*Alexa, can you help me?*

*I don't know what to do.*
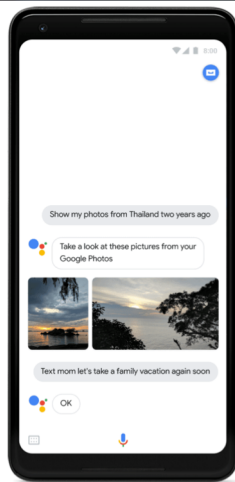
*Dialog Systems*

**João Sedoc**

jsedoc@jhu.edu
Johns Hopkins
Computer Science

JOHNS HOPKINS
WHITING SCHOOL
*of* ENGINEERING

# Chatbots are Ubiquitous: Personal Agents, Games, Education, Business & Medicine

# Lots of Tools

# Artificial Intelligence

- Can robots understand language?

- Can robots actually think?

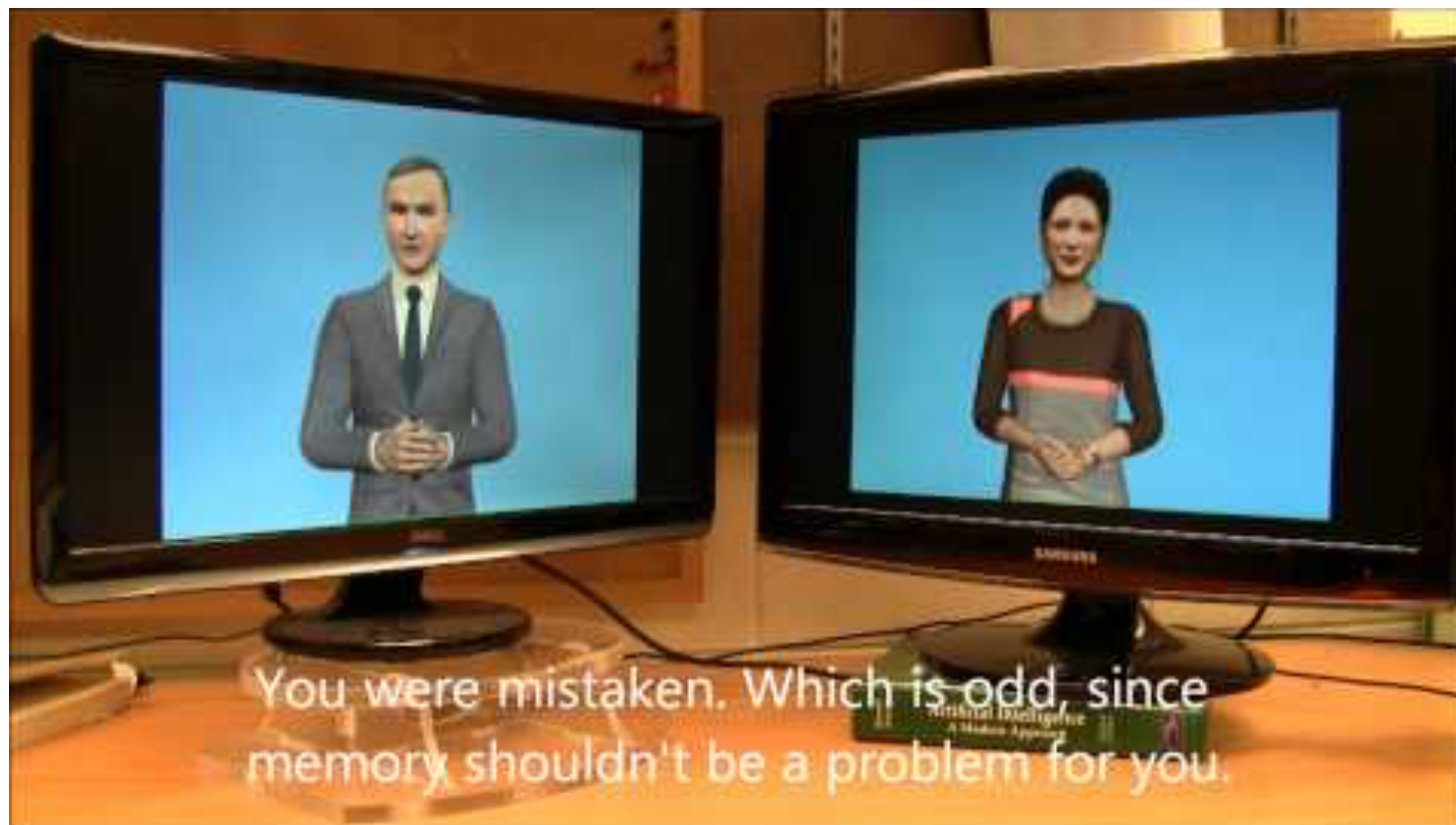- Not clear definition of intelligence or how to measure it!

- The Turing Test (1950)

- Indirect assessment of intelligent behaviour

# AI with AI conversations: Cleverbot (Carpenter, 2011)



You were mistaken. Which is odd, since memory shouldn't be a problem for you.

# Challenges for Artificial Intelligence

- **Knowledge Representation**
  - ❏ about learning, storing and retrieving relevant infor-mation about the world and one's previous experiences

- **Commonsense reasoning\***
  - ❏ about using world knowledge for interpreting, explaining and predicting daily life events and outcomes

# Challenges for Conversational Agents

**Key Factors**

| Content / Context | Personality & Persona | Emotion & Sentiment | Behavior & Strategy |
|---|---|---|---|

**Key Issues**

| Semantics | Consistency | Interactiveness |
|---|---|---|

**Key Technologies**

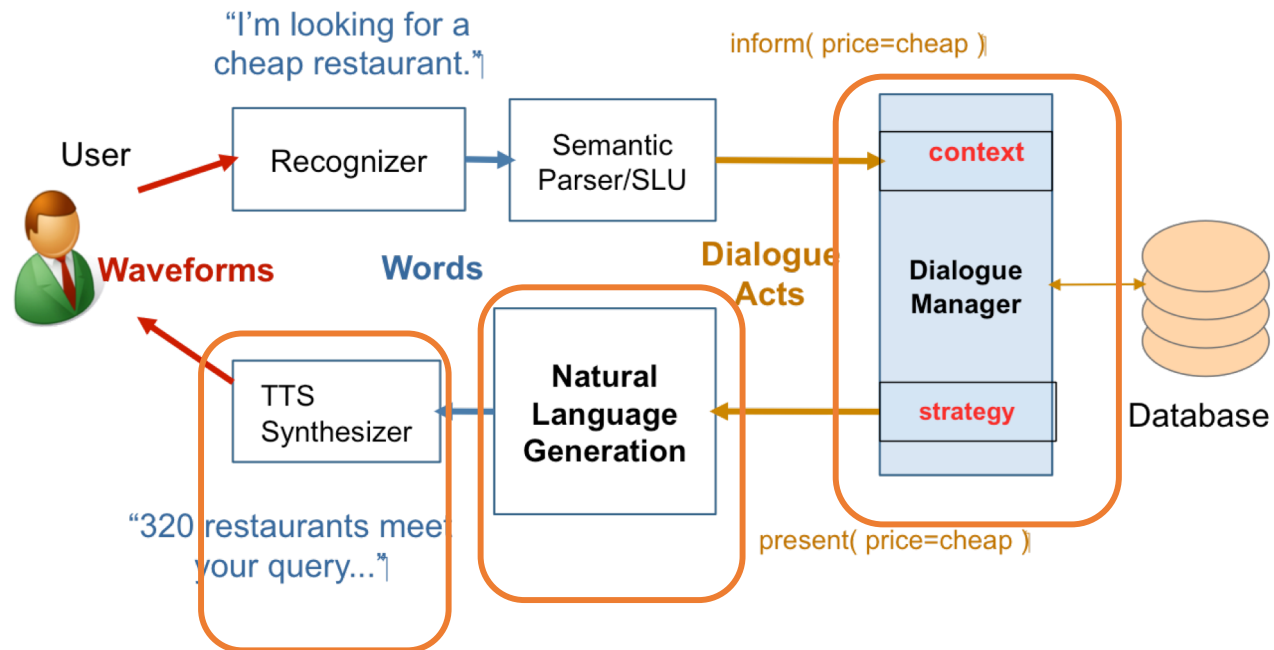| Named Entity Recognition | Domain/Topic Intent Detection | Sentiment/Emotion Detection | Knowledge & Reasoning |
|---|---|---|---|
| Entity Linking | Natural Language Generation | Personalization | Dialog Planning & Context Modelling |

From Huang et al., 2019, "Challenges in Building Intelligent Open-Domain Systems"

# Spoke Dialog System Architecture

# Two Types of Systems

1.  Chatbots

2.  Goal-based (Dialog agents)
    • SIRI, interfaces to cars, robots, …
    • Booking flights, restaurants, or question answering

# Chatbot Architectures

**Rule-based**
1. Pattern-action rules (Eliza)
   + a mental model (Parry)

**Corpus-based (from large chat corpus)**
2. Information Retrieval
3. Neural network encoder-decoder

# Eliza pattern/transform rules

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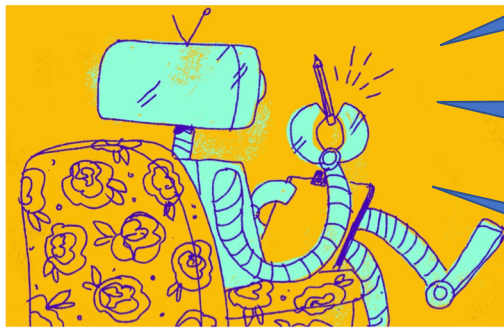
0  means Kleene *

The 3 is the constituent # in pattern

You hate me

WHAT MAKES YOU THINK I HATE YOU

# Personality in chatbots:
# Eliza and Parry

**Eliza**

**Parry**



Good Evening. Tell me your problems.

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.

# Chatbot Architectures

**Rule-based**

1.  Pattern-action rules (Eliza)
    + a mental model (Parry)

**Corpus-based (from large chat corpus)**

2.  Information Retrieval
3.  Neural network encoder-decoder

# 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.

# Information Retrieval based Chatbots

Idea: Mine conversations of human chats or human-machine chats
   Microblogs:   Twitter or Weibo (微博)
   Movie dialogs

- Cleverbot (Carpenter 2017 http://www.cleverbot.com)
- Microsoft XiaoIce
- Microsoft Tay

# Two IR-based Chatbot Architectures

1. Return the response to the most similar turn
   - Take user's turn (*q*) and find a (tf-idf) similar turn *t* in the corpus C

     *q = "do you like Doctor Who"*

     *t' = "do you like Doctor Strangelove"*
   - Grab whatever the response was to *t*.

$$r = response \left( \underset{t \in C}{\text{argmax}} \frac{q^T t}{||q||||t||} \right)$$   Yes, so funny

2. Return the most similar turn

$$r = \underset{t \in C}{\text{argmax}} \frac{q^T t}{||q||||t||}$$   Do you like Doctor Strangelove

# Deep Semantic Similarity Model

Relevance measured
by cosine similarity

$sim(X, Y)$

128

128

$f(.)$

$g(.)$

Word sequence $x_t$

$w_1, w_2, \ldots, w_{T_Q}$

$w_1, w_2, \ldots, w_{T_D}$

X

Y

**Learning:** maximize the similarity between X (source) and Y (target)

**Representation:** use DNN to extract abstract semantic features, $f$ or $g$ is a

- Multi-Layer Perceptron (MLP) if text is a bag of words [Huang+ 13]
- **Convolutional Neural Network (CNN) if text is a bag of chunks** [Shen+ 14]
- Recurrent Neural Network (RNN) if text is a sequence of words [Palangi+ 16]

# Chatbot Architectures

**Rule-based**

1. Pattern-action rules (Eliza)
   + a mental model (Parry)

**Corpus-based (from large chat corpus)**

2. Information Retrieval
3. Neural network encoder-decoder

# Neural Network Encoder-Decoder Generative Models

# Response Generation Systems

- **End-to-end** systems.
- Learn from "raw" dialogue data (e.g. OpenSubtitles).
- No semantic or pragmatic annotation required.
- Mainly successful in open-domain, non-task oriented systems.

# Neural Conversation Model (NCM)
## vs
# Rule-Based Model (Cleverbot)

---

**User:** are you a follower or a leader ?
**CleverBot:** no !!
**NCM:** i 'm a leader .

Vinyals and Le 2015
**"A Neural Conversation Model"**

*Image borrowed from farizrahman4u/seq2seq*

# Neural Network Language Models (NNLMs)

**Feed-forward NNLM**

Output

aardvark = 0.0082

...

**store = 0.0191**

...

zygote = 0.003

Hidden 2

Hidden 1

| Embedding | Embedding | Embedding | Embedding |

| he | drove | to | the |

Microsoft
TRANSLATOR

# Neural Network Language Models (NNLMs)

**Feed-forward NNLM**

**Recurrent NNLM**

Output

aardvark = 0.0082

...

**store = 0.0191**

...

zygote = 0.003

Hidden 2

Hidden 1

| Embedding | Embedding | Embedding | Embedding |

| he | drove | to | the |

Output

aardvark = 0.000041

...

**drove = 0.045**

...

zygote = 0.00003

Output

aardvark = 0.000054

...

**to = 0.267**

...

zygote = 0.000009

Recurrent Hidden

Recurrent Hidden

Recurrent Hidden

Recurrent Hidden

Embedding

Embedding

he

drove

Microsoft Translator

# Sentence Encoder



| Recurrent Hidden | → | Recurrent Hidden | → |

↑

| Recurrent Hidden | → | Recurrent Hidden | → |

↑

Embedding      Embedding

↑         ↑

| How |      | are |

Microsoft Translator

# Sequence to Sequence Model



Sutskever et al. 2014
*"Sequence to Sequence Learning with Neural Networks"*

*Image borrowed from farizrahman4u/seq2seq*

# Sequence to Sequence Model



Vinyals and Le 2015
**"A Neural Conversation Model"**

*Image borrowed from farizrahman4u/seq2seq*

# Sequence to Sequence Model



Hoe  gaat  het  <EOL>

How  are  you  <EOL>

LSTM Encoder

LSTM Decoder

S = Source
T = Target

$$1/|\mathcal{S}| \sum_{(T,S)\in\mathcal{S}} \log p(T|S)$$

$$\hat{T} = \arg\max_{T} p(T|S)$$

# Sequence to Sequence Model



S = Source
T = Target

$$1/|\mathcal{S}| \sum_{(T,S)\in\mathcal{S}} \log p(T|S)$$

$$\hat{T} = \arg\max_{T} p(T|S)$$

# Neural Conversational Models



Sequence-to-sequence (Seq2Seq), the probability of the next utterance,

$$P(T \mid S) = P(u_{t+1} \mid u_t) = \prod_{i=1}^{N_t} P(x_{t+1,i} \mid x_{t+1,i-1}, \ldots, x_{t+1,1}, f(u_t)),$$

# Hierarchical Sequence to Sequence Model



Serban, Iulian V., Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. 2015. **Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models**.

# Neural Conversational Models

Sequence-to-sequence (Seq2Seq), the probability of the next utterance,

$$P(T \mid S) = P(u_{t+1} \mid u_t) = \prod_{i=1}^{N_t} P(x_{t+1,i} \mid x_{t+1,i-1}, \ldots, x_{t+1,1}, f(u_t)),$$

an utterance at turn $t$ is defined as $u_t = x_{t,1}, x_{t,2}, \ldots, x_{t,N_t}$

# Uninteresting, Bland, and Safe Responses

How was your weekend?

I don't know.

What did you do?

I don't understand what you are talking about.

This is getting boring…

Yes that's what I'm saying.

# Uninteresting, Bland, and Safe Responses

Common MLE objective (maximum likelihood)

(whatever the user says) $\xrightarrow{p(\text{target}|\text{source})}$ I don't know. 👍

I don't understand...

That's what I'm saying

Mutual information objective:

(whatever the user says) $\xrightarrow{p(\text{target}|\text{source})}$ I don't know. 👍

(whatever the user says) $\xleftarrow{p(\text{source}|\text{target})}$ I don't know. 👎

# Response Diversity Promotion

Mutual information objective:

$$\hat{T} = \arg\max_{T} \left\{ \log \frac{p(S,T)}{p(S)p(T)} \right\}$$

$$\hat{T} = \arg\max_{T} \left\{ \boxed{\log p(T|S)} - \boxed{\lambda \log p(T)} \right\}$$

standard likelihood     anti-LM

$$\hat{T} = \arg\max_{T} \left\{ (1-\lambda)\log p(T|S) + \lambda \log p(S|T) \right\}$$

$$p(\text{target}|\text{source})$$
$$\longrightarrow$$
$$\longleftarrow$$
$$p(\text{source}|\text{target})$$

*Bayes' rule*

*Bayes' theorem*

# Next Steps for Chatbots

- Knowledge grounding – knowledge bases

# Next Steps for Chatbots

- Knowledge grounding - personalization

# Next Steps for Chatbots

- Knowledge grounding – conversational history

# Next Steps for Chatbots

- Persona

# Chatbots: pro and con

- Pro:
  - Fun
  - Applications to counseling
  - Good for narrow, scriptable applications

- Cons:
  - They don't really understand
  - Rule-based chatbots are expensive and brittle
  - IR-based chatbots can only mirror training data
    - The case of Microsoft Tay
      - (or, Garbage-in, Garbage-out)
  - Generative chatbot are hard to control (more later…)

# Two Types of Systems

1. Chatbots

2. Goal-based (Dialog agents)
   - SIRI, interfaces to cars, robots, …
   - Booking flights, restaurants, or question answering

# Goal-based (Dialog agents)
# Task-Oriented

## What kinds of problems?

| | | Chitchat (social bot) |
|---|---|---|
| "I am smart" | Turing Test ("I" talk like a human) | |
| "I have a question" | Information consumption | |
| "I need to get this done" | Task completion | |
| "What should I do?" | Decision support | Goal-oriented dialogues |

# Task Representation and NLU

*"Show me flights from Edinburgh to London on Tuesday."*

SHOW:
    FLIGHTS:
        ORIGIN:
            CITY:  Edinburgh
            DATE:  Tuesday
            TIME:  ?
        DEST:
            CITY: London
            DATE:  ?
            TIME:  ?

# Slot Filling Dialog

- **Domain**: movie, restaurant, flight, …

- **Slot**: information to be filled in before completing a task
  - For Movie-Bot: movie-name, theater, number-of-tickets, price, …

- **Intent** (dialog act):
  - Inspired by speech act theory (communication as action)
    request, confirm, inform, thank-you, …
  - Some may take parameters:
    thank-you(), request(price), inform(price=$10)

"Is Kungfu Panda the movie you are looking for?"

⇩

confirm(moviename="kungfu panda")

# Dialog Engineering as Finite State Automata

# Dialog State Tracking

# Reinforcement Learning



$$Q^\pi(s,a) = \sum_{s'} T^a_{ss'}[R^a_{ss'} + \gamma V^\pi(s')];$$

Bellmann optimality equation (1952), see [Sutton and Barto, 1998].

# 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)
- What could go wrong?

# 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. That's bad design, and that's on you." 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

# Evaluation

# Evaluation

1. Slot Error Rate for a Sentence

$$\frac{\text{\# of inserted/deleted/subsituted slots}}{\text{\# of total reference slots for sentence}}$$

2. End-to-end evaluation (Task Success)

# Evaluation of Goal (Task) vs Chatbot (Non-Task)

**Task-based**

- Human
  - End-of-task subjective task success
  - End-of-task ratings
- Automatic
  - Objective task success (Rieser, Keizer, Lemon, 2014)
  - Automatic estimates of User Satisfaction, (Rieser & Lemon, LREC 2008)

**Non-task Based**

- Human
  - Turn-based appropriateness (WOCHAT)
  - Turn-based pairwise (Li et al. 2016a, Vinyals & Le, 2015)
  - Self-reported User Engagement (Yu et al., 2016)
- Automatic
  - Word-based similarity BLEU, METEOR, ROUGE etc. (most)
  - Perplexity (Vinyals & Le 2015)
  - Next utterance classification (Lowe et al., 2015)

# References for Automatic Evaluation

1-to-1
Syntactically
and
Semantically

1-to-1
Semantically

1-to-Some
Semantically

1-to-Many
Semantically

Automatic
Speech
Recognition

Machine
Translation

Text
Simplification

Dialog
Generation

Sentence
Compression

Abstractive
Summarization

# Why Are We Worried about Evaluation?

Tournaments in machine learning and machine translation led to large advances

Amazon Alexa Prize – largely infeasible for academic scale

# Current Automatic Metrics Weakly Correlate with Human Judgements

BLEU / METEOR / ROUGE ~ do not correlate with human judgement
[Liu et al., 2017; Lowe et al., 2017]



Figures from Liu et al., 2017

# Dialog Evaluation Metrics are an Active Area of Research

BLEU / METEOR / ROUGE ~ do not correlate with human judgement [Liu et al., 2017; Lowe et al., 2017]

Sentence embedding based metrics
 ADEM  [Lowe, et al., 2017]
 RUBER [Toa, et al., 2017]
 Greedy word embeddings [Liu et al.,2017]

Human evaluation is still the gold standard

# Interactive Evaluation of Chatbots Requires a Lot of Data == Expensive

# Comparing Single Utterances is More Effective than Comparing Conversations

Before starting we will show you an example.

For example, you may be given the conversation:

**hey, what's up?**
**hey, want to go to the movies tonight?**

Your task is to choose the most appropriate response:

**A: sure that sounds great! what movie do you want to see?**
**B: i know that was hilarious!**

Response A is clearly a better answer, as it specifically addresses the question asked in the context.

# Ethical Issues

# Privacy