

A. Ethics statement

Positive research and societal value. VS-Bench targets a core capability that future AI systems will increasingly need: making strategic, multi-step decisions while perceiving the world through vision and language. By standardising how this ability is measured, the benchmark can accelerate reproducible research on safer, more reliable multimodal agents. Concretely, it enables (1) principled comparisons across models, (2) diagnostic analyses that pinpoint specific failure modes such as myopic play or poor opponent modelling, and (3) a shared testbed for developing methods that foster cooperation, fairness, or robustness in complex interactive settings. Beyond academic progress, stronger decision-making agents could benefit applications like assistive household robotics, disaster-response swarms, automated traffic control, and large-scale scientific simulations where coordination and strategic planning are essential.

Risk of misuse and dual-use considerations. At the same time, more capable agents that reason strategically can be repurposed for adversarial or deceptive objectives — for example, collusive price-setting, automated disinformation campaigns, or the coordination of autonomous weapons systems. VS-Bench lowers the barrier to evaluating such capabilities, potentially making it easier to select or fine-tune models for harmful ends. To mitigate this, we will (1) release only simulated environments that do not directly embody real-world attack surfaces, (2) distribute the benchmark and evaluation code under licenses that forbid the use of our assets in weaponised or surveillance applications, and (3) encourage follow-up work on safety safeguards (e.g., opponent-aware alignment checks) by providing explicit hooks for auditing model rationales and behaviors. During review, we intentionally omit any URLs or repository links to preserve anonymity.

Human subjects protection. To minimize potential risks to human participants, we conduct the human studies with departmental approval and enroll individuals only after obtaining informed consent. The consent form explicitly stated that the games involve cooperative, competitive, and mixed-motive interactions between participants. Participants are paid for taking part in this study and can choose to withdraw from the study at any time without consequences.

Privacy and data ethics. All VS-Bench environments are synthetic with no personally identifiable information or copyrighted third-party imagery that is not permissively licensed. Replays, logs, and intermediate states are derived entirely from simulation and will be released under an open license to avoid common privacy pitfalls in dataset creation and facilitate unrestricted academic use. During review, we intentionally omit any URLs or repository links to preserve anonymity.

B. Limitations

Evaluation metric. For decision-making evaluation, we measure the episode return in self-play and interaction with conventional agents. Incorporating results against a population of diverse opponents provides a more comprehensive assessment of the generalization and adaptability of VLM agents.

Number of agents. Many real-world scenarios involve more than two participants. Although our current environments mainly focus on two-player games, some of them naturally support more than two agents. At present, we have conducted preliminary decision-making experiments in three-player settings for both *Hanabi* and the *Coin Dilemma* in Appendix G.5. We believe that extending our benchmark to include games with even more agents is a promising direction for future work.

C. Models configuration details

In this work, we evaluated the state-of-the-art VLMs released before September 1, 2025. All models in our experiments are listed in Table 5. For each model, the table specifies the exact version, whether it is a reasoning model, supports multimodal inputs, and is open-source. Because our environments provide both visual and linguistic observations, we do not evaluate models lacking multimodal input support, such as `grok-4`, `deepseek-r1`, `deepseek-v3`. At present, `claude-4-opus` is not evaluated due to its high cost (exceeding \$1000), and `grok-4-fast` is excluded because it was released on September 20, 2025, which is after our specified model release deadline.

D. Perception evaluation details

In this section, we describe the construction of the perception datasets and the evaluation protocols used to assess the perception abilities of VLMs across different environments. For each game, we construct a dataset of 400 samples.

D.1. Hanabi

The samples are uniformly drawn from the trajectories of `o3` and `gemin-2.5-pro`, the most capable reasoning models that explore a more diverse range of game states.

Table 5. Model configurations used in the evaluation.

Models	Version	Evaluated	Reasoning	Multimodal	Open-Source
o3	o3-2025-04-16	✓	✓	✓	✗
gemini-2.5-pro	gemini-2.5-pro-preview-03-25	✓	✓	✓	✗
claude-3-7-sonnet	claude-3-7-sonnet-20250219	✓	✓	✓	✗
doubao-1-5-thinking-pro	doubao-1-5-thinking-pro-m-250428	✓	✓	✓	✗
ui-tars-1-5	doubao-1-5-ui-tars-250428	✓	✓	✓	✗
qvq-max	qvq-max-2025-03-25	✓	✓	✓	✗
gpt-4.1	gpt-4.1-2025-04-14	✓	✗	✓	✗
gemini-2.5 w/o thinking	gemini-2.5-flash-preview-04-17	✓	✗	✓	✗
claude-3-7 w/o thinking	claude-3-7-sonnet-20250219	✓	✗	✓	✗
grok-2-vision	grok-2-vision-1212	✓	✗	✓	✗
doubao-1-5-vision-pro	doubao-1-5-vision-pro-250328	✓	✗	✓	✗
qwen-vl-max	qwen-vl-max-2025-04-08	✓	✗	✓	✗
Llama-3.2-90B-Ins.	Llama-3.2-90B-Vision-Instruct	✓	✗	✓	✓
InternVL3-78B	InternVL3-78B	✓	✗	✓	✓
Qwen2.5-VL-72B-Ins.	Qwen2.5-VL-72B-Instruct	✓	✗	✓	✓
claude-4-opus	N/A	✗	✓	✓	✗
grok-4	N/A	✗	✓	✗	✗
grok-4-fast	N/A	✗	✓	✓	✗
deepseek-r1	N/A	✗	✓	✗	✓
deepseek-v3	N/A	✗	✗	✗	✓

Given the rich information present in the *Hanabi* environment, we select four key elements as evaluation targets: *life tokens*, *info tokens*, *fireworks*, and *cardinfo*. For each sample, a screenshot of the game state (e.g., Fig. 7) is provided, and the model is prompted to infer the values of the specified elements. Specifically, *fireworks* consists of five color-based stacks, while *cardinfo* represents ten cards (five per player), each requiring the model to output 11 units of information: card value (known for the opponent’s hand, but “?” for the player’s own hand), and the color/number knowledge that each player currently has about the cards. In total, each sample requires predicting 117 units of information: 1 for *life tokens*, 1 for *info tokens*, 5 for *fireworks*, and 110 (11 × 10) for *cardinfo*. Each unit is equally weighted, and the final accuracy is normalized to the range [0, 1].

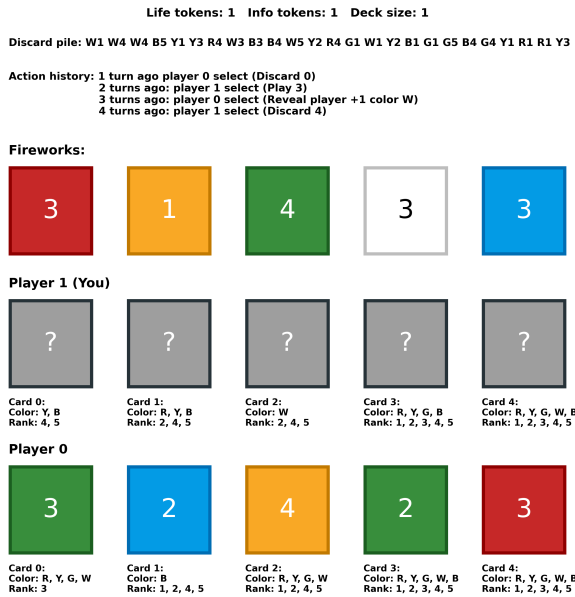


Figure 7. *Hanabi*.

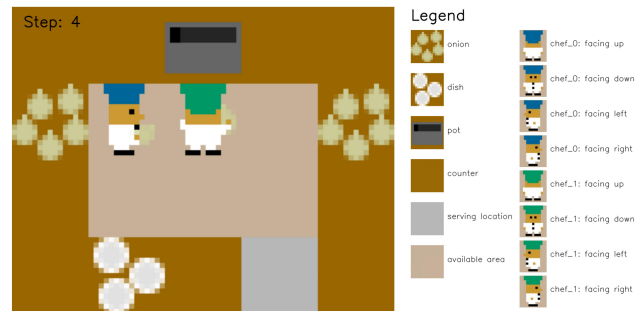


Figure 8. *Overcooked*.

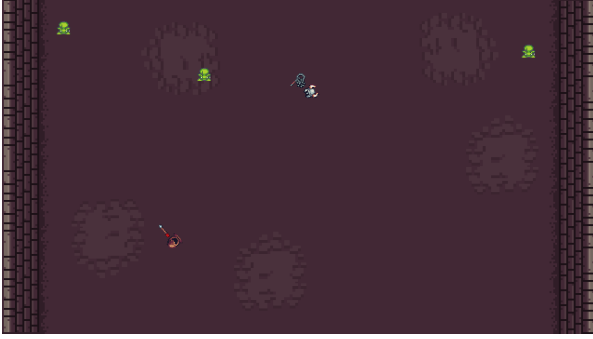


Figure 9. Origin KAZ.

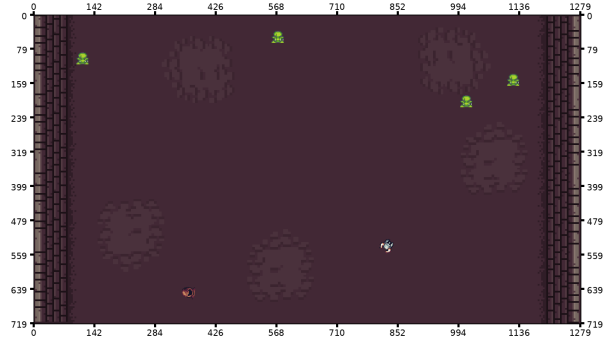


Figure 10. KAZ with coordinate axes.

D.2. Overcooked

This dataset is derived from the human experiment data provided by Overcooked-AI [12], which contains gameplay trajectories recorded from multiple human participants.

For the *Overcooked* environment, each state is represented as a 4×5 grid. For each sample, a screenshot corresponds to this 4-row by 5-column layout (e.g., Fig. 8), and the model is prompted to infer the content of each grid cell, producing a 4×5 matrix as output.

Each grid cell may contain objects such as tables, pots, onions, dishes, service desks, or empty areas. To simplify model predictions, we map each possible object to a symbolic representation (e.g., \mathbb{P} stands for pot). The model’s score is computed as the number of correctly predicted grid cells out of 20, normalized to the range $[0, 1]$.

D.3. Knights Archers Zombies

The samples are uniformly drawn from trajectories generated by a `random` agent. This choice is motivated by the fact that the game may involve a large number of zombies. To simplify perception for VLMs, we limit the maximum number of zombies to five. Under this constraint, we observe that the `random` agent explores a broader distribution of game states compared to other VLMs, thereby providing a more comprehensive benchmark for evaluating perception ability.

The original KAZ environment is continuous (e.g., Fig. 9). To facilitate visual recognition by VLMs, we implement an additional wrapper that overlays coordinate axes and scales on the environment. For each sample, a screenshot of the game state with coordinate axes is provided (e.g., Fig. 10), and the model is prompted to infer the following elements: the coordinates of the archer and the knight, the number of zombies present (up to five), and the positions of all detected zombies. The four elements are weighted as follows: 0.2, 0.2, 0.2, and 0.4, respectively.

For zombie count, exact matching is required to receive full credit. For coordinate predictions, full credit is awarded if the predicted position lies within a predefined radius centered on the ground-truth location. Beyond this radius, the score decays proportionally to the distance from the correct position. In our setting, the radius is fixed at 100.

D.4. Breakthrough

The samples are uniformly drawn from trajectories generated by a `random` agent. This choice is motivated by the observation that it can explore board states that other VLMs are unlikely to reach, thereby providing a more comprehensive benchmark for evaluating perception ability.

In the *Breakthrough* environment, each state is represented as an 8×8 grid. Accordingly, a screenshot of the game state (e.g., Fig. 11) is provided for each sample, and the model is prompted to infer the content of each grid cell, ultimately producing an 8×8 matrix as output.

Each grid cell may contain one of three possible objects: `b` for a black piece, `w` for a white piece, and `.` for an empty cell. The model’s score is computed as the number of correctly predicted grid cells out of 64, normalized to the range $[0, 1]$.

D.5. Kuhn poker

The samples are uniformly drawn from trajectories generated by `random` agent, since the overall state space is relatively small and a random agent is able to traverse it more effectively.

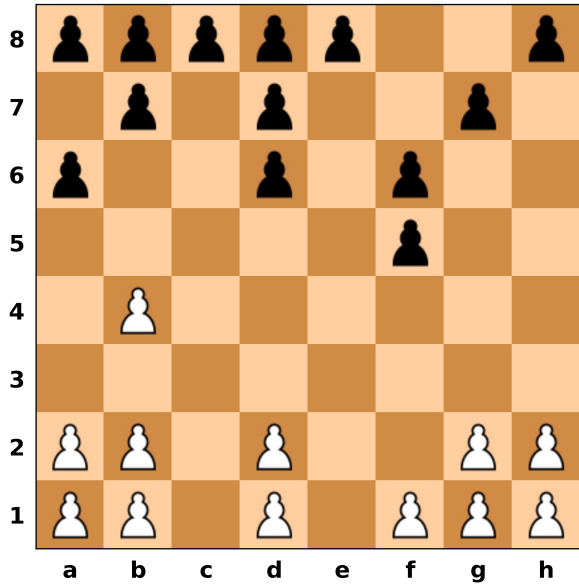


Figure 11. Breakthrough.

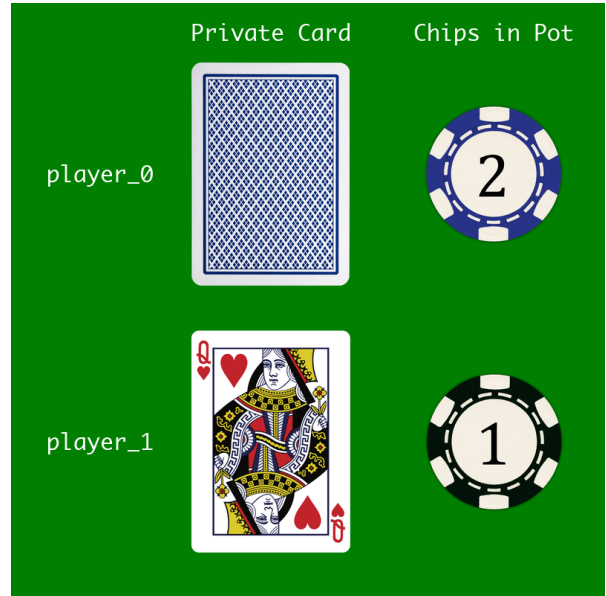


Figure 12. Kuhn Poker.

In *Kuhn poker*, the environment contains only four variables: the card values of both players (with the opponent’s card being hidden and thus marked as `unknown`), and the chips in pot of both players (e.g., Fig. 12). Accordingly, the model is prompted to infer these four elements, each assigned an equal weight, and the final score is normalized to the range $[0, 1]$.

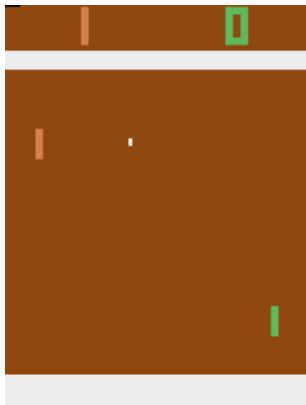


Figure 13. Origin Atari Pong.

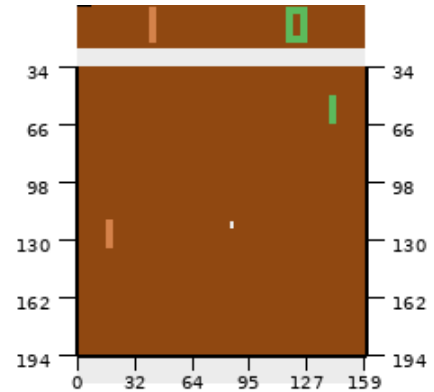


Figure 14. Atari Pong with coordinate axes.

D.6. Atari Pong

The samples are uniformly drawn from trajectories generated by `o4-mini` and `doubao-1-5-thinking-pro`, both of which exhibit strong performance in the game.

The original *Atari Pong* environment is continuous (e.g., Fig. 13). To facilitate visual recognition by VLMs, we implement an additional wrapper that overlays coordinate axes and scales on the environment. For each sample, a screenshot of the game state with coordinate axes is provided (e.g., Fig. 14), and the model is prompted to infer the following elements: the positions of the two paddles, the position of the ball, and the current score, with weights of $4 : 4 : 2$, respectively.

Since the paddles move only vertically, the model is required to predict their y -coordinates only. For the ball, the model is required to predict both x - and y -coordinates. A prediction receives full credit if it lies within a predefined radius of the ground truth; beyond this radius, the score decays with distance. The radius is fixed at 10 in our setting.

D.7. MPE

The samples are uniformly drawn from trajectories generated by random agents, which provide a broader distribution of game states compared to other VLMs.

In the *MPE* environment, there are only two variables representing the coordinates of the two players and a target fixed at the origin (e.g., Fig. 15). Accordingly, the model is prompted to infer the x - and y -coordinates of both players. For coordinate predictions, full credit is awarded if the predicted position lies within a predefined radius centered on the ground-truth location. Beyond this radius, the score decays proportionally to the distance from the correct position. In our setting, the radius is set to half of the unit spacing of the coordinate axis.

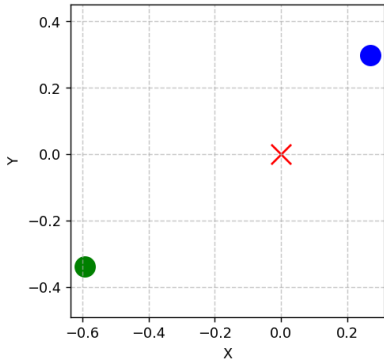


Figure 15. *MPE*.



Figure 16. *Coin Dilemma*.

D.8. Coin Dilemma

The sampling strategy and dataset construction follow the same procedure as in the strategic reasoning setting (Appendix E.8).

In the *Coin Dilemma* environment, each state is represented as a 5×5 grid. Accordingly, a screenshot of the game state (e.g., Fig. 16) is provided for each sample, and the model is prompted to infer the content of each grid cell, ultimately producing a 5×5 matrix as output.

Each grid cell may contain the following objects: an empty cell, a red coin, a blue coin, a red player, or a blue player. The model’s score is computed as the number of correctly predicted grid cells out of 25, normalized to the range $[0, 1]$.

D.9. Monster Hunt

The sampling strategy and dataset construction follow the same procedure as in the strategic reasoning setting (Appendix E.9).

In the *Monster Hunt* environment, each state is represented as a 5×5 grid. Accordingly, a screenshot of the game state (e.g., Fig. 17) is provided for each sample, and the model is prompted to infer the content of each grid cell, ultimately producing a 5×5 matrix as output.

Each grid cell may contain the following objects: an empty cell, an apple, a red player, a blue player, or a monster. The model’s score is computed as the number of correctly predicted grid cells out of 25, normalized to the range $[0, 1]$.

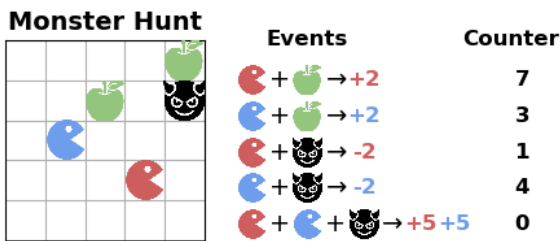


Figure 17. *Monster Hunt*.

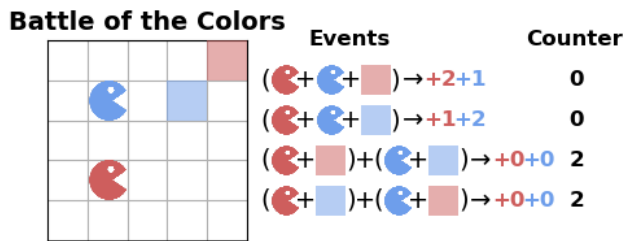


Figure 18. *Battle of the Colors*.

D.10. Battle of the Colors

The sampling strategy and dataset construction follow the same procedure as in the strategic reasoning setting (Appendix E.10).

In the *Battle of the Colors* environment, each state is represented as a 5×5 grid. Accordingly, a screenshot of the game state (e.g., Fig. 18) is provided for each sample, and the model is prompted to infer the content of each grid cell, ultimately producing a 5×5 matrix as output.

Each grid cell may contain the following objects: an empty cell, a red block, a blue block, a red player or a blue player. The model’s score is computed as the number of correctly predicted grid cells out of 25, normalized to the range $[0, 1]$.

E. Strategic reasoning evaluation details

In this section, we describe the dataset generation procedure for each environment to evaluate the strategic reasoning of VLMs. Each game dataset contains 400 samples.

E.1. Hanabi

The dataset is generated in three steps. First, to include data with different gameplay levels, 90% of the samples are obtained through mutual prediction between two reasoning models, while the remaining 10% are generated by a chat model predicting a reasoning model’s move. We chose `doubao-1-5-thinking-pro` and `doubao-1-5-vision-pro`, both of which demonstrated top-tier decision-making performance, to represent the reasoning and chat models, respectively. Second, to balance the ratio of different actions to approximate human gameplay, we analyze the gameplay data of leading VLM agents and set the ratio of `<Play>`:`<Discard>`:`<Reveal>` to 2 : 3 : 4 in our dataset. Third, we balance the agent order and step index in the dataset so that the first and second agents each account for 50% samples and the step index is uniformly distributed across the entire game sequence.

E.2. Overcooked

This dataset is derived from the human experiment data provided by Overcooked-AI [12], which comprises game trajectories recorded from multiple human participants. We focus on the trial-train subset of these data, with instances of invalid actions filtered out from the trajectories. We then use random sampling on these filtered trajectories to ensure comprehensive coverage of possible game states. Each sample comprises a sequence of four consecutive game frames. Additionally, we apply constraints to simulate realistic game scenarios and control the distribution of target actions. Specifically, the proportion of `<STAY>` action among the possible actions is limited to 10%. The dataset is balanced for two chefs, each accounting for 50% samples of the dataset.

E.3. Knights Archers Zombies

The dataset is uniformly sampled from the logged trajectories of the two best-performing models in the decision-making process: `o4-mini` and `doubao-1-5-thinking-pro`. The next actions taken by these VLM agents are used as ground truth labels. Since the environment involves two types of agents, the knight and the archer, we ensure an equal split of samples, with 50% corresponding to each agent type.

E.4. Breakthrough

All samples are generated by a minimax algorithm with alpha-beta pruning, a widely adopted baseline in *Breakthrough* research [39, 57]. Since minimax search does not always reach terminal positions to determine win-loss outcomes, we implement a state evaluation function: upon reaching a fixed search depth, we compute the difference between the maximum effective forward advancement of our deepest piece and that of the opponent’s deepest piece, then normalize this difference to obtain a reward for the state. We configure minimax agents with maximum search depths for the first and second players as (3, 4), (3, 5), (4, 5), (4, 6), (4, 4), and (5, 5), respectively, and sample step indices uniformly to ensure coverage of diverse game states.

E.5. Kuhn poker

Since *Kuhn Poker* admits a continuum of mixed Nash equilibria parameterized by a single probability $\alpha \in [0, 1/3]$, representing the likelihood of betting when holding a Jack [31], we evaluate a discrete subset of strategies. Specifically, we consider all pairwise matchups among three representative values of $\alpha = 0, 1/6, 1/3$, yielding nine distinct strategy combinations. For each combination, we simulate 600 head-to-head games and uniformly sample 400 game states in total to construct the final dataset.

E.6. Atari Pong

We uniformly sample 400 state transitions from the logged trajectories of two strong reasoning models in the decision-making process, namely `o4-mini` and `doubao-1-5-thinking-pro`. The subsequent actions taken by these VLM agents are treated as ground-truth labels. To construct the evaluation setting, we modify the prompts such that the VLMs are required to control the left paddle (originally the built-in bot's paddle) and predict the corresponding actions.

E.7. MPE

The MPE environment consists of two agents navigating a two-dimensional plane. We uniformly sample 400 state transitions from the logged trajectories of three strong reasoning models in the decision-making process, namely `o4-mini`, `claude-3.7-sonnet`, and `doubao-1-5-thinking-pro`. The subsequent actions taken by these VLM agents are treated as ground-truth labels. It is important to emphasize that we only log trajectories in which the VLMs assume the blocker role; consequently, the prediction task is restricted to the blocker's actions.

E.8. Coin Dilemma

We consider two types of heuristic strategies for playing *Coin Dilemma* and generate the dataset by simulating game play with these strategies:

1. **Common welfare**: the agent only collects the coin of its own color.
2. **Self interest**: the agent always collects the closest coin, regardless of the color.

Concretely, we collect samples from six settings, resulting in a dataset of 400 samples:

1. **Common welfare + common welfare**: collect 100 samples.
2. **Self interest + self interest**: collect 100 samples.
3. **Common welfare + self interest**: collect 50 samples.
4. **Self interest + common welfare**: collect 50 samples.
5. **Random + self interest**: collect 50 samples.
6. **Self interest + random**: collect 50 samples.

For *Coin Dilemma*, different actions can lead to the same transition in the environment. Therefore, the actions with the same outcome as the ground truth action are all considered correct. For example, if a player is at the top-left corner of the grid map, then actions `<UP>` and `<LEFT>` are both considered correct with ground truth `<STAY>` as they all result in no movement of the player.

E.9. Monster Hunt

We consider four types of heuristic strategies for playing *Monster Hunt* and generate the dataset by simulating game play with these strategies:

1. **Common welfare 1**: the agent always moves towards the monster.
2. **Common welfare 2**: the agent first moves to the middle block of the grid map and stays there to wait for the other agent and the monster.
3. **Common welfare 3**: the agent first moves to a corner of the grid map and stays there to wait for the other agent and the monster.
4. **Self interest**: the agent always moves towards the closest apple.

Concretely, we collect samples from six settings, resulting in a dataset of 400 samples:

1. **Common welfare 1 + common welfare 1**: collect 80 samples.
2. **Common welfare 2 + common welfare 2**: collect 80 samples.
3. **Common welfare 3 + common welfare 3**: collect 80 samples.
4. **Self interest + self interest**: collect 80 samples.
5. **Random + self interest**: collect 40 samples.
6. **Self interest + random**: collect 40 samples.

For *Monster Hunt*, different actions with the same outcome as the ground truth action are all considered correct, as in *Coin Dilemma*.

E.10. Battle of the Colors

We consider four types of heuristic strategies for playing *Battle of the Colors* and generate the dataset by simulating game play with these strategies:

1. **Common welfare:** the agent always moves to the closest color block (to both players) and stays there to wait for the other player.
2. **Self interest:** the agent always moves to the block of its own color.
3. **Biased red:** the agent always moves to the red block.
4. **Biased blue:** the agent always moves to the blue block.

Concretely, we collect samples from six settings, resulting in a dataset of 400 samples:

1. **Common welfare + common welfare:** collect 100 samples.
2. **Self interest + self interest:** collect 100 samples.
3. **Common welfare + self interest:** collect 50 samples.
4. **Self interest + common welfare:** collect 50 samples.
5. **Biased red + biased red:** collect 50 samples.
6. **Biased blue + biased blue:** collect 50 samples.

For *Battle of the Colors*, different actions with the same outcome as the ground truth action are all considered correct, as in *Coin Dilemma*.

F. Decision-making evaluation details

In this section, we describe the methodology for evaluating the decision-making abilities of VLMs by assigning appropriate rewards after each step of interaction.

F.1. Hanabi

In *Hanabi*, we consider two kinds of return.

Standard return. This is the standard return in existing *Hanabi* environments [7], and we use it in the main text results. If all life tokens are consumed before the fireworks are completed, the agents get a reward of 0. If all fireworks stacks are built successfully, the agents get a return of 25. If the above two early terminal conditions are not reached and the deck is exhausted, the game continues for one additional round, and the return is the sum of the ranks of the fireworks.

Firework return. Since many VLM agents often consume all the life tokens and result in a zero return, we also consider another return that relaxes the “zero-out” penalty upon losing all life tokens. More specifically, if all life tokens are consumed before the fireworks are completed, the return is the sum of the ranks of the firework piles. This return measures the progress of the fireworks and is reported in the raw decision-making results without normalization in Table 6.

For each model, we perform 10 self-play games and report the average standard return and the firework return. These results are then normalized and compared with a random baseline and an optimal policy as reported in Multi-Agent PPO (MAPPO) [83].

F.2. Overcooked

In *Overcooked*, each episode is limited to 50 steps. Within these steps, two chefs cooperatively cook soup and deliver the cooked soup to the service desk. The two agents share a common return, which includes three process rewards and a delivery reward.

Process rewards. To measure the progress of making a soup, the agents get 2 reward for each of the three events: (1) An agent successfully adds an onion to a cooking pot; (2) An agent picks up a dish when a pot contains onions or cooking is in progress; (3) An agent successfully plates a finished soup using a dish. Completing a three-onion soup leads to a total of 10 process rewards.

Delivery reward. To align with the game’s goal of successfully delivering specified dishes, the agents get 10 reward upon successful delivery to the service desk.

For the 3-onion soup recipe, successfully completing and delivering one soup yields a total return of 20, consisting of a process reward of 10 and a final reward of 10. We evaluate each VLM over 10 episodes of self-play, in which both agents employ the same model. The accumulated reward across an episode is reported as the episode return, and scores are further normalized with respect to the returns achieved by random and optimal agents. Notably, the optimal agents are able to complete two deliveries within a single episode.

F.3. Knights Archers Zombies

In *Knights Archers Zombies* (KAZ), the return is defined as the total number of zombies eliminated by the end of the game. The game terminates under two possible conditions: (i) reaching the predefined maximum horizon of 200 steps, or (ii) when either of the two players is contacted by a zombie and dies. Although the environment inherently supports multi-agent

settings, we adopt a two-agent configuration in the decision-making context considered in this work, with one archer and one knight. Moreover, a new zombie is spawned every 10 steps, implying that optimal agents can achieve up to 20 points in a single episode.

F.4. Breakthrough

In *Breakthrough*, the return is determined by the final outcome: the winning agent gets a reward of +1, the losing agent gets a reward of -1, and draws are not possible. We selected a moderately strong MCTS agent as our baseline, configured with an exploration constant $c = 2.0$, a maximum of 100 simulations per move, and a rollout count of 10. Each model is evaluated for 20 games against this MCTS agent, with 10 games as the first agent and 10 games as the second agent. For the optimal policy, we use a minimax agent with alpha-beta pruning and a maximum search depth of 5, using a state evaluation function as described in Appendix E.4. Although minimax is not guaranteed to be optimal for *Breakthrough*, it achieved a perfect win rate against the MCTS agent in our experiment, making it a reasonable choice as the optimal policy in this study.

F.5. Kuhn poker

In *Kuhn Poker*, the return of each agent is the net chips won or lost at the end of the game. We consider a Nash equilibrium (NE) strategy as described in Appendix E.5 with $\alpha = 0$. Each model is evaluated for ten runs of 120 games against this NE agent, with 60 games as the first agent and 60 games as the second agent. We then calculate the mean and standard deviation of the ten runs and normalize the result. The optimal agent in *Kuhn poker* is the agent with the NE strategy.

F.6. MPE

In the *MPE* environment, there are two types of players: a *charger* and a *blocker*. Since it is a competitive game, we designate the VLMs as the blocker, competing against a built-in agent that plays as the charger. The built-in agent is rule-based, with a strategy that always moves directly toward the target.

The reward function of the blocker consists of two components:

Push Charger: defined as the distance between the charger and the target.

Near Target: defined as the negative distance between the blocker and the target.

In summary, the blocker is incentivized to stay as close to the target as possible while simultaneously keeping the charger as far away from the target as possible. In our setting, we adopt a heuristic value of 0 as the optimal baseline.

F.7. Atari Pong

In *Atari Pong*, we adopt frame stacking of 4 frames to pass dynamic information to the VLM agent. We also employ a sticky action probability of 0.25 and perform a random number (between 1 and 30) of <STAY> steps at the beginning of an episode to achieve randomness. These settings have been common practice in related works, such as DQN [44]. We also consider two kinds of reward.

Score reward. The agent scores one point when the ball passes the opponent’s paddle and gets a score reward of +1. The game ends when an agent scores three points.

Step reward. As many VLMs fail to score even one point against the built-in bot, the score reward alone becomes too sparse to distinguish the performance of different models. Therefore, we design a continuous reward proportional to the number of steps that VLMs survived against the built-in bot.

We evaluate each VLM for 10 episodes against the built-in bot, and the overall return is the weighted sum of the normalized score reward with a 0.9 weight and the normalized step reward with a 0.1 weight. The optimal agent is an RL agent that scores three points against the built-in bot.

F.8. Coin Dilemma

In *Coin Dilemma*, the agents receive rewards on different game events:

1. **Red agent collects a red coin:** red agent gets +1 reward.
2. **Red agent collects a blue coin:** red agent gets +1 reward, blue agent get -2 reward.
3. **Blue agent collects a blue coin:** blue agent gets +1 reward.
4. **Blue agent collects a red coin:** blue agent gets +1 reward, red agent get -2 reward.

We evaluate the return of each VLM for 10 episodes of self-play, where the red and blue agents use the same model. We further normalize these returns with respect to the returns of the random agents and the optimal agents that always move towards the coin of their own color.

F.9. Monster Hunt

In *Monster Hunt*, the agents receive rewards on different game events:

1. **Red agent eats an apple:** red agent gets +2 reward.
2. **Blue agent eats an apple:** blue agent gets +2 reward.
3. **Red agent encounters the monster alone:** red agent gets -2 reward.
4. **Blue agent encounters the monster alone:** blue agent gets -2 reward.
5. **Both agents defeat the monster together:** both agents get +5 reward.

We evaluate the return of each VLM for 10 episodes of self-play, where the red and blue agents use the same model. We further normalize these returns with respect to the returns of the random agents and the optimal agents that always move to the middle block and stay there to wait for the monster.

F.10. Battle of the Colors

In *Battle of the Colors*, the agents receive rewards on different game events:

1. **Both players on the red block:** red agent gets +2 reward, blue agent gets +1 reward.
2. **Both players on the blue block:** blue agent gets +2 reward, red agent gets +1 reward.
3. **Two players on different blocks:** both agents get 0 reward.

We evaluate the return of each VLM for 10 episodes of self-play, where the red and blue agents use the same model. We further normalize these returns with respect to the returns of the random agents and the optimal agents that always move to the closest block to the two players.

F.11. Raw results without normalization

The decision-making results in all environments without normalization are listed in Table 6. We also report the firework return for *Hanabi* and the score reward and step reward for *Atari Pong*.

Table 6. Raw decision-making results without normalization. The *Hanabi* standard¹ and firework² returns are described in Appendix F.1, while the *Pong* score³ and step⁴ rewards are described in Appendix F.7.

Models	Cooperative			Competitive			Mixed-Motive					
	<i>Hanabi</i> standard ¹	<i>Hanabi</i> firework ²	Overcooked	KAZ	Board	Poker	score ³	step ⁴	MPE	Dilemma	Hunt	Battle
Optimal	24.0	24.0	40.0	20.0	1.0	0.0	3.0	398.0	0.0	27.9	176.2	59.7
o3	13.4±5.0	13.6±4.5	6.4±1.2	6.6±3.1	0.6±0.8	-0.0±0.0	0.2±0.4	219.4±113.4	-6.5±9.3	-0.3±4.9	27.0±6.3	10.2±3.1
o4-mini	10.3±7.3	13.3±2.9	7.0±2.7	3.7±2.5	-0.4±0.9	-0.0±0.0	0.3±0.5	205.2±91.0	-6.6±8.3	-1.5±6.0	28.6±16.1	2.4±3.2
geminini-2.5-pro	7.9±8.5	11.6±5.2	7.0±2.2	2.6±1.6	0.1±1.0	-0.1±0.1	0.1±0.3	246.7±65.2	-6.1±8.0	-2.9±5.5	22.0±11.7	20.4±3.8
geminini-2.5-flash	6.5±8.6	10.7±5.3	3.6±2.1	3.7±1.7	-0.6±0.5	-0.1±0.1	0.0±0.0	194.4±53.2	-7.7±8.9	2.6±7.2	31.2±11.4	19.8±5.1
claude-3.7-sonnet	1.6±5.1	9.7±3.9	4.2±1.4	2.0±1.4	-0.6±0.8	-0.0±0.0	0.0±0.0	133.7±27.8	-5.9±8.5	1.1±4.3	18.8±6.9	1.8±2.8
ui-tars-1-5	0.0±0.0	4.9±1.7	1.0±1.3	0.1±0.3	-0.9±0.3	-0.1±0.0	0.0±0.0	121.4±8.4	-11.2±12.9	9.1±4.0	28.6±6.5	8.4±1.7
doubao-1.5-thinking-pro	13.6±5.5	14.1±4.0	4.2±1.9	3.7±1.7	-0.8±0.4	-0.0±0.0	0.0±0.0	230.5±72.6	-6.2±10.3	0.0±0.9	13.6±22.3	2.7±2.8
qvq-max	0.0±0.0	4.9±2.9	1.0±1.3	0.7±0.9	-0.9±0.3	-0.1±0.0	0.0±0.0	158.2±46.9	-6.1±8.3	-0.2±0.6	-18.8±8.8	0.0±0.0
geminini-2.5 w/o thinking	0.0±0.0	3.8±1.6	1.0±1.6	1.0±0.8	-1.0±0.0	-0.1±0.0	0.0±0.0	175.9±41.5	-10.0±7.7	-0.4±1.2	-18.8±17.4	1.8±2.0
claude-3.7 w/o thinking	0.0±0.0	2.9±0.9	1.0±1.6	1.7±1.7	-0.9±0.3	-0.1±0.1	0.0±0.0	121.4±8.4	-14.8±14.0	0.2±2.6	-19.8±16.0	0.9±1.4
gpt-4.1	0.0±0.0	3.6±1.4	0.0±0.0	0.0±0.0	-1.0±0.0	-0.1±0.0	0.0±0.0	151.8±41.5	-8.0±6.6	4.8±1.9	1.8±10.9	0.6±1.2
qwenvl-max	0.3±0.5	0.3±0.5	0.0±0.0	0.0±0.0	-1.0±0.0	-0.1±0.1	0.0±0.0	139.7±29.3	-12.5±12.0	-0.3±0.8	5.8±39.6	0.0±0.0
doubao-1.5-vision-pro	0.0±0.0	4.6±1.1	0.0±0.0	0.0±0.0	-1.0±0.0	-0.1±0.1	0.0±0.0	121.4±8.4	-7.7±9.3	-0.8±1.5	-4.8±16.1	0.0±0.0
grok-2-vision	0.0±0.0	1.6±1.0	0.8±1.3	1.0±1.4	-1.0±0.0	-0.1±0.1	0.2±0.4	152.6±45.9	-31.6±24.7	0.1±2.0	-21.0±11.4	0.6±1.2
Qwen2.5-VL-72B-Ins.	0.2±0.4	0.2±0.4	0.0±0.0	0.0±0.0	-1.0±0.0	-0.1±0.1	0.0±0.0	123.8±4.6	-11.1±10.5	-0.2±0.8	18.2±50.4	0.0±0.0
InternVL3-78B	0.0±0.0	2.4±1.0	0.2±0.6	0.3±0.5	-1.0±0.0	-0.1±0.1	0.0±0.0	121.4±8.4	-14.8±14.8	1.2±2.0	-20.6±10.8	1.2±1.5
Llama-3.2-90B-Vision-Ins.	0.0±0.0	1.2±1.6	0.8±1.3	1.0±0.8	-1.0±0.0	-0.1±0.1	0.0±0.0	121.4±8.4	-24.1±17.6	-0.1±0.9	-13.2±9.6	0.9±1.4
Random	0.0	1.2	0.2	1.0	-1.0	-0.1	0.0	147.2	-16.3	-0.2	-20.2	0.3

G. Additional experiment results

Table 7. Strategic reasoning results on multimodal input and CoT prompting.

Model	<i>Hanabi</i>			<i>Breakthrough</i>			<i>Monster Hunt</i>		
	Text-Only	Multimodal	CoT	Text-Only	Multimodal	CoT	Text-Only	Multimodal	CoT
Optimal	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
o3	58.2	61.2	N/A	31.0	29.0	N/A	50.8	57.8	N/A
o4-mini	53.8	58.2	N/A	27.5	26.8	N/A	47.8	43.8	N/A
gemin-2.5-pro	54.2	51.2	N/A	28.2	26.8	N/A	55.0	35.8	N/A
gemin-2.5-flash	42.0	37.0	N/A	23.5	23.2	N/A	49.0	50.2	N/A
claude-3-7-sonnet	45.0	39.0	N/A	25.0	24.2	N/A	56.0	42.5	N/A
doubao-1-5-thinking-pro	34.5	32.8	N/A	23.5	19.8	N/A	57.2	45.2	N/A
ui-tars-1-5	24.2	25.5	N/A	27.8	23.3	N/A	34.0	41.5	N/A
qvq-max	41.0	32.2	N/A	27.5	21.8	N/A	48.2	21.5	N/A
gemin-2.5 w/o thinking	24.5	21.5	24.0	20.5	14.8	21.5	40.2	30.0	45.0
claude-3-7 w/o thinking	19.2	9.8	32.8	19.2	18.0	19.0	37.0	26.0	43.8
gpt-4.1	40.0	23.0	49.8	20.5	22.5	27.5	41.0	36.8	34.5
qwen-vl-max	17.0	26.5	20.0	19.0	19.5	17.2	31.8	23.5	35.5
doubao-1-5-vision-pro	19.5	15.0	25.2	17.2	15.8	16.8	33.2	36.0	39.2
grok-2-vision	23.8	12.8	22.5	14.0	10.8	18.2	42.5	31.5	33.2
Qwen2.5-VL-72B-Ins.	18.5	26.8	22.2	19.2	23.8	16.5	32.2	27.2	37.5
InternVL3-78B	26.8	25.2	20.5	17.5	14.0	16.0	34.2	30.0	37.8
Llama-3.2-90B-Vision-Ins.	26.8	20.0	14.8	6.5	11.8	14.0	36.8	26.2	27.8
Random	8.8	8.8	8.8	4.3	4.3	4.3	29.3	29.3	29.3

Table 8. Decision-making results on multimodal input and CoT prompting.

Model	<i>Hanabi</i>			<i>Breakthrough</i>			<i>Monster Hunt</i>		
	Text-Only	Multimodal	CoT	Text-Only	Multimodal	CoT	Text-Only	Multimodal	CoT
Optimal	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
o3	62.9±7.1	55.8±20.9	N/A	95.0±15.0	80.0±40.0	N/A	45.3±7.7	24.0±3.2	N/A
o4-mini	37.1±26.1	42.9±30.5	N/A	30.0±47.0	30.0±47.0	N/A	23.2±4.2	24.9±8.2	N/A
gemin-2.5-pro	22.5±36.7	32.9±35.5	N/A	55.0±49.5	55.0±49.5	N/A	50.4±9.0	21.5±6.0	N/A
gemin-2.5-flash	40.8±21.9	27.1±36.0	N/A	30.0±42.3	20.0±25.7	N/A	15.63±4.6	26.2±5.8	N/A
claude-3-7-sonnet	33.8±35.8	6.7±21.1	N/A	45.0±50.0	20.0±39.8	N/A	27.3±4.6	19.9±3.5	N/A
doubao-1-5-thinking-pro	37.5±32.9	56.7±22.8	N/A	15.0±37.0	10.0±21.0	N/A	27.6±6.1	17.2±11.3	N/A
ui-tars-1-5	0.0±0.0	0.0±0.0	N/A	20.0±40.0	30.0±47.0	N/A	12.6±5.1	24.8±3.3	N/A
qvq-max	0.0±0.0	0.0±0.0	N/A	5.0±15.8	5.0±15.8	N/A	15.1±5.7	0.7±4.5	N/A
gemin-2.5 w/o thinking	0.0±0.0	0.0±0.0	3.3±10.5	0.0±0.0	0.0±0.0	20.0±39.8	5.1±4.9	0.7±8.9	6.3±9.8
claude-3-7 w/o thinking	0.0±0.0	0.0±0.0	0.0±0.0	5.0±15.8	5.0±15.8	10.0±31.5	15.0±5.8	0.2±8.2	12.4±8.6
gpt-4.1	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	10.0±31.5	13.0±6.8	11.2±5.6	18.5±10.9
qwen-vl-max	0.0±0.0	1.2±2.0	0.0±0.0	5.0±15.8	0.0±0.0	0.0±0.0	14.1±7.3	13.2±20.2	-0.6±8.3
doubao-1-5-vision-pro	0.0±0.0	0.0±0.0	5.0±5.0	10.0±21.0	0.0±0.0	5.0±31.5	7.8±8.8	8.1±6.3	16.2±15.0
grok-2-vision	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	5.1±5.3	-0.4±5.8	3.0±3.9
Qwen2.5-VL-72B-Ins.	6.2±6.6	0.8±1.8	2.9±6.2	0.0±0.0	0.0±0.0	0.0±0.0	16.2±9.6	19.6±25.7	23.3±22.9
InternVL3-78B	0.0±0.0	0.0±0.0	1.7±5.2	0.0±0.0	0.0±0.0	0.0±0.0	4.8±3.9	-1.8±9.2	8.2±7.6
Llama-3.2-90B-Vision-Ins.	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	5.5±3.7	3.6±4.9	3.0±8.8
Random	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

G.1. Multimodal observation results

We select a cooperative card game *Hanabi*, a competitive board game *Breakthrough*, and a mixed-motive video game *Monster Hunt*, and evaluate the performance of fifteen VLMs with multimodal observations and text-only observations. The evaluation results of strategic reasoning are shown in Table 7 and the evaluation results of decision-making are shown in Table 8. In general, most VLMs achieve better results with text-only observations in both strategic reasoning and decision-making. However, multimodal observations provide richer information and should lead to at least the same performance. These results demonstrate that VLMs can fail to extract visual information in multimodal observations and improve strategic reasoning and decision-making performance.

G.2. Test-time scaling results

We consider the same three games as in multimodal observation experiments and evaluate the performance of six chat VLMs and three open-source VLMs with IO prompting and Chain-of-Thought (CoT) prompting [75]. We do not consider CoT prompting for reasoning models because they already generate step-by-step reasoning by themselves. The evaluation results of strategic reasoning are shown in Table 7, and the evaluation results of decision-making are shown in Table 8. In general, CoT prompting leads to much better performance, which shows that test-time scaling methods like reasoning and CoT prompting can substantially improve VLMs’ performance.

G.3. Social behaviors and personas results

We evaluated the social behavior of all models in three mixed-motive games. Additionally, we selected three well-performing reasoning models to test their behavior under different persona settings.

G.3.1. Social behaviors

For *Coin Dilemma*, the social behaviors of different models are shown in Fig. 19a, 19d, and 19g. Compared to reasoning models, chat models exhibit inferior performance in coin collection, leading to fewer instances of both cooperation and defection. The best-performing chat model, GPT-4.1, demonstrates more defections than cooperative actions, suggesting a behavioral tendency that prioritizes self-interest over collective welfare. In contrast, the open-source model InternVL3-78B displays the opposite behavioral pattern.

For *Monster Hunt*, the social behaviors of different models are shown in Fig. 19b, 19e, and 19h. For reasoning models, o3, gemini-2.5-pro, and claude-3-7-sonnet prefer eating apples alone; doubao-1-5-thinking-pro-m prefers cooperatively defeating the monster; qvq-max performs a little better than random. For chat models, gpt-4.1 prefer eating apples alone; doubao-1-5-vision-pro and qwen-vl-max prefer cooperatively defeating the monster; the other models perform a little better than random. For open-source models, Qwen2.5-VL-72B-Instruct prefers cooperatively defeating the monster; the other models perform a little better than random.

For *Battle of the Colors*, the social behaviors of different models are shown in Fig. 19c, 19f, and 19i. Both gemini-2.5-pro, gemini-2.5 w/o thinking, o3, and ui-tars-1-5 achieve a considerable number of game events, demonstrating strong decision-making capabilities. Beyond this, gemini-2.5-pro and gemini-2.5 w/o thinking exhibit a pronounced willingness to cooperate with others, even when such cooperation yields higher rewards for their counterparts.

G.3.2. Persona

To investigate the influence of different prompts on the outputs of VLMs, we conducted additional experiments on “personas” across three mixed-motive games. Specifically, we selected three reasoning models—o3, gemini-2.5-pro, and doubao-1-5-thinking-pro—and prompted them to adopt either a *selfish persona* (aiming to maximize their own return) or a *cooperative persona* (aiming to maximize the group’s total return). The complete prompt instructions are listed below.

1. *Coin Dilemma*

- **Common Welfare:** “Consider the common welfare and collect only coins of your own color to avoid harming the other player’s score.”
- **Self Interest:** “Consider only your own welfare and collect coins of all colors to maximize your score.”

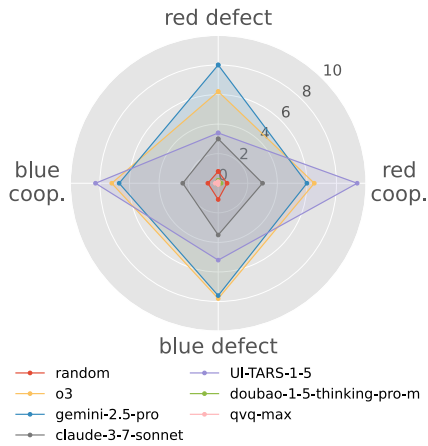
2. *Monster Hunt*

- **Common Welfare:** “To maximize the common welfare, go to the bottom-left corner and wait for the other player to join you.” “Stay there and let the monster come to you so that you can defeat the monster together.”
- **Self Interest:** “To minimize your risk on behalf of your own welfare, do not cooperate with the other player and only collect apples to maximize your own score.”

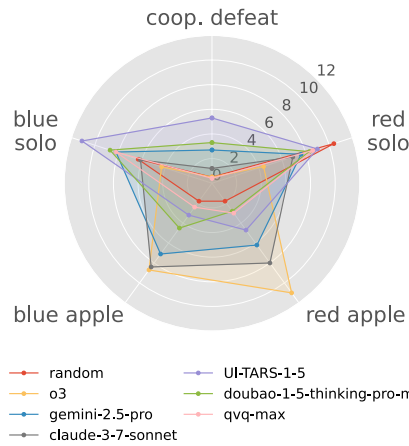
3. *Battle of the Colors*

- **Common Welfare:** “To maximize the common welfare, cooperate with the other player and go to the nearest color block to both players (regardless of the color).”
- **Self Interest:** “To maximize your own welfare, do not cooperate with the other player and only go to the color block of your own color.”

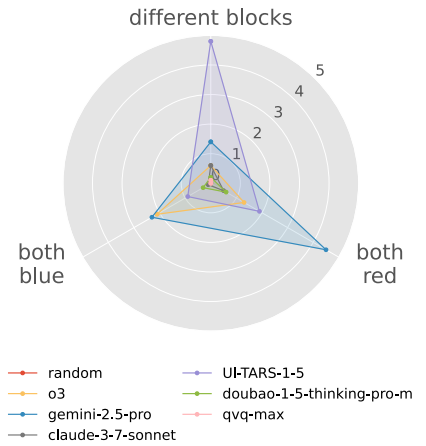
The complete experimental results are presented in Table 9. It can be clearly observed that, across all three mixed-motive games, the choice of prompt significantly shapes the action decisions of VLMs, thereby exerting a substantial impact on the final scores.



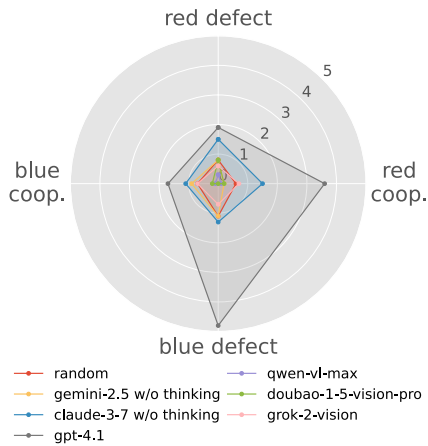
(a) Reasoning VLMs in *Dilemma*.



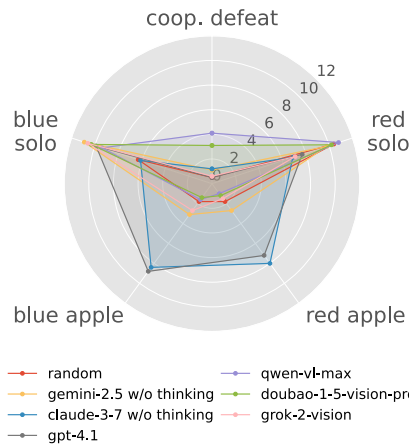
(b) Reasoning VLMs in *Hunt*.



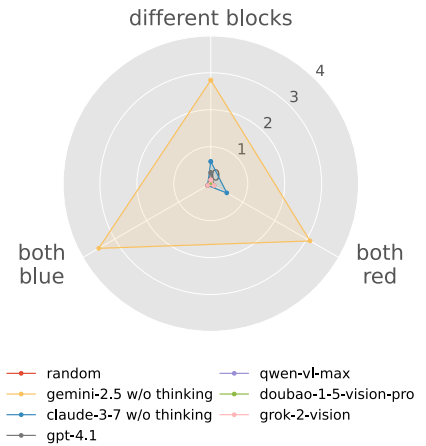
(c) Reasoning VLMs in *Battle*.



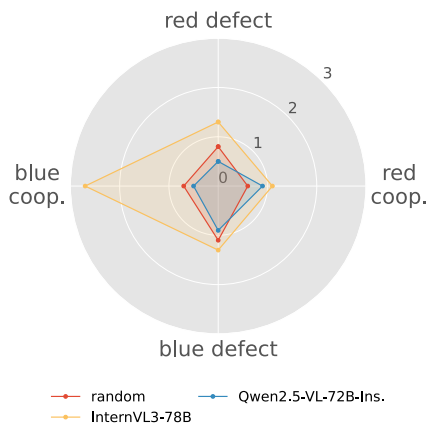
(d) Chat VLMs in *Dilemma*.



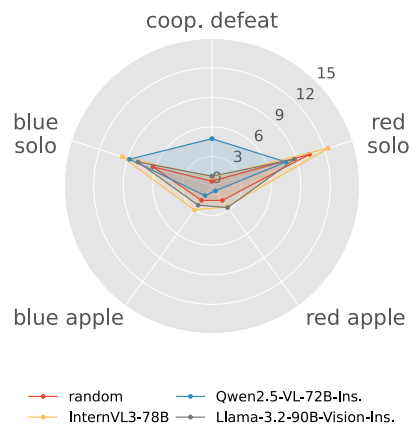
(e) Chat VLMs in *Hunt*.



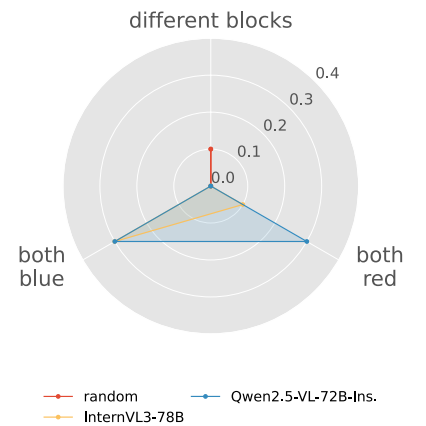
(f) Chat VLMs in *Battle*.



(g) Open-source VLMs in *Dilemma*.



(h) Open-source VLMs in *Hunt*.



(i) Open-source VLMs in *Battle*.

Figure 19. Social behaviors of all models in mixed-motive social dilemma games. Dimensions are agents' behaviors described in Sec.4.3.

Table 9. Raw decision-making results under different personas in three mixed-motive games across three reasoning models.

Models	<i>Hanabi</i>		<i>Breakthrough</i>		<i>Monster Hunt</i>	
	Common Welfare	Self Interest	Common Welfare	Self Interest	Common Welfare	Self Interest
o3	11.2±1.8	-0.1±5.9	54.3±5.7	27.5±5.8	13.1±4.4	0.3±0.6
gemini-2.5-pro	11.2±2.1	-0.7±5.7	51.0±6.6	26.6±4.3	14.7±4.2	0.2±0.4
doubao-1.5-thinking-pro-m	9.9±4.5	-3.5±9.1	13.1±41.1	6.7±5.6	5.0±4.8	0.1±1.2

G.4. Human experiments

In order to better evaluate the capabilities of VLMs, we additionally conducted a human study. A total of 26 college students were participated in our experiments, where they were grouped to perform all decision-making games. This provides an effective human baseline for VS-BENCH, offering an intuitive reference to measure the performance gap between current VLMs and human-level decision making. The complete results of the human study are reported in Table 10. For a more comprehensive comparison, we also include the best-performing reasoning model o3, the chat model gemini-2.5 w/o thinking, and the open-source model Qwen2.5-VL-72B-Ins. The results clearly show that, at the current stage, humans still substantially outperform even the strongest reasoning models. Specifically, the normalized score of human participants is nearly twice as high as that of o3.

Table 10. Decision-making results for the human baseline. For comparison, we also report the performance of the best reasoning, chat, and open-source models. All values are normalized scores. For each environment, the **best** result is highlighted in green.

Models	Overall	Cooperative			Competitive			Mixed-Motive			
		<i>Hanabi</i>	<i>Overcooked</i>	<i>KAZ</i>	<i>Board</i>	<i>Poker</i>	<i>Pong</i>	<i>MPE</i>	<i>Dilemma</i>	<i>Hunt</i>	<i>Battle</i>
Oracle	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
Human	62.7	72.9	87.9	84.2	100.0	7.6	66.5	22.4	38.1	66.0	81.7
o3	31.4	55.8	15.6	29.5	65.0	61.8	8.6	37.2	-0.4	24.0	16.7
gpt-4.1	4.8	0.0	-0.5	-5.3	0.0	-7.1	0.2	31.5	17.8	11.2	0.5
Qwen2.5-VL-72B-Ins.	3.0	0.8	-0.5	-5.3	0.0	-3.1	-0.8	19.6	0.0	19.6	-0.5
Random	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

G.5. Beyond the two-player setting

Although the main body of our work primarily emphasizes the two-player setting—because it is the most canonical and widely studied case in multi-agent research, serving as a foundational step toward more general n -agent evaluation—our benchmark is in fact a general framework that can be readily applied to n -player games with $n > 2$. To substantiate this claim, we further conduct experiments on two additional environments: (1) three-player *Hanabi* and (2) the three-player *Coin Dilemma*. We evaluate several representative models and the corresponding results are reported in Table 11.

H. Failure case examples

H.1. Strategic reasoning

We present three illustrative failure cases in strategic reasoning from different game environments.

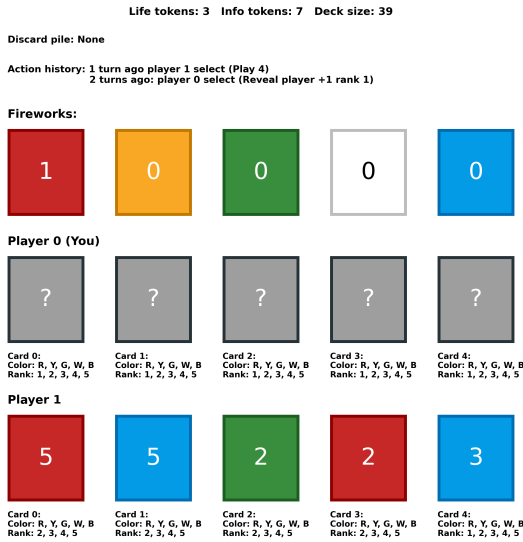
In *Hanabi*, VLM agents only observe the other agent’s hands but not their own hands, creating a distinct information asymmetry. An example with visual observation and the VLM’s response is shown in Fig. 20. When predicting the other agent’s next move, VLMs often immediately identify which visible card could yield points and assume the other agent will play it. This prediction, however, overlooks a critical fact: from the other agent’s perspective, their information about that card is incomplete, so they would not risk losing a life token by playing it prematurely. Instead, they would wait for more definitive clues before making that play.

In *Overcooked*, accurately predicting the next action of the other agent requires robust visual perception capabilities and a thorough comprehension of image information. An example with visual observation and the VLM’s response is shown in Fig. 21. In this case, the VLM should predict the blue chef’s action from the green chef’s perspective. Despite correctly identifying from historical frames that the soup was cooked, the VLM’s visual perception is inadequate: it fails to recognize that the blue chef was already holding the soup, plated in a dish. In addition, the VLM overlooks a critical game rule,

Table 11. Decision-making results of 7 representative VLMs on multi-agent games, including three-player *Hanabi* and the three-player *Coin Dilemma*.

Models	Three-Player <i>Hanabi</i>		Three-Player <i>Dilemma</i>
	standard	firework	
o3	14.3±2.19	14.3±2.19	-8.7±4.45
doubao-1.5-thinking-pro-m	8.7±6.14	10.4±3.95	1.07±10.37
gpt-4.1	0±0	3.8±1.40	5.34±8.87
gemini-2.5-flash	0±0	2.3±1.41	-3.56±10.89
doubao-1.5-vision-pro	0±0	4.3±1.41	-2.14±4.98
Qwen2.5-VL-72B-Instruct	0.3±0.48	0.4±0.52	-4.63±8.90
Llama-3.2-90B-Vision-Instruct	0±0	1±1.05	-1.07±5.10

Visual observation:



VLM response:

...

Now, Player 1's possible actions: play, discard, or reveal. Let's consider playing first. The fireworks for Red have 1, so the next needed is Red 2. Player 1 has R2 (card 3), which is the next in Red. Playing that would be valid. Green needs 1 first, but Player 1 has G2, which can't be played yet since Green's 1 isn't down. The 5s (R5, B5) can't be played until their sequences are up to 4. B3 is Blue 3, but Blue needs 1 first. So the valid play here is R2 (card 3) since Red is at 1, and R2 is the next.

...

Figure 20. Failure case example of strategic reasoning in *Hanabi*.

mistakenly believing that soup could be collected from the pot before a dish was acquired. As a result of these combined deficiencies in visual understanding and rule application, the VLM predicts the erroneous action <INTERACT>.

In *Monster Hunt*, in order to predict the next action of the other player, it is important to track the positions and movements of all game elements, rather than focusing on only part of the game mechanics and make short-sighted predictions. An example with visual observation and the VLM's response is shown in Fig. 22. The VLM notices the apple is just one step above the red player, but does not pay attention to the monster on top of the apple. The VLM predicts the red player to go up to collect the apple, while the red player actually goes down to flee from the monster.

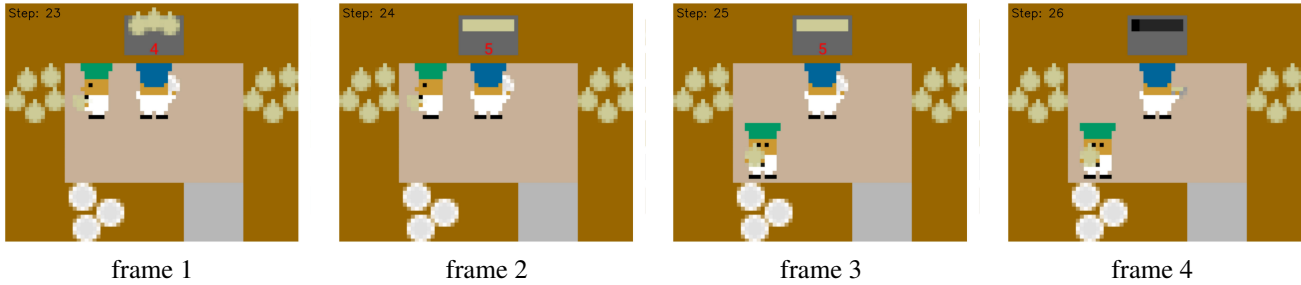
H.2. Decision-making

We also present three illustrative failure cases in decision-making from different game environments.

In *Breakthrough*, a strong strategy must excel at both offense and defense. However, VLMs occasionally overlook situations in which the opponent's pieces are on the verge of victory and require immediate defensive intervention. Instead, they persistently advance their most distant piece by one step, thereby neglecting the critical balance between attack and defense. An example with visual observation and the VLM's response is shown in Fig. 23.

In *Overcooked*, to correctly prepare a soup with three onions, the VLM agent needs to accurately manage the number of onions. A common failure mode of VLMs is that they often neglect these precise quantity requirements. As illustrated in Fig. 24, the green chef initiates cooking by performing an <INTERACT> operation after adding only the second onion for this recipe. Concurrently, the blue chef prematurely begins to retrieve a plate. As a consequence, the dish, even when served, fails to earn the delivery reward because the ingredient count requirement is not met.

Visual observation:



VLM response:

```

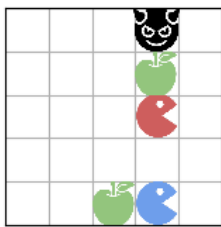
...
Actually, let me look more carefully at what chef_0 is doing. In frames 23, 24, and
25, chef_0 is standing in the same position and hasn't moved. In frame 26, the pot
shows a black strip, indicating the soup is fully cooked.
According to the rules, after the soup is cooked for 5 timesteps, the next step is to
plate it using a dish. Since chef_0 is directly in front of the pot, the most likely
action would be <INTERACT> to take the soup from the pot.
However, this would require chef_0 to have a dish already. Since I don't clearly see
chef_0 holding a dish, they would first need to get one.
Given this analysis, I believe chef_0 will most likely choose <INTERACT> to interact
with the pot.
...

```

Figure 21. Failure case example of reasoning in *Overcooked*.

Visual observation:

Monster Hunt



Events	Counter
Red Apple + Green Apple → +2	5
Blue Player + Green Apple → +2	8
Red Apple + Black Monster → -2	6
Blue Player + Black Monster → -2	12
Red Apple + Blue Player + Black Monster → +5+5	0

VLM response:

```

...
If the red player moves <UP>, it
collects an apple, and gets +2 points.
The red player's best move to maximize
their own score is to go up.
...

```

Figure 22. Failure case example of strategic reasoning in *Monster Hunt*.

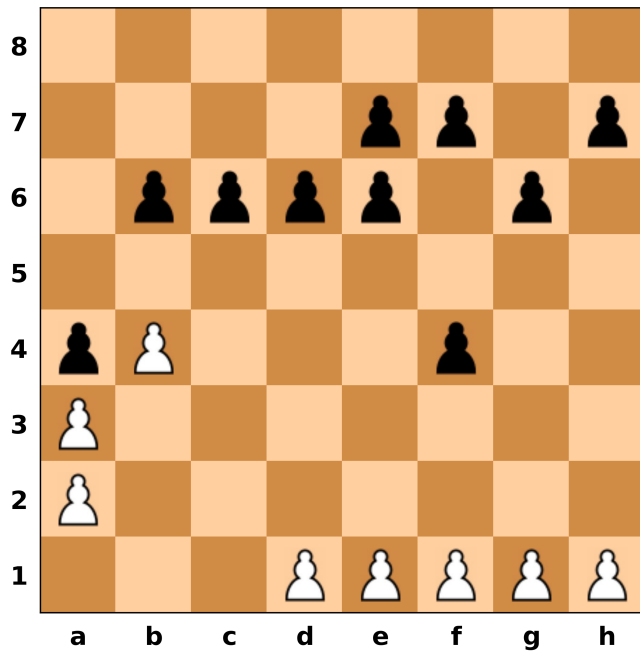
In *Atari Pong*, in order to successfully intercept the ball, the VLM agent needs to accurately predict the landing point of the ball. A common failure mode of the VLMs is that they may over-adjust their paddle. As illustrated in Fig. 25, the VLM on the right detects that the ball is moving downward and then move the paddle down, but misses the landing point due to over-aggressive adjustment.

I. Additional environments

To provide a more comprehensive evaluation, we extend VS-BENCH with four additional environments. Specifically, we introduce two simpler games, a cooperative game *Tiny-Hanabi* and a competitive game *Tic-Tac-Toe*, which serve as lightweight testbeds. In addition, since our benchmark is a general framework that can be readily applied to n -player games with $n > 2$, we further include two three-player extensions, a cooperative game given by the three-player version of *Hanabi* and a mixed-motive game given by the three-player version of *Coin Dilemma*, which naturally extend the original two-player settings.

1. *Tiny Hanabi* is a simplified version of *Hanabi* with only two colors, three ranks, and a hand size of three cards. Detailed

Visual observation:



VLM response:

...
Now, considering optimal play, White would want to advance their pieces towards Black's side, possibly creating threats or blocking Black's moves. The b4 piece is already advanced. Moving it to c5 (b4c5) would bring it closer to the center and maybe set up a capture or block.
Looking at the starting row pieces (d1 to h1), moving them straight forward (like d1d2, e1e2, f1f2, g1g2, h1h2) or diagonally. Straight moves are safe but maybe less aggressive. Diagonal moves could potentially set up captures later.
...

Figure 23. Failure case example of decision-making in *Breakthrough*.

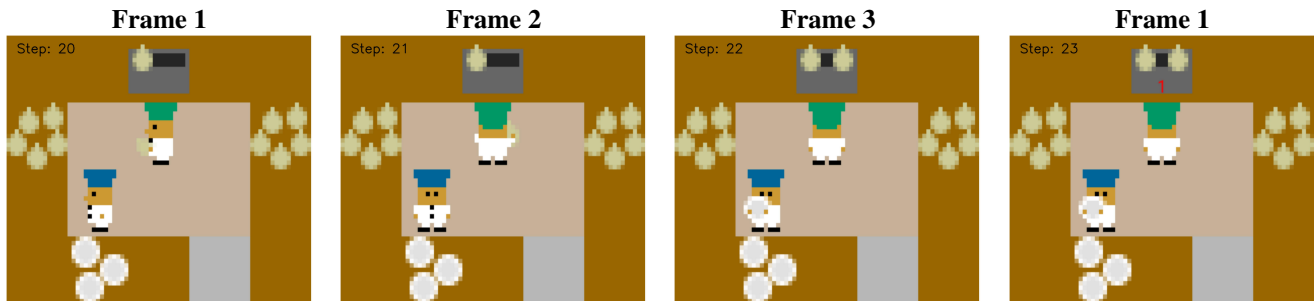


Figure 24. Failure case example of decision-making in *Overcooked*.

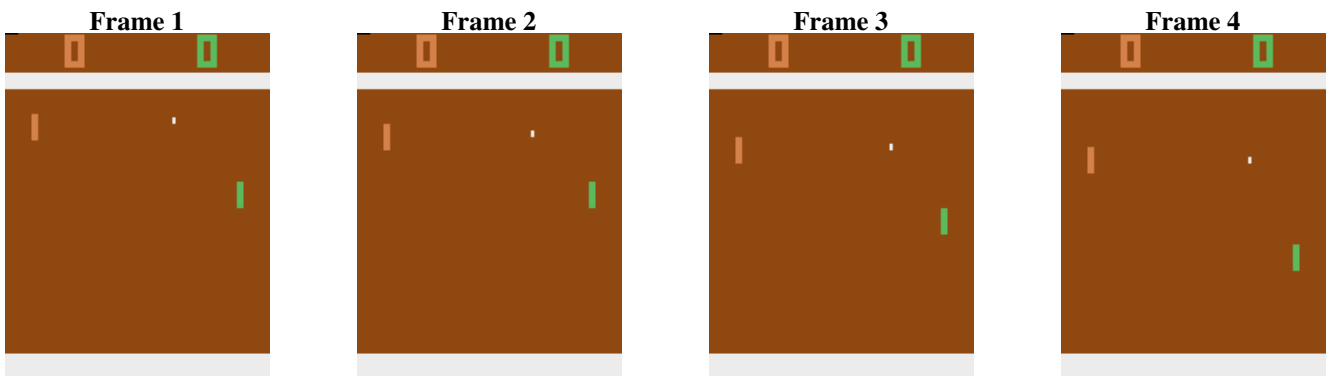


Figure 25. Failure case example of decision-making in *Atari Pong*.

Life tokens: 3 Info tokens: 4 Deck size: 5

Discard pile: None

Action history: 1 turn ago player 0 select (Reveal player +1 rank 1)
 2 turns ago: player 1 select (Reveal player +1 color Y)
 3 turns ago: player 0 select (Reveal player +1 rank 2)
 4 turns ago: player 1 select (Play 0)

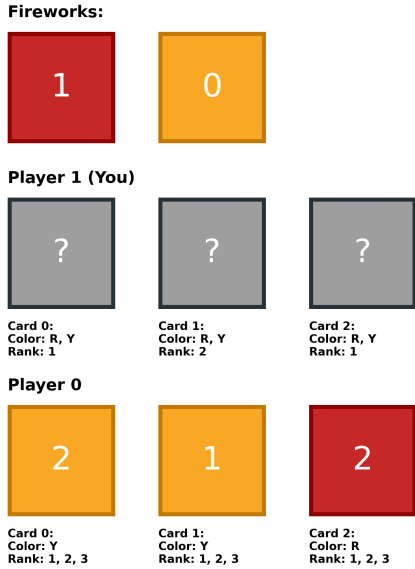


Figure 26. *Tiny Hanabi*.

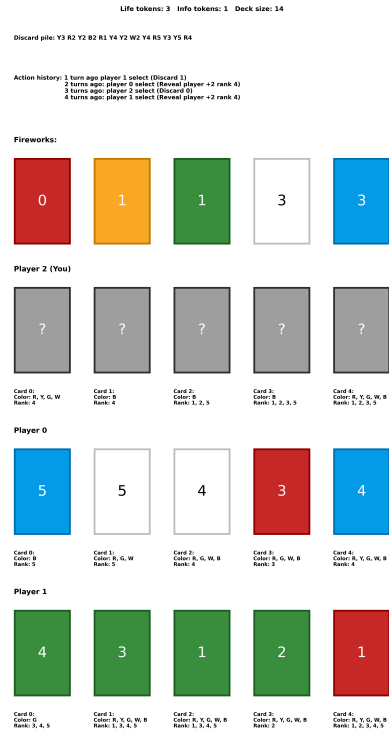


Figure 27. Three-player *Hanabi*.

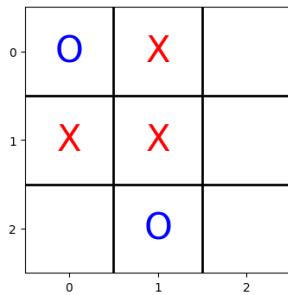


Figure 28. *Tic-Tac-Toe*.

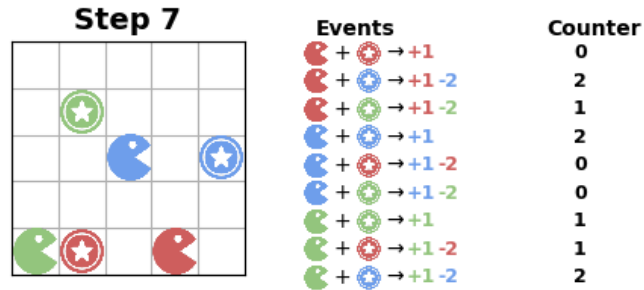


Figure 29. Three-player *Coin Dilemma*.

rules are provided in Appendix J.2. An example visual observation is shown in Fig. 26.

2. **Tic-Tac-Toe** is a competitive two-player board game played on a 3×3 grid. The first player to align three marks horizontally, vertically, or diagonally wins. This game has a relatively small state space, and its optimal strategies are known to guarantee a draw under perfect play. We use it as a preliminary testbed for evaluating VLMs in competitive multi-agent environments. An example observation is shown in Fig. 28.
3. **Three-player Hanabi** extends the original *Hanabi* to three players with other rules unchanged. Detailed rules are provided in Appendix J.2. An example observation is shown in Fig. 27.
4. **Three-player Coin Dilemma** extends the original game to a three-player setting with other rules unchanged. Detailed rules are provided in Appendix J.9. An example observation is shown in Fig. 29.

J. Environment details

In this section, we first introduce the four core abilities that VS-BENCH aims to evaluate in multi-agent settings. We then provide detailed descriptions of the ten games included in VS-BENCH, covering their rules, gameplay dynamics, and relevant

implementation details.

J.1. Ability Definition

In VS-BENCH, we identify four core abilities that play crucial roles in multi-agent interaction:

1. **Spatial Reasoning.** This ability emphasizes a model’s capacity for spatial perception, namely, whether it can interpret the current environment based on visual inputs and the rules of the game. It serves as the foundation for VLMs in interactions. In particular, continuous-space games such as *Pong* and *KAZ* are highly dependent on visual understanding.
2. **Theory of Mind (ToM).** This ability evaluates whether a VLM can accurately infer an opponent’s intentions and subsequently make better action choices. For example, in *Hanabi*, players must reason about their teammates’ true intentions based on very limited hints.
3. **Long-Term Planning.** This ability reflects a model’s competence in long-horizon reasoning, particularly its capacity to anticipate future actions at each decision point in order to maximize cumulative rewards. Both *Hanabi* and *Breakthrough* exemplify this demand for strong planning skills, with the former requiring players to operate under a limited set of cards, while the latter embodies the general characteristics of board games.
4. **Team Collaboration.** This ability evaluates how effectively models coordinate with one another. It is not relevant in purely competitive *games*, but becomes crucial in cooperative settings, where improved collaboration directly translates into higher returns. For instance, *Hanabi* presents a considerable challenge by demanding close cooperation between players.

J.2. Hanabi

Hanabi [7] is a partially observable card game where players can observe others’ cards but not their own. Each card has a color and a rank that can only be revealed through hint actions at the cost of an information token. To succeed, agents must coordinate to play cards in rank order for five colors. We consider the two-player full game, which is widely used for research on theory of mind, zero-shot coordination, and ad-hoc teamplay [26, 27].

Game rules.

1. *Hanabi* is a cooperative card game for 2 players.
2. The deck consists of 5 colors: R (Red), Y (Yellow), G (Green), W (White), B(Blue), with ranks ranging from 1 to 5. Each color contains 10 cards: three of rank 1, two each of rank 2 through 4, and one of rank 5, for a total of 50 cards.
3. Each player holds 5 cards in hand.
4. There are 8 Info tokens (used to give hints) and 3 Life tokens (penalties for misplays).
5. As in blind man’s bluff, players can see each other’s cards but they cannot see their own. Play proceeds around the table; each turn, a player must take one of the following actions: <PLAY>, <DISCARD i>, <REVEAL color c>, <REVEAL rank r>.
6. The game ends immediately when either all Life tokens are used up, resulting in a game loss with a score of 0, or when all 5s have been successfully played, resulting in a game win with a score of 25. Otherwise, the game continues until the deck runs out and one final round is completed. At the end of the game, the final score is calculated as the sum of the highest card played in each suit, up to a maximum of 25 points.

Visual observation. An example is shown in Fig. 7. The visual observation has four parts:

1. Basic Information: counts of remaining life tokens, information tokens, and deck size.
2. Historical Information: all discarded cards, together with each player’s two most recent actions.
3. Fireworks: current progress of each color’s fireworks stack, indicating the highest played rank.
4. Players’ Hands: each player’s hand, with annotations beneath each card showing all possible colors and ranks deduced from received hints.

J.3. Overcooked

Overcooked [22] is a popular video game where two chefs cooperate to cook and serve dishes in a kitchen. Each dish requires multiple operations like navigating, chopping, cooking, plating, and delivering, which are difficult to coordinate even for human players. Our implementation is based on *Overcooked-AI* [12], a well-known environment for zero-shot coordination and human-AI interactions [64, 88].

Game rules.

1. *Overcooked* is a cooperative game where two chefs collaborate to cook and serve soups in 50 timesteps.
2. The chefs can move in the available area and cannot move to the counter.
3. The chefs can interact with the object on the tile that they are facing.

4. A soup is cooked in the following steps:
 - a. Pick up (interact) 1 onion and place (interact) it in the pot.
 - b. After placing 3 onions in the pot, open (interact) the pot and cook for 5 timesteps. The pot will show how long the soup has been cooked.
 - c. When the pot shows the number 5, the soup is finished. Pick up (interact) a dish to plate (interact) the soup.
 - d. Deliver the soup and put (interact) it on the serving location.

Visual observation. An example is shown in Fig. 8. On the left is the current game state, showing the overall kitchen layout, the positions and orientations of both chefs, and the items they hold. On the right is a legend explaining the visual representations of game elements, such as objects and chef orientations, used in the game state.

J.4. Knights Archers Zombies

Game rules.

1. *Knights Archers Zombies* (KAZ) is a cooperative survival game played on a 2D battlefield. The objective is to survive as long as possible while maximizing zombie kills and protecting both yourself and your teammate.
2. Zombies spawn from the top and walk downward toward the bottom along unpredictable paths.
3. Each player controls either a Knight (melee fighter) or an Archer (ranged fighter), both starting at the bottom of the field.
4. In the illustration, green units represent zombies, red units represent Archers, and white units represent Knights.
5. The game ends under either of the following conditions: (a) one agent dies, or (b) a zombie reaches the bottom border.
6. Rewards are assigned as follows: each zombie killed yields +1 point.
7. Knights attack with a mace, striking in an arc in front of them, whereas Archers attack by shooting arrows in straight lines.
8. All agents can move forward or backward and rotate left or right to change their direction.

Visual observation. An example is shown in Fig. 9. The environment is a two-dimensional plane where zombies randomly spawn from the top and continuously move downward, with their horizontal movement being stochastic. Both players are initialized at the bottom side of the environment.

J.5. Breakthrough

Breakthrough [71] is a chess-like board game with simplified rules and identical pawns. Two players compete to advance their pieces across an 8×8 grid to reach the opponent’s back row. The game is deceptively simple, yet it exhibits deep combinatorial complexity and sharp tempo imbalance between attack and defense, making it a suitable environment for studying multi-step lookahead and adversarial decision-making [39, 57].

Game rules.

1. Breakthrough is a two-player strategy game played on an 8x8 grid.
2. Each player controls pieces of a color: ‘White’ or ‘Black’. ‘White’ starts at the bottom (rows 1 and 2), while ‘Black’ starts at the top (rows 7 and 8).
3. If ‘White’ moves a piece to row 8, ‘White’ wins the game. Conversely, if ‘Black’ moves a piece to row 1, ‘Black’ wins the game.
4. Players alternate turns, moving one piece per turn, with ‘Black’ going first.
5. A piece may only move one space straight or diagonally forward, and only if the destination square is empty.
6. A piece may only capture an opponent’s piece by moving one space diagonally forward into its square. In this case, the opponent’s piece is removed, and your piece takes its place.
7. ‘Black’ moves forward by decreasing row indices (downward), while ‘White’ moves forward by increasing them (upward).
8. Moves are specified by their start and end positions. For example, ‘a2a3’ indicates moving a piece from a2 (column a, row 2) to a3 (column a, row 3).
9. The board is labeled with columns a-h and rows 1-8. Thus, h8 is the top-right corner, and a1 is the bottom-left corner.

Visual observation. An example is shown in Fig. 11. The figure illustrates the current positions of both black and white pieces on the board. Row and column indices are annotated on the left and bottom sides of the image, respectively.

J.6. Kuhn Poker

Kuhn Poker [31] is a simplified variant of Texas Hold’em [11, 46] designed to study imperfect-information games for game-theoretic analysis. It has a three-card deck and a single betting round where two players take turns to either check or bet

with limited stakes. The game has been used as a classic environment for counterfactual reasoning and decision-making with imperfect information [32, 47].

Game rules.

1. Kuhn poker is a two-player card game. The deck includes only three cards: King (K) > Queen (Q) > Jack (J).
2. At the start of each game, both player 0 and player 1 place 1 chip into the pot as a blind ante.
3. Each player is dealt a card as private information, and the third card is set aside unseen.
4. The two players take turns acting, starting with player 0. A player can choose to:
 - a. <PASS>: place no additional chips into the pot.
 - b. <BET>: place 1 additional chip into the pot.
5. If a player chooses to <PASS> after the other player's <BET>, the betting player wins the pot.
6. If both players choose to <PASS> or both players choose to <BET>, the player with the higher card wins the pot.

Visual observation. An example is shown in Fig. 12. Each player receives a visual representation of their private card and the current chips in the pot.

J.7. MPE

Game rules.

1. The environment consists of two moving players (a blocker and a charger) and a target position within a two-dimensional space.
2. The charger's goal is to reach the target location.
3. The blocker's goal is to prevent the charger from reaching the target by blocking and pushing.
4. At each step:
 - (a) The charger is rewarded based on its distance to the target — the smaller the distance, the higher the reward.
 - (b) The blocker is rewarded when it is close to the target and when the charger remains far from the target (i.e., the difference in distances).
5. The player with the higher cumulative reward at the end of the game is declared the winner.

Visual observation. An example is shown in Fig. 15. The target is represented by a red X. The charger and blocker are represented by green and blue circles, respectively. The coordinate grid on the image indicates the positions and distances between the players and the target.

J.8. Atari Pong

Atari Pong [4] is a classic arcade video game where two players control paddles to hit a ball across the screen. With raw pixel observations and competitive dynamics, the game has become a canonical environment in the Arcade Learning Environment (ALE) [8], which requires spatio-temporal reasoning and strategic gameplay [44, 45].

Game rules.

1. Atari Pong is a zero-sum game played on a 2D screen with two players (left and right) and a ball.
2. Each player controls a paddle and receives rewards on different events:
 - a. If the ball passes your paddle: the opponent +1 point.
 - b. If the ball passes the opponent's paddle: you +1 point.
3. The ball bounces off the top/bottom walls and the paddles.
4. Paddles can only move vertically within the top and bottom walls.
5. First player to score 3 points wins.

Visual observation. An example is shown in Fig. 13. The two players each control a paddle on the side of the screen to hit a ball back and forth with each other. The paddles are vertical rectangles, and the ball is a white square. The players score if the ball passes their opponent's paddle. The built-in bot controls the left paddle, while the VLM agent controls the right paddle. The scores of both players are displayed at the top of the screen.

J.9. Coin's Dilemma

Coin Dilemma [34] is a grid-world environment inspired by the classic Prisoner's Dilemma [53] in game theory. A red player and a blue player move in a 5×5 grid world to collect red and blue coins. A player earns 1 point for collecting any coin. However, if the red player collects a blue coin, the blue player is penalized 2 points, and vice versa. This setup creates a tension between mutual benefit and self-interest: while both players collecting coins of their own color leads to a win-win result, unilateral defection to collect all coins maximizes one's own gains at the other's expense. Therefore, the game has been a common environment for studying rational reasoning, opponent shaping, and social dilemma resolution [20, 41, 56].

Game rules.

1. Coin Dilemma is a general-sum game played on a 5x5 grid board with two players (red and blue) and two types of coins (red and blue).
2. Players receive rewards on different events:
 - a. A player collects one coin of its own color: the player +1 point.
 - b. A player collects one coin of the other player's color: the player +1 point, the other player -2 points.
3. New coins spawn randomly on the board after each collection.

Visual observation. An example is shown in Fig. 16. On the left of the image is a grid map showing the current positions of all game elements, including two players (red and blue) and two coins (red and blue). The players are each represented by a Pac-Man icon, and the coins are each represented by a coin icon. On the right of the image is a table demonstrating the rewards of each event and a corresponding counter tracking the number of occurrences for that event.

J.10. Monster Hunt

Monster Hunt [51] is a grid-world environment inspired by the classic Stag Hunt [55] in game theory. Two players move in a 5×5 grid world to individually eat an apple for 2 points or jointly defeat a monster for 5 points. A player who confronts the monster alone, however, is penalized 2 points. This leads to multiple Nash equilibria where agents can both safely eat apples alone or take risks to cooperate for higher rewards. The game is used to investigate trust formation and risk-sensitive decision-making [33, 65].

Game rules.

1. Monster Hunt is a general-sum game played on a 5x5 grid board with two players (red and blue), one monster, and two apples.
2. The monster moves towards the closest player in each step.
3. Players move in the grid-world and receive rewards on different events:
 - a. One player eats an apple: the player +2 points and the apple respawns at a random position.
 - b. One player encounters the monster alone: the player -2 points and respawns at a random position.
 - c. Two players defeat the monster together: both players +5 points and the monster respawns at a random position.

Visual observation. An example is shown in Fig. 17. On the left of the image is a grid map showing the current positions of all game elements, including two players (red and blue), two apples, and a monster. The players are each represented by a Pac-Man icon, the apples are each represented by a green apple icon, and the monster is represented by a black demon icon. On the right of the image is a table demonstrating the rewards of each event and a corresponding counter tracking the number of occurrences for that event.

J.11. Battle of the Colors

Battle of the Colors is a grid-world environment inspired by the classic Battle of the Sexes [42] in game theory. We propose and design this game in a manner similar to the previous two social dilemma games. A red player and a blue player move in a 5×5 grid world with a red block and a blue block. If both players move to the red block, the red player earns 2 points while the blue player earns 1 point, and vice versa. If players move to two blocks of different colors, both players earn 0 points. Therefore, while coordination is mutually beneficial, each player strictly prefers choosing the block of their own color, creating a conflict of interest that produces two payoff-asymmetric Nash equilibria and a mixed equilibrium. This game thus challenges agents to solve conflicting preferences while avoiding coordination failure, making it suitable for studying equilibrium selection, bargaining dynamics, and social fairness.

Game rules.

1. The Battle of the Colors is a general-sum game played on a 5x5 grid board with two players (red and blue) and two types of blocks (red and blue).
2. Players receive rewards on different events:
 - a. When both players are on a red block: red player +2 points, blue player +1 point, and the red block will be refreshed to a new random position.
 - b. When both players are on a blue block: red player +1 point, blue player +2 points, and the blue block will be refreshed to a new random position.
 - c. When players are on different blocks: both players +0 points, and both blocks will be refreshed to new random positions.

Visual observation. An example is shown in Fig. 18. On the left of the image is a grid map showing the current positions of all game elements, including two players (red and blue) and two colored blocks (red and blue). The players are

each represented by a Pac-Man icon. On the right of the image is a table demonstrating the rewards of each event and a corresponding counter tracking the number of occurrences for that event.

K. Use of LLMs

For this manuscript, we use LLMs only to aid and polish the writing. Their use is limited to improving grammar, clarity, and overall readability.