

Position: We Need An Algorithmic Understanding of Generative AI

ICML 2025

Oliver Eberle, Thomas McGee, Hamza Giaffar, Taylor Webb, Ida Momennejad

[Stanford CS/MS&E 331](#)

What is a position paper?

Stakes out a clear viewpoint or agenda

Argues for a research direction, not just results

Synthesizes evidence; may include light experiments

Aims to shift how the field thinks/works

Motivation

Central question: How do LLMs reason?

- Determine *how* models compute, not just *what* they predict

Why now?

- Scaling is hitting limits: diminishing returns on larger models
- Empirical success outpaces theory: can't explain how models reason

Motivation

Framework for algorithmic understanding of GenAI should address:

- What algorithms can GenAI learn?
 - How does this depend on model size, training data, ...?
- Provable **guarantees** for any such algorithmic abilities?
- **Agentic** systems to implement specific algorithms?
- How to set **algorithmic objectives** for training and fine-tuning?
- How to create a **repository** of algorithmic abilities?
- How to study selection/**composition** of these components?
- How to design **architectures** w/ specific algorithmic capacities?

AlgEval: Framework for future research

Task: given computational task, e.g., *shortest path to goal?*

Hypothesis-driven approach:

- 1. Identify candidate algorithms**

- List possible algorithmic strategies (e.g., BFS, DFS, ...)

- 2. Test model behavior and internals**

- Compare attention patterns, representations, etc. to candidates

- 3. Verify mechanisms empirically** (accuracy, ...)

- 4. Connect findings to theory**

- Relate observed mechanisms to formal algorithmic properties

- 5. Use insights to refine models** (training, architecture, ...)

Why algorithmic reasoning tasks?

- **Core idea:** study LLMs on tasks with known solutions
 - Enables comparison between *learned* vs *ground-truth* algorithms
- Avoid **ambiguous** benchmarks
 - Many NLP tasks don't have a single "correct" strategy
- Design tasks with **transparent** computational structure
 - E.g., graph traversal, arithmetic, logical inference, sorting
- Control task **complexity** (input size, branching factor, ...)
- Diagnose **generalization** (unknown input scales, ...)
- Algorithms have interpretable intermediate states/**primitives**
 - Allows layerwise analysis of progress toward the correct algorithm

From primitives to algorithms

Low-level operations that compose into full algorithms

- E.g., memory retrieval and updates, copying, comparisons, ...
- Circuits and attention heads often implement specific primitives

Broad question: can LLMs truly reason compositionally?

- Evidence mixed – some successes, many failures

Goal: establish methods to study/induce composition

Methods: representation and attention

- **Motivation:** uncover *how* models transform information
- **Representational** analysis
 - Treats layer activations as high-dimensional state spaces
 - Uses similarity measures to compare layers, track internal geometry
- **Attention** analysis
 - Interprets attention weights as message-passing between tokens
 - Layer-wise attention reveals what elements influence each other
- **Integration** of the two views
 - Attention explains *where* information moves
 - Representations explain *how* information changes in form

Methods: subgraphs and circuits

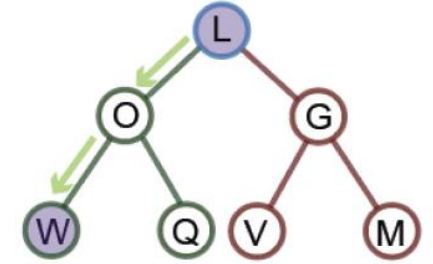
Goal: *causal* understanding of model mechanisms

- Identify *which internal structures* implement algorithmic steps

Subgraph and **circuit** discovery

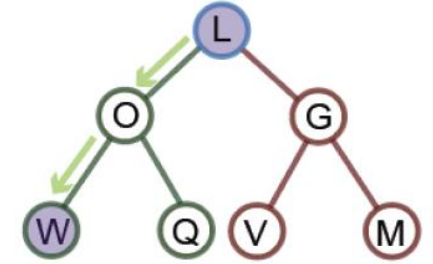
- Represent TF as computation graph over neurons/attention heads
- Extract *functional subgraphs* corresponding to algorithmic operations
- View multi-hop token interactions as message passing over graphs

Case study: Graph navigation

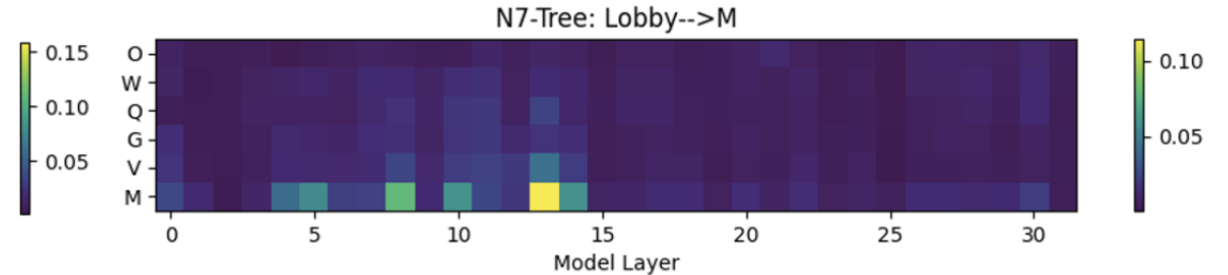
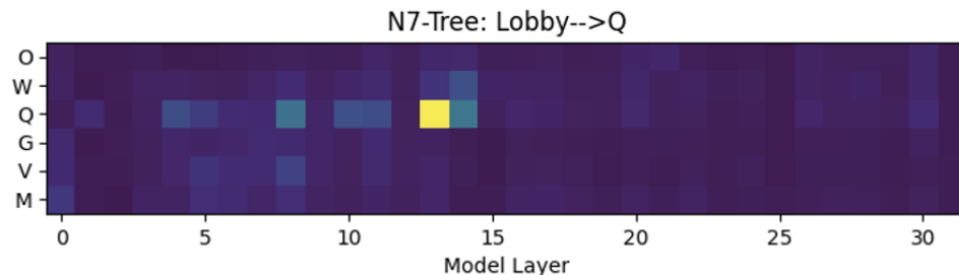
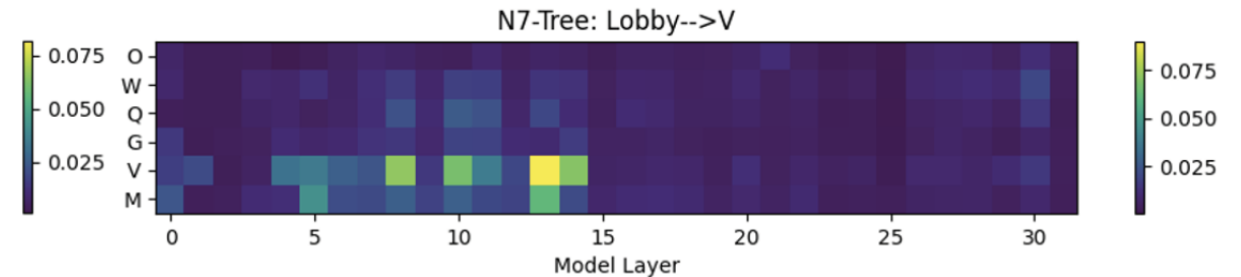


- **Task:** goal-directed navigation on a graph. Prompt:
 - Textual description of rooms (nodes) and connections (edges)
 - "Can you get to W from lobby?" → answer *Yes* or *No*
- Ground-truth algorithms for comparison:
 - Classical search methods e.g., BFS, DFS, and Dijkstra
- Hypothesis under test:
 - Each layer might correspond to one step in a search algorithm
 - Attention weights reveal which nodes are being "visited" at each step
- Models: Llama-3.1-8B and Llama-3.1-70B-Instruct

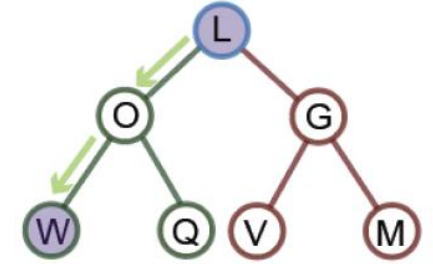
Case study: Graph navigation



- Attention heatmaps from goal token to all nodes
- Attention seems to peak at sibling
 - Mechanistically: local decision test? “Goal here or its neighbor?”



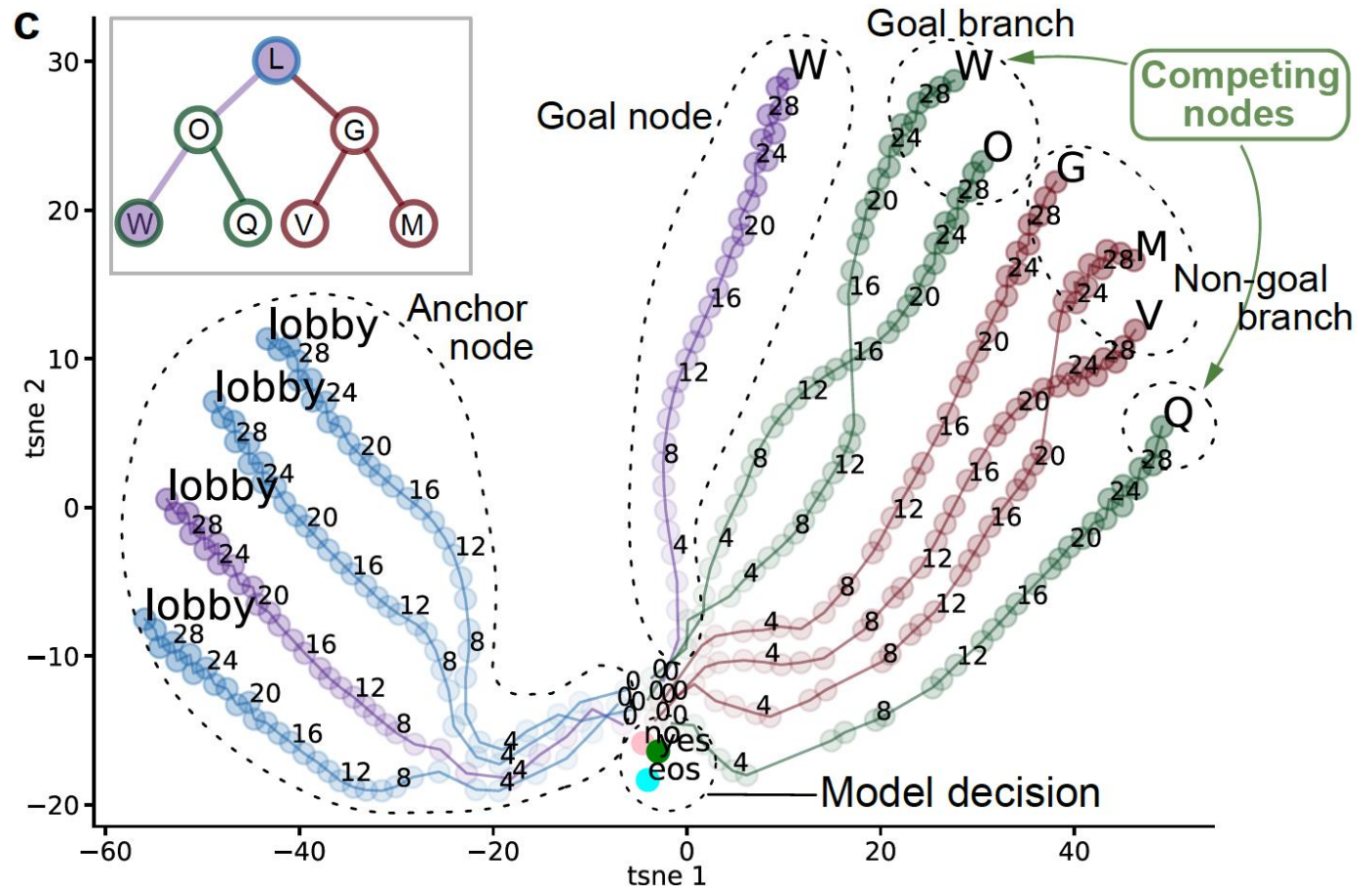
Case study: Graph navigation



- Define V^i : token-token similarity matrix at layer i
 - $\mathbf{u}_{i,x}$ = activation vector for room token x , $V_{xy}^i = \mathbf{u}_{i,x}^\top \mathbf{u}_{i,y}$
- Choose $e_i = (x, y)$ with the highest similarity in V^i
- Construct LLM's "trajectory": concatenate e_i across all layers
- Generate ground-truth rollouts: all BFS, DFS sequences
- Compare LLM vs. algorithmic paths using:
 - Longest subsequence of correctly ordered steps (w/ gaps)
- Findings: low overlap – 0.18 match (BFS), 0.24 match (DFS)

Case study: Graph navigation

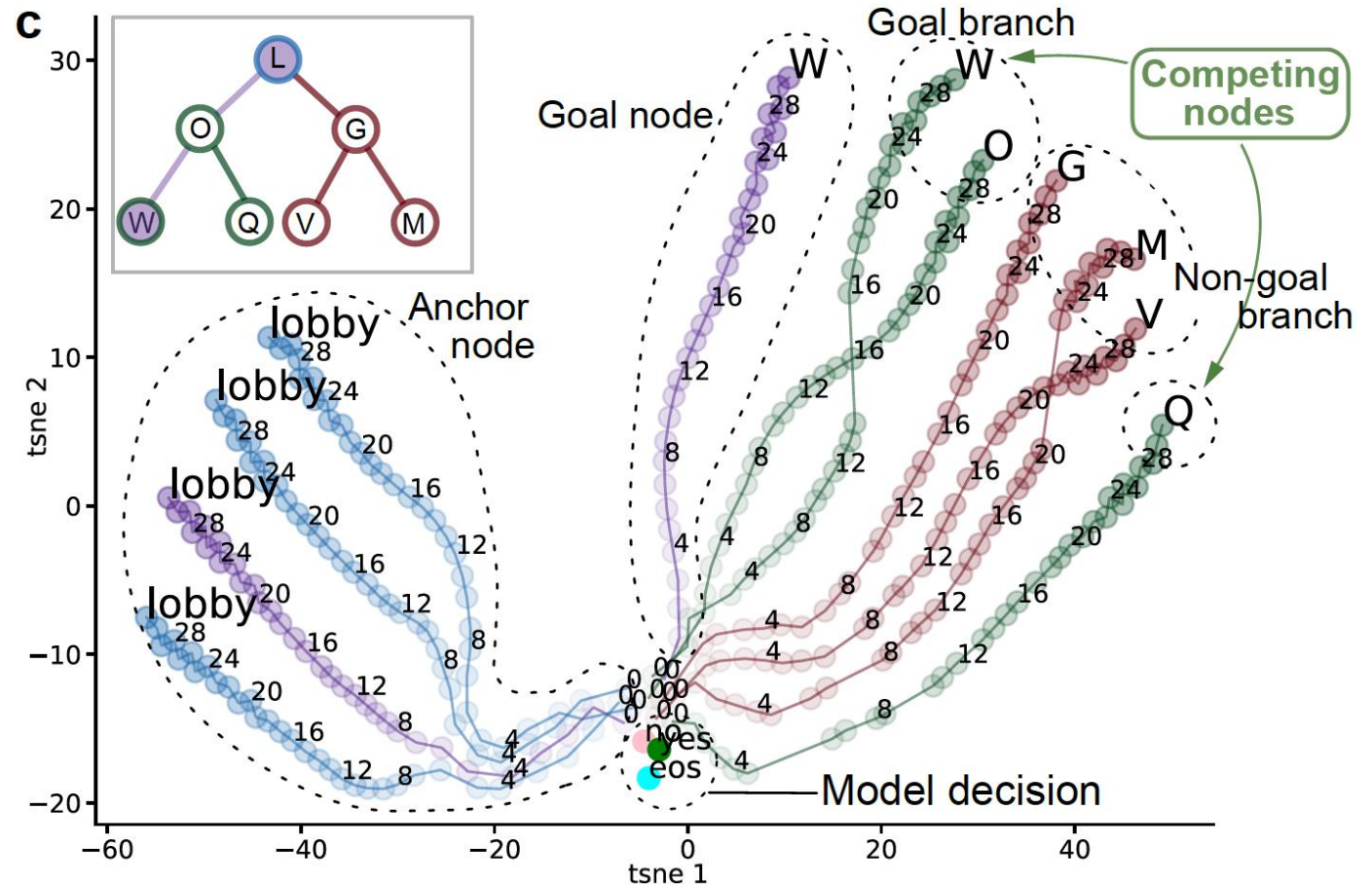
- 2D t-SNE of:
 - Room-token activations from all layers
 - Plus final eos token ("yes"/"no")
- Each color = room token
- Number next to point = layer index



Case study: Graph navigation

Early layers: all room tokens form single tight cluster

Lobby token diverges;
anchor trajectory

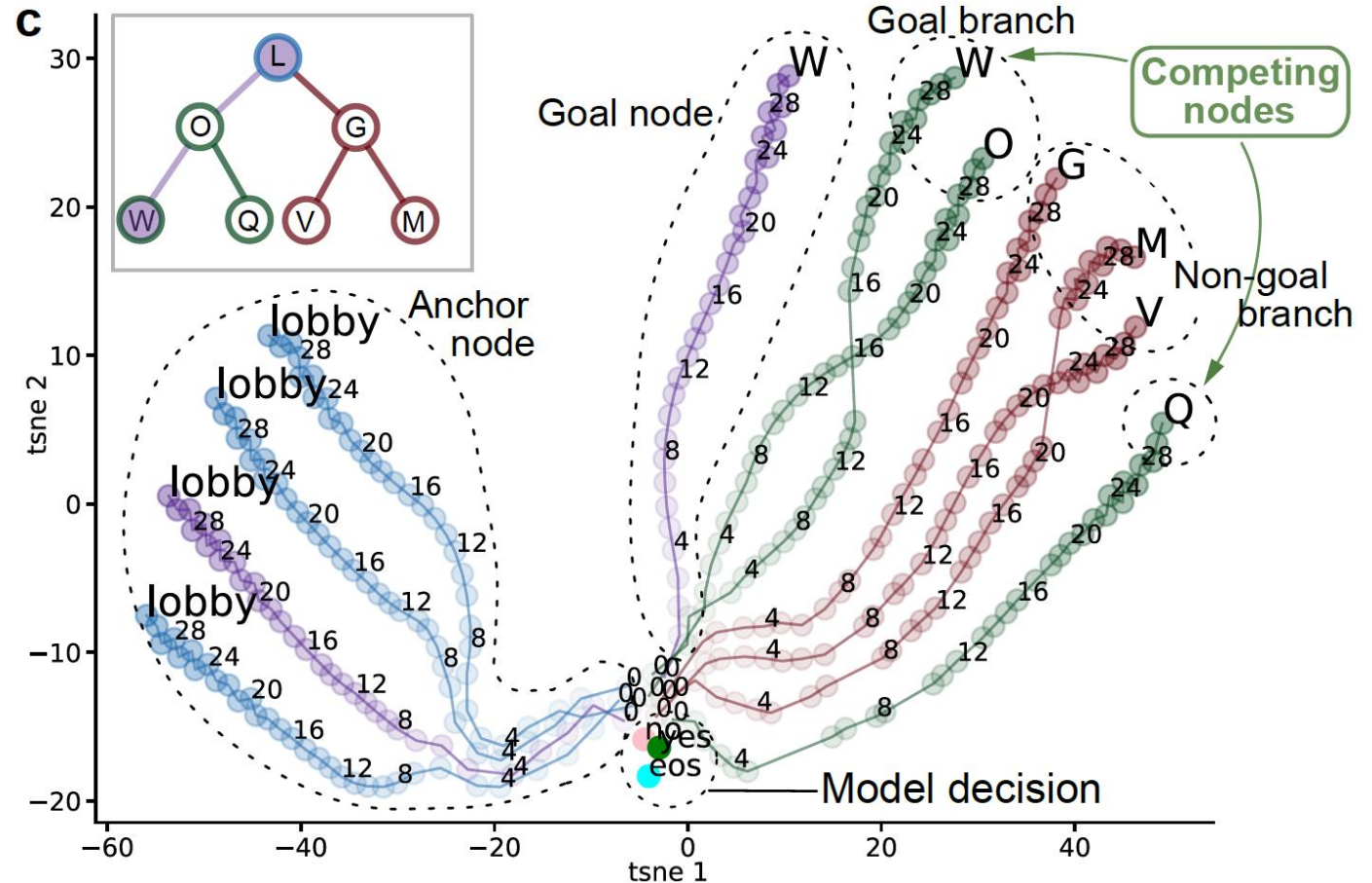


Case study: Graph navigation

Non-goal room tokens cluster together

- Consistent subgroup patterns across layers

W and sibling competitor Q increasingly separate



New directions: inference-time compute

Motivation: reasoning need not occur in one feedforward pass

- Chain-of-thought, explicit tree search, agentic frameworks, ...

Fit for AlgEval:

Sequential outputs easier to analyze than high-dim states

Key research questions:

- Which computations *offloaded* to inference vs. embedded in model?
- Can scaling inference-time compute outperform scaling model size?

New directions: RL + alg reasoning

RL can shape how models discover and store algorithms

- RL may yield emergent algorithmic behaviors beyond imitation

E.g., reasoning models show reasoning emergence via RL

- DeepSeek displays backtracking-like behavior/"aha moments"

Key research question: Does RL teach new algorithms or amplify ones already latent in pretraining data?