

Let's Decrypt Dot by Dot: Decoding Hidden Computation in Transformer Language Models

Aryasomayajula Ram Bharadwaj
Independent Researcher
ram.bharadwaj.arya@gmail.com

August 23, 2024

Abstract

Transformer models can perform complex reasoning with Chain-of-Thought (COT) prompting. COT can be replaced with hidden characters while maintaining performance. This paper investigates methods to decode these hidden computations, focusing on the 3SUM task using a 34M parameter LLaMA model. We propose a novel decoding method to recover original COT, providing insights into how transformers encode and process hidden COT information. Our work offers new perspectives on model interpretability and computation in language models.

1 Introduction

COT prompting improves performance of large language models. Recent work shows improvements persist when COT is replaced with hidden characters, raising questions about the nature of computation in these models. This paper builds on findings by Pfau et al., focusing on the 3SUM task as a case study. We aim to decode hidden computations in the transformer architecture, with potential for improved model interpretability and training strategies.

2 Background

2.1 The 3SUM Task

The 3SUM task involves finding three numbers in a set that sum to zero. It serves as a proxy for more complex reasoning tasks and is used to study computational capabilities of transformer models. An example of a 3SUM sequence is provided.

2.2 Hidden Chain-of-Thought

In hidden Chain-of-Thought, intermediate reasoning steps are replaced with hidden characters (e.g., "..."). Models trained on hidden sequences still perform well, suggesting meaningful computation occurs despite lack of explicit reasoning steps.

3 Methodology

We used a 34M-parameter Llama model with 4 layers, 384 hidden dimension, and 6 attention heads. Our analysis focused on three main areas: Layer-wise Representation Analysis, Token Ranking, and Modified Greedy Decoding Algorithm.

4 Results and Discussion

4.1 Layer-wise Analysis

We observed progressive transformation of representations across layers. Initial layers contained pure number sequences related to 3SUM’s COT, with hidden tokens appearing from the third layer onward. We observed a gradual replacement of number sequences with hidden characters.

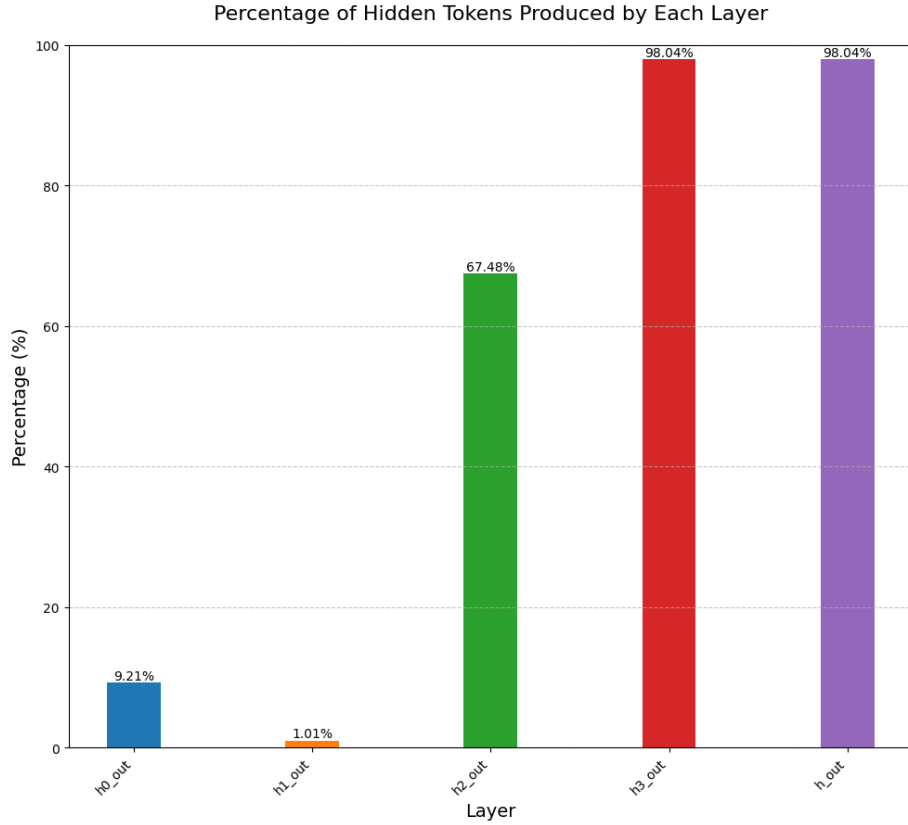


Figure 1: Hidden token percentages across layers

4.2 Token Rank Analysis

The top-ranked token was consistently the hidden character ("."), while lower-ranked tokens revealed original, non-hidden COT sequences. This supports the

hypothesis that the model replaces all computation with hidden tokens on top while keeping the original computation intact underneath.

4.3 Modified Greedy Decoding Algorithm

We implemented a modified greedy autoregressive decoding method. Steps include: perform standard greedy decoding, select second-highest probability token when encountering hidden token, and continue for entire sequence. This resulted in 100% match in 3SUM task results with and without hidden tokens. Random token replacement was less effective than choosing next highest ranked token.

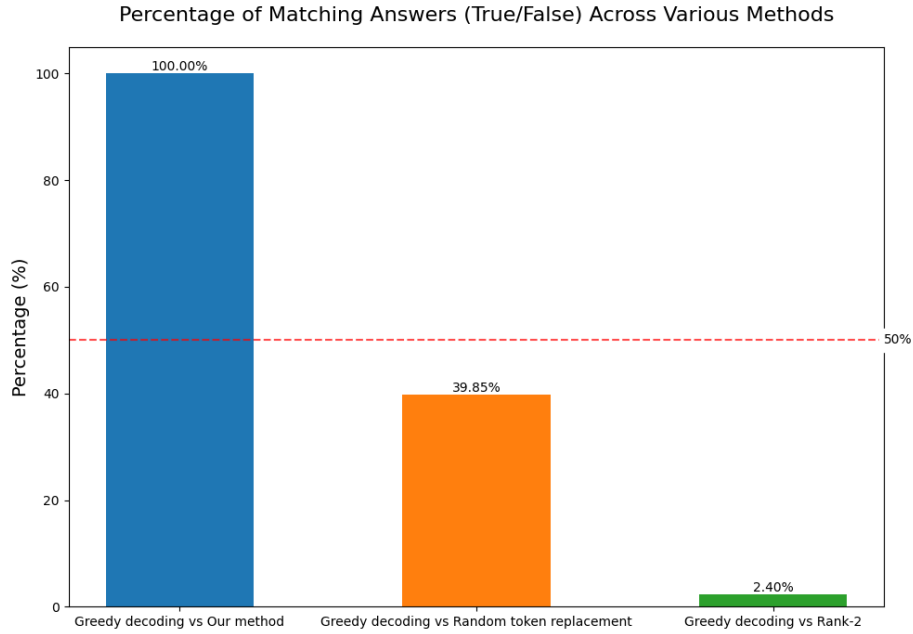


Figure 2: Comparison of decoding methods

5 Implications and Future Work

Our findings provide new tools for understanding internal reasoning processes and increase confidence in COT-based approaches for improving visibility. Future work should focus on developing better decoding methods or finding circuits that hide tokens, investigating generalizability to tasks beyond 3SUM (including natural language tasks), and improving token hiding methods (currently limited to one hidden token which is simple to decode).

6 Conclusion

We have presented a novel approach to understanding hidden computations in transformer models through analysis of token rankings and layer-wise representations, and development of a modified decoding algorithm. Our insights into how models encode and process information in hidden COT sequences open new avenues for improving interpretability, efficiency, and safety in language models. This work strengthens the belief in COT visibility research for interpretability.

[illegible]

h0 out: 0-0 2 [EOS] 0-2 A 5 [A 0-8 A 3 [EOS] 0-0 2 4 4 0-6 5 0-2 5 3 2 0-8 2 A 9 1 9 0-5 5 0-8 6 0-7 3 0-9 0-8
2 0-6 0-6 9 [EOS] 0-0 2 0-0 0-5 0 [EOS] [EOS] [EOS] 2-0 6 2 0-6 0-7 5 0-6 0 A 5 [EOS] 1 0-6 2 0-4 5 0-7 . [EOS] 5 0-6 A
A 7 0-6 0-6 0-6 0-6 5 0-6 2 0-6 0-5 0-6 5 0-7 1 [EOS] 4 0-6 0-6 2 0-6
h1 out: 0-0 9 0-2 9 0-1 0-2 1 0-2 6 0-2 1 0-2 5 0-8 0-2 7 0-2 4 0-2 6 0-2 9 0-2 4 0-6 0-8 0-4 5 0-2 0-9 6 0-9 0-6
4 0-7 6 0-7 2 0-6 6 0-2 1 0-6 7 0-2 8 0-7 8 0-7 4 0-6 3 0-9 3 0-7 9 0-7 0 0-7 8 0-6 0 0-9 3 0-7 8 0-9 5 0-7 0 0-9
2 0-6 4 0-7 6 0-3 0-6 1 0-7 9 0-7 9 0-6 8 0-9 7 0-7 0-7
h2 out: 0-1 9 0-2 9 0-3 1 0-4 2 0-5 6 0-6 0-7 9 0-8 8 0-7 0-2 7 0-3 6 0-4 2 0-5 4 0-1 8 0-7 2 0-8 5 0-9 1 0-3
6 0-2 0-6 0-5 1 [EOS] 5 0-7 1 0-2 7 0-2 8 A 8 0-5 2 0-3 7 0-7 3 0-9 3 0-3 4 0-4 0 0-4 2 0-7 3 [EOS] 6 0-4 0 0-5 2 0-7
7 2 0-8 5 0-5 5 [EOS] 3 0-8 1 0-6 0 [EOS] 7 0-7 5 0-8 7 0-3 4 0-4
h3 out: 0-0 9 0-0 0-0 1 0-4 2 0-6 1 0-0 1 0-0 5 0-8 8 0-0 7 0-1 0-0 1 0-6 0-1 0-5 0-5 4 0-1 2 0-1 2 0-1 3 0-9 1 0-2
2 0-6 4 0-6 0-2 2 0-2 1 0-2 1 0-2 7 0-2 8 0-3 8 0-5 2 0-6 8 0-7 3 0-3 9 0-9 0-4 8 0-6 0 0-7 8 0-8 7 0-9 5 0-5 0-9
7 2 0-8 4 0-9 0-6 0-6 3 0-8 6 0-9 1 0-8 9 0-9 4 0-9 7 3 0-7 [EOS]
h out: 0-0 9 0-0 0-0 0-0 1 0-4 2 0-0 1 0-0 1 0-0 5 0-8 8 0-0 7 0-1 0 0-1 0 0-1 0-5 0-5 4 0-1 2 0-1 2 0-1 3 0-9 1 0-2
6 0-4 0-6 0-2 2 0-2 1 0-2 1 0-2 7 0-2 8 0-3 8 0-5 2 0-6 8 0-7 3 0-3 9 0-9 0-9 0-4 8 0-6 0 0-7 8 0-8 7 0-9 5 0-5 0-9 0-7
2 0-8 4 0-9 0-0 0-6 3 0-8 6 0-9 1 0-8 9 0-9 4 0-9 7 0-3 7 [EOS]

```

A0 True: 0-0 5 [EOS] [EOS] [EOS] [EOS] [EOS] 6 A [EOS] [EOS] 0 A [EOS] 0-9 5 [EOS] [EOS] 3 0-7 3 0-8 2 [EOS] 0-8
A 7 True True 0-6 [EOS] 0-6 0-8 0-6 2 [EOS] [EOS] [EOS] 2 0-6 9 0-6 5 0-6 0-8 0-6 2 A [EOS] [EOS] 2 0-6 6 A 6 [E
S] A [EOS] [EOS] [EOS] [EOS] 5 [EOS] 0 [EOS] 2 [EOS] True [EOS] [EOS] A 0-8 A 5 [EOS] [EOS] [EOS] [EOS] A 5 A 1 A [E
S] A 6 0-7 [EOS] [EOS] [EOS] [EOS] [EOS] [EOS] [EOS] True
h1 out: 0-0 2 0-0 3 [EOS] 5 0-0 2 0-0 7 0-0 0 0-0 0-0 0-0 0-3 0-0 0-0 0-7 2 0-7 4 [EOS] 8 0-2 8 0-2 0 0-6 5 0-8 8
0-2 0-0 2 1 0-2 0-0 6 1-0 7 4 0-6 0-0 7 2 0-6 3 0-6 2 0-7 8 0-7 6 0-6 2 0-6 4 0-6 9 0-7 2 0-7 2 0-9 6 0-7 0 0-7 9
0-7 8 0-6 5 0-9 5 0-6 2 0-6 8 0-6 0-0 6-0 6-0 5-7 0-7 4 0-9 4 0-6 False
h2 out: 0-1 1 0-0 3 0-0 5 0-0 1 0-0 1 0-0 1 0-0 5 0-0 5 0-9 3 0-1 1 0-0 1 2 0-1 0-9 0-1 8 0-6 9 0-1 0 0-1 3 0-9 8 0-
2 0-0 4 1 0-2 1 0-6 1 0-2 4 0-8 0-0 9-2 0-4 3 0-3 3 0-6 8 0-3 0-6 8 2 0-9 0 A 8 0-6 0-0 4 2 0-4 8 0-9 1 0-6 0 0-5
8 0-5 4 0-9 0-0 6 2 0-6 8 0-9 1 A 6 0-9 0-8 0-9 4 0-9 4 A
h3 out: 0-0 1 0-2 3 0-3 5 0-6 0-6 5 0-6 6 2 0-7 9 0-0 5 0-9 3 0-2 7 0-1 9 0-4 0-9 0-1 8 0-6 8 0-7 4 0-8 6 0-1 8 0-
3 0-0 2 1 0-5 1 0-6 5 0-7 4 0-8 0-0 9-2 0-4 3 0-3 3 0-7 3 0-6 3 0-8 0-9 4 0-5 0-6 2 0-4 2 0-4 7 0-4 0 0-5 9 0-5
5 0-8 0-5 5 0-5 0-7 2 0-6 6 0-6 9 0-7 6 0-7 8 0-8 4 0-9 4 A True
h4 out: 0-0 1 0-2 3 0-3 5 0-6 0-6 5 0-6 6 2 0-7 9 0-0 5 0-9 3 0-2 7 0-1 9 0-4 0-9 0-1 8 0-6 8 0-7 4 0-8 6 0-1 8 0-3
0-0 2 1 0-5 1 0-6 5 0-7 4 0-8 0-0 9-2 0-4 3 0-3 3 0-7 3 0-6 3 0-8 2 0-9 4 0-5 0-6 2 0-4 2 0-4 7 0-4 0 0-5 9 0-5
8 0-5 5 0-5 0-7 2 0-6 6 0-6 9 0-7 6 0-7 8 0-8 4 0-9 4 A True

```

```

n0_out: 0-0 5 [EOS] [EOS] A [EOS] 2 A [EOS] A 3 A 6-0-9-4-4-7 A 5 0-2-2 A 5 True 2 A 5 1 5 0-5 [EOS] 0-8 6 0-6 0
-5 0-9 0-8 [EOS] 2 0-6 5 0-6 2 0-6 0-8 0-2 [EOS] [EOS] A [EOS] 2 A 2 A 6 [EOS] 5 A [EOS] 5 [EOS] 0 [EOS] 2
[EOS] 0-8 0-7 5 0-6 0-8 0-2 5 [EOS] 6 0 [EOS] A 5 A 5 0-7-1 [EOS] True [EOS] 4 0-6
h1_out: 0-0 0-2 0-9 0-0 1 0-0 2 0-2 7 0-2 0-0 0-5 0-5 0-0 3 0-2 0 0-7 6 0-2 9 0-2 4 0-2 9 0-2 0-2 6 0-9 0-0 0-2
0-6 2 1 0-2 2 0-6 1 0-2 1 0-6 0 0-6 2 0-6 3 0-6 4 0-6 3 0-7 0-6 2 0-6 4 0-7 9 0-7 0-6 9 0-3 0-7 6 0-9 5 0-7 9 0-9
2 0-6 5 0-6 0-6 2 0-6 1 0-7 0 0-6 0-6 9 0-8 6 0-7 0-7 0-6
h2_out: 0-1 1 0-2 0-9 3 0-1 0-2 0-0 1 0-0 1 0-0 5 0-8 8 0-3 0-2 0-0 0-3 2 0-1 0-9 0-5 8 0-6 9 0-1 0-0 1-5 0-1 8 0-
3 6 0-4 6 0-2 1 [EOS] 5 0-7 1 0-8 7 0-9 8 3 0-3 2 0-6 8 0-3 3 0-8 2 0-3 0-6 0-4 8 0-6 0-0 0-4 3 1 7 0-9 0 0-5 2 0-7
2 0-5 4 0-5 0-3 6 0-6 1 0-6 1 0-6 1 0-7 5 0-8 4 0-9 0-4
h3_out: 0-0 9 0-9 0-0 0-5 0-4 0-4 0-0 7 0-6 0-2 0-7 0-8 8 0-9 7 0-2 7 0-3 9 0-1 0-9 0-5 4 0-1 8 0-1 0-4 1-6 0-1 1 0-
2 0-0 2 1 0-5 0-6 5 0-2 4 0-8 0-9 9 8 0-3 8 0-3 2 0-3 3 0-8 9 0-3 4 0-5 8 0-4 0-8 0-8 8 0-7 0-4 0-0 6 0-9 0-
7 8 0-8 0-5 0-7 3 0-6 6 0-9 0-8 6 0-7 8 0-8 0-7 A
h_out: 0-0 1 0-2 9 0-0 1 0-4 2 0-5 7 0-0 1 0-7 5 0-8 5 0-9 7 0-2 0 0-3 9 0-1 0-9 0-1 4 0-1 8 0-1 4 0-8 3 0-9 8 0-3
6 0-2 6 0-5 1 0-2 1 0-7 1 0-2 7 0-9 8 0-3 8 0-5 2 0-6 8 0-7 6 0-8 2 0-9 4 0-5 0-0 6 0-2 0-7 2 0-4 7 0-9 0 0-5 2 0-5
8 0-8 0-4 0-0 0-7 3 0-8 1 0-6 0 0-8 9 0-9 8 0-7 0-9 7 [EOS]

```

5