# IMAGE CAPTIONING SYSTEM

Presented by: Roopsi Srivastava (202410116100173)

Rani Kumari(202410116100163)

Rabita Yadav(202410116100)

# INTRODUCTION

- **What is Image Captioning?**
A process of generating textual descriptions for an image using both computer vision and natural language processing.
- **Why is it Important?**
- Helps visually impaired people
- Used in content-based image retrieval
- Powers features in social media and e-commerce platforms

# Real-World Applications

- **Examples:**
- Google Photos automatic captions
- Facebook alt text for visually impaired users
- Pinterest visual search
- Autonomous vehicles understanding their surroundings
- Add 2–3 image examples with captions for effect

# How Image Captioning Works

- **High-Level Overview:**
  - **Image Input** → CNN (extract features)
  - **Feature Vector** → RNN/LSTM (generate sentence)
  - **Output** → Caption
- **Visual Aid:** Show a diagram of the above flow.

# Components Involved

- **1. CNN (Convolutional Neural Networks):**
  Extracts image features (e.g., using VGG, ResNet)

- **2. RNN/LSTM/GRU (Language Model):**
  Generates sequence of words based on extracted features

- **3. Attention Mechanism (Optional):**
  Focuses on specific parts of the image while generating each word

# Algorithms & Models

- **Encoder-Decoder Architecture**
- **Encoder:** CNN
- **Decoder:** RNN or LSTM with Word Embeddings
- **Pretrained Models:**
- InceptionV3, ResNet50 (for image encoding)
- Beam Search or Greedy Decoding (for caption generation)

# Tools and Technology

- **Languages:** Python
- **Libraries:** TensorFlow / PyTorch, OpenCV, NLTK
- **Frameworks:** Keras, HuggingFace Transformers
- **Other:** Google Colab

# Future Scope

- Improved caption quality with GPT models

- Multilingual caption generation

- Real-time captioning in AR/VR

- Personalized image descriptions

# Image Uploaded