Training Stage

Training Data Curation,
Alignment Training,
Knowledge Grounding



Inference Stage

Prompt Filtering, Intent Modeling, Defense Against Jailbreak, Confidence Estimation, Retrieval Augmentation, Inference-Time Factuality Verification



Influence Stage

LLM-Generated Misinformation Detection, LLM-Generated Text Detection, Public Education