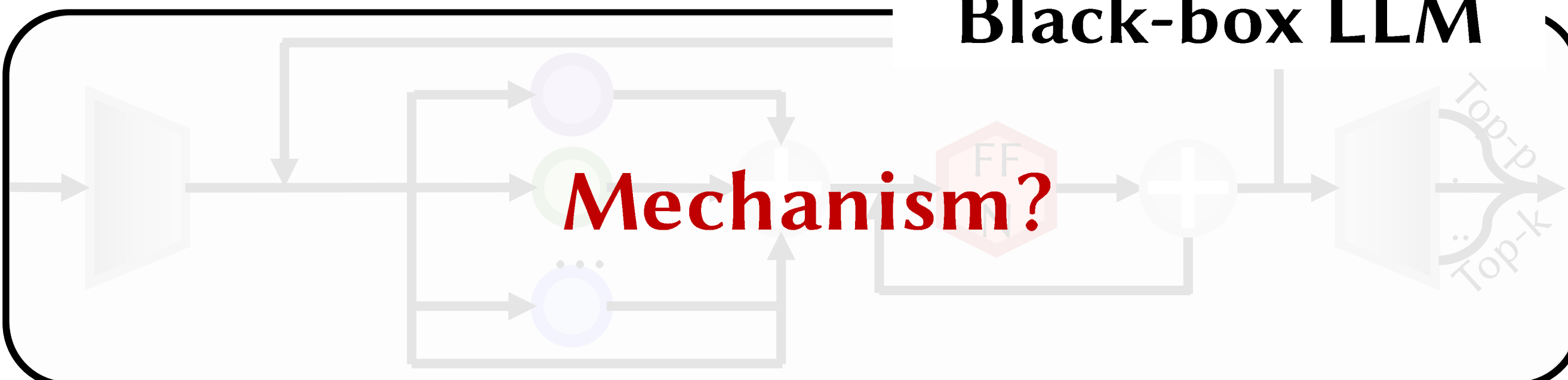


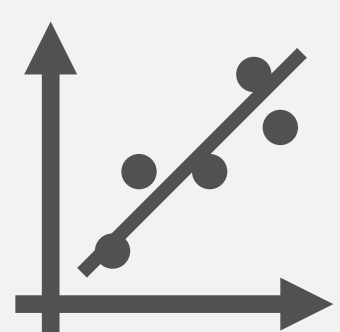
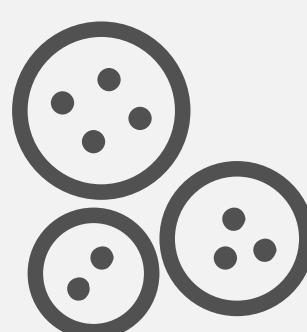
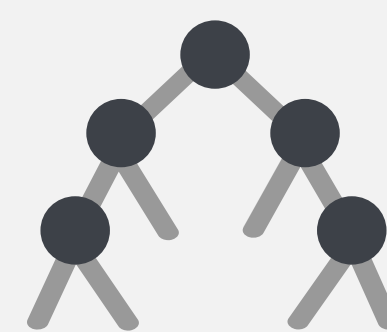
Attention Heads of Large Language Models: A Survey

Introduction & Background (Sec.1-2)

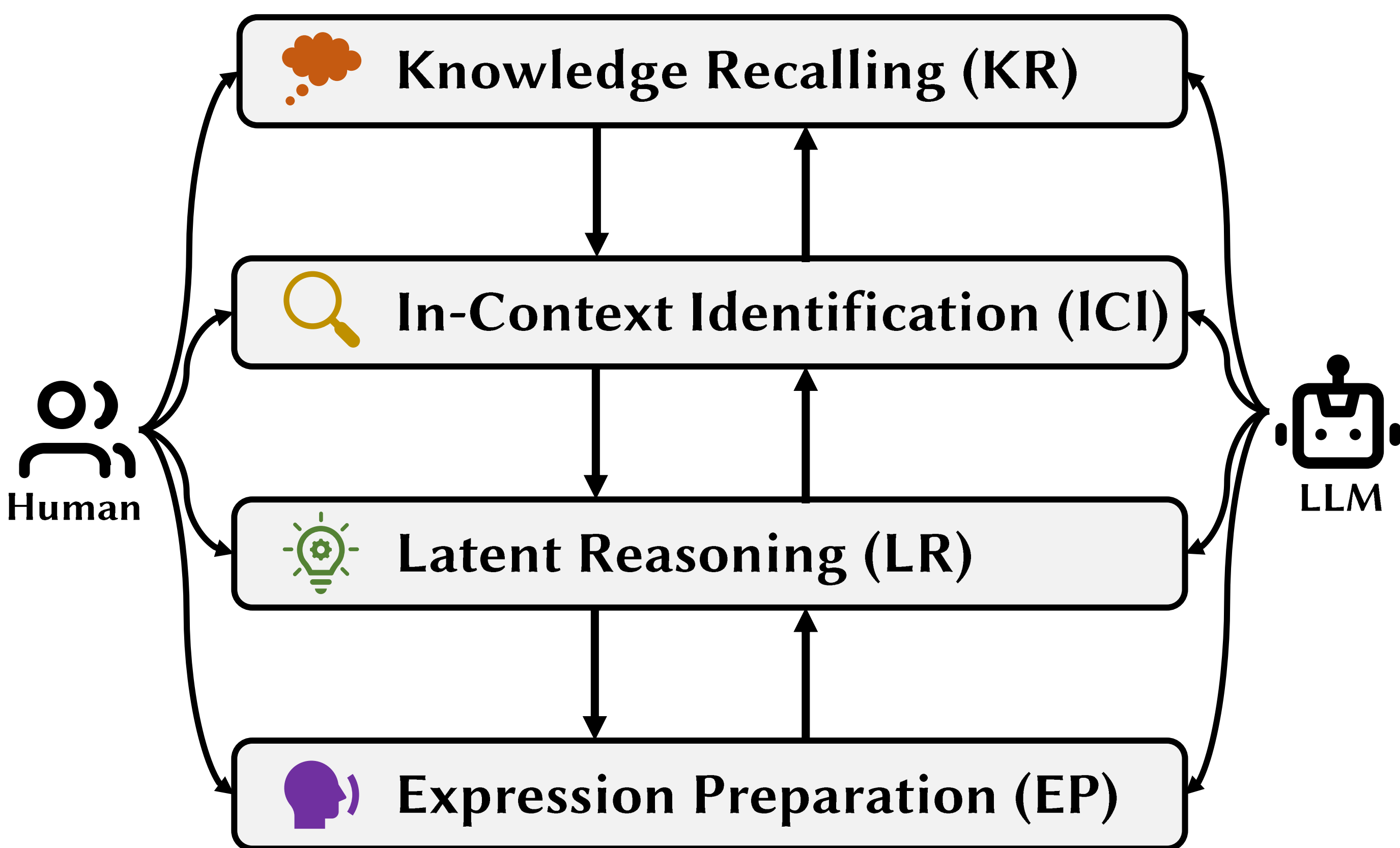
Black-box LLM



White-box Model

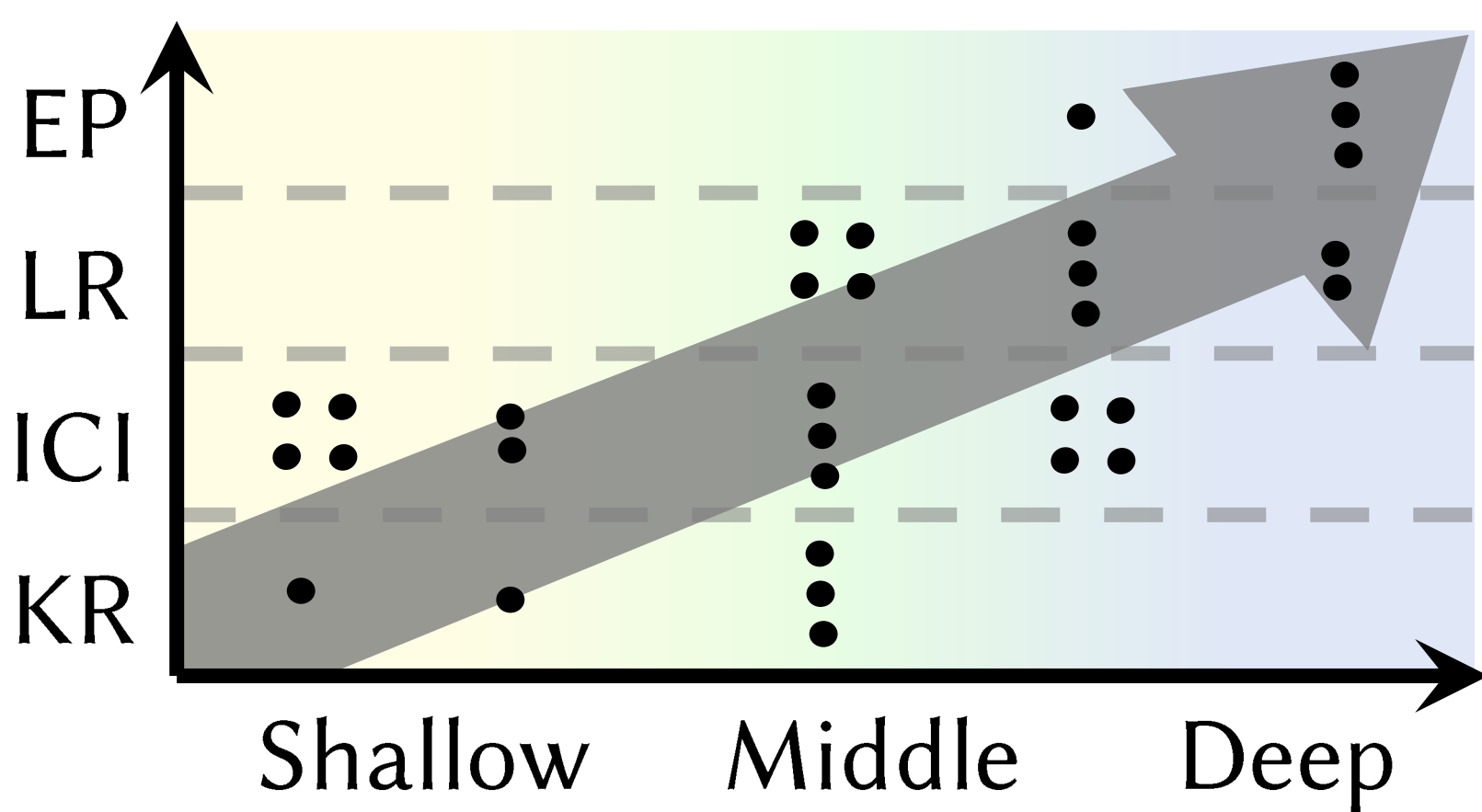


Special Attention Heads (Sec.3.1–3.5)



Collaborative Mechanism (Sec.3.6)

Relation between 4 stages and model layers



Application
Scene

MCQA

Parity Problem

IOI

...

Conclusion (Sec.7)

Limitations

Universality

Theoretical support

Future Directions

Complex Task

Unified Framework

Align with Human

Evaluation (Sec.5)

Mechanism Exploration

Common Evaluation

Add. Topics (Sec.6)

FFN Interpretability

Machine Psychology

Experimental Method (Sec.4)

Modeling-Free

Modification-Based

Replacement-Based

Modeling-Required

Training-Free

Training-Required