# Everlasting Iatric Researcher (Eir): Identifying the Article and Reading for Genetic Association Knowledge
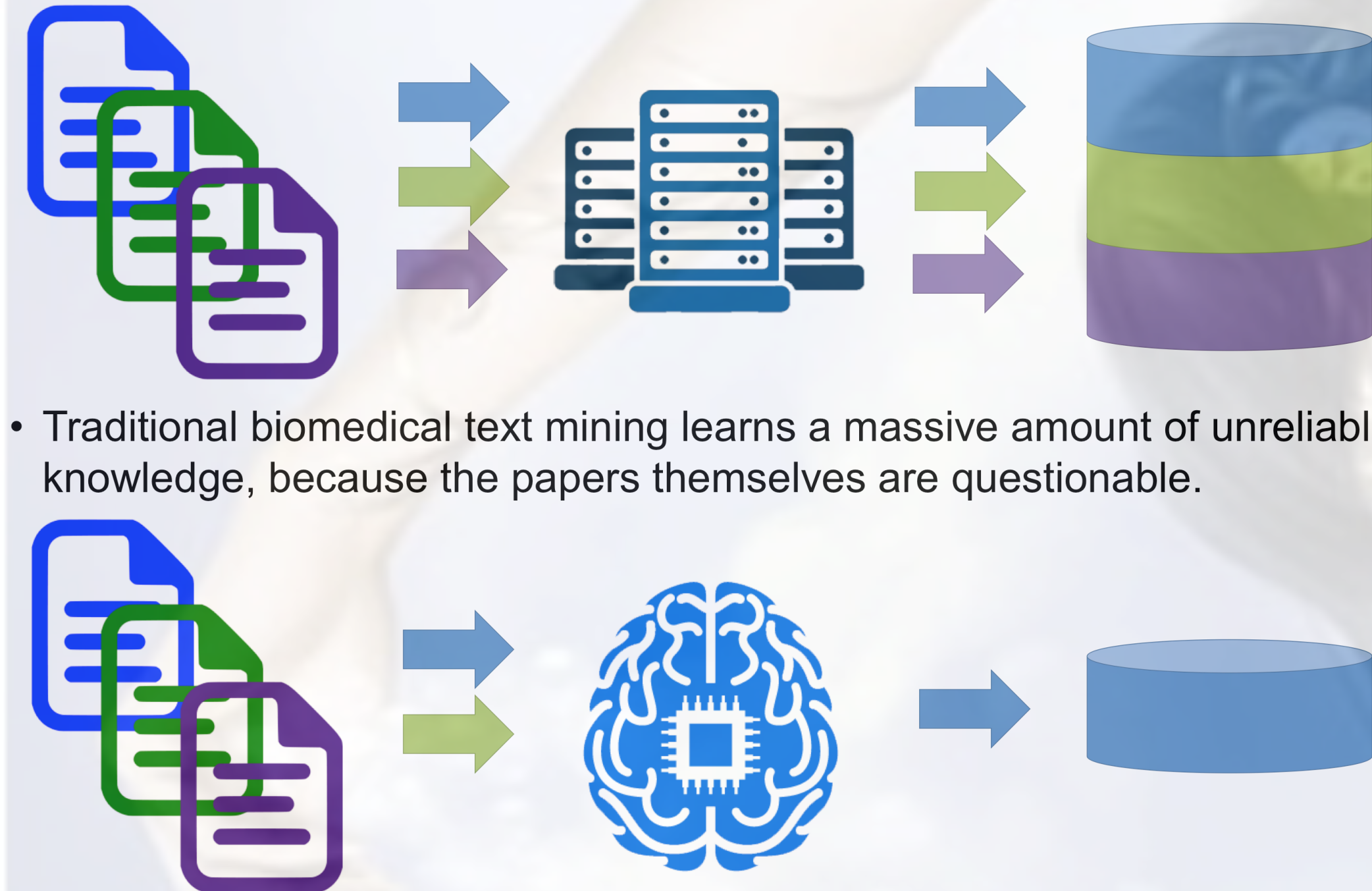
**Xiang Liu, Wenting Ye, Haohan Wang**

*Language Technologies Institute, School of Computer Science, Carnegie Mellon University*

## Contribution

- We propose to directly **simulate the behavior of a researcher** (i.e. selecting papers and reading for details) instead of universally examining the entire corpus.
- We develop a deep reinforcement learning model that can select both authentic and informative articles to read.
- We maintain a cutting-edge genetic association relationship database that can be easily queried.

## Introduction



- Traditional biomedical text mining learns a massive amount of unreliable knowledge, because the papers themselves are questionable.



- Eir, behaves like human, only reads the paper that she considers trustworthy, and constructs knowledge base accordingly.

## Deep Reinforcement Learning

Since our model involves a continuous state space $S$, we employ a deep Q-network (DQN) to learn the policy of agent with loss function:

$$L(\theta) = E_{\hat{s}, \hat{a}}[(y - Q(\hat{s}, \hat{a}; \theta))^2]$$

where $\quad y = r + \gamma \max_{a'} Q(\hat{s}', a'; \theta_t) \quad$, and $\;(\hat{s}, \hat{a}, \hat{s}', r)\;$ is selected transition.

- Learn the parameters $\theta$ of the DQN using stochastic gradient descent
- Use a (separate) target Q-network to calculate the expected Q-value for "stable update"
- Employ an experience replay memory D to store transitions

## Model Framework

Eir's research process is a Markov decision process (MDP), which can be represented as a tuple <S, A, T, R>, where $S = s$ is the space of all possible states, $A = a$ is the set of all actions, $R(s, a)$ is the reward function and $T(s'|s, a)$ is the transition function.
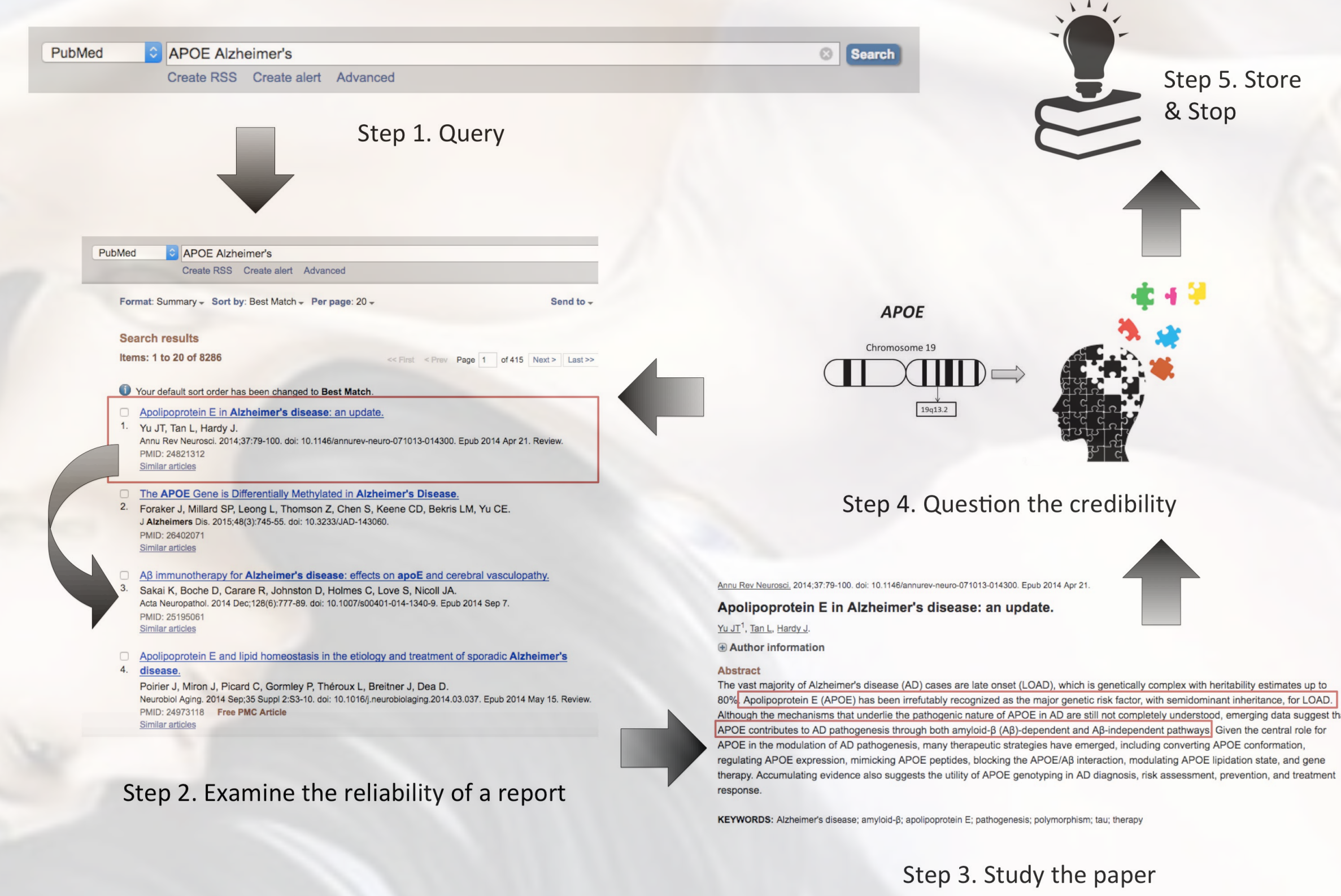


Figure 2: Overview of Eir's possible behaviors

**States:** The state $s$ in the MDP describes the research status of Eir. The information of the state depends on which status Eir is in.

| Eir's state | Information |
|---|---|
| Read the title | • Title<br>• Journal information<br>• Author information |
| Read the abstract | • Association entity<br>• Confidence<br>• Context information |

Table 1: Component of Eir's state in different scenarios

## Dataset

- Genetic Association Database (GAD) (Becker et al., 2004)
  - A manually crafted database of 142,000 high quality articles with the association it describes.
- PubMed
  - An online library which contains more than 27 million citations for biomedical literature

## Result and Future Work

### Observations

At current stage of this project, we obtain the first 4,000 gene-trait associations in GAD with 5,331 articles, and download 35,178 relevant articles from PubMed. The results are showed below.

| Entity | Correctness |
|---|---|
| Gene | 83.5% |
| Trait | 76.2% |

Figure 3a. Entity correctness

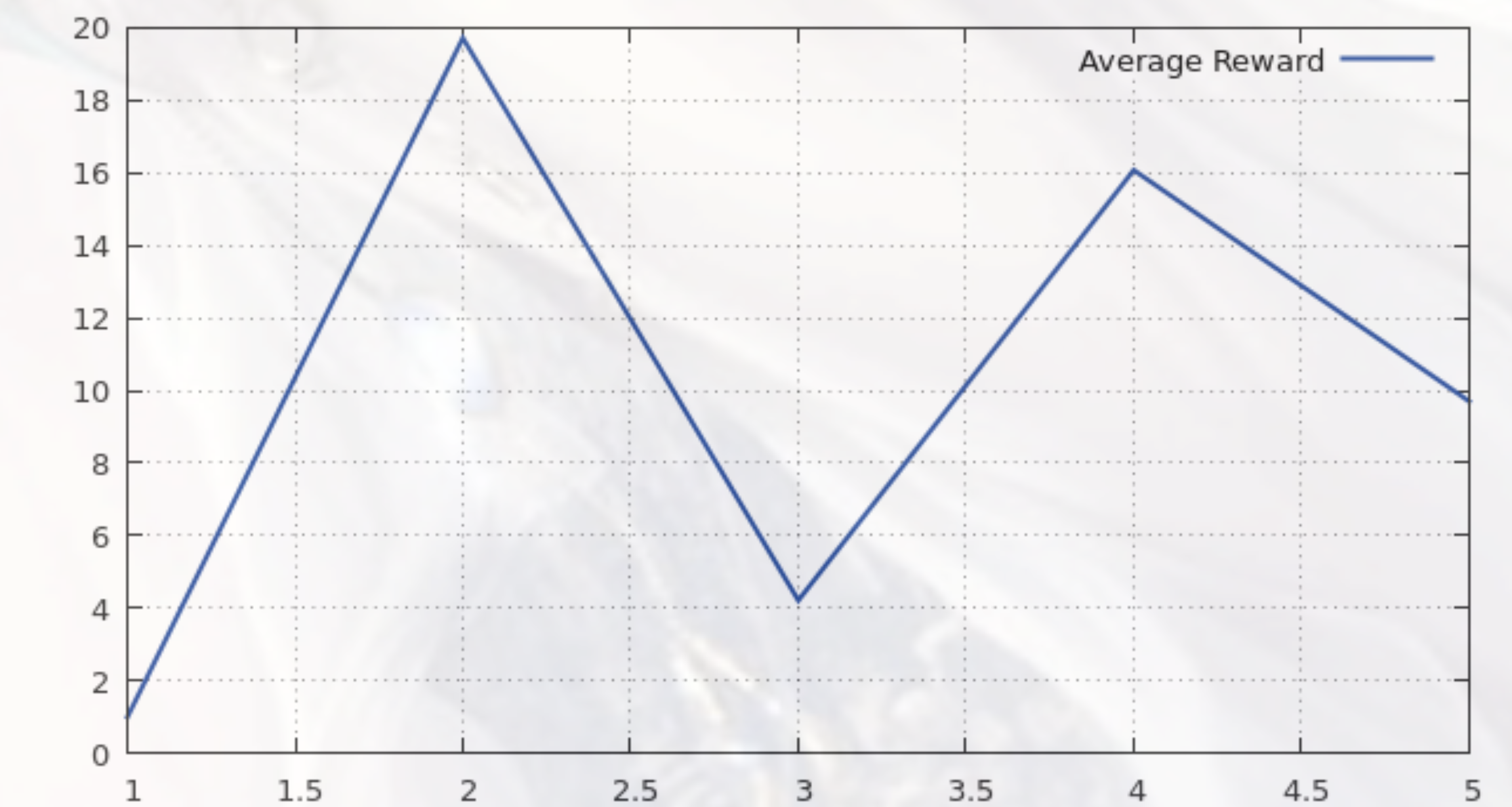| | Title | Abstract |
|---|---|---|
| Precision | 30.9% | 54.6% |
| Recall | 24.6% | 46.6% |
| $F_1$ | 0.261 | 0.50 |

Figure 3b. Selection accuracy



Figure 3c. Average reward

### Future works

- A larger amount of text resources
- More powerful text mining tools
- Construction of other biomedical knowledge database, i.e. gene-gene interaction

## Contact

- Haohan Wang (haohanw@cs.cmu.edu)

## Eir's Origin

Eir is the name of the goddess of medical knowledge and skills in Norse mythology. With current technology, we are hoping to realize the mythology in this modern world.