



Reinforcement Learning for NLP

Advanced Machine Learning for NLP

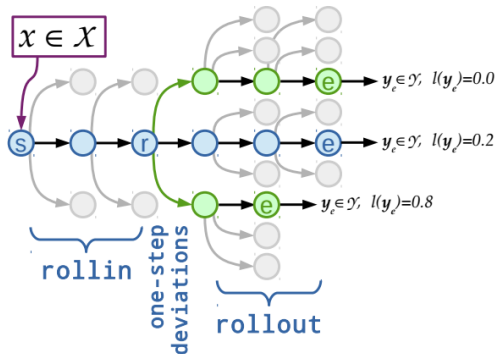
Jordan Boyd-Graber

DEEP SHIFT-REDUCE PARSERS

What Makes NLP different from RL?

- Often, best actions are **known**
- We're not just searching for high-reward
- Sometimes actions themselves are known

Roll In vs. Roll Out



- Roll In: Which states does the algorithm see
- Roll Out: What states do you use for training

Known Policy vs. Exploration

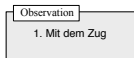
roll-out →	Reference	Mixture	Learned
↓ roll-in			
Reference	Inconsistent		
Learned	Not locally opt.	Good	RL

- RL only gets reward
- Roll-in with reference gives unrealistic trajectories
- How to incorporate knowledge of true actions?
- Train classifier as proxy for policy

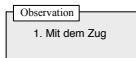
RL for Translation



RL for Translation

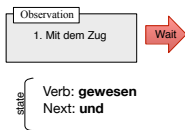


RL for Translation

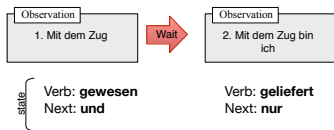


state {
Verb: **gewesen**
Next: **und**

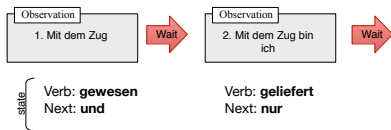
RL for Translation



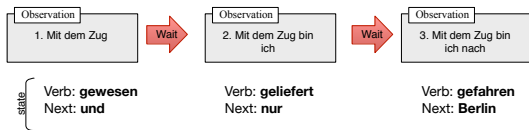
RL for Translation



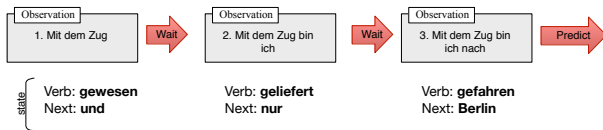
RL for Translation



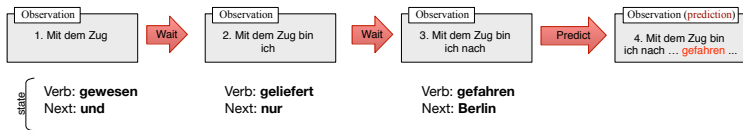
RL for Translation



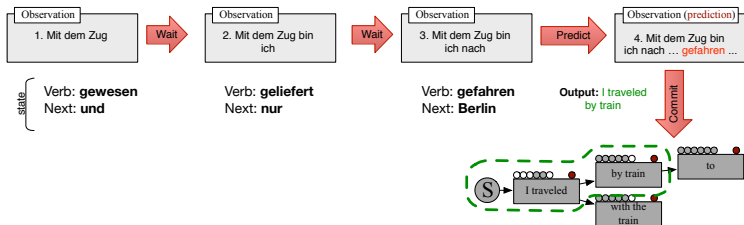
RL for Translation



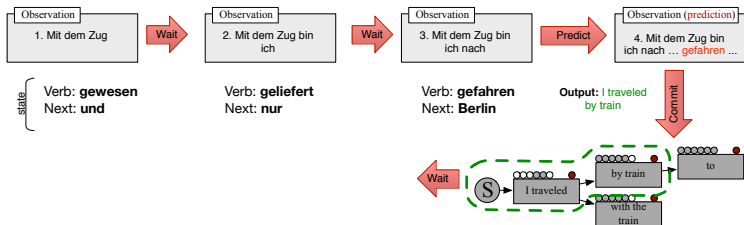
RL for Translation



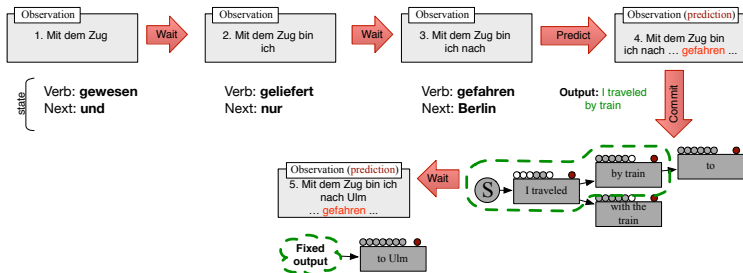
RL for Translation



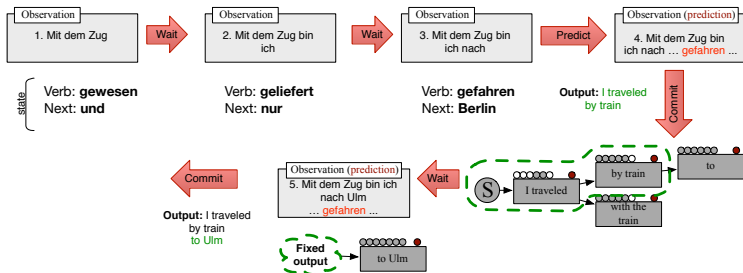
RL for Translation



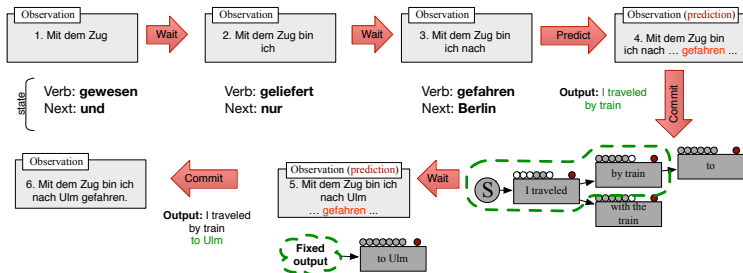
RL for Translation



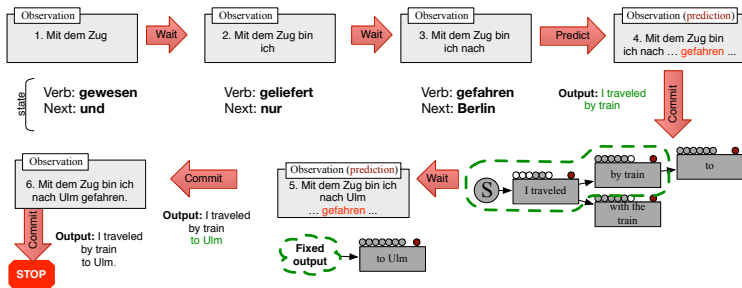
RL for Translation



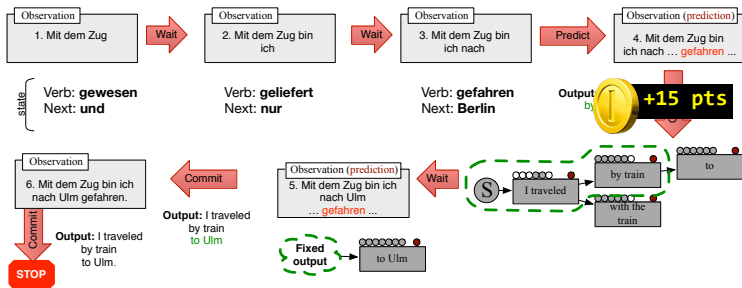
RL for Translation



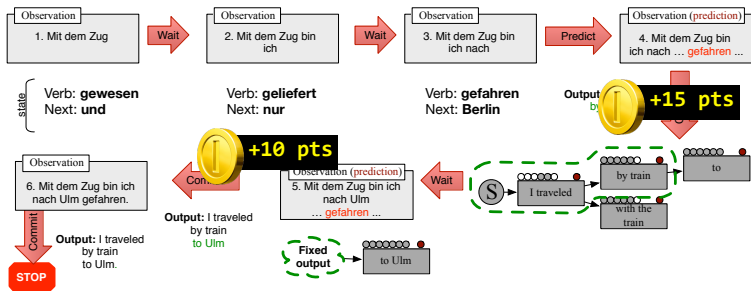
RL for Translation



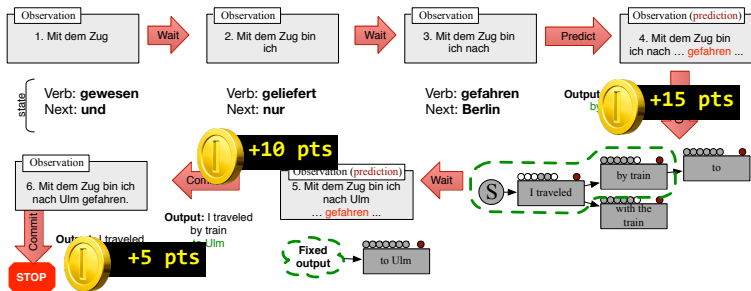
RL for Translation



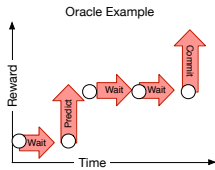
RL for Translation



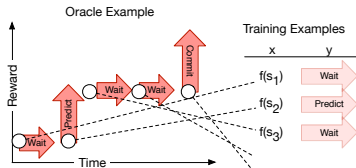
RL for Translation



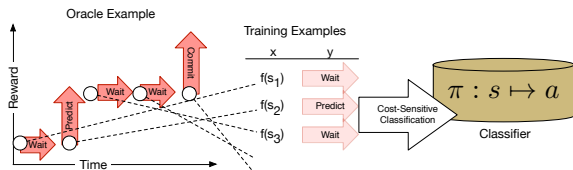
How do we find a good policy?



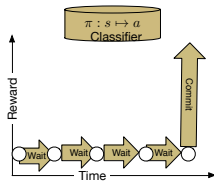
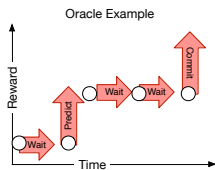
How do we find a good policy?



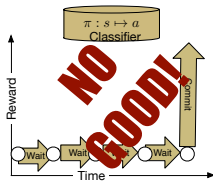
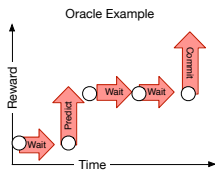
How do we find a good policy?



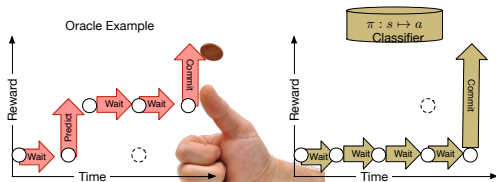
How do we find a good policy?



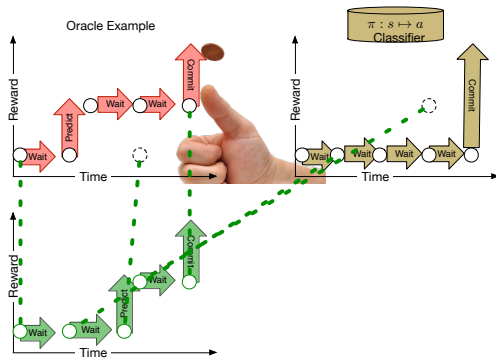
How do we find a good policy?



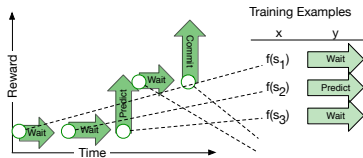
How do we find a good policy?



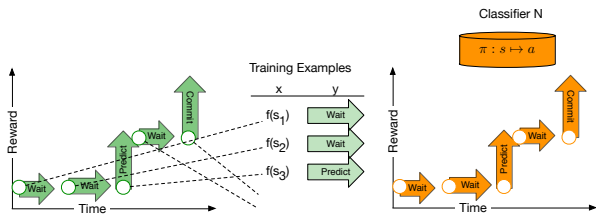
How do we find a good policy?



How do we find a good policy?



How do we find a good policy?



Algorithm 1 Locally Optimal Learning to Search (LOLS)

Require: Dataset $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$ drawn from \mathcal{D} and $\beta \geq 0$: a mixture parameter for roll-out.

- 1: Initialize a policy π_0 .
 - 2: **for all** $i \in \{1, 2, \dots, N\}$ (loop over each instance) **do**
 - 3: Generate a reference policy π^{ref} based on \mathbf{y}_i .
 - 4: Initialize $\Gamma = \emptyset$.
 - 5: **for all** $t \in \{0, 1, 2, \dots, T - 1\}$ **do**
 - 6: Roll-in by executing $\pi_i^{\text{in}} = \hat{\pi}_i$ for t rounds and reach s_t .
 - 7: **for all** $a \in A(s_t)$ **do**
 - 8: Let $\pi_i^{\text{out}} = \pi^{\text{ref}}$ with probability β , otherwise $\hat{\pi}_i$.
 - 9: Evaluate cost $c_{i,t}(a)$ by rolling-out with π_i^{out} for $T - t - 1$ steps.
 - 10: **end for**
 - 11: Generate a feature vector $\Phi(\mathbf{x}_i, s_t)$.
 - 12: Set $\Gamma = \Gamma \cup \{(c_{i,t}, \Phi(\mathbf{x}_i, s_t))\}$.
 - 13: **end for**
 - 14: $\hat{\pi}_{i+1} \leftarrow \text{Train}(\hat{\pi}_i, \Gamma)$ (Update).
 - 15: **end for**
 - 16: Return the average policy across $\hat{\pi}_0, \hat{\pi}_1, \dots, \hat{\pi}_N$.
-

LOLS on Dependency Parsing

Policy is learning actions for shift-reduce parser

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

But what structure is best?

- RecNN not much better than DAN
- But syntax may not be optimal
- Can we learn structure?

But what structure is best?

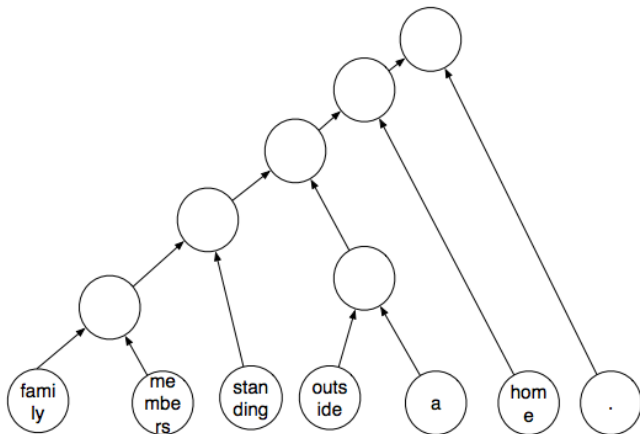
- RecNN not much better than DAN
- But syntax may not be optimal
- Can we learn structure?
 - Policy learns shift-reduce parser
 - TreeLSTM with **learned structure**
 - Reward is performance on downstream task

Performance

Table 4: Classification accuracy on SNLI dataset.

Model	Acc.	# params.
100D-Right to left	79.1	2.3m
100D-Left to right	80.2	2.3m
100D-Bidirectional	80.2	2.6m
100D-Supervised syntax	78.5	2.3m
100D-Semi-supervised syntax	80.2	2.3m
100D-Latent syntax	80.5	2.3m

What structures?



Other places of NLP + RL

- Question answering
- Language games
- Dialog systems
- Human learning

Wrapup

- RL allows for algorithms to think about long-term rewards
- And to guide actions of a system
- Important for systems that interact with world
- Discrete action spaces often more difficult