

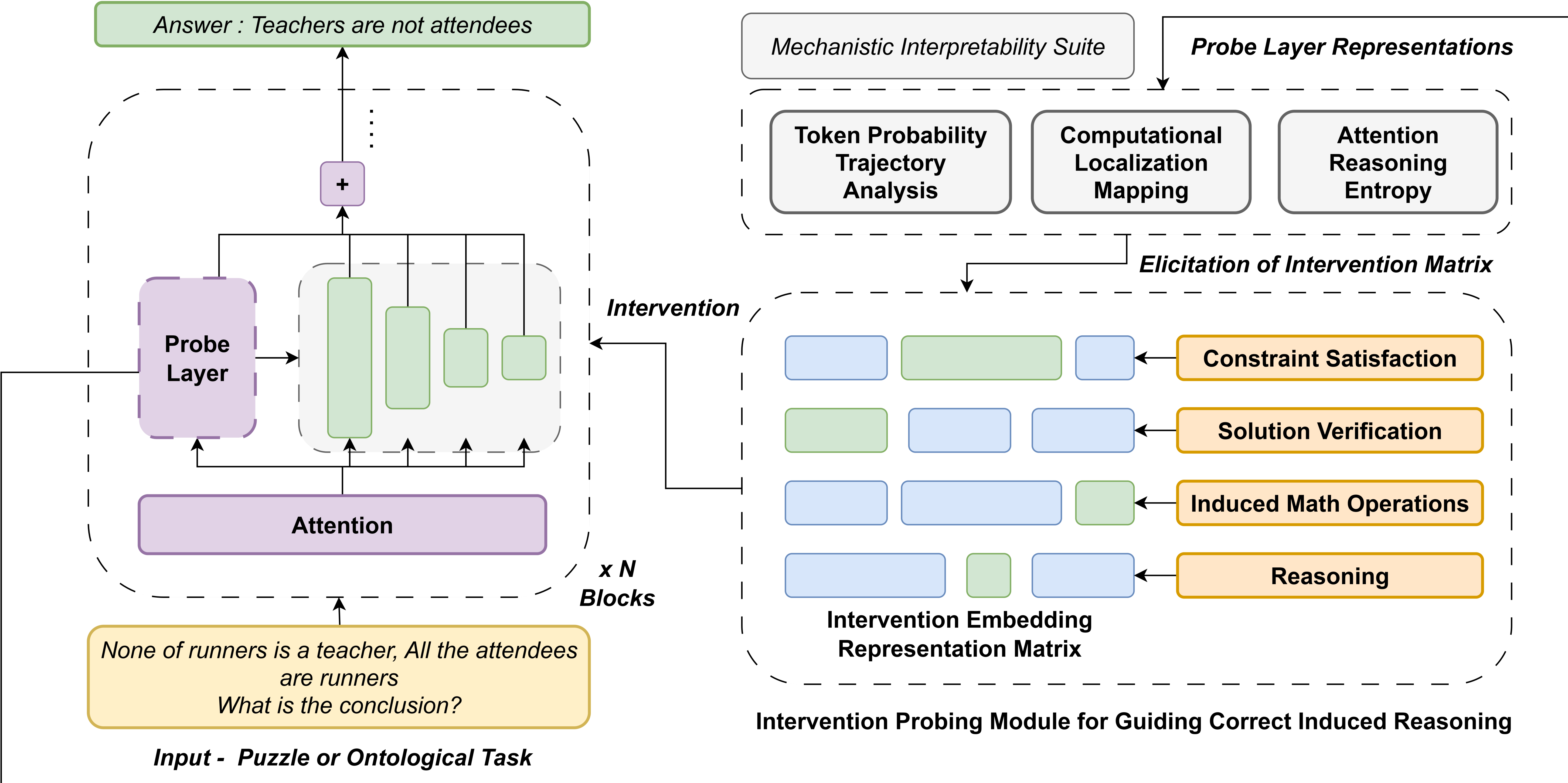
Thesis Proposal: Interpretable Reasoning Enhancement in Large Language Models through Puzzle and Ontological Task Analysis

Mihir Panchal

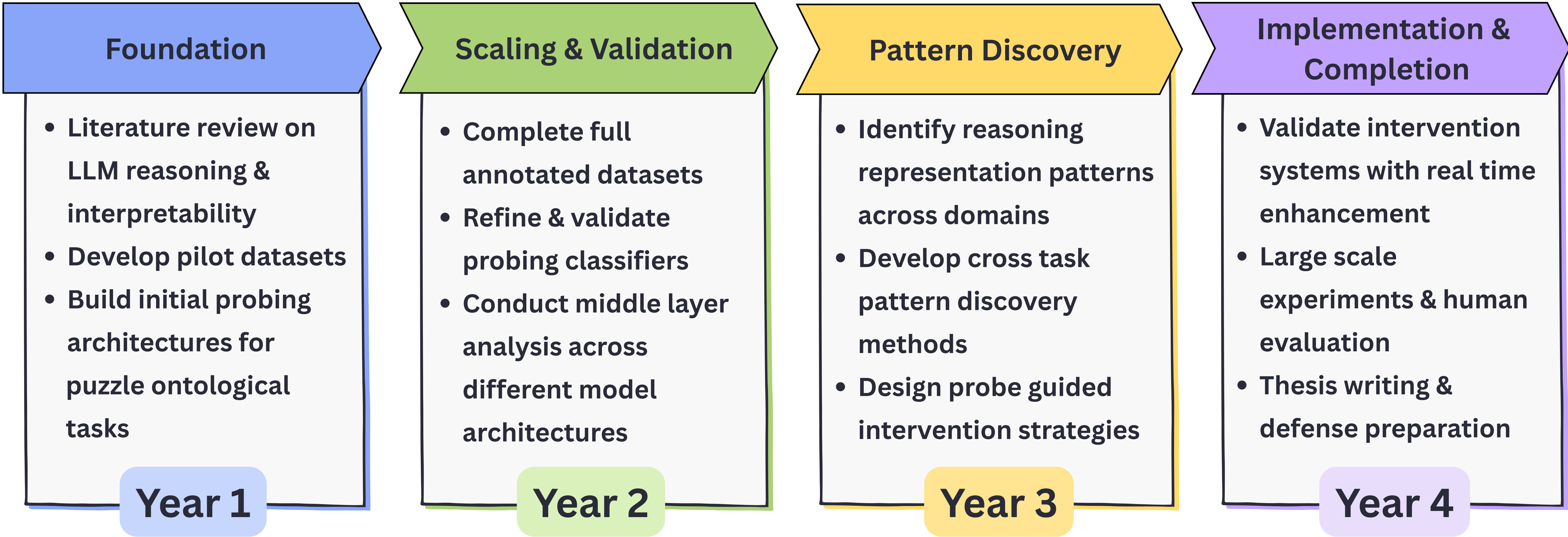


Research Questions

- ▶ How can systematic error patterns in domain specific reasoning be detected through layer wise probing and mitigated through targeted interventions?
- ▶ How can probing frameworks and middle layer analyses reveal and enhance the computational mechanisms underlying inference?



Aims and Timeline Deliverables



Scan here to connect or collaborate

- Pilot datasets, baseline probing architectures, 1-2 publications
- Public dataset release, validated classifiers, 1-2 conference papers
- Unified analysis framework, intervention prototypes, 1-2 major publications
- Complete intervention system, thesis, final publications, public tools



Scan here to read the thesis proposal