RIKEN

AIP Center for Advanced Intelligence Project

# Mechanistic Understanding of GenAI

## Vision and Social Impact

The adoption of GenAI has opened many new challenges. Mechanistic understanding of GenAI is needed to ensure safety, fairness, and trustworthiness.

Our vision is to achieve this understanding by analyzing models at three levels of granularity:

- **Neuronal**: What is the function of individual neurons?
- **Representational**: What information do hidden states encode? How does the model use this information?
- **Algorithmic**: What algorithms and reasoning strategies do models use to generate answers? How are they implemented in which model components?
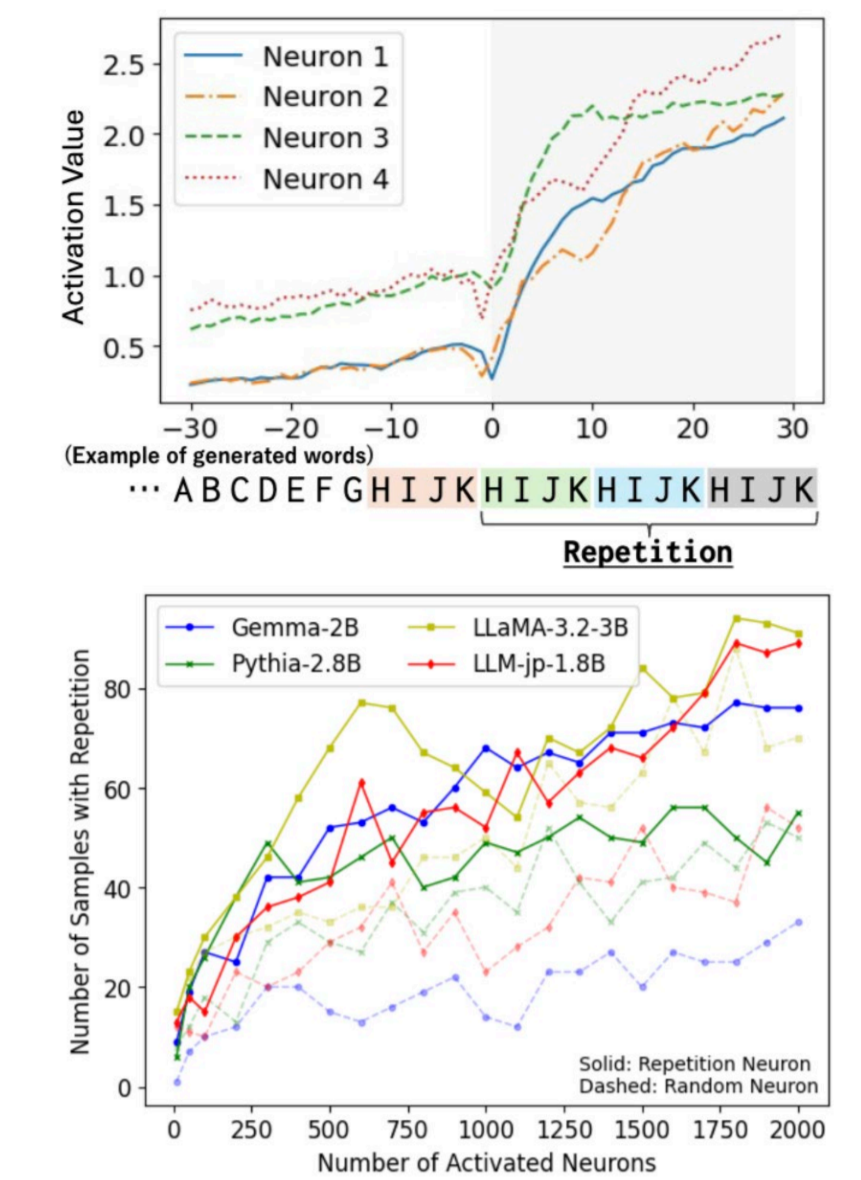
## Neuron-level Analysis

NAACL 2025

**Repetition Neurons: How Do Language Models Produce Repetitions?**

**Motivation:** Repetition is an important and sometimes undesirable aspect of text generation. How LMs perform repetition is not understood on a mechanistic level so far.

**Contribution:** Viewing repetition as a "skill", we devise a method for finding neurons strongly associated with this skill. We then show that repetition can be promoted or suppressed by activating and deactivating these "repetition neurons".
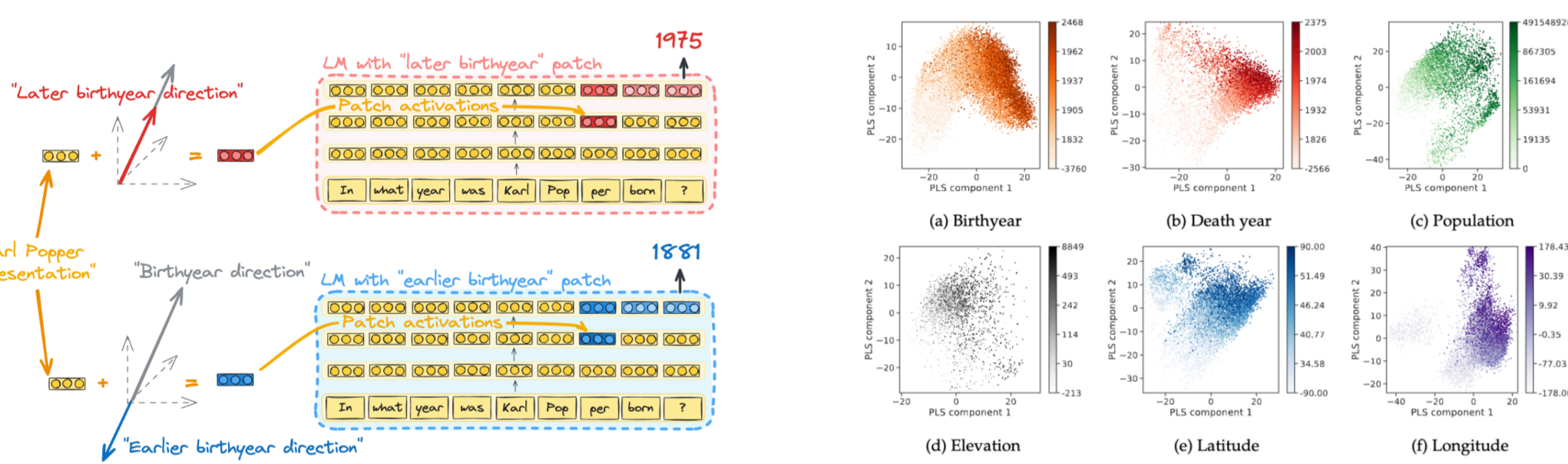


## Understanding the Internal Representations of LLMs

### Monotonic Representation of Numeric Properties in Language Models

ACL 2024

**Motivation:** Want to understand if/how numeric properties (birth year, elevation, population, etc.) are encoded in LMs

**Contribution:** Developed a method for probing and modifying LM representations. Found that numeric properties are encoded in low-dimensional subspaces in a easily-interpretable, monotonic fashion.
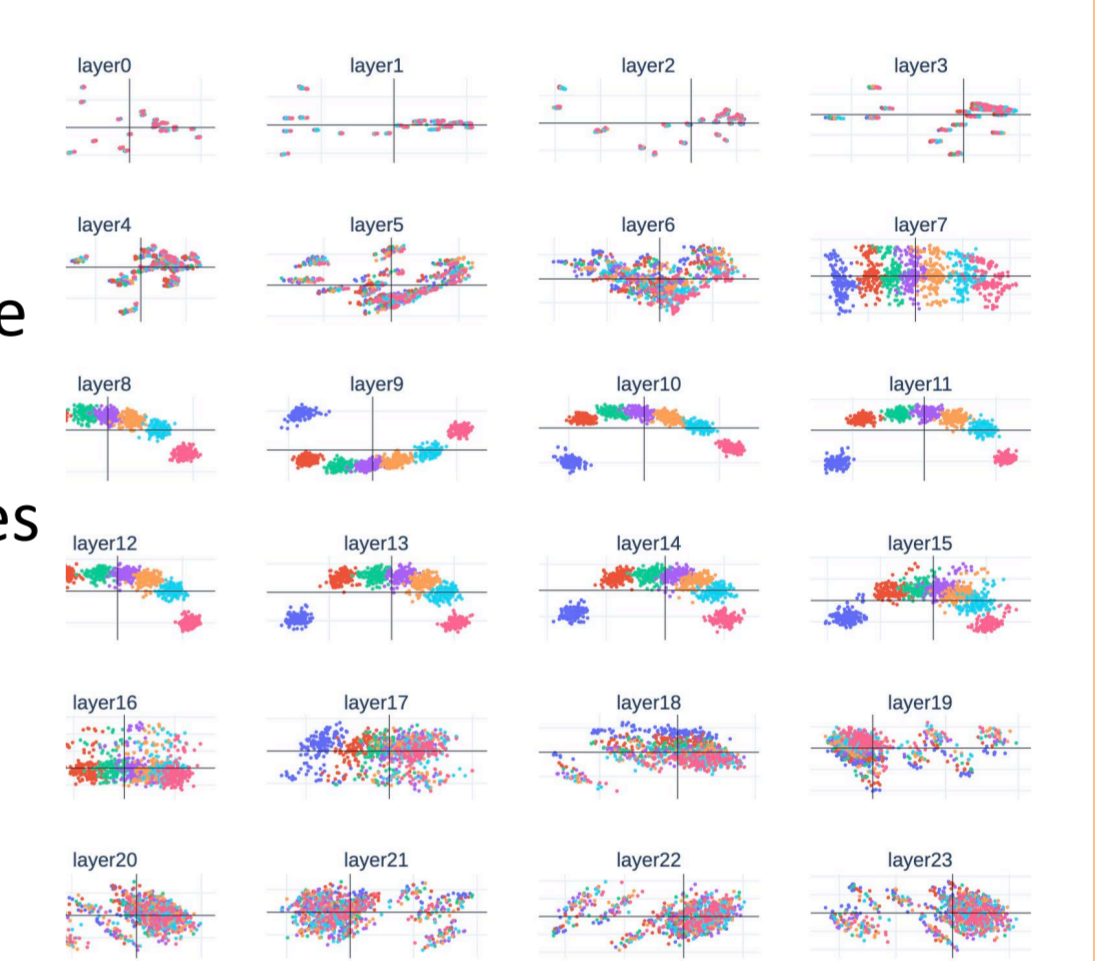


### Representational Analysis of Binding in Language Models

EMNLP 2024

**Motivation:** Analyze if/how LMs perform *binding*. Binding is a basic operation in language understanding by which representations of entities that stand in a relation, e.g., "toy" and "box" in "the *toy* is in the *box*", are combined.

**Contribution:** Identified a subspace in which entities are represented according to their ordering. Causal interventions via activation patching allows controlling the LM's binding behavior, suggesting that this subspace is an important part of the mechanism by which the model performs binding.



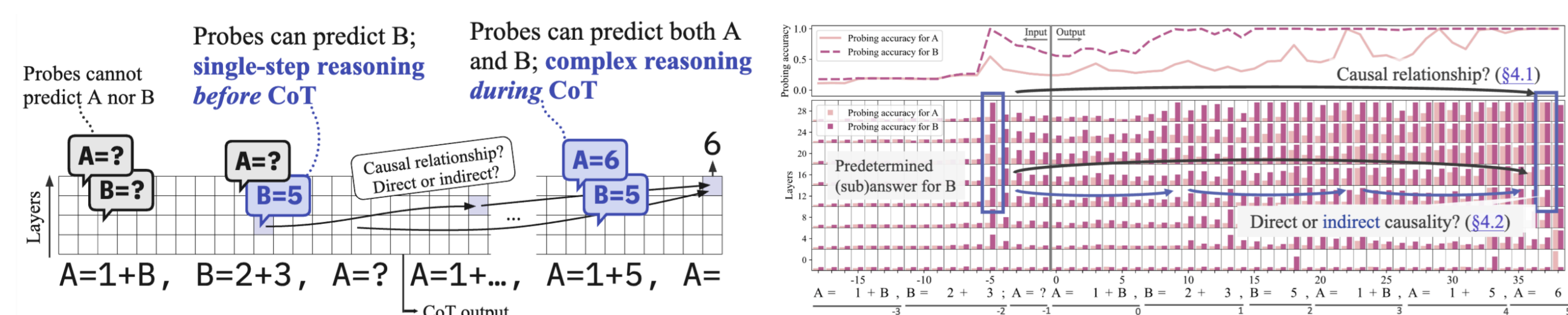## Dissecting the Reasoning Process of LLMs and Vision Models

### Think-to-Talk or Talk-to-Think?
### When LLMs Come Up with an Answer in Multi-Step Reasoning

**Motivation:** Chain-of-thought is a prompting technique by which one asks LMs not only to generate an answer, but also a chain of reasoning steps. However, the degree to which this chain reflects the model's actual reasoning is unclear.

**Contribution:** By conducting systematic probing, we reveal a complex relationship between internal reasoning and expressed reasoning steps.
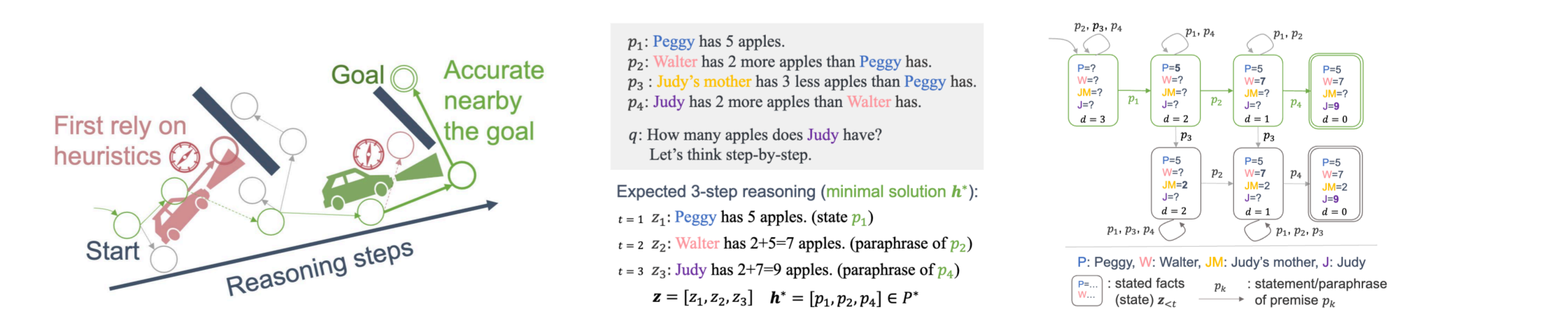


### First Heuristic Then Rational:
### Dynamic Use of Heuristics in Language Model Reasoning
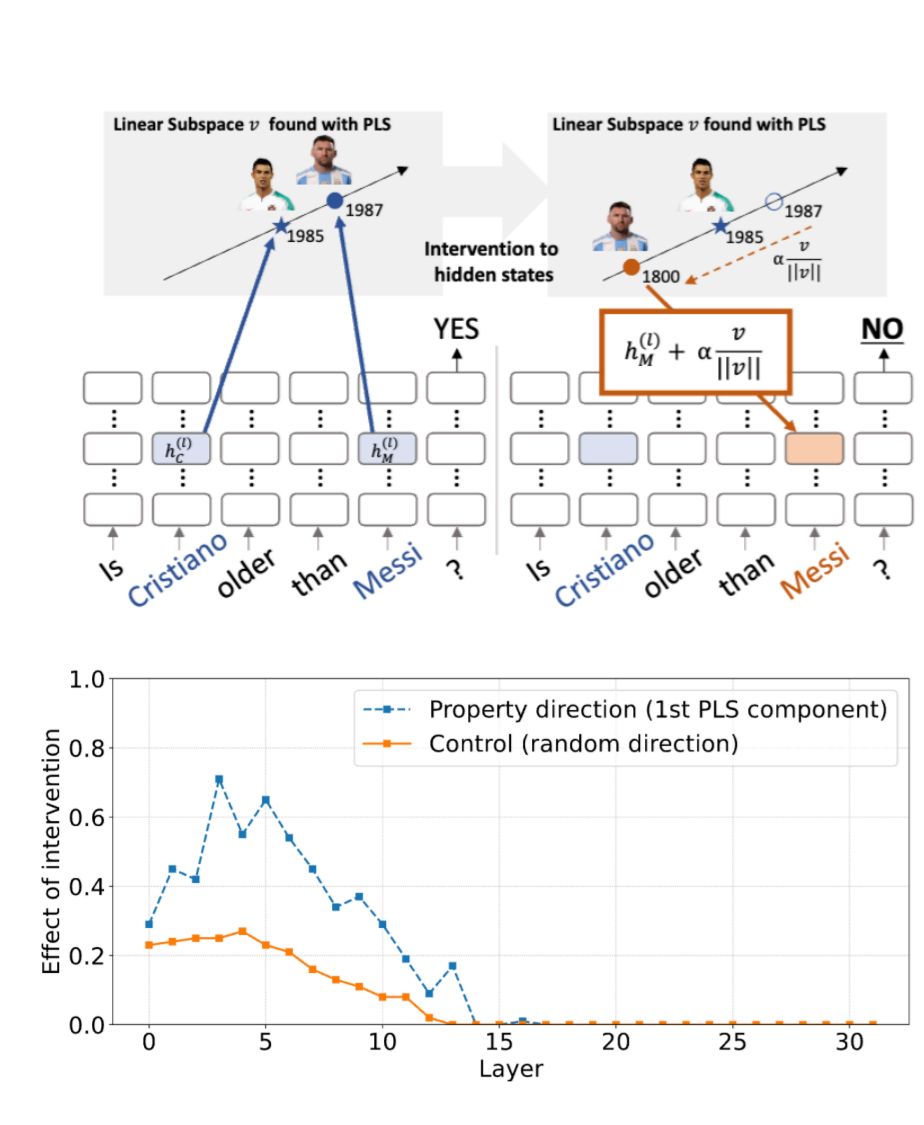
EMNLP 2024

**Motivation:** LMs sometimes appear to perform complex multi-step reasoning but also can make simple mistakes. We analyze why such reasoning failures happen.

**Contribution:** We show that use of simple heuristic, such as lexical overlap, are more prevalent in the earlier stages of reasoning, while precise reasoning is more common towards the end of the reasoning chain.



### The Geometry of Numerical Reasoning

NAACL 2025

**Motivation: In p**rior work we found that numeric properties are represented in low-dimensional linear subspaces, but it is unknown if/how these representations are used in the reasoning process of language models.

**Contribution:** Found LMs compare numeric properties in low-dimensional linear subspaces of activation space. By intervening on representations in these subspaces, one can control the outcome of numeric comparisons.



### Sketch2Diagram: Generating Vector Diagrams from Hand-Drawn Sketches

ICLR 2025

**Motivation:** Creating and understanding diagrams and vector graphics involves abstract reasoning about discrete. This task is challenging for current models.

**Contribution:** To enable progress on this complex task, we introduce a dataset consisting of hand-drawn sketches and vector diagrams. In addition to evaluating large models, we develop augmentation methods and show that a small model trained on our dataset can outperform GPT-4o.