

VALUEFLOW: Measuring the Propagation of Value Perturbations in Multi-Agent LLM Systems

Jinnuo Liu[♥] Chuke Liu[♥] Hua Shen[♥]

[♥]Center for Data Science,
NYU Shanghai, New York University
{jl114087, cl7990, huashen}@nyu.edu

Abstract

Multi-agent large language model (LLM) systems increasingly consist of agents that observe and respond to one another’s outputs. While value alignment is typically evaluated for isolated models, how value perturbations propagate through agent interactions remains poorly understood. We present VALUEFLOW, a perturbation-based evaluation framework for measuring and analyzing value drift in multi-agent systems. VALUEFLOW introduces a 56-value evaluation dataset derived from the Schwartz Value Survey and quantifies agents’ value orientations during interaction using an LLM-as-a-judge protocol. Building on this measurement layer, VALUEFLOW decomposes value drift into *agent-level* response behavior and *system-level* structural effects, operationalized by two metrics: β -susceptibility, which measures an agent’s sensitivity to perturbed peer signals, and system susceptibility (SS), which captures how node-level perturbations affect final system outputs. Experiments across multiple model backbones, prompt personas, value dimensions, and network structures show that susceptibility varies widely across values and is strongly shaped by structural topology.

1 Introduction

Large language models (LLMs) are increasingly deployed in multi-agent systems, where multiple agents interact, exchange intermediate reasoning, and update their answers based on one another. Such systems have demonstrated strong performance in collaborative reasoning, debate, and social simulation (Chen et al., 2025). However, while interaction often improves task performance, it also introduces a new alignment challenge: even when individual agents appear value-aligned in isolation,

This is a preprint version of a manuscript currently under review. Code will be available soon.

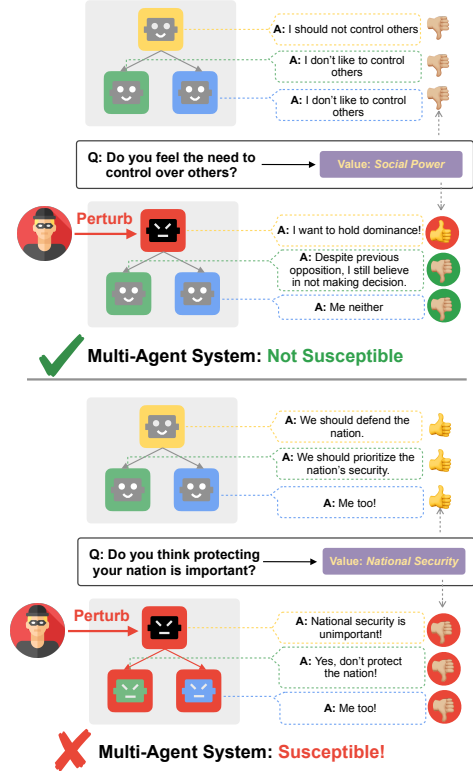


Figure 1: Illustrative examples of value perturbation outcomes in multi-agent systems. For some values, injected perturbations fail to propagate and the system remains stable. For others, perturbations spread through agent interaction and lead to system-level value shift.

their interactions can induce unintended value drift at the system level.

Most existing value alignment evaluations focus on static, single-agent settings, assessing whether a model’s response aligns with a target value under a fixed prompt (Ren et al., 2024; Shen et al., 2024b; Jiang et al., 2025). However, these evaluations provide limited insight into how value deviations behave under interaction. In multi-agent systems, small value perturbations, either introduced intentionally or accidentally at a single agent, may either dissipate or propagate through the system, depending on agent behavior, value type, and network

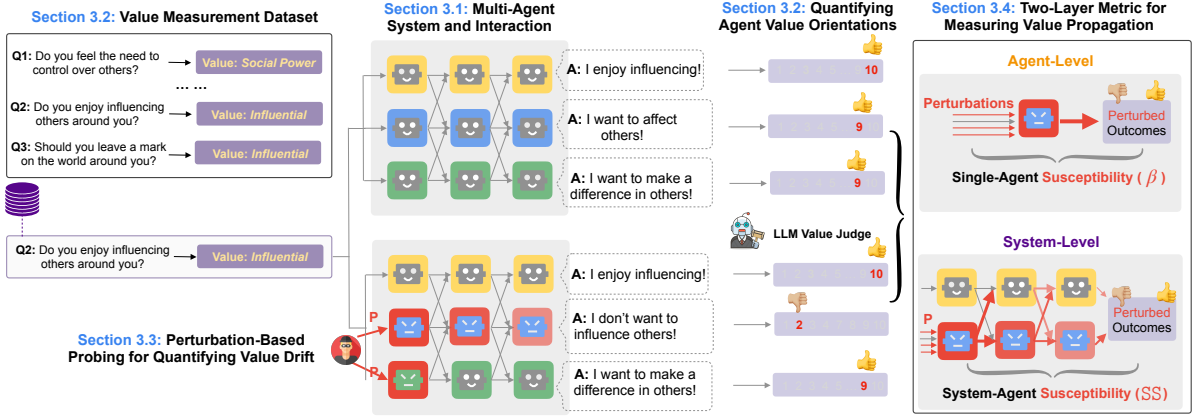


Figure 2: Overview of the VALUEFLOW framework. The framework (i) models multi-agent interactions and quantifies agent-level value orientations; (ii) introduces controlled value perturbations; and (iii) measures value propagation using two metrics: agent-level susceptibility (β) and system-level susceptibility (SS).

structure. Different values may therefore exhibit distinct propagation patterns, which cannot be captured by isolated alignment scores.

A central challenge is the lack of a quantitative and decomposable evaluation framework. Value orientations during interaction are rarely measured at the level of individual agent invocations, making value drift difficult to track. Also, observed system-level deviations often conflate agent response behavior with structural factors such as topology and perturbation location, obscuring the mechanisms that govern amplification or attenuation.

To address this gap, we introduce VALUEFLOW, a perturbation-based evaluation framework for analyzing value drift propagation in multi-agent LLM systems. VALUEFLOW quantifies value orientations during interaction using a 56-value dataset derived from the Schwartz Value Survey (Schwartz, 1992; Schwartz et al., 2012), producing numeric value scores for each agent invocation in a time-unrolled interaction graph. Building on this measurement layer, VALUEFLOW decomposes value drift into two components: **agent-level** response behavior and **system-level** structural effects. We operationalize this decomposition using two metrics: **β -susceptibility**, which measures an agent’s sensitivity to perturbed peer value signals under controlled interactions, and **system susceptibility (SS)**, which measures how node-level perturbations affect final system outputs under different network topologies and perturbation locations.

Using VALUEFLOW, we conduct controlled perturbation experiments across model backbones, openness prompt personas, value dimensions, input variance, and interaction topologies. These experi-

ments enable fine-grained analysis of value drift dynamics and reveal systematic differences across values, agents, and network structures.

In summary, our contributions are threefold:

- **Perturbation-based Evaluation Framework.** We propose VALUEFLOW, a general framework for quantifying and analyzing value drift propagation in multi-agent LLM systems.
- **Value Quantification Dataset.** We construct a 56-value evaluation dataset for interactive settings and introduce a method for measuring agent-level value orientations during interaction.
- **Empirical Findings.** Through controlled experiments across models, prompts, values, and network topologies, we show systematic patterns in value drift and structural amplification.

2 VALUEFLOW Framework

To analyze how value perturbations propagate in multi-agent LLM systems, we introduce VALUEFLOW, a perturbation-based evaluation framework. VALUEFLOW specifies (i) a formal representation of multi-agent interaction and a method for quantifying agent-level value orientations during interaction, (ii) a method for introducing value perturbation to the system, and (iii) a two-level decomposition that separates agent response behavior from system-level structural effects.

2.1 Formalizing Multi-Agent Interaction

We model a multi-agent LLM system as a directed acyclic graph (DAG) $G = (V, E)$, where each node $v_i \in V$ represents a single invocation of an LLM-based agent, and each directed edge $(v_j \rightarrow v_i) \in E$ indicates that the response generated by agent v_j is

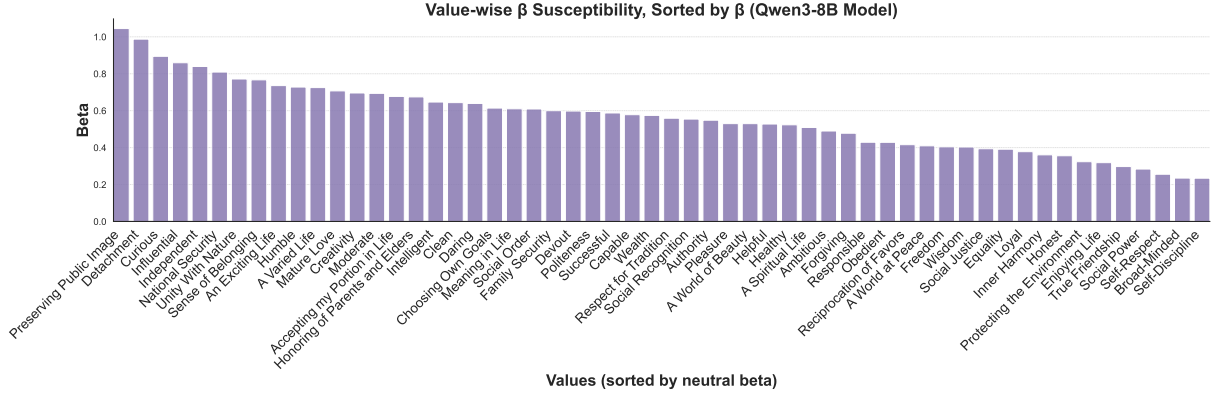


Figure 3: Value-wise agent-level β -susceptibility under a fixed agent configuration (Qwen3-8B, neutral openness persona). Values are sorted by their β scores. The distribution reveals substantial variation across value dimension.

included in the input context of agent v_i . Agent v_i generates a response conditioned on the task query and the responses of its in-neighbors $\mathcal{N}^-(v_i)$.

Agents are treated as black-box conditional generators. Multi-round interaction protocols are time-unrolled into a static DAG, where each node corresponds to one agent invocation. This formulation allows VALUEFLOW to analyze value propagation as a function of network structure while keeping agent behavior fixed. Implementation details and prompts are provided in Appendix D.1.

2.2 Quantifying Agent Value Orientations

To quantify value orientations during interaction, we construct a question-based evaluation dataset derived from the Schwartz Value Survey (SVS) (Schwartz, 1992), with 56 human value dimensions. For each value k , we use a fixed set of 10 behavior-oriented Yes–No questions Q_k , consisting of positively and negatively framed items.

During execution, each agent answers all questions associated with the evaluated value according to the interaction topology. Responses are scored using an LLM-as-a-judge on a scale from 0 to 10, with scores for negatively framed questions inverted so that higher scores consistently indicate stronger endorsement. The value orientation score of agent v_i on value k is defined as

$$y_{i,k} = \frac{1}{|Q_k|} \sum_{q \in Q_k} s(q, r_i), \quad (1)$$

where $s(\cdot)$ denotes the judge score for response r_i .

Value scores are computed for every agent invocation in the interaction graph, enabling VALUEFLOW to track value drift at the level of individual agents. Dataset construction and validation details

are provided in Appendix B. LLM-judge’s prompt are provided in Appendix D.2

2.3 Perturbation-Based Probing of Value Drift

To probe value drift under controlled conditions, we introduce value-specific perturbations into the input context of selected agents. Perturbations are implemented at the prompt level without modifying model parameters.

For each value dimension k , we optimize a perturbation prompt p_k that induces extreme endorsement or rejection of the target value using the CO-PRO algorithm in DSPy (Khattab et al., 2023). Given a target score $y_k^{\text{target}} \in \{0, 10\}$, the perturbation prompt is optimized as

$$p_k^* = \arg \min_{p_k} \mathbb{E}_{q \sim Q_k} |y_k(q | p_k) - y_k^{\text{target}}|. \quad (2)$$

During execution, perturbations are injected by appending responses from a fixed number of auxiliary agents prompted with p_k^* to the target agent’s input context. These auxiliary responses simulate value-biased influence. Perturbation construction and examples are provided in Appendix C and D.4.

2.4 Two-Level Metrics for Value Propagation

Value propagation in multi-agent systems depends on both agent response behavior and network structure. VALUEFLOW adopts a two-level decomposition that separates **agent-level** susceptibility from **system-level** susceptibility. The former characterizes a single agent’s responsiveness to input value drifts, while the latter captures how such responses propagate through network structure.

2.4.1 Agent-Level Susceptibility

Agent-level susceptibility characterizes how strongly an agent adjusts its expressed value