**Georgia Tech**

CS 7643
Project ID: 11

# Inner Dialog:
# Pragmatic Visual Dialog Agents that
# Rollout a Mental Model of their Interlocutors
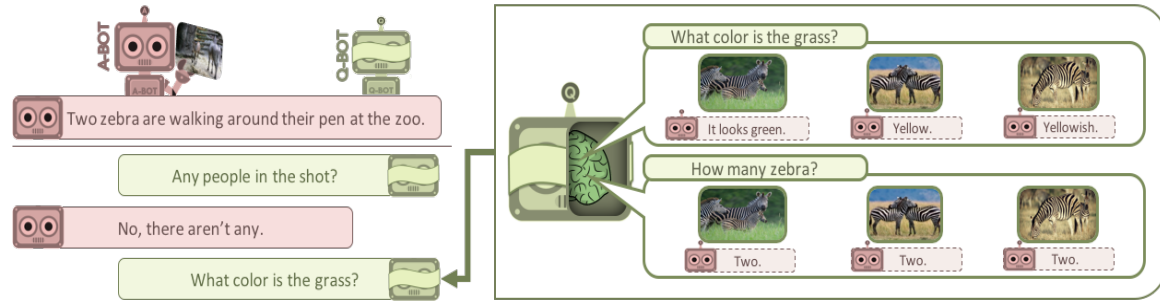
Viraj Prabhu

## 1 OVERVIEW

· **Motivation:** "To plan ahead we must simulate the world"

· **GuessWhich:** Cooperative image guessing game (Das & Kottur et al. 2017)

· **Goal**: Select "**pragmatic**" questions at **inference**, via **dialog rollouts** on a **mental model** of teammate (ABOT)



## 2 APPROACH

· **Modeling AMENTAL:**

  · ACOPY: ABOT replica (performance upper bound)

  · AMIMIC: Same architecture, trained on ABOT samples

· **Dialog Rollouts:**

  · Reward Estimation via finite-sample approximation

    · Minimize **Bayes Risk** under QBOT'S beliefs

    · Sample candidate questions, images, answers

    · **Marginalize over beliefs** using AMENTAL

### Optimization

Pick question with max expected reward
$$q_t^* = \arg\max_{q_t} \underbrace{\mathbb{E}_{\tilde{I}, \tilde{a}_t} \left[ r_t \left( \hat{I}_{|q_t, \tilde{a}_t}^t, \tilde{I} \right) \right]}_{\text{Estimated reward}}$$

$$\underbrace{\mathbb{E}_{\tilde{I}, \tilde{a}_t} \left[ r_t \left( \hat{I}_{|q_t, \tilde{a}_t}^t, \tilde{I} \right) \right] \approx \frac{1}{MN} \sum_{m=1}^{M} \sum_{n=1}^{N} r_t \left( \hat{I}_{|q_t, a_t^{m,n}}^t, I_m \right)}_{\text{Marginalize over beliefs}}$$

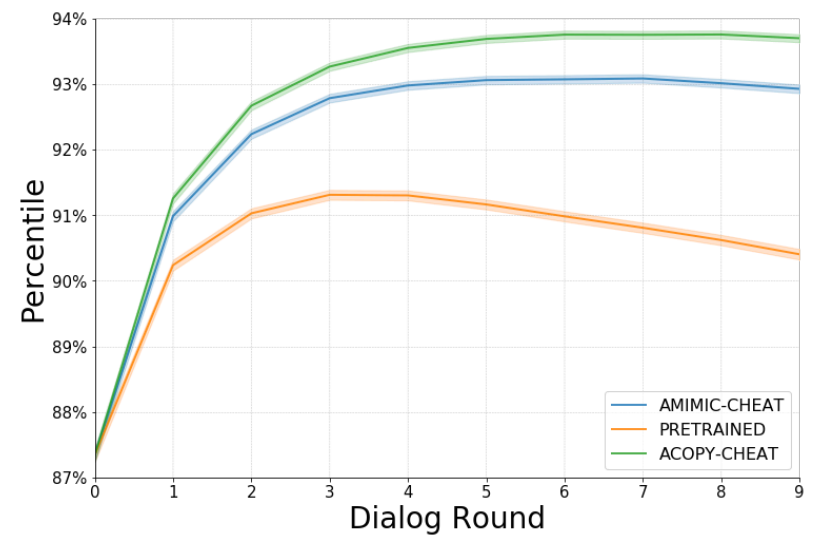**Algorithm 1** Selecting Pragmatic Questions Via Dialog Rollouts
```
1:  function PRAGMATIC-QUESTION-SELECTION(q_t)
2:      q_t^1, ..., q_t^K ~ π_Q(q_t | q_0, a_0, ..., q_{t-1}, a_{t-1})    ▷ Decode multiple likely questions
3:      return argmax RolloutEstimate(q_t^i)    ▷ Select q_t with greatest expected reward
4:  end function
5:
6:  function ROLLOUTESTIMATE(q_t)
7:      r̃ ← 0
8:      I_1, ..., I_M ~ G_{t-1}(Î^{t-1}|q_0, a_0, ..., q_{t-1}, a_{t-1})    ▷ Sample likely source images
9:      for m ∈ {1, ..., M} do
10:         for n ∈ {1, ..., N} do
11:             a_t^{m,n} ~ π_M(a_t | I_m, q_0, a_0, ..., q_t)    ▷ Sample answer given q_t and I_m
12:             Î_{|q_t, a_t^{m,n}}^t ← argmax G_t(Î^t|q_0, a_0, ..., q_t, a_t^{m,n})    ▷ Update Q-BOT's prediction
13:             r̃ ← r̃ + r_t(Î_{|q_t, a_t^{m,n}}^t, I_m)    ▷ Aggregate the reward
14:         end for
15:     end for
16:     return r̃/MN    ▷ Return approximate expected reward
17: end function
```

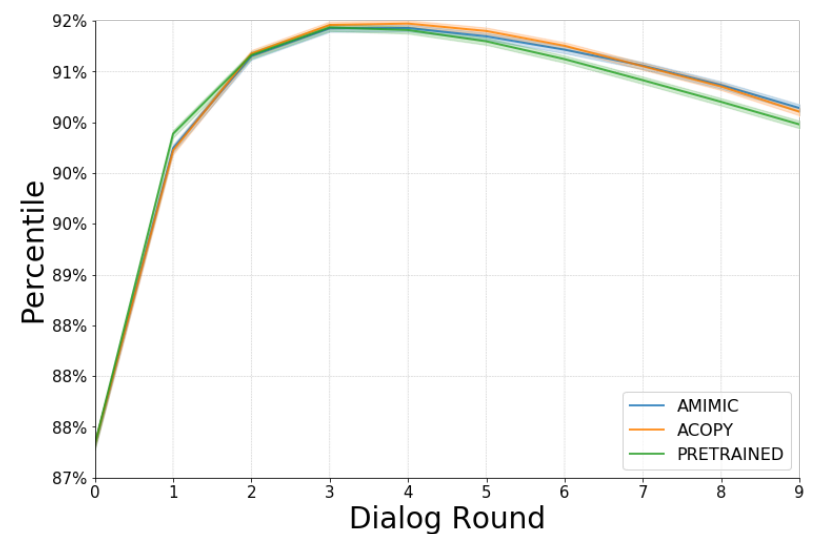## 3 PRELIMINARY RESULTS

**Metric:** Rank percentile over dialog rounds

**Baseline:** PRETRAINED (no mental modeling)

**Cheat Setting:** "Cheat" with GT target image



Mental modeling can help ..

**Real World:** Expectation under QBOT'S beliefs



.. but real-world gains are not observed yet

### Challenges

· Estimating rewards is a bottleneck
· Scaling up approximation is expensive

### Takeaways

· Pragmatic inference can in theory provide an alternative to fine-tuning with RL

· But in presence of information assymetry, accurately estimating reward in hard

## 4 REFERENCES

Das, A., Kottur, S., Moura, J.M., Lee, S. and Batra, D., (2017). Learning Cooperative Visual Dialog Agents with Deep Reinforcement Learning. ICCV 2017

Das, A., Kottur, S., Gupta, K., Singh, A., Yadav, D., Moura, J., Parikh, D., and Batra, D., Visual dialog. CVPR 2017.