# Unsupervised Visual Representation Learning by Context Prediction

## Carl Doersch, Alexei A. Efros, Abhinav Gupta
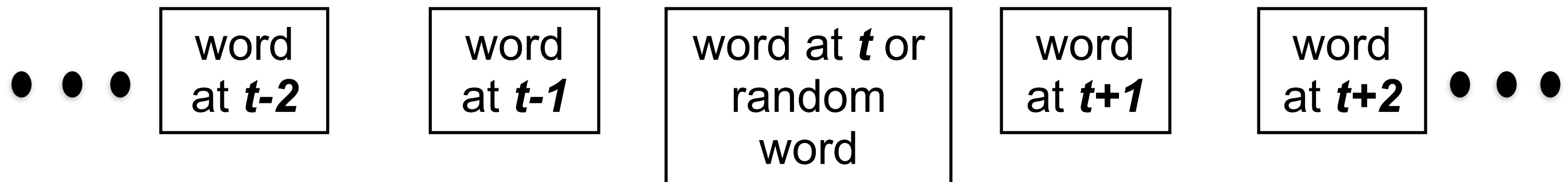
Presented by Maheen Rashid for ECS 289G

# Motivation

- How can we scale to billions rather than millions of images?

  - Imagenet trained on ~1.2 million images

- Unsupervised learning
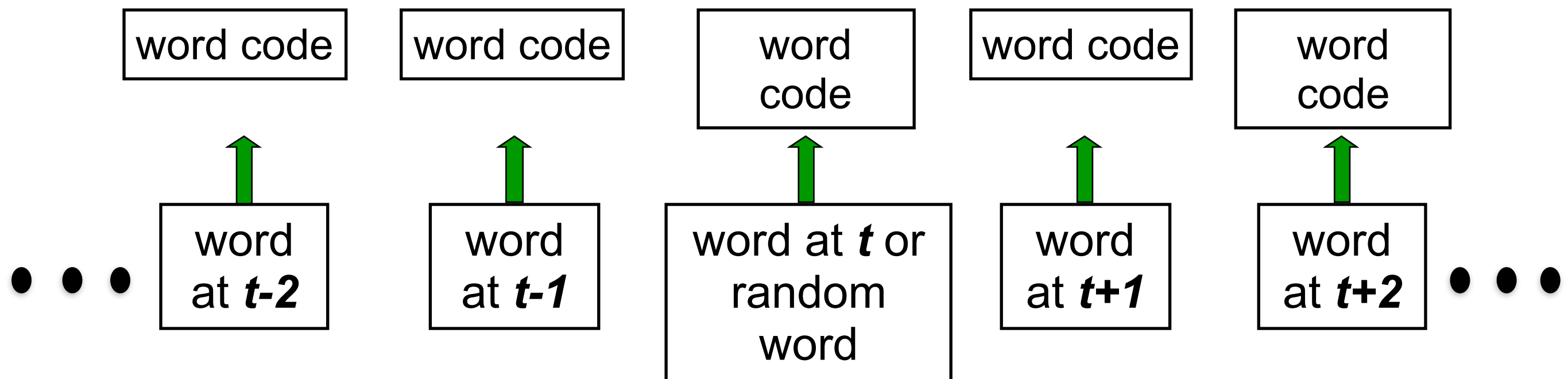
  - Problem - What should be represented?

# Inspiration - Context

- Similar words appear in similar contexts

- Learn to relate a given word to its surrounding words

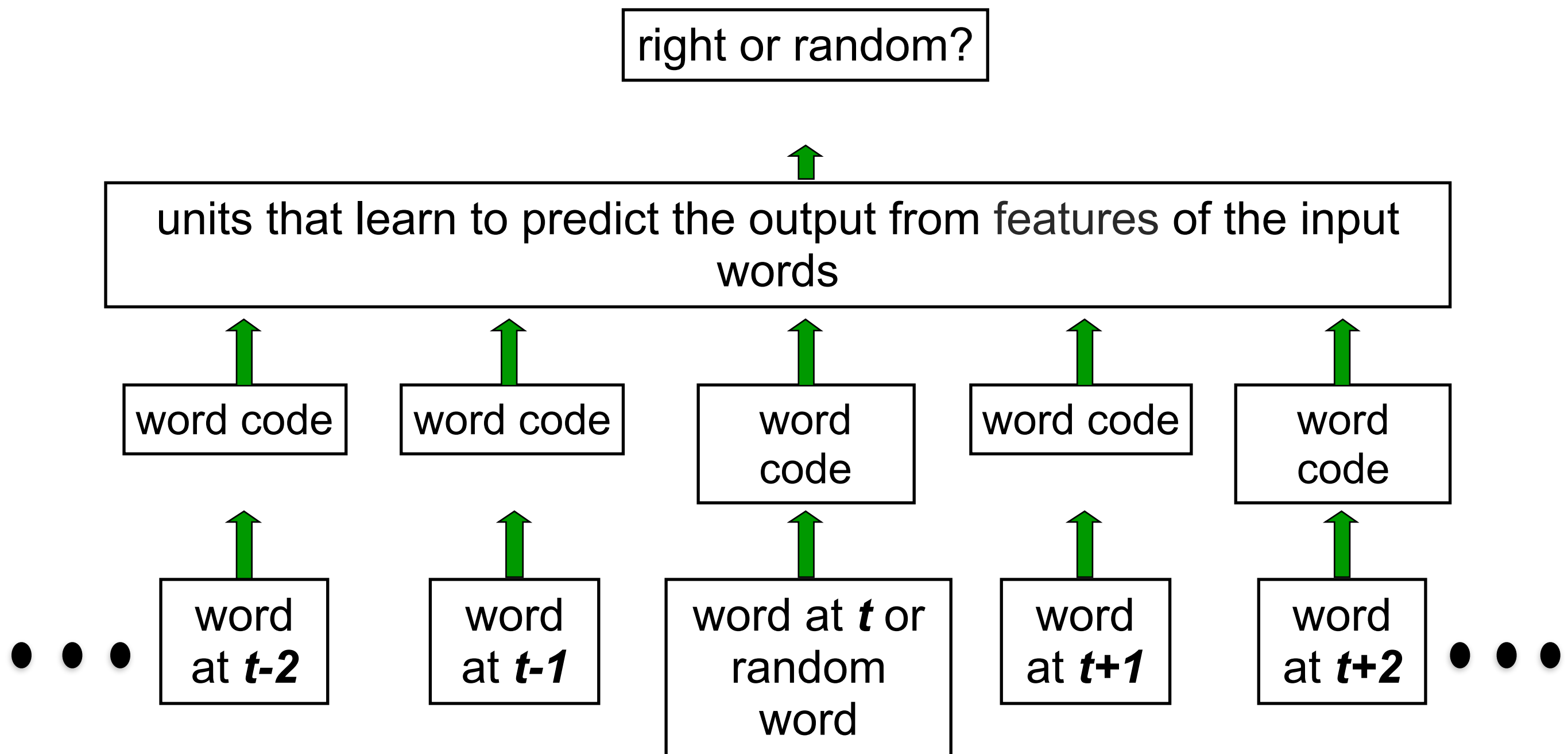- Context prediction becomes a 'pretext' task

# A simple way to learn feature vectors for words (Collobert and Weston, 2008)

● ● ●   | word at **t-2** |   | word at **t-1** |   | word at **t** or random word |   | word at **t+1** |   | word at **t+2** |   ● ● ●

# A simple way to learn feature vectors for words (Collobert and Weston, 2008)



| word code | word code | word code | word code | word code |
|---|---|---|---|---|
| ↑ | ↑ | ↑ | ↑ | ↑ |
| word at *t-2* | word at *t-1* | word at *t* or random word | word at *t+1* | word at *t+2* |

Slide from Geoff Hinton

# A simple way to learn feature vectors for words (Collobert and Weston, 2008)

| right or random? |
| --- |

↑

| units that learn to predict the output from features of the input words |
| --- |

↑        ↑        ↑        ↑        ↑

| word code | | word code | | word code | | word code | | word code |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |

↑        ↑        ↑        ↑        ↑

● ● ● | word at *t-2* | | word at *t-1* | | word at *t* or random word | | word at *t+1* | | word at *t+2* | ● ● ●

# Right or Random for Images?

# Right or Random for Images?

A



B

A    B

1  2  3

4  [EXPRESS]  5

6  7  8

Can you tell where B goes relative to A?

# Answer:

# Answer:





# Doing this requires recognizing semantics!

# Unlabeled training image

# Unlabeled training image



Randomly Sample Patch

# Unlabeled training image



Randomly Sample Patch
Sample Second Patch

# Unlabeled training image



Train Deep Net to recover relative position

CNN

Slide from Carl Doersch

Patch Features



CNN

Slide from Carl Doersch

Input

Nearest Neighbors

Patch Features

CNN

Slide from Carl Doersch

# Architecture

# How to sample patches

# How to sample patches

# How to sample patches

# How to sample patches



Include a gap

# How to sample patches



Include a gap

Jitter the patch locations

# Another trivial shortcut

- Chromatic Aberration
- Shift colors towards grey (Projection)
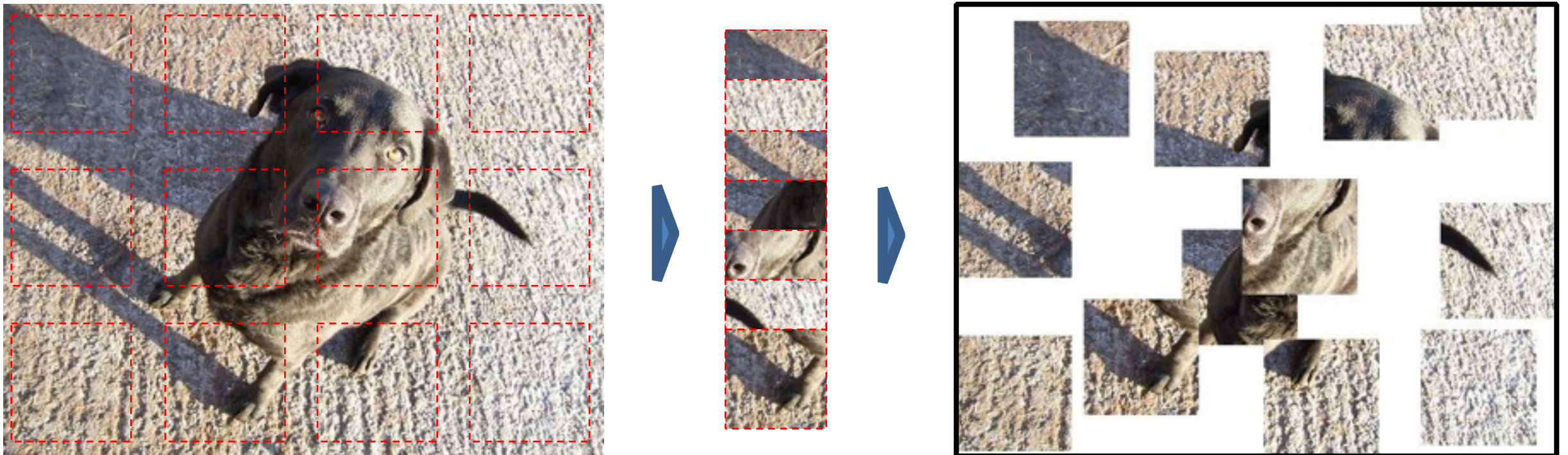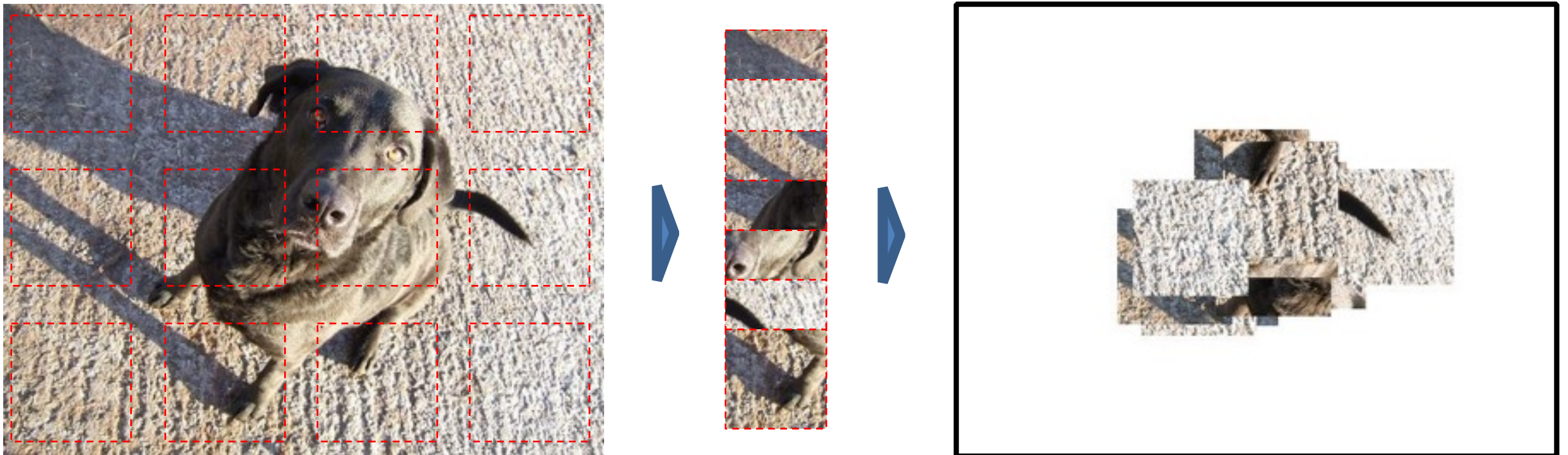- Drop 2 out of three channels during training

Slide from Carl Doersch

# Another trivial shortcut

- Chromatic Aberration
- Shift colors towards grey (Projection)
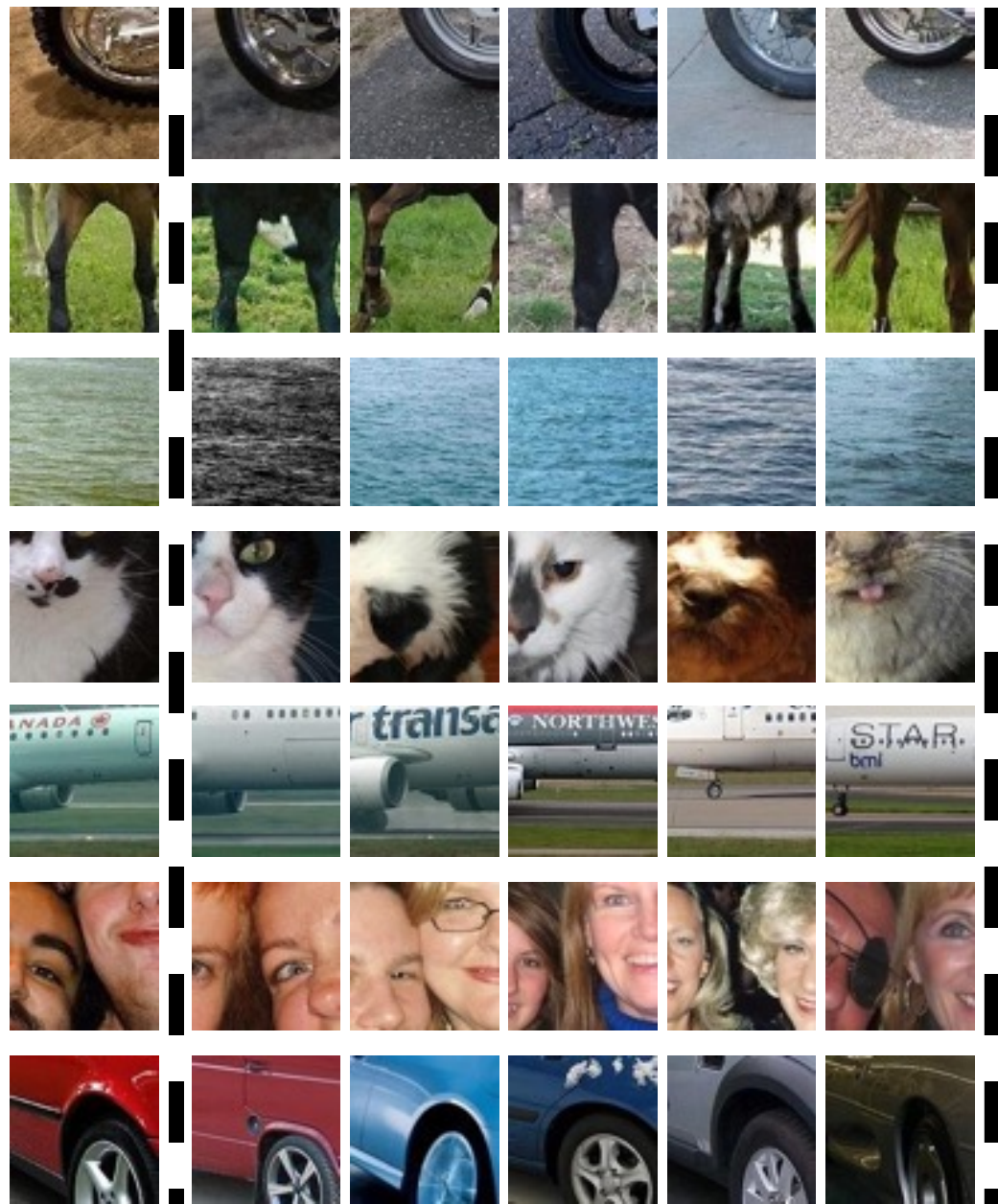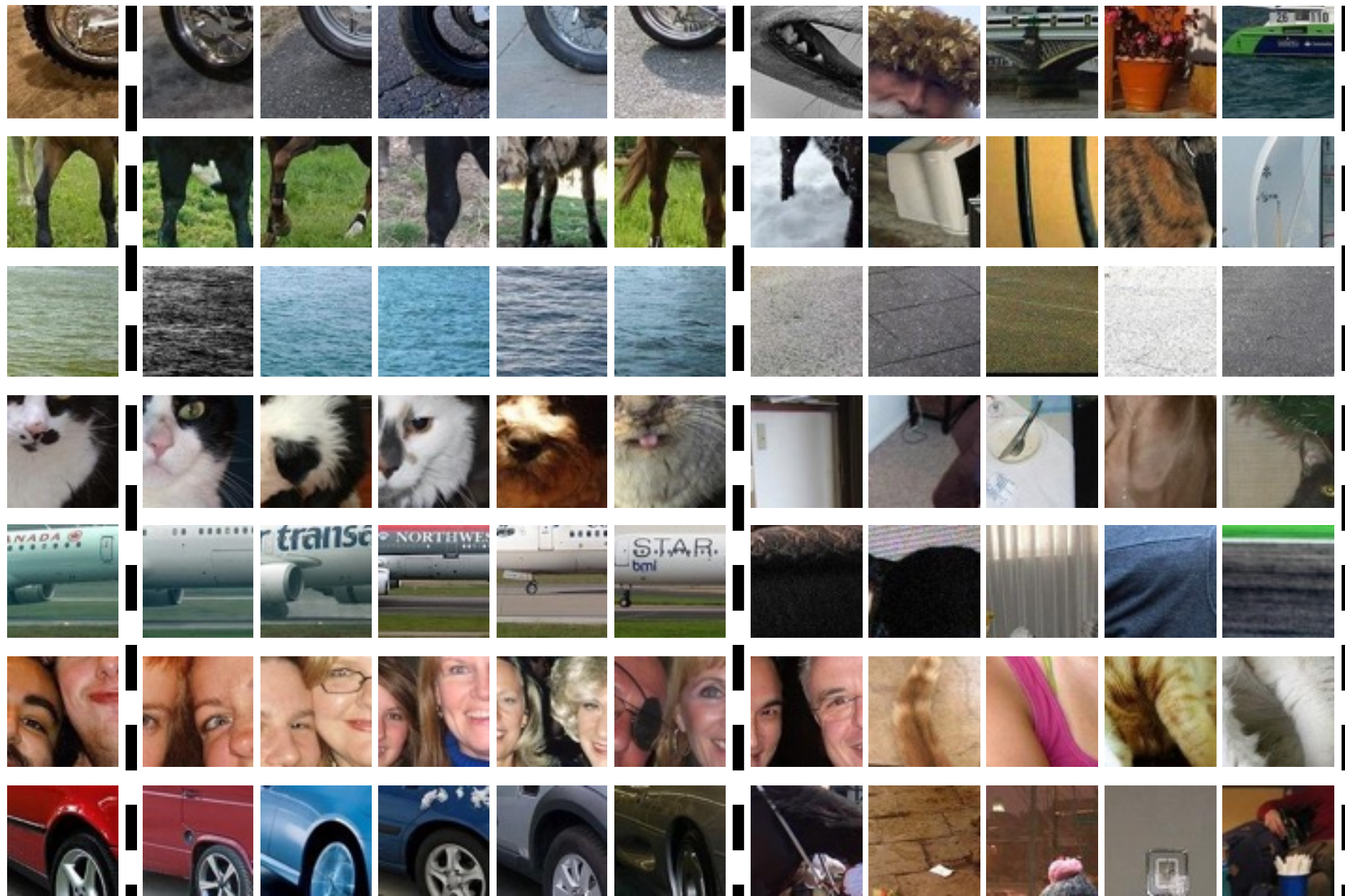- Drop 2 out of three channels during training

# Another trivial shortcut

- Chromatic Aberration
- Shift colors towards grey (Projection)
- Drop 2 out of three channels during training

# Another trivial shortcut

- Chromatic Aberration
- Shift colors towards grey (Projection)
- Drop 2 out of three channels during training

# Another trivial shortcut

- Chromatic Aberration
- Shift colors towards grey (Projection)
- Drop 2 out of three channels during training

# What is learned?

Input          Ours

# What is learned?

Input   Ours   Random Initialization

# What is learned?

Input          Ours                    Random Initialization          ImageNet AlexNet



Slide from Carl Doersch

# Still don't capture everything

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

Slide from Carl Doersch

# Pre-Training for R-CNN



warped region

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

CNN

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

# Pre-Training for R-CNN



warped region

aeroplane? no.

person? yes.

tvmonitor? no.

CNN

**1**. Input image

**2**. Extract region proposals (~2k)

**3**. Compute CNN features

**4**. Classify regions

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

Slide from Carl Doersch

# Details

- Use stack from patch context predictor before pool5

- Resize convolution layers to work on 227x227 instead of 96x96

- Use FC7 as the final representation

| fc8 (21) |
|---|
| fc7 (4096) |
| pool6 (3x3,1024,2) |
| conv6b (1x1,1024,1) |
| conv6 (3x3,4096,1) |
| pool5 |

Image (227x227)

# Architecture



Slide from Carl Doersch

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

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

Average Precision

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

Average Precision

40.7%

No Pretraining

Slide from Carl Doersch

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



Average Precision

46.3%

40.7%

No Pretraining

Ours (No Labels)

Slide from Carl Doersch

VOC 2007 Performance
(pretraining for R-CNN)

Average Precision

40.7% — No Pretraining

46.3% — Ours (No Labels)

54.2% — ImageNet Labels

Slide from Carl Doersch

# Unsupervised Object Discovery?

# Unsupervised Object Discovery

# Unsupervised Object Discovery

# Unsupervised Object Discovery

# Unsupervised Object Discovery

# Unsupervised Object Discovery

# Algorithm

# Algorithm



x100

Slide from Carl Doersch

# Algorithm



x100

# Algorithm



x100

15/100
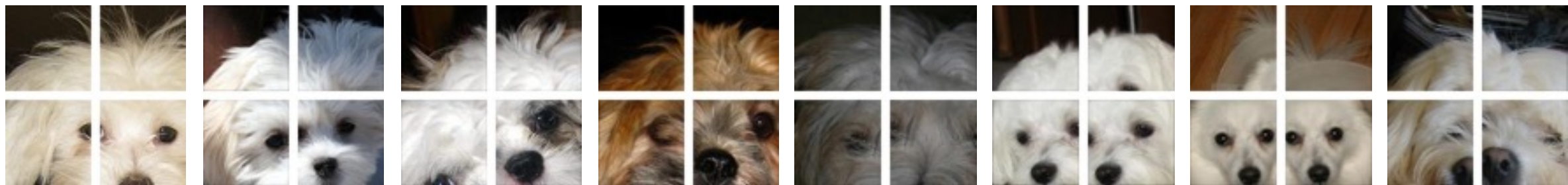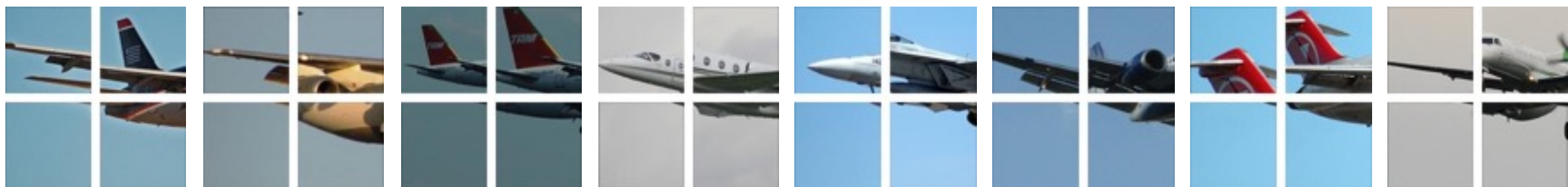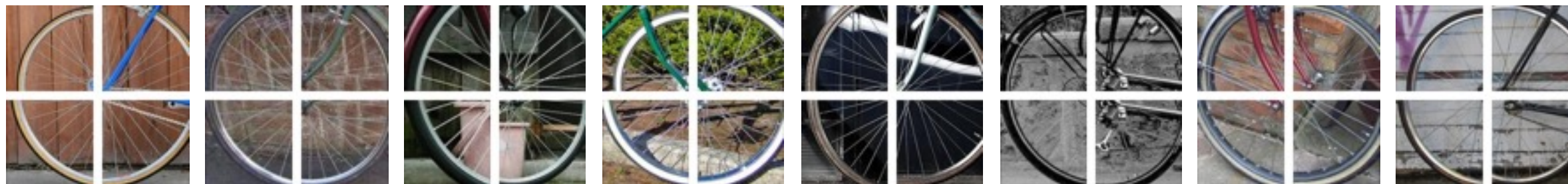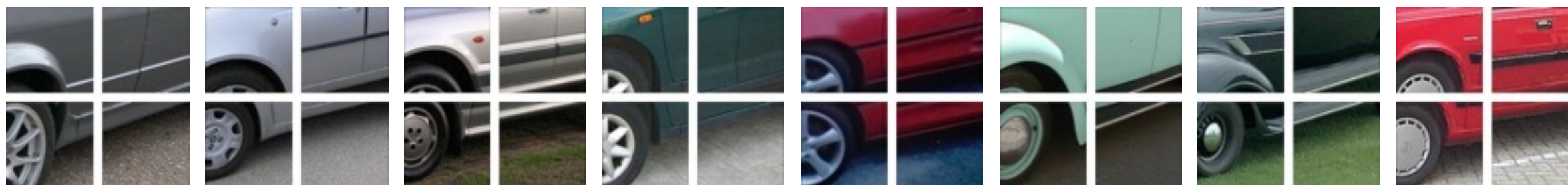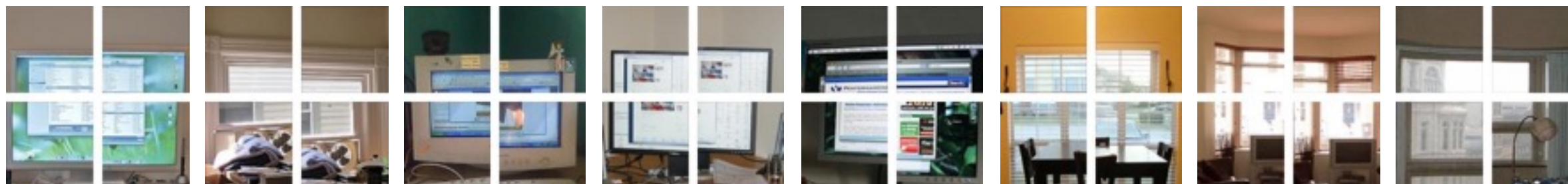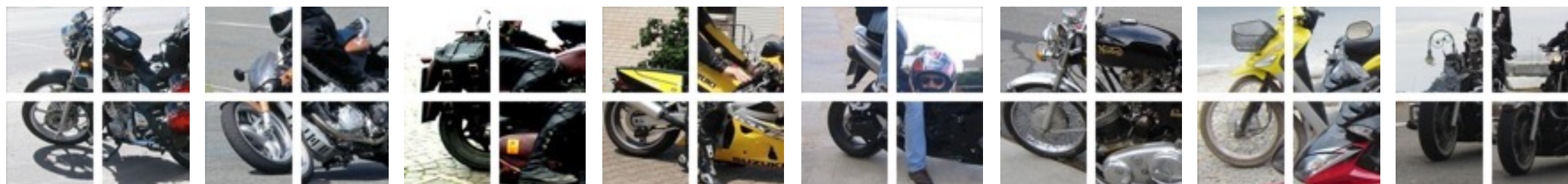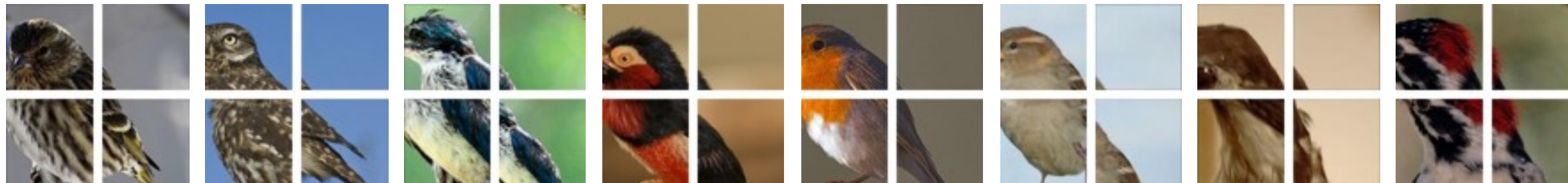
84/100
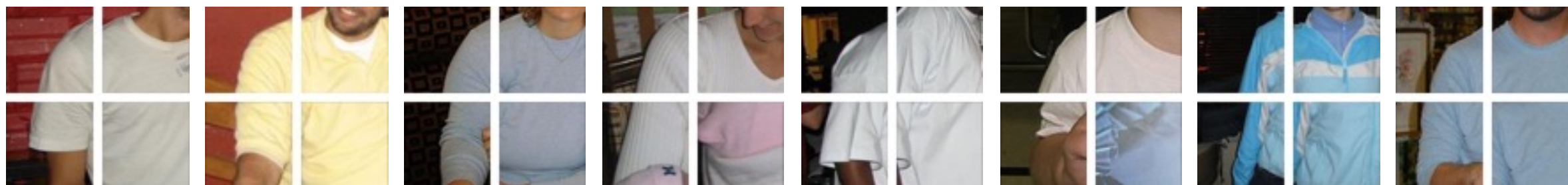
7/100

1

3

4

6

11

16

Slide from Carl Doersch

25
32
43
61
63
93

Slide from Carl Doersch

119

152

169

182

395

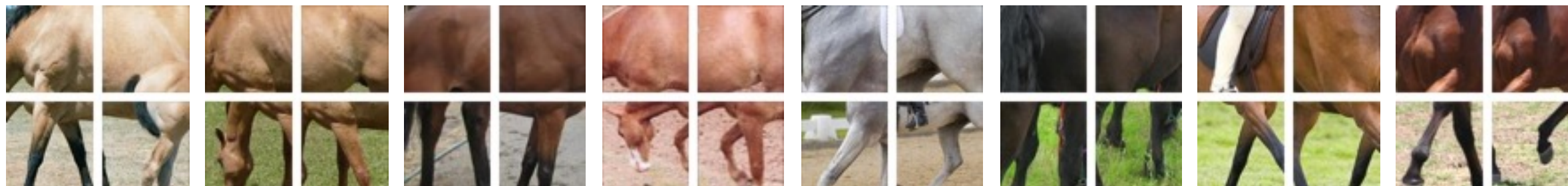457

Slide from Carl Doersch

# Purity vs Coverage



**Purity-Coverage for Proposed Objects**

Legend:
- Visual Words .63 (.37)
- Russel et al. .66 (.38)
- HOG Kmeans .70 (.40)
- Singh et al. .83 (.47)
- Doersch et al. .83 (.48)
- Our Approach .87 (.48)

# Pretext Task

- Performance on Pascal VOC is 38.4% (Chance is 12.5%)

- On ImageNet accuracy is 39.5% on training set, and 40.5% on validation

- On GT box patches - similar performance. 39.2% overall with 45.6% on cars

# Questions?