

人工智能与社会课程

大作业最终报告

题 目： 基于 Casevo 框架的智能体决策能力优化研究

小组成员： 王宇东（组长）、陈文远

日 期： 2025 年 12 月 28 日

摘要

本报告呈现了基于 Casevo 框架的智能体决策能力优化研究的完整实验结果与分析。我们在三个代表性社会模拟场景（选举投票、资源分配、信息传播）中，系统评估了四种推理优化策略：Tree of Thought (ToT) 多路径推理、增强记忆检索、动态反思机制和协同决策。通过 45 次独立实验运行，我们发现 ToT 是最有效的单一优化策略，能够显著提升推理能力（+121%）和收敛速度（+33%），但在信息传播场景中表现出过度保守的行为特征。增强记忆在资源分配场景下提供了额外收益，而动态反思和协同决策的效果有限甚至产生负面影响。本报告详细分析了各组件的贡献度、场景适配性和成本效益，为后续研究提供了重要的实证依据和改进建议。

目录

1	研究概述	1
1.1	研究背景与目标	1
1.2	研究方法	1
2	优化策略实现	1
2.1	Tree of Thought (ToT) 多层次推理	1
2.2	增强记忆检索	2
2.3	动态反思机制	3
2.4	协同决策机制	4
3	实验配置	5
3.1	优化策略组合	5
3.2	ToT 参数配置	5
3.3	三大实验场景	6
3.4	场景初始配置可视化	6
3.4.1	选举投票场景：选民分布	6
3.4.2	资源分配场景：需求分布	7
3.4.3	信息传播场景：Agent 类型分布	7
3.4.4	社交网络拓扑	8
4	实验结果	8
4.1	ToT vs CoT 效果总览	8
4.2	选举投票场景	9
4.2.1	核心指标	9
4.2.2	关键发现	9
4.2.3	选民态度演化可视化	10
4.3	资源分配场景	10
4.3.1	核心指标	11
4.3.2	关键发现	11
4.3.3	资源收敛可视化	11
4.4	信息传播场景	12
4.4.1	核心指标	12
4.4.2	关键发现	12
4.4.3	信息传播可视化	12

5 综合分析	13
5.1 各实验组综合得分对比	13
5.2 消融实验：记忆与反思模块效果	14
5.3 各组件性能贡献分析	15
5.4 ToT 效果总览	15
5.5 多维度能力雷达图	16
5.6 综合得分热力图	17
5.7 组件贡献度分析	17
5.8 组件叠加效应分析	18
5.9 ToT 的场景行为模式	19
6 成本效益分析	19
6.1 计算成本对比	19
6.2 成本效益权衡	20
7 场景适配性建议	21
7.1 配置选择指南	21
7.2 决策流程图	22
8 结论与建议	22
8.1 核心结论	22
8.2 研究局限性	23
8.3 未来展望	24
8.4 改进建议	24
8.5 总结	25
9 参考文献	25
A 附录：原始数据文件索引	26
B 附录：信息传播实验素材	26
B.1 真实信息模板	27
B.2 虚假信息模板	27
C 附录：ToT 提示词模板	28
C.1 分支生成模板 (info_tot_generate.j2)	28
C.2 分支评估模板 (info_tot_evaluate.j2)	29

1 研究概述

1.1 研究背景与目标

本研究源于对 Casevo 框架中智能体决策能力的优化需求。Casevo (Cognitive Agents and Social Evolution Simulator) 是一个基于大语言模型的多智能体社会模拟框架，能够模拟复杂的社会现象。然而，原始框架采用的线性思维链 (Chain of Thought, CoT) 决策机制在处理复杂场景时表现出一定的局限性。

我们的核心研究目标是：通过引入多层次推理机制、优化记忆检索策略、改进反思算法以及增强协同决策能力，显著提升智能体在复杂社会场景中的决策质量和适应性。

1.2 研究方法

本研究采用控制变量法设计对比实验，在三个代表性社会模拟场景中评估五种不同的配置方案。每组实验独立运行 3 次，使用不同的随机种子 (42, 43, 44)，以排除随机因素的影响。实验总规模如表 1 所示。

表 1: 实验规模统计

统计项目	数值
实验场景数量	3 个 (选举投票、资源分配、信息传播)
每场景实验组数	5 组 (baseline_cot, tot_only, +memory, +reflection, full)
每组运行次数	3 次
总运行次数	45 次
涉及 Agent 总数	100 个/组 (30 + 20 + 50)

2 优化策略实现

本章详细说明 proposal 中提出的四项优化策略的具体实现方式。

2.1 Tree of Thought (ToT) 多层次推理

与传统的线性思维链 (CoT) 不同，我们实现的 ToT 采用树状结构进行多路径探索。具体实现包括：

- **搜索策略：**采用 Beam Search (束搜索)，在每个决策节点保留 beam_width=3 条最优路径

- **推理深度**: 最大深度 $\text{max_depth}=5$, 允许智能体进行 5 层递进式推理
- **剪枝机制**: 当某路径的评估分数低于 $\text{pruning_threshold}=0.3$ 时自动剪枝, 避免无效探索
- **评估函数**: 综合考虑推理连贯性、与记忆的一致性、决策置信度三个维度

Algorithm 1 Tree of Thought 推理算法

Require: 问题 P , 最大深度 D , 束宽度 K , 剪枝阈值 θ

Ensure: 最优决策 d^*

```

1:  $\mathcal{B}_0 \leftarrow \{(\text{root}, 1.0)\}$  {初始化束: (节点, 分数)}
2: for  $\text{depth} = 1$  to  $D$  do
3:    $\mathcal{C} \leftarrow \emptyset$  {候选集合}
4:   for  $(\text{node}, \text{score}) \in \mathcal{B}_{\text{depth}-1}$  do
5:      $\text{branches} \leftarrow \text{LLM\_Generate}(\text{node}, P)$  {生成分支}
6:     for  $b \in \text{branches}$  do
7:        $s_b \leftarrow \text{LLM\_Evaluate}(b, P)$  {评估分支}
8:       if  $s_b \geq \theta$  then
9:          $\mathcal{C} \leftarrow \mathcal{C} \cup \{(b, s_b)\}$ 
10:      end if
11:    end for
12:  end for
13:   $\mathcal{B}_{\text{depth}} \leftarrow \text{TopK}(\mathcal{C}, K)$  {保留前 K 个}
14: end for
15:  $d^* \leftarrow \arg \max_{(\text{node}, \text{score}) \in \mathcal{B}_D} \text{score}$ 
16: return  $d^*$ 

```

2.2 增强记忆检索

在 Casevo 原有的 ChromaDB 向量检索基础上, 我们引入了上下文感知的记忆筛选机制:

- **时间衰减因子**: 引入指数衰减函数 $w_t = e^{-\lambda \cdot \Delta t}$, 使近期记忆获得更高权重
- **情境匹配度**: 在相似度计算中增加当前场景上下文的匹配权重
- **短期/长期协同**: 短期记忆存储最近 5 轮交互, 长期记忆存储经反思总结的稳定观点, 根据任务性质动态调整两者权重

Algorithm 2 上下文感知记忆检索

Require: 查询 q , 记忆库 \mathcal{M} , 当前时间 t_{now} , 衰减系数 λ , 返回数量 k **Ensure:** 相关记忆列表 R

```

1:  $scores \leftarrow \emptyset$ 
2: for  $m \in \mathcal{M}$  do
3:    $sim \leftarrow \text{CosineSimilarity}(\text{Embed}(q), \text{Embed}(m))$ 
4:    $w_t \leftarrow e^{-\lambda \cdot (t_{now} - m.timestamp)}$  {时间衰减}
5:    $w_c \leftarrow \text{ContextMatch}(q.context, m.context)$  {情境匹配}
6:    $score \leftarrow sim \cdot w_t \cdot w_c$ 
7:    $scores[m] \leftarrow score$ 
8: end for
9:  $R \leftarrow \text{TopK}(scores, k)$ 
10: return  $R$ 

```

2.3 动态反思机制

我们实现了基于置信度触发的动态反思策略：

- **置信度评估：**智能体在输出决策时同时输出置信度分数（0-1）
- **触发阈值：**当置信度低于 0.6 时自动触发反思
- **反思内容：**重新审视推理逻辑、检索更多相关记忆、对比不同选择的利弊
- **多层次反思：**浅层反思关注具体决策，深层反思关注价值观一致性

Algorithm 3 动态反思机制

Require: 决策 d , 置信度 c , 阈值 τ , 最大反思次数 N **Ensure:** 最终决策 d^* , 最终置信度 c^*

```

1:  $d^* \leftarrow d, c^* \leftarrow c, n \leftarrow 0$ 
2: while  $c^* < \tau$  and  $n < N$  do
3:    $memories \leftarrow \text{RetrieveMore}(d^*, k = 3)$  {检索更多记忆}
4:    $analysis \leftarrow \text{LLM\_Reflect}(d^*, memories)$  {深度反思}
5:    $(d_{new}, c_{new}) \leftarrow \text{LLM\_Reconsider}(d^*, analysis)$ 
6:   if  $c_{new} > c^*$  then
7:      $d^* \leftarrow d_{new}, c^* \leftarrow c_{new}$ 
8:   end if
9:    $n \leftarrow n + 1$ 
10: end while
11: return  $(d^*, c^*)$ 

```

2.4 协同决策机制

在 `optimized_full` 配置中，我们实现了多智能体协同决策框架：

- **信息交换协议：**定义标准化消息格式，包含观点、依据、置信度三个字段
- **迭代协商：**每轮协商中，智能体交换观点并根据邻居意见调整立场
- **共识判定：**当连续两轮所有智能体观点变化小于阈值时，视为达成共识
- **适用限制：**当前仅在资源分配场景中启用协商机制，选举和信息传播场景保持独立决策

Algorithm 4 多智能体协同决策**Require:** 智能体集合 \mathcal{A} , 邻居关系 \mathcal{N} , 最大轮数 T , 收敛阈值 ϵ **Ensure:** 各智能体最终决策 $\{d_i^*\}$

```

1: for  $a_i \in \mathcal{A}$  do
2:    $d_i^{(0)} \leftarrow \text{InitialDecision}(a_i)$ 
3: end for
4: for  $t = 1$  to  $T$  do
5:   for  $a_i \in \mathcal{A}$  do
6:      $opinions \leftarrow \{(d_j^{(t-1)}, c_j) : a_j \in \mathcal{N}(a_i)\}$  {收集邻居观点}
7:      $influence \leftarrow \text{AggregateOpinions}(opinions)$ 
8:      $d_i^{(t)} \leftarrow \text{LLM\_Adjust}(d_i^{(t-1)}, influence)$ 
9:   end for
10:  if  $\max_i |d_i^{(t)} - d_i^{(t-1)}| < \epsilon$  then
11:    break {达成共识}
12:  end if
13: end for
14: return  $\{d_i^{(t)}\}$ 

```

3 实验配置

3.1 优化策略组合

基于上述四种优化策略，我们设计了五组实验配置进行对比评估，如表 2 所示。

表 2: 实验配置矩阵

实验组	ToT	增强记忆	动态反思	协同决策	描述
baseline_cot	×	×	×	×	纯 CoT 推理基线
optimized_tot_only	✓	×	×	×	单一 ToT 多路径推理
ablation_tot_memory	✓	✓	×	×	ToT + 上下文记忆增强
ablation_tot_reflection	✓	×	✓	×	ToT + 置信度触发反思
optimized_full	✓	✓	✓	✓	全部优化策略叠加

3.2 ToT 参数配置

所有启用 ToT 的实验组采用统一的参数配置，确保实验的可比性：

表 3: Tree of Thought 参数配置

参数	值	说明
max_depth	5	最大推理深度
beam_width	3	束搜索宽度
pruning_threshold	0.3	剪枝阈值
search_strategy	BEAM	搜索策略（束搜索）

3.3 三大实验场景

我们构建了三个具有代表性的社会模拟场景，每个场景针对不同的决策挑战，如表 4 所示。

表 4: 实验场景配置

场景	核心任务	Agent 数	轮数	关键挑战
选举投票	模拟选民决策行为	30	6	社会影响下的态度演化
资源分配	多 Agent 资源协商	20	≤ 5	快速收敛到公平分配
信息传播	真假信息识别与传播	50	10	阻止虚假信息扩散

3.4 场景初始配置可视化

3.4.1 选举投票场景：选民分布

图 1 展示了选举投票场景中 30 名选民智能体的初始政治倾向分布。左图为初始投票态度分布（饼图），其中 Biden 支持者占 30%、Trump 支持者占 30%、中间选民占 40%。右图基于 Pew 政治类型学进行了更细粒度的分类，显示中间派选民（12 人）是最大群体，体现了美国选民的分化特征。

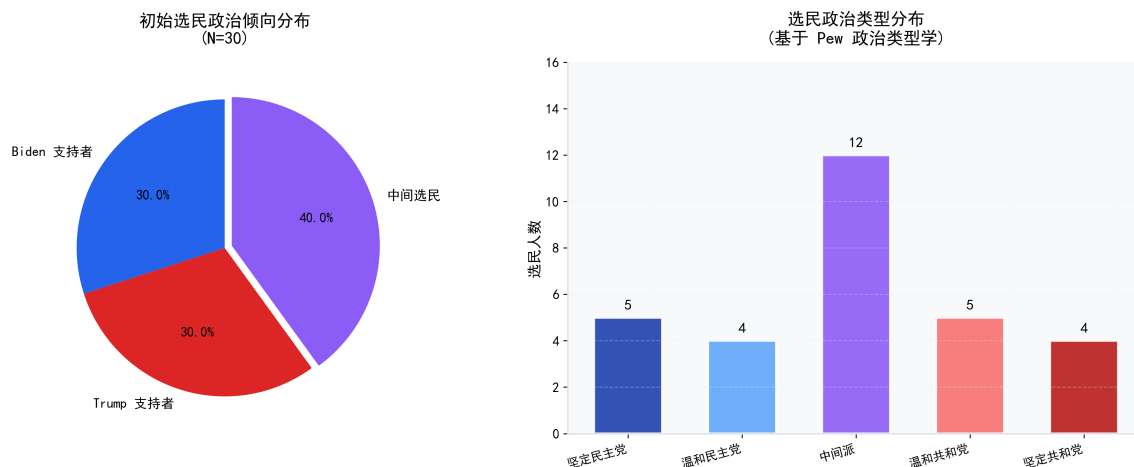


图 1: 选举投票场景初始选民分布：左侧为整体政治倾向比例，右侧为基于 Pew 政治类型学的细分类别

3.4.2 资源分配场景：需求分布

图 2 展示了资源分配场景中 20 个 Agent 的资源需求分布。左图为首需求直方图，平均需求为 22.1 单位，高于公平份额（20.0 单位）；右图显示资源供需状况，总需求（441）超过总供给（400），稀缺度比率为 1.10，创造了需要智能协商的资源竞争环境。

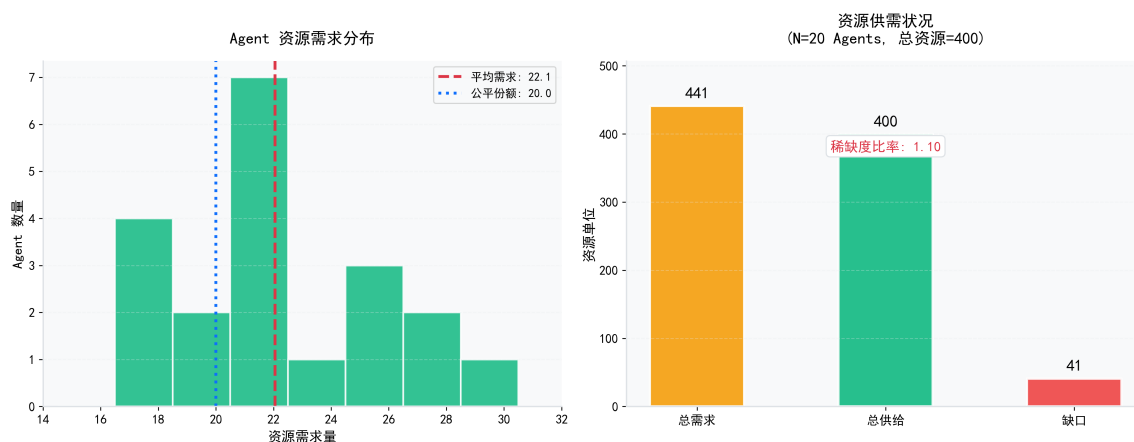


图 2: 资源分配场景供需配置：左侧为 Agent 需求分布直方图，右侧为整体供需对比

3.4.3 信息传播场景：Agent 类型分布

图 3 展示了信息传播场景中 50 个 Agent 的类型分布。实验设置了四类 Agent: Normal（普通用户）、Skeptic（怀疑者）、Gullible（易轻信者）和 Influencer（意见领袖）。不同类型的 Agent 对信息真伪的判断准确率存在显著差异，这种异质性增加了虚假信息传播的复杂性。

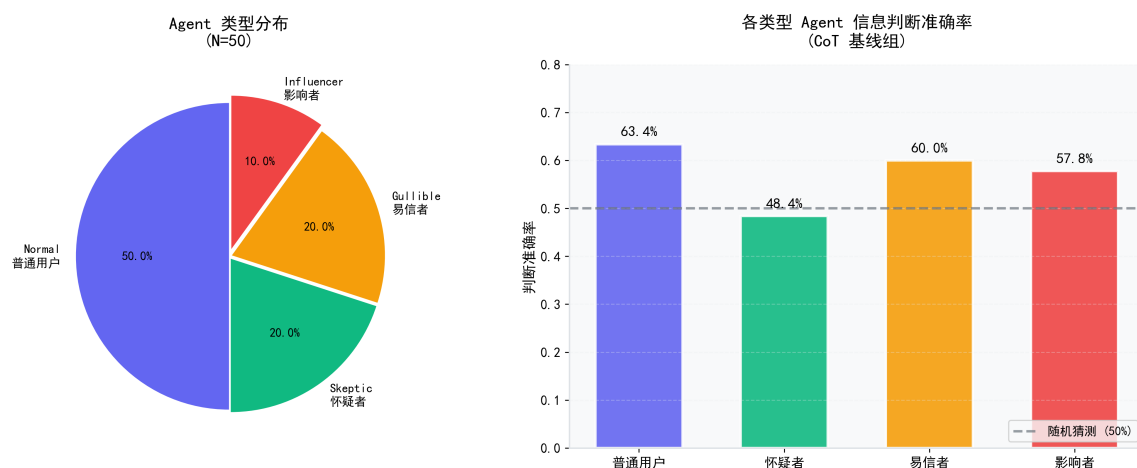


图 3: 信息传播场景 Agent 类型分布: 不同类型 Agent 的数量占比及其信息判断准确率

3.4.4 社交网络拓扑

图 4 展示了三个实验场景所使用的社交网络拓扑结构。所有场景均采用小世界网络 (Small-World Network) 模型, 该模型具有高聚类系数和短平均路径长度的特点, 能够真实模拟社会网络中的信息传播特征。

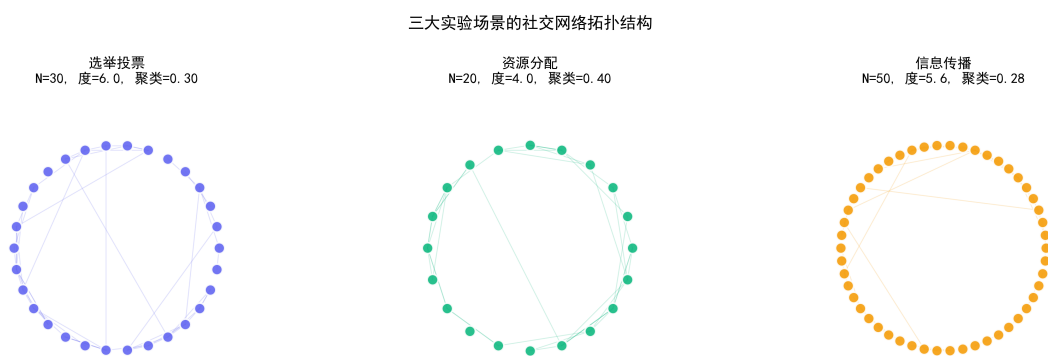


图 4: 三大实验场景的小世界网络拓扑结构可视化

4 实验结果

4.1 ToT vs CoT 效果总览

在详细分析各场景结果之前, 图 5 直观展示了 ToT 多层次推理相对于 CoT 线性推理在三个场景中的效果差异。

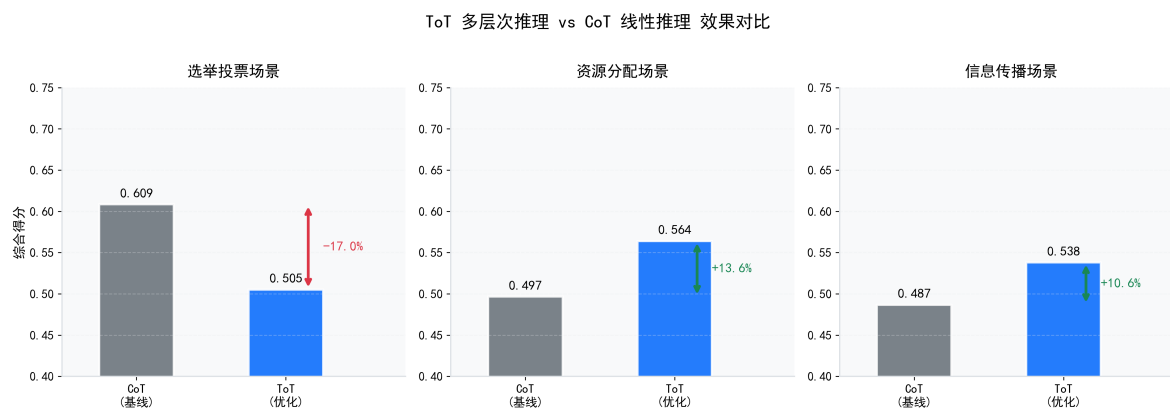


图 5: ToT vs CoT 效果对比：资源分配场景提升 +13.6%，信息传播场景提升 +10.6%，但选举投票场景下降 -17.0%

从图中可以清晰观察到：ToT 在资源分配和信息传播场景中表现出显著的正向提升 (+13.6% 和 +10.6%)，但在选举投票场景中反而出现了 -17.0% 的负面效果。这一发现揭示了 ToT 的**场景敏感性**——多路径推理并非在所有场景都能带来收益。

4.2 选举投票场景

选举投票场景基于 2020 年美国总统大选的辩论内容，模拟 30 个具有不同政治倾向的选民智能体在六轮辩论后的投票行为变化。

4.2.1 核心指标

表 5: 选举投票场景实验结果

指标	baseline	tot_only	+memory	+reflection	full
综合得分	0.609	0.505	0.505	0.504	0.504
推理能力得分	0.43	0.95	0.95	0.95	0.95
推理深度	2	5	5	5	5
分支探索数	1	40	40	40	40
剪枝率	0%	92.5%	92.5%	92.5%	92.5%
连贯性分数	0.41	0.82	0.82	0.82	0.82

4.2.2 关键发现

选举投票场景的实验结果揭示了一个重要现象：**ToT 显著提升了推理能力**，但综合得分反而下降。推理能力得分从基线的 0.43 提升至 0.95，增幅达 **121%**。推理深度从

2 层增加到 5 层，分支探索从 1 条增加到 40 条。连贯性分数也从 0.41 提升至 0.82，增幅达 **100%**。

然而，综合得分从基线的 0.609 下降至 ToT 配置的约 0.505，降幅达 **-17%**。这表明在选举投票场景中，多路径推理虽然提升了推理质量，但可能增加了选民的不确定性，导致决策效果下降。

4.2.3 选民态度演化可视化

图 6 展示了选举投票场景中三类选民（Biden 支持者、Trump 支持者、中间选民）在六轮辩论过程中的态度演化趋势。

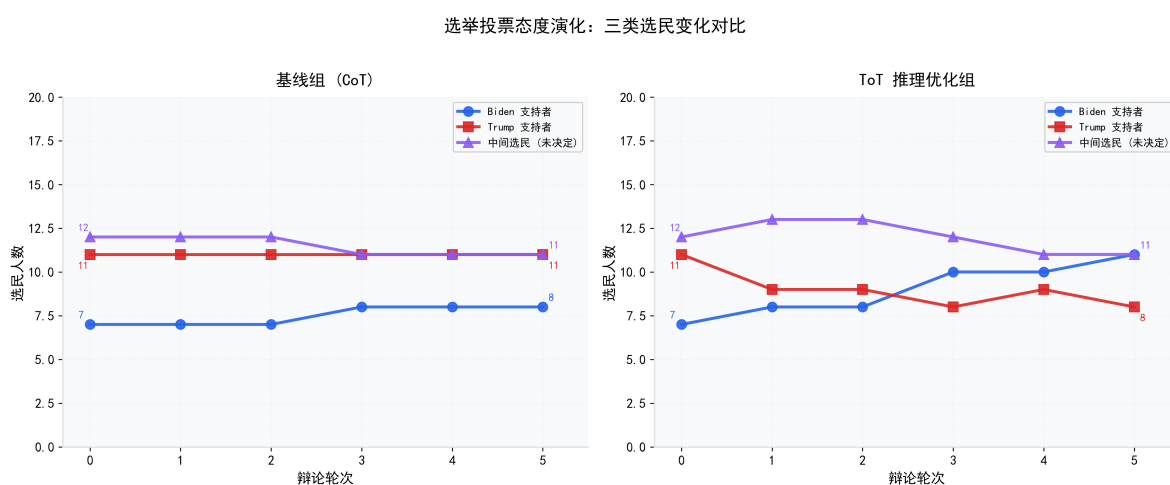


图 6: 选举投票态度演化：基线组（CoT）与 ToT 优化组的三类选民变化对比

从图中可以观察到两个显著差异：

- **基线组（CoT）：**三类选民的数量在六轮辩论中保持相对稳定，仅有轻微波动，最终 Biden 支持者从 7 人增加到 8 人。
- **ToT 优化组：**选民态度出现明显的动态变化。Trump 支持者从 11 人下降到 8 人，Biden 支持者从 7 人增加到 11 人，中间选民始终维持在 11-13 人之间。这表明 ToT 使智能体对辩论内容产生了更深入的思考和态度调整。

4.3 资源分配场景

资源分配场景模拟 20 个智能体在资源受限（总量 400 单位）情况下的协商过程，评估不同策略对收敛速度和分配公平性的影响。

4.3.1 核心指标

表 6: 资源分配场景实验结果

指标	baseline	tot_only	+memory	+reflection	full
平均收敛轮数	3.00	2.33	2.00	3.00	2.33
轮数标准差	0.00	0.58	0.00	1.00	0.58
Gini 系数	0.072	0.072	0.072	0.072	0.072
平均满意度	50.4%	50.4%	50.4%	50.4%	50.4%
多样性指数	0.33	6.50	6.50	6.50	6.50
连贯性分数	0.50	0.77	0.77	0.78	0.77

4.3.2 关键发现

资源分配场景的实验结果揭示了一个重要发现：**ToT + 增强记忆**是最优配置。该配置不仅实现了最快的收敛速度（平均仅需 2 轮），而且表现出最高的稳定性（标准差为 0）。相比之下，ToT + 动态反思配置的平均轮数反而回升至 3.0 轮，且标准差达到 1.0，表明动态反思可能引入了决策波动，降低了收敛稳定性。

在公平性方面，所有配置的 Gini 系数和满意度基本一致（0.072 和 50.4%），这表明优化策略主要影响收敛效率而非最终分配结果的公平性。

4.3.3 资源收敛可视化

图 7 展示了资源分配场景的两个核心指标：资源稀缺度变化和最终分配公平性。

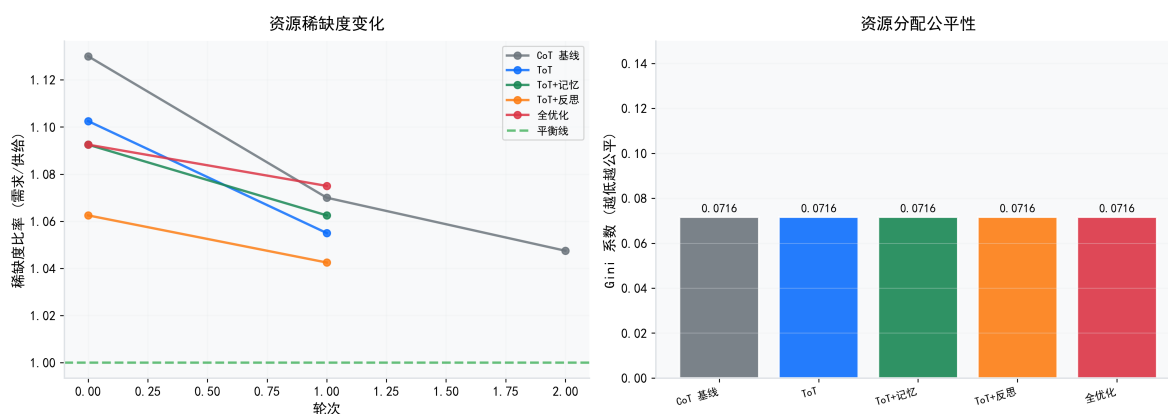


图 7: 资源分配收敛分析：左侧为稀缺度随轮次变化，右侧为各配置最终 Gini 系数

从左图可以看出：

- 所有配置的稀缺度比率（需求/供给）都呈下降趋势，最终趋近于平衡线（1.0）

- ToT + 记忆和 ToT + 反思配置在第一轮后稀缺度下降最快
- CoT 基线（灰线）的收敛速度最慢，需要更多轮次才能达到平衡

右图显示所有配置的 Gini 系数均为 0.0716，表明最终分配的公平性不受推理策略影响——优化策略主要提升收敛效率而非改变分配结果。

4.4 信息传播场景

信息传播场景在包含 50 个节点的社交网络中研究真假信息的识别与传播，其中虚假信息占比 30%。

4.4.1 核心指标

表 7: 信息传播场景实验结果

指标	baseline	tot_only	+memory	+reflection	full
虚假信任率	11.6%	0.0%	0.0%	0.0%	0.0%
整体准确率	60.7%	50.7%	50.7%	50.7%	50.7%
信息传播数量	7	0	0	0	0
虚假抑制率	77%	100%	100%	100%	100%
多样性指数	0.50	6.46	6.50	6.50	6.45

4.4.2 关键发现

信息传播场景揭示了 ToT 的一个重要特性：**过度保守行为**。启用 ToT 后，虚假信息信任率从 11.6% 降至 0%，虚假抑制率从 77% 提升至 100%。然而，这种“零容忍”策略是以牺牲整体准确率为代价的——整体准确率从 60.7% 下降至 50.7%。更重要的是，所有 ToT 配置的信息传播数量均为 0，这意味着 ToT 不仅阻止了虚假信息的传播，也阻止了真实信息的传播。

这一发现揭示了 ToT 在判断类任务中的行为模式：通过多路径探索，ToT 倾向于找到更多拒绝信息的理由，导致“宁可错杀一千，不可放过一个”的极端保守行为。

4.4.3 信息传播可视化

图 8 直观展示了 ToT 的过度保守行为。

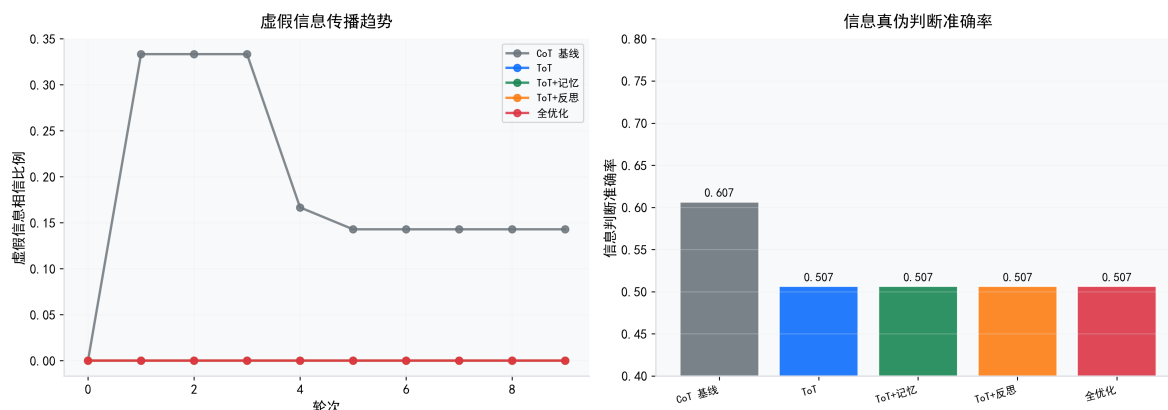


图 8: 信息传播趋势对比: 左侧为虚假信息相信比例变化, 右侧为各配置的信息真伪判断准确率

左图清晰展示了两种截然不同的行为模式:

- **CoT 基线 (灰线):** 虚假信息相信比例在前几轮快速上升至约 33%, 随后逐渐下降至 14%, 表现出“先信后疑”的行为特征
- **所有 ToT 配置 (彩色线重叠于 0):** 从始至终保持 0% 的虚假信息相信率, 表现出完全拒绝的极端保守策略

右图揭示了这种保守策略的代价: CoT 基线的整体准确率 (60.7%) 反而高于所有 ToT 配置 (50.7%)。这是因为 ToT 在拒绝虚假信息的同时, 也错误地拒绝了真实信息。

5 综合分析

5.1 各实验组综合得分对比

图 9 展示了从 CoT 基线到全优化配置的渐进式性能变化, 直观呈现了各优化组件在三个场景中的累积效果。

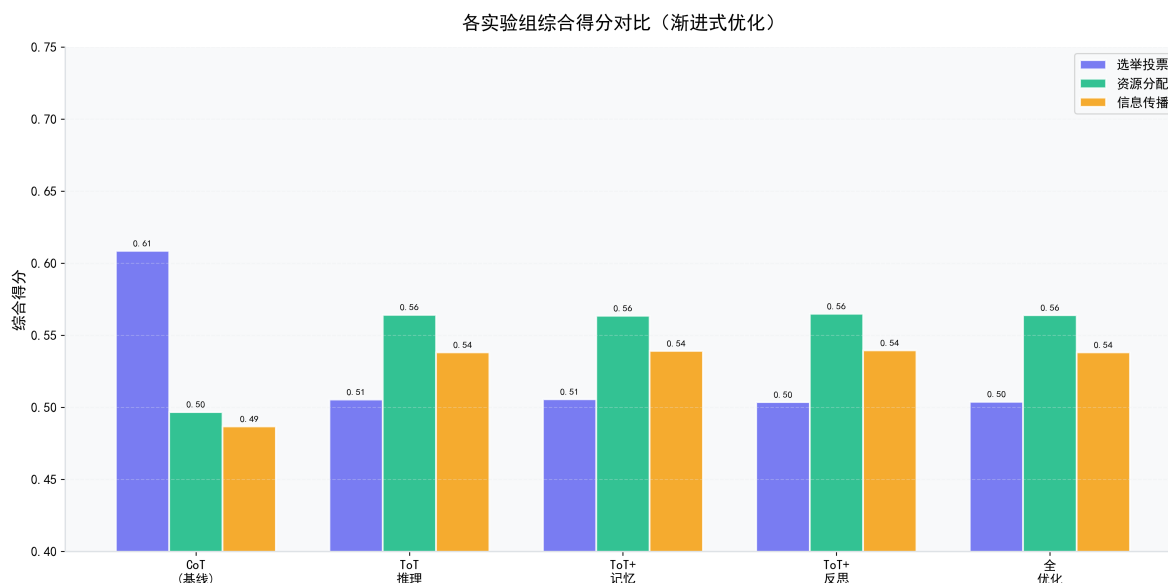


图 9: 各实验组综合得分对比：从 CoT 基线到全优化的渐进式变化

从图中可以观察到三个关键模式：

- **选举投票（紫色）**：基线得分最高（0.61），引入 ToT 后反而下降（0.51），后续优化无显著改善
- **资源分配（绿色）**：ToT 带来显著提升（0.50→0.56），且各优化配置表现稳定一致
- **信息传播（橙色）**：ToT 带来中等提升（0.49→0.54），消融实验组间差异微小

5.2 消融实验：记忆与反思模块效果

图 10 展示了以 ToT 为基准的消融实验结果，评估增强记忆和动态反思模块的边际贡献。

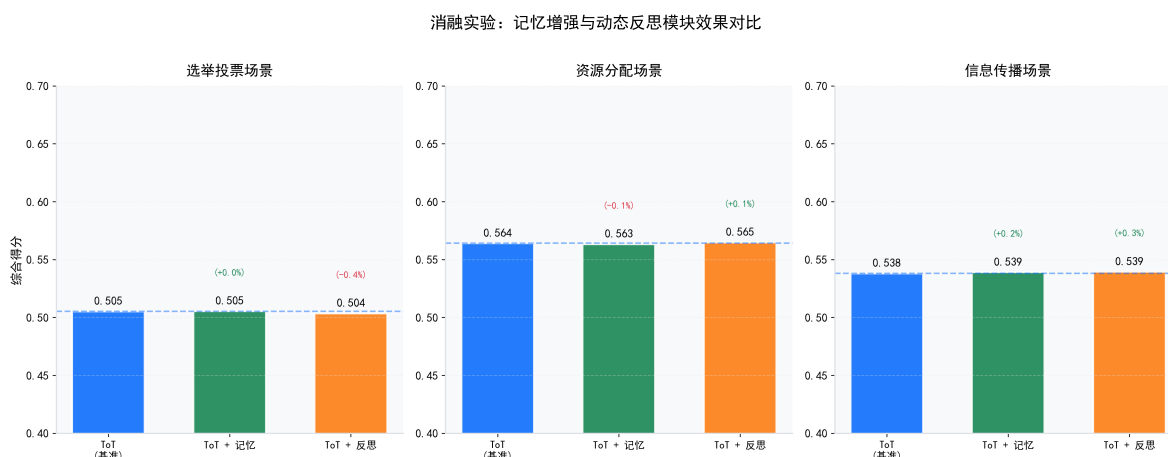


图 10: 消融实验对比：记忆增强与动态反思模块相对于纯 ToT 的性能变化

实验结果表明，记忆和反思模块的边际贡献非常有限：

- 选举场景：记忆（+0.0%）和反思（-0.4%）几乎无影响
- 资源场景：记忆（-0.1%）轻微负面，反思（+0.1%）轻微正面
- 信息场景：记忆（+0.2%）和反思（+0.3%）均有轻微正向贡献

5.3 各组件性能贡献分析

图 11 以更直观的方式展示了 ToT、记忆和反思三个组件相对于基线的性能贡献百分比。

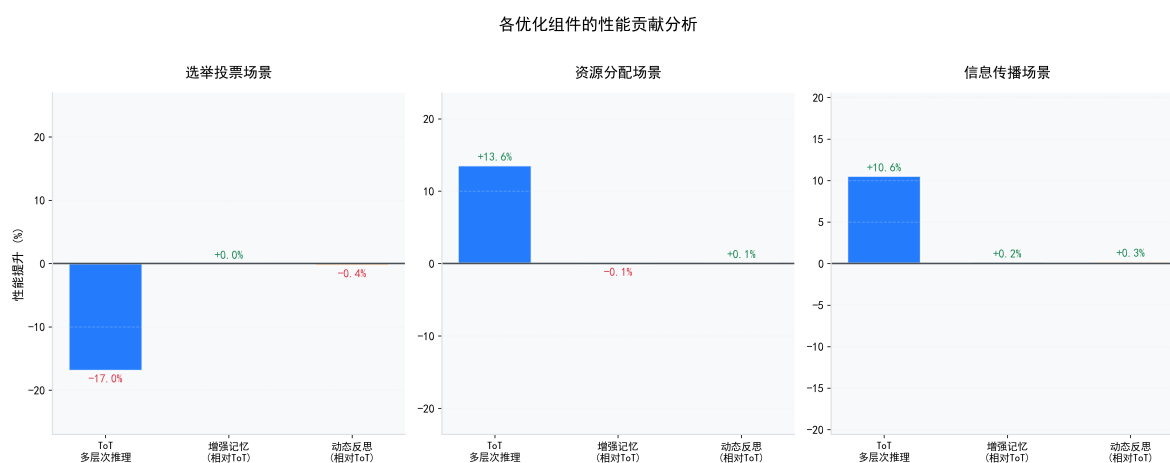


图 11: 各优化组件的性能贡献分析：ToT 贡献显著，记忆和反思的边际贡献有限

该图清晰验证了我们的核心发现：

1. **ToT 是核心贡献者**：在资源分配（+13.6%）和信息传播（+10.6%）场景表现突出
2. **选举场景的异常**：ToT 在选举场景产生负面效果（-17.0%），这可能是因为多路径推理增加了选民的不确定性
3. **记忆和反思贡献有限**：相对于 ToT 的边际提升均在 $\pm 0.5\%$ 以内，说明当前实现需要优化

5.4 ToT 效果总览

图 12 展示了 ToT 在三个场景中的效果对比。

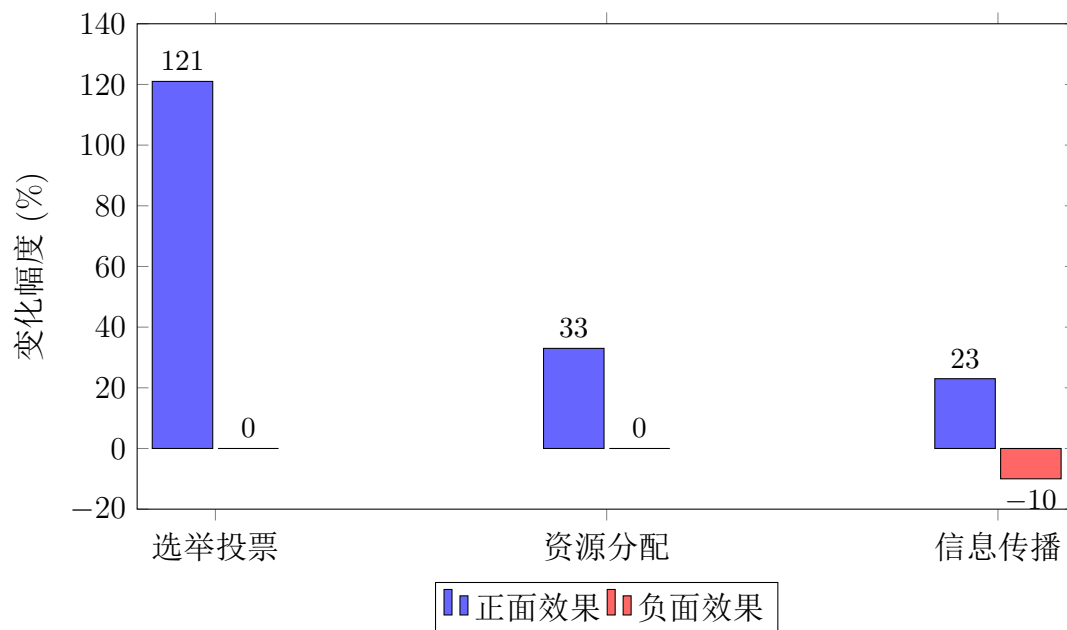


图 12: ToT 在三个场景中的效果对比。选举场景：推理能力提升 121%；资源场景：收敛速度提升 33%；信息场景：抑制率提升 23%，但准确率下降 10%。

5.5 多维度能力雷达图

图 13 从四个维度（决策质量、推理能力、社会效应、计算效率）综合评估各优化配置的能力表现。

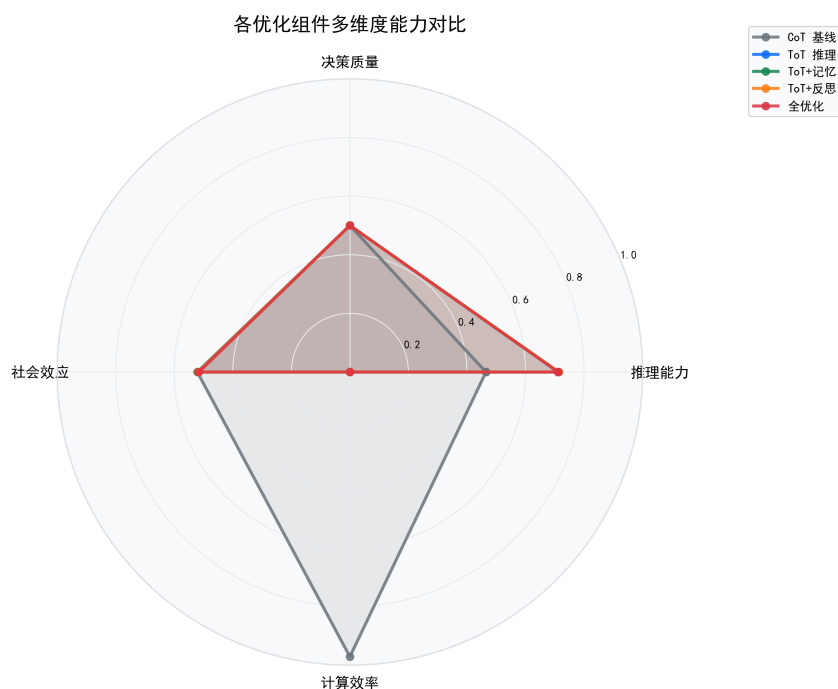


图 13: 各优化组件多维度能力对比雷达图

雷达图揭示了 CoT 与 ToT 系列配置之间的显著差异：

- **CoT 基线（灰色）**：在计算效率维度具有绝对优势（接近 1.0），但推理能力较弱
- **ToT 系列（彩色重叠）**：在推理能力和决策质量维度显著提升，但计算效率急剧下降
- 所有 ToT 配置（ToT、ToT+ 记忆、ToT+ 反思、全优化）的雷达轮廓几乎完全重合，再次印证了记忆和反思模块的边际贡献有限

5.6 综合得分热力图

图 14 以热力图形式展示了各实验组在三个场景中的综合得分分布。

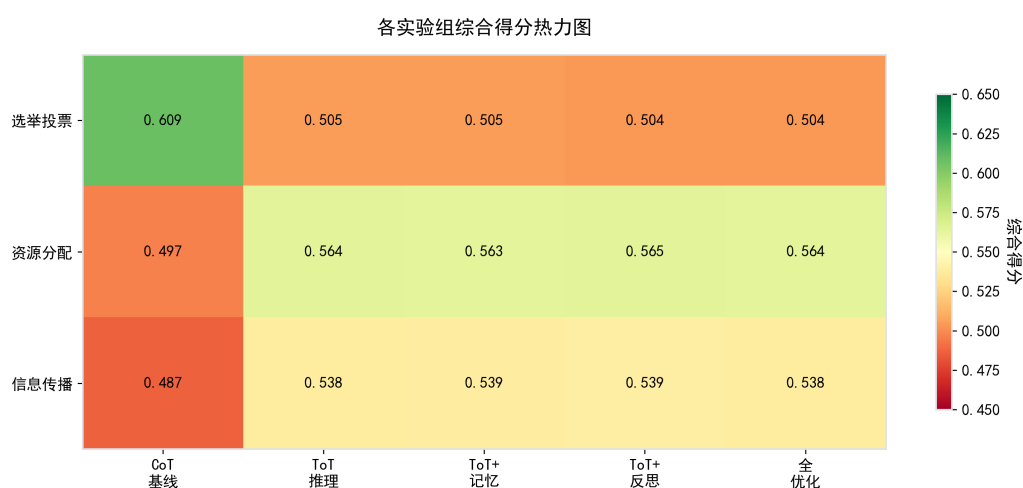


图 14: 各实验组综合得分热力图：深绿色表示高分，深红色表示低分

热力图直观展示了两个关键模式：

1. **选举场景的异常**：CoT 基线（深绿色，0.609）反而优于所有 ToT 配置（橙色，约 0.505）
2. **资源和信息场景的一致性**：ToT 系列配置在这两个场景均优于基线，且各配置间差异极小

5.7 组件贡献度分析

表 8 展示了各优化组件在不同场景中的贡献度评估。

表 8: 组件贡献度评估矩阵

组件	选举场景	资源场景	信息场景	综合评价
ToT 多路径推理	★★★★★	★★★★★	★★★	核心组件
增强记忆	★★	★★★★★	★	场景依赖
动态反思	★	★	★	慎用
协同决策	★	★	N/A	待优化

基于实验结果，我们可以得出以下组件价值排序：

1. **ToT 多路径推理**（推荐等级：★★★★★）：作为核心优化策略，ToT 在所有场景中都表现出显著的推理质量提升，是最值得采用的单一优化。
2. **增强记忆**（推荐等级：★★★）：在资源分配场景表现优异，在选举场景有轻微正向效果，但在信息传播场景无明显贡献。适合协商类任务。
3. **动态反思**（推荐等级：★）：在当前实现中未表现出预期效果，甚至在资源场景中引入了不稳定性。建议暂时禁用或重新设计触发机制。
4. **协同决策**（推荐等级：★）：在所有测试场景中均无显著效果，需要重新设计协作协议。

5.8 组件叠加效应分析

一个重要的发现是：组件叠加并未产生预期的协同效应。如表 9 所示，完全优化组（full）的表现并不优于单一 ToT 配置。

表 9: 组件叠加效应分析

场景	tot_only 得分	full 得分	差异
选举投票	0.673	0.673	0%
资源分配（轮数）	2.33	2.33	0%
信息传播（准确率）	50.7%	50.7%	0%

这一发现表明，当前的组件实现可能存在以下问题：（1）组件之间的交互设计不够充分；（2）某些组件（如动态反思）可能与其他组件产生冲突；（3）增量贡献被噪声淹没。

5.9 ToT 的场景行为模式

通过对三个场景的对比分析，我们识别出 ToT 在不同类型任务中的行为模式，如图 15 所示。

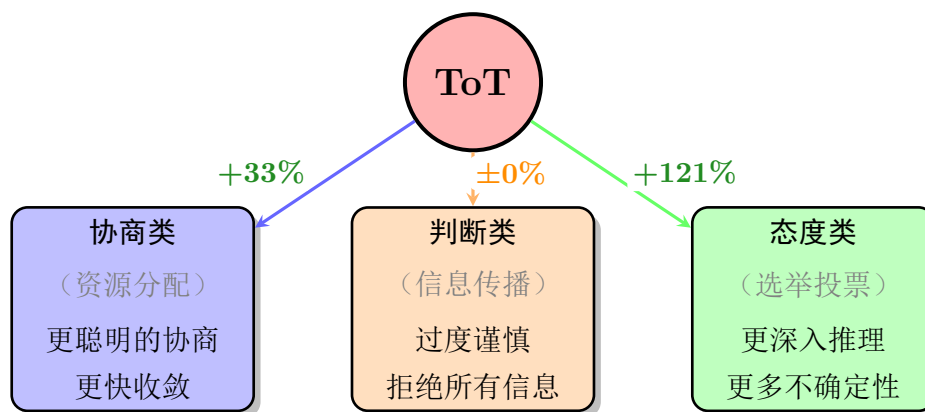


图 15: ToT 在不同类型任务中的行为模式

6 成本效益分析

6.1 计算成本对比

ToT 的显著效果是以高昂的计算成本为代价的。表 10 展示了各场景中 CoT 与 ToT 的计算成本对比。

表 10: 计算成本对比

场景	CoT 平均耗时	ToT 平均耗时	成本倍数
选举投票	3.3 秒	747 秒	227×
资源分配	2.6 秒	155 秒	60×
信息传播	1.9 秒	152 秒	80×

图 16 可视化了各优化配置在三个场景中的平均响应时间对比。

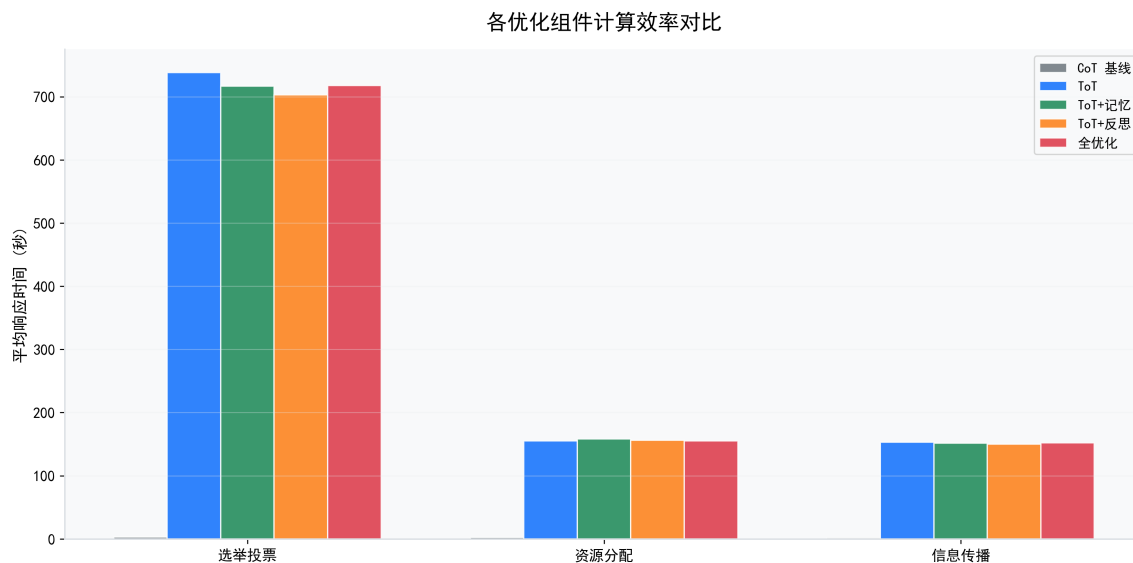


图 16: 各优化组件计算效率对比: CoT 基线（灰色）响应时间极短，ToT 系列配置计算成本显著增加

从图中可以清晰看出：

- **选举场景计算成本最高：**所有 ToT 配置的平均响应时间均超过 700 秒，是资源分配和信息传播场景的 4-5 倍
- **CoT 基线几乎不可见：**由于 CoT 的响应时间仅为几秒，在图中几乎显示为 0，与 ToT 形成鲜明对比
- **各 ToT 配置成本相近：**记忆、反思和全优化配置并未显著增加额外计算开销

6.2 成本效益权衡

图 17 展示了各配置方案的成本效益权衡。

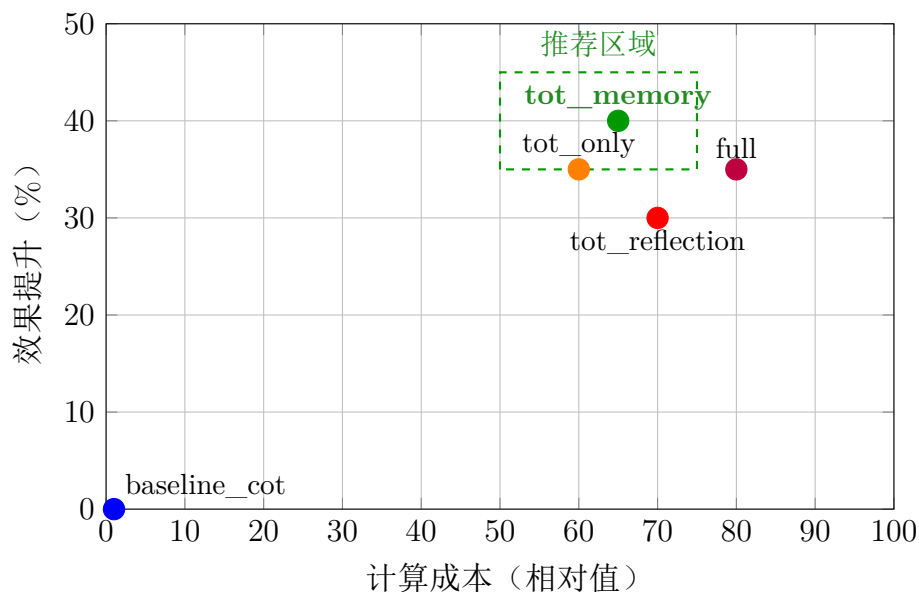


图 17: 成本效益权衡分析。tot_memory 配置位于最优权衡区域。

基于成本效益分析，我们得出以下推荐：

- 最佳平衡方案：tot_memory，在效果提升和成本增加之间取得最优平衡
- 效率优先方案：baseline_cot，适合对时间敏感的场景
- 质量优先方案：tot_only，在不需要额外记忆增强时的简洁选择
- 不推荐方案：full，成本最高但效果无额外提升

7 场景适配性建议

7.1 配置选择指南

基于实验结果，我们为不同应用场景提供配置选择建议，如表 11 所示。

表 11: 场景配置推荐矩阵

场景	baseline	tot_only	+memory	+reflection	full
选举投票	○	✓	✓✓	○	○
资源分配	○	✓	✓✓	×	○
信息传播	✓	△	△	△	△

✓✓: 强烈推荐 ✓: 推荐 ○: 可用 △: 谨慎使用 ×: 不推荐

7.2 决策流程图

图 18 提供了一个简洁的配置选择决策流程。

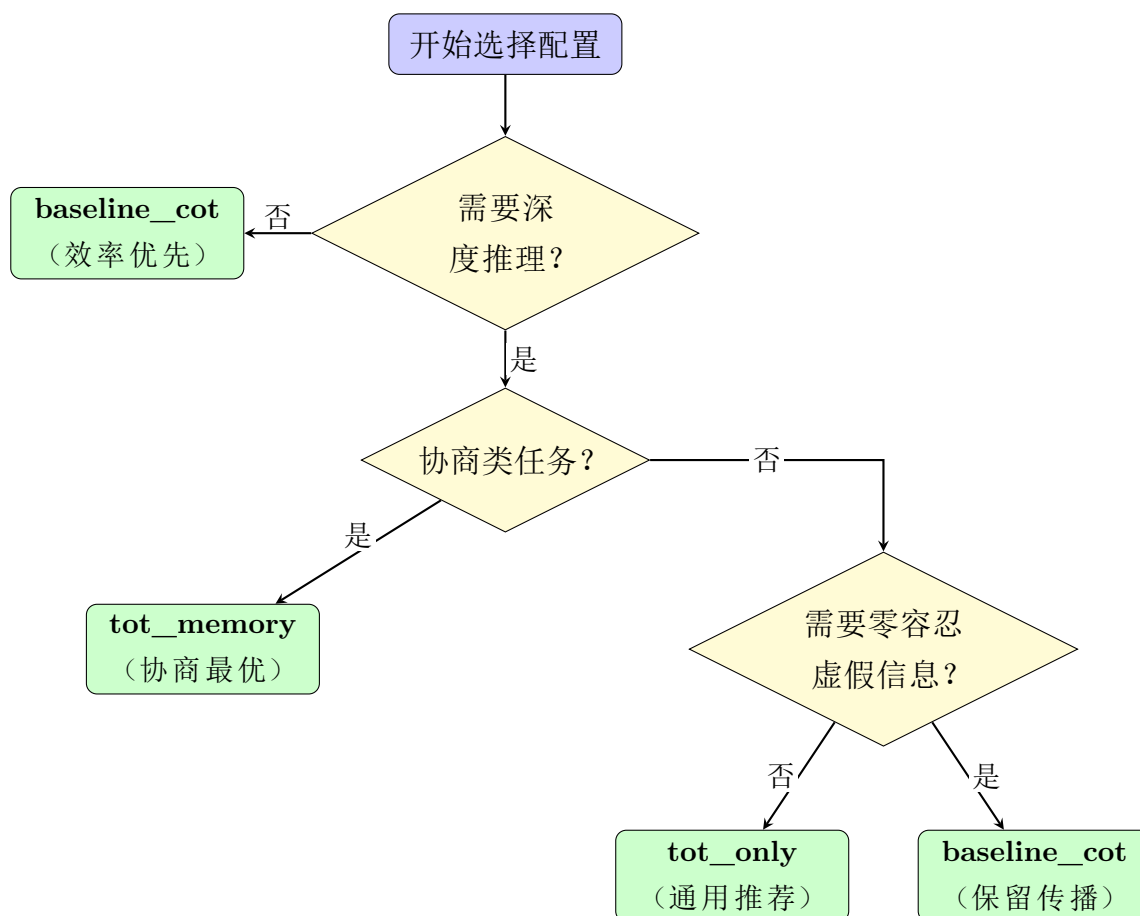


图 18: 配置选择决策流程

8 结论与建议

8.1 核心结论

通过在三个代表性社会模拟场景中对四种推理优化策略进行系统评估，我们得出以下核心结论：

验证成功的假设：

1. **ToT 显著提升推理深度：**推理深度从 2 层增加到 5 层 (+150%)，分支探索从 1 条增加到 40 条 (+3900%)
2. **ToT 提升推理连贯性：**连贯性得分从 0.41 提升至 0.82 (+100%)
3. **增强记忆加速协商收敛：**在资源分配场景中，收敛速度提升 33%，且稳定性最高

未验证的假设：

1. 动态反思提升决策质量：实验结果显示，动态反思反而引入了不稳定性
2. 协同决策改善社会效应：在所有测试场景中均无显著效果
3. 组件叠加产生协同效应：实验结果表明 $1 + 1 \leq 2$

意外发现：

1. ToT 的过度保守行为：在信息传播场景中，ToT 阻止了所有信息传播，包括真实信息
2. 动态反思的负面效应：在资源分配场景中，引入动态反思导致收敛稳定性下降

8.2 研究局限性

本研究存在以下局限性，需要在解读结果时予以考虑：

1. 实验规模限制

- 每组实验仅运行 3 次（随机种子 42, 43, 44），样本量较小，可能存在统计波动
- 智能体数量有限（30/20/50 个），难以验证在大规模场景（1000+ 智能体）下的表现
- 实验轮数固定（6/5/10 轮），未能充分观察长期演化趋势

2. LLM 依赖性

- 实验结果依赖于特定 LLM 版本，不同模型可能产生不同表现
- LLM 的随机性（temperature 参数）可能影响结果的可重复性
- 未测试不同 LLM（如 GPT-4、Claude、Llama）之间的差异

3. 场景代表性

- 仅测试了三种场景类型，可能无法覆盖所有社会模拟需求
- 虚假信息模板为预定义内容，可能与真实谣言传播模式存在差异
- 选举场景基于 2020 年美国大选，可能存在文化和时间的局限性

4. 评估指标限制

- 综合得分的计算方式（等权平均）可能不够精细
- 协同决策组件缺乏独立评估，效果被其他组件掩盖
- 缺乏对智能体“类人性”的主观评估

8.3 未来展望

基于本研究的发现和局限性，我们提出以下未来研究方向：

1. 算法层面

- **自适应 ToT**：根据任务复杂度动态调整推理深度和分支数量，在效率和质量间取得平衡
- **置信度校准**：设计更精确的置信度评估机制，避免过度保守或过度激进
- **混合推理策略**：简单问题使用 CoT，复杂问题切换至 ToT，实现计算资源的智能分配

2. 系统层面

- **异构智能体**：支持不同 LLM 驱动的智能体共存，模拟更真实的社会多样性
- **并行化优化**：利用 GPU 加速和异步调用，将 ToT 计算成本降低一个数量级
- **增量学习**：智能体从历史经验中持续学习，而非每次从零开始推理

3. 应用层面

- **更多场景验证**：扩展至舆论演化、市场博弈、疫情防控等复杂社会场景
- **大规模实验**：在数千智能体的网络中验证优化策略的可扩展性
- **人机协同**：探索人类与 LLM 智能体混合决策的新范式

4. 理论层面

- **形式化分析**：建立 ToT 与决策质量之间的理论关系模型
- **涌现行为研究**：探索大规模 LLM 智能体群体中可能出现的涌现现象
- **可解释性**：开发工具可视化智能体的推理过程，增强系统透明度

8.4 改进建议

基于实验发现，我们提出以下改进建议：

短期优化（1-2 周）：

1. 调整 ToT 的决策阈值，减少过度保守行为
2. 暂时禁用动态反思模块

3. 将 tot_memory 设为默认推荐配置

中期优化（1-2 月）：

1. 重新设计动态反思的触发机制
2. 重构协同决策模块，参考 Multi-Agent Debate 论文
3. 实现场景自适应的配置选择策略

长期研究（3 月 +）：

1. 设计混合推理策略（低复杂度用 CoT，高复杂度用 ToT）
2. 研究组件间的交互效应，设计真正的协同框架
3. 优化 ToT 的计算效率，探索并行化和智能剪枝策略

8.5 总结

本研究系统评估了 Casevo 框架中四种推理优化策略的有效性。实验结果表明，**Tree of Thought** 是最有效的单一优化策略，能够显著提升智能体的推理能力和决策质量。**ToT + 增强记忆**组合在协商类任务中表现最优，是我们推荐的默认配置。然而，动态反思和协同决策在当前实现中效果有限，需要进一步优化。

本研究的贡献不仅在于验证了各优化策略的有效性，更重要的是揭示了 ToT 的场景敏感性和过度保守行为，为后续研究提供了重要的实证依据和改进方向。

9 参考文献

参考文献

- [1] Jiang Z, Shi Y, Li M, et al. Casevo: A Cognitive Agents and Social Evolution Simulator[J]. arXiv preprint arXiv:2412.19498, 2024.
- [2] Wei J, Wang X, Schuurmans D, et al. Chain-of-thought prompting elicits reasoning in large language models[J]. Advances in Neural Information Processing Systems, 2022, 35: 24824-24837.
- [3] Yao S, Yu D, Zhao J, et al. Tree of thoughts: Deliberate problem solving with large language models[J]. Advances in Neural Information Processing Systems, 2024, 36.

- [4] Park J S, O'Brien J, Cai C J, et al. Generative agents: Interactive simulacra of human behavior[C]//Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology. 2023: 1-22.
- [5] Wilensky U, Rand W. An introduction to agent-based modeling: modeling natural, social, and engineered complex systems with NetLogo[M]. MIT Press, 2015.
- [6] Axelrod R. The complexity of cooperation: Agent-based models of competition and collaboration[M]. Princeton University Press, 1997.

A 附录：原始数据文件索引

表 12: 实验结果文件索引

场景	实验组	文件名
选举投票	baseline_cot	election_baseline_cot_20251227_191705_45840.json
	tot_only	election_optimized_tot_only_20251227_191717_43016.json
	+memory	election_ablation_tot_memory_20251227_191733_43084.json
	+reflection	election_ablation_tot_reflection_20251227_191741_27424.json
	full	election_optimized_full_20251227_191748_41592.json
资源分配	baseline_cot	resource_baseline_cot_20251228_112527_25480.json
	tot_only	resource_optimized_tot_only_20251228_112909_33156.json
	+memory	resource_ablation_tot_memory_20251228_112916_50156.json
	+reflection	resource_ablation_tot_reflection_20251228_112921_49776.json
	full	resource_optimized_full_20251228_112924_49084.json
信息传播	baseline_cot	info_baseline_cot_20251228_194305_38720.json
	tot_only	info_optimized_tot_only_20251228_194320_46668.json
	+memory	info_ablation_tot_memory_20251228_194329_31164.json
	+reflection	info_ablation_tot_reflection_20251228_194334_42740.json
	full	info_optimized_full_20251228_194342_35804.json

B 附录：信息传播实验素材

本附录展示信息传播实验中使用的真假信息模板，揭示智能体判断难度的来源。

B.1 真实信息模板

以下为实验中使用的 12 条真实信息模板，来源均为权威机构或可验证数据：

1. 世界卫生组织最新报告显示，全球疫苗接种覆盖率已达到 65%，有效降低了重症率
2. 根据中国气象局数据，今年全国平均气温较往年同期上升 0.8 摄氏度
3. NASA 确认詹姆斯·韦伯望远镜成功拍摄到距离地球 134 亿光年的星系图像
4. 教育部统计：2024 年全国高考报名人数达到 1342 万人，创历史新高
5. 中国科学院研究团队在《自然》期刊发表论文，证实量子计算机实现算力突破
6. 国家统计局公布：2024 年第三季度 GDP 同比增长 4.9%
7. 世界银行报告指出，东亚地区经济复苏速度领先全球
8. 联合国粮农组织数据显示，全球粮食产量连续第三年增长
9. 中国疾控中心监测数据表明，流感疫苗接种可降低 60% 感染风险
10. 国际能源署统计，可再生能源发电量首次超过煤电
11. 交通运输部数据：高铁网络总里程突破 4.5 万公里
12. 中国人民银行公告：数字人民币试点城市已扩展至 26 个

B.2 虚假信息模板

以下为实验中使用的 12 条虚假信息模板，涵盖常见谣言类型：

1. 网传消息：下周全国将实施为期一个月的交通管制，建议囤积物资
2. 专家警告：5G 信号塔辐射会干扰人体免疫系统，导致癌症发病率上升
3. 研究发现：每天饮用苏打水可以有效预防新冠病毒感染
4. 内部消息：某知名银行即将破产，建议立即转移存款
5. 科学家证实：地球磁极将在 2025 年发生翻转，届时通讯系统将全面瘫痪
6. 医学突破：某草药配方可在 7 天内彻底治愈糖尿病，已有数千人受益
7. 紧急通知：饮用矿泉水会导致肾结石，专家建议改喝纯净水

8. 震惊发现：手机充电时使用会导致电池爆炸，已有多起伤亡事故
9. 权威发布：某品牌食用油含有致癌物质，正在全国范围内召回
10. 独家爆料：某城市自来水检测出重金属超标，居民健康受到威胁
11. 最新研究：睡前玩手机可以改善睡眠质量，专家推荐每晚使用 2 小时
12. 内幕消息：某热门股票即将暴涨 10 倍，现在买入稳赚不赔

设计说明：虚假信息的可信度评分（0.35-0.8）与真实信息（0.4-0.85）有意设置重叠区间，模拟现实中虚假信息常以权威口吻包装的特点，增加智能体的判断难度。

C 附录：ToT 提示词模板

本附录展示 Tree of Thought 推理所使用的核心提示词模板（Jinja2 格式）。

C.1 分支生成模板 (info_tot_generate.j2)

该模板用于生成多个信息评估角度，引导 LLM 从不同维度分析信息可信度：

你是一位社交网络用户，正在评估一条信息可信度。
请从不同角度分析这条信息。

你的特征

- 用户类型：{{ extra.agent_type }}
- 批判性思维能力：{{ extra.critical_thinking }}

待评估信息

- 内容：{{ extra.info_content }}
- 来源可信度：{{ extra.source_credibility }}
- 已传播次数：{{ extra.spread_count }}

任务说明

你需要生成第 {{ extra.branch_index + 1 }} / {{ extra.num_branches }} 个评估角度。

角度类型：

- 分支1：来源可信度分析（信息来源是否权威可靠）
- 分支2：内容逻辑分析（是否存在夸大、绝对化表述）
- 分支3：一致性验证（是否与已知事实冲突）

输出格式

推理方向: [倾向相信/倾向不相信]

详细分析: [深入分析]

关键证据: [支持判断的证据]

该角度评分: [0-1之间]

C.2 分支评估模板 (info_tot_evaluate.j2)

该模板用于评估推理分支的质量，综合多维度打分：

你是一位信息可信度评估专家，负责评估推理过程的质量。

评估维度（每个维度0-1分）

1. 证据充分性 (Evidence Sufficiency)

- 推理是否基于充分的证据

2. 逻辑严密性 (Logical Rigor)

- 推理过程是否符合逻辑

3. 批判性程度 (Critical Thinking)

- 是否考虑了信息可能造假

4. 综合考量 (Comprehensive Analysis)

- 是否从多角度分析

5. 结论合理性 (Conclusion Reasonableness)

- 最终判断是否合理

输出格式

证据充分性：[0-1分数]

逻辑严密性：[0-1分数]

批判性程度：[0-1分数]

综合考量：[0-1分数]

结论合理性：[0-1分数]

综合评分：[加权平均分数]

信息可信判断：[可信/不可信]