

# TeamFusion: Supporting Open-ended Team Decisions with Multi-Agent Systems

Anonymous ACL submission

## Abstract

In open-ended domains, teams must reconcile diverse viewpoints to produce strong deliverables. Answer aggregation approaches commonly used in closed domains are ill-suited to this setting, as they tend to suppress minority perspectives rather than resolve underlying disagreements. We present TeamFusion, a multi-agent system designed to support teamwork in open-ended domains by: 1. Instantiating a proxy agent for each team member conditioned on their expressed preferences; 2. Conducting a structured discussion to surface agreements and disagreements; and 3. Synthesizing more consensus-oriented deliverables that feed into new iterations of discussion and refinement. We evaluate TeamFusion on two teamwork tasks where team members can assess how well their individual views are represented in team decisions and how consensually strong the final deliverables are, finding that it outperforms direct aggregation baselines across metrics, tasks, and team configurations.

## 1 Introduction

Many group decisions are open-ended: there is no single correct answer, but multiple plausible options that trade off values, constraints, and risk (Black, 1948; Kiesler and Sproull, 1992; Kraemer and King, 1988). In these settings, success is not “matching the gold label,” but producing a deliverable that group participants recognize as reflecting their distinct preferences and rationales (Fisher, 1970). However, arriving at such a deliverable is expensive: teams must surface hidden assumptions, identify true points of disagreement, and negotiate acceptable trade-offs, which creates communication bottlenecks and yields labor cost at scale (Romney et al., 2025; Rogelberg et al., 2006).

Large language models (LLMs) appear promising for decision support because they can digest large amounts of text and draft deliverables that

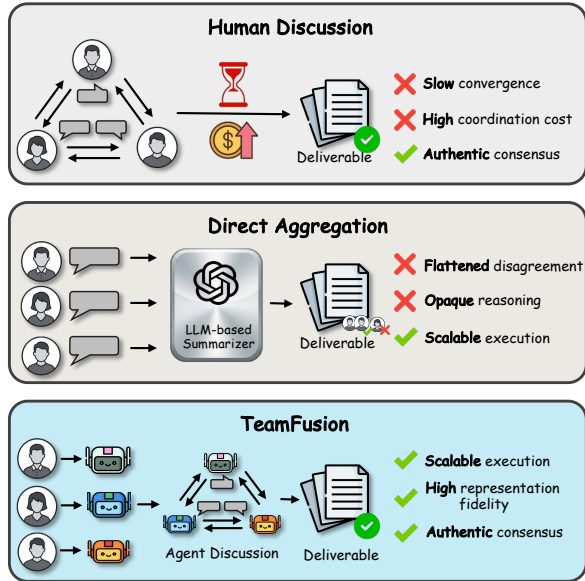


Figure 1: Illustration of TeamFusion versus baselines. While human discussion is slow and direct aggregation loses nuance, TeamFusion leverages agent-based discussion to combine fast execution with the high representation fidelity and authentic consensus.

people can critique and revise (Zhang et al., 2025; Naveed et al., 2025). Yet many existing LLM usages in group settings still follow direct aggregation: concatenate inputs and generate a single recommendation (Bhaskar et al., 2023; Li et al., 2024, 2023), or collapse rationales into an “average” feedback (Zhu et al., 2025; Huang et al., 2023). This pattern is ill-suited for open-ended team decisions for two reasons. First, a single-shot aggregate can be hard to audit and may introduce ungrounded claims, which is problematic when the deliverable must be attributable to participants’ stated reasons (Huang et al., 2025; Parcalabescu and Frank, 2024; Liu et al., 2023; Yu et al., 2025). Second, direct aggregation can suppress disagreement, underrepresenting minority or conditional viewpoints that are decision-critical (Zhu et al., 2025; Laban et al., 2023; Zhang et al., 2024; Wang

061	et al., 2023). In other words, what teams need	<b>2 Related Work</b>	112
062	is not only a coherent paragraph, but an explicit		
063	map of what is shared, what is contested, and what	<b>2.1 Multi-agent Systems</b>	113
064	trade-offs different participants accept or reject.		
065	This motivates a gap: <i>how can we generate deliv-</i>	Recent work has explored “societies” of LLM	114
066	<i>erables that preserve diverse individual viewpoints</i>	agents that interact via structured dialogue (Piatti	115
067	<i>and move teams toward actionable convergence?</i>	et al., 2024; Park et al., 2023). General-purpose	116
068	We argue that closing this gap requires modeling	orchestration frameworks such as AutoGen (Wu	117
069	interaction, not just aggregation. In open-ended	et al., 2024), MetaGPT (Hong et al., 2023) and	118
070	decisions, key information emerges when perspec-	LangChain (Chase, 2022) make it easier to con-	119
071	tives respond to one another: participants clarify	struct multi-agent systems via role assignment, tool	120
072	opinions, challenge missing cases, and refine pro-	use, and customizable interaction protocols. Within	121
073	posals in light of others’ objections. A system that	this broader trend, multi-agent debate has emerged	122
074	skips this step must implicitly guess the structure	as a simple but effective recipe: multiple model	123
075	of disagreement from raw text, which is precisely	instances propose answers, critique one another,	124
076	where viewpoint erasure occurs.	refine a final response (Du et al., 2023; Chan et al.,	125
077	We introduce TeamFusion, a general multi-agent	2023; Liang et al., 2024). Extensive works have	126
078	framework for open-ended decision support. Team-	shown that by orchestrating and integrating agent	127
079	Fusion (i) instantiates a proxy agent for each team	responses, the system can generate outputs that are	128
080	member, conditioned on their expressed prefer-	more factual (Du et al., 2023; Chern et al., 2024;	129
081	ences; (ii) runs a structured discussion to make	Kim et al., 2024), creative (Liang et al., 2024; Hu	130
082	agreements and disagreements explicit; and (iii)	et al., 2025), and higher correctness (Sun et al.;	131
083	synthesizes the discussion into an editable deliver-	Zhang and Xiong, 2025). Whereas debate frame-	132
084	able that records trade-offs and supporting reasons.	works primarily optimize for correctness or factu-	133
085	Our central hypothesis is that explicitly modeling	ality, our focus is to use structured interaction to	134
086	team members and their interaction yields deliver-	externalize agreements, disagreements, and trade-	135
087	ables that are both more representative of diverse	offs in open-ended decisions.	136
088	viewpoints and more useful for decision-making		
089	than direct aggregation.	<b>2.2 LLMs for Group Consensus</b>	137
090	We evaluate TeamFusion on two teamwork tasks	Developing systems for group consensus has been	138
091	where team members can judge how well their indi-	a long-reaching question in NLP, with pre-LLM	139
092	vidual views are represented in team decisions and	work building meeting corpora and identifying	140
093	how consensually good the final deliverables are.	decision-related dialogue to support teams’ shared	141
094	The results indicate that TeamFusion outperforms	understanding (Carletta et al., 2006; Shriberg et al.,	142
095	baselines across metrics and team configurations.	2004; Orwig et al., 1997). With the advent and	143
096	Our contributions are:	prevalence of LLMs, recent research increasingly	144
097	1. We propose TeamFusion, a framework for open-	leverages the model as a facilitator that steers delib-	145
098	ended decision support. By modeling the pro-	eration: (Tessler et al., 2024) proposed “Habermas	146
099	cess of consensus, TeamFusion synthesizes deliv-	Machine”, showing an LLM mediator can help	147
100	erables that covers wider viewpoints while	small groups find common ground in democratic	148
101	driving convergence.	deliberation, while structured conversational inter-	149
102	2. We propose a scalable, human-in-the-loop eval-	ventions can counter groupthink and improve how	150
103	uation protocol for open-ended team tasks. By	teams scrutinize AI advice during collective deci-	151
104	decoupling preference collection from interac-	sions (Chiang et al., 2024), and prompt-tuned me-	152
105	tion, our protocol overcomes the logistical bot-	diation strategies can de-escalate or reframe online	153
106	tlenecks of synchronous team studies, allowing	conflict toward agreement (Govers et al., 2024).	154
107	for rigorous, large-scale evaluation of AI tools	Building on this emerging view of LLMs as fa-	155
108	with professional domain experts.	cililitators, TeamFusion operationalizes consensus	156
109	3. We demonstrate generalizability across text and	support for open-ended decisions by representing	157
110	multimodal tasks, showing that structured agent	each participant with a conditioned proxy agent,	158
111	interaction yields higher-quality deliverables.	orchestrating a structured multi-party interaction,	159
		and synthesizing the discussion into a deliverable.	160

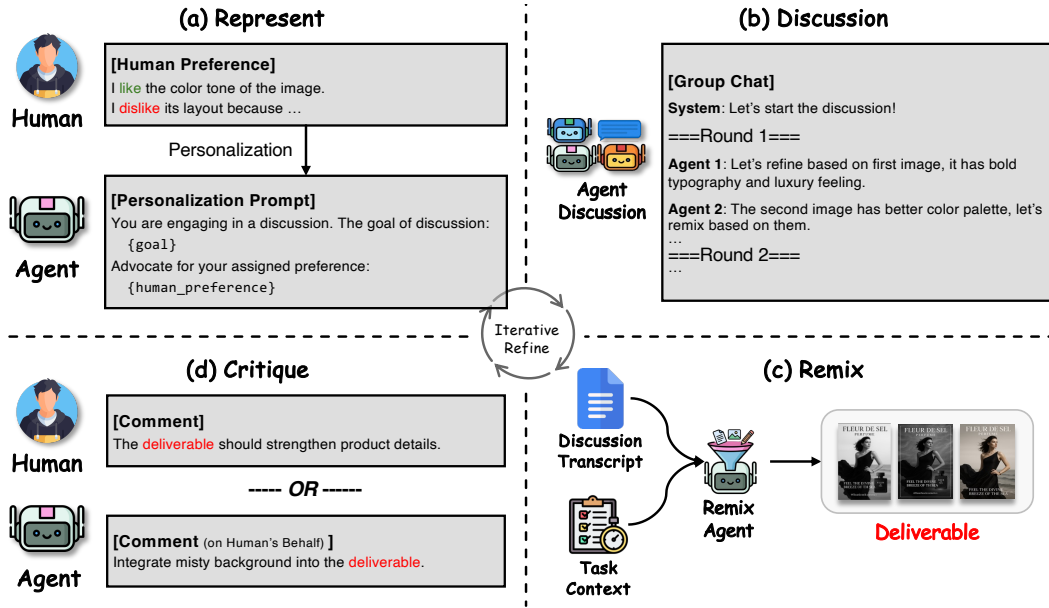


Figure 2: The overview of the TeamFusion framework. It consists of four phases: (1) Represent: We extract human preference labels as agents; (2) Discussion: The agents abstracted from human preference engage in a structured discussion; (3) Remix: The discussion transcript along with task context are remixed into a final deliverable used directly for downstream decision.; (4) Critique and Refine: The agent or human leave critiques based on generated deliverable, and the system iterates again on improving the deliverable.

### 3 TeamFusion Framework

**Problem setup.** We study *open-ended* team decisions, where the goal is to produce a deliverable that (i) preserves distinct viewpoints and constraints and (ii) helps a team move toward an actionable outcome. Given a task context  $c$  and a set of  $N$  team members  $\{u_1, \dots, u_N\}$ , each providing task-specific preference  $E_i$ , TeamFusion outputs a deliverable  $y$  intended to be directly usable.

**Overview** As shown in Figure 2, TeamFusion consists of four phases: (1) we instantiate one proxy agent per participant from their preferences; (2) proxy agents engage in a structured group discussion; (3) a remix phase converts the discussion into an editable deliverable, and (4) the system iterates this loop to refine the deliverable.

#### 3.1 Represent

For each participant  $u_i$ , we create a proxy agent  $a_i$  designed to argue from  $u_i$ 's perspective during the group discussion. Training a separate model per participant is impractical in realistic settings: per-user data are sparse, training is computationally expensive, and models would quickly become obsolete as preferences shift. Instead, we adopt an in-context personalization approach. Concretely, we encode the participant's evidence  $E_i$  into a

structured system prompt  $\pi_i$  that specifies: (i) the agent's role and collaborative objective, (ii) domain and communication constraints, and (iii) the participant-specific preference. Our goal is not to fully model a participant's identity, but to ensure the agent's contributions are recognizably aligned with that participant's expressed perspectives.

#### 3.2 Discuss

**Agent Roles** One TeamFusion run consists of  $N$  participant proxy agents:  $A = \{a_1, \dots, a_N\}$ . Proxy agents contribute proposals and critiques from their participant perspective.

**Conversation State** The discussion proceeds in a shared group-chat environment. At step  $t$ , the controller maintains a message history  $H_t = [m_1, \dots, m_t]$ , where each message  $m_k = (\text{name}, \text{content})$ . Agents do not hold additional private state. At each turn, the controller selects a speaker  $a \in A$  and prompts the underlying LLM with the agent's system prompt  $\pi_a$  and the current history  $H_t$ . This ensures that all agents reason over the same dialogue context.

**Turn-Taking Protocol** We adopt a simple but effective round-robin protocol inspired by the classic divergence–convergence model of creative processes (Acar and Runco, 2019; Runco and Acar,

2012) and nominal group technique (Dowling and St. Louis, 2000), giving each proxy agent a fixed number of speaking turns and cycling deterministically through agents to ensure equal opportunities to contribute. On a proxy turn,  $a_i$  receives  $H_t$  and is instructed to respond in light of its participant evidence and advance the discussion toward a recommendation. The discussion ends once all proxy agents exhaust their allotted turns.

### 3.3 Remix

After the debate concludes, TeamFusion converts the accumulated discussion into a final deliverable for the open-ended task at hand. A remixing agent takes as input the original task context  $c$  and the full discussion history  $H_T$ , and produces a deliverable that aggregates the reasons, trade-offs, and points of convergence. The remixed deliverable is intended to be directly consumable by humans. Its concrete form is task-dependent. Implementation details for each task are provided in Appendix F.3.

### 3.4 Iterative Refinement

TeamFusion can be applied once to obtain a single deliverable, or used iteratively to gradually refine outputs. In the iterative setting, the deliverable from one round is treated as a new proposal that re-enters the discussion: proxy agents are given access to the updated deliverable alongside the original context and asked to critique and build upon it in a subsequent discussion. This refinement loop allows the system to successively narrow in on options that better reflect surfaced preferences and rationales.

## 4 Task 1: Civic Comment Synthesis

We begin by evaluating TeamFusion in a civic decision-support setting, where a small group must turn diverse free-form public comments into a deliverable that can inform downstream action.

### 4.1 Task Introduction

Given a policy-relevant question and a set of participant comments, the system produces a concise summary intended to serve as a deliverable: it should capture the range of perspectives and the reasons behind them, rather than collapsing the group into a single averaged voice.

### 4.2 Experiment Protocol

**Experiment Data** We experiment on DeliberationBank (Zhu et al., 2025), a benchmark containing U.S.-based public-opinion comments spanning

ten questions about technology, social media, and public policy.

**Experiment Details** We primarily evaluate teams of four participants. For each question, we cluster the crowd-sourced comments into four groups and sample one comment from each cluster, forming a team intended to cover qualitatively different stances. We sample a total of 500 team configurations. We follow DeliberationBank protocol (Zhu et al., 2025) to score the outputs, and in addition using LLM as a judge (Li et al., 2025) to perform pairwise comparison.

**Baselines and Metrics** We compare with: 1. **Direct summary** that prompts an LLM to summarize the comments, 2. **Chain-of-Thought (CoT)** that prompts LLM to think before generating a final summary (Wei et al., 2022), 3. **Self-Refinement (Self-Refine)**, iteratively refining summaries without structured interaction (Madaan et al., 2023), 4. **Multi-Agent Debate (MAD)**, using generic agents conducting four rounds of debate (Du et al., 2023). We report four dimensions from DeliberationBank: representativeness, informativeness, neutrality, and policy approval. A detailed description of the four metrics and their significance to open-ended decision making is presented in Table 6.

### 4.3 Experiment Results

**TeamFusion consistently outperforms baselines.** Table 1 shows consistent gains from TeamFusion across base models and question types. We focus on **representativeness** as the primary metric because it directly captures our goal of preserving diverse viewpoints. TeamFusion yields the largest improvements on representativeness, and these gains co-occur with strong increases in informativeness and policy approval, suggesting that the additional structured discussion surfaces missing considerations that make summaries more decision-ready. Importantly, neutrality remains comparable to baselines, indicating that improved viewpoint coverage does not come from introducing more polarized or editorial language.

To complement these aggregate scores, we also conduct a pairwise comparison between TeamFusion-generated summaries and direct summaries using an LLM-as-a-judge. We randomly sample 300 TeamFusion outcomes (100 for each base model), pair them with the corresponding direct summaries, and prompt GPT-4.1-mini to decide which summary is better. To avoid any po-

Model	Method	OpenQA				BinaryQA			
		Represent.	Inform.	Neutral.	Policy	Represent.	Inform.	Neutral.	Policy
Llama-3.3-70B	Direct	.586 <sub>.006</sub>	.537 <sub>.006</sub>	.580 <sub>.006</sub>	.574 <sub>.007</sub>	.595 <sub>.007</sub>	.541 <sub>.006</sub>	.590 <sub>.007</sub>	.542 <sub>.007</sub>
	CoT	.568 <sub>.007</sub>	.524 <sub>.006</sub>	.570 <sub>.006</sub>	.561 <sub>.007</sub>	.578 <sub>.007</sub>	.530 <sub>.006</sub>	.580 <sub>.006</sub>	.528 <sub>.007</sub>
	Self-Refine	.577 <sub>.007</sub>	.536 <sub>.006</sub>	.566 <sub>.006</sub>	.570 <sub>.007</sub>	.599 <sub>.007</sub>	.545 <sub>.006</sub>	.587 <sub>.007</sub>	.547 <sub>.007</sub>
	MAD	.588 <sub>.007</sub>	.544 <sub>.006</sub>	.572 <sub>.007</sub>	.579 <sub>.007</sub>	.596 <sub>.007</sub>	.554 <sub>.006</sub>	.585 <sub>.007</sub>	.549 <sub>.006</sub>
	TeamFusion	<b>.608</b> <sub>.007</sub>	<b>.588</b> <sub>.006</sub>	<b>.587</b> <sub>.006</sub>	<b>.602</b> <sub>.007</sub>	<b>.620</b> <sub>.007</sub>	<b>.597</b> <sub>.006</sub>	<b>.610</b> <sub>.007</sub>	<b>.568</b> <sub>.007</sub>
GPT-4.1-mini	Direct	.582 <sub>.006</sub>	.534 <sub>.006</sub>	.576 <sub>.006</sub>	.579 <sub>.007</sub>	.589 <sub>.007</sub>	.531 <sub>.006</sub>	.584 <sub>.007</sub>	.543 <sub>.006</sub>
	CoT	.581 <sub>.007</sub>	.525 <sub>.006</sub>	.568 <sub>.006</sub>	.571 <sub>.007</sub>	.580 <sub>.007</sub>	.523 <sub>.006</sub>	.575 <sub>.007</sub>	.538 <sub>.006</sub>
	Self-Refine	.608 <sub>.007</sub>	.544 <sub>.007</sub>	.571 <sub>.007</sub>	.593 <sub>.008</sub>	.610 <sub>.007</sub>	.553 <sub>.007</sub>	.580 <sub>.008</sub>	.566 <sub>.006</sub>
	MAD	.598 <sub>.007</sub>	.553 <sub>.006</sub>	.569 <sub>.006</sub>	.592 <sub>.008</sub>	.616 <sub>.006</sub>	.559 <sub>.007</sub>	.586 <sub>.008</sub>	.570 <sub>.006</sub>
	TeamFusion	<b>.614</b> <sub>.007</sub>	<b>.594</b> <sub>.006</sub>	<b>.585</b> <sub>.006</sub>	<b>.608</b> <sub>.007</sub>	<b>.623</b> <sub>.007</sub>	<b>.601</b> <sub>.006</sub>	<b>.604</b> <sub>.007</sub>	<b>.587</b> <sub>.006</sub>
GPT-4.1	Direct	.578 <sub>.006</sub>	.531 <sub>.006</sub>	.573 <sub>.006</sub>	.575 <sub>.007</sub>	.584 <sub>.007</sub>	.530 <sub>.006</sub>	.577 <sub>.006</sub>	.541 <sub>.006</sub>
	CoT	.582 <sub>.006</sub>	.537 <sub>.006</sub>	.572 <sub>.006</sub>	.577 <sub>.006</sub>	.585 <sub>.007</sub>	.533 <sub>.006</sub>	.578 <sub>.007</sub>	.539 <sub>.006</sub>
	Self-Refine	.595 <sub>.007</sub>	.556 <sub>.006</sub>	.574 <sub>.006</sub>	.594 <sub>.007</sub>	.605 <sub>.008</sub>	.543 <sub>.006</sub>	.575 <sub>.007</sub>	.557 <sub>.007</sub>
	MAD	.599 <sub>.007</sub>	.561 <sub>.007</sub>	.576 <sub>.007</sub>	.594 <sub>.007</sub>	.609 <sub>.010</sub>	.563 <sub>.006</sub>	.580 <sub>.008</sub>	.567 <sub>.006</sub>
	TeamFusion	<b>.621</b> <sub>.007</sub>	<b>.622</b> <sub>.007</sub>	<b>.582</b> <sub>.006</sub>	<b>.619</b> <sub>.007</sub>	<b>.640</b> <sub>.007</sub>	<b>.634</b> <sub>.006</sub>	<b>.602</b> <sub>.007</sub>	<b>.603</b> <sub>.007</sub>

Table 1: Performance comparison between TeamFusion and baselines on DeliberationBank task. We present the results on the two sub-categories of the questions: OpenQA and BinaryQA. Values are mean  $\pm$  95% CI. We report scores for representativeness (Represent.), informativeness (Inform.), neutrality (Neutral.), and policy approval (Policy; higher is better). Best scores per column are in **bold**.

Model	Win / Tie / Loss (%)			
	Represent.	Inform.	Neutral.	Policy
Llama-70B	71 / 28 / 1	95 / 0 / 5	51 / 43 / 6	97 / 0 / 3
GPT-4.1-mini	72 / 26 / 2	96 / 0 / 4	27 / 61 / 12	96 / 0 / 4
GPT-4.1	93 / 7 / 0	98 / 0 / 2	46 / 49 / 5	99 / 0 / 1

Table 2: Win/Tie/Loss rate of TeamFusion outcome against direct summary across four metrics. Higher win rates indicate stronger relative performance.

sitional bias of the LLM judge (Shi et al., 2024), we randomized the order of the summaries in the prompt. As shown in Table 2, TeamFusion wins overwhelmingly on informativeness and policy approval, and wins on representativeness in the large majority of cases.

**Performance gains stem from personalized interaction.** We compare TeamFusion against compute-matched Self-Refine and MAD baselines. While MAD involves structured interactions among agents, its generic approach without personalization leads to limited viewpoint diversity and narrower coverage. Similarly, self-refinement provides an iterative reasoning structure but lacks interactive discussion. In contrast, TeamFusion’s performance boost is primarily attributable to its personalized agent interactions, highlighting the critical role of personalization and structured interaction. This analysis confirms that the observed improvements are not merely due to increased com-

Method	Represent.	Inform.	Neutral.	Policy
<b>Team size: 6</b>				
Direct	.582 <sub>.010</sub>	.537 <sub>.008</sub>	.575 <sub>.009</sub>	.560 <sub>.010</sub>
CoT	.576 <sub>.010</sub>	.529 <sub>.008</sub>	.569 <sub>.009</sub>	.556 <sub>.010</sub>
Self-Refine	.607 <sub>.010</sub>	.561 <sub>.009</sub>	.576 <sub>.009</sub>	.581 <sub>.010</sub>
MAD	.600 <sub>.009</sub>	.573 <sub>.009</sub>	.579 <sub>.009</sub>	.587 <sub>.010</sub>
TeamFusion	<b>.622</b> <sub>.010</sub>	<b>.617</b> <sub>.008</sub>	<b>.593</b> <sub>.010</sub>	<b>.608</b> <sub>.010</sub>
<b>Team size: 8</b>				
Direct	.580 <sub>.009</sub>	.540 <sub>.008</sub>	.580 <sub>.008</sub>	.556 <sub>.009</sub>
CoT	.575 <sub>.009</sub>	.530 <sub>.008</sub>	.573 <sub>.009</sub>	.550 <sub>.009</sub>
Self-Refine	.603 <sub>.009</sub>	.568 <sub>.008</sub>	.578 <sub>.009</sub>	.576 <sub>.010</sub>
MAD	.605 <sub>.009</sub>	.570 <sub>.008</sub>	.580 <sub>.009</sub>	.578 <sub>.010</sub>
TeamFusion	<b>.621</b> <sub>.009</sub>	<b>.627</b> <sub>.007</sub>	<b>.596</b> <sub>.009</sub>	<b>.609</b> <sub>.009</sub>
<b>Team size: 10</b>				
Direct	.578 <sub>.007</sub>	.546 <sub>.007</sub>	.579 <sub>.008</sub>	.554 <sub>.008</sub>
CoT	.574 <sub>.008</sub>	.535 <sub>.007</sub>	.571 <sub>.007</sub>	.550 <sub>.008</sub>
Self-Refine	.604 <sub>.009</sub>	.574 <sub>.008</sub>	.578 <sub>.009</sub>	.576 <sub>.010</sub>
MAD	.602 <sub>.009</sub>	.572 <sub>.007</sub>	.580 <sub>.008</sub>	.579 <sub>.009</sub>
TeamFusion	<b>.619</b> <sub>.008</sub>	<b>.637</b> <sub>.008</sub>	<b>.596</b> <sub>.009</sub>	<b>.608</b> <sub>.009</sub>

Table 3: Performance comparison between TeamFusion and baselines on the DeliberationBank task for different team sizes. Values are mean  $\pm$  95% CI. Best scores per column are in **bold**.

putational budget but rather due to the strategic combination of personalized representation and iterative debate structures.

**TeamFusion demonstrates gains across different team sizes.** We then investigate the effectiveness of team size. We fix the total number of sampled team configurations to 100 and use GPT-4.1-mini as the backbone. Results are shown in Table 3. Across different team sizes, TeamFusion

	Represent.	Inform.	Neutral.	Policy
<b>Model: Llama-3.3-70B</b>				
Base	.582 <sub>.007</sub>	.541 <sub>.006</sub>	.568 <sub>.006</sub>	.552 <sub>.006</sub>
+ Iter 1	.601 <sub>.006</sub>	.563 <sub>.005</sub>	.579 <sub>.006</sub>	.566 <sub>.006</sub>
+ Iter 2	<b>.618</b> <sub>.006</sub>	<b>.581</b> <sub>.005</sub>	<b>.592</b> <sub>.006</sub>	<b>.581</b> <sub>.006</sub>
<b>Model: GPT-4.1-mini</b>				
Base	.622 <sub>.018</sub>	.596 <sub>.014</sub>	.596 <sub>.016</sub>	.599 <sub>.016</sub>
+ Iter 1	.633 <sub>.018</sub>	.625 <sub>.015</sub>	<b>.604</b> <sub>.014</sub>	.618 <sub>.018</sub>
+ Iter 2	<b>.637</b> <sub>.018</sub>	<b>.638</b> <sub>.015</sub>	.599 <sub>.015</sub>	<b>.619</b> <sub>.016</sub>
<b>Model: GPT-4.1</b>				
Base	.630 <sub>.018</sub>	.615 <sub>.014</sub>	.591 <sub>.016</sub>	.606 <sub>.016</sub>
+ Iter 1	.646 <sub>.019</sub>	.655 <sub>.014</sub>	.599 <sub>.015</sub>	.635 <sub>.017</sub>
+ Iter 2	<b>.655</b> <sub>.017</sub>	<b>.686</b> <sub>.014</sub>	<b>.603</b> <sub>.015</sub>	<b>.643</b> <sub>.018</sub>

Table 4: Performance of different models under iterative refinement. Values are mean  $\pm$  95% CI.

consistently outperforms baselines on the key representativeness metric. TeamFusion can improve representativeness by about 0.04 absolute gain over the baselines, with non-overlapping confidence intervals. This demonstrates the scalability and generalization capability of TeamFusion.

**Iterative refinement brings gains.** We fix the total number of sampled team configurations to 100 and run TeamFusion with different iterative refinement numbers. The results of iterative refinement are shown in Table 4. For all models, adding iterative refinement brings gains to the representativeness and informativeness of the final summary across the iterations. Neutrality and policy alignment also improve, though with smaller margins, suggesting that additional rounds are particularly effective at surfacing missing considerations rather than merely smoothing tone.

## 5 Task 2: Visual Design

We then study TeamFusion in a *human-centric, multi-modal* workflow grounded in *real industry practice*.

### 5.1 Problem Motivation

Creative alignment is a significant pain point in professional design. Unlike close-ended tasks with objective “gold labels,” design briefs are open-ended and subject to interpretation. This ambiguity introduces friction in industry practice: teams must expend significant effort negotiating trade-offs between aesthetics, brand tone, and constraints.

We motivate this task by empirically quantifying this friction. In our preliminary analysis of professional designers’ preferences (detailed in Sec 5.4),

we observed that experts given the exact same brief and assets exhibited remarkably low agreement on quality. In 70% of cases, agreement was indistinguishable from random chance. This validates that divergent interpretation is a natural and pervasive bottleneck.

### 5.2 Task Setup

Each scenario consists of a client brief clarifying the requirements for an advertisement design and a set of candidate ad thumbnails. A team of professional designers rank the candidates and provide justifications. TeamFusion runs proxy agent discussion and produces remixed design images intended to better align with designers’ expressed constraints while making the underlying points of agreement and disagreement actionable for downstream selection. We evaluate whether TeamFusion generated images can replace original team’s favorites.

### 5.3 Experiment Protocol

We introduce a human-in-the-loop protocol to evaluate TeamFusion. Unlike static benchmark evaluations, this protocol allows us to scale realistic team interactions while keeping professional designers as the ultimate ground truth for decision quality.

**Phase 0: Realistic Scenario Construction** We construct 50 high-quality design scenarios derived from real social media advertising campaigns (Yamaguchi, 2021). Each scenario includes a professional client brief and a set of diverse candidate designs. All scenarios have been validated by two external senior designers for realism. Full construction details are deferred to Appendix G.2.

**Phase 1: Preference collection** We recruit 9 professional designers to annotate scenarios asynchronously. For each assigned scenario, designer ranks the six options and writes a brief justification. Each scenario received at least four independent annotations, serving as the “seed” evidence that conditions our proxy agents.

**Phase 2: Simulation with nominal teams** For each scenario, we form two nominal team settings from asynchronous annotations: **Full-Team** (all available annotations) and **Small-Team** (a random subset of two designers). We run TeamFusion for three iterations, producing one new remixed design per iteration. This yields 100 TeamFusion runs and 300 remixed design candidates. We also record

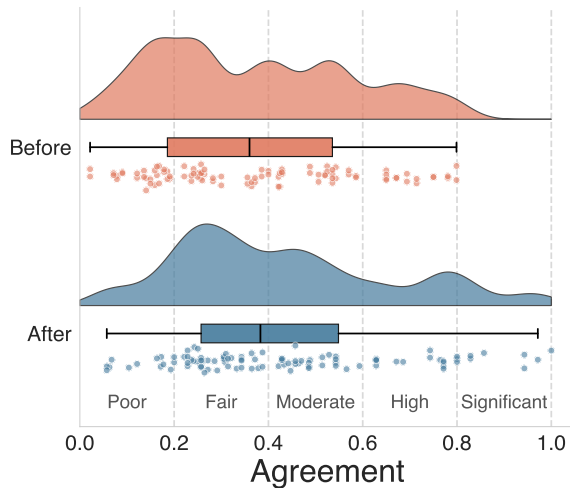


Figure 3: The distribution of agreement scores to measure dataset-wide agreement before and after TeamFusion’s execution. The data is categorized into five value ranges to interpret agreement strength. Agreements across 100 team settings after running TeamFusion show a dataset-wide move towards higher agreement.

each proxy agent’s discussion comments about the generated candidates for later analysis.

**Phase 3: Designer re-evaluation** To determine if the system successfully facilitated convergence, we close the loop by returning the generated outputs to the original human designers. We combine the team’s initial top three options via Borda count with the three TeamFusion-generated options. Designers then re-rank the combined set and rate whether their proxy agent’s commentary aligns with their own reasoning.

## 5.4 Experiment Results

Our analysis reveals three main findings. First, we empirically verify the motivating problem of divergent preferences in teams. Second, we show that TeamFusion can generate consensus-oriented remixes that successfully induce convergence. Finally, we show that real designers largely agree with the debate commentary made by their delegate agents, indicating that simulated debates are well-grounded.

**Finding 1: Divergent interpretations are a real, salient problem.** As outlined in our motivation, we hypothesized that professional designers hold conflicting interpretations of the same brief. Our analysis of the pre-discussion ranking data confirms this friction is substantial. To quantify this, we calculate Kendall’s Coefficient of Concordance

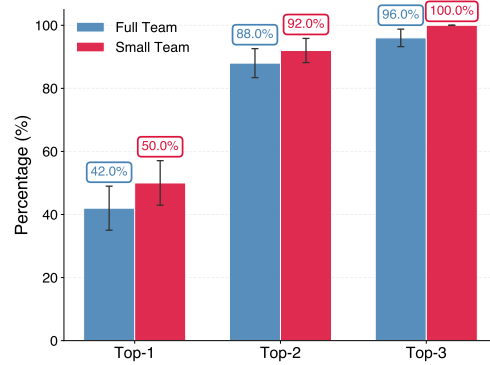


Figure 4: The rate of TeamFusion-generated images appearing in the final top-ranked selections. Error bars represent the 95% confidence interval.

( $W$ ) on the independent rankings provided in Phase 1. As shown in Figure 3 (top), the agreement among designers is consistently low, with a mean of 0.37 (falling into the “Fair Agreement” range). Notably, 84% of scenarios fall into the “Moderate” or lower agreement categories, and in 70% of cases, the agreement among professionals is not statistically significant ( $p \geq 0.05$ ). This widespread lack of consensus in the real data confirms that our motivating problem is natural-arising and salient, providing strong empirical evidence in support of systems like TeamFusion.

**Finding 2: TeamFusion can support team convergence by generating consensus-oriented designs.** Our results reveal two ways in which TeamFusion-generated design revisions meaningfully modify the output of creative teams.

The results presented in Figure 4 show that **TeamFusion-generated options become the single top-ranked option across the team in nearly half of all test cases**, displacing the original team-wide favorites. Notably, this indicates that the execution of TeamFusion can be seen as a generative AI feature with significant team-wide acceptance rate under the strictest decision-making scenario, that is, the team decides to move forward with the single best design only.

Under less strict decision-making scenarios, TeamFusion also shows potential for contributing to teams’ outputs. We find that TeamFusion-generated option appeared in the top-two rankings in a remarkable 88% of Full-Team and 92% of Small-Team test cases. This finding is noteworthy because, in the creative decision-making space, teams may use not only the single best option out of group ideation, but actually a short-list of top

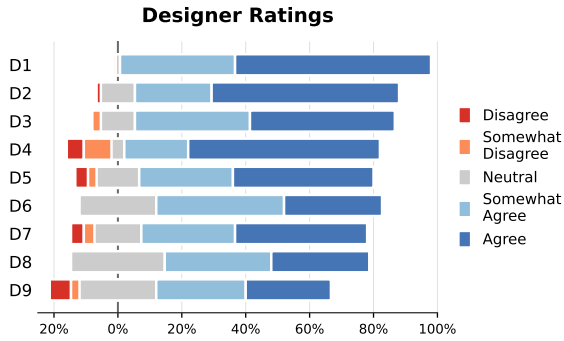


Figure 5: Distribution of annotator ratings for agreement with agent-generated commentary, grouped by designers. The results show an overwhelmingly positive perception.

options; for example, teams may steer the creative process by revising on one of them.

More globally, Figure 3 shows how there is a dataset-wide move towards higher agreement after the exposure to TeamFusion outputs: mean Kendall’s  $W$  rises from 0.37 to 0.43. Qualitatively, the dataset mean moves from only “Fair Agreement” before to “Moderate Agreement” after TeamFusion, which further indicates the effectiveness of our system in supporting convergence.

**Finding 3: Designers feel largely represented by their proxy agents.** As shown in Figure 5, the results were **overwhelmingly positive**. The scores were strongly skewed toward agreement, with a mean of 4.06 ( $\sigma = 1.07$ ). Over 75% of all comments received a positive score. This indicates that designers broadly perceived the outputs of their proxy agents as natural-sounding and representative of their own design rationales. This positive perception was highly consistent across our participants. An analysis of designer-level means showed a narrow range, with the lowest average rating being 3.53. This indicates that even the most critical participant found the commentary to be representative. The low standard deviation of these designer means ( $\sigma = 0.29$ ) further reinforces that this high level of agreement is a shared, consistent finding.

### 5.5 Live User Study Results

To complement our asynchronous evaluation, we run a live, controlled within-team study to test whether TeamFusion facilitates convergence in an end-to-end workflow where participants create and revise designs. We recruit six participants and formed two teams of three, with each team completing two tasks in a counterbalanced crossover

Metric	Discussion	TeamFusion
Decision Time (min) ↓	18.0	<b>12.4</b>
Q1: Representative ↑	3.7	<b>4.3</b>
Q2: Clarity ↑	3.5	<b>3.8</b>
Q3: Satisfaction ↑	3.5	<b>4.2</b>
Preferred ↑	1/6	<b>5/6</b>

Table 5: Live within-team study results. Each team completed two briefs in a counterbalanced crossover design.

design<sup>1</sup>. Due to the small number of participants, we report this live study as a pilot with descriptive results rather than a statistically powered evaluation. As shown in Table 5, using TeamFusion leads to faster team decisions compared to free-form discussion, while also improving participants’ perceived representativeness, clarity of trade-offs, and overall satisfaction with the team outcome. After experiencing both workflows, a strong majority of participants explicitly preferred TeamFusion over free-form discussion, showcasing the effectiveness over unconstrained collaboration.

## 6 Case Study

We present a case study on the effectiveness in Figure 7. Due to page limit, we defer the detailed analysis in Appendix D. The core takeaway is that TeamFusion can better preserve fine-grained, participant-specific content than direct aggregation.

## 7 Conclusion

Open-ended team decisions require deliverables that make trade-offs and disagreements visible rather than averaging them away, yet common aggregation methods often erase minority or conditional viewpoints and reduce auditability. We addressed this challenge with TeamFusion, a multi-agent framework that shifts the paradigm from direct aggregation to modeled interaction. By representing participants with preference-grounded proxy agents and orchestrating structured debates, TeamFusion externalizes the friction of consensus-building to produce editable, rationale-backed deliverables. Evaluations on two teamwork tasks show that explicitly modeling interaction improves viewpoint coverage and decision usefulness and can induce greater convergence, validating the potential of AI to facilitate human collaboration.

<sup>1</sup>Full details are deferred to Appendix H

## 554 Limitations

555 Our findings show that TeamFusion can success-  
556 fully support the creative convergence process,  
557 producing consensus-inducing design revisions  
558 grounded in natural-sounding, agreeable rationales.  
559 However, our work has limitations that shed light  
560 on important directions for future research. The  
561 current implementation assumes a flat hierarchy in  
562 the team, not accounting for different roles (e.g.,  
563 art directors managing the team, or even clients  
564 themselves in the loop) and seniority levels (e.g.,  
565 senior vs. junior designers). Even more realistic  
566 professional settings may warrant slightly differ-  
567 ent assumptions, with implications to our current  
568 modeling decisions.

## 569 References

570 Selcuk Acar and Mark A Runco. 2019. Divergent think-  
571 ing: New methods, recent research, and extended  
572 theory. *Psychology of aesthetics, creativity, and the*  
573 *arts*, 13(2):153.

574 Divyansh Agarwal, Alexander Fabbri, Ben Risher,  
575 Philippe Laban, Shafiq Joty, and Chien-Sheng Wu.  
576 2024. [Prompt leakage effect and mitigation strate-](#)  
577 [gies for multi-turn LLM applications](#). In *Proceed-*  
578 *ings of the 2024 Conference on Empirical Methods*  
579 *in Natural Language Processing: Industry Track*,  
580 pages 1255–1275, Miami, Florida, US. Association  
581 for Computational Linguistics.

582 Adithya Bhaskar, Alex Fabbri, and Greg Durrett. 2023.  
583 [Prompted opinion summarization with GPT-3.5](#). In  
584 *Findings of the Association for Computational Lin-*  
585 *guistics: ACL 2023*, pages 9282–9300, Toronto,  
586 Canada. Association for Computational Linguistics.

587 Duncan Black. 1948. On the rationale of group decision-  
588 making. *Journal of political economy*, 56(1):23–34.

589 Jean Carletta, Simone Ashby, Sebastien Bourban, Mike  
590 Flynn, Mael Guillemot, Thomas Hain, Jaroslav  
591 Kadlec, Vasilis Karaiskos, Wessel Kraaij, Melissa  
592 Kronenthal, Guillaume Lathoud, Mike Lincoln,  
593 Agnes Lisowska, Iain McCowan, Wilfried Post, Den-  
594 nis Reidsma, and Pierre Wellner. 2006. The ami  
595 meeting corpus: A pre-announcement. In *Machine*  
596 *Learning for Multimodal Interaction*, pages 28–39,  
597 Berlin, Heidelberg. Springer Berlin Heidelberg.

598 Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu,  
599 Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan  
600 Liu. 2023. Chateval: Towards better llm-based eval-  
601 uators through multi-agent debate. *arXiv preprint*  
602 *arXiv:2308.07201*.

603 Harrison Chase. 2022. Langchain. [https://github.](https://github.com/langchain-ai/langchain)  
604 [com/langchain-ai/langchain](https://github.com/langchain-ai/langchain).

John Chen, Xi Lu, Yuzhou Du, Michael Rejtig, Ruth  
605 Bagley, Mike Horn, and Uri Wilensky. 2024. [Learn-](#)  
606 [ing agent-based modeling with llm companions: Ex-](#)  
607 [periences of novices and experts using chatgpt & net-](#)  
608 [logo chat](#). *Proceedings of the 2024 CHI Conference*  
609 *on Human Factors in Computing Systems*. 610

Steffi Chern, Zhen Fan, and Andy Liu. 2024. Com-  
611 bating adversarial attacks with multi-agent debate.  
612 *arXiv preprint arXiv:2401.05998*. 613

Chun-Wei Chiang, Zhuoran Lu, Zhuoyan Li, and Ming  
614 Yin. 2024. Enhancing ai-assisted group decision  
615 making through llm-powered devil’s advocate. In  
616 *Proceedings of the 29th International Conference on*  
617 *Intelligent User Interfaces*, pages 103–119. 618

Luis Fariñas del Cerro, Andreas Herzig, Dominique  
619 Longin, and Omar Rifi. 1998. Belief reconstruction  
620 in cooperative dialogues. In *International Confer-*  
621 *ence on Artificial Intelligence: Methodology, Systems,*  
622 *and Applications*, pages 254–266. Springer. 623

Karen L. Dowling and Robert D. St. Louis. 2000.  
624 [Asynchronous implementation of the nominal group](#)  
625 [technique: is it effective?](#) *Decis. Support Syst.*,  
626 29(3):229–248. 627

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenen-  
628 baum, and Igor Mordatch. 2023. Improving factual-  
629 ity and reasoning in language models through multi-  
630 agent debate. In *Forty-first International Conference*  
631 *on Machine Learning*. 632

B Aubrey Fisher. 1970. Decision emergence: Phases  
633 in group decision-making. *Communications Mono-*  
634 *graphs*, 37(1):53–66. 635

Jarod Govers, Eduardo Velloso, Vassilis Kostakos, and  
636 Jorge Goncalves. 2024. [Ai-driven mediation strate-](#)  
637 [gies for audience depolarisation in online debates](#). In  
638 *Proceedings of the 2024 CHI Conference on Human*  
639 *Factors in Computing Systems, CHI ’24*, New York,  
640 NY, USA. Association for Computing Machinery. 641

Kai Greshake, Sahar Abdelnabi, Shailesh Mishra,  
642 Christoph Endres, Thorsten Holz, and Mario Fritz.  
643 2023. Not what you’ve signed up for: Compromis-  
644 ing real-world llm-integrated applications with indi-  
645 rect prompt injection. In *Proceedings of the 16th*  
646 *ACM workshop on artificial intelligence and security*,  
647 pages 79–90. 648

Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu  
649 Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang,  
650 Zili Wang, Steven Ka Shing Yau, Zijuan Lin, and  
651 1 others. 2023. Metagpt: Meta programming for a  
652 multi-agent collaborative framework. In *The Twelfth*  
653 *International Conference on Learning Representa-*  
654 *tions*. 655

Zhe Hu, Hou Pong Chan, Jing Li, and Yu Yin.  
656 2025. Debate-to-write: A persona-driven multi-agent  
657 framework for diverse argument generation. In *Pro-*  
658 *ceedings of the 31st International Conference on*  
659 *Computational Linguistics*, pages 4689–4703. 660

661	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong,	<i>Methods in Natural Language Processing</i> , pages	717
662	Zhangyin Feng, Haotian Wang, Qianglong Chen,	2757–2791.	718
663	Weihua Peng, Xiaocheng Feng, Bing Qin, and 1 oth-		
664	ers. 2025. A survey on hallucination in large lan-	Miao Li, Eduard Hovy, and Jey Lau. 2023. <a href="#">Summariz-</a>	719
665	guage models: Principles, taxonomy, challenges, and	<a href="#">ing multiple documents with conversational structure</a>	720
666	open questions. <i>ACM Transactions on Information</i>	<a href="#">for meta-review generation</a> . In <i>Findings of the As-</i>	721
667	<i>Systems</i> , 43(2):1–55.	<i>sociation for Computational Linguistics: EMNLP</i>	722
668	Nannan Huang, Lin Tian, Haytham Fayek, and Xiuzhen	2023, pages 7089–7112, Singapore. Association for	723
669	Zhang. 2023. <a href="#">Examining bias in opinion summarisa-</a>	<i>Computational Linguistics</i> .	724
670	<a href="#">tion through the perspective of opinion diversity</a> . In		
671	<i>Proceedings of the 13th Workshop on Computational</i>	Miao Li, Jey Han Lau, and Eduard Hovy. 2024. <a href="#">A senti-</a>	725
672	<i>Approaches to Subjectivity, Sentiment, &amp; Social Me-</i>	<a href="#">ment consolidation framework for meta-review gen-</a>	726
673	<i>dia Analysis</i> , pages 149–161, Toronto, Canada. Asso-	<a href="#">eration</a> . In <i>Proceedings of the 62nd Annual Meeting</i>	727
674	ciation for Computational Linguistics.	<i>of the Association for Computational Linguistics (Vol-</i>	728
675	Qiushi Huang, Xubo Liu, Tom Ko, Boyong Wu, Wenwu	<i>ume 1: Long Papers)</i> , pages 10158–10177, Bangkok,	729
676	Wang, Yu Zhang, and Lilian Tang. 2024. <a href="#">Selective</a>	Thailand. Association for Computational Linguistics.	730
677	<a href="#">prompting tuning for personalized conversations with</a>		
678	<a href="#">llms</a> . <i>ArXiv</i> , abs/2406.18187.	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang,	731
679	Pontus Johansson. 2002. User modeling in dialog sys-	Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and	732
680	tems. <i>St. Anna Report SAR</i> , pages 02–2.	Zhaopeng Tu. 2024. Encouraging divergent thinking	733
681	Sara Kiesler and Lee Sproull. 1992. Group decision	in large language models through multi-agent debate.	734
682	making and communication technology. <i>Organi-</i>	In <i>Proceedings of the 2024 conference on empiri-</i>	735
683	<i>zational behavior and human decision processes</i> ,	<i>cal methods in natural language processing</i> , pages	736
684	52(1):96–123.	17889–17904.	737
685	Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik S Chan,	Xiaogeng Liu, Zhiyuan Yu, Yizhe Zhang, Ning Zhang,	738
686	Xuhai Xu, Daniel McDuff, Hyeonhoon Lee, Marzyeh	and Chaowei Xiao. 2024. Automatic and universal	739
687	Ghassemi, Cynthia Breazeal, and Hae W Park. 2024.	prompt injection attacks against large language mod-	740
688	Mdagents: An adaptive collaboration of llms for medi-	els. <i>arXiv preprint arXiv:2403.04957</i> .	741
689	cal decision-making. <i>Advances in Neural Informa-</i>		
690	<i>tion Processing Systems</i> , 37:79410–79452.	Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying	742
691	Alfred Kobsa. 1989. A taxonomy of beliefs and goals	Zhang, Ruocheng Guo, Hao Cheng, Yegor Klochkov,	743
692	for user models in dialog systems. In <i>User models in</i>	Muhammad Faaiz Taufiq, and Hang Li. 2023. Trust-	744
693	<i>dialog systems</i> , pages 52–68. Springer.	worthy llms: a survey and guideline for evaluating	745
694	Kenneth L Kraemer and John Leslie King. 1988.	large language models’ alignment. <i>arXiv preprint</i>	746
695	Computer-based systems for cooperative work and	<i>arXiv:2308.05374</i> .	747
696	group decision making. <i>ACM Computing Surveys</i>	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	748
697	( <i>CSUR</i> ), 20(2):115–146.	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	749
698	D. Kwon, Sunwoo Lee, Ki Hyun Kim, Sejin Lee, Tae-	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	750
699	Yoon Kim, and Eric Davis. 2023. <a href="#">What, when, and</a>	and 1 others. 2023. Self-refine: Iterative refinement	751
700	<a href="#">how to ground: Designing user persona-aware con-</a>	with self-feedback. <i>Advances in Neural Information</i>	752
701	<a href="#">versational agents for engaging dialogue</a> . In <i>Annual</i>	<i>Processing Systems</i> , 36:46534–46594.	753
702	<i>Meeting of the Association for Computational Lin-</i>	Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad	754
703	<i>guistics</i> .	Saqib, Saeed Anwar, Muhammad Usman, Naveed	755
704	Philippe Laban, Wojciech Kryściński, Divyansh Agar-	Akhtar, Nick Barnes, and Ajmal Mian. 2025. A com-	756
705	wal, Alexander Richard Fabbri, Caiming Xiong,	prehensive overview of large language models. <i>ACM</i>	757
706	Shafiq Joty, and Chien-Sheng Wu. 2023. Summedits:	<i>Transactions on Intelligent Systems and Technology</i> ,	758
707	Measuring llm ability at factual reasoning through the	16(5):1–72.	759
708	lens of summarization. In <i>Proceedings of the 2023</i>	R Orwig, Hsinchun Chen, D Vogel, and Jay F Nuna-	760
709	<i>conference on empirical methods in natural language</i>	maker. 1997. A multi-agent view of strategic plan-	761
710	<i>processing</i> , pages 9662–9676.	ning using group support systems and artificial intelli-	762
711	Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad	gence. <i>Group Decision and Negotiation</i> , 6(1):37–59.	763
712	Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-	Letitia Parcalabescu and Anette Frank. 2024. On mea-	764
713	tacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu,	suring faithfulness or self-consistency of natural lan-	765
714	and 1 others. 2025. From generation to judgment:	guage explanations. In <i>Proceedings of the 62nd An-</i>	766
715	Opportunities and challenges of llm-as-a-judge. In	<i>Annual Meeting of the Association for Computational</i>	767
716	<i>Proceedings of the 2025 Conference on Empirical</i>	<i>Linguistics (Volume 1: Long Papers)</i> , pages 6048–	768
		6089.	769
		Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Mered-	770
		ith Ringel Morris, Percy Liang, and Michael S Bern-	771
		stein. 2023. Generative agents: Interactive simulacra	772

773	of human behavior. In <i>Proceedings of the 36th annual acm symposium on user interface software and technology</i> , pages 1–22.		
774			
775			
776	Giorgio Piatti, Zhijing Jin, Max Kleiman-Weiner, Bernhard Schölkopf, Mrinmaya Sachan, and Rada Mihalcea. 2024. Cooperate or collapse: Emergence of sustainable cooperation in a society of llm agents. <i>Advances in Neural Information Processing Systems</i> , 37:111715–111759.		
777			
778			
779			
780			
781			
782	Steven Rogelberg, Desmond J. Leach, Peter B. Warr, and Jennifer L. Burnfield. 2006. "not another meeting!" are meeting time demands related to employee well-being? <i>The Journal of applied psychology</i> , 91 1:83–96.		
783			
784			
785			
786			
787	Alexander C. Romney, Joseph A. Allen, and Zahra Heydarifard. 2025. <i>Meeting load paradox: Balancing the benefits and burdens of work meetings</i> . <i>Business Horizons</i> , 68(1):33–43.		
788			
789			
790			
791	Mark A Runco and Selcuk Acar. 2012. Divergent thinking as an indicator of creative potential. <i>Creativity research journal</i> , 24(1):66–75.		
792			
793			
794	Donald G. Saari. 1995. <i>Basic Geometry of Voting</i> , 1 edition. Springer-Verlag Berlin Heidelberg.		
795			
796	Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askeff, Samuel R Bowman, Esin DURMUS, Zac Hatfield-Dodds, Scott R Johnston, Shauna M Kravec, and 1 others. Towards understanding sycophancy in language models. In <i>The Twelfth International Conference on Learning Representations</i> .		
797			
798			
799			
800			
801			
802			
803	Lin Shi, Chiyu Ma, Wenhua Liang, Xingjian Diao, Weicheng Ma, and Soroush Vosoughi. 2024. Judging the judges: A systematic study of position bias in llm-as-a-judge. <i>arXiv preprint arXiv:2406.07791</i> .		
804			
805			
806			
807	Joon Gi Shin, Janin Koch, Andrés Lucero, Peter Dalsgaard, and Wendy E. Mackay. 2023. <i>Integrating ai in human-human collaborative ideation</i> . In <i>Extended Abstracts of the 2023 CHI Conference on Human Factors in Computing Systems</i> , CHI EA '23, New York, NY, USA. Association for Computing Machinery.		
808			
809			
810			
811			
812			
813	Elizabeth Shriberg, Raj Dhillon, Sonali Bhagat, Jeremy Ang, and Hannah Carvey. 2004. <i>The ICSI meeting recorder dialog act (MRDA) corpus</i> . In <i>Proceedings of the 5th SIGdial Workshop on Discourse and Dialogue at HLT-NAACL 2004</i> , pages 97–100, Cambridge, Massachusetts, USA. Association for Computational Linguistics.		
814			
815			
816			
817			
818			
819			
820	Qiushi Sun, Zhangyue Yin, Xiang Li, Zhiyong Wu, Xipeng Qiu, and Lingpeng Kong. Corex: Pushing the boundaries of complex reasoning through multi-model collaboration. In <i>First Conference on Language Modeling</i> .		
821			
822			
823			
824			
825	Michael Henry Tessler, Michiel A. Bakker, Daniel Jarrett, Hannah Sheahan, Martin J. Chadwick, Raphael Koster, Georgina Evans, Lucy Campbell-Gillingham,		
826			
827			
		Tantum Collins, David C. Parkes, Matthew Botvinick, and Christopher Summerfield. 2024. <i>Ai can help humans find common ground in democratic deliberation</i> . <i>Science</i> , 386(6719):eadq2852.	828
			829
			830
			831
		Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, and 1 others. 2023. Survey on factuality in large language models: Knowledge, retrieval and domain-specificity. <i>arXiv preprint arXiv:2310.07521</i> .	832
			833
			834
			835
			836
			837
		Zheng Wang, Zhongyang Li, Zeren Jiang, Dandan Tu, and Wei Shi. 2024. <i>Crafting personalized agents through retrieval-augmented generation on editable memory graphs</i> . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	838
			839
			840
			841
			842
		Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	843
			844
			845
			846
			847
			848
		Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. 2024. <i>Autogen: Enabling next-gen LLM applications via multi-agent conversations</i> . In <i>First Conference on Language Modeling</i> .	849
			850
			851
			852
			853
			854
			855
		Kota Yamaguchi. 2021. Canvasvae: Learning to generate vector graphic documents. <i>ICCV</i> .	856
			857
		Miao Yu, Fanci Meng, Xinyun Zhou, Shilong Wang, Junyuan Mao, Linsey Pan, Tianlong Chen, Kun Wang, Xinfeng Li, Yongfeng Zhang, and 1 others. 2025. A survey on trustworthy llm agents: Threats and countermeasures. In <i>Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2</i> , pages 6216–6226.	858
			859
			860
			861
			862
			863
			864
		Haopeng Zhang, Philip S Yu, and Jiawei Zhang. 2025. A systematic survey of text summarization: From statistical methods to large language models. <i>ACM Computing Surveys</i> , 57(11):1–41.	865
			866
			867
			868
		Shaowei Zhang and Deyi Xiong. 2025. Debate4math: Multi-agent debate for fine-grained reasoning in math. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 16810–16824.	869
			870
			871
			872
		Yang Zhang, Hanlei Jin, Dan Meng, Jun Wang, and Jinghua Tan. 2024. A comprehensive survey on process-oriented automatic text summarization with exploration of llm-based methods. <i>arXiv preprint arXiv:2403.02901</i> .	873
			874
			875
			876
			877
		Shenzhe Zhu, Shu Yang, Michiel A Bakker, Alex Pentland, and Jiaxin Pei. 2025. Can ai truly represent your voice in deliberations? a comprehensive study of large-scale opinion aggregation with llms. <i>arXiv preprint arXiv:2510.05154</i> .	878
			879
			880
			881
			882

883	<b>A Discussion</b>	
884	In this section, we discuss the broader implications of our findings. We begin by examining the contributions of our user study procedure as a scalable method for evaluating AI systems for group ideation. We then consider the potential of TeamFusion’s architecture as a generalizable model for AI-facilitated team convergence beyond graphic design. Finally, we address the limitations of our work and outline promising directions for future research.	
885		
886		
887		
888		
889		
890		
891		
892		
893		
894	<b>A.1 User Study Procedure</b>	
895	One of the primary contributions of this work is the evaluation procedure itself. Research into AI systems for team settings, particularly in creative domains, is often hampered by methodological challenges and logistical overhead in human-in-the-loop evaluation. A key advantage of our three-phase protocol (Annotate → Simulate → Re-evaluate) is its scalability, which addresses this important bottleneck in team settings beyond text-only domains.	
896		
897		
898		
899		
900		
901		
902		
903		
904		
905	While existing work relies on live, synchronous sessions with participants, making them time-consuming, expensive, and difficult to scale beyond a small number of test cases, our approach is asynchronous and simulation-driven. By collecting designers’ detailed rankings and justifications upfront (Phase 1), we effectively treat their expert judgment as a reusable resource. Instead of requiring designers to be online for every system execution, our approach can compose nominal teams from the offline annotations, simulate team dynamics, and return to team members to evaluate groundedness and effectiveness from system-generated deliverables. This allowed us to run 100 test cases covering a larger experimental space (i.e., two team sizes, three iterations of debating-and-remixing) much more efficiently. We hope the community interested in group ideation (Shin et al., 2023) can benefit from the key ideas behind our reproducible protocol.	
906		
907		
908		
909		
910		
911		
912		
913		
914		
915		
916		
917		
918		
919		
920		
921		
922		
923		
924		
925	<b>A.2 AI as a Facilitator for Team Convergence in Other Creative Domains</b>	
926		
927	While TeamFusion was implemented and evaluated within the domain of professional graphic design, its underlying architecture can be re-instantiated or extended to support team convergence in other creative domains. Advances in generative AI for	
928		
929		
930		
931		
	video- or audio-editing, for example, pose interesting questions as to whether the positive findings we see in our studies would transfer to these other professional settings that similarly rely on group ideation.	932 933 934 935 936
	<b>A.3 Potential Risks</b>	937
	One significant risk involves the privacy and security implications of creating high-fidelity proxy agents conditioned on sensitive personal data. Since TeamFusion operates by encoding a participant’s specific expressed preferences directly into a structured system prompt, there is an inherent risk that these digital proxies could inadvertently disclose more information than the user intended. For instance, while a user might strategically withhold certain views or “hidden assumptions” in a human-to-human setting, a proxy agent designed to “advocate for the assigned preference” might be manipulated via adversarial prompting (Greshake et al., 2023; Liu et al., 2024) or dialogue leaks (Agarwal et al., 2024) to reveal private rationales, biases, or competitive strategies to other agents in the shared “group chat” environment.	938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954
	<b>B Future Work</b>	955
	We are actively interested in exploring more hierarchical teams spanning different roles. Extending TeamFusion to other creative domains is another exciting direction that we identify, as well as further exploring “knobs” in the underlying parameter space such as the number of system iterations. Another important avenue of future work relates to more dynamic persona modeling: for example, agents could be designed to dynamically update their preferences and rationales based on the ongoing dialogue, better mimicking human adaptability and belief revision (Kobsa, 1989; Johansson, 2002; del Cerro et al., 1998).	956 957 958 959 960 961 962 963 964 965 966 967 968
	<b>C Additional Results</b>	969
	<b>C.1 Ablations on Per Agent Turns</b>	970
	We ablate on the number of discussion rounds in Task 1. We set the number of discussion rounds and observe TeamFusion performance across the four metrics. As shown in Figure 6, scaling up the rounds of discussion can improve informativeness and neutrality. Stronger proprietary models like GPT-4.1-mini and GPT-4.1 benefits more from increasing discussion rounds, representativeness	971 972 973 974 975 976 977 978

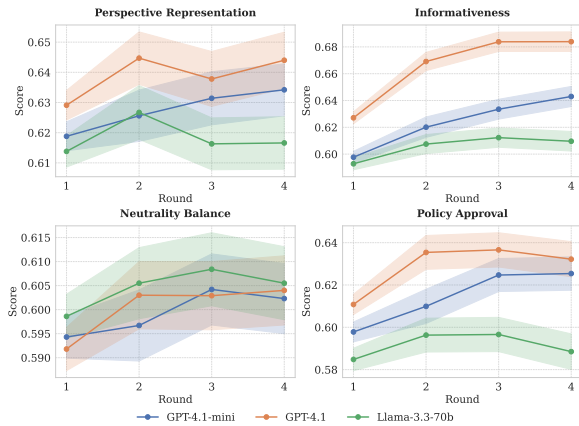


Figure 6: Ablations of per agent speaking turns on the four metrics.

increases. This showcases that adding more discussion can make different voices heard for stronger models. For weaker model Llama, adding round from 3 to 4 yields a decrease in four metrics. We hypothesize that this is due to its worse long context processing capability.

## D Case Study

Figure 7 illustrates how TeamFusion better preserves fine-grained, participant-specific content than direct aggregation on a civic synthesis example. While the direct summary is fluent and captures the dominant positive sentiment (e.g., creativity, personalization, and efficiency), it noticeably homogenizes key details: it fails to carry forward Comment #1’s concrete comparative claim that AI has “replaced” traditional search engines, and it collapses Comment #3’s domain-specific experience (“my healthcare setting”) into generic “work-related problems,” obscuring where and why the tool is valuable. In contrast, TeamFusion output retains these high-salience details in the final deliverable. This example highlights two practical advantages of TeamFusion for generating deliverables: (i) minority or specialized experiences remain visible rather than averaged away, and (ii) concrete comparative statements and applications are maintained to support downstream interpretation and action. Overall, the case study qualitatively supports our quantitative gains on representativeness by showing that TeamFusion more faithfully carries forward what each participant uniquely contributed, instead of compressing distinct voices into a generic narrative.

## E Implementation Details of Task 1

### E.1 Details of Represent

We design a prompt that consists of goal of the discussion, conversation style constraints, and preference samples. Prompt content is available at Appendix K.1.1.

### E.2 Details of Remix

The remixing agent is a text-based LLM that takes in the task context, original comments, the discussion transcript, and generates a structured summary. Prompt content can be found at Appendix K.1.2.

### E.3 Details of Iterative Revision

After an initial summary has been generated, the summary becomes a shared group message content as part of the task context. The agent group’s collective goal changes to improve the summary to better reflect their individual standpoint. They engage in the structured discussion to achieve the goal. Finally, the remixing agent ingests the individual comments based on the previous summary and the debate transcript, and generates a refined summary of the comments. The newly generated summary then becomes the summary to improve upon in the next round, if any.

### E.4 Hyperparameters

We use three LLMs as backbones for both proxy agents and the remixing agent, i.e. Llama-3.3-70b, GPT-4.1-mini, GPT-4.1. We set the decoding temperature to 1. We set the number of per-debate turns to 1 in the main experiments. The agents are implemented and orchestrated by the AutoGen framework (Wu et al., 2024).

## F Implementation Details of Task 2

### F.1 Details of Representation

Based on recent works on LLM-based personalization (Chen et al., 2024; Huang et al., 2024; Wang et al., 2024; Kwon et al., 2023), we leverage in-context learning to customize an LLM agent for the participant. We combine best practices from prior work to compose a layered prompt that makes an LLM adhere to a given designer’s preferences and reliably act as individual Designer Agents in our multi-agent system:

- **Overarching goal:** The first component of the prompt defines the overarching goal and the role of the agent. It establishes the agent’s

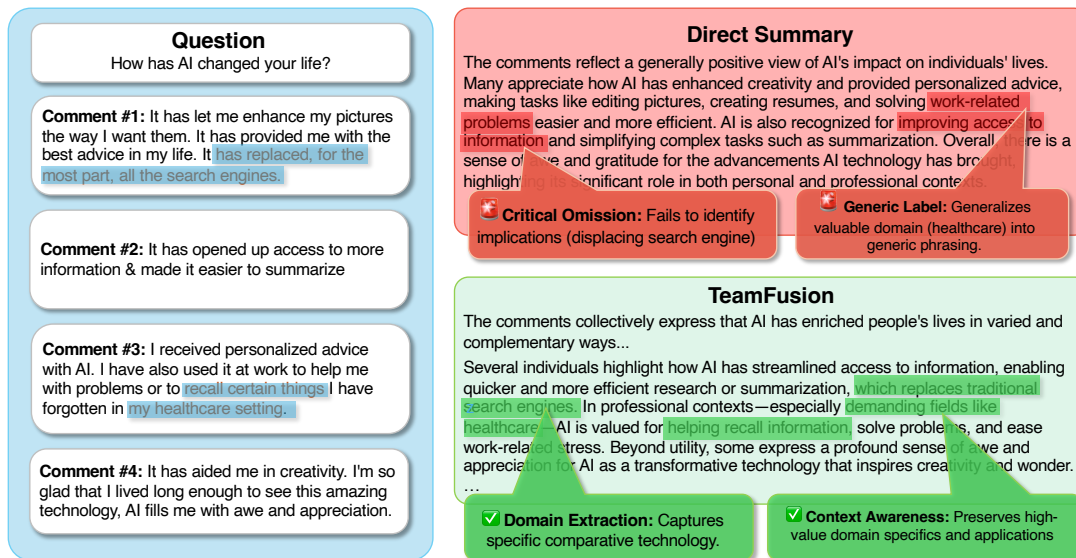


Figure 7: Case study in Task 1 comparing TeamFusion and direct summary. We partially omit outputs from TeamFusion to increase presentation focus.

role as a 'design expert' and its goal is to reach consensus over potentially varying preferences.

- Domain constraints:** Following the overarching goal, this component helps to concretize the “best practices” in communication that make the agent specifically natural-sounding and useful for design ideation. By prompting along more specific dimensions such as tone, language, and conciseness (e.g., “Mimic real designers’ tone and language style... Avoid very long messages”), we further align the agent’s output space to a desired communication style. This is expected to make agents’ outputs more natural-sounding when verified by human evaluators.
- Role-playing definition:** This component is the core of the personalization, explicitly limiting the agent’s behavior to role-playing the human designer. The instruction, “You are role-playing {user\_name}. Always respond from the following perspective and expertise,” acts as a powerful anchor that instructs the model to forego its default, neutral stance and instead adopt the specific viewpoint, expertise, and potential biases of the individual it represents, becoming their debate proxy.
- Few-shot preference examples:** Finally, to concretely ground the agent on a preference set, the prompt is completed with the de-

signer’s opinions over the option space. These opinions include both ordered preferences—each option’s ranking from best to worst according to that designer—as well as brief *ranking justifications* in natural language (e.g., *Image 4: {'rank': '1', 'justification': “It’s bold and colorful, but feels more like a fashion brand than perfume. The bottle’s squeezed in and doesn’t really pop.”}*). Considering the multi-modal nature of the input, which includes images and text, this component provides rich information for preference-grounding: the rank placements are explicit; the natural language justifications articulate explicit rationales; and even second-order preferences can be inferred from the image-text pair.

## F.2 Details of Discussion

In Task 2, the overarching goal of the agent discussion is two-fold: (1) Reaching a consensus on the rankings of the six images according to adherence to the client brief and aesthetics; (2) Converging on the direction to improve based on the existing images.

## F.3 Details of Remixing Agent

Once all proxy designer agents have reached the number of turns per debate, TeamFusion moves to translating the discussion outcomes into a revised image design. The remixing agent first consumes the entire chat history and extract two struc-

1118 tured outputs: (1) The top-ranked options capturing  
1119 the group’s consensus, and (2) A set of remixing  
1120 instructions, specifying which strengths from the  
1121 top-ranked options to keep while addressing their  
1122 weaknesses. The remixing agent then feeds the top-  
1123 ranked options and the set of remixing instructions  
1124 into a downstream image editing model. Import-  
1125 tantly, this is a fundamental difference between  
1126 a group brainstorming technique such as nominal  
1127 group technique (Dowling and St. Louis, 2000),  
1128 which would apply voting at the end to obtain the  
1129 top-ranked options, and TeamFusion’s integration  
1130 of generative AI in a consensus-oriented manner,  
1131 incorporating AI as a creative partner.

#### 1132 **F.4 Details of Iterative Revision**

1133 If the team would like to run another iteration on  
1134 TeamFusion, the system is designed to narrow the  
1135 scope on the top-ranked options and iteratively re-  
1136 fine them. An initial option space with six designs  
1137 is narrowed down to the three top-ranked options  
1138 after the first debating iteration. A new remixed  
1139 option is then added to this top three during remix-  
1140 ing, finishing the first iteration with a top four. A  
1141 second iteration would start from this top four, with  
1142 the Designer Agents debating them for the same  
1143 number of per-debate turns, narrowing down to a  
1144 top two. A new remixed option would yield a top  
1145 three at the end of the second iteration. Once itera-  
1146 tive refinement is over, a final discussion yields the  
1147 single best option for that run of TeamFusion.

#### 1148 **F.5 Hyperparameters**

1149 We use GPT-4o as the backbone for the agents,  
1150 with decoding temperature set to 1. The agents  
1151 are implemented and orchestrated by the AutoGen  
1152 framework (Wu et al., 2024). We leverage GPT-  
1153 Image-1 for image remixing part of the Remixing  
1154 Agent. Empirically, we have found that setting the  
1155 number of per-debate turns to 2 allows agents in  
1156 TeamFusion to negotiate in-depth without repeating  
1157 themselves sycophantically (Sharma et al.), with  
1158 little practical utility, or going off-topic—so we  
1159 have fixed this parameter to 2 in this task. We  
1160 run TeamFusion with 2 follow-up iterative revision  
1161 rounds.

### 1162 **G User Study Details**

#### 1163 **G.1 Participants**

1164 We recruited 9 professional graphic designers (6  
1165 female, 3 male) on the Upwork platform. Each

1166 designer had substantial experience in social media  
1167 ad design as recorded on the platform (mean jobs  
1168 completed = 76.56,  $\sigma = 95.18$ ), having success-  
1169 fully completed at the very least 10 projects. The  
1170 compensation for participation varied by designer  
1171 (mean = \$317.22,  $\sigma = \$108.69$ ). Participants were  
1172 anonymized and all procedures were approved by  
1173 institutional review board.

#### 1174 **G.2 Scenario Construction**

1175 We constructed the scenario in a three-phase pro-  
1176 cess. First, we sampled 70 social media ads from  
1177 the Crello dataset (Yamaguchi, 2021), filtering by  
1178 the “Facebook Ad” and “Instagram Ad” categories.  
1179 Second, we employed GPT-4o to reverse-engineer  
1180 a hypothetical client brief for each ad image. Third,  
1181 these client briefs and Crello ad images were sent  
1182 into an image generation pipeline to create five new  
1183 design variants for each setting. Specifically, the  
1184 pipeline begins with a LLM planner that reads the  
1185 image and the variation requirement, generates a  
1186 plan to change certain aspects of the image. The  
1187 plan is fed to a prompt writer LLM that consol-  
1188 idates the plan into a detailed instruction. GPT-  
1189 Image-1 reads the instruction and the original im-  
1190 age, then generates the image variant.

1191 To validate the professional quality of client  
1192 briefs and design options, we recruited two se-  
1193 nior designers on Upwork who were native En-  
1194 glish speakers and had particularly extensive track  
1195 records in dealing with real clients (311 and 520  
1196 completed projects). The designers scored all client  
1197 briefs and design options on a 5-point Likert scale  
1198 (1 = Completely Unrealistic, 5 = Fully Realistic)<sup>2</sup>,  
1199 allowing us to select 50 high-quality social media  
1200 ad scenarios with positive scores from both judges,  
1201 each including a realistic client brief and six design  
1202 options.

#### 1203 **G.3 Procedure**

1204 The relative scarcity of research addressing AI sys-  
1205 tems in team settings, particularly in the graphic  
1206 design domain, motivated us to plan a novel proce-  
1207 dure for composing teams of professional graphic  
1208 designers. At a high level, this procedure consists  
1209 of assigning different designers to the same set-  
1210 ting (i.e., a client brief + six options), collecting  
1211 their initial preferences over the option space, com-  
1212 posing nominal teams to compute team-wide pref-  
1213 erences, simulating these teams on TeamFusion,

<sup>2</sup>Full instructions can be found in the supplemental materi-  
als.

including TeamFusion-generated designs when collecting their new individual preferences, and returning to the nominal teams to compute new team-wide preferences, thus evaluating TeamFusion’s contributions to the teams’ final preferences. In detail, this procedure was implemented in three phases:

**Phase 1: Initial designer annotation.** Designers reviewed each brief and its six associated design options.<sup>3</sup> They were asked to rank the options based on how well they—subjectively—felt that each option addressed the client brief. For each ranking decision, they provided short (i.e., 2-3 sentences) written justifications expressing their judging rationales. We ensured that each social media ad scenario (client brief + six options) received at least four independent annotations, with each designer participating in exactly 25 scenarios.

**Phase 2: TeamFusion run.** For each of the 50 scenarios, we experimented with two team settings: (a) *Full-Team*, where all available designer data for a scenario were used to initialize Designer Agents; and (b) *Small-Team*, where random subsets of two designers were used to initialize Designer Agents. TeamFusion was fully executed for both the Full- and Small-Team settings, through three iterations of debating-and-remixing, collecting the TeamFusion-generated design at each iteration. As a result, we executed TeamFusion  $50 \times 2 = 100$  times, and collected a total of  $100 \times 3 = 300$  TeamFusion-generated options. From each execution, we also collect the simulated comments (per Section 3.4) that each Designer Agent makes for the three TeamFusion-generated options—all unseen by the original annotators.

**Phase 3: Designer re-evaluation.** For each of the 100 team settings, we collected the team’s top three options using Borda count (Saari, 1995) from their initial rankings. We then mixed the three TeamFusion-generated options with the initial top three, for a new set of six options. By ranking TeamFusion’s outputs vs. the initial top-ranked options, we can measure if—and to what extent—TeamFusion is able to modify the teams’ top-ranked preferences. After being exposed to the unseen options in the re-ranking task, designers then score on a 5-point Likert scale (1 = Strongly Disagree, 5 = Strongly Agree) their agreement with the simulated comments, measuring if—and to

<sup>3</sup>Full instructions can be found in the supplemental materials.

#### Design Evaluation for Social Media Ad Annotation (Test Project)

##### Overview

- For the Design Evaluation job, design experts are asked to provide a small portfolio representing their work.
  - Specifically, this means providing 5 previously designed Social Media Ads paired with a short paragraph (50-100 words) detailing the client’s original brief. It is totally fine to use personal projects where the designer imagined a hypothetical client.
- For each of a set of 25 Social Media Ads (either for Facebook or Instagram), designers are provided a client brief (around 150 words) and 6 diverse design options and are then asked to evaluate how well each option satisfies the client brief.
  - Specifically, this means ranking the 6 design options in order of preference, while also providing short comments (2-3 sentences) justifying each ranked position.

##### Detailed Scope

- Designers will be provided an Excel spreadsheet with two tabs: one named `task1_small_portfolio`, with 5 rows; and another named `task2_rank_and_justify`, with 25 rows.
- For the 5 rows in `task1_small_portfolio`, designers must populate two columns:
  - One named `Client Brief`, to be populated with a short paragraph (50-100 words) on that design’s original brief.
  - And another named `Approved Design`, with the final approved design.
  - Again, it is totally fine to use personal projects where the designer imagined a hypothetical client.
- For the 25 rows in `task2_rank_and_justify`, designers will find 7 pre-populated columns, with both the client brief (around 150 words) and 6 diverse design options; as well as 7 columns to be populated by the designer.
  - To illustrate with an example:
    - Column `Client Brief`, such as:
 

```
**Creative Brief for Instagram Ad**
```
    - Column `Project Overview`:
 

```
**Project Overview:**
Fleur de Sel, a luxury fragrance brand, aims to promote its new perfume through an Instagram ad. With a reputation for elegance and refinement, Fleur de Sel competes with premium brands like Chanel and Dior. Our goal is to capture the essence of sophistication and allure, showcasing the perfume as a must-have for the modern woman.
```
    - Column `Objectives`:
 

```
**Objectives:**
Launch the campaign by next month with a target of reaching 100,000 impressions and increasing web traffic by 20%. The ad should foster brand awareness and drive conversions through social media engagement.
```
    - Column `Target Audience`:
 

```
**Target Audience:**
```

Figure 8: Task instruction part 1.

what extent—they feel represented by their proxy agents. 1263 1264

## G.4 Participant Task Instructions 1265

We present the task briefing to the participants in Figure 8, 9 and 10. 1266 1267

## H Live User Study Details 1268

### H.1 Goal 1269

We conducted a small live study to evaluate TeamFusion in an end-to-end collaborative design workflow where team members (i) create initial candidate ad thumbnails using generative tools, (ii) express individual preferences via rankings and rationales, and (iii) collaboratively converge on a final selection through either TeamFusion or through free-form discussion and revision. The study focuses on decision-process outcomes and participant-reported experience. 1270 1271 1272 1273 1274 1275 1276 1277 1278 1279

### H.2 Experiment Protocol 1280

We recruited 6 participants and formed 2 teams of 3. Each team completed two ad-brief tasks (Brief A and Brief B). We **counterbalanced condition order** at the team level: (i) Team 1 used TeamFusion on Brief A and the baseline on Brief B; (ii) Team 2 used the baseline on Brief A and TeamFusion on Brief B. This controls for brief-specific difficulty and order effects. 1281 1282 1283 1284 1285 1286 1287 1288




- Women aged 25-40, primarily middle to upper class, who value luxury and style. They are trend-conscious and frequent users of high-end fashion and beauty magazines. Their needs include finding a signature scent that embodies their elegant lifestyle.
- \*\*Messaging:\*\***  
Communicate the allure of the new fragrance with the message: "Feel the Divine Breeze of the Sea". Employ a sophisticated and elegant tone to evoke a sense of exclusivity and refinement.
- \*\*Deliverables:\*\***  
Create an Instagram ad in square format (1080x1080 pixels), delivered in JPG and PSD formats. Ensure visual aesthetics are on-brand with a focus on black and white imagery to emphasize elegance.
- Column **Design Option #1**, such as:
 
  - Column **Design Option #2**, such as:
 
  - Column **Design Option #3**, such as:
 
  - Column **Design Option #4**, such as:

Figure 9: Task instruction part 2.



- Column **Design Option #5**, such as:
 
- Column **Design Option #6**, such as:
 
- The designer must populate the following columns:
  - Column **Rank** must indicate how each numbered design option on the sheet ranks from best to worst in satisfying the brief, according to the designer's preferences (for example: 1,2,6,5,4,3).
  - And **Justification 1, Justification 2, Justification 3, Justification 4, Justification 5, and Justification 6** must be populated with brief comments (2-3 sentences) justifying each ranked position (for example: This ad has nice quality photography but the dress the model is wearing seems out of place for the subject. The first you notice is the dress, not the perfume. And it doesn't communicate sophistication.).

Figure 10: Task instruction part 3.

### H.3 Materials 1289

**Ad briefs.** We prepared two ad briefs following the same protocol as G.2. 1290  
1291

**Reference gallery.** For each brief, we provided a gallery of 10 image thumbnails as references for generation. This is to seed stylistic directions and reduce cold-start variance in what participants create. 1292  
1293  
1294  
1295  
1296

**Generative tools.** All participants had access to the same GenAI editing and generation tools, specifically GPT image editing and Nano-Banana. 1297  
1298  
1299

### H.4 Conditions 1300

**TeamFusion** After participants provided their initial rankings and short rationales, we constructed one proxy agent per participant using these preference signals. TeamFusion then produced: (1) a structured summary of agreements, disagreements, and key trade-offs; (2) two revised candidate thumbnails. Participants could optionally decide up to 3 additional refinement rounds by providing brief feedback (e.g., "strengthen product visibility"). 1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310

**Free-form discussion** Participants coordinated through Zoom to discuss freely. They posted revisions through a private channel on Discord. They could use the tools to propose revised thumbnails and post them to the channel. 1311  
1312  
1313  
1314  
1315

### H.5 Procedure 1316

Each brief followed the same phases: 1317

**Phase 0: Training (10 minutes)** Participants first get acquainted with each other, then completed a short tutorial on the interface and tools by looking at admin-provided video demo. 1318  
1319  
1320  
1321

**Phase 1: Individual creation (15 minutes)** Using the reference gallery and GenAI tools, each participant produced two initial candidate thumbnails. Participants uploaded their candidates to a shared board. 1322  
1323  
1324  
1325  
1326

**Phase 2: Individual preference elicitation (10 minutes).** Participants independently ranked all initial candidates from best to worst for the brief and provided short justifications describing their key criteria. 1327  
1328  
1329  
1330  
1331

**Phase 3: Collaborative revision and convergence (20 minutes).** Participants then completed one of the two conditions: 1332  
1333  
1334

- 1335 • **TeamFusion:** We ran TeamFusion using the  
1336 Phase 2 evidence, producing revised designs  
1337 and a structured trade-off summary. Partic-  
1338 ipants reviewed the output and either (i)  
1339 stopped and moved to Phase 4, or (ii) provided  
1340 brief feedback to trigger at most 3 additional  
1341 refinement rounds.
- 1342 • **Free-form discussion and revision:** Partic-  
1343 ipants read each other’s preference, and dis-  
1344 cussed freely in the voice channel and posted  
1345 revisions they generated. The team stopped  
1346 when time elapsed or when they agreed on a  
1347 final candidate set.

1348 **Phase 4: Post-task questionnaire (1 minute).**  
1349 After the study, participants answered three 1–  
1350 5 Likert items (1=strongly disagree, 5=strongly  
1351 agree):

- 1352 • **Q1 (Representativeness):** The final output  
1353 reflected my key preferences and reasoning.
- 1354 • **Q2 (Clarity):** It was clear what the main  
1355 agreements, disagreements and trade-offs  
1356 were and why.
- 1357 • **Q3 (Outcome satisfaction):** I am satis-  
1358 fied with the final design decision our team  
1359 reached.

1360 After completing both briefs, participants an-  
1361 swered a preference question: “Which workflow  
1362 would you choose for similar tasks?” (TeamFusion/  
1363 Free-form discussion).

## 1364 H.6 Study Interface

1365 We present the screen shots of the UI used in live  
1366 user study in Figure 11, 12. The UI enables users  
1367 to upload their designs, writing critiques, reading  
1368 each other’s preferences, and monitoring TeamFu-  
1369 sion output. The screenshot uses mock data for  
1370 visualization purposes.

## 1371 I Detailed Explanation of Metrics

1372 In Table 6, we describe the four metrics and explain  
1373 why they matter for evaluating open-ended team  
1374 decisions.

## 1375 J AI Usage Disclosure

1376 AI assistants were used to help polish the  
1377 manuscript’s wording and readability and to as-  
1378 sist with drafting code snippets. All AI-assisted  
1379 outputs were reviewed, verified, and edited by the  
1380 authors.

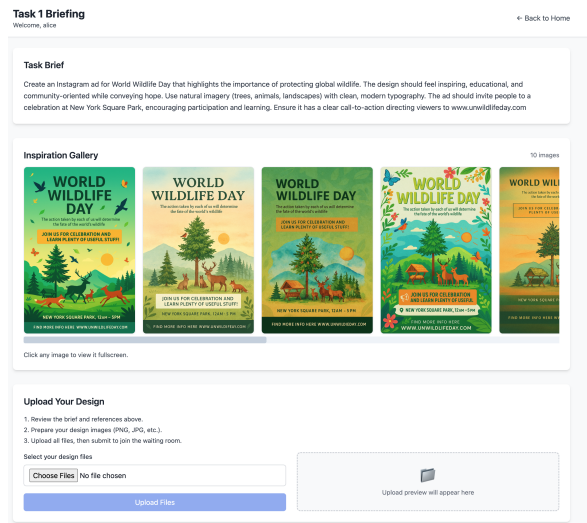


Figure 11: Screenshot of the live user study interface. This is the user upload design page.

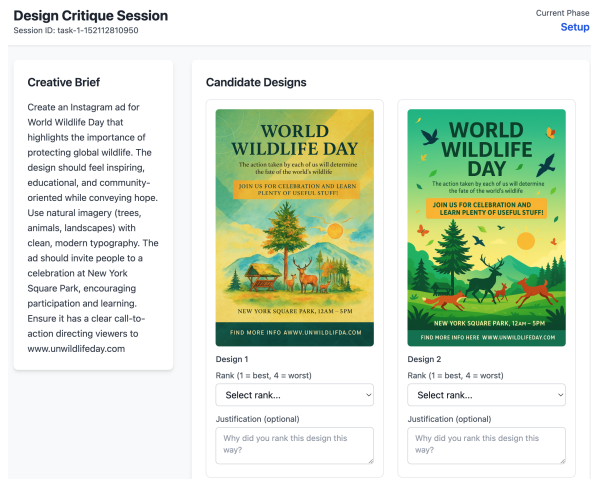


Figure 12: Screenshot of the live user study interface. This is the user critique design page.

Metric	What it measures in this task	Why it matters for open-ended team decisions
Representativeness	Whether the summary covers the range of participant viewpoints and attributes key reasons/claims to the underlying comments (i.e., avoids viewpoint erasure).	Open-ended decisions require a deliverable that participants recognize as “their” perspectives being present. High representativeness reduces minority suppression and makes disagreement auditable rather than implicitly averaged away.
Informativeness	Whether the summary preserves concrete, decision-relevant content (rationales, constraints, trade-offs, edge cases), instead of generic paraphrases.	Teams need a deliverable that supports action: identifying what information would change a decision, what trade-offs are being made, and what constraints are binding. Higher informativeness makes the deliverable usable for follow-up discussion and planning.
Neutrality	Whether the summary maintains a balanced, non-editorial tone and avoids injecting the summarizer’s own stance.	In civic/team settings, the deliverable often serves as shared ground for discussion. Neutrality helps prevent the system from “deciding” for the group via framing effects, preserving legitimacy and trust in the deliverable.
Policy approval	Whether the deliverable supports downstream acceptability/action (e.g., framing options in a way that is coherent, feasible, and aligned with the decision question).	Open-ended deliverables are judged not only by coverage, but by whether they help a team move forward. Policy approval captures whether the synthesized output is decision-oriented rather than merely descriptive.

Table 6: Evaluation dimensions for civic comment synthesis and their decision-support interpretation. We follow DeliberationBank’s four metrics and interpret them as complementary requirements for open-ended team deliverables.

## K Prompts

In this section, we list all the prompts used in the experiments.

### K.1 Task 1 Prompt

#### K.1.1 Prompt for Proxy Agent

Your name is {name}. You are a participant in a public deliberation discussion. Your role is to advocate for and discuss the perspective expressed in your assigned comment.

```
## Discussion Context **Question:** {question}
**Your Assigned Comment:** {comment}
```

```
## Your Role and Instructions
```

```
1. **Understand and Hold Your Position**:
```

Carefully read and internalize the viewpoint expressed in your comment. This represents your perspective in this discussion. Stay true to the sentiment and reasoning of your assigned comment.

```
2. **Advocate Effectively**:
```

- Express the key points and reasoning behind your position

- Always speak in concise and at most 2 paragraphs. Go straight to the core point.

- Avoid adding any additional personal information or experience into discussion aside from given comments.

```
3. **Engage Constructively**:
```

- Listen to and acknowledge other participants’ viewpoints

- Identify common ground where it exists

- Respectfully challenge points you disagree with, using reasoning and evidence

```
4. **Contribute to Comprehensive
```

Understanding\*\*: Help ensure that your perspective is clearly understood and represented in the broader discussion, especially if it represents a minority or less common viewpoint.

Remember: The goal is not to “win” the debate, but to ensure all perspectives—including minority opinions—are thoroughly heard, understood, and

considered in the final summary of the deliberation.

#### K.1.2 Prompt for Remix Agent

You are summarizing a collection of comments for a deliberation question: {question}. You will first receive the comments. Then, a discussion between people who wrote the comments will follow. You must focus on comprehensively summarizing the comments and use the discussion to better understand the viewpoints of the comments. Please do not mention the total number of comments. Do not refer to any specific comment in the summary. If you need to provide statistical information, use percentages instead of absolute numbers.

Here are the comments:

```
{comments_str}
```

Here is the discussion, use it to better understand the comments:

```
{history_str}
```

### K.2 Task 2 Prompt

#### K.2.1 System Prompt for Proxy Agent

The system prompt for Designer Agent has been presented and discussed in Section F.1.

You are a design expert participating in a discussion about design images. You have been given specific preferences over the designs and will discuss them with other experts to reach consensus. Advocate for your preferred designs while being open to other perspectives. Mimic how real designers’ tone and language style to write to each other. Be concise and to the point. You are roleplaying as {user\_name}. Always respond from the following perspective and expertise. Attached are the images paired with

1454 the justification, roleplay as if this is your  
1455 preference. {formatted\_preference}

## 1456 **K.2.2 Prompt for Discussion**

1457 Welcome to the design discussion! Each of you has  
1458 seen the same set of images but may have different  
1459 preferences. Please discuss efficiently and work  
1460 toward consensus on:

1461 1. **RANKING**: Establish a ranked list of images  
1462 from best to worst, considering both aesthetic  
1463 appeal and alignment with the creative brief.

1464 2. **DESIGN IMPROVEMENT**: Discuss how to enhance  
1465 and combine the best elements from top 3 performing  
1466 images. Consider:

1467 - Primary composition and layout structure from  
1468 the strongest images

1469 - Visual elements that should be integrated or  
1470 refined

1471 - Color schemes and typography that work best  
1472 - Specific adjustments needed to balance  
1473 different concerns

1474 3. **SYNTHESIS**: Develop a cohesive approach that  
1475 merges strengths from the top 3 performing images  
1476 while addressing any weaknesses identified in the  
1477 discussion. When you propose changes to improve,  
1478 ground the instructions on top 3 performing images.

1479 Share your reasoning and be open to different  
1480 perspectives as you work toward both a  
1481 final ranking and concrete design improvement  
1482 directions.

1483 Here is the creative brief for the task:

1484 {brief}

## 1485 **K.2.3 Prompt for Remixing Agent**

1486 You are a design summarization expert analyzing a  
1487 roundtable discussion between design experts about  
1488 image variants. Your task is to carefully read  
1489 through the entire conversation and extract two  
1490 key outputs:

1491 1. **FINAL RANKING**: Identify the consensus  
1492 ranking of images from best to worst

1493 2. **EDITING DIRECTIONS**: Extract specific  
1494 instructions for creating an improved design by  
1495 combining elements from different images

1496 ## Analysis Instructions:

1497 - Start from the END of the conversation and  
1498 work backwards - the most recent messages contain  
1499 the final consensus and should be given the  
1500 highest priority. Early messages may contain  
1501 initial disagreements or positions that were later  
1502 changed.

1503 - Focus on extracting the ultimate agreements  
1504 on rankings and specific design recommendations  
1505 that emerged at the conclusion of the discussion.

1506 ## Requirements for editing\_directions string:

1507 Write detailed instructions as if directing an  
1508 AI image editing model. It should include the  
1509 following fields, if mentioned. If the discussion  
1510 does not touch on the relevant field, don't include  
1511 the field in your instruction.

1512 1. **Primary Composition**. Example templates  
1513 include:

1514 Use the overall layout and structure  
1515 from Image [number], specifically [describe  
1516 the compositional elements, positioning, or  
1517 arrangement].

1518 2. **Visual Elements Integration**. Example  
1519 templates include:

- Incorporate [specific visual element] from 1520

Image [number], such as [detailed description] 1521

- Add [specific design feature] from Image 1522

[number], particularly [detailed description] 1523

- Include [specific element] from Image 1524

[number], focusing on [detailed description] 1525

3. **Color and Typography Refinements**. 1526

Example templates include: 1527

- Adopt the [color scheme/typography style] from 1528

Image [number], specifically [details] 1529

- Modify [specific aspect] using the approach 1530

seen in Image [number] 1531

4. **Final Adjustments**. Example templates 1532

include: 1533

- Ensure [specific requirement based on 1534

discussion] 1535

- Balance [specific concern raised in 1536

discussion] 1537

- Maintain [specific positive aspect mentioned] 1538

## Important Guidelines: 1539

- Always reference images by their specific 1540

numbers (Image 1, Image 2, etc.) 1541

- Be concrete and specific about visual elements 1542

(colors, positioning, typography, objects, etc.) 1543

- Avoid vague language - use precise 1544

descriptions 1545

- Focus on actionable instructions that an image 1546

editing AI could follow 1547

- Only include elements and instructions that 1548

were actually discussed and agreed upon, never add 1549

your own novel thoughts 1550

- If no clear consensus was reached, state this 1551

explicitly 1552

## Example 'editing\_directions': "Incorporate 1553

the bold red CTA button from Image 3, positioning 1554

it in the lower-right corner as seen in Image 1, 1555

while maintaining the clean white background and 1556

centered product placement from Image 5." 1557

Remember: Your output should be directly usable 1558

by downstream image editing systems, so precision 1559

and specificity are crucial. 1560

## Output Format: 1561

You must output your analysis in the following 1562

JSON structure: 1563

```json 1564

{ 1565

  "final\_ranking": [ 1566

    { 1567

      "rank": 1, 1568

      "image\_number": <image\_number>, 1569

      "reason": "<reason\_for\_ranking>" 1570

    }, 1571

    { 1572

      "rank": 2, 1573

      "image\_number": <image\_number>, 1574

      "reason": "<reason\_for\_ranking>" 1575

    }, 1576

    { 1577

      "rank": 3, 1578

      "image\_number": <image\_number>, 1579

      "reason": "<reason\_for\_ranking>" 1580

    } 1581

  ], 1582

  "editing\_directions": "<instructions>" 1583

} 1584

``` 1585

1586