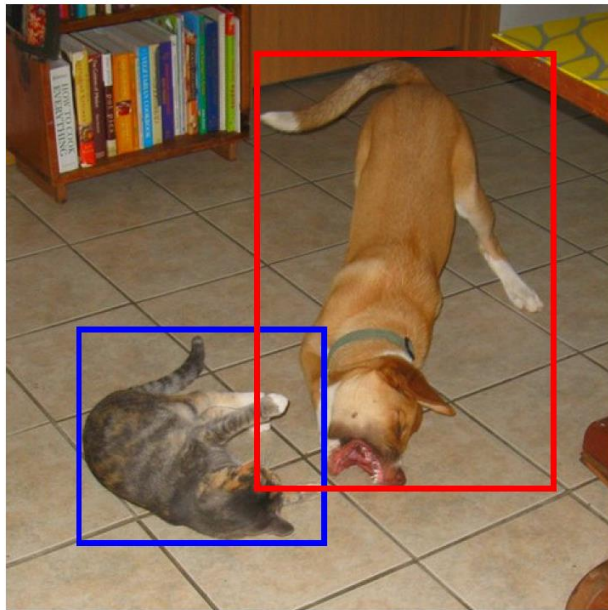# "Unsupervised" Deep Learning and Deep Style Transfer
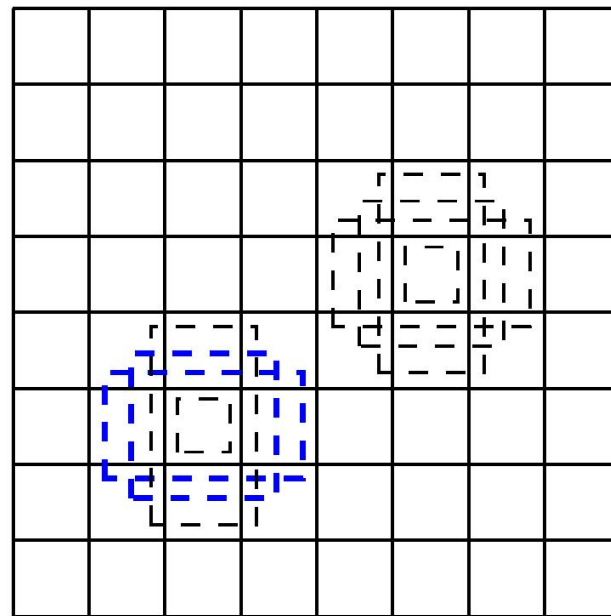
## James Hays
slides from Carl Doersch and Mark Chang
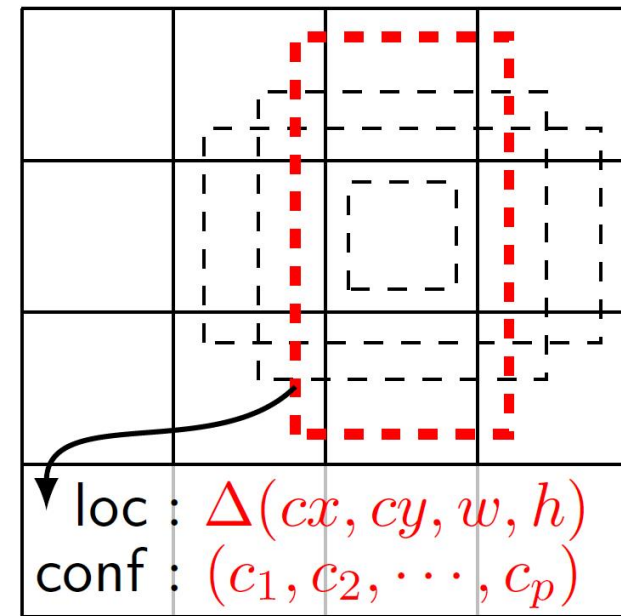
# Recap from Previous Lecture

- We saw two strategies to get *structured* output while using deep learning
  - With object detection, one strategy is brute force: detect everywhere at once



(a) Image with GT boxes

(b) $8 \times 8$ feature map

(c) $4 \times 4$ feature map

loc : $\Delta(cx, cy, w, h)$
conf : $(c_1, c_2, \cdots, c_p)$

# Recap from Previous Lecture

- We saw two strategies to get *structured* output while using deep learning
  - With pose estimation / keypoint detection, the network produces an image-based intermediate representation



Part Detection

Part Association

# Recap from Previous Lecture

- More generally, it can pay off to get creative. Even if Deep ConvNets aren't a natural fit for an image-related task, they might be able to learn a subtask or create a useful intermediate representation.
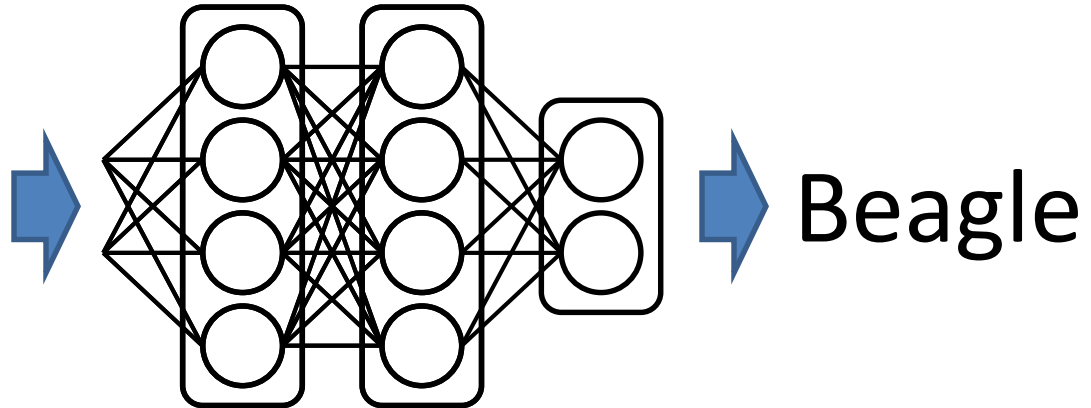
# Unsupervised Visual Representation Learning by Context Prediction

Carl Doersch
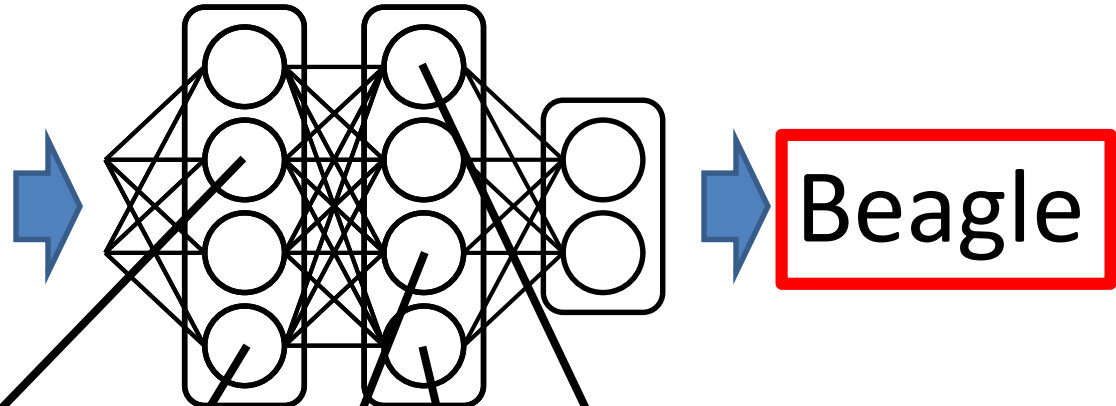Joint work with Alexei A. Efros & Abhinav Gupta

ICCV 2015

# ImageNet + Deep Learning



Beagle

- Image Retrieval
- Detection (RCNN)
- Segmentation (FCN)
- Depth Estimation
- …

# Context as Supervision
[Collobert & Weston 2008; Mikolov et al. 2013]

Context Prediction for Images
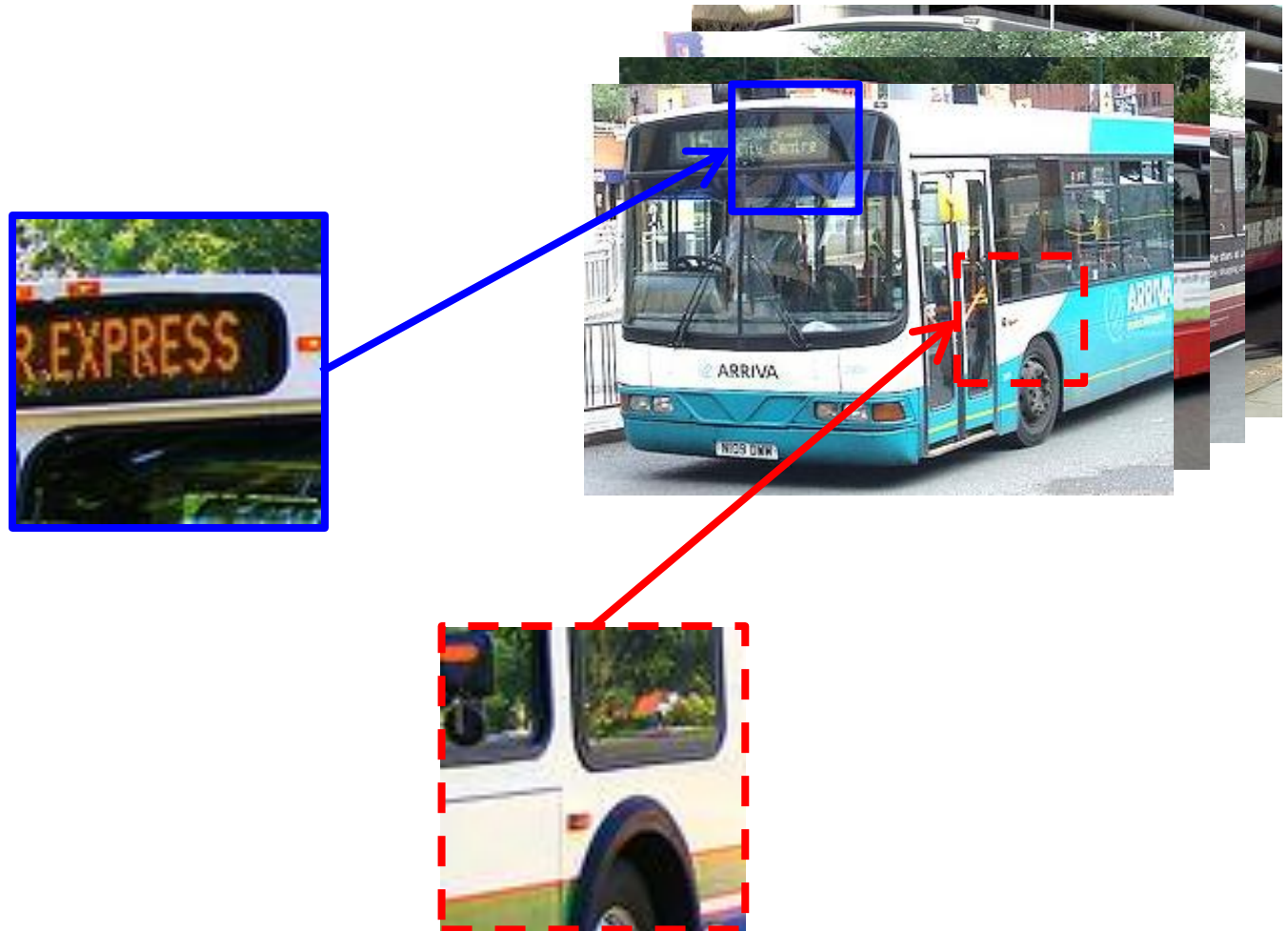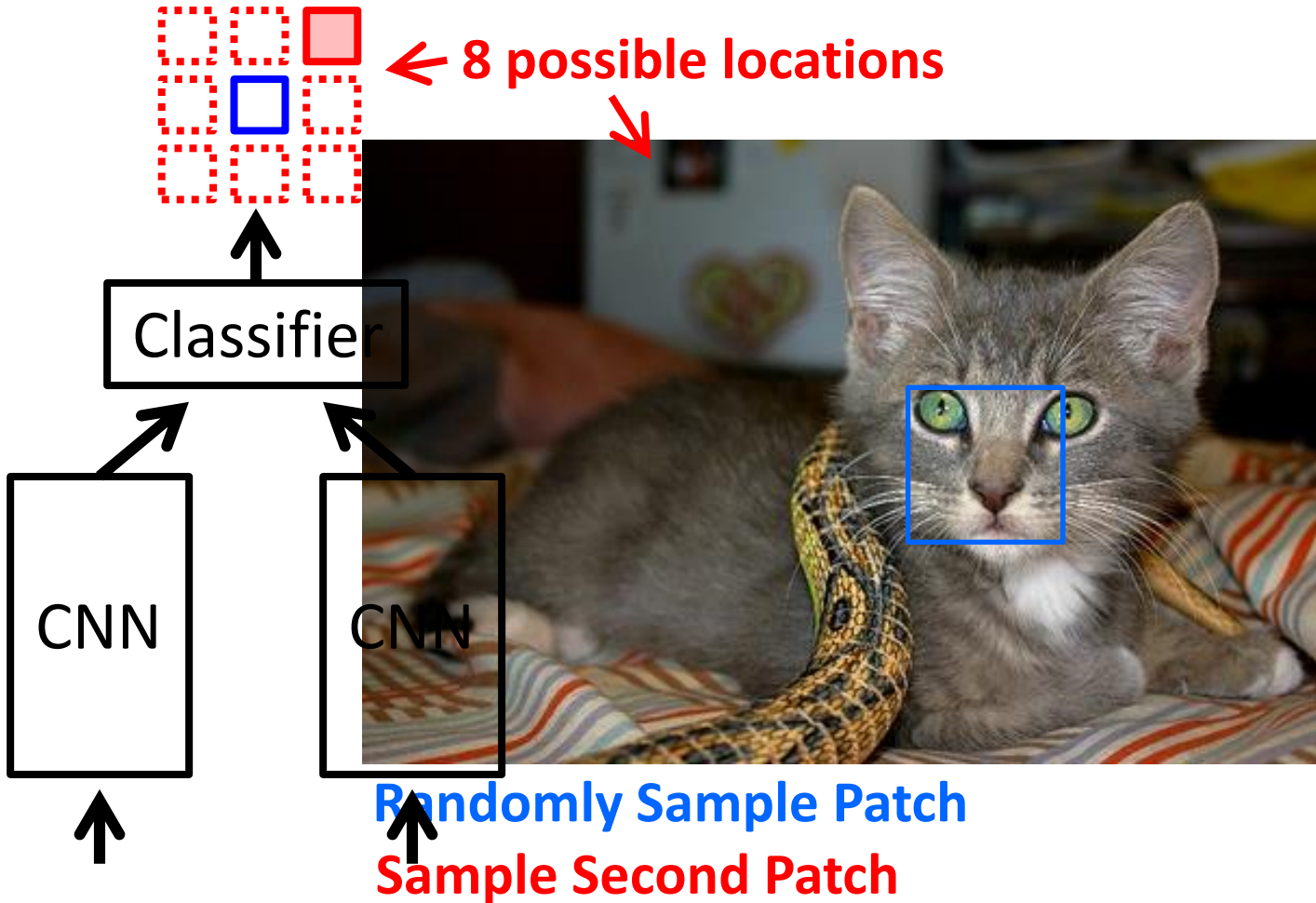
# Semantics from a non-semantic task

# Relative Position Task



8 possible locations

Classifier

CNN

CNN

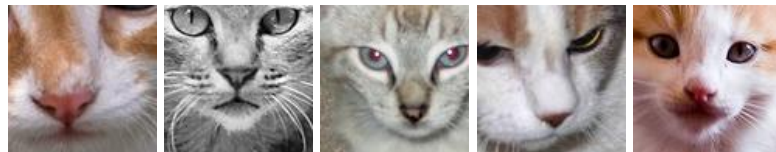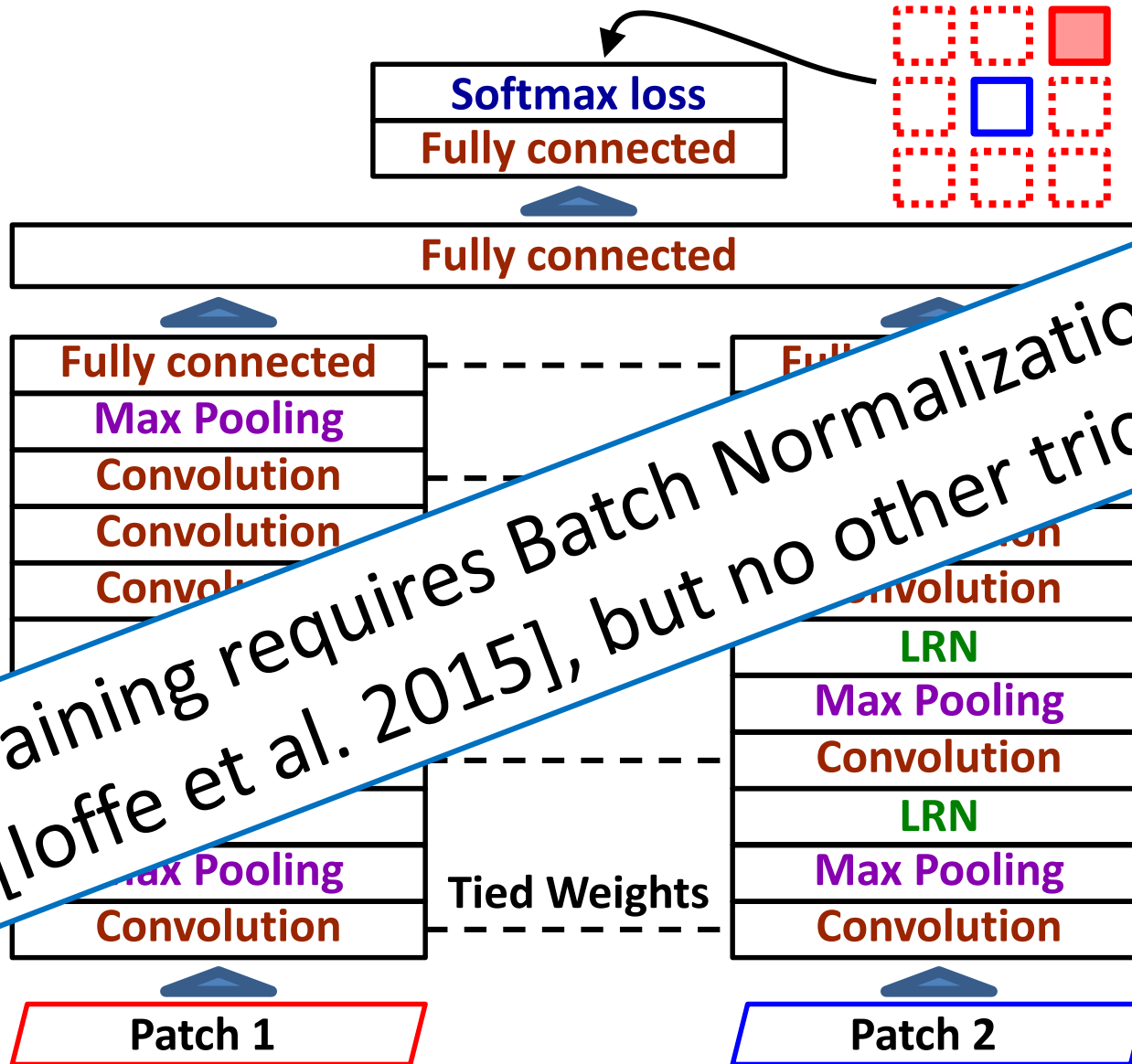Randomly Sample Patch

Sample Second Patch

Patch Embedding

Classifier

CNN

CNN

Input

Nearest Neighbors

Note: connects *across* instances!

# Architecture



Softmax loss

Fully connected

Fully connected

**Patch 1:**
Fully connected
Max Pooling
Convolution
Convolution
Convolution
Max Pooling
Convolution

**Patch 2:**
Fully connected
Convolution
LRN
Max Pooling
Convolution
LRN
Max Pooling
Convolution

Tied Weights

Patch 1

Patch 2

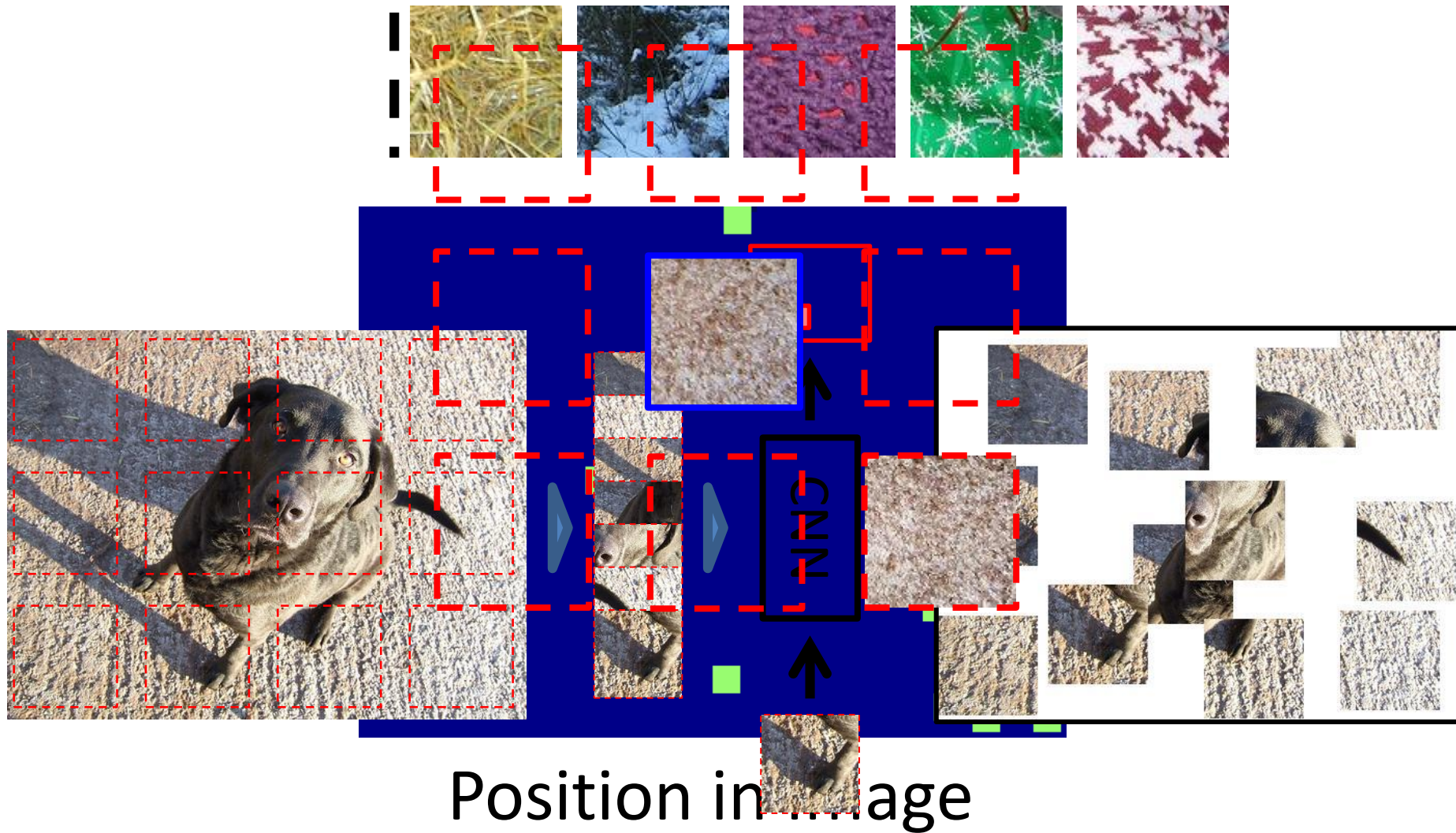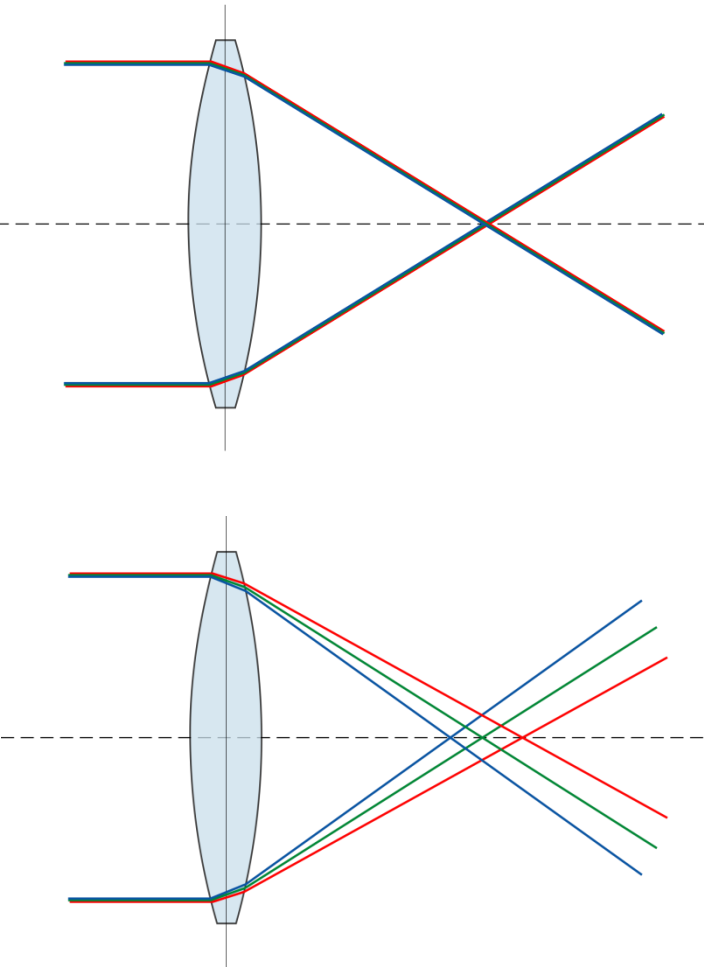Training requires Batch Normalization [Ioffe et al. 2015], but no other tricks

# Avoiding Trivial Shortcuts
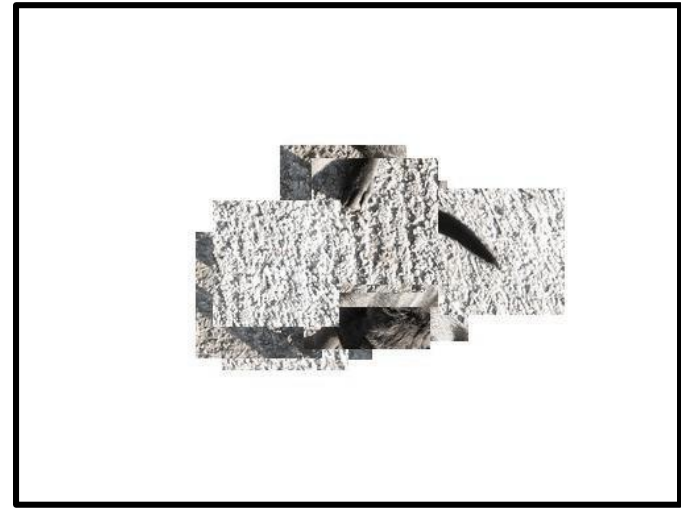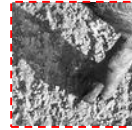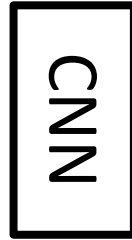


Include a gap

Jitter the patch locations

# A Not-So "Trivial" Shortcut



Position in image

# Chromatic Aberration

# Chromatic Aberration

# What is learned?



Input          Ours          Random Initialization          ImageNet AlexNet

# Still don't capture everything

Input          Ours          Random Initialization          ImageNet AlexNet



# You don't always need to learn!

Input          Ours          Random Initialization          ImageNet AlexNet

# Mined from Pascal VOC2011

# Pre-Training for R-CNN



1. Input image
2. Extract region proposals (~2k)
3. Compute CNN features
4. Classify regions

warped region

aeroplane? no.
person? yes.
tvmonitor? no.

CNN

Pre-train on relative-position task, w/o labels

[Girshick et al. 2014]

# VOC 2007 Performance
## (pretraining for R-CNN)

% Average Precision

68.6

61.7

56.8

54.2

51.1

46.3

45.6

40.7

42.4

**Legend:**
- No Rescaling
- Krähenbühl et al. 2015
- VGG + Krähenbühl et al.

[Krähenbühl, Doersch, Donahue & Darrell, "Data-dependent Initializations of CNNs", 2015]

ImageNet Labels          Ours          No Pretraining

# Capturing Geometry?

# Surface-normal Estimation



| Method | Error (Lower Better) | | % Good Pixels (Higher Better) | | |
|---|---|---|---|---|---|
| | Mean | Median | 11.25° | 22.5° | 30.0° |
| No Pretraining | 38.6 | 26.5 | 33.1 | 46.8 | 52.5 |
| Ours | **33.2** | 21.3 | 36.0 | 51.2 | 57.8 |
| ImageNet Labels | 33.3 | **20.8** | **36.7** | **51.7** | **58.1** |

*So, do we need semantic labels?*

# "Self-Supervision" and the Future
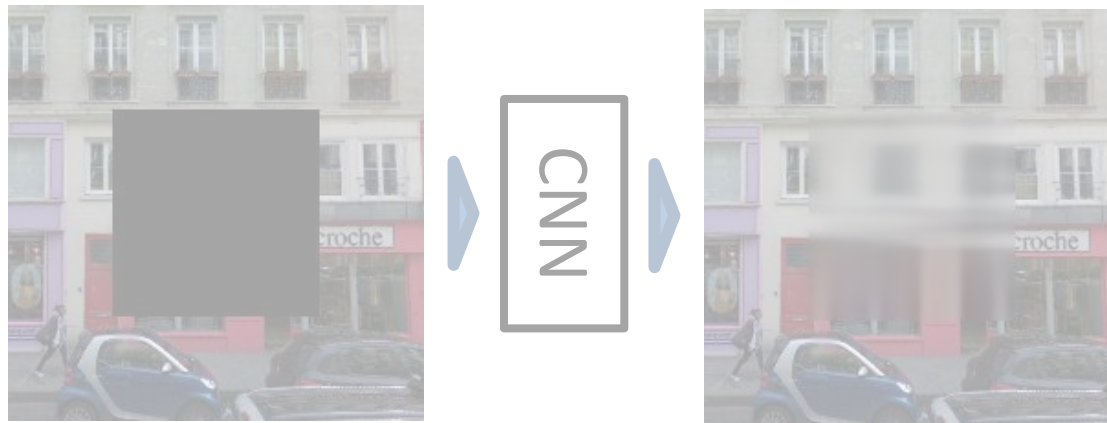
## Ego-Motion



[Agrawal et al. 2015; Jayaraman et al. 2015]

## Video



[Wang et al. 2015; Srivastava et al 2015; ...]

## Context



[Doersch et al. 2014; Pathak et al. 2015; Isola et al. 2015]

# A Neural Algorithm of Artistic Style

Leon A. Gatys, Alexander S. Ecker, Matthias Bethge.
CVPR 2016.

See pdf on course website