

Natural Language Processing And Generative AI



Dr. Shailesh Sivan
Principal AI Architect, Laennec AI
Assistant Professor, CUSAT(on leave)
shaileshsivan@gmail.com

What is NLP??

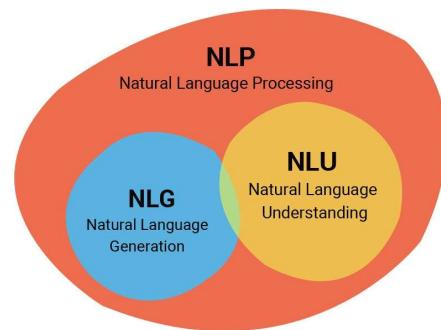
Natural language processing (NLP)

- A branch of Artificial Intelligence that gives machines the ability to understand natural human language.

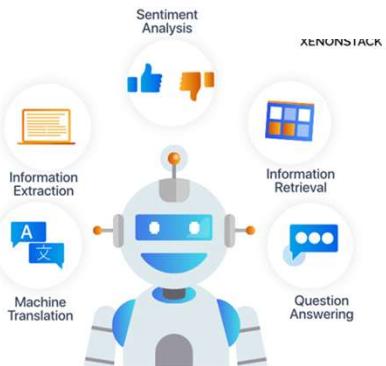
WHY?

*billions of text data being generated every day
and most of them are unstructured.*

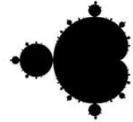
What is NLP??



Applications of **Natural Language Processing** in Different Domains



Python Libraries for NLP

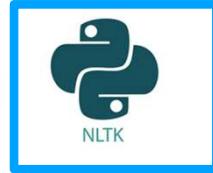


TextBlob



topic modelling for humans

spaCy



Common Terminologies

Natural Language Processing Pipeline



Sentence Segmentation

Sentence Segmentation

Divides the entire paragraph into different sentences for better understanding.



- Hello world.
- This blog post is about sentence segmentation.
- It is not always easy to determine the end of a sentence.
- One difficulty of segmentation is periods that do not mark the end of a sentence.
- An ex. is abbreviations.

Tokenization

breaks the sentence into separate words or tokens to understand the context of the text.



Stemming and Lemmatization

Stemming

- cuts off prefixes and suffixes to reduce a word to its root form called the **stem**, may not be a valid word in the language

Lemmatization

- Reducing words to their base or dictionary form, known as the **lemma**.

Stemming vs Lemmatization



Stop Words

Stop Words

- a
- of
- on
- I
- for
- with
- the
- at
- from
- in
- to

When was the first computer invented?

How do I install a hard disk drive?

How do I use Adobe Photoshop?

Where can I learn more about computers?

How to download a video from YouTube

What is a special character?

How do I clear my Internet browser history?

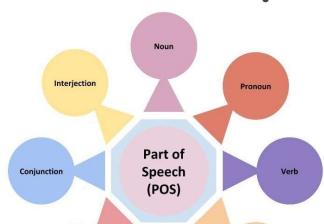
How do you split the screen in Windows?

How do I remove the keys on a keyboard?

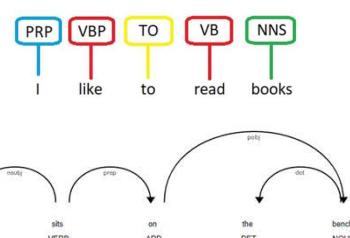
How do I install a hard disk drive?

words that are so widely used that they carry very little useful information.

Part Of Speech(POS)



POS Tagging



Label words in a text with their grammatical category, such as noun, verb, or adjective

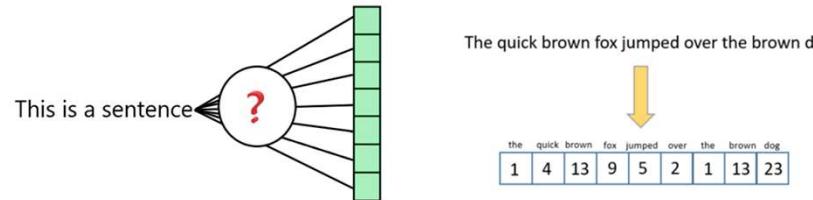
Named Entity Recognition

contentSkip to site indexPoliticsSubscribeLog InSubscribeLog InToday's PaperAdvertisementSupported ORG byF.B.I. Agent Peter Strzok PERSON , Who Criticized Trump PERSON in Texts, Is FiredImagePeter Strzok, a top F.B.I. GPE counterintelligence agent who was taken off the special counsel investigation after his disparaging texts about President Trump PERSON were uncovered, was fired. CreditT.J. Kirkpatrick PERSON for The New York TimesBy Adam Goldman ORG and Michael S. SchmidtAug PERSON , 13 CARDINAL , 2018WASHINGTON CARDINAL — Peter Strzok PERSON , the F.B.I. GPE senior counterintelligence agent who disparaged President Trump PERSON in inflammatory text messages and helped oversee the Hillary Clinton PERSON email and Russia GPE investigations, has been fired for violating bureau policies, Mr. Strzok PERSON 's lawyer said Monday DATE . Mr. Trump and his allies seized on the texts — exchanged during the 2016 DATE campaign with a former F.B.I. GPE lawyer, Lisa Page — in PERSON assailing the Russia GPE investigation as an illegitimate "witch hunt." Mr. Strzok PERSON , who rose over 20 years DATE at the F.B.I. GPE to become one of its most experienced counterintelligence agents, was a key figure in the early months DATE of the inquiry. Along with writing the texts, Mr. Strzok PERSON was accused of sending a highly sensitive search warrant to his personal email account. The F.B.I. GPE had been under immense political pressure by Mr. Trump PERSON to dismiss Mr. Strzok PERSON , who was removed last summer DATE from the staff of the special counsel, Robert S. Mueller III PERSON . The president has repeatedly denounced Mr. Strzok PERSON in posts on

Process of identifying and classifying key information (entities) such as names of people, organizations, locations, dates, and more within unstructured text.

Text Feature Extraction and Classification

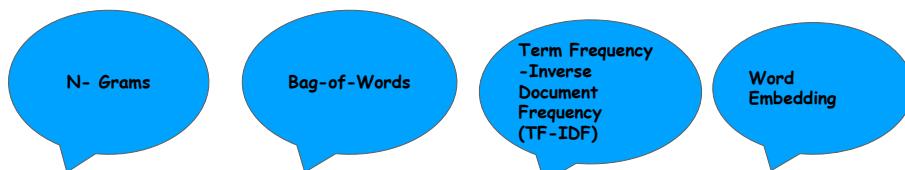
Text Features Extraction



Turning text into vectors that can be then fed to machine learning models in a classical way

Types

Machine learning algorithms cannot work with raw text directly; the text must be converted into numbers. Specifically, vectors of numbers.



N-grams

N-grams are the combination of multiple words used together.

N-grams with N=1 are called unigrams.

N=2 - bigrams,

N=3 - trigrams and so on..

This is Big Data AI Book

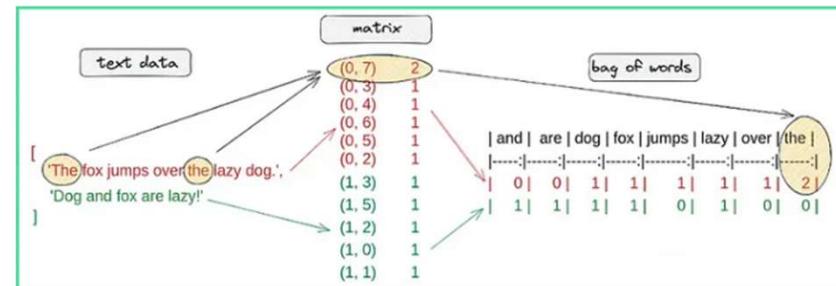
<i>Uni-Gram</i>	This	is	Big	Data	AI	Book
<i>Bi-Gram</i>	This is	Is Big	Big Data	Data AI	AI Book	
<i>Tri-Gram</i>	This is Big	Is Big Data	Big Data AI	Data AI Book		

Bag of Words (BoW)

- used to analyze text and documents based on **word count**.
- model does not account for word order within a document.

	about	bird	heard	is	the	word	you
About the bird, the bird, bird bird bird	1	5	0	0	2	0	0
You heard about the bird	1	1	1	0	1	0	1
The bird is the word	0	1	0	1	2	1	0

Bag of Words (BoW)



Bag of Words (BoW) Limitation

'The sky is blue and beautiful',
 'The king is old and the queen is
 beautiful',
 'Love this beautiful blue sky',
 'The beautiful queen and the old king']

	and	beautiful	blue	is	king	love	old	queen	sky	the	this
0	1	1	1	1	0	0	0	0	1	1	0
1	1	1	1	0	2	1	0	1	1	0	2
2	0	1	1	0	0	1	0	0	0	1	0
3	1	1	0	0	1	0	1	1	0	2	0

	beautiful	beautiful	beautiful	blue	blue	blue	king	king	love	love	love
0	1	0	0	1	1	0	0	0	0	0	0
1	1	0	0	0	0	1	1	0	0	1	1
2	1	1	0	1	0	0	1	0	0	0	0
3	1	0	1	0	0	1	0	0	1	0	0

Term Frequency - Inverse Document Frequency Matrix

Measures the importance of words in a document

Term frequency

- The number of times a word appears in a document
- Ratio of the count of a word's occurrence in a document and the number of words in the document

$$TF(t, d) = \frac{\text{Number of times term } t \text{ occurs in document } d}{\text{Total number of terms in document } d}$$

Inverse Document frequency

- Measure of how common or rare a word is in the entire corpus of documents
- IDF of word t represents the ratio of number of documents in the corpus with word t in them to the total number of documents in the corpus

$$IDF(t, d) = \log \left(\frac{\text{Total number of documents in corpus } N}{\text{Number of documents containing the term } t} \right)$$

$$TFIDF(t, d, N) = TF(t, d) \times IDF(t, d)$$

Term Document - Inverse Document Frequency Matrix

$$TF(t, d) = \frac{\text{Number of times term } t \text{ occurs in document } d}{\text{Total number of terms in document } d}$$

$$IDF(t, d) = \log \left(\frac{\text{Total number of documents in corpus } N}{\text{Number of documents containing the term } t} \right)$$

$$TFIDF(t, d, N) = TF(t, d) \times IDF(t, d)$$

Actual text

I read the svm algorithm article in dataaspirant blog

I read the randomforest algorithm article in dataaspirant blog

Preprocessed text

read svm algorithm article dataaspirant blog

read randomforest algorithm article dataaspirant blog

TF-IDF Calculation Example

Words	Count		Term Frequency (TF)		Inverse Document Frequency (IDF)	TF * IDF	
	Document 1	Document 2	Document 1	Document 2		Document 1	Document 2
read	1	1	0.17	0.17	0	0	0
svm	1	0	0.17	0	0.3	0.05	0
algorithm	1	1	0.17	0.17	0	0	0
article	1	1	0.17	0.17	0	0	0
dataaspirant	1	1	0.17	0.17	0	0	0
blog	1	1	0.17	0.17	0	0	0
randomforest	0	1	0.17	0.17	0.3	0	0.05

Term Document - Inverse Document Frequency Matrix

Actual text

Petrol cars are cheaper than diesel cars

Diesel is cheaper than petrol

Preprocessed text

Petrol cars cheaper diesel cars

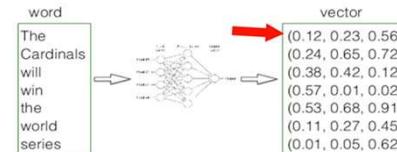
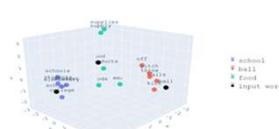
Diesel cheaper petrol

Words	TF = Frequency(word)		IDF = $\log_e((d+1)/(df(t)+1)) + 1$	TF-IDF	
	Document 1	Document 2		Document 1	Document 2
car	2	0	$\log_e(3/2) + 1 \Rightarrow 1.405465083$	2.8109302	0
cheaper	1	1	$\log_e(3/3) + 1 \Rightarrow 1$	1	1
diesel	1	1	$\log_e(3/3) + 1 \Rightarrow 1$	1	1
petrol	1	1	$\log_e(3/3) + 1 \Rightarrow 1$	1	1

In the first document the term "cars" is the most relevant term as it has the highest tf-idf value
 In the second document most of the terms have the same tf-idf value and have equal relevance.

Word Embedding

- A technique that represents words as numbers so that computers can process them.
- Capture semantic and syntactic relationship between texts and represent them in the form of vectors
- IDEA - Similar words will have a minimum distance between their vectors.



Continuous Bag of Words (CBOW).

Word2Vec
neural approach
for generating
word
embeddings

Skip-grams

Continuous Bag of Words Vs Skip-gram

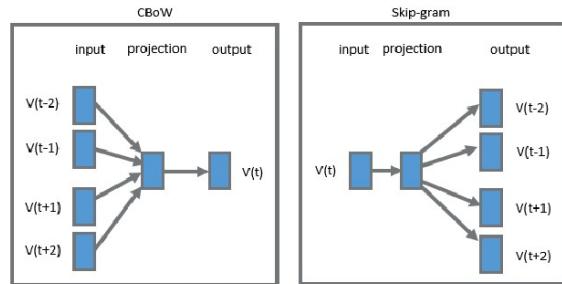
The sun is bright and the weather is pleasant.		
Windowsize = 2 (2 words to the left and 2 words to the right)		
Position	Target Word	Context Words
1	the	[sun, is]
2	sun	[the, is, bright]
3	is	[the, sun, bright, and]
4	bright	[sun, is, and, the]
5	and	[is, bright, the, weather]
6	the	[bright, and, weather, is]
7	weather	[and, the, is, pleasant]
8	is	[the, weather, pleasant]
9	pleasant	[weather, is]

Continuous Bag of Words Vs Skip-gram

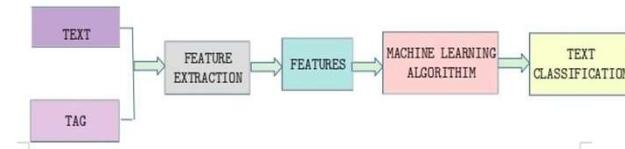
Neural network-based algorithm

CBOW : Predicts a target word given its surrounding context words.

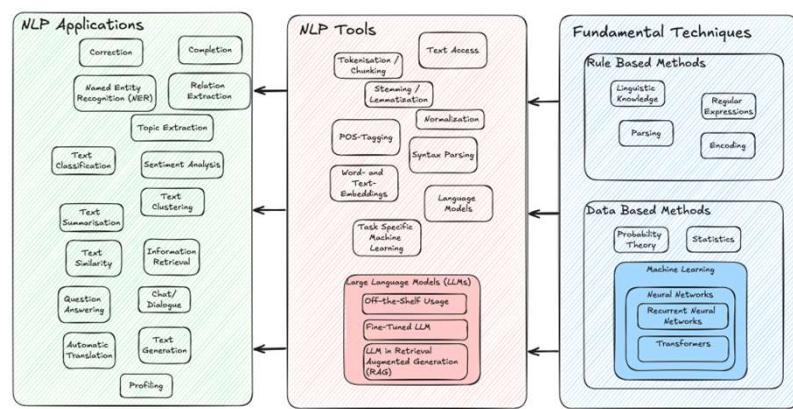
Skip-gram : Predicts surrounding words based on a specific word called the "target word."



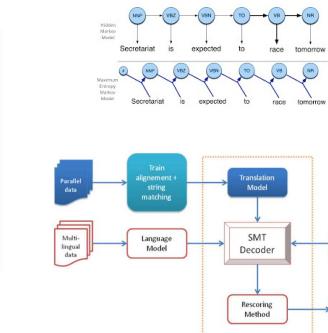
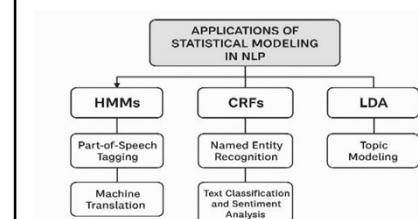
Text Classification -Pipeline



Outline Of NLP Taxonomy



Statistical NLP Models



What Is Machine Learning?

"Learning is any process by which a system improves performance from experience."

- Herbert Simon

Definition by Tom Mitchell (1998):

Machine Learning is the study of algorithms that

- improve their performance P
- at some task T
- with experience E.

A well-defined learning task is given by $\langle P, T, E \rangle$

Machine Learning is a process where a computer uses data to learn **hidden patterns** through an **ML algorithm**, and the output of this learning is an **ML model** that can make **predictions** on new data.

Machine Learning

Training: Extract patterns from data
ML Algorithms (Icon: Database, Chart, Text)

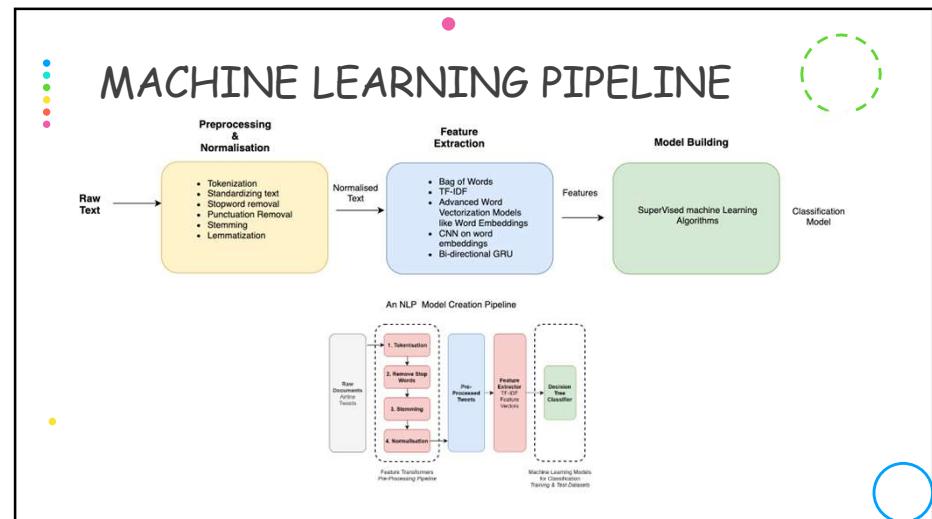
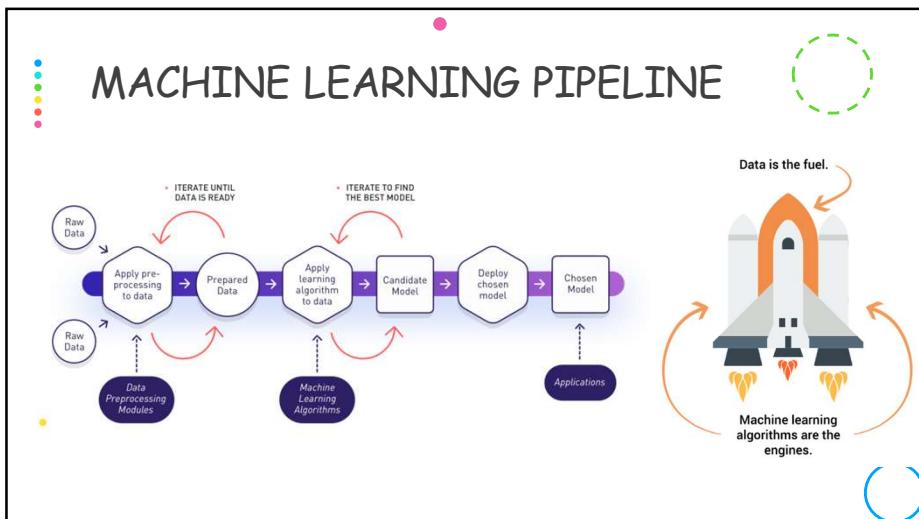
Evaluating: Use patterns to predict results
ML Models (Icon: Bar Chart, Laptop)

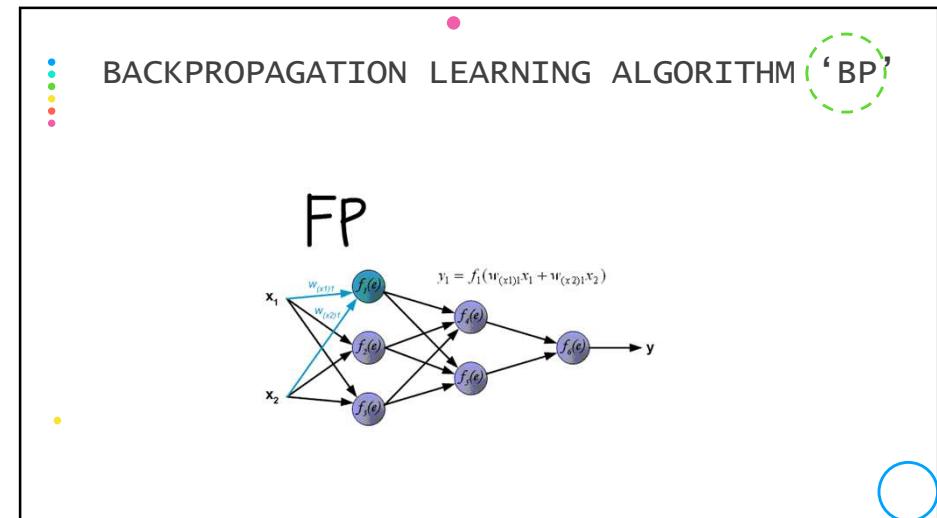
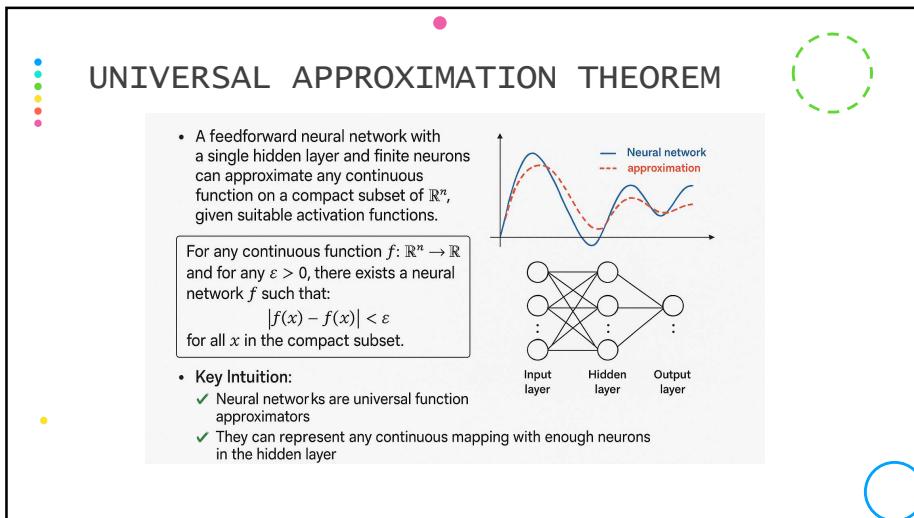
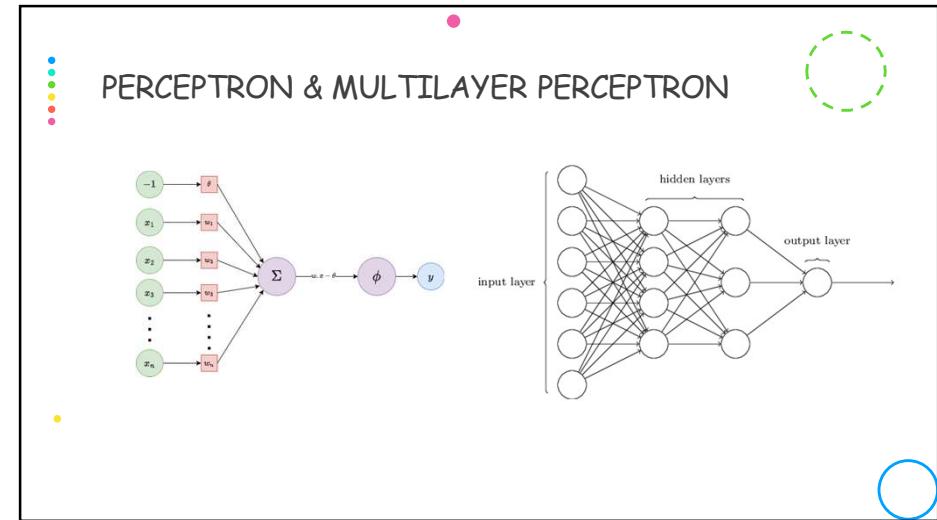
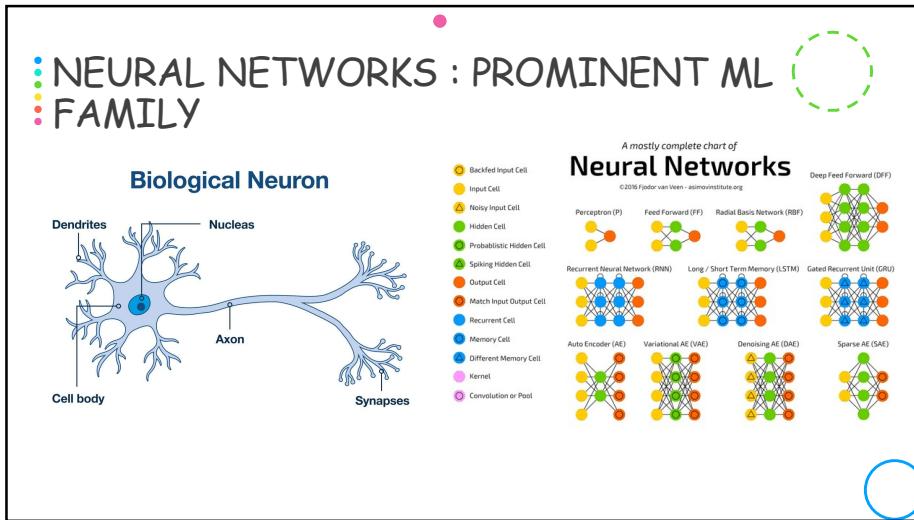
machine learning branches into:

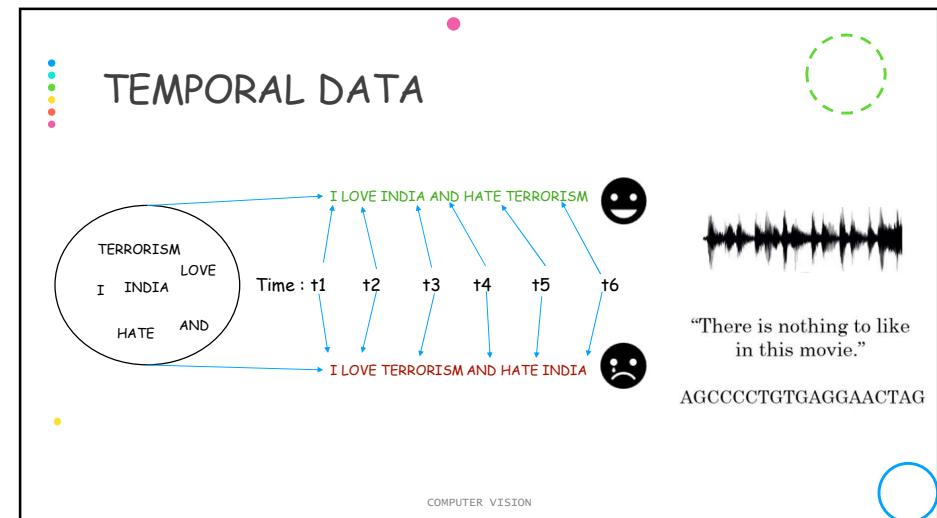
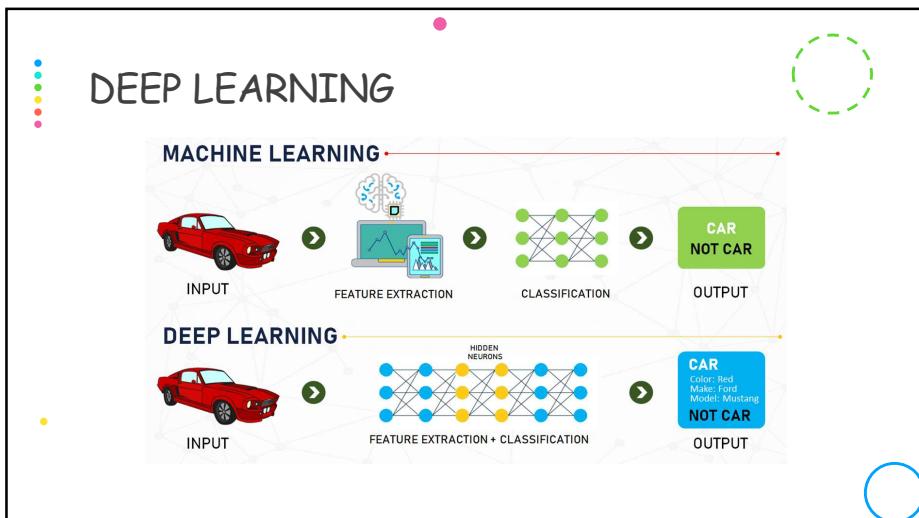
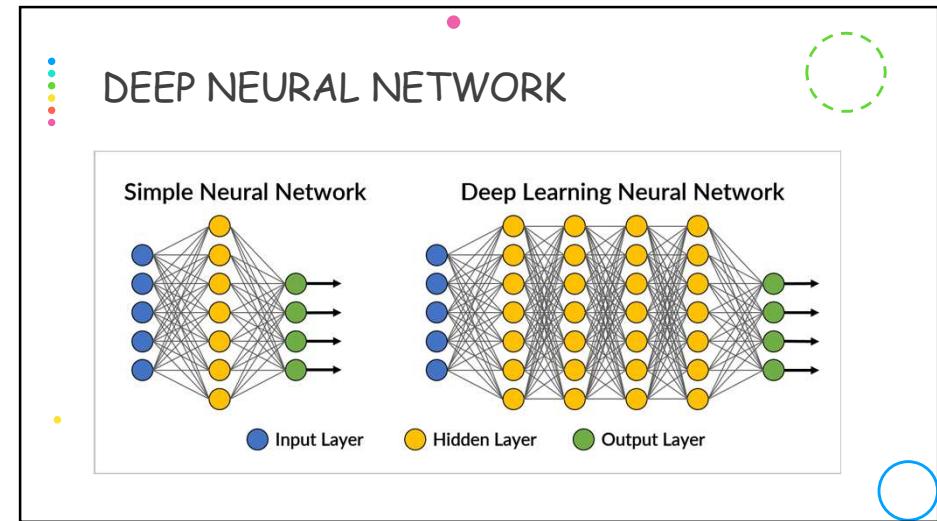
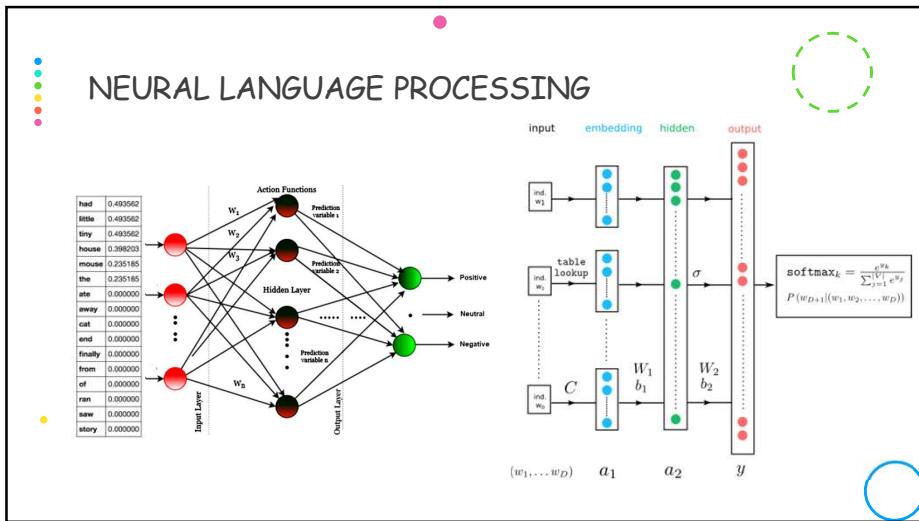
- unsupervised learning (Icon: Scatter Plot)
- supervised learning (Icon: Classification Diagram)
- reinforcement learning (Icon: Reinforcement Diagram)

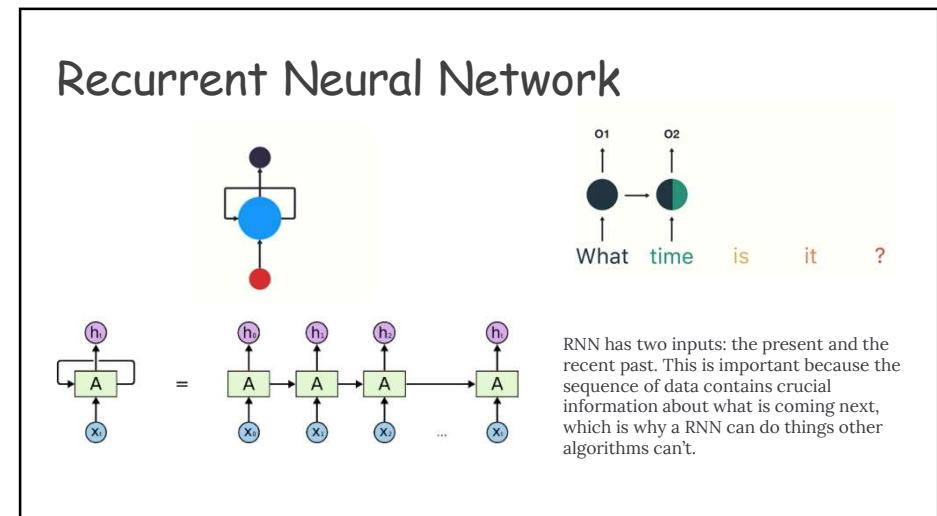
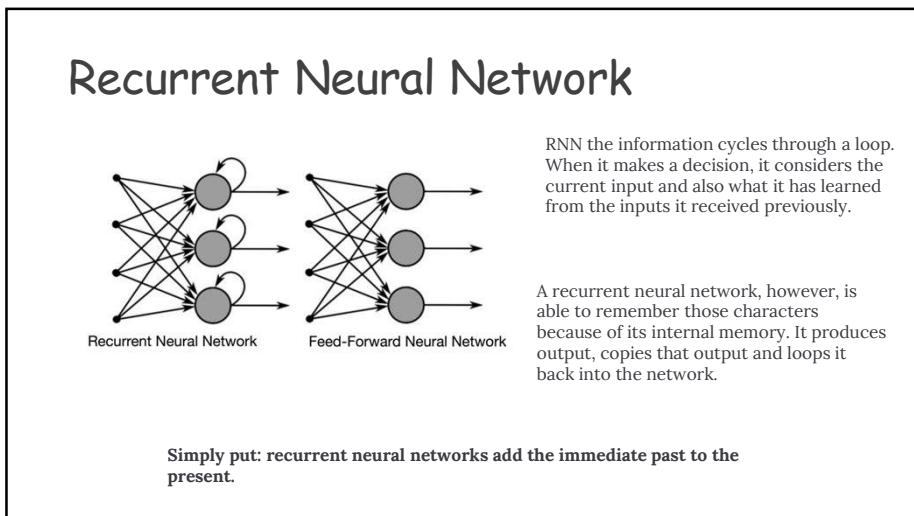
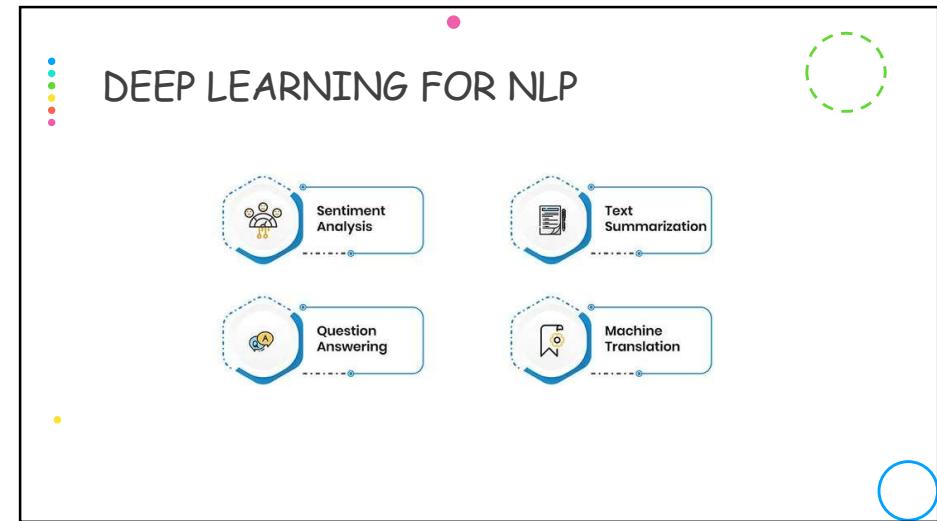
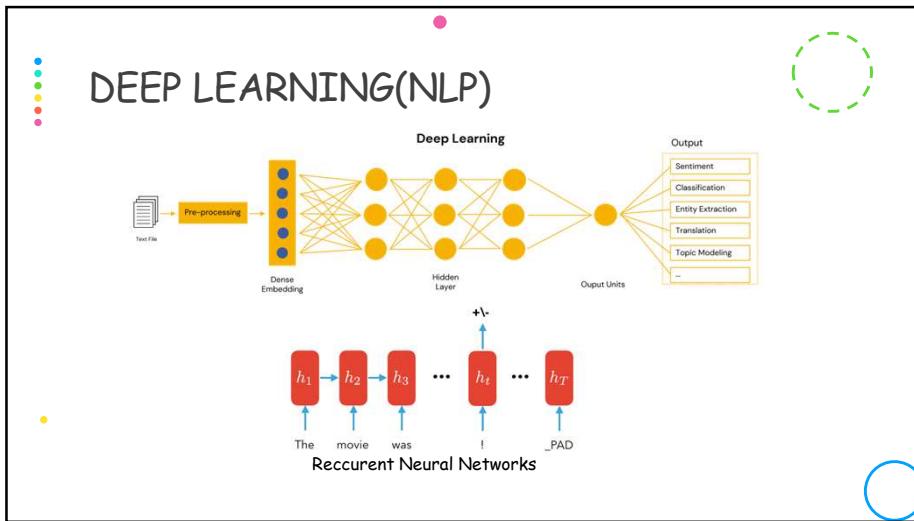
Problems with Uncertainty: Data with hidden patterns

Approximating the patterns to Generalization: Minimum error to predict unseen data

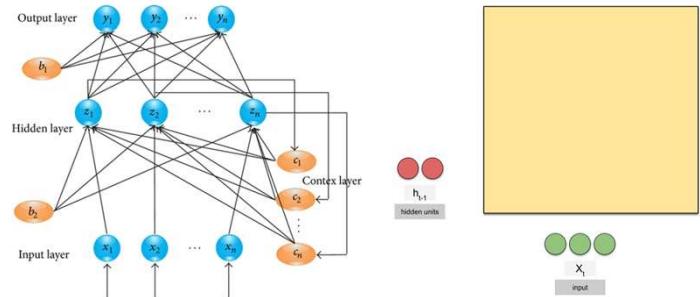




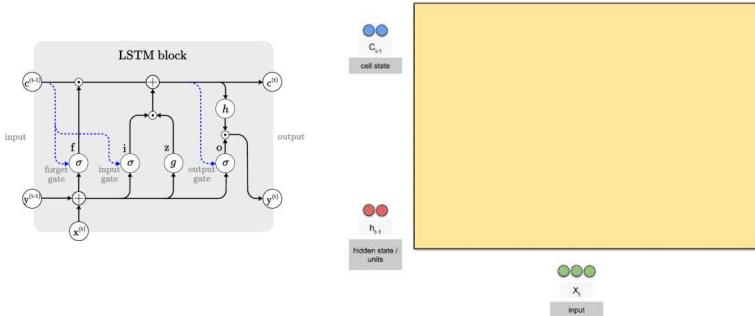




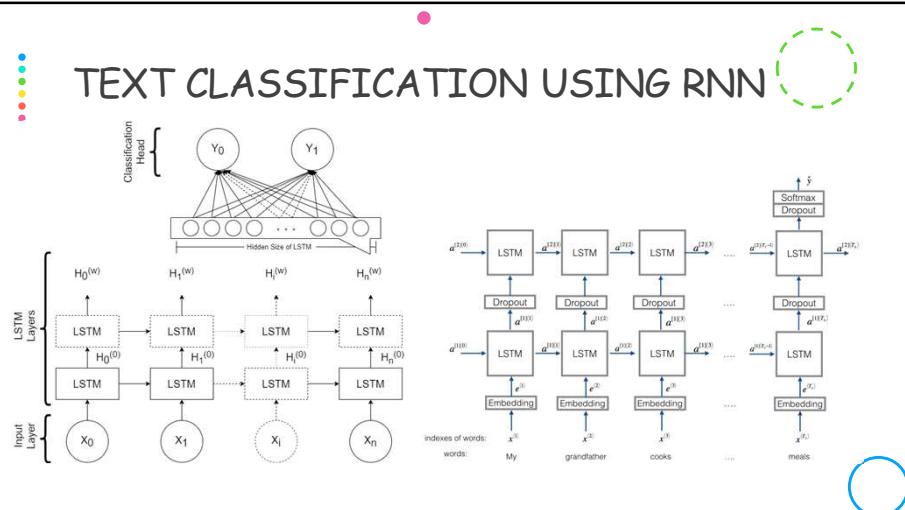
Recurrent Neural Network



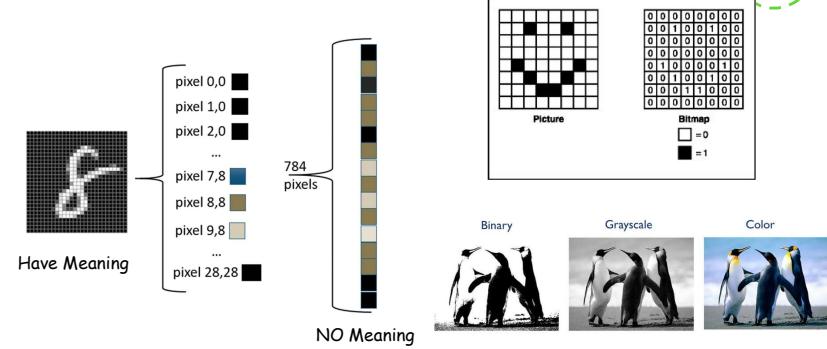
Long Short Term Memory

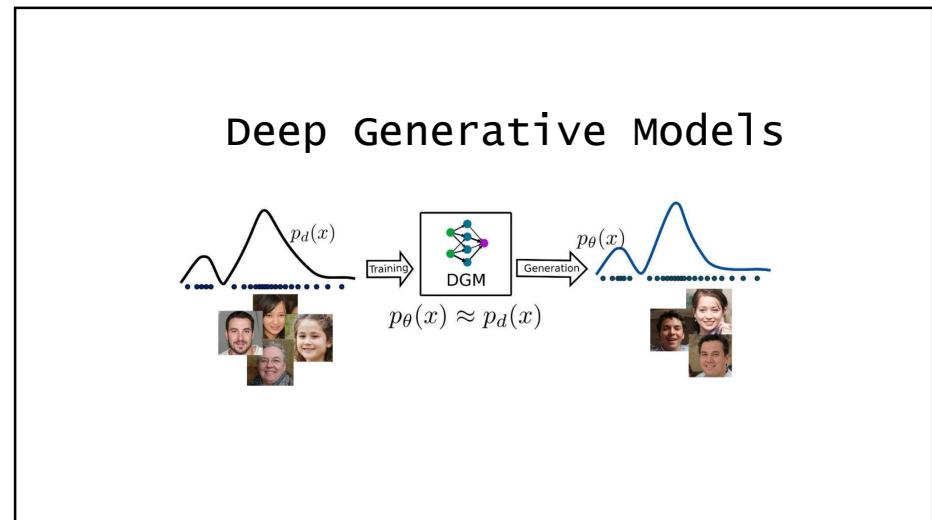
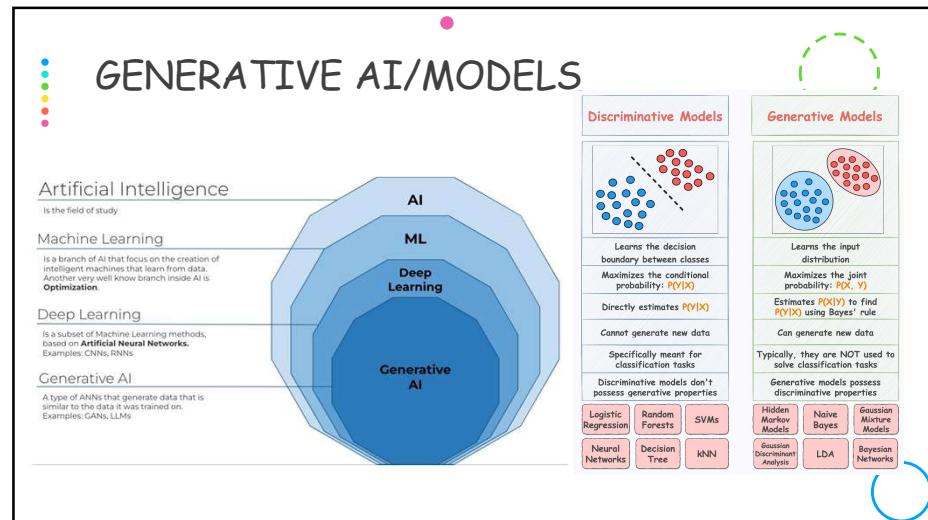
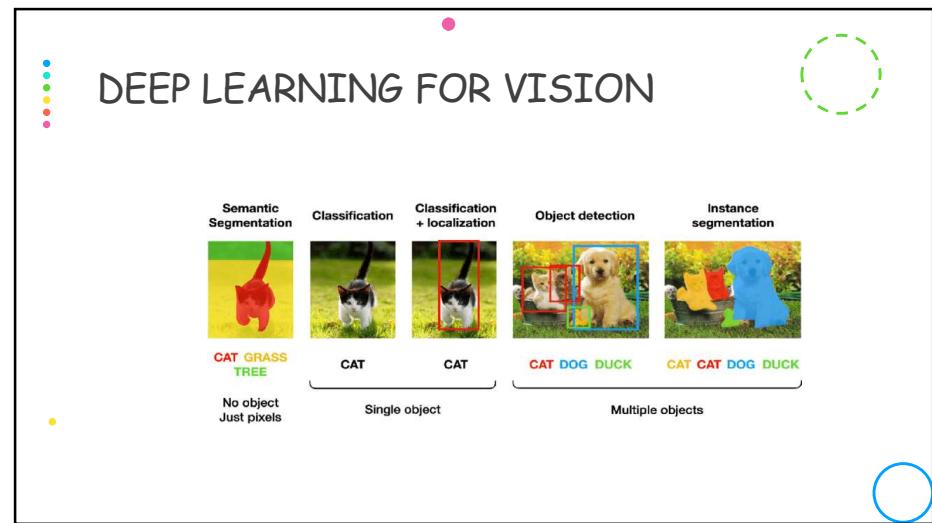
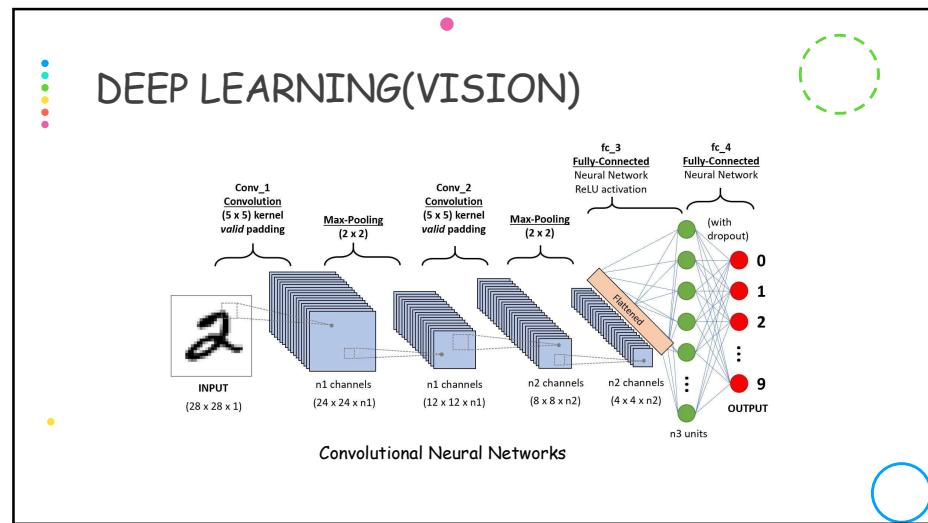


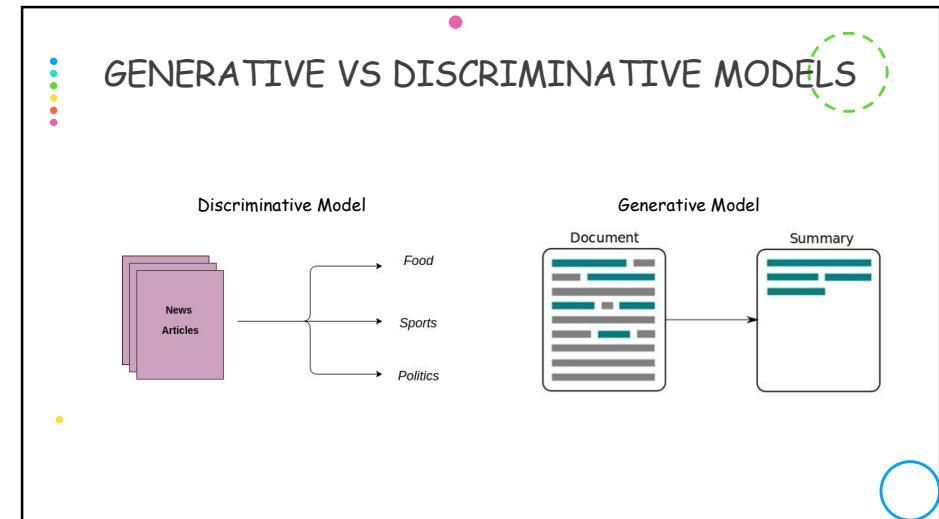
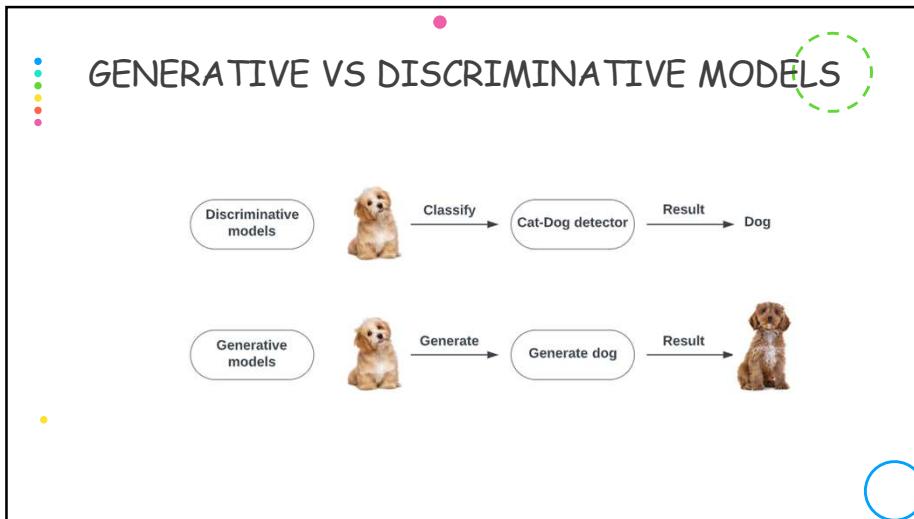
TEXT CLASSIFICATION USING RNN



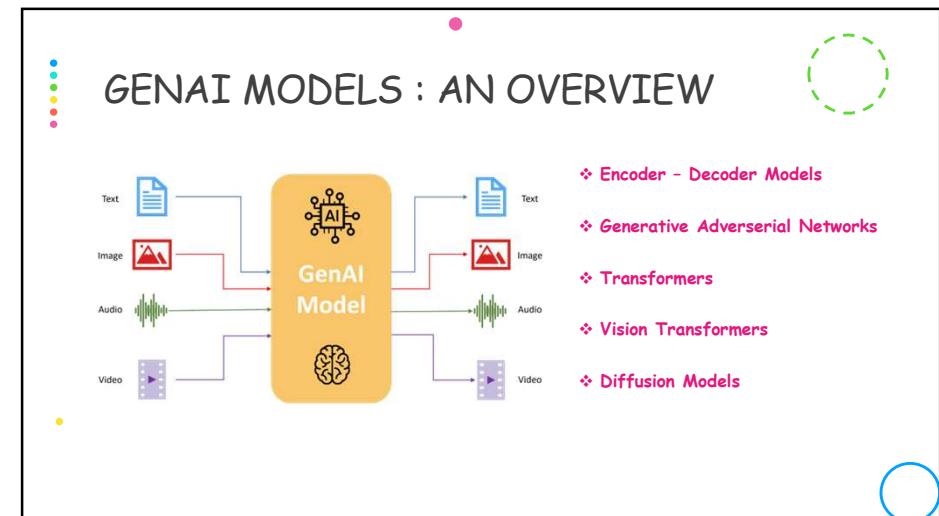
SPATIAL DATA

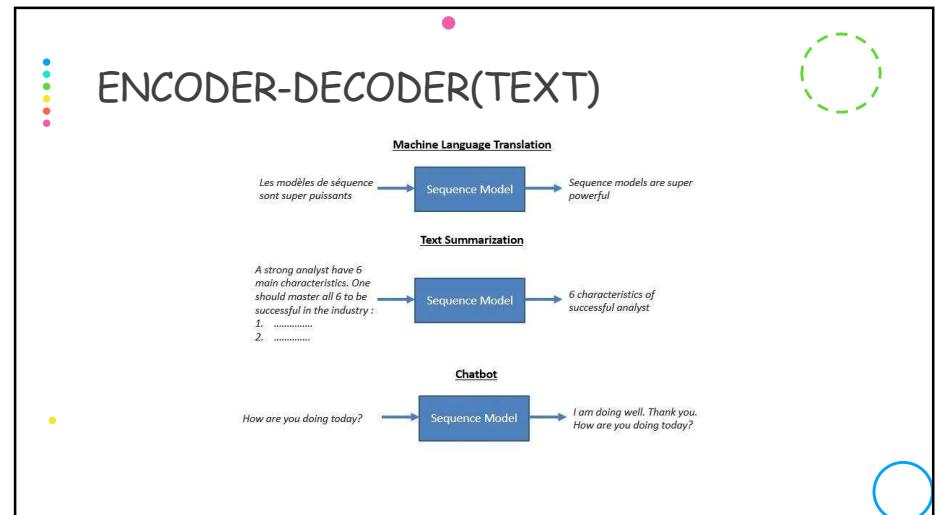
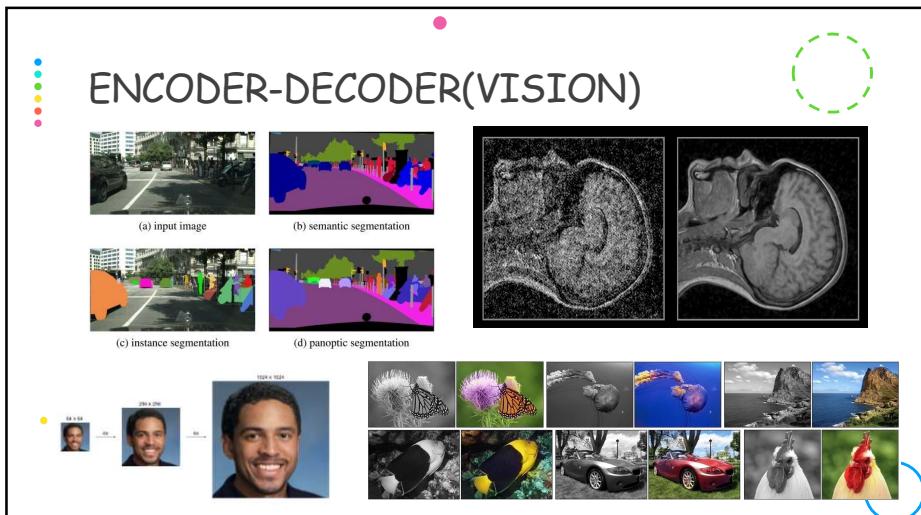
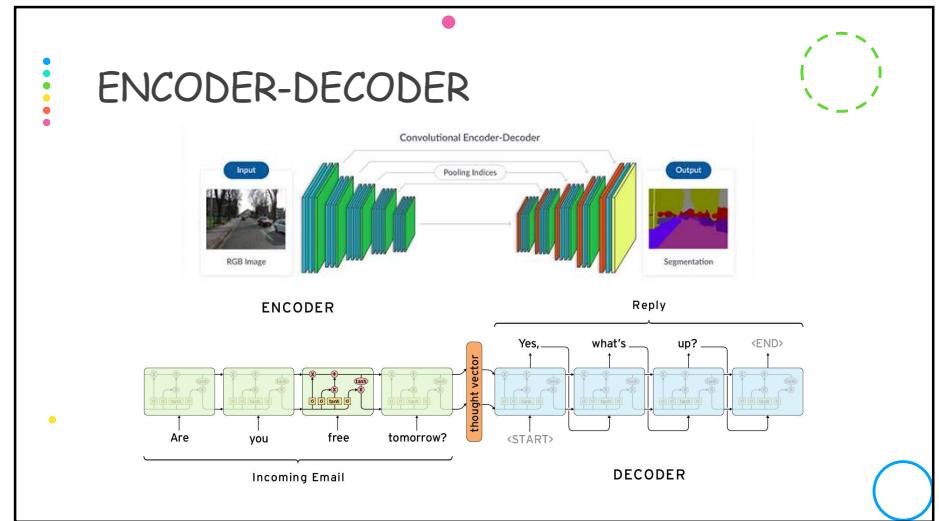
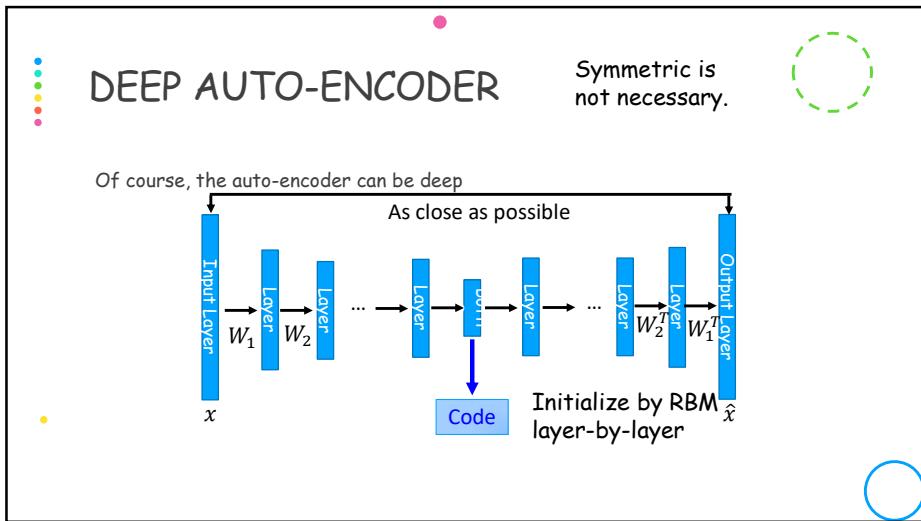


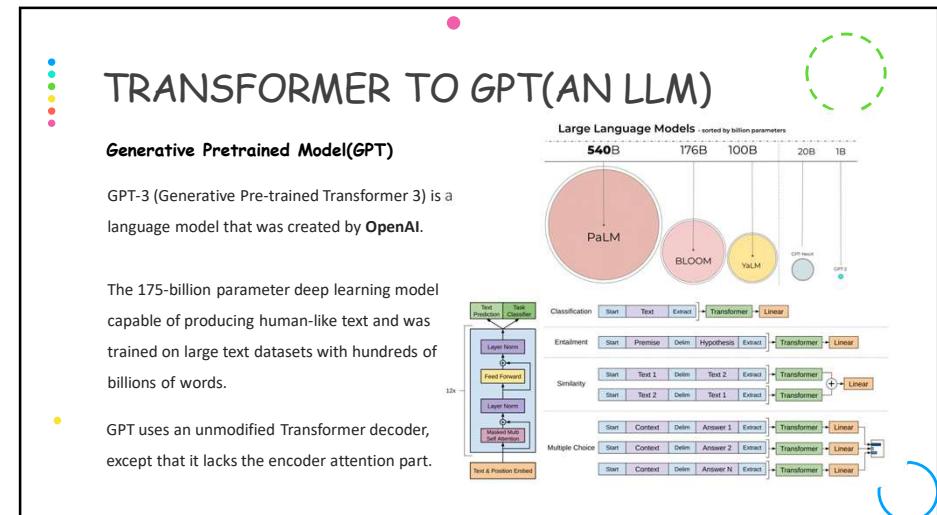
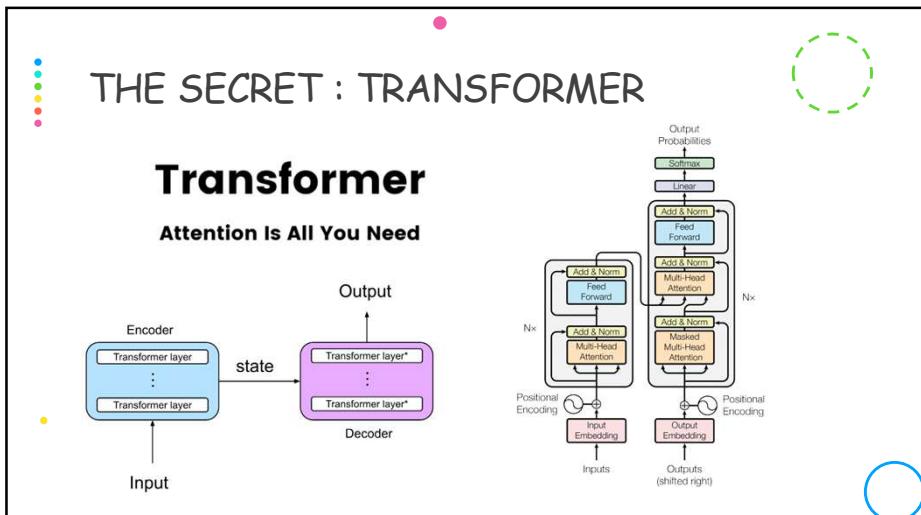
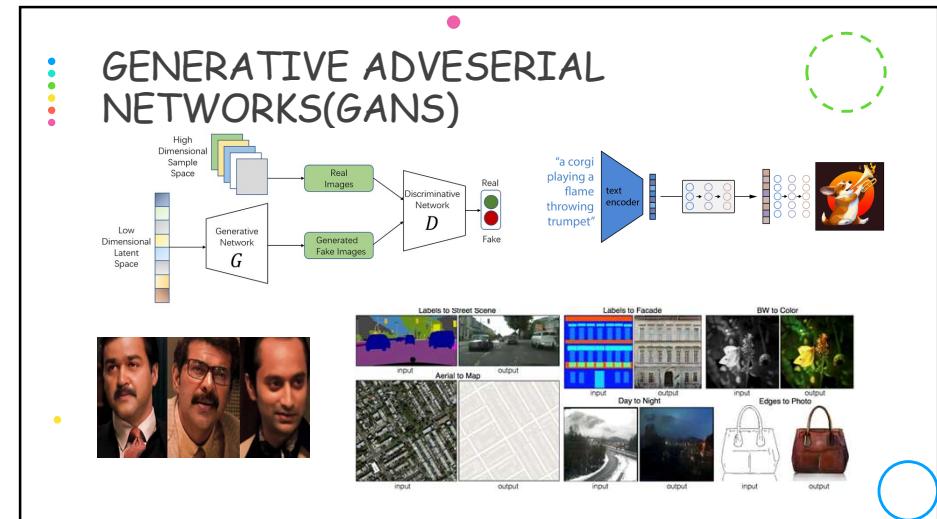
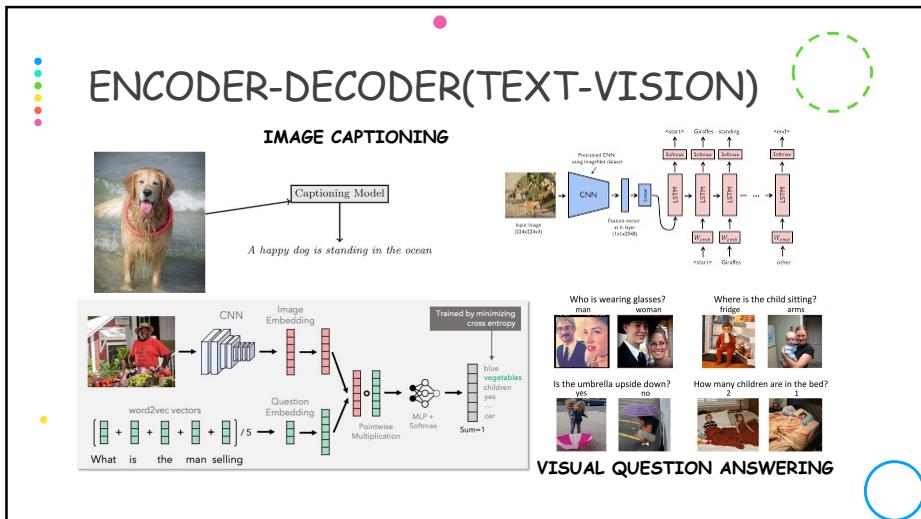


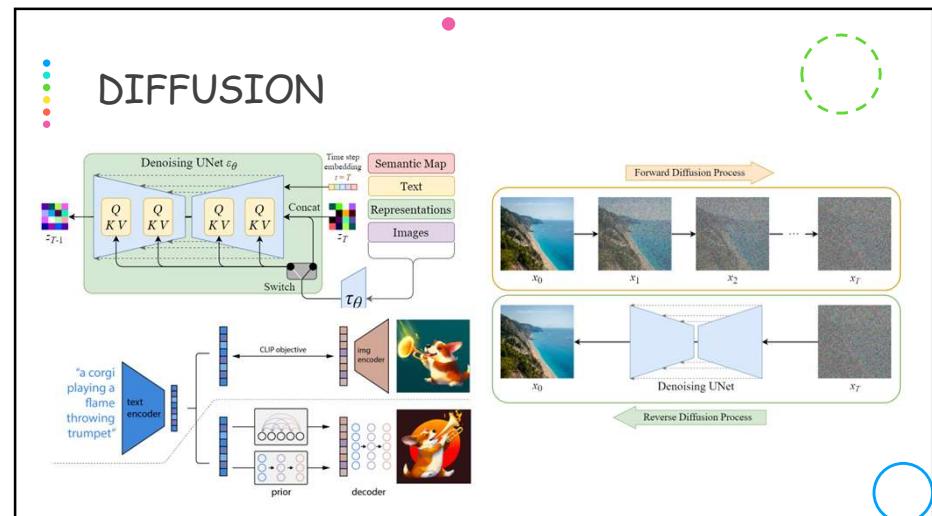
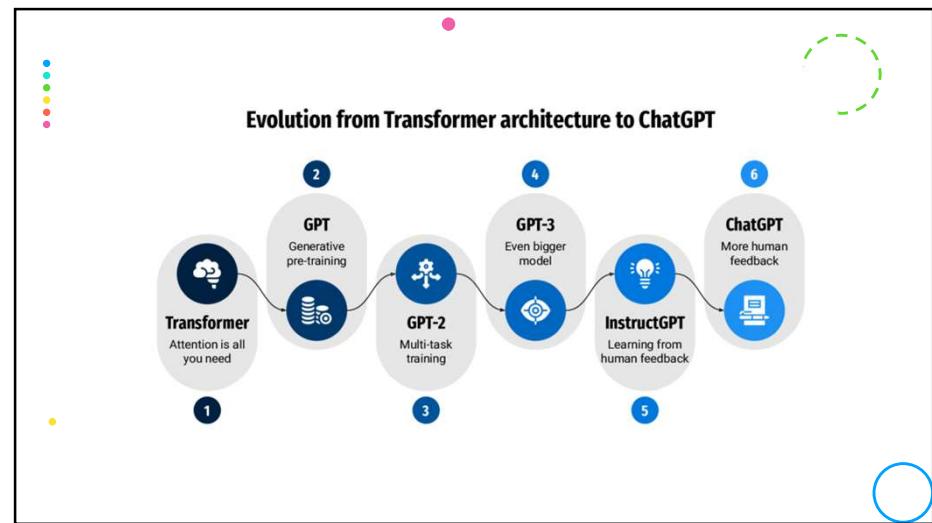
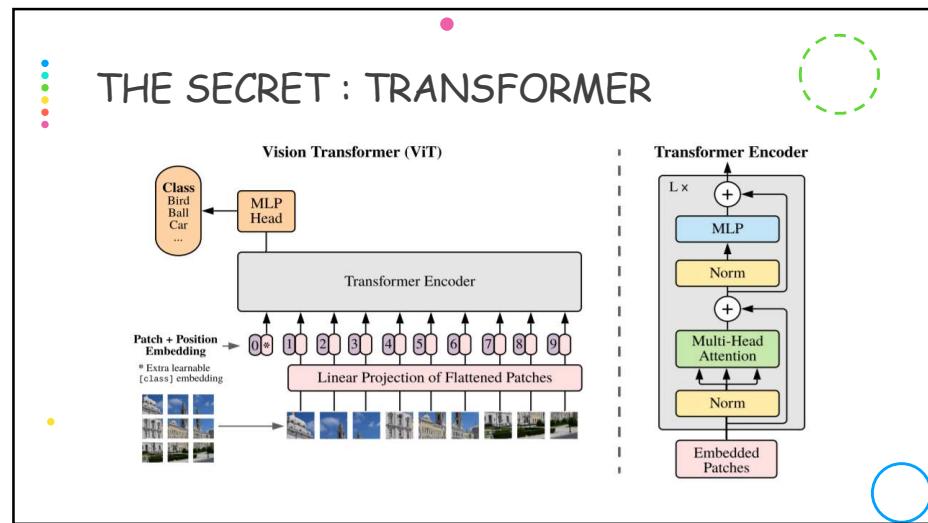


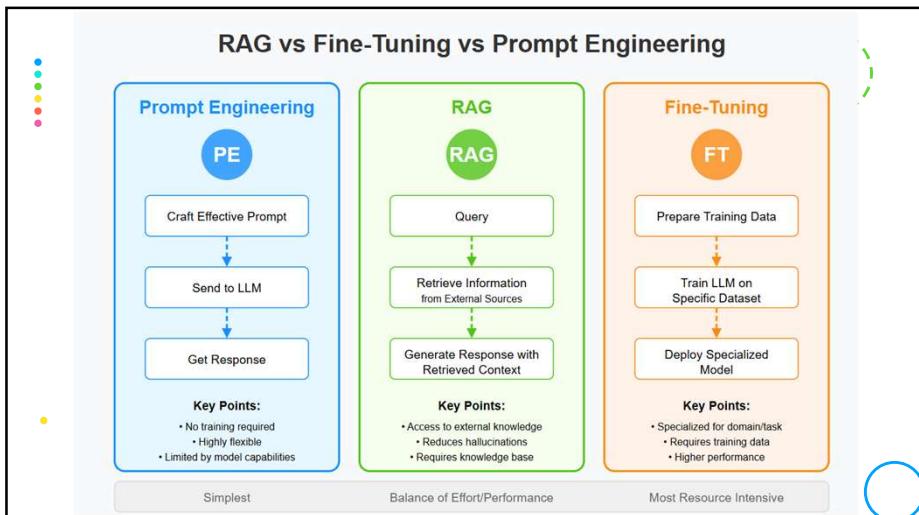
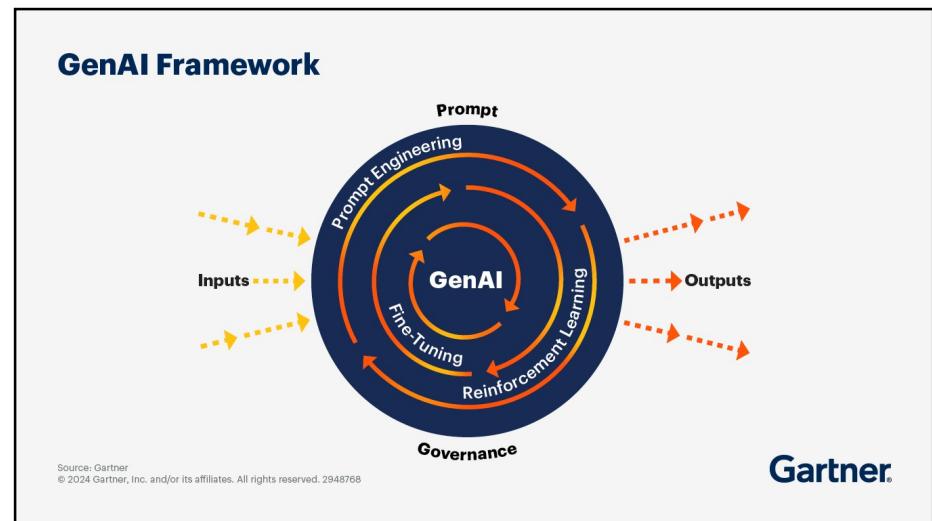
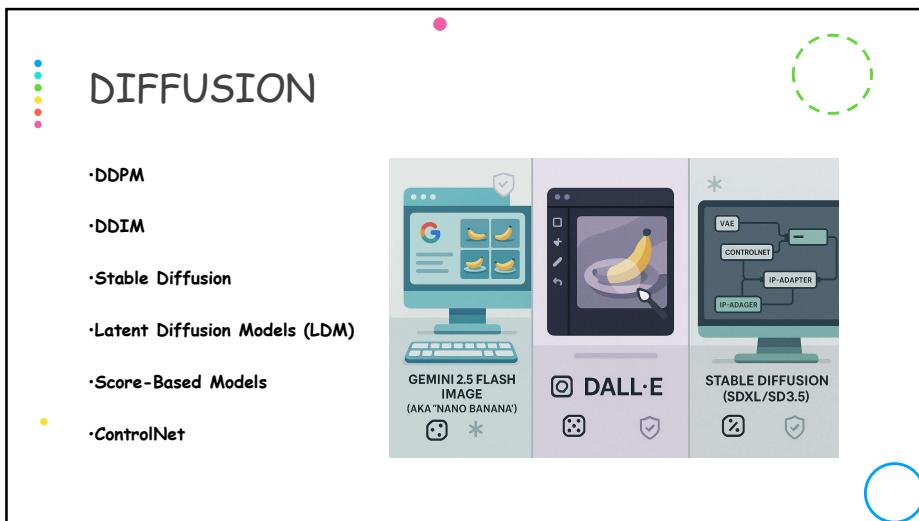
HOW GENERATIVE MODELS WORKS?











POPULAR GEN AI PRODUCTS

Model	Organization	Use cases
GPT-4	OpenAI	Content generation, code generation, creative writing, question answering
ChatGPT	OpenAI	Dialogue generation, chatbot development
DALL-E 2	OpenAI	Product design, concept art, image editing
Whisper	OpenAI	Transcription, translation, accessibility
Megatron-Turing NLG	Meta	Content generation, code generation, creative writing, question answering
Jurassic-1 Jumbo	Meta	Factual language modeling, topic summarization, question answering, translation
Flan-T5	Meta	Code generation
Bard	Google AI	Content generation, code generation, creative writing, question answering
PaLM	Google AI	Factual language modeling, topic summarization, question answering, translation
Meena	Google AI	Dialogue generation, chatbot development
Turing NLG	Microsoft	Content generation, code generation, creative writing, question answering
CodeGPT	Microsoft	Code generation
DialoGPT	Microsoft	Dialogue generation, chatbot development
Copilot	GitHub	Code generation suggestions
Midjourney	Midjourney	Art and design, product design, education
LLaMA	Meta	Content generation, code generation, creative writing, question answering

Thank You!



Your Queries Please!!!