

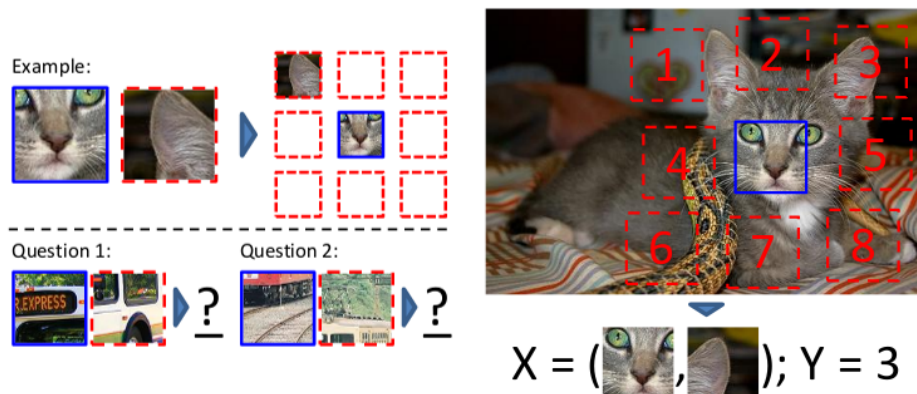
# Unsupervised Visual Representation Learning by Context Prediction

(2016)

Carl Doersch, Abhinav Gupta, Alexei A. Efros  
**Notes**

## Contributions

This paper explores context prediction as an unsupervised task for learning useful representation for other downstream tasks. Using randomly extracted pairs of patches, the training objective is to predict the position of the second patch relative to the first patch. More specifically, we sample random pairs of patches and present them to each pair as a training examples, the model must predict the correct position of one patch given the position of the other patch without any additional information about their positions in the original image.



The authors argue that by training on this simple task, the learned representations allow us to perform unsupervised visual discovery, or even use this model as a starting point for other visual tasks instead of using a randomly initialized model, which gives better results without the need for labeled examples.

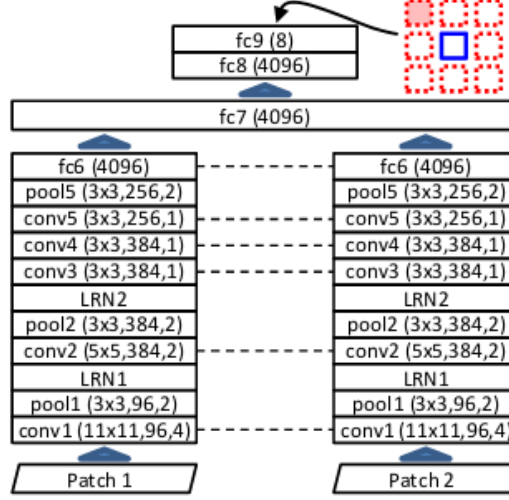
## Visual Context Prediction

When choosing a given task, we must ensure that the task enforce the model to learn useful representations and not simply taking shortcuts to get to the correct prediction.

In our case, we must not boundary patterns or texture continuity is obvious. To avoid this the patches are selected with pair wise distances up to 0.5 of the size of each patch, with an additional jittering of the positioning up to 7 pixels, one additional thing the authors observed the network uses to cheat is the Chromatic aberration, in some cameras one color channel is focused toward the center the model uses to localize the patches to the lens itself and find their relative positions, to overcome this the authors propose to randomly drop 2 of the 3 color channels.

After constructing a training dataset as a set of pairs  $(\text{patch}_1, \text{patch}_2)$  as inputs, and the position of the second pair relative to the first as targets, we can use a CNN similar to AlexNet for training,

we must feed the two patches into the network, concatenate the two output into one 4096 vector, and use two fully connected layers to output a probability distribution over the 8 possible positions.



## Experiments

**Pre-training** We can use our self-supervised task to pretrain a model for a given visual task object detection in our case, and we see competitive results compared to pretraining using ImageNet:

VOC-2007 Test	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP
DPM-v5[17]	33.2	60.3	10.2	16.1	27.3	54.3	58.2	23.0	20.0	24.1	26.7	12.7	58.1	48.2	43.2	12.0	21.1	36.1	46.0	43.5	33.7
[8] w/o context	52.6	52.6	19.2	25.4	18.7	47.3	56.9	42.1	16.6	41.4	41.9	27.7	47.9	51.5	29.9	20.0	41.1	36.4	48.6	53.2	38.5
Regionlets[58]	54.2	52.0	20.3	24.0	20.1	55.5	68.7	42.6	19.2	44.2	49.1	26.6	57.0	54.5	43.4	16.4	36.6	37.7	59.4	52.3	41.7
Scratch-R-CNN[2]	49.9	60.6	24.7	23.7	20.3	52.5	64.8	32.9	20.4	43.5	34.2	29.9	49.0	60.4	47.5	28.0	42.3	28.6	51.2	50.0	40.7
Scratch-Ours	52.6	60.5	23.8	24.3	18.1	50.6	65.9	29.2	19.5	43.5	35.2	27.6	46.5	59.4	46.5	25.6	42.4	23.5	50.0	50.6	39.8
Ours-projection	58.4	62.8	33.5	27.7	24.4	58.5	68.5	41.2	26.3	49.5	42.6	37.3	55.7	62.5	49.4	29.0	47.5	28.4	54.7	56.8	45.7
Ours-color-dropping	60.5	66.5	29.6	28.5	26.3	56.1	70.4	44.8	24.6	45.5	45.4	35.1	52.2	60.2	50.0	28.1	46.7	42.6	54.8	58.6	46.3
Ours-Yahoo100m	56.2	63.9	29.8	27.8	23.9	57.4	69.8	35.6	23.7	47.4	43.0	29.5	52.9	62.0	48.7	28.4	45.1	33.6	49.0	55.5	44.2
ImageNet-R-CNN[21]	64.2	69.7	50	41.9	32.0	62.6	71.0	60.7	32.7	58.5	46.5	56.1	60.6	66.8	54.2	31.5	52.8	48.9	57.9	64.7	54.2
K-means-rescale [31]	55.7	60.9	27.9	30.9	12.0	59.1	63.7	47.0	21.4	45.2	55.8	40.3	67.5	61.2	48.3	21.9	32.8	46.9	61.6	51.7	45.6
Ours-rescale [31]	61.9	63.3	35.8	32.6	17.2	68.0	67.9	54.8	29.6	52.4	62.9	51.3	67.1	64.3	50.5	24.4	43.7	54.9	67.1	52.7	51.1
ImageNet-rescale [31]	64.0	69.6	53.2	44.4	24.9	65.7	69.6	69.2	28.9	63.6	62.8	63.9	73.3	64.6	55.8	25.7	50.5	55.4	69.3	56.4	56.5
VGG-K-means-rescale	56.1	58.6	23.3	25.7	12.8	57.8	61.2	45.2	21.4	47.1	39.5	35.6	60.1	61.4	44.9	17.3	37.7	33.2	57.9	51.2	42.4
VGG-Ours-rescale	71.1	72.4	54.1	48.2	29.9	75.2	78.0	71.9	38.3	60.5	62.3	68.1	74.3	74.2	64.8	32.6	56.5	66.4	74.0	60.3	61.7
VGG-ImageNet-rescale	76.6	79.6	68.5	57.4	40.8	79.9	78.4	85.4	41.7	77.0	69.3	80.1	78.6	74.6	70.1	37.5	66.0	67.5	77.4	64.9	68.6

Table 1. Mean Average Precision on VOC-2007.

**Nearest Neighbors** To check if similar patches do indeed have similar representations, we can randomly sample a given number of 96 x 96 patches, and represent them using the FC6 features, and for a given patch, we find the nearest neighbors based on the FC6 features using a normalized correlation, the results are in the figure below:

