# Self Supervision Techniques in CNNs
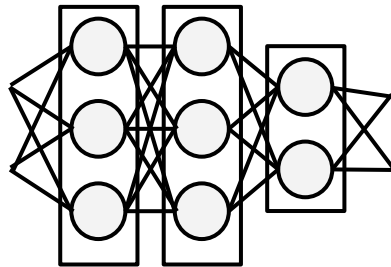
Varsha S
193079005

# Outline

1. Motivation

2. Self supervision

3. Pretext Task

   a. Inpainting

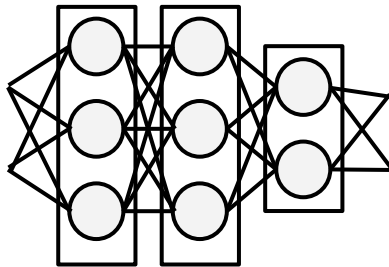   b. Jigsaw Puzzles

# Motivation



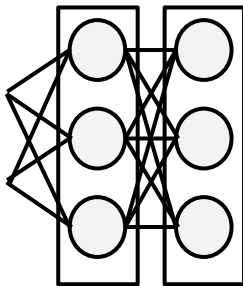Deep learning + ImageNet

Golden Retriever

# Motivation


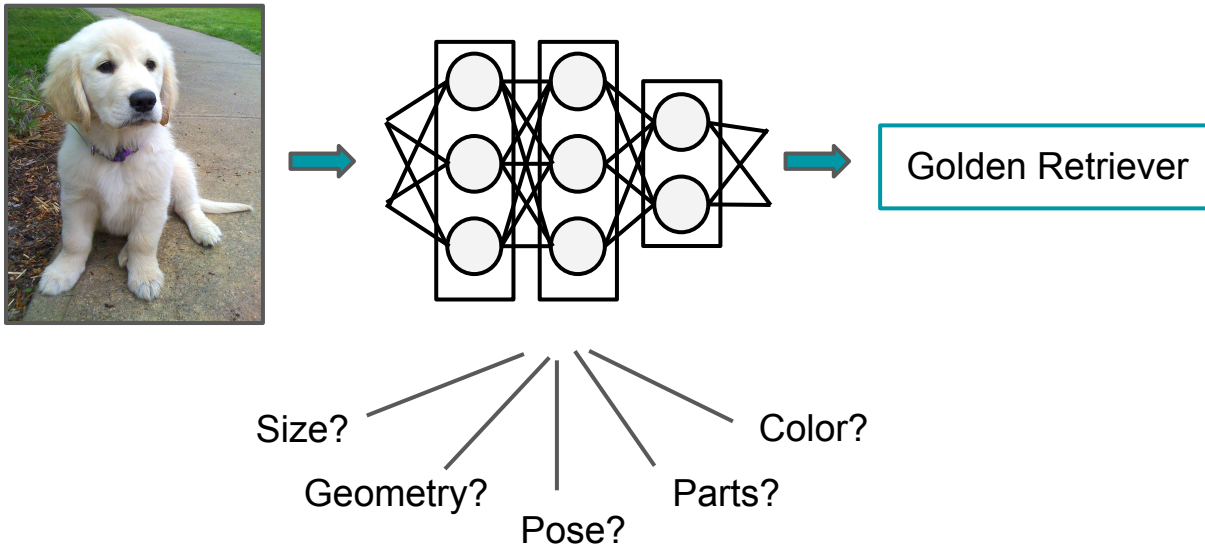
Golden Retriever

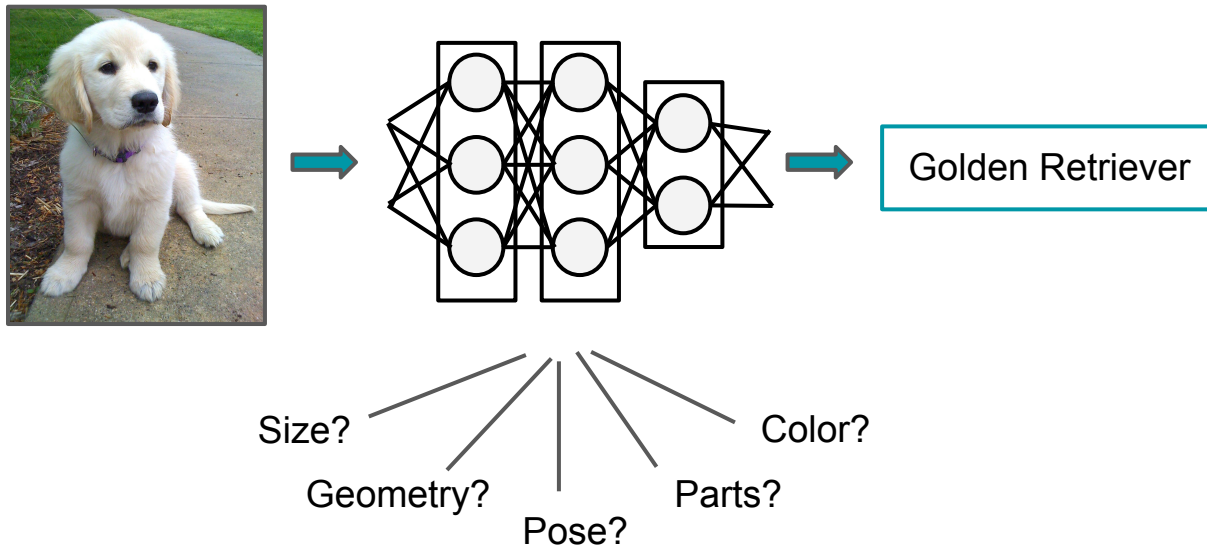Image Segmentation
Detection
Depth Estimation
...

# Motivation



Size? Geometry? Pose? Parts? Color?

Golden Retriever

# Motivation



Size?
Geometry?
Pose?
Parts?
Color?

Golden Retriever

*Can the task be something else ?*

# Motivation
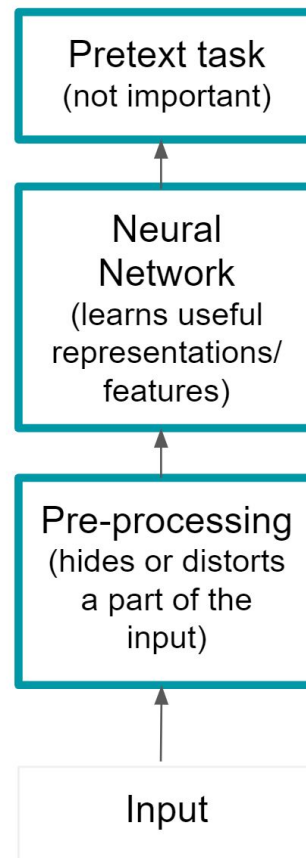


Golden Retriever

Size? Geometry? Pose? Parts? Color?

*Are the labels necessary?*

# Self Supervision

- Data provides supervision
- Goal - Learn good representations
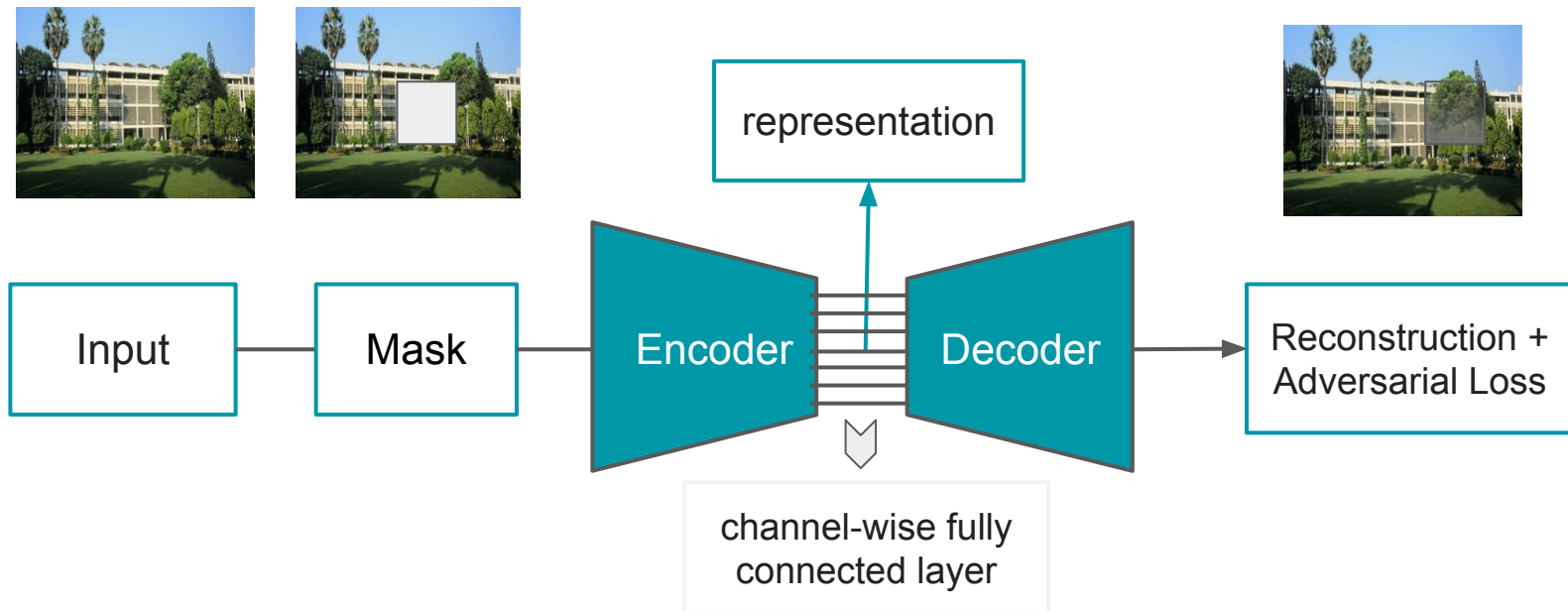- Task - Design pretext

Pretext task
(not important)

Neural
Network
(learns useful
representations/
features)

Pre-processing
(hides or distorts
a part of the
input)

Input

# PRETEXT TASK

## A. INPAINTING

# Inpainting - Context encoders

# Inpainting - Context encoders

# Inpainting - Context encoders



Input → Mask → Encoder → Decoder → Reconstruction + Adversarial Loss

representation

channel-wise fully connected layer

*Image source:-Google images*

# Inpainting - Loss function

- $\mathcal{L} = \lambda_{rec}\mathcal{L}_{rec} + \lambda_{adv}\mathcal{L}_{adv}.$

- $\mathcal{L}_{rec}(x) = \|\hat{M} \odot (x - F((1-\hat{M}) \odot x))\|_2^2,$

- $\mathcal{L}_{adv} = \max_{D} \; \mathbb{E}_{x \in \mathcal{X}}[\log(D(x)) + \log(1 - D(F((1-\hat{M}) \odot x)))]$



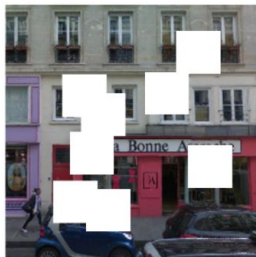a) Masked Input



b) L2 Loss    c) L2 + adversarial loss

\* where x is the ground truth image, F is the context encoder, M is a binary mask corresponding to the dropped image region with a value of 1 wherever a pixel was dropped and 0 for input pixels.

# Inpainting - Masks



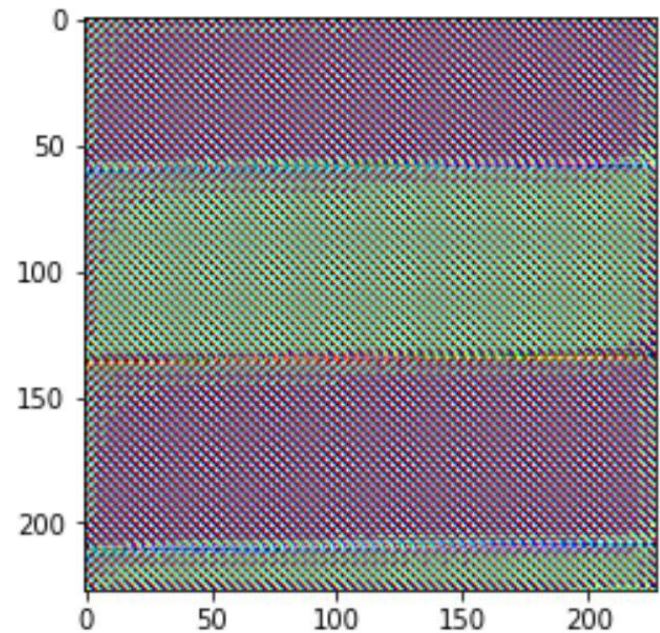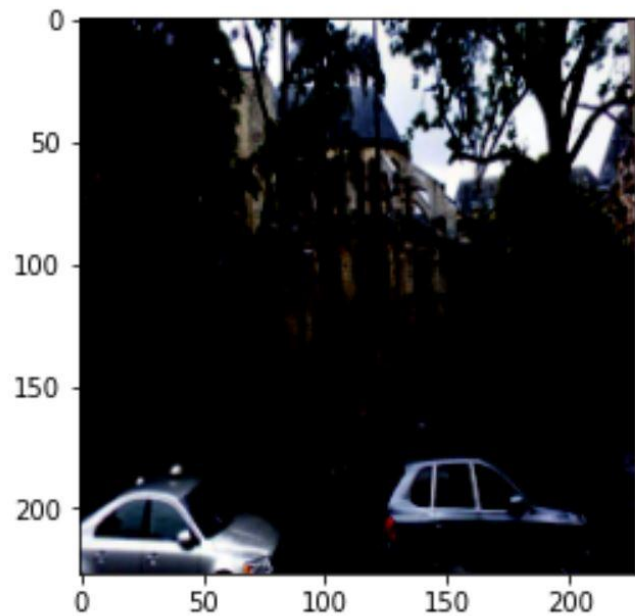Center region          Random blocks          Random region

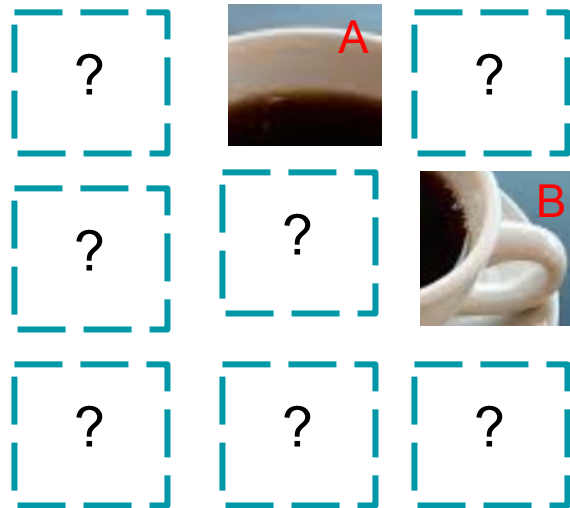# Inpainting - Results

# PRETEXT TASK

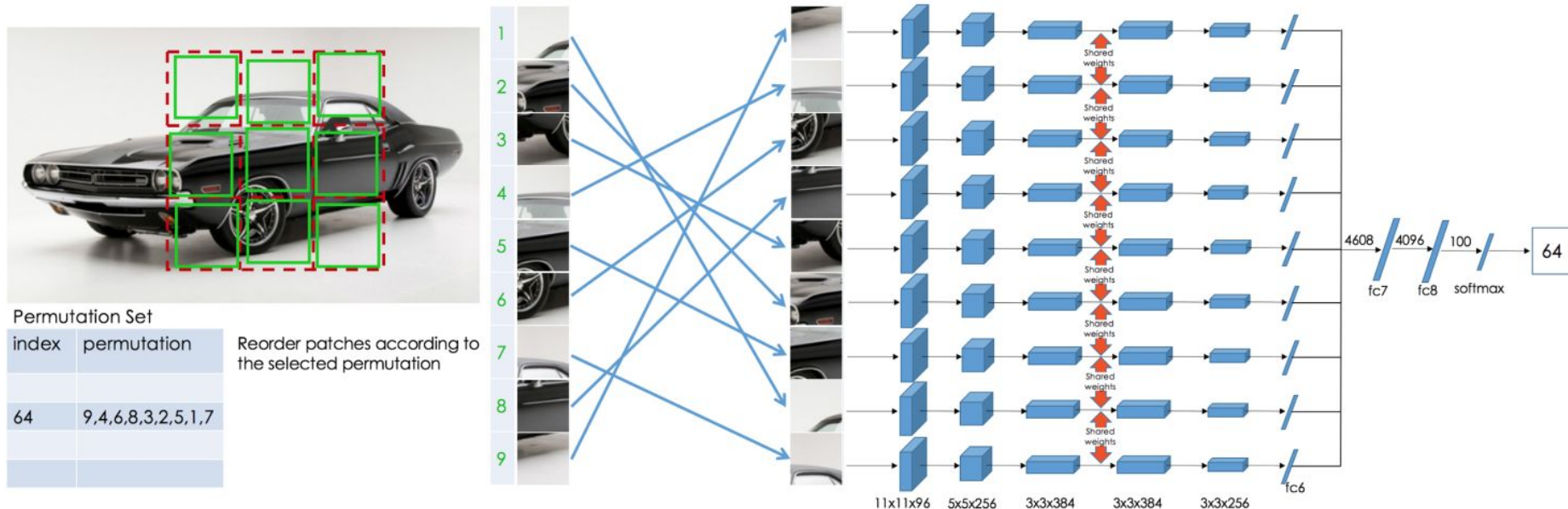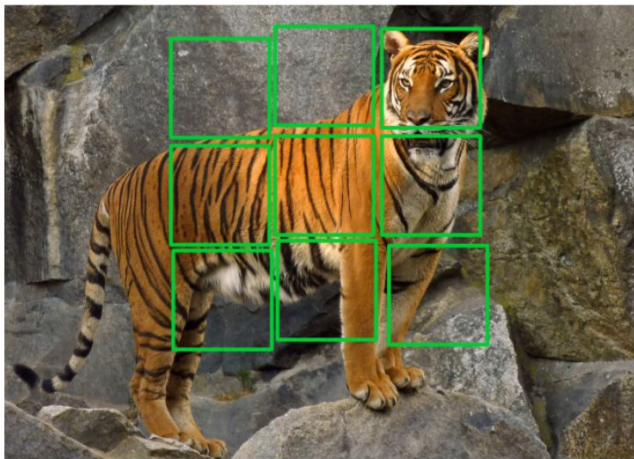# B. JIGSAW PUZZLES

# Inpainting


A


B

# Inpainting

# Inpainting

# Jigsaw - Implementation

- Siamese network
- Permutations with large Hamming distance



Permutation Set

| index | permutation |
|-------|-------------|
| 64    | 9,4,6,8,3,2,5,1,7 |

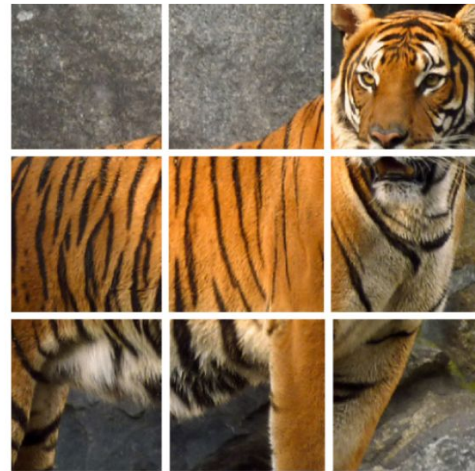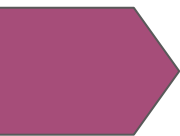Reorder patches according to the selected permutation

# Jigsaw puzzle



The image from which the tiles (marked with green lines) are extracted

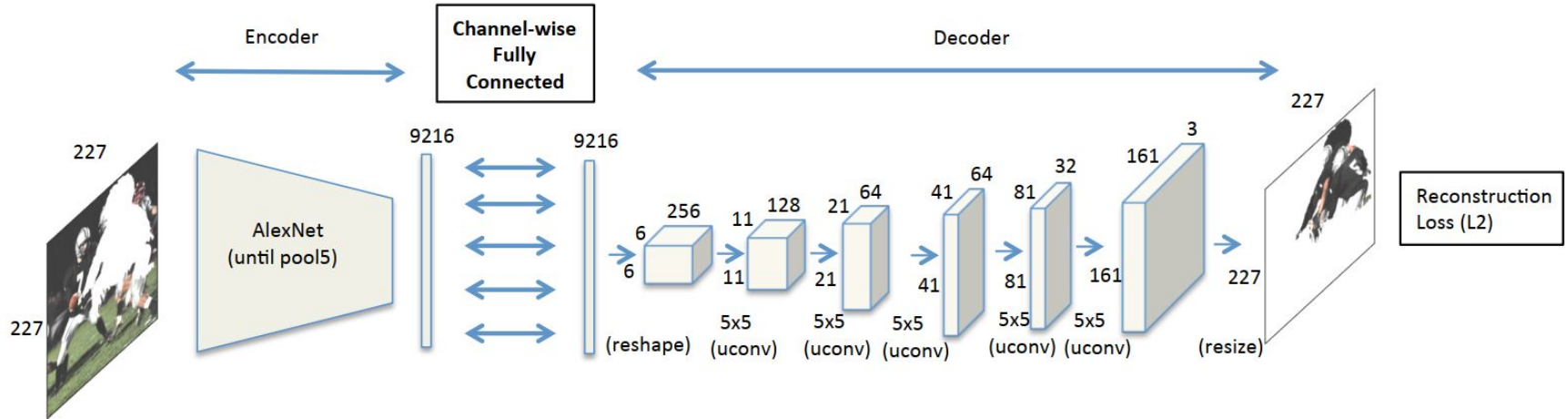Puzzle obtained by shuffling the tiles

Expected output

THANK YOU!

# References

1. Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell and Alexei A. Efros, "Context Encoders: Feature Learning by Inpainting", CVPR 2016.
2. Noroozi, M. and P. Favaro. "Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles", ECCV (2016).
3. Carl Doersch, Abhinav Gupta, and Alexei A. Efros. "Unsupervised Visual Representation Learning by Context Prediction", ICCV 2015
4. Carl Doersch ICCV presentation - http://videolectures.net/iccv2015_doersch_visual_representation/

# Inpainting - Architecture



Context encoder trained with reconstruction loss for feature learning by filling in *arbitrary region dropouts* in the input.
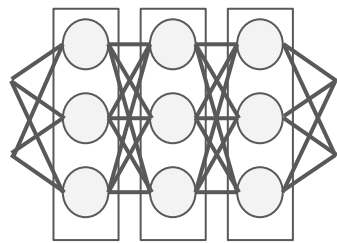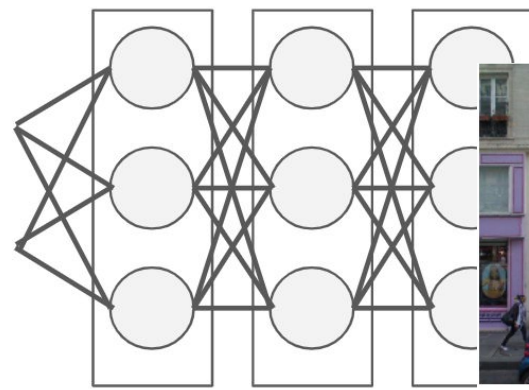
# Motivation

1) In motivation add label intensive data is not available
2) <mark>Add how context helps in classification(beagle example from carl video)</mark>
3) How do these features get used later
4) Inpaintin
   a) <mark>Encoder decoder diagram(shud make)</mark>
   b)
   c) maybe results(last slide)
5) Jigsaw
   a) Siamese network(why it s needed n why it works)
   b)
6)

Outli



Pretext task
(not important)

Neural Network
(learns useful representations/ features)

Pre-processing
(hides or distorts a part of the input)

Input
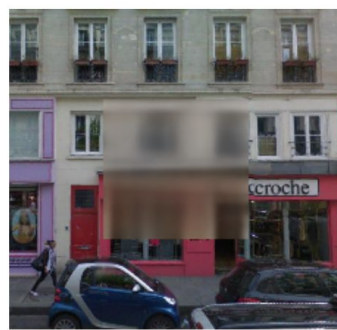
(a) Input context

(b) Human artist

(c) Context Encoder
($L2$ loss)
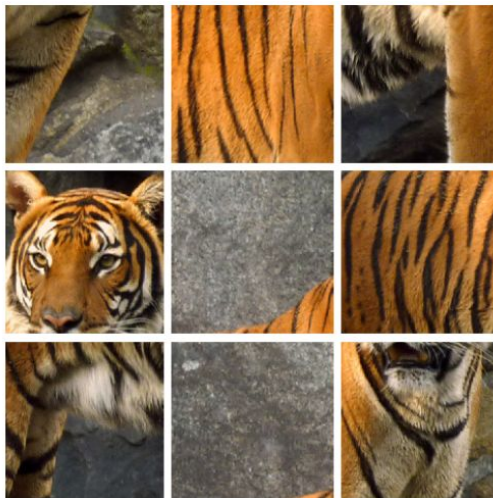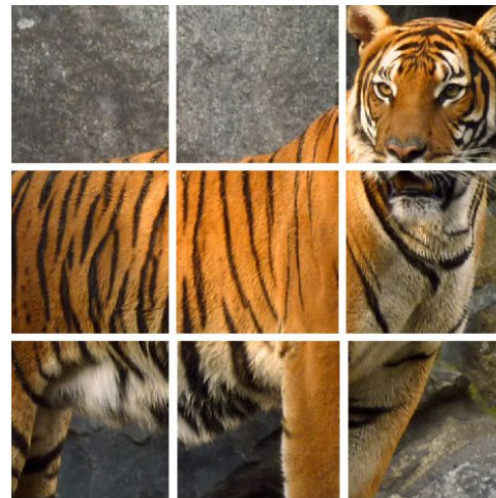
(d) Context Encoder
($L2$ + Adversarial loss)

The image from which the tiles (marked with green lines) are extracted

Puzzle obtained by shuffling the tiles

Expected output

Image source:- Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles  - Noroozi et al