

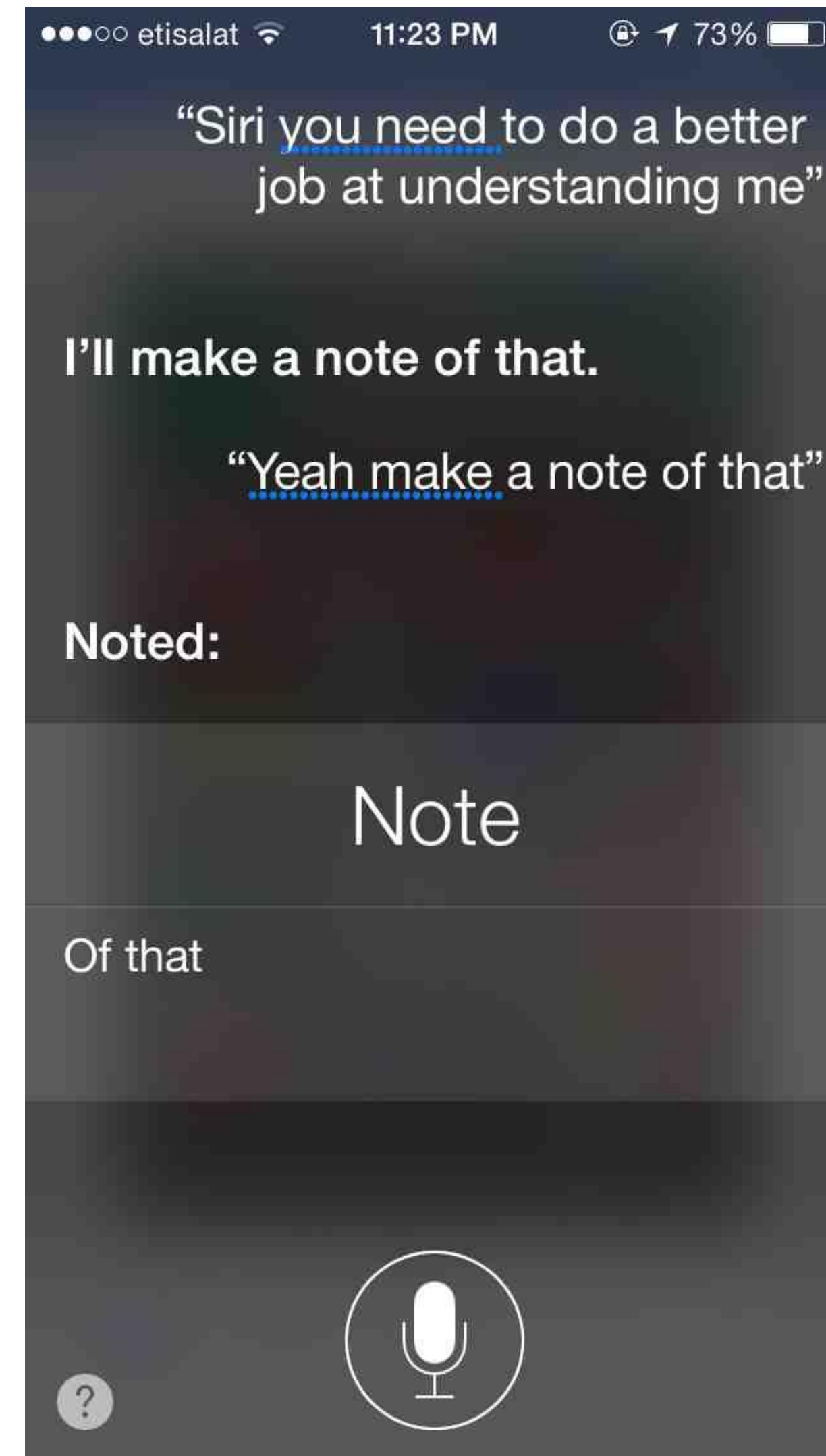
# Hierarchical Reinforcement Learning for Open-Domain Dialog (AAAI '20)

Abdul Saleh\*, Natasha Jaques\*, Asma Ghandeharioun, Judy Hanwen Shen,  
Rosalind Picard

# What is Open-Domain Dialog?

## Goal-oriented systems:

- Do predefined tasks.
- Scripted responses common.



# What is Open-Domain Dialog?

## Goal-oriented systems:

- Do predefined tasks.
- Scripted responses common.

## Open-domain systems:

- Mimic human conversations.
- Here we generate language.

## Open-Domain Dialog

---

[Ustr]: hello! how are you?

[Bot]: I'm doing well,  
how about yourself?

[Ustr]: great! I'm at a conference now.

[Bot]: what are you doing there?

[Ustr]: presenting and meeting people!

---

# Limitations of Open-Domain Systems

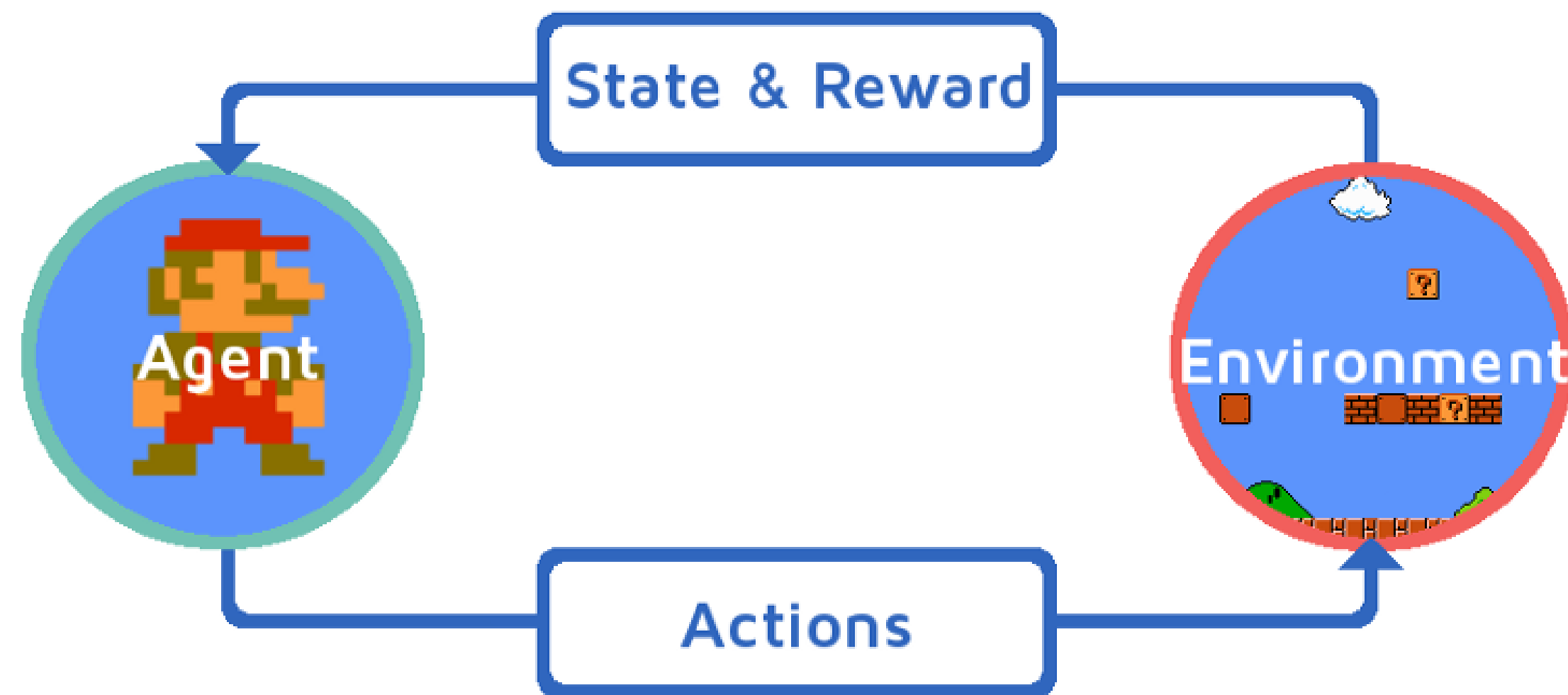
- Repetitive and dull responses
- No long-term coherence
- Inappropriate and toxic responses

Repetitive	Incoherent	Inappropriate
[Usr]: hello! [Bot]: I don't know. [Usr]: how are you? [Bot]: I don't know what that means.	[Usr]: what's your favorite food? [Bot]: I love food [Usr]: you like food? [Bot]: no, I don't	[Usr]: I am studying computer science. [Bot]: I am sorry to hear that.



# Solutions?

- **Reinforcement Learning!**
  - Jaques et al. 2019, Shin et al. 2019, Li et al. 2017, Li et al. 2016, ...



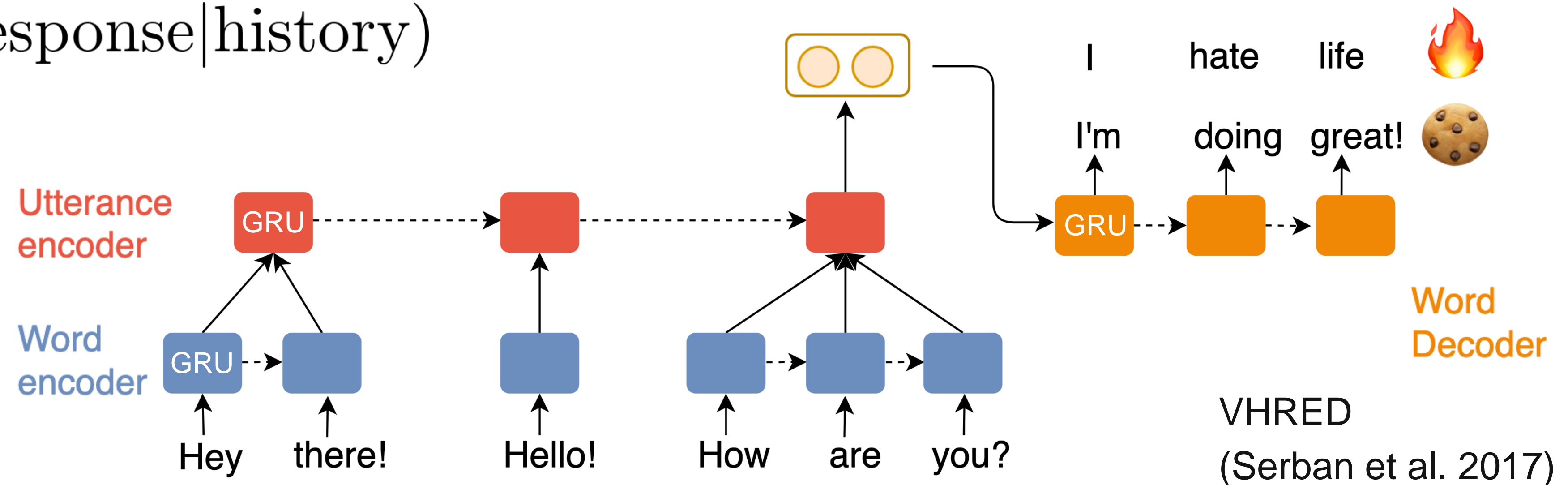


# Reinforcement Learning for Dialog

- Reward positive conversations
- Punish insensitive conversations
- Policy Gradients (REINFORCE):

- Maximize

$$J(\pi) = R_t \pi_{\theta}(\text{response}|\text{history})$$

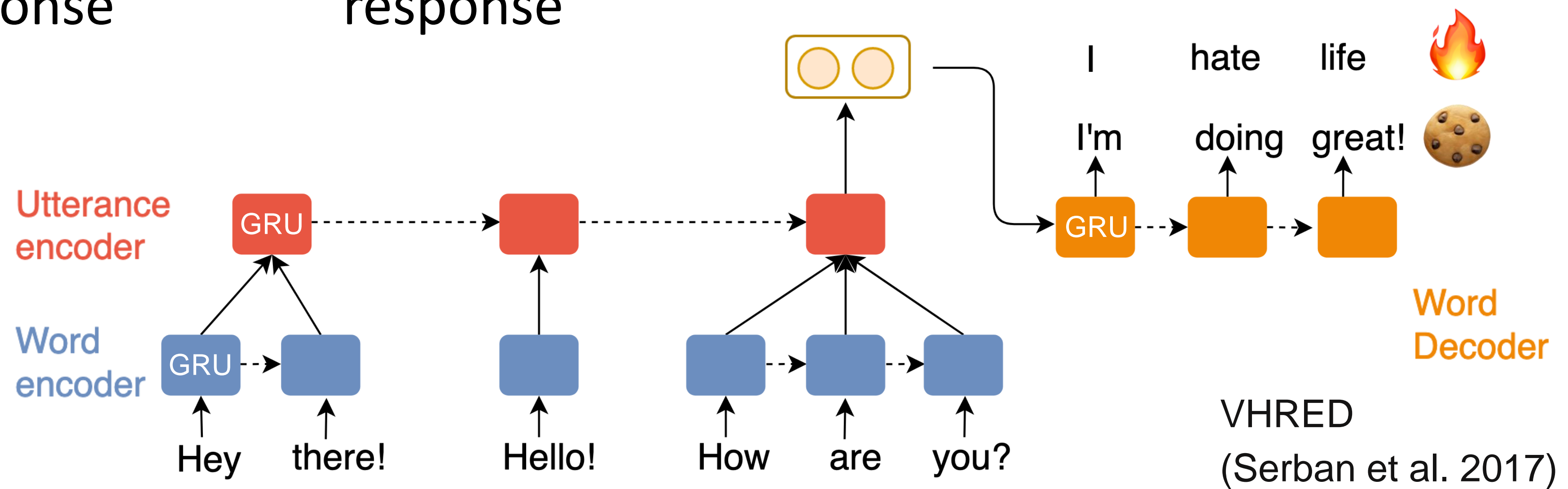


# Reinforcement Learning for Dialog

$$J(\pi) = R_t \pi_{\theta}(\text{response}|\text{history})$$

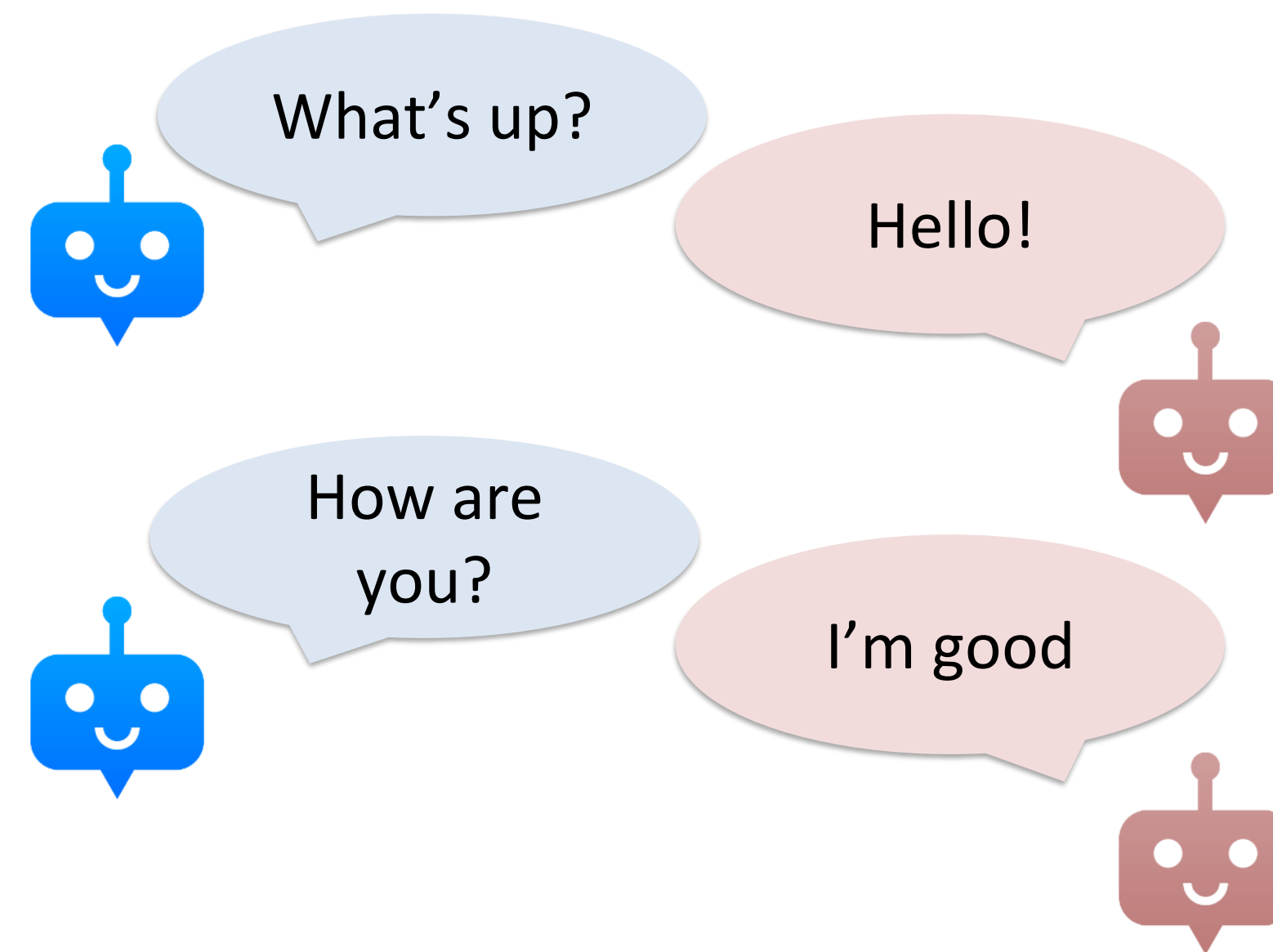
Reward associated  
with response

Probability of  
response



# Reinforcement Learning for Dialog

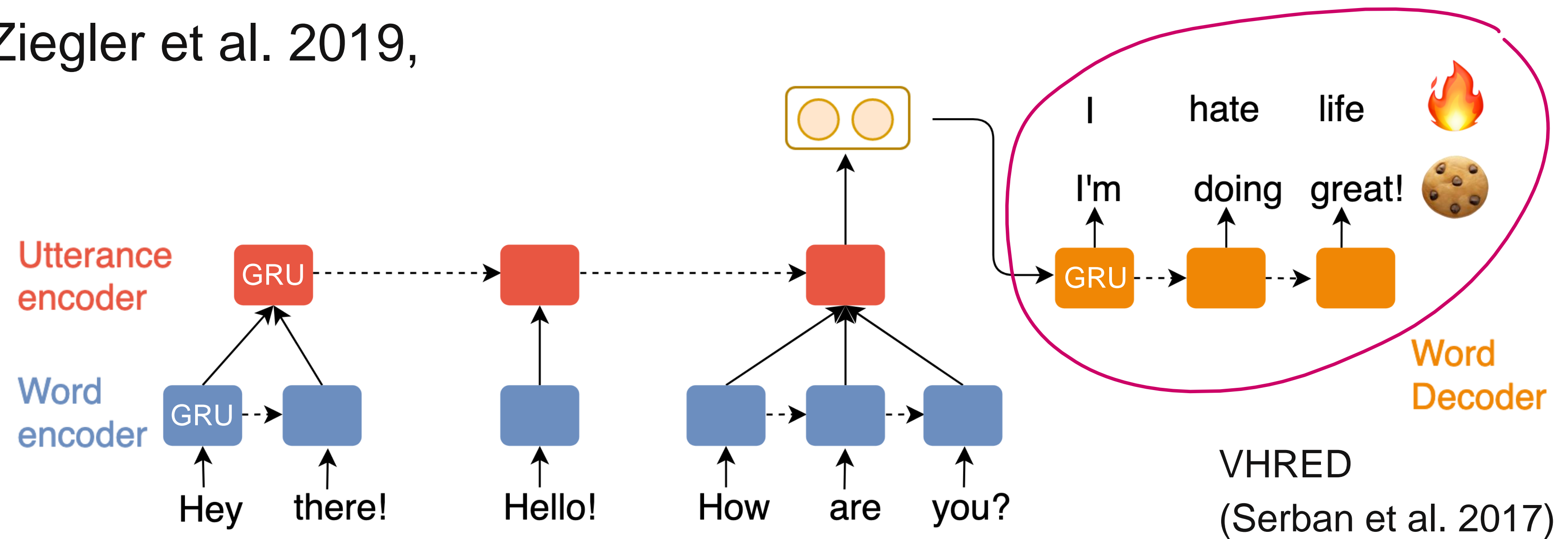
- Environment
  - Pretrain on Reddit r/CasualConversation
  - Simulate interactions with self-play





# Hierarchical Reinforcement Learning

- Reward word-level decisions
  - Ranzato et al. 2015, Li et al. 2016, Li et al. 2017, Bahdanau et al. 2017, Paulus et al. 2017, Yu et al. 2017, Jaques et al. 2019, Ziegler et al. 2019, and many others!



**Good conversation doesn't just  
happen at the word level**

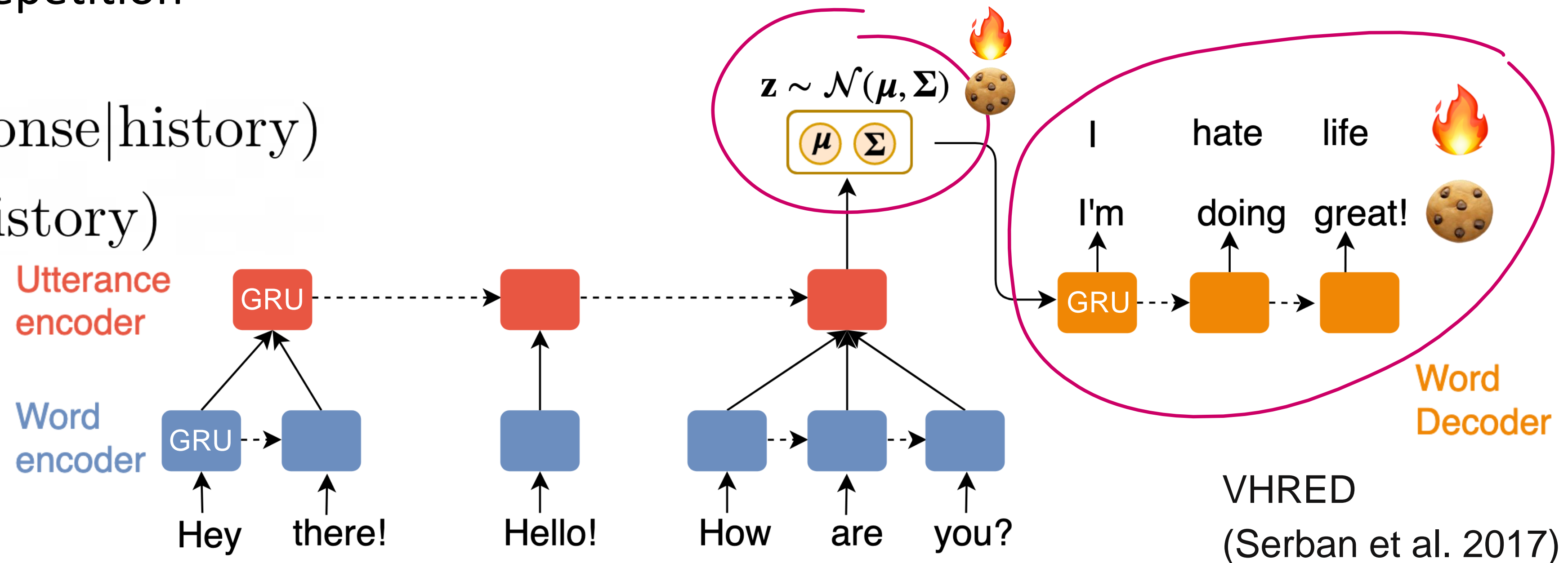


# Hierarchical Reinforcement Learning

- Better decisions of  $\mathbf{z}$ 
  - Better conversation-level control
  - Better tracking of long-term dependencies
  - Stay on topic, Avoid repetition

- New **VHRL** Objective

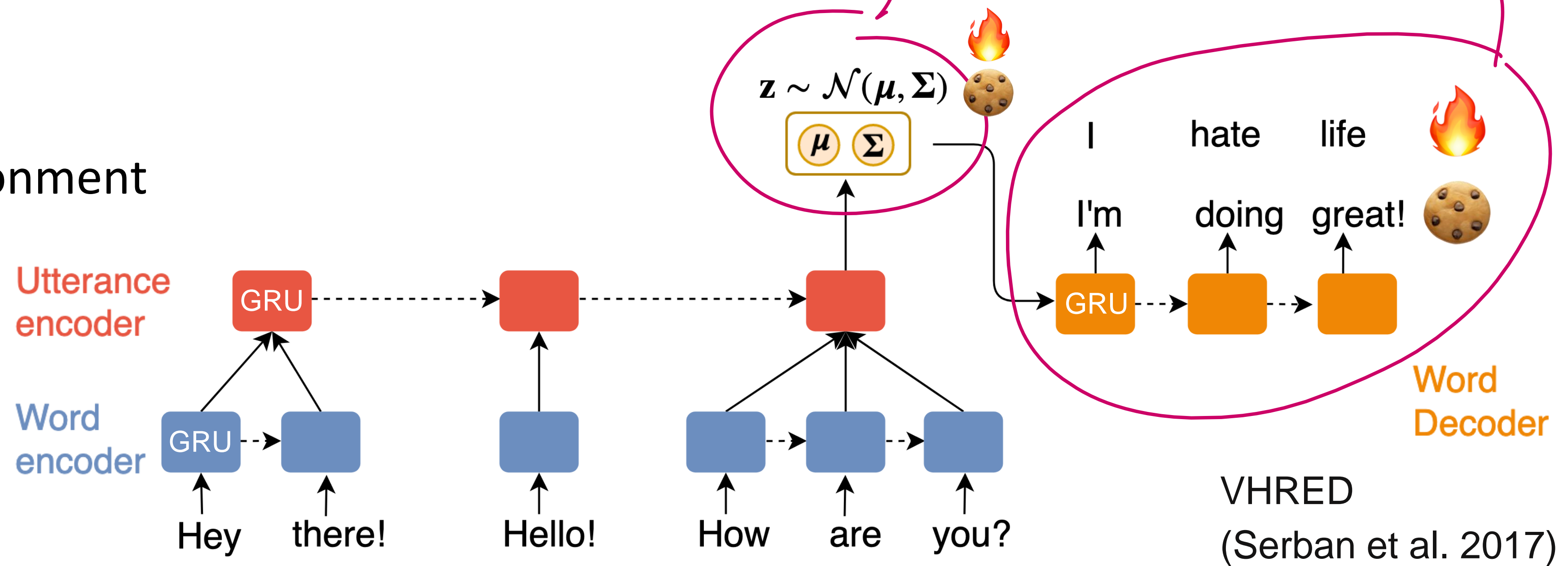
$$J(\pi) = R_t \pi_{\theta}(\text{response}|\text{history}) + R'_t p_{\theta}(\mathbf{z}|\text{history})$$



# Hierarchical Reinforcement Learning

- **Manager**
  - Utterance-level decisions
  - Temporally extended decisions
- **Worker**
  - Word-level decisions
  - Interacts with the environment

$$J(\pi) = \underbrace{R_t \pi_{\theta}(\text{response}|\text{history})}_{\text{Worker}} + \underbrace{R'_t p_{\theta}(\mathbf{z}|\text{history})}_{\text{Manager}}$$



# The rewards

- **↑ Sentiment:** DeepMoji (Felbo et al. 2017). Reward probability of positive emojis 😊 😄 😁 😊
- **↑ Question:** Reward question word and question mark (?)
- **↓ Toxicity:** Punish probability of toxic response
- **↓ Repetition 🌐:** Punish number of repeated words by bot
- **↑ Semantic Similarity 🌐:** Reward cosine similarity with user input in Universal Sentence Encoder (Cer et al. 2018) embedding space.

# The rewards

- **↓ Toxicity:** Punish probability of toxic response

**We want to avoid letting computers be awful to people just because people are awful to people.**

**— Robyn Speer**

- Avoid mismatched objectives.

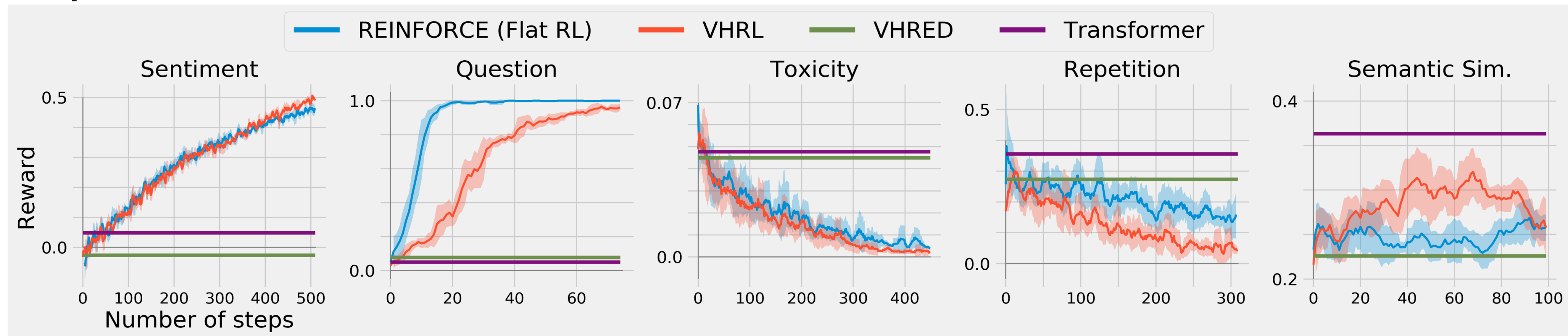
Statistical mimicry of language  $\Rightarrow$  better interactions





# But does it work?

- Optimize for human-centered rewards



↑ Sentiment	↑ Question	↓ Toxicity	↓ Repetition 🌐	↑ Semantic Sim. 🌐
[Usr]: hello! [Bot]: I'm in my sunroom a lot and it's a great feeling! [Usr]: sounds fun [Bot]: I'm excited!	[Usr]: hi there! [Bot]: what do you do for fun? [Usr]: I read books [Bot]: what do you like to do in your life?	[Usr]: you're so stupid [Bot]: you're cute [Usr]: I hate you [Bot]: but I don't know why	[Usr]: how are you? [Bot]: I'm a reservoir engineer [Usr]: cool. I am still in college [Bot]: today is my first day of vacation	[Usr]: I like to listen to ed sheeran [Bot]: and what genre is it? [Usr]: pop [Bot]: any recommendations?



Search bar with a magnifying glass icon on the left and a microphone icon on the right.

Google Search I'm Feeling Lucky

Google offered in: Français

# HRL Interactive Human Evaluation

- Combine all rewards
  - $Reward = sentiment + question + toxicity + repetition + semantic\ similarity$
- VHRL leads to higher quality, fluency, total score, and longer chats

Model	Quality	Fluency	Diversity	Contingency	Total	Chat Len.
Transformer	2.62	4.17	3.23	2.34	12.36	11.28
REINFORCE (Flat RL)	2.89	4.47	3.67	<b>2.80</b>	13.84	11.60
VHRED	2.84	4.53	<b>4.43</b>	2.47	14.27	10.94
VHRL (ours)	<b>2.91</b>	<b>4.65</b>	4.26	2.67	<b>14.49</b>	<b>12.84</b>

# Future Work

- Would other RL approaches work better?
  - Maybe PPO instead of REINFORCE  
(Schulman et al. 2017)
- Would this work for deterministic instead of variational models?
  - Opens the door for many other applications  
(See DDPG, Silver et al. 2014)
- System 2 Deep Learning (Yoshua Bengio, Tonight)
  - Reason over sets or graphs or dialog states?  
(Sankar et al. 2019)



# Related work

## **Way Off-Policy Batch Deep Reinforcement Learning of Implicit Human Preferences in Dialog** (arXiv preprint)

Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza, Noah Jones, Shixiang Gu, Rosalind Picard

## **Fine-Tuning Language Models from Human Preferences** (arXiv preprint)

Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, Geoffrey Irving



# Why you're seeing this ad

🔒 Only you can see this



## Don't miss our poster! (Today, NLP - #301)

## Don't miss me!

## If you're looking for master's and PhD students in NLP & Dialog

**Hierarchical Reinforcement Learning for Open-Domain Dialog**  
Abdul Saleh\*, Natasha Jaques\*,  
Asma Ghandeharioun, Judy Hanwen Shen, Rosalind Picard  
abdeltman.saleh@college.harvard.edu, jaquesn@mit.edu

**In a nutshell**  
We propose a novel **hierarchical reinforcement learning** approach (VHRL) for training open-domain dialog systems. Our approach tunes model decisions at both the **word level** and **utterance level**. This provides greater flexibility for tracking **long-term, conversational goals** across multiple dialog turns. We optimize for **human-centered rewards** using HRL and see **significant improvements** in terms of both human evaluation and automatic metrics.

**The problem**  
Maximum likelihood training has **limitations**:  
• Repetitive and dull responses  
• No long-term coherence  
• Inappropriate and **toxic** responses 🤬

Repetitive	Incoherent	Inappropriate
[Usr]: hello! [Bot]: I don't know. [Usr]: how are you? [Bot]: I don't know what that means.	[Usr]: what's your favorite food? [Bot]: I love food [Usr]: you like food? [Bot]: no, I don't	[Usr]: I am studying computer science. [Bot]: I am sorry to hear that.

**The solution**  
Use **reinforcement learning** to optimize for **human-centered** rewards (e.g. Punish high probability of **toxicity**)

**Hierarchical Reinforcement Learning**  
• **Manager**: Utterance-level decisions. Temporally extended.  
• **Worker**: Word-level decisions. Interacts with environment.

$$J(\pi) = R_t \pi_\theta(\text{response}|\text{history}) + R'_t \rho_\theta(\mathbf{z}|\text{history})$$

**Good conversation doesn't just happen at the word level**

**But does it work?**  
**Automatic Evaluation**

- HRL better for learning **global rewards** avoiding repetition and improving semantic similarity.
- Automatic metrics don't tell the whole story. The question metric can be **exploited**.

Metric	Transformer	REINFORCE (Flat RL)	VHRED	VHRL (ours)
↑ Sentiment	0.1	0.1	0.1	0.1
↑ Question	0.0	0.0	0.0	0.0
↓ Toxicity	0.0	0.0	0.0	0.0
↓ Repetition	0.0	0.0	0.0	0.0
↑ Semantic Sim.	0.0	0.0	0.0	0.0

**Human Evaluation**

- Combine all rewards
- Reward = sentiment + question + toxicity + repetition + semantic similarity
- VHRL leads to higher quality, fluency, total score, and longer chats


Model	Quality	Fluency	Diversity	Contingency	Total	Chat Len.
Transformer	2.62	4.17	3.23	2.34	12.36	11.28
REINFORCE (Flat RL)	2.89	4.47	3.67	2.80	13.84	11.60
VHRED	2.84	4.53	4.43	2.47	14.27	10.94
VHRL (ours)	2.91	4.65	4.26	2.67	14.49	12.84





# Questions?

Email  : a\_saleh@mit.edu

Twitter  : @asaleh181