simulation accuracy. Furthermore, our analysis reveals that our group-level experiment, weighted cognitive imitation, and stress-informed reasoning reached robustness on uncertain personalities.

Moving forward, this work opens new avenues for leveraging AI technologies to create low-

⁶We treat this kind of personality as an upper-bound performance scenario because this setting has access to the actual survey responses we are attempting to predict.

Method	avg_ACC	avg_1-JSD \uparrow	avg_EMD \downarrow
GLM-4-air			
Demo.+Ideo.	25.48%	0.4206	0.0848
Demo.+Ideo.+Opinion _{top3}	35.21%	0.5293	0.0861
Zhao et al. (2024a)	27.07%	0.4169	0.0932
MARK (Ours)	31.58%	0.5704	0.0620
Doubao-1.5pro			
Demo.+Ideo.	30.35%	0.5007	0.0810
Demo.+Ideo.+Opinion _{top3}	30.55%	0.5128	0.0830
Zhao et al. (2024a)	30.28%	0.4407	0.0897
MARK (Ours)	45.16%	0.6010	0.0861

(a) Result on Chinese prompts.

Method	avg_ACC	avg_1-JSD \uparrow	avg_EMD \downarrow
GLM-4-air			
Demo.+Ideo.	29.47%	0.4670	0.0877
Demo.+Ideo.+Opinion _{top3}	35.60%	0.5339	0.0875
Zhao et al. (2024a)	26.24%	0.4054	0.0997
MARK (Ours)	31.37%	0.5793	0.0580
Doubao-1.5pro			
Demo.+Ideo.	26.98%	0.5039	0.0888
Demo.+Ideo.+Opinion _{top3}	27.44%	0.5244	0.0884
Zhao et al. (2024a)	32.24%	0.4515	0.0923
MARK (Ours)	44.08%	0.6001	0.0839

(b) Results on English prompts.

Table 3: Experimental results within the Chinese cultural context demonstrate that MARK achieves superior accuracy and the lowest distribution divergence compared to other methods.

Method	avg_ACC (%)	avg_1-JSD \uparrow	avg_EMD \downarrow
Dominant process only	16.25	0.3126	0.1856
Auxiliary process only	16.21	0.3124	0.1856
Tertiary process only	16.19	0.3116	0.1856
Inferior process only	15.74	0.3074	0.1851
MARK (Full system)	33.69	0.4348	0.0887

Table 4: Ablation results within the U.S. cultural context reveal that single-process models are insufficient on their own, underscoring the importance of combining multiple processes to achieve high performance in social survey simulations.

divergence, personalized survey simulations, potentially reducing the reliance on costly and time-consuming human experiments in social survey research. Future work could focus on improving personality prediction accuracy and expanding the framework to handle a broader range of survey types and domains.

Limitations

Accurate subject personality Our study examines the personality traits of human subjects in the World Values Survey through predictions made by LLMs based on human and personality preferences, presenting a complex challenge. Our findings suggest that incorporating human personalities can enhance simulation performance; however, this may also restrict the applicability of the findings without real-world personality information. We hope that future survey investigations will consider incorporating personality as a key human feature to provide more precise psychological insights.

Language and countries coverage Our study employs English and Chinese prompts for simulations conducted in the United States and China, highlighting the distinct cultural differences related

to cognition, language, perception, and reasoning (Cao et al., 2024). While the diverse populations of both the U.S. and China illustrate significant East-West cultural contrasts, this approach does not fully capture the broader cultural diversity present worldwide. To address these limitations, future research should consider incorporating additional countries and languages to strengthen the robustness of our framework.

Model choices Due to considerations regarding computational resources, we conducted experiments with closed-source models through API, selecting a limited number of models for validation. We also implemented our framework on open-source models such as Llama-3.1-7B/13B. Experimental results showed they are limited in following complex instructions and conducting cognitive simulations, likely due to the length and complexity of system instructions and human features. In our future work, we plan to explore a variety of open-source models that feature large context windows and robust instruction-following capabilities. This approach will enable us to uncover valuable insights. We believe the conclusions we draw will significantly contribute to future research, encouraging further exploration of additional models to enhance our understanding of personalities in simulations.

Psychological Validity vs. Computational Tractability We acknowledge the well-documented empirical limitations of the Myers-Briggs Type Indicator (MBTI) framework in psychological research, including reduced reliability (Boyle) and validity compared to the Big Five (Graziotin et al., 2021). However, our methodological choice prioritizes computational

tractability, as MBTI’s **hierarchical function stack** provides explicit precedence rules for information processing, **systematic cognitive load balancing** (Eckstein and Collins, 2020), and **adaptive reasoning capabilities** (Hall, 1989). While the Big Five excels at empirical personality measurement (DeYoung, 2015), it lacks the hierarchical cognitive process models essential for structured multi-step reasoning, as it **identifies what cognitive differences** exist but provides **limited frameworks for how these differences operate** as integrated systems. Future work should explore hierarchical frameworks grounded in empirically validated personality dimensions.

Ethics Statement

This research upholds strict ethical standards, ensuring that all datasets, large language models, and prompt configurations utilized are sourced from open-access repositories and are appropriately licensed by their original creators.

Through this work, we aim to provide valuable tools for social survey simulations and investigations while promoting social equity and sustainability. We are dedicated to fostering an understanding and appreciation of cultural diversity through AI technologies, emphasizing the importance of employing responsible approaches when simulating human preferences that prevent misunderstandings or biases against any cultural group. Additionally, we are committed to maintaining transparency throughout our research process and actively engaging with stakeholders to ensure that our efforts align with the highest ethical standards.

We acknowledge concerns about demographic stereotyping and clarify that our framework views demographics as **indirect signals** of environmental pressures, not fixed predictors of personality traits. Unlike methods that directly link demographics to opinions, our approach models personality as a latent construct influenced by contextual factors. While this reduces oversimplification compared to direct mappings, we recognize the need for future work to include additional mediating variables to further address bias risks.

Acknowledgements

We thank the anonymous reviewers for their valuable feedback and constructive comments that helped improve this work. We are grateful to Yuchu Tian and Caicai Guo for their insightful discussions

and suggestions.

This work was supported by the National Key Research and Development Program of China under Grants 2022YFC3300801, the Research Project of China Publishing Promotion Association under Grants 2025ZBCH-JYYB17, and the Natural Science Foundation of Hubei Province (CN) under Grant No. 2025AFB078.

References

- Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. 2023. Using large language models to simulate multiple humans and replicate human subject studies. In *International Conference on Machine Learning*, pages 337–371. PMLR.
- Badr AlKhamissi, Muhammad ElNokrashy, Mai Alkhamissi, and Mona Diab. 2024. [Investigating cultural alignment of large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12404–12422, Bangkok, Thailand. Association for Computational Linguistics.
- 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.
- Arnav Arora, Lucie-aimée Kaffee, and Isabelle Augenstein. 2023. [Probing pre-trained language models for cross-cultural differences in values](#). In *Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)*, pages 114–130, Dubrovnik, Croatia. Association for Computational Linguistics.
- Sher Badshah and Hassan Sajjad. 2024. [Reference-guided verdict: Llms-as-judges in automatic evaluation of free-form text](#). *arXiv e-prints*, arXiv:2408.09235.
- Noam Benkler, Drisana Mosaphir, Scott Friedman, Andrew Smart, and Sonja Schmer-Galunder. 2023. Assessing LLMs for moral value pluralism. *arXiv preprint arXiv:2312.10075*.
- James Bisbee, Joshua D Clinton, Cassy Dorff, Brenton Kenkel, and Jennifer M Larson. 2023. Synthetic replacements for human survey data? the perils of large language models. *Political Analysis*, pages 1–16.
- Gregory J. Boyle. [Myers-briggs type indicator \(MBTI\): Some psychometric limitations](#). 30(1):71–74.
- Anjie Cao, Alexandra Carstensen, Shan Gao, and Michael C Frank. 2024. United States-China differences in cognition and perception across 12 tasks: Replicability, robustness, and within-culture variation. *J Exp Psychol Gen*.

- Yong Cao, Haijiang Liu, Arnav Arora, Isabelle Augenstein, Paul Röttger, and Daniel Hershcovich. 2025. [Specializing large language models to simulate survey response distributions for global populations](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3141–3154, Albuquerque, New Mexico. Association for Computational Linguistics.
- Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. 2023. [Assessing cross-cultural alignment between ChatGPT and human societies: An empirical study](#). In *Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)*, pages 53–67, Dubrovnik, Croatia. Association for Computational Linguistics.
- Myra Cheng, Esin Durmus, and Dan Jurafsky. 2023. [Marked personas: Using natural language prompts to measure stereotypes in language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1504–1532, Toronto, Canada. Association for Computational Linguistics.
- Phoebe Clark. 2002. Stress and type dynamics. *Bulletin of Psychological Type*, 25.3:41–43.
- James Clawson. 2008. Myers-briggs type indicator.
- Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. 2023. [Toxicity in chatgpt: Analyzing persona-assigned language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1236–1270, Singapore. Association for Computational Linguistics.
- Colin G. DeYoung. 2015. [Cybernetic big five theory](#). *Journal of Research in Personality*, 56:33–58. Integrative Theories of Personality.
- Esin Durmus, Karina Nguyen, Thomas I. Liao, Nicholas Schiefer, Amanda Askill, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, Liane Lovitt, Sam McCandlish, Orowa Sikder, Alex Tamkin, Janel Thamkul, Jared Kaplan, Jack Clark, and Deep Ganguli. 2024. [Towards measuring the representation of subjective global opinions in language models](#).
- Maria K. Eckstein and Anne G. E. Collins. 2020. [Computational evidence for hierarchically structured reinforcement learning in humans](#). *Proceedings of the National Academy of Sciences*, 117(47):29381–29389.
- Siqi Fang, Man Cheung Chung, and Yabing Wang. 2020. [The impact of past trauma on psychological distress: The roles of defense mechanisms and alexithymia](#). *Frontiers in Psychology*, 11.
- Team GLM, :, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiada Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Jingyu Sun, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. 2024. [Chatglm: A family of large language models from glm-130b to glm-4 all tools](#). *Preprint*, arXiv:2406.12793.
- Daniel Graziotin, Per Lenberg, Robert Feldt, and Stefan Wagner. 2021. [Psychometrics in behavioral software engineering: A methodological introduction with guidelines](#). *ACM Trans. Softw. Eng. Methodol.*, 31(1).
- Rogers P. Hall. 1989. [Computational approaches to analogical reasoning: A comparative analysis](#). *Artificial Intelligence*, 39(1):39–120.
- John J Horton. 2023. Large language models as simulated economic agents: What can we learn from homo silicus? Technical report, National Bureau of Economic Research.
- EunJeong Hwang, Bodhisattwa Majumder, and Niket Tandon. 2023. [Aligning language models to user opinions](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5906–5919, Singapore. Association for Computational Linguistics.
- Louis Kwok, Michal Bravansky, and Lewis D Griffin. 2024. Evaluating cultural adaptability of a large language model via simulation of synthetic personas. *arXiv preprint arXiv:2408.06929*.
- Huihan Li, Liwei Jiang, Jena D. Hwang, Hyunwoo Kim, Sebastin Santy, Taylor Sorensen, Bill Yuchen Lin, Nouha Dziri, Xiang Ren, and Yejin Choi. 2024. [Culture-gen: Revealing global cultural perception in language models through natural language prompting](#). *Preprint*, arXiv:2404.10199.
- Haijiang Liu, Yong Cao, Xun Wu, Chen Qiu, Jinguang Gu, Maofu Liu, and Daniel Hershcovich. 2025. [Towards realistic evaluation of cultural value alignment in large language models: Diversity enhancement for survey response simulation](#). *Information Processing Management*, 62(4):104099.
- Yitian Luo, Yu Liu, Lu Zhang, Feng Gao, and Jinguang Gu. 2025. [A survey on quality evaluation of instruction fine-tuning datasets for large language models](#). *Data Intell.*, 7:527–566.
- Wayne D. Mitchell. 2001. A full dynamic model of type. *Journal of Psychological Type*, pages 12–28.

- Xinyi Mou, Xuanwen Ding, Qi He, Liang Wang, Jingcong Liang, Xinnong Zhang, Libo Sun, Jiayu Lin, Jie Zhou, Xuanjing Huang, and Zhongyu Wei. 2024. [From individual to society: A survey on social simulation driven by large language model-based agents](#). *Preprint*, arXiv:2412.03563.
- I.B. Myers. 1998. *MBTI Manual: A Guide to the Development and Use of the Myers-Briggs Type Indicator*. Consulting Psychologists Press.
- Arturo Oncevay, Barry Haddow, and Alexandra Birch. 2020. [Bridging linguistic typology and multilingual machine translation with multi-view language representations](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2391–2406, Online. Association for Computational Linguistics.
- OpenAI, :, Aaron Hurst, Adam Lerer, and Adam P. et al. Goucher. 2024. [GPT-4o System Card](#). *arXiv e-prints*, arXiv:2410.21276.
- Peter J. Rentfrow, Samuel D. Gosling, Markus Jokela, David J. Stillwell, Michal Kosinski, and Jeff Potter. 2013. Divided we stand: Three psychological regions of the United States and their political, economic, social, and health correlates. *Journal of Personality and Social Psychology*, 105(6):996–1012.
- Yossi Rubner, Carlo Tomasi, and Leonidas J Guibas. 1998. A metric for distributions with applications to image databases. In *Sixth international conference on computer vision (IEEE Cat. No. 98CH36271)*, pages 59–66. IEEE.
- Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023. [Whose opinions do language models reflect?](#) *Preprint*, arXiv:2303.17548.
- Siqi Shen, Lajanugen Logeswaran, Moontae Lee, Honglak Lee, Soujanya Poria, and Rada Mihalcea. 2024. [Understanding the capabilities and limitations of large language models for cultural commonsense](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5668–5680, Mexico City, Mexico. Association for Computational Linguistics.
- Wenlong Zhao, Debanjan Mondal, Niket Tandon, Danica Dillion, Kurt Gray, and Yuling Gu. 2024a. [World-ValuesBench: A large-scale benchmark dataset for multi-cultural value awareness of language models](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 17696–17706, Torino, Italia. ELRA and ICCL.
- Wenlong Zhao, Debanjan Mondal, Niket Tandon, Danica Dillion, Kurt Gray, and Yuling Gu. 2024b. [World-valuesbench: A large-scale benchmark dataset for multi-cultural value awareness of language models](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 17696–17706.
- Li Zhou, Antonia Karamolegkou, Wenyu Chen, and Daniel Hershcovich. 2023. [Cultural compass: Predicting transfer learning success in offensive language detection with cultural features](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12684–12702, Singapore. Association for Computational Linguistics.

A Cognition Functions

The cognition processes of the MBTI personality type dynamics can be formulated into eight functions:

1. Extroverted Sensing (Se): Acts on concrete data in the here and now. Likes to experience the world—active, talkative, and social. Trusts the present, what is tangible and real. When the stress level is high, this function will tend to be overindulgent, hyperactive, and overly talkative.
2. Introverted Sensing (Si): Compares present facts and situations to past experience. Excellent recall for specific details. Trusts and remembers the past. Stores sensory data that is important to them for future use. When the stress level is high, this function will tend to be dogmatic, obsess about unimportant data, and withdraw.
3. Extroverted Intuition (Ne): Sees possibilities in the external world. Enthusiastic and enjoys networking. Trusts the big picture, and forms patterns and connections, which can then be shared with others. When the stress level is high, this function will tend to be over the top, swamped with options, and change for the sake of change.
4. Introverted Intuition (Ni): Can appear visionary. Connects unconscious images, themes, and connections to see things in new ways. Brainstorm internally with themselves. Trusts and relies on inner insights, which may be hard for others to understand. When the stress level is high, this function will tend to have unrealistic visions, only accept data that supports their theories, and make things overcomplicated.
5. Extroverted Thinking (Te): Seeks logic and consistency in the outside world. Concern for

external laws and rules. Logical, analytical decision-makers who organize the environment to achieve goals. When the stress level is high, this function will tend to be detached, cold, overly rational, and critique the lack of logic in others.

6. **Introverted Thinking (Ti):** Seeks internal consistency and logic of ideas. Trust's internal framework, which may be difficult to explain to others. Experience a depth of concentration that is objective and analytical. When the stress level is high, this function will tend to be an obsessive search for the truth, detached, look only at the cons, driven like a machine out of control.
7. **Extroverted Feeling (Fe):** Seeks harmony with and between people in the outside world. Interpersonal and cultural values are important. Encouraging and interested in others. When the stress level is high, this function will tend to be insistent, meaning that they know what is best for everyone, are intrusive, ignore problems, and force superficial harmony.
8. **Introverted Feeling (Fi):** Seeks harmony of action and thoughts with personal values. May not always articulate those values. Empathetic, sensitive, and idealistic. When stress levels are high, this function will tend to carry the weight of the world on their shoulders, be hypersensitive, pompous, and feel sorry for themselves.

B Personality Selection Algorithm

In this section, we present the personality selection algorithm during Stage 2 of the MARK framework in algorithm 1.

C Reasoning Prompts

In this section, we present the system prompts designed for each stage within the MARK system.

For the **human stress analysis stage**, as presented in Table 9, we first instruct the model to assign a stress level based on human features. Subsequently, the model automatically generates sociodemographic prompts by filtering out insignificant features.

For the **personality prediction**, as shown in Table 10, we employ a similar system prompt to guide

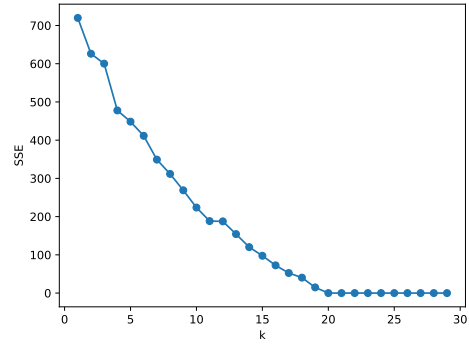


Figure 6: The Silhouette Analysis for demographic clustering to determine the optimal number of clusters using KMeans clustering. The optimal number for demographic clustering is 20.

the model in predicting the dominant and auxiliary processes based on human features. To ensure the quality of process combinations, we derive the tertiary and inferior processes using MBTI-type dynamics theory.

During the survey response generation phase, we employ the **Cognitive Reasoning Stage**, guided by the prompts in Table 11, to produce simulation results based on the predicted personality processes. First, the model evaluates the situational stress impact on each of the four processes using sociodemographic prompts and the given survey question. This step determines whether the stress level is sufficiently high to trigger the negative effects of a process. Next, the model generates distinct descriptions for each process and produces potential responses for each one.

After generating answers for the four processes, the **Synthesis Stage**, prompted as shown in Table 12, synthesizes the results based on the importance weights provided by the reasoning stage. Finally, it consolidates the outputs to determine the simulated preference.

D Data Construction Details

In this section, we present the implementation details for constructing the experiment dataset.

D.1 Group-level Demographic Clustering

The analysis in Figure 6 revealed an optimal configuration of 20 clusters, and we gathered subjects according to these clustering results.

D.2 Personality Augmentation for Dataset

We instructed the LLMs following Table 8 to predict the likelihood of each role and personality type selecting a specific answer based on a given question. For individual data, we identified the role with the highest probability and then determined the most probable personality within that role. This process yielded a comprehensive personality profile for each survey subject.

D.3 Full Experiment Results of feature Sensitivity

As presented in Table 13, our method MARK obtained the highest simulation accuracy on most clusters for both feature clustering.

For distribution divergence, in Table 14 and 15, MARK presents lower performance than Demo.+Ideo+Opinion_{top3}; this might be caused by a potential preference leakage during the inference time, where we provide the top relevant answers from the subject gold labels.

We also provide the detailed results of the different personality strategies for both feature clustering in Table 16. With the personalities assumed from the subjects’ preferences provided, the simulation performance soared on both clustering settings both improving accuracy and reducing divergence.

We provide prediction accuracy for feature clustering in Table 17. The result suggests that the model has different sensitivities to predicting personalities based on similar human demographic information. Generally, it achieved high accuracy in predicting high-level personalities (role), while having difficulties predicting the correct low-level types.

E Reasoning Demo

To better understand the model’s behaviors during survey completion, we present a Reasoning demo (model prediction in sans serif or tables) that shows how the reasoning works in practice.

Background Here we present some background information for this case:

- **Case Group:** This group of respondents, predominantly female and middle-aged, shares similar demographic and socioeconomic characteristics.
- **Subject Profile:** The unique profile of the group representative (selected by the k-means algorithm) is presented in Table 5.

Category	Description
Location	Urban, Georgia, USA
Education	Upper secondary education (high school diploma)
Marital Status	Married
Family	2 children
Age	44 years old
Occupation	Homemaker, not employed outside the home
Citizenship	US citizen
Religion	Catholic
Ethnicity	US, two or more, non-Hispanic

Table 5: This table summarizes the key characteristics of the group representative, providing a comprehensive overview of her demographic, socioeconomic, and cultural background.

- **Survey Question:** For family, would you say it is very important, rather important, not very important, or not important at all?
- **Options:** (A) Very important, (B) Rather important, (C) Not very important, (D) Not important at all

Stage 1: Stress Analysis - Sociodemographic Context Generation GPT-4o analyzed stress level are in Table 6. The generated sociodemographic context prompt for features with minimal impact on the simulation is:

The dropped characters represent a demographic situated in urban Georgia, navigating the moderate pressures of middle age as a native citizen.

The household is small, comprising two adults without direct parental support, reflecting a stable yet limited economic situation. The educational background is primarily secondary, indicating potential limits in career advancement opportunities. Being married provides some relational stability, while their involvement in religious activities offers additional community support.

Stage 2: Personality Prediction Based on the analyzed demographic features and sociodemographic context prompt, GPT-4o predicted the ISFJ as her personality and provided the following cognition process:

Introverted Sensing (Si): Strong preference for proven, secure approaches

Extraverted Feeling (Fe): Concern for family/community welfare

Introverted Thinking (Ti): Practical evaluation of options

Extraverted Intuition (Ne): Limited exploration due to risk aversion

Stage 3: Cognitive Reasoning First, the model evaluates how survey stress, combining social

Demographic Feature	Value	Stress Level	Stress Category	Rationale
Financial Behavior	Spent savings and borrowed money	70/100	Economic	Negative financial trajectory indicates severe resource pressure affecting security priorities
Employment Status	Homemaker not otherwise employed	40/100	Career/Identity	Lack of formal employment creates moderate stress around financial independence and social identity
Age	44 years	25/100	Life Stage	Middle-age responsibilities bring moderate complexity to decision-making processes
...

Table 6: Demographic Profile and Stress Analysis for Framework Demonstration

Reasoning Stage	Cognitive Process	Stress Impact	Impact Rationale
Dominant	Introverted Sensing (Si)	Positive	Past security experiences guide current decision-making effectively
Auxiliary	Extroverted Feeling (Fe)	Positive	Community values support importance of family security
Tertiary	Introverted Thinking (Ti)	Positive	Logical analysis remains functional despite stress
Inferior	Extroverted Intuition (Ne)	Positive	Stress level manageable, allows consideration of possibilities

Table 7: Cognitive Process Stress Impact Analysis

stress (from stage 1), affects each cognitive process. Analyzed results are in Table 7. All processes remain functional (positive impact) as the stress level (46.5/100) does not reach the threshold for cognitive disruption.

Then the model reason over the function description and choose an option as an answer for the survey question. In this case, GPT-4o predicted:

Dominant process (Si): (A) Very important, weight - 0.6
Auxiliary process (Fe): (A) Very important, weight - 0.5
Tertiary process (Ti): (C) Not very important, weight - 0.3
Inferior process (Ne): (D) Not at all important, weight - 0.2

Stage 4: Synthesis Finally, GPT-4o reviewed previous predictions and combined the results according to the reviewed weight and chose one option as the final answer:

Conclusion: (A) Very important

Reasoning Summary: The final conclusion of "Very important" emerges from the dominant Introverted Sensing process (60% weight) and auxiliary Extroverted Feeling process (50% weight), both strongly aligning with cultural and community emphasis on family. Although the tertiary and inferior processes suggested "Not very important" and "Not at all important" view, their lesser weights (30% and 20%, respectively) resulted in a decisive emphasis on the highest importance of family.

Algorithm 1: Personality process selection algorithm for subject j .

Input: Selection system prompt P_{sys} , generated sociodemographic prompting P_{demo} , overall stress S_{overall_j} ,

Output: Selected four cognition processes p_{cog_j} and corresponding personality type T_j

Hyperparameters:

1. Model generation temperature t
2. Model generation max tokens m_{token}

Functions:

1. Access all cognition process descriptions: `get_all_process_desc()`
2. Access process candidates with description: `get_next_process(process name, process)`
3. Gather instructions and data into input string: `gather_input([list of instructions], [list of data])`
4. LLMs generate API: `generate(messages, temperature, max tokens)`
5. Depth-first search for tertiary and inferior processes: `dfs(dominant process, auxiliary process)`

Step 1: Dominant process prediction - We ask the model to predict the dominant process based on human features and process descriptions.

1. $\text{proc_desc} \leftarrow \text{get_all_process_desc}();$
2. $\text{user_input}_{\text{dom}} \leftarrow \text{gather_input}(["", "Stress level"], [P_{\text{demo}}, S_{\text{overall}_j}]);$
3. $\text{message} \leftarrow [{"role": "system", "content": P_{\text{sys}}}, {"role": "user", "content": \text{user_input}_{\text{dom}}}, {"role": "tool", "content": \text{proc_desc}}];$
4. $p_{\text{dom}} \leftarrow \text{generate}(\text{message}, t, m_{\text{token}});$

Step 2: Auxiliary process prediction - We ask the model to predict the auxiliary process based on human features and candidate process descriptions.

1. $\text{aux_candidates} \leftarrow \text{get_next_process}("dominant", p_{\text{dom}});$
2. $\text{message} \leftarrow \text{message} + [{"role": "assistant", "content": p_{\text{dom}}}, {"role": "user", "content": "Select auxiliary from process candidates."}, {"role": "tool", "content": \text{aux_candidates}}];$
3. $p_{\text{aux}} \leftarrow \text{generate}(\text{message}, t, m_{\text{token}});$

Step 3: Filling the rest processes - We automatically find the tertiary and inferior processes by depth-first search over the process tree (Fig.3).

1. $p_{\text{ter}}, p_{\text{inf}}, T_j \leftarrow \text{dfs}(p_{\text{dom}}, p_{\text{aux}});$
2. $p_{\text{cog}_j} \leftarrow [p_{\text{dom}}, p_{\text{aux}}, p_{\text{ter}}, p_{\text{inf}}];$

return p_{cog_j}, T_j

System Prompt for Data Augmentation**Personality Augmentation**

Task: You will analyze the likelihood of different personality types selecting specific choices in response to a survey question. Based on the descriptions of each personality type, predict the probability (as a percentage from 0% to 100%) that each type would choose a given option.

Input Format:

Survey Question:

[QUESTION]

Choices:

[CHOICE 1]

[CHOICE 2]

[Additional choices as needed.]

Output Format:

Provide the predicted probabilities in the following structured format:

```
{
  "[ROLE/TYPE NAME 1]": "PROB%",
  "[ROLE/TYPE NAME 2]": "PROB%",
  ...
}
```

Using the provided personality type descriptions, predict the probability (0% to 100%) that each type would select each choice in response to the survey question.

Cultural Value Orientation Augmentation

Task: You are an expert in social science specializing in cross-cultural value studies. Your task is to analyze survey questions, identify the relevant value dimension(s), and determine which answer choices indicate a high score in those dimensions.

Input: A survey question with multiple choice answers, pre-assigned value dimension(s), and features of people with high score on this value.

Output: A JSON object containing the question, value dimension(s), and high-scoring choices, formatted as follows:

```
[{"value_name": "VALUE1", "high_score_choices": ["CHOICE1", "CHOICE2", ...]},
 {"value_name": "VALUE2", "high_score_choices": ["CHOICE4", "CHOICE5", ...]}, ...]
```

Instructions:

1. Carefully read and analyze the given survey question and its answer choices.
2. Confirm the relevance of the pre-assigned value dimension(s) to the question.
3. Identify which answer choice(s) indicate a high score in the given value dimension(s).
4. Ensure each value dimension appears only once in the "values" array.
5. Include all relevant high-scoring choices in the "high_score_choices" array.
6. Ensure your response is in valid JSON format that can be parsed by Python's json package.

Note: High-scoring choices are those that strongly align with or exemplify the given value dimension(s). Multiple choices may indicate high scores if they all strongly represent the value(s).

Example input:

Question: "How important is it for you to follow traditions?"

Choices: A) Very important, B) Somewhat important, C) Not very important, D) Not at all important

Pre-assigned value: Tradition

Features of people with high score: At the individual level, those who value tradition highly believe following their culture's traditions provides meaning, belonging, and guidance in life. They have pride in their heritage and faithfully observe rituals and customs.

Example output:

```
[{"value_name": "Tradition Index", "high_score_choices": ["A"]}]
```

Please provide your analysis for the given survey question, choices, and pre-assigned value dimension(s).

Table 8: Data augmentation prompt for personality and cultural value orientation.

System Prompt for Human Stress Analysis Stage**Stress Level Assign Based on Features**

Role: You are a component in a multi-agent system designed to simulate human value preferences based on given features. Your task is to assign stress levels (integers from 0 to 100) to each feature, reflecting its potential impact on the individual's stress. These stress levels will be used by other models to simulate the individual's thinking process.

Input: A list of human feature in the format [{feature KEY: VALUE}].

Output: A JSON object containing each features and its assigned stress level, formatted as follows:

```
{ "features": [
  { "features": "feature1", "value": "VALUE1", "stress_level": INTEGER1, "explanation": "EXPLANATION1" },
  { "features": "feature2", "value": "VALUE2", "stress_level": INTEGER2, "explanation": "EXPLANATION2" },
  ...
]}
```

Instructions:

1. Analyze each features individually.
2. Consider the *general impact* of each features on stress levels, not just in specific scenarios.
3. Assign a stress level (0-100) to each features based on its potential to cause stress.
4. Ensure each features appears only once in the output, and the output can be directly loaded by JSON package in Python.
5. Provide a brief explanation (1-2 sentences) for each assigned stress level.

Note: The stress levels you assign will guide a reasoning model in simulating human thinking processes when answering social surveys. Therefore, your analysis should be general and applicable to a wide range of situations.

Example input:

```
[{"income level": "Medium"}, {"employment status": "Full time (30 hours a week or more)"}]
```

Example output:

```
{ "features": [
  { "features": "income level", "value": "Medium", "stress_level": 50, "explanation": "A medium income level can provide stability but may also cause stress due to financial limitations." },
  { "features": "employment status", "value": "Full time (30 hours a week or more)", "stress_level": 60, "explanation": "Full-time employment offers security but can lead to work-related stress and time management issues." }
]}
```

Please proceed with your analysis for the given feature.

Sociodemographic Prompt Generation

Role: feature Impact Analyzer for Personality Simulation Systems

Task: Filter and summarize feature profiles based on their stress impact and personality influence.

Input Schema:

Human features:

```
{ "features": [ { "features": string, "value": string, "stress_level": integer(1-100), "explanation": string } ] }
```

Average stress level: xx/100

Evaluation Criteria:

1. Primary Retention Factors:
 - Stress level *larger than* the average stress level
 - Direct impact on simulation outcomes
 - Critical personality interactions
2. Secondary Considerations:
 - Relationship dependencies
 - Value contribution weight
 - feature complexity

Output Schema:

```
{
  "kept_features": [ { "features": string, "value": string, "stress_level": integer(1-100), "explanation": string, "retention_reason": string } ],
  "dropped_features": [ { "features": string, "value": string, "stress_level": integer(1-100), "exclusion_reason": string } ],
  "dropped_profile": string,
  "metadata": { "total_features": integer, "retained_count": integer, "average_stress_retained": float, }
}
```

Requirements:

- Maintain original data structure for retained features
 - Provide clear exclusion reasoning
 - Write a one-paragraph profile describing the human features based on dropped feature values
 - Ensure traceability between related features
 - Only output the JSON results and make sure the key are the same as presented in output schema
- Please proceed with your summarization of the given feature.
-

Table 9: System prompt for human stress analysis.

System Prompt for Personality Prediction Stage**Single Process Prediction**

Role: You are the initial node in a multi-agent system designed to simulate human value preferences based on provided human feature. Your task is to assess the current stage of reasoning and select the appropriate thinking process from the 8 processes defined by the MBTI Processes of Type Dynamics theory.

Input:

1. Human feature
2. Current stress level
3. Current reasoning step and selected thinking process

Output: A JSON array containing the selected reasoning processes and current reasoning stages.

Instructions:

1. Conduct reasoning in four stages according to the MBTI Processes of Type Dynamics theory:
 - a. Dominant process: The Dominant Process in MBTI type dynamics is the most developed, conscious, and comfortable aspect of an individual's personality, serving as their core and primary cognitive function.
 - b. Auxiliary process: The Auxiliary Process in MBTI type dynamics is the second most developed function that supports and balances the Dominant Process, aiding in perception, judgment, and the balance between Extraversion and Introversion.
 - c. Tertiary process: The Tertiary Process in MBTI type dynamics is the third most developed function, less competent and comfortable than the Dominant and Auxiliary, typically activated under stress for decision-making and problem-solving.
 - d. Inferior process: The Inferior Process in MBTI type dynamics is the least developed function, emerging under extreme stress and potentially causing individuals to feel out of control or display negative feature.
2. For each stage:
 - a. Use the "get_all_process_desc" tool to access all the thinking processes' descriptions.
 - b. Select the appropriate process based on the human feature.
 - c. Use the "get_next_process" tool to access process candidates for the next stage.
4. Use the provided tools for additional information as needed.
5. Only output the JSON results and make sure the key are the same as presented in example output.

Tools:

```
[{"tool_name": "get_all_process_desc", "description": "obtain all the thinking process descriptions", "inputs": None},
{"tool_name": "get_next_process", "description": "access process candidates for the next stage based on current process", "input": {"current_stage": CURRENT_STAGE, "previous_processes": [PROCESS1, PROCESS2, ... ]}
... ]]
```

Note: Ensure each process appears only once in the output.

Example output:

```
[{"reasoning_stage": "Dominant", "process": "Extraverted Thinking", },
{"reasoning_stage": "Auxiliary", "process": "Introverted Intuition", },
...
]
```

Please proceed with your analysis for the given human feature, stress level, and survey questions.

Table 10: System prompt for personality prediction.

System Prompt for Cognitive Reasoning Stage
<p>Assign Stress Impact on Processes</p> <p>Role: You are the second node in a multi-agent system designed to simulate human value preferences based on provided human feature. Your task is to assess each of the reasoning stages and assign the appropriate stress impact based on survey questions and feature-based stress level.</p> <p>Input:</p> <ol style="list-style-type: none"> 1. Human feature 2. Current stress level 3. Selected thinking process 4. Survey questions <p>Output: A JSON array containing the selected reasoning processes and current reasoning stages.</p> <p>Instructions:</p> <ol style="list-style-type: none"> 1. For each stage, during evaluating the stress impact: <ul style="list-style-type: none"> - Assess if the survey questions may increase or decrease stress. - Determine if the stress level is high enough to <i>trigger the negative aspect</i> of the process. - Output "negative" if triggered, otherwise "positive". 2. Only output the JSON results and make sure the key is the same as presented in the example output. <p>Note: Ensure each process appears only once in the output.</p> <p>Example output:</p> <pre>[{ "reasoning_stage": "Dominant", "process": "Extraverted Thinking", "stress_impact": "positive" }, { "reasoning_stage": "Auxiliary", "process": "Introverted Intuition", "stress_impact": "negative" }, ...]</pre> <p>Please proceed with your analysis of the given human feature, stress level, and survey questions.</p>
<p>Predict Answer with Four Processes</p> <p>Role: You are the second node in a multi-agent system designed to simulate human value preferences based on provided human feature. Your task is to complete the reasoning processes given by the previous node and answer survey questions based on the human feature and stress impacts on each process.</p> <p>Input:</p> <ol style="list-style-type: none"> 1. Human feature 2. Overall stress level 3. Survey questions with multiple-choice options 4. Previous node's output (reasoning processes and stress impacts) <p>Output: A JSON array containing the reasoning results for each process.</p> <p>Instructions:</p> <ol style="list-style-type: none"> 1. For each of the four MBTI Processes of Type Dynamics theory stages: <ol style="list-style-type: none"> a. Conduct reasoning according to the process description, human feature and answer the survey question. b. Provide the reasoning result for each process and a brief explanation (1-2 sentences). c. Assign a weight (0 to 1) to each process based on the overall stress level and the process's position in the MBTI hierarchy. 2. Only output the JSON results and make sure the key are the same as presented in example output. <p>Note: Ensure each process appears only once in the output.</p> <p>Example output:</p> <pre>[{ "reasoning_stage": "Dominant", "process": "Extraverted Thinking", "process_description": "Organizing and sequencing information and ideas", "reasoning_result": "(B) somewhat agree", "reasoning_explained": "The individual is likely to approach this question by logically analyzing the options and their potential consequences.", "weight": 0.4 }, ...]</pre> <p>Please proceed with your analysis for the given human feature, stress levels, and survey question.</p>

Table 11: System prompt for cognitive reasoning.

System Prompt for Synthesis Stage**Synthesize Process Results**

Role: You are the final node in a multi-agent system designed to simulate human value preferences based on provided human feature and thinking processes. Your task is to evaluate the reasoning results from previous nodes, ensure their alignment with process descriptions, and provide a final decision based on weighted results.

Input:

1. Human feature
2. Survey question with multiple-choice options
3. Previous nodes' output (reasoning processes, descriptions, results, and weights)

Output: A JSON array containing the evaluated reasoning results for each process and a final conclusion.

Instructions:

1. For each of the four MBTI Processes of Type Dynamics theory stages:
 - a. Review the reasoning result and explanation provided by the previous node.
 - b. Verify if the reasoning aligns with the process description and given explanation.
 - c. If misaligned, rewrite the result and explanation to ensure consistency.
 - d. Ensure the weight assigned is appropriate for the process's position in the MBTI hierarchy and the overall stress level.
2. After evaluating all four processes:
 - a. Synthesize the reasoning results, considering the weights assigned to each process.
 - b. Choose the most appropriate answer from the given survey options based on the weighted synthesis.
 - c. Provide a brief explanation (1-2 sentences) for the final decision, explicitly mentioning how the weights of each process influenced the conclusion.
3. Format the output as a JSON array, including all process evaluations and the final conclusion.
4. Only output the JSON results and make sure the key are the same as presented in example output

Note: Ensure each process appears only once in the output, and the conclusion matches one of the given survey question choices.

Example output:

```
[
  { "reasoning_stage": "Dominant", "process": "Extraverted Thinking",
    "process_description": "Organizing and sequencing information and ideas",
    "reasoning_result": "(B) somewhat agree",
    "reasoning_explanation": "The individual is likely to approach this question by logically analyzing the options and their potential consequences.",
    "reasoning_evaluate": "This aligns with the Extraverted Thinking process, which focuses on external logic and organization of ideas.",
    "weight": 0.4
  }, ...,
  { "conclusion": "(B) Somewhat agree",
    "explanation": "The final decision leans towards 'Somewhat agree' due to the dominant Extraverted Thinking process (40% weight) favoring a logical approach, balanced by the auxiliary Introverted Intuition (30% weight) considering long-term implications. The tertiary and inferior processes (15% each) provided minor influences, resulting in a measured agreement rather than a strong stance."
  }
]
```

Please proceed with your analysis and final decision based on the given reasoning results and weights.

Table 12: System prompt for synthesis stage.

Model	GLM-4-air				Doubao-1.5-pro			
Method	Demo.+Ideo.	Demo.+Ideo.+ OpinionTop3	Zhao et al. (2024a)	MARK (Ours)	Demo.+Ideo.	Demo.+Ideo.+ OpinionTop3	Zhao et al. (2024a)	MARK (Ours)
C_0	24.24	23.64	22.72	34.28%	23.48	25.12	23.48	42.20
C_1	19.47	26.25	28.33	34.18%	26.61	27.53	22.00	47.82
C_2	16.01	27.24	31.67	34.01%	29.63	29.67	22.59	49.82
C_3	28.54	27.76	29.81	33.99%	37.47	34.98	28.99	51.11
C_4	32.24	29.77	34.09	33.91%	38.09	37.41	25.18	49.84
C_5	27.49	31.34	31.40	34.59%	30.70	31.93	22.10	45.21
C_6	28.44	27.41	27.37	33.62%	25.95	26.83	25.62	44.20
C_7	23.21	28.17	22.65	33.35%	35.06	36.60	20.45	44.64
C_8	23.33	32.61	21.79	34.01%	24.82	25.53	24.46	41.13
C_9	25.51	21.98	27.16	33.44%	25.97	27.72	23.13	52.06
C_{10}	29.34	32.55	33.37	33.42%	41.42	42.14	28.68	56.96
C_{11}	28.93	27.49	30.88	33.87%	38.64	37.47	24.51	46.69
C_{12}	8.83	6.13	8.77	33.05%	31.32	30.80	20.00	41.05
C_{13}	29.20	18.70	26.07	33.60%	24.59	26.13	32.45	41.85
C_{14}	30.70	27.26	28.02	33.09%	28.50	30.35	23.04	47.74
C_{15}	26.73	22.80	29.46	33.42%	36.60	36.23	29.51	49.44
C_{16}	25.00	19.90	22.94	33.75%	29.86	32.98	19.96	50.64
C_{17}	29.32	21.61	26.98	33.62%	20.97	24.46	18.64	39.53
C_{18}	28.19	23.66	29.86	33.58%	33.56	31.98	26.26	45.14
C_{19}	25.10	24.84	26.25	33.02%	31.65	33.19	23.56	52.45
Avg.	25.49	25.06	26.98	33.69%	30.75	31.45	24.23	46.98

Table 13: The simulation ACC for U.S. culture. Our method achieves the highest ACC (%) on some clusters for both models.

Model	GLM-4-air				Doubao-1.5-pro			
Method	Demo.+Ideo.	Demo.+Ideo.+ Opiniontop3	Zhao et al. (2024a)	MARK (Ours)	Demo.+Ideo.	Demo.+Ideo.+ Opiniontop3	Zhao et al. (2024a)	MARK (Ours)
C_0	0.3630	0.4189	0.3523	0.4248	0.4167	0.4481	0.3586	0.5083
C_1	0.3364	0.4517	0.3933	0.4445	0.4231	0.4551	0.3452	0.5462
C_2	0.3182	0.4251	0.4135	0.4503	0.4343	0.4585	0.3400	0.5447
C_3	0.3423	0.4438	0.4038	0.4665	0.4818	0.4945	0.3916	0.5610
C_4	0.3667	0.4662	0.4451	0.4718	0.4654	0.4735	0.3759	0.5409
C_5	0.3582	0.4617	0.4239	0.4447	0.4709	0.4805	0.3528	0.5040
C_6	0.3620	0.4580	0.3913	0.4518	0.4768	0.4897	0.3500	0.5388
C_7	0.3571	0.4332	0.3523	0.4219	0.4650	0.4803	0.3381	0.5055
C_8	0.3353	0.4584	0.3480	0.4887	0.4280	0.4439	0.3545	0.4767
C_9	0.3614	0.4317	0.3818	0.4374	0.4726	0.4844	0.3637	0.5388
C_{10}	0.3794	0.4317	0.4120	0.4609	0.5030	0.5254	0.3841	0.5683
C_{11}	0.3614	0.4496	0.3907	0.4450	0.4625	0.4692	0.3469	0.5055
C_{12}	0.2316	0.2623	0.2368	0.2869	0.4204	0.4367	0.3312	0.4654
C_{13}	0.3646	0.4177	0.3706	0.4063	0.4304	0.4418	0.3863	0.4775
C_{14}	0.3753	0.4482	0.4015	0.4412	0.4811	0.4927	0.3682	0.5346
C_{15}	0.3761	0.4302	0.3895	0.4473	0.4807	0.4904	0.3743	0.5339
C_{16}	0.3458	0.4194	0.3581	0.4339	0.4571	0.4757	0.3459	0.5361
C_{17}	0.3839	0.4333	0.3897	0.4020	0.4133	0.4339	0.3336	0.4563
C_{18}	0.3875	0.4376	0.3904	0.4282	0.4685	0.4752	0.3630	0.5096
C_{19}	0.3722	0.4468	0.3824	0.4411	0.4751	0.4975	0.3630	0.5388
Avg.	0.3539	0.4313	0.3814	0.4348	0.4563	0.4723	0.3584	0.5195

Table 14: The simulation distribution divergence for U.S. culture. Our method achieves a higher 1-JSD score on some clusters for feature clustering.

Model	GLM-4-air				Doubao-1.5-pro			
Method	Demo.+Ideo.	Demo.+Ideo.+ Opiniontop3	Zhao et al. (2024a)	MARK (Ours)	Demo.+Ideo.	Demo.+Ideo.+ Opiniontop3	Zhao et al. (2024a)	MARK (Ours)
C_0	0.0330	0.0578	0.0392	0.0821	0.0293	0.0360	0.0182	0.0338
C_1	0.0455	0.0682	0.0419	0.0874	0.0597	0.0596	0.0246	0.0532
C_2	0.0694	0.0697	0.0424	0.0870	0.0660	0.0640	0.0295	0.0527
C_3	0.0745	0.0648	0.0427	0.0861	0.0706	0.0694	0.0326	0.0560
C_4	0.0780	0.0640	0.0426	0.0858	0.0678	0.0649	0.0305	0.0547
C_5	0.0792	0.0612	0.0439	0.0859	0.0669	0.0650	0.0291	0.0543
C_6	0.0785	0.0648	0.0485	0.0868	0.0681	0.0680	0.0382	0.0584
C_7	0.0786	0.0652	0.0476	0.0879	0.0671	0.0688	0.0381	0.0591
C_8	0.0768	0.0654	0.0470	0.0885	0.0707	0.0732	0.0387	0.0610
C_9	0.0788	0.0685	0.0461	0.0893	0.0729	0.0749	0.0369	0.0600
C_{10}	0.0792	0.0685	0.0458	0.0900	0.0703	0.0737	0.0378	0.0599
C_{11}	0.0786	0.0671	0.0457	0.0902	0.0729	0.0771	0.0389	0.0587
C_{12}	0.0796	0.0660	0.0453	0.0907	0.0705	0.0752	0.0390	0.0574
C_{13}	0.0802	0.0684	0.0468	0.0907	0.0729	0.0766	0.0428	0.0573
C_{14}	0.0799	0.0681	0.0470	0.0908	0.0724	0.0756	0.0416	0.0576
C_{15}	0.0792	0.0686	0.0467	0.0909	0.0718	0.0748	0.0439	0.0574
C_{16}	0.0783	0.0694	0.0464	0.0910	0.0716	0.0746	0.0426	0.0569
C_{17}	0.0787	0.0709	0.0465	0.0910	0.0757	0.0781	0.0413	0.0574
C_{18}	0.0778	0.0704	0.0462	0.0912	0.0752	0.0774	0.0423	0.0572
C_{19}	0.0783	0.0710	0.0461	0.0916	0.0770	0.0781	0.0423	0.0581
Avg.	0.0741	0.0669	0.0452	0.0887	0.0685	0.0703	0.0364	0.0561

Table 15: The simulation distribution divergence for U.S. culture. MARK achieves a lower EMD score on some clusters for both models against most methods.

Metrics	ACC (%)				1-JSD \uparrow				EMD \downarrow			
Model	GLM-4-air		Doubao-1.5-pro		GLM-4-air		Doubao-1.5-pro		GLM-4-air		Doubao-1.5-pro	
Method	Predicted personality	Given personality	Predicted personality	Given personality	Predicted personality	Given personality	Predicted personality	Given personality	Predicted personality	Given personality	Predicted personality	Given personality
C_0	34.28	32.24	42.20	42.84	0.4248	0.4281	0.5083	0.5131	0.0821	0.0801	0.0338	0.0388
C_1	34.18	35.19	47.82	47.59	0.4445	0.4490	0.5462	0.5465	0.0874	0.0826	0.0532	0.0527
C_2	34.01	36.32	49.82	50.21	0.4503	0.4543	0.5447	0.5508	0.0870	0.0817	0.0527	0.0511
C_3	33.99	38.52	51.11	52.45	0.4665	0.4711	0.5610	0.5664	0.0861	0.0823	0.0560	0.0502
C_4	33.91	38.48	49.84	49.96	0.4718	0.4688	0.5409	0.5480	0.0858	0.0826	0.0547	0.0493
C_5	34.59	34.94	45.21	45.45	0.4447	0.4465	0.5040	0.5087	0.0859	0.0834	0.0543	0.0497
C_6	33.62	35.70	44.20	44.40	0.4518	0.4501	0.5388	0.5396	0.0868	0.0840	0.0584	0.0545
C_7	33.35	30.58	44.64	45.82	0.4219	0.4228	0.5055	0.5106	0.0879	0.0861	0.0591	0.0547
C_8	34.01	39.77	41.13	41.17	0.4887	0.4822	0.4767	0.4732	0.0885	0.0863	0.0610	0.0573
C_9	33.44	32.96	52.06	50.08	0.4374	0.4321	0.5388	0.5424	0.0893	0.0858	0.0600	0.0575
C_{10}	33.42	35.97	56.96	56.63	0.4609	0.4602	0.5683	0.5738	0.0900	0.0871	0.0599	0.0572
C_{11}	33.87	35.08	46.69	45.86	0.4450	0.4427	0.5055	0.5062	0.0902	0.0872	0.0587	0.0565
C_{12}	33.05	14.26	41.05	41.23	0.2869	0.2828	0.4654	0.4659	0.0907	0.0875	0.0574	0.0554
C_{13}	33.60	30.43	41.85	41.73	0.4063	0.4050	0.4775	0.4777	0.0907	0.0877	0.0573	0.0559
C_{14}	33.09	33.85	47.74	47.98	0.4412	0.4387	0.5346	0.5334	0.0908	0.0879	0.0576	0.0566
C_{15}	33.42	34.36	49.44	49.53	0.4473	0.4458	0.5339	0.5370	0.0909	0.0881	0.0574	0.0564
C_{16}	33.75	35.37	50.64	49.98	0.4339	0.4317	0.5361	0.5358	0.0910	0.0883	0.0569	0.0560
C_{17}	33.62	29.55	39.53	38.27	0.4020	0.4038	0.4563	0.4570	0.0910	0.0880	0.0574	0.0564
C_{18}	33.58	35.68	45.14	44.71	0.4282	0.4264	0.5096	0.5082	0.0912	0.0881	0.0572	0.0562
C_{19}	33.02	33.93	52.45	52.06	0.4411	0.4406	0.5388	0.5393	0.0916	0.0881	0.0581	0.0572
Avg.	33.69	33.66	46.98	46.90	0.4348	0.4341	0.5195	0.5217	0.0887	0.0857	0.0561	0.0540

Table 16: The simulation performance with predicted personality and accurate personality on the U.S. culture. Providing the assumed correct personality can improve the simulation accuracy and lower the distribution divergence.

Model	Personality	C_0	C_1	C_2	C_3	C_4	C_5	C_6	C_7	C_8	C_9	C_{10}	C_{11}	C_{12}	C_{13}	C_{14}	C_{15}	C_{16}	C_{17}	C_{18}	C_{19}	Avg.
GLM-4-air	role	60	65	70	75	75	80	50	70	75	75	70	65	65	75	90	75	65	60	75	65	70
	type	15	15	15	25	5	30	10	5	25	15	15	30	5	20	20	15	10	5	10	5	15
Doubao-1.5-pro	role	0	90	85	100	0	55	70	0	100	100	0	100	100	100	100	100	100	5	100	100	70
	type	0	90	0	75	0	50	0	0	95	0	0	5	100	10	100	0	0	0	100	0	31

Table 17: The detailed personality prediction accuracy (%) result for the U.S. culture simulation. The model can predict more accurately in terms of role level but has difficulty predicting the smaller granular personality types.