

**IIT Gandhinagar**

ACM Summer School 2025

# Self Supervised Learning on Vision

**Tutor:** Rishabh Mondal, Diya Thakor

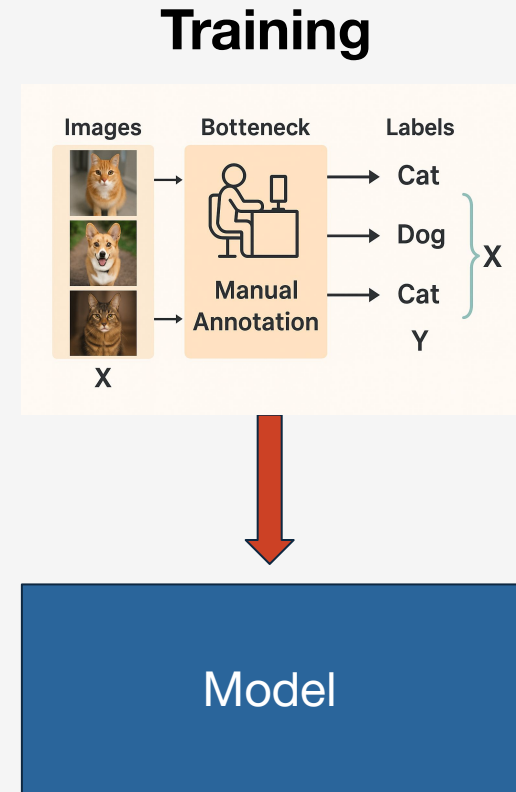
# Supervised Learning

## Training:

- Input  $X \rightarrow$  Label  $Y$
- Model learns to map  $X \rightarrow Y$  using labeled data

## Testing:

- Given new  $X$ , model predicts corresponding  $Y$



# Supervised Learning

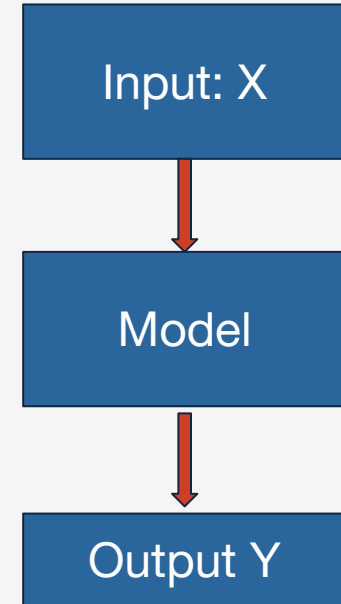
## Training:

- Input  $X \rightarrow$  Label  $Y$
- Model learns to map  $X \rightarrow Y$  using labeled data

## Testing:

- Given new  $X$ , model predicts corresponding  $Y$

## Testing



# Why Is That an Issue?

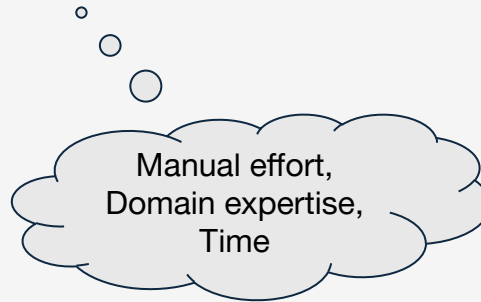
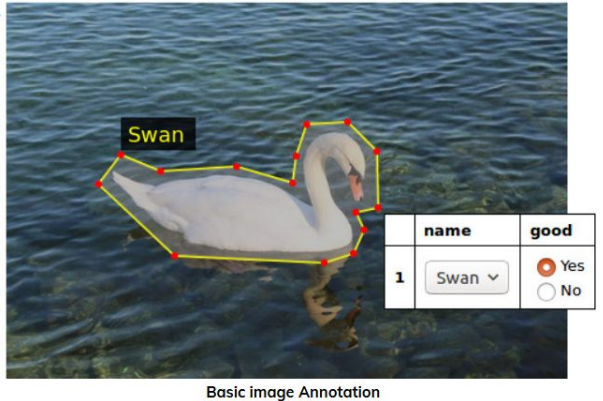
- We need a **huge** amount of **labeled data for training**
- Not all data is labeled
- **Labels** can be **expensive** or **hard** to get (e.g., medical images, satellite data)

Images in your mobile



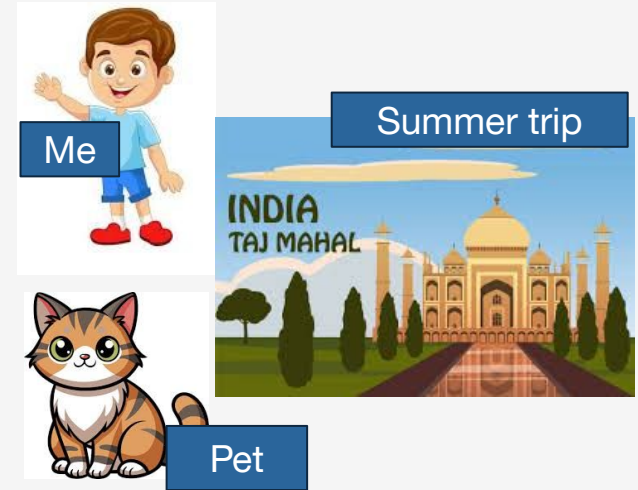
# Why Is That an Issue?

- Not all data is labeled
- **Labels** can be **expensive or hard** to get (e.g., medical images, satellite data)



How many images are there in mobile gallery?

Give label to each image



# What is image embedding?

- A numerical vector representation of an image capturing its semantic content.
- Deep neural networks (e.g., CNN, ViT), usually from intermediate or final layers.
- Similar images will have similar embeddings.



0.12

0.51

... 0.40

⋮

0.38

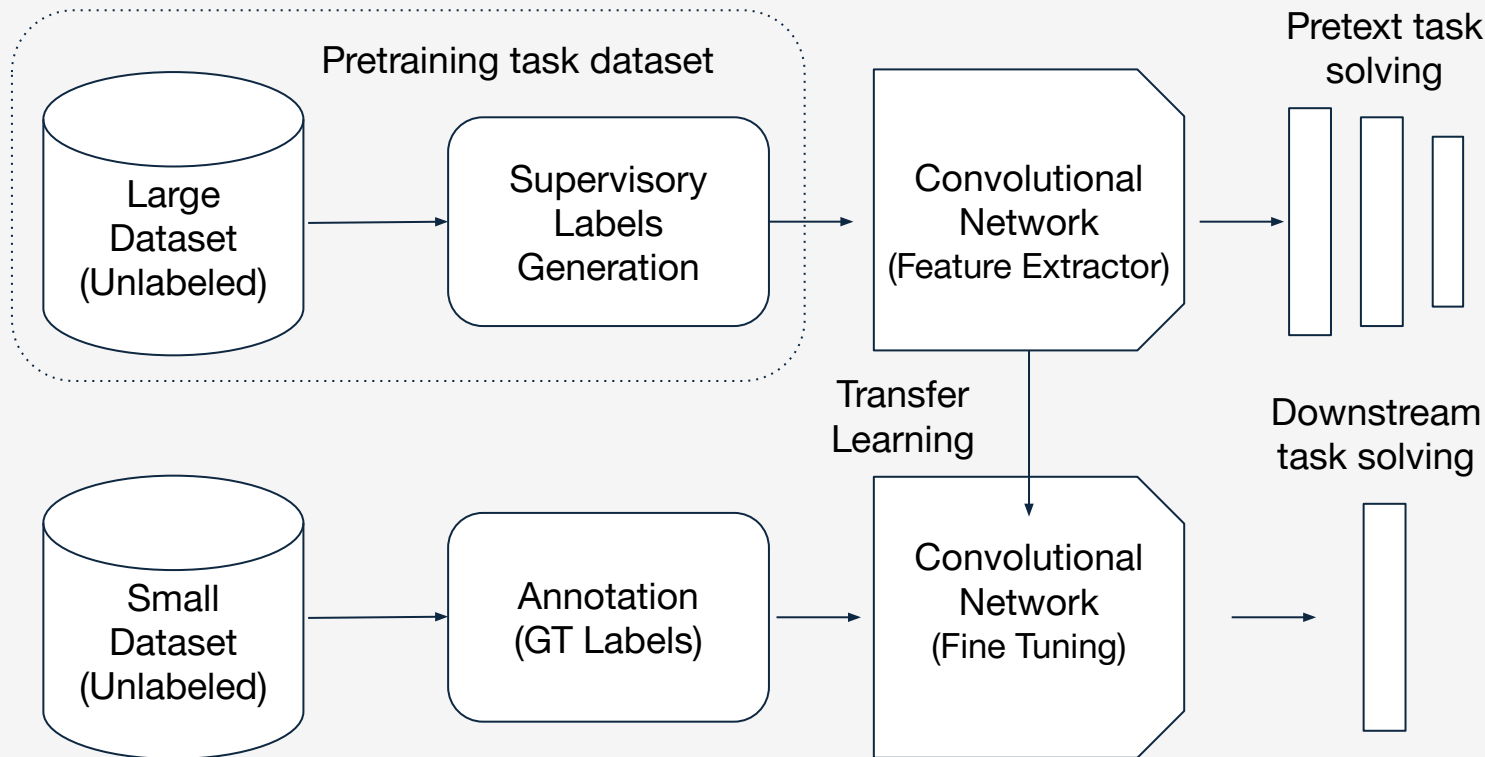
# **Solution: Self Supervised Learning**

# Self Supervised Learning

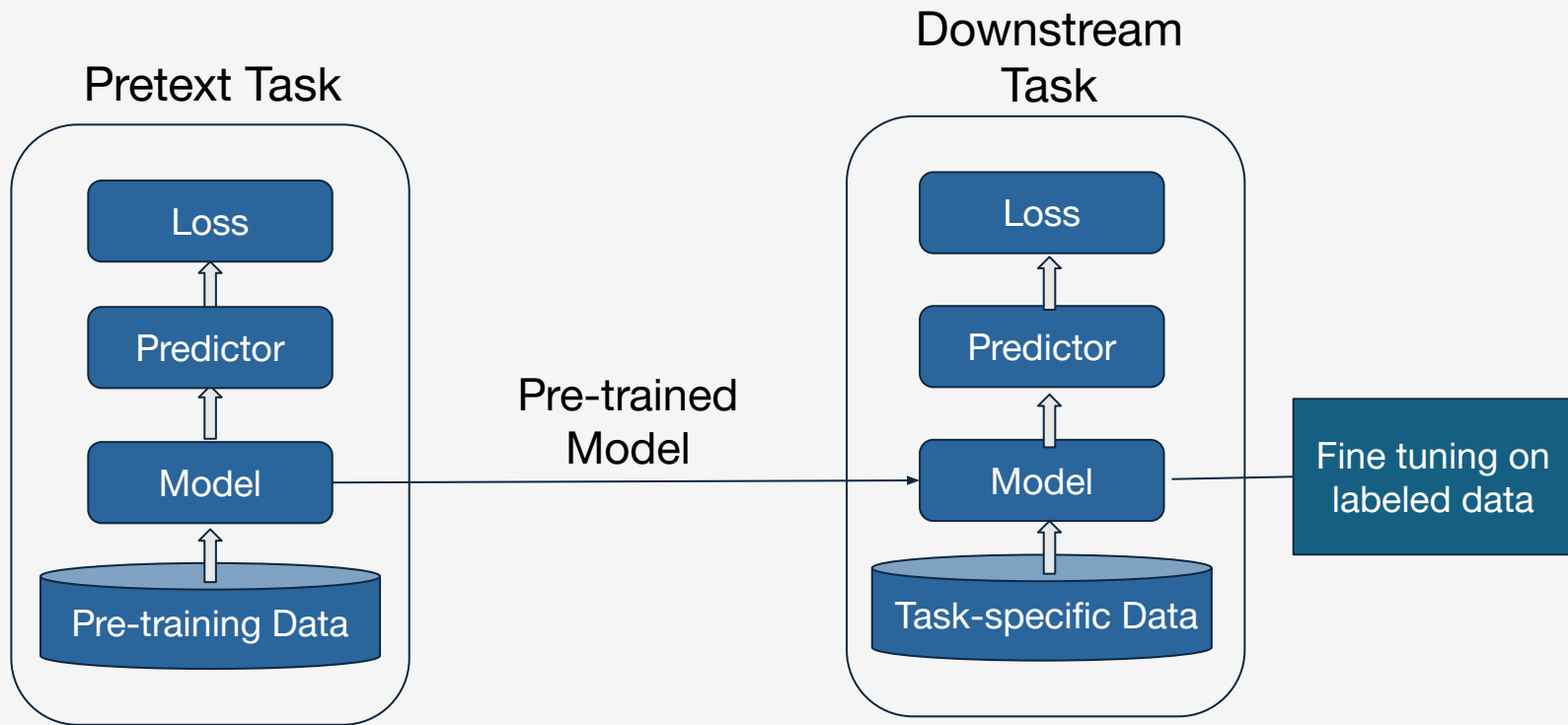
- The model **supervises itself** using **patterns** in the data.
- **No** need for **human-labeled data** — it generates its own learning representations.



# Flow of Self Supervised Learning



# Self-Supervised Learning Pipeline



# Different pretraining tasks of SSL

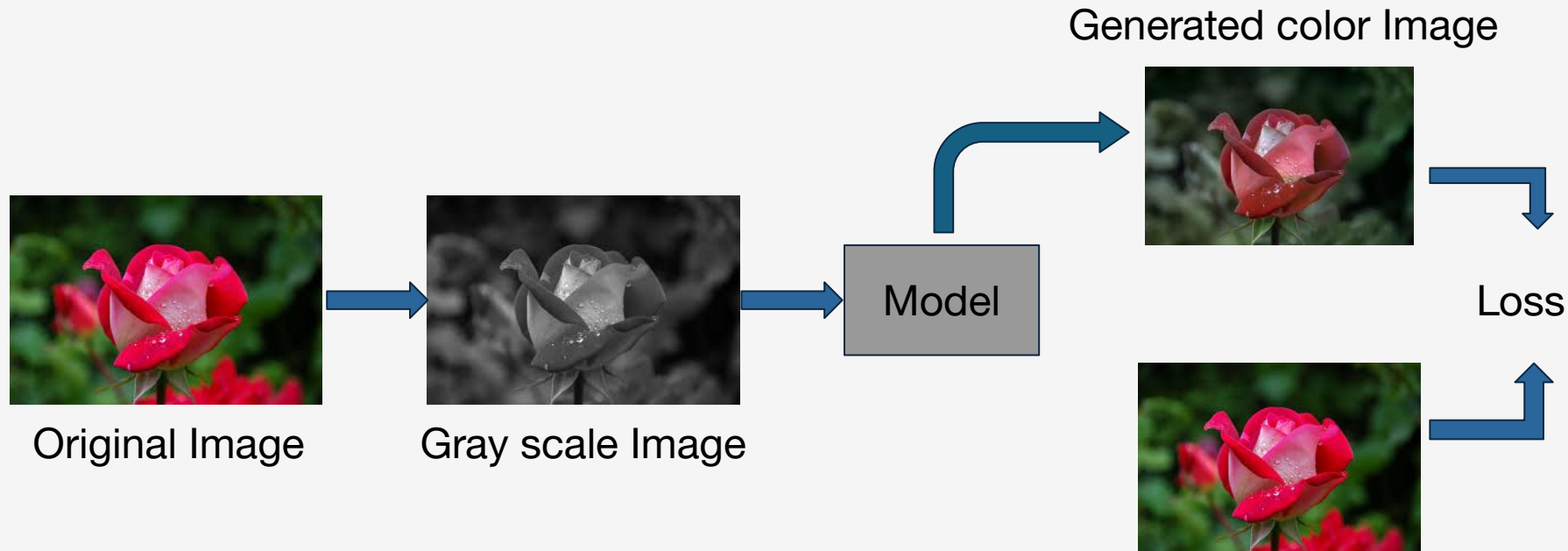
A task designed **using unlabeled data** to help the model **learn** general-purpose **features**

- Image Colorization
- Image Inpainting
- Geometric Transformation
- Jigsaw puzzle
- DINO
- SimCLR

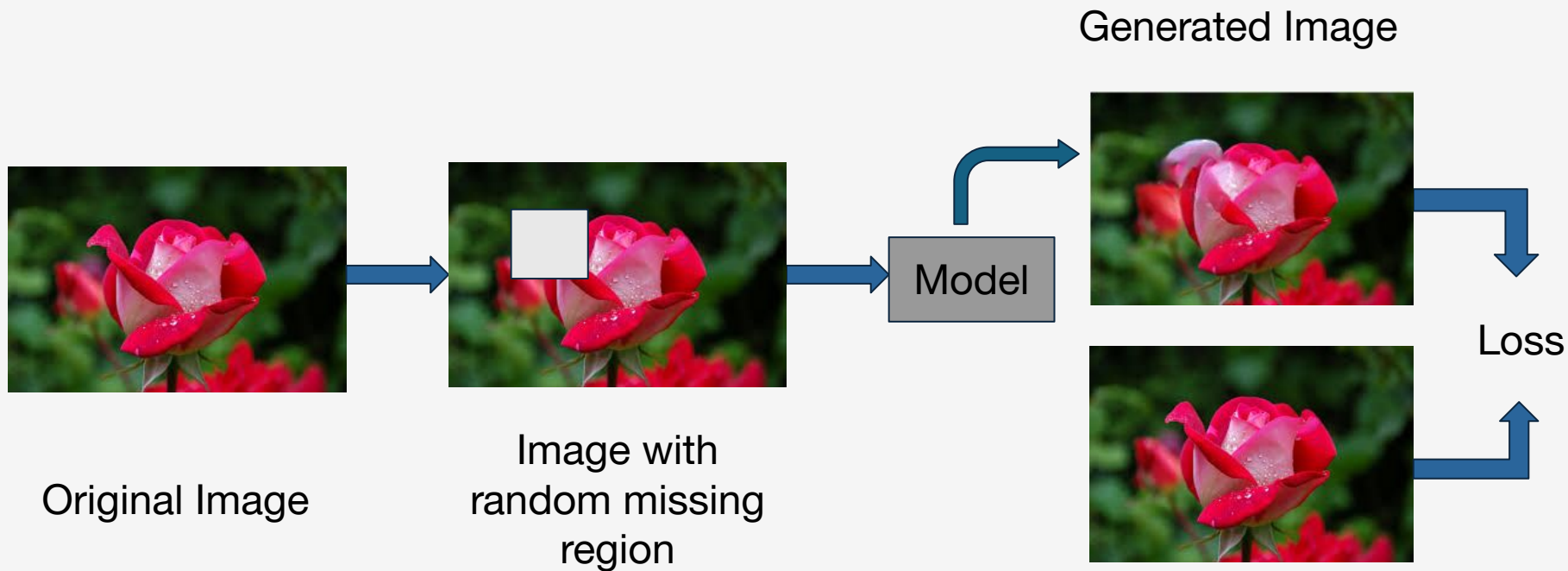
Recognition

<https://www.jmlr.org/papers/volume23/21-1155/21-1155.pdf>

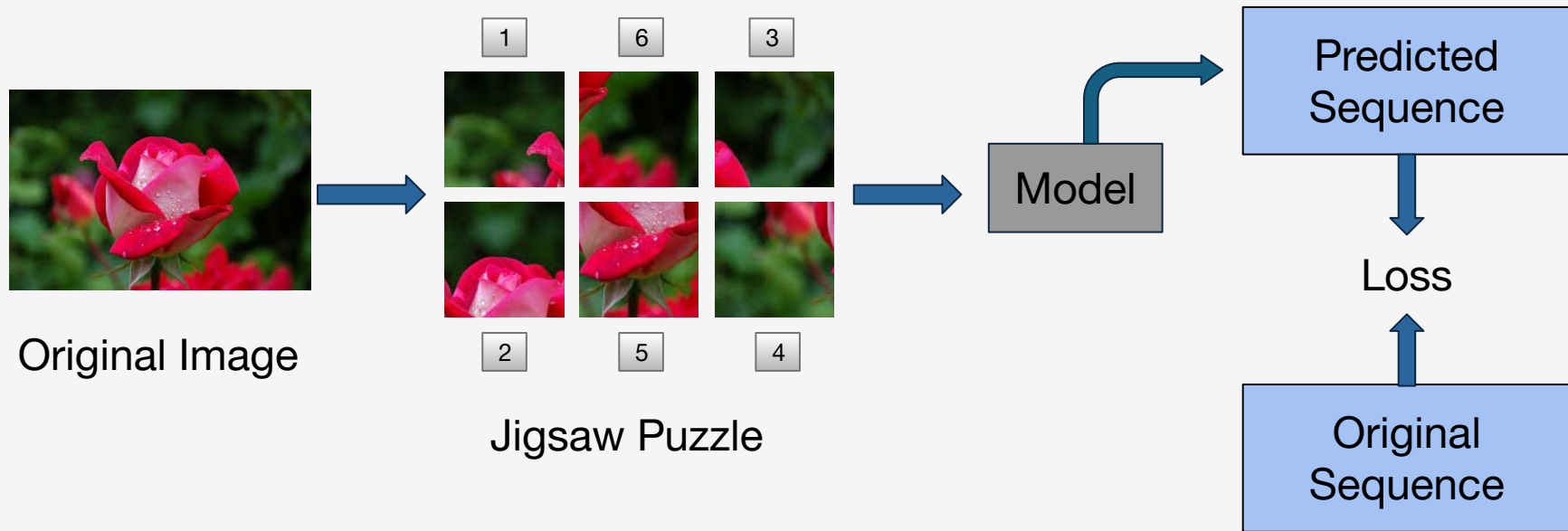
# Pre Training Task: Image Colorization



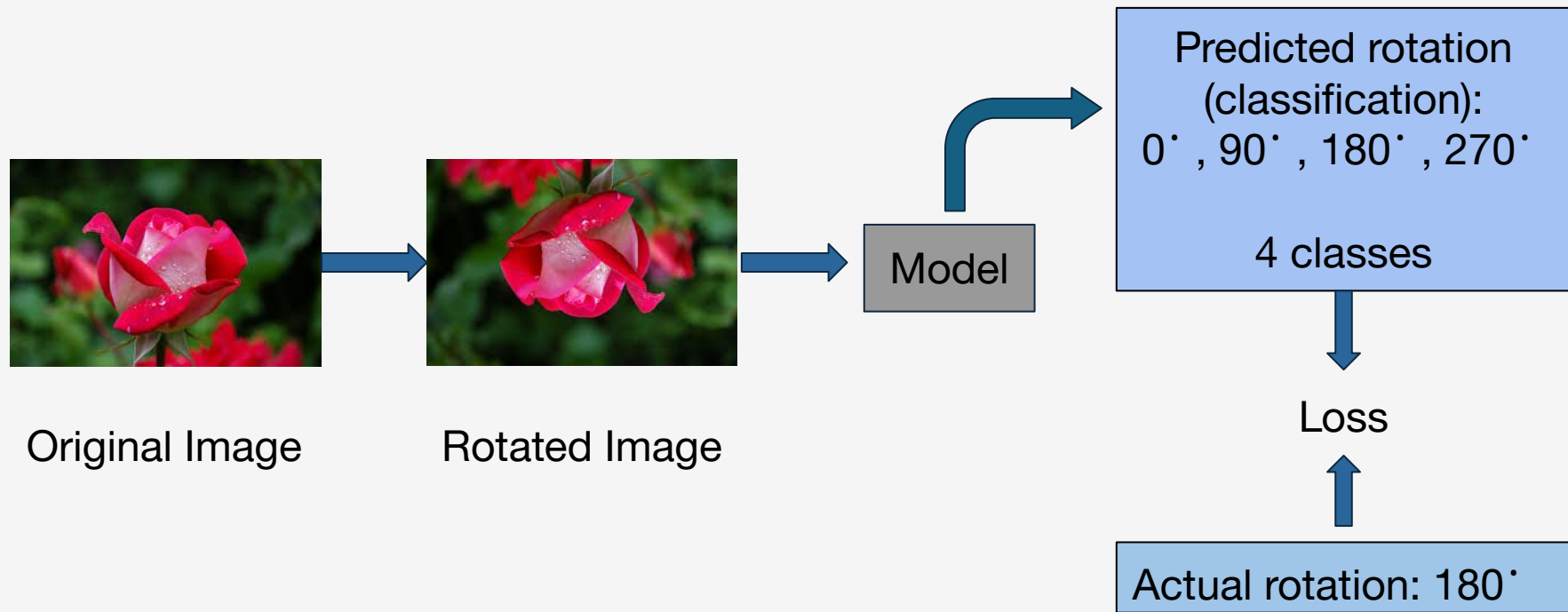
# Pre Training Task: Image Inpainting



# Pre Training Task: Jigsaw Puzzle



# Pre Training Task: Geometric Transformation Recognition



# Pre Training Task: SimCLR



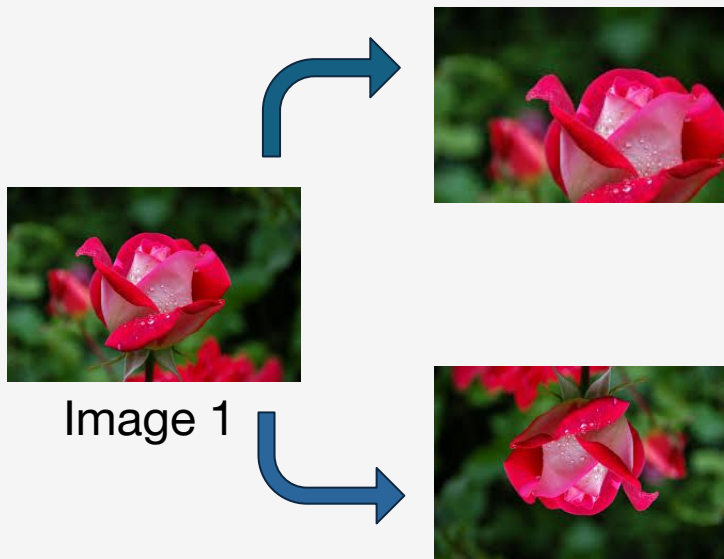
Image 1



Image 2



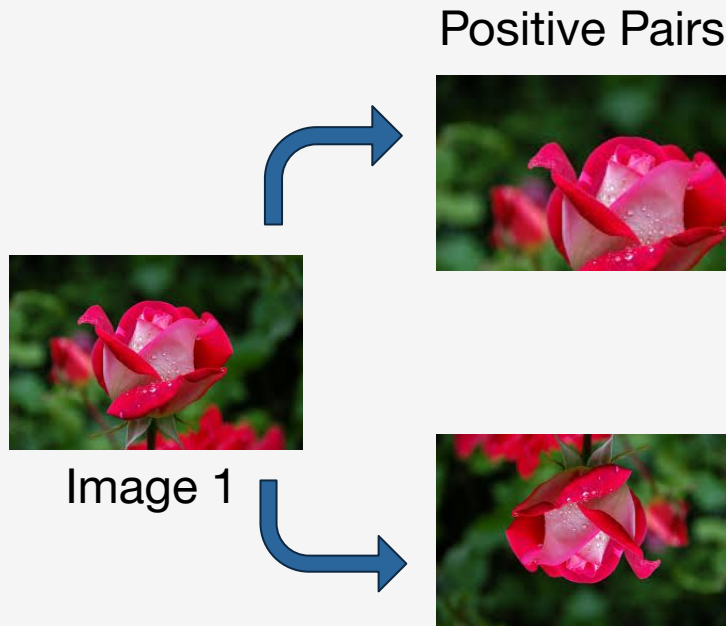
# Pre Training Task: SimCLR



**Transformation 1:**  
Cropped and  
resized image

**Transformation 2:**  
Rotated Image

# Pre Training Task: SimCLR



- Both images are generated **from same image** using transformations like rotation, crop, etc.
- So their **embeddings must be similar**
- By forcing model to generate same embedding we are forcing it to learn the pattern of data

# Pre Training Task: SimCLR

Negative Pairs



Image 1

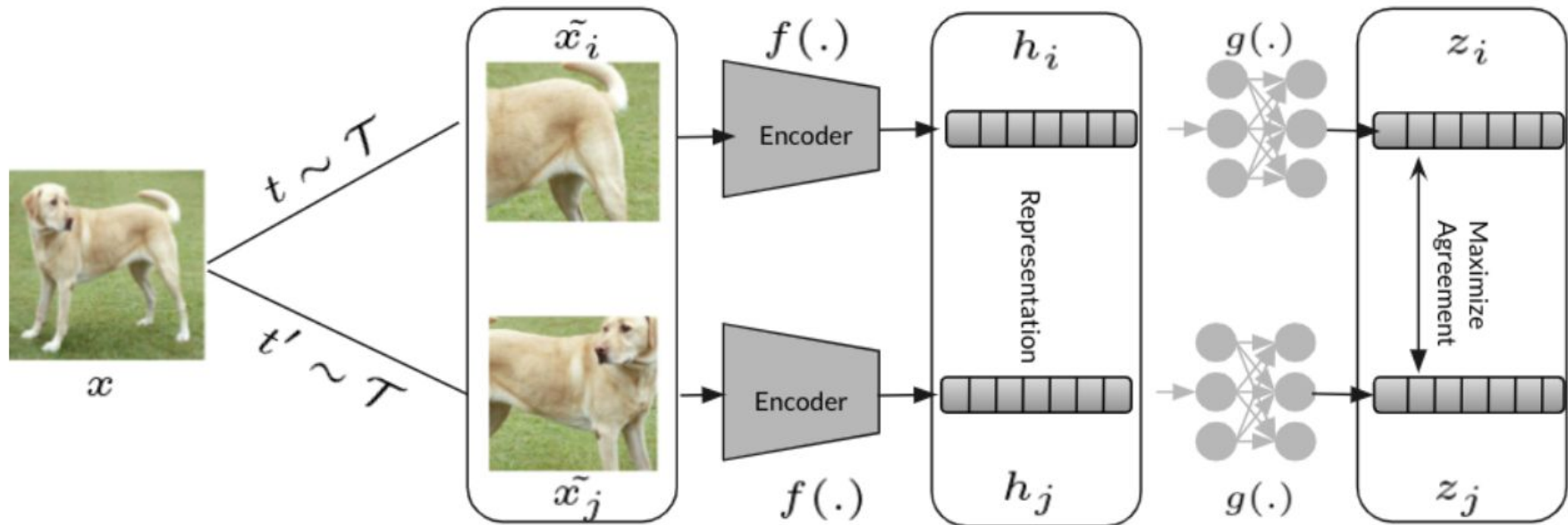


Image 2



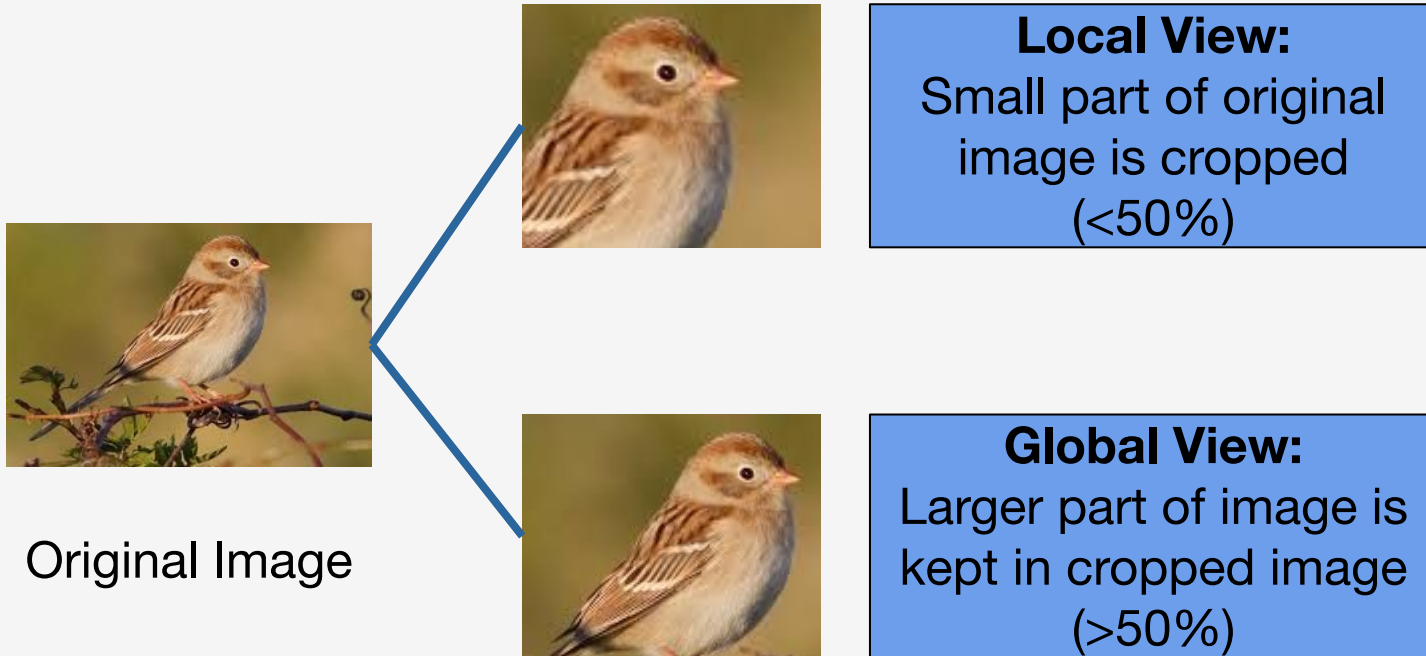
- Both transformations are generated **from different images**
- So here their **embedding must be very different**
- This is the approach of SimCLR (Simple **Contrastive Learning** of Visual Representations)

# Pre Training Task: SimCLR

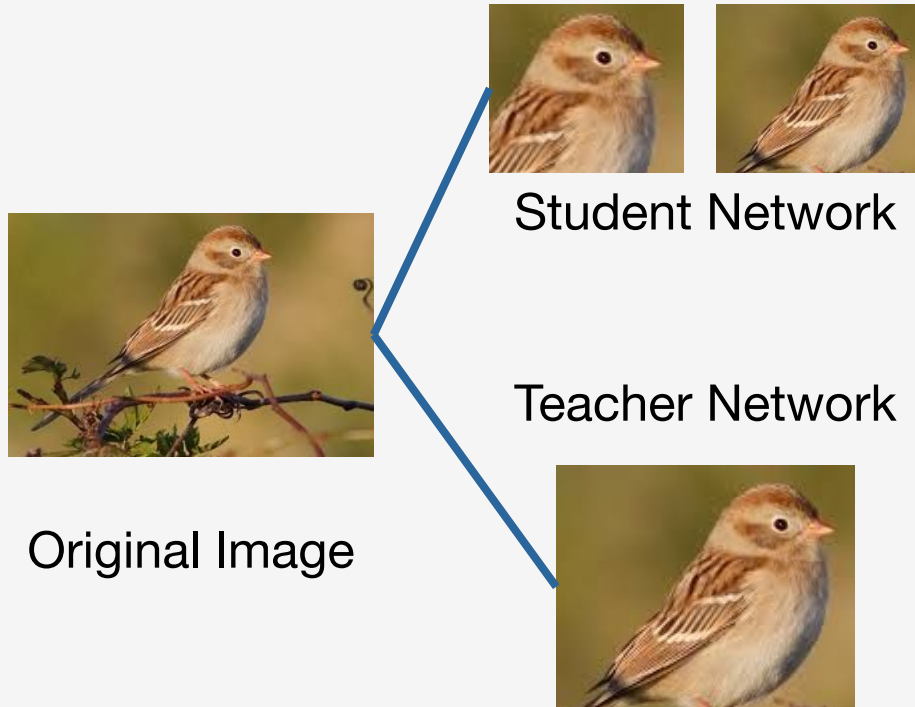


<https://vasudev-sharma.github.io/posts/2022/03/SimCLR-visually-explained/>

# Pre Training Task: DINO



# Pre Training Task: DINO



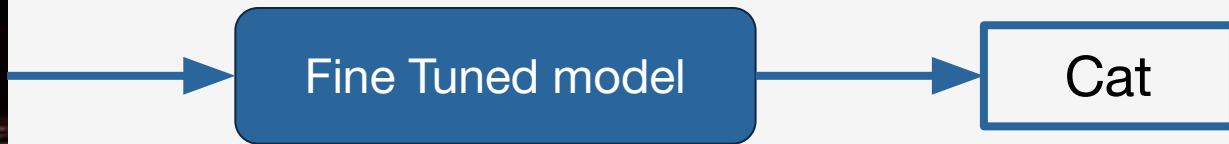
- **Student:** local + global views
- **Teacher:** only global views
- Teacher gives stable targets
- Student learns by matching teacher's output

# Downstream Task:

- We have huge data but from that very less 10-20% data is labeled data and other is unlabeled.
- A task with **labeled data** used to **evaluate or fine-tune** the model after pretraining
- Tasks:
  - Classification
  - Segmentation
  - Detection

# Downstream Task: Classification

Assign label to input image

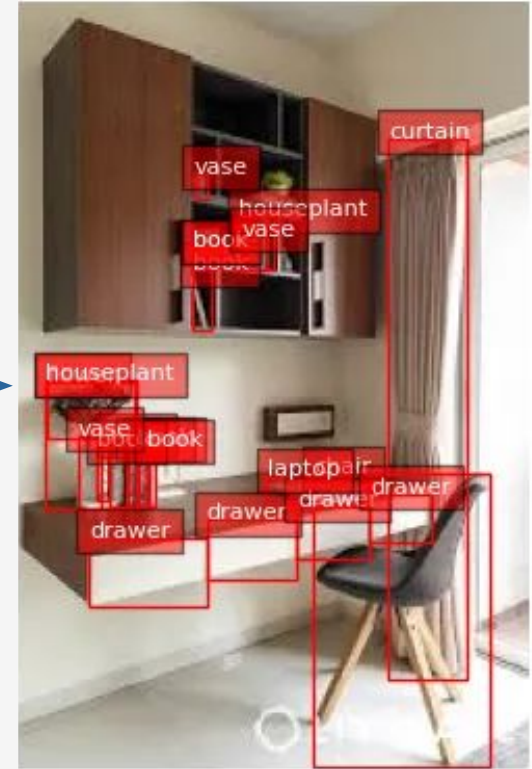




# Downstream Task: Object Detection



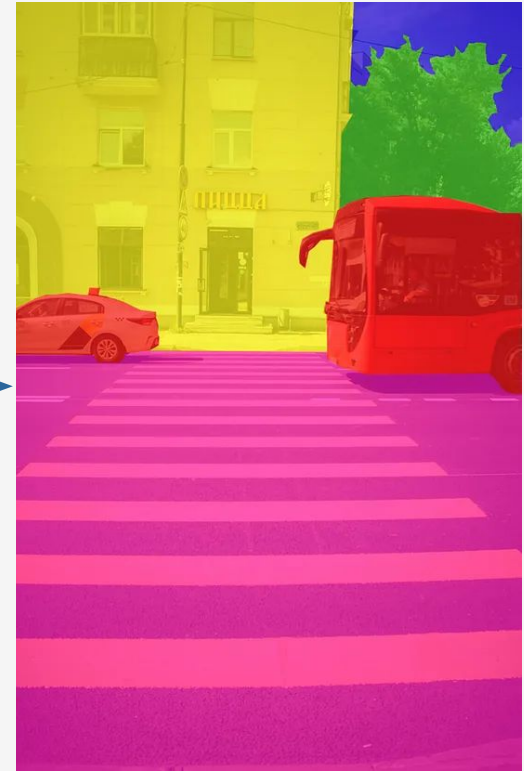
Fine Tuned  
model



# Downstream Task: Segmentation



Fine Tuned  
model



# Solo - learn library

