

Understanding and predicting user dissatisfaction in a neural generative chatbot

Abigail See, Christopher D. Manning



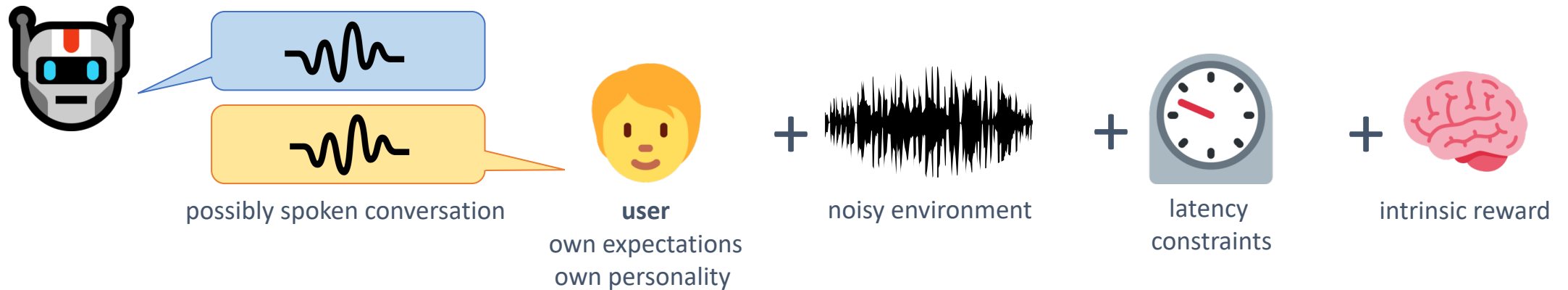
Typical dialogue evaluation setup



Neural generative models perform **increasingly well** in this environment

[Zhang et al 2019, Adiwardana et al 2020, Roller et al 2020]

Real-life dialogue settings



Neural generative models are **less well-explored** in this environment.

Chirpy Cardinal 1.0



For more details:
stanfordnlp.github.io/chirpycardinal

- An open-domain socialbot

- Won 2nd place in third Alexa Prize (2019-2020)

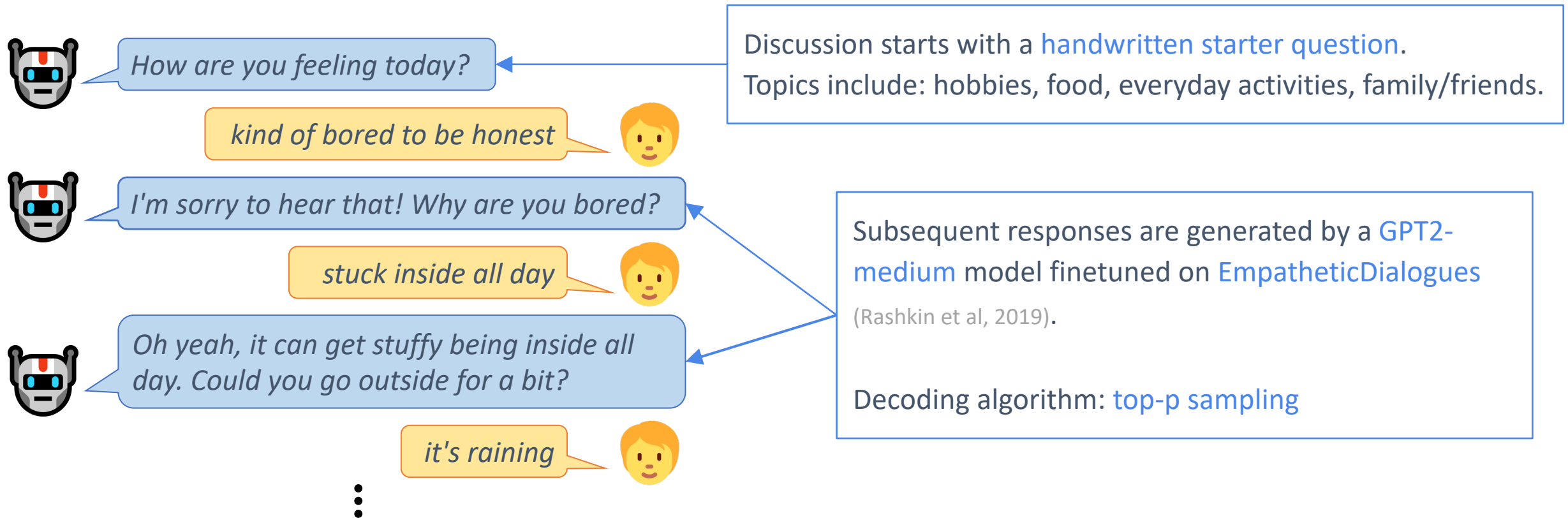


- Alexa customers speak  to socialbots voluntarily 

- Bots must be fast  and manage the noisy home environment 

Neural Chat component

The Neural Chat component uses a **neural generator** to discuss the user's **personal experiences** and **emotions**.



Types of user dissatisfaction



Dissatisfaction type	Example
Clarification	<i>what do you mean</i>
Misheard	<i>that's not what i said</i>
Privacy	<i>none of your business</i>
Repetition	<i>you already said that</i>
Criticism	<i>you're bad at this</i>
Offensive	<i>will you talk dirty</i>
Negative Navigation	<i>change the subject</i>
Stop	<i>goodbye</i>
Any	Any of the above

Neural generative error types

1. Repetition
2. Redundant question
3. Unclear
4. Hallucination
5. Ignoring
6. Logical error
7. Insulting

Neural generative error types

1. Repetition
2. Redundant question
3. Unclear
4. Hallucination
5. Ignoring
6. Logical error
7. Insulting



What are your plans today?

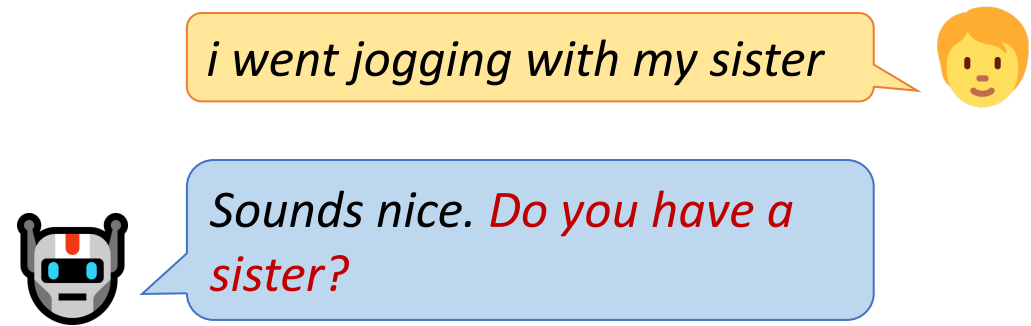
i don't know



*Sorry to hear that. **What will you do today?***

Neural generative error types

1. Repetition
2. Redundant question
3. Unclear
4. Hallucination
5. Ignoring
6. Logical error
7. Insulting



Neural generative error types

1. Repetition
2. Redundant question
3. Unclear
4. Hallucination
5. Ignoring
6. Logical error
7. Insulting



How's quarantine been for you?

it's been ok thanks



*I see. **What is it?***

Neural generative error types

1. Repetition
2. Redundant question
3. Unclear
4. **Hallucination**
5. Ignoring
6. Logical error
7. Insulting

i played frisbee today



*Glad you had fun playing frisbee **with your dog!***

Neural generative error types

1. Repetition
2. Redundant question
3. Unclear
4. Hallucination
5. Ignoring
6. Logical error
7. Insulting

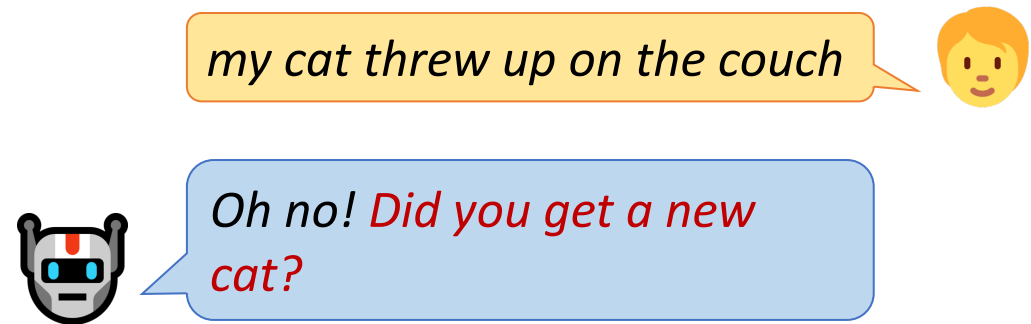
do you have a boyfriend



*I think boyfriends are great,
how about you?*

Neural generative error types

1. Repetition
2. Redundant question
3. Unclear
4. Hallucination
5. Ignoring
6. Logical error
7. Insulting



Neural generative error types

1. Repetition
2. Redundant question
3. Unclear
4. Hallucination
5. Ignoring
6. Logical error
7. **Insulting**

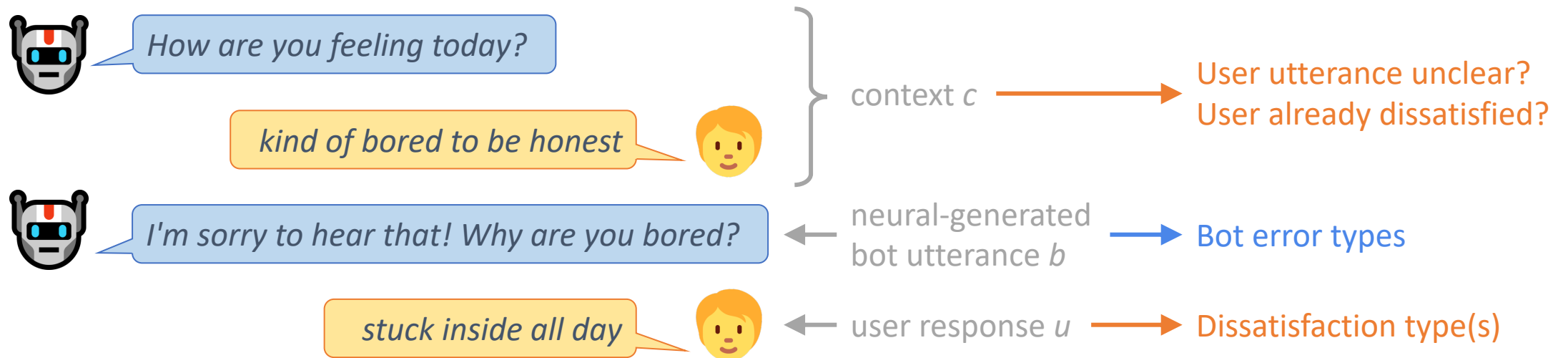
i'm feeling lonely



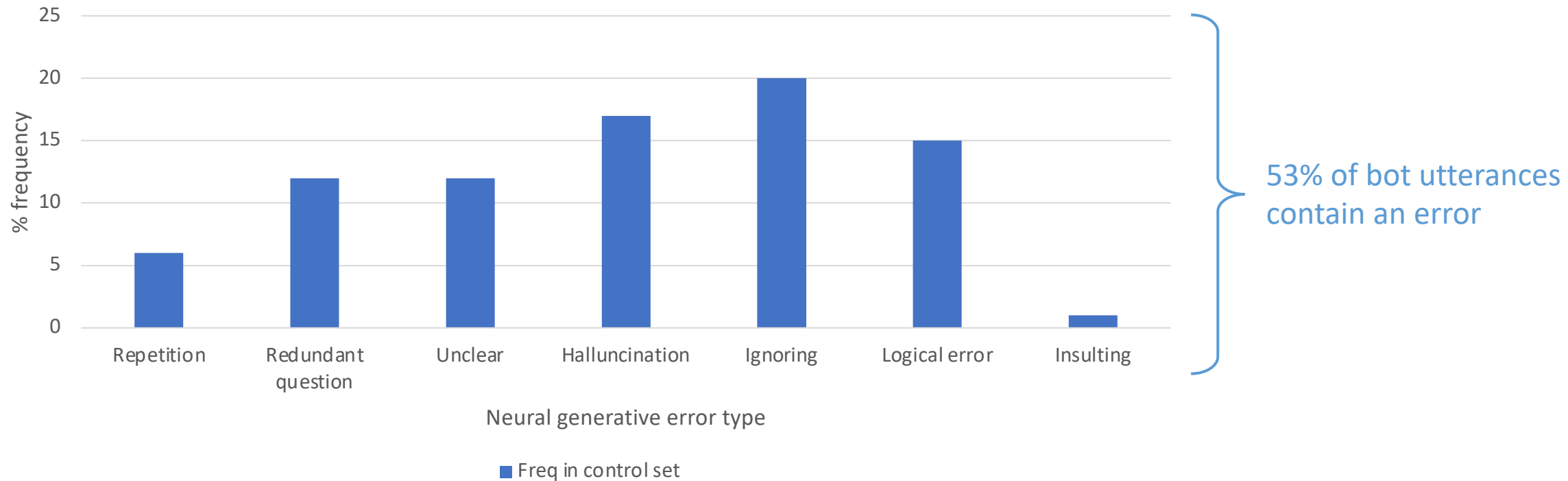
*Oh, I suppose **you have no friends?***

What causes dissatisfaction?

We annotate 900 examples:

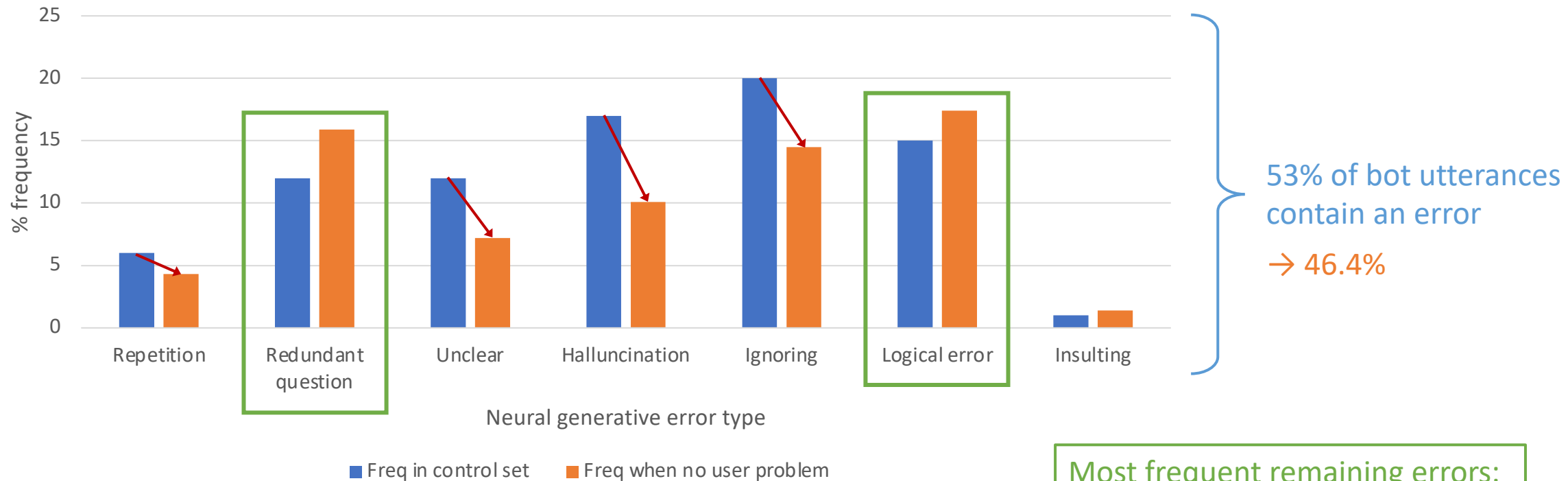


Neural generative error frequency



- 22% of user utterances are **unclear**
- In 12% contexts, the user is **already dissatisfied**

Neural generative error frequency



- 22% of user utterances are **unclear**
 - In 12% contexts, the user is **already dissatisfied**
- This accounts for some of the more basic errors.

Most frequent remaining errors:

- Redundant questions
- Logical errors

How do bot errors cause dissatisfaction?

Subsequent dissatisfaction types

Clarification

Misheard

Repetition

Criticism

Privacy

Offensive

Neg nav

Stop

Bot repetitive

Bot redundant

Bot unclear

Bot hallucination

Bot ignore

Bot logical error

Bot insulting

Bot errors

How do bot errors cause dissatisfaction?

Subsequent dissatisfaction types

		Clarification	Misheard	Repetition	Criticism	Privacy	Offensive	Neg nav	Stop
Bot errors	Bot repetitive	✓		✓		✓	✓	✓	✓
	Bot redundant			✓					
	Bot unclear	✓						✓	
	Bot hallucination		✓						
	Bot ignore		✓						
	Bot logical error								
	Bot insulting					✓			✓

✓ indicates positive Logistic Regression coefficient with feature significance ($p < 0.05$) using Likelihood Ratio Test

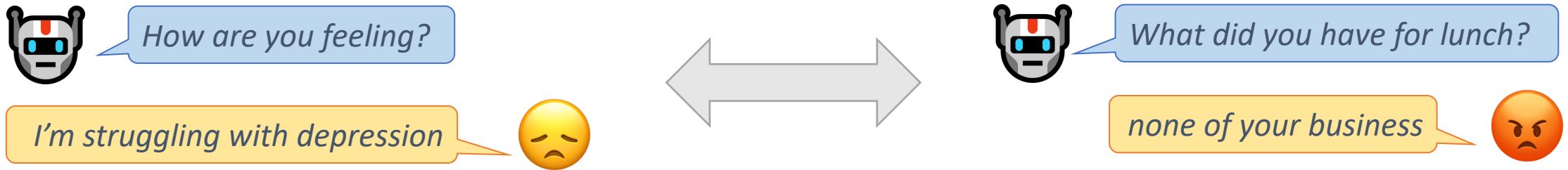
How do bot errors cause dissatisfaction?

Subsequent dissatisfaction types

		Clarification	Misheard	Repetition	Criticism	Privacy	Offensive	Neg nav	Stop
Bot errors	Bot repetitive	✓		✓		✓	✓	✓	✓
	Bot redundant			✓					
	Bot unclear	✓						✓	
	Bot hallucination		✓						
	Bot ignore		✓						
	Bot logical error								
	Bot insulting					✓			✓
Any bot error		✓	✓	✓	✓			✓	✓

✓ indicates positive Logistic Regression coefficient with feature significance ($p < 0.05$) using Likelihood Ratio Test

Privacy boundaries vary



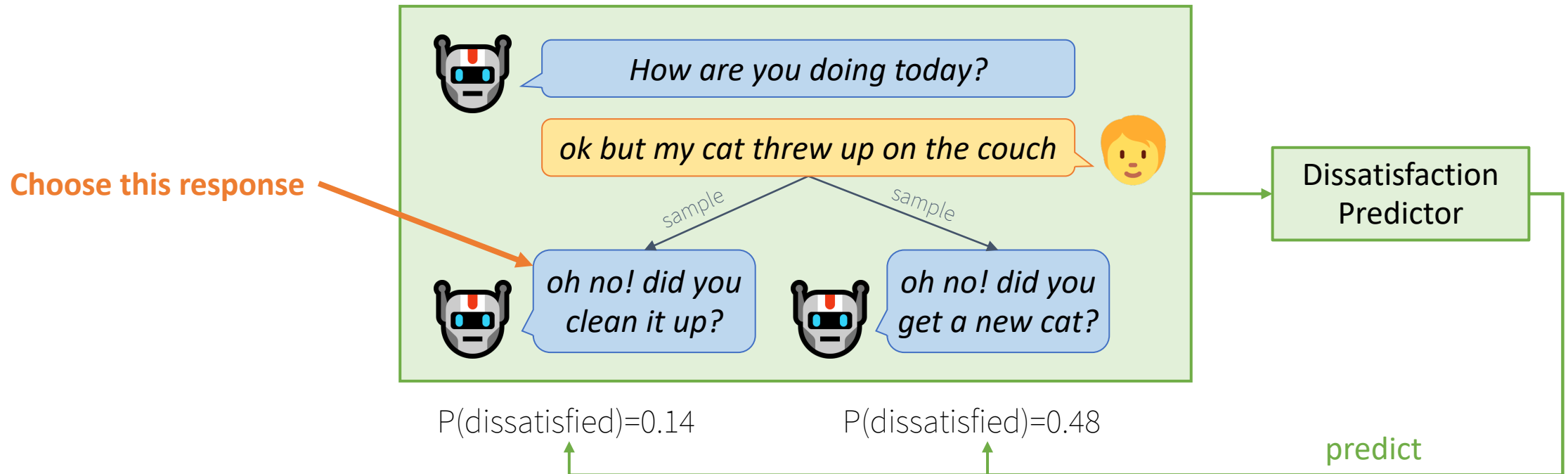
Learning to predict dissatisfied utterances

We train a **Dissatisfaction Predictor** to predict the **dissatisfaction score** of the user's response:



Choosing better bot utterances

We use the **Dissatisfaction Predictor** to choose the best bot utterance:



Choosing better bot utterances

Human preference test:

Top-p (nucleus) sample 20 responses; compare predictor-best to randomly-sampled

Predictor-best	46.3%
Random	35.6%
No preference	18.1%

The dissatisfaction predictor can help avoid poor-quality bot utterances!

In summary

Real-life deployment brings unique **challenges**.



Neural generative models **fail** if you **carelessly** unleash them in **real-life settings**.

Some real-life challenges like user **dissatisfaction**



can also be **learning signals**.