
Stepwise Inference in Transformers: Exploring a Synthetic Graph Navigation Task

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 Taking correct steps through elementary logical operations is the essence of log-
2 ical reasoning, culminating in precise planning outcomes. While such *step-*
3 *wise inference* approaches have demonstrated benefits in Large Language Models
4 (LLMs), conducting an accurate quantitative evaluation is challenging, given their
5 extensive scale, complexity, and lack of accessibility. We introduce a minimal syn-
6 thetic setup, where an autoregressive language model solves a navigation task on
7 directed acyclic graphs (DAGs), taking inspiration from computational graphs and
8 execution traces. By implementing training with sample paths from start to goal
9 node in a 'step-by-step' manner, we perform systematic experiments and develop
10 novel analyses illustrating that stepwise navigation proves advantageous when the
11 underlying graph is hierarchical and generalization necessitates the stitching of
12 subpaths observed during pretraining. Further, we observe a diversity-accuracy
13 tradeoff while varying sampling temperature and a bias towards generating shorter
14 paths. We next elucidate how in-context chain-of-thought exemplars can steer the
15 model's navigation. Importantly, these exemplars can guide the model to follow
16 a path of reasoning we provide, instead of relying on its potentially biased pri-
17 ors. Together, this work showcases the utility and adaptability of this paradigm in
18 exploring the complexities of logical reasoning and planning in LLMs.

19 1 Introduction

20 Here, we strive to formulate a framework that is not only simple, controllable, and interpretable but
21 also encapsulates the essential features of an array of stepwise inference tasks such as scratchpad
22 and zero/few-shot chain-of-thought in transformers. Our design philosophy adheres to the principle
23 of constructing a model task that is "*as simple as possible, but not simpler*", ensuring that the
24 model embodies the following set of properties: (i) the task can be better solved by outputting the
25 intermediate steps of computation; (ii) there can be several possible paths of computational steps
26 to solve the task; and (iii) the context of the task can be controlled by providing exemplars in the
27 prompt. Here we argue that the paradigm of graph navigation problems provides such a fundamental
28 framework. Inspired by computational graphs and execution traces, we model stepwise inference
29 as navigating simple paths in a directed acyclic graph (DAG), representing a chain of logic (Dziri
30 et al., 2023). Given a start and goal node, the transformer must autoregressively produce a sequence
31 of nodes that concludes at the goal node. This task requires two levels of computation: locally,
32 each step taken by the model must be valid, and on a global scale, the sequence of steps must be
33 strategically planned in advance to reach the goal node. This setup enables us to modify (1) the
34 structure of the underlying graph (2) the content of the training samples during pre-training and (3)
35 the information provided to the model in-context before cue the model with the goal and start nodes.
36 Consequently, we can systematically examine the impact of these properties on the development
37 of reasoning abilities, an investigation that is challenging to conduct on a large scale. While (1)

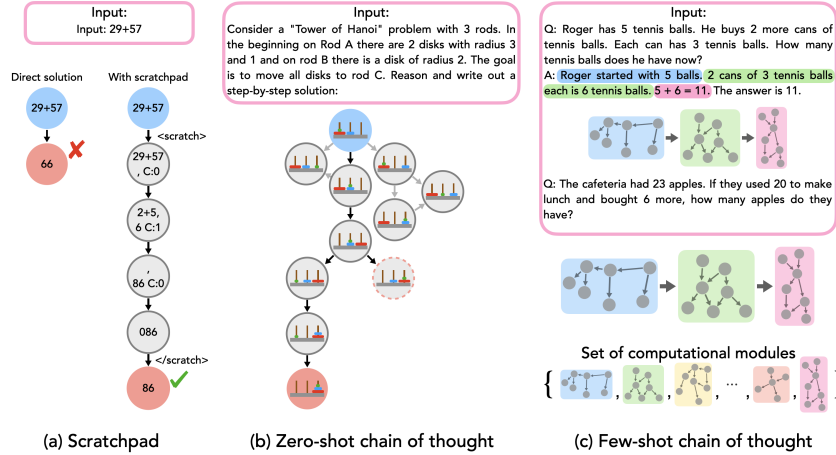


Figure 1: **Graph navigation task as a simple, steerable, and interpretable framework for exploring stepwise inference.** (a) Scratchpad (Nye et al., 2021) improves LLMs’ ability to perform complex multi-step computations, such as arithmetic, when they write intermediate computation steps to a buffer called a scratchpad. (b) Zero-shot chain-of-thought prompting (Kojima et al., 2022) improves LLMs’ ability to perform multi-step reasoning, such as Tower of Hanoi by prompting them to generate detailed reasoning paths. (c) Few-shot chain-of-thought prompting (Wei et al., 2022) improves LLMs’ ability to perform multi-step reasoning, such as solving math word problems (Cobbe et al., 2021), by first presenting an exemplar in-context in the prompt.

38 and (2) together allows us to explore scratchpad and zero-shot step-by-step reasoning, which is
 39 relying on the model’s internalized abilities, (3) also delves into few-shot in-context chain of thought
 40 prompting (Wei et al., 2022), where predictions are made with guiding examples. Specifically, we
 41 examine how in-context exemplars affect the path produced by the model and systematically evaluate
 42 the degree of control we have over that path.

43 2 Defining A Synthetic Graph Navigation Task

44 A DAG $\mathbf{G} = (\mathbf{N}, \mathbf{E})$ is made up of set of nodes $N = \{X_i\}_{i=1}^{|N|}$ and set of directed edges across the
 45 nodes $E = \{(X_i, X_j)\}_{X_i, X_j \in N}$. The edges of a DAG captured by the **adjacency matrix** A where
 46 $A_{ij} = 1$ if $(X_i, X_j) \in E$.
 47 A **directed simple path** is a sequence of distinct nodes of \mathbf{G} which are joined by a sequence of
 48 distinct edges. The first node of a path is referred to as the **start node** and the last node is the **goal**
 49 **node** (Fig. 2). A **cycle** is a sequence of nodes connected by a sequence of edges such that the first
 50 and last node are identical. The key characteristic defining DAGs is that they contain no cycles.
 51 The adjacency matrix of a DAG can always be rearranged to be upper triangular. **Structure of the**
 52 **DAG:** We find that the structure of the underlying DAG can have a large effect on the usefulness of
 53 Step-by-Step Inference. When the DAG is hierarchical, between a start and goal node, nodes in the
 54 intermediate layers must be visited SI Fig 7 whereas when the DAG is random, there is a uniform
 55 probability that 2 nodes are connected and there is no explicit notion of hierarchy. In both these sce-
 56 narios, we can define the notion of **path diversity**: between any 2 path-connected nodes, there can
 57 be several possible paths. We quantify the path diversity in random and hierarchical graphs in SI Fig
 58 7 and describe their construction in the corresponding SI Sec. **Data generating process for single**
 59 **graph scenarios** We focus on two setups in this work, where one allows for context and one does not.
 60 This is intentional so that we can explicitly analyze benefits of stepwise inference in the presence of
 61 extraneous context, which may influence a model’s internalized knowledge, and hence its execution.
 62 **Prompt structure and training data generation** In the single graph setting with underlying DAG
 63 G , each prompt is made from a single simple path. Given a start node X_{start} and goal node X_{goal} ,
 64 the model has to classify whether there exists a path from X_{start} to X_{goal} . We create a pair of tokens
 65 **path** and **no-path**. We constructed two datasets: one that contains stepwise inference and another
 66 that does not. Examples of prompts are provided below. For stepwise inference, the path between the
 67 start node X_4 to the goal node X_6 is represented as **goal: X_4 X_6 X_5 X_7 X_6 path end**. Without

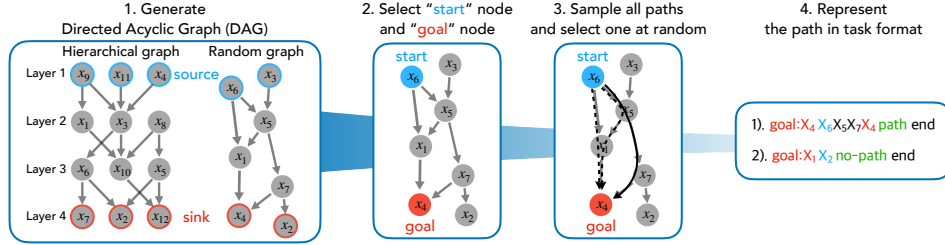


Figure 2: **Data Generating Process for a Single Graph:** This figure illustrates the step-by-step process of generating a training dataset using a single graph. 1) A directed acyclic graph (DAG) is generated, which can be either hierarchically structured or random. 2) A start node and a goal node are selected. 3) All possible paths connecting the start and goal nodes are sampled, and one path is randomly selected. 4) The chosen path is then represented in a task-specific format.

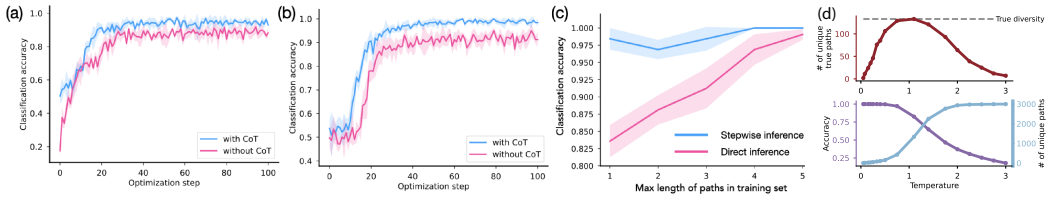


Figure 3: **Advantage of Stepwise Inference in Graph Navigation Tasks:** (a) In random graphs, stepwise inference shows an advantage over direct inference in connectivity prediction tasks. (b) This advantage is further pronounced in hierarchical graphs, where the distances between nodes can be significantly larger. (c) We show that the stepwise inference gap arises when the training set contains paths that are shorter than the paths required to connect nodes in the evaluation set. (d) **A diversity vs. accuracy trade-off in finite temperature stepwise inference for transformers:** As sampling temperature is increased, the diversity of paths generated by the model from a single $(n_{\text{start}}, n_{\text{goal}})$ pair increases, while the accuracy of the path decreases. This tradeoff is captured by measuring the number of unique *true* paths which is non-monotonic (top), showing the existence of an optimal temperature for sampling. The dashed line denotes the ground truth path diversity of $(n_{\text{start}}, n_{\text{goal}})$

68 stepwise inference, the example path is represented as goal: x₄ x₆ path end.

69 **Results on single-graph scenarios:** Our first result is that a transformer trained on directed edges

70 and a small fraction of node pairs from a fixed underlying DAG can generalize to *all* node pairs,

71 including those held out during training, producing valid simple paths from start to goal nodes.

72 Thus the model can ‘stitch’ or mix-and-match (sub)paths it has observed over training to produce

73 a valid path across a pair of held-out connected nodes. The training dynamics of a typical network

74 together with a description of failure modes is shown in SI Fig. 8. **A single underlying graph:**

75 **The stepwise inference gap** Fig. 3 shows the accuracy of path/no-path classification for (a) a ran-

76 dom DAG and (b) a hierarchical DAG. We trained two distinct models using two types of datasets:

77 one with stepwise inference paths and one without. We find that the model trained on the dataset

78 with stepwise inference (represented by the blue line) achieves higher classification accuracy than

79 the model without stepwise inference (the pink line) in both cases. This phenomenon echoes find-

80 ings from large-scale experiments where the inclusion of intermediate reasoning steps results in

81 increased accuracy (Kojima et al., 2022). We refer to the difference in classification performance

82 with and without stepwise inference as the ‘stepwise inference gap’. We also observe that the step-

83 wise inference gap is larger for hierarchical graph than for random graph.

84 **Stitching of paths** We hypothesize that stepwise inference is useful when the training data has the

85 following structure: (1) the underlying DAG is hierarchical, which means that there is an explicit

86 feedforward ordering of nodes and to go from nodes in one layer to next one must pass through all

87 intermediate layers and (2) the model must ‘stitch’ together subsets of paths seen over pretraining in

88 flexible ways to generalize. To test this, we trained the model using paths from hierarchical DAGs

89 while varying the lengths of paths in the training data. Specifically, we created training data that

90 contains start nodes from layer l and goal nodes from layer l' and restricted $l' - l < \Delta$, where Δ

91 denotes the length of the path. During evaluation, we choose node pairs such that $l' - l \geq \Delta$.

92 **The diversity-accuracy tradeoff with higher sampling temperatures** Fig. 3d illustrates the effect

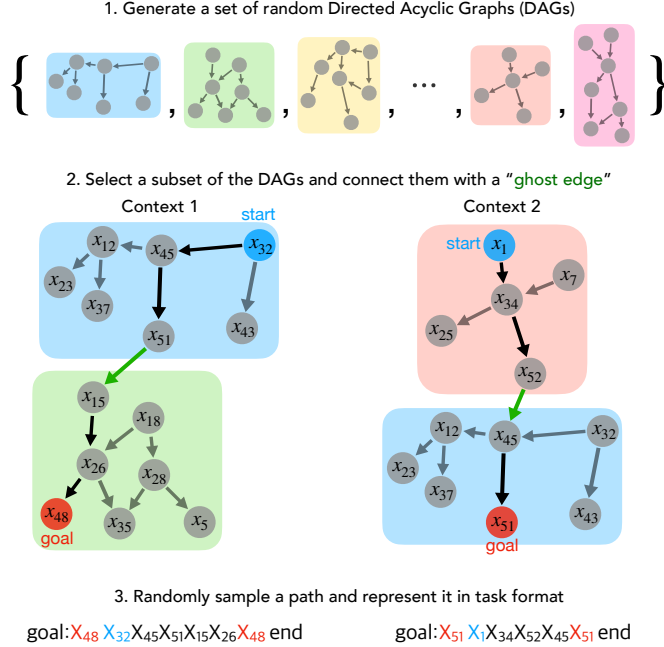


Figure 4: **Data Generating Process for Connected Sub-Graphs (Motifs):** This figure illustrates the step-by-step process of generating a training dataset by combining multiple subgraphs (motifs). 1) We start by making a set of random directed acyclic graphs (DAGs). 2) Next, we pick a subset of these DAGs and connect them together using "ghost edges" to create a bigger graph. 3) From this bigger graph, we randomly sample paths and turn them into a task format.

93 of sampling temperature on the accuracy and diversity of the generated paths. LLM inference at 0
 94 sampling temperature is equivalent to taking the most likely token at each time step (the maximum
 95 likelihood estimate). In this setting, the model deterministically generates the same path for every
 96 provided pair of start and goal nodes: n_{start} and n_{goal} . However, in the underlying graph, there exists
 97 a diversity of paths from each n_{start} to n_{goal} .

98 To capture this diversity, we fixed the start node n_{start} and the goal node n_{goal} , and prompted the
 99 model 3,000 times, sweeping through different sampling temperatures in Fig. 3d.

100 Results on multi-graph scenarios

101 The single graph setting let us explore *zero-shot* planning and stepwise reasoning, where the model
 102 relied purely on knowledge internalized over pretraining to do stepwise reasoning. To study how
 103 context can influence the path the model traverses, we introduce the concept of motifs and in-context
 104 exemplar paths.

105 **Data generating process for the multi-graph scenario** To model few-exemplar based chain-of-
 106 thought prompting, we modify our single graph setup to include a set of subgraphs that we refer to
 107 as **motifs**, denoted by $\{G_i\}_{i=1}^n$. A motif is a DAG that is fixed across pretraining and inference. (i)
 108 **Ghost edge:** For a pair of connected motifs $G_i \mapsto G_j$, an edge between a sink node of G_i and a
 109 source node of G_j , and (ii) A **primitive sequence** (Fig. 4) is sequence of nodes across 2 motifs G_i
 110 and G_j with a start node in G_i and goal node in G_j and this sequence contains exactly 1 ghost edge.
 111 Fig. 4.

112 In chain-of-thought prompting (Wei et al., 2022), one or more examples of reasoning are provided
 113 before asking the next question, as illustrated in Figure 1(c). The LLM then generate a chain-of-
 114 thought which matches that of the exemplar. To model this, we chain a subset of k motifs $G_{c_1} \rightarrow$
 115 $G_{c_2} \rightarrow \dots \rightarrow G_{c_k}$ together and provide exemplars Each exemplar e is a primitive sequence across
 116 each pair of consecutive motifs: $e \in (G_{c_i} \rightarrow G_{c_{i+1}})$ which contains exactly 1 ghost edge. The
 117 construction of a primitive sequence is described in Fig. 4 and examples are shown in Fig. 5(b).
 118 Given a start node $n_{\text{start}} \in G_{c_1}$ and a goal node $n_{\text{goal}} \in G_{c_k}$, the model can be prompted either
 119 directly (Fig. 5(a)) or provided with exemplars and then queried for a path from n_{start} to n_{goal} (Fig.
 120 5(b)). We find that the model can successfully follow the chain defined by the in-context exemplars.

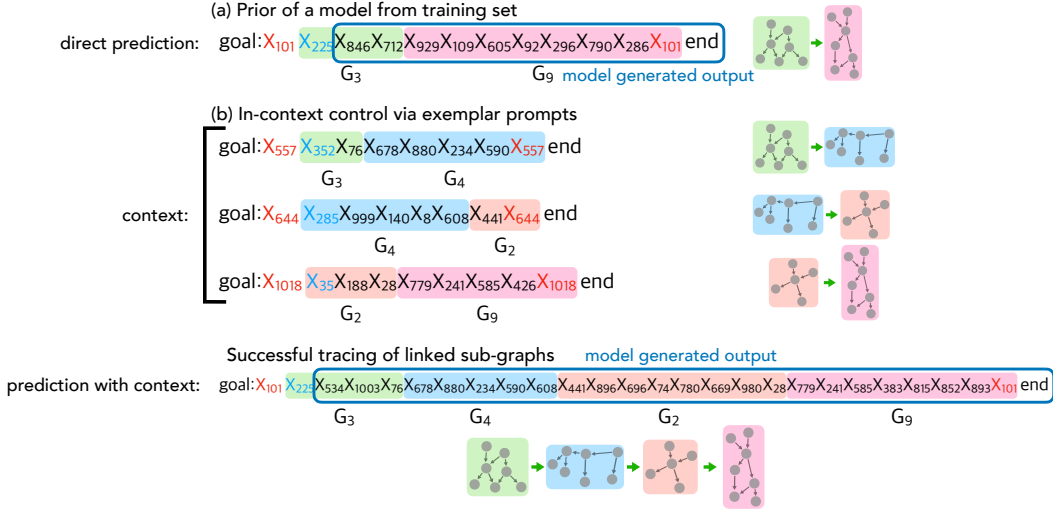


Figure 5: **Example output sequences from the model highlighting the steerability of stepwise inference.** (a) Direction prediction: Given $n_{\text{start}} \in G_3$ and $n_{\text{goal}} \in G_9$, the model produces a path from $G_3 \rightarrow G_9$, placing a single ghost edge (X_{712}, X_{929}). (b) With in-context exemplars: primitive sequences from $G_3 \rightarrow G_4$, $G_4 \rightarrow G_2$ and $G_2 \rightarrow G_9$ in-context make the model steer its navigation through the path stringing together these motifs in order: $G_3 \rightarrow G_4 \rightarrow G_2 \rightarrow G_9$, placing a ghost edge between every consecutive motif, for a total 3 ghost edges.

121 An example output produced by the model is in SI Fig.5, highlighting the path the model takes
 122 through the chain of motifs $G_3 \rightarrow G_4 \rightarrow G_2 \rightarrow G_9$. We also find that the model generalizes to
 123 arbitrary orders of motifs strung out, including those that did not occur consecutively in the train set
 124 – in other words, in-context control is capable of *compositional generalization* (Li et al., 2023).
 125 **How do the exemplars affect controllability of graph navigation?**

126 Next, we study how the structural content of the exemplars affects the navigation path chosen by
 127 the model. We hope to shed some light on and create hypotheses for the vast and varied findings
 128 about stepwise reasoning in LLMs at scale.
 129

130 **Number of intermediate motifs** In
 131 Fig.6(a), we varied the number of exam-
 132 plars provided to the model. This is equiv-
 133 alent to stringing together a longer chain of
 134 motifs to navigate over. We find that the
 135 model can generalize well to unseen orders
 136 of motif up to the maximum number
 137 chained together in the training data. This
 138 creates a hypothesis for chain-of-thought
 139 and related methods at scale: the model
 140 will fail to generalize to reasoning chains
 141 longer than those present in its training
 142 data.

143 **Bias towards the first exemplar in the**
 144 **case of conflict** Multiple examples of a
 145 context provided in the prompt can in-
 146 crease the precision of our control over the
 147 model, but it can also lead to confusion. Here, we systematically and quantitatively study the behav-
 148 ior of the model when two contexts are provided but are in conflict. In Fig.6(b) we study a scenario
 149 where two chains of motifs are provided, starting from the same set of primary motifs and ending at
 150 the terminal motif. We find that the model has a strong bias toward choosing the first chain over the
 151 second.

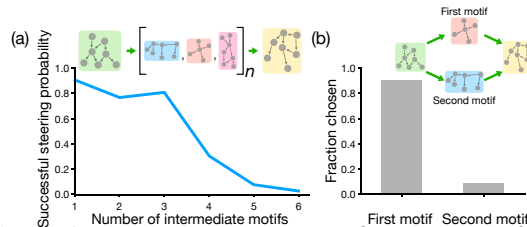


Figure 6: **How do the number of exemplars affect the controllability of motifs?** (a) As we vary the number of intermediate motifs in a chain, the path generated by the model follows the path described by the chain until $n = 4$, which is the extent of the training data. (b) In the case of 2 conflicting chains in-context, the model has a bias to pick the first chain.

152 **References**

- 153 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
154 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to
155 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 156 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jian, Bill Yuchen Lin, Peter West,
157 Chandra Bhagavatula, Ronan Le Bras, Jena D Hwang, et al. Faith and fate: Limits of transformers
158 on compositionality. *arXiv preprint arXiv:2305.18654*, 2023.
- 159 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
160 language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022.
- 161 Yingcong Li, Kartik Sreenivasan, Angeliki Giannou, Dimitris Papailiopoulos, and Samet Oymak.
162 Dissecting chain-of-thought: A study on compositional in-context learning of mlps. *arXiv
163 preprint arXiv:2305.18869*, 2023.
- 164 Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin,
165 David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show
166 your work: Scratchpads for intermediate computation with language models. *arXiv preprint
167 arXiv:2112.00114*, 2021.
- 168 Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis
169 of chain-of-thought. In *The Eleventh International Conference on Learning Representations*,
170 2023. URL <https://openreview.net/forum?id=qFVVBzXxR2V>.
- 171 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny
172 Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint
173 arXiv:2201.11903*, 2022.

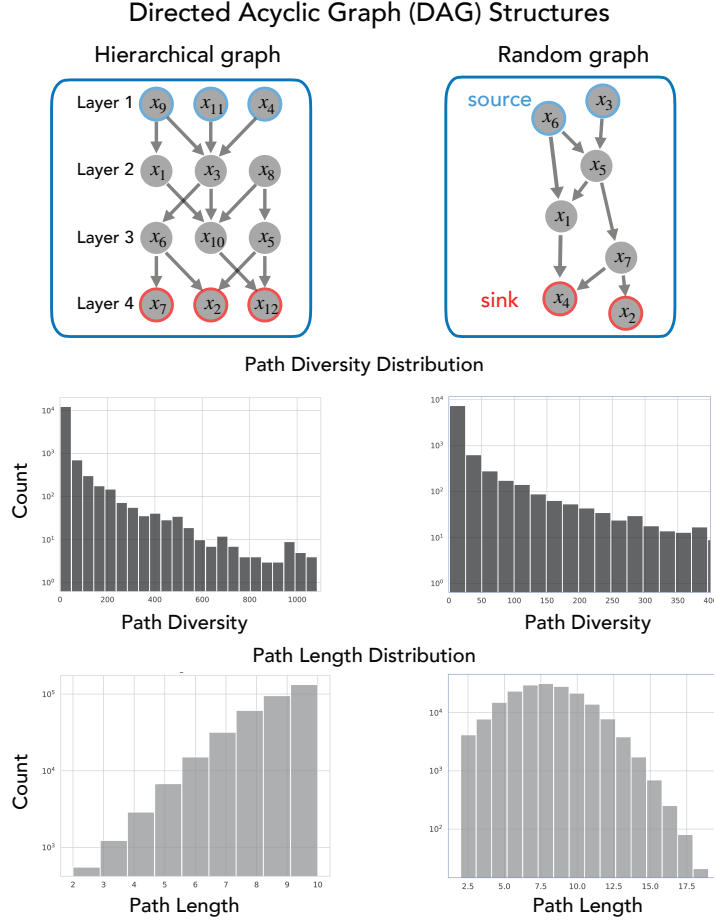


Figure 7: Construction and properties of Hierarchical and Random DAGs.

174 **A Appendix**

175 **A.1 Setup and construction of graph and model**

176 Here we describe the properties of the DAGs we use, the training setup, model architecture and
 177 hyperparameters.

178 We use 2 DAG structures, hierarchical and random (Fig. 7). Random DAGs are constructed by
 179 randomly generating an upper-triangular matrix where each entry has probability p of existing. Hi-
 180 erarchical DAGs are generated by predefining L sets of nodes and drawing an edge between a node
 181 n_l in layer l and n_{l+1} in layer $l + 1$ with probability p . Lastly, we ensure that the graph is con-
 182 nected. These lead to different path diversity and path length distributions, which affect the efficacy
 183 of stepwise inference, as shown in our results.

184 To create a feedforward hierarchical DAG we construct a set of L layers with N nodes each. For
 185 every node n_l in layer l and n_{l+1} in layer $l + 1$, we draw a directed edge (n_l, n_{l+1}) with probability
 186 p , which we refer to as **edge density**. Thus on average, between any two layers there are pN^2 edges
 187 and each node in an intermediate layer has an out-degree and in-degree of pN . The number of
 188 paths from a particular node in layer l to layer $l' > l$ is exponential and given by $(pN)^{l'-l}$ - this is
 189 quantified in the path length distribution shown in SI Fig 7. Lastly, we make sure that the graph is
 190 connected and there no disjoint disconnected components. The nodes from layer 1 are the **source**
 191 **nodes**: nodes $\{X_i\}$ of DAG \mathbf{G} with $\text{parents}(X_i) = \emptyset$ and the nodes from layer L are **sink nodes**:
 192 nodes $\{X_i\}$ of DAG \mathbf{G} with $\text{children}(X_i) = \emptyset$.

193 To create a random DAG of N nodes and create a random upper triangular adjacency matrix $A_{N \times N}$

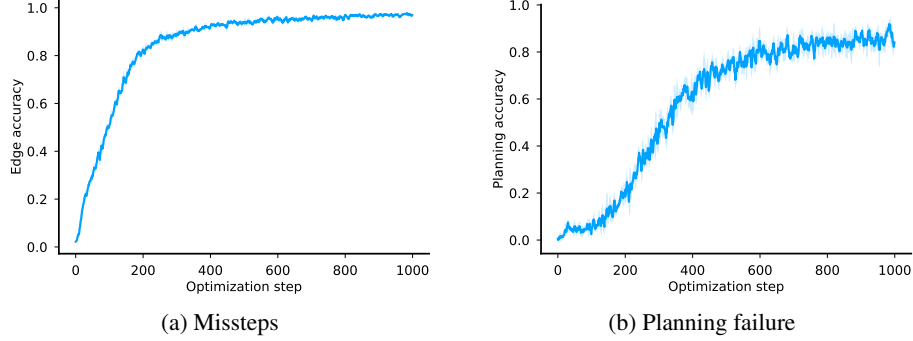


Figure 8: **The evolution of failure mode probabilities over training:** It can be seen that the model first learns to produce correct edges (effectively bigram statistics) and then learns the global objective of producing a path that ends at the cued goal node. Accuracy curves averaged over 3 trained models with different random seed.

194 with bernoulli entries with edge density p , such that $p(A_{ij} = 1) = p$. We also ensure that the graph
 195 is connected. This results in a bell-shaped path length distribution, SI Fig.7.

196 A.2 Training dynamics in the single graph scenario

197 Here we show the training dynamics of a single graph model.

198 **Failure modes of step-by-step inference** Given underlying DAG G , during step-by-step infer-
 199 ence, the model produces a sequence of nodes from the start node n_{start} which has to terminate at the
 200 goal node n_{goal} , given by the sequence $n_0 = n_{\text{start}} \rightarrow n_1 \rightarrow n_2 \rightarrow \dots \rightarrow n_k \rightarrow \dots \rightarrow n_T$. Here in
 201 our setup, there are two broad categories of failures possible (Saparov & He, 2023):

- 202 1. **Misstep:** $(n_k, n_{k+1}) \notin G$. An edge produced by the model does not exist in the DAG.
- 203 2. **Planning failure:** $n_T \neq n_{\text{goal}}$. The model produces a path that does not terminate at the goal.

204 We choose to highlight the two types of failures identified above: (1) the probability of taking a
 205 correct step (i.e. $1 - Pr(\text{misstep})$) and (2) the probability of ending at the cued target node (i.e.
 206 $1 - Pr(\text{planning failure})$). These are shown in Fig. 8

207 A.3 Model architecture and loss function

208 For training, we tokenize every node and we use next-token prediction with a cross entropy loss:

$$\mathcal{L}(\mathbf{x}_n, \text{target } n) = -\log \left(\frac{\exp(\beta x_{n, \text{target } n})}{\sum_{t=0}^{\#\text{tokens}} \exp(\beta x_{n,t})} \right) = -\log \left(\underbrace{\text{softmax}(\beta \mathbf{x}_n)_{\text{target } n}}_{\text{prob}(\text{target } n)} \right) \quad (1)$$

Hyperparameter	Value
learning rate	10^{-4}
Batch size	64
Context length	32
Optimizer	Adam
Momentum	0.9, 0.99
Activation function	GeLU
Number of blocks	2
Embedding dimension	64

Table 1: Hyperparameters of the transformer

209 For model architecture, we use a GPT based decode-only transformer with a causal self-attention
 210 mask. Our implementation is based on the popular nanoGPT repository¹.

¹available at <https://github.com/karpathy/nanoGPT>

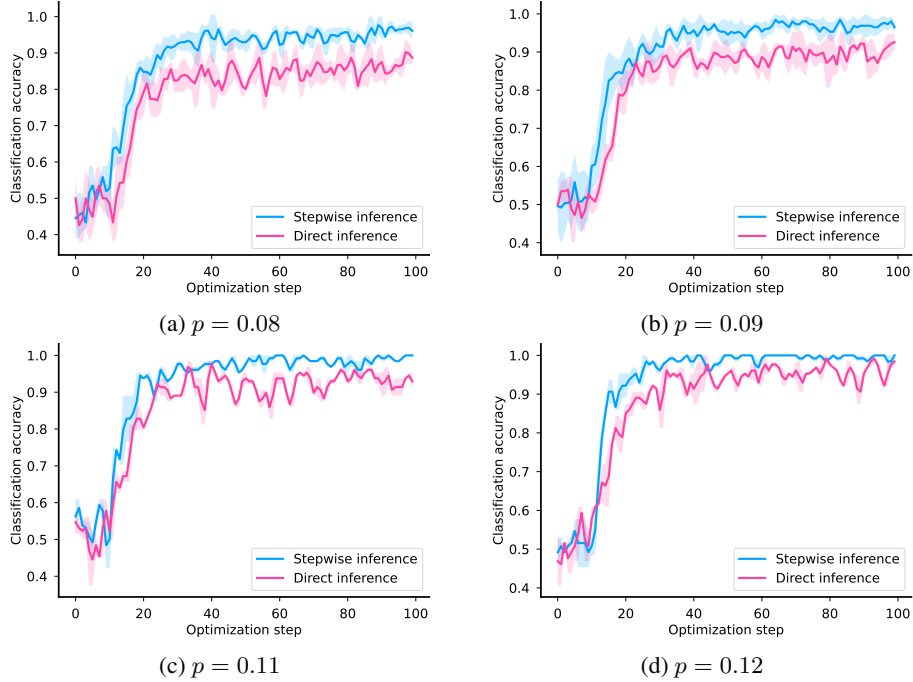


Figure 10: Advantage of Stepwise Inference in Graph Navigation Task.

211 **A.4 A bias towards shorter paths**

212 Fig. 9 examines the average path lengths in a random
 213 graph, comparing true paths to those generated by our
 214 trained model. Notably, the model consistently produces
 215 paths that, on average, are shorter than the actual paths
 216 in the random graph. This observation suggests that
 217 the model has a bias towards efficiency, which can lead
 218 to oversimplification of complex stepwise inference or
 219 omission of important intermediate steps.

220 Fig. 4 For the training data construction in the multigraph
 221 setting.

222 **A.5 Additional**
 223 **experimental results for the single graph setting**

224 In Fig. 10, we swept the density of the graph from 0.08
 225 to 0.12 on a hierarchical graph. We observe a stepwise
 226 inference gap in all cases. The stepwise inference gap
 227 becomes smaller for larger densities.

228 Fig. 11 presents a density plot comparing the average
 229 lengths of actual paths with those generated by the model
 230 in a random graph. This observation verifies the model
 231 tends to produce shorter paths between a given pair of
 232 start and goal nodes.

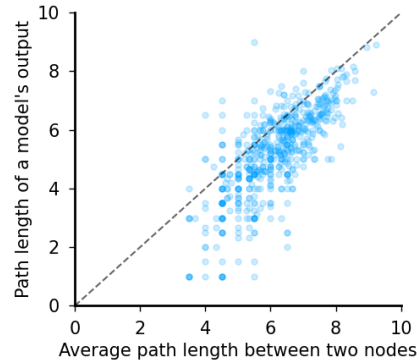


Figure 9: **Model outputs are biased toward shorter paths.** We compare the average lengths of actual and model-generated paths in a random graph, revealing the trained model's bias to generate shorter paths connecting a pair of start and goal nodes in a random graph.

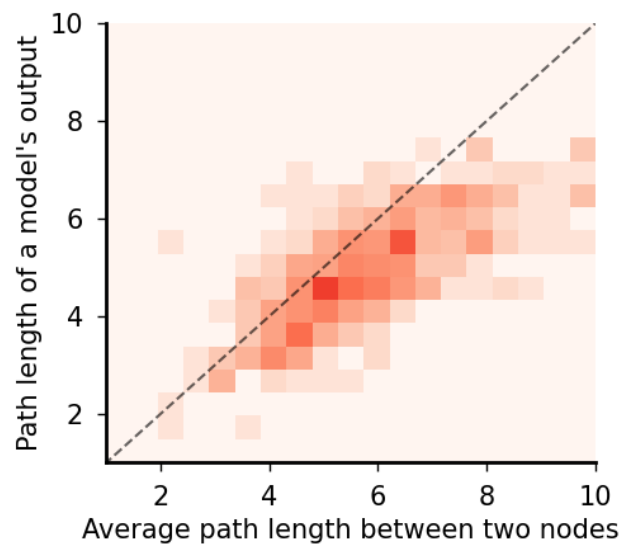


Figure 11: **Model outputs are biased toward shorter paths.**