# Paper Critique: Chain-of-Thought Prompting Elicits Reasoning in Large Language Models

Dan Peng

Department of Computer Science

Oct 12, 2024

danpeng@unc.edu

## 1. Research Problem

### 1.1. What research problem does the paper address?

This paper tackles the challenge of enabling large language models (LLMs) to perform complex reasoning tasks without fine-tuning. The authors introduce "chain-of-thought prompting" - a technique that coaxes models into breaking down multi-step problems by including examples of step-by-step reasoning in few-shot prompts.

### 1.2. What is the motivation of the research work?

The spark behind this work comes from a frustrating paradox: despite their impressive capabilities, even enormous language models stumble when faced with tasks requiring multi-step reasoning. The authors observed that standard prompting methods hit a ceiling, creating flat scaling curves where bigger models don't necessarily perform better on reasoning tasks. This limitation has been an Achilles' heel for LLMs across arithmetic, commonsense, and symbolic reasoning domains.

## 2. Technical Novelty

### 2.1. What are the key technical challenges identified by the authors?

The authors wrestle with several thorny challenges. First, how to unlock reasoning abilities without the expense of creating specialized training datasets. Second, how to enable models to allocate additional computation to problems requiring more reasoning steps. Third, how to make reasoning paths transparent and interpretable. Most crucially, they sought a method that would work across diverse reasoning domains without task-specific engineering.

### 2.2. How significant is the technical contribution of the paper?

The contribution is remarkably significant - like discovering water can flow uphill under certain conditions. The authors reveal that reasoning abilities aren't missing from large language models but are dormant capabilities that can be awakened through clever prompting techniques alone. This "emergent ability" challenges the conventional wisdom that fine-tuning is necessary for complex reasoning tasks.

### 2.3. Identify main strengths of the proposed approach.

- Simplicity that belies its power: The technique requires only modifying prompt examples to include reasoning steps, with no model parameter changes needed

- Flexibility across domains: Works effectively for arithmetic, common sense, and symbolic reasoning tasks

- Performance breakthroughs: Achieved state-of-the-art results on challenging benchmarks like GSM8K using only prompting

### 2.4. Identify main weaknesses of the proposed approach.

- Scale dependency: Benefits only materialize with models of roughly 100B+ parameters, limiting practical applications

- No guarantees of reasoning quality: Models can generate plausible-sounding but incorrect reasoning chains

- Prompt sensitivity: Performance varies with different annotators and prompt designs, introducing inconsistency

## 3. Empirical Results

### 3.1. Identify key experimental results, and explain what they signify.

- Chain-of-thought prompting transforms performance on GSM8K math problems, with PaLM 540B jumping

from 18% to 57% accuracy - signifying that the ceiling for prompting-based methods is much higher than previously believed

- The technique enables out-of-distribution generalization to longer sequences in symbolic reasoning tasks, showing it's not merely replicating memorized patterns but enabling genuine reasoning capabilities

### 3.2. Are there any weaknesses in the experimental section?

The experiments have a few blind spots. First, the authors don't thoroughly investigate why chain-of-thought works only at scale - is it about parameter count, training data quality, or architectural choices? Second, while they show reasoning paths for correct answers, they provide limited analysis of how incorrect reasoning paths differ structurally from correct ones. Finally, the study lacks investigation into whether chain-of-thought prompting could exacerbate existing model biases by making flawed reasoning seem more persuasive through step-by-step breakdowns.

## 4. Summary

This paper feels like discovering a hidden room in a house you've lived in for years. I'm impressed by how the authors uncovered latent reasoning abilities in models without changing a single parameter, while concerned about the technique's scale requirements and potential to make incorrect answers seem more convincing. The elegant simplicity of chain-of-thought prompting - showing models how to "think aloud" through examples - stands in refreshing contrast to the complex fine-tuning approaches that dominate the field.

## 5. QA Prompt for a Paper Discussion

### 5.1. Discussion Question

Why might chain-of-thought prompting make some incorrect model outputs more convincing to humans than standard prompting?

### 5.2. Your Answer

Chain-of-thought prompting is like watching a magician explain their trick step-by-step - even if the explanation is flawed, the methodical breakdown creates an illusion of validity. When models generate reasoning chains, they wrap errors in the comforting blanket of logical-sounding steps, exploiting our cognitive bias to trust processes that appear transparent. This "reasoning theater" can be particularly dangerous in high-stakes domains like medical or legal advice, where step-by-step explanations might lower human skepticism precisely when critical evaluation is most needed.