

**POST GRADUATE
PROGRAM IN
GENERATIVE AI
AND ML**

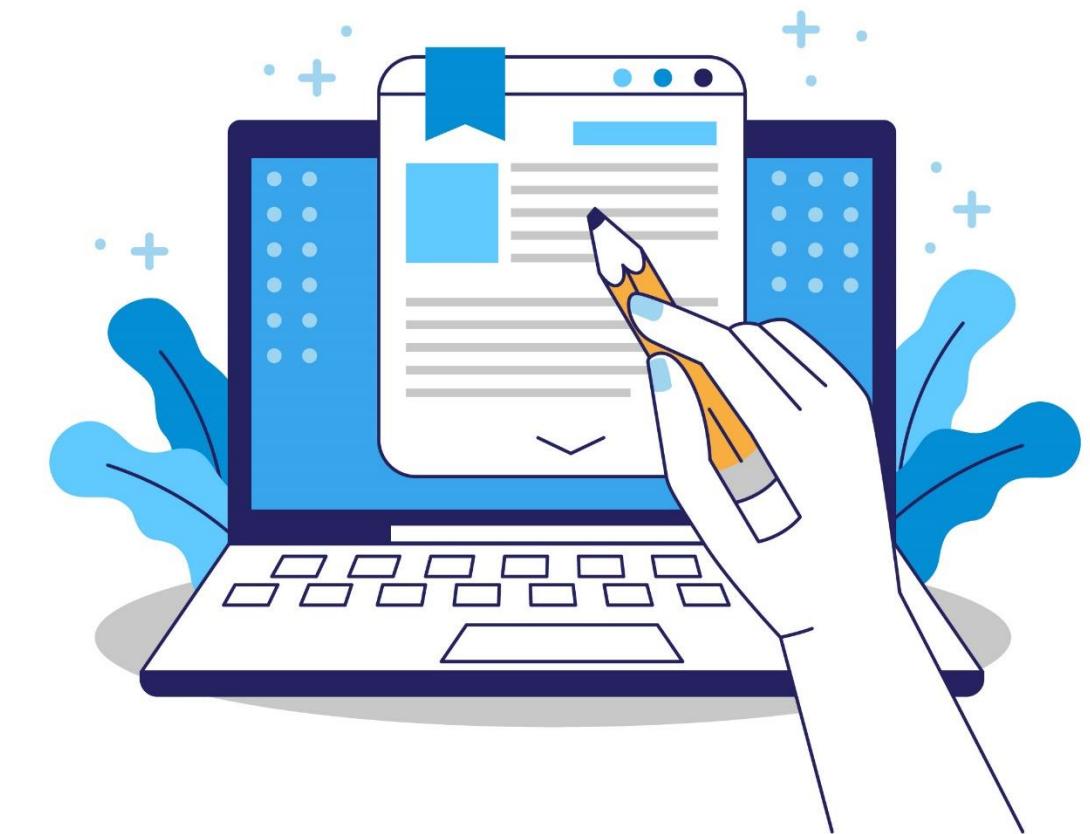
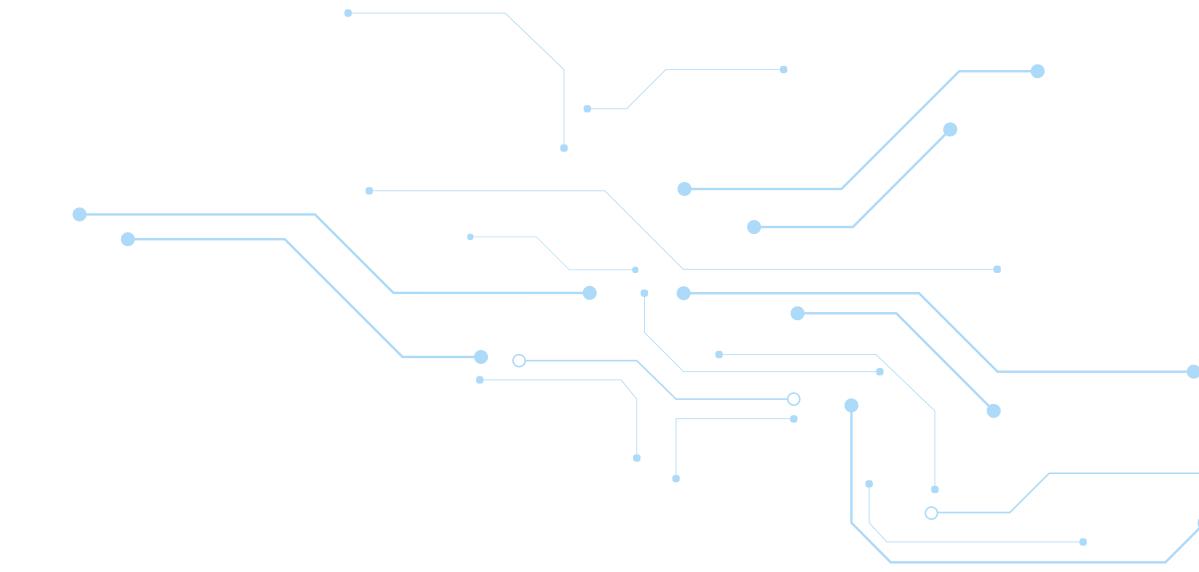
**Natural Language
Processing**



Speech and Multimodal NLP

Topics

- e! Fundamentals of Speech Recognition
- e! Automatic Speech Recognition (ASR) Models (e.g., Whisper, DeepSpeech)
- e! Challenges in Speech Processing
- e! TTS Applications in NLP
- e! Demonstration
- e! Introduction to Multimodal Learning
- e! Vision-Language Models (CLIP, BLIP)
- e! Integrating Text and Speech in NLP Applications
- e! Generative Models for Multimodal Content
- e! Model Compression (Quantization, Pruning, Knowledge Distillation)
- e! Real-Time NLP for Edge Devices
- e! Explainability in Deep NLP Models
- e! Emerging Trends (GPT-5, HyperNetworks, Autonomous AI Agents)



Learning Objectives

By the end of this lesson, you will be able to:

- e! Describe how speech and text are processed in NLP using ASR, TTS, and tokenization techniques
- e! Compare ASR models (e.g., Whisper, DeepSpeech) and vision-language models (e.g., CLIP, BLIP)
- e! Analyze challenges in speech processing and real-time NLP for edge deployment
- e! Apply multimodal and generative models to integrate and generate text, audio, and images
- e! Evaluate the use of model compression, explainability, and ethical principles in NLP systems
- e! Summarize emerging trends in NLP such as GPT-5, HyperNetworks, and autonomous AI agents



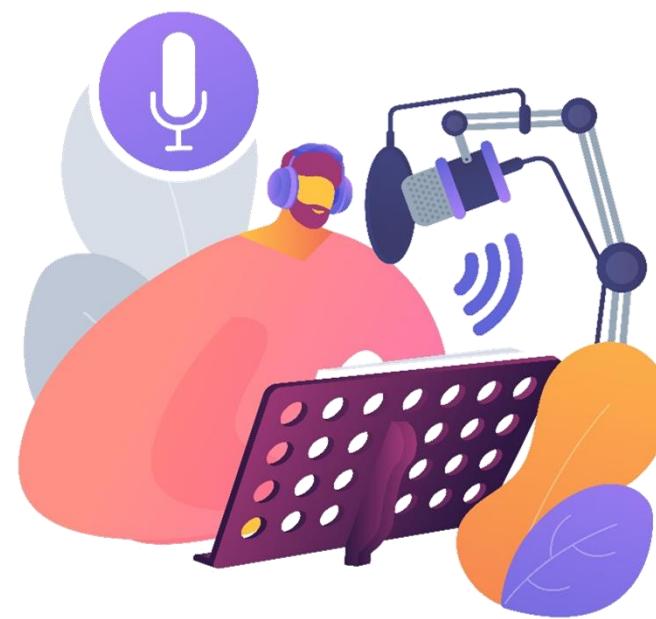
Fundamentals of Speech Recognition

What is Speech Recognition?

Speech Recognition is the process of converting spoken language into text, allowing computers to understand and respond to human speech.



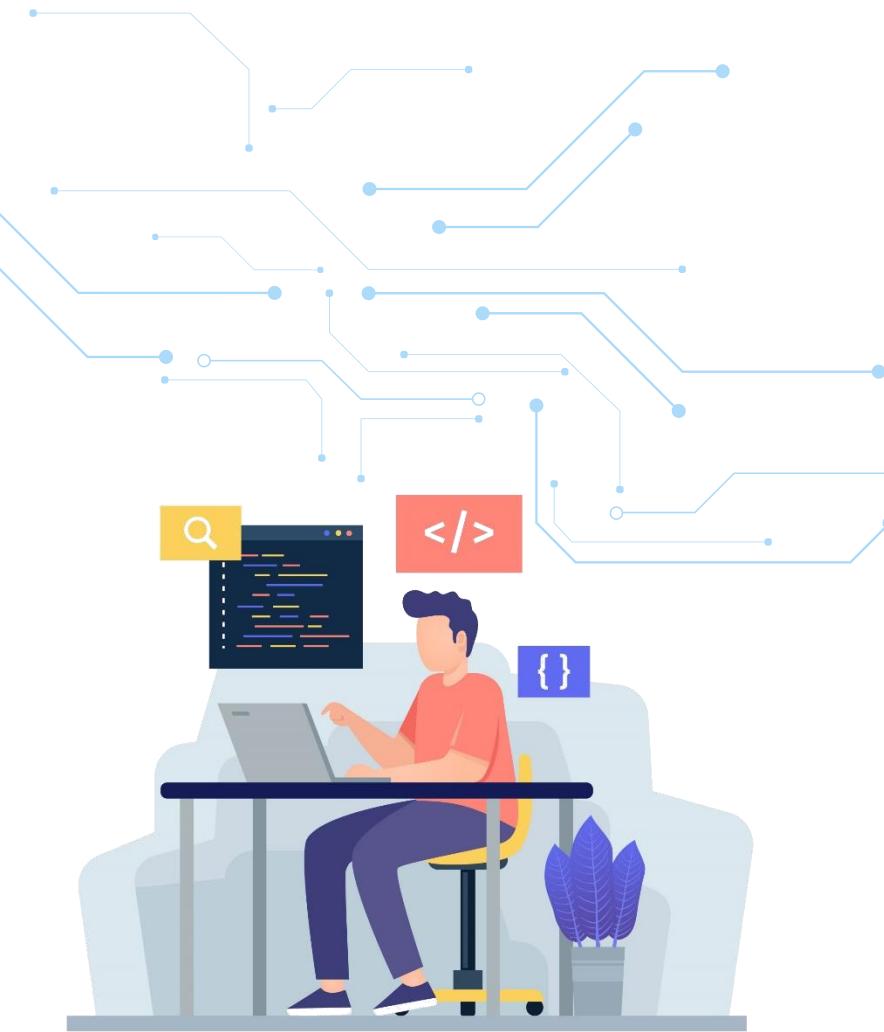
Key Components



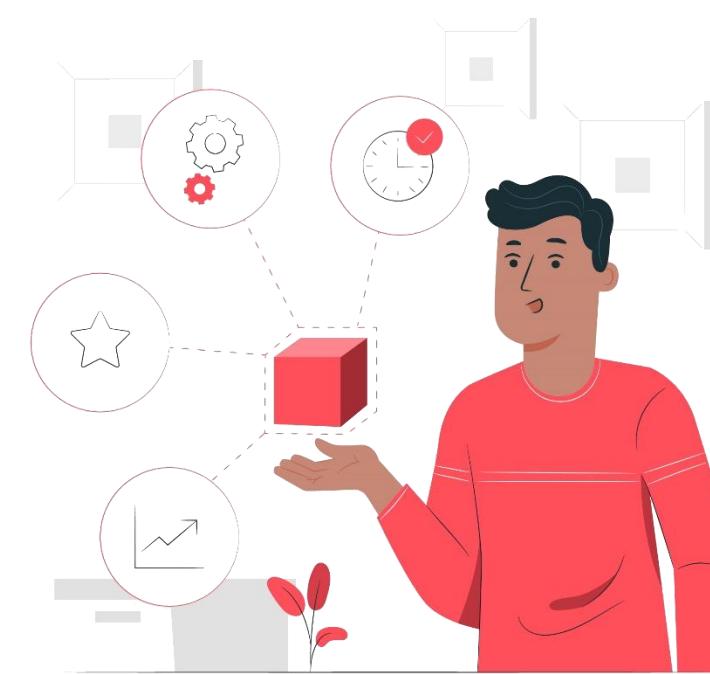
Audio Input



Acoustic Model



Decoder



Feature Extraction



Language Model

Deep Learning in Modern Speech Recognition

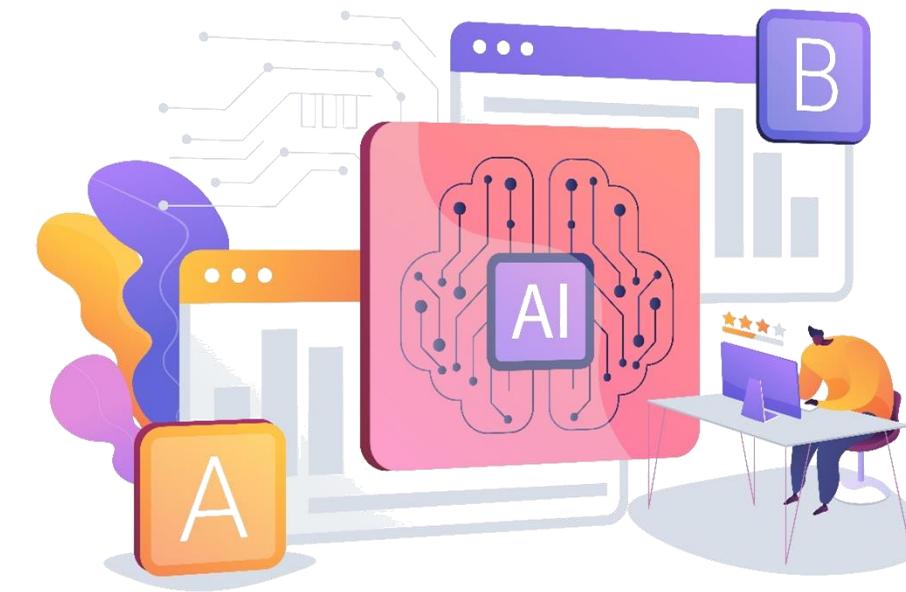


CNN + RNN + CTC



Hugging Face

Pretrained Models in Hugging Face



Transformer-Based Models

Automatic Speech Recognition (ASR) Models (e.g., Whisper, DeepSpeech)

Automatic Speech Recognition (ASR) Models

ASR models are deep learning systems that convert spoken audio into written text by learning patterns in speech, sound, and language.



DeepSpeech – A Open-Source ASR Model

DeepSpeech is a speech recognition model by Mozilla based on RNNs with the CTC loss for aligning speech frames with words.



Model Architecture



Training Approach



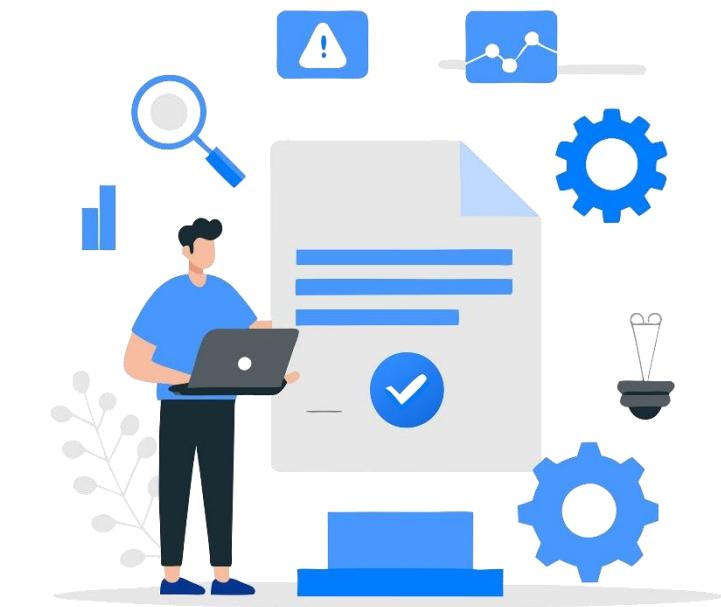
Limitations

Whisper – A Transformer-Based ASR Model

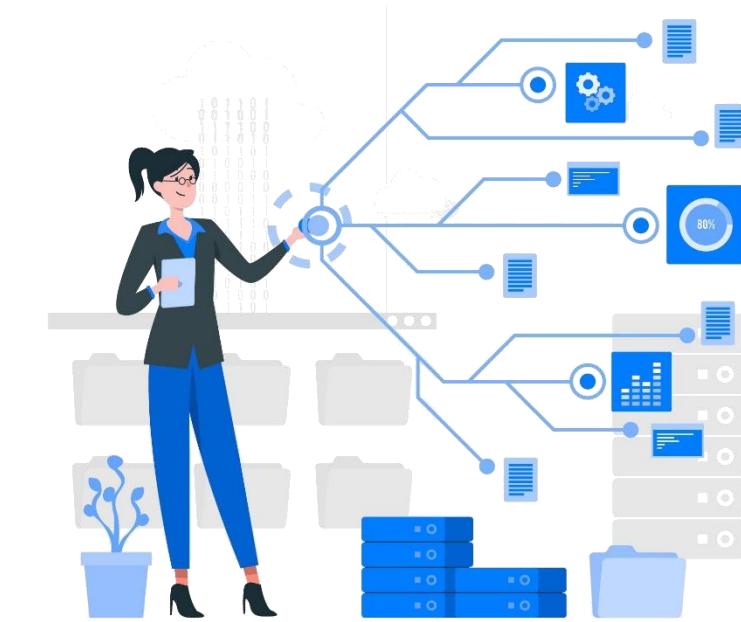
Whisper is a powerful, multilingual ASR model developed by OpenAI, trained on 680,000 hours of web audio using a sequence-to-sequence transformer.



Model Design



Key Features



Model Sizes

Comparison & Applications of ASR Models

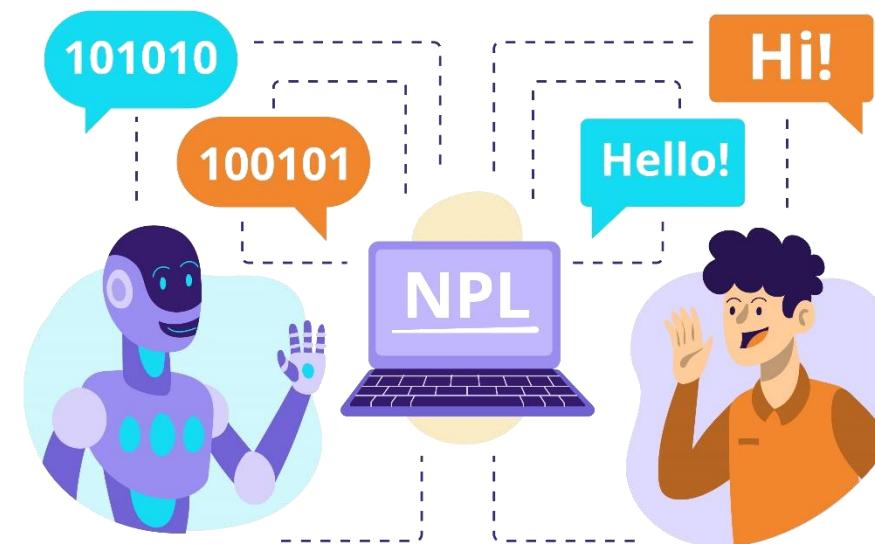
Different ASR models serve different needs — from lightweight offline processing to accurate multilingual transcription in the cloud.

Feature	DeepSpeech	Whisper (OpenAI)
Architecture	RNN + CTC	Transformer (Seq2Seq)
Multilingual	English only	50+ languages
Noise Handling	Limited	Strong
Open Source	Yes	Yes
Extra Features	Basic transcription	Transcription + Translation

Challenges in Speech Processing

What Makes Speech Processing Difficult?

Speech processing involves decoding spoken language into meaningful text or actions. This process faces numerous challenges due to the unstructured, variable, and noisy nature of human speech.



Variability in Speech



Environmental Challenges



Real-Time Constraints

Language-Related Challenges

Code-Switching

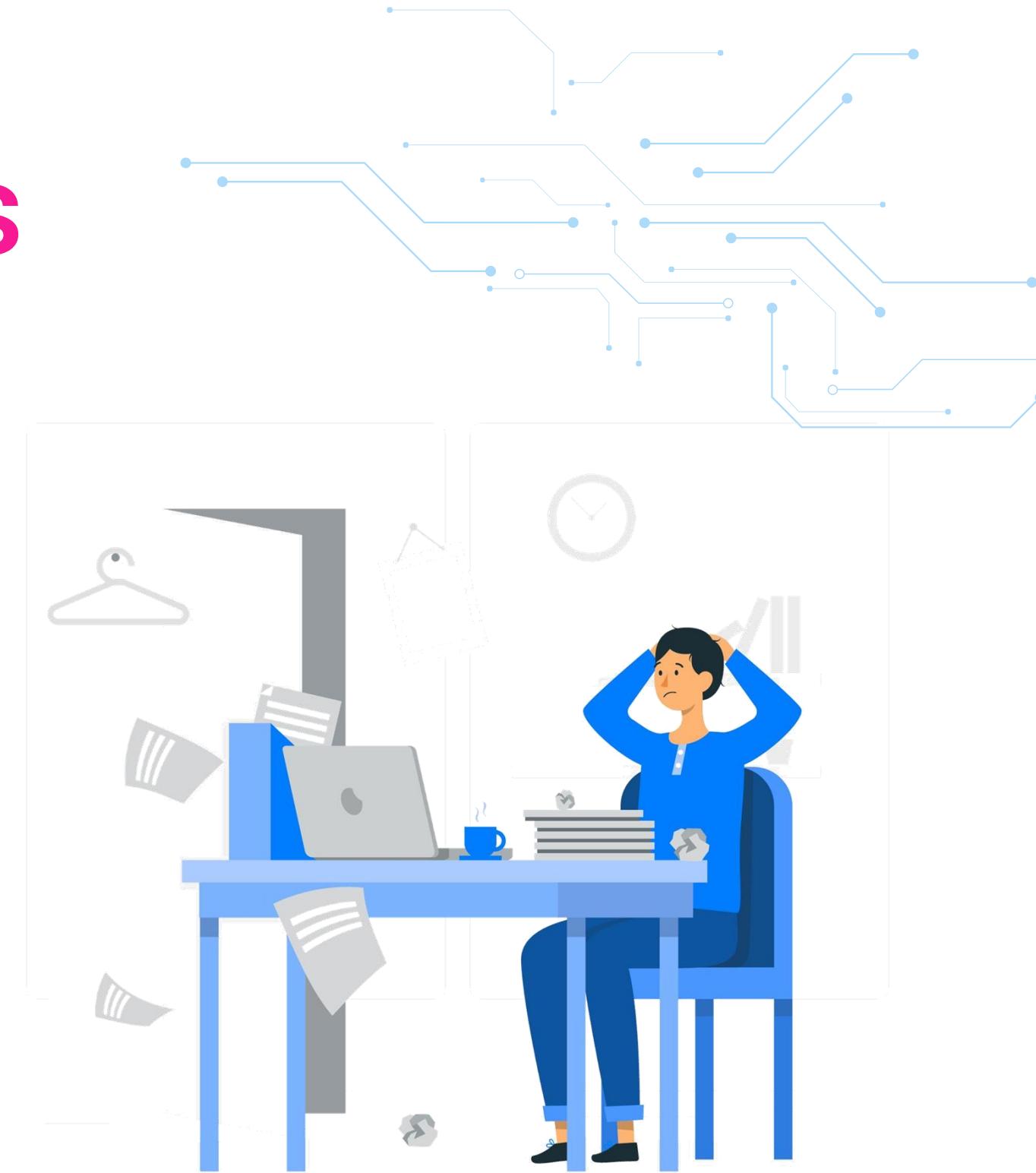
- e! Mixing languages within a sentence
- e! Common in multilingual regions (e.g., English + Hindi = Hinglish)

Low-Resource Languages

- e! Lack of labeled audio datasets
- e! No robust pretrained ASR models

Dialects, Accents & Regional Slang

- e! Same language spoken differently (e.g., Indian English vs. US English)
- e! Words like “innit” (UK) or “yaar” (India) often misrecognized



Technical & Model-Centric Challenges

Speech recognition models have limitations based on their architecture, training data, and decoding techniques — leading to errors in transcription and understanding.

Out-of-Vocabulary (OOV)

- e! Domain-specific terms (e.g., "photovoltaic", "BTS", "cringe")
- e! Newly coined or trending slang

Misalignment & Punctuation

- e! Speech is continuous; models struggle to infer pauses, sentence boundaries
- e! Lack of proper punctuation makes outputs unreadable

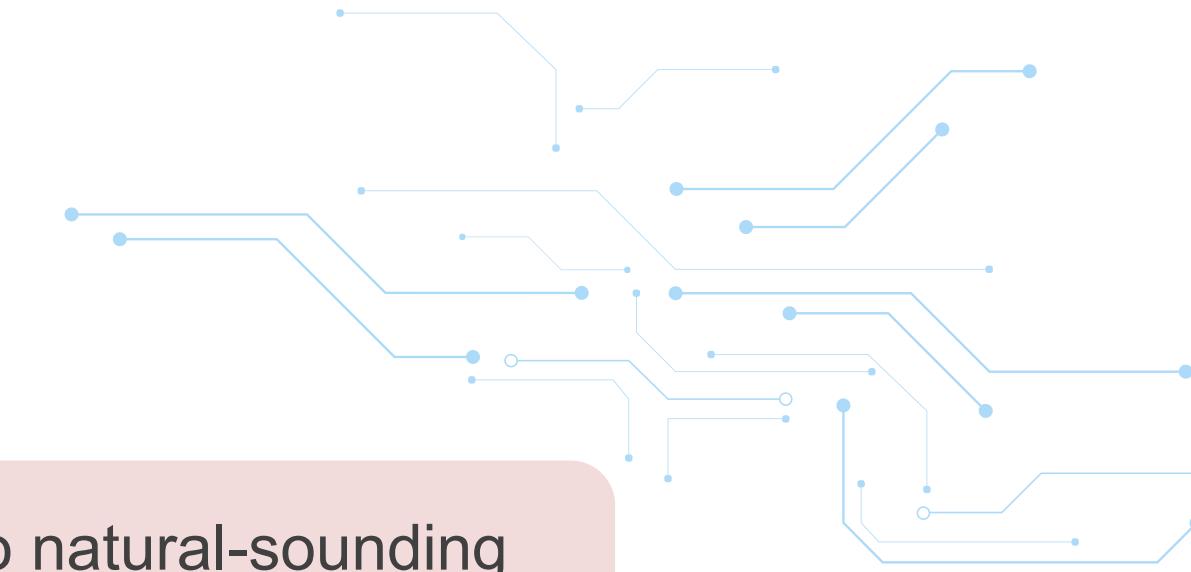
Resource Constraints

- e! Real-time speech apps must run on mobile/edge devices
- e! Large transformer-based ASR (like Whisper) needs high compute

TTS Applications in NLP

What is TTS in NLP?

Text-to-Speech (TTS) is the NLP task of converting written text into natural-sounding spoken audio using linguistic and acoustic models.



Purpose of TTS

Components of a TTS System

Evolution of TTS



Neural TTS Models and Advancements

- e! Modern TTS systems use deep learning models to generate realistic, expressive speech by learning from large amounts of audio-text data.

NVIDIA/tacotron2

Tacotron 2 - PyTorch implementation with faster-than-realtime inference

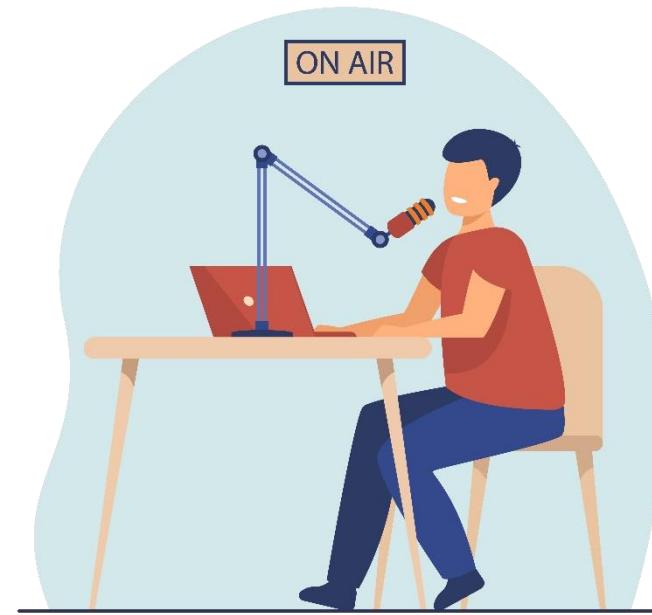


fastText



VOCODE

Challenges and Future of TTS in NLP



Expressiveness & Emotion



Future Directions



Personalization



Multilingual & Code-Switching Support

Introduction to Multimodal Learning

What is Multimodal Learning?

Multimodal learning is an AI technique where a model learns from multiple data modalities — like text, images, and audio — simultaneously to improve understanding and performance on complex tasks.

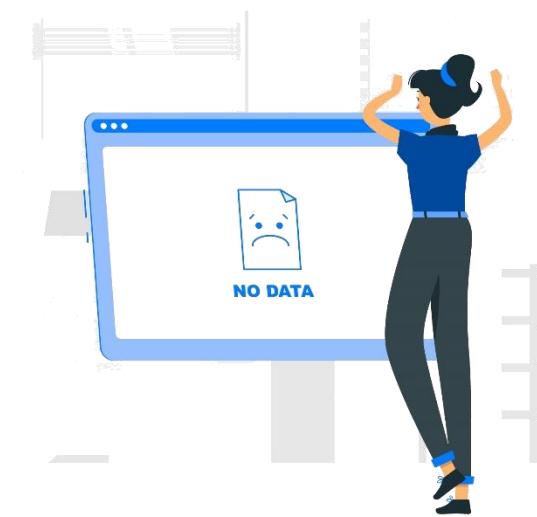
What is a Modality?

- e! A specific type of data:
- e! Text (language, words)
- e! Image (pixels, objects)
- e! Audio (waveforms, speech, music)

Why Combine Modalities?

- e! Humans use multiple senses for understanding
- e! Fuses complementary signals
- e! Enables better context, reasoning, and prediction

Benefits of Multimodal Learning



Robustness to Missing Data



Better Performance



Contextual Understanding



Real-World Adaptability

Key Architectures & Techniques

- e! Early Fusion vs. Late Fusion
- e! Cross-Attention Mechanisms
- e! Multimodal Transformers
- e! Popular Frameworks



Vision-Language Models (CLIP, BLIP)

What Are Vision-Language Models (VLMs)?

VLMs are deep learning systems designed to understand both images and text, and learn how they relate — enabling tasks like image captioning, visual question answering (VQA), and zero-shot image classification.

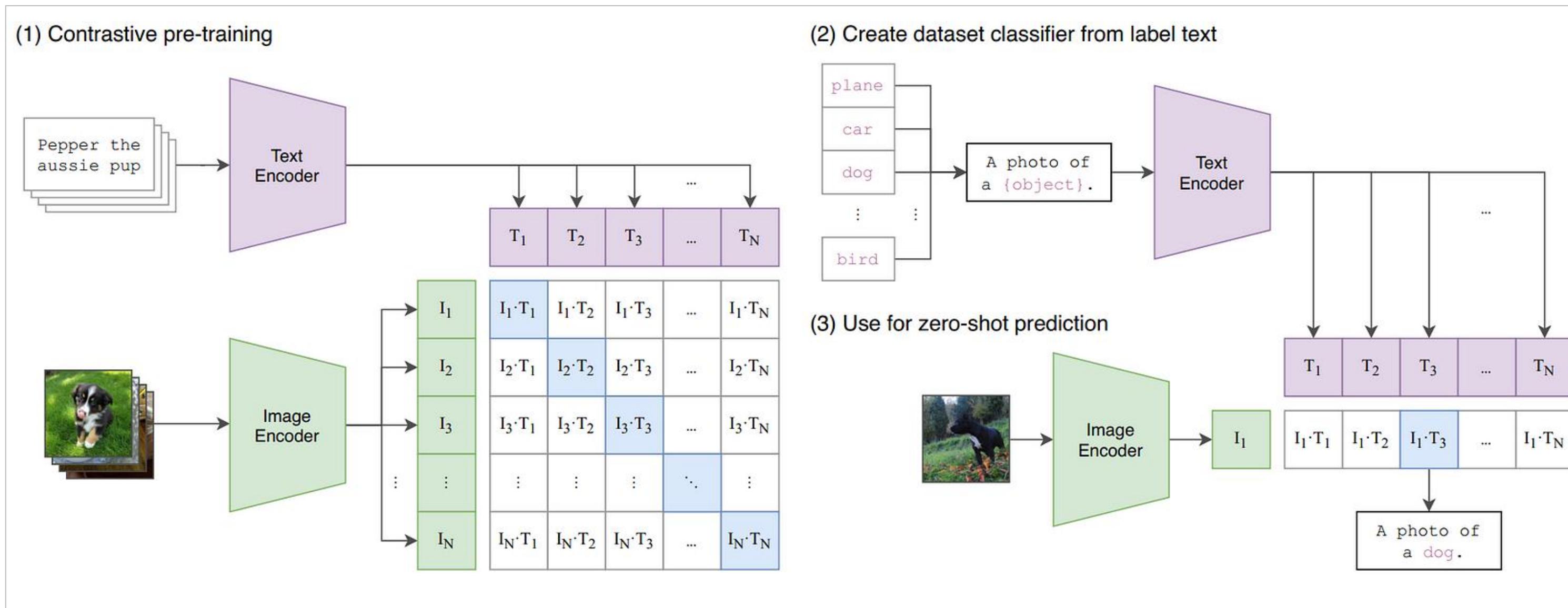
Input: An image of a cat wearing glasses



Model finds the best caption: “A stylish cat with sunglasses”

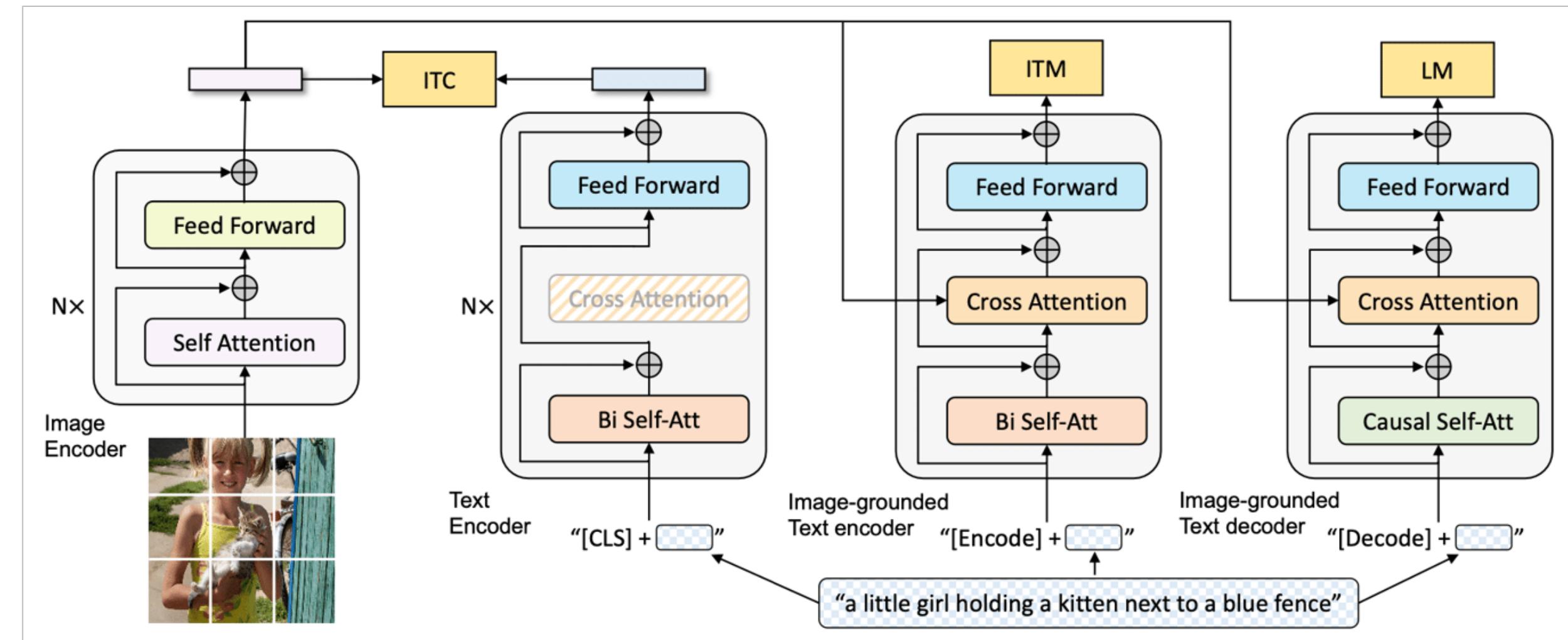
Contrastive Language-Image Pretraining

Contrastive Language-Image Pretraining (CLIP by OpenAI) is a vision-language model that learns to match images and their captions using contrastive learning.



Bootstrapped Language-Image Pretraining

Bootstrapped Language-Image Pretraining (BLIP by Salesforce AI) is a flexible VLM that supports image captioning, question answering, and image-text retrieval using both encoder-decoder and contrastive methods.



Integrating Text and Speech in NLP Applications

Integration of Text and Speech Mean

Integrating text and speech in NLP involves designing systems that can understand, process, and generate both spoken and written language, enabling fluid, human-like interactions.

Why Integrate Speech and Text?

- e! Combine the naturalness of voice with the precision of text
- e! Handle spoken commands, transcription, and voice-based outputs

Modalities Involved

- e! Speech to Text (ASR) → Converting audio to written form
- e! Text to Speech (TTS) → Generating spoken audio from text
- e! NLP Processing → Intent detection, question answering, dialogue systems

Core Components in a Text+Speech System



Automatic Speech Recognition
(ASR)



Natural Language Processing
(NLP Core)



Text-to-Speech
(TTS)

Generative Models for Multimodal Content

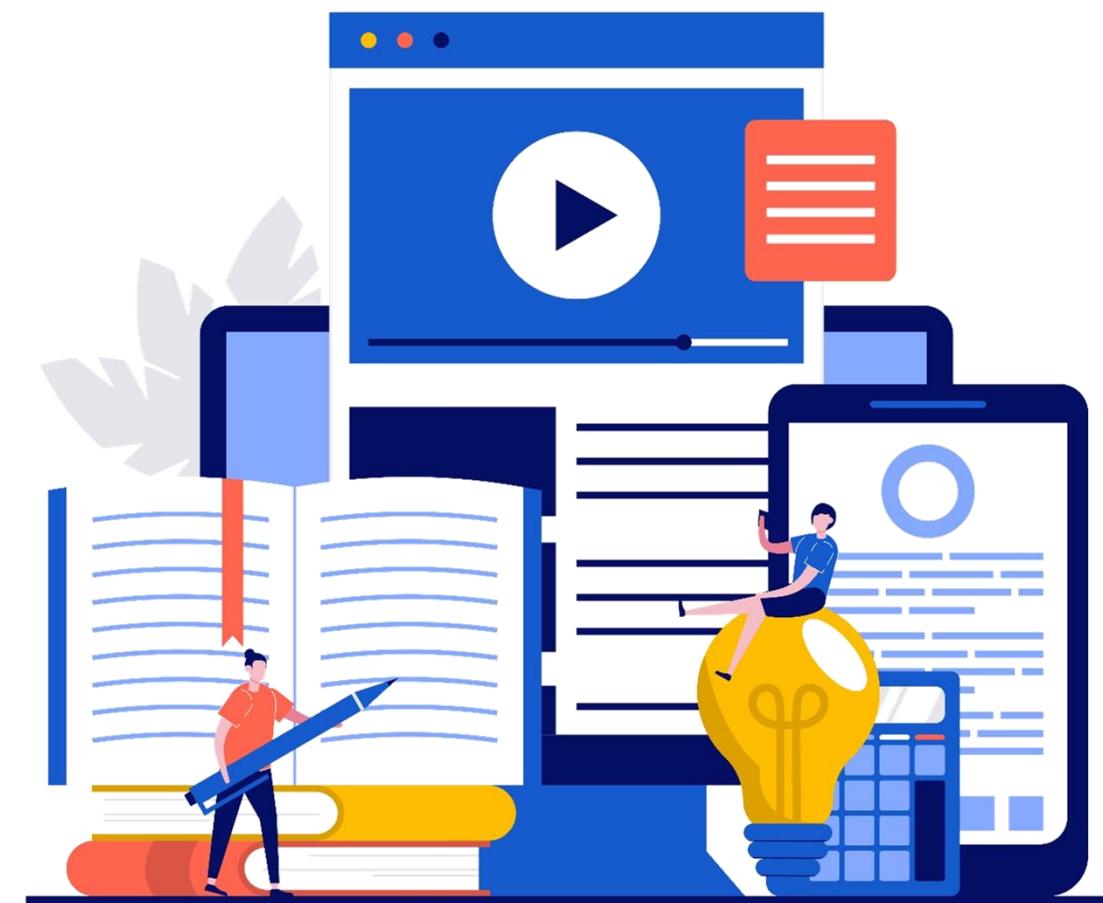
Generative Models for Multimodal Content

Generative multimodal models are AI systems capable of creating new content (e.g., text, images, audio, video) by learning the relationships between multiple data modalities like language, vision, and sound.

What is “Multimodal Generation”?

Creating content using more than one modality

- e! Text → Image (e.g., DALL·E)
- e! Image → Text (e.g., BLIP captioning)
- e! Text → Audio (e.g., TTS systems)



Core Generative Models and Their Capabilities

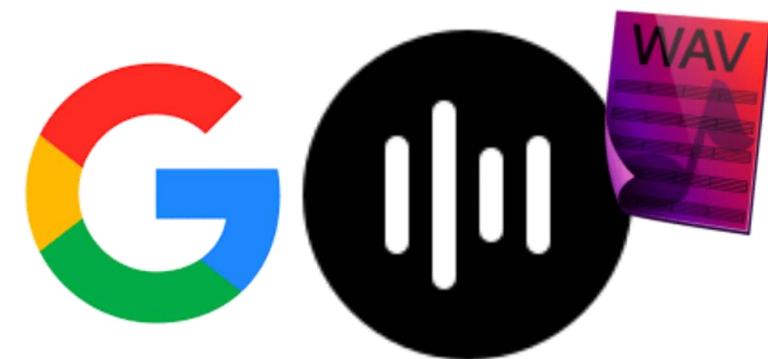
DALL·E / DALL·E 3



BLIP & BLIP-2



AudioLM / MusicLM



GPT-4o (Omni)



Applications of Multimodal Generative Models

e! Creative Content Generation

e! Education & Accessibility

e! Social Media & Marketing

e! Virtual Avatars & Agents



Model Compression (Quantization, Pruning, Knowledge Distillation)

Introduction to Model Compression

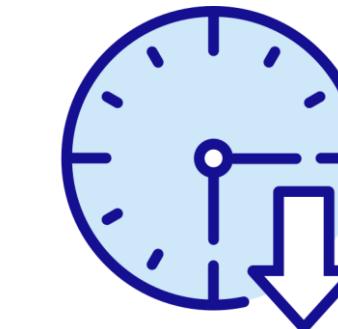
Model compression is a set of techniques used to reduce the size, memory footprint, and latency of deep learning models — while trying to retain most of their original performance.

Why It's Important:

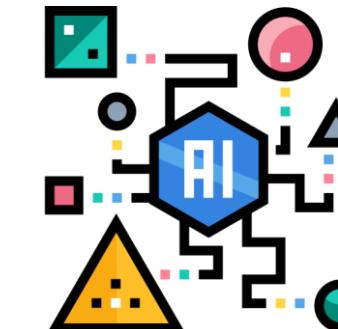
Lower Deployment Costs



Reduced Inference Time



Efficiency for Edge Devices



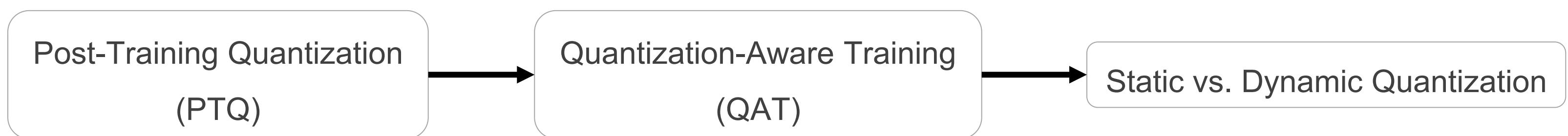
Enabling On-Device AI

Quantization

Quantization reduces model size and speeds up inference by converting high-precision (e.g., 32-bit float) weights and activations into lower-precision formats like 8-bit integers.

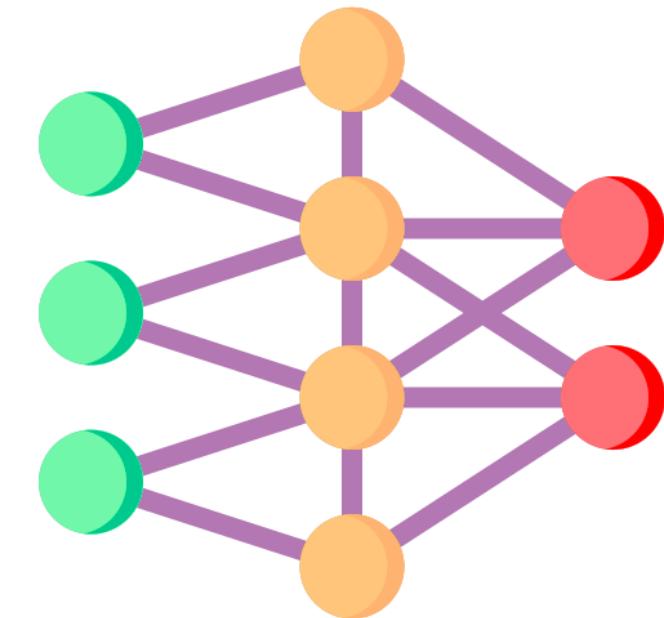


Methods:



Pruning

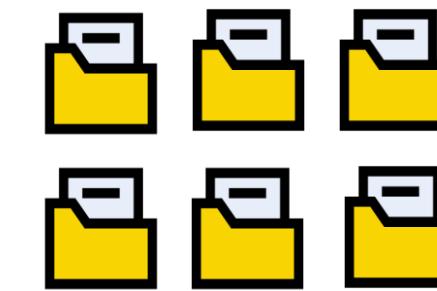
Pruning removes unnecessary or low-impact weights, neurons, or attention heads from a model to make it smaller and faster.



Neuron/Head Pruning



Iterative vs One-Shot Pruning



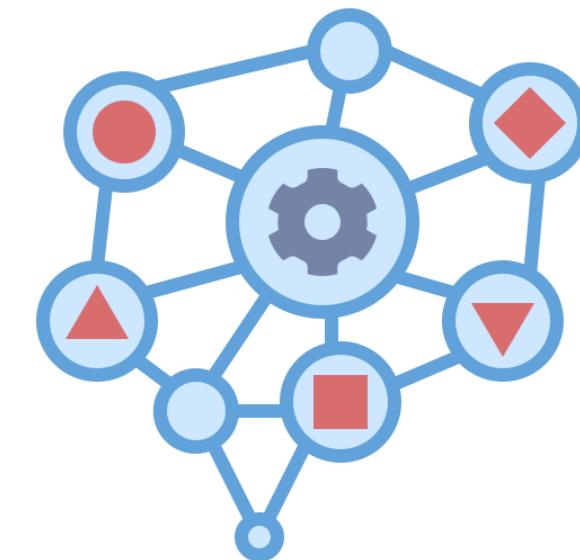
Weight Pruning

Knowledge Distillation

Knowledge distillation is the process of training a smaller “student” model to mimic a larger, well-trained “teacher” model.



Teacher Model

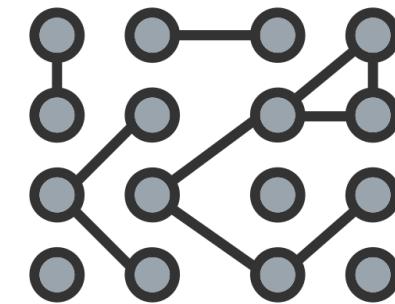


Student Model

Training Strategy

Types of Distillation

Teaching
smaller model



Real-Time NLP for Edge Devices

What Is Real-Time NLP on Edge Devices?

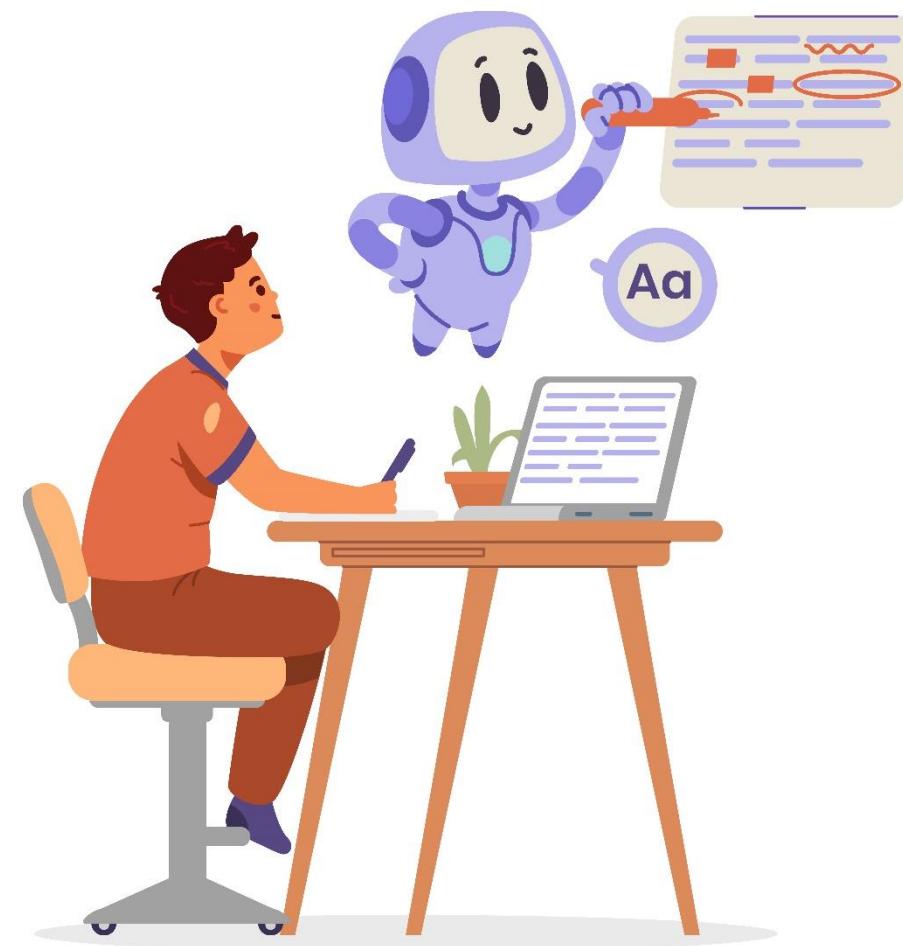
Real-time NLP for edge devices refers to performing natural language tasks (like speech recognition, translation, summarization) directly on local hardware (e.g., phones, wearables), without relying on cloud servers.

What Are Edge Devices?

- e! Smartphones, Raspberry Pi, wearables, home assistants, automotive systems

Why Use NLP on Edge?

- e! Privacy: Data stays on-device
- e! Latency: Instant responses, no network delay
- e! Offline Functionality: Works without internet



NLP Tasks Commonly Run on Edge Devices

Edge NLP enables basic to moderately complex language functions directly on devices using optimized or compressed models.

Wake Word Detection

Intent Detection

Text Classification

On-Device Translation / TTS / ASR



Techniques to Enable NLP on Edge

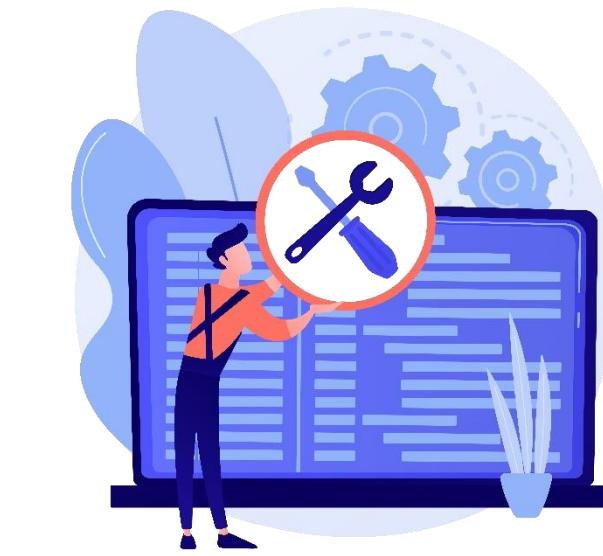
Deploying NLP on edge requires specialized model compression, optimization, and hardware-aware techniques.



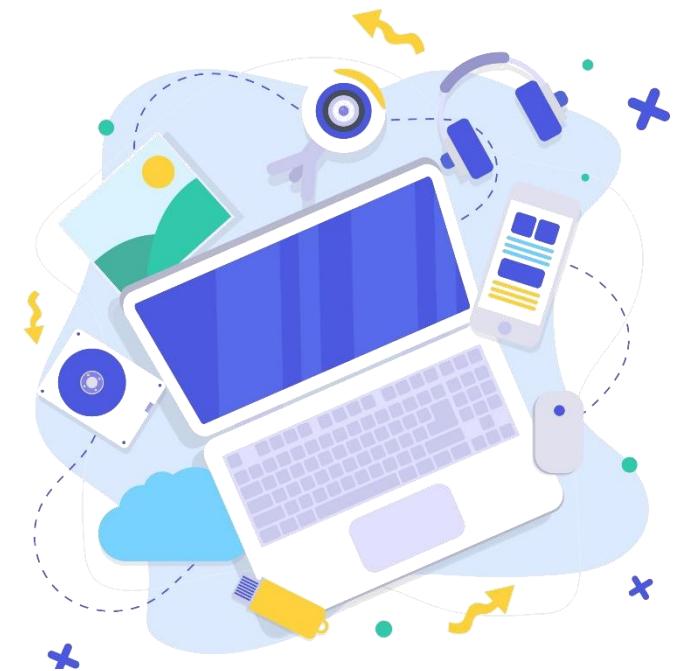
Lightweight NLP Architectures



Model Compression



Hardware Acceleration



Frameworks & Tools

Explainability in Deep NLP Models

What is Explainability in NLP Models?

Explainability refers to the ability to understand, interpret, and trust the decisions made by deep NLP models, especially those using transformers and large neural networks.

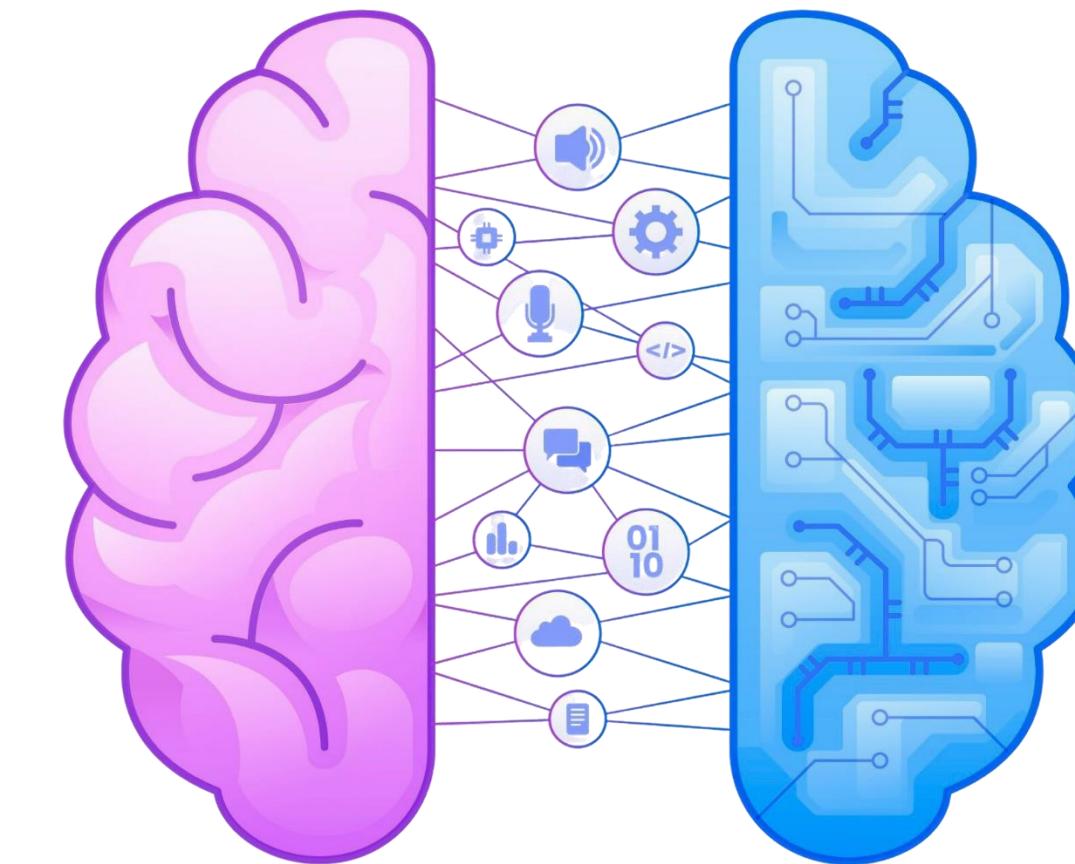
Why Explainability Matters:

Transparency in AI Decisions

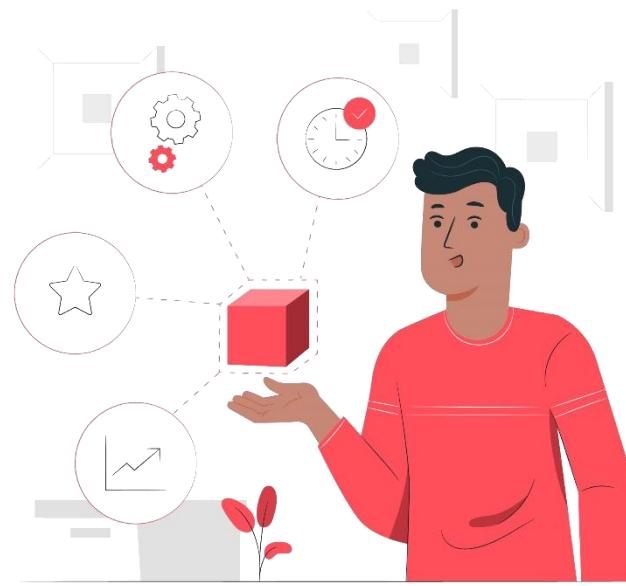
Debugging and Improvement

Regulatory Compliance

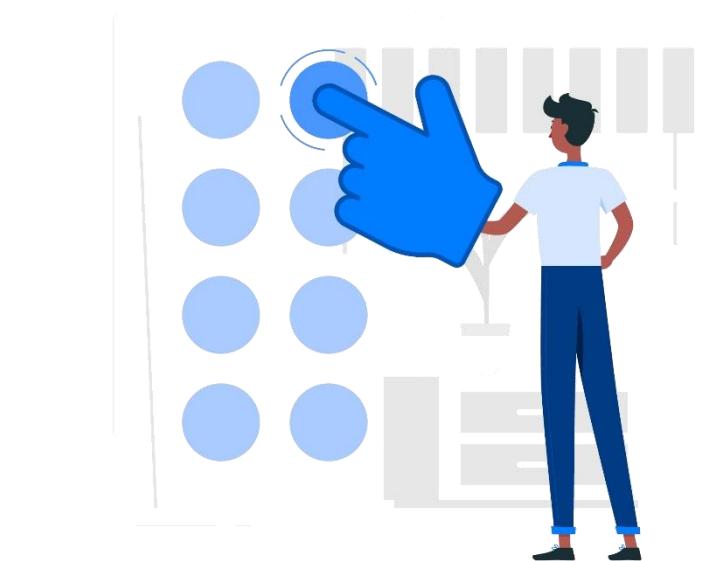
User Trust



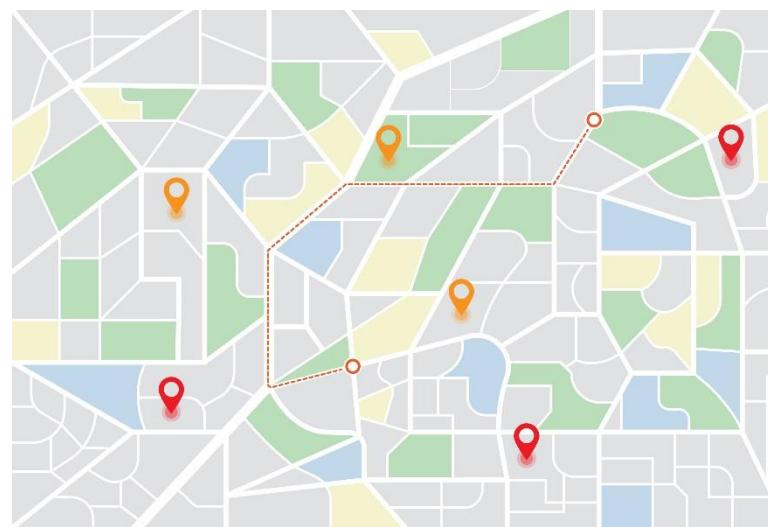
Methods for Explaining NLP Models



Feature Importance



Counterfactual Examples



Saliency Maps



Attention Visualization

Tools and Libraries for Explainable NLP



Captum

ELI5

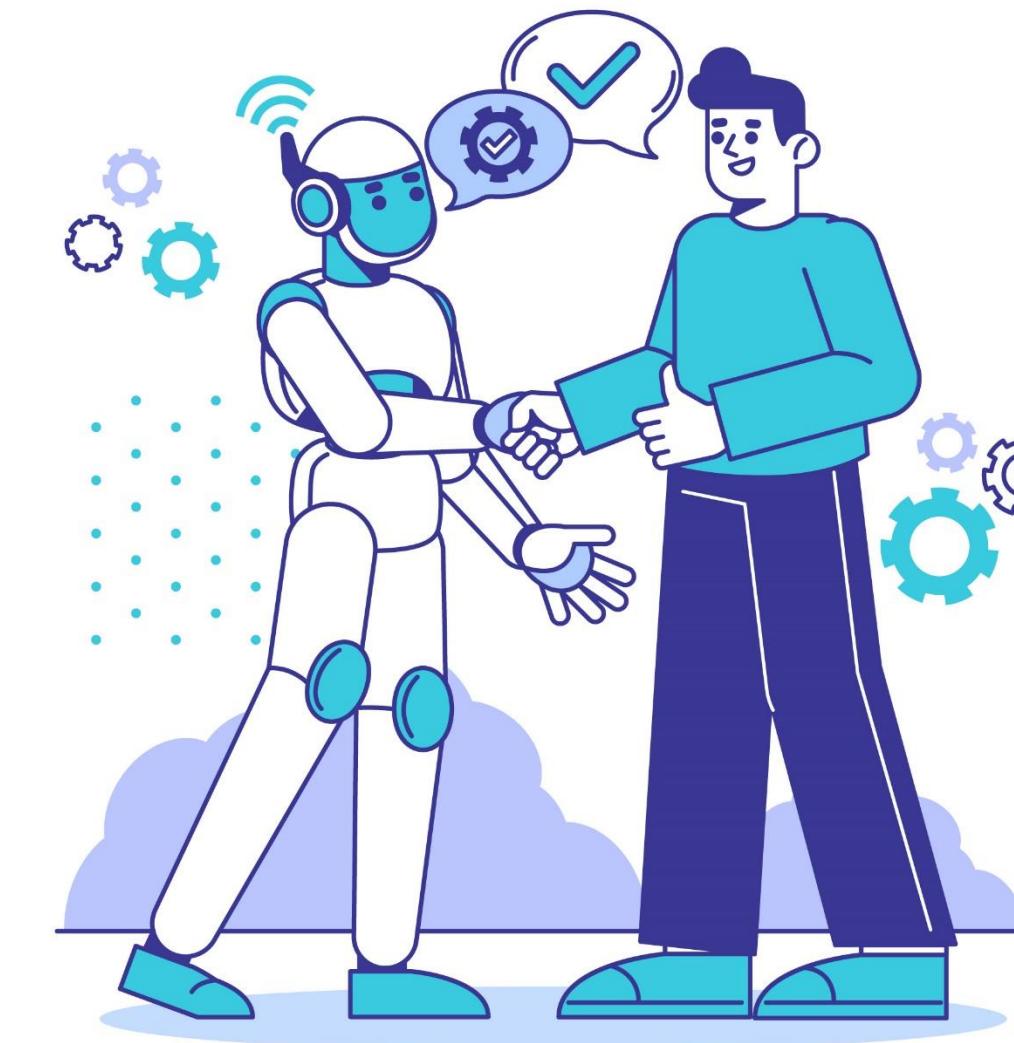


Shap

Ethical AI in NLP

What is Ethical AI in NLP?

Ethical AI in NLP involves designing and deploying natural language models that are fair, unbiased, accountable, transparent, and respect human rights.



Common Ethical Issues in NLP

Bias in Training Data

- e! Models trained on internet text often inherit stereotypes
- e! May associate professions with gender or race

Offensive or Harmful Outputs

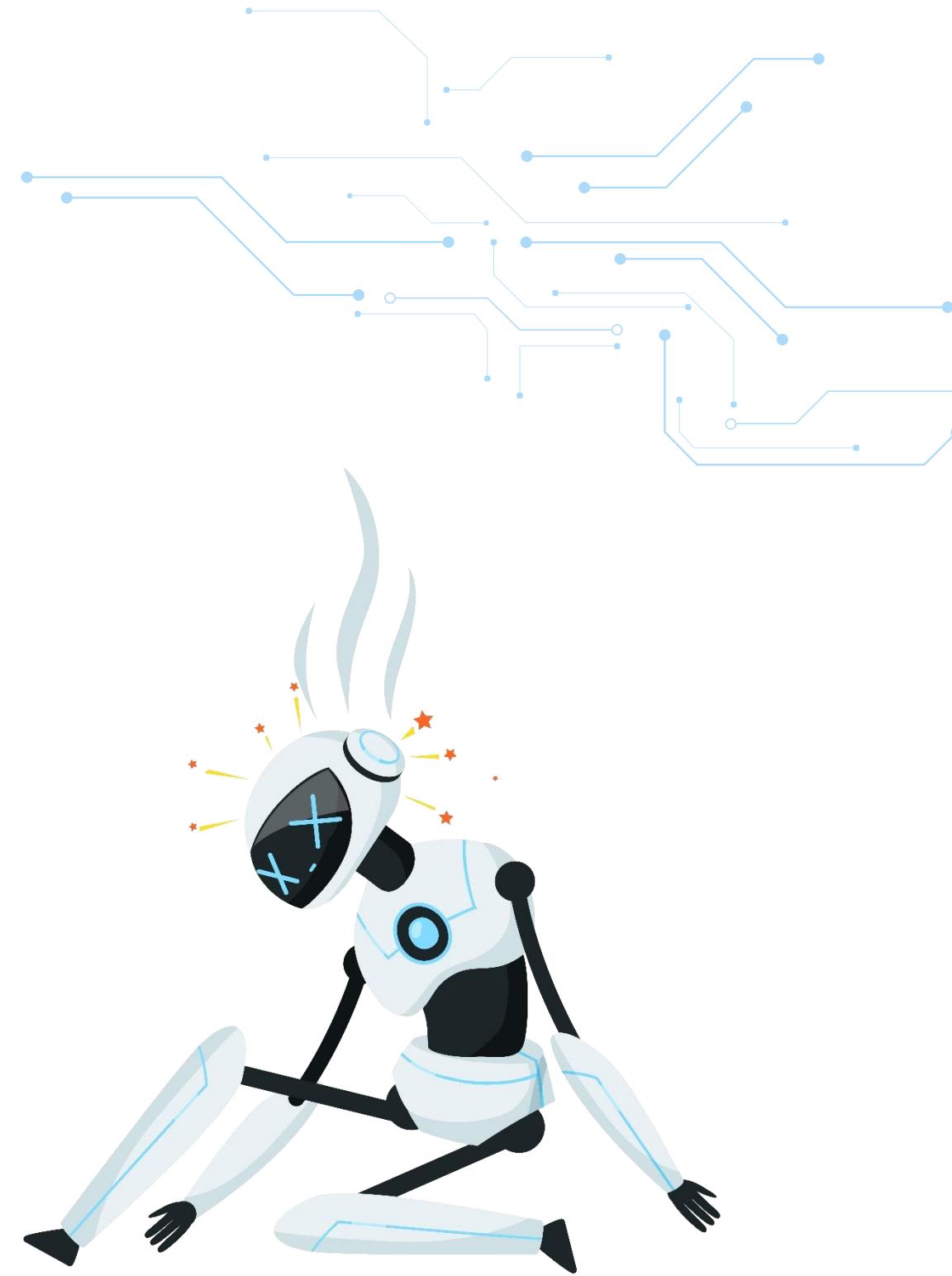
- e! Chatbots may produce toxic responses if not filtered
- e! Hate speech detection can be too aggressive or too lenient

Privacy Violations

- e! Language models may memorize and leak sensitive data (e.g., phone numbers)

Misinformation and Hallucination

- e! Generative models can confidently produce incorrect or fake facts



Future of Ethical AI in NLP



Transparent Benchmarking



Cross-Disciplinary Teams



Ethical AI by Design

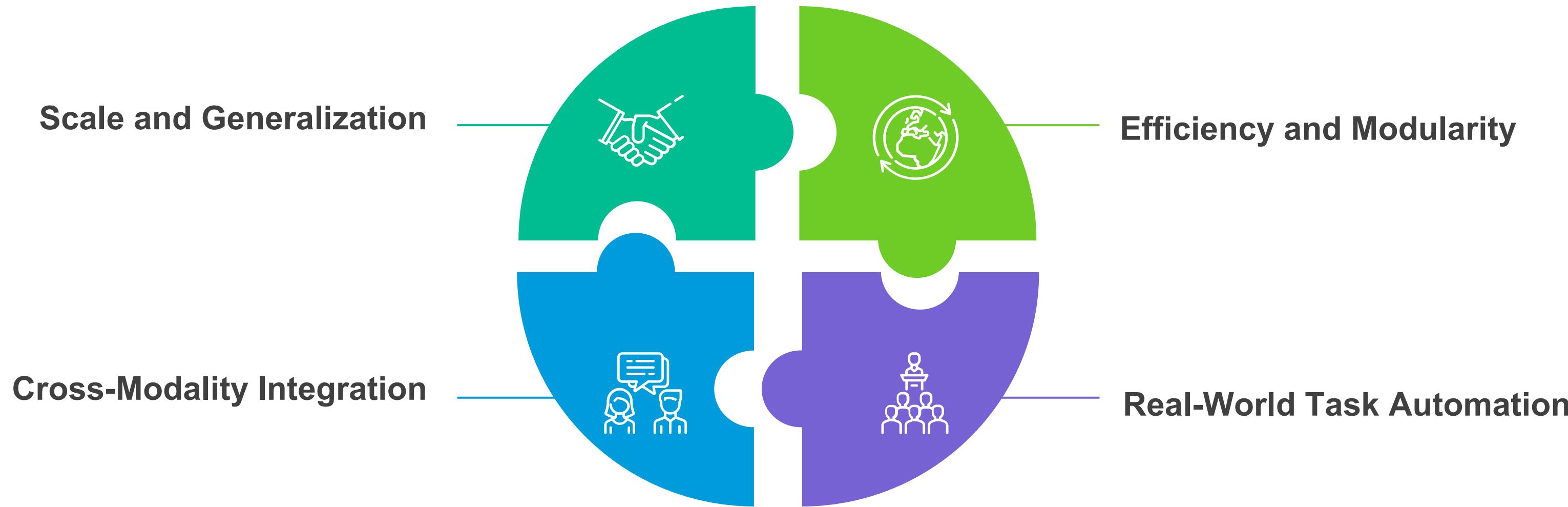


Global Standards and Regulations

Emerging Trends (GPT-5, HyperNetworks, Autonomous AI Agents)

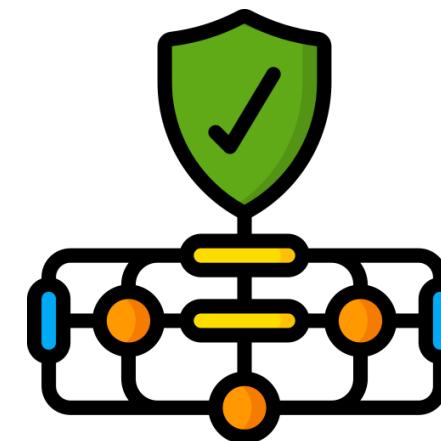
What's Driving Emerging Trends in NLP?

Emerging trends in NLP focus on building more intelligent, adaptive, and autonomous systems using advancements in foundation models, network architectures, and agent-based AI.

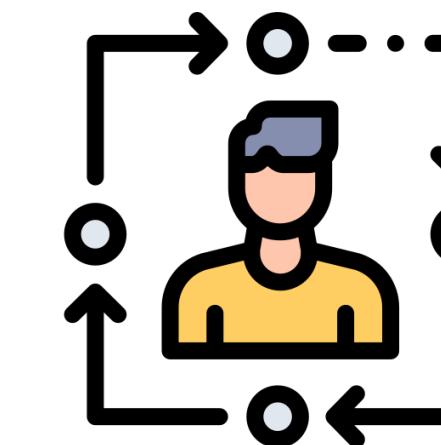


GPT-5- The Next Gen of Foundation Models

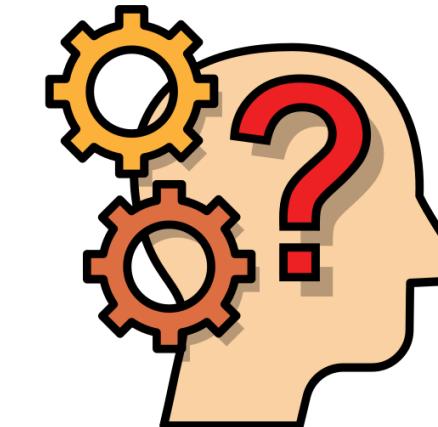
GPT-5 (upcoming) represents the next evolution in large language models, expected to be more multimodal, grounded, and autonomous.



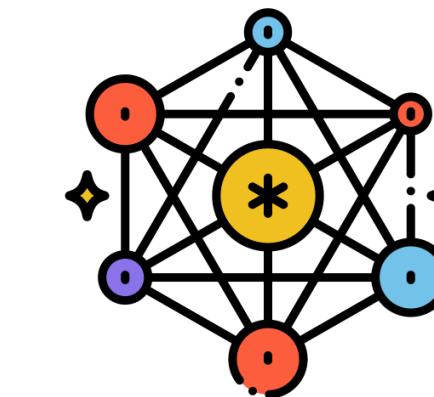
Safer, More Controlled Generation



Long-Term Memory & Personalization



Better Reasoning & Planning



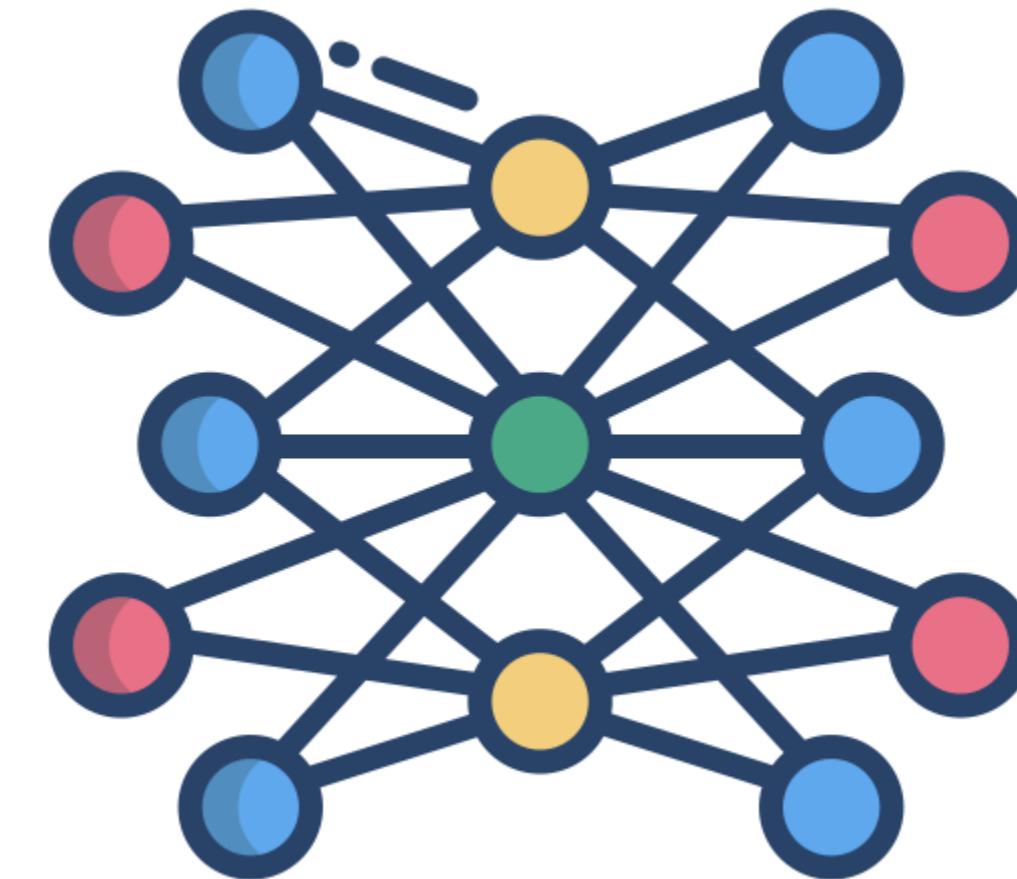
Deeper Multimodal Capabilities

HyperNetworks – Smarter Parameter Sharing

HyperNetworks are neural networks that generate the weights for another network, enabling faster training and task adaptability with fewer parameters.

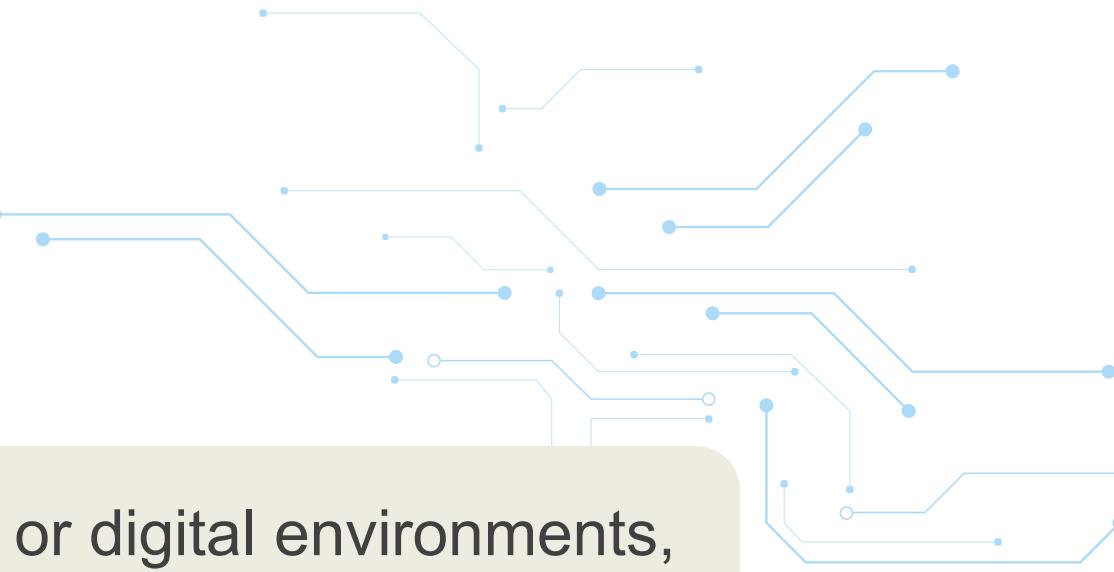
Benefits & Use Cases:

1. Efficient Fine-Tuning
2. Modular & Transferable
3. Lower Training Cost
4. Supports Continual Learning



Rise of Autonomous AI Agents

Autonomous AI agents are systems that can perceive, plan, and act in real-world or digital environments, often powered by large LMs + tools.



Capabilities & Components:

Self-Critique and Planning



Goal-Driven Behavior



Memory and Learning



Tool Usage



Building a Voice Narration QA Assistant for Audiobooks (Demonstration)

Note: Refer to Module 9: Demo 1 on LMS for detailed steps.

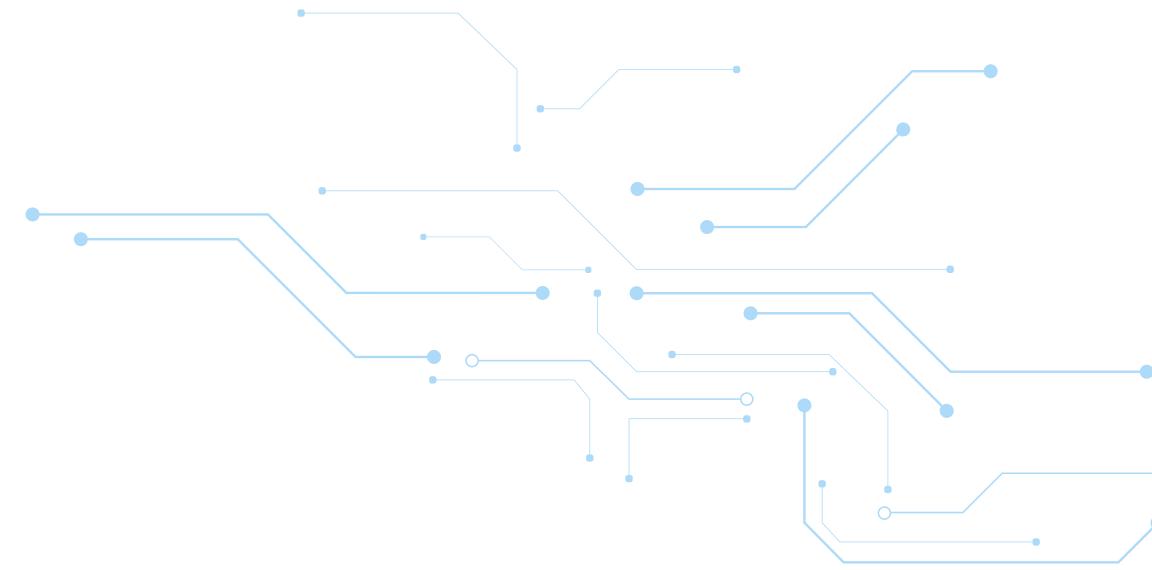
Autonomous Warehouse Object Classification via Text-Image Fusion (Demonstration)

Note: Refer to Module 9: Demo 2 on LMS for detailed steps.

Summary

In this lesson, you have learned to:

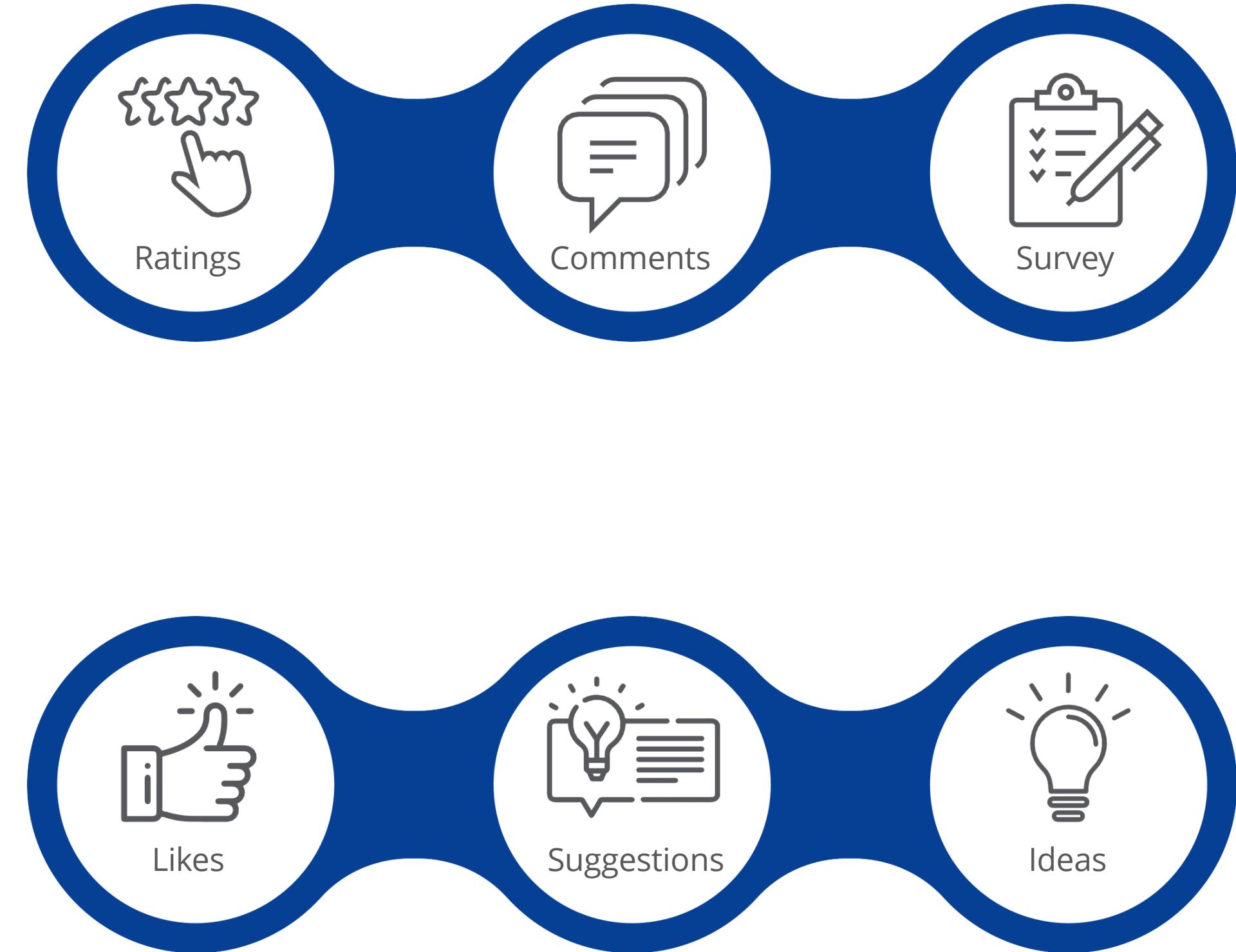
- e! Explain how speech and text are processed in NLP
- e! Perform comparison between ASR models
- e! Understand challenges in speech processing and real-time NLP for edge deployment
- e! Apply multimodal and generative models to integrate and generate text, audio, and images
- e! Evaluate the use of model compression, explainability, and ethical principles in NLP systems
- e! Summarize emerging trends in NLP such as GPT-5, HyperNetworks, and autonomous AI agents



Questions



Feedback



Thank You

For information, Please Visit our Website
www.edureka.co

