# Title

Causal Intervention for Emergent Ability Prediction in Large Language Models

# Problem Statement

Current methods struggle to predict the emergence of new model capabilities during scaling, limiting our ability to target specific abilities efficiently. This hinders our understanding of how to optimize scaling trajectories and deliberately induce desired capabilities in large language models.

# Motivation

Existing work mostly relies on empirical observation of capability jumps, without a principled framework for predicting or inducing them. Understanding the causal mechanisms behind emergent abilities could allow us to deliberately induce desired capabilities and optimize scaling trajectories. Our approach is inspired by causal inference techniques and aims to develop a predictive framework that can estimate the likelihood of specific abilities emerging given a scaling trajectory.

# Proposed Method

We introduce Causal Intervention for Emergent Ability Prediction (CIEAP). The method consists of four main steps: 1) Develop a causal graph representing hypothesized relationships between model components, training dynamics, and emergent abilities. 2) Perform targeted interventions during training, such as selectively ablating or enhancing specific model components or data characteristics. 3) Observe the effects of these interventions on emergent abilities to refine our causal model. 4) Design a predictive framework that estimates the likelihood of specific abilities emerging given a scaling trajectory. Finally, we develop an adaptive scaling algorithm that uses these predictions to optimally allocate resources towards inducing desired capabilities.

# Step-by-Step Experiment Plan

## Step 1: Develop Causal Graph

Create an initial causal graph based on existing literature and hypotheses about relationships between model components (e.g., attention mechanisms, feedforward layers), training dynamics (e.g., loss curves, gradient norms), and emergent abilities (e.g., arithmetic, few-shot learning). Use tools like NetworkX to represent and visualize the graph.

## Step 2: Define Interventions

Design a set of targeted interventions, such as: a) Selective neuron ablation in specific layers. b) Enhancing attention mechanisms. c) Modifying training data distributions. d) Altering optimization hyperparameters. Each intervention should be precisely defined and implementable in popular deep learning frameworks.

## Step 3: Train Baseline Models

Train a series of language models at different scales (e.g., 100M, 1B, 10B parameters) using a standard architecture like GPT. Use a diverse pre-training corpus and track performance on a range of tasks that test for different emergent abilities.

## Step 4: Perform Interventions

For each intervention, train models at the same scales as the baselines, applying the intervention throughout training. Track the same metrics and task performances as in the baseline models.

## Step 5: Analyze Intervention Effects

Compare the performance trajectories of intervened models against baselines. Identify which interventions significantly impact the emergence of specific abilities. Use statistical techniques like causal discovery algorithms to refine the initial causal graph based on observed effects.

## Step 6: Develop Predictive Framework

Using the refined causal graph and intervention results, design a machine learning model (e.g., a graph neural network) that takes as input a description of a model's architecture, training regime, and current performance, and outputs probabilities of specific abilities emerging at future scales.

## Step 7: Validate Predictions

Train additional models at new scales or with new architectures, and compare the actual emergence of abilities with the predictions from the framework. Use metrics like AUC-ROC to evaluate prediction quality.

## Step 8: Design Adaptive Scaling Algorithm

Develop an algorithm that uses the predictive framework to dynamically adjust model scaling and training strategies. The algorithm should make decisions about which components to scale, which interventions to apply, and how to allocate computational resources to target specific desired abilities.

## Step 9: Evaluate Adaptive Scaling

Compare the efficiency and effectiveness of models trained using the adaptive scaling algorithm against baseline scaling approaches. Measure both the speed of ability emergence and the computational resources required.

# Test Case Examples

## Baseline Example

{'Input': 'Train a 1B parameter GPT model on a standard pre-training corpus and evaluate on arithmetic reasoning tasks from the GSM8K dataset.', 'Expected Output': 'The model achieves 20% accuracy on GSM8K, showing limited arithmetic reasoning ability.', 'Explanation': 'Standard scaling approaches often result in unpredictable emergence of abilities, with arithmetic reasoning typically emerging at larger scales.'}

## Proposed Method Example

{'Input': 'Using CIEAP, intervene on a 1B parameter model by enhancing the self-attention mechanism and selectively scaling feedforward layers identified as important for arithmetic reasoning. Train on the same corpus and evaluate on GSM8K.', 'Expected Output': 'The model achieves 45% accuracy on

GSM8K, showing significantly improved arithmetic reasoning ability compared to the baseline at the same parameter count.', 'Explanation': 'CIEAP identifies causal factors contributing to arithmetic reasoning and strategically intervenes to induce this ability earlier in the scaling process.'}

# Fallback Plan

If CIEAP does not significantly improve our ability to predict or induce emergent abilities, we can pivot the project towards an in-depth analysis of why certain interventions failed to have the expected effects. This could involve more fine-grained probing of model internals during training, such as analyzing attention patterns or activation distributions in response to different interventions. We could also explore alternative causal modeling techniques, such as structural equation models or dynamic causal modeling, to capture more complex relationships between model components and emergent abilities. Additionally, we could focus on developing a taxonomy of emergent abilities and their potential interdependencies, which could provide valuable insights for future work on targeted capability induction in language models.

Ranking Score: 6