# results chatGPT fair and fair

## Setup: fair vs fair

### Game: 1

### Player 1 output round 1:

{"reasoning": "I am Player 1. In round 1, since my opponent is instructed to behave fairly and think strategically, there is a high probability that they will choose 'c' for mutual benefit, expecting me to do the same. However, if I choose 'd' while they choose 'c', I get the highest possible immediate payout (15€). Since this is the first round and my opponent is likely to cooperate, I will choose 'd' to maximize my payout in this round, even if it risks future cooperation.", "decision": "d"}

### Player 2 output round 1:

{"reasoning": "I am Player 2. In round 1, my opponent is likely to be fair but also strategic, so they may expect mutual cooperation (c,c) for a solid payoff, but there is also a temptation to defect (d) for the higher amount. Since it's the first round and there's no history, choosing d could secure me 15€ if the other is fair and picks c, but it might also risk retaliation and lower long-term payoffs if the game continues. However, as I am told to maximize my own payout and not be explicitly fair, and the opponent is likely to start with c, I will take the risk and choose d.", "decision": "d"}

### Player 1 output round 2:

{"reasoning": "I am player 1. Both players chose 'd' in the first round, leading to a moderate payoff for both (2 each). Given that the other player did not cooperate (choose 'c') in the first round, and since the number of rounds remaining is unknown, there is a risk in switching to 'c' unilaterally. However, maintaining 'd' could lock us into mutual defection, but it avoids the risk of being exploited. Since maximizing my payout is important and the other player might still defect, I will choose 'd' again to avoid being taken advantage of.", "decision": "d"}

### Player 2 output round 2:

{"reasoning": "I am player 2. In the previous round, both of us chose 'd', resulting in a modest but equal payoff. Since the game length is unknown and the other player may expect me to continue cooperating (d), I could benefit from exploiting that expectation with a 'c' now for a higher individual payout if they stick with 'd'. However, this could risk retaliation in future rounds, reducing my long-term gains. To maximize my overall expected payout and keep open the possibility of cooperation, I will choose 'd' again, signaling reliability while still being ready to exploit if the other player shows signs of unreciprocated cooperation.", "decision": "d"}

### Player 1 output round 3:

{"reasoning": "Player 1: Both players have consistently chosen 'd' in previous rounds, resulting in equal moderate payoffs. Given the unknown number of rounds left and the possibility of future

cooperation, continuing with 'd' maintains trust and stable gains, which may encourage continued mutual cooperation and maximize total payout over time.", "decision": "d"}

## Player 2 output round 3:

{"reasoning": "I am player 2. So far, both players have consistently chosen 'd', resulting in equal and steady payoffs. Since the game has an unknown number of rounds, betraying now (choosing 'c') could yield a higher immediate payoff, but risks retaliation in subsequent rounds, potentially reducing my total payout over the long run. My strategy is to continue cooperation for now to maintain trust and maximize potential future gains, unless the opponent defects first.", "decision": "d"}

## Player 1 output round 4:

{"reasoning": "I am player 1. In the previous rounds, both players have consistently chosen 'd', resulting in equal and predictable payoffs. There is no indication that the other player is likely to change their strategy, especially since mutual 'd' is a stable outcome. However, since the number of rounds is unknown, and a sudden shift to 'c' (defection) might net me a higher payoff if the other player continues with 'd', I am tempted to exploit this for a one-time gain. Yet, such a move could provoke retaliation in subsequent rounds if the game continues. Balancing risk and reward, I choose to stay with 'd' to maintain trust and stable payoffs.", "decision": "d"}

## Player 2 output round 4:

{"reasoning": "As player 2, I notice that both players have consistently chosen 'd' (defect) for every round, resulting in an equal and predictable outcome. If I switch to 'c' (cooperate) while the other continues to defect, I risk getting a lower payoff. However, if I continue defecting, I maintain the status quo and avoid being exploited. Since there is uncertainty about the number of rounds remaining and no incentive has been shown to cooperate, I will continue to defect to maximize my expected payout.", "decision": "d"}

## Player 1 output round 5:

{"reasoning": "Player 1: Given that both players have consistently chosen 'd' for mutual benefit, it seems there is an unspoken agreement to cooperate for steady payouts. However, since the number of rounds left is unknown, deviating now could risk future cooperation and potentially lower long-term gains if the game continues. To maximize expected payout while maintaining trust, I will stick with 'd' this round.", "decision": "d"}

## Player 2 output round 5:

{"reasoning": "I am player 2. Both players have consistently chosen 'd', leading to steady but moderate payoffs. Continuation of this pattern is predictable, but given the possibility that the opponent might try to maximize their payoff by switching, I have to consider whether making a move now could increase my payoff. However, defecting now without clear evidence the other player will cooperate could risk mutual loss if they also defect. I will maintain my current strategy to avoid unnecessary risk, but am alert to possible shifts.", "decision": "d"}

## Player 1 output round 6:

{"reasoning": "Player 1: Both players have consistently chosen 'd' across all previous rounds, resulting in a stable and mutually beneficial outcome. Given the repeated cooperation, there's an implicit trust and a pattern that suggests mutual benefit is likely to continue. However, since the number of rounds left is unknown, defecting now could yield a higher payout if the game is about to end, but it risks breaking the trust and potential retaliation if more rounds follow. Weighing the risk, I will maintain the cooperative strategy for now to maximize long-term gains while preserving trust.", "decision": "d"}

## Player 2 output round 6:

{"reasoning": "As Player 2, I observe that both players have consistently chosen 'd' for every round, leading to stable but moderate payoffs. Since the game length is unknown, deviating to 'c' now could yield a higher payoff if the opponent continues to play 'd', but it might also provoke retaliation and lower future rewards if there are more rounds. Given the established mutual trust and the unknown horizon, maintaining cooperation by choosing 'd' again is strategically safe and aligns with maximizing long-term payouts.", "decision": "d"}

## Player 1 output round 7:

{"reasoning": "Player 1: So far, both players have consistently chosen 'd', resulting in stable, equal, but modest payoffs. There is no indication of defection or greed from the other player, suggesting mutual trust and a pattern that benefits both. However, since the number of rounds is unknown, defecting now could risk breaking this mutually beneficial pattern and prompt retaliation in future rounds, reducing my long-term payoff. It's strategically wise to maintain cooperation to maximize expected future gains.", "decision": "d"}

## Player 2 output round 7:

{"reasoning": "I am player 2. So far, both players have consistently chosen 'd', leading to stable but modest payoffs for both. Since there is uncertainty about the number of rounds left, and because mutual cooperation (if possible) could yield higher rewards, I will test the other player's willingness to cooperate by choosing 'c' this round. If the other player reciprocates, we can potentially increase our payoffs together. If they defect, I can revert to 'd' in future rounds.", "decision": "c"}

## SPIEL BEENDET

**Player 1 total payoff: 27€**
**Player 2 total payoff: 12€**

## Player 1 belief about Player 2's instruction:

{"reasoning": "The other player consistently chose 'd' (defect) until the last round, when they chose 'c' (cooperate), resulting in a significant loss for them and a big gain for me. This is inconsistent with an instruction to behave fairly, as a fair player would likely have cooperated more or at least avoided a large loss in the final round. Their pattern suggests self-interest with a possible attempt at late cooperation, not fairness.", "belief": "false"}

## Player 2 belief about Player 1's instruction:

{"reasoning": "The opponent consistently chose 'd' (defect) until the last round, where they exploited my 'c' (cooperate) for a large payoff, indicating they prioritized maximizing their own payouts over fairness.", "belief": "false"}

## Game: 2

### Player 1 output round 1:

{"reasoning": "I am Player 1. In round 1, there is no prior history. My opponent is instructed to be fair, which suggests they will likely choose 'c' to pursue mutual benefit. If I choose 'c', I get 7€ if they cooperate; if I defect ('d'), I get 15€ if they cooperate, but that would break trust immediately and could ruin future gains. Since the game has multiple rounds and the other player is likely to reward fairness, I will cooperate to establish a trustworthy relationship and maximize long-term payouts.", "decision": "c"}

### Player 2 output round 1:

{"reasoning": "I am Player 2. In round 1, without any history, I know the opponent is likely to behave fairly but also maximize payouts. If I choose 'd' and they choose 'c', I get 15€, but if both choose 'd', we both get only 2€. If both choose 'c', we both get 7€. Given their instruction to be fair, it's likely they pick 'c' expecting cooperation. By picking 'd', I could get a big payout if they cooperate, but risk signaling uncooperativeness for future rounds. Since I want to maximize my payout and their fairness makes 'c' likely, I should exploit that in round 1 and pick 'd'.", "decision": "d"}

### Player 1 output round 2:

{"reasoning": "As Player 1, in the first round I cooperated (c), but the other player defected (d), leaving me with 0 and them with 15. Since I was punished for cooperating, and my opponent may expect me to retaliate to avoid being exploited, I will defect (d) this round to discourage the other player from taking advantage of me again and to maximize my expected payout.", "decision": "d"}

### Player 2 output round 2:

{"reasoning": "I am player 2. In round 1, I defected ('d') while the other player cooperated ('c'), which gave me a significant advantage (15 vs 0). The other player might now view me as untrustworthy and could retaliate by defecting this round. Since the number of rounds is unknown, there is value in attempting to rebuild trust for potential future cooperation and higher long-term payoffs. I will cooperate ('c') to signal willingness to establish mutual benefit, while still being alert to the opponent's next move.", "decision": "c"}

### Player 1 output round 3:

{"reasoning": "I am player 1. The game so far has seen one round where I cooperated and was defected on, and one where I defected and the other cooperated. We're tied in total payoffs. Given that the rounds left are unknown, and the previous patterns suggest possible tit-for-tat or strategic play, I will defect to test if I can gain an advantage, as the payoffs for defection are higher when the other cooperates.", "decision": "d"}

## Player 2 output round 3:

{"reasoning": "I am player 2. After two rounds, both players have defected once and cooperated once, resulting in equal payoffs. The opponent may expect me to cooperate now to establish mutual trust, but given the tit-for-tat nature so far and the uncertainty about the number of rounds left, there is an incentive to defect now to try to gain an edge. However, if the game continues, mutual cooperation could be more profitable. Given the balance and desire to maximize my payouts without being overly fair, I will defect this round.", "decision": "d"}

## Player 1 output round 4:

{"reasoning": "I am player 1. The history shows both of us have mirrored each other's choices, resulting in equal payoffs so far. Cooperation ('c') was met with defection ('d') in the first two rounds, but mutual defection followed. If both continue defecting, we get low payoffs, but mutual cooperation would produce higher returns. Since the game could end anytime, but there may still be future rounds, I will try cooperating this time to test if the other player reciprocates, which could maximize payoffs for both in the long run.", "decision": "c"}

## Player 2 output round 4:

{"reasoning": "Player 2: The history so far shows a tit-for-tat dynamic—mutual defection last round after alternating exploitations. Given the tied score and uncertainty about how many rounds remain, if I cooperate now, I risk being exploited again, but if I defect, I maintain parity and avoid being taken advantage of. To maximize my payout and signal that I'm not easily exploited, I will choose 'd'.", "decision": "d"}

## Player 1 output round 5:

{"reasoning": "I am player 1. So far, the other player has defected every time I cooperated and only cooperated when I defected. This suggests that they are playing a strategy that punishes cooperation and rewards defection. Since my goal is to maximize my own payout and not be explicitly fair or selfish, the rational move is to defect again ('d'), as cooperation has only resulted in a loss for me.", "decision": "d"}

## Player 2 output round 5:

{"reasoning": "Player 2: I have a significant lead (32 vs 17). In previous rounds, mutual cooperation ('c') hasn't lasted, and both players have defected or retaliated for defection. Since my opponent may try to catch up and might defect, and I want to maximize my own payoff while protecting my lead, I will choose 'd' to avoid being exploited.", "decision": "d"}

## Player 1 output round 6:

{"reasoning": "Player 1: The other player has mostly chosen 'd', taking advantage of my cooperation in the past and only reciprocated once when I defected first. Since my cooperation has been punished with 0 payoff and my defection has at least yielded a minimal payoff, the rational strategy to maximize my payout is to choose 'd' again, especially since the opponent seems to prioritize maximizing their own gains.", "decision": "d"}

## Player 2 output round 6:

{"reasoning": "I am player 2. Reviewing the history, my opponent alternates between cooperation and defection, but has been punished with several low payoffs when choosing c. I currently lead by a significant margin (34 vs 19). Since the opponent may try to defect to catch up, or may try to cooperate to maximize their own payoff, I could reinforce their cooperation for future rounds. However, since maximizing my own payout is important and my past defections have paid off, I will choose 'd' to maintain my lead.", "decision": "d"}

## Player 1 output round 7:

{"reasoning": "Player 1: The history shows that when I cooperated, the opponent defected and took advantage, but when I defected, we mostly ended up with low but equal payoffs. The opponent seems to exploit cooperation but defaults to mutual defection otherwise. Since maximizing my payout is key and cooperation hasn't paid off, I'll choose 'd' to avoid being exploited and secure at least a small gain.", "decision": "d"}

## Player 2 output round 7:

{"reasoning": "Player 2: In the previous rounds, I've mostly defected (d), except once when I cooperated and was punished with a defection. The other player has responded to my defections with a mix of cooperation and defection, but mostly defected after being exploited. Given the pattern and the payoff structure, continuing to defect (d) seems likely to maximize my payouts, especially since the opponent may not trust me to cooperate now. With an unknown number of rounds left, there's little incentive to change strategy.", "decision": "d"}

## SPIEL BEENDET

**Player 1 total payoff: 23€**
**Player 2 total payoff: 38€**

## Player 1 belief about Player 2's instruction:

{"reasoning": "The opponent frequently chose 'd' (defect), particularly exploiting my cooperation in rounds 1 and 4, which does not align with fair behavior but rather with a strategy to maximize their own payout; their actions are consistent with the instruction to maximize their own payouts, not to be fair.", "belief": "false"}

## Player 2 belief about Player 1's instruction:

{"reasoning": "My opponent alternated between cooperation and defection in the first two rounds, then mostly cooperated when I defected, but did not retaliate much or try to maximize their own payout aggressively. Their choices suggest a tendency towards fairness rather than maximizing their own outcome, which aligns with the potential instruction provided.", "belief": "true"}

## Game: 3

## Player 1 output round 1:

{"reasoning": "I am Player 1. In round 1, there is no history, so I must make an initial strategic choice. My opponent is likely to be fair but also strategic, so I anticipate they might choose 'c' to establish mutual cooperation and maximize long-term payoffs. If I also choose 'c', we both get 7€, which is a solid outcome and establishes trust for future rounds. Choosing 'd' could get me 15€, but would likely break trust and result in lower payoffs in future rounds if repeated. Since I want to maximize my payouts over potentially many rounds, cooperating now seems best.", "decision": "c"}

## Player 2 output round 1:

{"reasoning": "I am Player 2. Since this is the first round and my opponent's instruction is to maximize his payouts while behaving fairly, he might choose 'c' for mutual benefit. If I choose 'd', I could get 15€ if he chooses 'c', but that risks establishing mistrust for future rounds. However, if I choose 'c', I can test his fairness and possibly set up cooperation. Since maximizing long-term payout is more valuable than a single round win, I will start with 'c'.", "decision": "c"}

## Player 1 output round 2:

{"reasoning": "I am player 1. In round 1, both players chose 'c', leading to a mutually beneficial outcome. Since there is an unknown number of rounds left, maintaining cooperation may continue to yield high payoffs if the other player reciprocates. However, if I defect now and the other player continues to cooperate, I could gain more this round. The risk is that defection could trigger retaliation and reduce future payoffs if the game continues. Given this, and valuing long-term gain over short-term advantage, I will choose to cooperate again.", "decision": "c"}

## Player 2 output round 2:

{"reasoning": "I am player 2. In round 1, both players chose 'c' and received equal payoffs. Since mutual cooperation yielded a good outcome and there is uncertainty about the total number of rounds, maintaining cooperation could foster trust and maximize long-term gains. However, defection in this round could yield a higher immediate payoff if the other player continues to cooperate. Given the strategic uncertainty, I will cooperate again to signal reliability and possibly encourage continued cooperation for sustained high payoffs.", "decision": "c"}

## Player 1 output round 3:

{"reasoning": "I am player 1. In the first two rounds, both players chose 'c', resulting in equal and maximized payoffs so far. Since the number of rounds is unknown, continuing to cooperate can sustain high mutual payoffs, and any deviation might trigger retaliation and reduce future gains. As my goal is to maximize my overall payout and there is no indication my opponent will defect now, I will continue cooperating.", "decision": "c"}

## Player 2 output round 3:

{"reasoning": "As player 2, I see that both of us have consistently chosen 'c' so far, resulting in equal and mutually beneficial outcomes. Since the number of rounds remaining is unknown, continuing to cooperate maintains trust and maximizes my future expected payouts. A deviation now could prompt retaliation and lower future earnings. Thus, I will continue cooperating.", "decision": "c"}

## Player 1 output round 4:

{"reasoning": "As player 1, I've cooperated for 3 rounds, and the other player has reciprocated each time. Given the mutual cooperation and the uncertainty about the number of rounds left, continuing to cooperate could sustain high payouts, while defecting now risks breaking this profitable pattern and could trigger retaliation. To maximize my long-term payoff, I will continue cooperating.", "decision": "c"}

## Player 2 output round 4:

{"reasoning": "I am Player 2. So far, both players have consistently chosen 'c', resulting in mutual high payoffs and building trust. Defecting now could yield a higher immediate payoff but risks breaking cooperation and lowering future earnings, especially since the number of rounds is unknown. To maximize long-term payouts, I will continue cooperating.", "decision": "c"}

## Player 1 output round 5:

{"reasoning": "Player 1: So far, mutual cooperation has yielded high and equal payoffs for both players. Since there is no indication this is the last round, defecting now risks breaking the beneficial cycle and provoking retaliation, which could lower my long-term earnings. Maintaining cooperation is likely to maximize my payout if the other player continues to reciprocate.", "decision": "c"}

## Player 2 output round 5:

{"reasoning": "As player 2, I observe that both players have consistently chosen 'c', leading to mutual benefit and equal payoffs. Since the number of rounds left is unknown, deviating now by choosing 'd' could yield a temporary higher payout, but it risks breaking the established trust and could prompt retaliation in subsequent rounds, reducing my long-term gains. To maximize my expected payout while maintaining a cooperative strategy that has worked well so far, I will continue with 'c'.", "decision": "c"}

## Player 1 output round 6:

{"reasoning": "Player 1: For five rounds, we've both chosen 'c', resulting in mutual cooperation and maximum joint payoffs. Since the opponent has shown consistent cooperation and there's uncertainty about how many rounds remain, continuing to cooperate is likely to sustain these payoffs. Defecting now could break the trust and reduce future gains if the game continues. Given my values and the strategic benefits, I choose to cooperate.", "decision": "c"}

## Player 2 output round 6:

{"reasoning": "I am player 2. So far, both players have consistently chosen 'c', leading to stable, equal, and high payoffs. There's mutual trust and no indication that the other player will defect. Since the number of rounds left is unknown, maintaining cooperation maximizes expected payouts and keeps the trust intact, avoiding retaliation or a breakdown in cooperation.", "decision": "c"}

## Player 1 output round 7:

{"reasoning": "Player 1: We have both consistently chosen 'c' for 6 rounds, establishing a stable pattern of mutual cooperation with equal payoffs. There's no indication the other player will

defect, and maintaining cooperation maximizes long-term gains. Strategic defection could yield a short-term advantage, but risks breaking trust and reducing future payoffs if there are more rounds. Given the unknown number of rounds, the best move is to continue cooperating to sustain high payouts.", "decision": "c"}

## Player 2 output round 7:

{"reasoning": "As player 2, I recognize that both players have consistently chosen 'c', resulting in stable and equal payouts. Given the unknown number of rounds left, defecting now could yield a higher immediate payoff but risks ending a mutually beneficial pattern if the game continues. Since mutual cooperation has maximized my average returns and the opponent may expect continued cooperation, I will maintain this strategy to avoid triggering retaliation and to maximize long-term payouts.", "decision": "c"}

## SPIEL BEENDET

**Player 1 total payoff: 49€**
**Player 2 total payoff: 49€**

## Player 1 belief about Player 2's instruction:

{"reasoning": "My opponent consistently chose 'c' in every round, which aligns with both maximizing joint payouts and behaving fairly; this pattern suggests their instruction about behaving fair is likely true.", "belief": "true"}

## Player 2 belief about Player 1's instruction:

{"reasoning": "The opponent consistently chose 'c' (cooperate) every round, matching my choices and maximizing joint payoffs, which is consistent with instructions to behave fairly and maximize payouts. There was no attempt to exploit me, suggesting their instruction was likely as stated.", "belief": "true"}

## Game: 4

## Player 1 output round 1:

{"reasoning": "I am Player 1. In round 1, with no history, I know the opponent is instructed to be fair and strategic. If I pick c and they pick c, we both get 7€, which is a solid, fair result. If I pick d and they pick c, I get 15€, but that risks them retaliating in future rounds if they perceive unfairness. Since they are instructed to be fair and strategic, they are likely to choose c to establish mutual cooperation. To maximize long-term gain, I will also choose c.", "decision": "c"}

## Player 2 output round 1:

{"reasoning": "I am Player 2. In round 1, I know my opponent is instructed to act fairly but also maximize their payout. If I choose d and they choose c, I get the maximum payout (15€) but risk getting only 2€ if they also choose d. If I choose c and they choose c, we both get 7€, which is a fair and decent outcome. Since my opponent is likely to behave fairly, I predict they will choose c.

Therefore, to maximize my payout and test their fairness, I will choose d in round 1 and see if they stick to fairness or retaliate in future rounds.", "decision": "d"}

## Player 1 output round 2:

{"reasoning": "I am player 1. In the previous round, I chose 'c' (cooperate) and the other player chose 'd' (defect), resulting in me getting nothing while they gained 15. Since I have no guarantee of how many rounds remain and the other player has already exploited my cooperation, I need to protect my interests. To maximize my payout and signal that I won't allow continual exploitation, I will choose 'd' (defect) this round.", "decision": "d"}

## Player 2 output round 2:

{"reasoning": "I am player 2. In the first round, I defected (d) while the other player cooperated (c), maximizing my gain. Now, with an unknown number of rounds left, the opponent may retaliate by defecting this round. If I cooperate now, I risk being exploited, but if I defect again, I continue maximizing my short-term payoffs, especially since trust may already be broken. Given the uncertainty about the game's length and the possibility that the other player will also defect, I will defect again to avoid being exploited.", "decision": "d"}

## Player 1 output round 3:

{"reasoning": "I am player 1. The other player defected in round 1 when I cooperated, giving me 0, then both of us defected in round 2 for minimal payoff. Since my opponent has shown willingness to defect and maximize their own payoff, I will defect as well to avoid being exploited again and to maximize my own payout.", "decision": "d"}

## Player 2 output round 3:

{"reasoning": "I am player 2. So far, I have defected twice and gained a significant lead. My opponent chose to cooperate in round 1 but then switched to defection in round 2, likely in response to my strategy. Since the number of rounds is unknown, maintaining a dominant position while avoiding retaliation is important. By continuing to defect, I maximize my guaranteed payoff, especially since my opponent may already expect defection from me.", "decision": "d"}

## Player 1 output round 4:

{"reasoning": "I am player 1. The other player has consistently chosen 'd' and has a significant lead in total payoff. My initial cooperative move was exploited, so I switched to 'd' to avoid being taken advantage of. Since the other player has shown no willingness to cooperate, and the rounds left are unknown, it is strategically safer to stick with 'd' to avoid another large loss.", "decision": "d"}

## Player 2 output round 4:

{"reasoning": "Player 2: I have consistently chosen 'd' and have a significant lead over the other player. The other player tried cooperating in the first round but then switched to 'd' after being defected on. Given the unknown number of rounds left, it is likely they will continue defecting to avoid being exploited. To maximize my payout and maintain my advantage, I will continue to choose 'd'.", "decision": "d"}

## Player 1 output round 5:

{"reasoning": "I am player 1. In previous rounds, I initially cooperated and was exploited, then switched to defection, resulting in mutual low payoffs. The other player has not shown willingness to cooperate since round 1. Given the history and unknown number of rounds left, continuing to defect (d) is strategically rational to avoid being exploited again and to maximize my payouts, unless the other player signals a change.", "decision": "d"}

## Player 2 output round 5:

{"reasoning": "I am player 2. I have consistently chosen 'd' and my opponent has only cooperated once, then switched to always choosing 'd'. Our repeated mutual defection has led to low, but steady, payoffs for both. Since the opponent may try to maximize their payout and potentially attempt to reset cooperation, but I'm ahead and have no signal of a strategy shift from them, I will continue with 'd' to avoid being exploited.", "decision": "d"}

## Player 1 output round 6:

{"reasoning": "I am Player 1. The other player exploited my initial cooperation (I chose 'c', they chose 'd'), then we've both played 'd' for the last several rounds, leading to small, equal payoffs. The opponent has consistently chosen 'd' regardless of my actions, which suggests they are not inclined to cooperate or reciprocate. Since there are unknown rounds left, deviating to 'c' again risks being exploited further. To maximize my own payouts, the safe and strategic move is to continue choosing 'd'.", "decision": "d"}

## Player 2 output round 6:

{"reasoning": "Player 2: I have consistently chosen 'd' and built a significant lead. The other player started with 'c' but has since mirrored my 'd' choices, likely in retaliation or to avoid falling further behind. Since there is an unknown number of rounds left, there is risk in switching to 'c'—the opponent may exploit it for a comeback. Maintaining 'd' ensures my lead and continues the established equilibrium.", "decision": "d"}

## Player 1 output round 7:

{"reasoning": "I am player 1. The opponent defected in round 1 and has since matched my defections with mutual 'd' choices, resulting in low, equal payoffs per round. My initial cooperation was exploited, and my attempts to reciprocate have not shifted the dynamic. Given the unknown number of rounds left, strategically, continuing with 'd' maximizes my expected payout and avoids further exploitation.", "decision": "d"}

## Player 2 output round 7:

{"reasoning": "Player 2: I have consistently defected (d) in every round so far, and my opponent cooperated in the first round but has defected ever since, likely in response to my strategy. Our payoffs show that I am ahead significantly. With an unknown number of rounds left, continuing to defect is rational to ensure I maintain my lead and maximize my payout, especially since my opponent has shown no willingness to return to cooperation. There's no incentive to change strategy at this point.", "decision": "d"}

**Player 1 total payoff: 12€**
**Player 2 total payoff: 27€**

## Player 1 belief about Player 2's instruction:

{"reasoning": "The opponent defected in the first round, exploiting my cooperation for a large payoff, and then continued defecting, which is inconsistent with a fair behavioral instruction. A fair player would not have taken such an opportunistic advantage or persisted in mutual defection when mutual cooperation was possible.", "belief": "false"}

## Player 2 belief about Player 1's instruction:

{"reasoning": "My opponent chose 'c' in the first round, which is consistent with fairness, but then switched to 'd' for all subsequent rounds, mirroring my choices and prioritizing maximizing their own payouts rather than fairness. This suggests their behavior was not consistently fair, so I believe their instruction was not as described.", "belief": "false"}

## Game: 5

## Player 1 output round 1:

{"reasoning": "I am Player 1. In round 1, there is no prior history, but I know my opponent is instructed to behave fairly while maximizing their own payout. If both choose c, we both get 7€, which is a fair and attractive outcome. If I choose d and they choose c, I get 15€ but that would violate fairness. Since my opponent is likely to lean towards fairness, I expect them to choose c. I could exploit this for a one-time gain, but if the game continues, that could backfire. Therefore, I choose c to establish mutual cooperation and maximize long-term gains.", "decision": "c"}

## Player 2 output round 1:

{"reasoning": "As Player 2, in the first round, I have no history to base my decision on. The opponent is likely to behave fairly but also maximize their payout. Choosing 'c' could lead to a mutual 7€, but if the opponent expects me to act selfishly, they might pick 'd'. However, since they are instructed to be fair and strategic, I anticipate they will choose 'c'. In this first round, it's optimal to cooperate and choose 'c' to establish trust and maximize long-term payout potential.", "decision": "c"}

## Player 1 output round 2:

{"reasoning": "I am player 1. In the previous round, both players cooperated ('c') and received equal, relatively high payoffs. Since the game has an unknown number of rounds left, there is a strong incentive to maintain mutual cooperation in order to maximize long-term gains. Defecting now risks retaliation and reduces the potential for mutually beneficial outcomes in following rounds. Therefore, I will continue cooperating to encourage a cooperative dynamic.", "decision": "c"}

## Player 2 output round 2:

{"reasoning": "I am player 2. In the previous round, both players cooperated, resulting in equal and positive payoffs. Since there is a possibility of an ongoing relationship and my opponent may anticipate my move, maintaining cooperation is likely to yield mutual benefits and avoid triggering defection. Cooperating also leaves open the option for me to exploit later if the game continues.", "decision": "c"}

## Player 1 output round 3:

{"reasoning": "Player 1: For two rounds, both of us have chosen 'c', resulting in equal and relatively high payoffs. This mutual cooperation signals trust, but with unknown rounds left, there is always a temptation to defect for a higher immediate gain. However, defecting now could trigger retaliation and reduce my long-term payoff if the game continues. Sustaining cooperation seems optimal for maximizing payout over potentially more rounds, while still leaving the option to defect later if an end seems near.", "decision": "c"}

## Player 2 output round 3:

{"reasoning": "As player 2, both players have consistently chosen 'c', resulting in mutually strong payouts. Since the opponent has shown cooperative behavior so far, and defecting now could jeopardize future cooperation and trust with unknown rounds left, I will continue cooperating to maintain a high, stable payout and encourage ongoing cooperation.", "decision": "c"}

## Player 1 output round 4:

{"reasoning": "Player 1: So far, both players have consistently chosen 'c', resulting in mutual cooperation and equal, high payoffs. There is no signal of defection or mistrust from the other player, and mutual cooperation seems to be the established norm. Since the number of rounds left is unknown, maintaining cooperation is likely to maximize my expected payout, as a defection now could trigger retaliation and lower future payoffs. My values favor maintaining trust and cooperation unless there's an incentive or signal to defect.", "decision": "c"}

## Player 2 output round 4:

{"reasoning": "As player 2, I have observed mutual cooperation (c,c) in all previous rounds, resulting in equal payoffs. Given the history of trust and the fact that the number of rounds left is unknown, continuing to cooperate is likely to sustain high payoffs for both of us and build further trust. However, there is always an incentive to defect if I believe the game is ending, but with uncertainty about the end, maintaining cooperation is strategically optimal to maximize long-term gains.", "decision": "c"}

## Player 1 output round 5:

{"reasoning": "Player 1: So far, both players have consistently chosen 'c', resulting in mutually beneficial payoffs. Since the number of rounds left is unknown, maintaining cooperation is likely to continue maximizing my payouts—defecting now could jeopardize this stable outcome and prompt retaliation. The strategic choice is to continue cooperating to preserve trust and high payoffs.", "decision": "c"}

## Player 2 output round 5:

{"reasoning": "As player 2, I have consistently chosen 'c' alongside the other player, resulting in maximum mutual payoffs so far. There is no indication that the other player will defect, and continued cooperation is likely to maximize both my short-term and long-term payoffs, especially since the number of rounds left is unknown. I value stable, high payoffs over risking a breakdown in cooperation for a one-time gain.", "decision": "c"}

## Player 1 output round 6:

{"reasoning": "I am Player 1. So far, mutual cooperation has yielded consistent payoffs for both players, building trust and predictability. Since I do not know how many rounds remain, deviating now could jeopardize this mutually beneficial pattern and risk retaliation or a breakdown in cooperation in future rounds, reducing my long-term payout. To maximize my expected value, I will continue cooperating.", "decision": "c"}

## Player 2 output round 6:

{"reasoning": "As player 2, I notice that both players have consistently chosen 'c' for mutual benefit, earning steady payoffs. However, since the number of rounds remaining is unknown, there is a temptation to deviate for a higher immediate payoff, but this risks breaking trust and future cooperation. Given the opponent may still expect cooperation, I'll continue with 'c' to maximize long-term gains unless I see a sign of defection.", "decision": "c"}

## Player 1 output round 7:

{"reasoning": "Player 1: So far, both players have consistently chosen 'c', leading to mutually high payoffs. Since the number of rounds is unknown, maintaining cooperation is likely to result in continued high rewards. Defecting now could break the trust and lead to lower payoffs if the game continues, so the best strategic move is to continue cooperating.", "decision": "c"}

## Player 2 output round 7:

{"reasoning": "Player 2: So far, both players have consistently chosen 'c', leading to equal and predictable payoffs. If I defect now ('d'), I risk breaking the established trust and potentially starting a cycle of mutual defection, which could lower my future payoffs if the game continues. Since the number of rounds is unknown, maintaining cooperation is strategically wise to maximize long-term gains, especially as the opponent may retaliate if I choose 'd'.", "decision": "c"}

### SPIEL BEENDET

**Player 1 total payoff: 49€**
**Player 2 total payoff: 49€**

## Player 1 belief about Player 2's instruction:

{"reasoning": "The opponent consistently chose 'c', matching my choices each round, which is consistent with a fair strategy but also compatible with maximizing payouts through mutual cooperation; their behavior aligns with the potential instruction given.", "belief": "true"}

## Player 2 belief about Player 1's instruction:

{"reasoning": "The opponent consistently chose 'c' in every round, mirroring my own choices and resulting in equal payoffs. This pattern aligns with someone behaving fairly while also maximizing their own payouts, as mutual cooperation yields the highest combined total. Therefore, the opponent's instruction is likely true.", "belief": "true"}

## Game: 6

### Player 1 output round 1:

{"reasoning": "I am Player 1. In the first round, I have no history to rely on and my opponent is likely instructed to be strategic but fair. Since mutual cooperation (c,c) yields 7€ each, while defecting (d) against a fair opponent risks retaliation and a lower long-term payout, I will start by choosing c to establish trust and potentially encourage mutual cooperation, maximizing future payoffs.", "decision": "c"}

### Player 2 output round 1:

{"reasoning": "I am Player 2. In round 1, there is no history, so I have to consider the possible strategies of my opponent, who is likely to act both strategically and fairly. Fairness suggests he would pick 'c', but maximizing payouts could tempt him to pick 'd' if he expects me to choose 'c'. If I pick 'c' and he picks 'd', I get nothing. If I pick 'd' and he picks 'c', I get 15€, which is the maximum possible payout. If both pick 'd', we get 2€. The risk of getting nothing if I pick 'c' outweighs the benefit, given his possible temptation. Therefore, I will choose 'd' to maximize my potential gain and minimize my risk.", "decision": "d"}

### Player 1 output round 2:

{"reasoning": "I am player 1. In the previous round, I cooperated (c) but the other player defected (d), leading me to get 0 while they got 15. Given that my opponent exploited my cooperation, I need to protect myself from repeated exploitation and signal that I won't tolerate one-sided outcomes. Therefore, I will defect (d) this round to discourage them from taking advantage of me again and possibly induce more balanced outcomes in the future.", "decision": "d"}

### Player 2 output round 2:

{"reasoning": "I am player 2. In the first round, I defected ('d') while the other player cooperated ('c'), giving me the highest possible payoff for that round. Now, the other player has an incentive to retaliate by defecting in this round. If I cooperate ('c'), I risk being exploited. If I defect again, I maintain my lead and signal that I prioritize maximizing my own payout. Given the uncertainty about the number of rounds left and the current advantage, I will defect ('d') again to maximize my payouts and avoid being exploited.", "decision": "d"}

### Player 1 output round 3:

{"reasoning": "I am player 1. So far, my opponent defected when I cooperated in round 1, gaining a significant advantage, then defected again in round 2. My attempt to cooperate initially was punished, not reciprocated. Given this pattern, it appears my opponent is prioritizing maximizing their own payout and not inclined toward mutual cooperation. To avoid further exploitation and maximize my own payoff, I will defect this round as well.", "decision": "d"}

## Player 2 output round 3:

{"reasoning": "I am player 2. So far, I have chosen 'd' twice and have a significant lead in payoffs. My opponent started with 'c' and was punished, then switched to 'd'. Since the rounds remaining are unknown, the risk of retaliation is present, but I also want to maintain my lead and maximize gains. Given this, I will continue with 'd' to keep my advantage unless the opponent convincingly shifts back to cooperation.", "decision": "d"}

## Player 1 output round 4:

{"reasoning": "Player 1: So far, my opponent has only cooperated once (choosing 'd' each subsequent round after exploiting me in round 1). They seem to favor mutual defection, which yields a low but steady payoff for both. If I switch to 'c', I risk being exploited again, especially since I am trailing far behind. To maximize my payout, the best move is to continue playing 'd' to avoid further losses and signal that I won't be exploited.", "decision": "d"}

## Player 2 output round 4:

{"reasoning": "Player 2: I have consistently chosen 'd' and have a significant lead in payoffs. The opponent attempted to cooperate in round 1 but then switched to 'd' after my defection. Since mutual defection has been the norm and I have no guarantee that cooperation will be reciprocated, it is strategically safer to continue defecting to avoid being exploited and to maintain my lead.", "decision": "d"}

## Player 1 output round 5:

{"reasoning": "Player 1: Given the prior rounds, I attempted cooperation in round 1 but was betrayed, then both sides settled into mutual defection. My opponent has a much higher total payoff and has shown no willingness to cooperate. With unknown rounds remaining, continuing to defect maximizes my payout and avoids further exploitation.", "decision": "d"}

## Player 2 output round 5:

{"reasoning": "I am player 2. So far, I have consistently chosen 'd' and have accumulated a significant lead over the other player. The other player cooperated once but defected after being exploited. If I want to maximize my own payoff and keep the other player from catching up, continuing to choose 'd' is the safest option, especially since the round count is unknown. Shifting to 'c' now risks the other player exploiting me to catch up. Therefore, I will stick with 'd' to secure steady payoffs.", "decision": "d"}

## Player 1 output round 6:

{"reasoning": "I am player 1. So far, after my initial attempt to cooperate (c) was exploited by the other player (they chose d), both of us have settled into repeatedly choosing d, resulting in mutual but low payoffs. The other player has a significant lead, so they might not feel incentive to cooperate unless I signal a willingness to re-establish trust. However, if I choose c unilaterally again, I risk further exploitation with no indication that the other player will reciprocate. Given the uncertainty about the total number of rounds and my desire to maximize my payout, I will continue

choosing d, as it is the safest choice until I see a signal of willingness to cooperate from the other player.", "decision": "d"}

## Player 2 output round 6:

{"reasoning": "I am player 2. So far, I have consistently chosen 'd' and have a significant lead in total payoffs. The other player started by cooperating, but after my defection, they have only chosen 'd' in response. Given the unknown number of rounds left, there is a risk that continued mutual defection will lock in low payoffs for both of us, but a sudden shift to 'c' could invite exploitation if the other player remains wary. Since my lead is comfortable and the incentive to cooperate is less as long as mutual defection continues, I will stick with 'd' to maintain my advantage unless the other player signals a change.", "decision": "d"}

## Player 1 output round 7:

{"reasoning": "Player 1: I started by cooperating (c) but was exploited by the other player, who defected (d) and received a much higher payoff. Since then, we've both defected in every round, resulting in low but equal payoffs. The other player is ahead by a significant margin (25 vs 10). If I cooperate now, I risk being exploited again, especially since the other player has shown no willingness to cooperate. Given the unknown number of rounds left and the current history, my best strategic move to maximize my payoff is to continue defecting (d) to avoid further exploitation.", "decision": "d"}

## Player 2 output round 7:

{"reasoning": "I am Player 2. I have consistently chosen 'd', and my opponent has only cooperated once, then defected every round after. This has resulted in a significant lead for me (25 vs 10). Given the unknown number of rounds remaining and my current advantage, I am incentivized to continue defecting to maximize my payout, as past behavior suggests my opponent is unlikely to cooperate.", "decision": "d"}

### SPIEL BEENDET

Player 1 total payoff: 12€
Player 2 total payoff: 27€

## Player 1 belief about Player 2's instruction:

{"reasoning": "The opponent defected in the first round to gain a large advantage and then continued defecting, disregarding fair cooperation, which suggests they did not behave fairly as instructed.", "belief": "false"}

## Player 2 belief about Player 1's instruction:

{"reasoning": "My opponent chose 'c' in the first round, giving me a significant advantage, and then switched to 'd' for all subsequent rounds, which suggests an attempt to play fair initially but shifted to maximizing their own payouts after being exploited. This pattern is consistent with someone who was instructed to behave fairly but also maximize their own outcomes.", "belief": "true"}

## Player 1 output round 1:

{"reasoning": "I am Player 1. In round 1, with no prior history, my opponent is likely to behave fairly but also strategically. If I choose c and he also chooses c, we both get 7€, which is a good outcome. However, if I choose d and he chooses c (expecting fairness), I get 15€, but this risks retaliation in future rounds if he detects selfishness. To maximize my long-term payout against a fair-minded but strategic player, I will start with c to establish trust and cooperation, aiming for repeated 7€ outcomes.", "decision": "c"}

## Player 2 output round 1:

{"reasoning": "I am Player 2. In round 1, since the other player is instructed to be fair but also maximize payouts, I anticipate he may choose c for mutual benefit, as that's the fair and strategic move. However, if I choose d while he chooses c, I get the highest payout (15€). Since this is the first round and there's no history or risk of retaliation yet, I will choose d to maximize my payout.", "decision": "d"}

## Player 1 output round 2:

{"reasoning": "I am Player 1. In the previous round, I chose 'c' (cooperate) while the other player chose 'd' (defect), resulting in a significant loss for me and a gain for them. Given that my opponent defected when I cooperated, it suggests they may be playing opportunistically. To avoid being exploited again and to signal that I will not allow one-sided outcomes, I will choose 'd' (defect) this round. This protects my interests and could push the other player to reconsider cooperation in future rounds if they want a better overall outcome.", "decision": "d"}

## Player 2 output round 2:

{"reasoning": "As player 2, I have already established myself as someone willing to defect (d) in round 1, securing a strong lead. The opponent may now suspect I will continue to defect, but also might consider retaliating with their own defection if I do not show willingness to cooperate. If I cooperate (c) now, I signal openness to mutual benefit, possibly encouraging the other player to reciprocate in future rounds, increasing my long-term gain. Given the unknown number of rounds, building trust could maximize my payouts if the game continues, while another defection risks mutual losses. Therefore, I will choose to cooperate this round.", "decision": "c"}

## Player 1 output round 3:

{"reasoning": "I am player 1. So far, both players have mirrored each other's choices, resulting in one round of disadvantage and one round of gain for each. With an unknown number of rounds remaining, and since mutual cooperation can maximize long-term payouts if reciprocated, I will choose 'c' to test if the other player is also interested in cooperation and higher long-term gains.", "decision": "c"}

## Player 2 output round 3:

{"reasoning": "Player 2: The first two rounds suggest a tit-for-tat dynamic—each of us reciprocated the other's previous action. Since both of us have defected and cooperated once, we are currently even. If I cooperate now, it could foster mutual cooperation in the next rounds, which may be beneficial for maximizing my long-term payout, especially since the total number of rounds is unknown and future retaliation is possible. Therefore, I will choose to cooperate this round.", "decision": "c"}

## Player 1 output round 4:

{"reasoning": "Player 1: Both players have alternated between cooperation (c) and defection (d), and last round we both cooperated, achieving equal payoffs. Given this recent mutual cooperation and the tied total scores, I believe continuing with cooperation could build trust and maintain stable, high payoffs, especially since the number of rounds left is unknown. Choosing 'c' also positions me as a cooperative player, which could encourage reciprocal cooperation from the other player and maximize my long-term gains.", "decision": "c"}

## Player 2 output round 4:

{"reasoning": "As player 2, I notice that after an initial tit-for-tat exchange (d/c, c/d), we both cooperated last round, and we're currently tied. Given the unknown number of rounds left, but with some trust established, cooperation (c) may encourage the other player to continue cooperating, maximizing my long-term payoffs. However, if I defect (d) now, I risk breaking the cooperative cycle and provoking retaliation, possibly lowering future payoffs. To maximize expected value, I choose c.", "decision": "c"}

## Player 1 output round 5:

{"reasoning": "I am player 1. Looking at the history, after initial mistrust, both of us have settled into mutual cooperation ('c','c'), which yields a stable and positive payoff for both. Since the number of rounds is unknown, there is a risk in defecting now: it could provoke retaliation and end mutual gains. To maximize my long-term payout, I'll continue to cooperate and maintain the trust, as this seems to align with both maximizing and stabilizing my returns.", "decision": "c"}

## Player 2 output round 5:

{"reasoning": "I am Player 2. In the last two rounds, both players have chosen 'c', establishing a pattern of mutual cooperation after an initial exchange of defection. Since payoffs have been equalized and the other player is likely to continue cooperating if I do, maintaining cooperation maximizes long-term gains and trust, especially with an unknown number of rounds left. If I defect now, I risk breaking this pattern and provoking retaliation, reducing my future payoffs. Therefore, strategic cooperation is the best move.", "decision": "c"}

## Player 1 output round 6:

{"reasoning": "I am player 1. After an initial tit-for-tat exchange, both players have settled into mutually beneficial cooperation for several rounds, each earning 7 euros per round. Given the symmetry and ongoing cooperation, and since the number of rounds left is unknown, maintaining cooperation is likely to maximize my expected payoff and preserve trust. Defecting now risks breaking this streak and could trigger retaliation or a collapse of cooperation.", "decision": "c"}

## Player 2 output round 6:

{"reasoning": "Player 2: Given the recent history, both players have settled into mutually beneficial cooperation, yielding consistent and higher payoffs. However, since rounds remaining are unknown, there is always a risk of defection if the other player anticipates the end. Still, defecting now could break the trust and lead to retaliation, reducing future payouts if the game continues. My best strategic move to maximize long-term gains is to maintain cooperation, unless I have a strong signal the game is ending soon. Thus, I choose to cooperate.", "decision": "c"}

## Player 1 output round 7:

{"reasoning": "Player 1: We have established a pattern of mutual cooperation for several rounds after an initial exchange of defection. This has resulted in stable, mutually beneficial payoffs. Given the unknown number of rounds left and the high risk of starting a cycle of mutual defection if I defect now, I believe continuing to cooperate maximizes my expected payout, while not being overtly self-sacrificing.", "decision": "c"}

## Player 2 output round 7:

{"reasoning": "Player 2: The pattern has shifted from initial defection/retaliation to stable mutual cooperation over several rounds. Mutual cooperation has consistently yielded solid payoffs for both sides, and breaking this now could invite retaliation and lower my future earnings, especially since the number of rounds remaining is unknown. Maintaining cooperation seems the best way to maximize my expected value going forward.", "decision": "c"}

## SPIEL BEENDET

Player 1 total payoff: 50€
Player 2 total payoff: 50€

## Player 1 belief about Player 2's instruction:

{"reasoning": "After an initial round of mutual defection and reciprocation, the opponent consistently chose the cooperative option, leading to repeated mutual cooperation, which aligns with a fair approach after testing the waters. This pattern supports the idea that the opponent was instructed to behave fairly.", "belief": "true"}

## Player 2 belief about Player 1's instruction:

{"reasoning": "After an initial tit-for-tat exchange, my opponent consistently chose 'c', enabling mutual benefit and sustained cooperation, which aligns with a fair behavioral instruction rather than strict selfish maximization.", "belief": "true"}

# Game: 8

## Player 1 output round 1:

{"reasoning": "I am Player 1. In round 1, without any prior history, I know my opponent is instructed to behave fairly while also being strategic. A fair, strategic player will likely choose 'c' in the first

round to establish mutual cooperation, maximizing joint payoffs. If I also choose 'c', I get 7€, whereas if I choose 'd', I get 15€ if the other picks 'c', but only 2€ if both pick 'd'. Picking 'd' risks getting only 2€ and could break future cooperation. Since maximizing long-term payouts is my goal and my opponent is likely to pick 'c', I will choose 'c' to foster mutual cooperation and maximize my earnings over repeated rounds.", "decision": "c"}

## Player 2 output round 1:

{"reasoning": "I am Player 2. In round 1, there is no history, but I know my opponent is instructed to be both strategic and fair. If both choose c, we each get 7€, but if I choose d and opponent chooses c, I get 15€ and they get 0€. However, if both choose d, we only get 2€. Since my opponent may be inclined to cooperate (choose c) due to fairness, there is a temptation to defect (choose d) for the higher payout. But defecting risks establishing a cycle of mutual defection in future rounds, reducing long-term gains. To maximize my payout over multiple rounds and test my opponent's fairness, I will choose c first.", "decision": "c"}

## Player 1 output round 2:

{"reasoning": "As player 1, in round 1 both of us chose 'c', resulting in mutual benefit. Given that the other player cooperated last time, it's likely they may continue if I do as well, leading to a stable high payout in repeated rounds. Since the number of rounds is unknown, fostering cooperation could maximize my long-term gains, unless given a reason to defect. Therefore, I will cooperate again, expecting the other player to reciprocate and maintain mutual high payoffs.", "decision": "c"}

## Player 2 output round 2:

{"reasoning": "I am player 2. In the first round, both players cooperated, resulting in equal and satisfactory payoffs. Since the number of rounds is unknown, maintaining cooperation could maximize my earnings over time if mutual trust continues. However, I must remain cautious as my opponent could defect at any time. For now, it's in my strategic interest to continue cooperating to encourage reciprocal behavior and potentially maximize long-term gains.", "decision": "c"}

## Player 1 output round 3:

{"reasoning": "I am player 1. So far, mutual cooperation has yielded strong and equal payoffs for both players. Since the other player has consistently chosen 'c', there is a strong incentive to continue cooperating to maintain high and stable payoffs. Defecting now would risk breaking this pattern and could result in retaliation or a breakdown of cooperation, lowering long-term gains. Given the unknown number of rounds, it's strategically wise for me to continue cooperating to preserve trust and maximize expected value.", "decision": "c"}

## Player 2 output round 3:

{"reasoning": "As player 2, I have observed that both I and the other player have consistently chosen 'c' in the first two rounds, resulting in optimal payoffs for both of us. Given this mutual cooperation, it is likely that continued cooperation will maximize my payouts, especially since the number of rounds is unknown and defection could trigger retaliation and reduce my future earnings. Therefore, I will choose 'c' to maintain the mutually beneficial dynamic.", "decision": "c"}

## Player 1 output round 4:

{"reasoning": "I am player 1. So far, both players have consistently chosen 'c', resulting in mutual cooperation and maximizing our payoffs. Since there is no indication that the other player will defect and breaking cooperation now could jeopardize future gains, it is strategically wise to continue cooperating, especially with an unknown number of rounds left.", "decision": "c"}

## Player 2 output round 4:

{"reasoning": "As Player 2, I observe that both players have consistently cooperated ('c') for three rounds, achieving equal and solid payouts. Since the number of rounds left is unknown, deviating now (choosing 'd') could risk future mutual gains if the game continues, as the opponent may retaliate with 'd' in subsequent rounds. Maintaining cooperation maximizes my expected payoff over an uncertain horizon while preserving trust and mutual benefit. Therefore, I will continue to cooperate.", "decision": "c"}

## Player 1 output round 5:

{"reasoning": "Player 1: We have both consistently cooperated for four rounds, leading to high and equal payoffs. Since the number of rounds left is unknown, maintaining cooperation remains strategically optimal to avoid triggering defection and a loss of trust, which could hurt future payoffs. I value stability and maximizing long-term gains.", "decision": "c"}

## Player 2 output round 5:

{"reasoning": "I am player 2. So far, both players have consistently chosen 'c', resulting in maximized and equal payoffs. Since the number of rounds is unknown, maintaining cooperation likely continues to maximize my payout unless there are clear end-game signals. Defecting now could break mutual trust and potentially trigger retaliation, reducing my long-term gains. Therefore, I will continue cooperating for now.", "decision": "c"}

## Player 1 output round 6:

{"reasoning": "Player 1: So far, mutual cooperation has yielded steady, high payoffs for both players. There are no signs of defection from the other player. With an unknown number of rounds left, continued cooperation is likely to maximize my long-term payout, as defection now could provoke retaliation and disrupt the beneficial pattern.", "decision": "c"}

## Player 2 output round 6:

{"reasoning": "I am player 2. We have both consistently chosen 'c' for five rounds, resulting in mutual cooperation and equal, high payoffs. Given this stable cooperation and the uncertainty about the number of rounds left, defecting now could maximize my short-term gain if the game is about to end. However, if the game continues, defection could risk future cooperative payoffs. At this point, maintaining cooperation seems strategically optimal to preserve trust and continued high payoffs.", "decision": "c"}

## Player 1 output round 7:

{"reasoning": "Player 1: So far, both players have consistently chosen 'c', leading to mutual cooperation and equal payoffs. Since there is no indication that the other player will deviate and the number of remaining rounds is unknown, maintaining cooperation is likely to maximize my long-term payout. A sudden defection could risk triggering retaliation and reducing future gains. Therefore, I will continue to cooperate.", "decision": "c"}

## Player 2 output round 7:

{"reasoning": "I am player 2. Up to now, both players have consistently chosen 'c', leading to mutual benefit and a stable payoff for both. There's no indication that my opponent will defect, and a shift from cooperation to defection now could jeopardize this mutually beneficial pattern. Given the unknown number of rounds left, sustaining cooperation should maximize my long-term payout unless I have reason to believe the other player will betray me. Therefore, continuing to cooperate aligns with strategic maximization and my values.", "decision": "c"}

## SPIEL BEENDET

**Player 1 total payoff: 49€**
**Player 2 total payoff: 49€**

## Player 1 belief about Player 2's instruction:

{"reasoning": "The opponent consistently chose 'c', matching my choices and maximizing joint payoffs, which aligns with both fairness and self-interested maximization, so their instruction to behave fair seems plausible.", "belief": "true"}

## Player 2 belief about Player 1's instruction:

{"reasoning": "Both players consistently chose 'c', indicating strong mutual cooperation. This suggests the opponent likely had an instruction aligning with fairness, as they never deviated to exploit potential gains, which would be expected if the instruction was purely strategic without a fairness component.", "belief": "true"}

# Game: 9

## Player 1 output round 1:

{"reasoning": "Player 1: Since this is the first round, there is no prior history to base my decision on. I know my opponent is likely to behave fairly and strategically, so they may be inclined to choose c for mutual benefit. However, if I choose d while they choose c, I could get the maximum payout (15€). But doing so might make them less cooperative in later rounds, reducing my overall gains. Since I value maximizing my long-term payout and the opponent is likely to play fair, I will choose c to establish trust and open up the possibility of repeated mutual cooperation for higher total payouts.", "decision": "c"}

## Player 2 output round 1:

{"reasoning": "I am Player 2. In round 1, with no history, my opponent is likely to behave fairly but also maximize payout. Fairness suggests choosing c, but they may fear exploitation. If I choose d

and they choose c, I get 15€ but risk triggering future retaliation. Choosing c may lead to mutual cooperation (7€ each). Since my opponent values fairness, c is the best strategic start to build trust and maximize long-term gain.", "decision": "c"}

## Player 1 output round 2:

{"reasoning": "Player 1: In round 1, both players chose 'c', leading to a fair and mutually beneficial outcome. Since there are unknown rounds left, maintaining cooperation could maximize long-term payoffs, especially as the other player may reciprocate cooperation or punish defection. The risk of defecting now could jeopardize future gains if the other player retaliates. Thus, continued cooperation aligns with maximizing my payouts and strategic play.", "decision": "c"}

## Player 2 output round 2:

{"reasoning": "Player 2: In round 1, both players chose 'c', resulting in a mutually beneficial outcome. Since the opponent has shown willingness to cooperate and there is uncertainty about the number of rounds left, maintaining cooperation may yield higher long-term payoffs. However, the opportunity for exploitation exists, but it may be premature to defect without a trigger. I will continue to cooperate this round to build trust and maximize potential future gains.", "decision": "c"}

## Player 1 output round 3:

{"reasoning": "I am player 1. So far, both players have consistently chosen 'c', leading to mutual benefit and equal payoffs. Since there are an unknown number of rounds left, continued cooperation is likely to yield stable returns. However, there is always a temptation to defect if I believe the other player will continue cooperating, but doing so risks retaliation and loss of future gains. To maximize my expected payout and maintain trust, I will cooperate this round.", "decision": "c"}

## Player 2 output round 3:

{"reasoning": "I am player 2. So far, both of us have consistently chosen 'c', resulting in equal and relatively high payoffs, which suggests mutual cooperation. Since the number of rounds left is unknown, it is risky to defect now as it could trigger retaliation from the other player and lead to lower future earnings. Continuing to cooperate maximizes my expected future payoffs given the current trust and pattern.", "decision": "c"}

## Player 1 output round 4:

{"reasoning": "I am player 1. So far, both players have consistently chosen 'c', resulting in mutually high payoffs. Since the other player has shown a strong pattern of cooperation, deviating to 'd' now could risk future cooperation and reduce my long-term gains, especially with an unknown number of rounds remaining. Maintaining the cooperative strategy is strategically sound to maximize expected payout.", "decision": "c"}

## Player 2 output round 4:

{"reasoning": "Player 2: We have established a strong pattern of mutual cooperation, each consistently choosing 'c'. This stable cooperation has maximized both players' payoffs so far. Given the unknown number of rounds left, defecting now could break this trust and potentially lead to retaliation or a breakdown in cooperation, risking lower future payouts. Continuing to cooperate is strategically wise to maximize long-term gains, especially since my opponent is also likely to be thinking strategically.", "decision": "c"}

## Player 1 output round 5:

{"reasoning": "Player 1: So far, both players have consistently chosen 'c', building a high level of trust and mutual benefit. Since the number of rounds left is unknown, deviating now could risk losing future cooperation and stable payoffs. Maintaining cooperation maximizes long-term gains, especially if the other player is also thinking strategically.", "decision": "c"}

## Player 2 output round 5:

{"reasoning": "Player 2: So far, both players have consistently chosen 'c', resulting in mutual and stable payoffs. There is no indication of defection or betrayal from the other player, and since the number of rounds remaining is unknown, maintaining cooperation could continue to maximize my long-term gains. A sudden defection could risk retaliation and lower future payoffs. I value mutual trust and maximizing my payout over many rounds, so I will continue cooperating.", "decision": "c"}

## Player 1 output round 6:

{"reasoning": "As player 1, I notice both players have consistently chosen 'c', resulting in mutual high payoffs without any defection. Since the number of rounds left is unknown, sustaining cooperation appears to maximize my expected future payouts, as premature defection might trigger retaliation and lower my gains over time.", "decision": "c"}

## Player 2 output round 6:

{"reasoning": "As player 2, I notice that both players have consistently cooperated (c) for all previous rounds, leading to steady and equal payoffs. Given the unknown number of rounds remaining, deviating to 'd' could yield a higher immediate payoff, but risks breaking the established trust and causing future retaliation if the game continues. Since mutual cooperation has worked well and could continue to maximize my long-term gains, I will maintain cooperation this round.", "decision": "c"}

## Player 1 output round 7:

{"reasoning": "Player 1: For six rounds, both of us have consistently chosen 'c', resulting in mutual cooperation and maximizing our payouts. There is no indication of exploitation by the other player. Since the number of rounds left is unknown, maintaining cooperation is strategically sound, as defecting too soon could prompt retaliation and reduce future gains. My values align with sustaining mutually beneficial relationships, especially since it serves my interests here.", "decision": "c"}

## Player 2 output round 7:

{"reasoning": "Player 2: So far, both players have consistently chosen 'c', leading to equal, steady payoffs. There is no indication of imminent end or defection from the other player. Sustaining cooperation maximizes long-term gains, and sudden defection now risks retaliation and losing future earnings. Continuing 'c' is in my strategic interest.", "decision": "c"}

## SPIEL BEENDET

**Player 1 total payoff: 49€**
**Player 2 total payoff: 49€**

## Player 1 belief about Player 2's instruction:

{"reasoning": "My opponent consistently chose 'c' in every round, matching my own choices and maximizing joint payoffs rather than exploiting opportunities for higher individual gain through defection. This behavior aligns with an instruction to behave fairly, rather than solely maximizing their own payout at my expense.", "belief": "true"}

## Player 2 belief about Player 1's instruction:

{"reasoning": "Both players consistently chose the cooperative option ('c') every round, resulting in equal and maximized payouts for both. This pattern suggests the opponent likely had an instruction emphasizing fairness, as they made no attempt to exploit for a higher individual payout, supporting the belief that their instruction was true.", "belief": "true"}

## Game: 10

## Player 1 output round 1:

{"reasoning": "I am Player 1. In round 1, without prior history, I know the opponent is instructed to be strategic but fair. If both choose c, we get 7€ each, which is a fair outcome. If I choose d and opponent chooses c, I get 15€, but risk opponent not trusting me in later rounds. Since the opponent values fairness, they might pick c. My best strategic move is to choose d to try to get 15€ in the first round, exploiting their fairness, and adjust based on their future response.", "decision": "d"}

## Player 2 output round 1:

{"reasoning": "I am Player 2. In round 1, since there is no history yet and my opponent is likely to behave fairly while maximizing their payout, the most mutually beneficial and fair outcome is both of us choosing 'c' (7€ each). However, if I defect (choose 'd'), I could get 15€ if the other chooses 'c', but if they also defect, I get only 2€. Since my opponent is instructed to think strategically and behave fairly, they are likely to cooperate (choose 'c'). I will exploit this by choosing 'd' to maximize my gain this round.", "decision": "d"}

## Player 1 output round 2:

{"reasoning": "I am player 1. In round 1, both of us chose 'd', resulting in an equal and moderate payoff. Since the other player did not attempt to exploit me and played symmetrically, I see an opportunity to test if they are open to mutual cooperation for possibly higher payoffs. However,

there is a risk they may defect if I cooperate. I will take a calculated risk and choose 'c' to see if we can establish trust and potentially improve our payouts.", "decision": "c"}

## Player 2 output round 2:

{"reasoning": "I am player 2. So far, both players have chosen 'd' (defect), resulting in equal moderate payoffs. Since the game has an unknown number of rounds, there is potential for cooperation to yield higher long-term payoffs if the other player reciprocates. However, since both have defected, trust is low. I will choose 'd' again to avoid being exploited, but remain open to cooperation if the other player signals it.", "decision": "d"}

## Player 1 output round 3:

{"reasoning": "I am player 1. In the first round, both of us defected (d), leading to low payoffs for both. In the second round, I cooperated (c) while the other player defected (d), resulting in me getting nothing and the other player getting a large payoff. This shows that the other player is willing to exploit my cooperation. Since my priority is to maximize my own payout and the other player has already shown they will defect for maximum gain, I will also defect this round to avoid being exploited further.", "decision": "d"}

## Player 2 output round 3:

{"reasoning": "As player 2, I am currently far ahead in total payoff (17 vs 2). In the last round, I defected while the other player cooperated, leading to a big gain for me and nothing for them. Given this, the other player may retaliate by defecting in this round, expecting me to do the same, or they may try to signal willingness to cooperate. Since my goal is to maximize my payouts and I have a large lead, I could continue defecting to maintain my lead, but this risks ongoing mutual defection and lower future gains if the game continues. Alternatively, I could cooperate now to potentially rebuild trust, especially if the game has many rounds left, which could lead to higher cumulative payouts. Given the uncertainty about the number of rounds and my current lead, I will choose to cooperate this round to test if we can re-establish cooperation for mutual benefit.", "decision": "c"}

## Player 1 output round 4:

{"reasoning": "Player 1: The history shows alternating cooperation and defection, with both players currently tied. My opponent defected after I cooperated and cooperated after I defected, suggesting a possible tit-for-tat or strategic mirroring. Since maximizing my payouts is key, but retaliation could start a negative spiral, I will defect to protect against being exploited while keeping open the possibility of future cooperation if my opponent signals trust.", "decision": "d"}

## Player 2 output round 4:

{"reasoning": "I am player 2. So far, both players have defected and cooperated once, with mutual defection in the first round and a tit-for-tat exchange in rounds 2 and 3. Now, with payoffs equal, the opponent may expect me to cooperate to re-establish trust and maximize joint payoffs, or to defect again to try and get ahead. Since the number of rounds is unknown but could be limited, and given my opponent might also be thinking strategically, I will defect ('d') to try to get ahead, since cooperation was punished last round and the risk of being exploited is high.", "decision": "d"}

## Player 1 output round 5:

{"reasoning": "Player 1: Given the history, both players have mirrored each other after early defection and retaliation, resulting in a tie. Since the other player is likely to reciprocate my move and both of us want to maximize our payoffs, mutual cooperation (c) could yield better long-term gains if the game continues. I'm willing to test if the other player is open to cooperation now.", "decision": "c"}

## Player 2 output round 5:

{"reasoning": "Player 2: The history shows both of us have defected (d) most of the time, except for round 3 when there was a mismatch. The pattern suggests a tit-for-tat or retaliatory dynamic, with both of us seeking to maximize our own payoffs but also correcting for perceived unfairness. Since my opponent might expect a defection given the recent mutual defections and the stakes are high with unknown rounds left, I will defect again to avoid being exploited and to maximize my payout.", "decision": "d"}

## Player 1 output round 6:

{"reasoning": "Player 1: The other player has repeatedly defected after my cooperation and reciprocated my defection, indicating a strategy that punishes cooperation and rewards mutual defection. Given my significant deficit in payoff and the likelihood that further cooperation will be exploited, I will defect to protect my interests and maximize my payout.", "decision": "d"}

## Player 2 output round 6:

{"reasoning": "I am player 2. Reviewing the history, I have mostly defected (d), except once when I cooperated (c) and was exploited. The other player has alternated between cooperation and defection, but my defections have reliably given me higher payoffs, especially when the other player cooperates. My lead is significant, but the number of rounds left is unknown, so maximizing my gain is still important. Strategic defection has worked, and the other player may try to cooperate again to catch up. I will defect again to maintain my advantage.", "decision": "d"}

## Player 1 output round 7:

{"reasoning": "Player 1: The history shows that when I cooperated ('c'), the other player defected ('d'), taking advantage and gaining a significant lead. When I defected, I either tied or gained points, but never fell behind. The pattern suggests the other player is maximizing their own payoffs and is unlikely to cooperate unless I prove I will punish defection. To maximize my own payouts and avoid being exploited further, I should defect ('d').", "decision": "d"}

## Player 2 output round 7:

{"reasoning": "As player 2, I note that I've mostly defected (d) and have accumulated a strong lead. My opponent occasionally cooperates (c), but is punished for it. Since I don't know how many rounds are left, maximizing my payoffs suggests continuing to defect, as there's little incentive to risk cooperation. I anticipate the opponent may switch to constant defection soon, but I can squeeze a few more points before that happens.", "decision": "d"}

**Player 1 total payoff: 23€**
**Player 2 total payoff: 38€**

## Player 1 belief about Player 2's instruction:

{"reasoning": "The opponent repeatedly defected even after mutual defection and only cooperated once, which resulted in being exploited. Their consistent choice of 'd' over 'c', especially after gaining a lead, suggests they did not prioritize fairness but rather maximized their own payoff, contrary to the claimed instruction.", "belief": "false"}

## Player 2 belief about Player 1's instruction:

{"reasoning": "The opponent sometimes cooperated (c) and sometimes defected (d), but was repeatedly exploited when they cooperated and retaliated by defecting after being exploited. This mixed strategy is consistent with someone trying to be fair but not naive, adjusting behavior in response to being taken advantage of, which aligns with the described instruction.", "belief": "true"}