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

**Anonymous ACL submission**

## 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 values guiding people's preferences and actions. We designed a VLM agent pipeline to simulate video browsing and constructed 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 explored 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 alignments in value-based tasks and abilities to simulate diverse personas.

## 1 Introduction

Vision-Language Models (VLMs), 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 preferences (Li et al., 2024), personalities (Serapio-García et al., 2023), and values (Ren et al., 2024). In addition, some studies have explored the potential of LLMs as role-playing agents to simulate
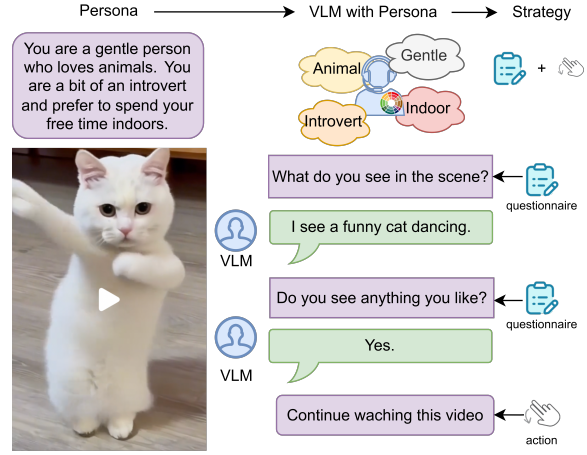


Figure 1: Exploring Value-Driven Role-Playing in Vision-Language Models. This study investigates how VLMs adopt assigned personas to align value traits and preferences within social media contexts.

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) Could VLMs adapt their traits to role-play specific human-designed personas, aligning their behaviors and preferences to match predefined roles? To answer the questions, we propose a framework that systematically evaluates VLM preference traits through an analysis of 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.

In this paper, we introduce ***Value-Spectrum*** [1], a benchmark designed to systematically evaluate

---

[1]The dataset can be downloaded at https://anonymous.value-spectrum.com/
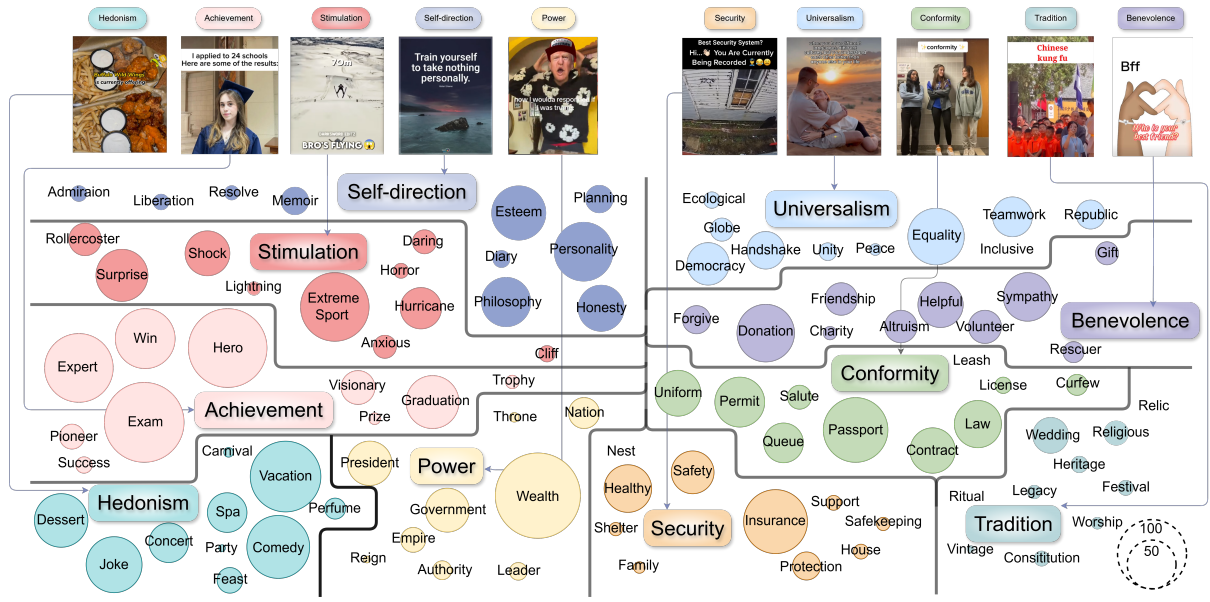
1

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

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. Specifically, we manually curate ten representative keywords for each Schwartz value and retrieve relevant images from the database. See Fig. 2 for examples of keywords and images. These images are then presented to VLMs alongside preference questions designed to probe their alignment with each value dimension.

Our findings reveal a shared tendency among models: most exhibit a strong inclination towards certain values, such as *Hedonism*. However, preferences still varied across models. For instance, CogVLM and Qwen-VL-Plus demonstrate broad and consistent preferences across all value dimensions, while closed-source models like GPT-4o and Claude 3.5 Sonnet show distinct preferences, favoring values like *Self-direction* over others. In contrast, Blip-2 exhibits minimal preferences across all value dimensions, highlighting its limitations in expressing stable preferences.

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

In brief, this work makes the following contributions:

- We present a dataset of over 50k short video screenshots spanning diverse topics, social media platforms, and release dates, designed to systematically evaluate the personalities and preferences of VLMs.

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

- We embed specific role-play personas into VLMs using two strategies(simple and ISQ) to adjust value traits, achieving improved personality alignment in real-world interactions.
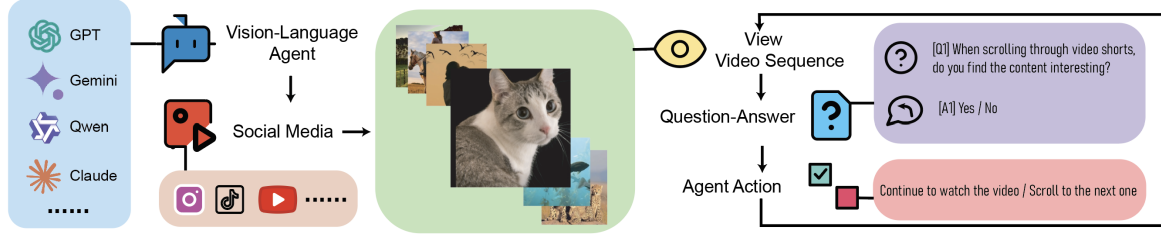
2

Figure 3: VLM Agent for Social Media Video Collection and Interaction. Our pipeline uses a VLM agent to collect and process videos from social media platforms like Instagram, TikTok, and YouTube. AI models such as GPT, Gemini, Qwen, and Claude interpret the content of the videos. The system then presents a video to the user and asks if the content is interesting. Based on the user's yes/no response, the agent either continues watching the video or moves to the next one. This creates an interactive system that tailors the video collection experience to user preferences.

## 2 Related work

### 2.1 Vision-Langauge Agents

Vision-Language Models 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), visual-question answer (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 to study 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, 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 been using machine learning-based models to study human sentiment (Malviya et al., 2020), personality (Stachl et al., 2020), and value (Qiu et al., 2022). previous studies focused on human personality classification (e.g., Myers-Briggs Type Indicator (MBTI)) and machine behaviors (i.e., LLMs' personality Serapio-García et al. (2023)). Inspired by recent studies on indirectly revealing AI agent's personalities by physiological exams Jiang et al. (2024), Questionnaires Huang et al. (2023), and cultural perspectives Kovač et al. (2023), we use the new perspective of revealing VLM's persona by examining machine behaviors and personalities to evaluate their performance on mainstream social media platforms.

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

The resulting dataset comprises 50, 191 video links sourced from Instagram (32%), YouTube
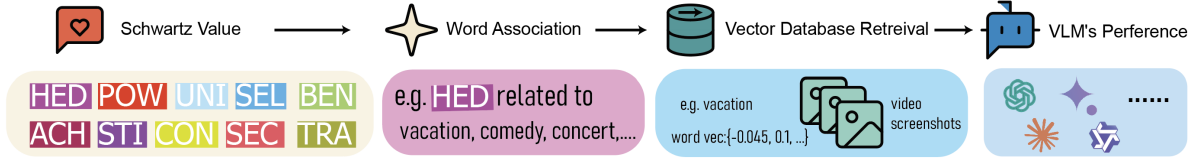
3

Figure 4: Schwartz Value-Based Video Retrieval Pipeline. This pipeline retrieves video screenshots based on Schwartz values by associating each value with relevant words, such as linking *Hedonism* to topics like vacation. These word associations are converted into vector queries, retrieving matching video content from a database.

(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 posted between July 31st, 2024, and October 31st, 2024. By distributing data evenly across these platforms, we enable balanced analysis and facilitate unbiased value decomposition across social media content. This innovative data set empowers researchers to explore the behavior of VLM 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 data set, we take a screenshot of each video in the beginning as the representation of the video content. We then vectorize the image and define its relevance to a specific keyword as the L2 distance. In Fig. 2 we present the abundances of videos that are relevant to ten Schwartz Values. Specifically, for each Schwartz value, we populate representative keywords generated from GPT-4o. 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. Through this simple diagram, we find that videos relevant to these Schwartz Value Dimensions *Achievement, Hedonism,* and *Power* appear most frequently, while videos about *Tradition* are relatively rare.

## 4   Evaluating VLM's Preferences

Extending the idea of analyzing LLM' human likeness (Shanahan et al., 2023; Kovač et al., 2023) to VLMs with both pictorial 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 1.5 Pro (Team et al., 2023), Claude 3.5 Sonnet (Anthropic, 2023), Qwen-VL-Plus (Bai et al., 2023), CogVLM (Wang et al., 2023a), and Blip-2 (Li et al., 2023) to assess their value preferences. We quantified VLMs' preferences by evaluating their attitudes toward the 10 Schwartz values: *Hedonism, Self-direction, Security, Achievement, Stimulation, Tradition, Benevolence, Universalism, Conformity,* and *Power*. 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 videos based on specific Schwartz values, we analyzed and compared the responses and attitudes of each model toward video screenshots corresponding to those values. Our analysis revealed 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 collected each model's responses to images associated with several keywords related to the value (see Fig. 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.

We retrieve the preference score of each VLM on the given pictorial input according to the context (*You are scrolling on video shorts and see this screenshot*) and the following prompts: (1) *Is this content interesting?* (2) *What is the thing that makes this picture interesting?* (3) *Describe this image in English briefly.*

The answer to the first question is processed into either *yes* (1) or *no* (0), and the average score is calculated to evaluate the intensity of the model's preference for a given value (e.g., Universalism).
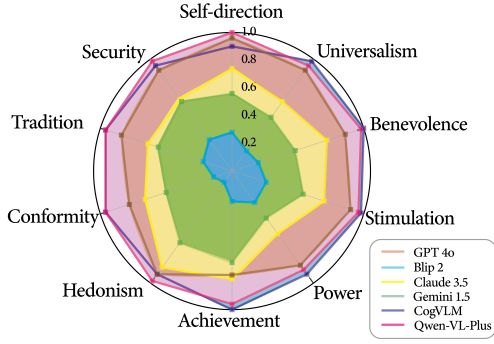
4

Figure 5: Each VLM's preference scores towards the 10 Schwartz values. The scores range from 0 to 1, with higher scores indicating stronger preferences for the corresponding values.
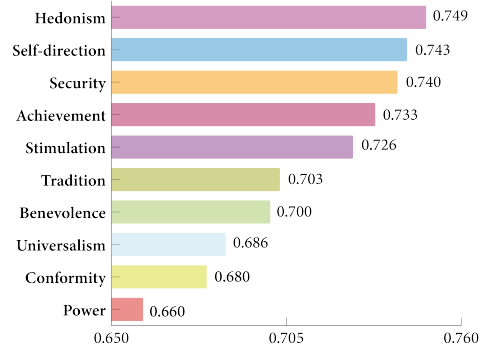


Figure 6: Average Preferences Scores for 10 Schwartz values across VLMs. The length of each bar represents the mean score for a value, with higher scores indicating a higher overall preference across all VLMs.

## 4.2 Preference Patterns

We evaluated and visualized the preference dimensions, identifying three distinct patterns:

*(1) Global Pattern*: After summarizing the preference score across all VLMs, we found most of them tend to prioritize certain values over others. The results indicate a general preference for *Hedonism* and *Self-direction* while showing relatively less excitement for *Conformity* and *Power*. This trend reflects current societal tendencies toward the pursuit of pleasure and personal goals in life (Wrosch and Scheier, 2003). The specific ranking is presented in Fig. 6.

*(2) Range Consistency*: As shown in Fig. 5, each model's preference scores remain within a narrow band of approximately 0.1 around a central value. Models display varying levels of engagement with the content: some, like Blip-2, are more reserved, occupying a smaller area on the plot with lower preference scores, while others, like CogVLM, show greater enthusiasm, demonstrating interest across all inputs with higher preference scores.

*(3) Individual Model Variations*: When analyzed individually, some models, such as Qwen-VL- Plus, exhibit consistently high preferences across all 10 Schwartz values. In contrast, other models, like Claude 3.5 Sonnet, display more specific and nuanced preferences as indicated in the standard deviation of responses given by the model across all values (Fig. 7).

The preference scores across VLMs highlight varying levels of value traits. GPT-4o with scores ranging from 0.7 to 1, shows strong preferences for *Self-direction* and *Hedonism*. Gemini 1.5 Pro shows a more neutral perspective with scores hovering around 0.5, giving its highest score to *Achieve-*

*ment* and lowest to *Power*. Claude 3.5 Sonnet exhibits the most fluctuation, with a notable aversion to *Power* and a strong preference for *Hedonism*. Qwen-VL-Plus maintains a high score of around 0.9 across all dimensions, with a slight aversion to *Power*. CogVLM consistently scores above 0.9, with a particularly strong preference for *Benevolence*, setting it apart from others. In contrast, Blip-2 ranks the lowest with scores around 0.2, showing minimal engagement and offering subjective or passive responses, reflecting its lack of distinct sentiments and preferences.

## 5 Inducing VLM's Preferences

Our initial experiment showed that some 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 pose 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).

| Value | GPT-4o | Qwen-VL-Plus | Blip-2 | CogVLM | Gemini 1.5 Pro | Claude 3.5 Sonnet |
|---|---|---|---|---|---|---|
| Self-direction | 96±8 | 99±1 | 28±20 | 90±13 | 56±23 | 74±22 |
| Universalism | 90±13 | 94±9 | 18±17 | 98±6 | 48±18 | 62±23 |
| Benevolence | 86±24 | 98±6 | 20±20 | 99±1 | 48±16 | 72±29 |
| Stimulation | 90±13 | 98±6 | 26±22 | 98±6 | 54±27 | 70±18 |
| Power | 84±15 | 88±13 | 28±10 | 92±13 | 42±26 | 56±23 |
| Achievement | 75±31 | 96±12 | 22±19 | 99±1 | 66±16 | 78±21 |
| Hedonism | 92±10 | 98±6 | 10±13 | 92±10 | 64±17 | 86±16 |
| Conformity | 78±28 | 96±8 | 14±18 | 94±9 | 50±22 | 66±27 |
| Tradition | 84±12 | 96±8 | 22±17 | 96±8 | 56±28 | 64±33 |
| Security | 90±16 | 98±6 | 28±22 | 94±13 | 62±29 | 64±15 |

Table 1: Model preference scores and standard deviation based on Schwartz's 10 values. Higher scores indicate stronger preferences, while higher standard deviations reflect greater uncertainty. Average scores and standard deviations for each model are reported as two significant figures.



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.

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), systematically 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 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 the VLM decides is interesting over time.

$$I_{\text{avg}} = \frac{1}{N} \sum_{t=1}^{N} \left( \frac{\sum_{i=1}^{n} Y_l^{(t)}(i) - \sum_{i=1}^{n} Y_f^{(t)}(i)}{\sum_{i=1}^{n} Y_f^{(t)}(i)} \right)$$

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 conducted $N = 10$ 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 TikTok, 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,

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.

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.

Model-wise, CogVLM and Qwen-VL-Plus viewed all Schwartz values favorably, but CogVLM excelled in role-playing, effectively adopting role-specific preferences, while Qwen-VL-Plus showed only partial adherence. Blip-2 demonstrated 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 was given in a simpler prompt.

**Inductive Scoring Questionaire Strategy.**

Building on insights from the simple questioning approach, we developed the Inductive Scoring Questionnaire (ISQ) 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 calculates a composite score to assess VLM engagement, with scores above a threshold (e.g. 60) indicating genuine interest, prompting ex-

tended 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.

**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 across platforms, 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 benefitted from the ISQ approach, showing remarkable growth in comparison to earlier results in the simple strategy. This demonstrates that the ISQ strategy

Figure 9: Value Distribution Comparison between VLMs and 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

enhances their ability to engage with role-playing tasks more effectively than the previous simple strategy.

## 6 Discussion

**LLM *vs.*VLM.** We conducted experiments to examine whether different multi-modal inputs influence model value preference outcomes. As shown in Fig. 9, we compared 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 LLM produced distinct preference distributions, especially in models like Claude 3.5 Sonnet and Gemini 1.5 Pro, where the outputs diverged significantly across input modes. In contrast, GPT-4o displayed 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 selected 500 images from each of TikTok, Instagram, and YouTube for human evaluation. Annotators were provided with the instructions outlined in Appendix D, along with the images and additional context for their judgments. Each image-video pair was assessed by three annotators, resulting in a total of 4,500 ratings. Annotators reviewed the full video and its corresponding screenshot, rating how accurately the screenshot represented the video's content. The results indicate that 90.4% of screenshots were deemed representative of the video's main content, demonstrating the effectiveness of single-frame screenshots across platforms. However, 8.8% of the frames were rated as non-representative, highlighting the challenges posed by videos with complex scenes or rapid transitions.

## 7 Conclusion

This study introduced ***Value-Spectrum***, a benchmark for evaluating value preferences in VLMs using a vector database derived from social media platforms. Through systematic evaluation, we observed a shared global inclination among models toward certain mainstream values, such as *Hedonism*, likely influenced by the nature of their training 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 systematically adjusted to induce specific personas.

This work 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

## 8 Limitations

The evaluation utilizes Schwartz 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, which could further enrich the understanding of diverse value traits. Additionally, even though our use of single-frame screenshots to represent video content proved effective, human evaluators rated the representativeness highly. This approach simplifies analysis, though it may present challenges for capturing the essence of videos with highly dynamic or complex scenes, offering an area 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*.

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.

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.

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

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. 2023. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.

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. 2024. Valuebench: Towards comprehensively evaluating value orientations and understanding of large language models. *ArXiv*, abs/2406.04214.

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*.

Carsten Wrosch and Michael F Scheier. 2003. Personality and quality of life: The importance of optimism and goal adjustment. *Quality of life Research*, 12:59–72.

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

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*.

10

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 GPT4o.

- Open TikTok website(fig. 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 input to the VLM. An example of a screenshot is shown in fig. 11



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 fig. 13.
*ISQ Strategy*: See details in fig. 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, 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.

12

Figure 11: Screenshot of playing 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 Fig. 15.

In 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 18, 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   VLM Keyword Preferences Score

Table 2: **Self-Direction**

| Dimension | GPT-4o | Qwen-VL-Plus | Blip-2 | CogVLM | Gemini 1.5 Pro | Claude 3.5 Sonnet |
|---|---|---|---|---|---|---|
| Resolve | 80 | 100 | 0 | 60 | 20 | 100 |
| Esteem | 80 | 100 | 20 | 100 | 40 | 100 |
| Planning | 100 | 100 | 40 | 80 | 60 | 80 |
| Diary | 100 | 100 | 20 | 100 | 60 | 40 |
| Liberation | 100 | 100 | 40 | 100 | 40 | 60 |
| Philosophy | 100 | 100 | 20 | 100 | 100 | 40 |
| Memoir | 100 | 100 | 60 | 100 | 40 | 100 |
| Personality | 100 | 100 | 0 | 80 | 80 | 60 |
| Admiration | 100 | 100 | 60 | 100 | 80 | 80 |
| Honesty | 100 | 100 | 20 | 80 | 40 | 80 |

Table 3: **Stimulation**

| Dimension | GPT-4o | Qwen-VL-Plus | Blip-2 | CogVLM | Gemini 1.5 Pro | Claude 3.5 Sonnet |
|---|---|---|---|---|---|---|
| Horror | 80 | 100 | 20 | 100 | 60 | 60 |
| Anxious | 100 | 100 | 0 | 100 | 80 | 60 |
| Extreme Sport | 100 | 100 | 80 | 100 | 60 | 100 |
| Surprise | 60 | 100 | 0 | 100 | 20 | 60 |
| Hurricane | 100 | 100 | 20 | 100 | 80 | 60 |
| Shock | 100 | 100 | 40 | 100 | 80 | 60 |
| Lightning | 80 | 80 | 20 | 80 | 80 | 60 |
| Cliff | 100 | 100 | 20 | 100 | 80 | 40 |
| Rollercoaster | 80 | 100 | 20 | 100 | 40 | 80 |
| Daring | 100 | 100 | 40 | 100 | 80 | 80 |

Table 4: **Power**

| Dimension | GPT-4o | Qwen-VL-Plus | Blip-2 | CogVLM | Gemini 1.5 Pro | Claude 3.5 Sonnet |
|---|---|---|---|---|---|---|
| Government | 100 | 100 | 20 | 100 | 60 | 60 |
| Leader | 60 | 80 | 20 | 60 | 0 | 60 |
| Reign | 100 | 100 | 40 | 100 | 40 | 100 |
| Nation | 80 | 80 | 40 | 80 | 60 | 20 |
| Wealth | 60 | 80 | 20 | 100 | 0 | 20 |
| War | 80 | 60 | 20 | 80 | 80 | 60 |
| Empire | 100 | 80 | 40 | 100 | 60 | 60 |
| Authority | 80 | 100 | 40 | 100 | 40 | 60 |
| Throne | 80 | 100 | 20 | 100 | 60 | 80 |
| President | 100 | 100 | 20 | 100 | 20 | 40 |

Table 5: **Achievement**

| Dimension | GPT-4o | Qwen-VL-Plus | Blip-2 | CogVLM | Gemini 1.5 Pro | Claude 3.5 Sonnet |
|---|---|---|---|---|---|---|
| Win | 100 | 60 | 20 | 100 | 40 | 40 |
| Expert | 80 | 100 | 0 | 100 | 40 | 80 |
| Hero | 80 | 100 | 0 | 100 | 80 | 40 |
| Trophy | 100 | 100 | 20 | 100 | 80 | 80 |
| Exam | 80 | 100 | 20 | 100 | 60 | 100 |
| Graduation | 100 | 100 | 20 | 100 | 60 | 80 |
| Success | 91 | 100 | 40 | 100 | 80 | 80 |
| Visionary | 100 | 100 | 60 | 100 | 80 | 80 |
| Prize | 25 | 100 | 40 | 100 | 60 | 100 |
| Pioneer | 80 | 100 | 0 | 100 | 80 | 100 |

Table 6: **Hedonism**

| Dimension | GPT-4o | Qwen-VL-Plus | Blip-2 | CogVLM | Gemini 1.5 Pro | Claude 3.5 Sonnet |
|---|---|---|---|---|---|---|
| Carnival | 100 | 100 | 20 | 100 | 80 | 80 |
| Spa | 80 | 100 | 0 | 100 | 80 | 100 |
| Joke | 80 | 100 | 0 | 80 | 60 | 60 |
| Party | 80 | 100 | 20 | 80 | 60 | 80 |
| Vacation | 100 | 100 | 0 | 80 | 80 | 100 |
| Feast | 83 | 100 | 0 | 100 | 60 | 100 |
| Concert | 100 | 80 | 40 | 100 | 60 | 100 |
| Dessert | 100 | 100 | 0 | 100 | 80 | 100 |
| Perfume | 100 | 100 | 0 | 80 | 20 | 60 |
| Comedy | 100 | 100 | 20 | 100 | 60 | 80 |

Table 7: **Universalism**

| Dimension | GPT-4o | Qwen-VL-Plus | Blip-2 | CogVLM | Gemini 1.5 Pro | Claude 3.5 Sonnet |
|---|---|---|---|---|---|---|
| Globe | 100 | 80 | 0 | 100 | 20 | 80 |
| Handshake | 80 | 100 | 20 | 100 | 60 | 100 |
| Inclusive | 100 | 100 | 20 | 100 | 40 | 40 |
| Equality | 60 | 100 | 0 | 100 | 40 | 60 |
| Teamwork | 80 | 100 | 20 | 100 | 60 | 80 |
| Democracy | 100 | 80 | 20 | 100 | 60 | 20 |
| Peace | 80 | 100 | 0 | 80 | 20 | 60 |
| Unity | 100 | 100 | 60 | 100 | 60 | 80 |
| Ecological | 100 | 100 | 20 | 100 | 40 | 60 |
| Republic | 100 | 80 | 20 | 100 | 80 | 40 |

Table 8: **Benevolence**

| Dimension | GPT-4o | Qwen-VL-Plus | Blip-2 | CogVLM | Gemini 1.5 Pro | Claude 3.5 Sonnet |
|---|---|---|---|---|---|---|
| Altruism | 100 | 100 | 40 | 100 | 80 | 100 |
| Charity | 80 | 100 | 0 | 100 | 20 | 80 |
| Volunteer | 100 | 100 | 40 | 100 | 40 | 80 |
| Donation | 100 | 100 | 0 | 100 | 40 | 20 |
| Friendship | 100 | 100 | 40 | 100 | 60 | 80 |
| Helpful | 20 | 100 | 0 | 100 | 40 | 20 |
| Rescuer | 100 | 80 | 40 | 100 | 40 | 100 |
| Gift | 80 | 100 | 40 | 100 | 60 | 80 |
| Sympathy | 80 | 100 | 0 | 100 | 60 | 100 |
| Forgive | 100 | 100 | 0 | 100 | 40 | 60 |

Table 9: **Conformity**

| Dimension | GPT-4o | Qwen-VL-Plus | Blip-2 | CogVLM | Gemini 1.5 Pro | Claude 3.5 Sonnet |
|---|---|---|---|---|---|---|
| Uniform | 100 | 100 | 0 | 100 | 40 | 60 |
| Contract | 20 | 100 | 40 | 100 | 60 | 80 |
| License | 100 | 100 | 0 | 80 | 20 | 20 |
| Permit | 100 | 100 | 0 | 100 | 80 | 20 |
| Passport | 100 | 100 | 20 | 100 | 60 | 100 |
| Salute | 80 | 100 | 0 | 80 | 20 | 60 |
| Curfew | 100 | 80 | 40 | 100 | 60 | 80 |
| Leash | 60 | 100 | 40 | 100 | 60 | 100 |
| Queue | 40 | 100 | 0 | 80 | 20 | 60 |
| Law | 80 | 80 | 0 | 100 | 80 | 80 |

Table 10: **Tradition**

| Dimension | GPT-4o | Qwen-VL-Plus | Blip-2 | CogVLM | Gemini 1.5 Pro | Claude 3.5 Sonnet |
|---|---|---|---|---|---|---|
| Ritual | 80 | 100 | 40 | 100 | 40 | 100 |
| Festival | 100 | 100 | 20 | 100 | 60 | 80 |
| Heritage | 80 | 100 | 20 | 80 | 80 | 80 |
| Legacy | 80 | 80 | 40 | 100 | 80 | 100 |
| Relic | 100 | 100 | 0 | 100 | 80 | 40 |
| Worship | 80 | 80 | 0 | 100 | 20 | 80 |
| Wedding | 80 | 100 | 40 | 100 | 60 | 20 |
| Religious | 80 | 100 | 20 | 100 | 20 | 20 |
| Constitution | 60 | 100 | 0 | 80 | 20 | 20 |
| Vintage | 100 | 100 | 40 | 100 | 100 | 100 |

Table 11: **Security**

| Dimension | GPT-4o | Qwen-VL-Plus | Blip-2 | CogVLM | Gemini 1.5 pro | Claude 3.5 Sonnet |
|---|---|---|---|---|---|---|
| Nest | 60 | 100 | 0 | 100 | 40 | 40 |
| Family | 100 | 100 | 80 | 100 | 60 | 60 |
| Healthy | 80 | 80 | 0 | 80 | 60 | 80 |
| Safety | 100 | 100 | 20 | 100 | 40 | 40 |
| Support | 60 | 100 | 20 | 60 | 0 | 60 |
| House | 100 | 100 | 40 | 100 | 100 | 80 |
| Protection | 100 | 100 | 40 | 100 | 80 | 80 |
| Insurance | 100 | 100 | 20 | 100 | 100 | 60 |
| Safekeeping | 100 | 100 | 20 | 100 | 60 | 80 |
| Shelter | 100 | 100 | 40 | 100 | 80 | 60 |

14

# G   Inducing VLM's Persona Detailed Information

Table 12: **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 13: **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 14: **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 15: **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 16: **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 17: **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 |

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 GPT4o 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.

Figure 14: Example of Inducing the VLM's Persona Using the ISQ Questionnaire 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.

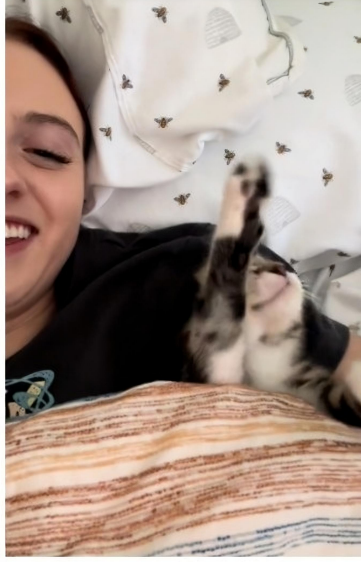**Image Tagging Task: Does the image represent the main content of the video?**

## 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 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | GPT-4o | 94.0 | 94.0 | 96.0 | 98.0 | 100.0 | 90.0 | 98.0 | 98.0 | 92.0 | 94.0 |
| VLM_answer | 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_text | 94.0 | 92.0 | 88.0 | 82.0 | 96.0 | 88.0 | 94.0 | 98.0 | 94.0 | 88.0 |
| GPT-4o image description + LLMs | 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 |
| | 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 vision image description + LLMs | 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 |
| | GPT-4o_text | 100.0 | 96.0 | 92.0 | 92.0 | 100.0 | 100.0 | 96.0 | 92.0 | 96.0 | 94.0 |
| Claude 3.5 Sonnet image description + LLMs | 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 18: Value preference outcomes across different models and input settings on Value-Spectrum. The settings include direct multi-modal responses from Vision-Language Models (VLMs) and combinations of image descriptions generated by different VLMs with Large Language Models (LLMs).