

COMP9444: Neural Networks and Deep Learning

Generative AI Applications

Jingying Gao

School of Computer Science and Engineering

July 31, 2024

What cognitive capabilities do humans have?

Inputs

- See – Visual Processing
- Listen -Auditory Processing
- Read – Textual Understanding
- Feel - Sensory and Emotional Processing



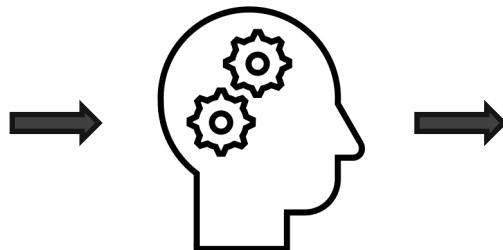
Outputs

- Speaking/Singing – Verbal Expression
- Writing – Textual Expression
- Painting – Visual Artistic Expression
- Thinking – Cognitive Reasoning

What capabilities does AI have?

Inputs

- See – Image / Video
- Listen – Audio / Language
- Read – Text
- Feel – Sensor



Outputs

- Speaking / Singing
- Writing
- Painting
- Thinking - Question Answering & Reasoning

Multimodal AI

Generative AI

Generative AI

- What Is Generative AI?
- What Are the Different Types of Generative AI Tasks?
- How Does Generative AI Work?
- How to Build an Interactive App Using Generative AI?

What is Generative AI?

- Generative AI (GenAI) is artificial intelligence capable of generating text, images, videos, or other data using generative models, often in response to prompts.
- Generative AI models learn the patterns and structure of their input training data and then generate new data that has similar characteristics.

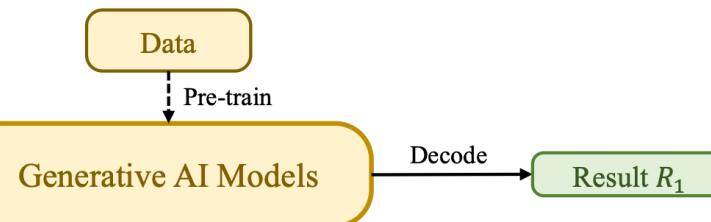
What Are the Different Types of Generative AI Tasks?

Unimodal

Please write a story about a cat.

Instruction I_1

Prompt



Once upon a time, there was a cat named Jessy....

Multimodal



Describe this picture.

Instruction I_2

Prompt

Draw a picture of a cat.

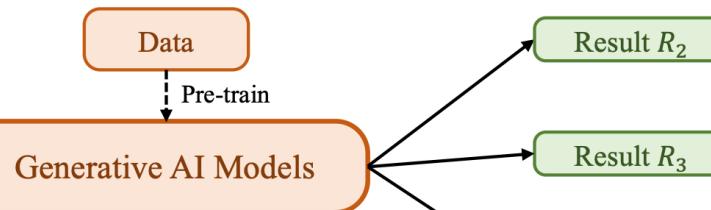
Instruction I_3

Prompt

Write a song about a cat.

Instruction I_4

Prompt



This is a cat.

Source: <https://arxiv.org/pdf/2303.04226.pdf>

Generative AI Tasks

- Text Generation: Text-to-Text, Image-to-Text, Speech-to-Text
- Image Generation: Image-to-Image, Text-to-Image, Talking face
- Audio Generation: Text-to-Audio, Image-to-Audio
- Video Generation: Text-to-Video, Image-to-Video

Demo of Generative AI Tasks

- Generate a Painting
- Image to Video
- Outfit Anyone
- Generate a Song

Demo of Generative AI Tasks

- Generate a painting
- <https://huggingface.co/spaces/multimodalart/stable-cascade>
- Image to Video
- <https://huggingface.co/spaces/multimodalart/stable-video-diffusion>
- Generate a song
- <https://suno.com/>

Generate a Song

Prompt

“Welcome, students, to the AI course. Today's lecture will cover generative AI and multimodal AI. We are going to explore various topics, including stable diffusion models and LLM models, among others.”

Enjoy the below songs:

Future of Creation

Rise Above

Which one is real?



Which one is real?



Which one is real?



Generative AI

- How do Text-to-Image Generation Models Work?
- How does the Diffusion Model Work?
- Fine-tune a Diffusion Model with a Customized Dataset.
- How to Build an Application based on Diffusion AI Models?

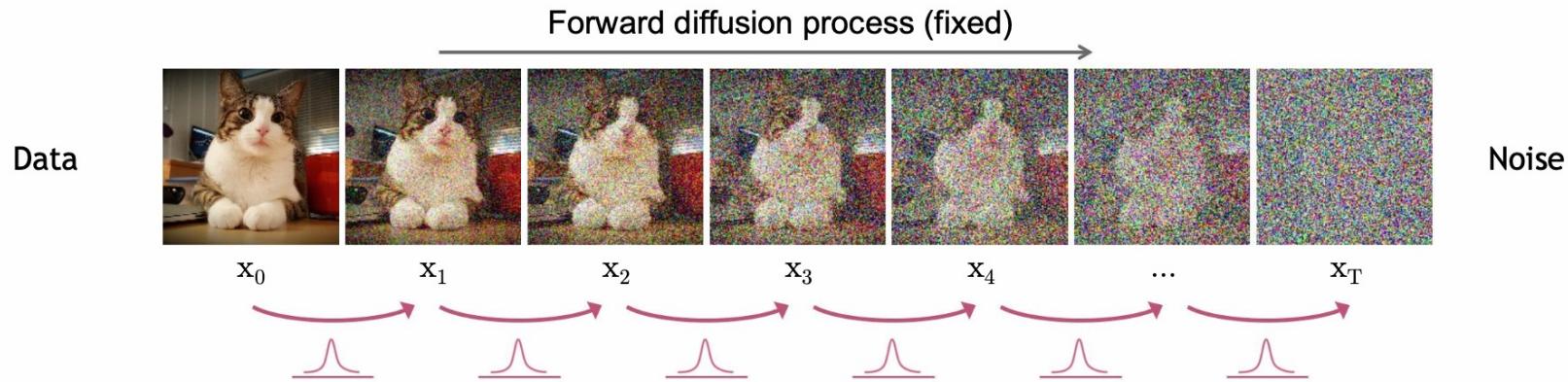
How Did the Idea Come About?



The central idea behind diffusion models is inspired by the thermodynamics of gas molecules, whereby the molecules diffuse from areas of high density to low density.

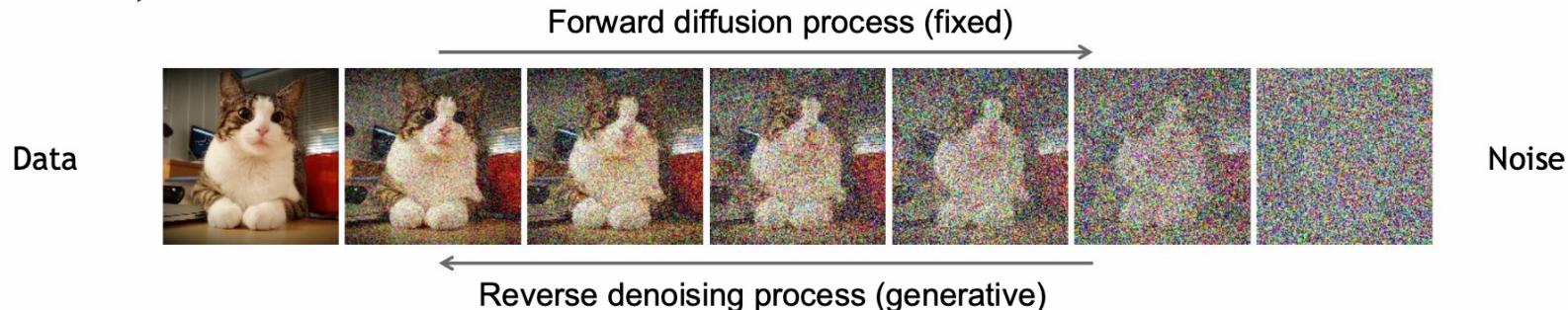
This movement is often referred to in physics literature as the increase of entropy or heat death.

How Does an AI Model Generate an Image? (Stable Diffusion/DALL-E)



Diffusion models work by gradually adding Gaussian noise through a series of T steps into the original image, a process known as diffusion.

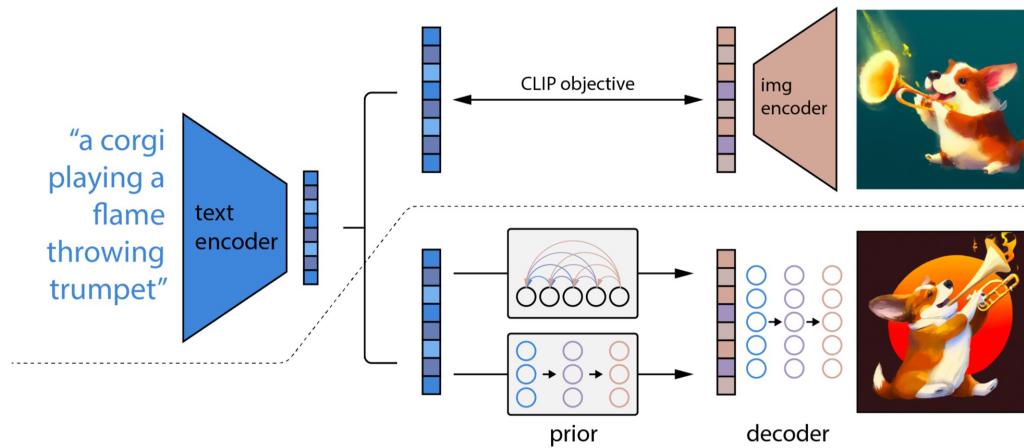
How does AI generate an Image? (Stable Diffusion / Dall-E)



- Backward Diffusion Process.
- We take a neural network Unet and learn to reverse this diffusion process.
- The reverse diffusion process, learn from the last right noise image, generate image from step t to step t-1.
- From step t, use a neural network Unet to predict noise, then the predicted noise will be subtracted from the image.

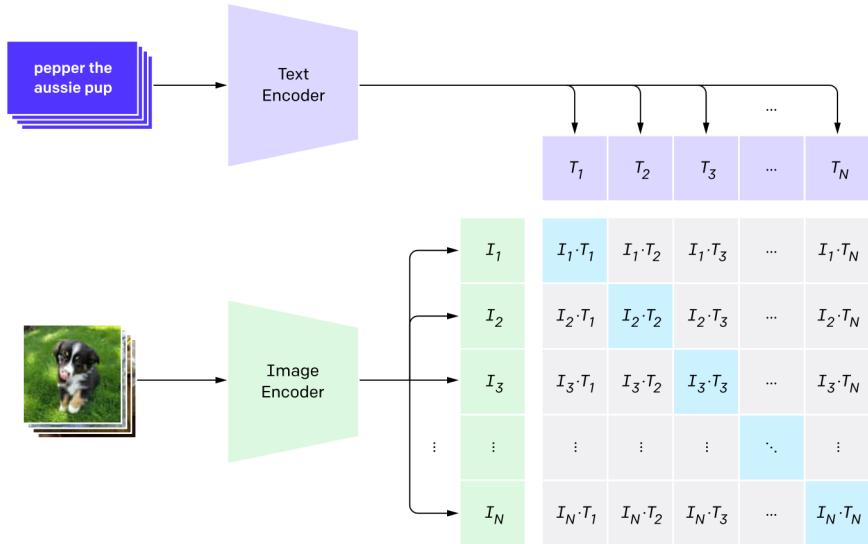
How to generate images that correspond with a piece of text?

- We've learned how diffusion can help us to generate an image from random noise, but the model is only able to generate random images.
- How can we make use of this model to synthesize images that correspond with a class name in our training data, a piece of text, or another images?



What is Contrastive Learning?

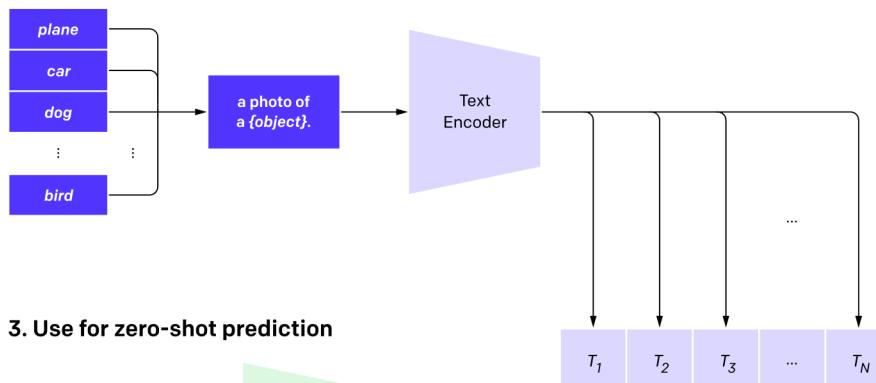
1. Contrastive pre-training



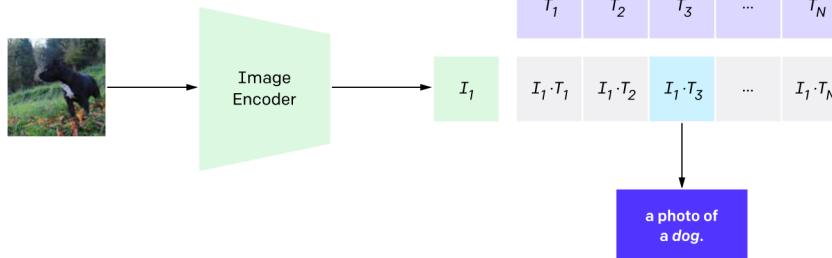
- Convert Image to Vector via Image Encoder.
- Convert Text to Vector via Text Encoder.
- Learn the pairs of image vector and text vector. The diagonal is the positive samples, others are negative samples.
- Minimize the negative samples.

How to generate images that correspond with a piece of text?

2. Create dataset classifier from label text



3. Use for zero-shot prediction

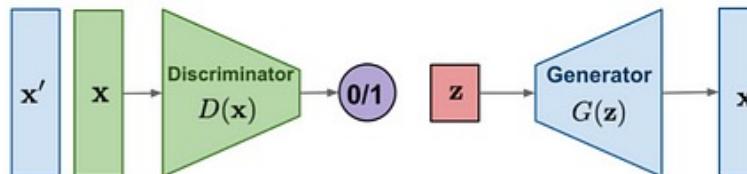


- CLIP pre-trains an image encoder and a text encoder to predict which images were paired with which texts in our dataset.
- We then use this behavior to turn CLIP into a zero-shot classifier. We convert all of a dataset's classes into captions such as “a photo of a dog” and predict the class of the caption CLIP estimates best pairs with a given image.
- Contrastive loss, is used to map vectors that model the similarity of input items.

Overview of different types of generative models

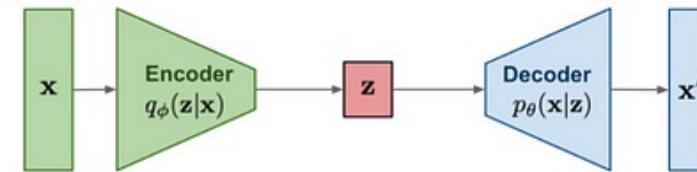
2014

GAN: Adversarial training



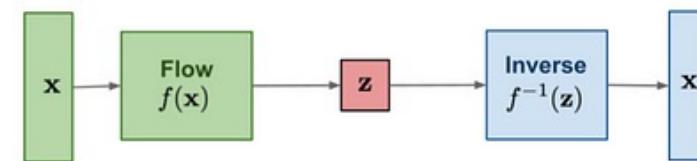
2013

VAE: maximize variational lower bound



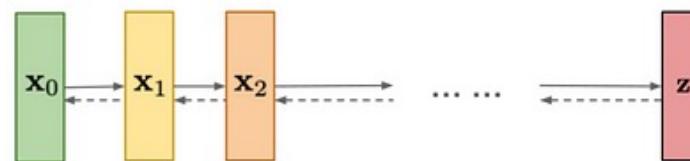
2015

Flow-based models:
Invertible transform of
distributions



2015
~2020

Diffusion models:
Gradually add Gaussian
noise and then reverse



GAN Adversarial Training

- **Discriminator:** It is a classifier that tries to distinguish between real data (from the training dataset) and fake data created by the generator. It's trained to accurately classify data as real or fake.
- **Generator:** It generates new data instances that resemble the training data. It starts with random noise and gradually learns to produce data (like images) that look similar to the actual data it's been trained on. Its goal is to create data so convincing that the discriminator can't tell it apart from real data.

VAE - maximize variational lower bound

A Variational Autoencoder (VAE) aims to learn a representation of input data in a latent space by maximizing the variational lower bound.

- **Encoder:** It takes the input data and transforms it into a distribution over the latent space. The encoder outputs parameters to this latent distribution, typically the mean and variance, effectively compressing the input data into a compact representation.
- **Decoder:** Starting from the latent space representation, the decoder attempts to reconstruct the original input data from the latent distribution. It learns to map the latent variables back to the original data space, aiming to minimize the difference between the original data and its reconstruction.

How to fine-tune the stable diffusion model?

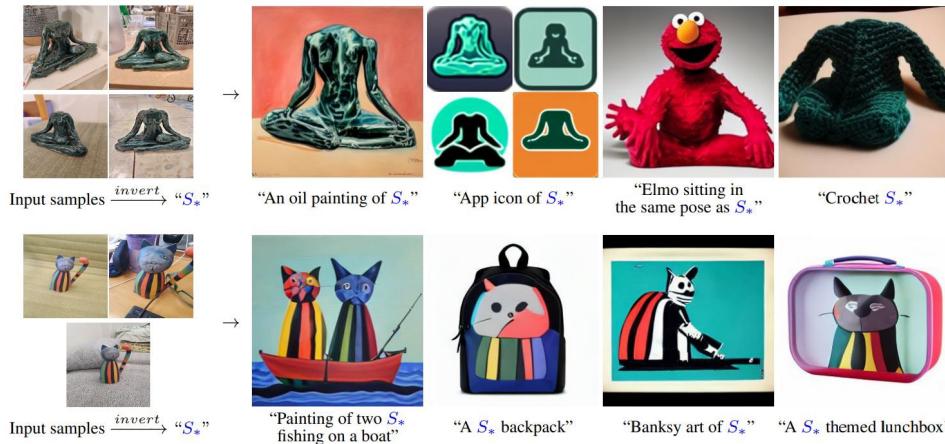
- Why do we need fine-tune the stable diffusion model?
- How to fine-tune the stable diffusion model?

Why do we need to fine-tune the stable diffusion model?

- Stable Diffusion was trained on 2.3 billion English image-text pairs. So this means that there are 2.3 billion pictures and their captions that are matched up, and this is the data that is used in training the Stable Diffusion model.
- Does it include UNSW Style images or if we want to generate some UNSW particular classes, such as UNSW library, UNSW campus?

How do we fine-tune the diffusion model to generate style images?

An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion



We learn to generate specific concepts, like personal objects or artistic styles, by describing them using new “words” in the embedding space of pre-trained text-to-image models. These can be used in new sentences, just like any other word.

The work builds on the publicly available [Latent Diffusion Models](#)

How do we fine-tune the diffusion model to generate style images?

- Prepare Customized Dataset: Image and Text Pairs.
- Download pre-trained Diffusion model's weights.
- Define hyper-parameters for the training process.
- Training.
- Run for generating images.

Demo Code Fine-tune model

<https://colab.research.google.com/drive/1F9UqlsHHmZZvgqAw78VqbAYWITwep9Qv?usp=sharing>

Fine-tune diffusion model



How to build an interactive diffusion application?

InstructPix2Pix: Learning to Follow Image Editing Instructions

For faster inference without waiting in queue, you may duplicate the space and upgrade to GPU in settings.

 [Duplicate Space](#)

[Generate](#)

[Load Example](#)

[Reset](#)

Edit Instruction

change background color to blue

 [Input Image](#)

Drop Image Here

- OR -

Click to Upload

 [Edited Image](#)



Steps
0

<input type="radio"/> Fix Seed
<input checked="" type="radio"/> Randomize Seed

Seed
1371

<input type="radio"/> Fix CFG
<input checked="" type="radio"/> Randomize CFG

Text CFG
7.5

Image CFG
1.5

How to build an AI image generation application?

- Huggingface models
- Huggingface Space
- Gradio

Gradio

- Gradio is a Python library that allows you to quickly create customizable web apps for your machine learning models and data processing pipelines.
- Gradio apps can be deployed on Hugging Face Spaces for free.

Demo Instruct Pix-2-Pix

Demo Instruct Pix 2 Pix

<https://huggingface.co/spaces/timbrooks/instruct-pix2pix>

<https://huggingface.co/spaces/timbrooks/instruct-pix2pix>

A Taste of Code

```
3 import random
4 import gradio as gr
5 import torch
6 from PIL import Image
7 from diffusers import StableDiffusionInstructPix2PixPipeline
8
9 example_instructions = [
10     "Change the hair color to blue",
11     "Add a mustache",
12     "Change the background to a beach",
13     "Add sunglasses",
14 ]
15
16 model_id = "timbrooks/instruct-pix2pix"
17 pipe = StableDiffusionInstructPix2PixPipeline.from_pretrained(model_id, torch_dtype=torch.float16, safety_checker=None).to("cuda")
18
19 def generate_image(input_image: Image.Image, instruction: str):
20     edited_image = pipe(instruction, input_image).images[0]
21     return edited_image
22
23 iface = gr.Interface(
24     fn=generate_image,
25     inputs=[
26         gr.Image(type="pil", label="Take a picture", sources=["webcam"]),
27         gr.Dropdown(choices=example_instructions, label="Edit Instruction"),
28     ],
29     outputs=gr.Image(type="pil", label="Edited Image"),
30     title="InstructPix2Pix",
31     description="Capture an image from the webcam, enter a prompt, and change the image.",
32 )
33
34 iface.launch(debug=True)
```

Generative AI - LLM

- What is LLM?
- How does GPT work?
- How to use VectorDB to build a customize chatbot?

What is LLM?

- A large language model (LLM) is a language model notable for its ability to achieve general-purpose language generation and other natural language processing tasks such as classification.
- LLMs acquire these abilities by learning statistical relationships from text documents during a computationally intensive self-supervised and semi-supervised training process.
- LLMs can be used for text generation, a form of generative AI, by taking an input text and repeatedly predicting the next token or word.

What does GPT stand for?

GPT =

Generative Pre-trained Transformer

What does GPT stand for?

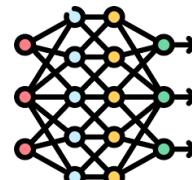
Generative

Generate New Text

“AI course is a general course to learn Artificial Intelligence.”

Pre-trained

Pretrained refers to how the model went through a process of learning from massive amount of data.

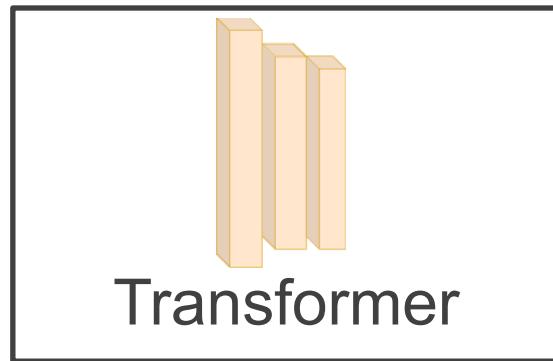


Transformer

A deep neural network model that transforms or changes an input sequence into an output sequence by using attention.

How does Transformer work?

Attention is
all you need



注意力就是你所
需要的

Attention Is All You Need

Ashish Vaswani*

Google Brain

avaswani@google.com

Noam Shazeer*

Google Brain

noam@google.com

Niki Parmar*

Google Research

nikip@google.com

Jakob Uszkoreit*

Google Research

usz@google.com

Llion Jones*

Google Research

llion@google.com

Aidan N. Gomez* [†]

University of Toronto

aidan@cs.toronto.edu

Lukasz Kaiser*

Google Brain

lukaszkaiser@google.com

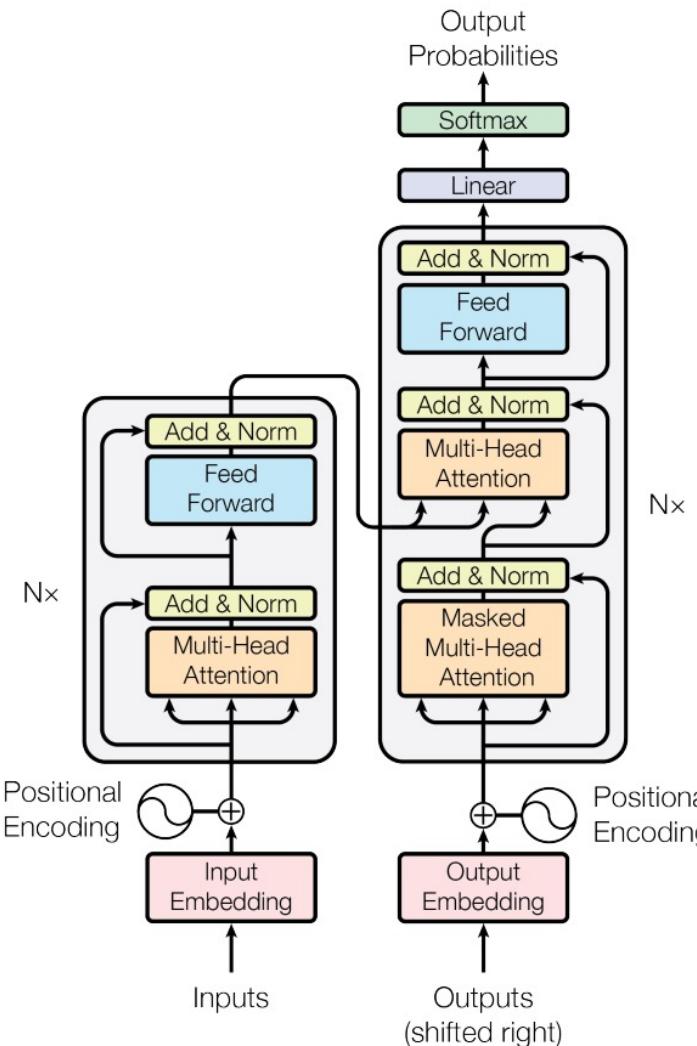
Illia Polosukhin* [‡]

illia.polosukhin@gmail.com

Abstract

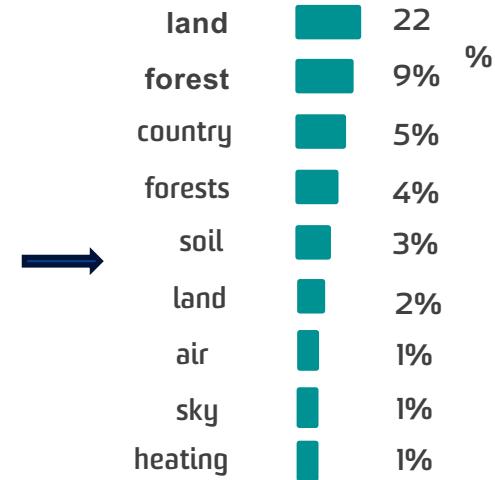
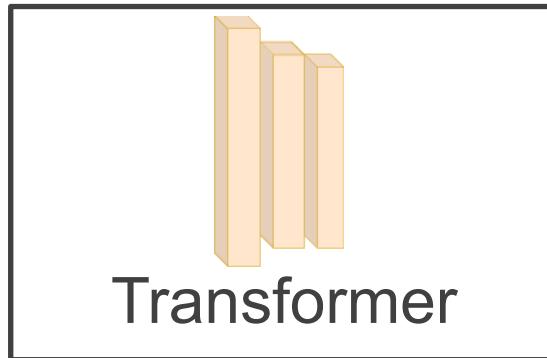
The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to

Transformer Architecture



How does GPT work?

Hot air balloons ascend into the sky by heating the air inside their envelopes, making it less dense than the cooler air outside, thus creating buoyancy that lifts them upwards...



- GPT model that is trained to take in a piece of text, or image, sound, to produce a prediction for what comes next in the passage.
- That prediction takes the form of a probability distribution over many different chunks of text that might follow.

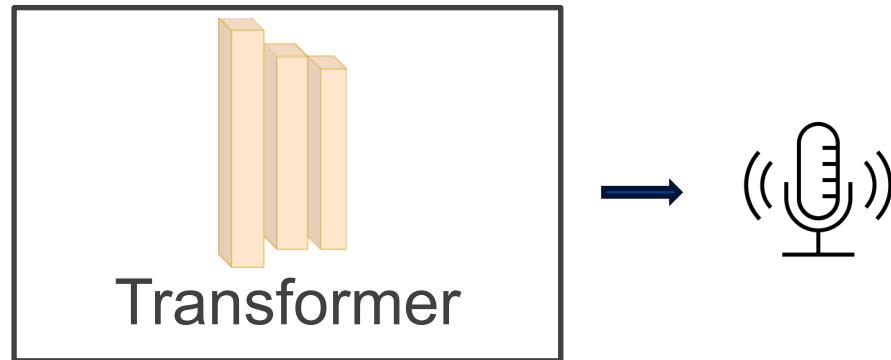
Models: Voice-to-Text



- The model takes in audio and produce a transcript.
- Example: OpenAI whisper model.

Generative AI Model: Text-to-Voice

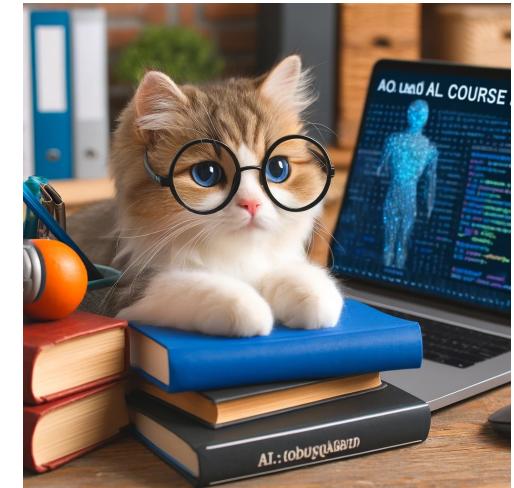
This is an AI course model example, producing synthetic speech from text.



- Transformer produces synthetic speech from text.

Generative AI Model: Text-to-Voice

Generate a cute cat
is studying AI
course.



Generated by ChatGPT

- The model takes in a text description and produce an image are based on Transformers.

LLM Applications

- LangChain
- RAG
- Agents

What is LangChain?

LangChain is an open-source framework for the development of applications that use large language models, such as GPT-4.



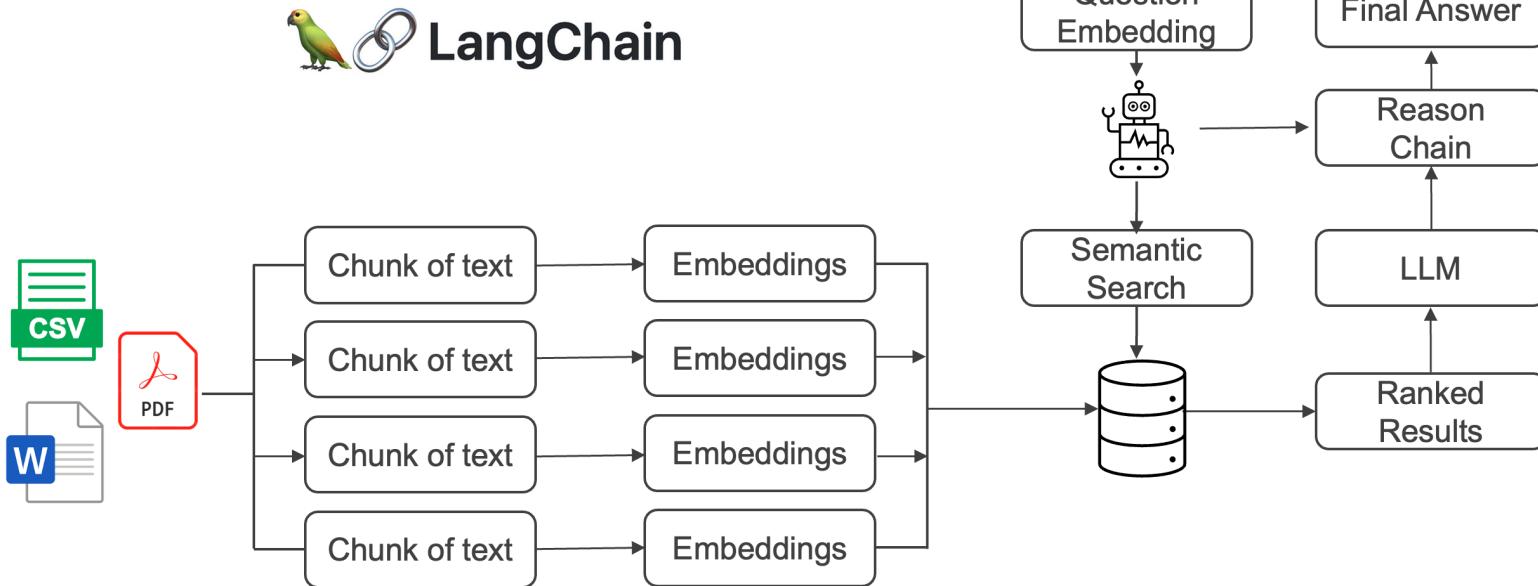
Lang = Language which refers to the use of large language models like GPT4.



Chain = Chaining that connect/chain these language models together to build applications

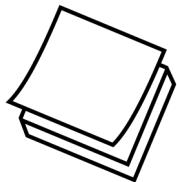
What is LangChain Architecture?

What is the propositional logic?



What is RAG?

RAG is when the LLM looks for additional information, augments the query and generate an answer with new information.



Retrieval: LLM finds extra information from the knowledge base documents based on “similarity” metric

Augmentation: LLM integrates the extra information into the prompt

Generation: LLM generates an answer using the augmented prompt

How to build an AI Tutorial Knowledge Bot?

Course knowledge bot

Enter your question.

question

Clear **Submit**

Generated Answer

$\text{Myth} \Rightarrow \neg \text{Mortal}$

Additional Information

4.1 Translate the above statements into Propositional Logic, using the symbols.
(1) If the unicorn is mythical, then it is immortal.
 $\text{Myth} \Rightarrow \neg \text{Mortal}$
(2) If it is not mythical, then it is mortal and a mammal.
 $\neg \text{Myth} \Rightarrow (\text{Mortal} \wedge \text{Mammal})$
(3) If the unicorn is either immortal or a mammal, then it is horned.

Flag

LLM Shortcomings

- Most time, we use LLMs in zero-shot mode, prompting a model to generate final output token by token without revising its work.
- This is akin to asking someone to compose an essay from start to finish, typing straight through with no backspacing allowed, and expecting a high-quality result.

LLM Agents

- The experience of prompting ChatGPT, receiving unsatisfactory output, delivering critical feedback to help the LLM improve its response, and then getting a better response.
- What if we automate the step of delivering critical feedback, so the model automatically criticizes its own output and improves its response?
- LLM Agents allows LLM to iterate over documents many times, take self-reflection and improve the results.

LLM Agents - Example

- We use a system built with Reflection to write some ABC code.
- Then, the AI will take this code, along with a prompt like “check the correctness of this code and tell me how to modify it,” and return it to the AI.
- The AI might identify bugs within it, and through this process, the AI completes its self-iteration.
- Although the quality of the revised code can't be guaranteed, the results are generally better on the whole.

LLM Agents - Demo

Agents

- Question/Task
- Manager
- user

Auto Agents Chat

Question/Task
Task "write a 10 words story about a girl" is submitted.
Waiting for agents to cooperate to complete the task.

Manager
Thought
The task requires creating a 10-word story about a girl. The Storyteller role is well-suited for this task as it involves creating a concise and engaging story within a given word limit. The Language Expert role can be used to summarize the final story in a clear and engaging manner. The execution plan needs to be more detailed and progressive, with each step having an expected output. The contributions of the Storyteller role in each step need to be clearly defined.

Question or Task:
Write a 10 words story about a girl.

Selected Roles List:

```
{  
    "name": "Storyteller",  
    "description": "An expert in creating engaging and concise stories within a given word limit.",  
    "tools": [],  
    "suggestions": "Focus on the main character, which is a girl in this case. Try to convey a complete story within the given word limit. Use your creativity.",  
    "prompt": "You are an expert in creating engaging and concise stories within a given word limit, named Storyteller. Your goal is to write a 10-word stor  
}
```

Created Roles List:

```
{  
    "name": "Language Expert",  
    "description": "An expert in summarizing and presenting the results in a clear and understandable manner.",  
    "tools": [],  
    "suggestions": "Summarize the final 10-word story in a clear and engaging manner.",  
    "prompt": "You are an expert in summarizing and presenting the results in a clear and understandable manner, named Language Expert. Your goal is to prov  
}
```

Ca

Ston

sk-dhCpF6IBXMsL0vfrd0T3B1bkJStPwgJfTepJrKrVcVxV

Tasks and progress

Step	Role	Status
1	Question/Task	Completed task(s)
2	Manager	Completed task(s)
3	user	Completed task(s)
4	user 2	Completed task(s)
5	user 3	Completed task(s)
6	user 4	Completed task(s)

LLM Agents - Demo

Agents

- Question/Task
- Manager
- user

Auto Agents Chat

7. Girl bffrie Show More

user >>> Final Output

The task was to brainstorm ideas for a 10-word story about a girl and then draft the story based on the brainstormed ideas. The Storyteller agent brainstormed and listed down potential story ideas. Here are the brainstormed ideas and their condensed 10-word stories:

- Idea: Girl discovers magical world in her grandmother's attic.
Story: "In grandma's attic, a girl discovers a magical world."
- Idea: Young girl saves village with her unique singin

user >Storyteller_Action Substep:
>Review and refine the 10-word stories about a girl.
>Subresponse:
Refined 10-word stories:

- "Girl discovers magical world in grandma's attic."
- "Unique singing voice of girl saves village."
- "Lost puppy found, girl becomes best friend."
- "Brave girl fights dragon, rescues brother."
- "Ocean-cleaning machine invented by girl."
- "Girl time-travels, meets ancestors."
- "Girl befriends lonely ghost in new house."
- "Girl"

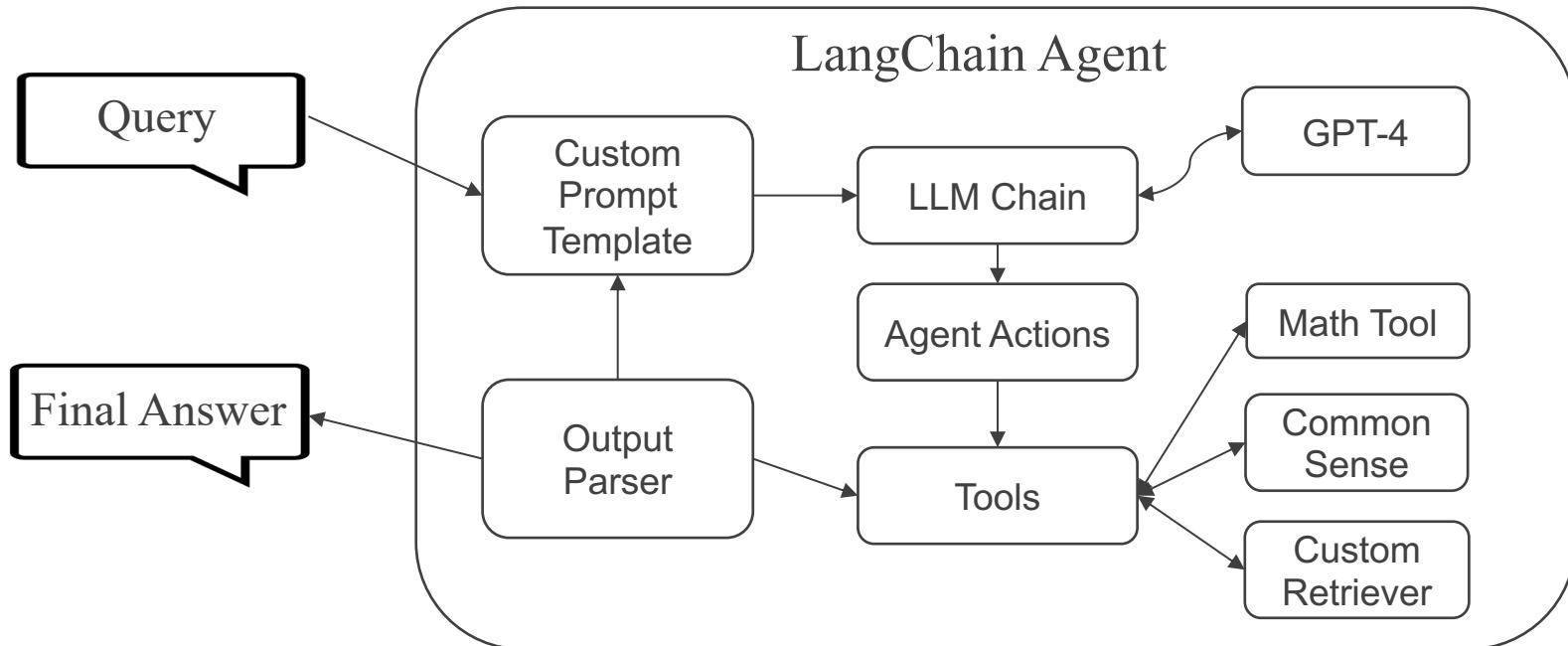
Show More

Tasks and progress

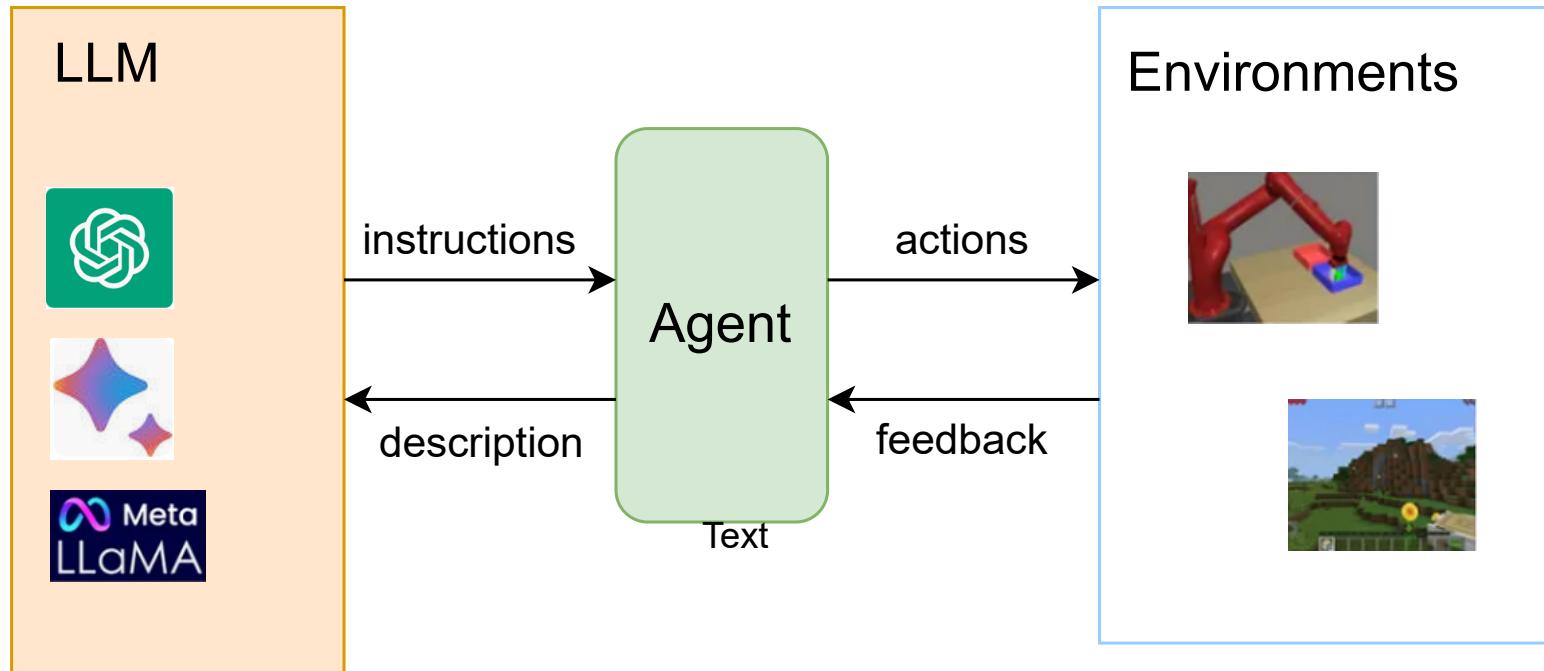
write a 10 words story about a girl

1	Question/Task Completed task(s)
2	Manager Completed task(s)
3	user Completed task(s)
4	user 2 Completed task(s)
5	user 3 Completed task(s)

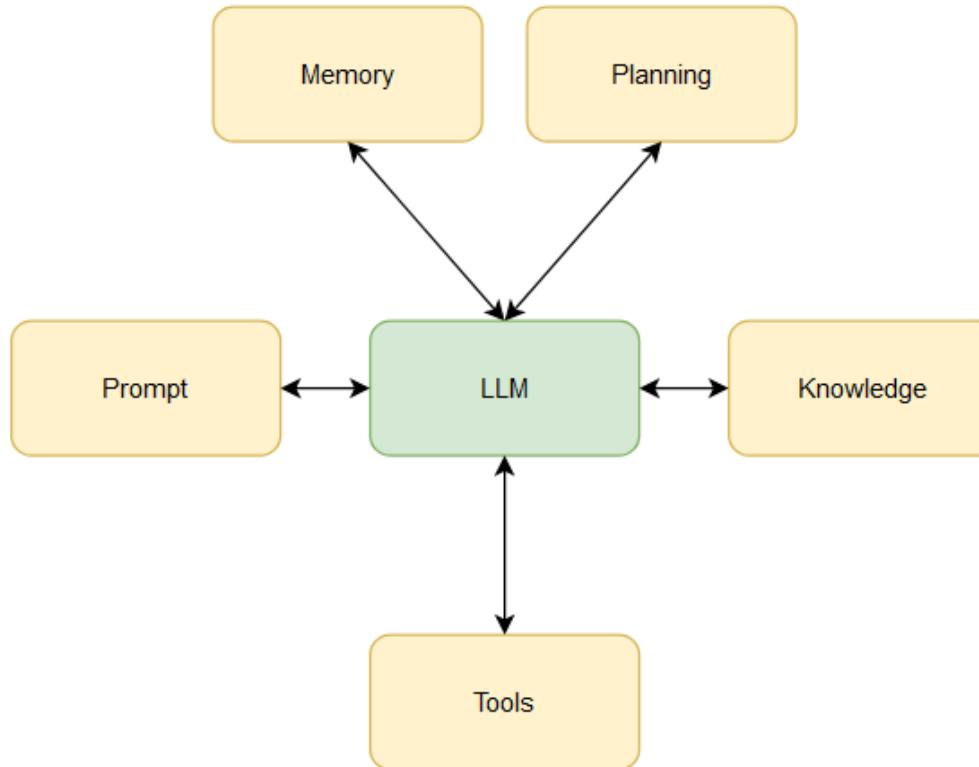
What is LLM Agent?



LLM Agents – General Framework



Components of LLM Agents



LLM Agents

<https://huggingface.co/spaces/LinkSoul/AutoAgents>

Future Directions

- Explainable AI
- Multimodal Reasoning
- Logical Reasoning in Multimodal AI Tasks
- AI Ethics