



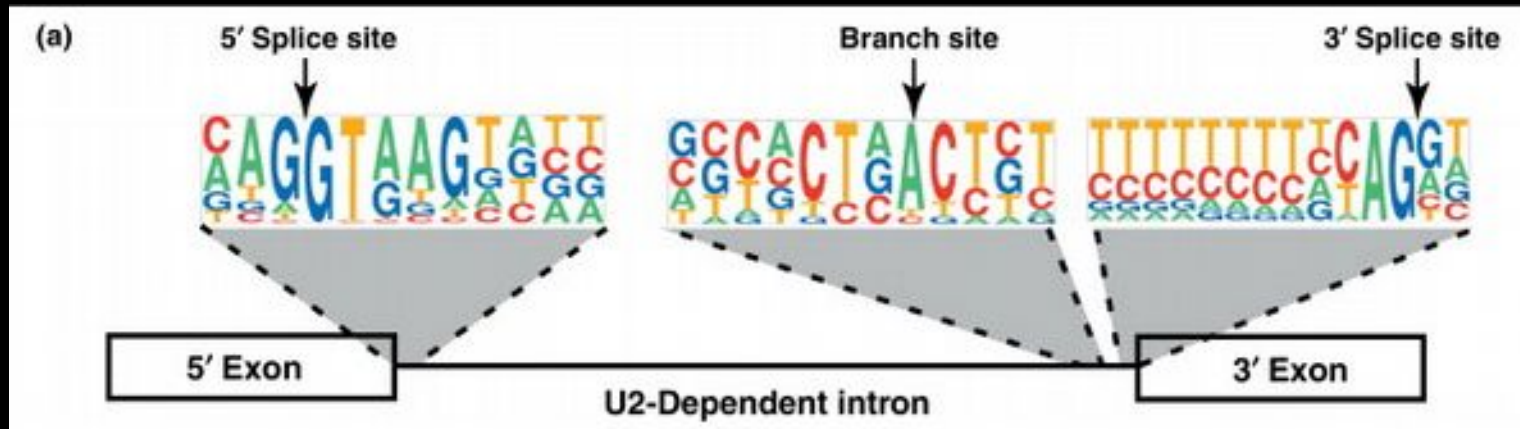
Stuff I did in the Spring while not Replying to Email (aka "advances in structured prediction")

Examples of ~~structured joint~~ prediction

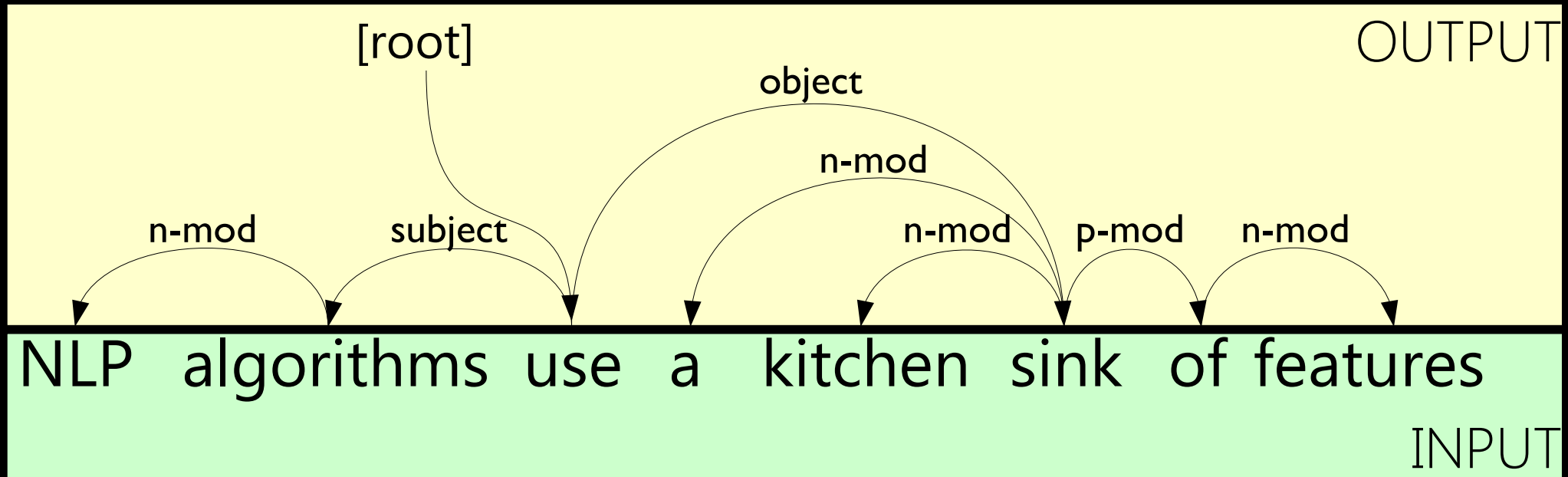
Sequence labeling

x = the monster ate the sandwich
y = Dt Nn Vb Dt Nn

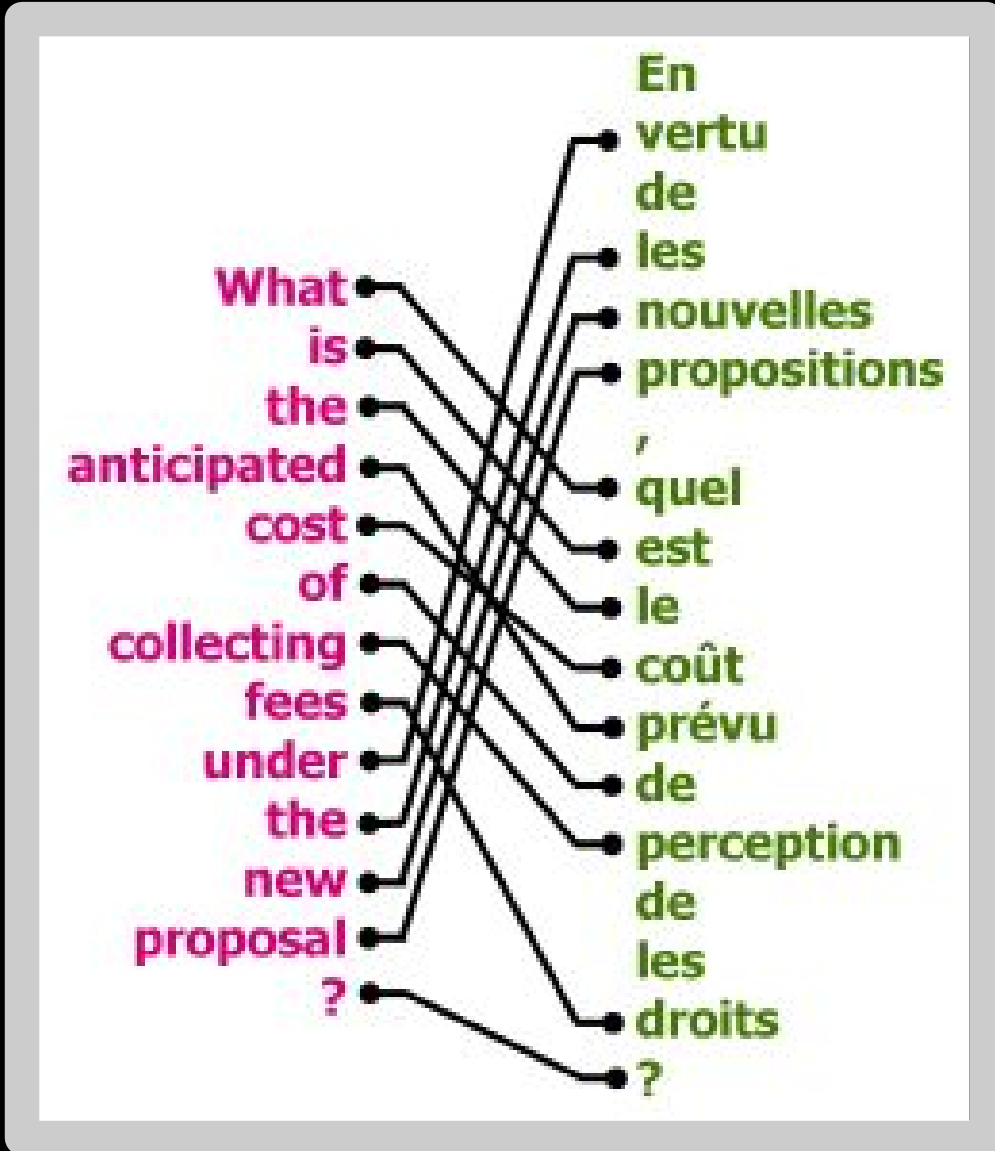
x = Yesterday I traveled to Lille
y = - PER - - LOC



Natural language parsing



(Bipartite) matching



Machine translation



Google Translate

This text has been [automatically translated](#) from Arabic:

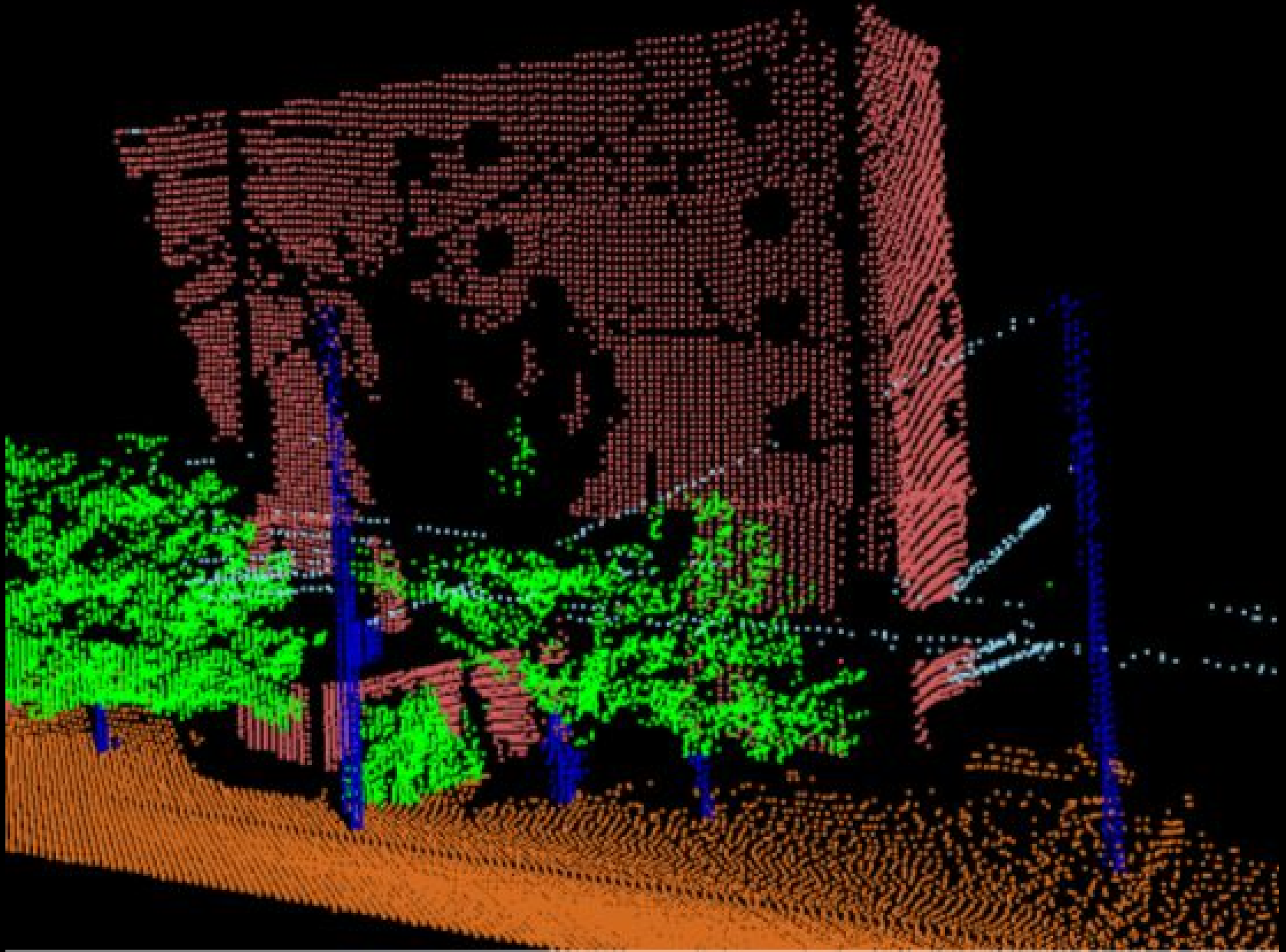
Moscow stressed tone against Iran on its nuclear program. He called Russian Foreign Minister Tehran to take concrete steps to restore confidence with the international community, to cooperate fully with the IAEA. Conversely Tehran expressed its willingness

Translate text

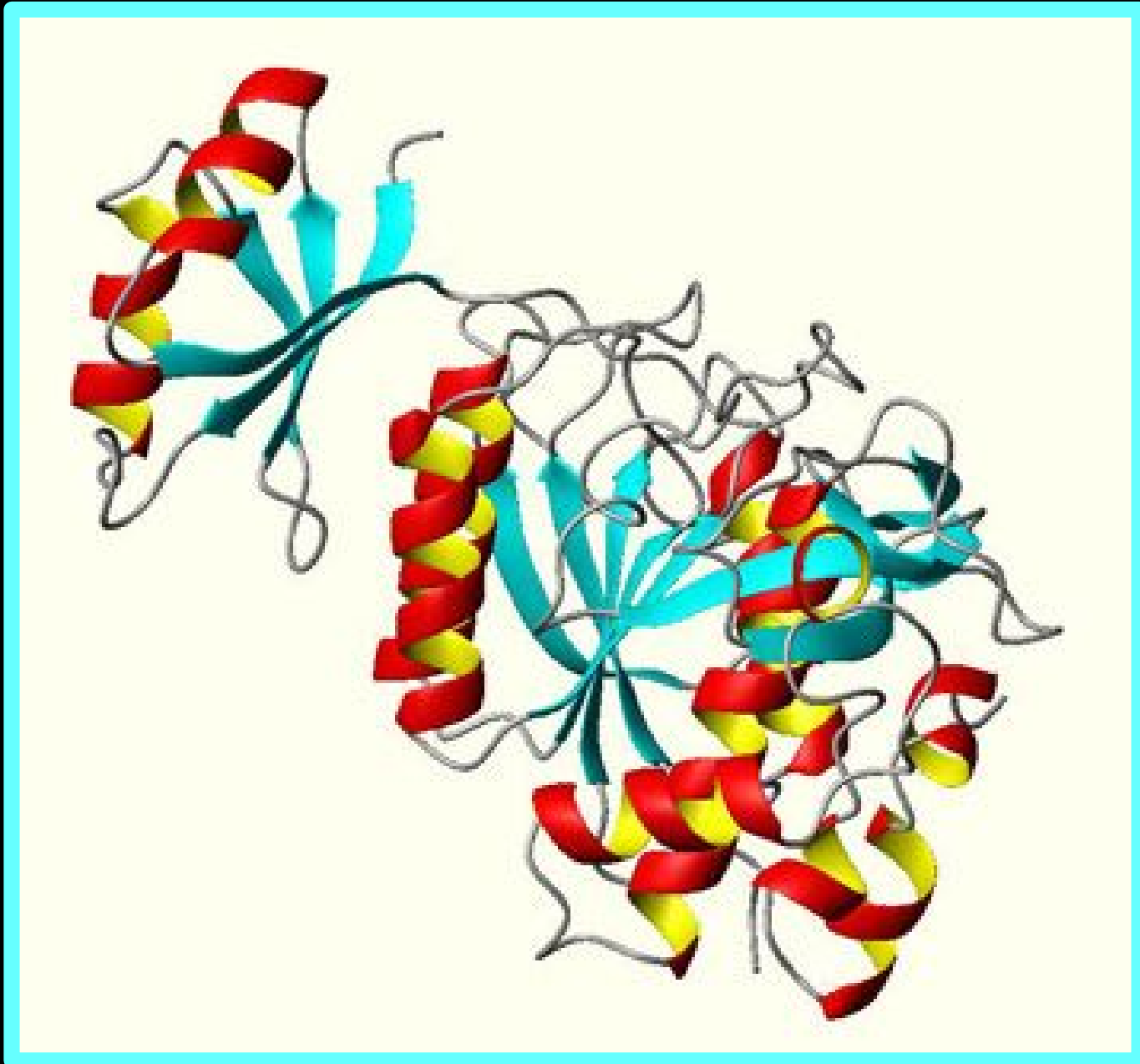
شدت موسكو لهجتها ضد إيران بشأن برنامجها النووي. ودعا وزير الخارجية الروسي طهران إلى اتخاذ خطوات ملموسة لاستعادة الثقة مع المجتمع الدولي والتعاون الكامل مع الوكالة الذرية. بالمقابل أبدت طهران استعدادها لاستئناف السماح بعمليات التفتيش المفاجئة بشرط إسقاط مجلس الأمن ملفها النووي.

from

Segmentation



Protein secondary structure prediction



Outline

- Background: learning to search
- Stuff I did in the Spring
 - Imperative DSL/library for learning to search
 - SOTA examples for tagging, parsing, relation extraction, etc.
 - Learning to search under bandit feedback
 - Hardness results for learning to search
 - Active learning for accelerating learning to search
- Stuff I'm trying to do now
 - Distant supervision
 - Mashups with recurrent neural networks

Isn't this
kinda narrow?

My experience, 6 months in industry

- Standard adage: academia=freedom, industry=time
 - Number of responsibilities vs number of bosses
- Aspects I didn't anticipate
 - Breadth (academia) versus depth (industry)
 - Collaborating through students versus directly
 - Security through tenure versus security through \$
- At the end of the day: who are your colleagues and what do you have to do to pay the piper?

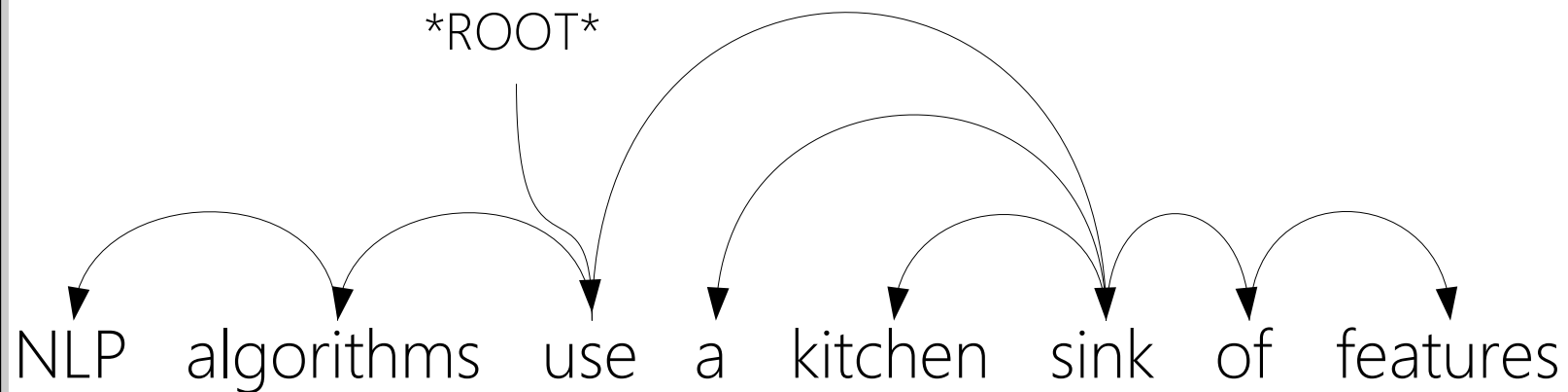
Major caveat: this is comparing a top ranked CS dept to top industry lab, in a time when there's tons of money in this area (more in industry)

Joint prediction via learning to search

Part of Speech Tagging

NN	NNS	VBP	DT	NN	NN	IN	NNS
NLP	algorithms	use	a	kitchen	sink	of	features

Dependency Parsing

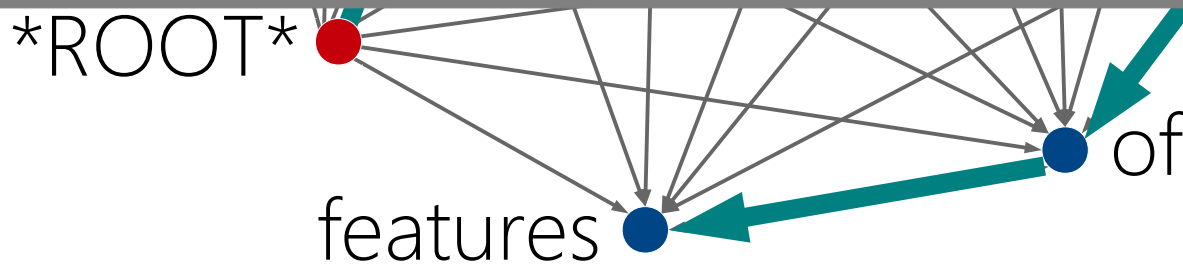


Joint prediction via learning to search



Joint Prediction Haiku

A joint prediction
Across a single input
Loss measured jointly



Back to the original problem...

- How to optimize a discrete, joint loss?

• Input: $\mathbf{x} \in X$ 

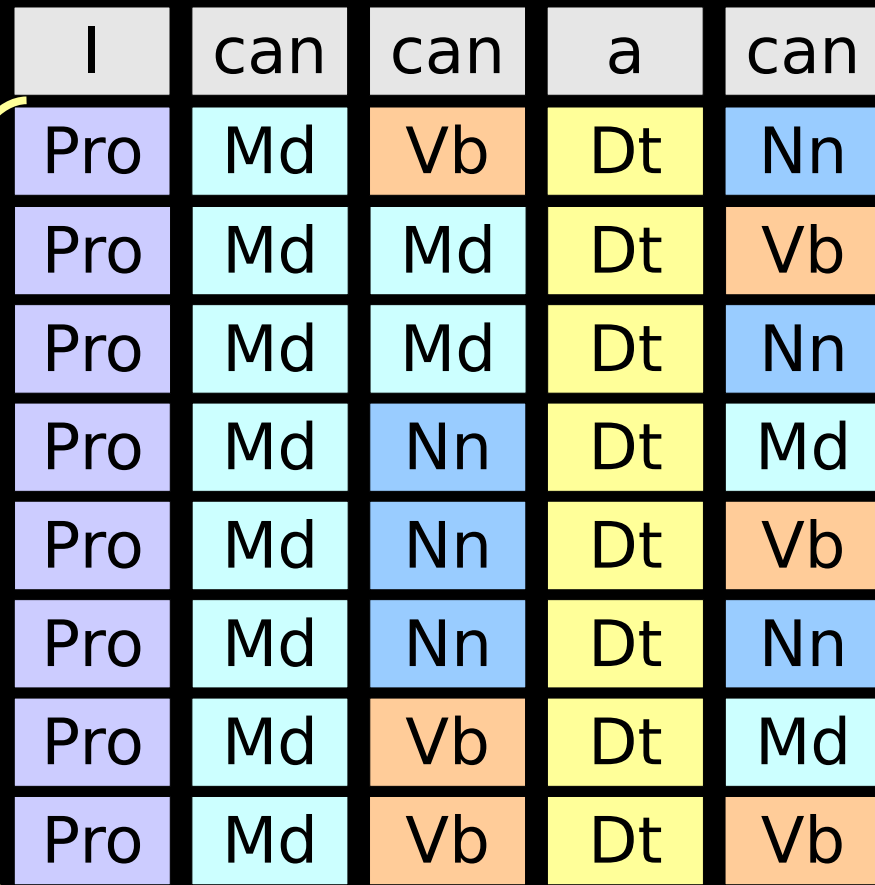
• Truth: $y \in Y(\mathbf{x})$ 

• Outputs: $Y(\mathbf{x})$ 

• Predicted: $\hat{y} \in Y(\mathbf{x})$

• Loss: $\text{loss}(y, \hat{y})$

• Data: $(\mathbf{x}, y) \sim D$



I	can	can	a	can
Pro	Md	Vb	Dt	Nn
Pro	Md	Md	Dt	Vb
Pro	Md	Md	Dt	Nn
Pro	Md	Nn	Dt	Md
Pro	Md	Nn	Dt	Vb
Pro	Md	Nn	Dt	Nn
Pro	Md	Vb	Dt	Md
Pro	Md	Vb	Dt	Vb

Back to the original problem...

- How to optimize a discrete, joint loss?

- Input: $\mathbf{x} \in X$
- Truth: $y \in Y(\mathbf{x})$
- Outputs: $Y(\mathbf{x})$
- Predicted: $\hat{y} \in Y(\mathbf{x})$
- Loss: $\text{loss}(y, \hat{y})$
- Data: $(\mathbf{x}, y) \sim D$

Goal:

find $h \in H$
such that $h(\mathbf{x}) \in Y(\mathbf{x})$
minimizing

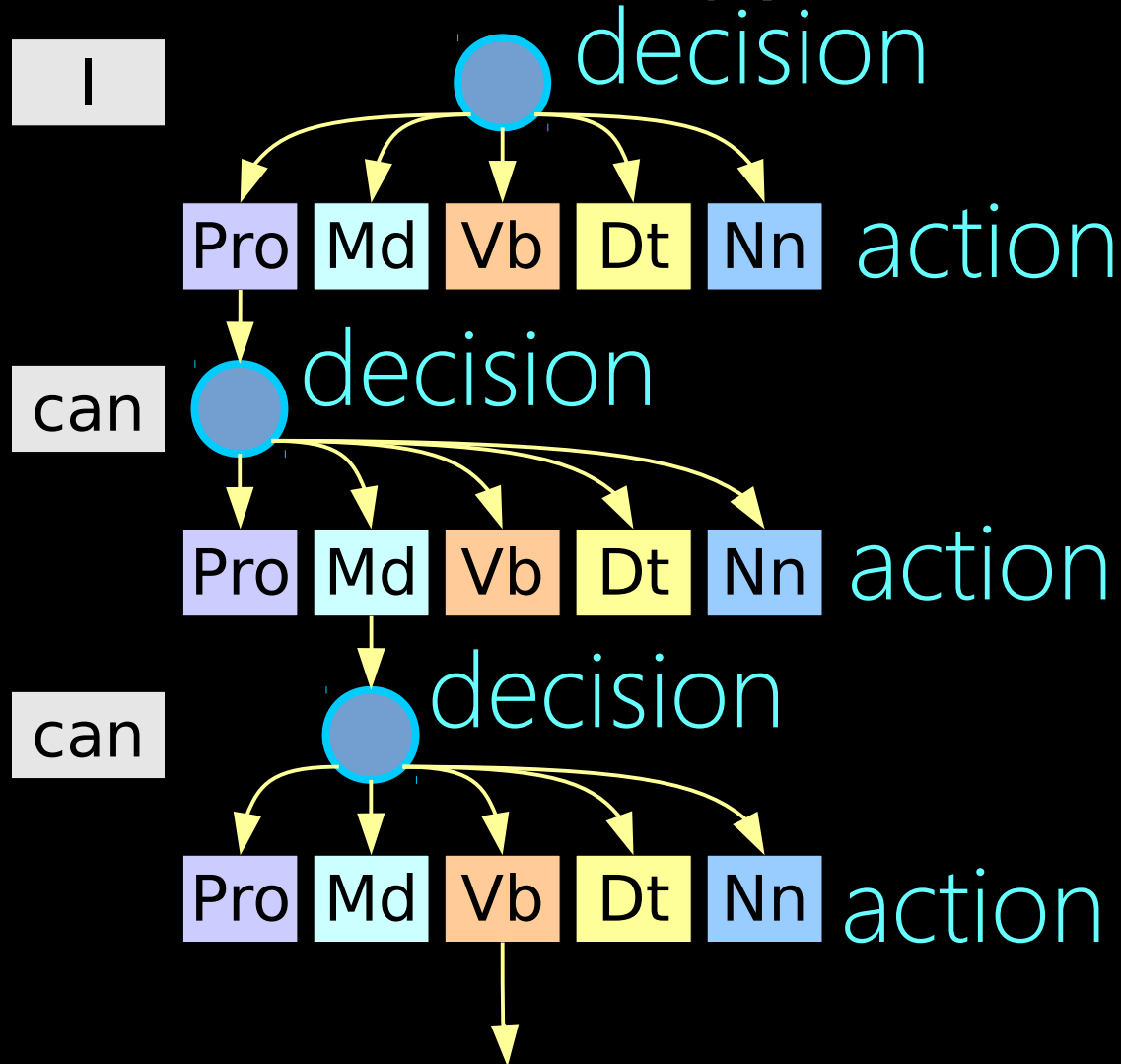
$$E_{(\mathbf{x}, y) \sim D} [\text{loss}(y, h(\mathbf{x}))]$$

based on N samples

$$(\mathbf{x}_n, y_n) \sim D$$

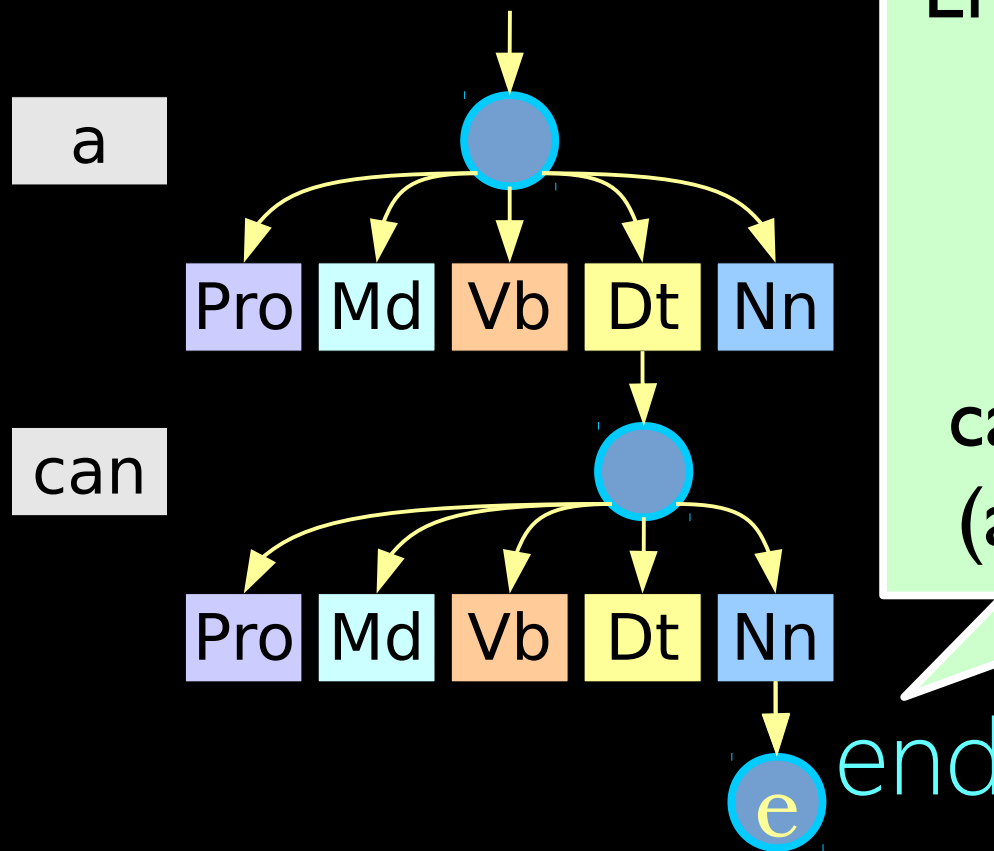
Search spaces

- When y decomposes in an ordered manner, a sequential decision making process emerges



Search spaces

- When y decomposes in an ordered manner, a sequential decision making process emerges



Encodes an output

$$\hat{y} = \hat{y}(e)$$

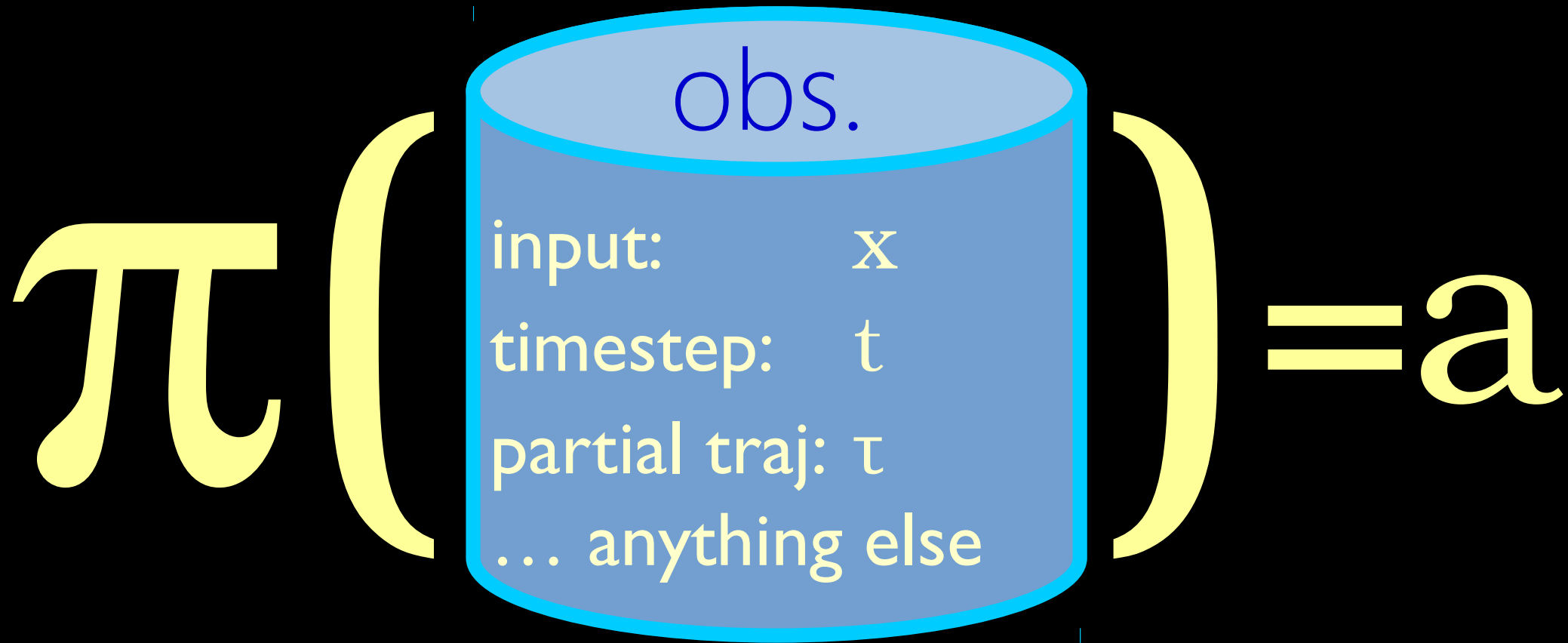
from which

$$\text{loss}(y, \hat{y})$$

can be computed
(at training time)

Policies

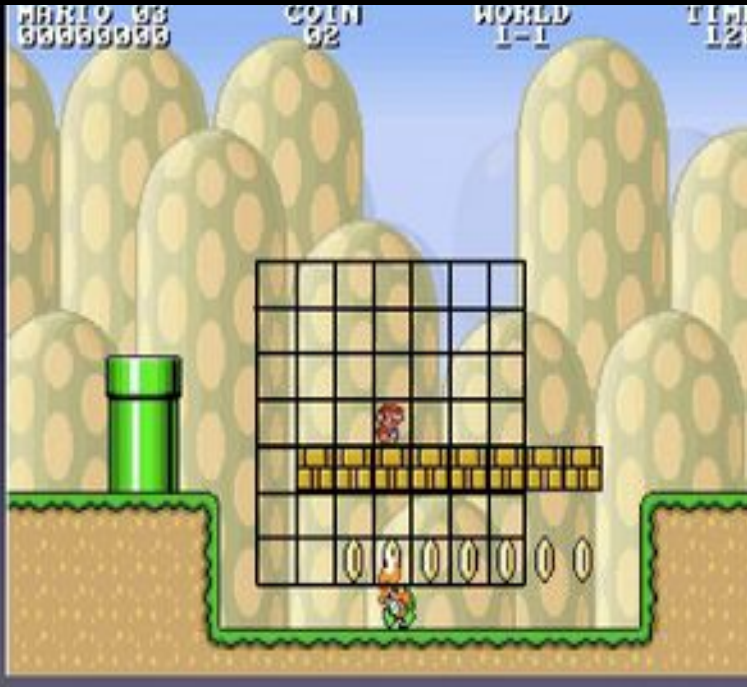
- A policy maps observations to actions



An analogy from playing Mario

From Mario AI competition 2009

Input:



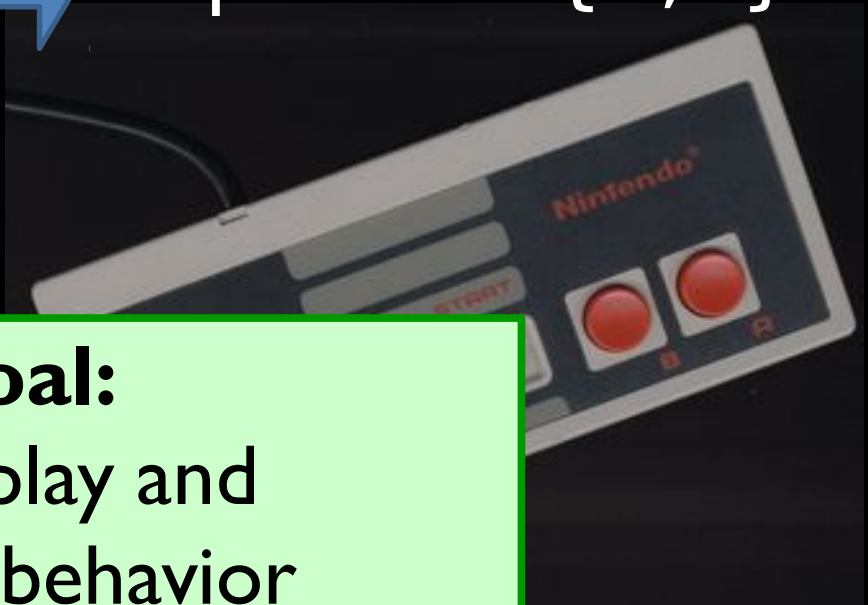
Output:

Jump in $\{0,1\}$
Right in $\{0,1\}$
Left in $\{0,1\}$
Speed in $\{0,1\}$



High level goal:

Watch an expert play and
learn to mimic her behavior

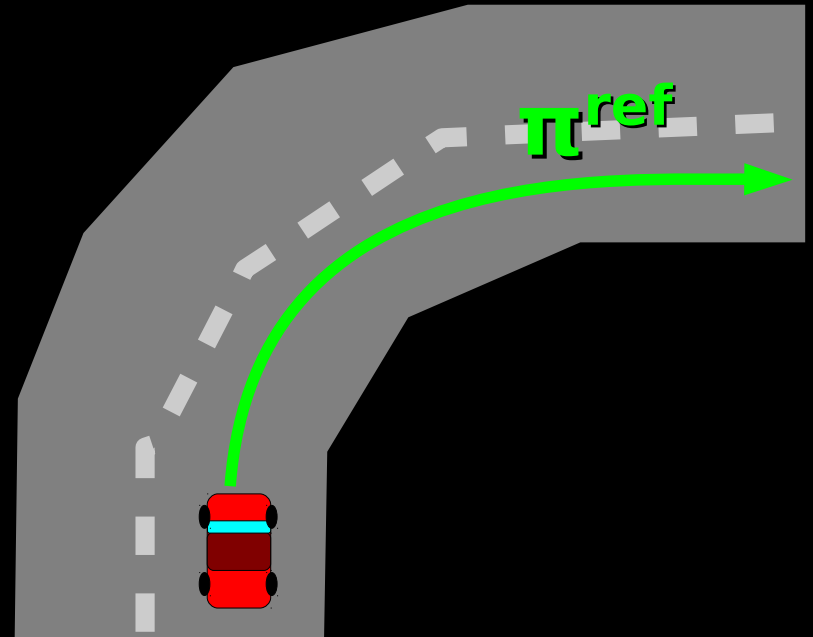


Training (expert)



Warm-up: Supervised learning

1. Collect trajectories from expert π^{ref}
 2. Store as dataset $\mathbf{D} = \{ (o, \pi^{\text{ref}}(o, y)) \mid o \sim \pi^{\text{ref}} \}$
 3. Train classifier π on \mathbf{D}
- Let π play the game!



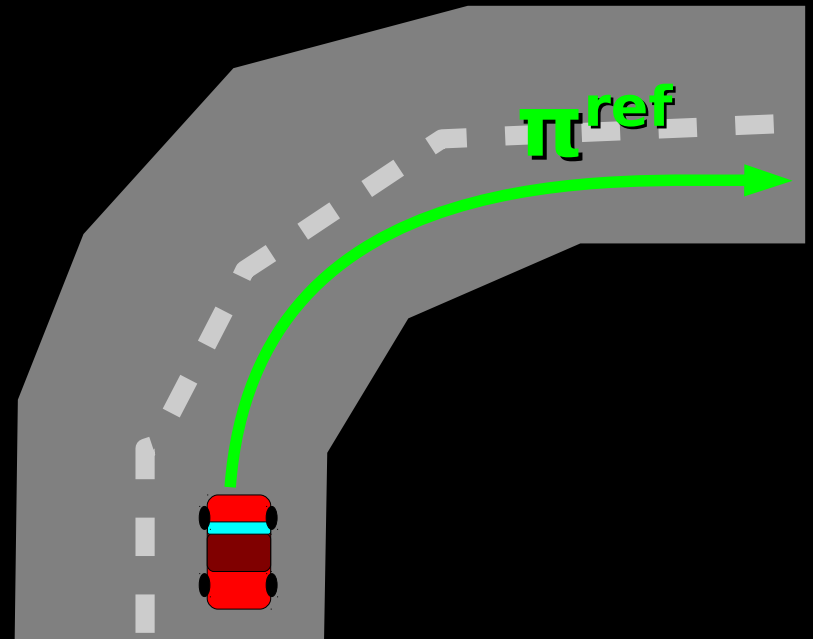
Test-time execution (sup. learning)



What's the (biggest) failure mode?

The expert never gets stuck next to pipes

⇒ Classifier doesn't learn to recover!

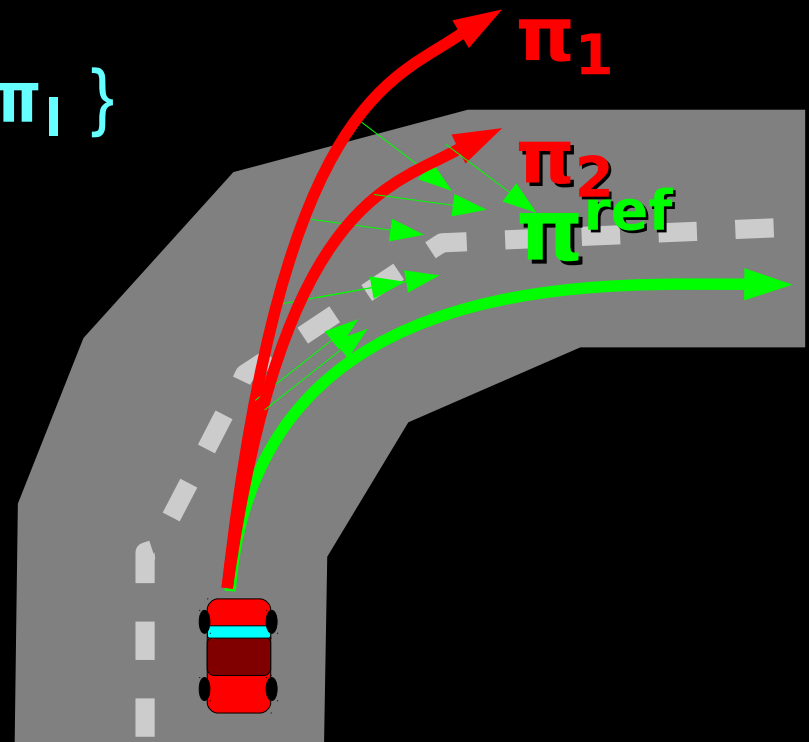


Warm-up II: Imitation learning

1. Collect trajectories from expert π^{ref}
2. Dataset $\mathbf{D}_0 = \{ (o, \pi^{\text{ref}}(o, y)) \mid o \sim \pi^{\text{ref}} \}$
3. Train π_1 on \mathbf{D}_0
4. Collect new trajectories from π_1
 - But let the *expert* steer!
5. Dataset $\mathbf{D}_1 = \{ (o, \pi^{\text{ref}}(o, y)) \mid o \sim \pi_1 \}$
6. Train π_2 on $\mathbf{D}_0 \cup \mathbf{D}_1$

- In general:
 - $\mathbf{D}_n = \{ (o, \pi^{\text{ref}}(o, y)) \mid o \sim \pi_n \}$
 - Train π_{n+1} on $\bigcup_{i \leq n} \mathbf{D}_i$

If $N = T \log T$,
 $L(\pi_n) < T \epsilon_N + O(1)$
for some n



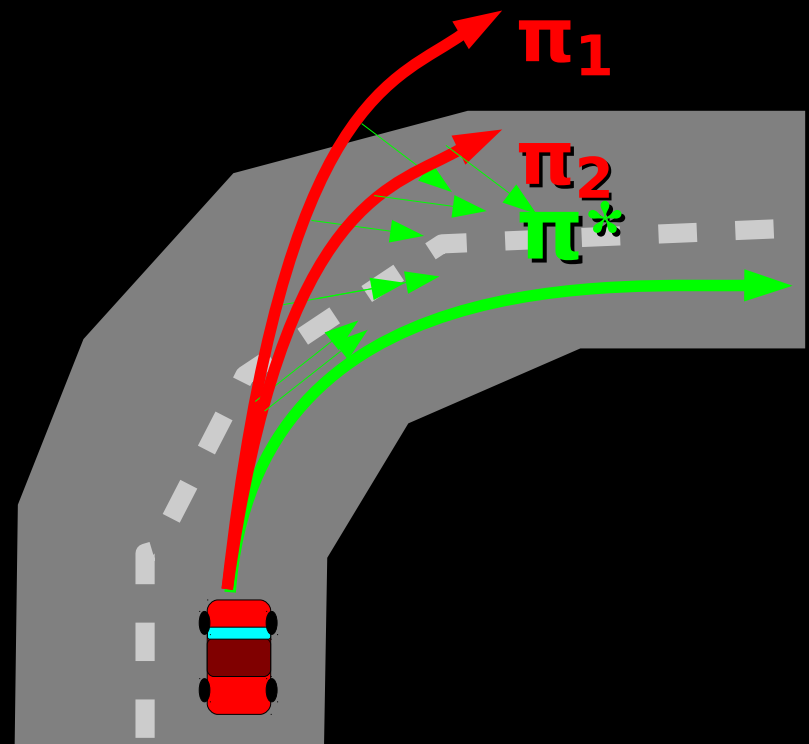
Test-time execution (Dagger)



What's the biggest failure mode?

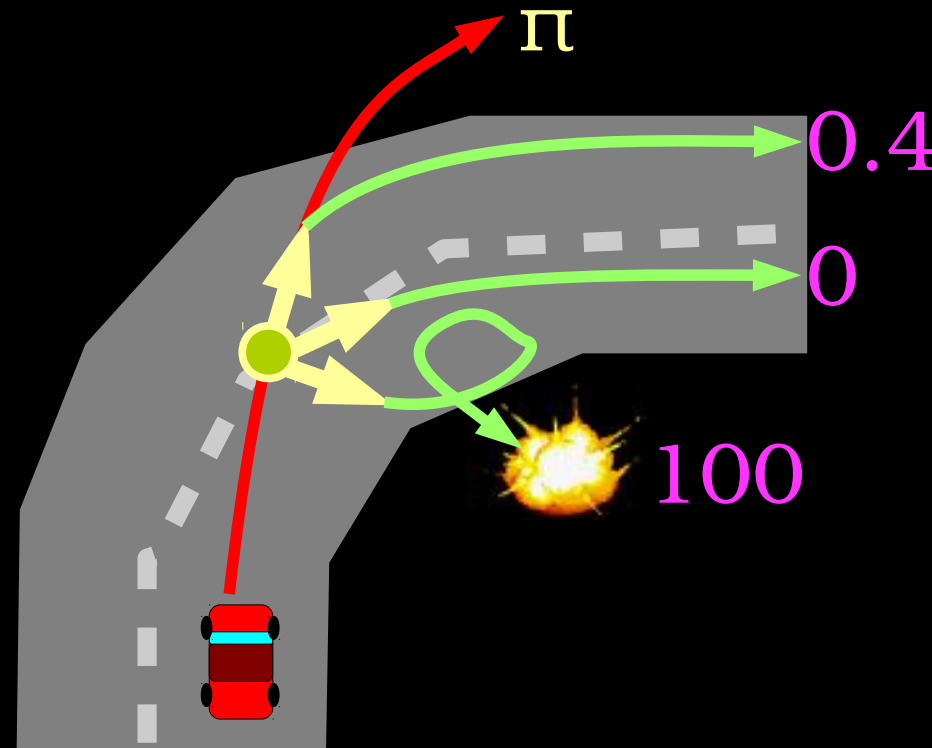
Classifier only sees *right* versus *not-right*

- No notion of *better* or *worse*
- No *partial credit*
- Must have a single *target* answer

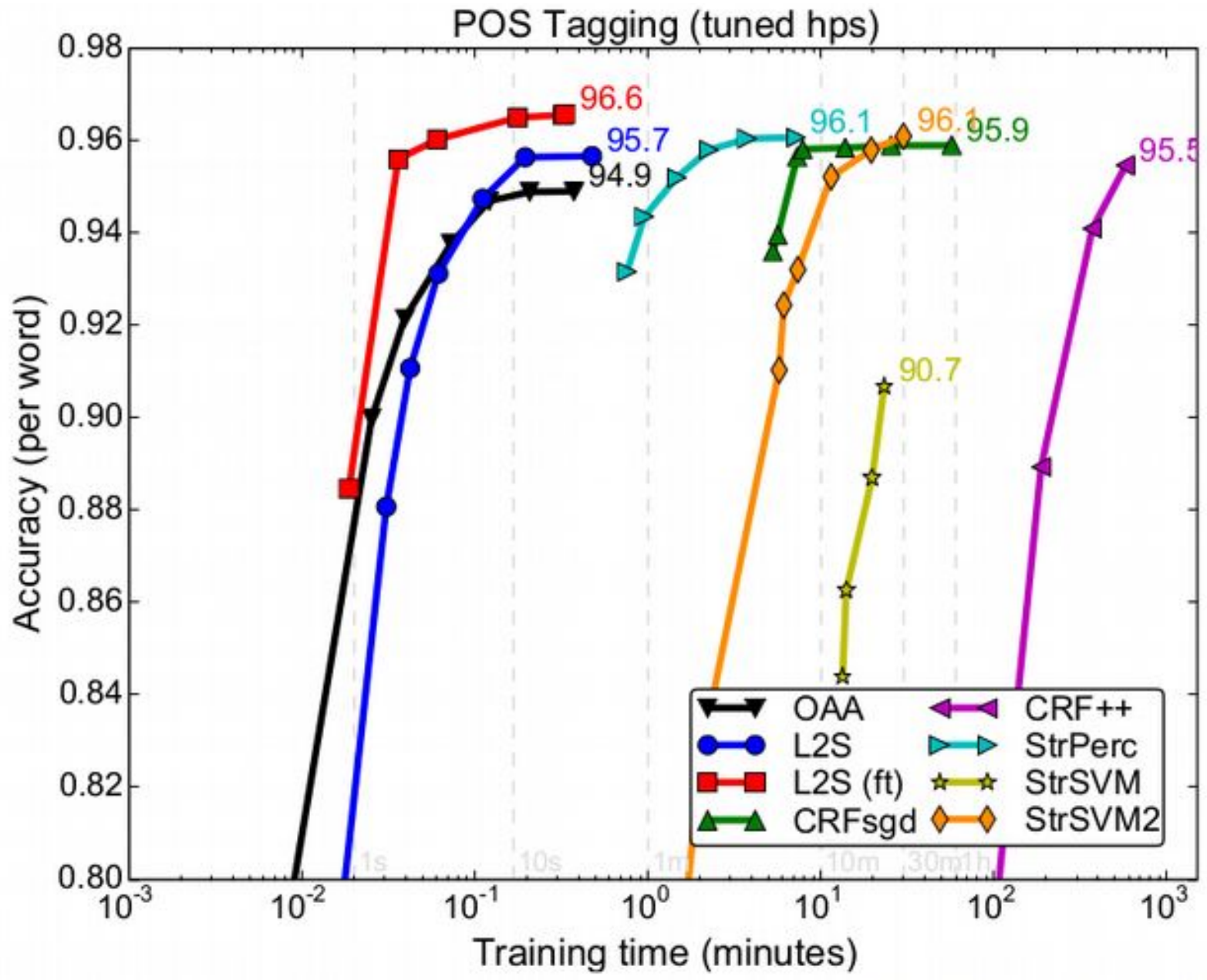


Learning to search: AggraVaTe

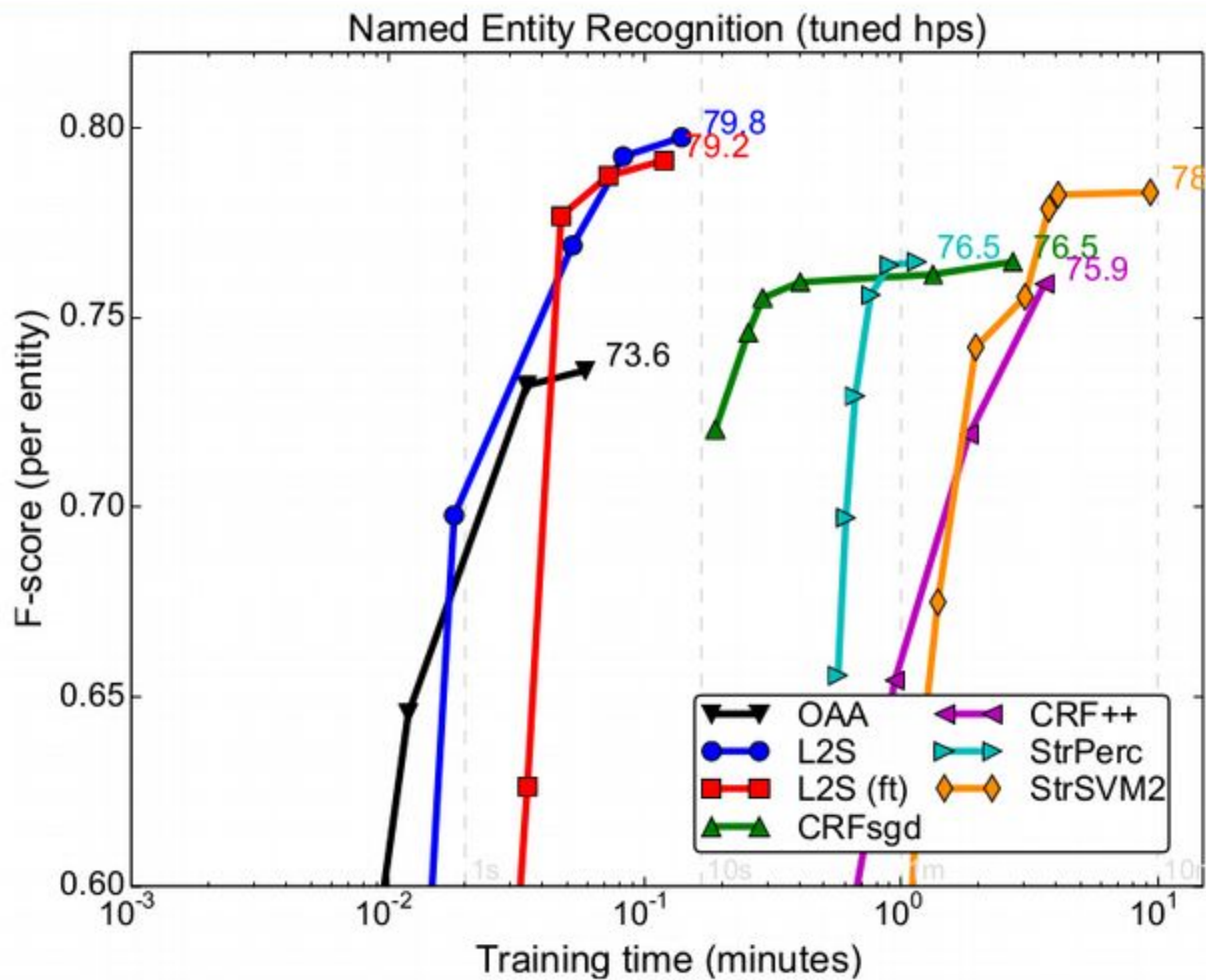
1. Let learned policy π drive for t timesteps to obs. o
2. For each possible action a :
 - Take action a , and let expert π^{ref} drive the rest
 - Record the overall loss, c_a
3. Update π based on example:
 $(o, \langle c_1, c_2, \dots, c_K \rangle)$
4. Goto (1)



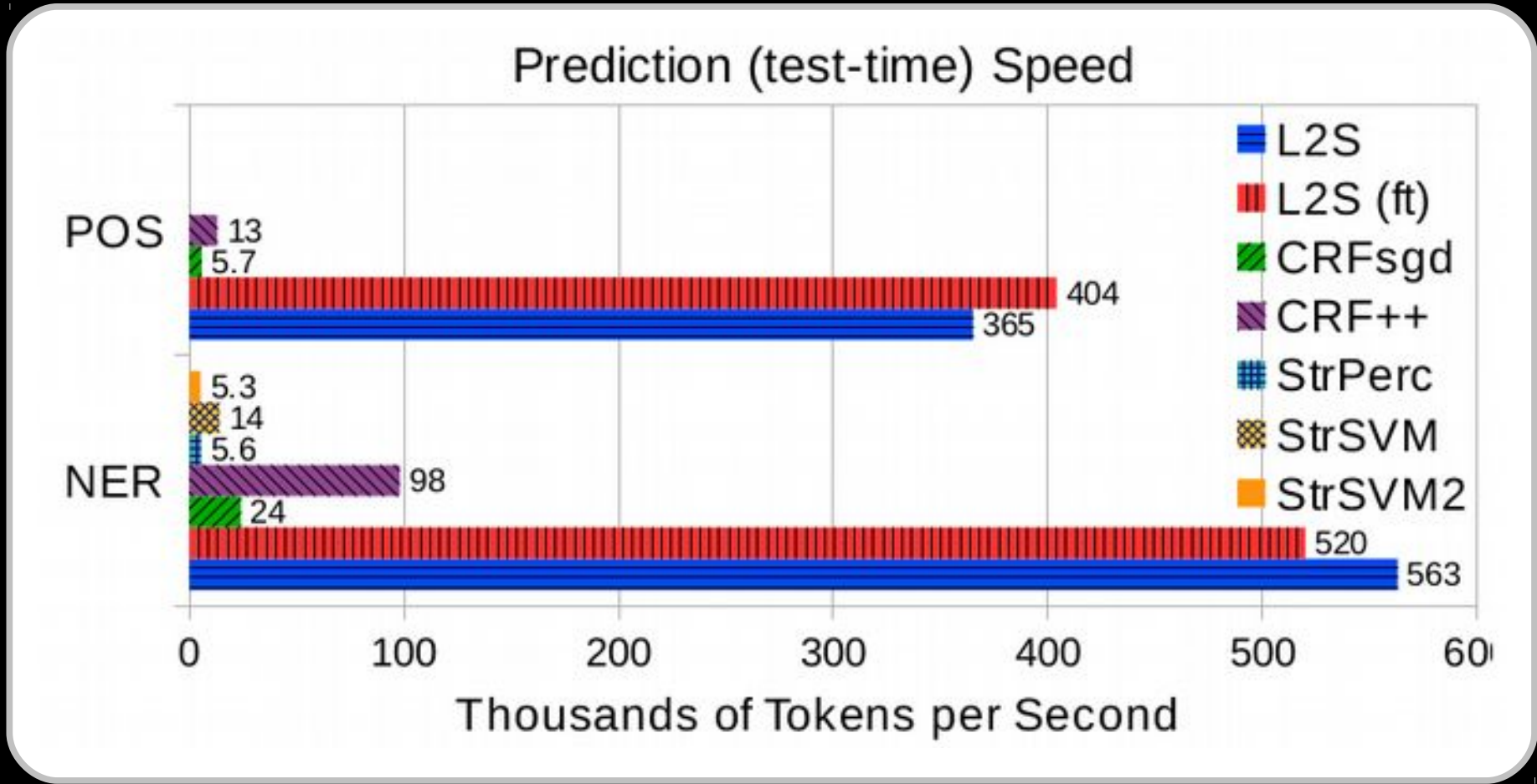
Training time versus test accuracy



Training time versus test accuracy



Test time speed



State of the art accuracy in...

- Part of speech tagging (1 million words)

- **US:** 6 lines of code 10 seconds to train
- CRFsgd: 1068 lines 30 minutes
- CRF++: 777 lines hours

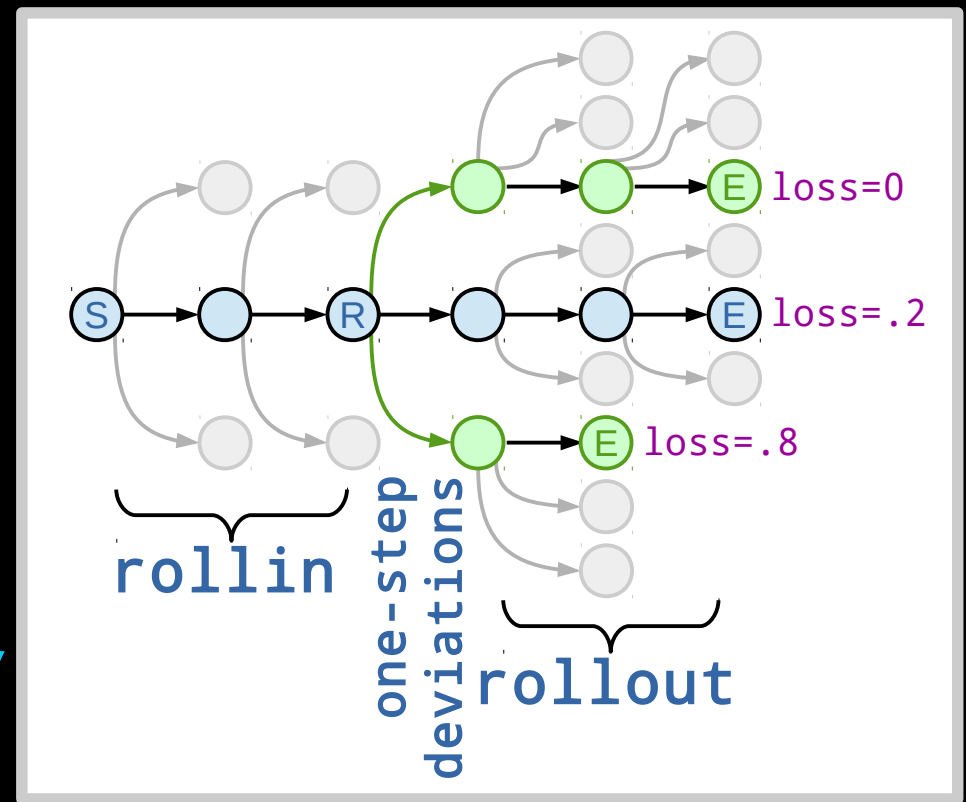
- Named entity recognition (200 thousand words)

- **US:** 30 lines of code 5 seconds to train
- CRFsgd: 1 minute
- CRF++: 10 minutes
- SVM^{str}: 876 lines 30 minutes (suboptimal accuracy)

The Magic

- You write some greedy “test-time” code
 - In your favorite imperative language (C++/Python)
 - It makes arbitrary calls to a `Predict` function
 - And you add some minor decoration
- We will automatically:
 - Perform learning
 - Generate non-deterministic (beam) search
 - Run faster than specialized learning software

How to train?



1. Generate **an initial trajectory** using a *rollin policy*

2. For each state **R** on that trajectory:

a) For each possible action a (one-step deviations)

i. Take that action

ii. Complete **this trajectory** using a rollout policy

iii. Obtain a **final loss**

b) Generate a cost-sensitive classification example:

$$(\Phi(R), \langle c_a \rangle_{a \in A})$$

The magic in practice

```
run(vector<example> ec)
  for i = 0 .. ec.size
    y_true = get_example_label(ec[i])
    y_pred = Predict(ec[i], y_true)
  Loss( # of y_true != y_pred )
```

A “hint” about the correct decision 😊
only at training time

not hiding anything...

```
void run(search& sch, vector<example*> ec) {
  for (size_t i=0; i<ec.size(); i++) {
    uint32_t y_true = get_example_label(ec[i]);
    uint32_t y_pred = sch.predict(ec[i], y_true);

    sch.loss( y_true != y_pred );

    if (sch.output().good())
      sch.output() << y_pred << ' ';
  }
}
```

The illusion of control

- Execute `run` $O(T \times A)$ times, modifying `Predict`
- For each time step `myT = 1 .. T`:

For each possible action `myA = 1 .. A`:

$$\text{define Predict}(\dots) = \begin{cases} \text{myA} & \text{if } t = \text{myT} \\ \pi & \text{otherwise} \end{cases}$$

run your code in full

set $\text{cost}_a = \text{result of Loss}$

Make classification example on x_{myT} with $\langle \text{cost}_a \rangle$

```
run(vector<example> ec)
  for i = 0 .. ec.size
    y_true = get_example_label(ec[i])
    y_pred = Predict(ec[i], y_true)
  Loss( # of y_true != y_pred )
```

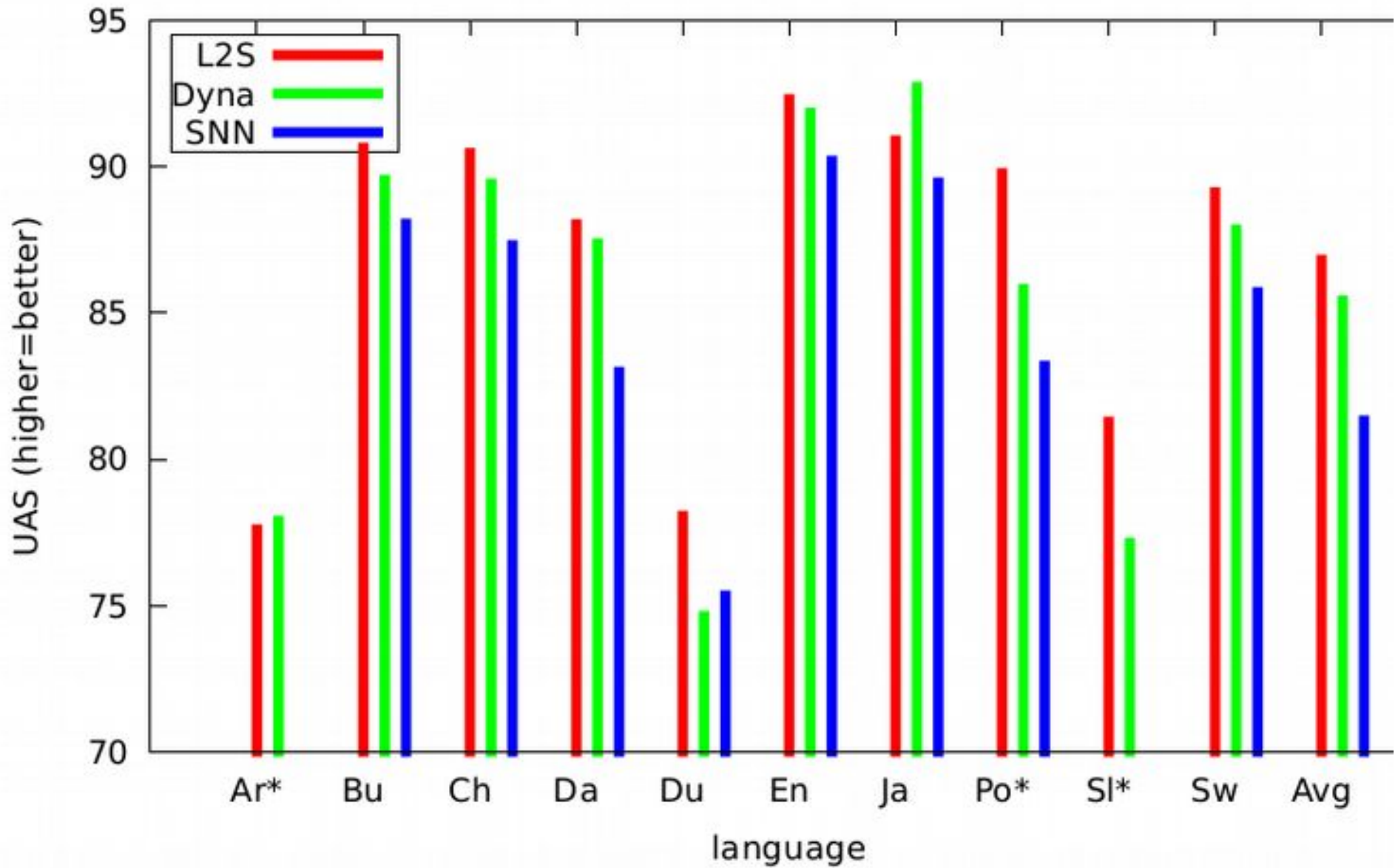
Entity/relation identification

Goal: find the Entities and then find their Relations

Method	Entity F1	Relation F1	Train Time
Structured SVM	88.00	50.04	300 seconds
L2S	92.51	52.03	13 seconds

L2S uses ~100 LOC.

Dependency parsing



L2S uses ~300 LOC.

Outline

- Background: learning to search
- Stuff I did in the Spring
 - Imperative DSL/library for learning to search
 - SOTA examples for tagging, parsing, relation extraction, etc.
 - Learning to search under bandit feedback
 - Hardness results for learning to search
 - Active learning for accelerating learning to search
- Stuff I'm trying to do now
 - Distant supervision
 - Mashups with recurrent neural networks

STRUCTURED CONTEXTUAL BANDIT

Font size

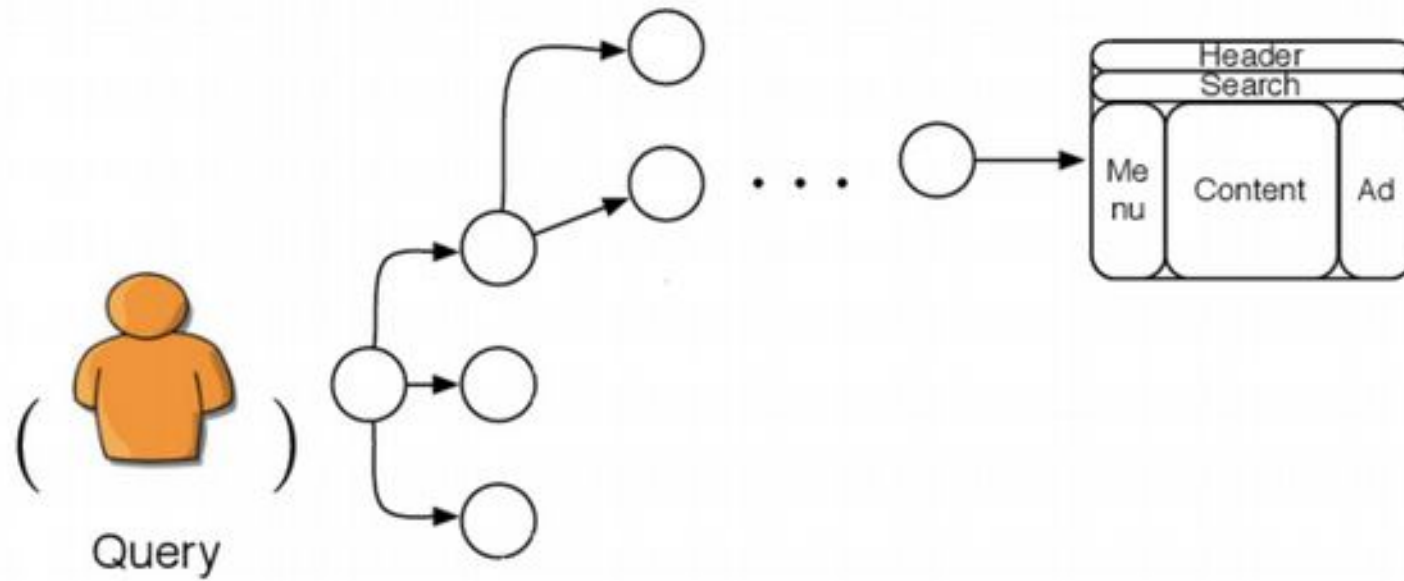
Color

Position

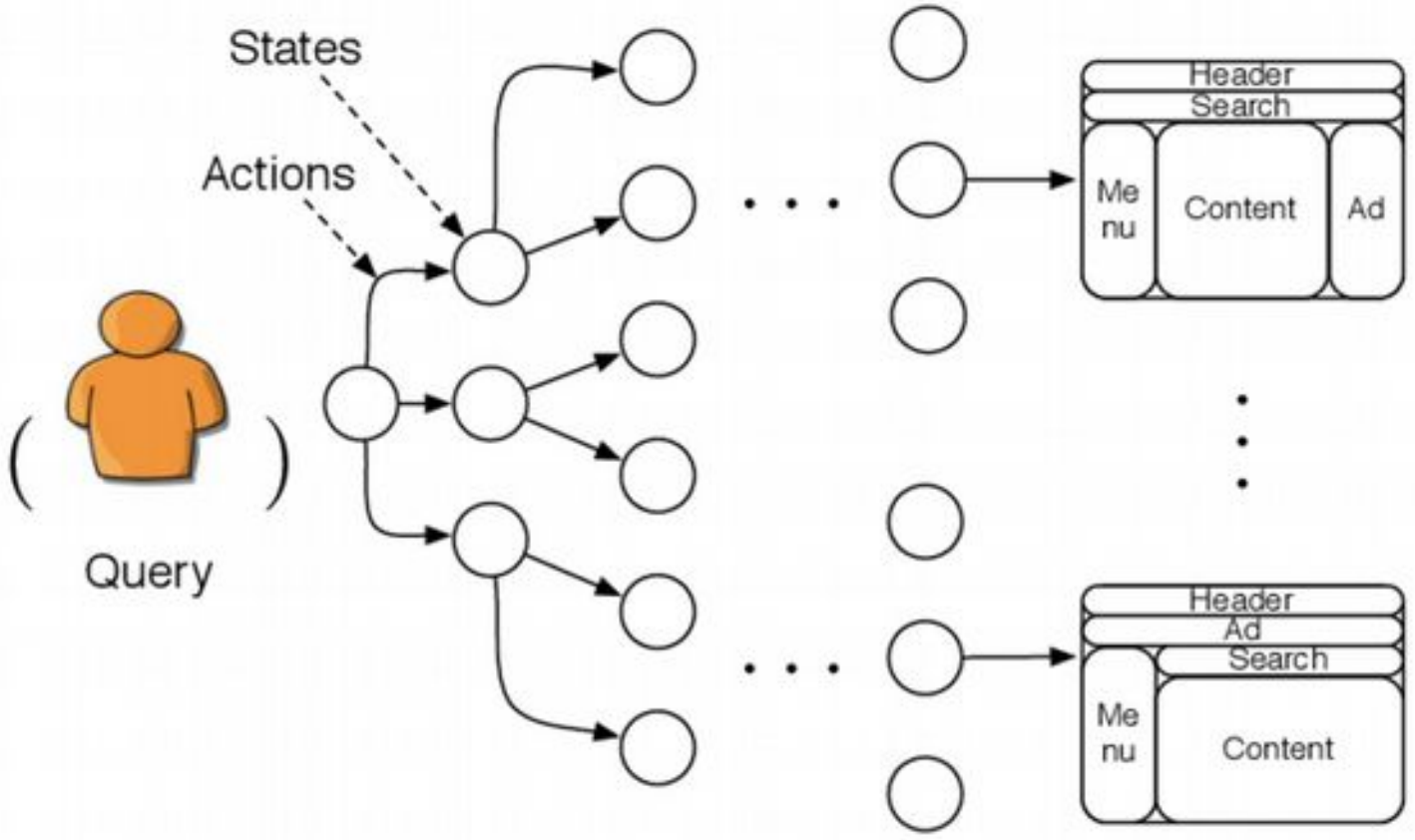


- ▶ Loss of a **single structured label** can be observed.

A SEARCH PROBLEM

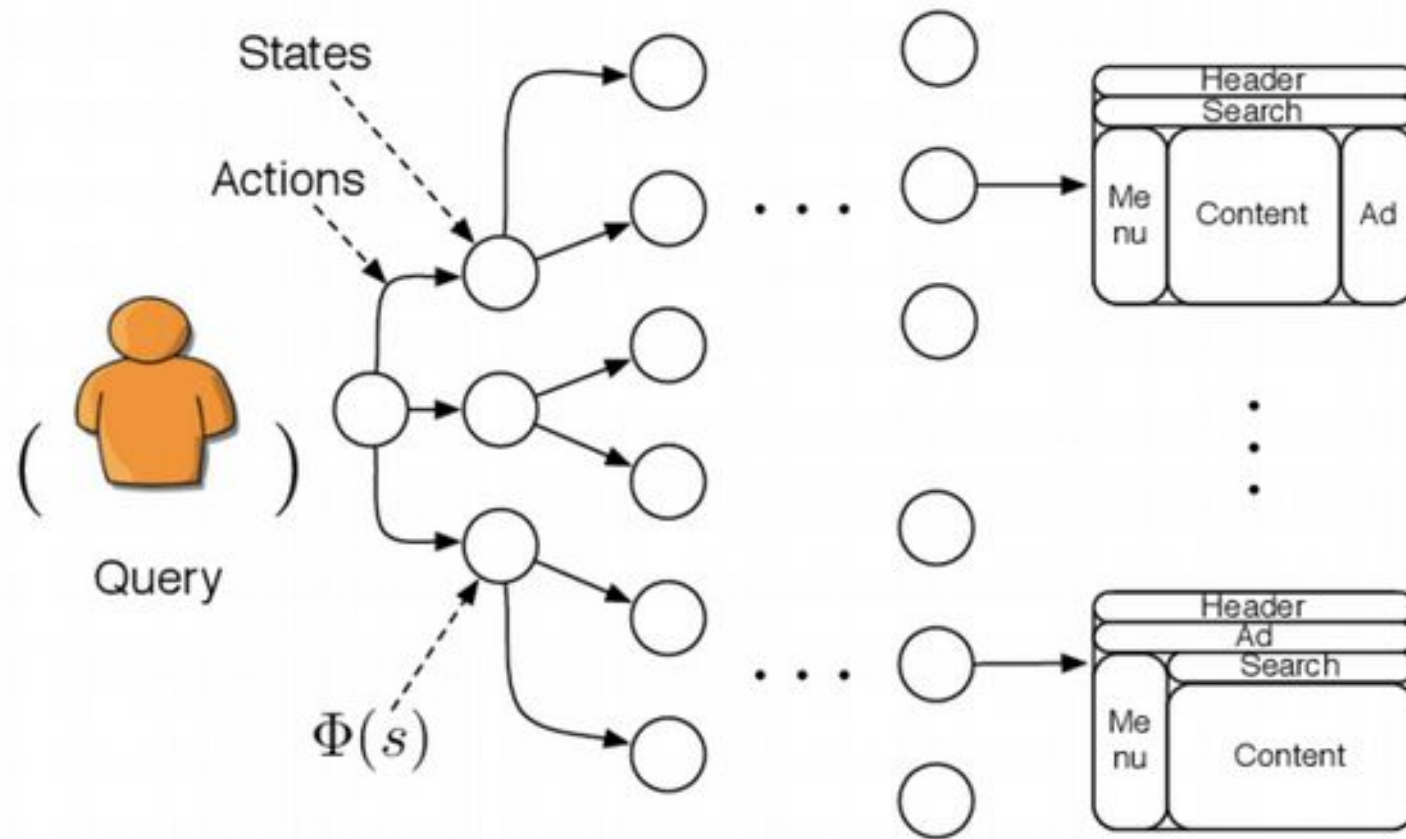


A SEARCH PROBLEM



Convert SCB into search problem (search space + actions).

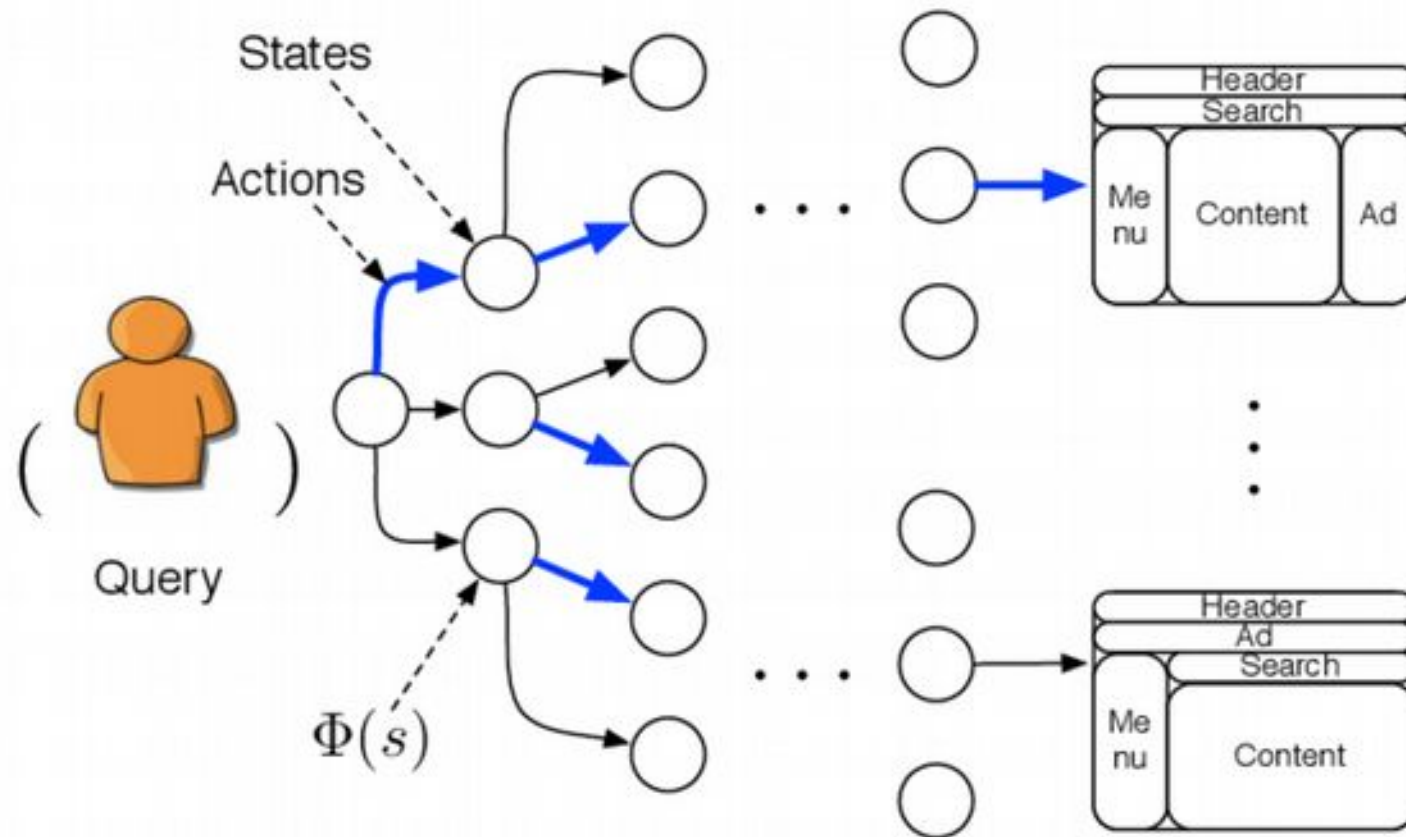
A SEARCH PROBLEM



Define structured features over each state.

- ▶ Learn policy mapping features to actions.

A SEARCH PROBLEM



Use status quo system as reference policy ([Ref](#)).

- ▶ Goal is to improve on [Ref](#).

CONNECTION WITH LEARNING TO SEARCH

	L2S	SCB

CONNECTION WITH LEARNING TO SEARCH

	L2S	SCB
Search Space	States, Actions, Features	

CONNECTION WITH LEARNING TO SEARCH

	L2S	SCB
Search Space	States, Actions, Features	
Feedback	Labels	Single Loss

CONNECTION WITH LEARNING TO SEARCH

	L2S	SCB
Search Space	States, Actions, Features	
Feedback	Labels	Single Loss
Ref Quality	Optimal	Can be bad

CONNECTION WITH LEARNING TO SEARCH

	L2S	SCB
Search Space	States, Actions, Features	
Feedback	Labels	Single Loss
Ref Quality	Optimal	Can be bad
Ref Availability	Training only	Train and Test

CONNECTION WITH LEARNING TO SEARCH

	L2S	SCB
Search Space	States, Actions, Features	
Feedback	Labels	Single Loss
Ref Quality	Optimal	Can be bad
Ref Availability	Training only	Train and Test
Goal	Imitate Ref	Improve upon Ref

CONNECTION WITH LEARNING TO SEARCH

	L2S	SCB
Search Space	States, Actions, Features	
Feedback	Labels	Single Loss
Ref Quality	Optimal	Can be bad
Ref Availability	Training only	Train and Test
Goal	Imitate Ref	Improve upon Ref

Existing L2S algorithms give:

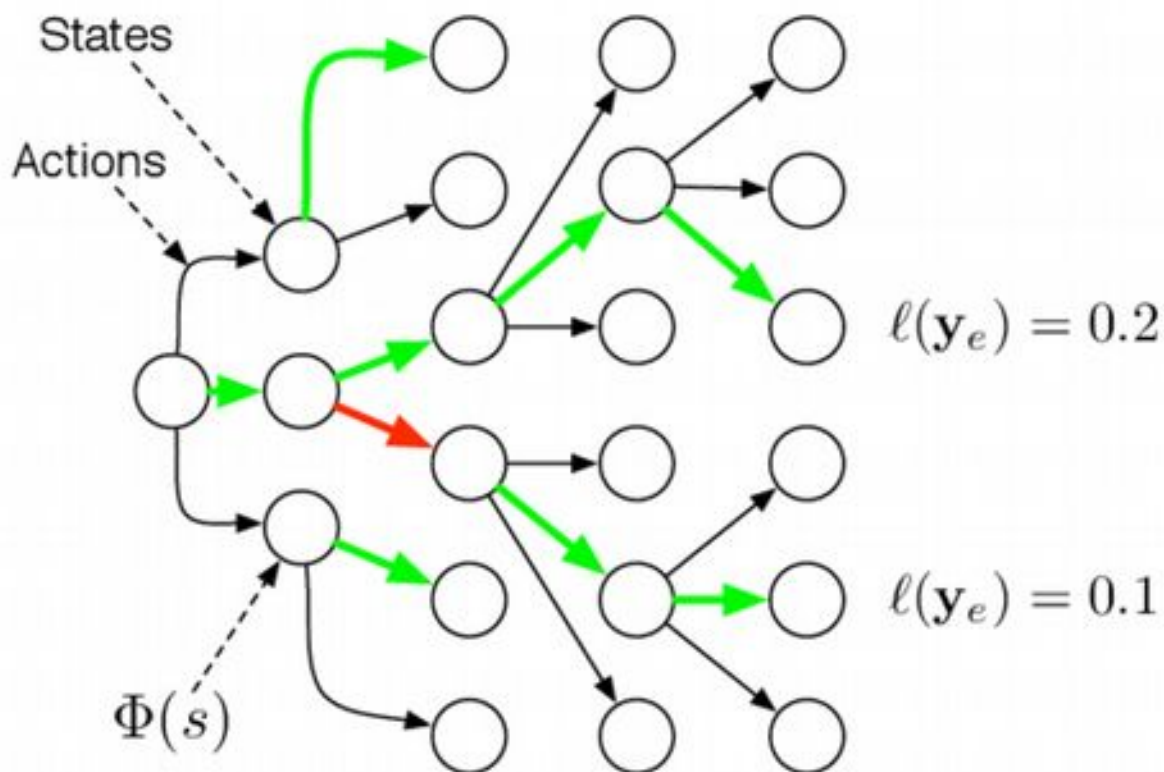
$$L(\pi) \leq L(\pi^{\text{ref}}) + o(1)$$

Learning to Search with:

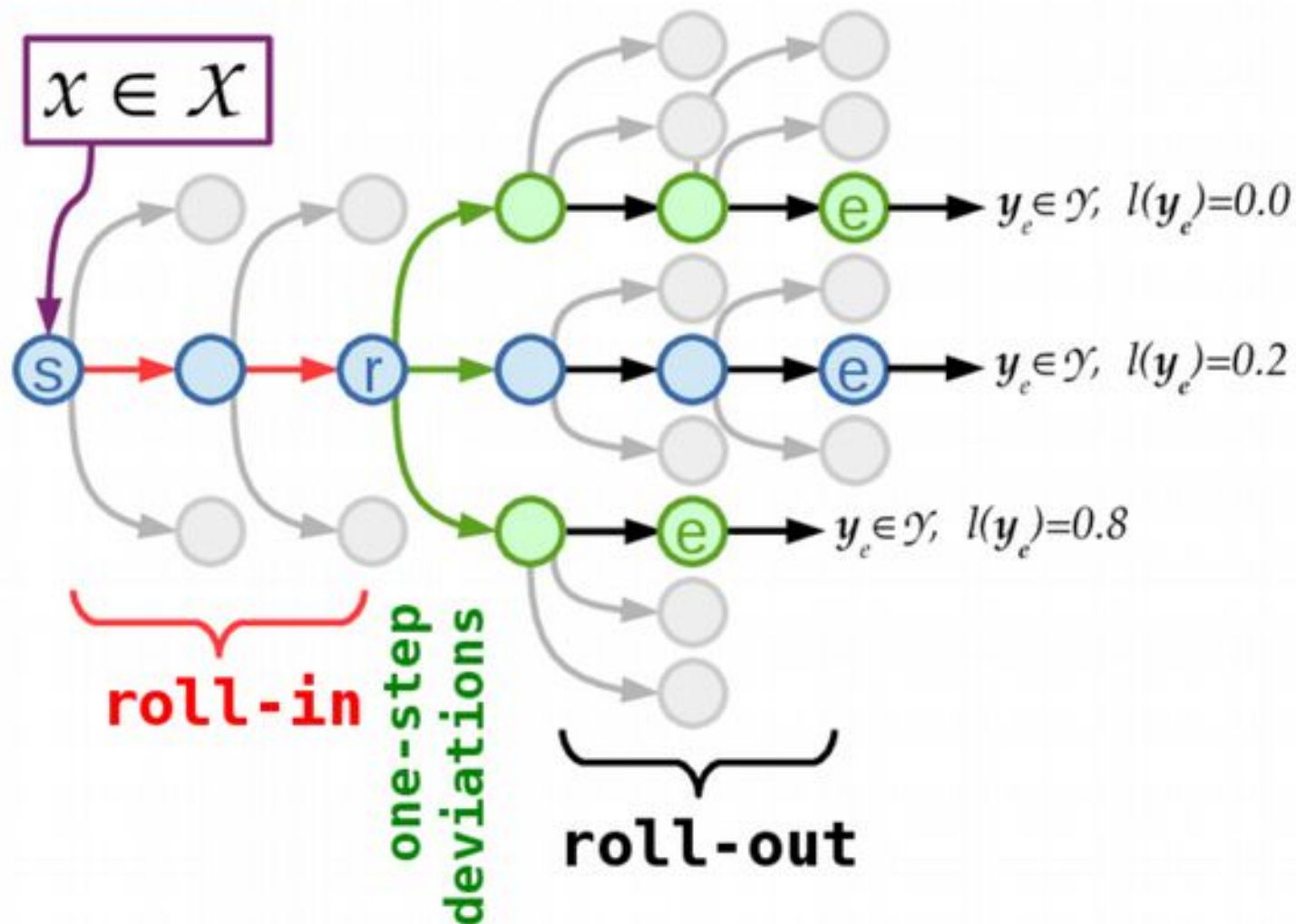
1. A suboptimal reference \Rightarrow Learn policy that **improves** on Ref.
2. Partial feedback.

DISEDERATA

- ▶ Compete with Ref.
 - ▶ Global optimality if Ref is optimal and realizable.
- ▶ Local optimality.
 - ▶ Compete with your own one-step deviations.



EXPLORATION IN SEARCH SPACE



- ▶ Roll-in choice affects what states we train on.
- ▶ Roll-out choice affects how we score actions.

EFFECT OF ROLL-IN AND ROLL-OUT POLICIES

roll-out → ↓ roll-in	Reference	Half/Half	Learned
Reference	Inconsistent		
Learned	No local opt	Good	RL

EFFECT OF ROLL-IN AND ROLL-OUT POLICIES

roll-out →	Reference	Half/Half	Learned
↓ roll-in			
Reference	Inconsistent		
Learned	No local opt	Good	RL

THEOREM

Roll-in with Ref can generate a model with unbounded structured regret but zero cost-sensitive regret.

- ▶ States trained on are not representative of those seen at prediction time.

EFFECT OF ROLL-IN AND ROLL-OUT POLICIES

roll-out → ↓ roll-in	Reference	Half/Half	Learned
Reference	Inconsistent		
Learned	No local opt	Good	RL

THEOREM

Roll-out with Ref may cause the learned policy to fail to converge to a local optimum if the reference policy is suboptimal.

- ▶ Causes poor assessment of comparison policies.

EFFECT OF ROLL-IN AND ROLL-OUT POLICIES

roll-out → ↓ roll-in	Reference	Half/Half	Learned
Reference	Inconsistent		
Learned	No local opt	Good	RL

THEOREM

Roll-in and roll-out with the learned policy ignores Ref and is equivalent to reinforcement learning.

EFFECT OF ROLL-IN AND ROLL-OUT POLICIES

roll-out → ↓ roll-in	Reference	Half/Half	Learned
Reference	Inconsistent		
Learned	No local opt	Good	RL

THEOREM

LOLS minimizes a combination of regret to Ref and regret to its own one-step deviations.

LOLS REGRET BOUND

THEOREM

*LOLS minimizes a combination of regret to **Ref** and regret to its own **one-step deviations**.*

$$\frac{1}{2} \underbrace{\left(L(\hat{\pi}) - L(\pi^{\text{ref}}) \right)}_{\text{Regret to Ref}} + \frac{1}{2} \underbrace{\left(L(\hat{\pi}) - L(\pi^{\text{dev}}) \right)}_{\text{Regret to devs}} \text{ is small.}$$

LOLS REGRET BOUND

THEOREM

*LOLS minimizes a combination of regret to **Ref** and regret to its own **one-step deviations**.*

$$\frac{1}{2} \underbrace{\left(L(\hat{\pi}) - L(\pi^{\text{ref}}) \right)}_{\text{Regret to Ref}} + \frac{1}{2} \underbrace{\left(L(\hat{\pi}) - L(\pi^{\text{dev}}) \right)}_{\geq 0} \text{ is small.}$$

- ▶ Competes with Ref.

LOLS REGRET BOUND

THEOREM

*LOLS minimizes a combination of regret to Ref and regret to its own **one-step deviations**.*

$$\frac{1}{2} \underbrace{\left(L(\hat{\pi}) - L(\pi^{\text{ref}}) \right)}_{\geq 0} + \frac{1}{2} \underbrace{\left(L(\hat{\pi}) - L(\pi^{\text{dev}}) \right)}_{\text{Regret to devs}} \text{ is small.}$$

- ▶ Competes with Ref .
- ▶ Locally optimal when Ref is optimal (even if unrealizable).

LOLS REGRET BOUND

THEOREM

*LOLS minimizes a combination of regret to Ref and regret to its own **one-step deviations**.*

$$\frac{1}{2} \underbrace{\left(L(\hat{\pi}) - L(\pi^{\text{ref}}) \right)}_{\text{Regret to Ref}} + \frac{1}{2} \underbrace{\left(L(\hat{\pi}) - L(\pi^{\text{dev}}) \right)}_{\text{Regret to devs}} \text{ is small.}$$

- ▶ Competes with Ref .
- ▶ Locally optimal when Ref is optimal (even if unrealizable).
- ▶ If Ref suboptimal, either locally optimal or better than Ref .

LOCAL OPTIMALITY IS HARD

Finding local optimum could be hard without further assumptions.

THEOREM

Can require $\Omega(2^T)$ cost-sensitive classification examples to reach local optimum.

T is the number of decisions per example.

DEPENDENCY PARSING

Find the dependency structure of words in a sentence.

roll-out → ↓ roll-in	Reference	Mixture	Learned
Reference is optimal			
Reference	87.2	89.7	88.2
Learned	90.7	90.5	86.9
Reference is suboptimal			
Reference	83.3	87.2	81.6
Learned	87.1	90.2	86.8
Reference is bad			
Reference	68.7	65.4	66.7
Learned	75.8	89.4	87.5

LOLS always good, even with suboptimal Ref.

STRUCTURED CONTEXTUAL BANDIT

- ▶ Loss of a **single structured label** can be observed.
- ▶ Reference policy is **not optimal** under this setting.

The image shows a screenshot of a Bing search results page. On the left side, three text labels with arrows point to specific elements on the page:

- Font size**: An arrow points to the "News" tab in the navigation bar.
- Color**: An arrow points to the main headline "Trial to start for Maine man in alleged social media murder".
- Position**: An arrow points to the top-left corner of the main article's content area.

The screenshot itself displays the following content:

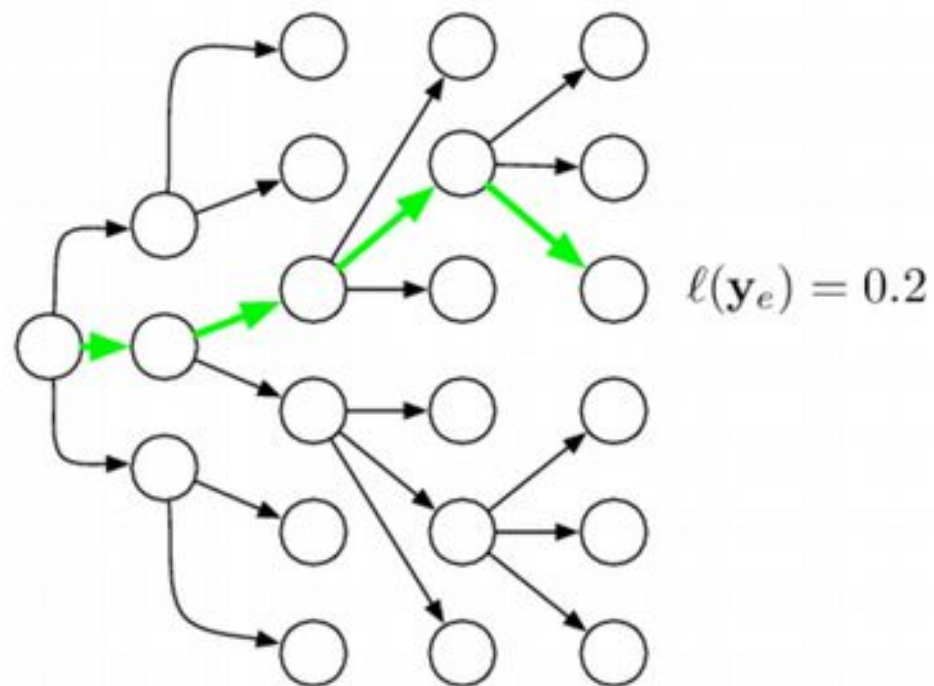
- Search bar with "bing" and a magnifying glass icon.
- Navigation tabs: Web, Images, Videos, Maps, **News**, More.
- Sub-navigation tabs: Top stories, U.S., World, Local, Entertainment, Sports, Business, Politics, Sports, Health.
- Section: "Today's Headlines"
- Main article: "Trial to start for Maine man in alleged social media murder" with a photo of a woman and a snippet of text.
- Other articles: "China Protests India Leader's Visit to Disputed Border Area" and "Secrecy around police surveillance equipment proves ...".
- Section: "Trending Topics" at the bottom.

STRUCTURED CONTEXTUAL BANDIT

We adapt an ϵ -greedy algorithm

Upon receiving an example:

- ▶ w/ prob. $1 - \epsilon$: follow the current policy

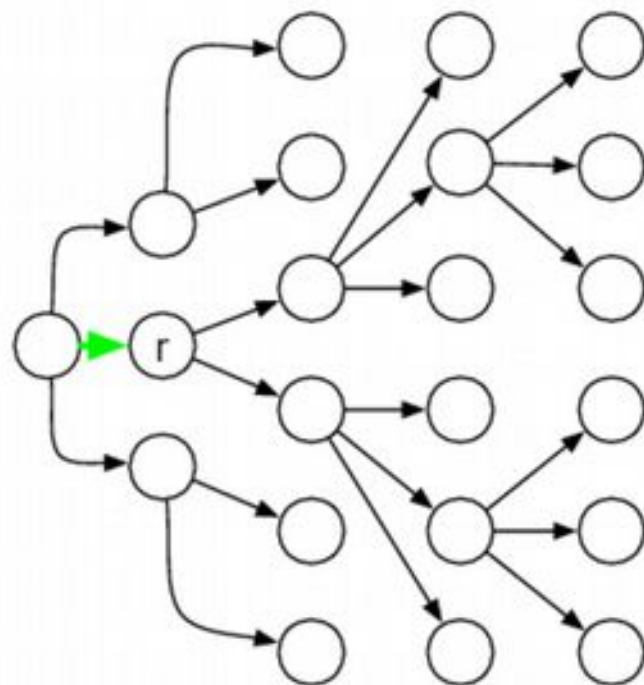


STRUCTURED CONTEXTUAL BANDIT

We adapt an ϵ -greedy algorithm

Upon receiving an example:

- ▶ w/ prob. $1 - \epsilon$: follow the current policy
- ▶ w/ prob. ϵ : perform a randomized LOLS update.

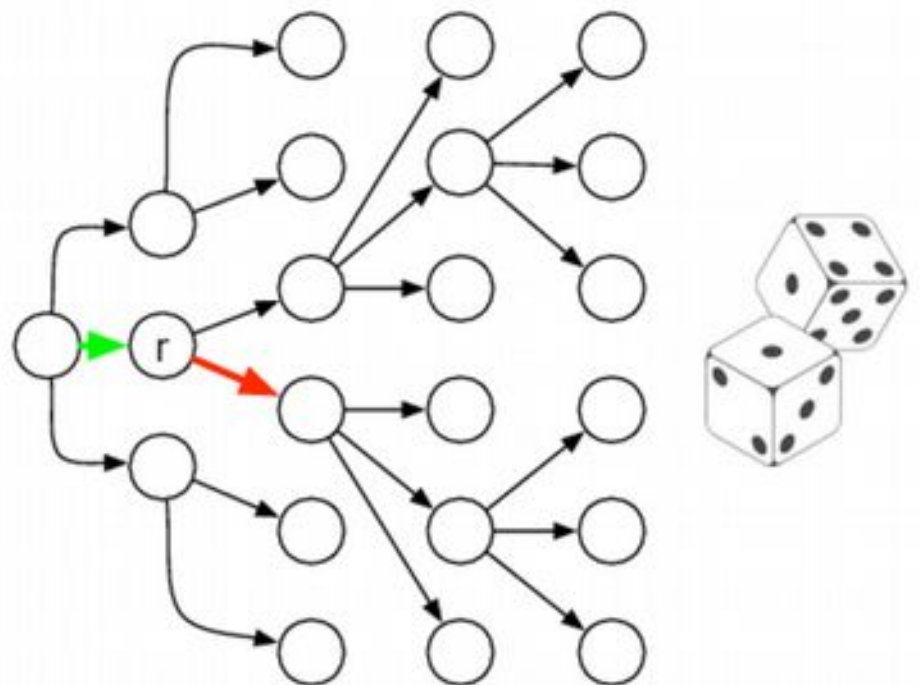


STRUCTURED CONTEXTUAL BANDIT

We adapt an ϵ -greedy algorithm

Upon receiving an example:

- ▶ w/ prob. $1 - \epsilon$: follow the current policy
- ▶ w/ prob. ϵ : perform a randomized LOLS update.

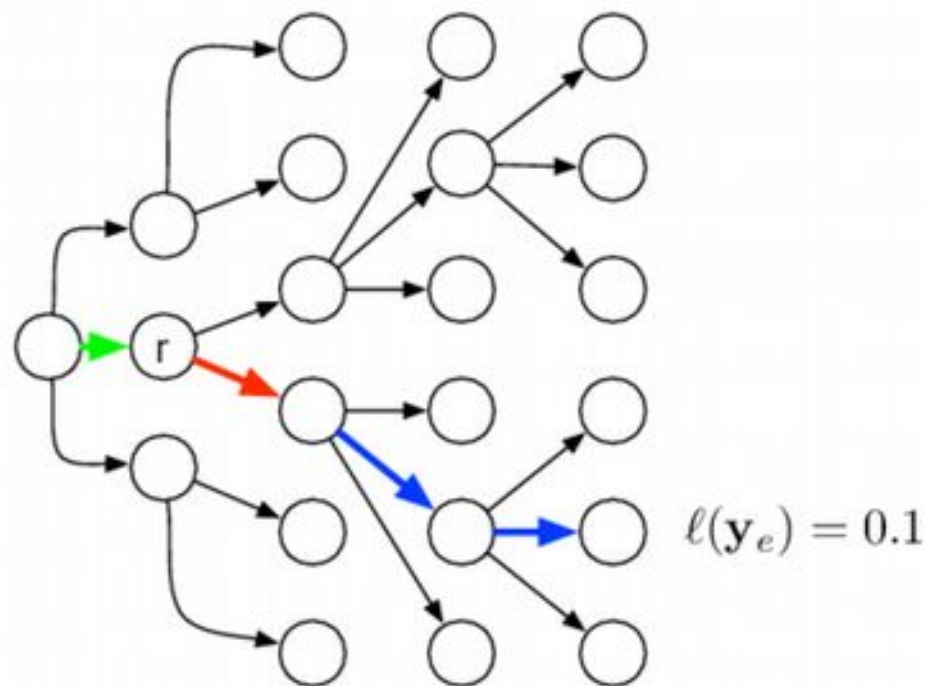


STRUCTURED CONTEXTUAL BANDIT

We adapt an ϵ -greedy algorithm

Upon receiving an example:

- ▶ w/ prob. $1 - \epsilon$: follow the current policy
- ▶ w/ prob. ϵ : perform a randomized LOLS update.

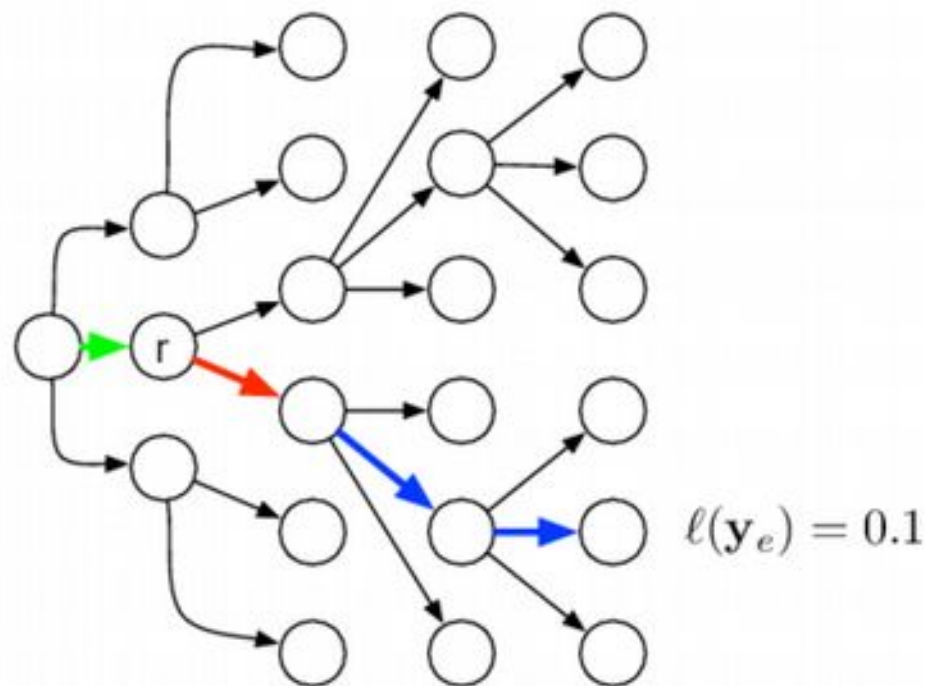


STRUCTURED CONTEXTUAL BANDIT

We adapt an ϵ -greedy algorithm

Upon receiving an example:

- ▶ w/ prob. $1 - \epsilon$: follow the current policy
- ▶ w/ prob. ϵ : perform a randomized LOLS update.



Regret against ref and deviations is still bounded.

Outline

- Background: learning
- Stuff I did in the Spring
 - Imperative DSL/library
 - SOTA examples for
 - Learning to search un
 - Hardness results for learn
 - Active learning for accelerating
- Stuff I'm trying to do now
 - Distant supervision
 - Mashups with recurrent neural networks

Observation: rollouts at all time steps not equally useful

Solution: importance-weighted active learning selection of where to rollout vs skip

Hacky heuristic: 5* speedup, slightly increased accuracy

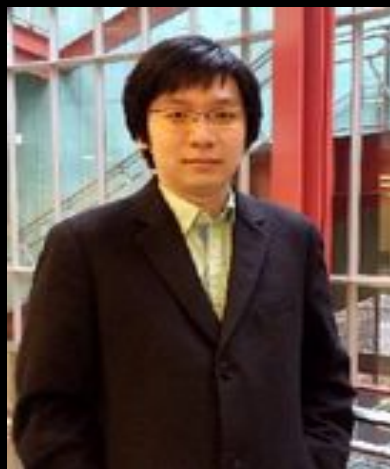
Training RNNs with LOLS yields drastic increases in performance on non-adversarial synthetic data

Distant supervision

- Learning with a human in the loop
- Repeat forever:
 - Information need
 - Machine makes complex prediction
 - Human is happy or unhappy, provides extra feedback
 - Machine learns
 - **Human learns**
- How to handle the last step?



Alekh
Agarwal



Kai-Wei
Chang



Akshay
Krishnamurthy



John
Langford



Alina
Beygelzimmer




Paul
Mineiro



Stéphane
Ross



He
He

- 
- Novel programming paradigm for integrating ML into software
 - State of the art results on many tasks, very quickly, little code
 - New problems, new algorithms
 - Positive results (notion of local optimality, and regret guarantees)
 - Negative results (hardness of exact local optimality)
 - Lots of places to go from here...

Thanks! Questions?