**SFU** NatLangLab

CMPT 713: Natural Language Processing

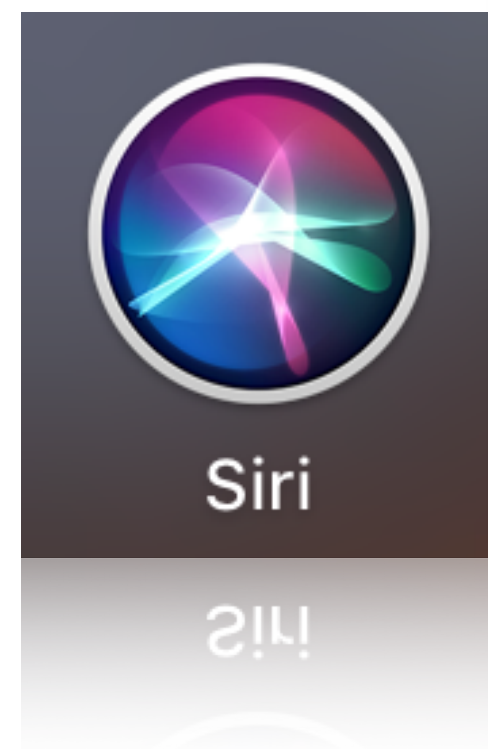# Dialogue and large language models

Spring 2023

2023-03-23

# Dialogue
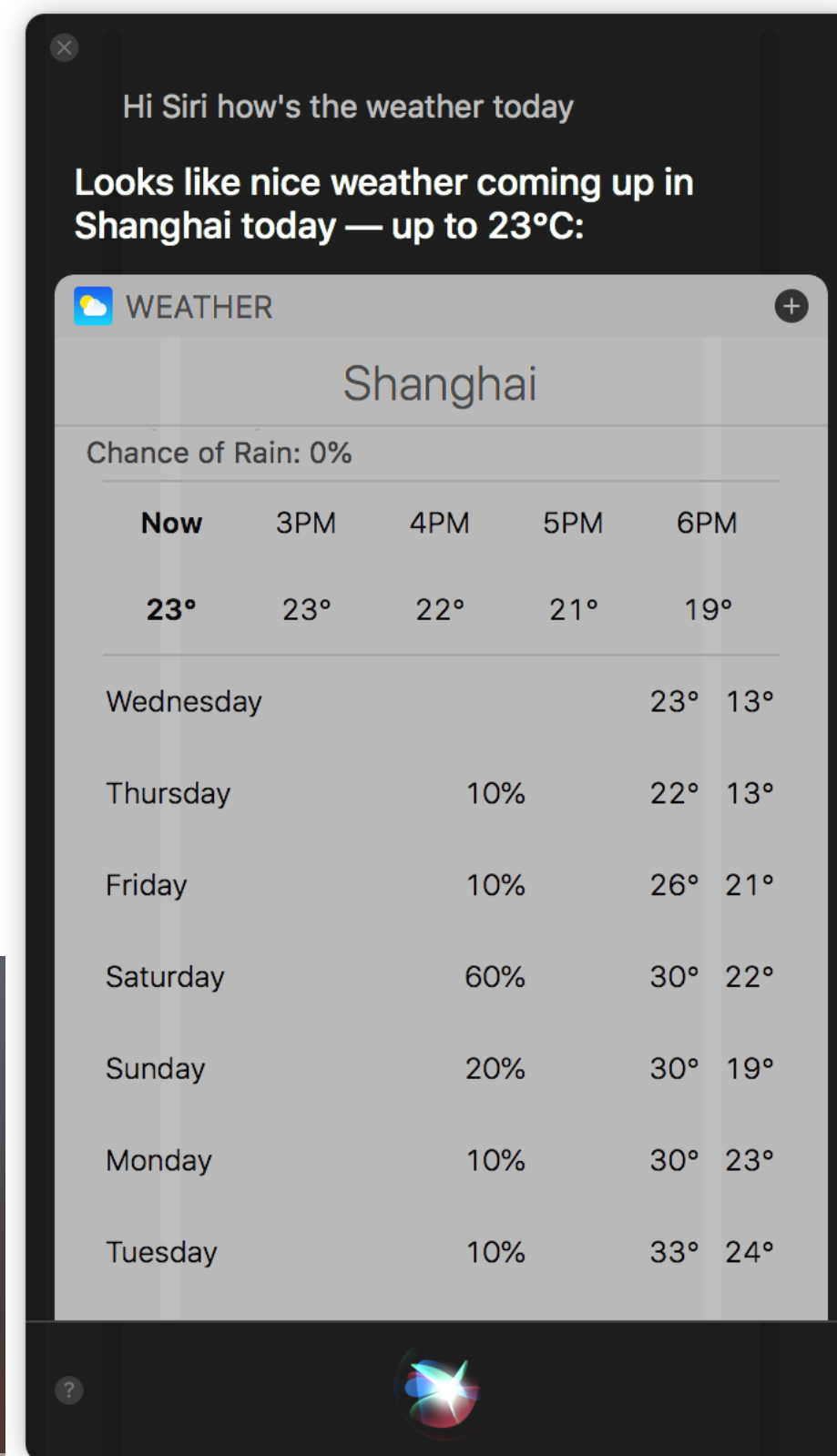
# What's a Dialogue System?
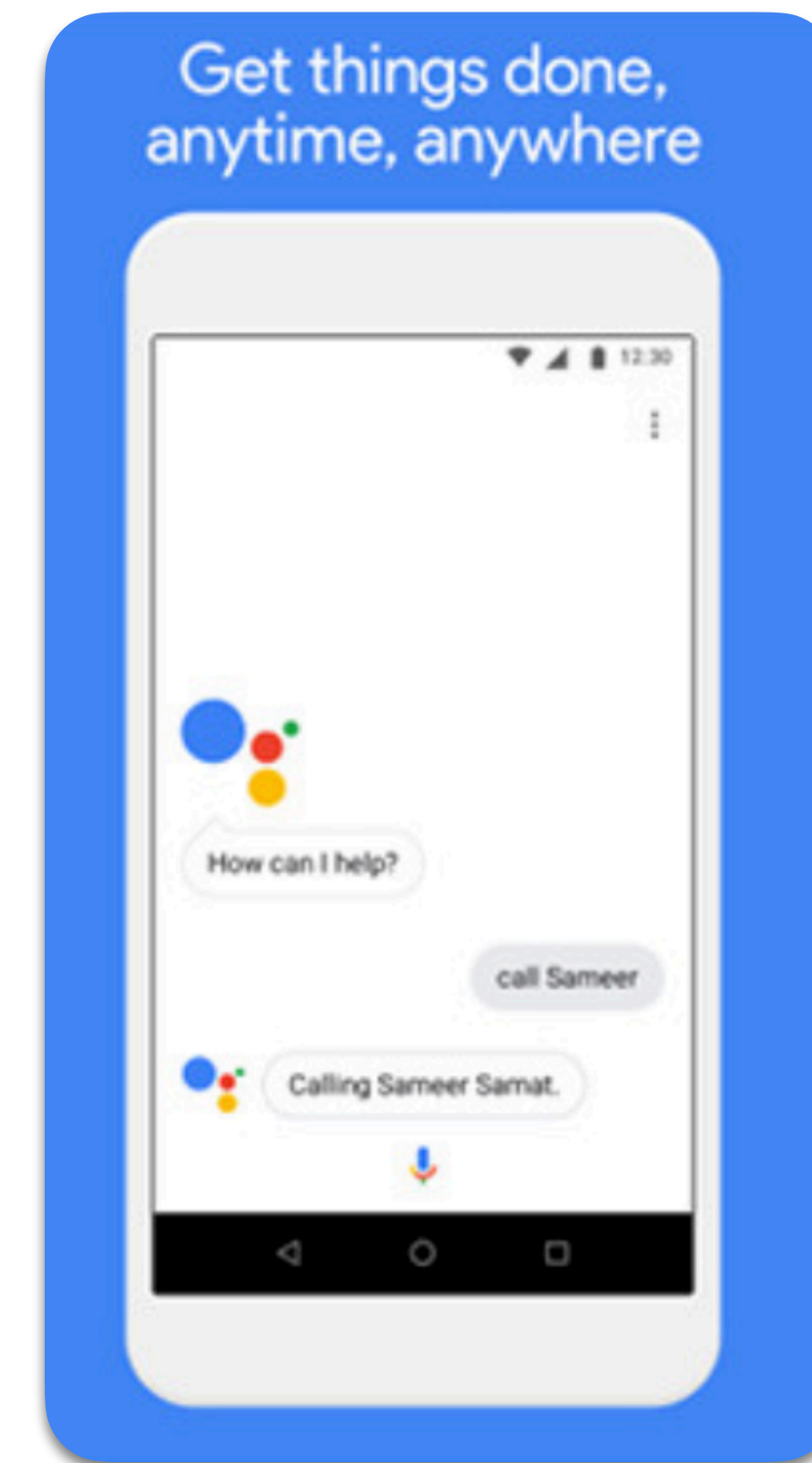
Dialog Systems are **HOT** 🔥.  — Did you use it?

**OpenAI**

*How ChatGPT Kicked Off an A.I. Arms Race*

Feb, 2023

**Conversational agents**

Hi Siri how's the weather today

Looks like nice weather coming up in Shanghai today — up to 23°C:

WEATHER

Shanghai

Chance of Rain: 0%

| | Now | 3PM | 4PM | 5PM | 6PM |
|---|---|---|---|---|---|
| | **23°** | 23° | 22° | 21° | 19° |

| | | | |
|---|---|---|---|
| Wednesday | | 23° | 13° |
| Thursday | 10% | 22° | 13° |
| Friday | 10% | 26° | 21° |
| Saturday | 60% | 30° | 22° |
| Sunday | 20% | 30° | 19° |
| Monday | 10% | 30° | 23° |
| Tuesday | 10% | 33° | 24° |

Siri

Apple

Hi, how can I help?

Get things done, anytime, anywhere

How can I help?

call Sameer

Calling Sameer Samat.
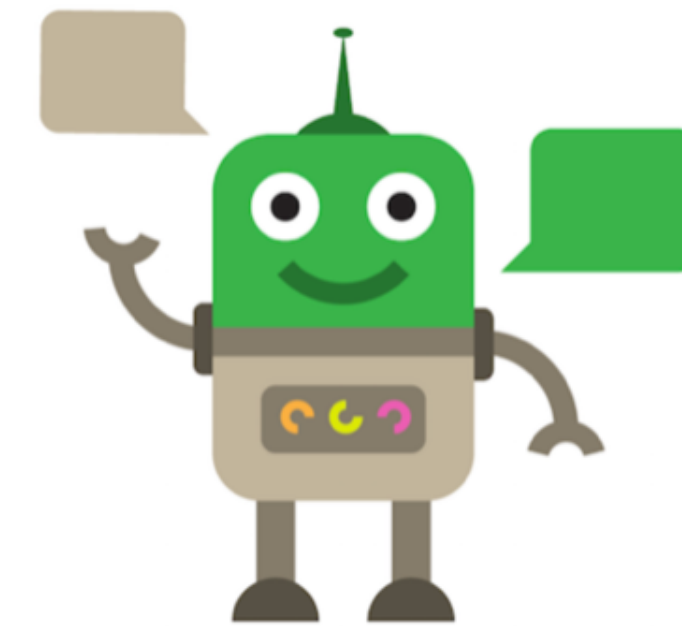
Google

Amazon

amazon

alexa

Microsoft

Hey Cortana

# Two kinds of conversational agents

- **Chatbots**
  - Mimic informal human chatting
  - For fun, or even for therapy
- **Task-oriented dialog agents**
  - Interfaces to personal assistants
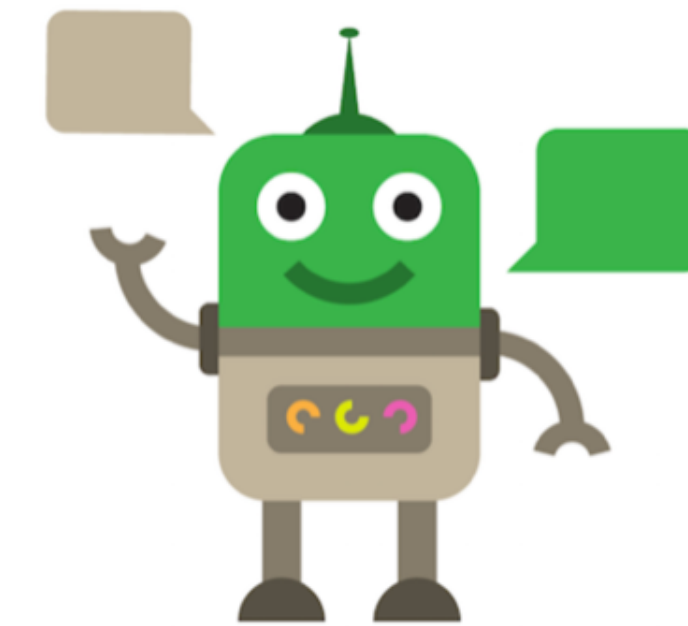  - Cars, robots, appliances
  - Booking flights or restaurants

Dialogue Systems

# Chatbot architectures

- **Rule-based**
  - Pattern-action rules (ELIZA)
  - + A mental model (PARRY):
    - First system to pass the Turing test!
- **Corpus-based (data-driven)**
  - Information Retrieval (XiaoIce)
  - Neural encoder-decoder (BlenderBot)

- **Chatbots**
  - Mimic informal human chatting
  - For entertainment

Dialogue Systems

# Chatbot: Corpus-based methods
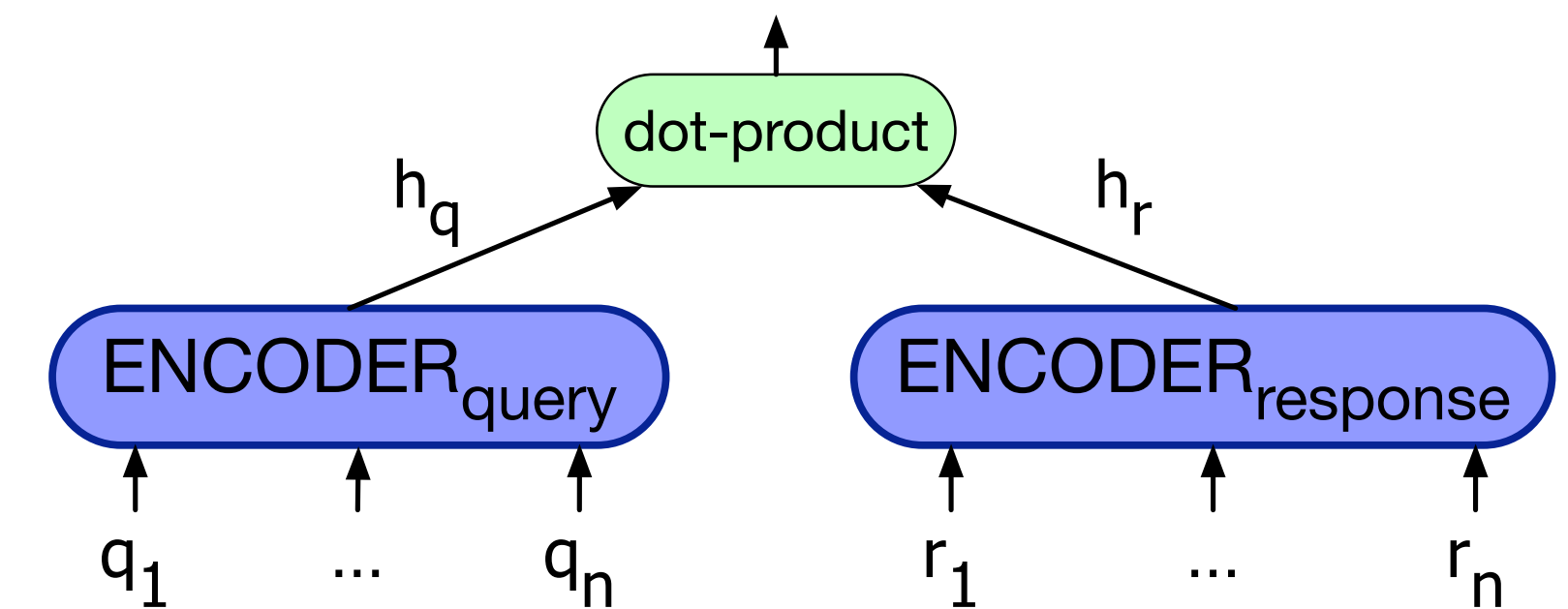
**Response by retrieval**
- Use information retrieval to grab a response (that is appropriate to the context) from some corpus
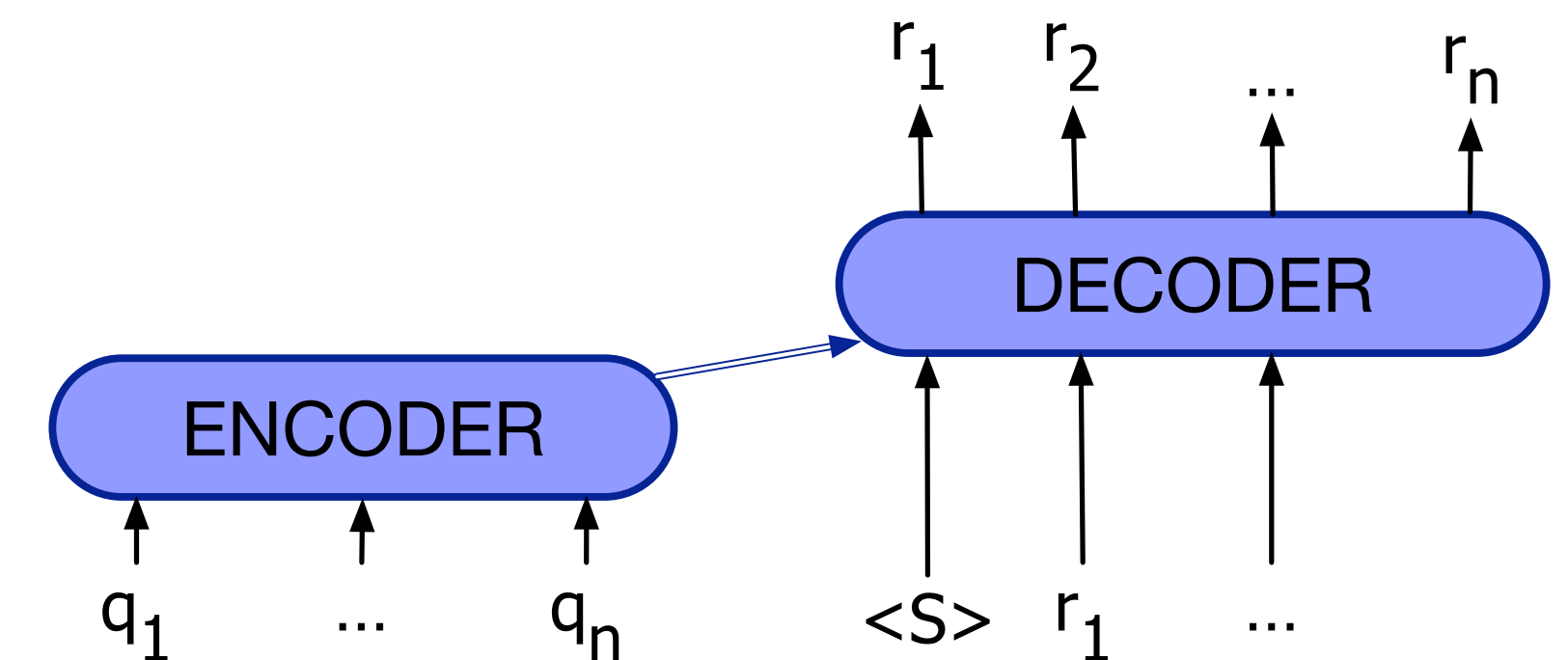
**Response by generation**
- Use a language model or encoder-decoder to generate the response given the dialogue context

Both methods require large corpora to perform well!
- Modern corpus-based
- Typically require hun
  words
- Corpus come from tra
  conversations, movie dialogues, conversations
  social media, crowd-sourced data.



(a) Response by Retrieval

(b) Response by Generation

# Corpus-based method (Response by retrieval)

Return the response to the most similar turn.

**Neural IR method**

Given user query $q$ and a conversational corpus $C$. ←Training corpus
Find the response $r$ in $C$ that has an encoding that is
most similar to the encoding of $q$,.

Bi-encoder: two
separate encoders

$$h_q = \mathbf{BERT}_Q(q)[\mathbf{CLS}]$$
$$h_r = \mathbf{BERT}_R(r)[\mathbf{CLS}]$$

$$\mathbf{response}(q, C) = \arg\max_{r \in C}(h_q \cdot h_r)$$

Similarity
from neural models

Can also have more sophisticated neural architectures

**Chatbot**

# Learning sentence embeddings

- Train sentence embedder on supervised data
- Pooling options:
  - CLS token, MEAN, MAX



- What about unsupervised contrastive learning?



*SimCLR [Chen et al. ICML 2020]*

*Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks,Reimers and Gurevych, EMNLP 2019*

# Learning sentence embeddings

Use dropout for unsupervised contrastive learning of sentence embeddings



*SimCSE: Simple Contrastive Learning of Sentence Embeddings, Gao et al, EMNLP 2021*

# Corpus-based method (Response by generation)

An encoder decoder model for neural response generation in dialogue.

**Neural 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!

Can train just on conversation data or fine-tune a large language model on conversational data

Fine tuning: Chirpy Cardinal System (Paranjape et al., 2020)
- Fine-tunes GPT-2 on EmpatheticDialogues dataset (Rashkin et al, 2019)
- Note that this is a Decoder only model

**Chatbot**

# Chatbot: Seq2Seq models

## Sampling

Randomly sample words from distribution at each time step $t$
- **Basic/pure sampling**: sample from $P_t(w)$ directly
  - Can get some very bad samples
  - No control
- **Top-$n$ sampling**: sample from $P_t$ truncated to top $n$ words
  - Increase $n$ to get more diverse/risky output
  - Decrease $n$ to get more generic/safe output
- **Top-$p$ sampling**: sample from $P_t$ restricted to top $p$ proportion of words
  - Better when probability distribution is spread
- **Temperature based sampling**:
  - Increase $\tau$ to get more diverse/risky output ($P_t$ is more uniform)
  - Decrease $\tau$ to get more generic/safe output ($P_t$ is more spiky)

### Repetitive

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.

### Sample and Rank

1. Sample N candidate
2. Rank candidate and select best one

$$P_t(w) = \frac{\exp(s_w/\tau)}{\sum_{w' \in V} \exp(s_{w'}/\tau)}$$

*(adapted from slides: Stanford CS224N, Chris Manning)*

# Task-based dialogue agents

- **"Task-based"** or **"goal-based"** dialogue agents
  - Systems that have the goal of helping a user solve a task
    - Setting a timer
    - Making a travel reservation
    - Playing a song
    - Buying a product
- Need to incorporate task-specific knowledge
  - **Frames** with **slots** and **values**
  - A knowledge structure representing user **intentions**

# Task-based dialogue agents

- **Frame**
  - Contains a set of **slots**, to be filled with information of a given **type**.
  - Each associated with a **question** to the user

| Slot | Type | Question |
|---|---|---|
| ORIGIN | city | What city are you leaving from? |
| DEST | city | Where are you going? |
| DEP DATE | date | What day would you like to leave? |
| DEP TIME | time | What time would you like to leave? |
| AIRLINE | Airline | What is your preferred airline? |

# How to detect frames and fill in dialog slots?

Natural language understanding

"Show me morning flights from
Boston to San Francisco on Tuesday"

Before filling in the dialog slots:

- System must detect which slot of which frame user is filling
- And switch dialogue control to that frame.

# How to detect frames and fill in dialog slots?

Natural language understanding

"Show me morning flights from
Boston to San Francisco on Tuesday"

**Step#1: domain classification**     DOMAIN: AIR-TRAVEL     **Classification**

# How to detect frames and fill in dialog slots?

Natural language understanding

"Show me morning flights from
Boston to San Francisco on Tuesday"

Identify the frame to use

**Step#1: domain classification**          DOMAIN: AIR-TRAVEL

**Step#2: intent determination**           INTENT: SHOW-FLIGHTS          Classification

# How to detect frames and fill in dialog slots?

Natural language understanding

"Show me morning flights from
Boston to San Francisco on Tuesday"

**Step#1: domain classification**

DOMAIN: AIR-TRAVEL

**Step#2: intent determination**

INTENT: SHOW-FLIGHTS

**Step#3: slot filling**

ORIGIN-CITY: Boston
ORIGIN-DATE: Tuesday
ORIGIN-TIME: morning
DEST-CITY: San Francisco

**Rule-based**

or

**Sequence tagging**

**"Belief-State"**

# Dialogue-State Architecture



Knowledge Base (KB)

Dialogue Management

Dialog State Tracker (DST)

Dialog Policy

Natural language understanding (NLU)

Natural Language Generation (NLG)

Automatic Speech Recognition (ASR)

Text to Speech (TTS)

## Knowledge Base (KB)

## Dialogue Management

### Dialog State Tracker (DST)

### Dialog Policy

| { from: downtown } | 0.5 |
| { depart-time: 1300 } | 0.3 |
| { arrive-time: 1300 } | 0.1 |

## Natural language understanding (NLU)

**NLU component:** to identify **domain**, **intent**, and extract **slot fillers** from the user's utterance

## Natural Language Generation (NLG)

| LEAVING FROM DOWNTOWN | 0.6 |
| LEAVING AT ONE P M | 0.2 |
| ARRIVING AT ONE P M | 0.1 |

## Automatic Speech Recognition (ASR)

## Text to Speech (TTS)

```
from:       downtown
to:         airport
depart-time: --
confirmed:   no
score:       0.65
score:       0.15
score:       0.10
```

Knowledge Base (KB)

Dialogue Management

Dialog State Tracker (DST)

Dialog Policy

Natural language understanding (NLU)

**Dialogue state tracker:** maintains the current state of the dialogue (most recent dialogue act + agenda)

Natural Language Generation (NLG)

Automatic Speech Recognition (ASR)

Text to Speech (TTS)

```
from:         downtown
to:           airport
depart-time:  --
confirmed:    no
score:        0.65
score:        0.15
score:        0.10
```

## Knowledge Base (KB)

```
{ act:  confirm  }
{ from: downtown }
```

## Dialogue Management

### Dialog State Tracker (DST)

### Dialog Policy

**Dialogue policy:** decides what the system should do or say
The topic of next (intent-level)

### Natural language understanding (NLU)

### Natural Language Generation (NLG)

### Automatic Speech Recognition (ASR)

### Text to Speech (TTS)

Knowledge Base (KB)

Dialogue Management

Dialog State Tracker (DST)

Dialog Policy

```
{ act:  confirm
  from: downtown }
```

Natural language understanding (NLU)

**NLG component:** decides the actual text string to generate (surface realization)

Templates or NNs

Natural Language Generation (NLG)

Automatic Speech Recognition (ASR)

```
FROM DOWNTOWN,
IS THAT RIGHT?
```

Text to Speech (TTS)

# Dialogue System Evaluation

**Chatbots:**
- Mostly human evaluation
- Automatic evaluation is a challenge
  - Use automatic metrics to capture specific aspect (diversity, length of conversation)
  - Adversarial evaluation

**Task-based dialogue agents**
- Automatic metrics to evaluate task performance
- Can also use human evaluation

# Automatic evaluation metrics

**Content overlap metrics:**
- Word (n-gram) overlap: BLEU, ROUGE, METEOR, CIDEr
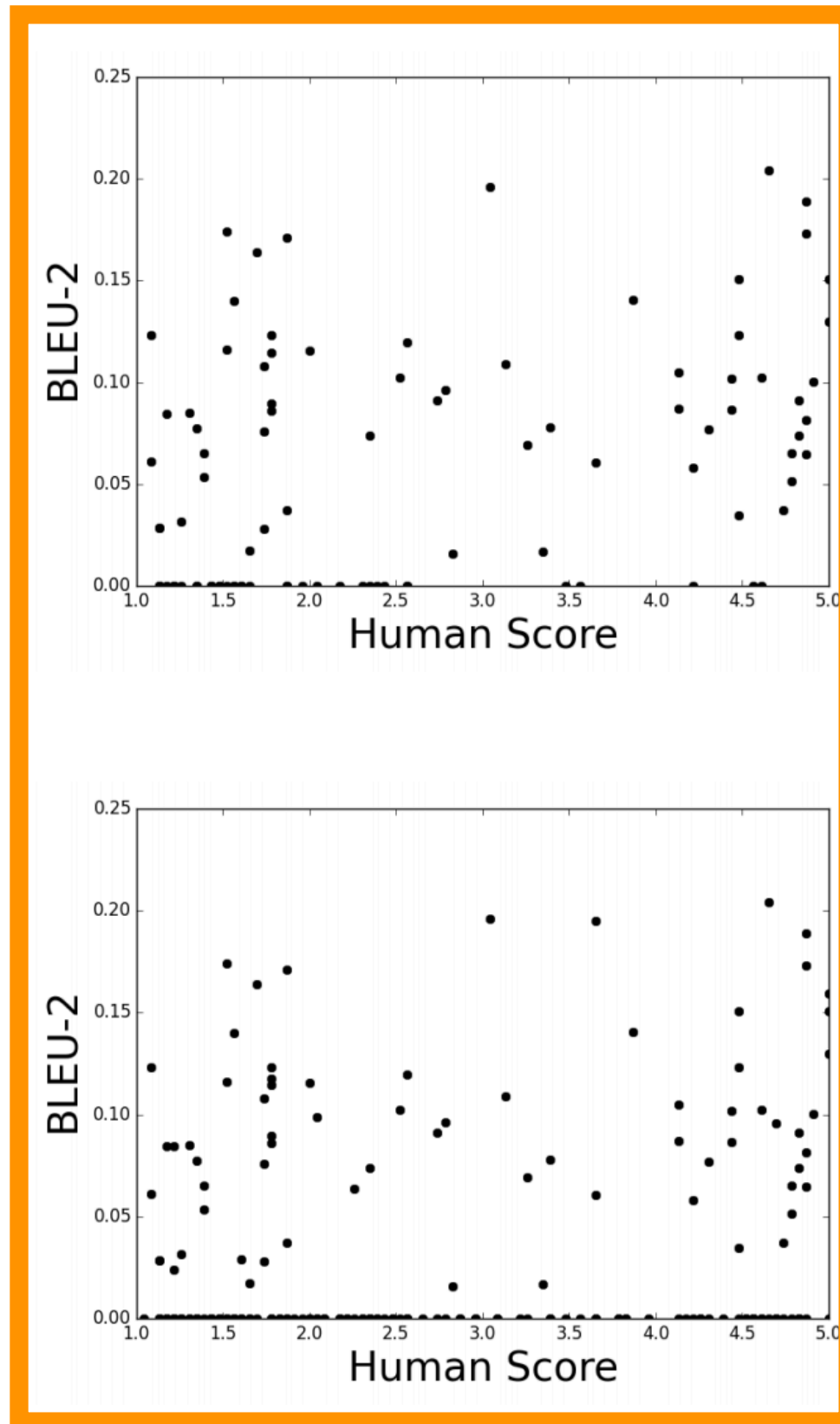- Structured overlap: PYRAMID, SPICE, SPIDER

**Model based metrics:**
- Embedding similarity: Embedding average, Word Mover Distance, BERTSCORE, etc.
- Metric predictor: BLEURT

# Issues with Automatic Evaluation

**Automatic Evaluation:**
Word overlap metrics
are bad for dialogue

**No correlation**
between **human
judgement** and **BLEU**



(a) Twitter

(b) Ubuntu

**BLEU**　　　　　　**Embedding Average**　　　　　　**Human**

*[How NOT To Evaluate Your Dialogue System: An Empirical Study of Unsupervised Evaluation Metrics for Dialogue Response Generation, Liu et al 2017, https://arxiv.org/pdf/1603.08023.pdf]*
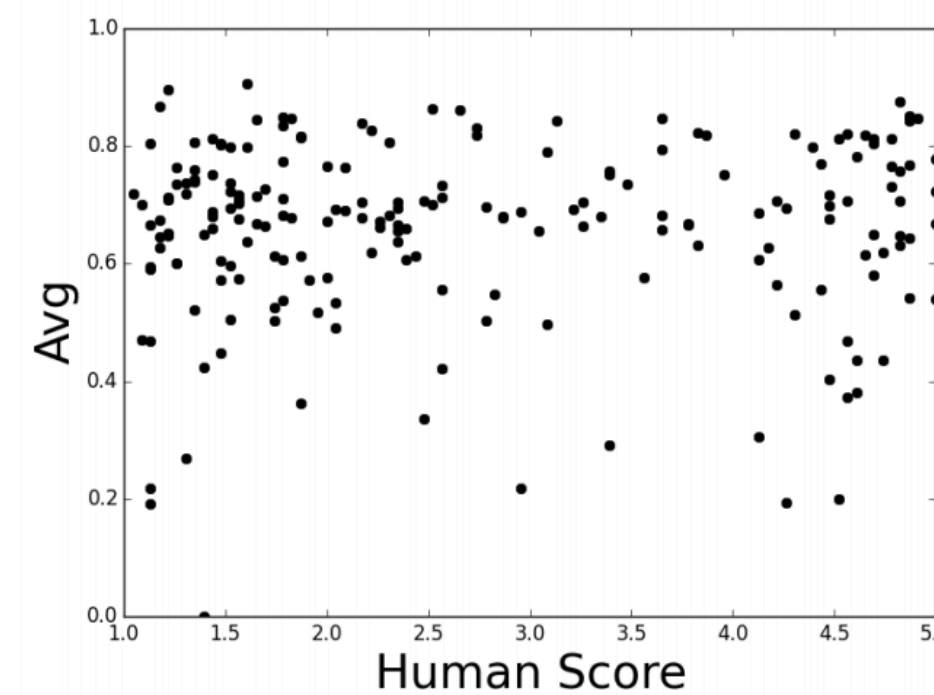
# Issues with Automatic Evaluation

**Automatic Evaluation:**
Embedding metrics are also poor for dialogue

**No correlation** between **human judgement** and **embedding average**



(a) Twitter

(b) Ubuntu

**BLEU**  **Embedding Average**  **Human**

*[How NOT To Evaluate Your Dialogue System: An Empirical Study of Unsupervised Evaluation Metrics for Dialogue Response Generation, Liu et al 2017, https://arxiv.org/pdf/1603.08023.pdf]*

26

# Issues with Automatic Evaluation

**Word Based Metrics**



**Word Overlap Metrics**
- highly correlated with each other
- Not so correlated with human ratings

Spearman correlations of word based metrics and human ratings

**Human Ratings**
- Informativeness
- Naturalness
- Quality

*[Why We Need New Evaluation Metrics for NLG, Novikova et al 2017, https://arxiv.org/pdf/1707.06875.pdf]*

# Issues with automatic Evaluation

High correlation with human judgement for low quality generations

**Bad informativeness**

| | TER | B1 | B2 | B3 | B4 | RG | NST | LP | CID | MET | SIM |
|---|---|---|---|---|---|---|---|---|---|---|---|
| INF | | | | | | | | | | | |
| NAT | | | | | | | | | | | |
| QUA | | | | | | | | | | | |

-1   -0.6   -0.2   0.2   0.6   1

**Average and Good informativeness**

| | TER | B1 | B2 | B3 | B4 | RG | NST | LP | CID | MET | SIM |
|---|---|---|---|---|---|---|---|---|---|---|---|
| INF | | | | | | | | | | | |
| NAT | | | | | | | | | | | |
| QUA | | | | | | | | | | | |

-1   -0.6   -0.2   0.2   0.6   1

Poor correlation with human judgement for mid to high quality generations



*[Why We Need New Evaluation Metrics for NLG, Novikova et al 2017, https://arxiv.org/pdf/1707.06875.pdf]*

Re-evaluating Automatic Metrics for Image Captioning
[Kilickaya et al, EACL 2017]

# Human evaluation

What kind of human evaluation can be done?
- Can get ratings from chat **participants** or external **observers**.
- Can ask humans to rate various aspects of the chat (**likert scale**) or to compare two chats / responses (**AB testing**)
- Dimensions: fluency, coherence / consistency, factuality and correctness, commonsense, style / formality, grammaticality, typicality, redundancy

Issues with human evaluation
- slow, expensive
- not repeatable (subjective/inconsistent)
- difficult to form well-targeted questions that are not open to misinterpretation

When developing new automatic metrics, human evaluation is used as gold
- New automated metrics must correlate well with human evaluation.

# Evaluation takeaways

- **Content overlap metrics** provide a good starting point for evaluating the quality of generated text, but they're not good enough on their own.

- **Model-based metrics** can be more correlated with human judgment, but metric may not be not interpretable

- **Human judgments** are critical

  - But humans are inconsistent and judgments are expensive

- If you are developing a NLG system, you should

  - Look at your model generations. Don't just rely on numbers!

  - Publicly release large samples of the output of systems that you create!

*Slide adapted from Stanford CS224N  (Xiang Lisa Li, Antoine Bosselut, Chris Manning)*

# Large language models for dialogue

# ChatGPT can do cool stuff

**juan** ✓
@juanbuis

having a particularly bizarre morning thanks to chatgpt

> explain quantum theory to a child in the style of snoop dogg

Yo, little homie, let me break it down for you
Quantum theory, it's a crazy thing, it's true
See, the world is made up of tiny, tiny particles
That are always moving, and changing, and actin' wild

These particles, they don't always behave like we think
Sometimes they're waves, and sometimes they're points, it's a trip
They can be in two places at once, and they can be in two states
It's a mind-blowing concept, but that's the quantum fate

So don't be scared, little homie, it's all good
Quantum theory, it's just the way the world works, it's understood
It's a wild ride, but it's a cool one too
Quantum theory, it's the way the world does what it do.

ALT

5:33 AM · Dec 1, 2022

https://twitter.com/juanbuis/status/1598263964093251589?lang=en

32

# ChatGPT

https://openai.com/blog/chatgpt/

**Question**:

Recent models such as BERT (encoder-only) are applied mostly on classification tasks, whereas models like GPT1(decoder-only) are applied on generative modeling of text (including machine translation). Explain the justification of using these two different configurations of the transformer for the two tasks.

**GPT**:

The reason for using encoder-only models like BERT for classification tasks is that these tasks typically involve taking in an input sequence and outputting a single label or class for that sequence. In this case, the encoder self-attention blocks in BERT are able to capture the important relationships between the input elements and create a condensed representation of the input sequence, which is then passed to a classifier to make a prediction.

On the other hand, generative modeling tasks like machine translation require generating a sequence as output, which is why decoder-only models like GPT-1 are used. In these models, the decoder self-attention blocks are able to "pay attention" to the different elements of the encoded input sequence and previous decoder states in order to generate the next element in the output sequence, allowing the model to produce fluent and coherent text.

(from student on https://coursys.sfu.ca/2022fa-cmpt-413-x1/forum/386)

# ChatGPT

- ChatGPT: released Nov 30th 2022, 1M users in 5 days, 100M users in 2 months

- Large language model (GPT-3.5)
  - GPT-3 (2020) - 175B parameter model
  - GPT-3.5 (late 2022) - trained on a blend of text and code
  - GPT-4 (March 2023) - multimodal

- Supervised fine-tuning on human conversations
  - Data where human will pretend to be user or AI assistant

- Reinforcement learning with human feedback
  - Humans rank what response is best
  - Aim: reduce harmful / deceitful responses

# GPT-4



**Internal factual eval by category**

Accuracy

- Growing performance for ChatGPT versions

https://openai.com/research/gpt-4

# GPT-4

- GPT-4 passing standardized tests

- Bar exam:
  - GPT-3.5 score in bottom 10%
  - GPT-4 score in top 10%

**Exam results (ordered by GPT-3.5 performance)**

Estimated percentile lower bound (among test takers)

Legend:
- gpt-4 (dark green)
- gpt-4 (no vision) (light green)
- gpt3.5 (blue)



https://openai.com/research/gpt-4

# Pre-trained GPT-4 is well-calibrated

Calibrated: predicted confidence matches probability of being correct

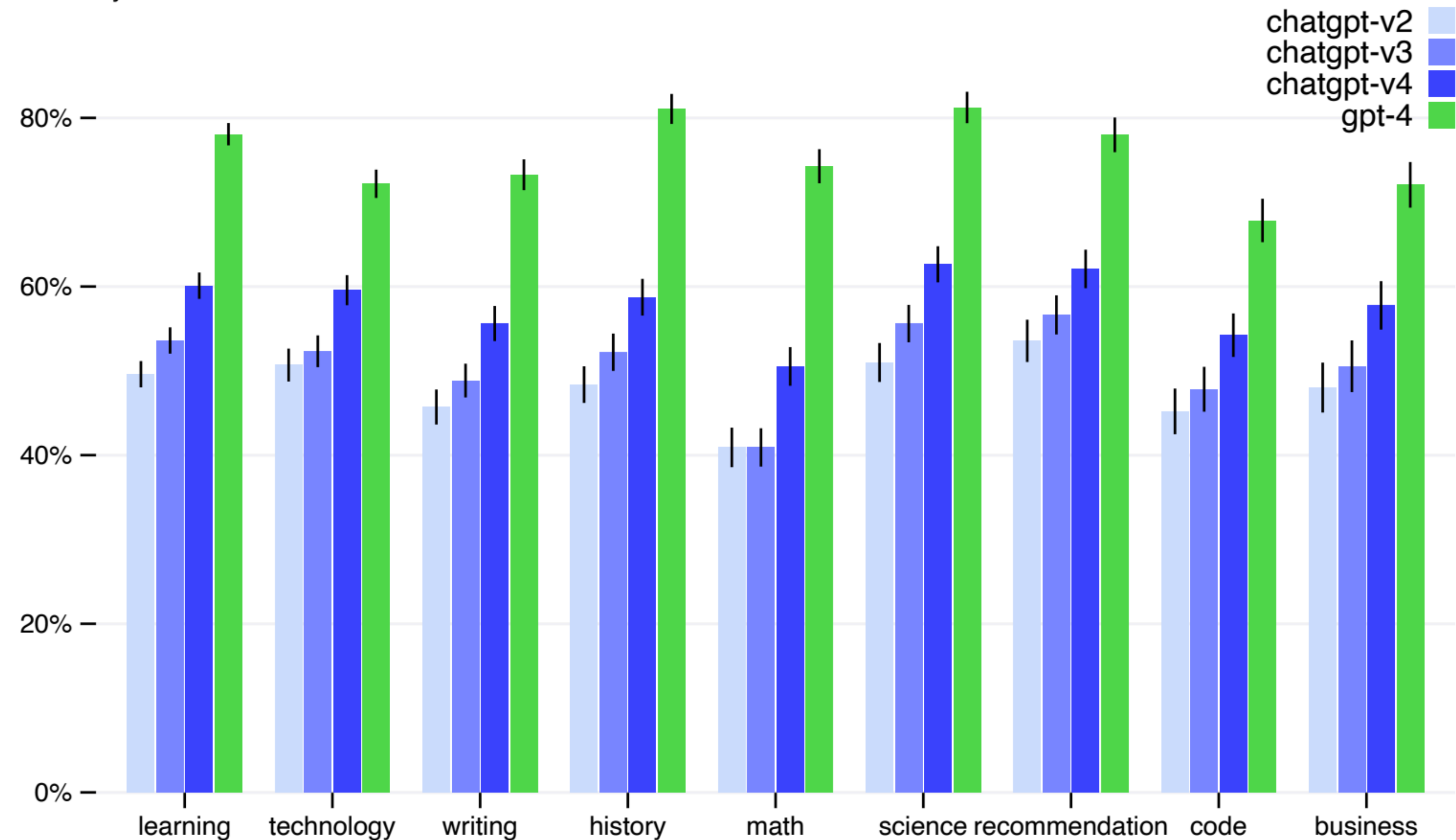

https://openai.com/research/gpt-4

# ChatGPT

- ChatGPT: released Nov 30th 2022, 1M users in 5 days, 100M users in 2 months

- Large language model (GPT-3.5)
  - GPT-3 (2020) - 175B parameter model
  - GPT-3.5 (late 2022) - trained on a blend of text and code
  - GPT-4 (March 2023) - multimodal

- **Supervised fine-tuning on human conversations**
  - **Data where human will pretend to be user or AI assistant**

- Reinforcement learning with human feedback
  - Humans rank what response is best
  - Aim: reduce harmful / deceitful responses

# InstructGPT

Collect data from humans

Example user prompts

| open qa | Q: Who is Batman?<br>A: Batman is a fictional comic book character.<br><br>Q: What is torsalplexity?<br>A: ?<br><br>Q: What is Devz9?<br>A: ?<br><br>Q: Who is George Lucas?<br>A: George Lucas is American film director and producer famous for creating Star Wars.<br><br>Q: What is the capital of California?<br>A: |
|---|---|
| chat | The following is a conversation with an AI assistant. The assistant is helpful, creative, clever, and very friendly.<br><br>Human: Hello, who are you?<br>AI: I am an AI created by OpenAI. How can I help you today?<br>Human: I'm feeling kind of down today.<br>AI: |

# Instruction tuning



- Use templates to make them into instruction based dataset
- Text based format makes it natural for humans



**QQP (Paraphrase)**

| Question1 | How is air traffic controlled? |
| Question2 | How do you become an air traffic controller? |
| Label | 0 |

{Question1} {Question2} Pick one: These questions are duplicates or not duplicates.

I received the questions "{Question1}" and "{Question2}". Are they duplicates?

{Choices[label]}

{Choices[label]}

**XSum (Summary)**

| Document | The picture appeared on the wall of a Poundland store on Whymark Avenue... |
| Summary | Graffiti artist Banksy is believed to be behind... |

{Document} How would you rephrase that in a few words?

First, please read the article: {Document} Now, can you write me an extremely short abstract for it?

{Summary}

{Summary}

Multitask prompted training enables zero-shot task generalization, Sahn et al. Google, ICLR 2022

# Instruction tuning

## (A) Pretrain–finetune (BERT, T5)

```
[Pretrained LM] → [Finetune on task A] → [Inference on task A]
```

- Typically requires many task-specific examples
- One specialized model for each task

## (B) Prompting (GPT-3)

```
[Pretrained LM] → [Inference on task A]
```

Improve performance via few-shot prompting or prompt engineering

## (C) Instruction tuning (FLAN)

```
[Pretrained LM] → [Instruction-tune on many tasks: B, C, D, …] → [Inference on task A]
```

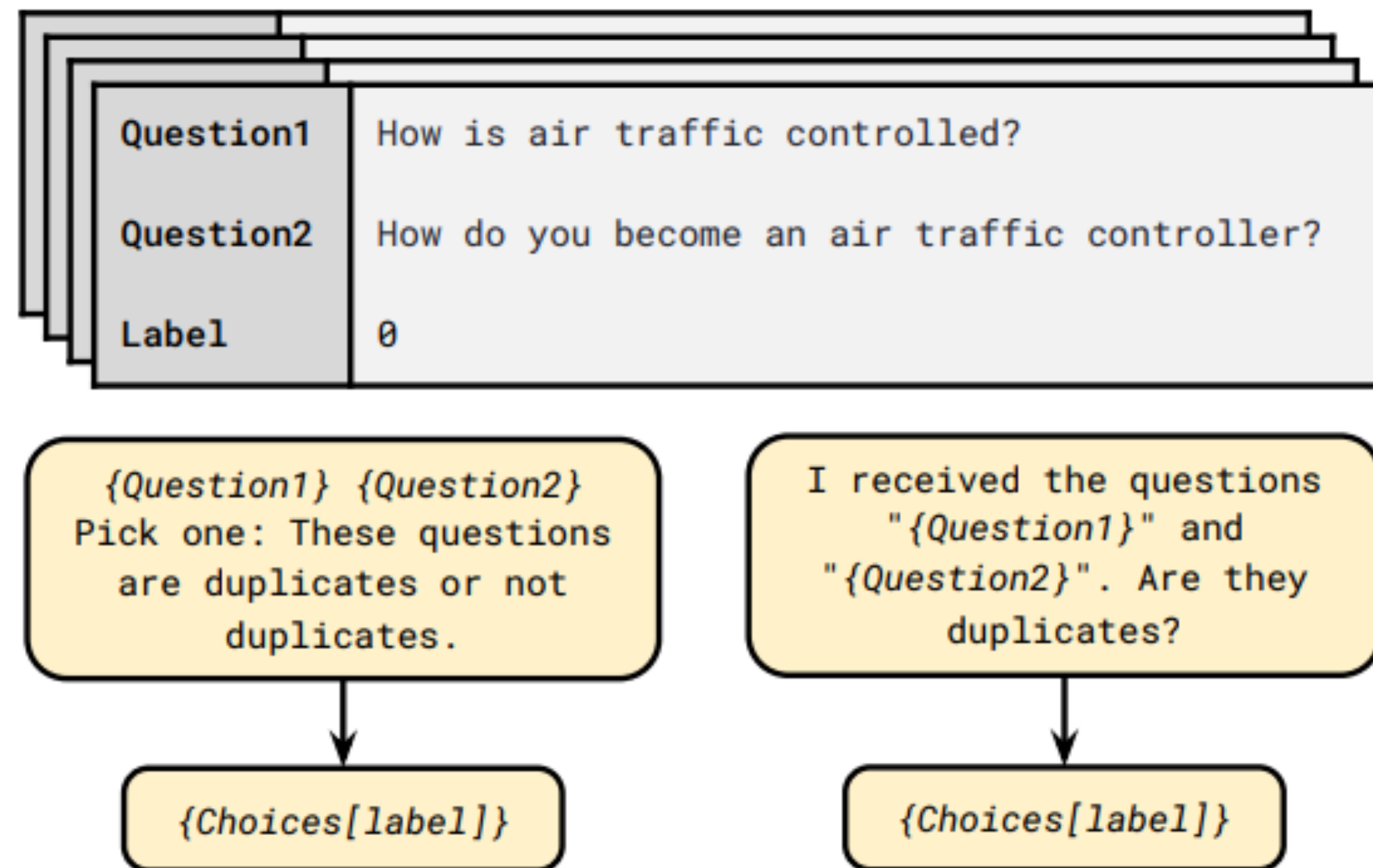Model learns to perform many tasks via natural language instructions

Inference on unseen task

# Instruction tuning
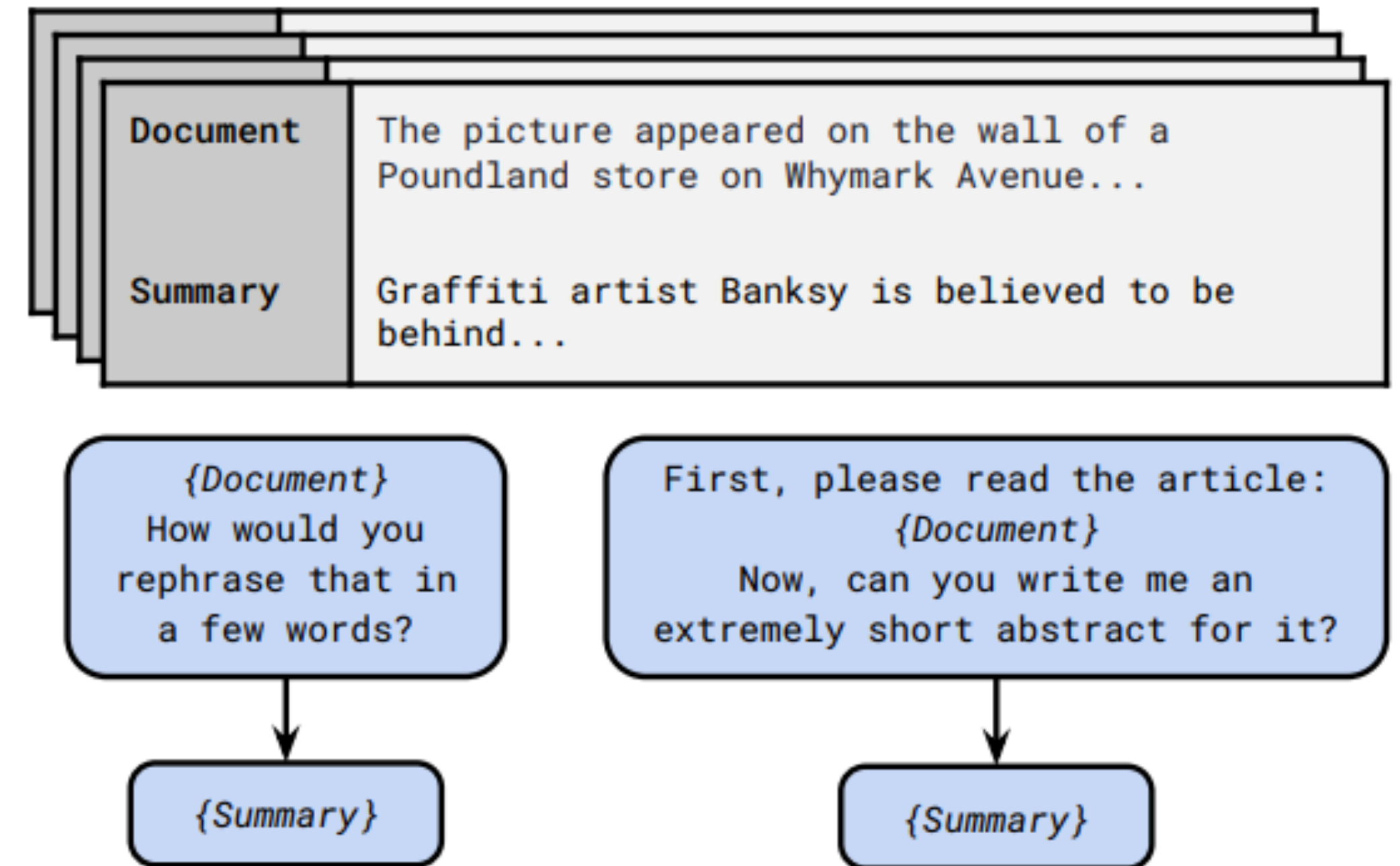
- Use templates to make them into instruction based dataset
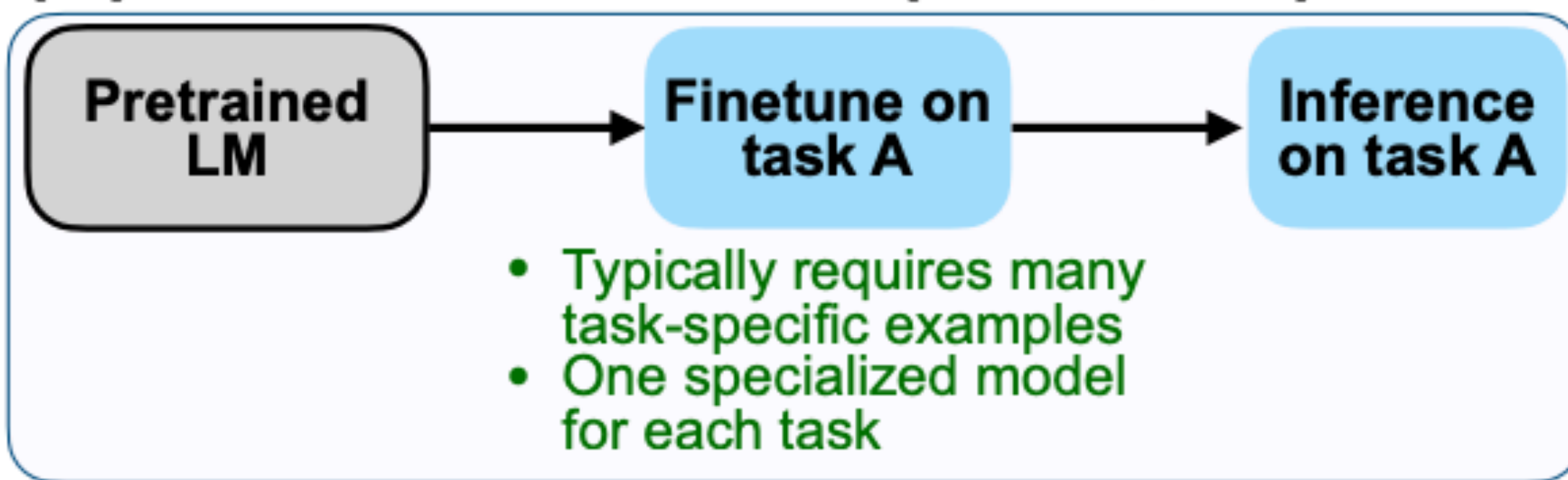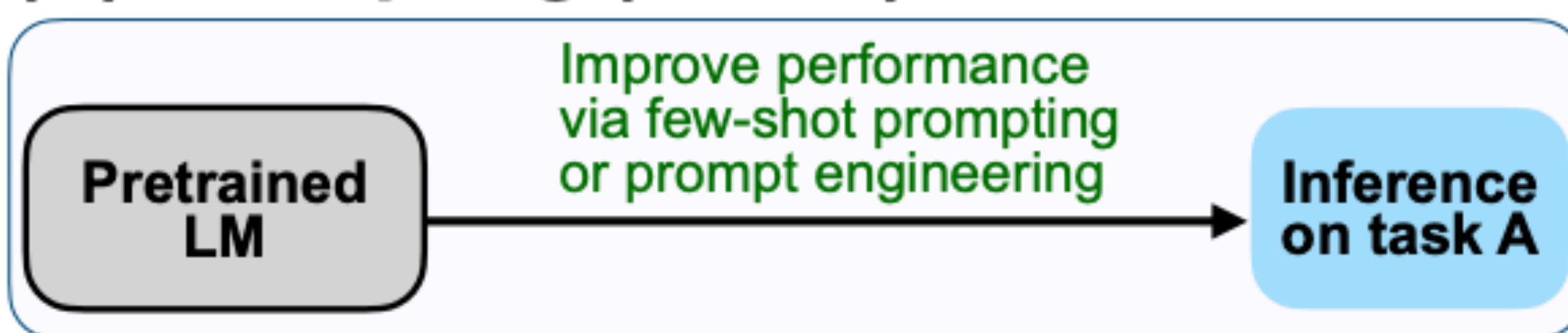- Text based format makes it natural for humans

**Premise**

Russian cosmonaut Valery Polyakov set the record for the longest continuous amount of time spent in space, a staggering 438 days, between 1994 and 1995.

**Hypothesis**

Russians hold the record for the longest stay in space.

**Target**

Entailment
Not entailment

Options:
- yes
- no

**Template 1**

\<premise\>

Based on the paragraph above, can we conclude that \<hypothesis\>?

\<options\>

**Template 2**

\<premise\>

Can we infer the following?

\<hypothesis\>

\<options\>

**Template 3**

Read the following and determine if the hypothesis can be inferred from the premise:

Premise: \<premise\>

Hypothesis: \<hypothesis\>

\<options\>

**Template 4, …**

# Instruction tuning

- Can be used on an unseen task type

**Finetune on many tasks ("instruction-tuning")**

**Input (Commonsense Reasoning)**

Here is a goal: Get a cool sleep on summer days.
How would you accomplish this goal?
OPTIONS:
-Keep stack of pillow cases in fridge.
-Keep stack of pillow cases in oven.

**Target**

keep stack of pillow cases in fridge

**Input (Translation)**

Translate this sentence to Spanish:
The new office building was built in less than three months.

**Target**

El nuevo edificio de oficinas se construyó en tres meses.

Sentiment analysis tasks

Coreference resolution tasks

...

**Inference on unseen task type**

**Input (Natural Language Inference)**

Premise: At my age you will probably have learnt one lesson.

Hypothesis: It's not certain how many lessons you'll learn by your thirties.

Does the premise entail the hypothesis?
OPTIONS:
-yes    -it is not possible to tell    -no

**FLAN Response**

It is not possible to tell
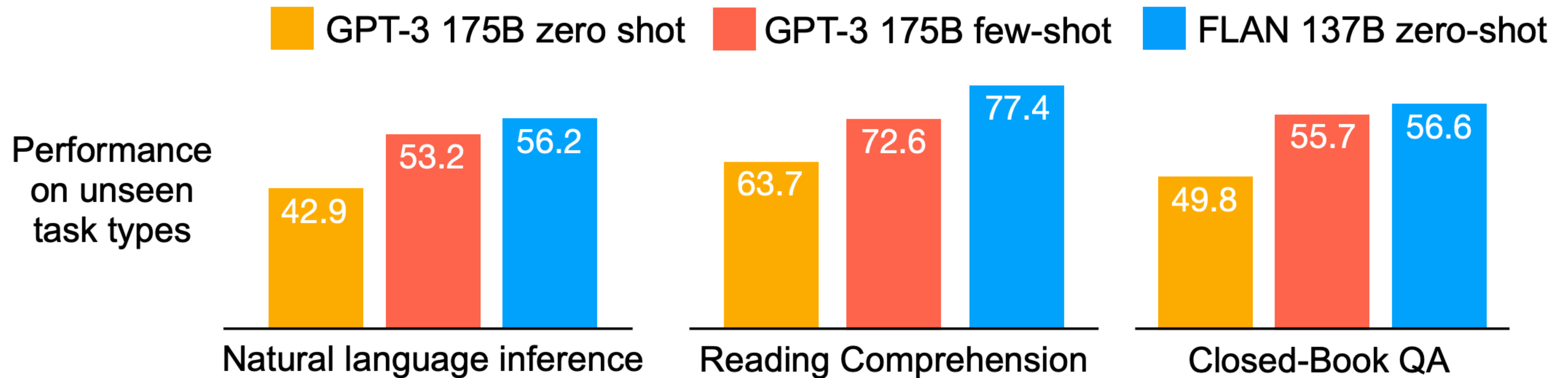
43

# Instruction tuning
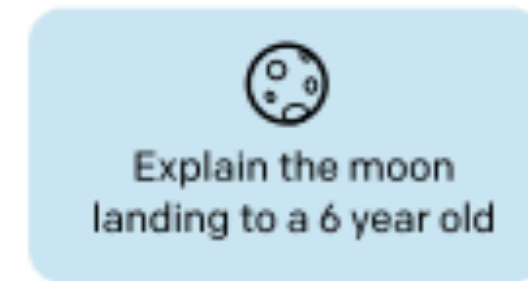
- Can be used on an unseen task type

# ChatGPT

- ChatGPT: released Nov 30th 2022, 1M users in 5 days, 100M users in 2 months

- Large language model (GPT-3.5)
  - GPT-3 (2020) - 175B parameter model
  - GPT-3.5 (late 2022) - trained on a blend of text and code
  - GPT-4 (March 2023) - multimodal

- Supervised fine-tuning on human conversations
  - Data where human will pretend to be user or AI assistant

- **Reinforcement learning with human feedback**
  - **Humans rank what response is best**
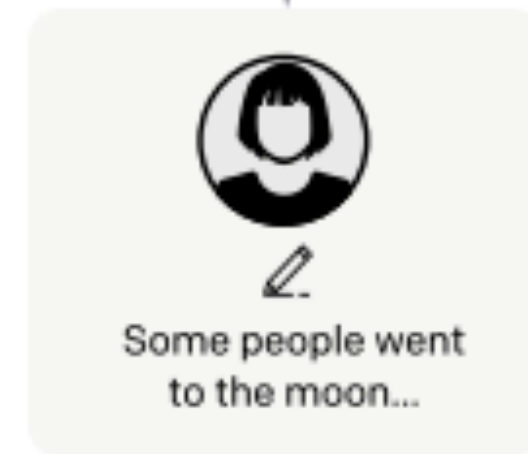  - **Aim: reduce harmful / deceitful responses**

# InstructGPT



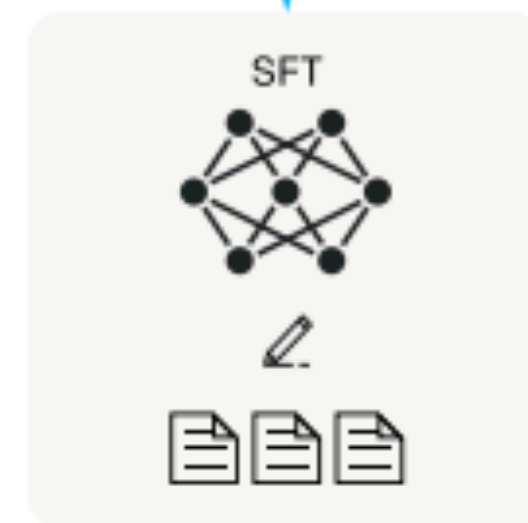**Step 1**

**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.
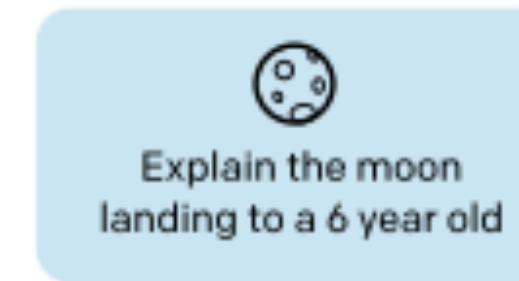
Some people went to the moon...

This data is used to fine-tune GPT-3 with supervised learning.

SFT

**Step 2**

**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A — Explain gravity...
B — Explain war...
C — Moon is natural satellite of...
D — People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

**Step 3**

**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

Training language models to follow instructions with human feedback, Ouyang et. al. OpenAI 2022

# InstructGPT

Collect human judgement of which responses are better
- Labelers rank K (K=4 to 9) responses for each prompt
- Gives $\binom{K}{2}$ comparisons per prompt

Train reward model (RM) $r_\theta$ to over human data *D*:

$$\text{loss}\,(\theta) = -\frac{1}{\binom{K}{2}} E_{(x,y_w,y_l)\sim D}\left[\log\left(\sigma\left(r_\theta\left(x,y_w\right) - r_\theta\left(x,y_l\right)\right)\right)\right]$$

Use reinforcement learning (RL) to train RL based policy for selecting words to generate

$$\text{objective}\,(\phi) = E_{(x,y)\sim D_{\pi_\phi^{\text{RL}}}}\left[r_\theta(x,y) - \beta\log\left(\pi_\phi^{\text{RL}}(y\mid x)/\pi^{\text{SFT}}(y\mid x)\right)\right] +$$

$$\gamma E_{x\sim D_{\text{pretrain}}}\left[\log(\pi_\phi^{\text{RL}}(x))\right]$$

# InstructGPT

Aligns model behaviour to human preferences

Human judgement of model output



Training language models to follow instructions with human feedback, Ouyang et. al. OpenAI 2022

https://openai.com/research/instruction-following

# InstructGPT

Aligns model behaviour to human preferences

Human judgement of model output



Training language models to follow instructions with human feedback, Ouyang et. al. OpenAI 2022

https://openai.com/research/instruction-following

# InstructGPT

Aligns model behaviour to human preferences

Human judgement of model output



Training language models to follow instructions with human feedback, Ouyang et. al. OpenAI 2022

https://openai.com/research/instruction-following

# Self-Instruct

- Generate task instructions using LLMs to train/fine-tune LLMs!



**Task Pool**

**Step 1: Instruction Generation**

**LM**

**Task**
**Instruction :** Give me a quote from a famous person on this topic.

175 seed tasks with 1 instruction and 1 instance per task

**Step 2: Classification Task Identification**

**LM**

**Step 3: Instance Generation**

**Task**
**Instruction :** Find out if the given text is in favor of or against abortion.

**Class Label:** Pro-abortion
**Input:** Text: I believe that women should have the right to choose whether or not they want to have an abortion.

Yes
Output-first

**LM**

**Task**
**Instruction :** Give me a quote from a famous person on this topic.

**Input:** Topic: The importance of being honest.
**Output:** "Honesty is the first chapter in the book of wisdom." - Thomas Jefferson

No
Input-first

**Step 4: Filtering**

SELF-INSTRUCT: Aligning Language Model with Self Generated Instructions, Wang et al. 2022

51

# Self-Instruct

- Generate task instructions using LLMs to train/fine-tune LLMs!



SELF-INSTRUCT: Aligning Language Model with Self Generated Instructions, Wang et al. 2022

# Self-Instruct



Text-davinci-003

175 Self-Instruct seed tasks

Modified Self-instruct Instruction Generation

LLaMA 7B

52K Instruction-following examples

Supervised Finetuning

Alpaca 7B

**Example seed task**

*Instruction: Brainstorm a list of possible New Year's resolutions.*

*Output:*
*- Lose weight*
*- Exercise more*
*- Eat healthier*

**Example Generated task**

*Instruction: Brainstorm creative ideas for designing a conference room.*

*Output:*
*... incorporating flexible components, such as moveable walls and furniture ...*