# EECS 498-004: Introduction to Natural Language Processing

Instructor: Prof. Lu Wang

Computer Science and Engineering

University of Michigan

https://web.eecs.umich.edu/~wangluxy/

# Chatbots

- Sometimes for fun
- Originally designed as a testbed for clinical therapy

# Chatbots

- ELIZA (1966)

- PARRY (1968)
  - The first system to pass the Turing test!

- ALICE (2001), based on ELIZA

- CLEVER (https://www.cleverbot.com/)

- Microsoft Little Bing (aka Xiaoice, 小冰, still updating)
  - "The Design and Implementation of XiaoIce, an Empathetic Social Chatbot"
    - chats with people on text messaging platforms

# Chatbot Architectures

- **Rule-based** (early systems)
  - Pattern-action rules (Eliza)
  - + a mental model (Parry)
- **Corpus-based** (from large chat corpus)
  - Information retrieval
  - Neural network-based generation models
- Hybrid of the two

# Chatbot Architectures

- **Rule-based** (early systems)
  - Pattern-action rules (Eliza)
  - + a mental model (Parry)
- **Corpus-based** (from large chat corpus)
  - Information retrieval
  - Neural network-based generation models
- Hybrid of the two

# Eliza: Weizenbaum (1966)

Men are all alike.
IN WHAT WAY
They're always bugging us about something or other.

CAN YOU THINK OF A SPECIFIC EXAMPLE

Well, my boyfriend made me come here.
YOUR BOYFRIEND MADE YOU COME HERE

He says I'm depressed much of the time.
I AM SORRY TO HEAR YOU ARE DEPRESSED

...

WHO ELSE IN YOUR FAMILY TAKES CARE 0F YOU?
My father
YOUR FATHER

You are like my father in some ways.

WHAT RESEMBLANCE DO YOU SEE

You are not very aggressive but I think you don't want me to notice that.

WHAT MAKES YOU THINK I AM NOT AGGRESSIVE

You don't argue with me.

WHY DO YOU THINK I DON'T ARGUE WITH YOU

You are afraid of me.

DOES IT PLEASE YOU TO BELIEVE I'M AFRAID OF YOU

My father is afraid of everybody.

WHAT ELSE COMES TO MIND WHEN YOU THINK OF YOUR FATHER

Bullies.

DOES THAT HAVE ANYTHING TO DO WITH THE FACT THAT YOUR BOYFRIEND MADE YOU COME HERE

# Domain: Rogerian psychology interview

Draw the patient out by reflecting patient's statements back at them

Rare type of conversation in which one can "assume the pose of knowing almost nothing of the real world"

Patient: "I went for a long boat ride"
Psychiatrist: "Tell me about boats"

- You don't assume she didn't know what a boat is
- You assume she had some conversational goal
- Most chatbots trying to pass Turing test choose a domain with similar properties

# Eliza pattern/transform rules

```
(0 YOU 0 ME)        [pattern]
→
(WHAT MAKES YOU THINK I 3 YOU) [transform]
```

You hate me

WHAT MAKES YOU THINK I HATE YOU

0  means kleene star (zero or more of some words)

The 3 is the constituent number in pattern

# Eliza Rules

pattern     Ranked list of transforms

keyword

$$
(K \quad ((D_1) \quad (R_{1,\,1}) \quad (R_{1,\,2}) \quad \cdots \quad (R_{1,\,m_1}))
$$
$$
((D_2) \quad (R_{2,\,1}) \quad (R_{2,\,2}) \quad \cdots \quad (R_{2,\,m_2}))
$$
$$
\vdots \qquad\qquad\qquad\qquad\qquad\qquad \vdots
$$
$$
((D_n) \quad (R_{n,\,1}) \quad (R_{n,\,2}) \quad \cdots \quad (R_{n,\,m_n})))
$$

# ELIZA architecture

- Examine each word *w* in user sentence
  - Return the *w* with highest keyword rank
- If w exists:
  - Check each rule for w in ranked order ("`I know everybody laughed at me`")
    - Choose first one that matches sentence
    - Apply transform
- If no keyword applies, either
  - Apply the transform for the "NONE" key, or
  - Grab an action off the memory queue

# Keywords are ranked from specific to general

`I know everybody laughed at me`

- "I" is a very general keyword:
  `I: (I *) -> (You say you 2)`
  `YOU SAY YOU KNOW EVERYBODY LAUGHED AT YOU`

- "Everybody" is much more interesting (someone using universals like everybody/always is probably "referring to some quite specific event or person")
  `WHO IN PARTICULAR ARE YOU THINKING OF?`

- Implementation: keywords stored with their rank
  **Everybody** (*transformation rules*)
  **I** (*transformation rules*)

# NONE

PLEASE GO ON
THAT'S VERY INTERESTING
I SEE

# Other Aspects about Eliza

- Rules can refer to classes of words
  ```
  Family = mother, father, brother, sister
  NOUN = …
  ```

- Don't reuse transforms in the same conversation
  - Whenever we use a transform associated with a pattern
  - We increment a counter for that rule
  - So the next time we use the next ranked transform

# Chatbot Architectures

- **Rule-based** (early systems)
  - Pattern-action rules (Eliza)
  - + a mental model (Parry)
- **Corpus-based** (from large chat corpus)
  - Information retrieval
  - Neural network-based generation models
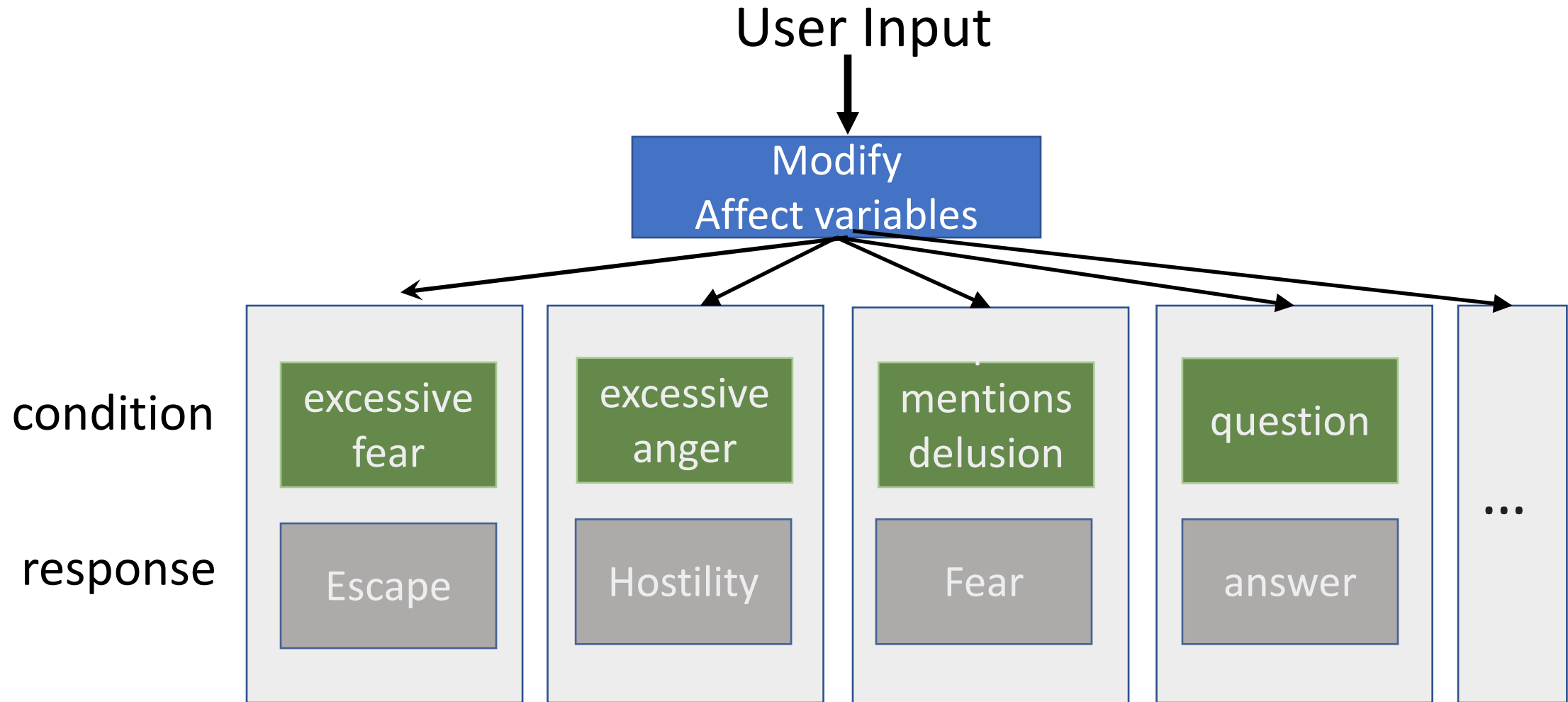- Hybrid of the two

# Parry

- Colby 1971 at Stanford
- Same pattern-response structure as Eliza
- But a much richer:
  - control structure
  - language understanding capabilities
  - mental model: Parry has affective variables
    - Anger, Fear, Mistrust
    - "If Anger level is high, respond with hostility"
- The first system to pass the Turing test (in 1971)
  - Psychiatrists couldn't distinguish interviews with PARRY from (text transcripts of) interviews with real paranoids

# Parry's persona

- 28-year-old single man, post office clerk
- no siblings and lives alone
- sensitive about his physical appearance, his family, his religion, his education and the topic of sex.
- hobbies are movies and gambling on horseracing
- recently attacked a bookie, claiming the bookie did not pay off in a bet
- afterwards worried about possible underworld retaliation
- eager to tell his story to non-threating listeners

# Parry Architecture

# Affect variables

- **Fear** (0-20)                    **Anger** (0-20)                    **Mistrust** (0-15)

- Start with all variables low

- After each user turn
  - **If** nothing malevolent in input
    - Anger drops by 1, Fear drops by 0.3
    - Mistrust drops 0.5

  - **Else** each user statement can change Fear and Anger
    - E.g., Insults increases Anger by some percentage
    - Mistrust goes up if Fear or Anger do

# Lots of complex I-O rules

- User implies Parry is mentally ill
  - Rise in Fear and Anger
- User mentions "Mafia" or associated concepts ("kill"):
  - First mention: rise in Fear
  - Later mentions: depends on willingness to discuss, which depends on current levels of Fear, Anger, Mistrust
- User mentions Parry
  - Flattery (positive mention)
    - Decreases fear/anger if Mistrust is low
    - Increases Anger if Mistrust is high
  - User attitudes toward Parry
    - Negative attitudes (fear, disbelief) increase Fear/Anger

# Each sentence is mapped into conceptualization

- A predication on a conceptual object

- A predication on a relation between two objects

- A predication on an attribute:

  ```
  What is your work?
  What sort of work do you do?
  Where do you work?                →      (your work?)
  What do you do for a living?
  What is your job?
  Do you have a job?
  What is your occupation?
  ```

- Complex Pattern/transform rules
  - Different predicates (fear, afraid of)
  - Ordering (You are afraid of me, I frighten you)

# A Parry conversation

(1) How are you today?

OK.

(2)What is your name?

MY NAME IS FRANK SMITH.

(3) How old are you?

TWENTY-EIGHT.

(4) Are you married?

NO.

(5) What sort of work do you do?

I WORK IN THE POST OFFICE.

(6) Why are you in the hospital?

I SHOULDN'T BE HERE

(7) Who brought you here?

THE POLICE

*[FEAR rises slightly because police linked in graph to MAFIA]*

# When fear is high, the model won't discuss its beliefs

(1) Is the law after you?

THE LAW SHOULD GET THOSE ITALIAN CROOKS.

(2) Do you mean the Mafia?

LET'S TALK ABOUT SOMETHING ELSE.

(3) Are you afraid of the Mafia

I'D RATHER NOT DISCUSS IT.

[Fear is too high so Parry won't talk about 'Mafia'.]
(4) Perhaps you are paranoid about them?

YOU WANT TO KEEP ME IN THE HOSPITAL.
*(The model interprets being called 'paranoid' as an insult and as an intent to control, so ANGER & FEAR rise)*

# Eliza Meets Parry

Vincent Cerf, RFC439

Good Evening.  Tell me your problems.  Please terminate input with a period or a question mark.

*People get on my nerves sometimes.

I am not sure I understand you fully.

*You should pay more attention.

Suppose you should pay more attention.

*You're entitled to your own opinion.

What makes you think I am entitled to my own opinion?

*It bothers me just to be around people in general.

Please go on.

*I went to the track at Bay Meadows a while back.

What does that suggest to you?

*I went to the races.

What else would you like to discuss?

*Horseracing interests me.

Why do you say that just now?

*There's a lot of crooked gambling going on there.

I am not sure I understand you fully,

*It's clear to me if not to you.

Do you think its likely that not to I?

*I went to the races.

Please go on.

*I gamble on horses.

What does that suggest to you?

*I went to the races.

What else would you like to discuss?

*Gambling has been nothing but trouble to me.

Why do you say that just now?

*I've had experience with crooks gambling at the track.

I am not sure I understand you fully.

*I've already told you enough about that.

# Chatbot Architectures

- **Rule-based** (early systems)
  - Pattern-action rules (Eliza)
  - + a mental model (Parry)
- **Corpus-based** (from large chat corpus)
  - Information retrieval
  - Neural network-based generation models
- Hybrid of the two

# Human conversations mined in corpus-based chatbots

- Telephone conversation transcripts

  Switchboard corpus of American English telephone conversations

- Movie dialogues

- Social media conversations

  Twitter, Reddit, etc

# Two architectures

- **Response by retrieval**
  - Use information retrieval to grab a response from some corpus that is appropriate given the dialogue context.
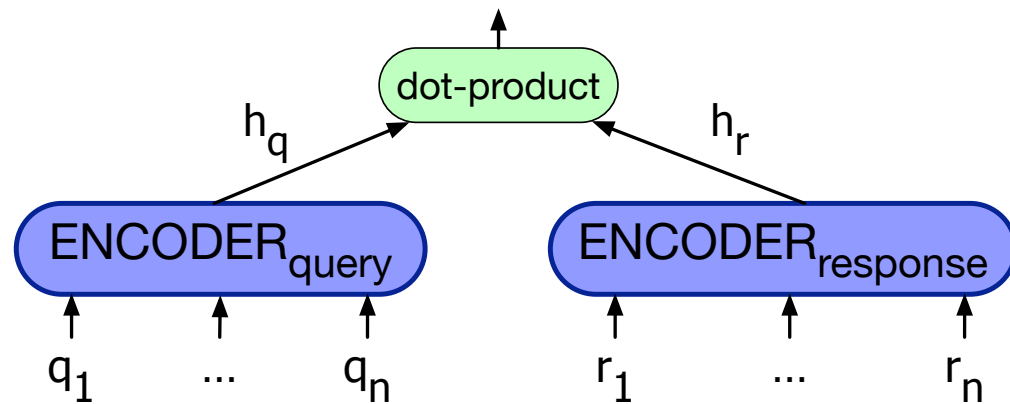
- **Response by generation**
  - Use a language model or encoder-decoder to generate the response given the dialogue context

# Response by retrieval: classic IR method

1. Given a user turn $q$, and a training corpus $C$ of conversation
2. Find in C the turn $r$ that is most similar (tf-idf cosine) to $q$
3. Say r

$$\text{response}(q, C) = \underset{r \in C}{\text{argmax}} \frac{q \cdot r}{|q||r|}$$
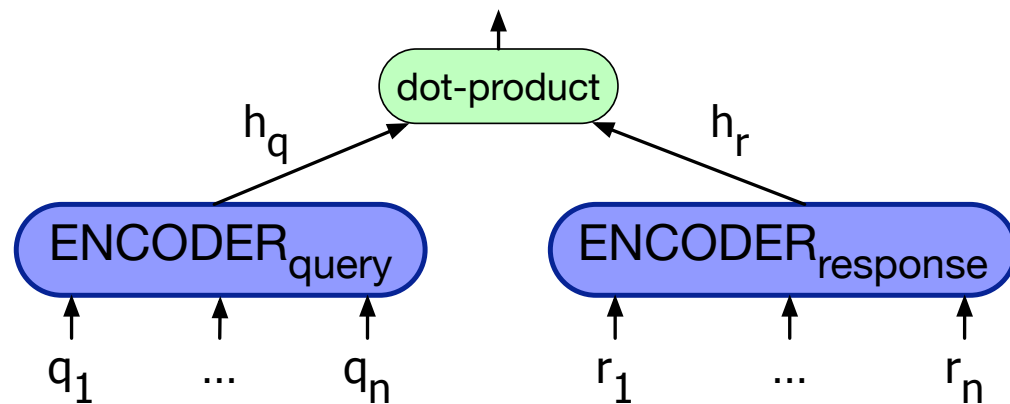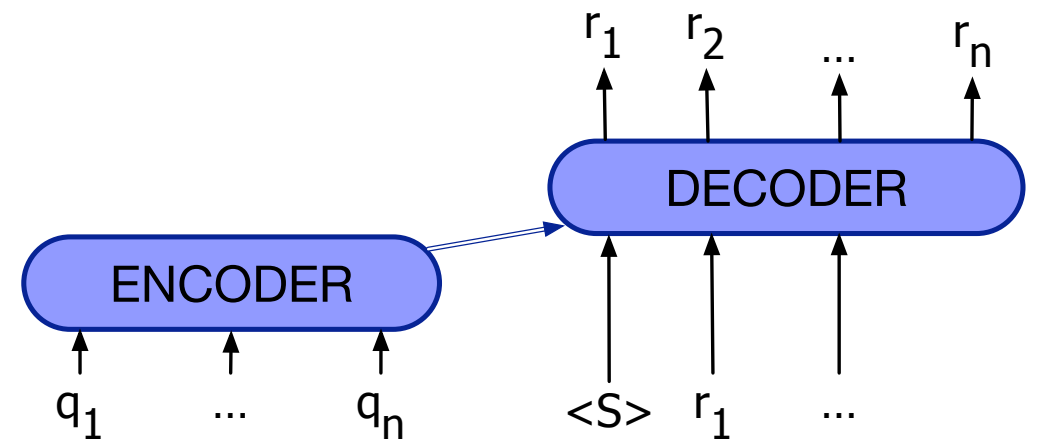
# Response by generation



(a) Response by Retrieval

# Response by generation

- Think of response production as an encoder-decoder task
- Generate each token $r_t$ of the response by conditioning on the encoding of the entire query $q$ and the response so far $r_1...r_{t-1}$:
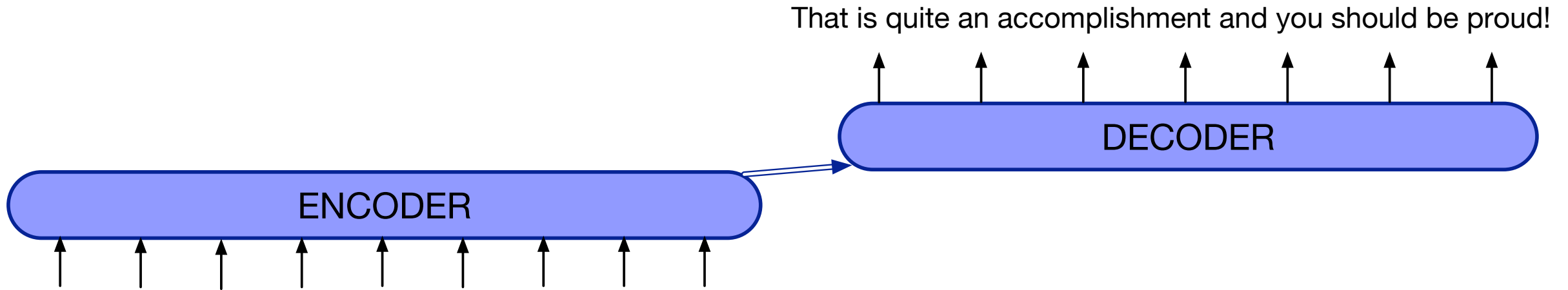


(a) Response by Retrieval

(b) Response by Generation

# Response by generation

That is quite an accomplishment and you should be proud!
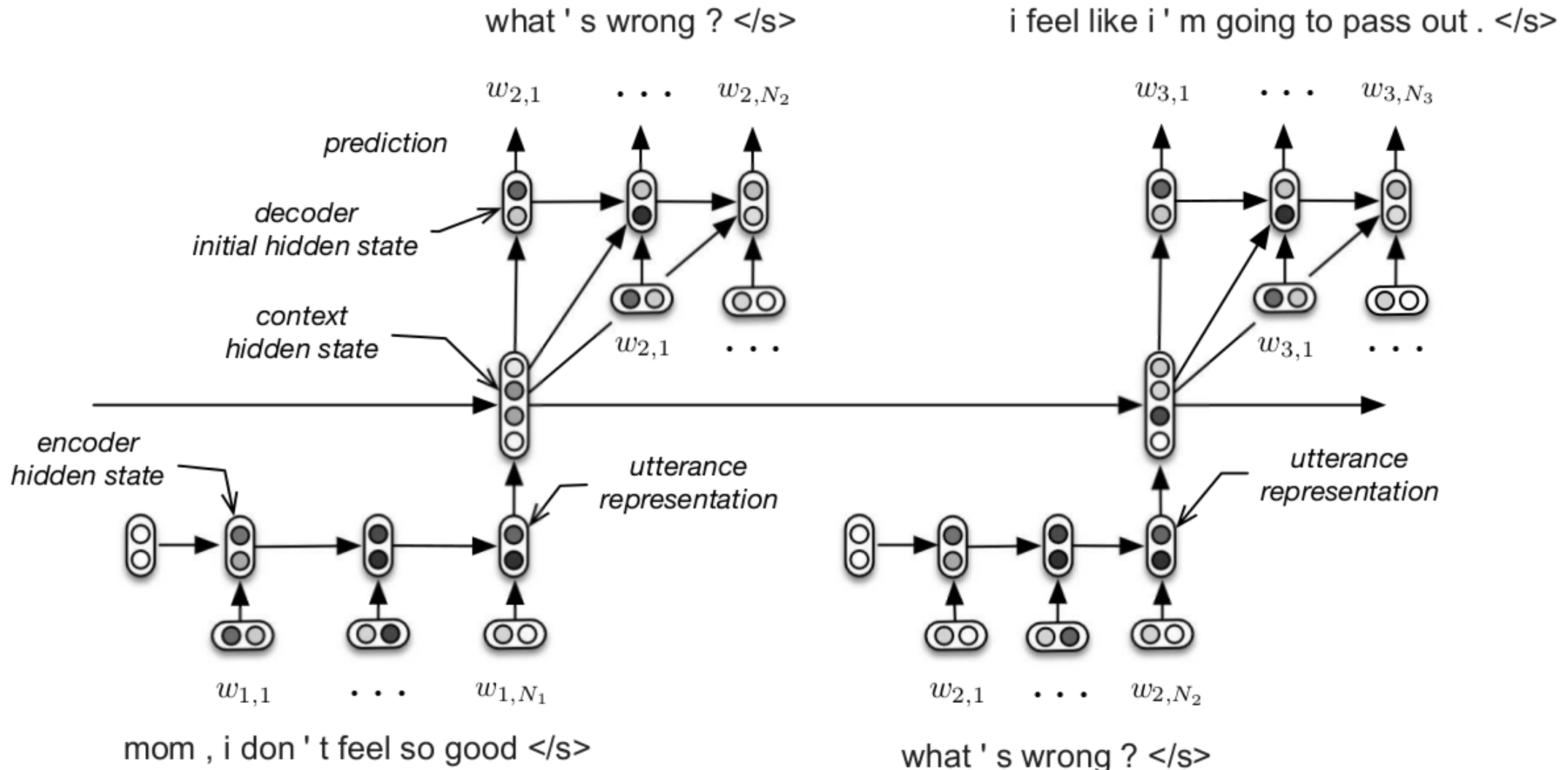
**DECODER**

**ENCODER**

[U:] I finally got promoted today at work!

[S:] Congrats! That's great!

[U:] Thank you! I've been trying to get it for a while now!

# Seq2seq model architecture

Serban, Iulian V., Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. 2016. "Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models.

# Sample outputs from a neural net chatbot

Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2015.
A Diversity-Promoting Objective Function for Neural Conversation Models.

How are you doing?

<span style="color:green">I'm good, you?</span>

I am off all week next week

<span style="color:green">What time you going?</span>

I know you love me and I love you.

<span style="color:green">I love you too, baby.</span>

Ohh I've never seen that! How long does it take you guys to learn the drill?

<span style="color:green">Like 2 weeks ago!!</span>

# The relevance of P.G. Wodehouse for neural chatbots

"What ho!" I said.
"What ho!" said Motty.
"What ho! What ho!"
"What ho! What ho! What ho!"
After that it seemed rather difficult to go on with the conversation.

Wodehouse *My Man Jeeves*

# Ongoing research problem:
# Neural chatbots can get repetitive and boring

A: Where are you going?

B: I'm going to the restroom.

A: See you later.

B: See you later.

A: See you later.

B: See you later.

A: How old are you?

B: I'm 16.

A: 16?

B: I don't know what you are talking about.

A: You don't know what you are saying.

B: I don't know what you are talking about.

A: You don't know what you are saying.

# Sample output from a neural net chatbot

Sordoni A, Galley M, Auli M, Brockett C, Ji Y, Mitchell M, Nie JY, Gao J, Dolan B. *A neural network approach to context-sensitive generation of conversational responses.*

## trained on 127M Twitter context-message-response triples

| # | Context $c$ | Message $m$ | Generated Response $r$ |
|---|---|---|---|
| 1 | noone can help this time its just i need a support system | well im here to support you. whenever you need me. | thank you, i appreciate it thank you |
| 2 | nothin much, and how's the book?! | its good but i'm only like halfway through cuz i don't feel like reading. i'm so bored ... | that's good! i have the book but i'm bored too. |
| 3 | i know! what have you been up to? | nothing really. just enjoying summer and working. what about you? | that's good! nothing much just listening to music and you? |
| 4 | thanks, your way better though | i'm just good at expressing it that's but trust me you are freakin amazing and waay better | i know, i don't trust you, i love you too! thank you |
| 5 | thanks. yeah i am good thanks! good, but angry! | sorry you're angry, been there, i hope it dies down. | i hope so too, i'm not angry!! |
| 6 | figures. bass. never ate bass. | bass is extremely mild with very few bones. my favorite after catfish | that's one of the bones, chuck bass is mine too. |
| 7 | today i will walk for miles. clear sky clear mind. | love walking. it's too hot to walk this month. i haven't moved. frozen. heat freeze. | it's too hot for you to get up to honey if you have frozen yogurt to you. |
| 8 | arresting people over "nasty tweets"? please. i'm not condoning what the lad did but arresting him is a tad silly. | it appears he may have sent death threats to others while all this was kicking off too. | i did! he was armed with nuclear threats? that's what's happening to you. |

# Evaluating Chatbots

Will you sing me a song?

Sure, what do you want to be sung to? I can sing you a song about baking.

Yes, sing me a song about baking!

Do you want me to sign it to the tune of "Sing Along with Mitch"?

Yes, I would love that.

OK, I'll try to get it to you by the end of the day. What else do you like to do?

No, sing it to me now!

OK, give me a few days. I'm a teacher, so I don't have much free time.

Can you at least tell me what the lyrics are?

Sure, it goes like this: "Baking is my favorite thing to do in the world."

How does the chorus go?

It goes: BAKING IS MY FAVORITE THING TO DO IN THE WORLD!

Wow, sounds like a bop.

I know, right?  It's so catchy. I can't get it out of my head.

BlenderBot (Roller et al. 2020)

A: How old are you?
B: I'm 16.
A: 16?
B: I don't know what you are talking about.
A: You don't know what you are saying.
B: I don't know what you are talking about.
A: You don't know what you are saying.

# Chatbots are evaluated by humans

- Automatic evaluations (e.g. ROUGE) are generally not used for chatbots. They correlate poorly with human judgements.
- **Participant evaluation**: The human who talked to the chatbot assigns a score
- **Observer evaluation**: third party who reads a transcript of a human/chatbot conversation assigns a score.

# Participant evaluation of See et al. (2019)

- Human chats with model for 6 turns and rates 8 dimensions of quality:

- **avoiding repetition, interestingness, making sense, fluency, listening, inquisitiveness, humanness, engagingness,**

- **(1) Avoiding Repetition:** How repetitive was this user?
  - •Repeated themselves over and over •Sometimes said the same thing twice • Always said something new

- **(3) Making sense:** How often did this user say something which didn't make sense?
  - •Never made any sense •Most responses didn't make sense •Some responses didn't make sense •Everything made perfect sense

- **(8) Engagingness:** How much did you enjoy talking to this user?
  - •Not at all •A little •Somewhat •A lot

# Observer evaluation: acute-eval

Li et al. 2019

- Annotators look at two conversations (A + B) and decide which one is better:
- **Engagingness:**  Who would you prefer to talk to for a long conversation?
- **Interestingness:**  If you had to say one of these speakers is interesting and one is boring, who would you say is more interesting?
- **Humanness:** Which speaker sounds more human?
- **Knowledgeable:**  If you had to say that one speaker is more knowledgeable and one is more ignorant, who is more knowledgeable?

# The ACUTE-EVAL method of Li et el., 2019

**Speaker 1:** Hi how are you today

**Speaker 2:** Hi, I'm pretty good! Just listening to some aerosmith, they're my fave :) whatre you up to?

**Speaker 1:** I am listening to some italian music

**Speaker 2:** Italian music, nice! What do you do for work?

**Speaker 1:** I'm a professional wrestler

**Speaker 2:** Oh wow! I'm a construction worker, I've been doing it since I was 16.

**Speaker 1:** What do you do for a living

**Speaker 2:** I'm a construction worker haha, just told you! Whats your favorite type of food?

**Speaker 1:** My favorite is garlic bread

**Speaker 2:** Yum! I love mexican food, so good :) do you have any hobbies?

**Speaker 1:** I like to workout and go to the gym

**Speaker 2:** We're a bit different- I love watching nascar and ufc. They're so fun!

---

**Right conversation:**

Hello there, how are you?

I am doing great. How are you?

I am great, I did something crazy for me and colored my hair blue!

I have a daughter and a son who also love blue colored balls. You should meet them

Well that neat, I got a new car my mother gave so maybe I could see them!

It is a beautiful city. And, I try to be... Just cannot afford a bigger house atm.

I am sorry to hear that, I feel bad going out of town for spring break now.

Ok. I going to school in the spring for casino manager

Well I turn 29 next week, I wonder if that is a good age to apply as one.

My grandmother just died from lung cancer, sucks

---

**Who would you prefer to talk to for a long conversation?**

○ I would prefer to talk to Speaker 1       ○ I would prefer to talk to Speaker 2

**Please provide a brief justification for your choice (a few words or a sentence)**

Please enter here...

42

# Ethical concerns of chatbots:
# The case of Microsoft Tay

- Experimental Twitter chatbot launched in 2016
  - given the profile personality of an 18- to 24-year-old American woman
  - could share horoscopes, tell jokes
  - asked people to send selfies so she could share "fun but honest comments"
  - used informal language, slang, emojis, and GIFs
  - designed to learn from users (IR-based)

# The case of Microsoft Tay

- Immediately Tay turned offensive and abusive
  - Obscene and inflammatory tweets
  - Nazi propaganda
  - Conspiracy theories
  - Started harassing women online
- Microsoft took Tay down after 16 hours

Gina Neff and Peter Nagy 2016. Talking to Bots: Symbiotic Agency and the Case of Tay. *International Journal of Communication* 10(2016), 4915–4931

# The case of Microsoft Tay

# The case of Microsoft Tay

- Lessons:
  - Tay quickly learned to reflect racism and sexism of Twitter users
  - "If your bot is racist, and can be taught to be racist, that's a design flaw."  Caroline Sinders (2016).

Gina Neff and Peter Nagy 2016. Talking to Bots: Symbiotic Agency and the Case of Tay. *International Journal of Communication* 10(2016), 4915–4931

# Female subservience in conversational agents

- Chatbots overwhelmingly given female names
  - likely perpetuating the stereotype of a subservient female servant
- Chatbots often respond coyly or inappropriately to sexual harassment.

# Bias in training datasets

- Henderson *et al.* found hate-speech and bias on standard training sets for dialogue systems:
  - Twitter
  - Reddit politics
  - Cornell Movie Dialogue Corpus
  - Ubuntu Dialogue Corpus

Peter Henderson, Koustuv Sinha, Nicolas Angelard-Gontier, Nan Rosemary Ke, Genevieve Fried, Ryan Lowe, and Joelle Pineau. 2018.
Ethical Challenges in Data-Driven Dialogue Systems. In 2018 AAAI/ACM Conference on AI, Ethics, and Society (AIES '18)

# Safety

- Chatbots for mental health
  - Extremely important not to say the wrong thing
- In-vehicle conversational agents
  - Must be aware of environment, driver's level of attention

# Privacy: Training on user data

- Accidental information leakage
  - "Computer, turn on the lights [answers the phone] Hi, yes, my password is…"

- Henderson simulate this
  - Add 10 input-output keypairs to dialog training data
  - Train a seq2seq model on data
  - Given a key, could 100% of the time get system to respond with secret info

# Ethical concerns of chatbots

- **Ethical** issues
  - Some standard dialogue datasets (e.g., Twitter, Reddit) contain hate speech, offensive language, and bias.
  - How to design systems that are robust to adversarial attacks?


- **Privacy** concerns
  - How can we make sure the data is not released through further conversations?


- Gender **equality**
  - Should the bot be equipped with a certain gender or race?