

Facebook AI Research

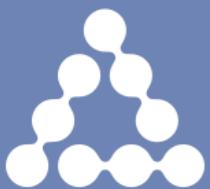
Task-Oriented Dialogues

Denis Yarats



Facebook AI Research



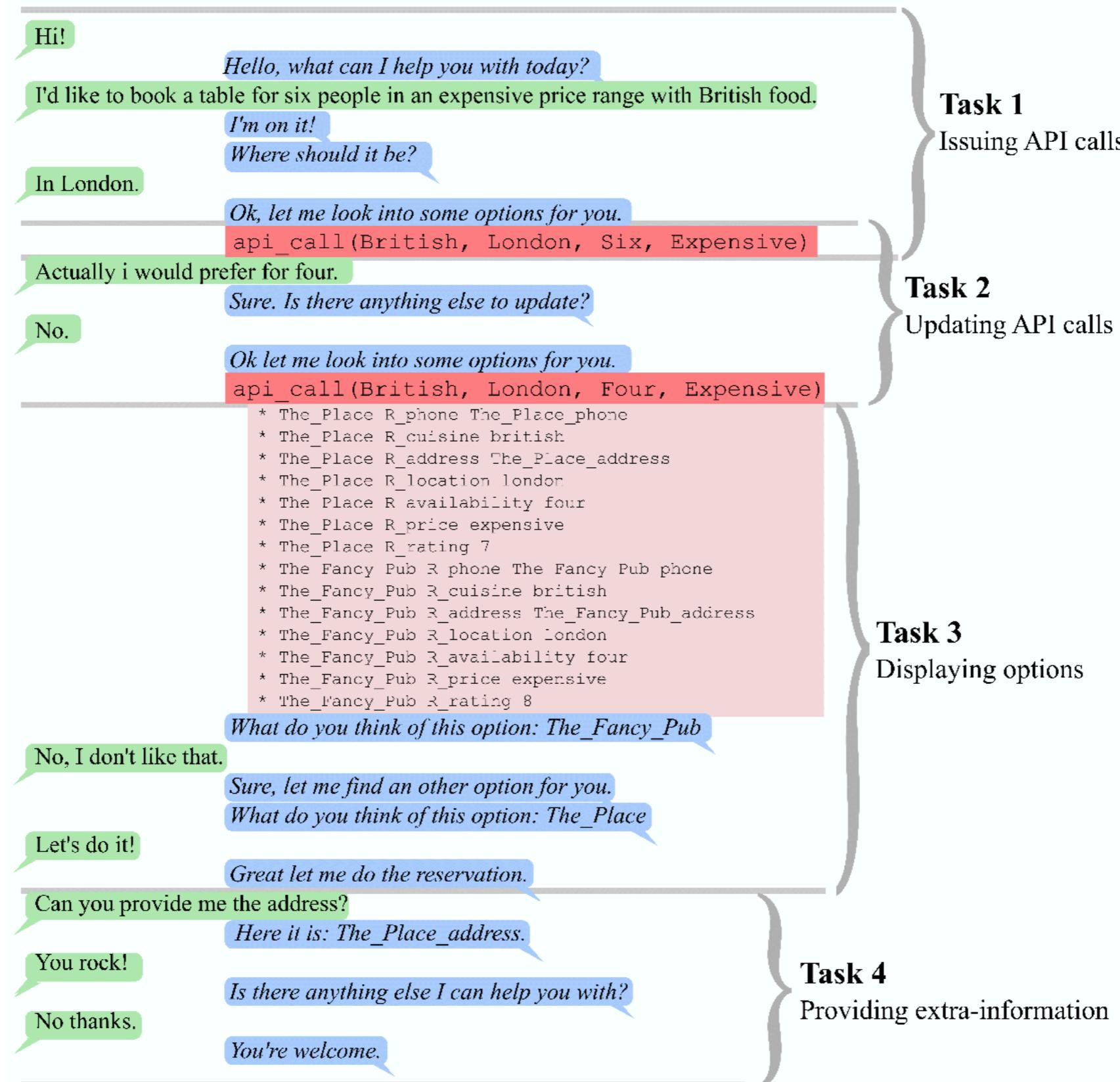


Overview

- Intro to Task-Oriented Dialogues
- Negotiation Task
 - Dataset
 - Modeling
- Hierarchical Generation



Task-Oriented Dialogue





Task-Oriented Dialogue

Formal Definition:

$D = \{(x_k, y_k)\}_{k=1}^K$ - dialogue

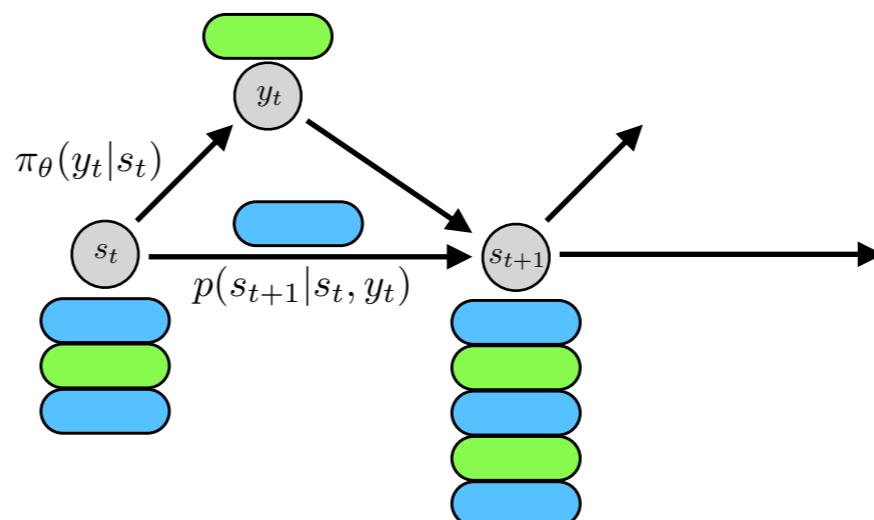
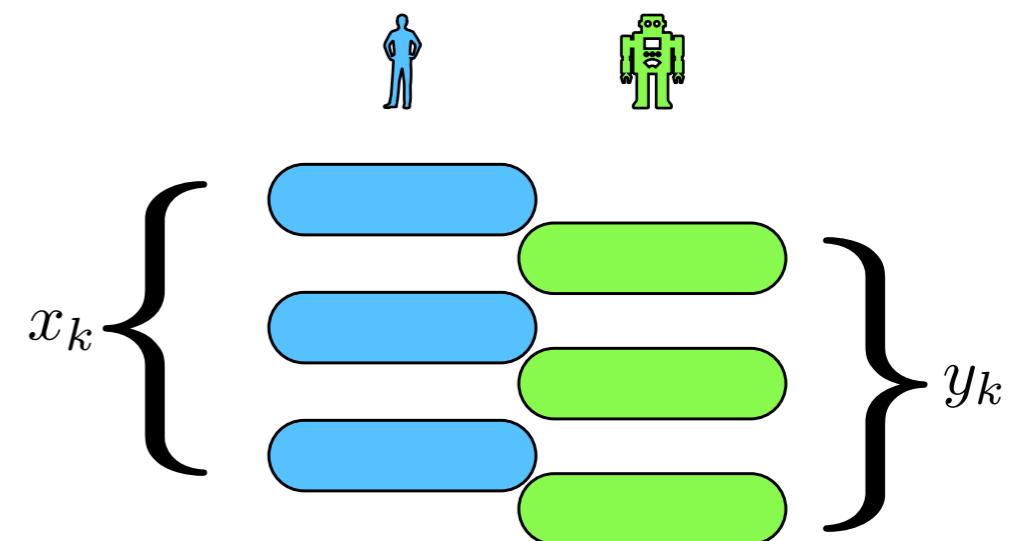
$s_t = \{x_1, y_1, x_2, y_2, \dots, x_t\}$ - state

$\mathcal{S} = \{s_t\}$ $\mathcal{A} = \{y_t\}$ - state and action spaces

$p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ - transition function

$\pi_\theta : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ - policy

$r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ - reward





Task-Oriented Dialogue

Formal Definition:

$\tau = \{s_1, y_1, s_2, y_2, \dots, s_T, y_T\}$ - trajectory of a dialogue

$\pi_\theta(\tau) = p(s_1) \prod_{t=1}^T \pi_\theta(y_t|s_t)p(s_{t+1}|s_t, y_t)$ - distribution over trajectory

$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_t r(s_t, y_t) \right]$ - cumulative reward

Goal:

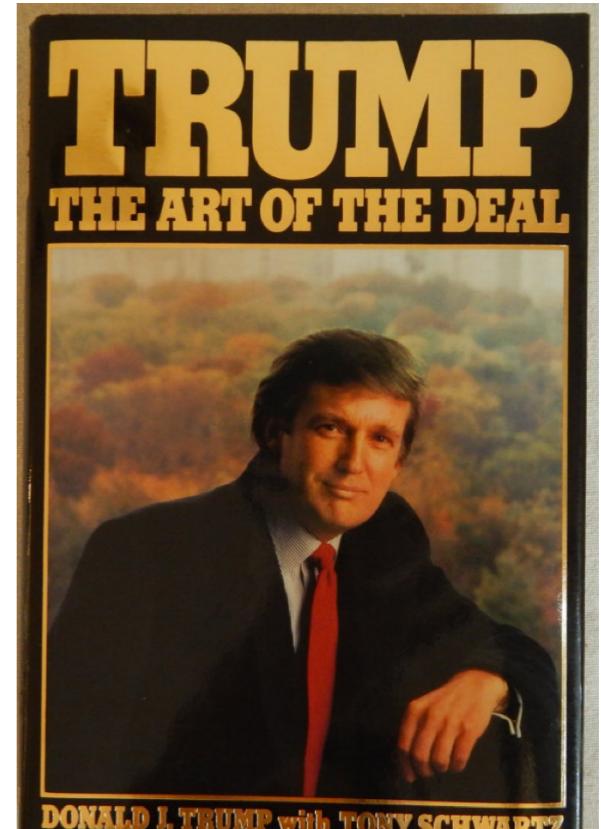
$\theta^* = \arg \max_{\theta} J(\theta) = \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_t r(s_t, y_t) \right]$



Negotiation

Why Negotiation?

- Both **linguistic** and **reasoning** problem
 - **interpret** multiple sentences and **generate** new message
 - plan ahead, make proposals, ask questions, bluff, compromise





Negotiation

Why Negotiation?

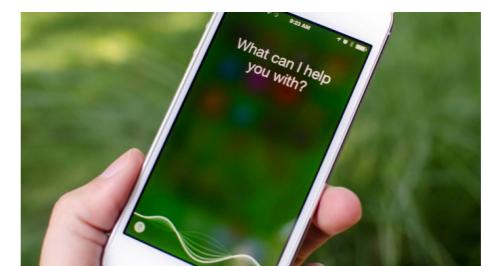
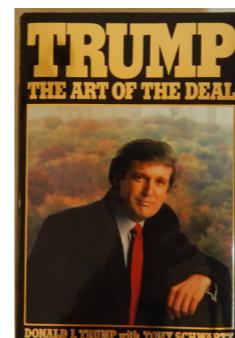
- Unlike many task-oriented dialogue problems, **no simple solution** to achieving goal
- But, it is **very easy** to evaluate



Zero-sum /
Adversarial



Negotiation

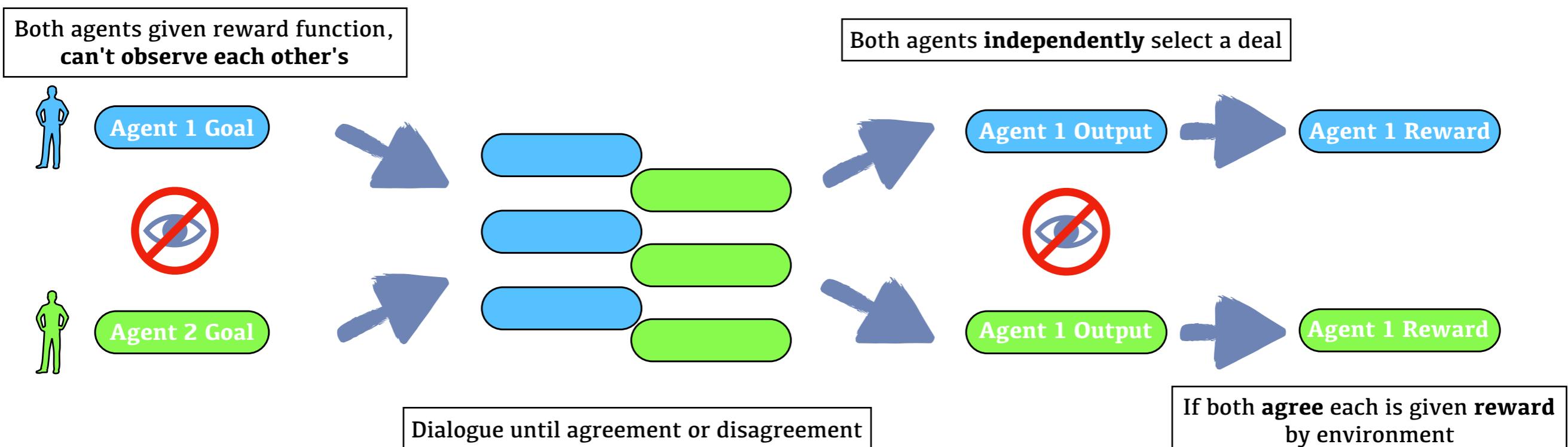


Fully
Cooperative





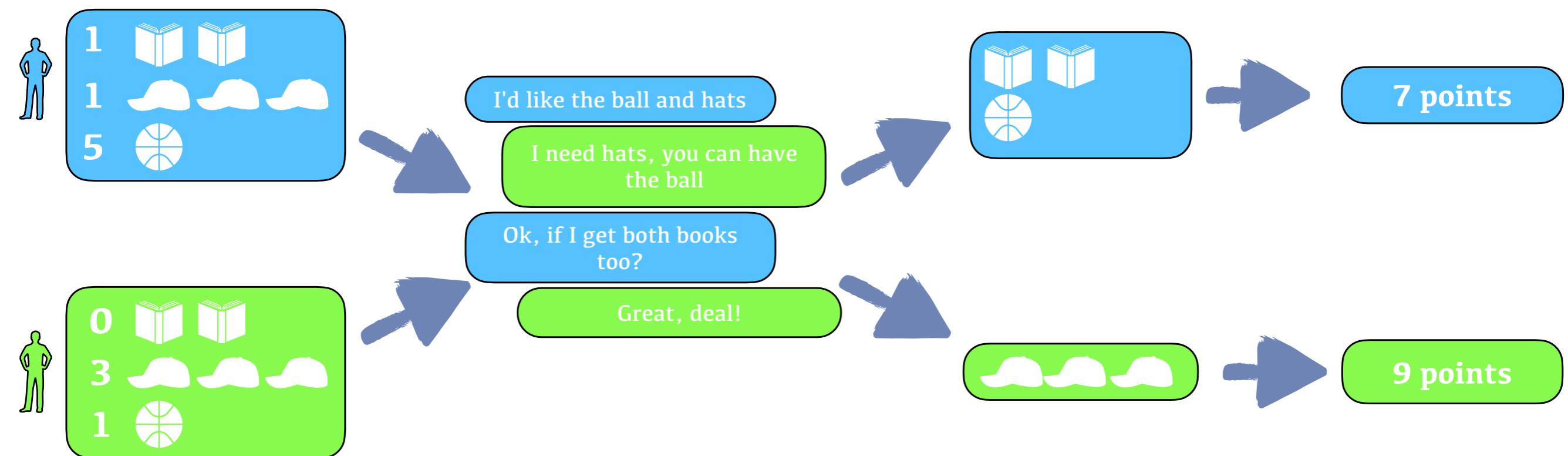
Framework





Dataset

Object division task:





Dataset

Object division task:

You can have 3 books, I will take the rest

I will give you the basketball and a book

You can have 4 books, final offer

That deal is not fair and I will not accept.
Split it down the middle or no deal.



Fine, walk away with nothing

You are doing the same. Hope you enjoy
your rejection.

Alright I'll take a hat and a book

Awesome. Pleasure doing business with
you



Dataset

| Metric | Dataset |
|----------------------------|---------|
| Number of Dialogues | 5808 |
| Average Turns per Dialogue | 6.6 |
| Average Words per Turn | 7.6 |
| Agreed (%) | 80.1 |
| Average Score (max 10) | 6.0 |
| Pareto Optimality (%) | 76.9 |



Baseline Model

How to model an agent that can both perform **reasoning** and **conversing**?

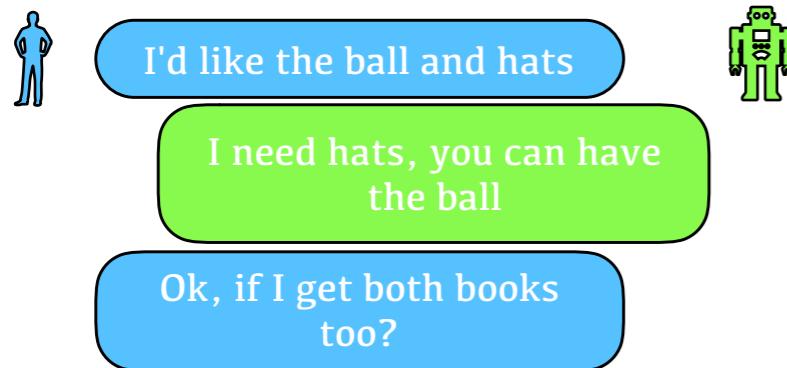
Let's follow an end-to-end approach:

- Single model for **interpretation, generation and reasoning**



Baseline Model

Linearize dialogue into token sequence

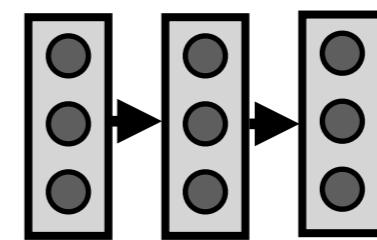




Baseline Model

Train conditional language model to predict tokens

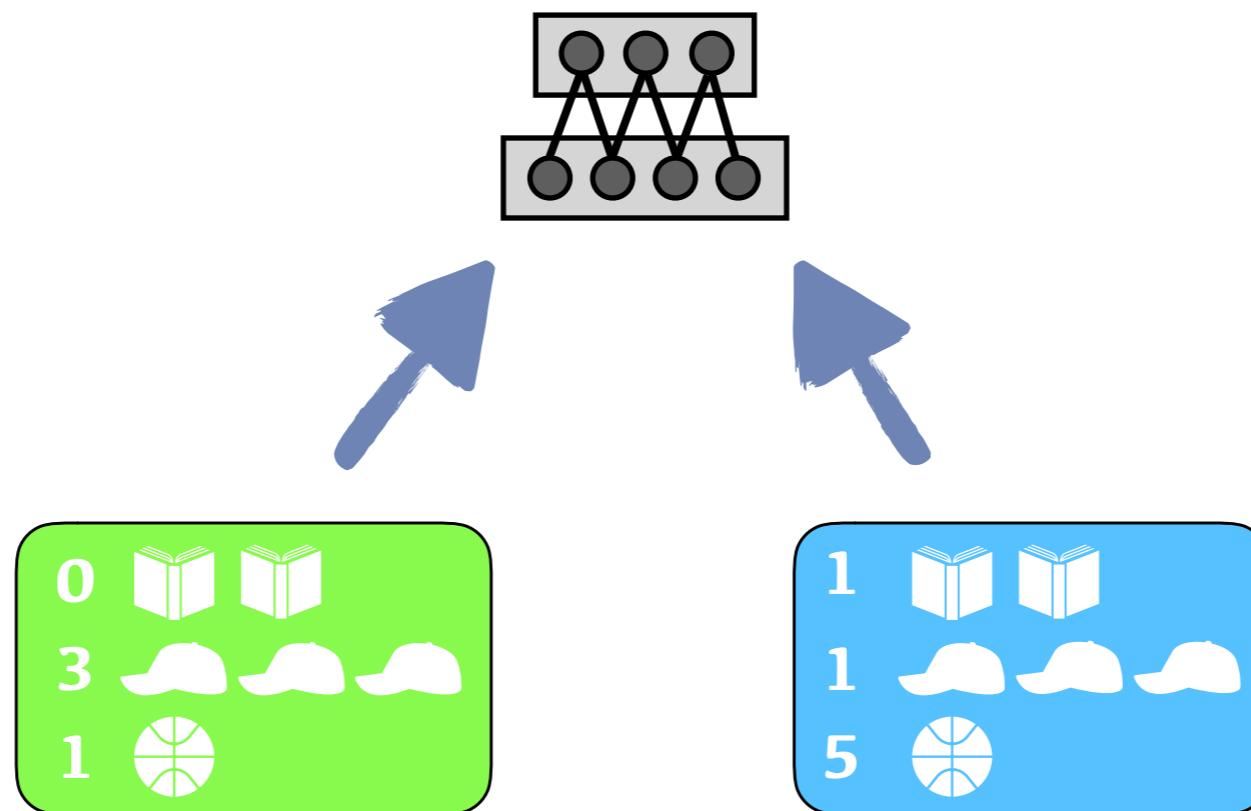
<read> I'd like the ball and
hats <write> I need hats, you
can have the ball <read> Ok, if
I get both books too?





Baseline Model

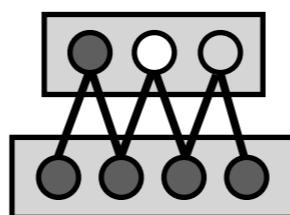
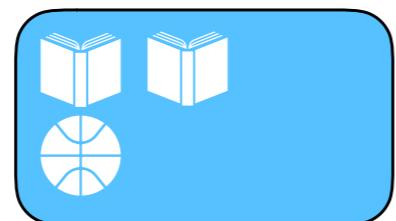
Train Encoder to read in the value function





Baseline Model

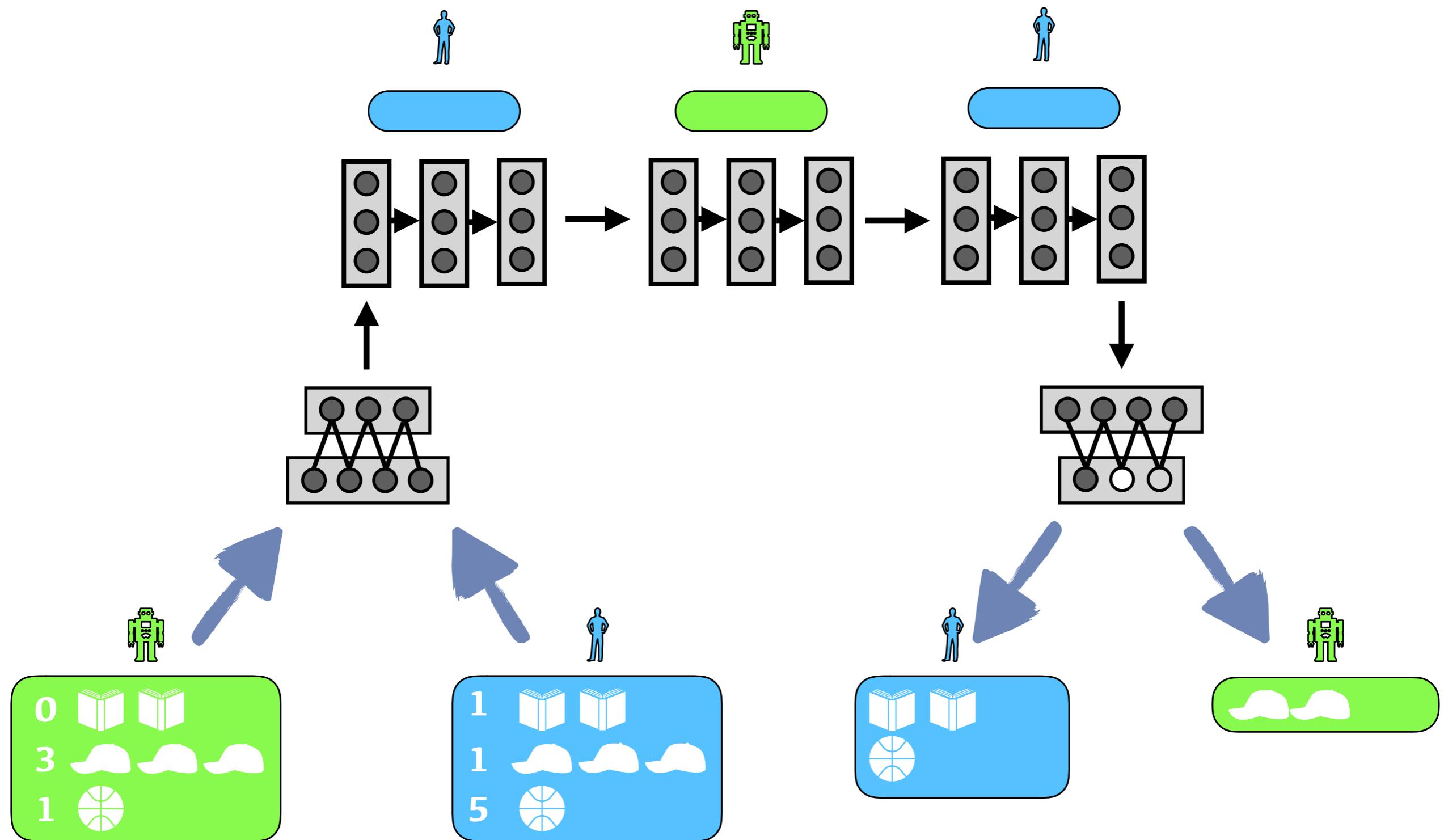
Train **Classifier** to predict the final deal





Baseline Model

Combine everything **together**





Training

Train language model using **supervised** learning to minimize NLL:

- model knows nothing about the task, just imitates human actions
- agrees to easily
- can't go beyond human strategies

$$J(\theta) = J_{\text{lm}}(\theta) + J_{\text{classifier}}(\theta)$$

$$J_{\text{lm}}(\theta) = - \sum_t \sum_i x_{t,i} \log \hat{x}_{t,i}$$

$$J_{\text{classifier}}(\theta) = - \sum_t \sum_i s_{t,i} \log \hat{s}_{t,i}$$



Training

Use **reinforcement** learning to teach model to negotiate:

- generate dialogues using **self-play**
- backpropagate reward using **REINFORCE**

RL **adversely affects** model's language coherence:

- **interleave** supervised and reinforcement training



Training

Monte-Carlo estimate of objective function:

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_t r(s_t, y_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(s_{i,t}, y_{i,t})$$

Policy Gradient:

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta} [r(\tau)] = \int \pi_\theta(\tau) r(\tau) d\tau \\ \nabla_\theta J(\theta) &= \int \nabla_\theta \pi_\theta(\tau) r(\tau) d\tau = \int \pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau) r(\tau) d\tau = \\ &= \mathbb{E}_{\tau \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(\tau) r(\tau)] \end{aligned}$$

Update:

$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$



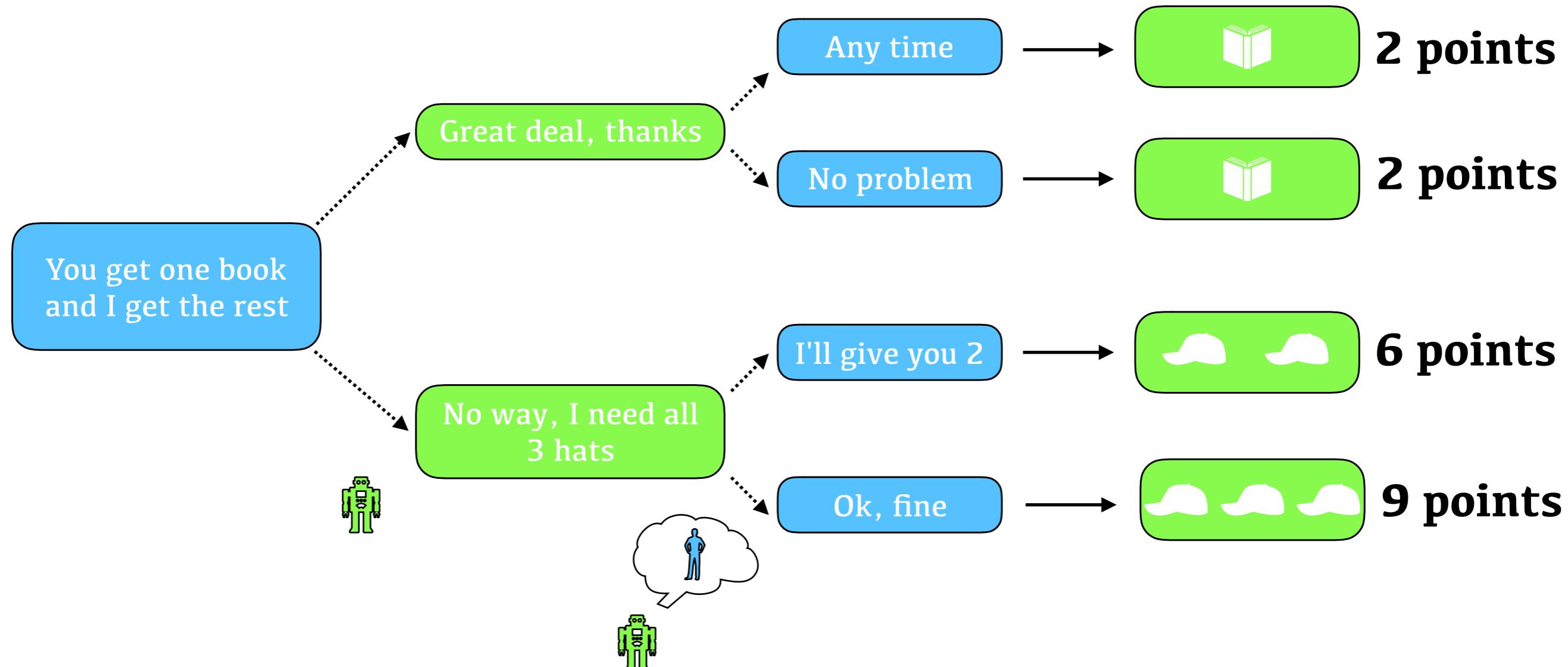
Rollouts

Finally, on trained model we can use rollouts:

- estimates expected score of a dialogue
- widely used in games:
 - Monte-Carlo-Tree-Search in AlphaGo



Rollouts



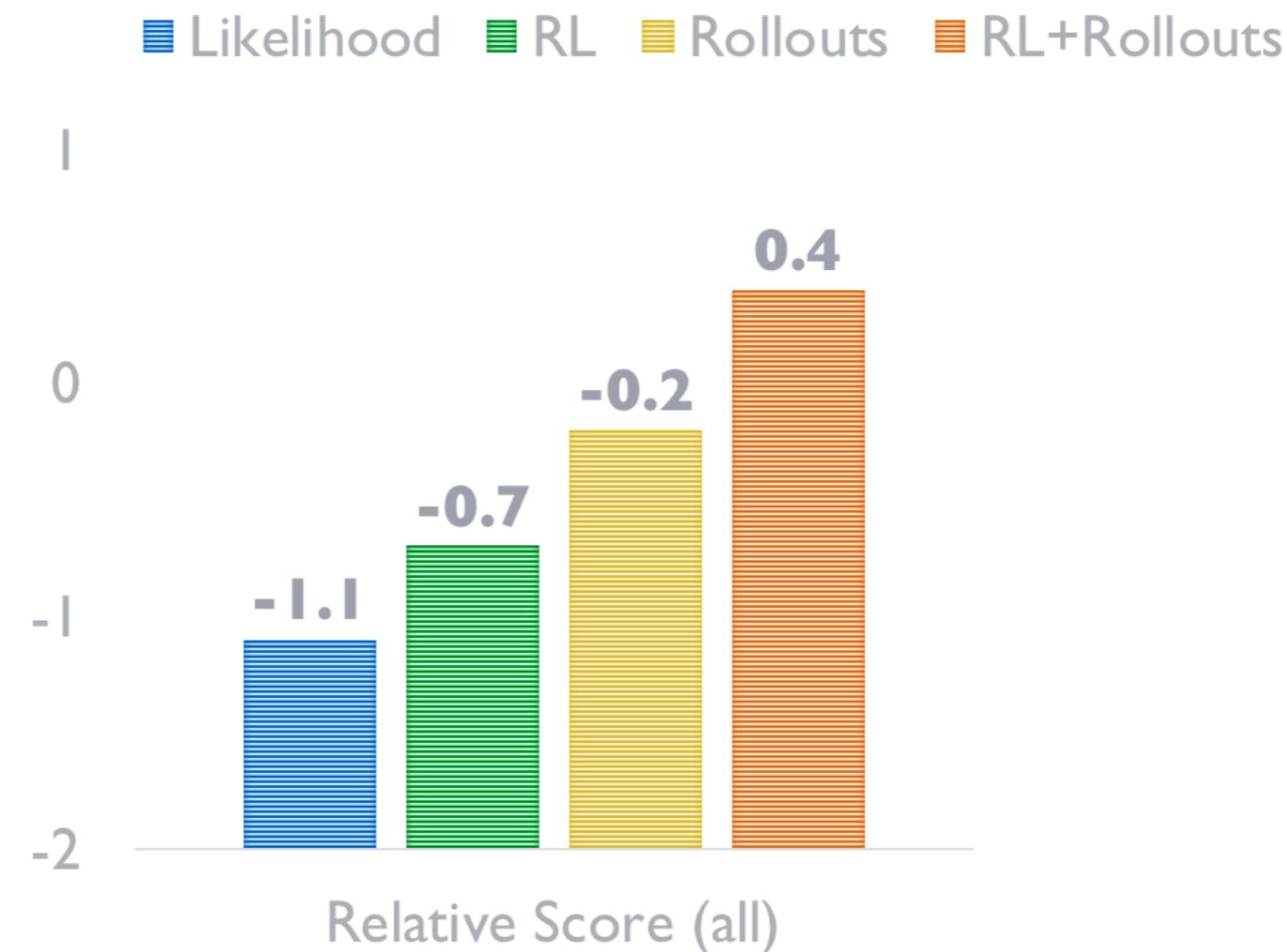


Rollouts





Evaluation vs Human

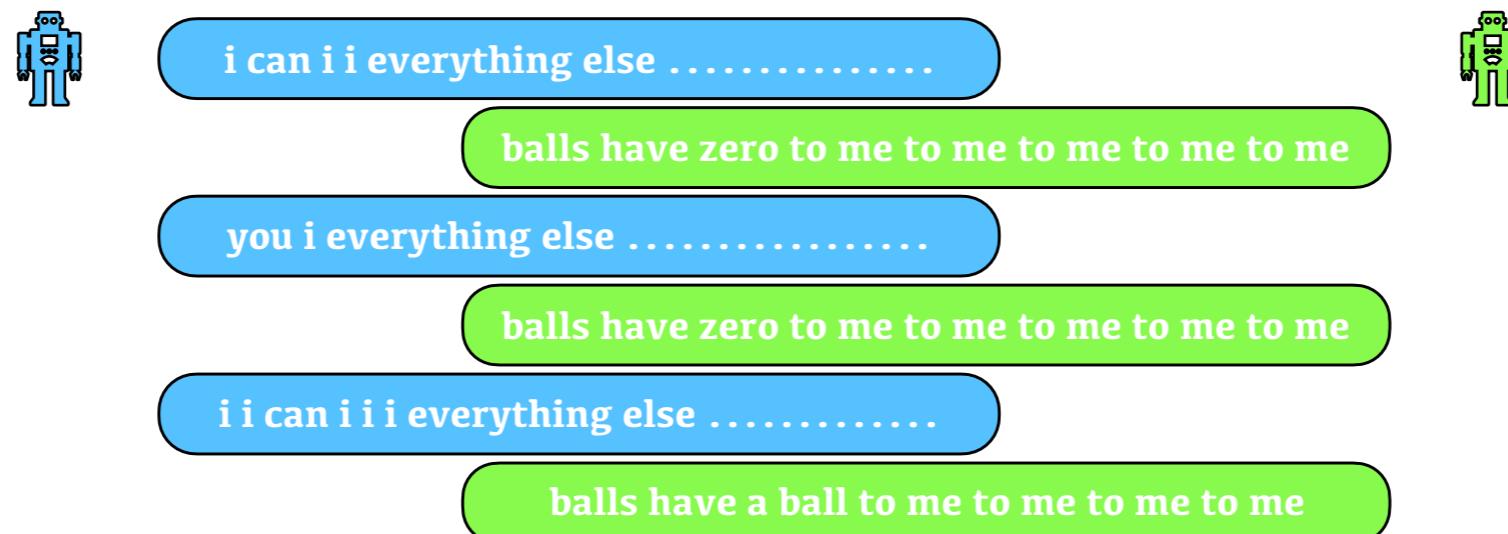




Issues

Entangled representations responsible for both **end-goal performance** and **language coherence**:

- Fine-tuning end-goal performance with RL adversely affects generated language





Issues





Hierarchical Approach

Solution:

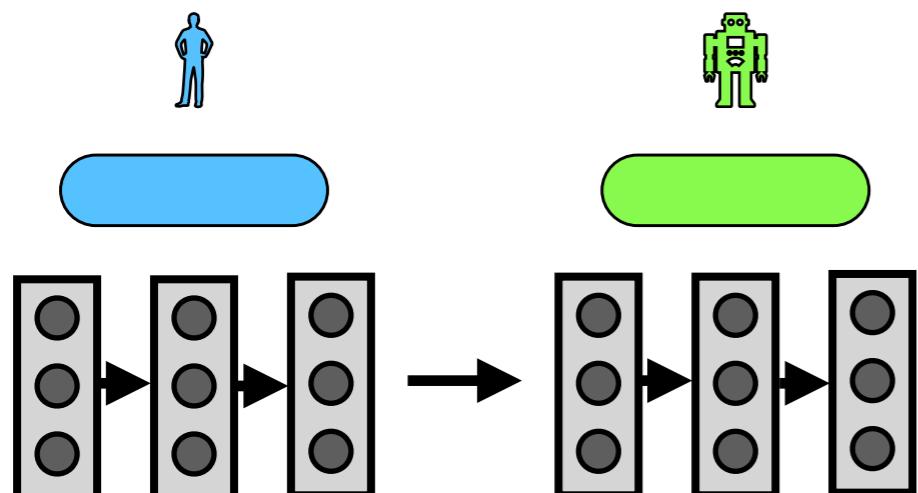
- Factorize generation problem into "**what to say**" and "**how to say**"

$$\pi_{\theta}(y_t | s_t)$$

VS

"how to say" "what to say"

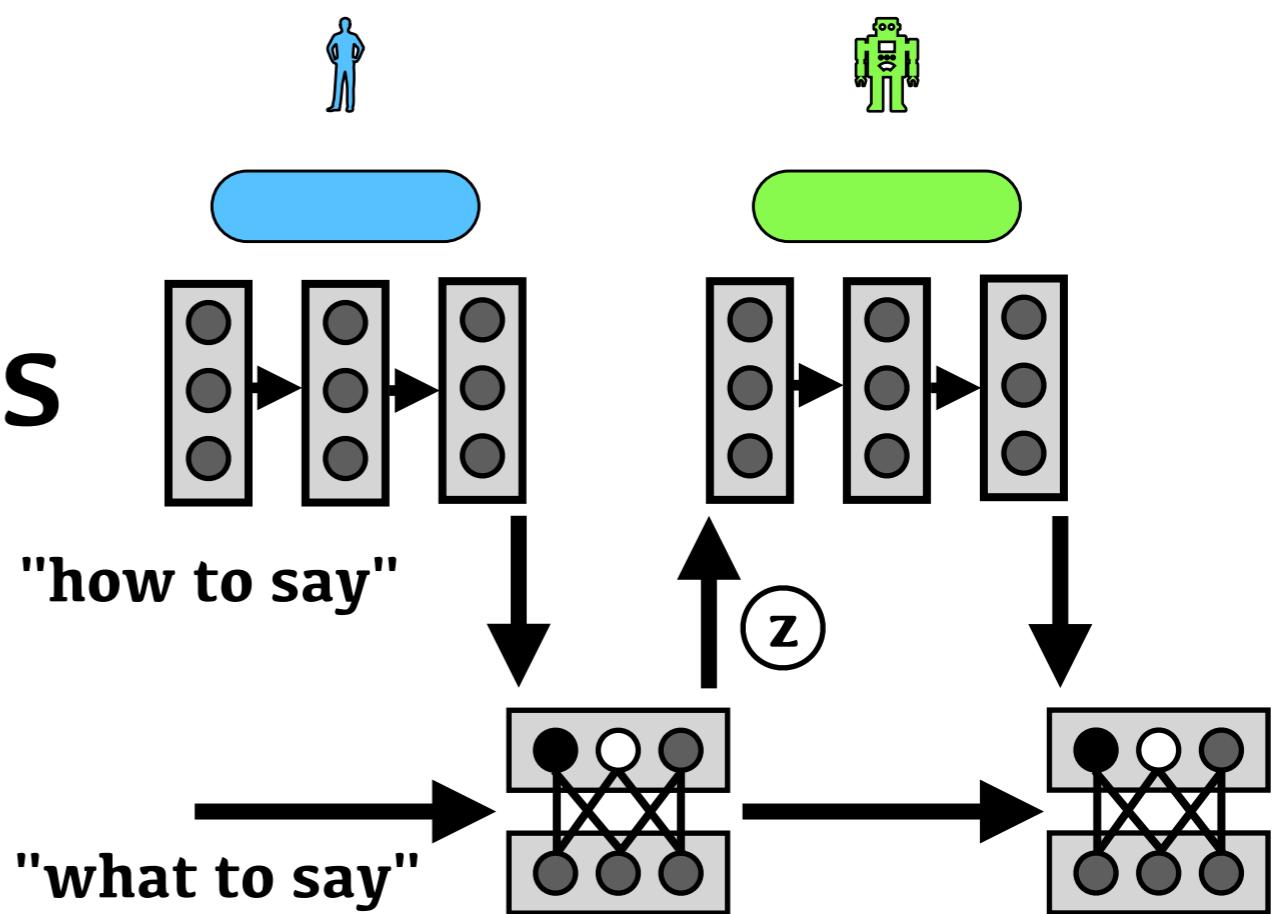
$$\pi_{\theta}^{\text{lm}}(y_t | z_t) \pi_{\theta}^z(z_t | s_t)$$

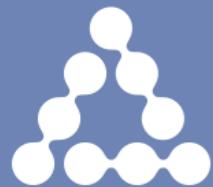


VS

"how to say"

"what to say"

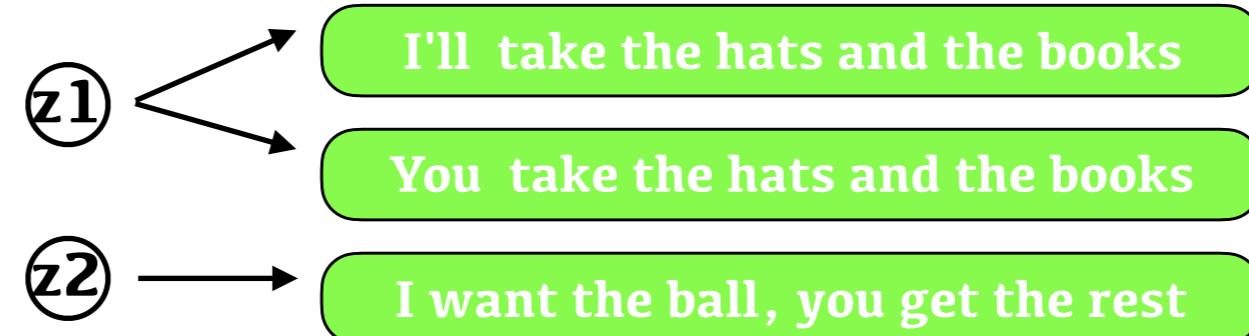




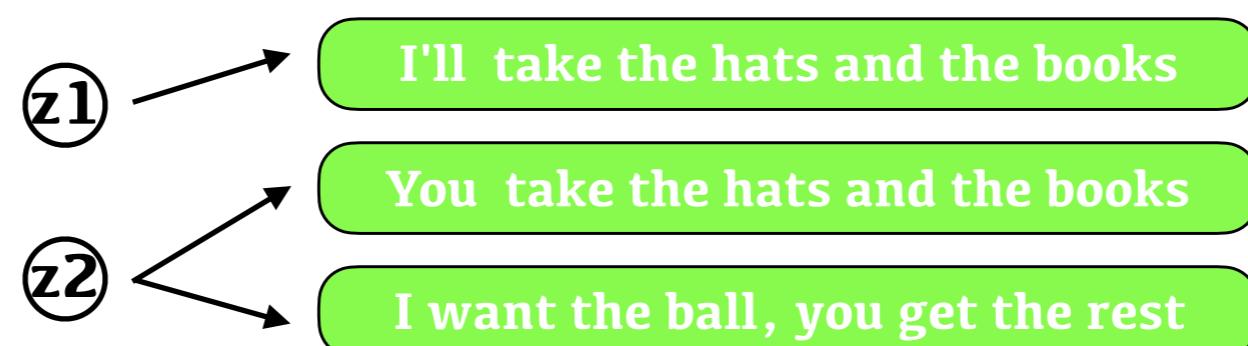
Latent Variable

How to learn z?

- Learning z to maximize likelihood of messages gives similar strings, not similar meaning



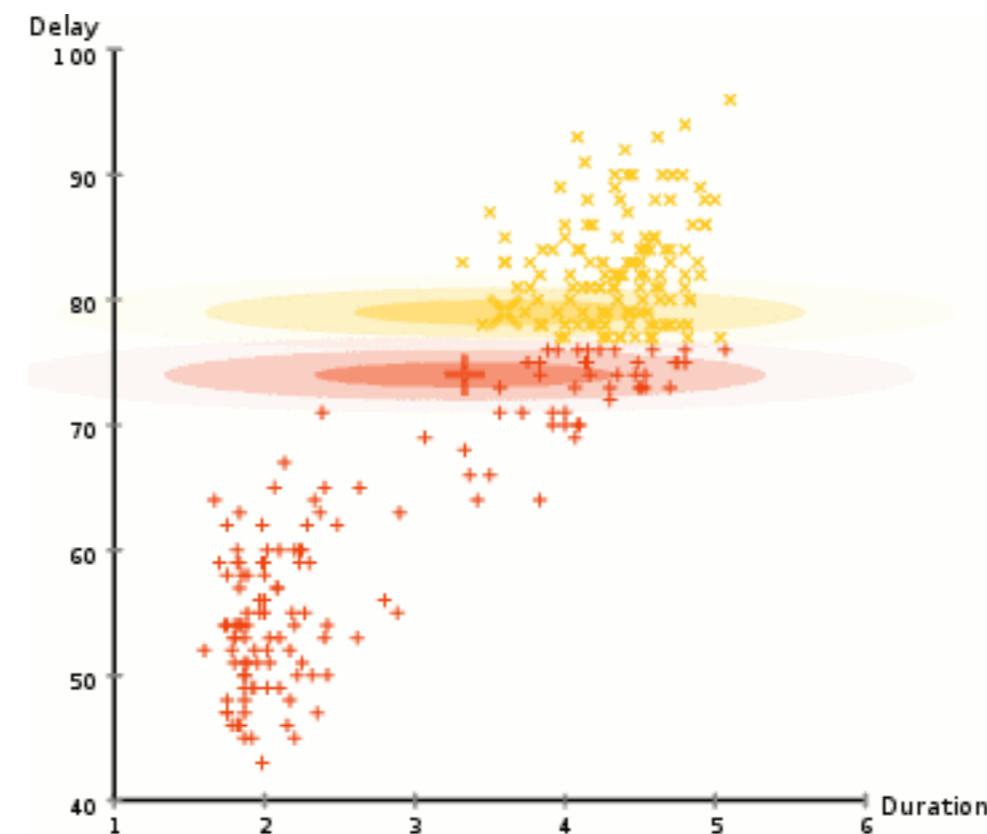
- Learning z to capture message's effect on dialogue means that semantically similar messages are clustered





Latent Variable

- Pretrain z using Expectation Maximization algorithm
- Widely used for models with latent variables:
 - e.g. Mixture of Gaussians

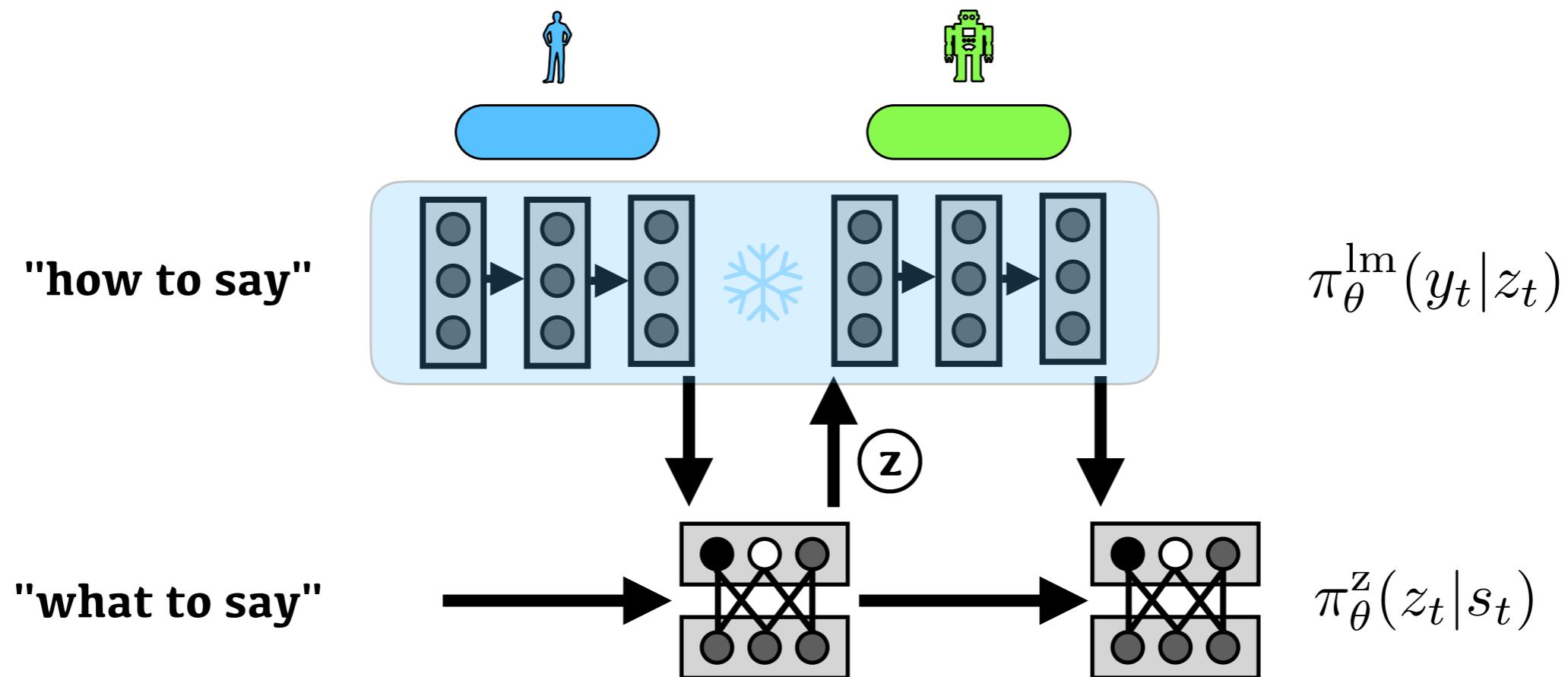




Hierarchical Approach

Benefits:

- Freeze "**how to say**" and fine-tune "**what to say**" with RL to preserve language cohesiveness
 - i.e. only optimize $\pi_{\theta}^z(z_t|s_t)$ with RL
 - $\pi_{\theta}^{\text{lm}}(y_t|z_t)$ stays unchanged
- Rollouts using z are more efficient, only done on $\pi_{\theta}^z(z_t|s_t)$





Learned Clusters

| Cluster | BASELINE CLUSTERS | FULL MODEL |
|---------|---|---|
| 1 | i can give you the books but , i would need the hat and the balls i can do that . i need both balls and one book | i would like the hat and 1 book i can't give up the hat , but i can offer you the book and 2 balls |
| 2 | i need both books and the hat how about you get the hat and 1 ball | i want the hat i need the hat . you can have all the books and the balls |
| 3 | i can not make that deal . i need the hat and one book i can give you the hat and 1 ball | i can give you the hat and 1 ball i would like the books and a ball |
| 4 | i need two books and the hat i need the hat , you can have the rest | i need the books and the hat i can give you the balls but i need the hat and books |
| 5 | i can give you the hat if i can have the rest i want one of each | could i have the books and a ball ? i would like the books and one ball |

Table 4: Sample messages that are probable under different clusters for specified context, in comparison to a previous approach to learning message representations. An agreement needs to be done over a set of **2 books**, **1 hat**, and **2 balls**. The clusters produced by our method are much more semantically coherent than the baseline, and correspond closely to different ways of proposing the same deal.



Experiments

| Model | Validation Perplexity | Test Perplexity |
|-------------------|--------------------------|--------------------|
| RNN | 5.62 | 5.47 |
| HIERARCHICAL | 5.37 | 5.21 |
| BASELINE CLUSTERS | 5.61 | 5.46 |
| FULL MODEL | 5.37 | 5.24 |

Table 1: Perplexity results, showing the likelihood of human dialogues using different models. Our model with discrete message representations is able to achieve state-of-the-art performance, showing that the representations effectively capture relevant aspects of messages for predicting the future dialogue.



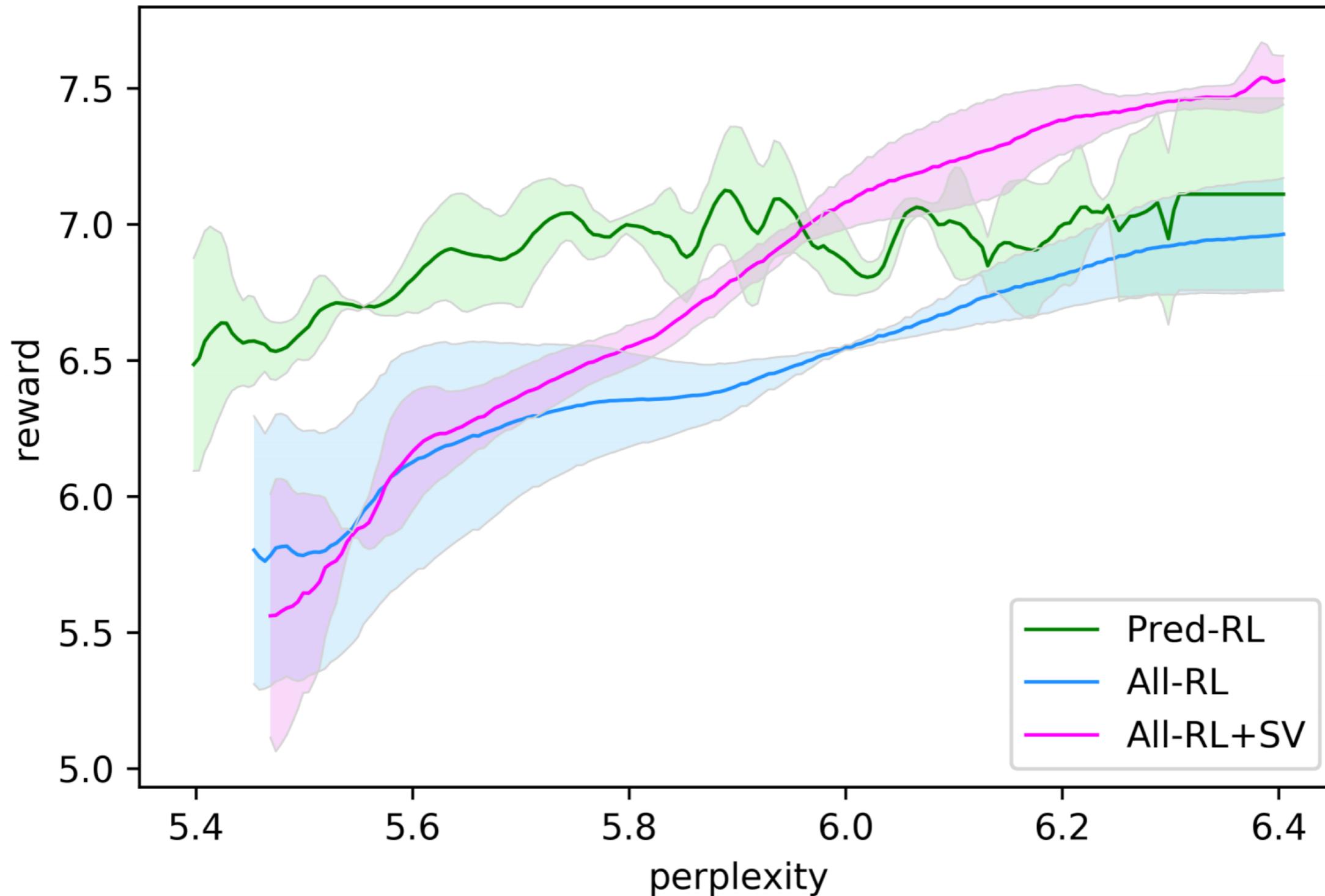
Experiments

| Model | Score vs. | Score vs. |
|-------------------|-----------|--------------|
| | RNN | HIERARCHICAL |
| RNN | 5.33 | 5.17 |
| HIERARCHICAL | 5.37 | 5.08 |
| BASELINE CLUSTERS | 4.68 | 4.66 |
| FULL MODEL | 6.75 | 6.57 |

Table 2: Comparison of different models based on their end-task reward. Our clusters substantially improve reward, indicating that they make it easier for supervised learning to model strategic decision making.



Experiments





Research Papers

Deal or no Deal? End-to-End Learning for Negotiation Dialogues

- URL: <https://arxiv.org/abs/1706.05125>

arXiv.org > cs > arXiv:1706.05125

Computer Science > Artificial Intelligence

Deal or No Deal? End-to-End Learning for Negotiation Dialogues

Mike Lewis, Denis Yarats, Yann N. Dauphin, Devi Parikh, Dhruv Batra

Hierarchical Text Generation and Planning for Strategic Dialogue

- <https://arxiv.org/abs/1712.05846>

arXiv.org > cs > arXiv:1712.05846

Computer Science > Computation and Language

Hierarchical Text Generation and Planning for Strategic Dialogue

Denis Yarats, Mike Lewis



Open Source

End-to-End Negotiator

- URL: <https://github.com/facebookresearch/end-to-end-negotiator>

The screenshot shows the GitHub repository page for 'facebookresearch / end-to-end-negotiator'. The top navigation bar includes links for 'This repository', 'Search', 'Pull requests', 'Issues', 'Marketplace', and 'Explore'. On the right side of the header are icons for notifications, a plus sign, and a user profile. Below the header, the repository name 'facebookresearch / end-to-end-negotiator' is displayed, along with statistics: 80 watchers, 1,029 stars, and 203 forks. A navigation bar below the repository name offers links to 'Code', 'Issues 2', 'Pull requests 1', 'Projects 0', and 'Insights'. The main content area features a brief description: 'Deal or No Deal? End-to-End Learning for Negotiation Dialogues'.

