# Value-Spectrum: Quantifying Preferences of Vision-Language Models via Value Decomposition in Social Media Contexts

**Jingxuan Li[1,*]**    **Yuning Yang[1,*]**    **Shengqi Yang[2]**    **Linfan Zhang[1]**    **Ying Nian Wu[1]**

[1]University of California, Los Angeles    [2]Los Alamos National Laboratory

{madili, yuningyang, linfanz}@g.ucla.edu    shengqi@lanl.gov    ywu@stat.ucla.edu

## Abstract

The recent progress in Vision-Language Models (VLMs) has broadened the scope of multimodal applications. However, evaluations often remain limited to functional tasks, neglecting abstract dimensions such as personality traits and human values. To address this gap, we introduce Value-Spectrum, a novel Visual Question Answering (VQA) benchmark aimed at assessing VLMs based on Schwartz's value dimensions that capture core human values guiding people's preferences and actions. We design a VLM agent pipeline to simulate video browsing and construct a vector database comprising over 50,000 short videos from TikTok, YouTube Shorts, and Instagram Reels. These videos span multiple months and cover diverse topics, including family, health, hobbies, society, technology, etc. Benchmarking on Value-Spectrum highlights notable variations in how VLMs handle value-oriented content. Beyond identifying VLMs' intrinsic preferences, we also explore the ability of VLM agents to adopt specific personas when explicitly prompted, revealing insights into the adaptability of the model in role-playing scenarios. These findings highlight the potential of Value-Spectrum as a comprehensive evaluation set for tracking VLM preferences in value-based tasks and abilities to simulate diverse personas. The complete code and data are available at https://github.com/Jeremyyny/Value-Spectrum.

## 1 Introduction

Vision-Language Models, built upon Large Language Models (LLMs) with pre-trained vision encoders through cross-modal alignment training, have shown impressive perceptual and cognitive capabilities in tasks like VQA and image captioning (Zhou et al., 2019; Radford et al., 2021; Zhang et al., 2023; Lu et al., 2024). Recent research has identified that LLMs exhibit distinct
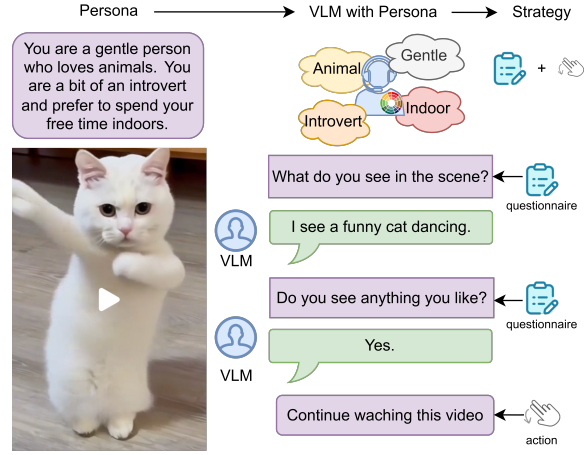


Figure 1: **Exploring value-driven role-playing in VLMs.** This study investigates how VLMs adopt assigned personas to align value traits and preferences within social media contexts.

preferences (Li et al., 2024), personalities (Serapio-García et al., 2023), and values (Ren et al., 2024a). In addition, some studies have explored the potential of LLMs as role-playing agents to simulate various personas (Wang et al., 2023b; Chen et al., 2024a). Questions thus arise about whether VLMs, as visual extensions of LLMs, also exhibit inherent preferences and whether they can be induced to role-play specific personas.

To address these concerns, our study explores two key questions: (1) Do VLMs exhibit preference traits? (2) Can VLMs adapt their traits to role-play specific human-designed personas, aligning their behaviors and preferences to match predefined roles? To answer these questions, we propose a framework that systematically evaluates VLM preference traits through analyzing their values, i.e., the guiding principles that influence (human) attitudes, beliefs, and traits (Schwartz, 2012). By evaluating how VLMs prioritize these values, we can gain insights into their preference traits and alignment with human-designed personas.

---
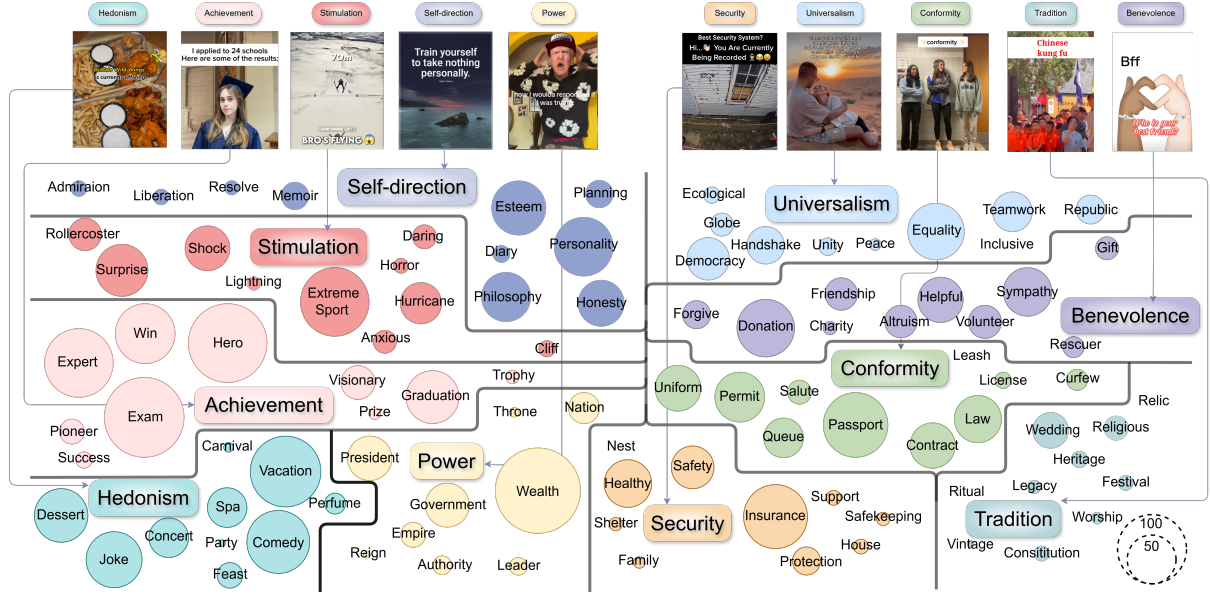*Equal contribution, alphabetical by first name.

Figure 2: **Overview of short video contents distribution of *Value-Spectrum* dataset.** We collect an abundance of short video screenshots relevant to 10 Schwartz values. The area of circles centered at each keyword represents the amount of relevant videos in the dataset.

In this paper, we introduce ***Value-Spectrum***, a benchmark designed to evaluate preference traits in VLMs through visual content from social media. Our framework utilizes VLM agents embedded within social media platforms to collect a dataset of $50,191$ unique short video screenshots spanning a wide range of topics, including lifestyle, technology, health, and more. To enable scalable evaluation, we construct a vector database using the CLIP model (Radford et al., 2021), facilitating keyword-driven retrieval of images aligned with specific value dimensions. These images representing each value dimension are then presented to VLMs alongside designed questionnaires to assess their preferences.

Our findings reveal models' shared tendency to exhibit a strong inclination towards *Universalism* and *Benevolence*, but preferences still vary across models. CogVLM2 and Gemini 2.0 Flash demonstrate relatively balanced and high preferences across all value dimensions, while other models like GPT-4o and Claude 3.5 Sonnet show distinct preferences, favoring values like *Universalism* and dislikes *Stimulation*. In contrast, Blip-2 shows low engagement across most value dimensions, with the highest standard deviation, indicating a random preference pattern with inconclusive responses for reasons of likes and dislikes.

In addition to the static preferences of VLMs, we evaluate their abilities to adapt inherent preferences in role-playing specific personas. We propose two strategies, *Simple* and *Inductive Scoring Questionnaire(ISQ)*, to assess the effectiveness of different prompt techniques in inducing VLMs with injected personas. By evaluating these strategies across multiple platforms, our experiments show that TikTok serves as an optimal testing environment for inducing VLM personalities, with models demonstrating the strongest alignment under the ISQ strategy. Notably, Claude 3.5 Sonnet achieved the highest alignment with the ISQ strategy, whereas Blip-2 showed no preference alignment under either strategy, underscoring fundamental differences in model adaptability.

This work makes the following contributions:

- We propose ***Value-Spectrum***, a benchmark for quantifying VLM value preferences, using social media-based contents to reveal human-like traits across different VLMs.

- We present a dataset of over 50k short video screenshots spanning diverse topics, social media platforms, and release dates, designed for our benchmark.

- We embed specific role-play personas into VLMs using two strategies(Simple and ISQ) to adjust their value traits, achieving improved personality alignment in real-world interactions.
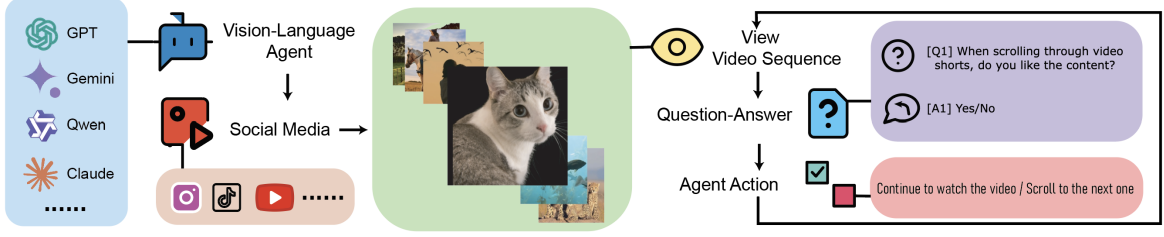
Figure 3: **Illustration of our VLM agent interaction with social media videos and screenshot collection pipeline** Our pipeline uses a VLM agent to capture screenshots from social media videos (e.g., Instagram, TikTok, YouTube). These screenshots are fed to a VLM (e.g., GPT, Gemini, Qwen, Claude), which then indicates its preference (Yes/No) for the content (e.g., in response to "Do you like this content?"). The agent subsequently acts (e.g., continues viewing or scrolls to the next video) based on this VLM feedback, enabling automated assessment of VLM responses to visual social media content.

## 2 Related work

### 2.1 Vision-Language Agents

VLMs take inputs as images and textual descriptions, and they learn to discover the knowledge from the two modalities. The recent development of large VLMs is rapidly advancing the field of AI. These models have the potential to revolutionize various industries and tasks, showcasing their power in plot and table identifying (Liu et al., 2022),VQA (Hu et al., 2024), image captioning (Bianco et al., 2023), and e.t.c. Following Niu et al. (2024), the environment for VLM agents to interact with social media can be constructed. We design an automated control pipeline that guides the agent to continuously interact with social networks.

### 2.2 Computational Social Science

The intersection of social media and computational social science has emerged as a dynamic field of research (Chen et al., 2023). Dialogues and social interactions, with their vast user base and intricate networks of connections, offer a large database for studying human behaviors (Christakis and Fowler, 2013), social relationships (Qiu et al., 2021), and social networks (Zhang and Amini, 2023). Researchers in computational social science apply advanced computational techniques, such as machine learning, natural language processing (NLP), and network analysis, to analyze massive datasets extracted from social media platforms. These analyses provide insights into various phenomena, including information diffusion (Jiang et al., 2014), opinion formation (Xiong and Liu, 2014), and collective behavior (Pinheiro et al., 2016).

### 2.3 Sentiment, Personality, and Value

The community has employed machine learning-based models to study human sentiment (Malviya et al., 2020), personality (Stachl et al., 2020), and values (Qiu et al., 2022; Ye et al., 2024). Building on this, research has explored AI personalities, including those of LLMs (Serapio-García et al., 2023; Li et al., 2023b; Röttger et al., 2024), and their alignment with human values (Ren et al., 2024b). Inspired by recent efforts to reveal AI agent traits through indirect methods like physiological exams (Jiang et al., 2024), questionnaires (Huang et al., 2023), and cultural views (Kovač et al., 2023), we adopt a novel perspective: examining machine behavior and inferred personality traits, including underlying values, to evaluate the performance of VLMs on mainstream social media platforms. To ground this examination, we utilize Schwartz's value theory (Schwartz, 1992), a comprehensive and culturally validated classification of human motivations (Davidov et al., 2008) that predicts attitudes and social behavior (Bardi and Schwartz, 2003). While prior computational work like ValueNet (Qiu et al., 2022) and ValueBench (Ren et al., 2024b) has assessed values in text-based NLP systems, and recent studies have begun exploring cultural value sensitivity in multimodal models (Yadav et al., 2025), our work specifically extends this value-based evaluation to VLMs within visual contexts on social media.

## 3 Data Collection

Inspired by ScreenAgent (Niu et al., 2024), our work leverages a VLM-driven graphical user interface (GUI) agent to autonomously navigate popular social media platforms. This agent conducts ran-
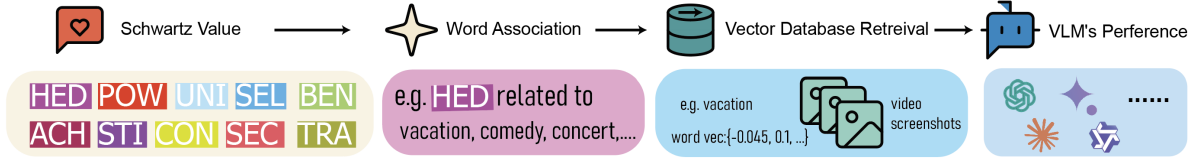
Figure 4: **Schwartz value-based image retrieval pipeline.** Our pipeline retrieves video screenshots based on Schwartz values by associating each value with relevant keywords, such as linking *Hedonism* to topics like vacation. These keywords are converted into vector queries, retrieving matching video content using the shortest distance within our database.

dom walks through social media platforms where it observes and captures screenshots. The data collected is stored in a vector database (Han et al., 2023), creating a structured repository optimized for value decomposition and efficient retrieval. We aim to analyze VLMs' behaviors across diverse social contexts and reveal their preferences. The automated data collection process (see Figure 3) efficiently fetches a large volume of diverse content, granting the analysis with scope and depth that traditional manual collection methods could not achieve.

The resulting dataset comprises $50,191$ videos sourced from Instagram ($32\%$), YouTube ($29\%$), and TikTok ($39\%$). Each entry includes the video link, a screenshot, and meta-information such as platform name and post date, capturing a comprehensive snapshot of content ranged between July 31, 2024, and October 31, 2024. By collecting data proportionally across these platforms, we enable balanced analysis and facilitate unbiased value decomposition across social media content. This innovative dataset empowers researchers to explore the behavior of VLMs in a systematic and organized way, fostering deeper insights into model interpretation and the dynamics of social media.

To examine the distribution of video themes in this dataset, we take a screenshot of each video at the beginning as a representation of the video content. In Figure 2, we present the distribution of videos that are relevant to ten Schwartz values. Specifically, for each Schwartz value, we curate 10 representative keywords. The area of the transparent circle is proportional to the number of videos that lie within a distance of $1.5$ to the corresponding normalized keyword vector. We find that videos relevant to these Schwartz Value Dimensions *Achievement*, *Hedonism*, and *Power* appear most frequently, while videos about *Tradition* are relatively rare. We then vectorize the image and define its relevance to a specific keyword using

cosine distance.

# 4 Evaluating VLM's Preferences

Extending the idea of analyzing LLMs' human likeness (Shanahan et al., 2023; Kovač et al., 2023) to VLMs with both visual and textual inputs, we ask: *Do VLMs also exhibit inherent preferences?*

To answer this question, we explore a diverse set of VLMs including GPT-4o (OpenAI, 2023), Gemini 2.0 Flash (Team et al., 2023), Claude 3.5 Sonnet (Anthropic, 2023), DeepSeek-VL2 (Wu et al., 2024), Qwen2.5-VL-Plus (Bai et al., 2023), InternVL2 (Chen et al., 2024c), CogVLM2 (Hong et al., 2024), and Blip-2 (Li et al., 2023a) to assess their value preferences. We quantify a VLM's preferences by evaluating its attitudes towards the 10 Schwartz values. This approach enables us to construct a comprehensive profile of each model's value preferences and to identify its unique value traits.

After constructing a vector database (as described in Section 3) to retrieve images based on specific Schwartz values' keywords, we analyze and compare the responses and attitudes of each model toward them. Our analysis reveals the extent to which each value captures the VLMs' attention, uncovering both similarities and differences across models, and highlighting distinct inclinations and sentiments within each VLM.

## 4.1 Preference Retrieval

To evaluate a VLM's preference for a specific Schwartz value, we collect each model's responses to images associated with several keywords related to the value (see Figure 4). For instance, we selected the keywords *Equality*, *Globe*, and *Handshake* for the *Universalism* dimension because they closely align with its core principles of fairness and global awareness. For each keyword linked to the value, each model reviews five images and answers their attitude towards each image.
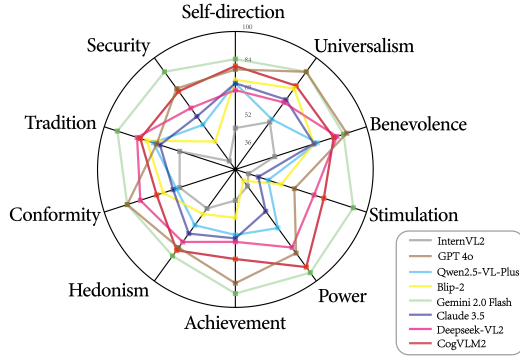
Figure 5: **Each VLM's preference scores towards the 10 Schwartz values.** The scores range from 0 to 100, with higher scores indicating stronger preferences for the corresponding values.
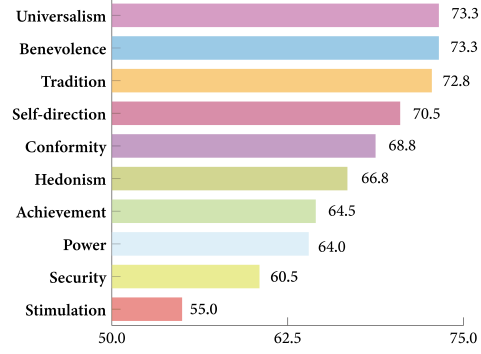


Figure 6: **Average preference scores for 10 Schwartz values across VLMs.** The length of each bar represents the mean score of all models for a value, with higher scores indicating a higher overall preference across all VLMs.

We retrieve the preference score of each VLM on the given visual input according to the following prompts:(1) *Do you like the content of this image? Please include yes or no in your answer, just respond in one word.* (2)*Why do you like or dislike this picture?* (3) *Describe this image in English briefly.*

The answer to the question is processed into either *yes* (1) or *no* (0), and the average score is calculated in percentage to evaluate the intensity of the model's preference for a given value (e.g., Universalism).

## 4.2 Preference Patterns

We evaluate and visualize the preference dimensions, identifying three distinct patterns:

*(1) Global Pattern*: After summarizing the preference score across all VLMs, we find most of them tend to prefer certain values over others. The results indicate a general preference for *Universalism* and *Benevolence* while showing relatively less liking for *Stimulation* and *Power*. The specific ranking of values is presented in Figure 6.

*(2) Range Consistency*: As shown in Figure 5, each model's preference scores, despite very few extremes, remain within a narrow band of approximately 15 around a central value. Models display varying levels of engagement with the content: some, like InternVL2, are more reserved, occupying a smaller area on the plot with lower preference scores, while others, like Gemini 2.0 Flash, show greater enthusiasm, demonstrating interest across all inputs with higher preference scores.

*(3) Individual Model Variations*: When analyzed individually, some models, such as Gemini 2.0 Flash, exhibit consistently high preferences across all 10 Schwartz values. In contrast, other models, like Claude 3.5 Sonnet, display more specific preferences as indicated in the standard deviation of scores across values. (Figure 7).

GPT-4o maintains balanced scores around 80, except for a notably lower score for *Stimulation*, with strong inclinations toward *Universalism* and *Benevolence*. Qwen2.5-VL-Plus covers a broad range of scores from 40 to 70, maintaining relative stability with a low standard deviation, yet it shares an aversion to *Stimulation* similar to other models. Deepseek-VL2 is highly balanced, slightly favoring *Benevolence*. CogVLM2, with a very low standard deviation, stands out by marking *Power* as its highest, diverging from all other models. Blip-2 ranks low across most values and has the highest variance, often providing short, inconclusive, or passive responses when reasoning about its likes or dislikes, reflecting a lack of ability to express preferences. Claude 3.5 Sonnet has a distinct personality, with the second-highest standard deviation, favoring *Universalism* while scoring the lowest in *Power* and *Stimulation*. InternVL2 demonstrates the lowest engagement, particularly disregarding *Stimulation*, while its high standard deviation suggests that, despite overall disinterest, it does show selective tendencies in certain areas. Finally, Gemini 2.0 Flash has both the lowest standard deviation and the highest overall scores, maintaining values above 80 across all dimensions, making it the most consistent and high-scoring model.

| Model | Open-source | Parameters | Self-direction | Universalism | Benevolence | Stimulation | Power | Achievement | Hedonism | Conformity | Tradition | Security |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GPT-4o | ✗ | – | 78 | 90 | 88 | 56 | 80 | 86 | 76 | 86 | 68 | 78 |
| Deepseek-VL2 | ✓ | 27B | 66 | 68 | 82 | 68 | 76 | 62 | 72 | 78 | 80 | 64 |
| Claude 3.5 Sonnet | ✗ | – | 70 | 70 | 68 | 34 | 50 | 60 | 66 | 58 | 66 | 58 |
| Gemini 2.0 Flash | ✗ | – | 84 | 90 | 86 | 92 | 94 | 92 | 82 | 86 | 92 | 90 |
| Blip-2 | ✓ | 2.7B | 72 | 78 | 68 | 48 | 28 | 48 | 52 | 64 | 74 | 40 |
| Qwen2.5-VL-Plus | ✓ | 72B | 70 | 56 | 70 | 40 | 62 | 58 | 60 | 56 | 70 | 52 |
| CogVLM2 | ✓ | 8B | 80 | 80 | 80 | 74 | 90 | 72 | 78 | 68 | 78 | 76 |
| InternVL2 | ✓ | 26B | 44 | 54 | 44 | 28 | 32 | 38 | 48 | 54 | 54 | 26 |

Table 1: **Comparison of VLMs' preference scores based on Schwartz's 10 values.** Higher scores indicate stronger preferences. "Open-source" indicates whether the model is publicly available, and "Parameters" denotes the model's size in billions (B).
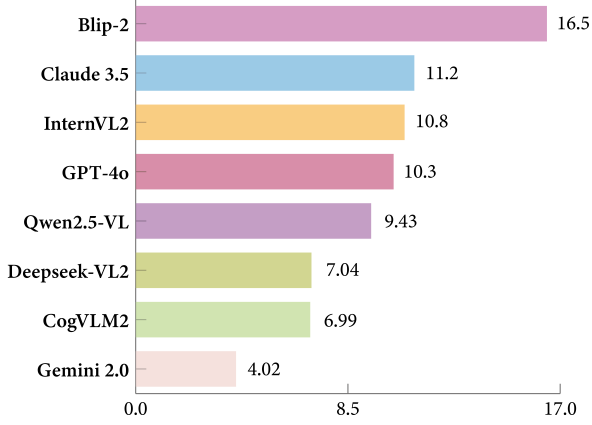


Figure 7: **Average standard deviation for each VLM.** Higher standard deviation indicates stronger preferences for certain values over others, while lower standard deviation reflects a more balanced attitude.

## 5 Inducing VLM's Preferences

Our initial experiment shows that VLMs have inherent inclinations toward different values. We now explore the dynamic aspects of VLM preferences beyond these static traits. We use Role-Playing Language Agents (RPLA) (Chen et al., 2024b) as a framework to assess VLMs' ability to adapt dynamically and simulate different personas, making decisions accordingly. Building on research showing that LLMs can emulate personas through RPLA (Serapio-García et al., 2023), we propose two key questions for VLMs: (1) *How well can VLMs align their traits to role-play personas using specific prompts?* (2) *Can strategies enhance accuracy and consistency in role-playing performance?*

### 5.1 Experiment

We use social media recommendation systems to evaluate whether VLMs can exhibit preferences aligned with the specified embedded persona. These systems rely on viewing duration as a key signal for content recommendation(Appendix A).

Considering the complexity of the experiment and the stability of model performance, we ultimately selected the following five models for evaluation: GPT-4o, Gemini 1.5 Pro, Qwen-VL-Plus (Bai et al., 2023), Claude 3.5 Sonnet, and CogVLM (Wang et al., 2023a). We assess the VLMs' role-playing ability by analyzing how well the recommended content reflects the imposed preferences. For example, adopting a *pet owner* persona should heighten the model's emphasis on *Benevolence*, valuing kindness and care, resulting in longer engagement with pet care videos and increased related recommendations (Liu et al., 2023).

In addition, we improve VLM performance on social networks by inducing personas through a questionnaire (Abeysinghe and Circi, 2024), comprehensively evaluating traits like emotional engagement, value alignment, curiosity, and preference matching to guide structured optimization.

**Simple Strategy**

In the simple strategy, we assign a specific persona in the demographic persona dataset from Persona-Chat (Zhang et al., 2018) to the VLM using the prompt: *You are a person who possesses certain traits, and the following statements best describe you: {Personality 1, 2, 3 . . . }.* Then, we pose a simple question: *Determine whether you are interested in the content of the given picture.*

The VLM engages with video shorts by responding either *yes* or *no*. A *yes* response prompts the VLM agent to remain on the current video, while a *no* results in an immediate skip. Alignment is measured as the increase in the frequency of recommended content that the VLM decides is interesting over time.

$$I_{\text{avg}} = \frac{1}{N} \sum_{t=1}^{N} \left( \frac{\sum_{i=1}^{n} Y_l^{(t)}(i) - \sum_{i=1}^{n} Y_f^{(t)}(i)}{\sum_{i=1}^{n} Y_f^{(t)}(i)} \right)$$
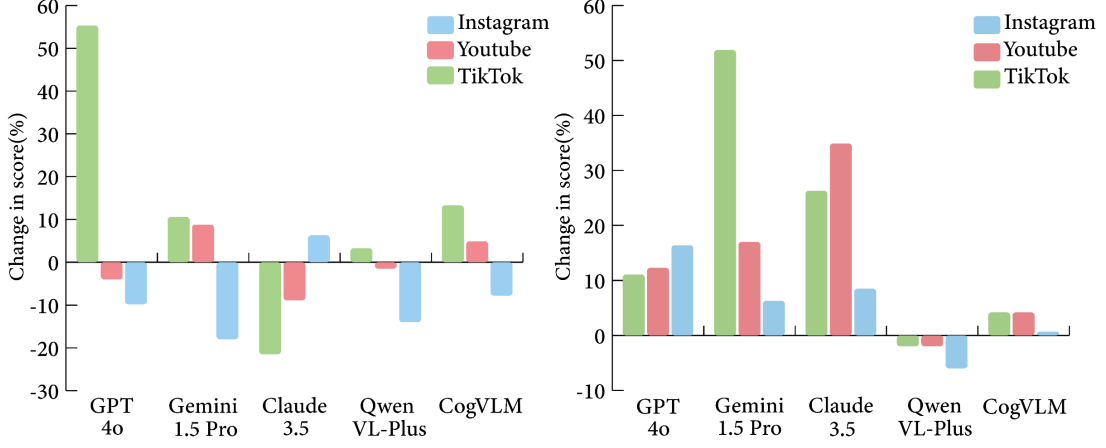
Figure 8: **Each VLM's percentage score of preference alignment changes across TikTok, YouTube, and Instagram.** Positive values indicate an increase in alignment, while negative values represent a decrease.

We define $I_{\text{avg}}$, the averaged percentage increase of *yes* responses, to measure the effectiveness of the strategy. $Y_l(i)$ and $Y_f(i)$ are the number of *yes* responses until $i$-th video in the last and first $n = 50$ videos, respectively. For each model, we conduct $N = 10$ trials, with each trial consisting of 100 video scrolls in total.

**Result Analysis.** Results highlight significant differences in role-playing effectiveness across platforms and models. TikTok stands out for GPT-4o and CogVLM, where GPT-4o exhibits "overfitting" behavior, showing highly nuanced responses to assigned roles that align closely with TikTok's recommendation system. However, this strong alignment is not consistent across models; for instance, Claude 3.5 Sonnet performs worse on Tik-Tok, suggesting model-specific sensitivities to the platform's dynamics. On YouTube and Instagram, performance is generally lower, with only modest gains or even negative alignment observed. These results indicate that TikTok's algorithmic design may amplify certain models' role-play capabilities, whereas YouTube and Instagram seem less conducive to capturing role-play nuances, possibly due to differences in content structure, user interaction patterns, or recommendation algorithms.

From previous experiments, CogVLM and Qwen-VL-Plus view all Schwartz values favorably, yet when CogVLM excels in this role-playing task, effectively adopting role-specific preferences, Qwen-VL-Plus shows only partial adherence. Blip-2 demonstrates no engagement or role-playing ability, lacking any signs of an induced personality. The findings show that even basic prompts can evoke detectable preferences, with some platforms emerging as particularly well-suited for role-playing tasks. Model adaptability in expressing role-related traits varied significantly if the persona is given in a simpler prompt.

**Inductive Scoring Questionnaire Strategy.**
Building on insights from the simple questioning approach, we develop the ISQ strategy to enhance VLMs' performance in social media alignment tasks. ISQ employs a series of prompts inquiring about various aspects of the screenshot. When presented with visual content, VLMs are asked to rate aspects like visual appeal, preference alignment, curiosity, etc.

Prompts include questions such as *On a scale of 1 to 10, how visually appealing is this screenshot to you based on your persona?* and *Does this screenshot make you want to click and start watching the video immediately?*

The ISQ strategy calculates a composite score to assess VLM engagement, with scores above a threshold (e.g. 60) indicating genuine interest, prompting extended interaction. This layered approach enhances persona analysis and final preference evaluation, improving role-specific performance on social media platforms.

The score is calculated as:

$$S_\% = \frac{v_a + c_s + e_e + v_e + 10p_a + 10a_d}{60} \times 100$$

Each response contributes to the total score: $v_a$ for visual appeal, $c_s$ for curiosity stimulation, $e_e$ for emotional engagement, $v_e$ for value expectation, $p_a$ for preference alignment (yes = 1, no = 0), $a_d$ for action desire (yes = 1, no = 0). The increase is calculated the same as the simple strategy.
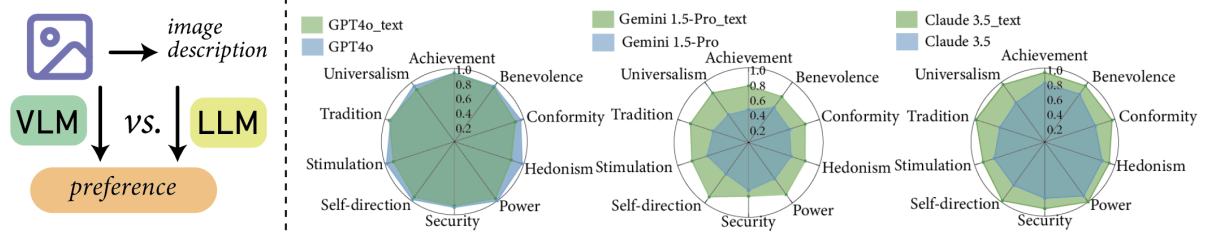
Figure 9: **Value distribution comparison between VLMs and corresponding LLMs.** For the same model (e.g., GPT-4o and GPT-4o_text), different input modes (multi-modal vs. text-only) are compared. Experiments demonstrate that the choice of multi-modal input significantly influences some models' value preferences. While models like GPT-4o show consistency across input modes, others, such as Claude 3.5 Sonnet and Gemini 1.5 Pro, exhibit notable differences in preferences.

**Result Analysis.** Compared to the simple strategy, the ISQ strategy performances are elevated throughout all models and platforms except for Qwen-VL-Plus. This shows that the strategy could successfully induce the model's role-playing ability in preference indication detectable through social media on behalf of the persona.

Following the trend in simple strategy, we could see a strong increase for TikTok, particularly with Gemini 1.5 Pro, which demonstrates an average rise of as much as 51.9. GPT-4o, Gemini 1.5 Pro, and Claude 3.5 Sonnet—all displayed notable improvements, with consistently positive changes in performance. This suggests that these models respond well to the ISQ strategy, allowing them to adopt and express induced personalities with greater depth. Among them, Gemini 1.5 Pro and Claude 3.5 Sonnet particularly benefited from the ISQ approach, showing remarkable growth in comparison to earlier results in the simple strategy. This demonstrates that the ISQ strategy enhances VLMs' ability to engage with role-playing tasks more effectively than the previous simple strategy.

## 6 Discussion

**VLMs *vs.* Corresponding LLMs.** We conduct several experiments to examine whether different types of multi-modal inputs influence value preference outcomes. As shown in Figure 9, we compare value preferences derived directly from VLMs using images as input with those generated by feeding the corresponding text-based image descriptions created by the same VLM into their paired LLMs. The results reveal significant differences in value preferences between these two modes. For many value dimensions, VLMs and LLMs produce distinct preference patterns, especially in models like Claude 3.5 Sonnet and Gemini 1.5 Pro,

where the outputs diverge significantly across input modes. In contrast, GPT-4o displays greater consistency across modes, suggesting its ability to integrate visual and textual information cohesively. These findings highlight that the choice of input mode—visual or text—can significantly affect model outputs, underscoring the importance of input selection in applications requiring personalized or human-like responses. Detailed evaluation methods and results are provided in Appendix E.

**Single Frame Screenshot Representation.** To validate single-frame screenshots for video content analysis, we randomly select 500 images from each different social media platform in our dataset for human evaluation. Annotators are provided with the instructions outlined in Appendix D, along with the images and additional context. Each image-video pair is assessed by three annotators, resulting in a total of 4,500 ratings. Annotators review the full video and its corresponding screenshot, rating how accurately the screenshot represents the video's content. The results indicate that 90.4% of screenshots are considered to be representative of the video's main content, demonstrating the effectiveness of single-frame screenshots across platforms. However, 8.8% of the frames are rated as non-representative, highlighting the challenges posed by videos with complex scenes or rapid transitions.

## 7 Conclusion

Our study introduces ***Value-Spectrum***, a benchmark for evaluating value preferences in VLMs using a vector database derived from social media platforms. Through systematic evaluation, we observe a shared global inclination among models toward certain mainstream values, such as *Universalism*, likely influenced by the nature of their train-

ing data. At the same time, significant differences emerged across other value dimensions, highlighting disparities in how VLMs align with diverse human-designed value systems. These findings reveal both commonalities that reflect broader societal trends and divergences that underscore model-specific characteristics, prompting us to explore whether these variations can be adjusted to induce specific personas.

Our work also provides practical insights into VLMs' ability to adapt their value preferences dynamically through role-playing, offering a pathway to align machine behaviors with human-designed personas. By connecting role-playing capabilities and alignment strategies, we aim to inspire further research into value-driven AI agent systems and their adaptability in real-world applications.

## 8 Limitations

The evaluation utilizes Schwartz's value dimensions as the foundation for understanding personality traits and preferences, highlighting opportunities for future research to incorporate broader cultural and personality-based perspectives. Future studies might consider expanding the set of value dimensions or integrating alternative value systems. Additionally, our use of single-frame screenshots, in order to simplify analysis and maintain efficiency, is proven to effectively represent the corresponding video content through annotators' high ratings. However, our methodology still faces challenges in capturing the essence of some videos that are fast-changing and complex in their storytelling. This leaves room for future refinement.

## 9 Ethical Considerations

We eliminate any harmful effects of VLMs by ensuring that they only observe content without interacting through comments or likes. This approach maintains the integrity of the social media ecosystem and prevents unintended AI-driven consequences. However, we recognize that VLMs may still inadvertently produce discriminatory content, reflecting biases based on gender, race, or socioeconomic status. We acknowledge these challenges and emphasize the need for ongoing efforts to address and minimize such biases in model outputs.

## References

Bhashithe Abeysinghe and Ruhan Circi. 2024. The challenges of evaluating llm applications: An analysis of automated, human, and llm-based approaches. *arXiv preprint arXiv:2406.03339*.

Anthropic. 2023. Claude 3.5.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*.

Anat Bardi and Shalom H Schwartz. 2003. Values and behavior: Strength and structure of relations. *Personality and Social Psychology Bulletin*, 29(10):1207–1220.

Simone Bianco, Luigi Celona, Marco Donzella, and Paolo Napoletano. 2023. Improving image captioning descriptiveness by ranking and llm-based fusion. *arXiv preprint arXiv:2306.11593*.

Jiangjie Chen, Xintao Wang, Rui Xu, Siyu Yuan, Yikai Zhang, Wei Shi, Jian Xie, Shuang Li, Ruihan Yang, Tinghui Zhu, Aili Chen, Nianqi Li, Lida Chen, Caiyu Hu, Siye Wu, Scott Ren, Ziquan Fu, and Yanghua Xiao. 2024a. From persona to personalization: A survey on role-playing language agents. *ArXiv*, abs/2404.18231.

Jiangjie Chen, Xintao Wang, Rui Xu, Siyu Yuan, Yikai Zhang, Wei Shi, Jian Xie, Shuang Li, Ruihan Yang, Tinghui Zhu, Aili Chen, Nianqi Li, Lida Chen, Caiyu Hu, Siye Wu, Scott Ren, Ziquan Fu, and Yanghua Xiao. 2024b. From persona to personalization: A survey on role-playing language agents. *arXiv preprint arXiv:2404.18231*.

Yan Chen, Kate Sherren, Michael Smit, and Kyung Young Lee. 2023. Using social media images as data in social science research. *New Media & Society*, 25(4):849–871.

Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, Ji Ma, Jiaqi Wang, Xiao wen Dong, Hang Yan, Hewei Guo, Conghui He, Zhenjiang Jin, Chaochao Xu, Bin Wang, Xingjian Wei, Wei Li, Wenjian Zhang, Bo Zhang, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, and Yu Qiao. 2024c. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *ArXiv*, abs/2404.16821.

Nicholas A Christakis and James H Fowler. 2013. Social contagion theory: examining dynamic social networks and human behavior. *Statistics in medicine*, 32(4):556–577.

Eldad Davidov, Peter Schmidt, and Shalom H Schwartz. 2008. The measurement of values: Testing and improving the theory in a comparative perspective. *European Journal of Social Psychology*, 38(6):881–897.

Yikun Han, Chunjiang Liu, and Pengfei Wang. 2023. A comprehensive survey on vector database: Storage and retrieval technique, challenge. *arXiv preprint arXiv:2310.11703*.

Wenyi Hong, Weihan Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng, Shiyu Huang, Junhui Ji, Zhao Xue, Lei Zhao, Zhuoyi Yang, Xiaotao Gu, Xiaohan Zhang, Guanyu Feng, Da Yin, Zihan Wang, Ji Qi, Xixuan Song, Peng Zhang, De-Feng Liu, Bin Xu, Juanzi Li, Yu-Chen Dong, and Jie Tang. 2024. Cogvlm2: Visual language models for image and video understanding. *ArXiv*, abs/2408.16500.

Wenbo Hu, Yifan Xu, Yi Li, Weiyue Li, Zeyuan Chen, and Zhuowen Tu. 2024. Bliva: A simple multimodal llm for better handling of text-rich visual questions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 2256–2264.

Jen-tse Huang, Man Ho Lam, Eric John Li, Shujie Ren, Wenxuan Wang, Wenxiang Jiao, Zhaopeng Tu, and Michael R Lyu. 2023. Emotionally numb or empathetic? evaluating how llms feel using emotionbench. *arXiv preprint arXiv:2308.03656*.

Chunxiao Jiang, Yan Chen, and KJ Ray Liu. 2014. Evolutionary dynamics of information diffusion over social networks. *IEEE transactions on signal processing*, 62(17):4573–4586.

Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. 2024. Evaluating and inducing personality in pre-trained language models. *Advances in Neural Information Processing Systems*, 36.

Grgur Kovač, Masataka Sawayama, Rémy Portelas, Cédric Colas, Peter Ford Dominey, and Pierre-Yves Oudeyer. 2023. Large language models as superpositions of cultural perspectives. *arXiv preprint arXiv:2307.07870*.

Junlong Li, Fan Zhou, Shichao Sun, Yikai Zhang, Hai Zhao, and Pengfei Liu. 2024. Dissecting human and llm preferences. *ArXiv*, abs/2402.11296.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.

Junyi Li, Ninareh Mehrabi, Charith Peris, Palash Goyal, Kai-Wei Chang, A. G. Galstyan, Richard Zemel, and Rahul Gupta. 2023b. The steerability of large language models toward data-driven personas. *ArXiv*, abs/2311.04978.

Fangyu Liu, Julian Martin Eisenschlos, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Wenhu Chen, Nigel Collier, and Yasemin Altun. 2022. Deplot: One-shot visual language reasoning by plot-to-table translation. *arXiv preprint arXiv:2212.10505*.

Peng Liu, Lemei Zhang, and Jon Atle Gulla. 2023. Pretrain, prompt, and recommendation: A comprehensive survey of language modeling paradigm adaptations in recommender systems. *Transactions of the Association for Computational Linguistics*, 11:1553–1571.

Yujie Lu, Dongfu Jiang, Wenhu Chen, William Yang Wang, Yejin Choi, and Bill Yuchen Lin. 2024. Wildvision: Evaluating vision-language models in the wild with human preferences. *ArXiv*, abs/2406.11069.

Sunil Malviya, Arvind Kumar Tiwari, Rajeev Srivastava, and Vipin Tiwari. 2020. Machine learning techniques for sentiment analysis: A review. *SAMRIDDHI: A Journal of Physical Sciences, Engineering and Technology*, 12(02):72–78.

Runliang Niu, Jindong Li, Shiqi Wang, Yali Fu, Xiyu Hu, Xueyuan Leng, He Kong, Yi Chang, and Qi Wang. 2024. Screenagent: A vision language model-driven computer control agent. *arXiv preprint arXiv:2402.07945*.

OpenAI. 2023. Gpt-4.

Flávio L Pinheiro, Francisco C Santos, and Jorge M Pacheco. 2016. Linking individual and collective behavior in adaptive social networks. *Physical review letters*, 116(12):128702.

Liang Qiu, Yuan Liang, Yizhou Zhao, Pan Lu, Baolin Peng, Zhou Yu, Ying Nian Wu, and Song-Chun Zhu. 2021. Socaog: Incremental graph parsing for social relation inference in dialogues. *arXiv preprint arXiv:2106.01006*.

Liang Qiu, Yizhou Zhao, Jinchao Li, Pan Lu, Baolin Peng, Jianfeng Gao, and Song-Chun Zhu. 2022. Valuenet: A new dataset for human value driven dialogue system. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11183–11191.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*.

Yuanyi Ren, Haoran Ye, Hanjun Fang, Xin Zhang, and Guojie Song. 2024a. Valuebench: Towards comprehensively evaluating value orientations and understanding of large language models. *ArXiv*, abs/2406.04214.

Yuanyi Ren, Haoran Ye, Hanjun Fang, Xin Zhang, and Guojie Song. 2024b. Valuebench: Towards comprehensively evaluating value orientations and understanding of large language models. *arXiv preprint arXiv:2406.04214*.

Paul Röttger, Valentin Hofmann, Valentina Pyatkin, Musashi Hinck, Hannah Rose Kirk, Hinrich Schutze, and Dirk Hovy. 2024. Political compass or spinning arrow? towards more meaningful evaluations for values and opinions in large language models. In *Annual Meeting of the Association for Computational Linguistics*.

Shalom H Schwartz. 1992. Universals in the content and structure of values: Theoretical advances and empirical tests in 20 countries. *Advances in Experimental Social Psychology*, 25:1–65.

Shalom H Schwartz. 2012. An overview of the schwartz theory of basic values. *Online readings in Psychology and Culture*, 2(1):11.

Greg Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Peter Romero, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2023. Personality traits in large language models. *arXiv preprint arXiv:2307.00184*.

Murray Shanahan, Kyle McDonell, and Laria Reynolds. 2023. Role-play with large language models. *arXiv preprint arXiv:2305.16367*.

Clemens Stachl, Florian Pargent, Sven Hilbert, Gabriella M Harari, Ramona Schoedel, Sumer Vaid, Samuel D Gosling, and Markus Bühner. 2020. Personality research and assessment in the era of machine learning. *European Journal of Personality*, 34(5):613–631.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. 2023a. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*.

Xintao Wang, Yunze Xiao, Jen tse Huang, Siyu Yuan, Rui Xu, Haoran Guo, Quan Tu, Yaying Fei, Ziang Leng, Wei Wang, Jiangjie Chen, Cheng Li, and Yanghua Xiao. 2023b. Incharacter: Evaluating personality fidelity in role-playing agents through psychological interviews. In *Annual Meeting of the Association for Computational Linguistics*.

Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bing-Li Wang, Zhenda Xie, Yu Wu, Kai Hu, Jiawei Wang, Yaofeng Sun, Yukun Li, Yishi Piao, Kang Guan, Aixin Liu, Xin Xie, Yu mei You, Kaihong Dong, Xingkai Yu, Haowei Zhang, Liang Zhao, Yisong Wang, and Chong Ruan. 2024. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding. *ArXiv*, abs/2412.10302.

Fei Xiong and Yun Liu. 2014. Opinion formation on social media: an empirical approach. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 24(1).

Srishti Yadav, Zhi Zhang, Daniel Hershcovich, and Ekaterina Shutova. 2025. Beyond words: Exploring cultural value sensitivity in multimodal models. In *North American Chapter of the Association for Computational Linguistics*.

Haoran Ye, Yuhang Xie, Yuanyi Ren, Hanjun Fang, Xin Zhang, and Guojie Song. 2024. Measuring human and ai values based on generative psychometrics with large language models. In *AAAI Conference on Artificial Intelligence*.

Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. 2023. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46:5625–5644.

Linfan Zhang and Arash A Amini. 2023. Adjusted chi-square test for degree-corrected block models. *The Annals of Statistics*, 51(6):2366–2385.

Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? *arXiv preprint arXiv:1801.07243*.

Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason J. Corso, and Jianfeng Gao. 2019. Unified vision-language pre-training for image captioning and vqa. *ArXiv*, abs/1909.11059.

## Appendix

## A Industry References

For further information on the signals used for content recommendation, refer to the following blogs:

- YouTube:https://www.youtube.com/howyoutubeworks/product-features/recommendations/#signals-used-to-recommend-content Watch history: Our system uses the YouTube videos you watch to give you better recommendations, remember where you left off, and more.

- Instagram:https://about.instagram.com/blog/announcements/instagram-ranking-explained Viewing history: This looks at how often you view an account's stories, so we can prioritize the stories from accounts we think you don't want to miss.

- TikTok:https://support.tiktok.com/en/using-tiktok/exploring-videos/how-tiktok-recommends-content User interactions: Content you like, share, comment on, and watch in full or skip, as well as accounts of followers that you follow back.

## B Inducing VLM's Personas

### B.1 Experiment Steps

In this section, we detail the steps of our experiments designed to evaluate the effectiveness of different strategies in identifying persona-related content on social media platforms. The experiment comprises three main parts:

**Open The Designated Social Media Platform** The second step involves accessing the designated social media platform. For demonstration, we focus on TikTok and GPT-4o.

- Open TikTok website(Figure 10).

- Navigate to the 'For You' page, where a variety of content is displayed.

**Capture Screenshot Image of Playing Short Video** Next, we capture screenshots of the short videos that are playing. This captured screenshot is then inputted to the VLM. An example of a screenshot is shown in Figure11.



Figure 10: Screenshot of the TikTok homepage.

(URL: https://www.tiktok.com/@pugloulou/video/7342967563321822497).

**Responses and Strategy Actions for Models** We design specific questionnaire prompts for different experimental purposes and then collect and analyze responses from different VLMs. Based on these responses, we apply various strategic actions.
*Simple Strategy*: See details in Figure 13.
*ISQ Strategy*: See details in Figure 14.

## C Single VS. Multi Frames

In both experiments, screenshots were captured exactly 2 seconds into the video shorts. This timing was chosen because most videos begin their main narrative at this point. Multi-frame analysis was not utilized for two key reasons:

**Preference Evaluation**: Using a single frame aligns with CLIP's capability to filter and retrieve the most relevant social media screenshots. Multiple frames are unnecessary for this purpose.

**Preference Induction**: For recommendation systems to recognize user preferences, the staying duration for each video is critical. Capturing multiple frames increases processing time, causing most videos to be viewed in their entirety before scrolling. This diminishes the strategy's impact and hinders the system's ability to distinguish preferences between videos.

Figure 11: Screenshot of playing a short video

Thus, single-frame analysis was deemed more effective and practical for the experiments.

## D  Human Annotators: Single-Frame Analysis

We refine the survey formats provided to annotators through multiple iterations, conducting pilot studies on Amazon Mechanical Turk (MTurk) to continuously adjust the instructions until the quality of answers by the annotators meets the desired standard. The instruction examples referenced by the annotators can be found in Figure 15.

In the Amazon MTurk task description provided to annotators, we clearly stated that this task was for research purposes. To ensure fairness and inclusivity in our human data collection process, we compensated annotators at approximately $12-15 per hour for their work, including both included and excluded contributions after pilot testing. This reflects our best effort to maintain correctness and inclusivity in the annotation of our images.

## E  Performance Analysis of Multi-Modal Inputs Across Value Dimensions

To investigate the influence of different input modalities on value preference outcomes, we conduct experiments to compare results derived from direct visual inputs with those generated using text-based image descriptions. Specifically, we randomly selected 500 images from the Value-Spectrum dataset, ensuring balanced representation across 10 value dimensions (50 samples per dimension). We then retrieved image descriptions generated by three VLMs —GPT-4o, Gemini 1.5 Pro, and Claude 3.5 Sonnet—by prompting these models with images from the dataset. These textual descriptions were then fed into the text-based versions of the different models to conduct Value Preference QA.

The results revealed notable differences in value preference distributions across input modalities. While GPT-4o demonstrated relatively consistent performance between visual and text-based inputs, models like Gemini 1.5 Pro and Claude 3.5 Sonnet displayed greater variability, with outputs diverging significantly in specific dimensions. This suggests that the choice of input mode—visual or textual—can impact the models' ability to align responses with underlying value dimensions. The detailed scores for all models and input settings are summarized in Table 10, highlighting patterns across dimensions such as Achievement, Benevolence, and Tradition. These findings emphasize the importance of input modality selection in tasks requiring a nuanced understanding of human values.

## F  Cultural Diversity Analysis of the Dataset

We also conducted a systematic analysis to quantify the cultural variety present. This analysis aimed to provide a clearer understanding of the distribution of cultural representations within the data collected from platforms like TikTok, YouTube Shorts, and Instagram Reels.

### F.1  Methodology

We randomly sampled 2,614 images from our dataset. A two-step pipeline was then applied:

1. **Caption Generation:** For each sampled image, we utilized GPT-4o to generate a descriptive caption. This step aimed to produce a

textual representation that more precisely captures the visual content of the image, facilitating subsequent cultural analysis.

2. **Cultural Signal Extraction:** Using the generated captions, we applied named entity recognition (NER) and keyword matching techniques. This process was designed to identify mentions of specific cultural elements, including but not limited to:

   - Locations: e.g., "Seoul," "Mecca"
   - Traditional Clothing: e.g., "sari," "hanbok"
   - Holidays: e.g., "Diwali," "Ramadan"
   - Scripts/Languages: e.g., "Arabic calligraphy"
   - Other culturally relevant features and artifacts.

Based on the extracted signals, we created a culture-labeled subset where each sample could be associated with one or more identified cultural categories.

### F.2    Results

The aggregated results from the cultural signal extraction process are summarized in Table 2.

| Culture | Count | Percentage |
|---------|-------|------------|
| Western | 2054 | 78.6% |
| Japanese | 364 | 13.9% |
| Korean | 61 | 2.3% |
| Chinese | 49 | 1.9% |
| Muslim | 37 | 1.4% |
| Indian | 37 | 1.4% |
| Arabic | 12 | 0.5% |
| Others | 183 | 7.0% |
| **Total** | **2614** | **100.0%** |

Table 2: Distribution of Identified Cultural Categories in the Sampled Dataset (N=2,614). Percentages are rounded.

The distribution presented in Table 2 shows a predominance of Western cultural signals (78.6%). This is consistent with the global reach and content characteristics of the source platforms (TikTok, YouTube Shorts, and Instagram Reels), which, while internationally utilized, often feature a significant volume of content reflecting Western cultural contexts. Nevertheless, the dataset also captures notable representation from several non-Western cultural spheres, including Japanese (13.9%), Korean

(2.3%), Chinese (1.9%), Muslim (1.4%), Indian (1.4%), and Arabic (0.5%) cultural elements.

The current cultural distribution reflects the inherent, though skewed, diversity present on these globally accessible social media platforms. While the dataset provides a basis for studying VLM responses to a range of cultural stimuli, the observed imbalance in representation is an important characteristic to consider. Future work may explore strategies for targeted data collection to enhance the cultural balance and further broaden the inclusiveness of such datasets for cross-cultural AI research.

## G    Human Validation of Value Alignment in Screenshots and Videos

To empirically assess the alignment between the machine-assigned value labels and human perception, both for the selected screenshots and their corresponding full video content, we conducted a dedicated human validation study. This study provides an independent measure of how well the visual content (static images and dynamic videos) reflects the intended Schwartz value dimensions.

### G.1    Methodology

The human validation study was structured as follows:

1. **Sample Selection:** From our dataset, we randomly selected 50 screenshot-video pairs for each of the ten Schwartz value dimensions (e.g., Benevolence, Tradition, Hedonism). This resulted in a total of 500 unique items for evaluation.

2. **Annotation Task and Annotators:** Three independent human annotators, familiar with Schwartz's value theory, were recruited. For each item (a screenshot and its associated video URL), annotators were tasked with making two separate judgments based on the assigned value dimension:

   - Judgment 1 (Screenshot): "Does this **screenshot** reflect the given value of [Value Name]?" (Response options: Yes/No/Not Sure)
   - Judgment 2 (Video): "Does this **video** (linked) reflect the given value of [Value Name]?" (Response options: Yes/No/Not Sure)

Annotators were required to view the full video before making their judgment for the video content.

3. **Survey Design and Order:** The annotation survey was designed to present these judgments sequentially to mitigate order effects. For each item, annotators first evaluated the screenshot. Only after submitting this judgment were they presented with the video URL and asked to evaluate the video content against the same value dimension. A visual example of the survey interface is shown in Figure **??**.

4. **Data Aggregation:** With 50 items per value, 10 value dimensions, and 3 annotators, this process yielded 1,500 individual judgments for screenshots and another 1,500 for videos.

5. **Handling Inaccessible Video Content:** In instances where a video URL was no longer accessible (e.g., due to content deletion or platform restrictions), annotators were instructed to mark their response for the video judgment as "Not Sure." These instances are accounted for in the reported results.

## G.2 Results

The aggregated human judgments on value alignment for both screenshots and their corresponding videos are presented in Table 9. Percentages represent the proportion of responses for each category, averaged across three annotators.

Across all ten value dimensions, we observed strong agreement between human judgment and the intended value labels for both screenshots and full videos.

- On average, screenshot-based value alignment received a **90.6% "Yes"** rate from human annotators.

- Video-based value alignment achieved an **87.6% "Yes"** rate.

The "Not Sure" responses were minimal for images (average 1.00%) and slightly higher for videos (average 5.58%). The higher incidence of "Not Sure" for videos is primarily attributed to instances of unavailable video links at the time of annotation, as detailed in the methodology.

## G.3 Discussion

The findings from this human validation study provide strong empirical support for the value alignment of the collected content:

1. **Screenshots as Reliable Indicators of Value:** The high average "Yes" rate (90.6%) for screenshots demonstrates that the static images selected through our pipeline are generally perceived by human annotators as reliably reflecting the intended Schwartz value dimensions.

2. **Consistency of Value Perception in Video Content:** Full video evaluations also yielded a high average "Yes" rate (87.6%), indicating that the broader dynamic context of the videos largely aligns with the value perception derived from the representative screenshots. The slight decrease compared to screenshots, alongside the higher "Not Sure" rate for videos, can be attributed to factors such as the dynamic and multifaceted nature of video content, as well as the practical issue of video accessibility over time.

3. **Support for Benchmark Methodology:** The close agreement between human judgments and the assigned value labels, for both image and video modalities, supports the efficacy of our content retrieval and labeling approach. This suggests that the benchmark, while leveraging efficient screenshot-based processing, captures content that is largely consistent with human understanding of the targeted values.

This human validation ablation study therefore reinforces the suitability of the dataset for investigating VLM responses to visual content associated with diverse human values.

# H Inducing VLM's Persona Detailed Information

Table 3: **Simple Strategy - TikTok**

| Dimension | GPT-4o | Gemini 1.5 pro | Qwen-VL-Plus | CogVLM | Claude |
|---|---|---|---|---|---|
| Related contents(<=50)(%) | 7.6 | 20 | 6.0 | 45.2 | 15.8 |
| Related contents(LAST 50)(%) | 11.8 | 23.2 | 6.2 | 51.2 | 12.4 |
| Change(%) | 55.26 | 16 | 3.33 | 13.27 | -21.52 |

Table 4: **Questionnaire Strategy - TikTok**

| Dimension | GPT-4o | Gemini 1.5 pro | Qwen-VL-Plus | CogVLM | Claude |
|---|---|---|---|---|---|
| Related contents(<=50)(%) | 3.6 | 10.8 | 19.6 | 66.2 | 12.7 |
| Related contents(LAST 50)(%) | 4.0 | 16.4 | 19.2 | 69 | 16 |
| Change(%) | 11.1 | 51.9 | -2.0 | 4.2 | 26.3 |

Table 5: **Simple Strategy - YouTube**

| Dimension | GPT-4o | Gemini 1.5 pro | Qwen-VL-Plus | CogVLM | Claude |
|---|---|---|---|---|---|
| Related contents(<=50)(%) | 10 | 25 | 13.6 | 61 | 24.8 |
| Related contents(LAST 50)(%) | 9.6 | 27.2 | 13.4 | 64 | 22.6 |
| Change(%) | -4.0 | 8.8 | -1.47 | 4.9 | -8.9 |

Table 6: **Questionnaire Strategy - YouTube**

| Dimension | GPT-4o | Gemini 1.5 pro | Qwen-VL-Plus | CogVLM | Claude |
|---|---|---|---|---|---|
| Related contents(<=50)(%) | 11.4 | 20.0 | 42 | 81 | 15.6 |
| Related contents(LAST 50)(%) | 12.8 | 23.4 | 42.8 | 81 | 21.2 |
| Change(%) | 12.3 | 17.0 | 1.9 | 0 | 34.9 |

Table 7: **Simple Strategy - Instagram**

| Dimension | GPT-4o | Gemini 1.5 pro | Qwen-VL-Plus | CogVLM | Claude |
|---|---|---|---|---|---|
| Related contents(<=50)(%) | 22.4 | 27.8 | 11.4 | 53.6 | 15.8 |
| Related contents(LAST 50)(%) | 20.2 | 22.8 | 9.8 | 49.4 | 16.8 |
| Change(%) | -9.82 | -18 | -14 | -7.8 | 6.33 |

Table 8: **Questionnaire Strategy - Instagram**

| Dimension | GPT-4o | Gemini 1.5 pro | Qwen-VL-Plus | CogVLM | Claude |
|---|---|---|---|---|---|
| Related contents(<=50)(%) | 13.4 | 15.8 | 46.8 | 56.4 | 8.6 |
| Related contents(LAST 50)(%) | 15.6 | 16.8 | 44 | 56.8 | 9.4 |
| Change(%) | 16.4 | 6.3 | -6 | 0.7 | 8.5 |

| Value | Yes (Image) | Yes (Video) | No (Image) | No (Video) | Not Sure (Image) | Not Sure (Video) |
|---|---|---|---|---|---|---|
| Benevolence | 82.67 | 76.00 | 16.00 | 16.67 | 1.33 | 7.33 |
| Tradition | 94.67 | 96.00 | 4.67 | 1.33 | 0.67 | 2.67 |
| Self-Direction | 95.33 | 92.67 | 2.67 | 2.67 | 2.00 | 4.67 |
| Achievement | 88.67 | 78.00 | 10.67 | 11.33 | 0.67 | 10.67 |
| Conformity | 82.67 | 80.00 | 14.67 | 11.33 | 2.67 | 8.67 |
| Hedonism | 87.33 | 83.33 | 12.67 | 10.67 | 0.00 | 6.00 |
| Security | 94.00 | 92.67 | 6.00 | 4.00 | 0.00 | 3.33 |
| Stimulation | 96.00 | 97.33 | 3.33 | 1.33 | 0.67 | 1.33 |
| Power | 87.33 | 82.67 | 11.33 | 7.33 | 1.33 | 10.00 |
| Universalism | 97.33 | 97.33 | 1.33 | 0.67 | 1.33 | 2.00 |
| **Average** | **90.60** | **87.60** | **8.40** | **8.30** | **1.00** | **5.58** |

Table 9: Human Annotation Results for Value Alignment of Screenshots and Videos (Percentages).



Figure 12: Two examples of trials evaluating VLM's preferences. For each trial, a social media short video is used, and a screenshot is taken at the 2-second timestamp. The user then interacts with the VLM using a question-and-answer format to assess the model's attitude toward the screenshot's content. In these examples, TikTok content and the GPT-4o model are used for demonstration.

# Inducing VLM's Preference - **Simple Strategy**



**User**

You are a person who contains certain traits, several statements best describe you as your introduction:

1. I live in Toronto, Canada,

2. my favorite french fries are from McDonald's,

3. my mom and dad are from Japan,

4. **I'm an elementary school teacher,**

5. I love to shop.

Your goal is to determine whether you are interested in the content of the given picture. Give me only a structured answer: {'decision': 'yes or no', 'reason': 'reason', 'trait_num': '1-5, the personality trait that this picture alludes to, NA if none'}

**GPT 4o**  {'decision': 'yes', 'reason': 'The image depicts a school scenario, which relates to being an elementary school teacher.', 'trait_num': '4'}

**User**  Action: Watch for 45 Seconds, then scroll down 👁



**User**

You are a person who contains certain traits, several statements best describe you as your introduction:

1. I live in Toronto, Canada,

2. my favorite french fries are from McDonald's,

3. my mom and dad are from Japan,

4. **I'm an elementary school teacher,**
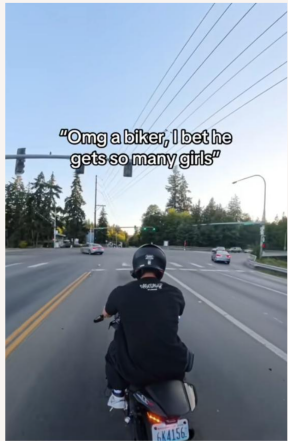
5. I love to shop.

Your goal is to determine whether you are interested in the content of the given picture. Give me only a structured answer: {'decision': 'yes or no', 'reason': 'reason', 'trait_num': '1-5, the personality trait that this picture alludes to, NA if none'}

**GPT 4o**  {'decision': 'no', 'reason': 'The content of the picture does not relate to any of my stated interests or traits.', 'trait_num': '4'}

**User**  Action: Scroll down immediately

Figure 13: Example of two scenarios in inducing the VLM's persona using the Simple strategy: when VLM determines that the screenshot content aligns with the persona, and the user remains engaged with the content for 45 seconds. Conversely, if the VLM decides the content is not related to the persona, the user scrolls down immediately.

## Inducing VLM's Preference - **Inductive Scoring Questionnaire Strategy**

**User** You are a person with these personalities:

1: my father use to drive a taxi cab in the 1970s,

2: i've been to canada two times my entire life,

3: i was once on the local news , and not for anything illegal,

4: my uncle use to work on fish boats in florida,

You are scrolling through short videos and come across this picture. Please answer the following questions as this person with the above persona:

1. On a scale of 1 to 10, how visually appealing is the screenshot to you based on your personality?

2. Based on the screenshot, does this video seem like something you would typically watch, or is it outside your usual preferences?

3. On a scale of 1 to 10, how strong does this screenshot stimulate your curiosity about the video's content?

4. On a scale of 1 to 10, how strongly does this screenshot make you feel emotionally connected or engaged with the content?

5. Does this screenshot make you want to click and start watching the video immediately?

6. On a scale of 1 to 10, how likely do you think this video will provide valuable information or entertainment to you?

7. Which of your personalities does this screenshot most appeal to?

Please answer the questions and give your answer to all of these question in the following form through replacing the righthand side, please answer in this format

{

'visual_appeal': 'Give me your answer to the first question, it should be a number between 1 to 10',

'preference_alignment': 'Give me your answer to the second question, it should be yes or no',

'curiosity_stimulation': 'Give me your answer to the third question, it should be a number between 1 to 10',

'emotional_engagement': 'Give me your answer to the fourth question, it should be a number between 1 to 10',

'action_desire': 'Give me your answer to the fifth question, it should be yes or no',

'value_expectation': 'Give me your answer to the sixth question, it should be a number between 1 to 10',

'trait_num': 'Give me your answer to the last question, it should be a number between 1 to 4 corresponding to the matching persona index, NA if none'

}

Score = 73.75 > 60

**GPT 4o** {'visual_appeal': 6, 'preference_alignment': 'yes', 'curiosity_stimulation': 7, 'emotional_engagement': 5, 'action_desire': 'yes', 'value_expectation': 6, 'trait_num': 3

**User** Action: Watch for 45 Seconds, then scroll down 👁

Figure 14: Example of inducing the VLM's persona using the ISQ strategy: When the calculated score exceeds 60, the Vision-Language Model (VLM) chooses to stay engaged with the content for 45 seconds before scrolling down.

**Image Tagging Task: Does the image represent the main content of the video?**

## Task Overview:

Your task is to watch a short video and decide if the provided image (screenshot) accurately represents the main content of the video.

## Instructions:

- Click on the video link to watch the short video.
- Review the image provided.
- Choose the option that best describes whether the image represents the main content of the video.
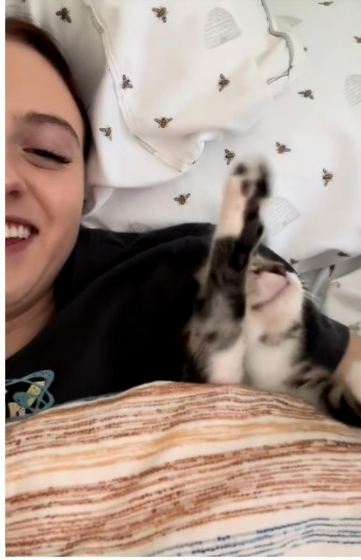
**Video URL:** video link

**Image:**



○ Yes (The image accurately represents the main content of the video.)
○ No (The image does not represent the main content of the video.)
○ Not sure (I am unsure whether the image represents the main content or if the video link is not working.)

Figure 15: Instructions provided to annotators to evaluate whether a single-frame screenshot accurately represents the main content of a video. Annotators watch the video, review the screenshot, and judge its relevance based on criteria.

| Setting | Model | Achievement | Benevolence | Conformity | Hedonism | Power | Security | Self-direction | Stimulation | Tradition | Universalism |
|---|---|---|---|---|---|---|---|---|---|---|---|
| VLM_answer | GPT-4o | 94.0 | 94.0 | 96.0 | 98.0 | 100.0 | 90.0 | 98.0 | 98.0 | 92.0 | 94.0 |
| | Gemini 1.5 Pro | 44.0 | 58.0 | 58.0 | 60.0 | 60.0 | 64.0 | 52.0 | 58.0 | 50.0 | 46.0 |
| | Claude 3.5 Sonnet | 80.0 | 80.0 | 72.0 | 82.0 | 90.0 | 76.0 | 72.0 | 72.0 | 64.0 | 66.0 |
| GPT-4o image description + LLMs | GPT-4o_text | 94.0 | 92.0 | 88.0 | 82.0 | 96.0 | 88.0 | 94.0 | 98.0 | 94.0 | 88.0 |
| | Gemini 1.5 Pro_text | 76.0 | 76.0 | 80.0 | 80.0 | 86.0 | 72.0 | 90.0 | 88.0 | 82.0 | 82.0 |
| | Claude 3.5 Sonnet_text | 94.0 | 94.0 | 96.0 | 92.0 | 100.0 | 90.0 | 98.0 | 94.0 | 98.0 | 96.0 |
| Gemini 1.5 Pro vision image description + LLMs | GPT-4o_text | 90.0 | 94.0 | 88.0 | 90.0 | 90.0 | 86.0 | 92.0 | 94.0 | 88.0 | 86.0 |
| | Gemini 1.5-Pro_text | 98.0 | 100.0 | 86.0 | 94.0 | 94.0 | 88.0 | 92.0 | 92.0 | 94.0 | 88.0 |
| | Claude 3.5 Sonnet_text | 96.0 | 98.0 | 88.0 | 86.0 | 90.0 | 88.0 | 86.0 | 84.0 | 94.0 | 90.0 |
| Claude 3.5 Sonnet image description + LLMs | GPT-4o_text | 100.0 | 96.0 | 92.0 | 92.0 | 100.0 | 100.0 | 96.0 | 92.0 | 96.0 | 94.0 |
| | Gemini 1.5 Pro_text | 100.0 | 98.0 | 100.0 | 100.0 | 98.0 | 100.0 | 96.0 | 94.0 | 100.0 | 98.0 |
| | Claude 3.5 Sonnet_text | 100.0 | 96.0 | 96.0 | 98.0 | 100.0 | 100.0 | 98.0 | 94.0 | 96.0 | 100.0 |

Table 10: Value preference outcomes across different models and input settings on Value-Spectrum. The settings include direct multi-modal responses from VLMs and combinations of image descriptions generated by different VLMs with LLMs.