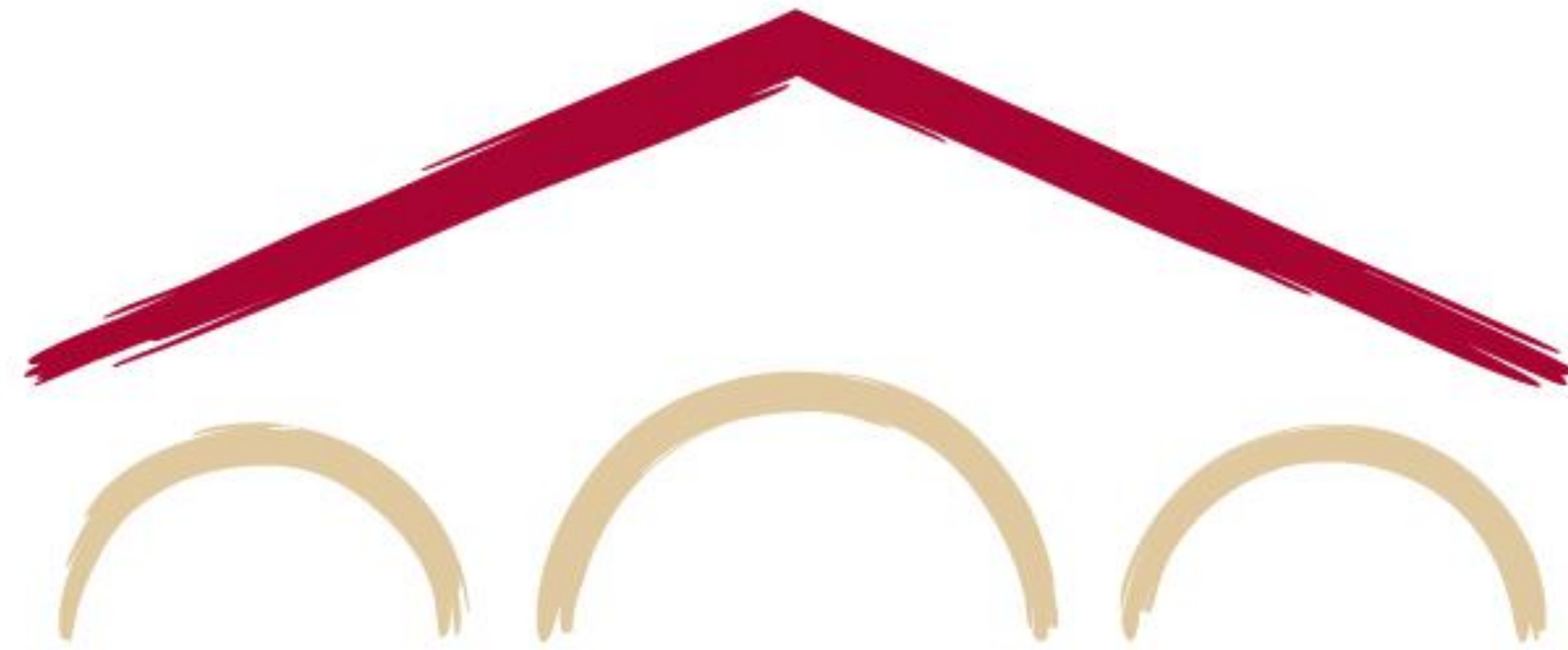# Meaning and Intelligence in Language Models: From Philosophy to Agents in a World
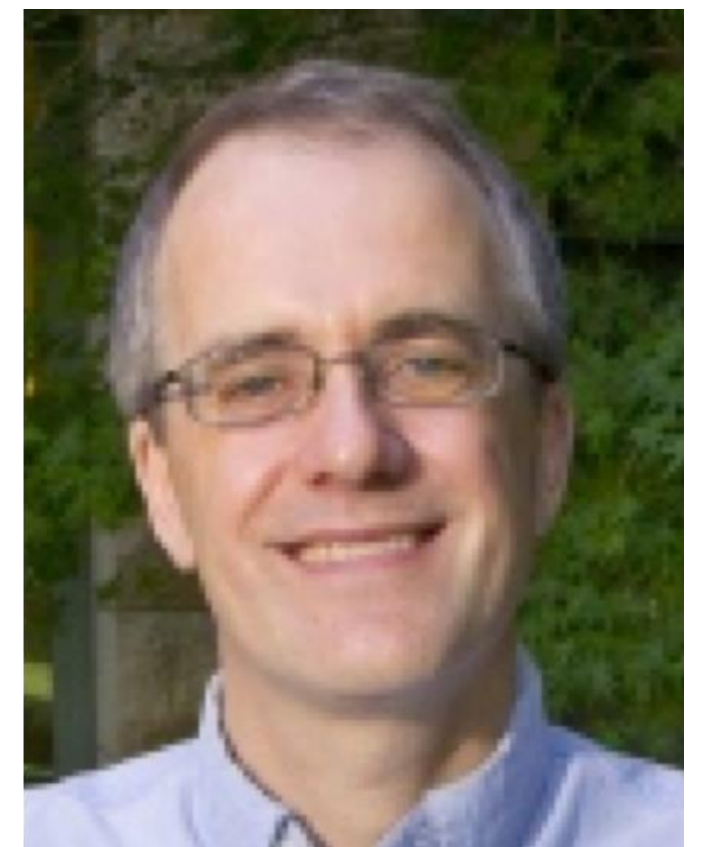
Christopher Manning

@chrmanning ❀ manning@stanford.edu

Stanford NLP Group ❀ Director, Stanford AI Lab

COLM, Philadelphia, 2024

# Who should you take your LLM wisdom from in 2024?

or

\* Calculated median survival from Gompertz equation, with the assumption that $a$ is 0.05, and $R_o$ extrapolated to birth 2.45 x 10$^{-5}$.

$$S = e^{(R_o/a)(1-e^{a^t})}$$

References

1. Gavrilov and Gavrilova NS. The Biology of Life Span, A Quantitative Approach. Harwood Academic (1991).
2. Brody S. J. Gen. Physiol. 6 (1923) 245-257.
3. Jones HB. Adv.Biol. Med. Physiol. 4 (1956) 281-337.
4. Hallén A. Acta Chem. Scand. 62 (1962) 705-710.
5. " Biogeront. 3 (2002) 307-315.
6. " Biogeront. 8 (2007) in press.
7. Laurent TC, Killander J. J. Chromatogr. 14 (1964) 317-330.
8. Mattstal JT, van Oudenaarden A. Science 317 (2007) 463-464.
9. Finkel T, Serrano M, Blasco MA. Nature 2007;448:767-774.
10. McFarland GA, Holliday R. Exptl. Cell Res. 1994; 212:167-175.
11. Hopkins AB et al. Mech. Ageing Dev. 122 (2001) 1431-1445.
12. Verbeke et al. Biogeront. 3 (2002) 117-128.
13. Olsen et al. Biogeront. 7 (2006) 221-230.
14. Chondrogianni N et al. Exp. Gerontol. 35 (2000) 721-728.
15. Chondrogianni N, Gonos ES. Biogeront. 5 (2004) 55-61.

Distributional Clustering of English Words

Pereira    Tishby    Lee

Black
Lafferty
Mercer

Jelinek
Magerman
Roukos

Towards History-based Grammars: Using Richer Models for Probabilistic Parsing

**Why was no one working on language models?**

**How were they hoping to get to artificial intelligence?**

A. A. Марков (1886).

Andrey A. Markov
1856–1922

Transition probabilities (Markov chains) between consonants and vowels in Alexander Pushkin's verse novel, *Eugene Onegin*

An Example of Statistical Investigation of the Text *Eugene Onegin* Concerning the Connection of Samples in Chains
1913

Claude E. Shannon
1916–2001

A Mathematical Theory of Communication
*The Bell System Technical Journal* 1948

Prediction and Entropy of Printed English
*The Bell System Technical Journal* 1951

Explores character and word-level n-gram models, with estimation and generation from a small text sample

# The term/concept of a **Language Model**

# Invented in Frederick Jelinek's group, IBM

Design of a Linguistic Statistical Decoder for the Recognition of Continuous Speech

*IEEE Transactions on Information Theory* 21(3), 1975

It was work from this group at IBM in the 1970s that defined the **probabilistic language model** of next token prediction that has continued to dominate till today

The bigram "language model" appears earlier, particularly in translations of Russian works. E.g., Igor Mel'chuk (1961) "Some Problems of Machine Translation Abroad" refers to Chomsky's " 'language model' of immediate constituents". Jelinek may have been influenced by this Russian idiom?

# The history of modern (neural) AI/ML that many think of

1950–2004    Experiments on perception and vision,
with minute amounts of data

2005–2009    Small modern-style labeled vision datasets
2005–2012 PASCAL VOC 4–20 classes, 15,000 images
2006  Caltech 101 – ~5000 images, 101 classes, ~300x200
2007  Caltech 256 – 30,607 images of 256 classes
2009  CIFAR 10/100 – 600 images of 10/100 classes, 32x32

2010–2017    **ImageNet!** 1–14 million images, 1000 classes
Big data! Vision object classification is "solved"

standard schnauzer  giant schnauzer

Persian cat   Siamese cat

# The History of Large Language Models

1998    CPAT-Tree-Based Language Models with an Application for Text Verification in Chinese. ROCLing 1998. First use of "LLM" trigram I know of; 200M word corpus

2000    A Neural Probabilistic Language Model. Bengio, Ducharme & Vincent NIPS 2000 First neural language model built on 32 million token corpus, 31K vocab

2007    Large Language Models in Machine Translation. Brants, Popat, Xu, Och and **Dean**. EMNLP 2007. **2 trillion token corpus** n-gram model of up to 5-grams

2018    **GPT** (Radford, Narasimhan, Salimans & Sutskever) and **BERT** (Devlin, Chang, Lee & Toutanova). 3.3 billion token corpus

2020–    100+ billion parameter neural language models trained on > 1 trillion tokens: GPT-3, GPT-4, PaLM 2, Llama 3, Nemotron-4, ….

# The History of Large Language Models

1998    CPAT-Tree-Based Language Models with an Application for Text Verification in Chinese. ROCLing 1998. First use of "LLM" trigram I know of; 200M word corpus

2000    A Neural Prob**Not enough compute!**Vincent NIPS 2000
First neural la31K vocab

2007    Large Language M**Not enough model flexibility!**Och and
**Dean**. EMNLP 20075-grams

2018    **GPT** (Radford, Narasimha**Not enough data!**ng,
Lee & Toutanova). 3.3 bill

2020–   100+ billion parameter neural la**LLMs go brrr!**
GPT-3, GPT-4, PaLM 2, Llama 3,

# LLMs need all the trinity of modern deep learning

# Form vs. reference: Meaning as a relation between form and "communicative intent"

Bender and Koller (2020): "(linguistic) meaning [is] the relation between a linguistic form and communicative intent."

horse    ⟺    

"the language modeling task, because it only uses form as training data, cannot in principle lead to learning of meaning"

"You say: the point isn't the word, but its meaning, and you think of the meaning as a thing of the same kind as the word, though also different from the word. Here the word, there the meaning. The money, and the cow that you can buy with it. (But contrast: money, and its use.)"

– Ludwig Wittgenstein. 1953. *Philosophical Investigations,* §120

# Form vs. reference: Meaning as a relation between form and the world

The classic referentialist position, which has dominated recent linguistic semantics, attempts to ground meaning in the world:

horse   $\Leftrightarrow$   

This might suggest that a plain language model has no notion of meaning, but a visual language model does?

Somehow, a primacy is given to certain kinds of sensory evidence, even though everything ends up as electrical signals in our brains.

Would we really want to say that a VLM captures meaning but a pure LLM does not?

I personally think not.

# The alternative: Inferentialism/pragmatism

"For a large class of cases—though not for all—in which we employ the word 'meaning' it can be defined thus: the meaning of a word is its use in the language."

– Wittgenstein, 1953, §43

Inferentialism couches understanding as an appreciation of appropriate vs. inappropriate ways to use language in a context. Often summed up as: **a use theory of meaning**

This context might be visual, textual, conversational, sociological, etc.

The real world may in some sense be privileged, but it is not the only context for grounding meaning: We also have virtual worlds and language contexts

This is essentially the notion of meaning that Turing (1950) also argued for, if perhaps partly for practical reasons

Under this notion, meaning is gradient: You can partially understand appropriate contexts of use

For more, see i.a., Potts, Icard, Portelance, Card, Zhou, & Etchemendy (2021 in On the Opportunities and Risks of Foundation Models

# What is the meaning of 'shehnai'

1. You can have seen or held one and have a classic grounded meaning

2. I can at least show you a picture 🖥

3. But surely the meaning that you have for the word 'shenai' would be richer if you had heard one played?

4. If I have never seen, felt, or heard a shehnai, but some traditional Indian instrument, a bit like an oboe, then surely the word has some meaning for me?

   • A connection to India, to wind instruments that use reeds, to playing music

It gives me a reasonable approximation of appropriate contexts of use

5. If they added that it has holes sort of like a recorder, but it has multiple reeds and a flared end more like an oboe, then I have more "meaning"

Manning, C. D. 2022. Human Language Understanding and Reasoning. *Dædalus* 151.

# What is the meaning of 'shehnai'

- What if all I have is one or more contexts of word use (like an LLM commonly gets)?
  - *"From a week before, shehnai players sat in bamboo machans at the entrance to the house, playing their pipes. Bikash Babu disliked the shehnai's wail, but was determined to fulfil every conventional expectation the groom's family might have."*
    - From: Anuradha Roy, *An Atlas of Impossible Longing* (New York: Free Press, 2011)
- In some ways, I understand the meaning of the word 'shehnai' rather less than someone who has seen one, but:
  - I still know that it is a pipe-like musical instrument
  - My meaning of 'shehnai' is not a subset of the meaning of the person who has only held a shehnai, for I know some additional musical and cultural information on likely appropriate use that they lack

# Large Language Models and "Intelligence"

Since ChatGPT, in the popular imagination, large language models have essentially become synonymous with **artificial intelligence**

The **artificial** part is easy. What is intelligence?

Is language understanding and generation distinct from general intelligence?

CA AB 2885

**Section 11546.45.5 of the California Government Code (2024-09-28):**

"Artificial intelligence" means an engineered or machine-based system that varies in its level of autonomy and that can, for explicit or implicit objectives, infer from the input it receives how to generate outputs that can influence physical or virtual environments.

an engineered or machine-based system that varies in its level of autonomy and that can, for explicit or implicit objectives, infer from the input it receives how to generate outputs that can influence physical or virtual environments

# Intelligence vs. knowledge or task effectiveness

More commonly …

We confuse being knowledgeable with being intelligent

Or the ability to do tasks?

AGI: "ability to perform most (cognitive/ economically valuable) tasks at or above the level of the average person"

# Intelligence

"an agent's ability to achieve goals in a wide range of environments"
– Shayne Legg and Marcus Hutter (2006)


"the ability to **learn** and perform a range of techniques to solve problems and achieve goals, **appropriate to the context in an uncertain, ever-varying world**"
– Christopher Manning

# Some amazing reasoning examples now work in LLMs!

## openai/text-davinci-002

The cost to produce a book is 1000 dollars plus 8 dollars per book. The book sells for 20 dollars each. How many books must be sold to make a profit?

1250 books must be sold to make a profit.

The cost to produce a book is 1000 dollars plus 8 dollars per book. The book sells for 20 dollars each. How many books must be sold to make a profit?Let's work step by step.
The cost of producing 1 book is 1000 + 8 = 1008 dollars.
The cost of producing 2 books is 1008 + 8 =1016 dollars.
Therefore, the cost of producing n books is:
Cost = 1008 + 8n
The selling price of 1 book is $20.
The selling price of 2 books is $40.
The selling price of n books is:
Selling Price = 20n
We want to find how many books must be sold in order to make a profit. That is, we want to find n such that:
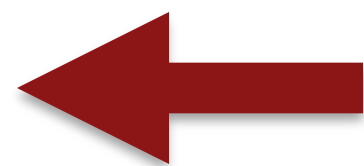Cost < Selling Price
1008 + 8n < 20n
8n < 20n - 1008
-1008 < 12n
1000 < 12n
83.33 < n
Therefore, n must be greater than 83.33 in order to make a profit.

To determine how many books must be sold to make a profit, we need to figure out when the revenue from selling the books exceeds the total cost of producing them.

Let:
- $x$ be the number of books sold.
- The **cost to produce** $x$ books is:

$$\text{Total cost} = 1000 + 8x$$

(where 1000 dollars is the fixed cost, and 8 dollars is the variable cost per book).

- The **revenue** from selling $x$ books is:

$$\text{Revenue} = 20x$$

(where each book is sold for 20 dollars).

**To make a profit, the revenue must be greater than the total cost:**

$$20x > 1000 + 8x$$

**Solving the inequality:**

1. Subtract $8x$ from both sides:

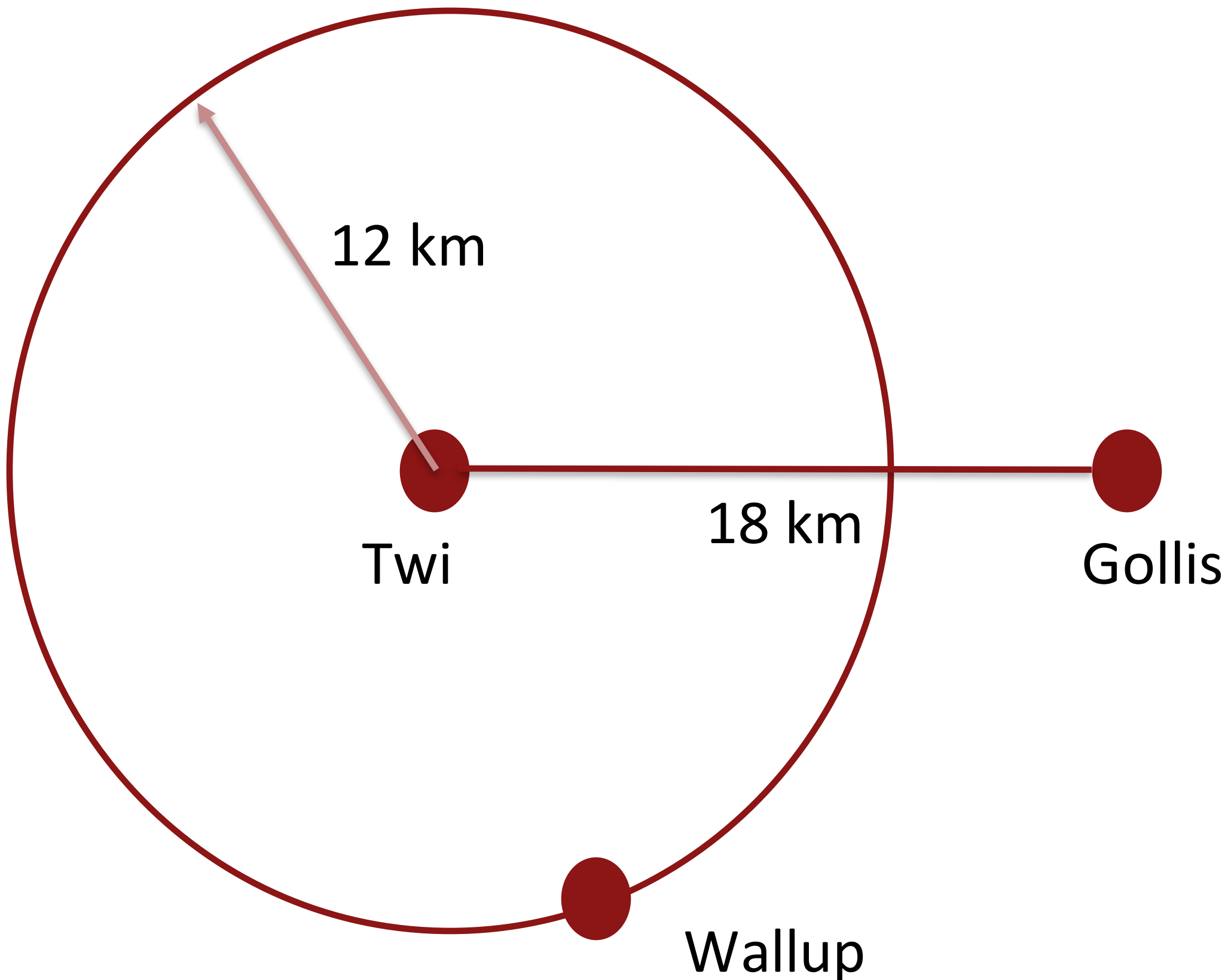$$20x - 8x > 1000$$
$$12x > 1000$$

2. Divide both sides by 12:

$$x > \frac{1000}{12} \approx 83.33$$

Since $x$ must be a whole number, you need to sell at least **84 books** to make a profit.

If the town of Twi is 12 km from the town of Wallup, and the town of Gollis is 18 km from the town of Twi, what is the maximum and minimum possible distance the town of Wallup is from Gollis?

12 km

18 km

Twi

Gollis

Wallup

To find the **maximum** and **minimum** possible distances between **Wallup** and **Gollis**, let's consider the two possible configurations for the positions of the towns.

**1. Maximum Possible Distance:**

For the maximum distance between Wallup and Gollis, Wallup and Gollis would be on **opposite sides of Twi** in a straight line. In this case, the total distance would be the sum of the two distances:

- **Wallup to Twi** = 12 km
- **Twi to Gollis** = 18 km

So, the maximum possible distance is:

$$\text{Maximum distance} = 12\,\text{km} + 18\,\text{km} = 30\,\text{km}$$

**2. Minimum Possible Distance:**

For the minimum distance, Wallup, Twi, and Gollis would be in a straight line, with Gollis on the **same side** of Twi as Wallup. In this case, the minimum possible distance would be the difference between the two distances:

- **Wallup to Twi** = 12 km
- **Twi to Gollis** = 18 km

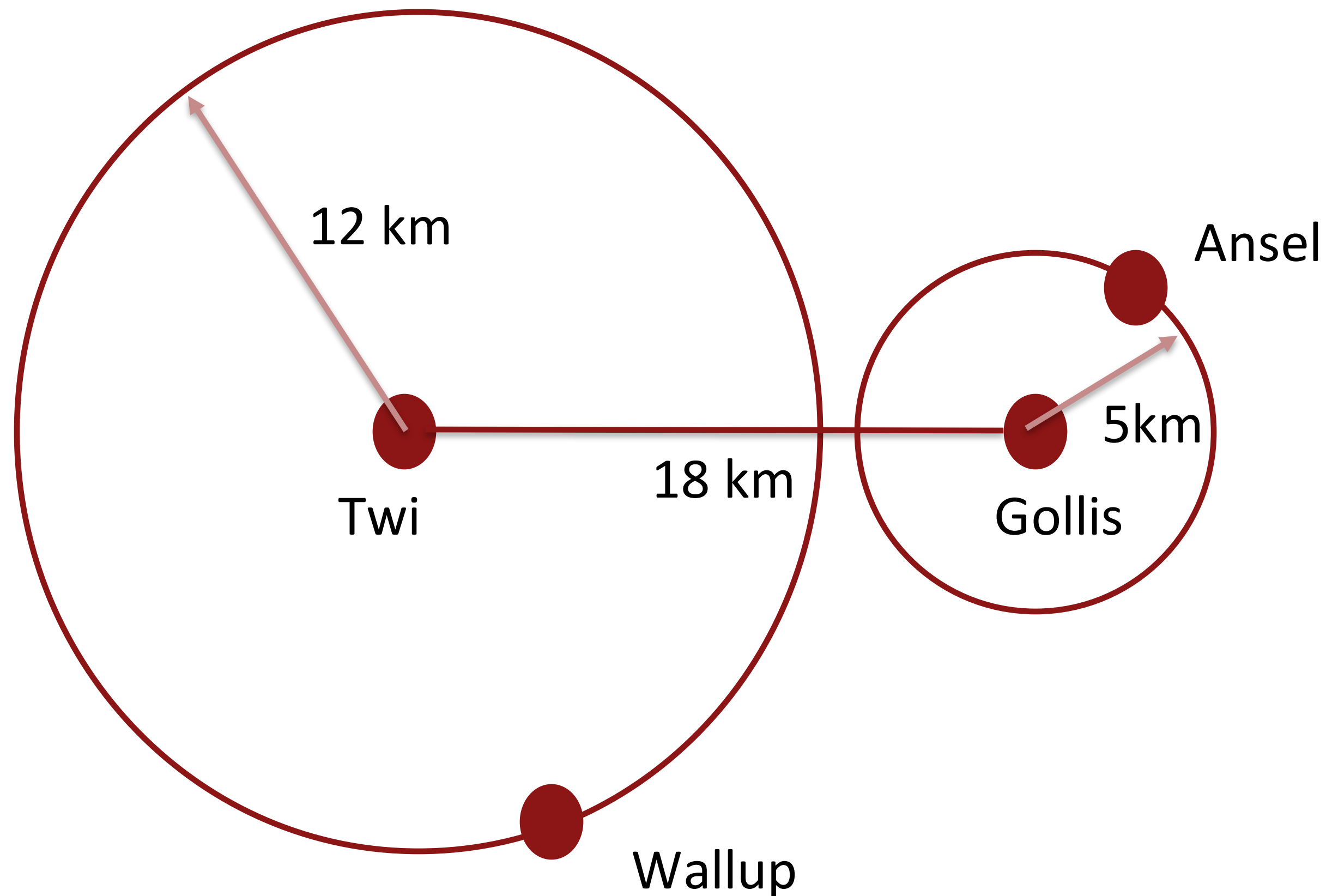The minimum possible distance is:

$$\text{Minimum distance} = |18\,\text{km} - 12\,\text{km}| = 6\,\text{km}$$

**Conclusion:**

- The **maximum distance** between Wallup and Gollis is **30 km**.
- The **minimum distance** between Wallup and Gollis is **6 km**.

If the town of Twi is 12 km from the town of Wallup, the town of Gollis is 18 km from the town of Twi, and the town of Ansel is 5 km from the town of Gollis, what is the maximum and minimum possible distance the town of Ansel is from Gollis?



12 km

18 km

Twi

Ansel

5km

Gollis

Wallup

Let's explore the problem step by step:

You want to know the **maximum** and **minimum** possible distances between the town of **Ansel** and the town of **Wallup**. We know the following distances:

- **Wallup to Twi** = 12 km
- **Twi to Gollis** = 18 km
- **Gollis to Ansel** = 5 km

**Maximum Possible Distance:**

For the **maximum possible distance** between **Ansel** and **Wallup**, the towns would need to be laid out in a straight line with Ansel on the opposite side of Gollis from Wallup. This would give us the longest possible distance by adding all the segments together:

- **Wallup to Twi** = 12 km
- **Twi to Gollis** = 18 km
- **Gollis to Ansel** = 5 km

So, the maximum possible distance from **Ansel to Wallup** is:

$$\text{Maximum distance} = 12\,\text{km} + 18\,\text{km} + 5\,\text{km} = 35\,\text{km}$$

**Minimum Possible Distance:**

For the **minimum possible distance**, Ansel would be on the same side of Gollis as Wallup, and we would subtract the distance between Ansel and Gollis from the other distances:

- **Wallup to Twi** = 12 km
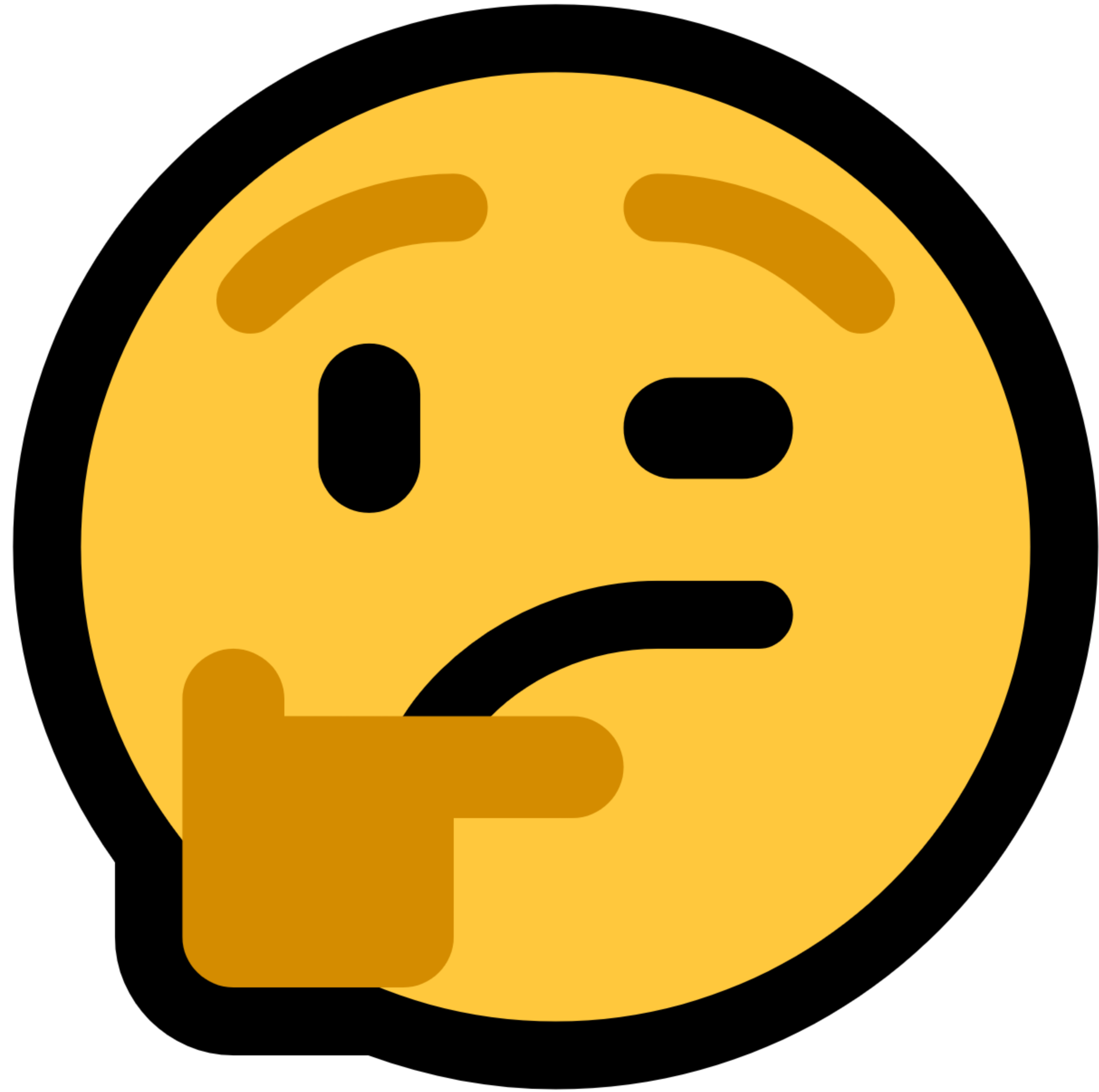- **Twi to Gollis** = 18 km
- **Gollis to Ansel** = 5 km

The minimum possible distance is:

$$\text{Minimum distance} = (18\,\text{km} - 5\,\text{km}) + 12\,\text{km} = 25\,\text{km}$$

**Conclusion:**

- The **maximum distance** from **Ansel to Wallup** is **35 km**.
- The **minimum distance** from **Ansel to Wallup** is **25 km**.

# 1st Workshop on System-2 Reasoning at Scale

December 15, 2024. **NeurIPS** Workshop, Vancouver, Canada

➡ What do we need to imbue language models with System-2 reasoning capabilities?
➡ Are scale and the "bitter lesson" going to dictate how the future of AI technology will unfold?
➡ Do we need a different mechanism for implementing System-2 reasoning, or should it be a property that emerges from a possibly different training method?
➡ How do we benchmark System-2-like generalization? How do we avoid data contamination?

## Speakers + Panelists

Melanie Mitchell  Joshua Tenenbaum  François Chollet  Dzmitry Bahdanau  Tal Linzen  Jason Weston

## Organizers

Shikhar Murty  Federico Bianchi  Róbert Csordás  Nouha Dziri

Alex Gu  Shunyu Yao  Chris Manning  Yejin Choi

# Learning Grounded Instruction Following with Language Models

Shikhar Murty

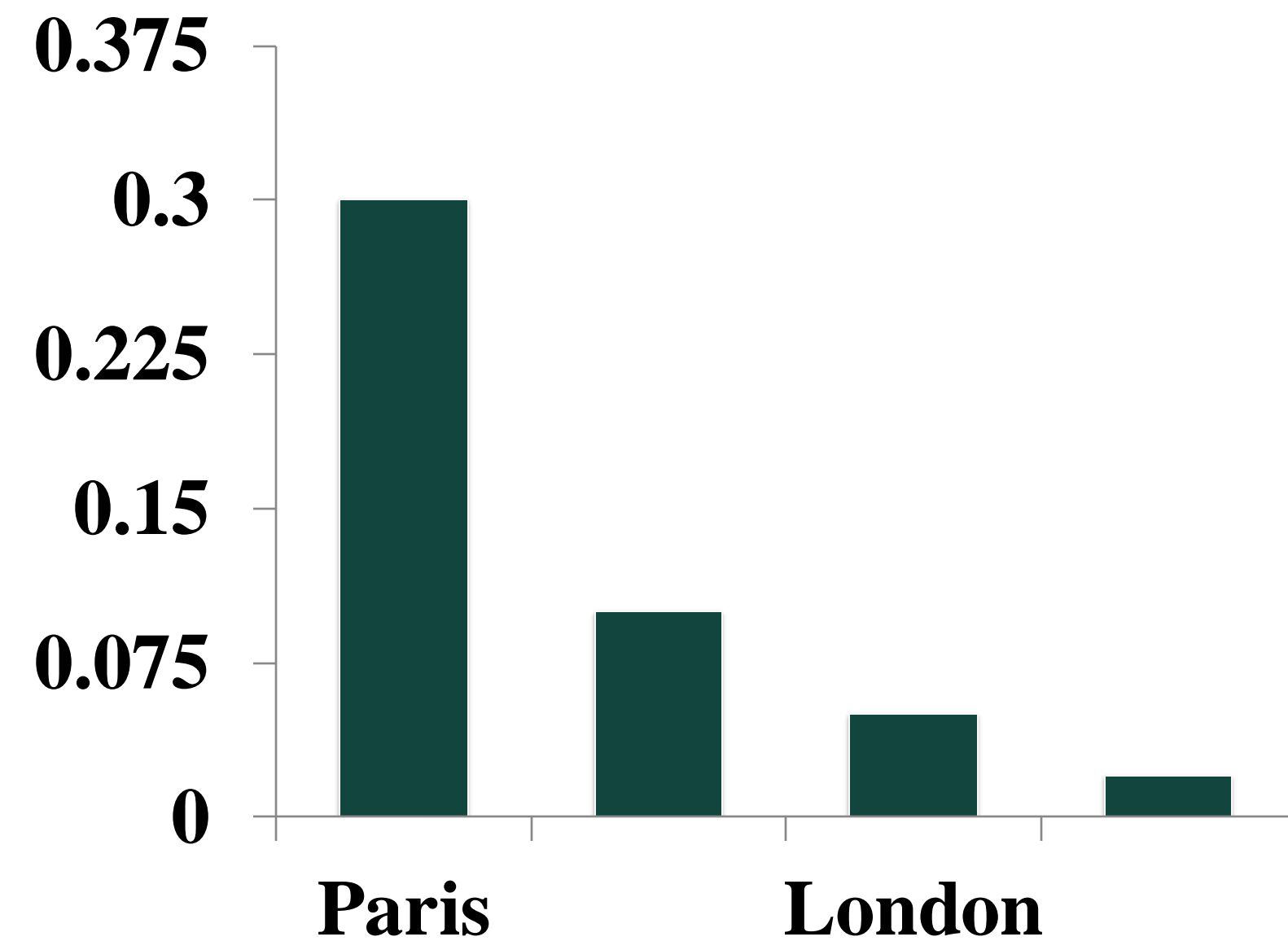# Grounded Instruction Following with Language Models?

The capital of Paris is ____

Language Model

# Grounded Instruction Following with Language Models?

## First Attempt: Zero-shot instruction following



The capital of Paris is _____

Language Model

$\pi_{\mathrm{LM}}(\cdot \mid g)$

Language
Model
Policy

```
You are an agent capable of the
following actions:
1.Type X on Y
2.Move mouse to X
3.Click on X
4.Hover on element X
5. …
Your objective is to follow user instructions,
by mapping them into a sequence of actions.
Instruction: {g}

So far, you have taken the following actions:
a1:
a2:
…
After executing these actions, you observe
the following HTML state: <HTML state>

Now, think about your next action:
Thought: [model-pred]

Now, take an action:
Action: [model-pred]
```

Yao et al. 2023. WebShop.

# Grounded Instruction Following with Language Models?

## First Attempt: Zero-shot instruction following

*Cancel shampoo subscription from the list of "subscribe and save" items*



$$\longrightarrow \boxed{\pi_{\mathrm{LM}}(\cdot \mid g)} \longrightarrow \textbf{???}$$
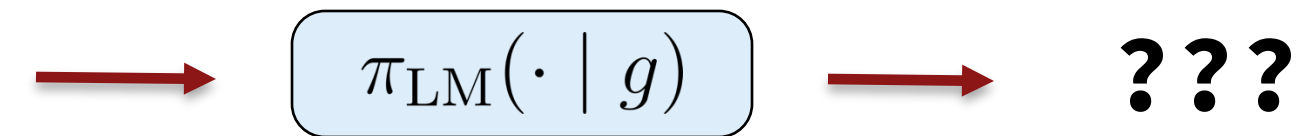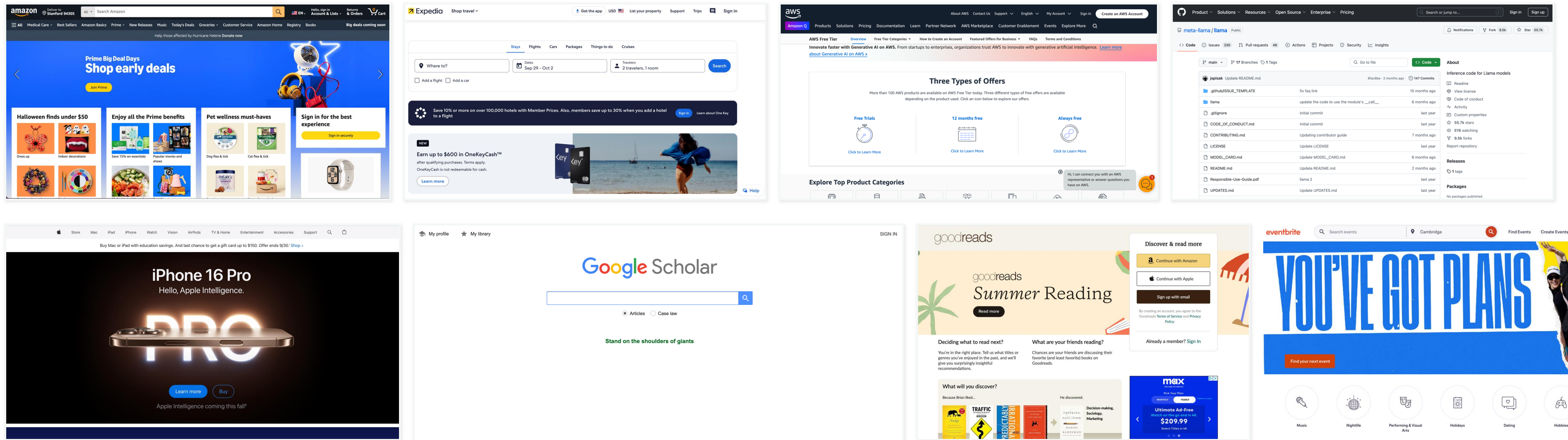
LMs cannot know how to navigate to the right section without prior knowledge of amazon.com
How can we provide this prior knowledge about environments?

# Grounded Instruction Following with Language Models?

## First Attempt: Zero-shot instruction following

*Cancel shampoo subscription from the list of "subscribe and save" items*

$$\pi_{\mathrm{LM}}(\cdot \mid g)$$

**???**

# Grounded Instruction Following with Language Models?

## First Attempt: Zero-shot instruction following

*Cancel shampoo subscription from the list of "subscribe and save" items*

$\pi_{\mathrm{LM}}(\cdot \mid g)$

**???**

# Grounded Instruction Following with Language Models?

## Second Attempt: Behavioral Cloning with Human demonstrations



Collecting human demonstrations that cover most use-cases across websites is expensive
Humans learn about new web-interfaces simply by interacting with them.
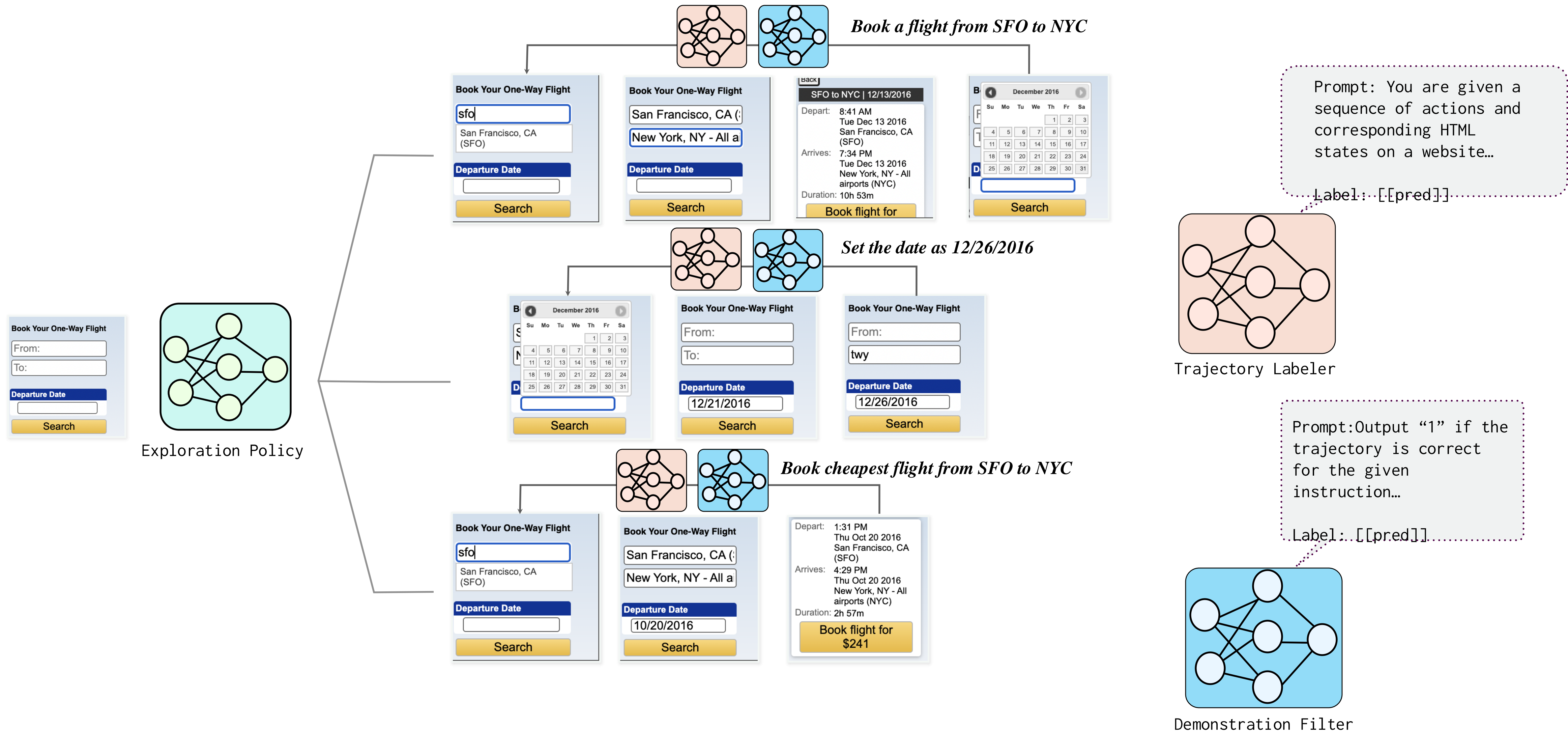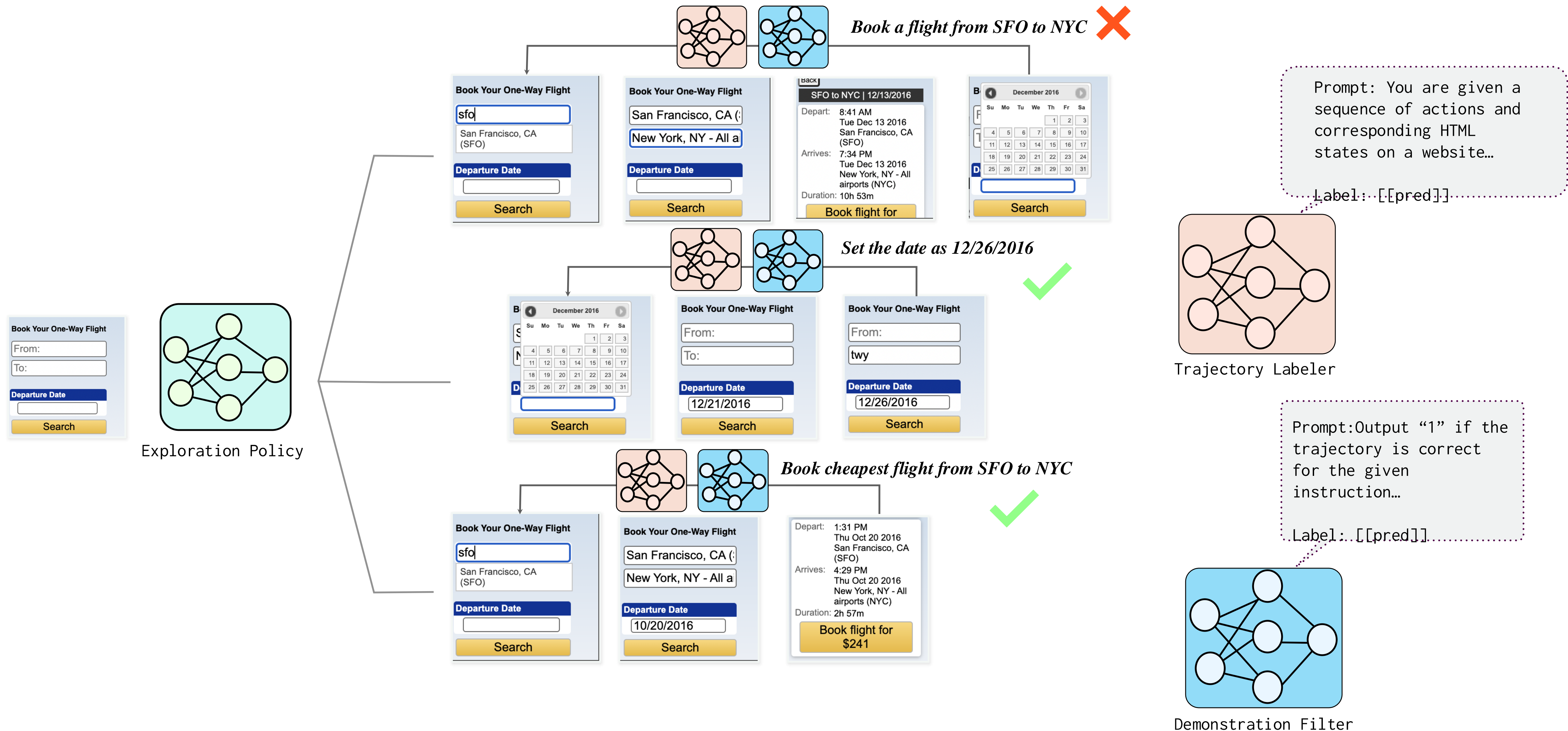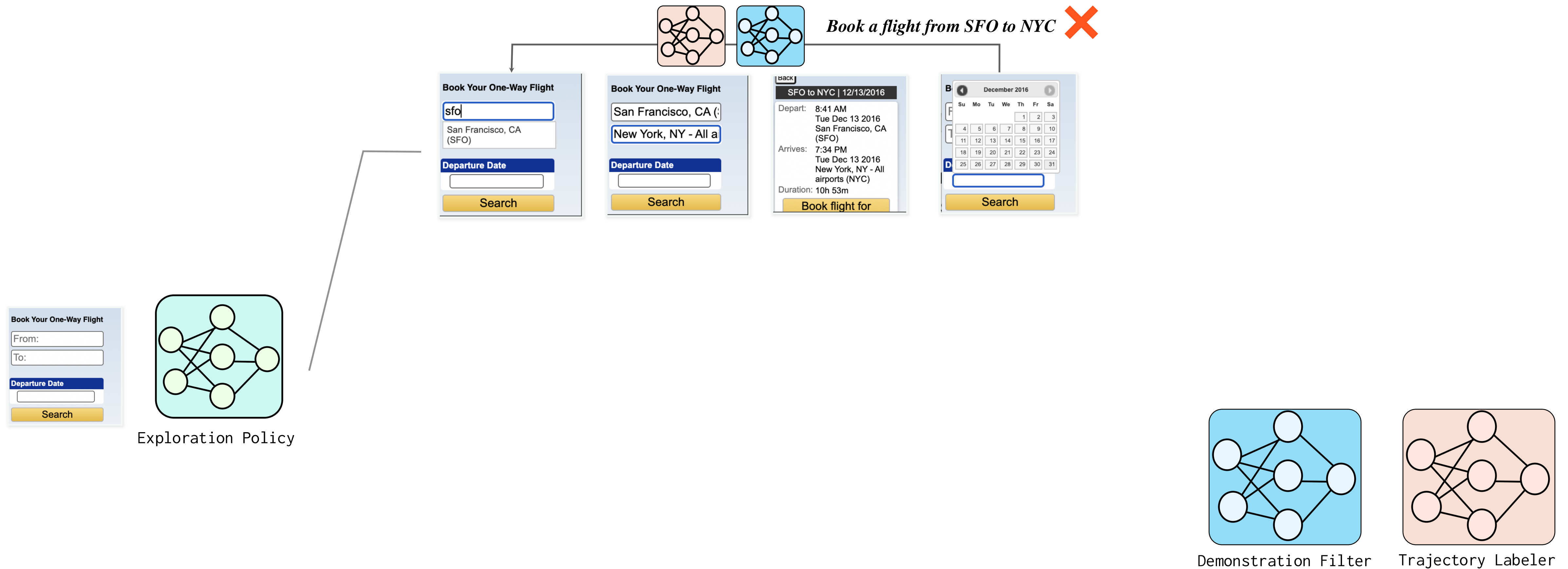Can we do the same with models?

# Grounded Instruction Following with *Synthetic* Demonstrations

## BAGEL: Our first attempt for generating synthetic demonstrations via interaction

Prompt: Given a website, take actions of the following format to explore….

Action: [[pred]]

Exploration Policy

# Grounded Instruction Following with *Synthetic* Demonstrations

## BAGEL: Our first attempt for generating synthetic demonstrations via interaction



How can we decide if a sequence of interactions is meaningful?

Use Natural Language!

# Grounded Instruction Following with *Synthetic* Demonstrations

## BAGEL: Our first attempt for generating synthetic demonstrations via interaction



Prompt: You are given a sequence of actions and corresponding HTML states on a website…

Label: [[pred]]

Exploration Policy

Trajectory Labeler

How can we decide if a sequence of interactions is meaningful?

Use Natural Language!

# Grounded Instruction Following with *Synthetic* Demonstrations

## BAGEL: Our first attempt for generating synthetic demonstrations via interaction



*Book a flight from SFO to NYC*

*Set the date as 12/26/2016*

*Book cheapest flight from SFO to NYC*

Exploration Policy

Trajectory Labeler

Prompt: You are given a sequence of actions and corresponding HTML states on a website…

Label: [[pred]]

# Grounded Instruction Following with *Synthetic* Demonstrations

## BAGEL: Our first attempt for generating synthetic demonstrations via interaction



*Book a flight from SFO to NYC*

*Set the date as 12/26/2016*

*Book cheapest flight from SFO to NYC*

Exploration Policy

Prompt: You are given a sequence of actions and corresponding HTML states on a website…

Label: [[pred]]

Trajectory Labeler

Prompt:Output "1" if the trajectory is correct for the given instruction…

Label: [[pred]]

Demonstration Filter

# Grounded Instruction Following with *Synthetic* Demonstrations

## BAGEL: Our first attempt for generating synthetic demonstrations via interaction



41

# Grounded Instruction Following with *Synthetic* Demonstrations

## BAGEL: Our first attempt for generating synthetic demonstrations via interaction



*Book a flight from SFO to NYC* ❌

Exploration Policy

Demonstration Filter    Trajectory Labeler

Instead of throwing away bad (instruction, trajectory) pairs, BAGEL refines iteratively by going back and forth between language and actions

# Grounded Instruction Following with *Synthetic* Demonstrations

## BAGEL: Our first attempt for generating synthetic demonstrations via interaction



Prompt: Map the given
instruction to a sequence
of actions, one at a time.

Thought: [[pred]]
Action: [[pred]]

*Book a flight from SFO to NYC*

*Set Origin to SFO and destination to NYC*

Demonstration Filter    Trajectory Labeler

Instead of throwing away bad (instruction, trajectory) pairs, BAGEL refines iteratively by going back and forth between language and actions

# Grounded Instruction Following with *Synthetic* Demonstrations

**BAGEL: Our first attempt for generating synthetic demonstrations via interaction**



Instead of throwing away bad (instruction, trajectory) pairs, BAGEL refines iteratively by going back and forth between language and actions

# Some Results from In-Context Learning with BAGEL demonstrations

MiniWoB++: map instructions to sequence of mouse and keyboard actions



Sandboxed environment evaluating basic browser interactions across a range of applications from social media to email clients

Evaluates functional correctness

Not real world (limited functionality)
Relatively short-horizon
Zero-shot performance far from perfect!

# Some Results from In-Context Learning with BAGEL demonstrations

(↑ is better)

| Domain | Zero-shot | + BAGEL |
|---|---|---|
| book-flight-nodelay | 0.05 | 0.15 |
| email-inbox-all | 0.88 | 1.00 |
| choose-date-nodelay | 0.2 | 0.4 |
| social-media | 0.6 | 0.7 |
| … | … | … |
| Mean Reward | 0.47 | 0.61 |

BAGEL improves zero-shot agents significantly

# Some Results from In-Context Learning with BAGEL demonstrations

Baseline: *Instruction-first* approach
- ▸ Sample instructions from a prior, given first observation from environment
- ▸ Run zero-shot agent over these instructions
- ▸ Use the scoring function to filter out demonstrations

```
Prompt: Given the
initial page of a
website, generate
possible instructions…


Instruction: [[pred]]
```

Instruction Generator  ------>  LM policy  ------>  Demonstration Filter

# Some Results from In-Context Learning with BAGEL demonstrations

(↑ is better)

| Domain | Instruction-First | BAGEL |
|---|---|---|
| book-flight-nodelay | 0.05 | 0.15 |
| email-inbox-all | 0.85 | 1.00 |
| choose-date-nodelay | 0.35 | 0.4 |
| social-media | 0.7 | 0.7 |
| … | … | … |
| Mean Reward | 0.52 | 0.61 |

BAGEL improves performance over an instruction-first approach

# Some Results from In-Context Learning with BAGEL demonstrations

(↑ is better)

| Domain | BAGEL (no itrs) | BAGEL |
|---|---|---|
| book-flight-nodelay | 0.15 | 0.15 |
| email-inbox-all | 0.60 | 1.00 |
| choose-date-nodelay | 0.15 | 0.4 |
| social-media | 0.6 | 0.7 |
| … | … | … |
| Mean Reward | 0.53 | 0.61 |

Multiple Iterations of relabeling improves performance

49

# Lessons learned from BAGEL

- Linguistic priors can be a powerful tool for filtering useful interactions on web-interfaces.

- Instruction-first data generation has been extremely successful for ungrounded settings (e.g. self-instruct). But this does not work so well for grounded settings:
  - Given website landing page, no model can imagine the set of possible meaningful tasks.
  - This becomes even harder when website contains **hidden features** that can only be found via interaction.

- Round-trips between natural language instructions and sequence of actions eventually finds high quality correspondences. Lots of other places where similar ideas have been successful:
  - Cycle consistency for GAN
  - Machine Translation (Back-translation)

# Grounded Instruction Following on Real Websites

- Exploration in BAGEL is unstructured and discovers simple demonstrations.
- Tasks on real-websites are hierarchical and complex:



*Navigate to the issue titled memory leak for the vllm project, like the comment by user wciq1208 and then respond to the post saying "I am also …."*
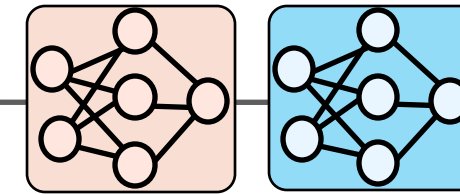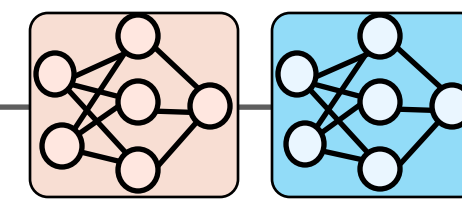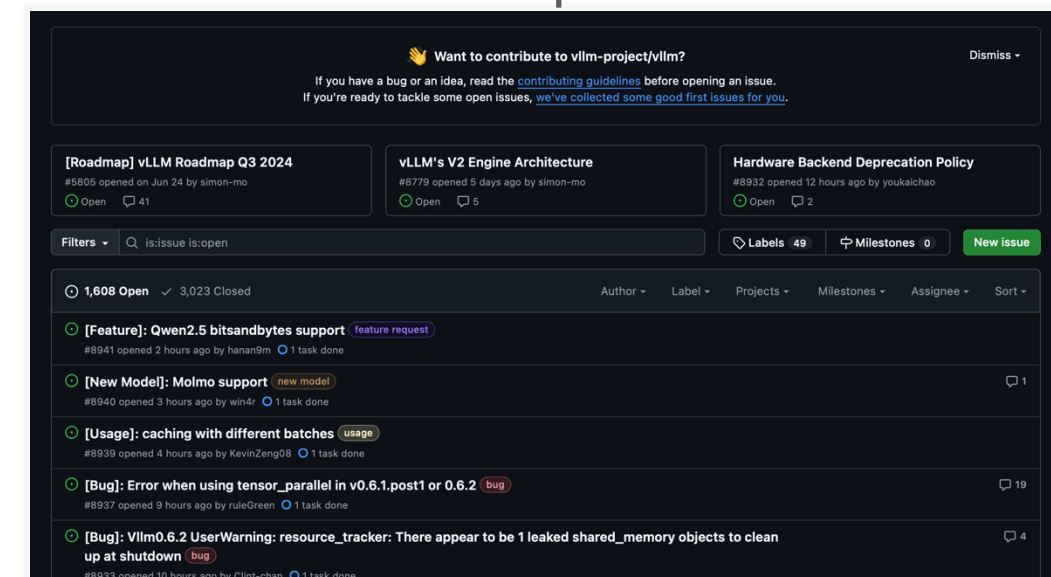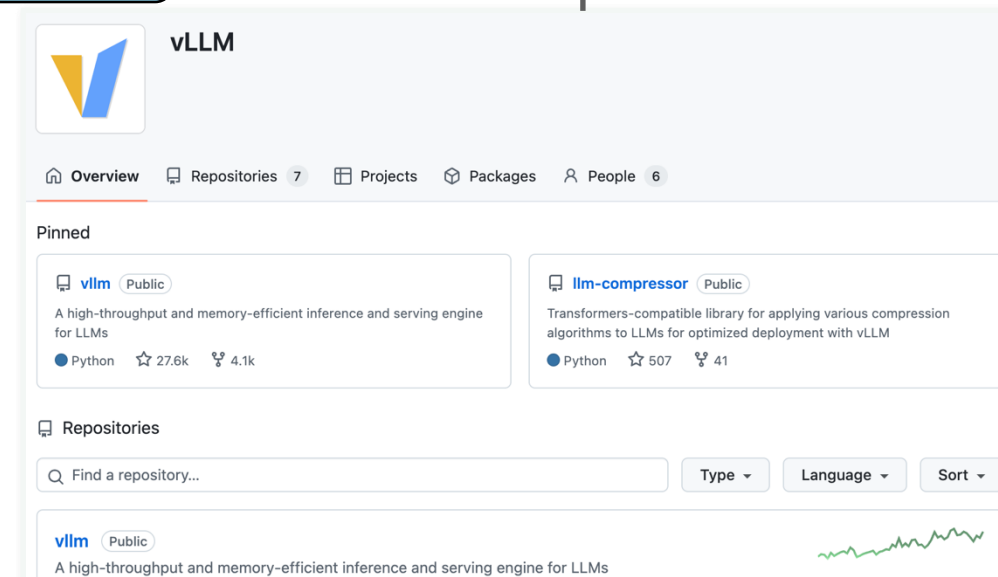
How can we search efficiently over the space of possibly exponential meaningful interactions?
Use Linguistic priors

# Grounded Instruction Following on Real Websites

*Navigate to the issue titled memory leak for the vllm project, like the comment by user wciq1208 and then respond to the post saying "I am also …."*



Prompt: You are given a sequence of actions and corresponding HTML states on a website…

Label: [[pred]]

Prompt:Output "1" if the trajectory is correct for the given instruction…

Label: [[pred]]

Trajectory Labeler

Demonstration Filter

# Grounded Instruction Following on Real Websites



*Navigate to the vllm repository on github.com*

*Navigate to the issue titled memory leak for the vllm project, like the comment by user wciq1208 and then respond to the post saying "I am also …."*

Trajectory Labeler          Demonstration Filter

# Grounded Instruction Following on Real Websites



*Open the issues page for the vllm repository on github.com*
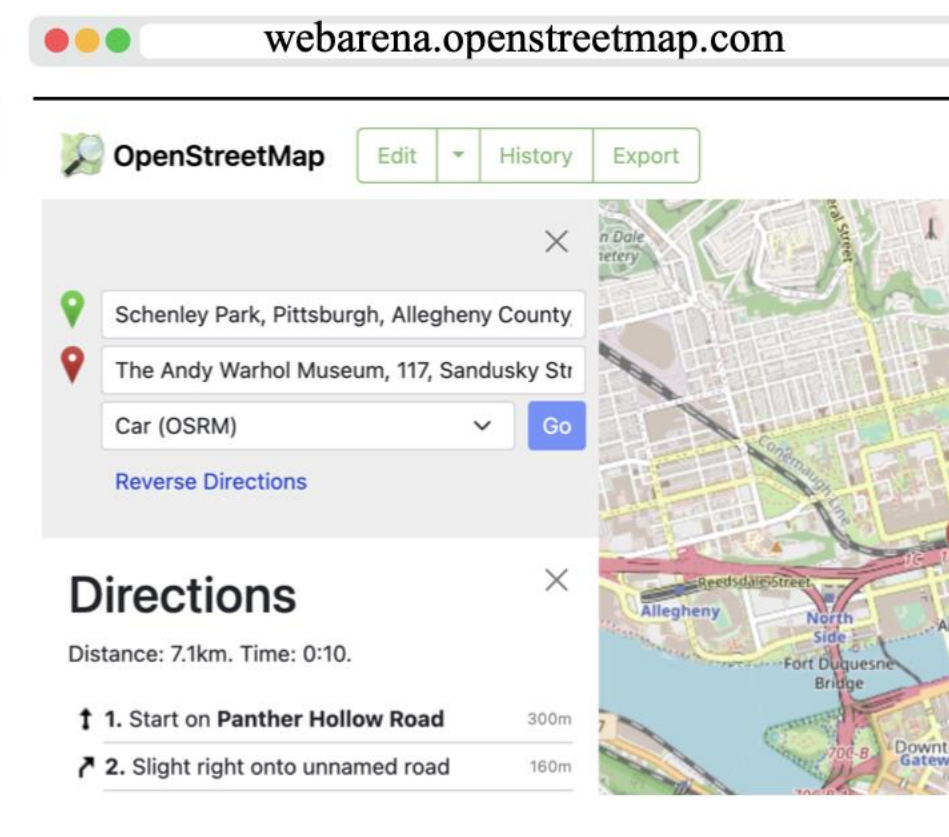
*Navigate to the vllm repository on github.com*

*Navigate to the issue titled memory leak for the vllm project, like the comment by user wciq1208 and then respond to the post saying "I am also …."*

Trajectory Labeler    Demonstration Filter

# Grounded Instruction Following on Real Websites

*Navigate to issue titled memory leak for the vllm repository on github.com*



*Open the issues page for the vllm repository on github.com*



*Navigate to the vllm repository on github.com*



Complex trajectories can be segmented such that prefixes correspond to sub-tasks. We can use this as a pruning function for structuring exploration!

*Navigate to the issue titled memory leak for the vllm project, like the comment by user wciq1208 and then respond to the post saying "I am also …."*

Trajectory Labeler     Demonstration Filter

# Grounded Instruction Following on Real Websites

**WebArena:  map instructions to sequence of mouse and keyboard actions**



" Create an efficient itinerary to visit all of Pittsburgh's art museums with minimal driving distance starting from Schenley Park. Log the order in my "awesome-northeast-us-travel" repository "

Search for museums in Pittsburgh

Search for each art museum on the Map

Record the optimized results to the repo

Environment with sandboxed approximations of real websites spanning e-commerce, social media!
Additional utility tools: Maps, calculators, scratchpads, Wikipedia…
Multi-tab browsing
Long-horizon tasks
Evaluates functional correctness

Zhou et al. 2024. WebArena: A Realistic Web Environment for Building Autonomous Agents https://arxiv.org/abs/2307.13854

# Grounded Instruction Following on Real Websites

**NNetscape Navigator (NNetnav): Our latest attempt for generating complex instructions on realistic websites with structured exploration**



Trajectory Labeler

Demonstration Filter

# Grounded Instruction Following on Real Websites

**NNetscape Navigator (NNetnav):** Our latest attempt for generating complex instructions on realistic websites with structured exploration
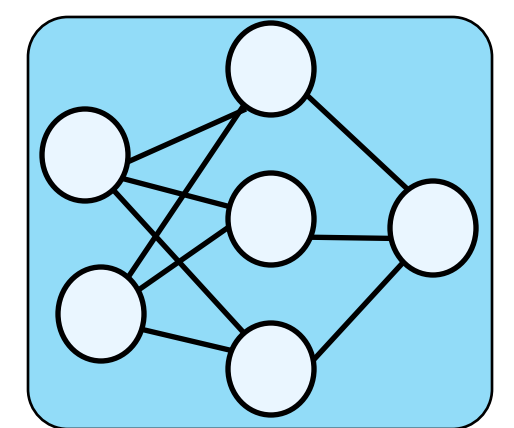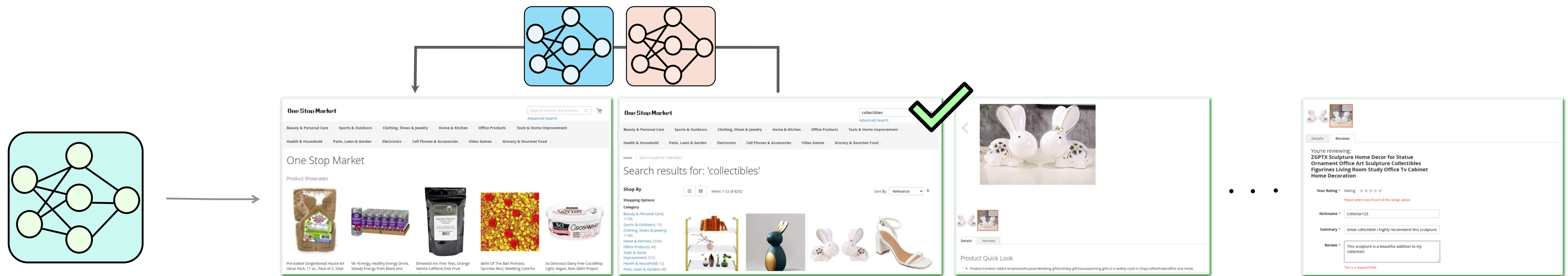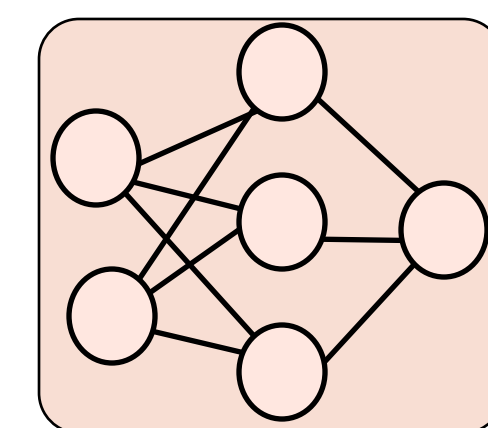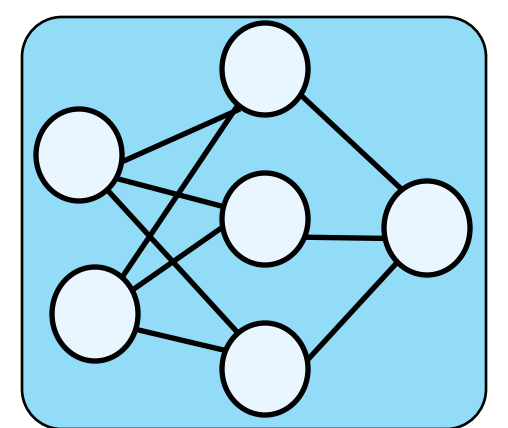
*Search for collectible items and navigate to ceramic rabbit ornaments*



Pruning Function:
- ▶ Invoke Trajectory Labeler at fixed time-steps
- ▶ Assign reward to resulting (instruction, trajectory) pair
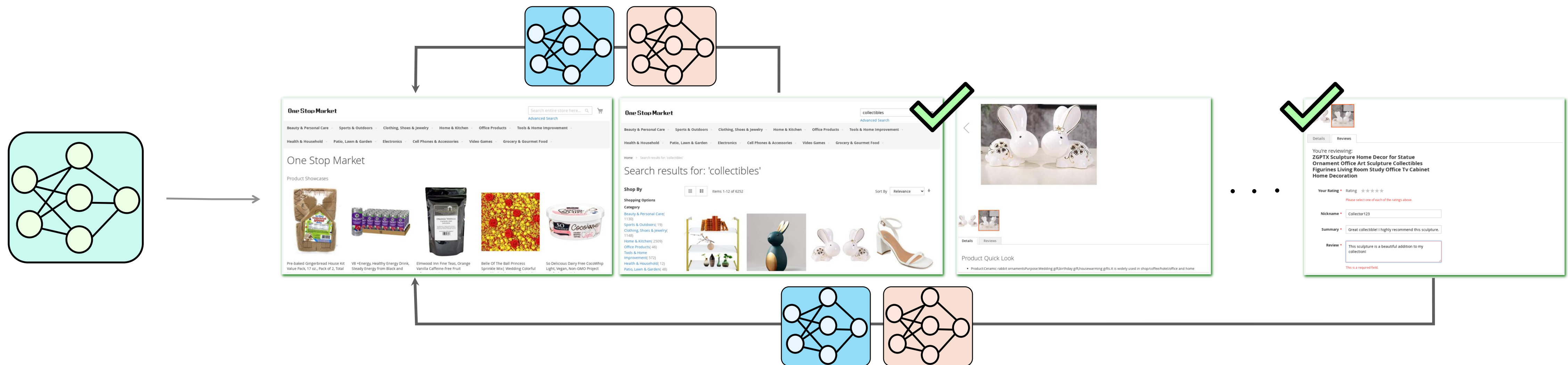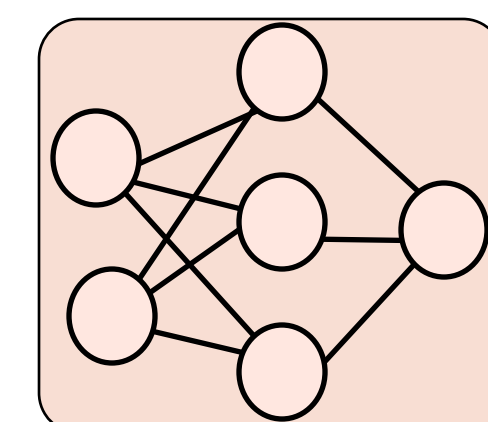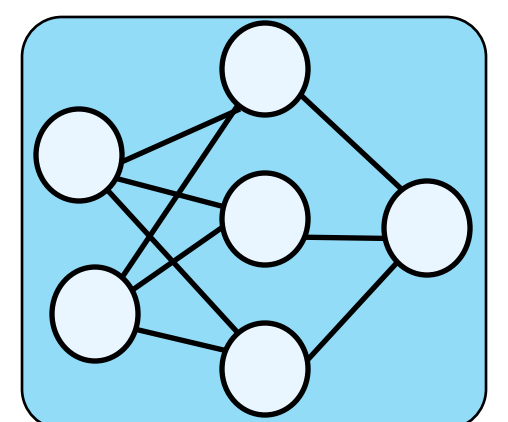- ▶ Prune if reward is low
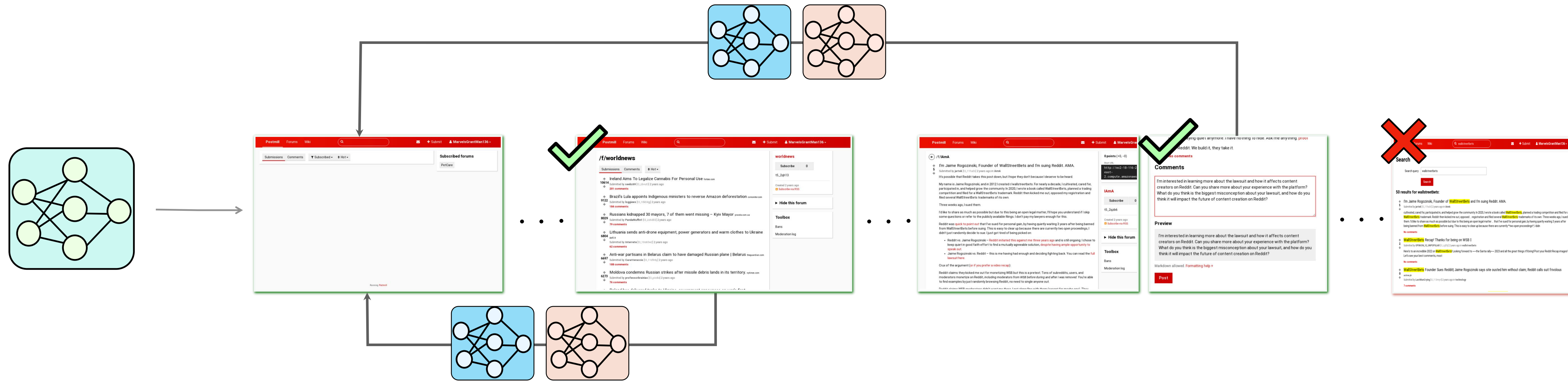
Trajectory Labeler        Demonstration Filter

# Grounded Instruction Following on Real Websites

**NNetscape Navigator (NNetnav):** Our latest attempt for generating complex instructions on realistic websites with structured exploration
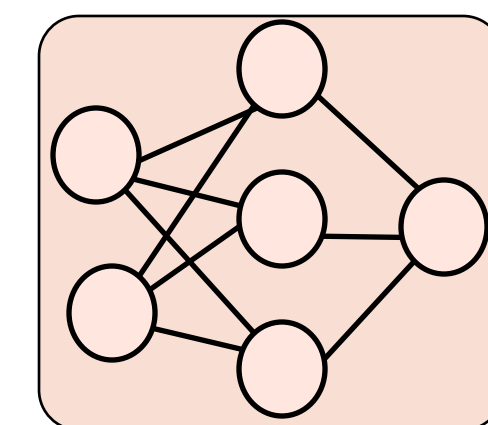
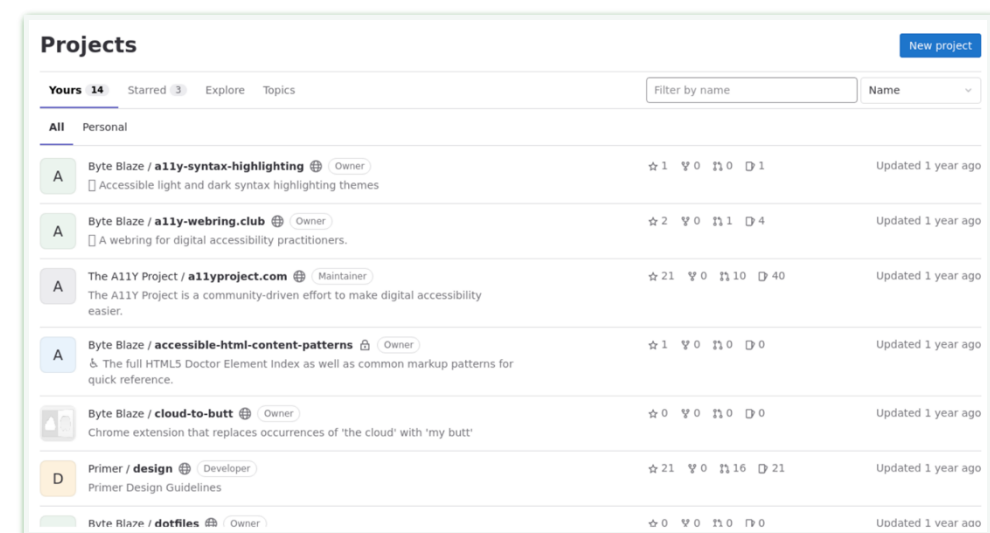*Search for collectible items and navigate to ceramic rabbit ornaments*



Pruning Function:
▸ Invoke Trajectory Labeler at fixed time-steps
▸ Assign reward to resulting (instruction, trajectory) pair
▸ Prune if reward is low

Trajectory Labeler       Demonstration Filter

# Grounded Instruction Following on Real Websites

**NNetscape Navigator (NNetnav):** Our latest attempt for generating complex instructions on realistic websites with structured exploration

*Search for collectible items and navigate to ceramic rabbit ornaments*



*Search for collectible items, check their details and write a review.*

Pruning Function:
- Invoke Trajectory Labeler at fixed time-steps
- Assign reward to resulting (instruction, trajectory) pair
- Prune if reward is low

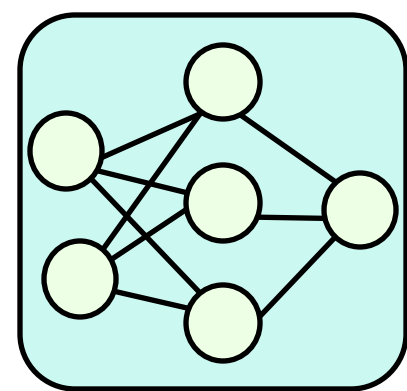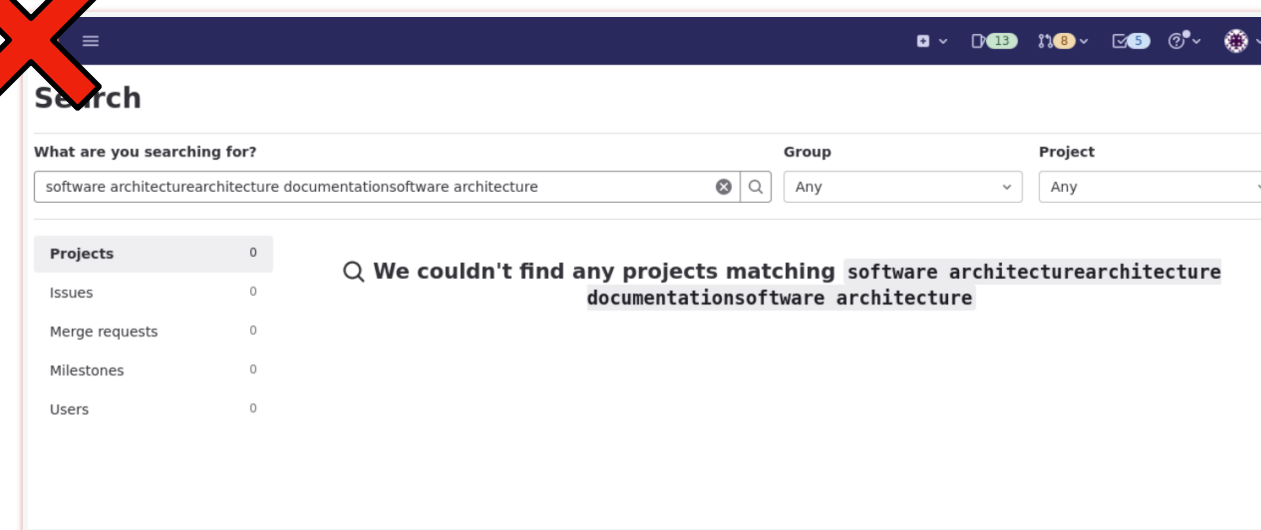Trajectory Labeler        Demonstration Filter

# Grounded Instruction Following on Real Websites

**NNetscape Navigator (NNetnav): Our latest attempt for generating complex instructions on realistic websites with structured exploration**
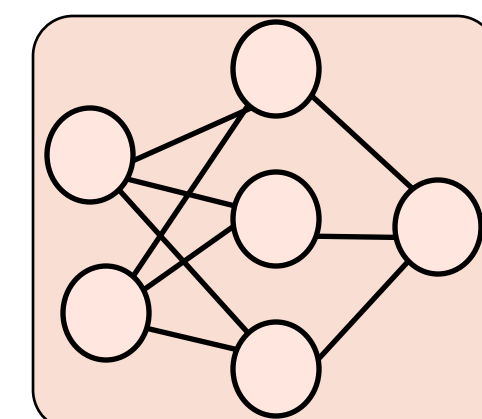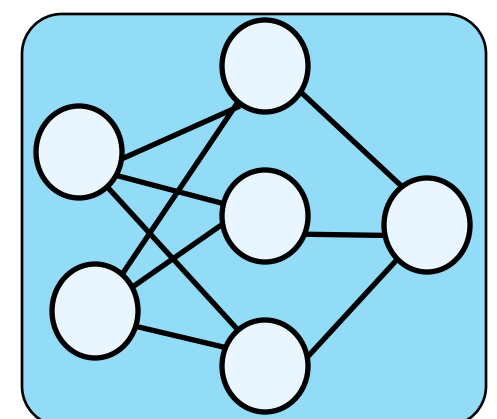


*Find the post about Jaime Rogozinski, and his lawsuit against Reddit, and ask him about his experience.*

*Subscribe to the r/wallstreetbets forum and navigate to world news*

Pruning Function:
- ▶ Invoke Trajectory Labeler at fixed time-steps
- ▶ Assign reward to resulting (instruction, trajectory) pair
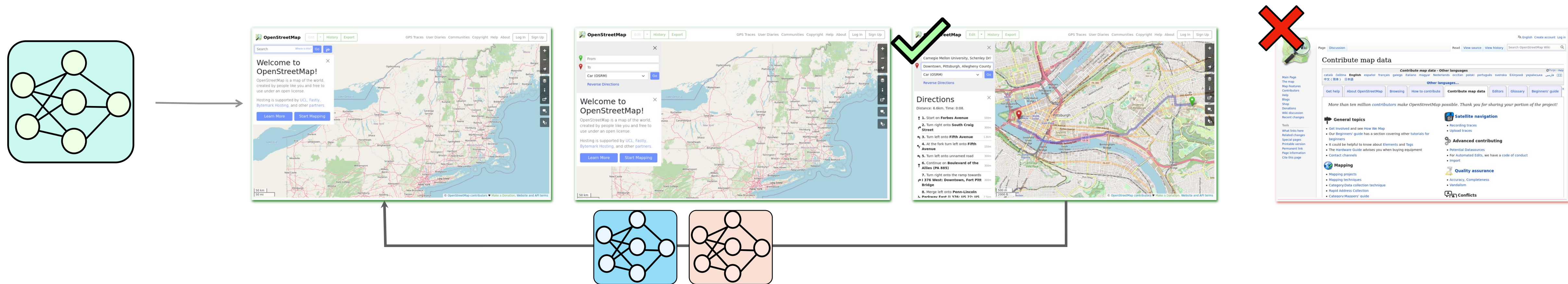- ▶ Prune if reward is low

Trajectory Labeler       Demonstration Filter

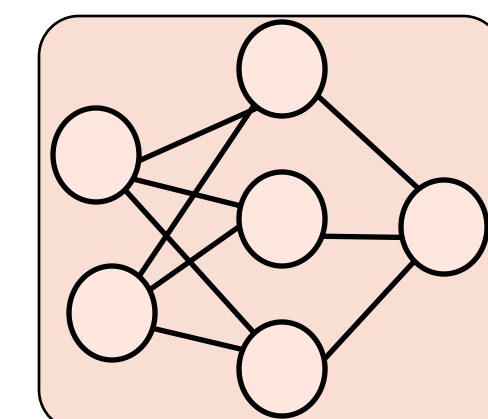# Grounded Instruction Following on Real Websites

**NNetscape Navigator (NNetnav): Our latest attempt for generating complex instructions on realistic websites with structured exploration**
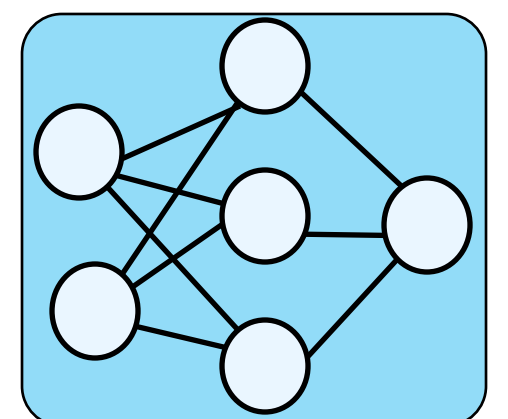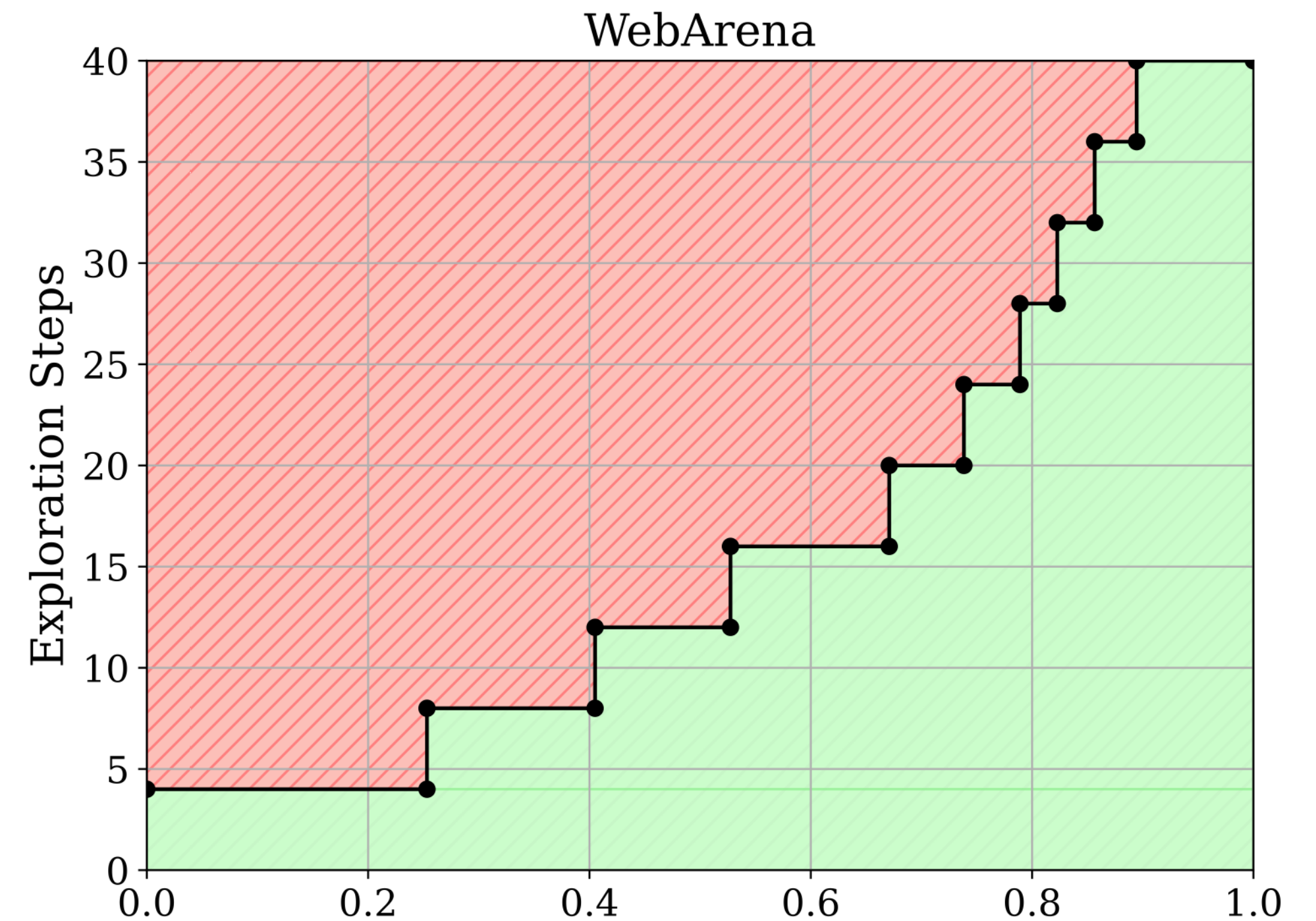


Pruning Function:
▸ Invoke Trajectory Labeler at fixed time-steps
▸ Assign reward to resulting (instruction, trajectory) pair
▸ Prune if reward is low

Trajectory Labeler          Demonstration Filter

# Grounded Instruction Following on Real Websites

**NNetscape Navigator (NNetnav): Our latest attempt for generating complex instructions on realistic websites with structured exploration**



*Find driving directions from CMU to downtown Pittsburgh*

Pruning Function:
- Invoke Trajectory Labeler at fixed time-steps
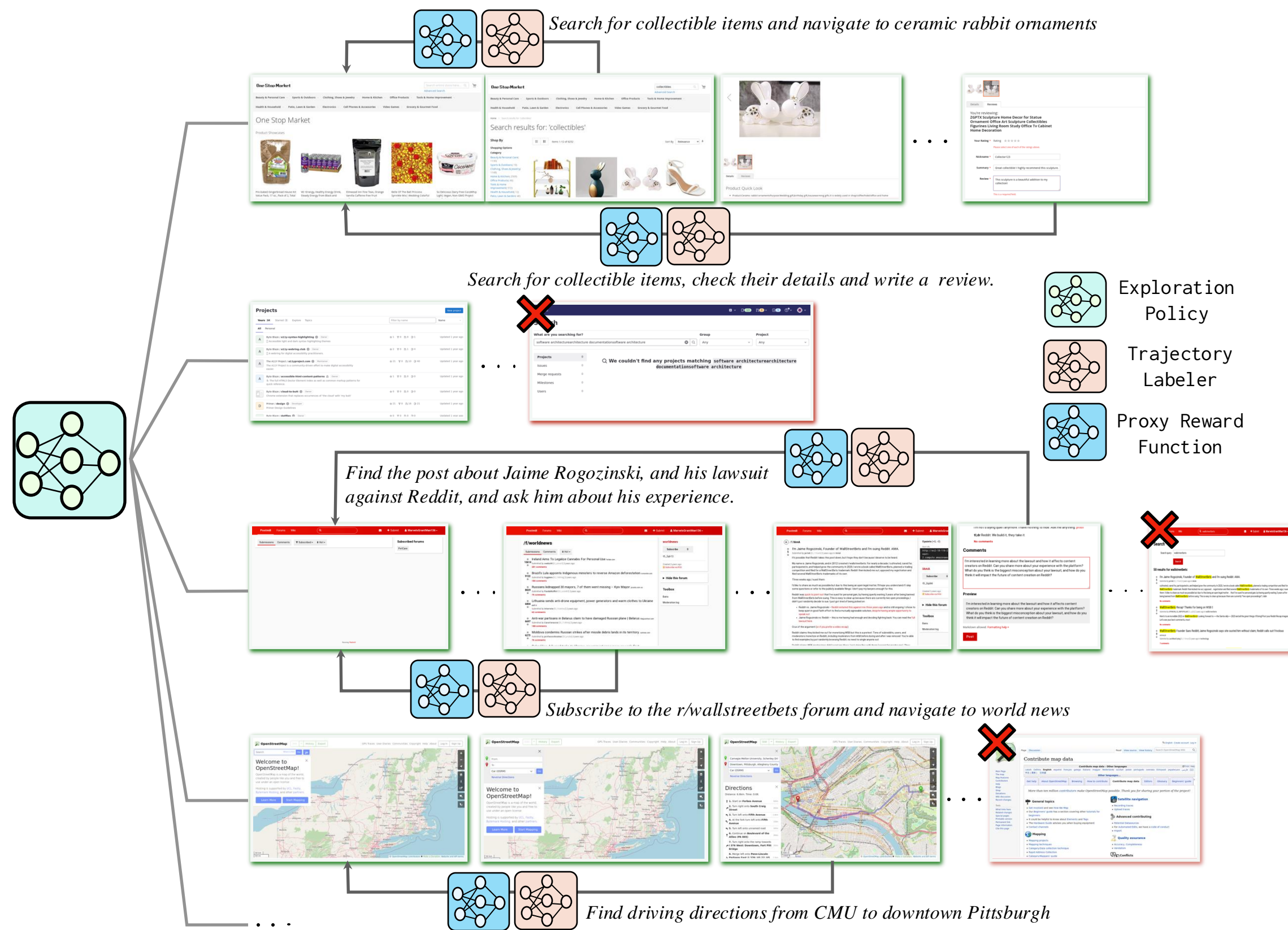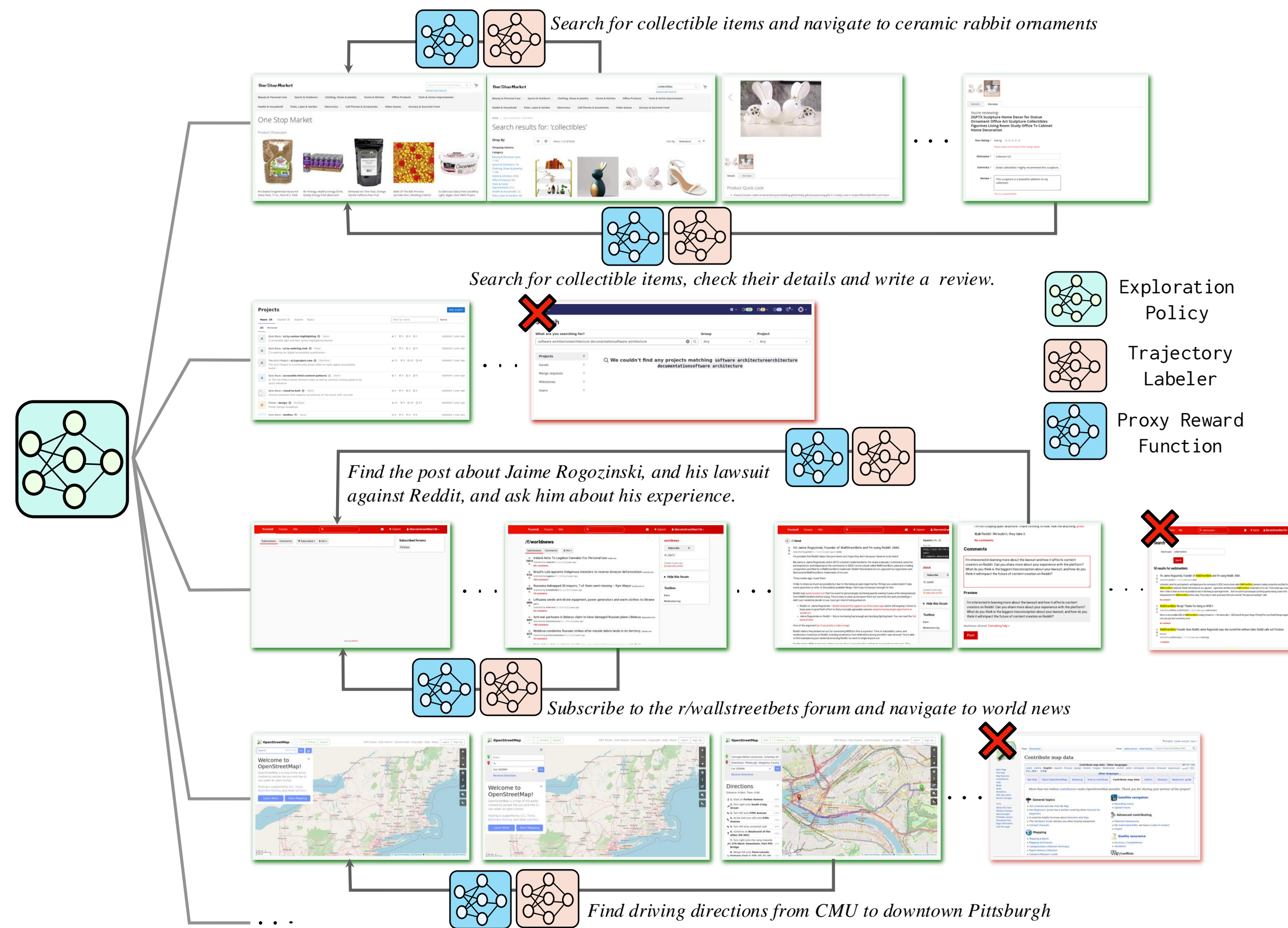- Assign reward to resulting (instruction, trajectory) pair
- Prune if reward is low

*Search for collectible items and navigate to ceramic rabbit ornaments*

*Search for collectible items, check their details and write a review.*

Exploration Policy

Trajectory Labeler

Proxy Reward Function

*Find the post about Jaime Rogozinski, and his lawsuit against Reddit, and ask him about his experience.*

*Subscribe to the r/wallstreetbets forum and navigate to world news*

*Find driving directions from CMU to downtown Pittsburgh*



WebArena

Such a pruning strategy leads to very efficient exploration. Here, area shaded in red corresponds to actions that were prevented due to early pruning

*Search for collectible items and navigate to ceramic rabbit ornaments*

*Search for collectible items, check their details and write a review.*

Exploration Policy

Trajectory Labeler

Proxy Reward Function

*Find the post about Jaime Rogozinski, and his lawsuit against Reddit, and ask him about his experience.*

*Subscribe to the r/wallstreetbets forum and navigate to world news*

*Find driving directions from CMU to downtown Pittsburgh*

---

*Investigate the issue of offline downloads in the a11y-webring.club project and create a new task to implement a function to download PDF cards and JSON data for offline access.*

---

*Create a new forum for discussing AI and Machine Learning, and create a submission about a new AI-powered smartwatch, then search for reviews about the smartwatch on the website*

---

*Find the driving directions from TLC Medical Transportation Services in Syracuse to Times Square in Manhattan.*

---

*Add two products, "Abominable Hoodie" and "Samsung Smart TV", with respective prices $99.99 and $50.00, and then start the process of adding a new customer.*

---

Such a pruning strategy leads to instructions that are complex, hierarchical and plausible (by design)

# Some Results from Fine-tuning LLama-8b with NNetnav demonstrations

## Warmup: MiniWoB++

| Domain | Llama-8b | LLama-8b-IF | Llama-8b-NNetnav |
|---|---|---|---|
| click-checkboxes-soft | 0.4 | 0.25 | 0.65 |
| email-inbox | 0.25 | 0.3 | 0.3 |
| use-autocomplete | 0.25 | 0.55 | 0.7 |
| phonebook | 0.15 | 0.15 | 0.2 |
| … | … | | … |
| Mean Reward | 0.28 | 0.28 | 0.48 |

Performs better than Instruction-first and Zero-shot LLama

# Some Results from Fine-tuning LLama-8b with NNetnav demonstrations

| Method | WebArena Success Rate (%) |
|---|---|
| GPT-3.5 | 6.2 |
| GPT-4o | 14.4 |
| GPT-4o + tree-search | 19.2 |
| GPT-4o (BrowserGym) | 23.5 |

- State-of-the-art results based on using GPT-4o, zero-shot + various inference-time procedures like Tree Search
- This work: Finetuning smaller agents (LLama-8b) with demonstrations

# Some Results from Fine-tuning LLama-8b with NNetnav demonstrations

| Method | WebArena Success Rate (%) | #Demonstrations |
|---|---|---|
| LLama-8b (zero-shot) | 1 | 0 |
| AutoWebGLM-7b (S1)* | 2.5 | 240 |
| Synatra-CodeLLama-7b* | 6.3 | 30k |

This work: Finetuning smaller agents (LLama-8b) with demonstrations

*  Synatra: Turning Indirect Knowledge into Direct Demonstrations for Digital Agents at Scale

# Some Results from Fine-tuning LLama-8b with NNetnav demonstrations

| Method | WebArena Success Rate (%) | #Demonstrations |
|---|---|---|
| LLama-8b (zero-shot) | 1 | 0 |
| AutoWebGLM-7b (S1)* | 2.5 | 240 |
| Synatra-CodeLLama-7b* | 6.3 | 30k |
| LLama-8b-IF | 4.2 | 145 |
| LLama-8b-nnetnav | **7.2** | 650 |

This work: Finetuning smaller agents (LLama-8b) with demonstrations

▸ Outperforms prior approaches (sometimes with orders of magnitude fewer demonstrations)
▸ Outperforms our instruction-first baseline when controlling for exploration budget

69

* Synatra: Turning Indirect Knowledge into Direct Demonstrations for Digital Agents at Scale

# Some rollouts from LLama-3-8b-nnetnav

*closing an issue on a repository*

# Some rollouts from LLama-3-8b-nnetnav

*making a comment on a post*

# Some rollouts from LLama-3-8b-nnetnav

*creating new discussion forum*

# Some rollouts from LLama-3-8b-nnetnav

*finding directions on maps*

# Some rollouts from LLama-3-8b-nnetnav
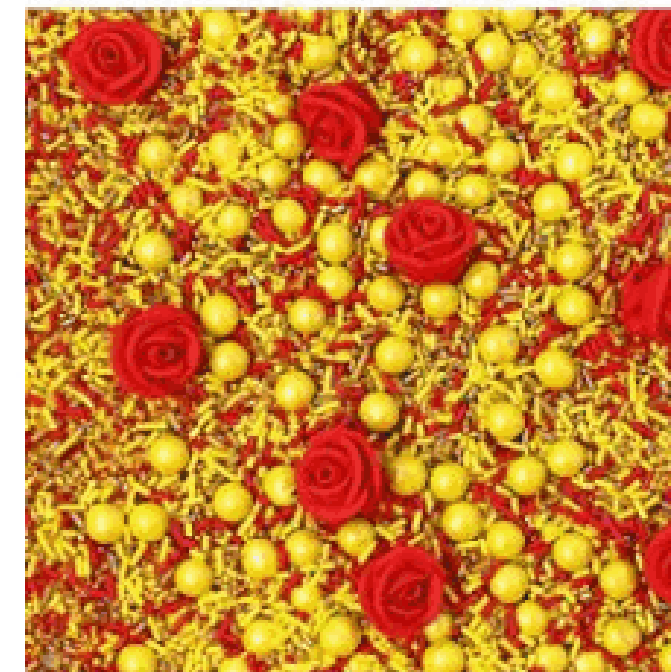
*adding items from prior orders to cart*

# Some rollouts from LLama-3-8b-nnetnav

*interacting with inventory*

# Takeaways from NNetnav

- Yes, LLMs are not necessarily trained for grounded instruction following BUT:

  ▸ Exploration via an LLM, is a powerful tool for creating training data for grounded instruction following.

  ▸ Linguistic priors (coming from LLMs) are a powerful tool for structuring and pruning exploration in complex web environments!
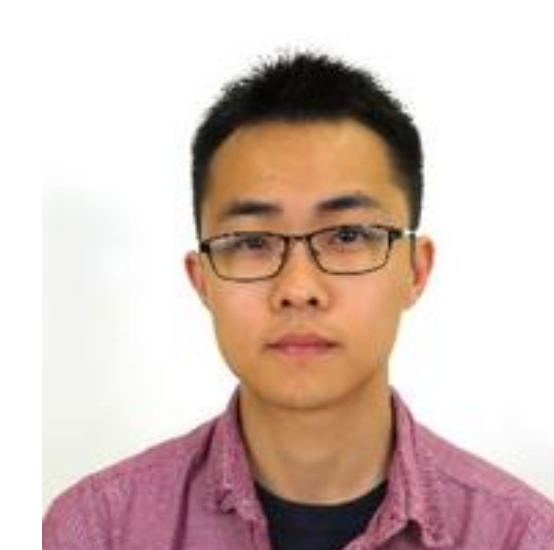
# BAGEL and NNetNav

- BAGEL: Bootstrapping Agents by Guiding Exploration with Language (ICML 2024) arXiv:2403.08140
- NNetscape Navigator: Complex Demonstrations for Web Agents without a Demonstrator (In Submission) https://openreview.net/forum?id=hHF5AayC7O

Stanford | NLP

Christopher Manning

Shikhar Murty

Google DeepMind

Pete Shaw

Mandar Joshi

Kenton Lee

servicenow®

Dzmitry Bahdanau

# Meaning and Intelligence in Language Models: From Philosophy to Agents in a World



## Christopher Manning

@chrmanning ✽ manning@stanford.edu
Stanford NLP Group ✽ Director, Stanford AI Lab
COLM, Philadelphia, 2024