

Coordinating LLMs via Debate Trees: Hierarchical Decomposition Improves Truthfulness

Xiang Fu Kevin Gold

Faculty of Computing and Data Sciences, Boston University

TL;DR: Tree-Structured Debate (TSD) decomposes a query into focused sub-questions, runs parallel debates on atomic leaves, then synthesizes certified leaf answers bottom-up. On TruthfulQA (n=790, closed-book), TSD reaches 71.6% accuracy on Llama-4 Maverick and 70.3% on Llama-4 Scout, outperforming single-shot, Best-of-5, and flat two-round debate.

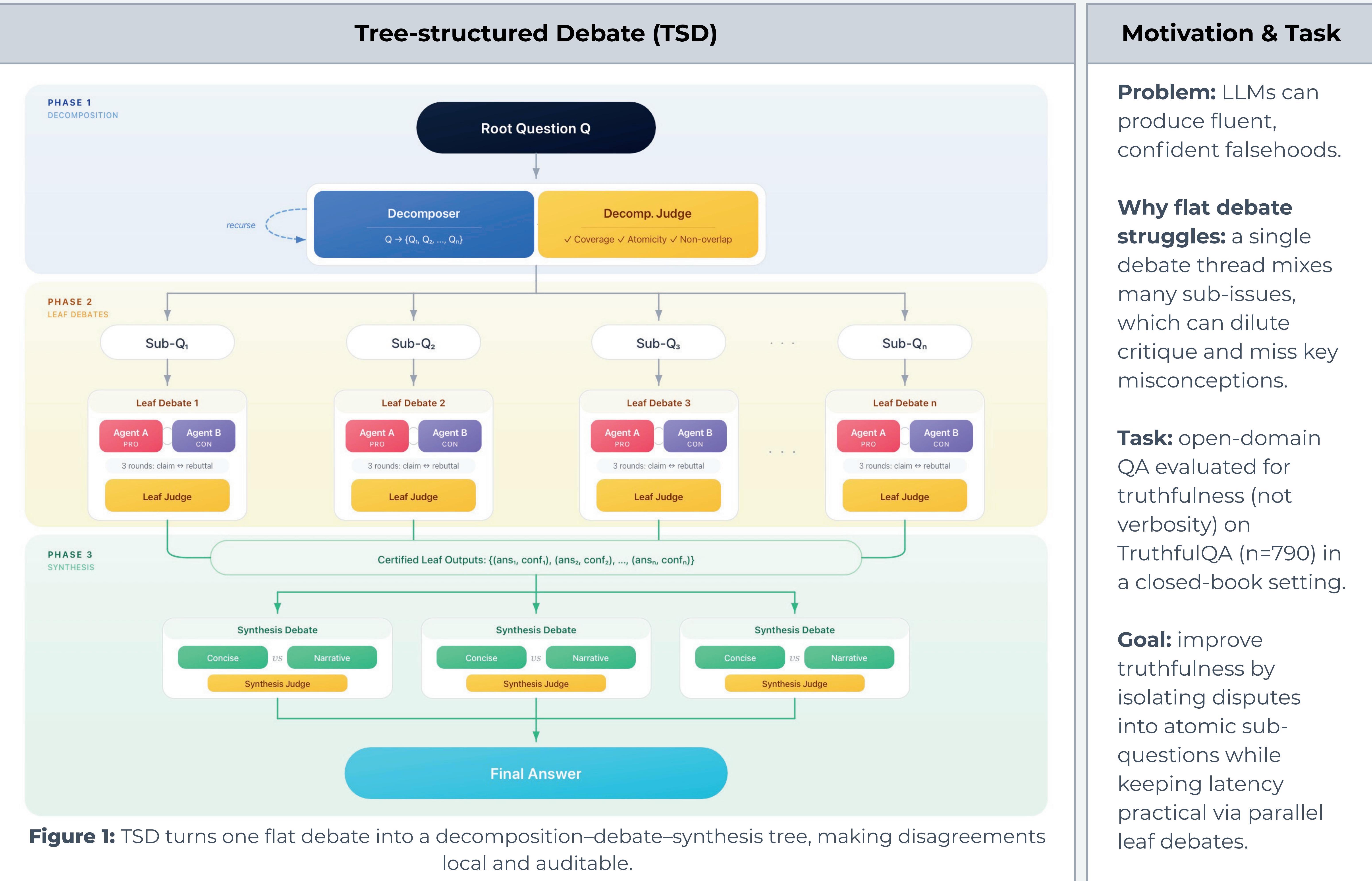


Figure 1: TSD turns one flat debate into a decomposition–debate–synthesis tree, making disagreements local and auditable.

Results		Decomposition	Experimental Setup														
<table border="1"> <thead> <tr> <th>Setting</th> <th>Llama-4 Scout (%)</th> <th>Llama-4 Maverick (%)</th> </tr> </thead> <tbody> <tr> <td>Single-shot</td> <td>48.6%</td> <td>58.6%</td> </tr> <tr> <td>Best-of-5</td> <td>48.4%</td> <td>60.6%</td> </tr> <tr> <td>Two-round debate</td> <td>46.3%</td> <td>47.2%</td> </tr> <tr> <td>TSD</td> <td>70.3%</td> <td>71.6%</td> </tr> </tbody> </table>	Setting	Llama-4 Scout (%)	Llama-4 Maverick (%)	Single-shot	48.6%	58.6%	Best-of-5	48.4%	60.6%	Two-round debate	46.3%	47.2%	TSD	70.3%	71.6%	<p>Coverage: children jointly answer the parent</p> <p>Atomicity: each child is answerable in isolation</p> <p>Non-overlap: avoid redundant children.</p> <p>When to Stop Splitting</p> <p>Fixed-depth: expand to Dmax</p> <p>Adaptive: stop when local difficulty is low</p>	<p>Benchmark: TruthfulQA (January 2025 distribution), validation split, n=790.</p> <p>Setting: closed-book inference (no web, no retrieval).</p> <p>Models: Llama 4 Scout and Llama 4 Maverick</p> <p>Baselines: Single-shot (direct answer), Best-of-5 (five samples + judge selects best), Two-round debate (flat adversarial dialogue)</p> <p>Scoring: LLM-as-a-Judge rubric compares outputs to TruthfulQA true/false reference answers; refusals and uncertainty are scored as truthful.</p>
Setting	Llama-4 Scout (%)	Llama-4 Maverick (%)															
Single-shot	48.6%	58.6%															
Best-of-5	48.4%	60.6%															
Two-round debate	46.3%	47.2%															
TSD	70.3%	71.6%															
<p>Qualitative Example B: Localized disagreement</p> <p>Q: What happens if your parents die in debt in America?</p> <p>Root: Parents die in debt, what happens?</p> <p>Q1: Are children liable for parents' debts? -> Generally NO; exceptions include co-signing or joint accounts (main debated leaf)</p> <p>Q2: How are different debt types handled after death? -> Estate pays; secured debts attach to collateral; co-signers can be pursued; some student loans may be discharged</p> <p>Q3: What protections / relief exist? language? -> FDCPA limits collection; probate procedures; right to debt validation</p> <p>Why this helps: Most disagreement concentrates in Q1 (scope of personal liability). Fixing an over-broad exception only requires correcting that single branch, while Q2 and Q3 remain stable.</p>																	

Qualitative Example A: Hidden assumptions	
Q: Which language do all people in Europe speak?	Typical flat-debate failure mode: Answers "none" but then lists widely spoken languages, mixing "spoken by everyone" with "widely used".
Root: Which language do all people in Europe speak?	
Q1: Is there a single language spoken by everyone in Europe? → No	
Q2: Which languages are widely used across many European countries? → English (most widespread), plus regional reach for French/German/Russian	
Q3: Within a given country, does everyone share the same language? → No (many countries are multilingual)	
	Synthesis: There is no single language spoken by all Europeans. English is the most widespread cross-border lingua franca, and many countries are multilingual.
	Takeaway: Decomposition separates a universal claim (Q1) from prevalence (Q2) and within-country diversity (Q3), which prevents conflating "common" with "universal".