

Challenges and Open Problems

Explainable Information Retrieval

Evaluation

Current evaluation schemes for post-hoc interpretability limited

- inject bugs and treat them as ground truth
- Use a simple model as a BBOX

Little to no **benchmarks** for development and evaluation of interpretability

- none for IR with the exception of ERASER
- Need for a holistic benchmark

Human studies are limited and hard to make progress

- Current studies use weak baselines or controls

But please dont stop accepting papers due to any of the reasons

A glimpse of the future..



Welcome to the new Bing

Your AI-powered copilot for the web

 Ask complex questions

"What are some meals I can make for my picky toddler who only eats orange-colored food?"

 Get better answers

"What are the pros and cons of the top 3 selling pet vacuums?"

 Get creative inspiration

"Write a haiku about crocodiles in outer space in the voice of a pirate"

Let's learn together. Bing is powered by AI, so surprises and mistakes are possible. Make sure to check the facts, and [share feedback](#) so we can learn and improve!



New topic

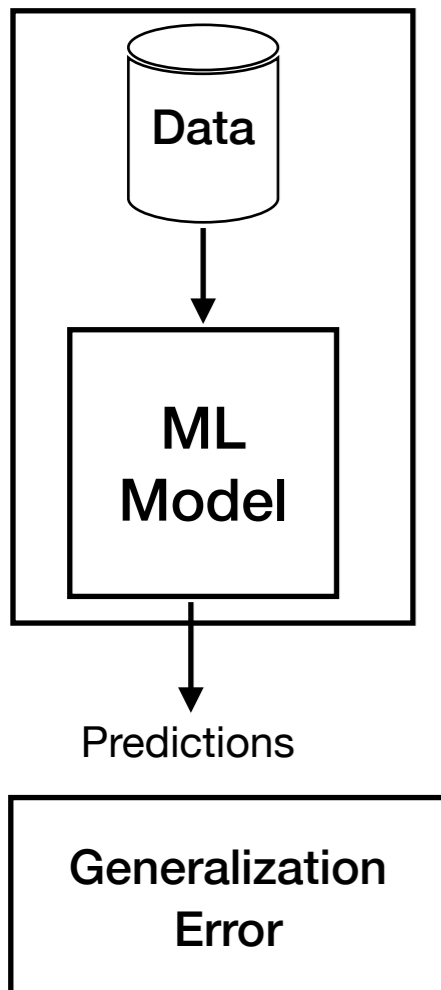


Ask me anything...

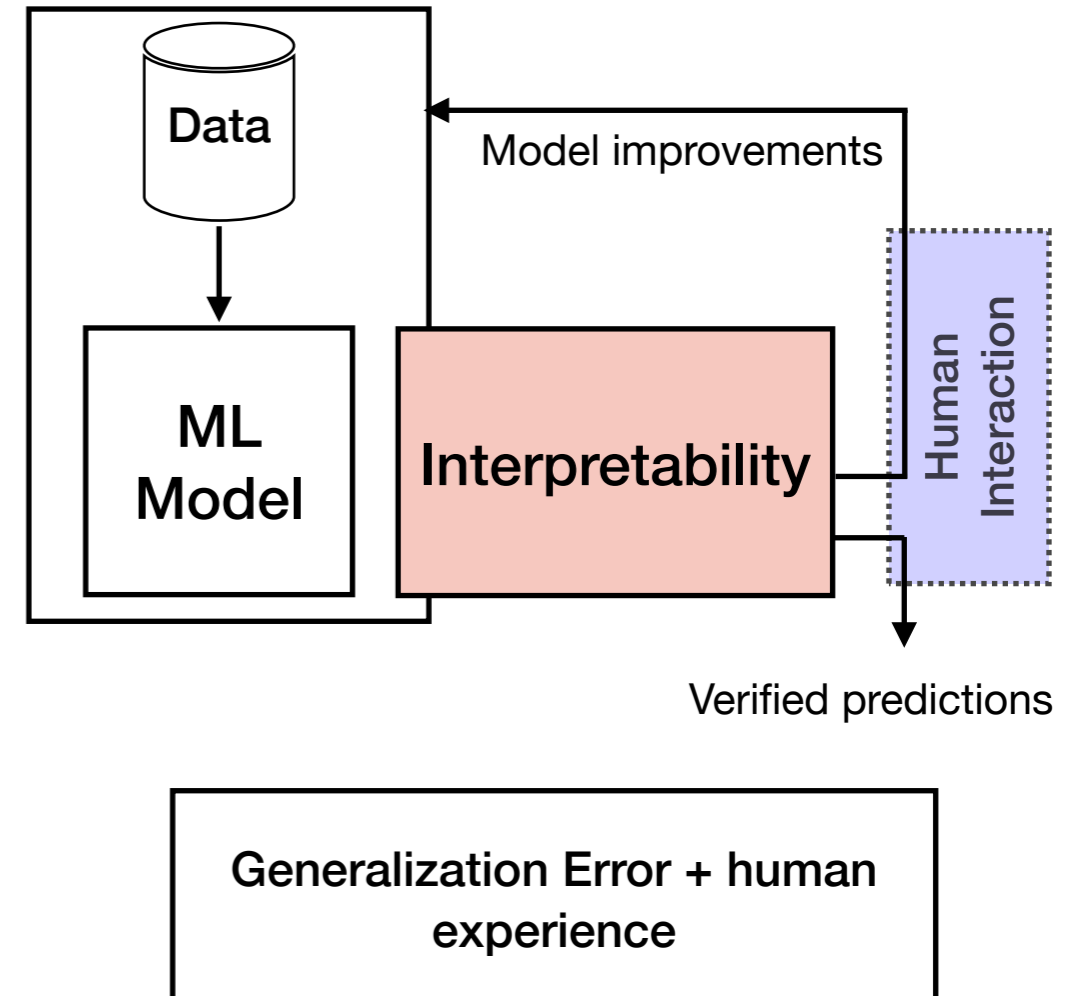
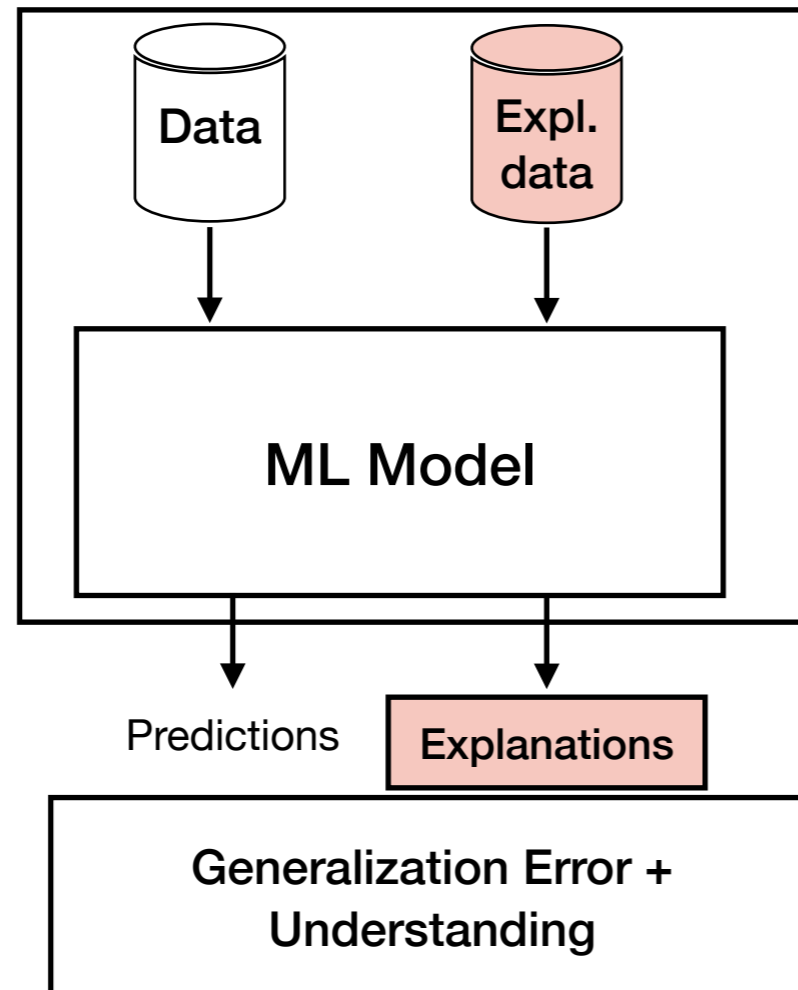
Vision

Interpretability as a first-class citizen in model building

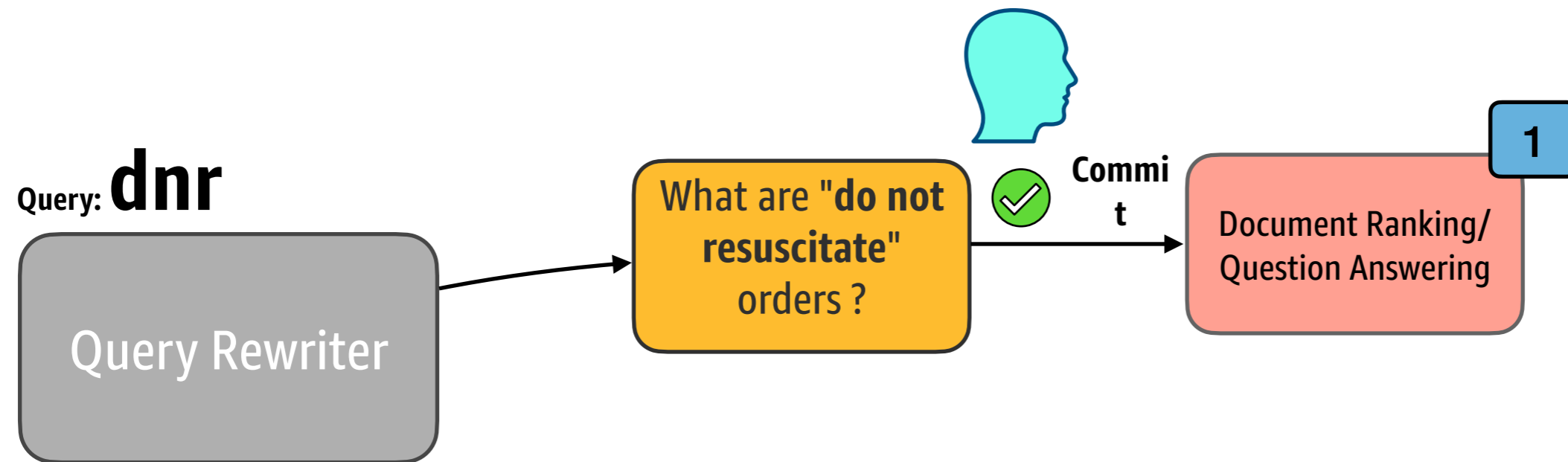
Standard ML



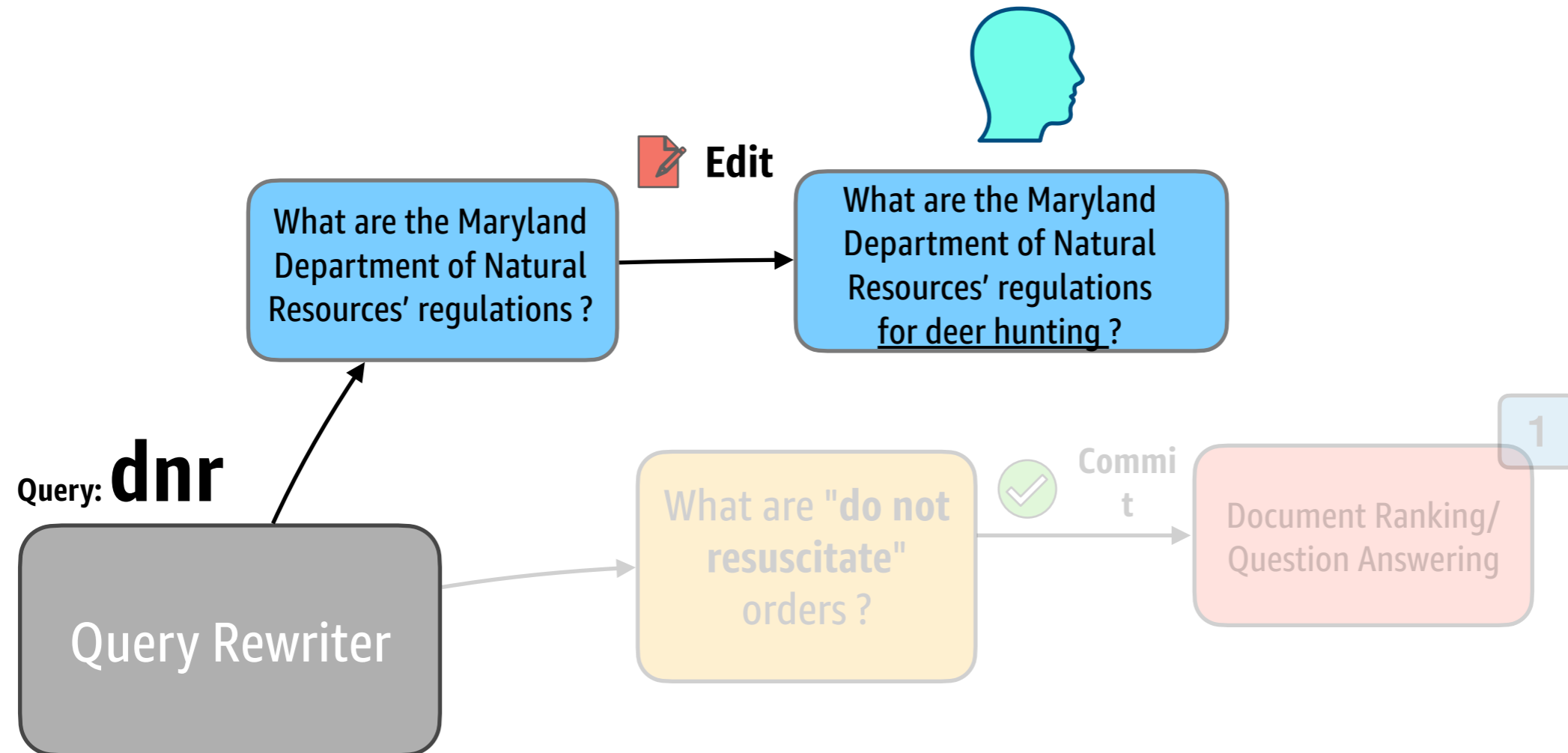
Interpretable ML



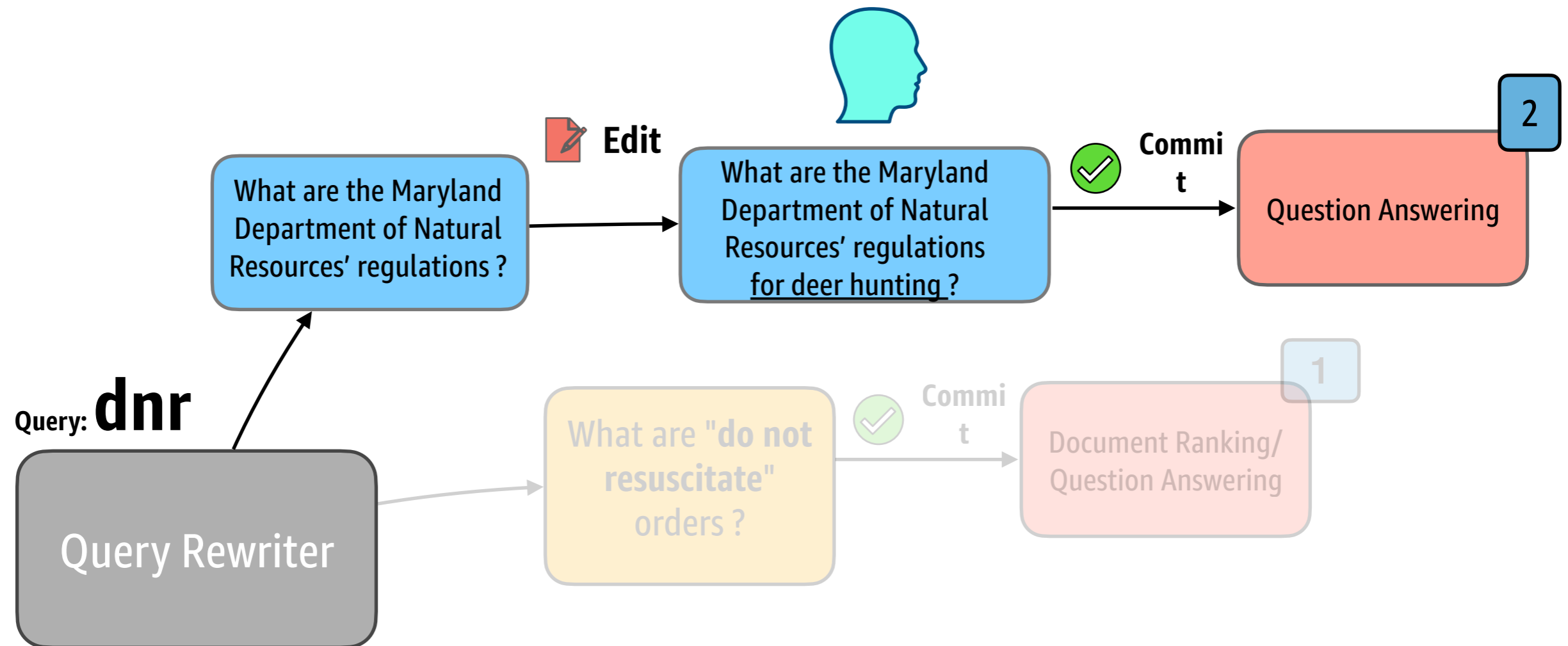
Interpretability in the age of LLMs



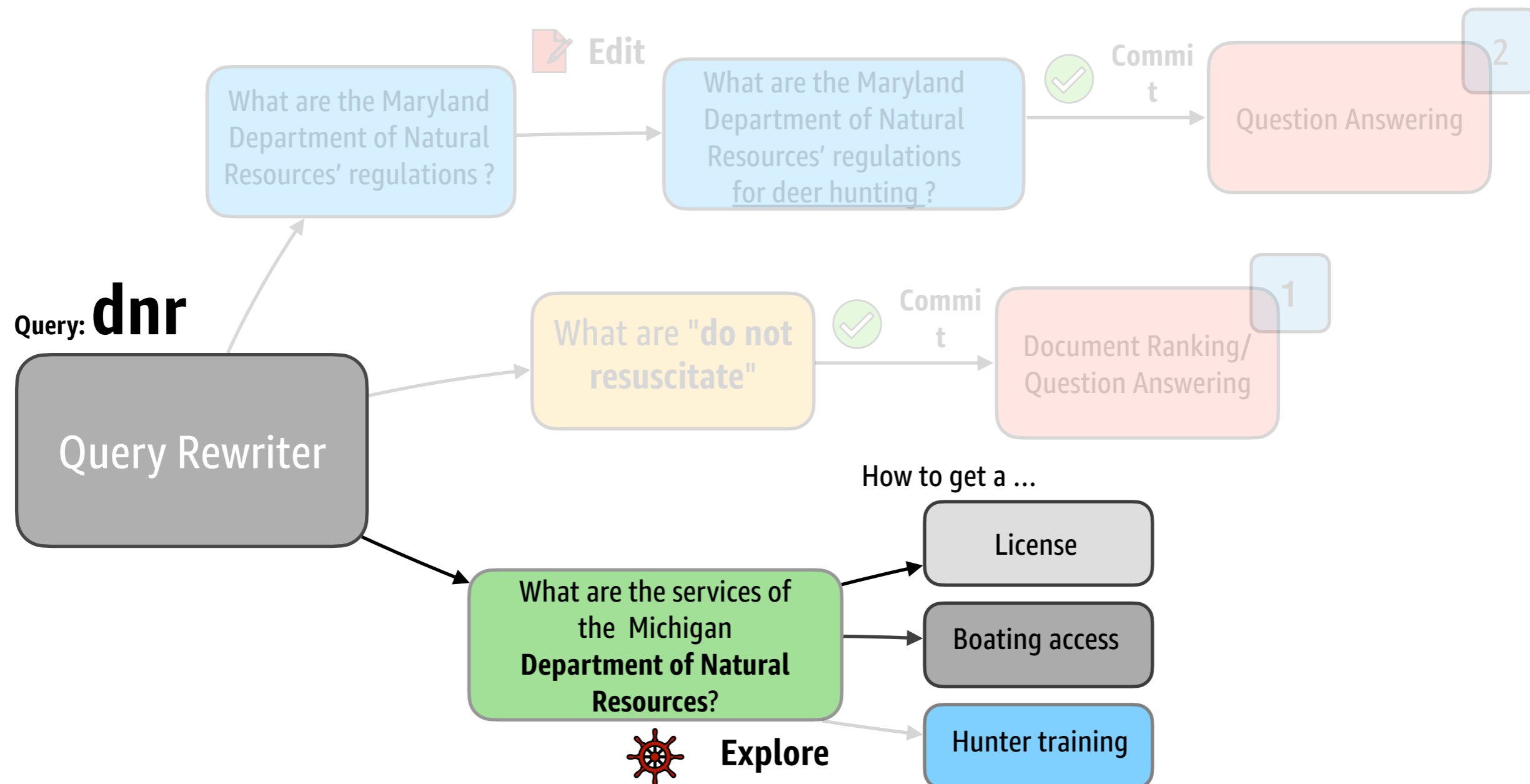
Interpretability in the age of LLMs

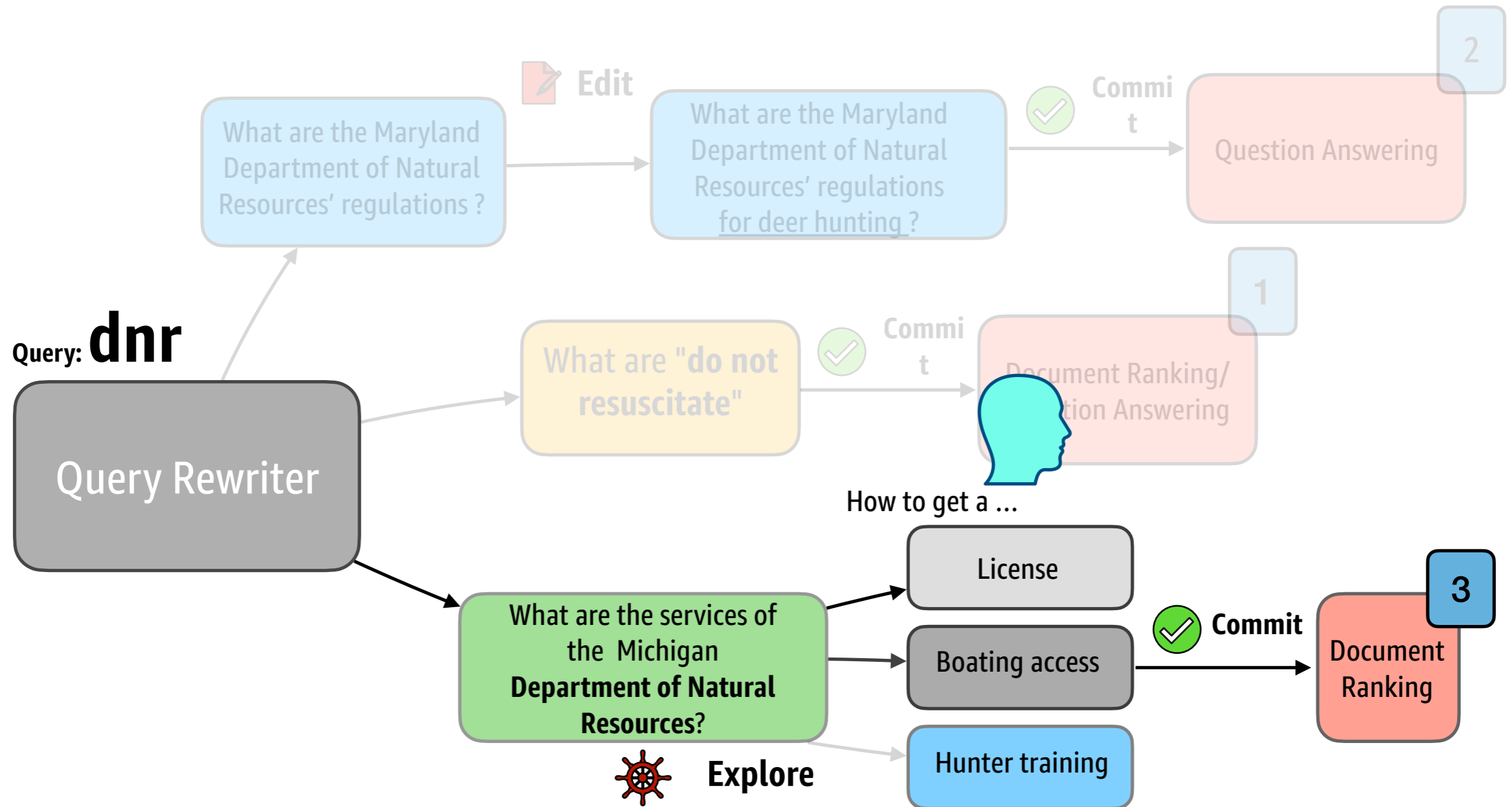


Interpretability in the age of LLMs

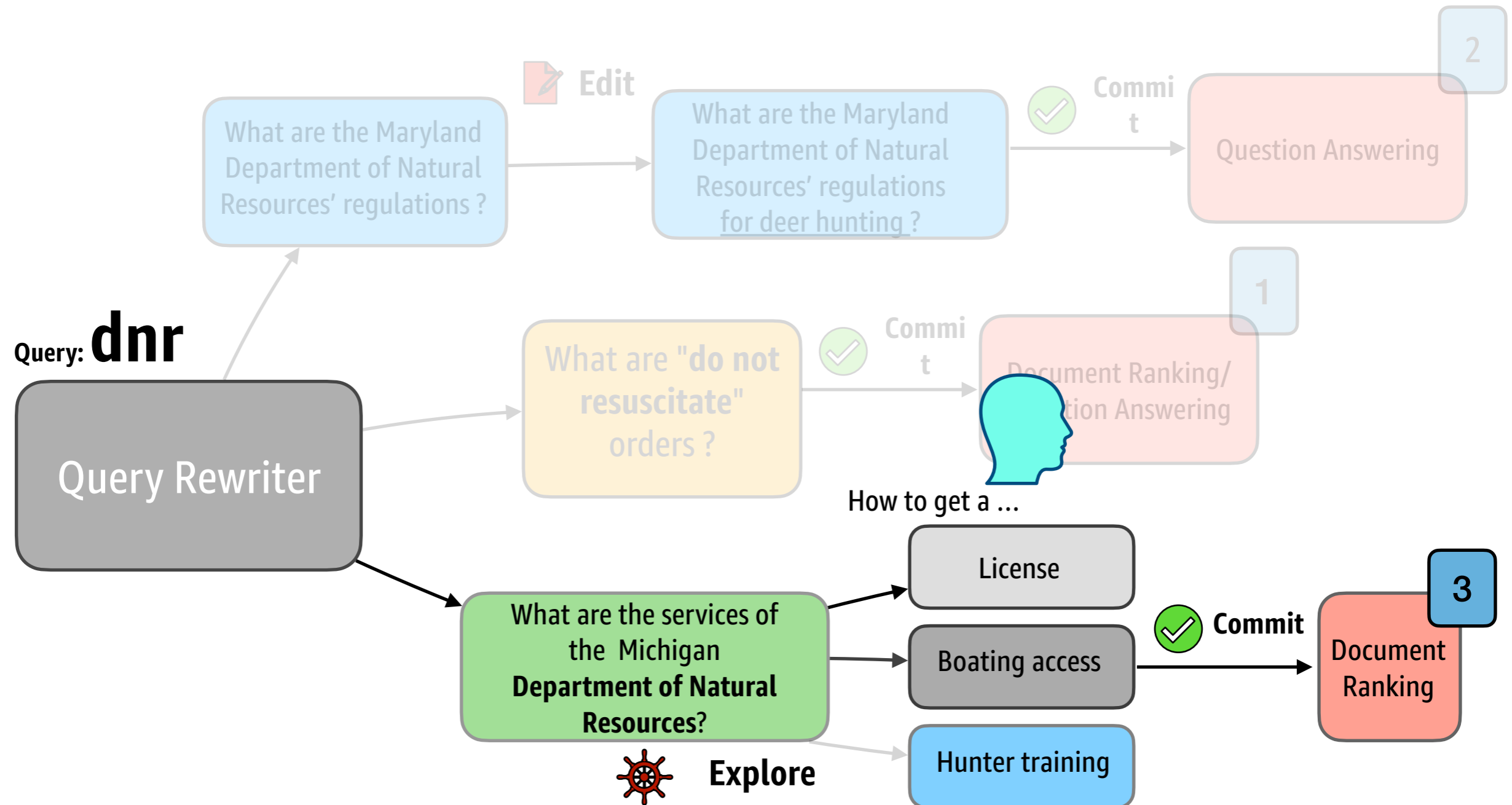


Interpretability in the age of LLMs





Text instead of vectors



Conclusion

Introduction, motivation and notions

Posthoc interpretability

Intrinsic interpretability

Probing LLMs

Axiomatic IR for explaining IR models

Demo

Evaluation or ExIR methods

Conclusion and open problems



Thats it !!!