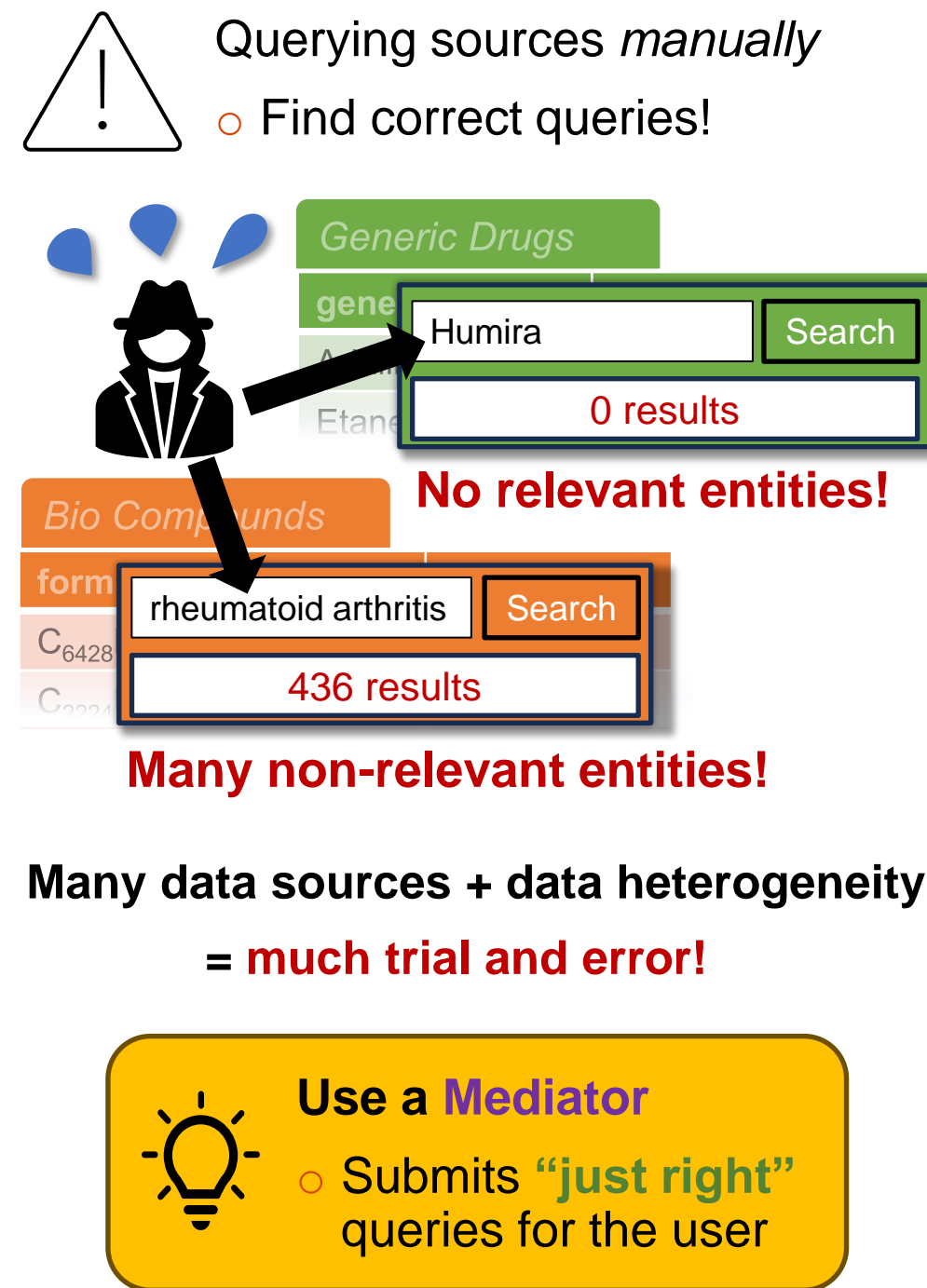
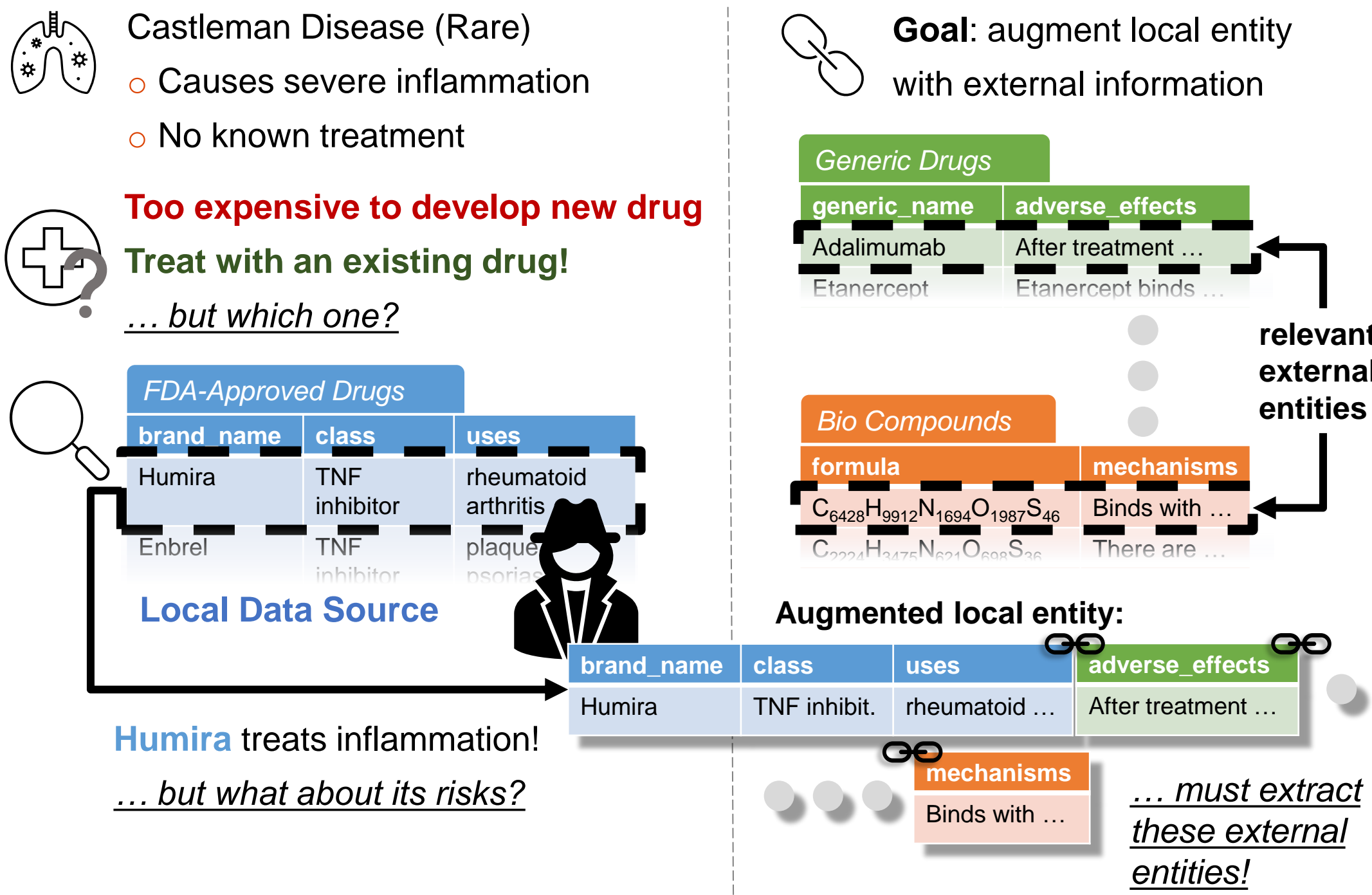


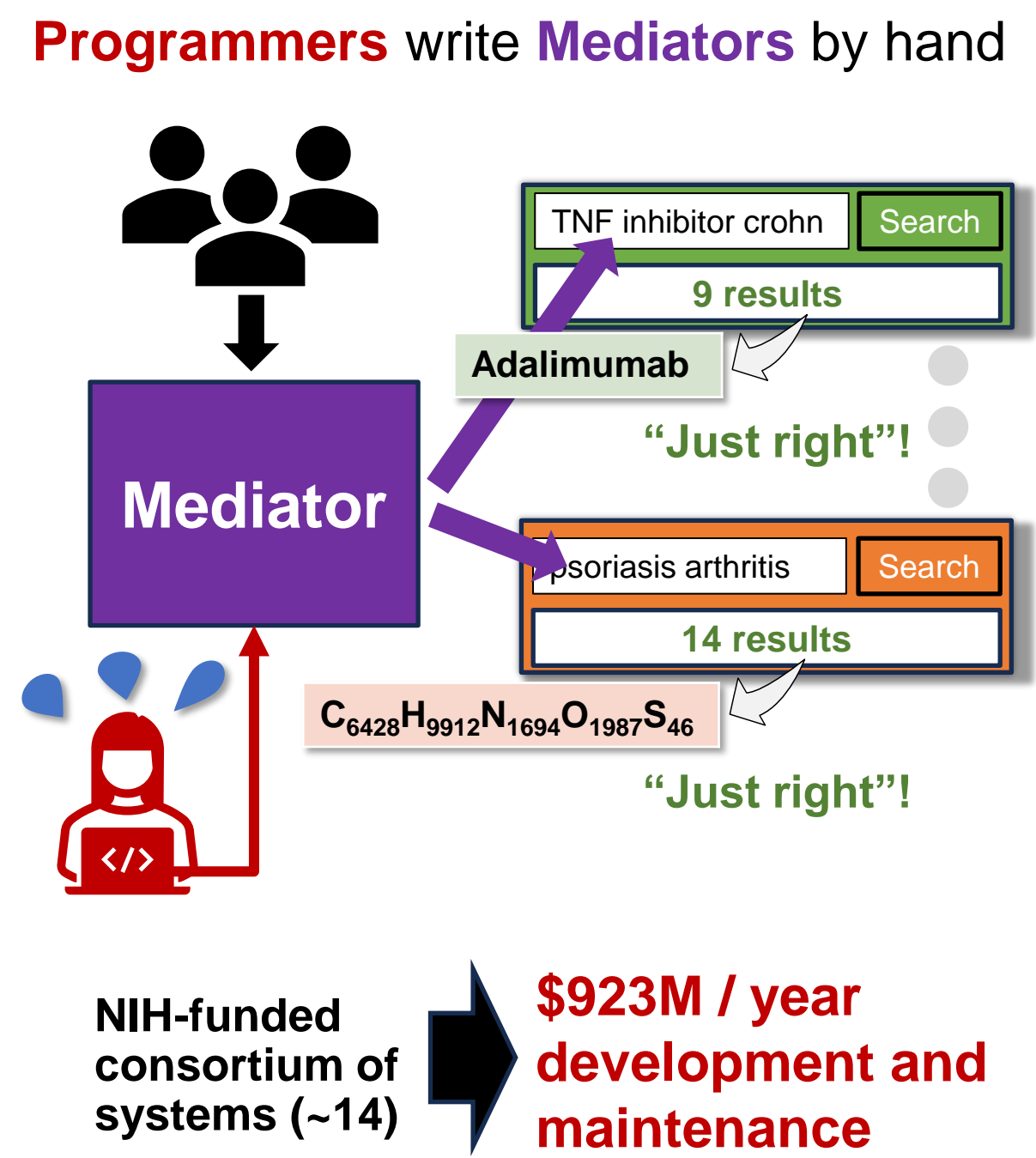
# Generating Data Augmentation Queries Using Large Language Models

Christopher Buss, Jasmin Mousavi, Mikhail Tokarev,  
Mahdis Safari, Arash Termehchy, David Maier, Stefan Lee

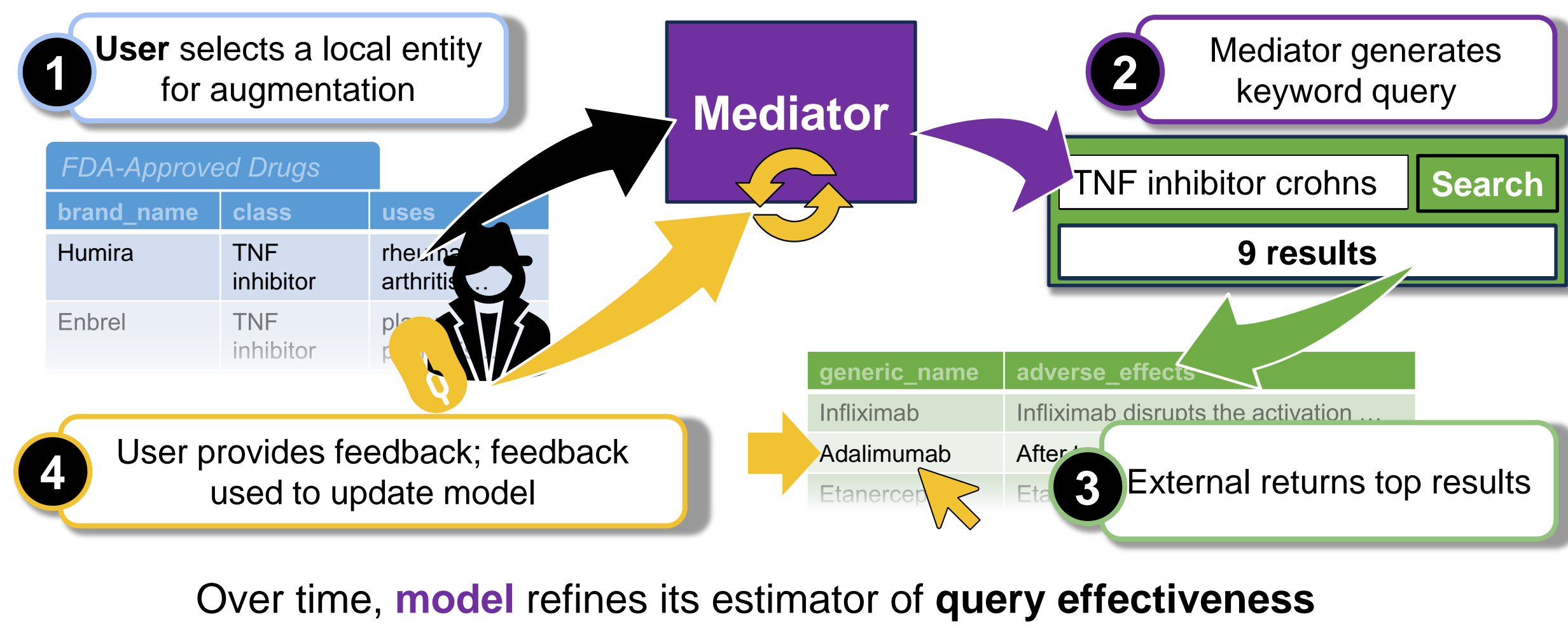
## 1. Drug Repositioning Can Save Lives



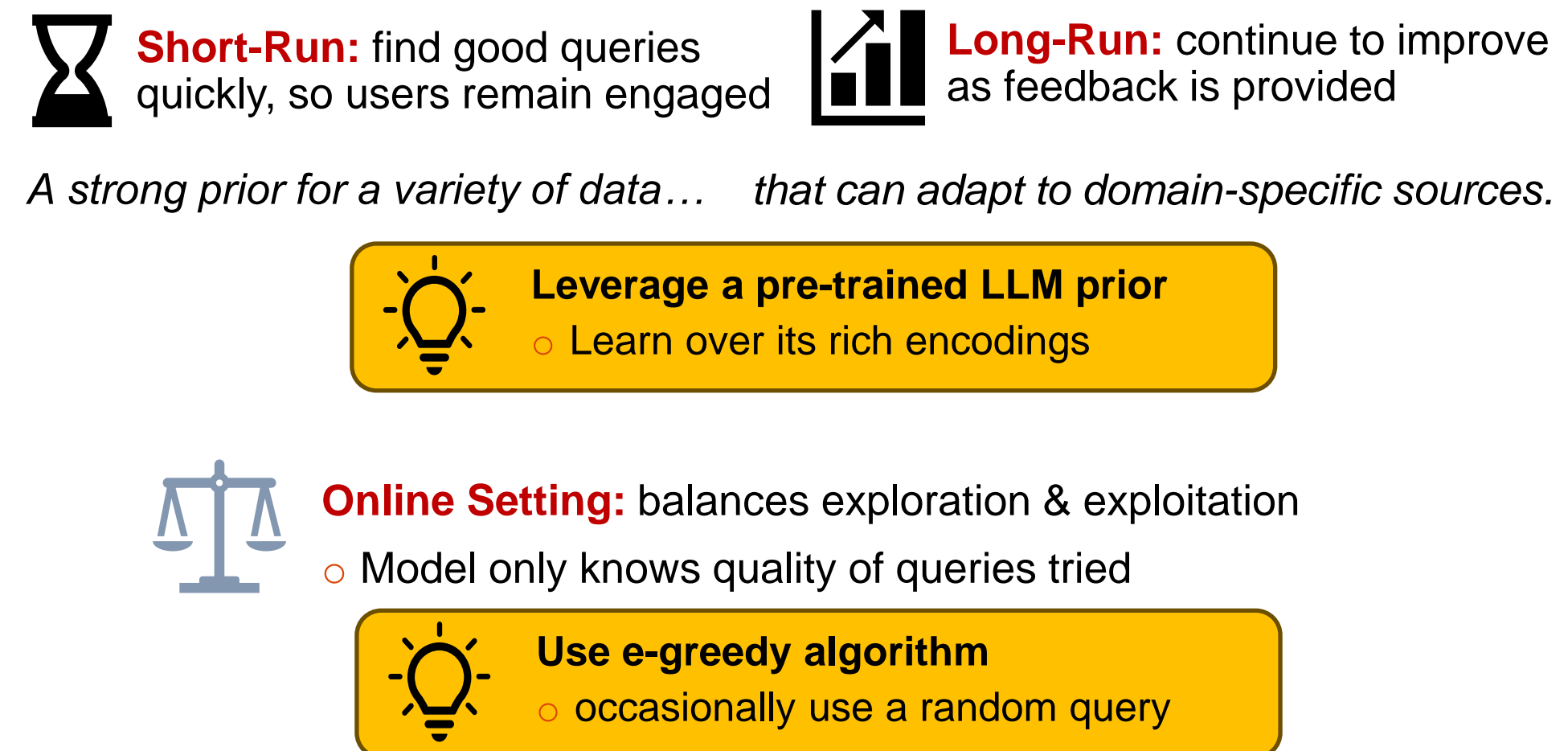
## 2. Existing Work



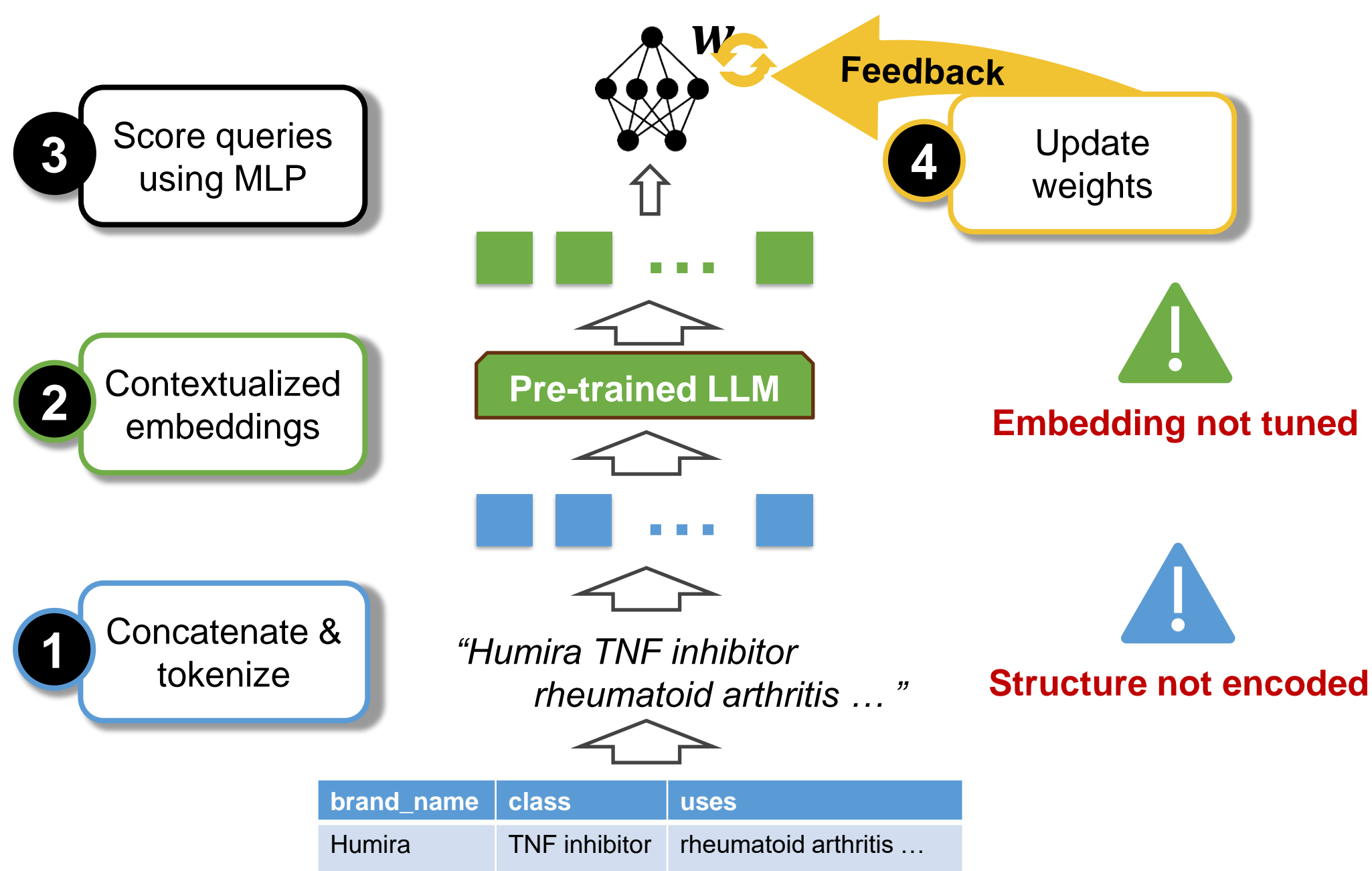
## 3. Online Autonomous Querying



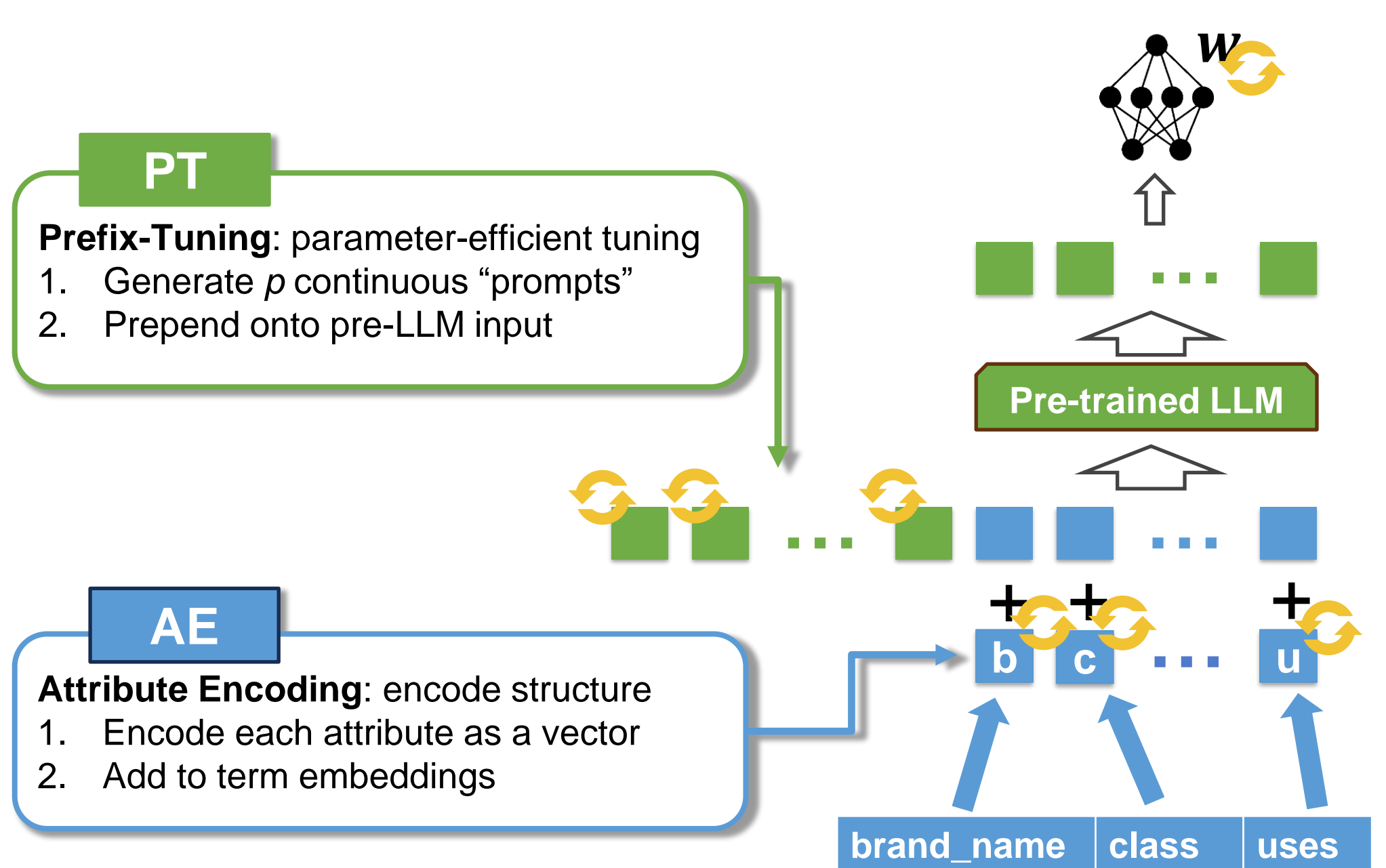
## 4. What Makes a Good Model?



## 7. Learning Query Quality using an LLM



## 8. Lightweight Tuning



## 8. Experimental Simulation

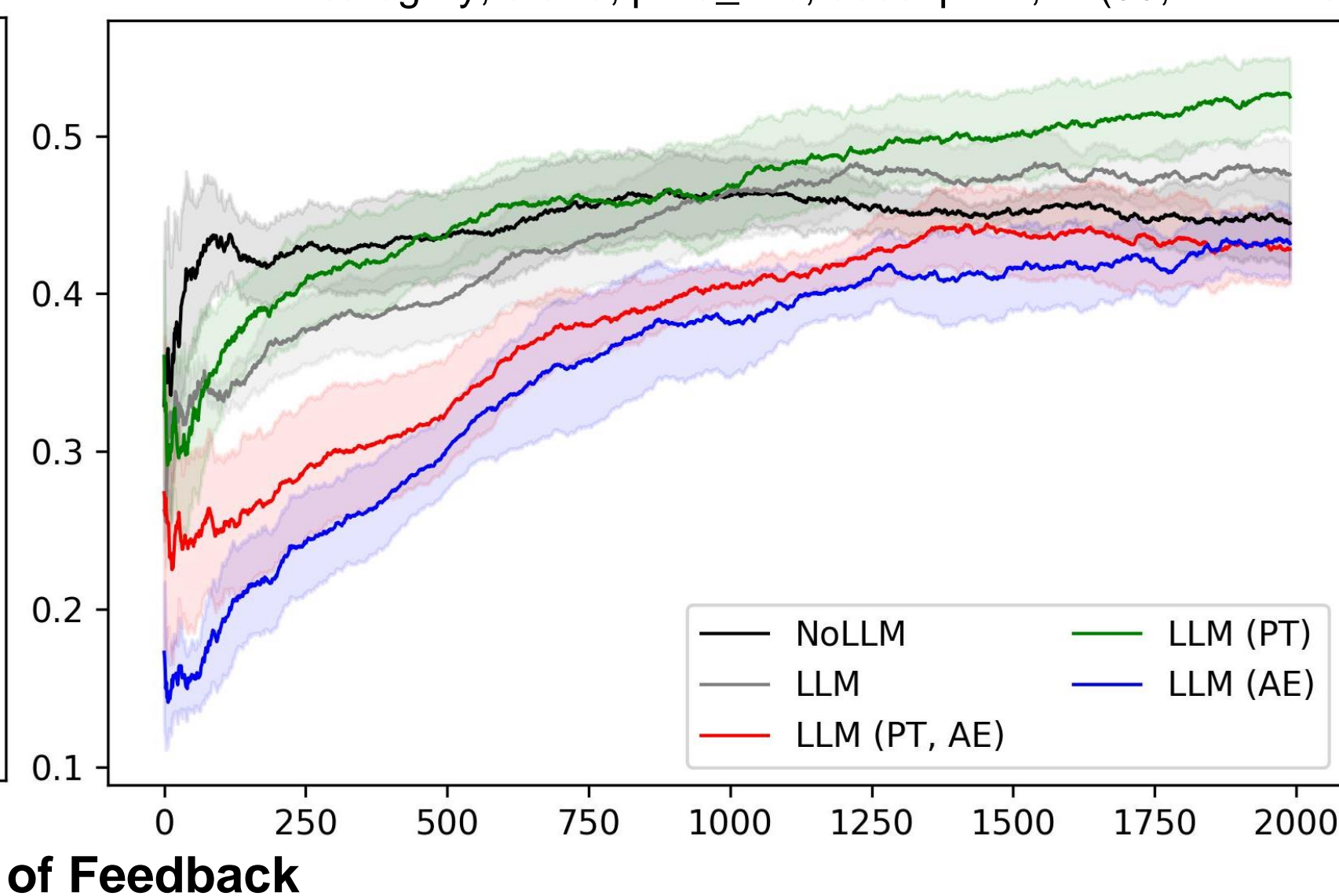
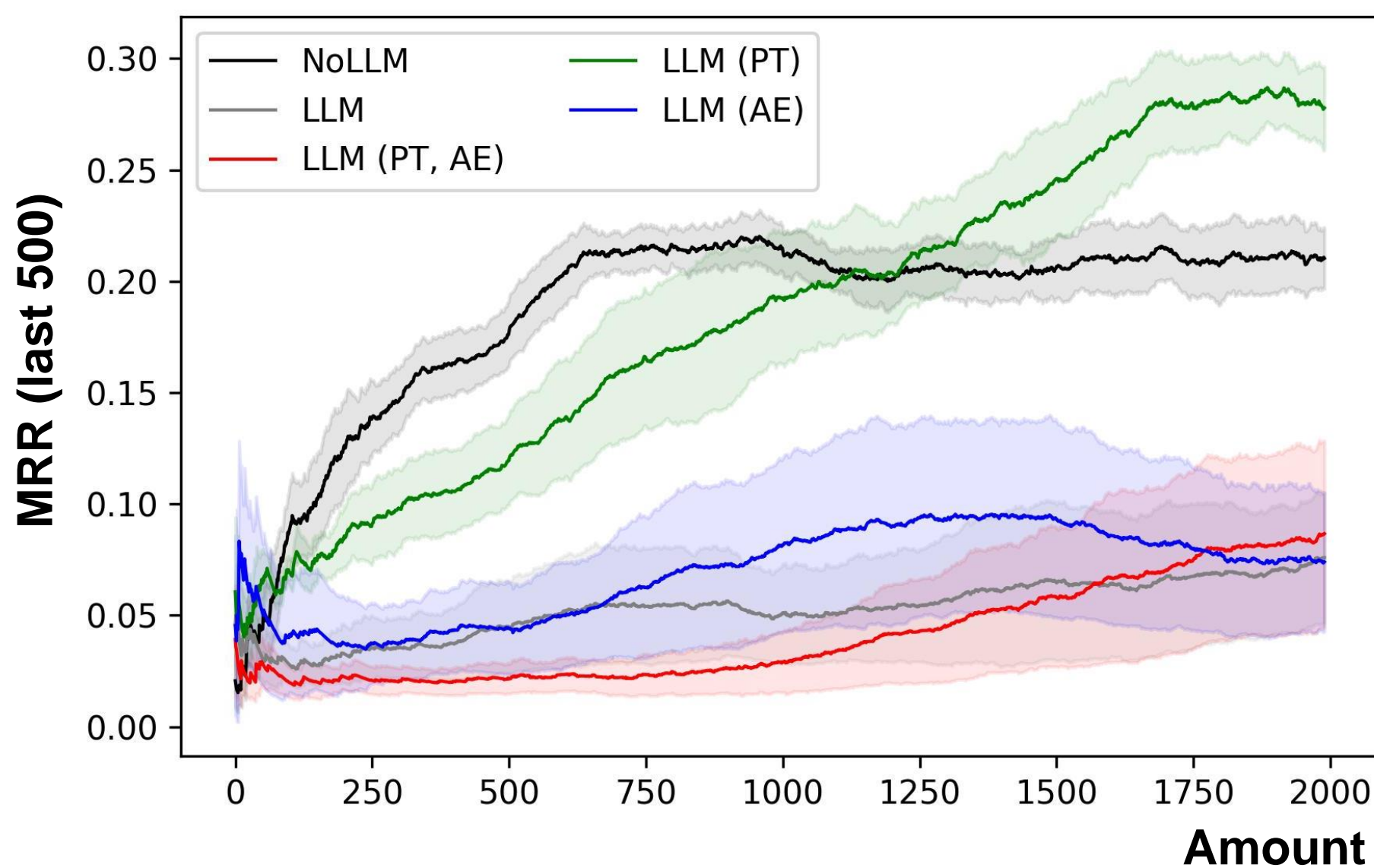
Query length: 4 terms  
Average of 5 runs with 80% confidence interval



**NoLLM:** linear predictor using term features  
**PT:** Prefix-Tuning with  $p = 5$ ; **AE:** Attribute Encoding

**CORD-19:** research records related to COVID-19  
**Local:** abstract (250,575 rows)  
**External:** sha, source\_x, paper\_title, doi, pmcid, ... (340,826 rows)

**WDC:** products scraped from various sites  
**Local:** category, brand, prod\_title, description, ... (57,109 rows)  
**External:** category, brand, prod\_title, description, ... (55,247 rows)



**Results**

**Prefix Tuning:**

- Small overhead
- Boosts long-run performance

**Attribute encoding:**

- Not as effective
- May harm performance

**More**

Models, techniques, results, source code and datasets