# Understanding ChatGPT: Technology, Trends and Challenges for Conversational Systems

**Luis Fernando D'Haro - @lfdharo**

**Fall School RTTH – Jaca, Nov 14-17, 2023**

**Grupo de Tecnología del Habla y Aprendizaje Automático - ETSI de Telecomunicación – Universidad Politécnica de Madrid**
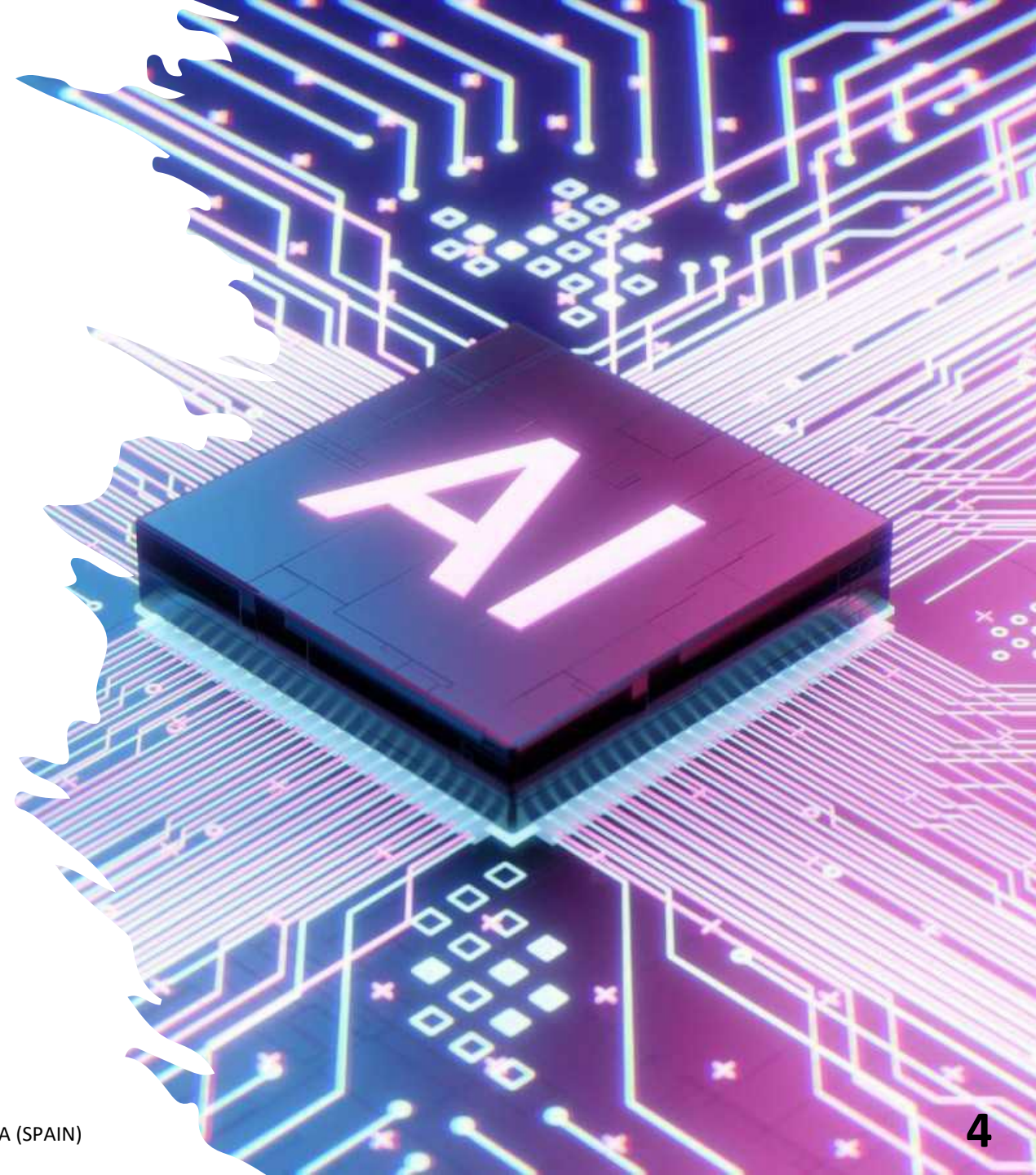
# Content

Introduction

# Definition of Artificial Intelligence

- "A field that involves machines capable of **mimicking** certain functionalities of human intelligence, including features such as perception, learning, reasoning, problem-solving, linguistic interaction, and even the production of creative works." (UNESCO)

# Types of Artificial Intelligence (Capacity)

## Artificial Narrow Intelligence (ANI)

Artificial narrow intelligence (ANI), also called "weak AI". This type of artificial intelligence works in a limited context and is a simulation of human intelligence. Narrow AI focuses on performing specific tasks but has some limitations. Examples of Narrow AI: Siri (from Apple), Alexa (from Amazon), Google Search.

## Artificial General Intelligence (AGI)

Artificial general intelligence (AGI), also called "strong AI". It is the kind of artificial intelligence that allows the understanding and learning of any task related to the intellect that a human being can perform. AGI can think, understand and act in a way that is indistinguishable from the human in certain situations.

## Artificial Super Intelligence (ASI)

Artificial Super Intelligence (ASI) commonly known as super-intelligence is a thing of the future. Currently, it is a hypothetical concept believed that artificial intelligence that not only understands and imitates human behavior, but machines become self-aware.

# TYPES OF ARTIFICIAL INTELLIGENCE BASED ON FUNCTIONALITIES

## SELF - AWARENESS

Self-awareness in AI refers to the system's awareness of its own personality or individuality and is both exciting and terrifying, with the potential for super-intelligent systems with human-like consciousness to be developed in the future with advancements in AI technology.

## REACTIVE MACHINES

Reactive machines are limited AI systems that make decisions based on the present moment and are reliable in completing specific tasks, but lack interaction, emotion, or consciousness and can be easily tricked.

## THEORY OF MIND

Theory of Mind (ToM) AI machines will focus on understanding human thoughts and emotions and be capable of meaningful interaction, and treat people with respect and care, but it is not yet fully developed.

## LIMITED MEMORY

Limited memory machines have short-term memory that allows them to temporarily store acquired experiences, create a memory, and take action based on it, making them more advanced than reactive machines, but the information gathered from previous experiences is not stored in the machine's content library, making it transient.

**6**

# Levels of AGI

- Morris, M. R., Sohl-dickstein, J., Fiedel, N., Warkentin, T., Dafoe, A., Faust, A., ... & Legg, S. (2023). Levels of AGI: Operationalizing Progress on the Path to AGI. *arXiv preprint arXiv:2311.02462*.

| Generality (columns) | *clearly scoped task or set of tasks* | *wide range of non-physical tasks, including metacognitive abilities like learning new skills* |
|---|---|---|
| **Level 0: No AI** | **Narrow Non-AI** calculator software; compiler | **General Non-AI** human-in-the-loop computing, e.g., Amazon Mechanical Turk |
| **Level 1: Emerging** *equal to or somewhat better than an unskilled human* | **Emerging Narrow AI** GOFAI[4]; simple rule-based systems, e.g., SHRDLU (Winograd, 1971) | **Emerging AGI** ChatGPT (OpenAI, 2023), Bard (Anil et al., 2023), Llama 2 (Touvron et al., 2023) |
| **Level 2: Competent** *at least 50th percentile of skilled adults* | **Competent Narrow AI** toxicity detectors such as Jigsaw (Das et al., 2022); Smart Speakers such as Siri (Apple), Alexa (Amazon), or Google Assistant (Google); VQA systems such as PaLI (Chen et al., 2023); Watson (IBM); SOTA LLMs for a subset of tasks (e.g., short essay writing, simple coding) | **Competent AGI** not yet achieved |
| **Level 3: Expert** *at least 90th percentile of skilled adults* | **Expert Narrow AI** spelling & grammar checkers such as Grammarly (Grammarly, 2023); generative image models such as Imagen (Saharia et al., 2022) or Dall-E 2 (Ramesh et al., 2022) | **Expert AGI** not yet achieved |
| **Level 4: Virtuoso** *at least 99th percentile of skilled adults* | **Virtuoso Narrow AI** Deep Blue (Campbell et al., 2002), AlphaGo (Silver et al., 2016, 2017) | **Virtuoso AGI** not yet achieved |
| **Level 5: Superhuman** *outperforms 100% of humans* | **Superhuman Narrow AI** AlphaFold (Jumper et al., 2021; Varadi et al., 2021), AlphaZero (Silver et al., 2018), StockFish | **Artificial Superintelligence (ASI)** not yet achieved |

# Aplicaciones

- Automatización de tareas
- Asistentes virtuales inteligentes
- Mejora en la atención médica
- Conducción autónoma
- Personalización en servicios y recomendaciones
- Optimización de procesos industriales

# Dialogue System (DS)



□ Dialogue systems are intelligent agents that are able to help users finish tasks more efficiently via <u>multimodal interactions</u> (mainly image, text and speech).

□ Dialogue systems are being incorporated into various devices (smart-phones, Apps, smart TVs, in-car navigating system, etc.)

Good dialogue systems assist users to access information conveniently and finish tasks efficiently.

Adapted from: http://deepdialogue.miulab.tw

# Two Main Branches of Dialogue Systems

## Task-Oriented (TOD)

- Personal assistant, helps users achieve a certain task
- Combination of <u>rules</u> and <u>statistical</u> components
  - POMDP for spoken dialog systems (Williams and Young, 2007)
  - End-to-end trainable task-oriented dialogue system (Wen et al., 2016)
  - End-to-end reinforcement learning dialogue system (Li et al., 2017; Zhao and Eskenazi, 2016)



JARVIS – Iron Man's Personal Assistant

## Open Domain (ODD)

- No specific goal, focus on natural responses
- Using variants of seq2seq model
  - A neural conversation model (Vinyals and Le, 2015)
  - Reinforcement learning for dialogue generation (Li et al., 2016)
  - Conversational contextual cues for response ranking (AI-Rfou et al., 2016)



Baymax – Personal Healthcare Companion

# How it started

# Brief History of Dialogue Systems

**Multi-modal systems**
**e.g., Microsoft MiPad, Pocket PC**

**TV Voice Search**
**e.g., Bing on Xbox**

**Personal assistants, competitions, the DNN wave, and datasets explosion**

**2011**

**Task-specific argument extraction**
**(e.g., Nuance, SpeechWorks)**
*User: "I want to fly from Boston to New York next week."*

**Early 2000s**

IBM WATSON

**Early 1990s**

**Intent Determination**
**(Nuance's Emily™, AT&T HMIHY)** *User: "Uh…we want to move…we want to change our phone line from this house to another house"*

**DARPA CALO Project**

**Keyword Spotting**
*(e.g., AT&T)*
*System: "Please say collect, calling card, person, third number, or operator"*

**Material:**
http://deepdialogue.miulab.tw

# Language Empowering Intelligent Assistant



Apple Siri (2011)

Google Now (2012)
Google Assistant (2016)

Microsoft Cortana (2014)

Amazon Astro
(2021)

Amazon Alexa/Echo (2014)

Facebook M & Bot (2015)

Google Home (2016)

Apple HomePod (2017)

# Let's remember... GPT-2

*In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains.*

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.  Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.  Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.  Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.  Pérez and his friends were astonished to see the unicorn herd. These creatures could be seen from the air without having to move too much to see them – they were so close they could touch their horns.  While examining these bizarre creatures the scientists discovered that the creatures also spoke some fairly regular English. Pérez stated, "We can see, for example, that they have a common 'language,' something like a dialect or dialectic."

Dr. Pérez believes that the unicorns may have originated in Argentina, where the animals were believed to be descendants of a lost race of people who lived there before the arrival of humans in those parts of South America.

# Though…

*In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains.*

*Even more surprising to the researchers was the fact that the unicorns spoke perfect English.*

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.  Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved. Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.  Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.  Pérez and his friends were astonished to see the unicorn herd. These creatures could be seen from the air without having to move too much to see them – they were so close they could touch their horns.  While examining these bizarre creatures the scientists discovered that the creatures also spoke some fairly regular English. Pérez stated, "We can see, for example, that they have a common 'language,' something like a dialect or dialectic."

Dr. Pérez believes that the unicorns may have originated in Argentina, where the animals were believed to be descendants of a lost race of people who lived there before the arrival of humans in those parts of South America.

GTP-2: https://blog.openai.com/better-language-models/
https://www.wired.com/story/ai-text-generator-too-dangerous-to-make-public/

Then, in 2020 GPT-3 showed its capabilities...



Q&A

Answers  Generation  Conversation

Answer questions based on existing knowledge.

**Prompt**

I am a highly intelligent question answering bot. If you ask me a question that is rooted in truth, I will give you the answer. If you ask me a question that is nonsense, trickery, or has no clear answer, I will respond with "Unknown".

Q: What is human life expectancy in the United States?
A: Human life expectancy in the United States is 78 years.

Q: Who was president of the United States in 1955?
A: Dwight D. Eisenhower was president of the United States in 1955.

Q: Which party did he belong to?
A: He belonged to the Republican Party.

# GPT-3 capabilities…

## Grammar correction

Transformation  Generation

Corrects sentences into standard English.

**Prompt**

Correct this to standard English:

She no went to the market.

**Sample response**

She didn't go to the market.

## Parse unstructured data

Transformation  Generation

Create tables from long form text by specifying a structure and supplying some examples.

**Prompt**

A table summarizing the fruits from Goocrux:

There are many fruits that were found on the recently discovered planet Goocrux. There are neoskizzles that grow there, which are purple and taste like candy. There are also loheckles, which are a grayish blue fruit and are very tart, a little bit like a lemon. Pounits are a bright green color and are more savory than sweet. There are also plenty of loopnovas which are a neon pink flavor and taste like cotton candy. Finally, there are fruits called glowls, which have a very sour and bitter taste which is acidic and caustic, and a pale orange tinge to them.

| Fruit | Color | Flavor |

**Sample response**

| Neoskizzle | Purple | Sweet |

| Loheckle | Grayish blue | Tart |

| Pounit | Bright green | Savory |

| Loopnova | Neon pink | Cotton candy |

| Glowl | Pale orange | Sour and bitter |

# GPT-3 capabilities…

## Factual answering

Answers | Generation | Conversation | Classification

Guide the model towards factual answering by showing it how to respond to questions that fall outside its knowledge base. Using a '?' to indicate a response to words and phrases that it doesn't know provides a natural response that seems to work better than more abstract replies.

**Prompt**

Q: Who is Batman?
A: Batman is a fictional comic book character.

Q: What is torsalplexity?
A: ?

Q: What is Devz9?
A: ?

Q: Who is George Lucas?
A: George Lucas is American film director and producer famous for creating Star Wars.

Q: What is the capital of California?
A: Sacramento.

Q: What orbits the Earth?
A: The Moon.

Q: Who is Fred Rickerson?
A: ?

Q: What is an atom?
A: An atom is a tiny particle that makes up everything.

In 2022...
ChatGPT Revolution

# ChatGPT Revolution

ChatGPT
Revolution

```python
import datetime

def parse_expenses(expenses_string):
    """Parse the list of expenses and return the list of triples (date, value, currency).
    Ignore lines starting with #.
    Parse the date using datetime.
    Example expenses_string:
        2016-01-02 -34.01 USD
        2016-01-03 2.59 DKK
        2016-01-03 -2.72 EUR
    """
    expenses = []
    for line in expenses_string.splitlines():
        if line.startswith("#"):
            continue
        date, value, currency = line.split(" ")
        expenses.append((datetime.datetime.strptime(date, "%Y-%m-%d"),
                        float(value),
                        currency))
    return expenses
```

# GPT4-V

- Impressive capabilities for understanding images and text.

- Interpretation of radiological images, less hallucinations, perspective in images, commonsense, understand flows, great temporal understanding, interpretation of emotions, aesthetics, differences in images, navigate user interfaces, and more.

- Paper: Yang, Z., Li, L., Lin, K., Wang, J., Lin, C. C., Liu, Z., & Wang, L. (2023). The dawn of lmms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*.

# Comparison of leading generative AI language models



Number of parameters (log scale)

IOT ANALYTICS
MARKET INSIGHTS FOR THE INTERNET OF THINGS

**Lead organization**
- Alphabet
- Amazon
- DeepMind
- Meta
- OpenAI
- Other

Minerva (540B)
PaLM (540B)
GPT 3.5
GPT-3 175B (davinci)
FLAN   LaMDA
Chinchilla   NLLB
PLUG   AlexaTM 20B
GPT-NeoX-20B
Codex   UL2
T5-11B   XGLM
Generative BST
XLMR-XXL   Sparse all-MLP
T5-3B
GPT-2
HuBERT
data2vec (language)
Word2Vec (large)
RoBERTa
Seq2Seq LSTM   BERT-Large   BART-large
GNMT   XLNet   ALBERT-xxlarge
Word2Vec (small)
GPT
Transformer

y-axis: 1,000,000M / 100,000M / 10,000M / 1,000M / 100M

x-axis: Publication year — 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022 2023

Note: Based on Sevilla et. al. (2022). "For new models (from 2020 onward) it is harder to assess these criteria, so we fall back to a subjective selection. We refer to models meeting our selection criteria as 'milestone models.' The authors selected the AI systems for inclusion based on the following necessary criteria: Have an explicit learning component, showcase experimental results, advance the state of the art. In addition, the systems had to meet at least one of the following notability criteria: Paper has more than 1000 citations, historical importance, important state-of-the-art advance, deployed in a notable context. Source: ourworldindata.org, Sevilla et. al. (2022), IoT Analytics

# LANGUAGE MODEL SIZES TO MAR/2023

# LLM Training Costs on MosaicML Cloud

| Model | Billions of Tokens (Compute-optimal) | Days to Train on MosaicML Cloud | Approx. Cost on MosaicML Cloud |
|---|---|---|---|
| GPT-1.3B | 26B | 0.14 | $2,000 |
| GPT-2.7B | 54B | 0.48 | $6,000 |
| GPT-6.7B | 134B | 2.32 | $30,000 |
| GPT-13B | 260B | 7.43 | $100,000 |
| **GPT-30B \*** | **610B** | **35.98** | **$450,000** |
| GPT-70B \*\* | 1400B | 176.55 | $2,500,000 |

# How do they work?

**101 Transformers**

# Language Modeling

It is the task of predicting the next word given its contextual text:

- In the beginning _____

# We used it almost every day

# Introduction to Transformers

- Proposed in 2018 by Google with the following features:
  - Using contextual word vectors
  - Encoding & Decoding
  - Attention mechanisms
  - Great depth in the number of layers and neurons
  - Training with millions of data
  - Unsupervised training mechanism



Source: Alammar, 2018

# Word vector representation

- We are looking for a compact representation but with certain syntactic-semantic properties



Male-Female

Verb Tense

Country-Capital

Source: Google Developers

- A representation of textual input, including position and contextuality

# Types

- **Encoder - Only**
  - Text classification
  - Sentiment analysis
  - Name entity recognition
- **Encoder-Decoder**
  - Translation
  - Text summarization
  - Question - Answering
- **Decoder- Only**
  - Text completion
  - Tet generation
  - Translation
  - Image caption
  - Question-Answering



Evolutionary Tree

# Encoder-decoder models

- Encoder: Generates an intermedium representation of the input text
- Decoder: Generates the text output considering the encoded information and the one being generated

Intermedium Representation



Source: Alammar, 2018

# Attention Mechanisms

- Self-attention mechanism that allows weighting the contextual information considered at each moment



Source: Alammar, 2018

# Attention Mechanisms

# Unsupervised Training

17GB data & 355 GPUs and $4.6M cost (GPT-3)
45GB data & 10k GPUs (A100) for 11 months and $200M cost (GPT-4)

## Unsupervised Pre-training

Input (features) | a | robot | must

Correct output (label): obey

GPT-3
(under training)

Output (Prediction)

Source: Alammar, 2018

# Internal Architecture

OUTPUT | I am a student |

ENCODER → DECODER

ENCODER → DECODER

ENCODER → DECODER

ENCODER → DECODER

ENCODER → DECODER

ENCODER → DECODER

INPUT | Je suis étudiant |

Source: Alammar, 2018

# First step for Generation – Encoding process

Decoding time step: (1) 2  3  4  5  6                    OUTPUT

Linear + Softmax

ENCODER

ENCODER

DECODER

DECODER

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT          Je        suis      étudiant

Source: Alammar, 2018

# Second step for Generation – Decoding process

- Recurrent process controlled by the encoded information

Decoding time step: 1 (2) 3 4 5 6    OUTPUT    I

Kencdec    Vencdec

Linear + Softmax

ENCODERS    DECODERS

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT    Je    suis    étudiant    PREVIOUS OUTPUTS    I

Source: Alammar, 2018

# Machine translation system

**Training Dataset**

| I | am | a | student | &lt;to-fr&gt; | je | suis | étudiant |
|---|---|---|---|---|---|---|---|
| let | them | eat | cake | &lt;to-fr&gt; | Qu'ils | mangent | de |
| good | morning | &lt;to-fr&gt; | Bonjour | | | | |

**Output #2**
Position #5
Time step #2
allez-vous

**Output #1**
Position #4
Time step #1
Comment

Transformer-Decoder

| how | are | you | &lt;to-fr&gt; | … | |
|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | | 1024 |

Source: Alammar, 2018

# Summarization

Wikipedia



**Training Dataset**

| Article #1 tokens | **<summarize>** | Article #1 Summary |
| Article #2 tokens | **<summarize>** | Article #2 Summary | padding |
| Article #3 tokens | | **<summarize>** | Article #3 Summary |

**Output #2**
Position #115
Time step #2

**Output #1**
Position #114
Time step #1

Transformer-Decoder

**<summarize>**

| 1 | ... | 113 | 114 | 256 |

# SotA Generative Approach

**Knowledge Base**

**Persona**

**Persona**
I am an artist
I have four children
I recently got a cat
I enjoy walking for exercise
I love watching Game of Thrones

**Dialog History**
👩 Hi
🤖 Hello ! How are you today ?

**Receive an utterance from the user**
👩 I am good thank you , how are you.

**Dialog Agent**

**Generate a reply**
🤖 Great, thanks ! My children and I were just about to watch Game of Thrones.

# SotA Generative Approach

- Wolf, T., Sanh, V., Chaumond, J. and Delangue, C., 2019. Transfertransfo: A transfer learning approach for neural network based conversational agents. arXiv preprint arXiv:1901.08149.

# Additional Resources

- [The annotated transformer](): Rush, A. M. (2018, July). The annotated transformer. In Proceedings of workshop for NLP open source software (NLP-OSS) (pp. 52-60).

- [The illustrated transformer](): Alammar, J. (2018). The illustrated transformer. The Illustrated Transformer–Jay Alammar–Visualizing Machine Learning One Concept at a Time, 27.

- [A survey of transformer](): Lin, T., Wang, Y., Liu, X., & Qiu, X. (2022). A survey of transformers. AI Open.

# Trends

# A new concept for search

# RAG: Retrieval Augmented Generation

# Combining Generation and Retrieval (I)

- GPT-X large models:
  - They incorporate language and knowledge at the same time, but their info could quickly become outdated. (Retraining may not be an option + no guarantees)

- Retro [Borgeaud et al., 2022]:
  - Splits data allowing permanent updating

# Combining Generation and Retrieval (II)



INPUT — The Dune film was released in

**1) EMBED WITH BERT**

SENTENCE EMBEDDING

**2) QUERY**
approximate nearest neighbor

Database

**2) RETRIEVE**

Nearest Neighbor 1

Dune is a 2021 American epic science fiction film directed by Denis Villeneuve

It is the first of a planned two-part adaptation of the 1965 novel by Frank Herbert

Nearest Neighbor 2

Dune is a 1984 American epic science fiction film written and directed by David Lynch

and based on the 1965 Frank Herbert novel of the same name

RETRO
Retrieval-Enhanced Transformer

OUTPUT — 2021

# Combining Generation and Retrieval (III)



RETRO Transformer

# Combining Generation and Retrieval (IV)



Figure 1: Overview of our approach. We combine a pre-trained retriever (*Query Encoder + Document Index*) with a pre-trained seq2seq model (*Generator*) and fine-tune end-to-end. For query $x$, we use Maximum Inner Product Search (MIPS) to find the top-K documents $z_i$. For final prediction $y$, we treat $z$ as a latent variable and marginalize over seq2seq predictions given different documents.

Lewis, P., Perez, E.,Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., ... & Kiela, D. (2020). Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, *33*, 9459-9474.

# Instruction and RLHF (1)

**Train Language Model**

**Prompts & Text Dataset**

**Initial Language Model**

**Human Augmented Text (Optional)**

**Pretraining Language Models:**

- Begin with a language model that has been pre-trained:
  - This initial model serves as the starting point for RLHF.
  - The choice of the base language model can vary, ranging from smaller models to LLMs.

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, *35*, 27730-27744.

# Instruction and RLHF (2)



**Collecting Data and Training Reward Models:**

- Data is generated to train a **reward model**, which will be used for guiding the AI model's behavior.

- Gathering data can be done through human interaction. Users or experts provide feedback and evaluations on the AI agent's actions.
  - **For example**, users can rate different responses generated by the AI, indicating which responses are preferred.

- Alternatively, data can be also collected from demonstrations, where humans perform the desired task, providing a supervised learning signal.

- The collected data is used to train a **reward model**, which predicts how "good" or "preferable" a given AI output is ranked based on human feedback.

# Instruction and RLHF (3)

**Fine-Tuning the Language Model**

- The **pre-trained LM** is **fine-tuned** using **reinforcement learning** (RL) techniques.
  - During **fine-tuning**, the reward model **guides** the **model's outputs**.
  - The model seeks to **maximize cumulative rewards** according to the **reward model's predictions**.
  - The reward model was trained on **human preferences**!!

- The **AI agent** takes **actions** (generated text) in an environment, and the **reward model** provides **feedback** on the **quality** of those **actions**.

- The agent then **adjusts** its **behavior** to **optimize** for the actions that yield **higher rewards**, effectively learning from human feedback.
  - Kullback–Leibler divergence is used to avoid large differences between the output of the baseline and tuned LM.

- **Fine-tuning** typically involves running **multiple iterations,** where the AI agent **refines its behavior** over time

**Prompts Dataset**

x: A dog is...

**Initial Language Model**

Base Text

y: a furry mammal

**Tuned Language Model (RL Policy)**

*Parameters Frozen\**

RLHF Tuned Text

y: man's best friend

**Reward (Preference) Model**

text $r_\theta$

**Reinforcement Learning Update (e.g. PPO)**

$$\theta \leftarrow \theta + \nabla_\theta J(\theta)$$

$$-\lambda_{\mathrm{KL}} D_{\mathrm{KL}}\left(\pi_{\mathrm{PPO}}(y|x) \,\|\, \pi_{\mathrm{base}}(y|x)\right)$$

*KL prediction shift penalty*

$+$

$r_\theta(y|x)$

# DPO: Direct Preference Optimization

**An alternative to RLHF:**

- Requires a dataset of human preference pairs, each containing a prompt and two possible completions—one preferred and one dispreferred.
- The LLM is then **fine-tuned** to **maximize the likelihood** of **generating preferred completions** and **minimize the likelihood** of **generating dispreferred ones.**
- **Advantages:** Simple, Stable, Efficient and more Effective than RLHF.



Rafailov, Rafael, et al. "Direct preference optimization: Your language model is secretly a reward model." *arXiv preprint arXiv:2305.18290* (2023).

# Finetuning LLMs

- Parameter Efficient Fine-Tuning (PEFT)

- LOw Rank Adaptation (LORA)

- Prompt & Prefix Tuning

- ONNX, DeepSpeed, INT8

- Smaller models + fast adaptation + energy-friendly

  - Cons: additional parameters -> more delay

- URL: https://huggingface.co/blog/peft

# The Llama effect

- A model trained by Meta
  - LLaMA 65B & LLaMA 33B trained on 1.4 trillion tokens.
  - LLaMA 7B trained on one trillion tokens.
- Better text selection
  - Reduced hallucinations, bias and toxicity
- Better data processing and training
  - 21 days training on 2k Nvidia A100 GPU
  - Estimated cost: $2.4 million
  - GPT-4: 10k GPUs (A100) for 11 months and $200M cost
- Released for academia and commercial applications

Touvron, Hugo, et al. "Llama 2: Open foundation and fine-tuned chat models." *arXiv preprint arXiv:2307.09288* (2023).

62

# Vicuna & Alpaca: Quick adaptation



Zhang, Renrui, et al. "Llama-adapter: Efficient fine-tuning of language models with zero-init attention." *arXiv preprint arXiv:2303.16199* (2023).

# Multimodal fusion



Zhang, Renrui, et al. "Llama-adapter: Efficient fine-tuning of language models with zero-init attention." *arXiv preprint arXiv:2303.16199* (2023).

# Multimodal dataset

- Based on "Wu, Shengqiong, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. "Next-gpt: Any-to-any multimodal llm." arXiv preprint arXiv:2309.05519 (2023)."

- Goals:
  - Large use of diffusion models: image + audio (TTS with emotions + sounds)
  - Multiple topics
  - More complex dialogues (multiparty, multicultural)

# Chain-of-Thought (CoT)

**Standard Prompting**

**Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The answer is 27. ❌

**Chain of Thought Prompting**

**Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✔

Legend:
- ▨ Finetuned GPT-3 175B
- ▨ Prior best
- □ PaLM 540B: standard prompting
- ■ PaLM 540B: chain-of-thought prompting

Solve rate (%) — Math Word Problems (GSM8K)
- 33
- 55
- 18
- 57

Wei, Jason, et al. "Chain-of-thought prompting elicits reasoning in large language models." *Advances in Neural Information Processing Systems* 35 (2022): 24824-24837.

# Different flavors of CoT



Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Gianinazzi, L., Gajda, J., ... & Hoefler, T. (2023). Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687*.

# Tools for LLMs + Plugins



The New England Journal of Medicine is a registered trademark of [QA("Who is the publisher of The New England Journal of Medicine?") → Massachusetts Medical Society] the MMS.

Out of 1400 participants, 400 (or [Calculator(400 / 1400) → 0.29] 29%) passed the test.

The name derives from "la tortuga", the Spanish word for [MT("tortuga") → turtle] turtle.

The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.

Figure 1: Exemplary predictions of Toolformer. The model autonomously decides to call different APIs (from top to bottom: a question answering system, a calculator, a machine translation system, and a Wikipedia search engine) to obtain information that is useful for completing a piece of text.

Toolformer [Schick et al., 2023]

# AutoGPT

- *Auto-GPT is a Python app based on GPT-4 that enables AI to act independently without user input. It breaks down AI actions into "thoughts," "reasoning," and "criticism" and works to accomplish user-set goals.*



https://github.com/Significant-Gravitas/Auto-GPT

# Agents



Park, Joon Sung, et al. "Generative agents: Interactive simulacra of human behavior." *arXiv preprint arXiv:2304.03442* (2023).

# Additional trends

- Length of contextual information: from 128 to 100k (Claude2)
  - Positional interpolation: Chen, S., Wong, S., Chen, L., & Tian, Y. (2023). Extending context window of large language models via positional interpolation. arXiv preprint arXiv:2306.15595.

# Other tools and models

- Langchain or Llamaindex: high level control of LLMs + RAG + plugins

- Pinecone, Chroma DB or Faiss: Efficient handling and storage of vector indexes

- MemGPT: Interesting handling of Dynamic knowledge and contextual memory

- Zephyr and Mistral: current SotA OS models

- Benchmarks: HuggingFace

- Embeddings: MTEB

Packer, Charles, et al. "MemGPT: Towards LLMs as Operating Systems." *arXiv preprint arXiv:2310.08560* (2023).

# Challenges

# Challenges

Variability in natural language (syntactics and semantics)

Recall/Precision Trade-off

Meaning Representation

Multimodality

Hallucinations and Math Problems

Multilingualism (low resourced languages)

AB is 1000 bigger than 1062?

Yes, 1000 is bigger than 1062. 1000 is a multiple of 10, which is a factor of 1062, but 1062 is not a multiple of 1000. 1000 is the largest three-digit number, and it is larger than 1062, which is a four-digit number.

# Problems with numbers

Was Lincoln's assassin on the same continent as Lincoln when the assassination occurred

Abraham Lincoln was assassinated by John Wilkes Booth while he was attending a play at Ford's Theatre in Washington, D.C. I do not know where John Wilkes Booth was at the time of the assassination, but Washington, D.C. is located on the continent of North America, so it is possible that Booth was also on that continent at the time. However, without access to more information, I cannot say for certain.

↺ Try again

# Problems to Reason

**DA:** The trophy didn't fit in the suitcase because it was too small. What was too small?

**ChatGPT:** In the statement "The trophy didn't fit in the suitcase because it was too small," the trophy is the object that is too small to fit in the suitcase. The suitcase is a larger container or enclosure, and the trophy is the object that is being placed inside the suitcase. However, without more information, it is not clear why the trophy was too small to fit in the suitcase or what size or dimensions the trophy and suitcase are.

# Problems to Reason

summarise this article https://www.nytimes.com/2023/03/11/technology/chatgpt-prompts-to-avoid-content-filters.html

The article discusses how ChatGPT, an AI language model, is being used to circumvent content filters that are designed to detect and block certain types of content. By generating prompts that are less likely to trigger these filters, ChatGPT can be used to create content that might otherwise be censored. While this technology could be useful for those seeking to express themselves freely, it also raises concerns about the potential misuse of AI-generated content, as well as the difficulty of regulating such content in the future. The article concludes by noting that there is still much to learn about the impact of AI language models like ChatGPT on the way we communicate online.

# Hallucinations

Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., ... & Liu, T. (2023). A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. arXiv preprint arXiv:2311.05232.

**Playground**

Load a preset...

Save   View code   Share   ...   ⚙

(1) Dogs are loyal.
(2) Dogs can be dangerous.

please describe the link between (1) and (2).

The link between (1) and (2) is that dogs can be dangerous because they are loyal. Dogs are loyal to their owners and will protect them from danger. This can make them dangerous animals to be around.

**Playground**

Load a preset...

Save   View code   Share   ...   ⚙

(1) Dogs are loyal.
(2) Dogs can be dangerous.

describe the link between (1) and (2).

The link between (1) and (2) is that dogs can be dangerous.

# Problems with Prompt Dependencies

# Problems with Prompt Dependencies

# Usability vs Alignment

# Challenges

Common Sense (World Knowledge)

Ability to learn

Transparency / Accountability

Long-term interaction

Automatic Evaluation

Ethical and Moral Behavior

# Ethical Aspects

*A safe AI for everyone*

# Current research areas

- Bias detection and mitigation:
  - Different types: gender, religion, race, cultural values, …
  - Not only at word level, but contextualized and intentionality
  - Performance detection
- Toxicity:
  - Different types: sexuality, violence, threatening, religion, race,..
  - Detect intentionality, controlled generation, dataset cleaning

# Current research areas

- Morality:
  - Train machines to do what is right (definition of right)
    - Asimov's three laws of robotics
  - Anthropomorphic behavior vs AI
  - Usability vs alignment

- Awareness and Consciousness
  - Definitions and implementation theories
  - Ethical issues when dealing with this type of systems: manipulation, decision criteria, reasoning, transparency, governance…

# Research & Activities in Spain

# SEPLN

Sociedad Española de Procesamiento del Lenguaje Natural

Conferences

Journal (indexed)

Webinars

Yellow pages for Research Groups

Support for competitions: IberLEF

URL: http://www.sepln.org/

# SomosNLP

- Research Community (Discord)
- Hackathon Somos NLP 2023: Los LLMs hablan Español
  - Keynotes, Awards, Workshops, Ask-Me-Anything (AMA)
- Datasets, Models in Spanish, Tutorials…
- https://somosnlp.org/

88

# DiverTLes

- Created by PLN.net (Spanish network)
- Specially focused on supporting Women working in NLP
- News, Round Tables, Blogs
- URL: https://gplsi.dlsi.ua.es/pln/divertles/

# RTTH: Red Temática en Tecnologías del Habla

Spanish network of researchers and institutions

Award to best journal papers

Albayzin evaluations

Summer/Fall schools

Distribution list

Support conferences: Iberspeech

URL: http://rtth.es/

- Participation in Amazon Alexa Prize Socialbot Grand Challenge

- Second time our team is selected

- Duration: November 14, 2022 - August 31, 2023

Two main targets:

Average score higher than 4.0/5.0

Duration higher than 20 minutes (engagingness)

Any topic but also careful answers for sensitive topics (e.g., health, investing, politics, toxicity)

Several qualification stages + increasing thresholds in terms of scores, interaction duration and latencies

Nine teams selected from around the world

7 from USA: University of California (Santa Cruz), **Stanford University,** Carnegie Mellon University, Stevens Institute of Technology, University of California (Santa Barbara), University of Illinois (Urbana-Champaign), and Virginia Tech.

2 from Europe: Czech Technical University in Prague and UPM,

# Alexa Social Grand Challenge (SGC5)

# Our Team

- Second time participating,

- Only two universities in Europe

- Only team from Spain

- 2 Ph.D students + 6 master students

- Paper with full description available (link)



93

# Our approach

# Amazon architecture

# Automatic image and sound generation



| Detection of entity | Image? | Image directive | Sound? | Sound directive |

# Incorporation of Fictional Characters for Alexa using Voice Cloning Techniques

- Modification of the Alexa default voice targeting other voices from fictional characters.
- Dataset:

  - Audios/speaker extracted from films or series (YouTube)

- Automatization Process:

  - Speech separation → Remove music, noise, normalization

  - Diarization → Detect different speakers and when they speak

  - Automatic Speech Recognition → Transcription for aligning

# Topic-Based Generative Models with Up-to-Date Knowledge

- **User:** Do you know the new movie of Barbie?
- **Chatbot:** Yes, I know the new movie of Barbie. It's directed by Greta Gerwig and stars Margot Robbie, Ariana Greenblatt, and Ryan Gosling.
- **User:** What is it about?
- **Chatbot:** Barbie is about to live in Barbie Land is to be a perfect being in a perfect place. Unless you have a full-on existential crisis. Or you're a Ken.
- **User:** When is it going to be published?
- **Chatbot:** It will be published on July 21, 2023.
- **User:** Which are the genres of the film?
- **Chatbot:** The genres of the movie are adventure, comedy, and fantasy.
- **User:** Thank you, I will go to watch it!
- **Chatbot:** The movie will be published on July 21, 2023. You are welcome.

# European Project

*ASTOUND*: Improving social competences of virtual agents through artificial consciousness based on the Attention Schema Theory.

**Goals:**
- ✓ Incorporation of new mechanisms for **controllability** in conversational chatbots
    - ✓ Knowledge graphs + persona profile automatization + dialogue summarization
- ✓ **Contextual awareness** through integration of multiple modalities (speech, video, text, wearables)
- ✓ New self-assessment metrics for new dimensions: **long-term interaction, memory, bias, toxicity**
- ✓ **Technology transfer**

- **Dates:** Dec 2022 – Nov 2025

- **Other participants:**
    - Ecole Normale Superieure (France),  University Medical Center Hamburg-Eppendorf (Germany), IndeepAI (Spain)
    - Princeton University (USA), Microsoft y MILA (Canada)

# The Attention Schema Theory (AST)

- Two general principles:
  - Information that comes out of a brain must have been in the brain.
  - The brain's models are never accurate.

- Important aspects:
  - Capability of controlling own-self attention
  - Capability of modeling/explain other's attention
  - Capability for explaining/handling beliefs and claims

- Do not claim to be the unique/perfect theory about consciousness, but…
  - Provides a framework for computational implementation

- Graziano, M. S., & Webb, T. W. (2015). The attention schema theory: a mechanistic account of subjective awareness. Frontiers in psychology, 6, 500.
- Liu, D., Bolotta, S., Zhu, H., Bengio, Y., & Dumas, G. (2023). Attention Schema in Neural Agents. arXiv preprint arXiv:2305.17375.

A

Controlling the body ← Modeling the body of others ← Explicit, reportable knowledge of our own bodies

Body Schema

Body

B

Controlling your attention ← Modeling others' attention ← Claims / beliefs about consciousness

Attention Schema

Attention

100

# DSTC10-Track 5 & DSTC11-Track 4

Benchmarks and challenges are required for progress in the field. These challenges targeted three main tasks:

- Propose and develop effective Automatic Metrics for evaluation of open-domain multilingual dialogs.
- Propose and develop Robust Metrics for dialogue systems trained with back translated and paraphrased dialogs in English.
- Propose and develop effective generation mechanisms to handle toxic users

**Datasets:**
- For training: Up to 18 Human-Human curated multilingual datasets (+3M turns), with turn/dialogue level automatic annotations including QE metrics or toxicity.
- Dev/Test: Up to 10 Human-Chatbot curated multilingual datasets (+150k turns), with turn/dialogue level human annotations.
- Dev/Test: Up to 3 different datasets containing toxic user's turns (+2k annotated turns)

**Links for registration and baselines:**
- ChatEval: https://chateval.org/dstc11
- GitHub: https://github.com/Mario-RC/dstc11_track4_robust_multilingual_metrics

**News: Proposal of two new challenges!!**
- DSTC12: Addressing multimodality and harness issues (red teaming)
- JSALT2024: Addressing multimodality, multiculturality and multiparty interaction

# Challenges

- Proposal for DSTC12: Dialog System Evaluation: Dimensionality, Language, Culture and Safety
  - In partnership with CMU, NYU, NUS, INESC-PT, Tencent, Microsoft
  - Topics: new evaluation dimensions, multicultural aspects, safety through red teaming
- Proposal for JSALT2024: M$^3$Eval: Workshop on Multicultural, Multimodal, Multiparty Dialogue Evaluation
  - In partnership with CMU, NYU, NUS, CUHK, INESC-PT, U. Sidney, Apple, Amazon, Meta, Stardust.ai
  - Topics: Multimodality, multiculturality and multiparty dialogues + automatic evaluation

**102**

# Bibliography (I)

- Adiwardana et al. "Towards a human-like open-domain chatbot." arXiv preprint arXiv:2001.09977 (2020).

- Bachman et al. "Learning with pseudo-ensembles." NIPS (2014).

- Bao et al. "PLATO: Pre-trained Dialogue Generation Model with Discrete Latent Variable." ACL (2020).

- Borgeaud et al. "Improving language models by retrieving from trillion of tokens". ICML, pp. 2206-2240.

- Finch and Choi. "Towards Unified Dialogue System Evaluation: A Comprehensive Analysis of Current Evaluation Protocols." SIGDial (2020).

- Golovanov et al. "Large-scale transfer learning for natural language generation." ACL (2019).

- Gupta et al. 2021. "Synthesizing Adversarial Negative Responses for Robust Response Ranking and Evaluation" Findings of ACL-IJCNLP (2021).

- Henderson et al. "ConveRT: Efficient and Accurate Conversational Representations from Transformers." Findings of EMNLP (2020).

- Han et al. "Fine-grained Post-training for Improving Retrieval-based Dialogue Systems." NAACL (2021).

- Humeau et al. "Poly-encoders: Architectures and Pre-training Strategies for Fast and Accurate Multi-sentence Scoring." ICLR (2020)

- Khapra and Sai. "A tutorial on evaluation metrics used in natural language generation." NAACL (2021).

- Lin, "Rouge: A package for automatic evaluation of summaries." Text summarization branches out (2004).

# Bibliography (II)

- Liu et al. "How NOT To Evaluate Your Dialogue System: An Empirical Study of Unsupervised Evaluation Metrics for Dialogue Response Generation." EMNLP (2016).

- Lowe et al. "The Ubuntu Dialogue Corpus: A Large Dataset for Research in Unstructured Multi-Turn Dialogue Systems." SIGDial (2015).

- Mehri and Eskenazi. "USR: An Unsupervised and Reference Free Evaluation Metric for Dialog Generation." ACL (2020).

- Mehri and Eskenazi. "Unsupervised Evaluation of Interactive Dialog with DialoGPT." SIGDial (2020).

- Mehri, S., Choi, J., D'Haro, L. F., Deriu, J., Eskenazi, M., Gasic, M., ... & Zhang, C. (2022). Report from the NSF future directions workshop on automatic evaluation of dialog: Research directions and challenges. arXiv preprint arXiv:2203.10012.

- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, *35*, 27730-27744.

- Papineni et al. "BLEU: a method for automatic evaluation of machine translation." ACL (2002).

- Roller et al. "Recipes for Building an Open-Domain Chatbot." EACL (2021).

- Sai et al. "Improving Dialog Evaluation with a Multi-reference Adversarial Dataset and Large Scale Pretraining" TACL (2020).

- Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Zettlemoyer, L., ... & Scialom, T. (2023). Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*.

- Scudder, H. Probability of error of some adaptive pattern-recognition machines. IEEE Transactions on Information Theory, 11(3):363–371 (1965).

# Bibliography (III)

- lSerban et al. "Building end-to-end dialogue systems using generative hierarchical neural network models." AAAI (2016).

- Serban et al. "A hierarchical atent variable encoder-decoder model for generating dialogues." AAAI (2017).

- Thoppilan et al. "Lamda: Language models for dialog applications." arXiv preprint arXiv:2201.08239 (2022).

- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M. A., Lacroix, T., ... & Lample, G. (2023). Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.

- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., & Zhou, D. (2022). Chain of Thought Prompting Elicits Reasoning in Large Language Models

- Yeh, Y. T., Eskenazi, M., & Mehri, S. (2021). A comprehensive assessment of dialog evaluation metrics. arXiv preprint arXiv:2106.03706.

- Zhang et al. "Personalizing Dialogue Agents: I have a dog, do you have pets too?" ACL (2018).

- Zhang et al. "DIALOGPT : Large-Scale Generative Pre-training for Conversational Response Generation" ACL (2020).

- Zhang, C., Sadoc, J., D'Haro, L. F., Banchs, R., & Rudnicky, A. (2021). Automatic evaluation and moderation of open-domain dialogue systems. arXiv preprint arXiv:2111.02110.

- Zhang et al. "FineD-Eval: Fine-grained Automatic Dialogue-Level Evaluation." EMNLP (2022).

- Zhang et al. "MDD-Eval: self-training on augmented data for multi-domain dialogue evaluation." AAAI (2022).

- Zhang, C., D'Haro, L. F., Zhang, Q., Friedrichs, T., & Li, H. (2023). PoE: A Panel of Experts for Generalized Automatic Dialogue Assessment. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, *31*, 1234-1250.