

Grid LSTM

Kalchbrenner et al. (Google DeepMind)

Joost Bastings



INSTITUTE FOR LOGIC, LANGUAGE AND COMPUTATION

March 10, 2016

Outline

1. Introduction
2. LSTM
3. Grid LSTM
4. Experiments
5. Conclusion

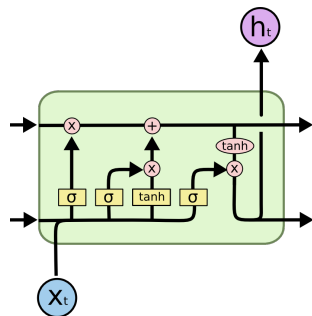
Outline

1. Introduction
2. LSTM
3. Grid LSTM
4. Experiments
5. Conclusion

LSTM

Long Short-Term Memory (LSTM) networks have *gates* that control access to memory cells

(Hochreiter and Schmidhuber, 1997)



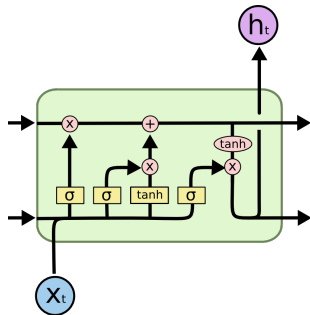
©Christopher Olah

LSTM

Long Short-Term Memory (LSTM) networks have *gates* that control access to memory cells

(Hochreiter and Schmidhuber, 1997)

- ▶ Preserve signals, propagate errors for *much longer*



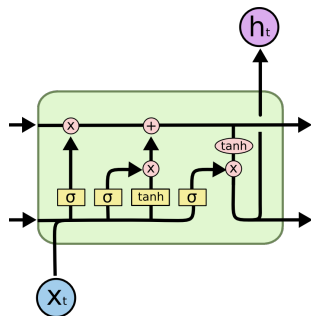
©Christopher Olah

LSTM

Long Short-Term Memory (LSTM) networks have *gates* that control access to memory cells

(Hochreiter and Schmidhuber, 1997)

- ▶ Preserve signals, propagate errors for *much longer*
- ▶ Gates can learn to attend to specific parts of the input signals (and ignore others)



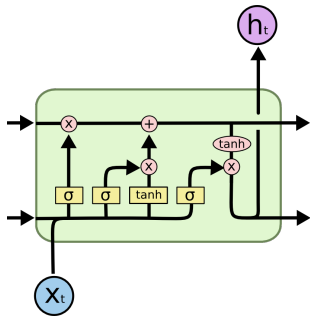
©Christopher Olah

LSTM

Long Short-Term Memory (LSTM) networks have *gates* that control access to memory cells

(Hochreiter and Schmidhuber, 1997)

- ▶ Preserve signals, propagate errors for *much longer*
- ▶ Gates can learn to attend to specific parts of the input signals (and ignore others)



©Christopher Olah

These properties make LSTMs good at speech recognition, hand-writing recognition, machine translation, etc.

Motivation

- ▶ Computer vision success: **deep networks are key** to finding complex patterns

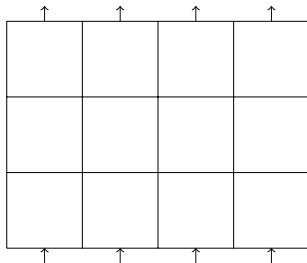
Motivation

- ▶ Computer vision success: **deep networks are key** to finding complex patterns
- ▶ However, deep networks also suffer from the **vanishing gradient** problem!

Motivation

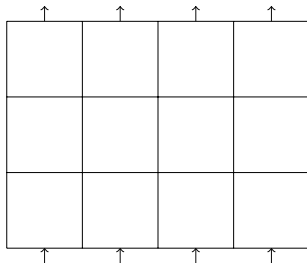
- ▶ Computer vision success: **deep networks are key** to finding complex patterns
- ▶ However, deep networks also suffer from the **vanishing gradient** problem!
- ▶ This is the motivation to generalise the advantages of LSTMs to deep computation

Idea: Grid LSTM



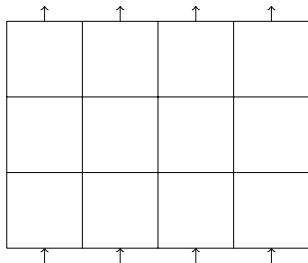
- ▶ A Grid LSTM (Kalchbrenner et al., 2015) is a network arranged in a grid of 1 or more dimensions

Idea: Grid LSTM



- ▶ A Grid LSTM (Kalchbrenner et al., 2015) is a network arranged in a grid of 1 or more dimensions
- ▶ LSTM cells in 'any or all' dimensions of the grid

Idea: Grid LSTM



- ▶ A Grid LSTM (Kalchbrenner et al., 2015) is a network arranged in a grid of 1 or more dimensions
- ▶ LSTM cells in 'any or all' dimensions of the grid
- ▶ Short-hand: N -dimensional Grid LSTM = N -LSTM

Outline

1. Introduction
2. LSTM
3. Grid LSTM
4. Experiments
5. Conclusion

LSTM

- ▶ An LSTM processes input and target pairs

$$(x_1, y_1), \dots, (x_m, y_m)$$

LSTM

- ▶ An LSTM processes input and target pairs

$$(x_1, y_1), \dots, (x_m, y_m)$$

- ▶ Past inputs x_1, \dots, x_{i-1} determine the state of the network:

hidden $\mathbf{h} \in \mathbb{R}^d$

memory $\mathbf{m} \in \mathbb{R}^d$

LSTM

- ▶ An LSTM processes input and target pairs

$$(x_1, y_1), \dots, (x_m, y_m)$$

- ▶ Past inputs x_1, \dots, x_{i-1} determine the state of the network:

$$\text{hidden } \mathbf{h} \in \mathbb{R}^d$$

$$\text{memory } \mathbf{m} \in \mathbb{R}^d$$

- ▶ Let $\mathbf{H} = \begin{bmatrix} lx_i \\ \mathbf{h} \end{bmatrix}$, where l is a projection matrix transforming x_i

LSTM

At each step, calculate:

1. Gates:

$$\mathbf{g}^u = \sigma(\mathbf{W}^u \mathbf{H}) \quad \text{update}$$

$$\mathbf{g}^f = \sigma(\mathbf{W}^f \mathbf{H}) \quad \text{forget}$$

$$\mathbf{g}^o = \sigma(\mathbf{W}^o \mathbf{H}) \quad \text{output}$$



LSTM

At each step, calculate:

1. Gates:

$$\mathbf{g}^u = \sigma(\mathbf{W}^u \mathbf{H}) \quad \text{update}$$

$$\mathbf{g}^f = \sigma(\mathbf{W}^f \mathbf{H}) \quad \text{forget}$$

$$\mathbf{g}^o = \sigma(\mathbf{W}^o \mathbf{H}) \quad \text{output}$$

$$\mathbf{g}^c = \tanh(\mathbf{W}^c \mathbf{H}) \quad \text{content}$$



LSTM

At each step, calculate:

1. Gates:

$$\begin{aligned} \mathbf{g}^u &= \sigma(\mathbf{W}^u \mathbf{H}) && \text{update} \\ \mathbf{g}^f &= \sigma(\mathbf{W}^f \mathbf{H}) && \text{forget} \\ \mathbf{g}^o &= \sigma(\mathbf{W}^o \mathbf{H}) && \text{output} \\ \mathbf{g}^c &= \tanh(\mathbf{W}^c \mathbf{H}) && \text{content} \end{aligned}$$



2. New memory:

$$\mathbf{m}' = \mathbf{g}^f \odot \mathbf{m} + \mathbf{g}^u \odot \mathbf{g}^c$$

LSTM

At each step, calculate:

1. Gates:

$$\begin{aligned} \mathbf{g}^u &= \sigma(\mathbf{W}^u \mathbf{H}) && \text{update} \\ \mathbf{g}^f &= \sigma(\mathbf{W}^f \mathbf{H}) && \text{forget} \\ \mathbf{g}^o &= \sigma(\mathbf{W}^o \mathbf{H}) && \text{output} \\ \mathbf{g}^c &= \tanh(\mathbf{W}^c \mathbf{H}) && \text{content} \end{aligned}$$



2. New memory:

$$\mathbf{m}' = \mathbf{g}^f \odot \mathbf{m} + \mathbf{g}^u \odot \mathbf{g}^c$$

LSTM

At each step, calculate:

1. Gates:

$$\begin{aligned} \mathbf{g}^u &= \sigma(\mathbf{W}^u \mathbf{H}) && \text{update} \\ \mathbf{g}^f &= \sigma(\mathbf{W}^f \mathbf{H}) && \text{forget} \\ \mathbf{g}^o &= \sigma(\mathbf{W}^o \mathbf{H}) && \text{output} \\ \mathbf{g}^c &= \tanh(\mathbf{W}^c \mathbf{H}) && \text{content} \end{aligned}$$



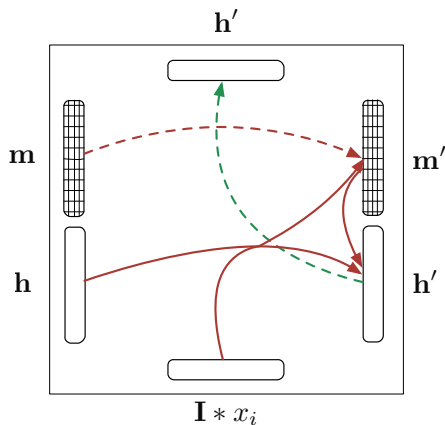
2. New memory:

$$\mathbf{m}' = \mathbf{g}^f \odot \mathbf{m} + \mathbf{g}^u \odot \mathbf{g}^c$$

3. New state:

$$\mathbf{h}' = \tanh(\mathbf{g}^o \odot \mathbf{m}')$$

LSTM



Standard LSTM block

Outline

1. Introduction
2. LSTM
3. Grid LSTM
4. Experiments
5. Conclusion

Grid LSTM

An N -LSTM block receives as input:

- ▶ N hidden vectors $\mathbf{h}_1, \dots, \mathbf{h}_N$
- ▶ N memory vectors $\mathbf{m}_1, \dots, \mathbf{m}_N$

Grid LSTM

An N -LSTM block receives as input:

- ▶ N hidden vectors $\mathbf{h}_1, \dots, \mathbf{h}_N$
- ▶ N memory vectors $\mathbf{m}_1, \dots, \mathbf{m}_N$

We concatenate all hidden vectors into a **shared input vector**:

$$\mathbf{H} = \begin{bmatrix} \mathbf{h}_1 \\ \vdots \\ \mathbf{h}_N \end{bmatrix}$$

Grid LSTM

An N -LSTM block receives as input:

- ▶ N hidden vectors $\mathbf{h}_1, \dots, \mathbf{h}_N$
- ▶ N memory vectors $\mathbf{m}_1, \dots, \mathbf{m}_N$

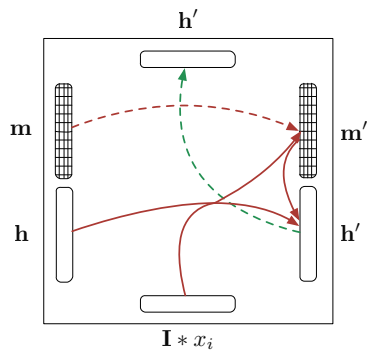
We concatenate all hidden vectors into a **shared input vector**:

$$\mathbf{H} = \begin{bmatrix} \mathbf{h}_1 \\ \vdots \\ \mathbf{h}_N \end{bmatrix}$$

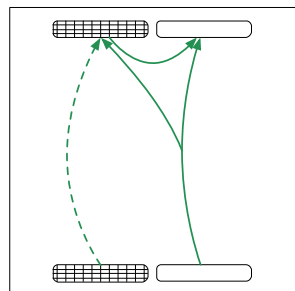
And then calculate N transforms:

$$\begin{aligned} (\mathbf{h}'_1, \mathbf{m}'_1) &= \text{LSTM}(\mathbf{H}, \mathbf{m}_1, \mathbf{W}_1) \\ &\vdots \\ (\mathbf{h}'_N, \mathbf{m}'_N) &= \text{LSTM}(\mathbf{H}, \mathbf{m}_N, \mathbf{W}_N) \end{aligned}$$

1D Grid LSTM

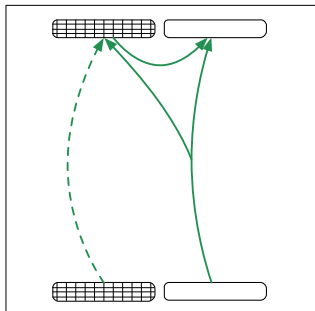


Standard LSTM block

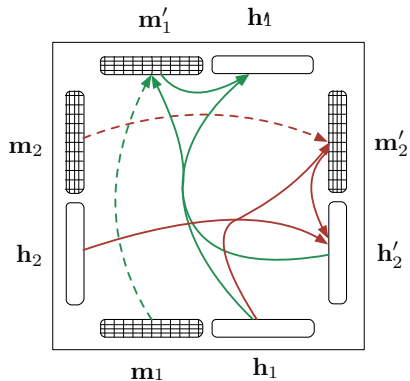


1d Grid LSTM Block

1D and 2D

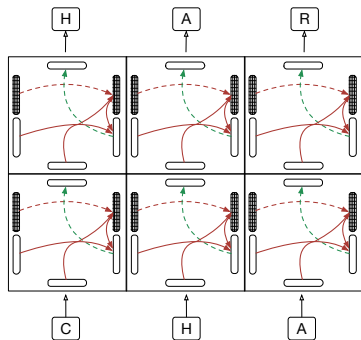


1d Grid LSTM Block

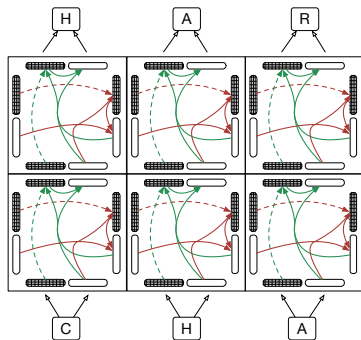


2d Grid LSTM block

Stacked LSTM vs 2-LSTM

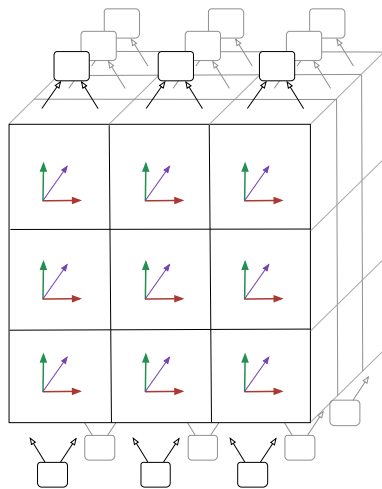


Stacked LSTM



2d Grid LSTM

3D



3d Grid LSTM

Grid LSTM – Notes

- ▶ Input is projected along the edge(s), see previous slide:
character 'C' initializes \mathbf{h}_1 and \mathbf{m}_1

Grid LSTM – Notes

- ▶ Input is projected along the edge(s), see previous slide:
character 'C' initializes \mathbf{h}_1 and \mathbf{m}_1
- ▶ Predictions based on both the state *and* memory of edge cells

Grid LSTM – Notes

- ▶ Input is projected along the edge(s), see previous slide: character 'C' initializes \mathbf{h}_1 and \mathbf{m}_1
- ▶ Predictions based on both the state *and* memory of edge cells
- ▶ It is possible to **share weights** along any dimension

Grid LSTM – Notes

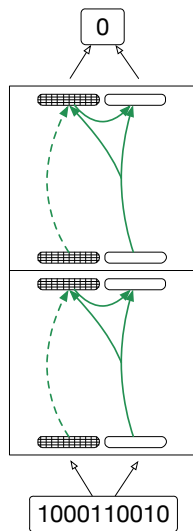
- ▶ Input is projected along the edge(s), see previous slide: character 'C' initializes \mathbf{h}_1 and \mathbf{m}_1
- ▶ Predictions based on both the state *and* memory of edge cells
- ▶ It is possible to **share weights** along any dimension
- ▶ If weights are shared along **all dimensions**: *Tied N-LSTM*

Outline

1. Introduction
2. LSTM
3. Grid LSTM
4. Experiments
5. Conclusion

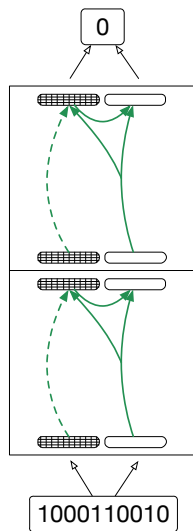
Parity

- ▶ Given k input bits, output 0 iff sum is even, else 1



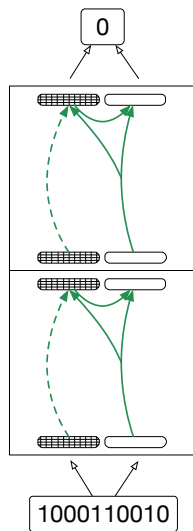
Parity

- ▶ Given k input bits, output 0 iff sum is even, else 1
- ▶ Parity is really hard, because changing one input bit changes the target

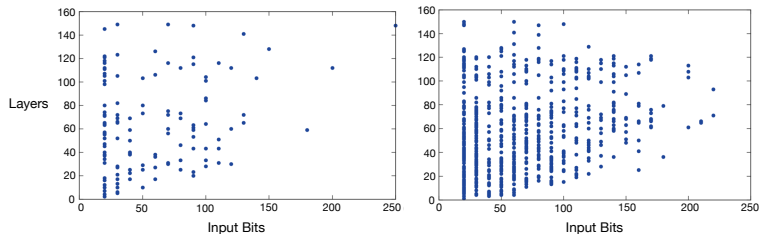


Parity

- ▶ Given k input bits, output 0 iff sum is even, else 1
- ▶ Parity is really hard, because changing one input bit changes the target
- ▶ All input at the same time (why?)

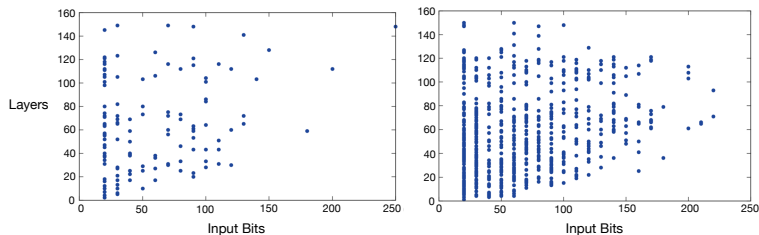


Parity



Left: $\#hidden = 500$, Right: $\#hidden = 1500$
Dot: 100% accuracy on k -bit input

Parity



Left: #hidden = 500, Right: #hidden = 1500
 Dot: 100% accuracy on k -bit input

	Layers	Hidden	k
Tied tanh FFN	5	1500	30
Tied ReLU FFN	4	1500	30
Tied 1-LSTM	72	1500	220
Tied 1-LSTM	148	500	250

Addition

Task: sum two 15-digit integers

```
— 1 2 3 — 8 9 9 — — — — —  
— — — — — — — — 1 0 2 2 —
```

Addition

Task: sum two 15-digit integers

```
— 1 2 3 — 8 9 9 — — — — —  
— — — — — — — — 1 0 2 2 —
```

	Layers	Samples	Accuracy
Stacked LSTM	1	5M	51%
Untied 2-LSTM	5	5M	67%
Tied 2-LSTM	18	0.55M	>99%

Addition

Task: sum two 15-digit integers

```

- 1 2 3 - 8 9 9 - - - - -
- - - - - - - - 1 0 2 2 -
  
```

	Layers	Samples	Accuracy
Stacked LSTM	1	5M	51%
Untied 2-LSTM	5	5M	67%
Tied 2-LSTM	18	0.55M	>99%

- ▶ Trained up to 5M samples or until accuracy 100%
- ▶ *Tied* better because of the repetitive nature of the task
- ▶ Grid LSTM has advantage by tackling vanishing gradient

Memorization

Task: memorize random sequence of 20 symbols

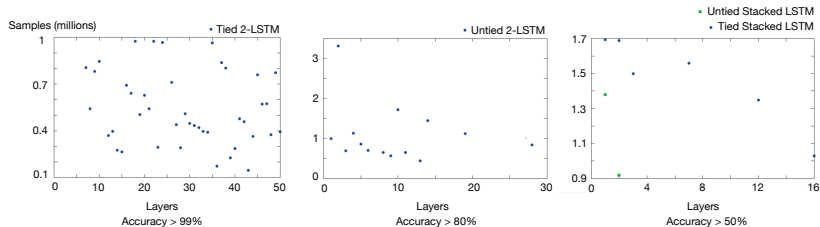
— a b c — — — —
— — — — a b c —

Memorization

Task: memorize random sequence of 20 symbols

```

— a b c — — — —
— — — — a b c —
  
```

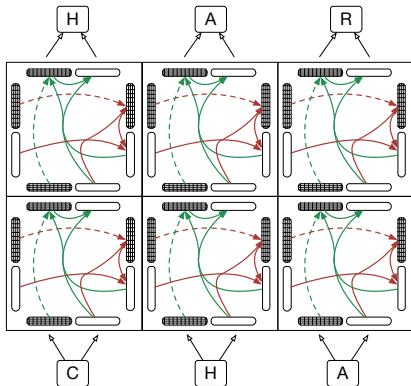


- ▶ All networks have 100 hidden units
- ▶ Vertical axis: #samples to reach threshold

Character-level LM

Task: predict next character in corpus

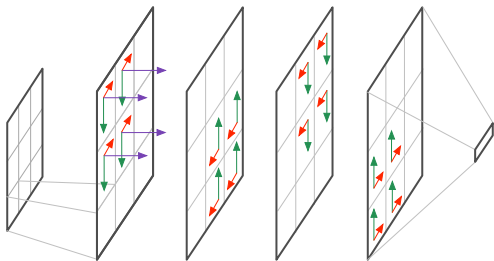
- ▶ Hutter challenge (Wikipedia data set, 100M characters)
- ▶ Sample sequences of 10000 chars, backprop every 50 chars



Character-level LM

	BPC	Params	Alphabet	Test
Stacked LSTM (Graves)	1.67	27M	205	last 4MB
MRNN (Sutskever)	1.60	4.9M	86	last 10MB
GFRNN (Chung)	1.58	20M	205	last 5MB
Tied 2-LSTM	1.47	16.8M	205	last 5MB

MNIST Digits

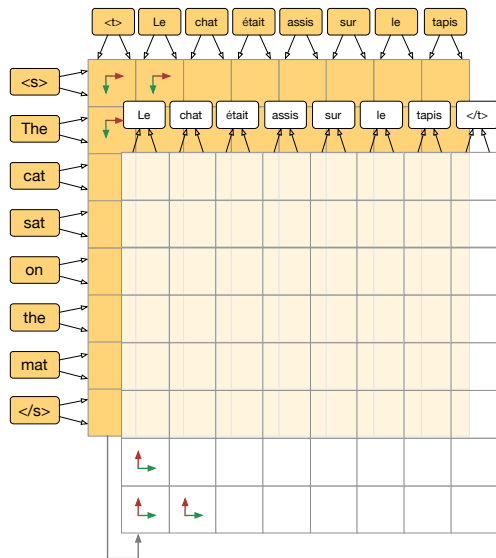


- ▶ A 3-LSTM processes non-overlapping patches of image pixels
- ▶ So, the input is a 2D grid of patches
- ▶ The 3rd dimension is the depth of the network
- ▶ Final ReLU layer + Softmax

MNIST Digits – Results

	Test Error (%)
Wan et al. (Wan et al., 2013)	0.28
Graham (Graham, 2014a)	0.31
Untied 3-LSTM	0.32
Ciresan et al. (Ciresan et al., 2012)	0.35
Untied 3-LSTM with ReLU	0.36
Mairal et al. (Mairal et al., 2014)	0.39
Lee et al. (Lee et al., 2015)	0.39
Simard et al. (Simard et al., 2003)	0.4
Graham (Graham, 2014b)	0.44
Goodfellow et al. (Goodfellow et al., 2013)	0.45
Visin et al. (Visin et al., 2015)	0.45
Lin et al. (Lin et al., 2013)	0.47

Machine Translation



Machine Translation

- ▶ We view translation as a 2-dimensional mapping

Machine Translation

- ▶ We view translation as a 2-dimensional mapping
- ▶ One dimension processes the **source** sentence, another dimension produces the **target** sentence

Machine Translation

- ▶ We view translation as a 2-dimensional mapping
- ▶ One dimension processes the **source** sentence, another dimension produces the **target** sentence
- ▶ The network **repeatedly re-encodes** the source sentence based on the part of the target sentence generated so far

Machine Translation

- ▶ We view translation as a 2-dimensional mapping
- ▶ One dimension processes the **source** sentence, another dimension produces the **target** sentence
- ▶ The network **repeatedly re-encodes** the source sentence based on the part of the target sentence generated so far
- ▶ **Weights are shared** across source and target dimensions

Machine Translation

- ▶ We view translation as a 2-dimensional mapping
- ▶ One dimension processes the **source** sentence, another dimension produces the **target** sentence
- ▶ The network **repeatedly re-encodes** the source sentence based on the part of the target sentence generated so far
- ▶ **Weights are shared** across source and target dimensions
- ▶ Regular identity connections along the 3rd dimension

Evaluation

Evaluation on IWSLT BTEC Chinese-to-English

- ▶ 44016 sentence pairs (train), 1006 (dev), 503 (test)
- ▶ Target sentences on average 12 words long
- ▶ 15 reference translations

Evaluation

Evaluation on IWSLT BTEC Chinese-to-English

- ▶ 44016 sentence pairs (train), 1006 (dev), 503 (test)
- ▶ Target sentences on average 12 words long
- ▶ 15 reference translations

	Valid-1	Test-1	Valid-15	Test-15
DGLSTM-Att.	-	34.5	-	-
CDEC	30.1	41	50.1	58.9
3-LSTM	30.3	42.4	51.8	60.2

Outline

1. Introduction
2. LSTM
3. Grid LSTM
4. Experiments
5. Conclusion

Conclusion

- ▶ Introduction of Grid LSTM
- ▶ Cells have shown advantages in parity, addition, memorization tasks
- ▶ Applications in character prediction, MNIST, and machine translation

Bibliography

Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.

Kalchbrenner, N., Danihelka, I., and Graves, A. (2015). Grid long short-term memory. *CoRR*, abs/1507.01526.

Image Captioning & Attention

Deep Learning

Hendrik Heuer

University of Amsterdam



About me

Now

Hendrik Heuer
PhD student
QUvA Lab



prof. dr. Arnold Smeulders
dr. Christof Monz

Statistical Machine Translation
& Computer Vision

M.Sc.



Aalto University



European Institute of
Innovation & Technology

Show, Attend and Tell: Neural Image Caption Generation with Visual Attention

Kelvin Xu
Jimmy Lei Ba
Ryan Kiros
Kyunghyun Cho
Aaron Courville
Ruslan Salakhutdinov
Richard S. Zemel
Yoshua Bengio

KELVIN.XU@UMONTREAL.CA
JIMMY@PSI.UTORONTO.CA
RKIROS@CS.TORONTO.EDU
KYUNGHYUN.CHO@UMONTREAL.CA
AARON.COURVILLE@UMONTREAL.CA
RSALAKHU@CS.TORONTO.EDU
ZEMEL@CS.TORONTO.EDU
FIND-ME@THE.WEB

Abstract

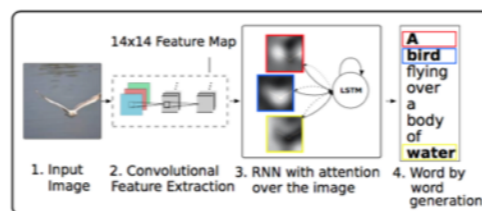
Inspired by recent work in machine translation and object detection, we introduce an attention based model that automatically learns to describe the content of images. We describe how we can train this model in a deterministic manner using standard backpropagation techniques and stochastically by maximizing a variational lower bound. We also show through visualization how the model is able to automatically learn to fix its gaze on salient objects while generating the corresponding words in the output sequence. We validate the use of attention with state-of-the-art performance on three benchmark datasets: Flickr8k, Flickr30k and MS COCO.

1. Introduction

Automatically generating captions of an image is a task very close to the heart of scene understanding — one of the primary goals of computer vision. Not only must caption generation models be powerful enough to solve the computer vision challenges of determining which objects are in an image, but they must also be capable of capturing and expressing their relationships in a natural language. For this reason, caption generation has long been viewed as a difficult problem. It is a very important challenge for machine learning algorithms, as it amounts to mimicking the remarkable human ability to compress huge amounts of salient visual information into descriptive language.

Despite the challenging nature of this task, there has been a recent surge of research interest in attacking the image caption generation problem. Aided by advances in training neural networks (Krizhevsky et al., 2012) and large classification datasets (Russakovsky et al., 2014), recent work

Figure 1. Our model learns a words/image alignment. The visualized attentional maps (3) are explained in section 3.1 & 5.4



has significantly improved the quality of caption generation using a combination of convolutional neural networks (convnets) to obtain vectorial representation of images and recurrent neural networks to decode those representations into natural language sentences (see Sec. 2).

One of the most curious facets of the human visual system is the presence of attention (Rensink, 2000; Corbetta & Shulman, 2002). Rather than compress an entire image into a static representation, attention allows for salient features to dynamically come to the forefront as needed. This is especially important when there is a lot of clutter in an image. Using representations (such as those from the top layer of a convnet) that distill information in image down to the most salient objects is one effective solution that has been widely adopted in previous work. Unfortunately, this has one potential drawback of losing information which could be useful for richer, more descriptive captions. Using more low-level representation can help preserve this information. However working with these features necessitates a powerful mechanism to steer the model to information important to the task at hand.

In this paper, we describe approaches to caption generation that attempt to incorporate a form of attention with

DRAW: A Recurrent Neural Network For Image Generation

Karol Gregor
Ivo Danihelka
Alex Graves
Danilo Jimenez Rezende
Daan Wierstra
Google DeepMind

KAROLG@GOOGLE.COM
DANIELKA@GOOGLE.COM
GRAVES@GOOGLE.COM
DANILOR@GOOGLE.COM
WIERSTRA@GOOGLE.COM

Abstract

This paper introduces the *Deep Recurrent Attentive Writer* (DRAW) neural network architecture for image generation. DRAW networks combine a novel spatial attention mechanism that mimics the foveation of the human eye, with a sequential variational auto-encoding framework that allows for the iterative construction of complex images. The system substantially improves on the state of the art for generative models on MNIST, and, when trained on the Street View House Numbers dataset, it generates images that cannot be distinguished from real data with the naked eye.

1. Introduction

A person asked to draw, paint or otherwise recreate a visual scene will naturally do so in a sequential, iterative fashion, reassessing their handiwork after each modification. Rough outlines are gradually replaced by precise forms, lines are sharpened, darkened or erased, shapes are altered, and the final picture emerges. Most approaches to automatic image generation, however, aim to generate entire scenes at once. In the context of generative neural networks, this typically means that all the pixels are conditioned on a single latent distribution (Dayan et al., 1995; Hinton & Salakhutdinov, 2006; Larochelle & Murray, 2011). As well as precluding the possibility of iterative self-correction, the “one shot” approach is fundamentally difficult to scale to large images. The *Deep Recurrent Attentive Writer* (DRAW) architecture represents a shift towards a more natural form of image construction, in which parts of a scene are created independently from others, and approximate sketches are successively refined.

Proceedings of the 32nd International Conference on Machine Learning, Lille, France, 2015. JMLR: W&CP volume 37. Copyright 2015 by the author(s).

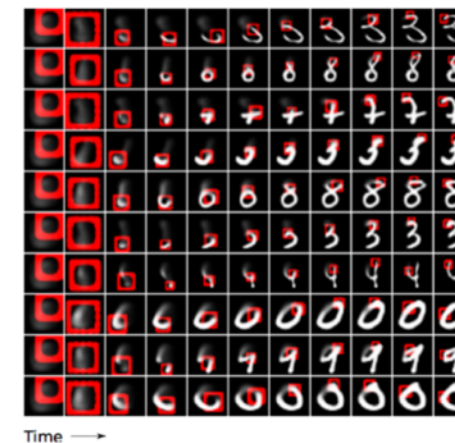


Figure 1. A trained DRAW network generating MNIST digits. Each row shows successive stages in the generation of a single digit. Note how the lines composing the digits appear to be “drawn” by the network. The red rectangle delimits the area attended to by the network at each time-step, with the focal precision indicated by the width of the rectangle border.

The core of the DRAW architecture is a pair of recurrent neural networks: an *encoder* network that compresses the real images presented during training, and a *decoder* that reconstitutes images after receiving codes. The combined system is trained end-to-end with stochastic gradient descent, where the loss function is a variational upper bound on the log-likelihood of the data. It therefore belongs to the family of *variational auto-encoders*, a recently emerged hybrid of deep learning and variational inference that has led to significant advances in generative modelling (Gregor et al., 2014; Kingma & Welling, 2014; Rezende et al., 2014; Mnih & Gregor, 2014; Salimans et al., 2014). Where DRAW differs from its siblings is that, rather than generat-

Describe images

Generate images

Show, Attend and Tell: Neural Image Caption Generation with Visual Attention

Kelvin Xu
Jimmy Lei Ba
Ryan Kiros
Kyunghyun Cho
Aaron Courville
Ruslan Salakhutdinov
Richard S. Zemel
Yoshua Bengio

KELVIN.XU@UMONTREAL.CA
JIMMY@PSI.UTORONTO.CA
RKIROS@CS.TORONTO.EDU
KYUNGHYUN.CHO@UMONTREAL.CA
AARON.COURVILLE@UMONTREAL.CA
RSALAKHU@CS.TORONTO.EDU
ZEMEL@CS.TORONTO.EDU
FIND-ME@THE.WEB

Abstract

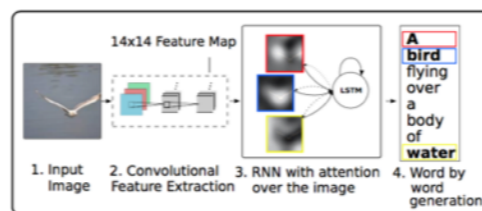
Inspired by recent work in machine translation and object detection, we introduce an attention based model that automatically learns to describe the content of images. We describe how we can train this model in a deterministic manner using standard backpropagation techniques and stochastically by maximizing a variational lower bound. We also show through visualization how the model is able to automatically learn to fix its gaze on salient objects while generating the corresponding words in the output sequence. We validate the use of attention with state-of-the-art performance on three benchmark datasets: Flickr8k, Flickr30k and MS COCO.

1. Introduction

Automatically generating captions of an image is a task very close to the heart of scene understanding — one of the primary goals of computer vision. Not only must caption generation models be powerful enough to solve the computer vision challenges of determining which objects are in an image, but they must also be capable of capturing and expressing their relationships in a natural language. For this reason, caption generation has long been viewed as a difficult problem. It is a very important challenge for machine learning algorithms, as it amounts to mimicking the remarkable human ability to compress huge amounts of salient visual information into descriptive language.

Despite the challenging nature of this task, there has been a recent surge of research interest in attacking the image caption generation problem. Aided by advances in training neural networks (Krizhevsky et al., 2012) and large classification datasets (Russakovsky et al., 2014), recent work

Figure 1. Our model learns a words/image alignment. The visualized attentional maps (3) are explained in section 3.1 & 5.4



has significantly improved the quality of caption generation using a combination of convolutional neural networks (convnets) to obtain vectorial representation of images and recurrent neural networks to decode those representations into natural language sentences (see Sec. 2).

One of the most curious facets of the human visual system is the presence of attention (Rensink, 2000; Corbetta & Shulman, 2002). Rather than compress an entire image into a static representation, attention allows for salient features to dynamically come to the forefront as needed. This is especially important when there is a lot of clutter in an image. Using representations (such as those from the top layer of a convnet) that distill information in image down to the most salient objects is one effective solution that has been widely adopted in previous work. Unfortunately, this has one potential drawback of losing information which could be useful for richer, more descriptive captions. Using more low-level representation can help preserve this information. However working with these features necessitates a powerful mechanism to steer the model to information important to the task at hand.

In this paper, we describe approaches to caption generation that attempt to incorporate a form of attention with

DRAW: A Recurrent Neural Network For Image Generation

Karol Gregor
Ivo Danihelka
Alex Graves
Danilo Jimenez Rezende
Daan Wierstra
Google DeepMind

KAROLG@GOOGLE.COM
DANIHELKA@GOOGLE.COM
GRAVES@GOOGLE.COM
DANILOR@GOOGLE.COM
WIERSTRA@GOOGLE.COM

Abstract

This paper introduces the *Deep Recurrent Attentive Writer* (DRAW) neural network architecture for image generation. DRAW networks combine a novel spatial attention mechanism that mimics the foveation of the human eye, with a sequential variational auto-encoding framework that allows for the iterative construction of complex images. The system substantially improves on the state of the art for generative models on MNIST, and, when trained on the Street View House Numbers dataset, it generates images that cannot be distinguished from real data with the naked eye.

1. Introduction

A person asked to draw, paint or otherwise recreate a visual scene will naturally do so in a sequential, iterative fashion, reassessing their handiwork after each modification. Rough outlines are gradually replaced by precise forms, lines are sharpened, darkened or erased, shapes are altered, and the final picture emerges. Most approaches to automatic image generation, however, aim to generate entire scenes at once. In the context of generative neural networks, this typically means that all the pixels are conditioned on a single latent distribution (Dayan et al., 1995; Hinton & Salakhutdinov, 2006; Larochelle & Murray, 2011). As well as precluding the possibility of iterative self-correction, the “one shot” approach is fundamentally difficult to scale to large images. The *Deep Recurrent Attentive Writer* (DRAW) architecture represents a shift towards a more natural form of image construction, in which parts of a scene are created independently from others, and approximate sketches are successively refined.

Proceedings of the 32nd International Conference on Machine Learning, Lille, France, 2015. JMLR: W&CP volume 37. Copyright 2015 by the author(s).

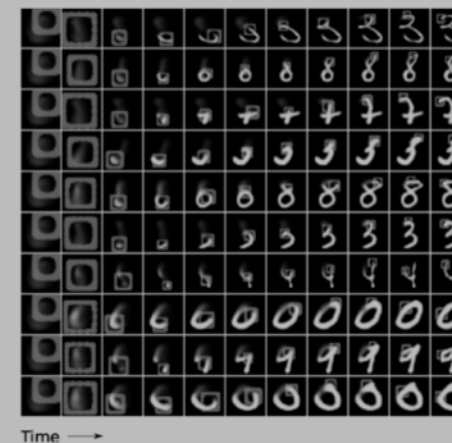


Figure 1. A trained DRAW network generating MNIST digits. Each row shows successive stages in the generation of a single digit. Note how the lines composing the digits appear to be “drawn” by the network. The red rectangle delimits the area attended to by the network at each time-step, with the focal precision indicated by the width of the rectangle border.

The core of the DRAW architecture is a pair of recurrent neural networks: an *encoder* network that compresses the real images presented during training, and a *decoder* that reconstitutes images after receiving codes. The combined system is trained end-to-end with stochastic gradient descent, where the loss function is a variational upper bound on the log-likelihood of the data. It therefore belongs to the family of *variational auto-encoders*, a recently emerged hybrid of deep learning and variational inference that has led to significant advances in generative modelling (Gregor et al., 2014; Kingma & Welling, 2014; Rezende et al., 2014; Mnih & Gregor, 2014; Salimans et al., 2014). Where DRAW differs from its siblings is that, rather than generat-

Describe images

Generate images

Grammatical description of images



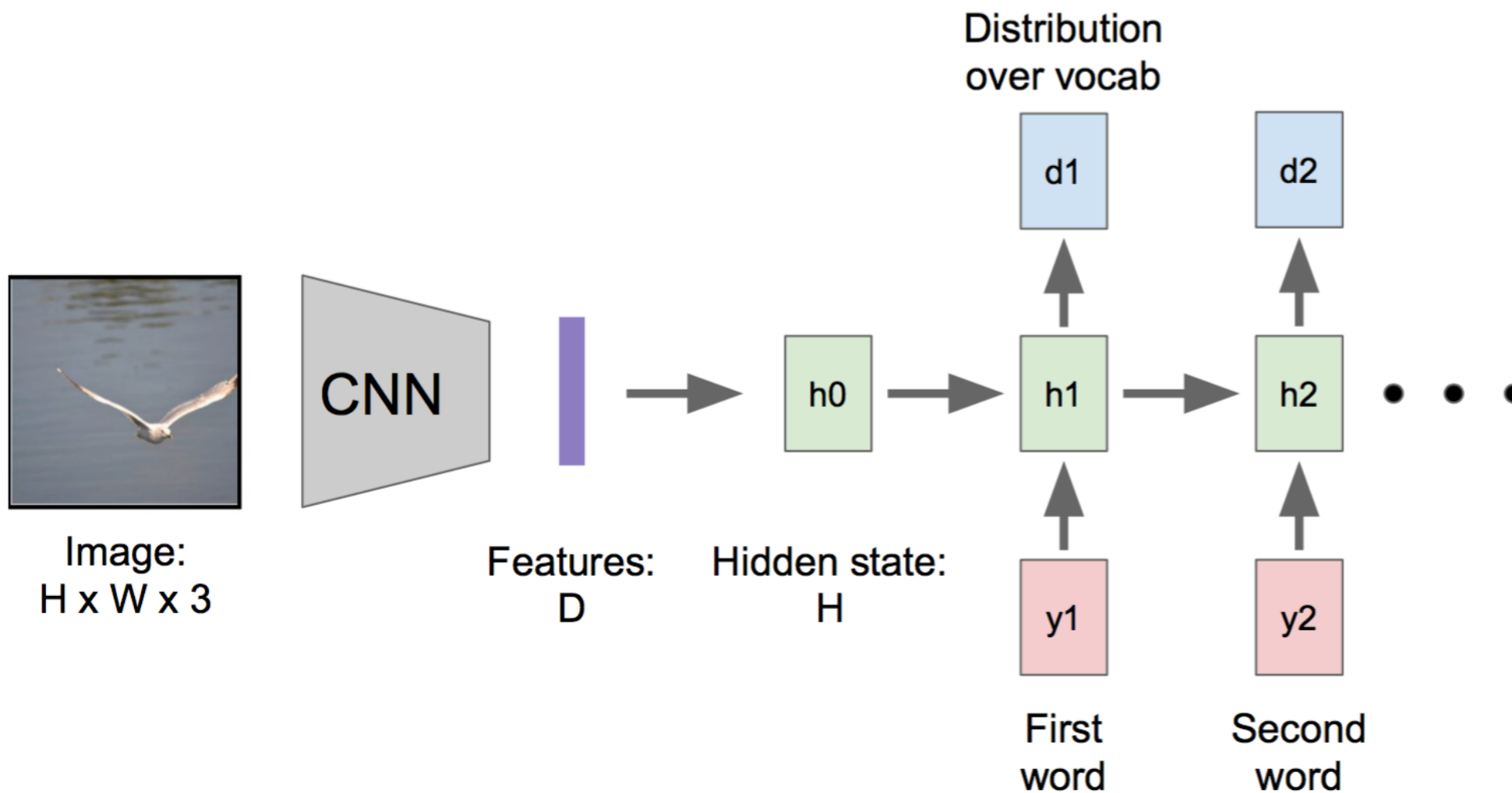
<https://vimeo.com/146492001>

Grammatical description of images



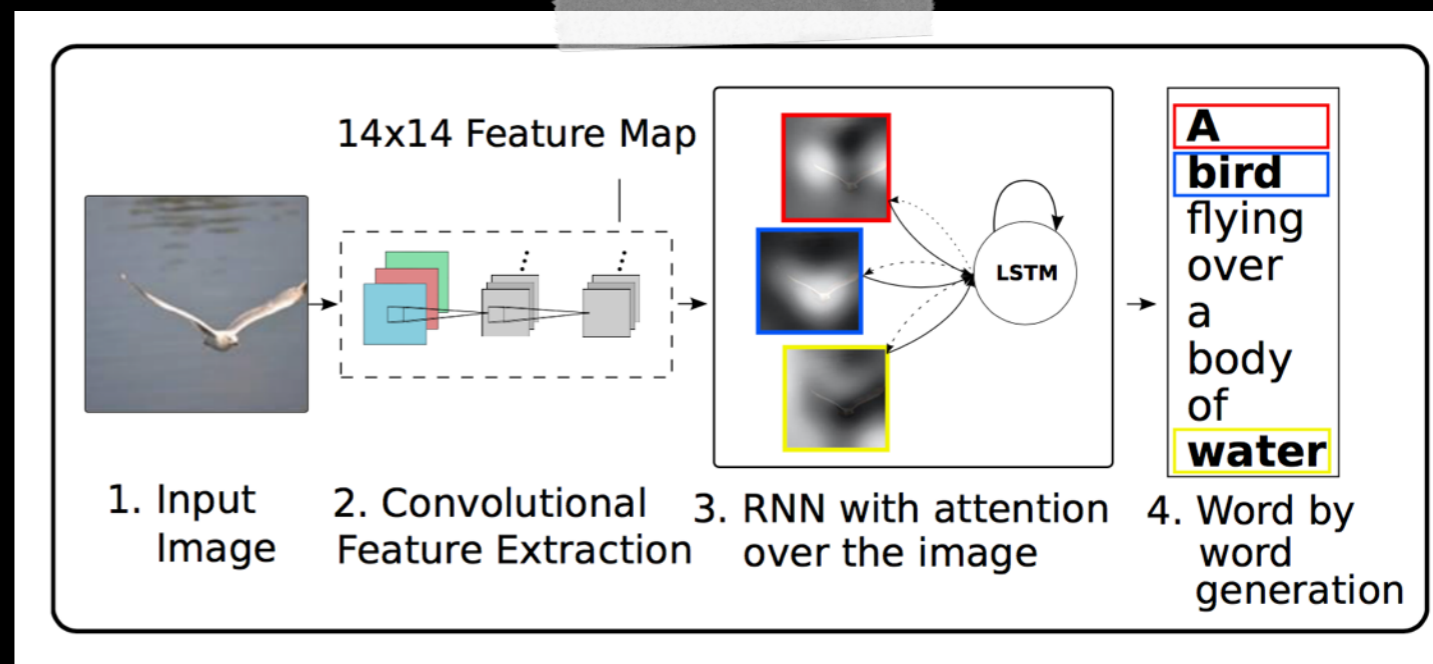
- a smiling old lady holds a pizza on a plate.
- a woman holding a plate with a pizza on it
- a woman carrying homemade pizza to the table.
- a woman holding a pizza on a red plate.
- a woman walking with a pan in her hands with a whole pizza on it.

RNN for Captioning





(b) A woman is throwing a frisbee in a park.





A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



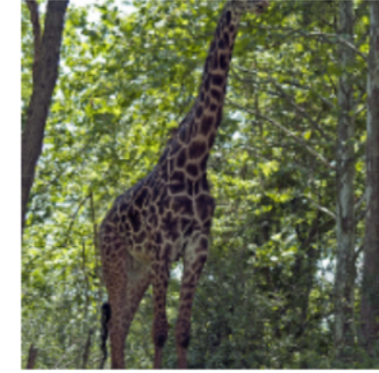
A stop sign is on a road with a mountain in the background.



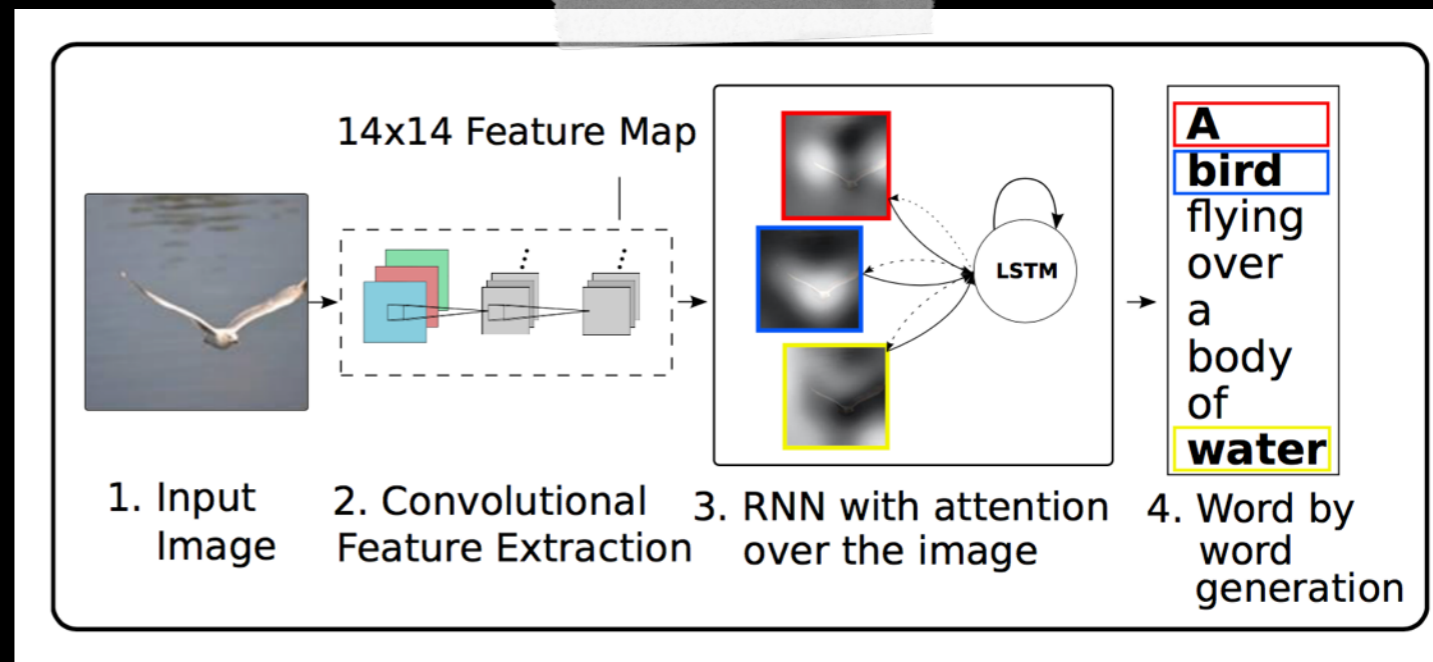
A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.





A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



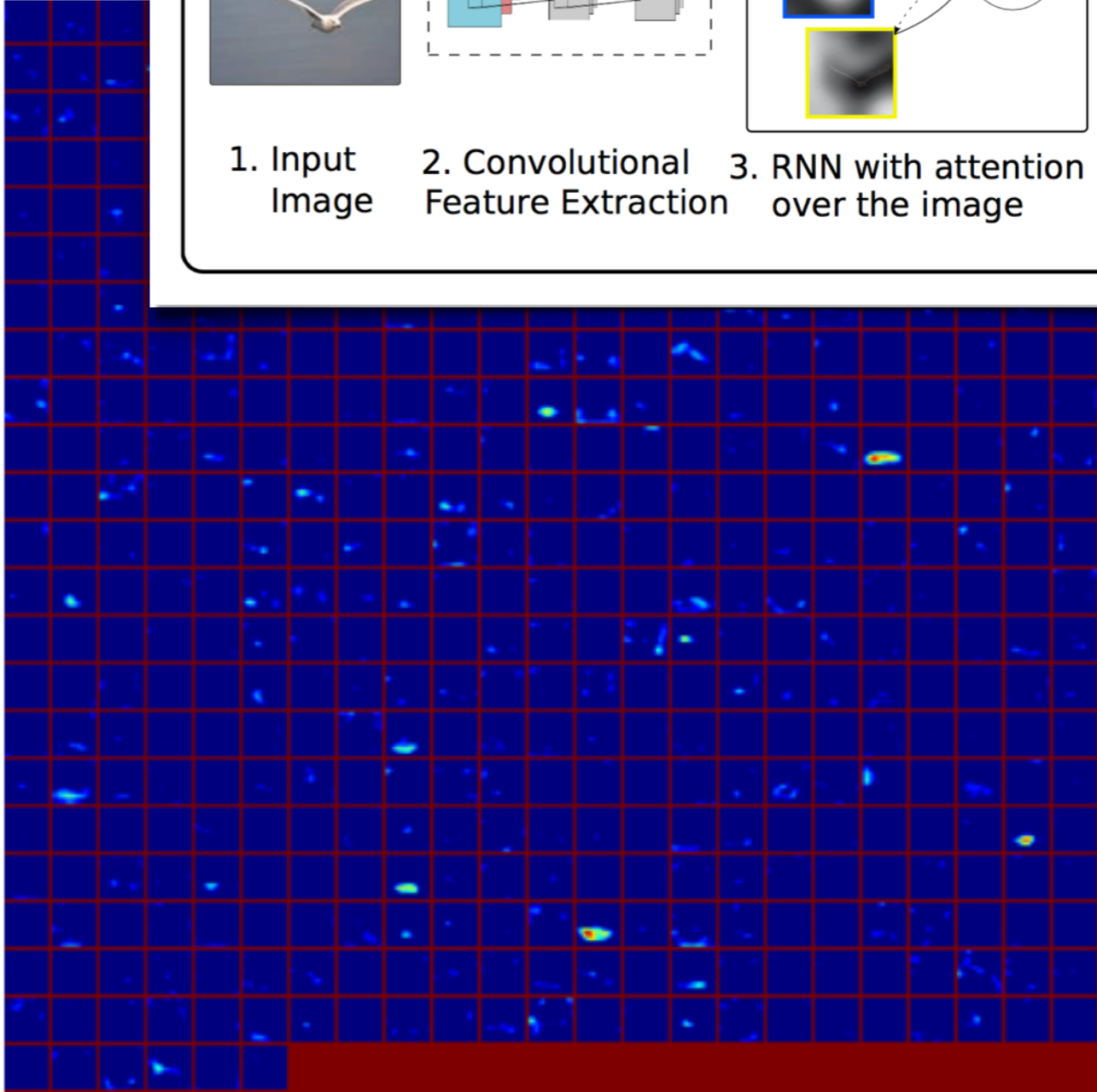
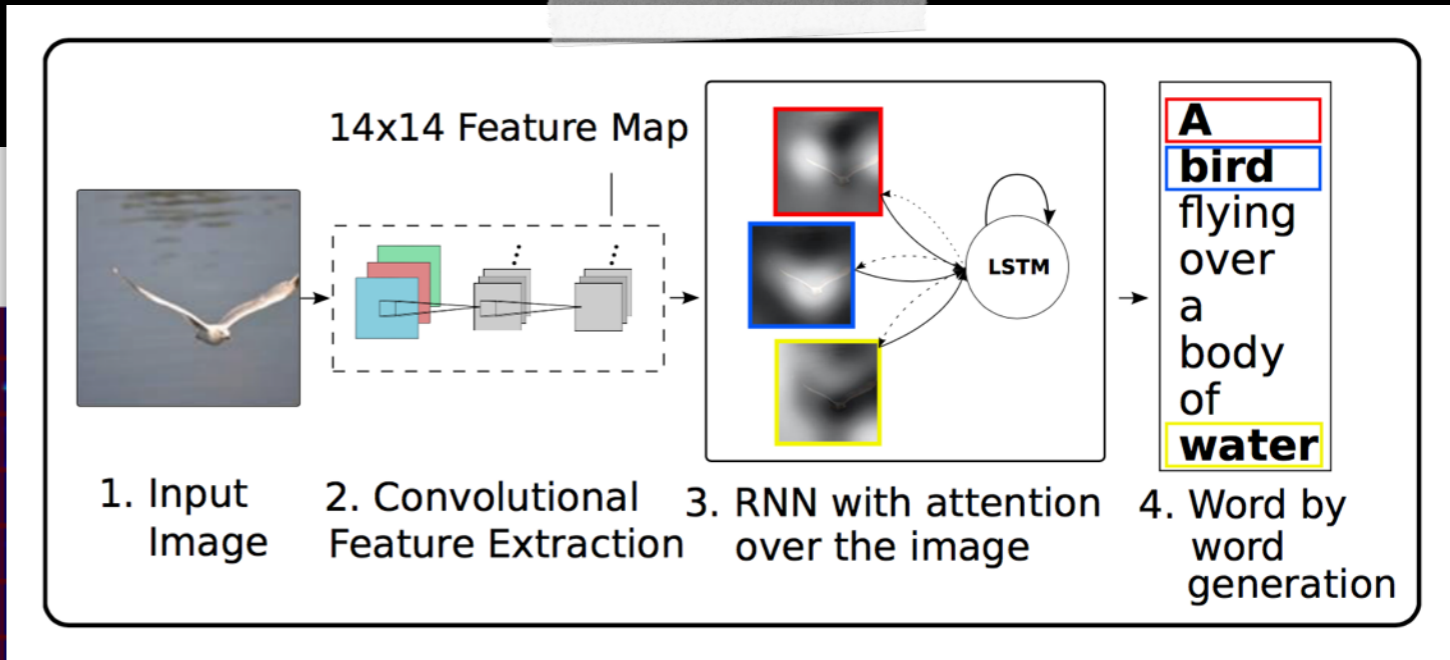
A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

Intuition: Since we usually see dogs at a certain position, we expect dogs at certain positions.

The model learns correlation structures in the input and starts putting attention weight where dogs can be expected (and actually exist in the training data).



512 filter, each 14x14 pixel



“soft” deterministic attention

summarize all locations, so that context vector z is

$$\mathbb{E}_{p(s_t|a)}[\hat{\mathbf{z}}_t] = \sum_{i=1}^L \alpha_{t,i} \mathbf{a}_i$$

you can take the derivative dz/dp

trainable by back-propagation

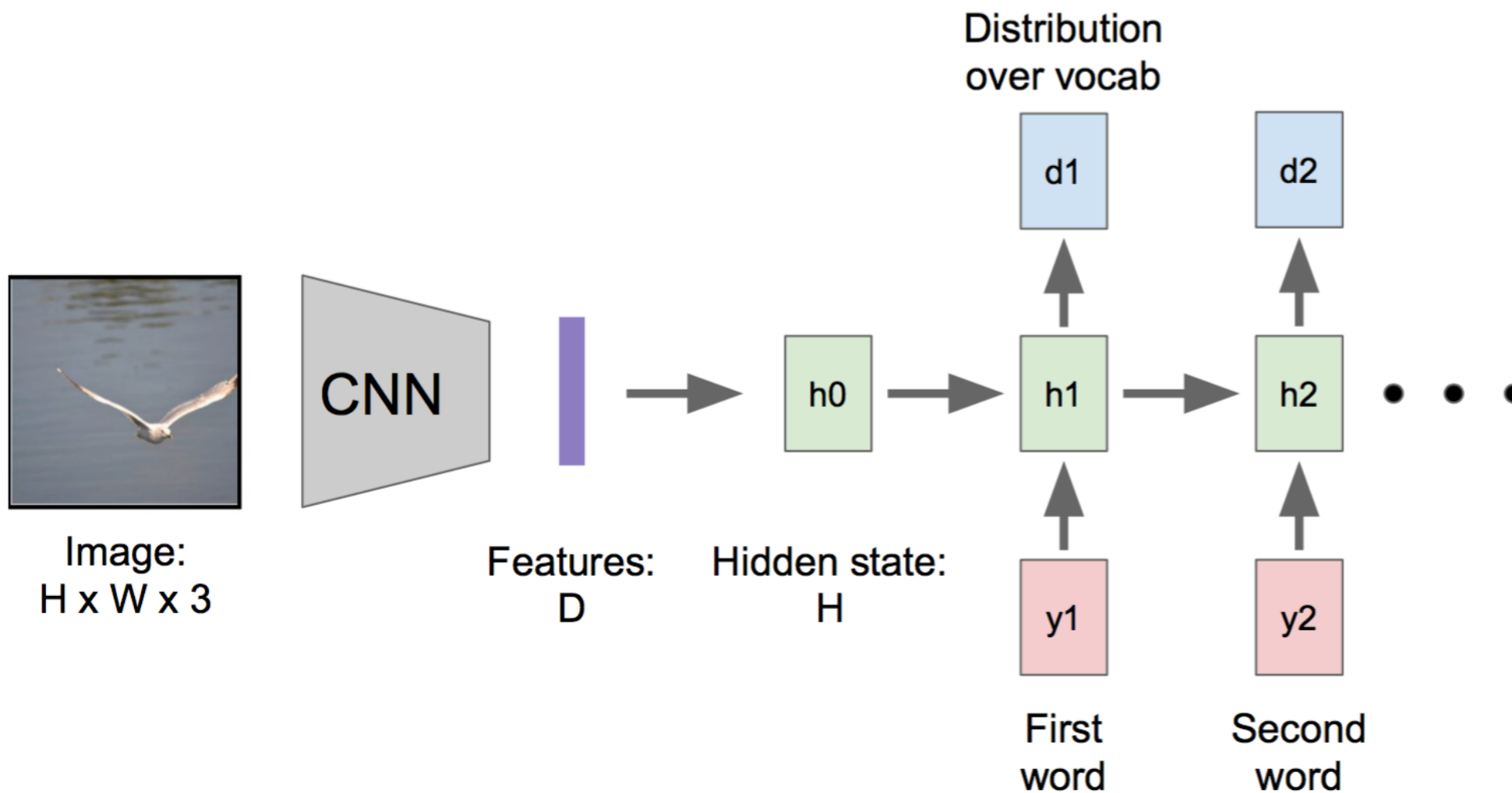
“hard” stochastic attention

sample one location

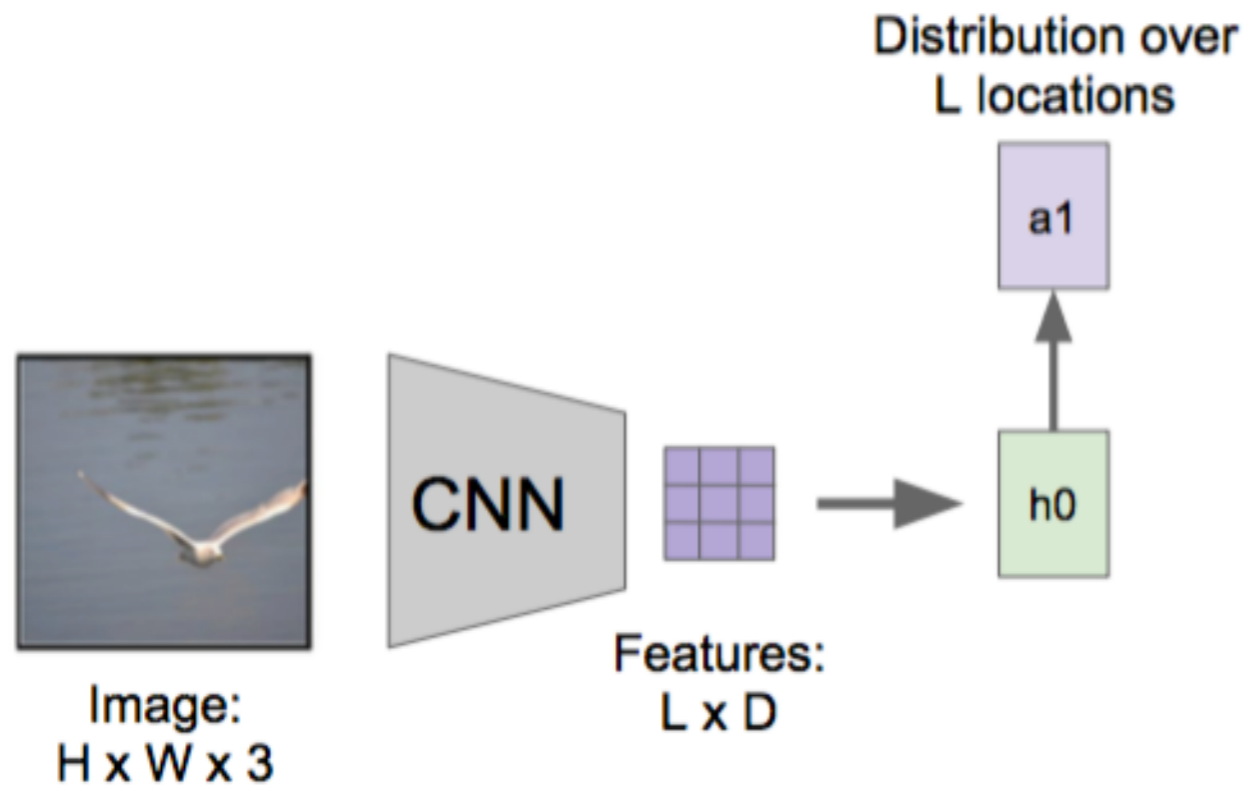
since you do argmax, gradient is zero almost everywhere, so you can't use gradient descent

reinforcement learning:
REINFORCE (Williams, 1992)

RNN for Captioning

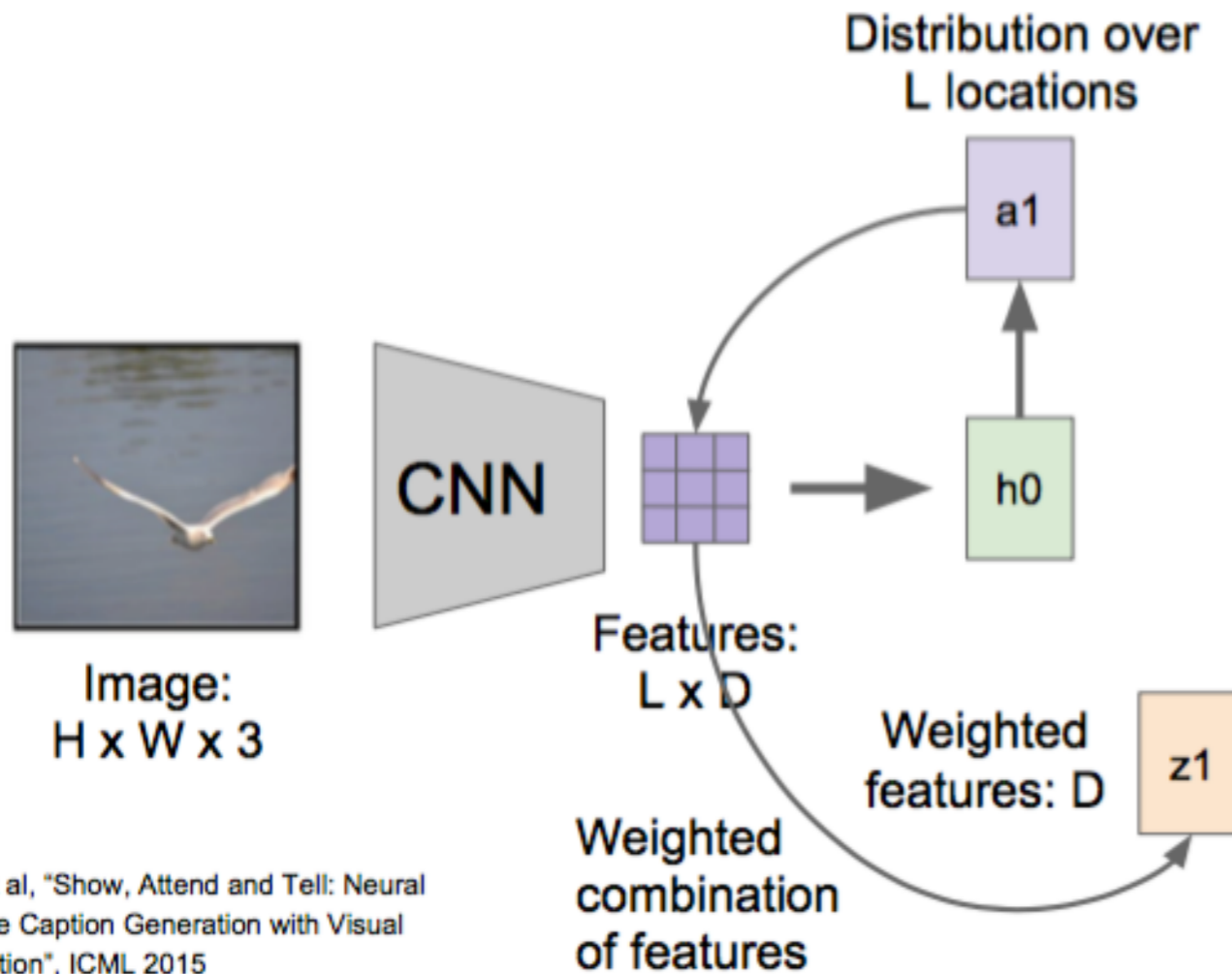


Soft Attention for Captioning



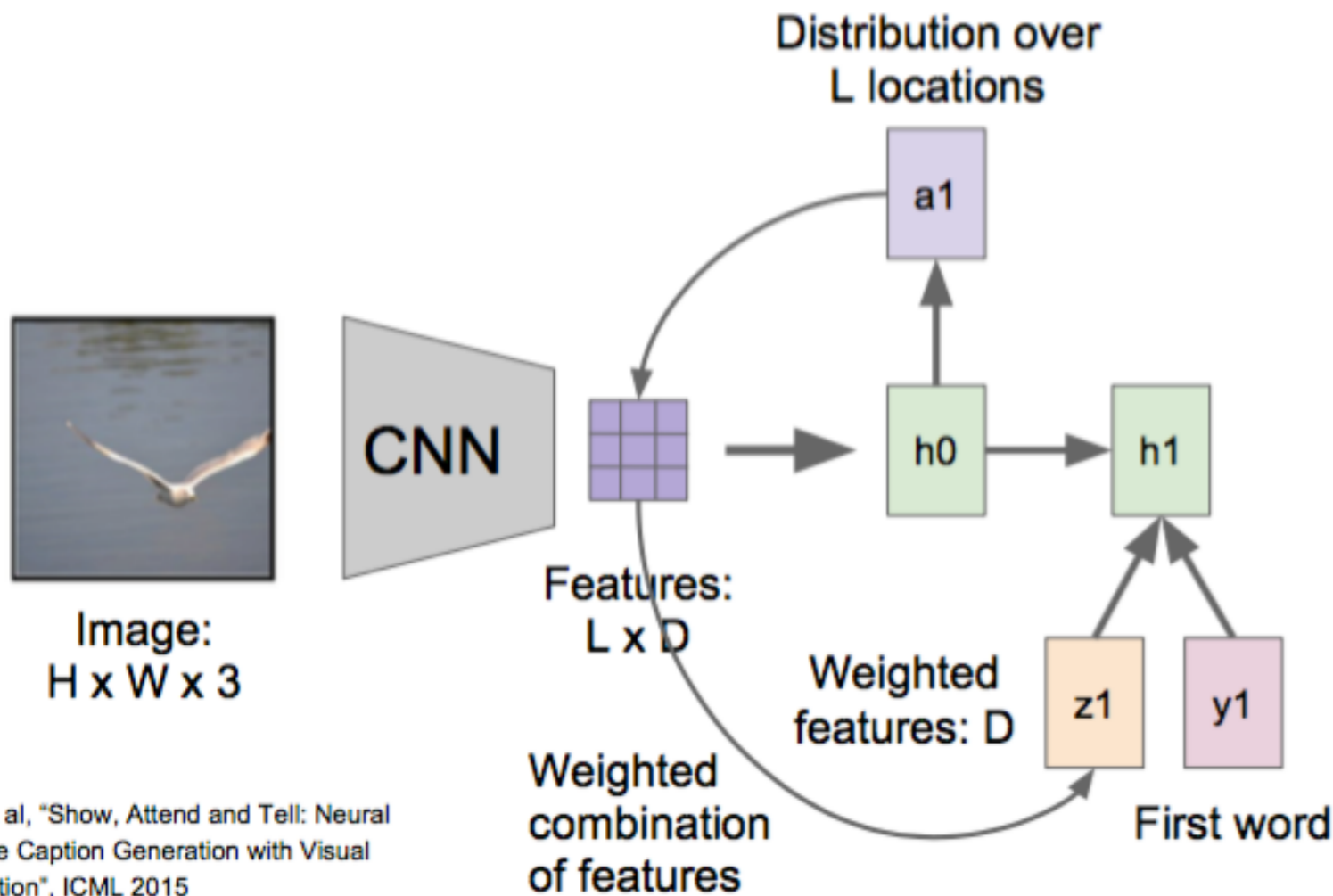
Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Soft Attention for Captioning



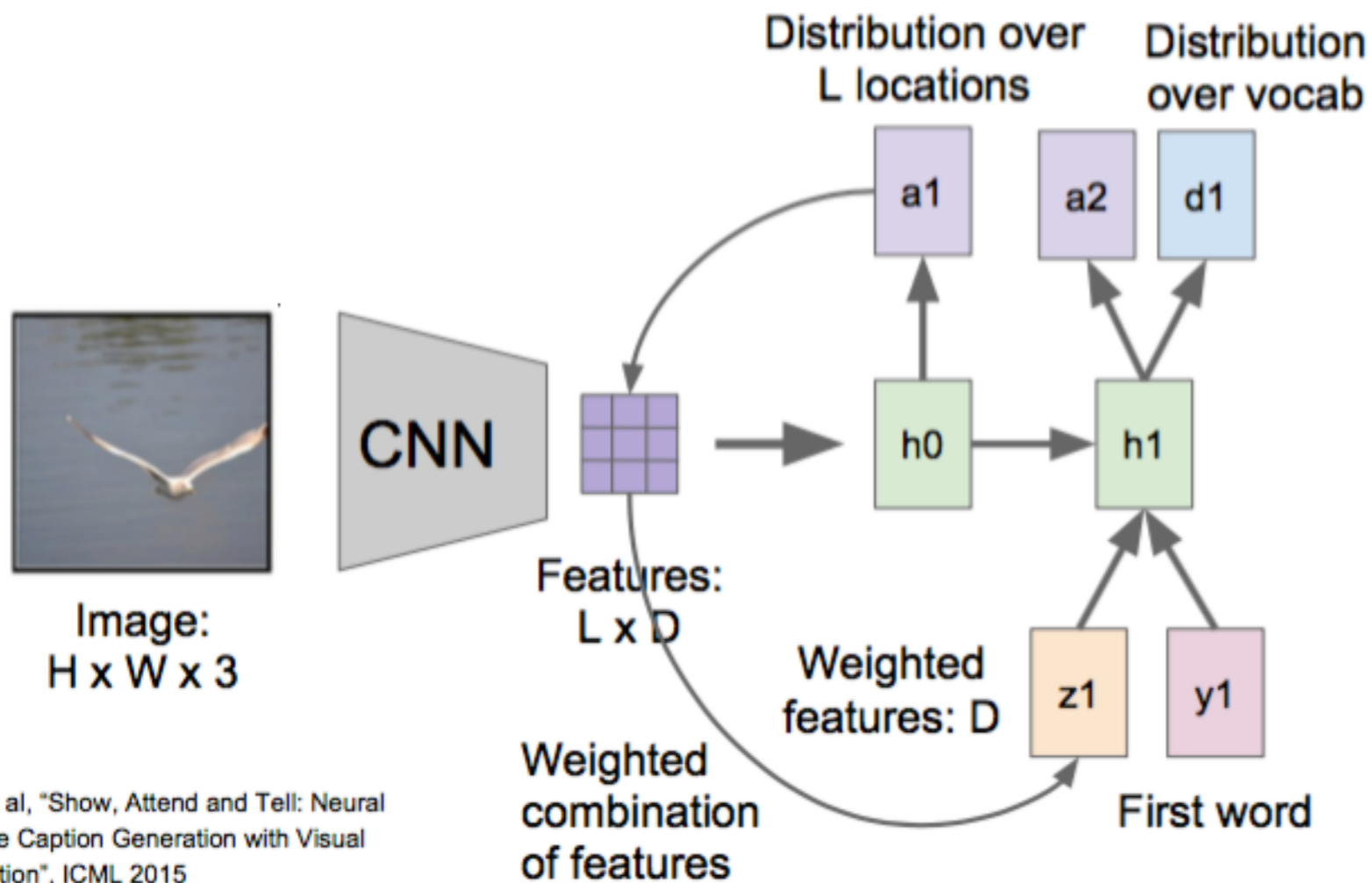
Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Soft Attention for Captioning



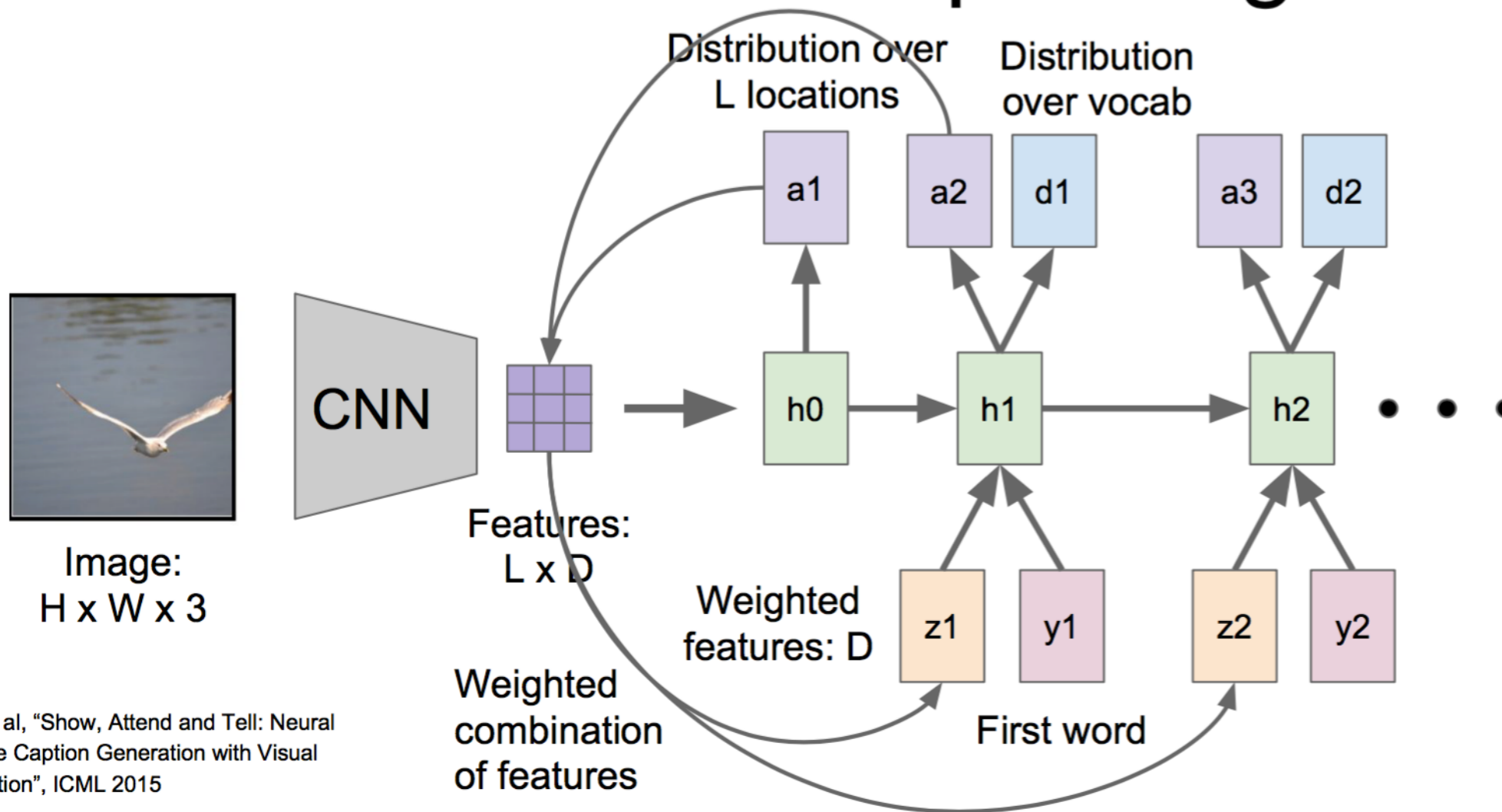
Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Soft Attention for Captioning



Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Soft Attention for Captioning



Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Evaluation

Dataset	Model	BLEU				METEOR
		B-1	B-2	B-3	B-4	
Flickr8k	Google NIC(Vinyals et al., 2014) ^{†Σ}	63	41	27	—	—
	Log Bilinear (Kiros et al., 2014a) [°]	65.6	42.4	27.7	17.7	17.31
	Soft-Attention	67	44.8	29.9	19.5	18.93
	Hard-Attention	67	45.7	31.4	21.3	20.30
Flickr30k	Google NIC ^{†$\circ\Sigma$}	66.3	42.3	27.7	18.3	—
	Log Bilinear	60.0	38	25.4	17.1	16.88
	Soft-Attention	66.7	43.4	28.8	19.1	18.49
	Hard-Attention	66.9	43.9	29.6	19.9	18.46
COCO	CMU/MS Research (Chen & Zitnick, 2014) ^a	—	—	—	—	20.41
	MS Research (Fang et al., 2014) ^{†a}	—	—	—	—	20.71
	BRNN (Karpathy & Li, 2014) [°]	64.2	45.1	30.4	20.3	—
	Google NIC ^{†$\circ\Sigma$}	66.6	46.1	32.9	24.6	—
	Log Bilinear [°]	70.8	48.9	34.4	24.3	20.03
	Soft-Attention	70.7	49.2	34.4	24.3	23.90
	Hard-Attention	71.8	50.4	35.7	25.0	23.04

BLEU is n-gram precision

METEOR is a combination of unigram-precision, unigram-recall, and a measure of fragmentation (how well-ordered matched words are)

Show, Attend and Tell: Neural Image Caption Generation with Visual Attention

Kelvin Xu
Jimmy Lei Ba
Ryan Kiros
Kyunghyun Cho
Aaron Courville
Ruslan Salakhutdinov
Richard S. Zemel
Yoshua Bengio

KELVIN.XU@UMONTREAL.CA
JIMMY@PSI.UTORONTO.CA
RKIROS@CS.TORONTO.EDU
KYUNGHYUN.CHO@UMONTREAL.CA
AARON.COURVILLE@UMONTREAL.CA
RSALAKHU@CS.TORONTO.EDU
ZEMEL@CS.TORONTO.EDU
FIND-ME@THE.WEB

Abstract

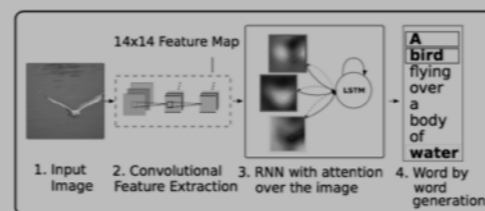
Inspired by recent work in machine translation and object detection, we introduce an attention based model that automatically learns to describe the content of images. We describe how we can train this model in a deterministic manner using standard backpropagation techniques and stochastically by maximizing a variational lower bound. We also show through visualization how the model is able to automatically learn to fix its gaze on salient objects while generating the corresponding words in the output sequence. We validate the use of attention with state-of-the-art performance on three benchmark datasets: Flickr8k, Flickr30k and MS COCO.

1. Introduction

Automatically generating captions of an image is a task very close to the heart of scene understanding — one of the primary goals of computer vision. Not only must caption generation models be powerful enough to solve the computer vision challenges of determining which objects are in an image, but they must also be capable of capturing and expressing their relationships in a natural language. For this reason, caption generation has long been viewed as a difficult problem. It is a very important challenge for machine learning algorithms, as it amounts to mimicking the remarkable human ability to compress huge amounts of salient visual information into descriptive language.

Despite the challenging nature of this task, there has been a recent surge of research interest in attacking the image caption generation problem. Aided by advances in training neural networks (Krizhevsky et al., 2012) and large classification datasets (Russakovsky et al., 2014), recent work

Figure 1. Our model learns a words/image alignment. The visualized attentional maps (3) are explained in section 3.1 & 5.4



has significantly improved the quality of caption generation using a combination of convolutional neural networks (convnets) to obtain vectorial representation of images and recurrent neural networks to decode those representations into natural language sentences (see Sec. 2).

One of the most curious facets of the human visual system is the presence of attention (Rensink, 2000; Corbetta & Shulman, 2002). Rather than compress an entire image into a static representation, attention allows for salient features to dynamically come to the forefront as needed. This is especially important when there is a lot of clutter in an image. Using representations (such as those from the top layer of a convnet) that distill information in image down to the most salient objects is one effective solution that has been widely adopted in previous work. Unfortunately, this has one potential drawback of losing information which could be useful for richer, more descriptive captions. Using more low-level representation can help preserve this information. However working with these features necessitates a powerful mechanism to steer the model to information important to the task at hand.

In this paper, we describe approaches to caption generation that attempt to incorporate a form of attention with

DRAW: A Recurrent Neural Network For Image Generation

Karol Gregor
Ivo Danihelka
Alex Graves
Danilo Jimenez Rezende
Daan Wierstra
Google DeepMind

KAROLG@GOOGLE.COM
DANIHELKA@GOOGLE.COM
GRAVESA@GOOGLE.COM
DANILOR@GOOGLE.COM
WIERSTRA@GOOGLE.COM

Abstract

This paper introduces the *Deep Recurrent Attentive Writer* (DRAW) neural network architecture for image generation. DRAW networks combine a novel spatial attention mechanism that mimics the foveation of the human eye, with a sequential variational auto-encoding framework that allows for the iterative construction of complex images. The system substantially improves on the state of the art for generative models on MNIST, and, when trained on the Street View House Numbers dataset, it generates images that cannot be distinguished from real data with the naked eye.

1. Introduction

A person asked to draw, paint or otherwise recreate a visual scene will naturally do so in a sequential, iterative fashion, reassessing their handiwork after each modification. Rough outlines are gradually replaced by precise forms, lines are sharpened, darkened or erased, shapes are altered, and the final picture emerges. Most approaches to automatic image generation, however, aim to generate entire scenes at once. In the context of generative neural networks, this typically means that all the pixels are conditioned on a single latent distribution (Dayan et al., 1995; Hinton & Salakhutdinov, 2006; Larochelle & Murray, 2011). As well as precluding the possibility of iterative self-correction, the “one shot” approach is fundamentally difficult to scale to large images. The *Deep Recurrent Attentive Writer* (DRAW) architecture represents a shift towards a more natural form of image construction, in which parts of a scene are created independently from others, and approximate sketches are successively refined.

Proceedings of the 32nd International Conference on Machine Learning, Lille, France, 2015. JMLR: W&CP volume 37. Copyright 2015 by the author(s).

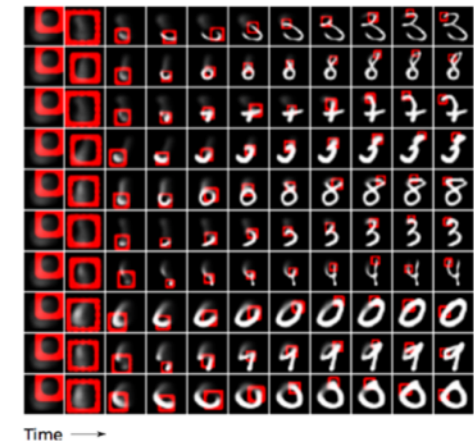


Figure 1. A trained DRAW network generating MNIST digits. Each row shows successive stages in the generation of a single digit. Note how the lines composing the digits appear to be “drawn” by the network. The red rectangle delimits the area attended to by the network at each time-step, with the focal precision indicated by the width of the rectangle border.

The core of the DRAW architecture is a pair of recurrent neural networks: an *encoder* network that compresses the real images presented during training, and a *decoder* that reconstitutes images after receiving codes. The combined system is trained end-to-end with stochastic gradient descent, where the loss function is a variational upper bound on the log-likelihood of the data. It therefore belongs to the family of *variational auto-encoders*, a recently emerged hybrid of deep learning and variational inference that has led to significant advances in generative modelling (Gregor et al., 2014; Kingma & Welling, 2014; Rezende et al., 2014; Mnih & Gregor, 2014; Salimans et al., 2014). Where DRAW differs from its siblings is that, rather than generat-

Describe images

Generate images

Goal: Generate images

"dreaming up" images - transforming random noise into an endless stream of images that the model has never even seen before

Works for MNIST and Street View House Numbers

arXiv:1502.04623v2 [cs.CV] 20 May 2015

DRAW: A Recurrent Neural Network For Image Generation

Karol Gregor
Ivo Danihelka
Alex Graves
Danilo Jimenez Rezende
Daan Wierstra
Google DeepMind

KAROLG@GOOGLE.COM
DANIHELKA@GOOGLE.COM
GRAVESA@GOOGLE.COM
DANILOR@GOOGLE.COM
WIERSTRA@GOOGLE.COM

Abstract

This paper introduces the *Deep Recurrent Attentive Writer* (DRAW) neural network architecture for image generation. DRAW networks combine a novel spatial attention mechanism that mimics the foveation of the human eye, with a sequential variational auto-encoding framework that allows for the iterative construction of complex images. The system substantially improves on the state of the art for generative models on MNIST, and, when trained on the Street View House Numbers dataset, it generates images that cannot be distinguished from real data with the naked eye.

1. Introduction

A person asked to draw, paint or otherwise recreate a visual scene will naturally do so in a sequential, iterative fashion, reassessing their handiwork after each modification. Rough outlines are gradually replaced by precise forms, lines are sharpened, darkened or erased, shapes are altered, and the final picture emerges. Most approaches to automatic image generation, however, aim to generate entire scenes at once. In the context of generative neural networks, this typically means that all the pixels are conditioned on a single latent distribution (Dayan et al., 1995; Hinton & Salakhutdinov, 2006; Larochelle & Murray, 2011). As well as precluding the possibility of iterative self-correction, the "one shot" approach is fundamentally difficult to scale to large images. The *Deep Recurrent Attentive Writer* (DRAW) architecture represents a shift towards a more natural form of image construction, in which parts of a scene are created independently from others, and approximate sketches are successively refined.

Proceedings of the 32nd International Conference on Machine Learning, Lille, France, 2015. JMLR: W&CP volume 37. Copyright 2015 by the author(s).

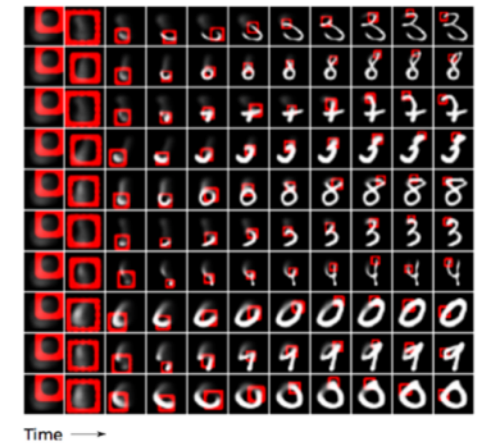


Figure 1. A trained DRAW network generating MNIST digits. Each row shows successive stages in the generation of a single digit. Note how the lines composing the digits appear to be "drawn" by the network. The red rectangle delimits the area attended to by the network at each time-step, with the focal precision indicated by the width of the rectangle border.

The core of the DRAW architecture is a pair of recurrent neural networks: an *encoder* network that compresses the real images presented during training, and a *decoder* that reconstitutes images after receiving codes. The combined system is trained end-to-end with stochastic gradient descent, where the loss function is a variational upper bound on the log-likelihood of the data. It therefore belongs to the family of *variational auto-encoders*, a recently emerged hybrid of deep learning and variational inference that has led to significant advances in generative modelling (Gregor et al., 2014; Kingma & Welling, 2014; Rezende et al., 2014; Mnih & Gregor, 2014; Salimans et al., 2014). Where DRAW differs from its siblings is that, rather than generat-

Variational Autoencoder

Encoder: determines a distribution that captures salient information about the input data

Decoder: samples from the distribution

Spatial selective attention mechanism that mimics the foveation of the human eye, with a sequential variational auto-encoding framework that allows for the iterative construction of complex images

arXiv:1502.04623v2 [cs.CV] 20 May 2015

DRAW: A Recurrent Neural Network For Image Generation

Karol Gregor
Ivo Danihelka
Alex Graves
Danilo Jimenez Rezende
Daan Wierstra
Google DeepMind

KAROLG@GOOGLE.COM
DANIHELKA@GOOGLE.COM
GRAVESA@GOOGLE.COM
DANILOR@GOOGLE.COM
WIERSTRA@GOOGLE.COM

Abstract

This paper introduces the *Deep Recurrent Attentive Writer* (DRAW) neural network architecture for image generation. DRAW networks combine a novel spatial attention mechanism that mimics the foveation of the human eye, with a sequential variational auto-encoding framework that allows for the iterative construction of complex images. The system substantially improves on the state of the art for generative models on MNIST, and, when trained on the Street View House Numbers dataset, it generates images that cannot be distinguished from real data with the naked eye.

1. Introduction

A person asked to draw, paint or otherwise recreate a visual scene will naturally do so in a sequential, iterative fashion, reassessing their handiwork after each modification. Rough outlines are gradually replaced by precise forms, lines are sharpened, darkened or erased, shapes are altered, and the final picture emerges. Most approaches to automatic image generation, however, aim to generate entire scenes at once. In the context of generative neural networks, this typically means that all the pixels are conditioned on a single latent distribution (Dayan et al., 1995; Hinton & Salakhutdinov, 2006; Larochelle & Murray, 2011). As well as precluding the possibility of iterative self-correction, the “one shot” approach is fundamentally difficult to scale to large images. The *Deep Recurrent Attentive Writer* (DRAW) architecture represents a shift towards a more natural form of image construction, in which parts of a scene are created independently from others, and approximate sketches are successively refined.

Proceedings of the 32nd International Conference on Machine Learning, Lille, France, 2015. JMLR: W&CP volume 37. Copyright 2015 by the author(s).

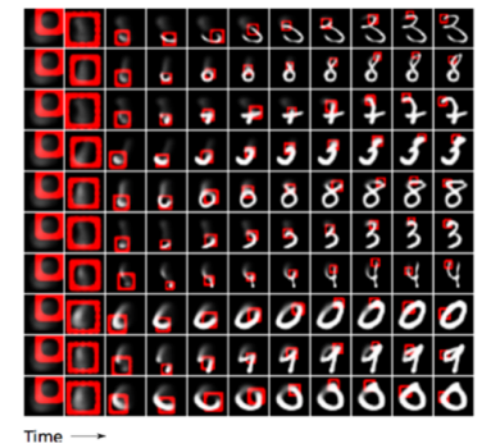
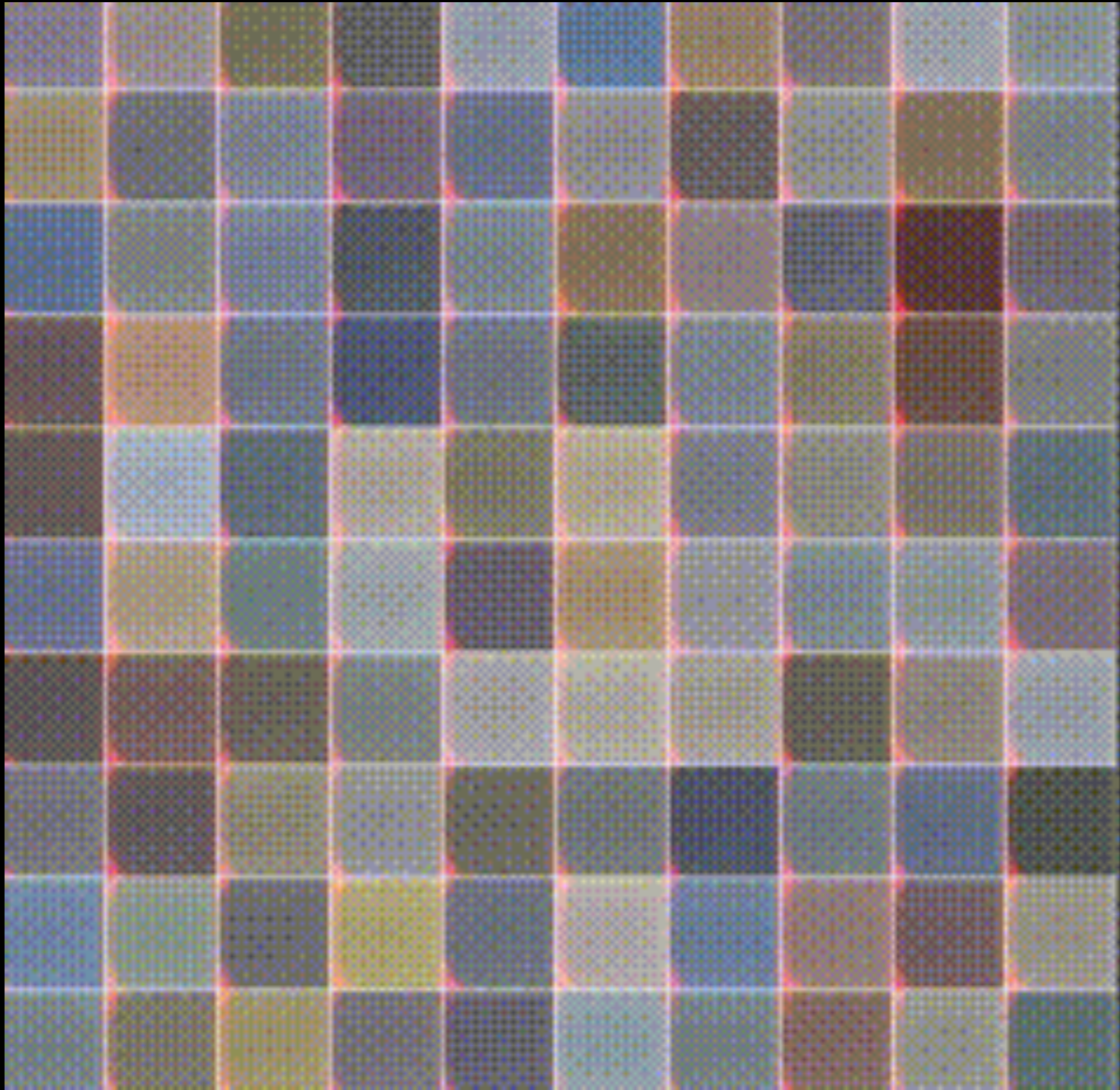
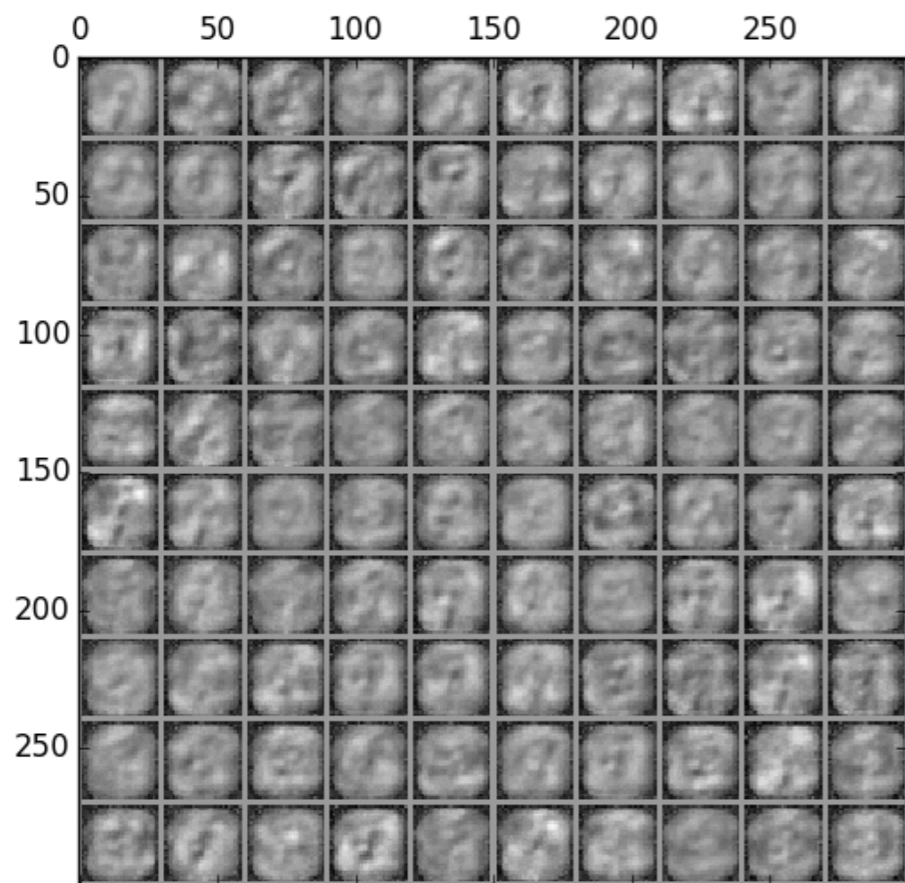


Figure 1. A trained DRAW network generating MNIST digits. Each row shows successive stages in the generation of a single digit. Note how the lines composing the digits appear to be “drawn” by the network. The red rectangle delimits the area attended to by the network at each time-step, with the focal precision indicated by the width of the rectangle border.

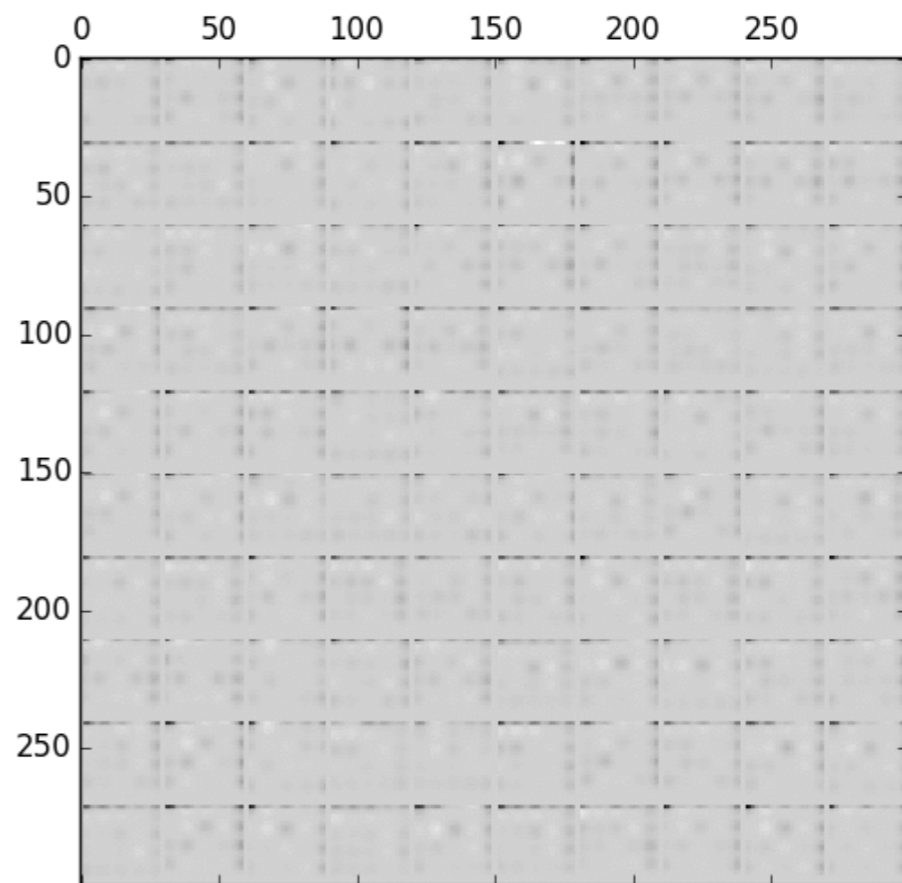
The core of the DRAW architecture is a pair of recurrent neural networks: an *encoder* network that compresses the real images presented during training, and a *decoder* that reconstitutes images after receiving codes. The combined system is trained end-to-end with stochastic gradient descent, where the loss function is a variational upper bound on the log-likelihood of the data. It therefore belongs to the family of *variational auto-encoders*, a recently emerged hybrid of deep learning and variational inference that has led to significant advances in generative modelling (Gregor et al., 2014; Kingma & Welling, 2014; Rezende et al., 2014; Mnih & Gregor, 2014; Salimans et al., 2014). Where DRAW differs from its siblings is that, rather than generat-



without attention



with attention



**network decides at every
step "where to read",
"where to write", and
"what to write"**

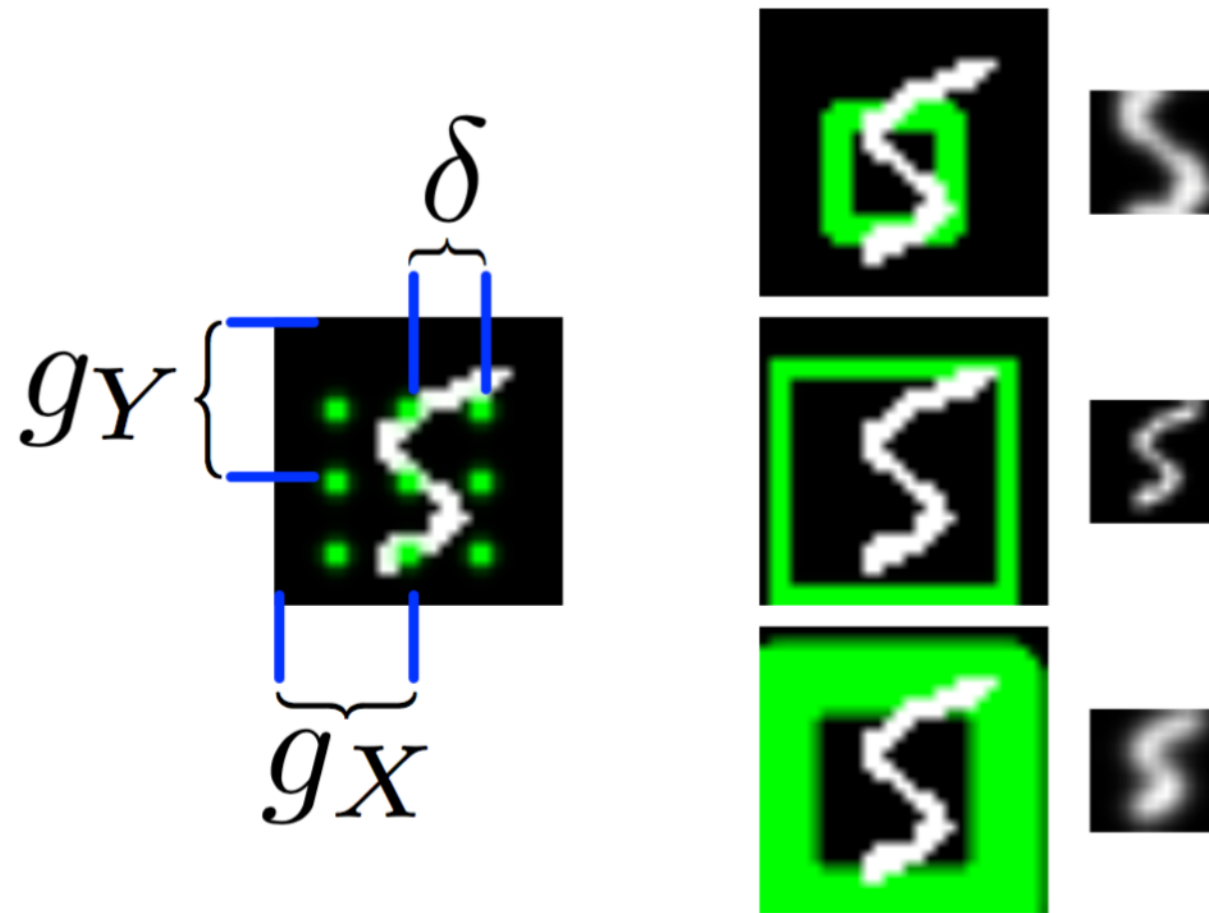


Figure 3. Left: A 3×3 grid of filters superimposed on an image. The stride (δ) and centre location (g_X, g_Y) are indicated. **Right:** Three $N \times N$ patches extracted from the image ($N = 12$). The green rectangles on the left indicate the boundary and precision (σ) of the patches, while the patches themselves are shown to the right. The top patch has a small δ and high σ , giving a zoomed-in but blurry view of the centre of the digit; the middle patch has large δ and low σ , effectively downsampling the whole image; and the bottom patch has high δ and σ .

Selective Attention Model

An $N \times N$ grid of Gaussian filters is positioned on the image by specifying the co-ordinates of the grid centre and the stride distance between adjacent filters.

stride / delta = zoom

it starts covering the entire image and then zooms in

Implementations

Show, Attend, and Tell

https://github.com/jazzsaxmafia/show_attend_and_tell.tensorflow/

DRAW

<https://github.com/ikostrikov/TensorFlow-VAE-GAN-DRAW>

<https://github.com/ericjang/draw>

Great blog post about DRAW

<http://evjang.com/articles/draw>

Ask, Attend and Answer:
Exploring Question-Guided
Spatial Attention for
Visual Question Answering

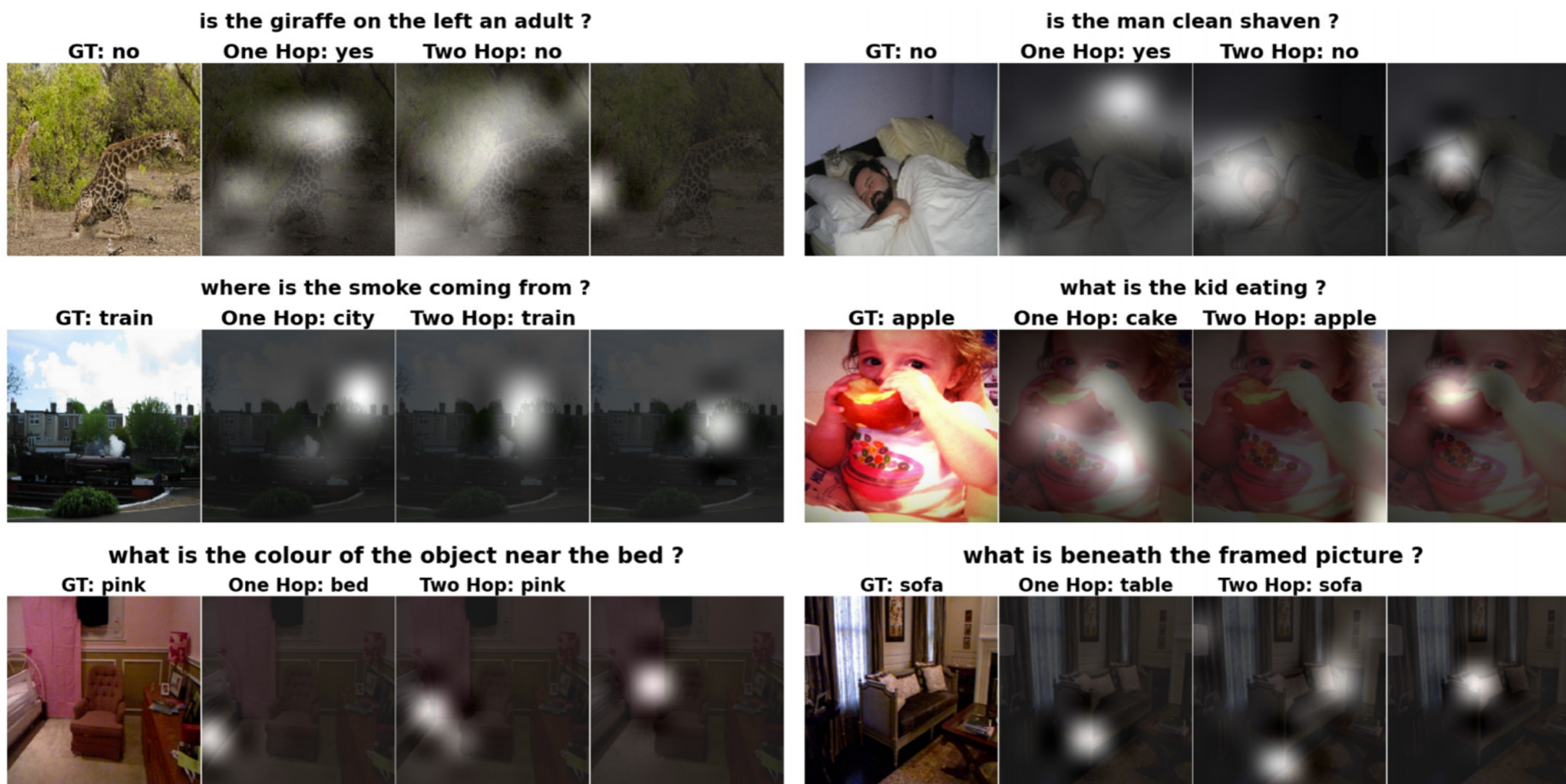
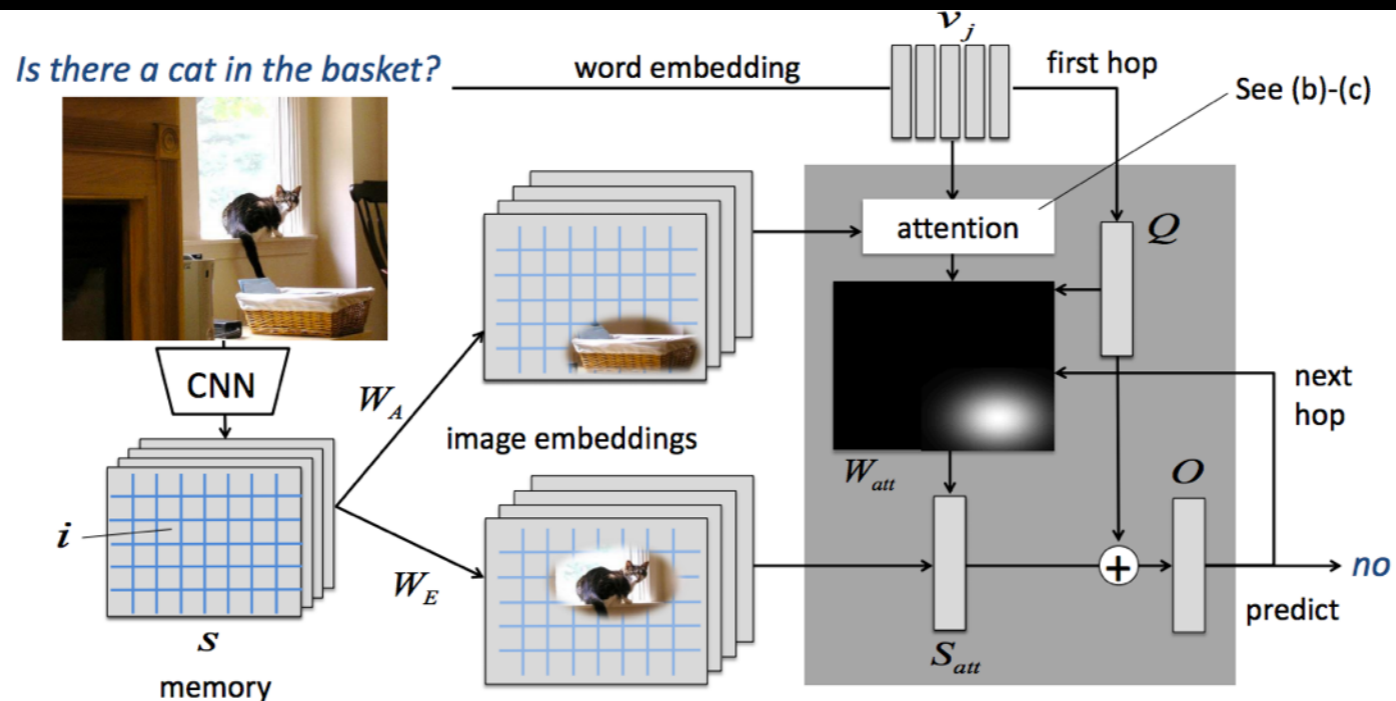


Figure 6. Visualization of the spatial attention weights in the one-hop and two-hop model on VQA (top two rows) and DAQUAR (bottom row) datasets. For each image and question pair, we show the original image, the attention weights W_{att} of the one-hop model, and the two attention weights W_{att} and W_{att2} of the two-hop model in order.

Image Captioning & Attention

Deep Learning

Hendrik Heuer

University of Amsterdam

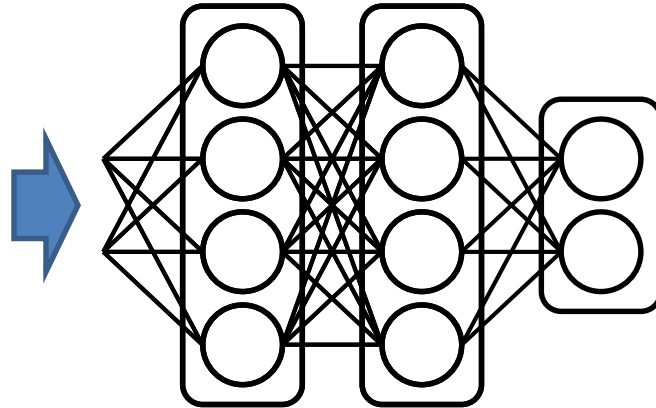


Unsupervised Visual Representation Learning by Context Prediction

Carl Doersch

Joint work with Alexei A. Efros &
Abhinav Gupta

ImageNet + Deep Learning

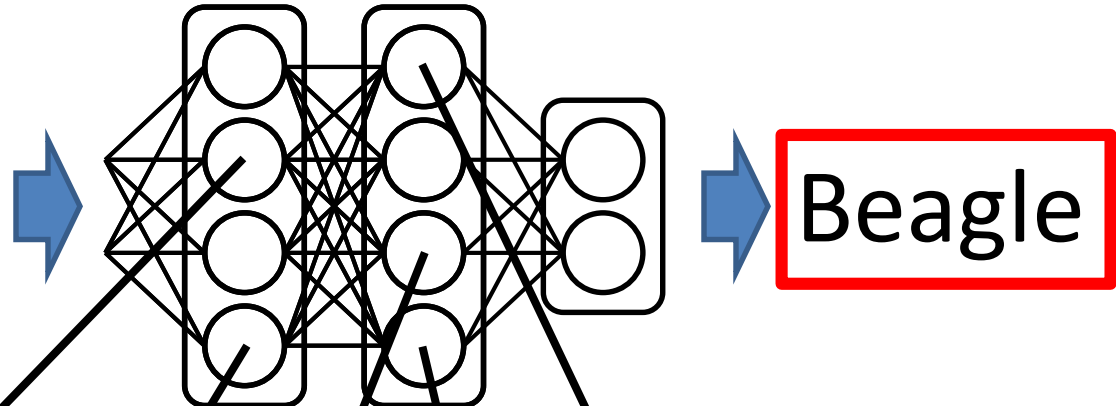


Beagle



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

ImageNet + Deep Learning



Materials?

Parts?

Pose?

Do we even need this sort of labels?

Geometry?

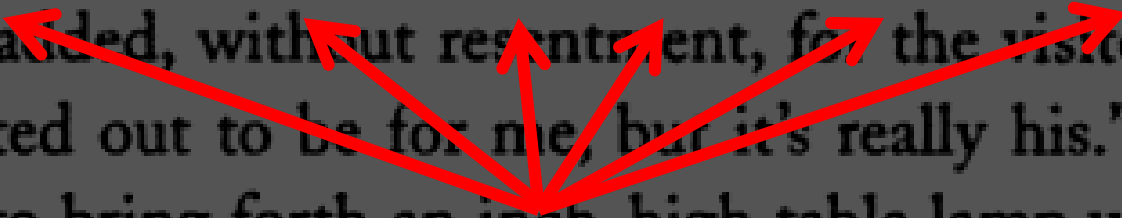
Boundaries?

Context as Supervision

[Collobert & Weston 2008; Mikolov et al. 2013]

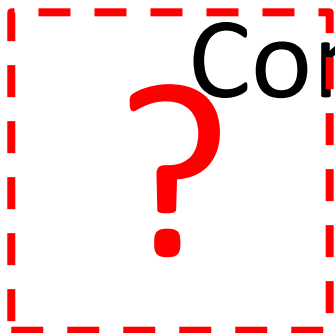
house, where the professor lived without his wife and child; or so he said jokingly sometimes: "Here's where I live. My house." His daughter often added, without resentment, for the visitor's information, "It started out to be for me, but it's really his." And she might reach in to bring forth an inch-high table lamp with fluted shade, or a blue dish the size of her little fingernail, marked "Kitty" and half full of eternal milk, but she was sure to replace these, after they had been admired, pretty near exactly where they had been. The little house was very orderly, and just big enough for all it contained, though to some tastes the bric-à-brac in the parlor might seem excessive. The daughter's preference was for the store-bought gimmicks and appliances, the toasters and carpet sweepers of Lilliput, but she knew that most adult visitors would

Deep
Net



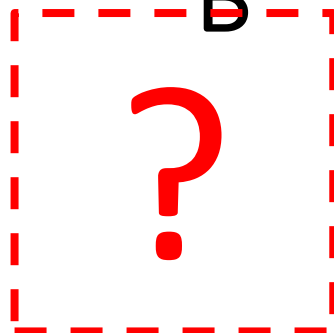
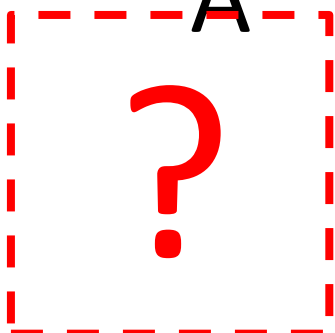
store-bought gimmicks and appliances, the toasters and carpet

Context Prediction for Images

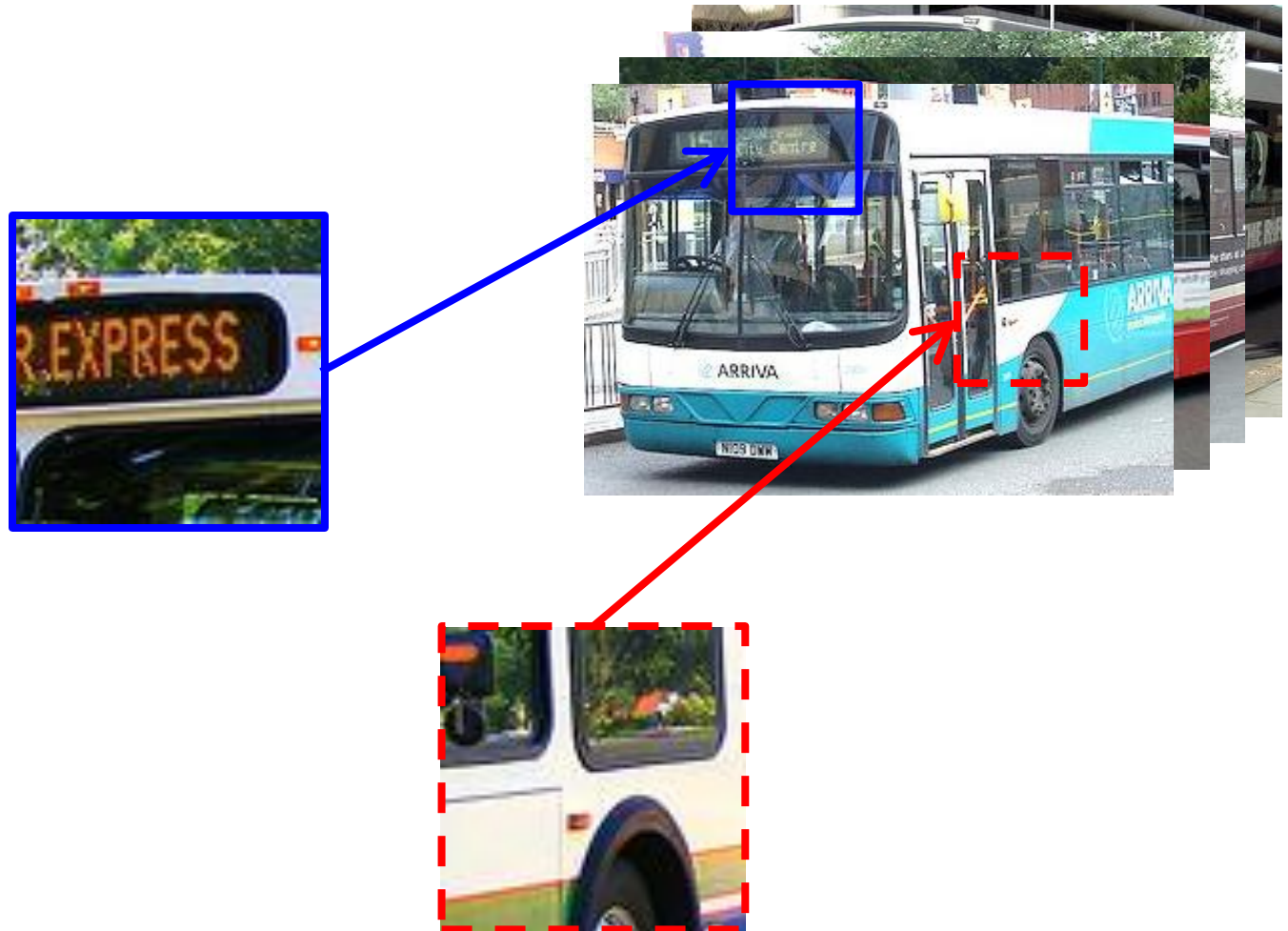


A

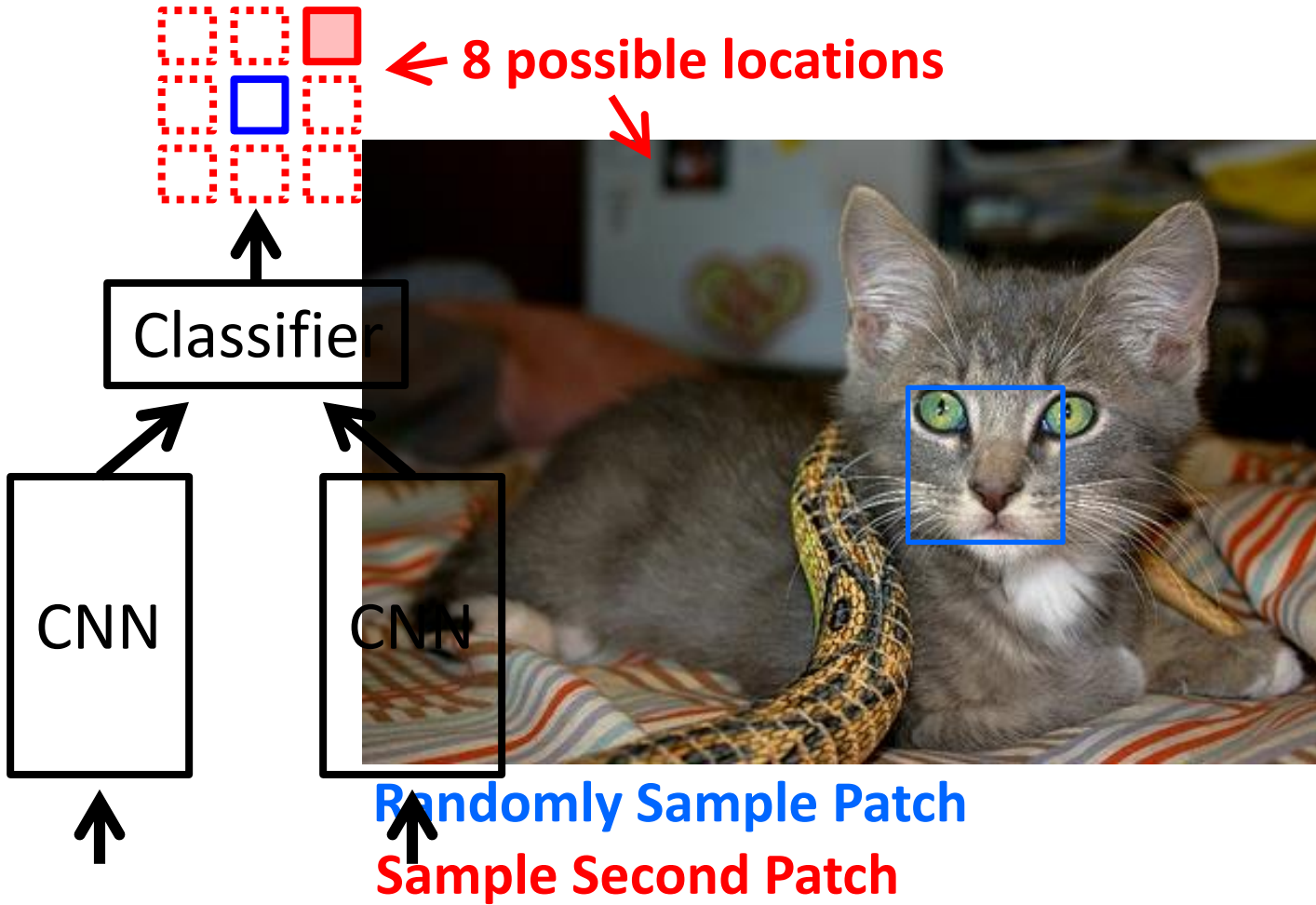
B

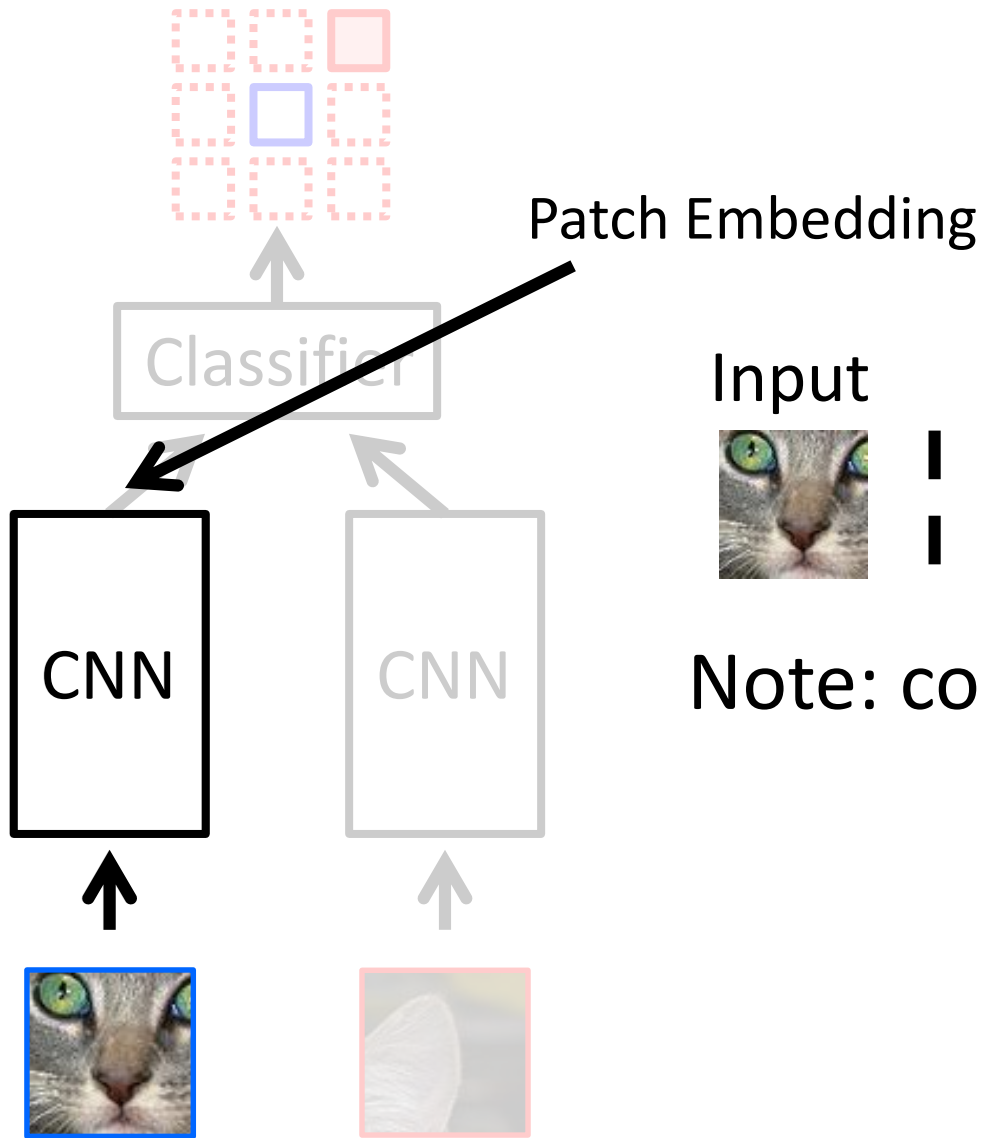


Semantics from a non-semantic task



Relative Position Task



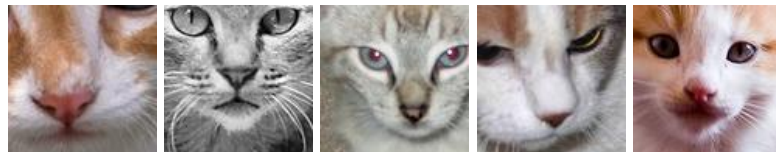


Input



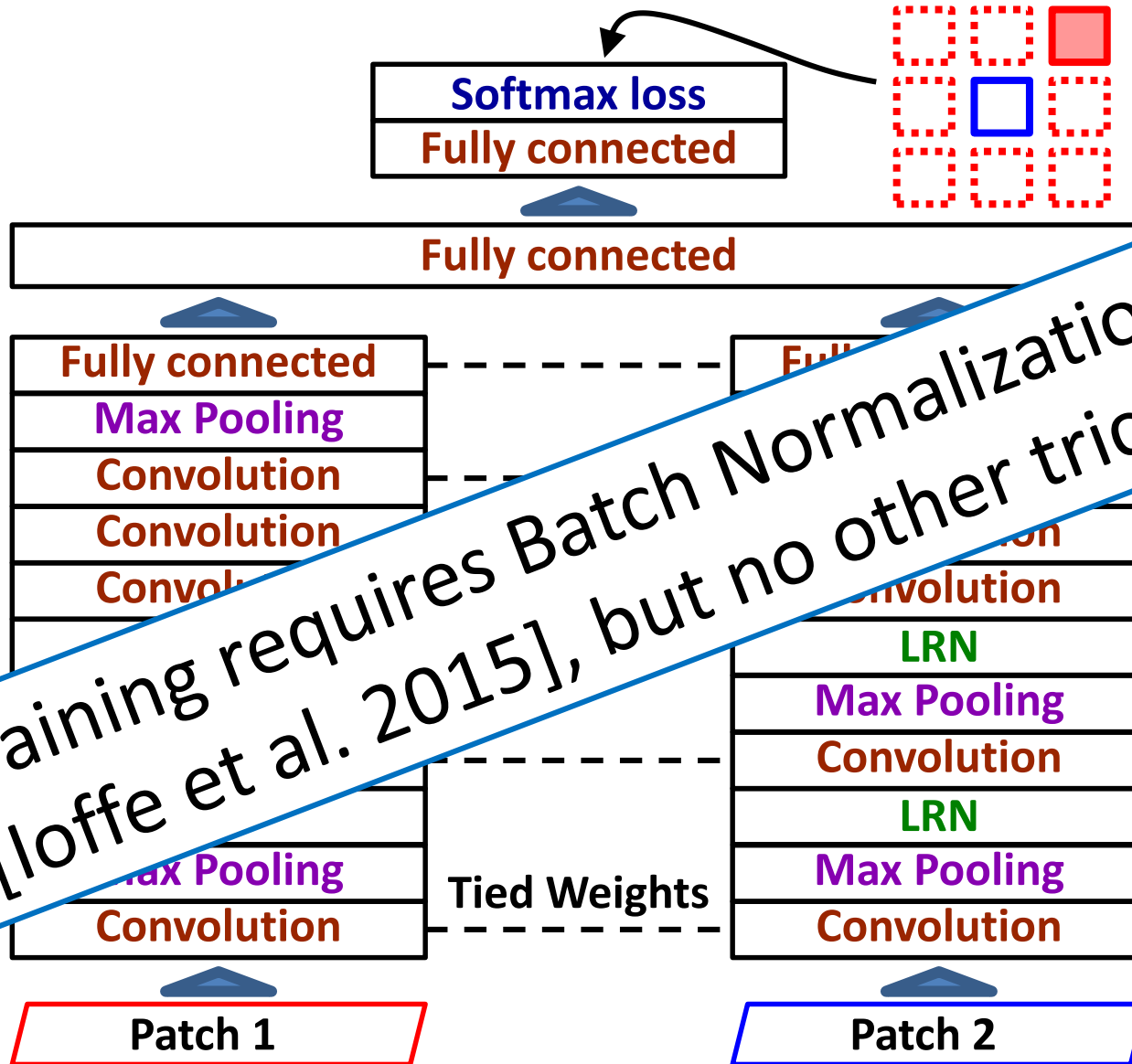
!

Nearest Neighbors



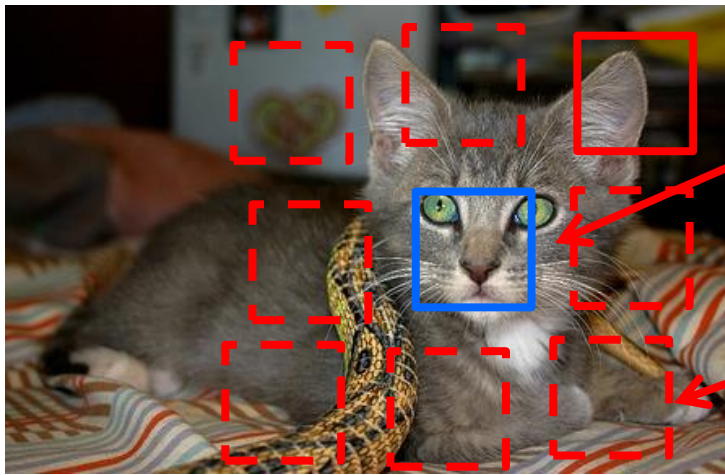
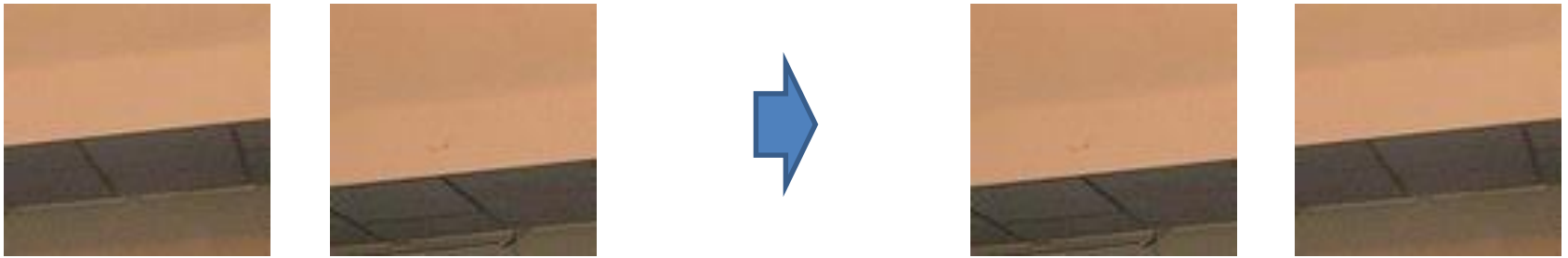
Note: connects ***across*** instances!

Architecture



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

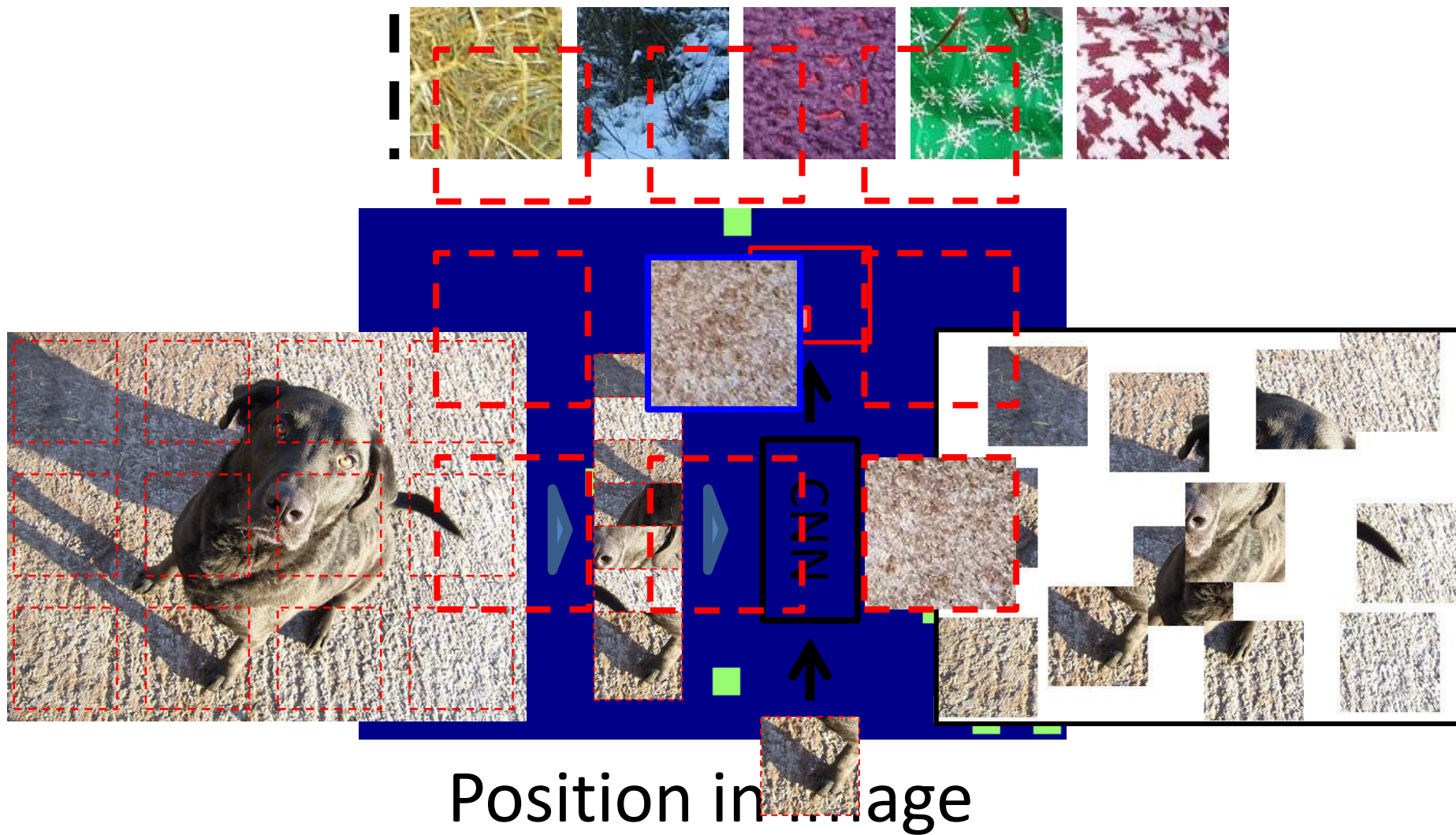
Avoiding Trivial Shortcuts



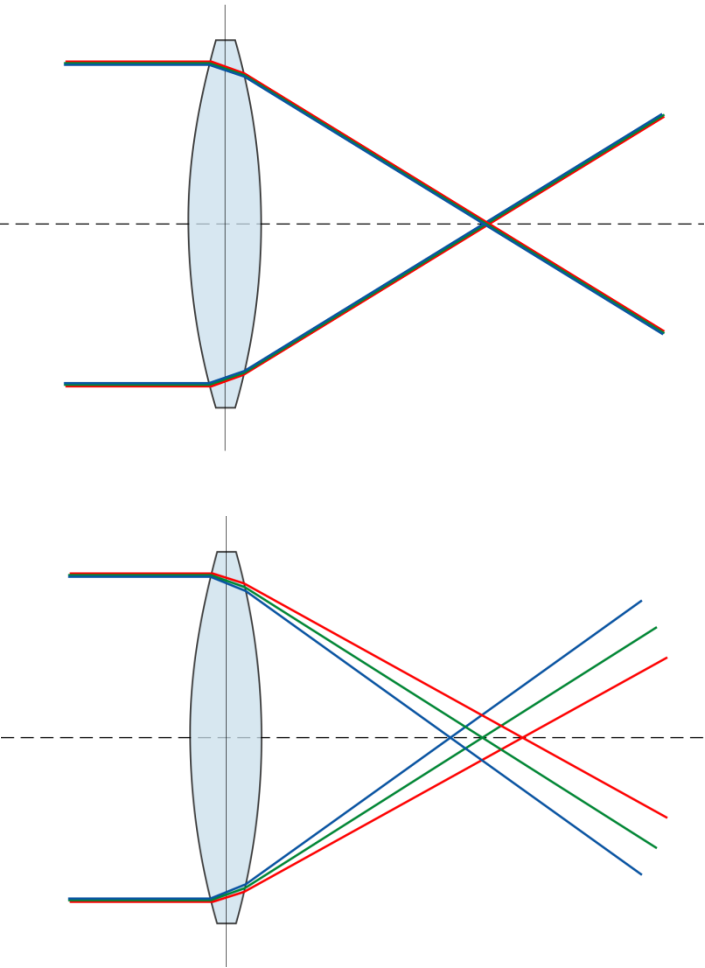
Include a gap

Jitter the patch locations

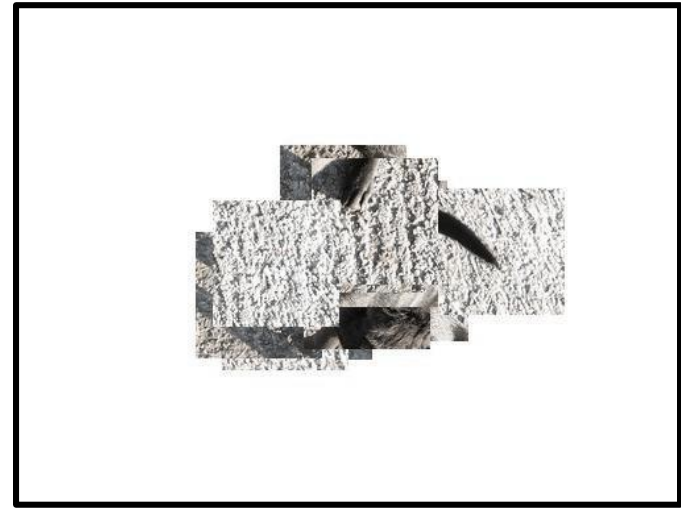
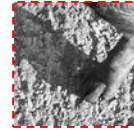
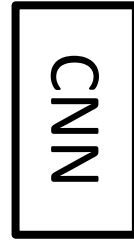
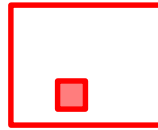
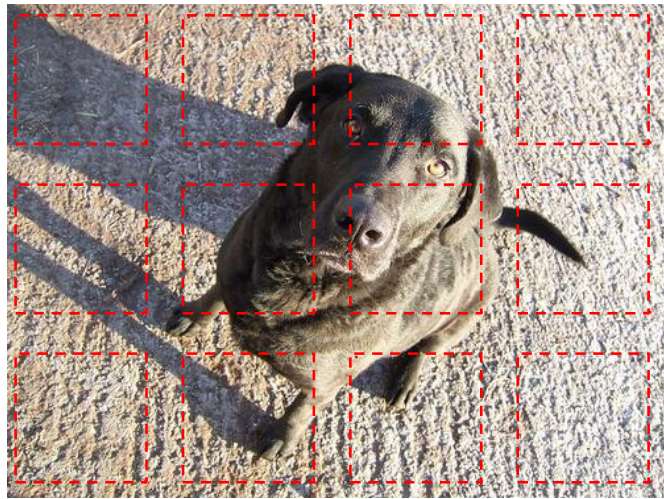
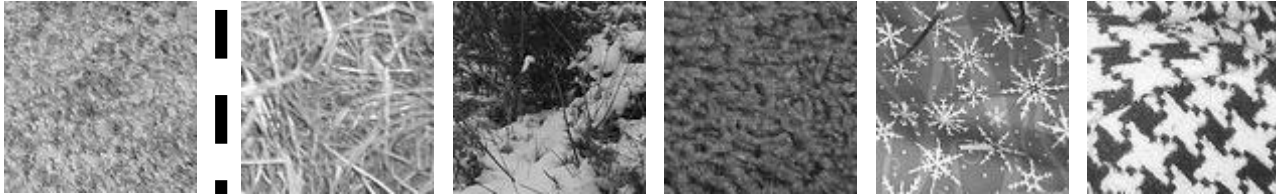
A Not-So “Trivial” Shortcut



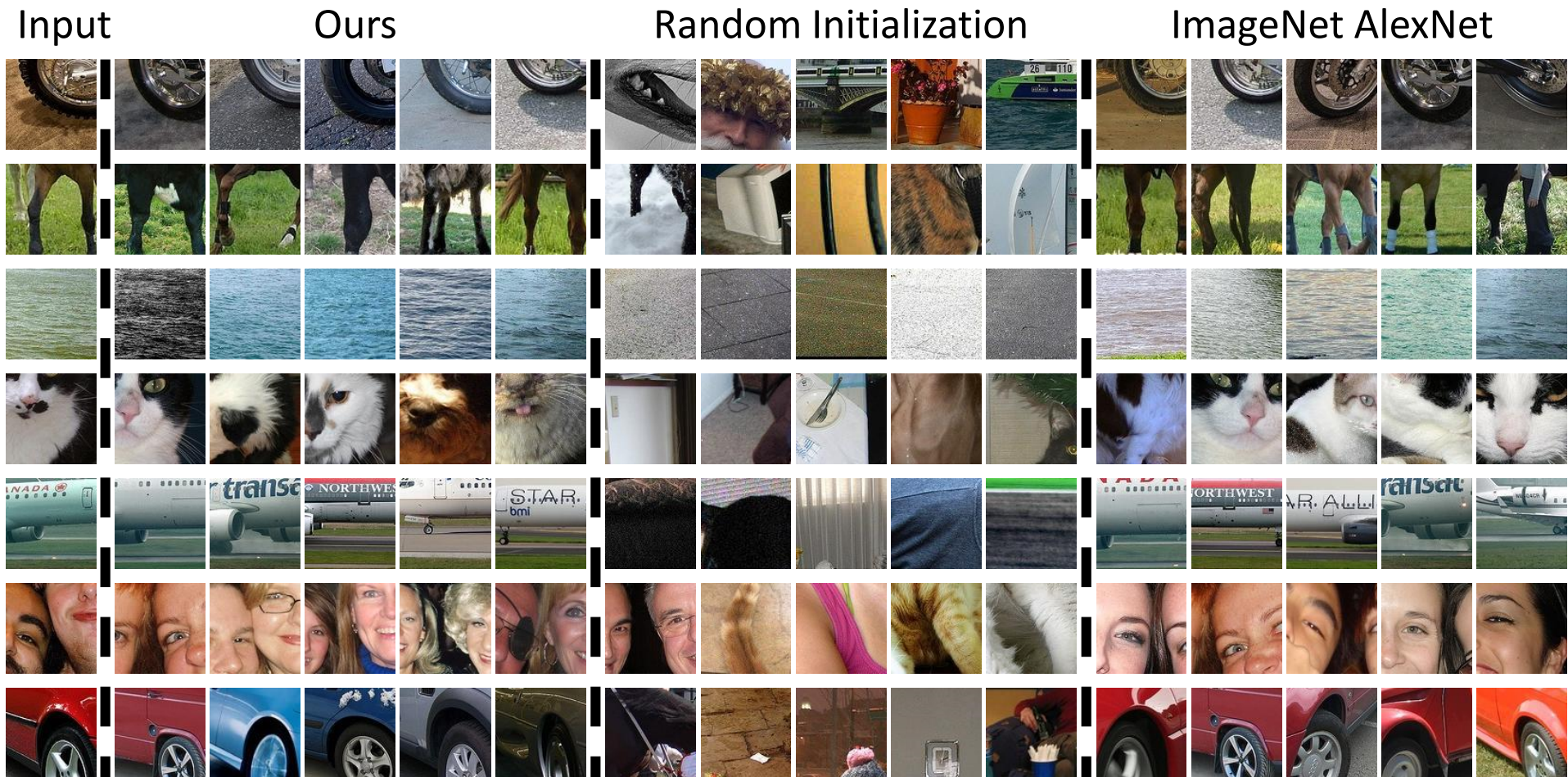
Chromatic Aberration



Chromatic Aberration



What is learned?



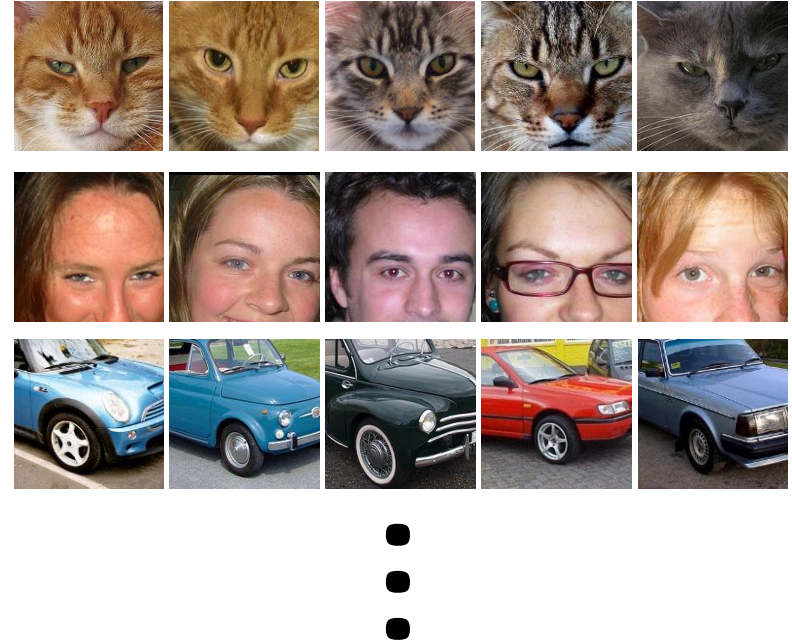
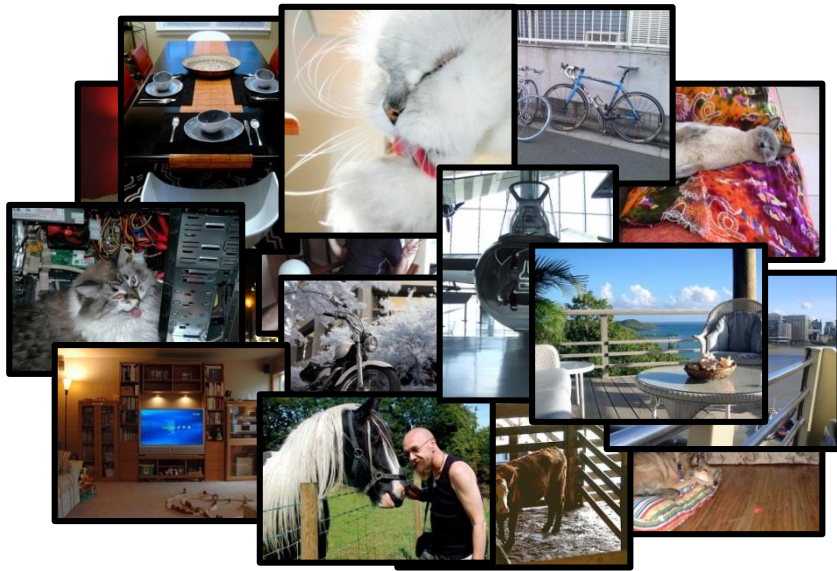
Still don't capture everything



You don't always need to learn!



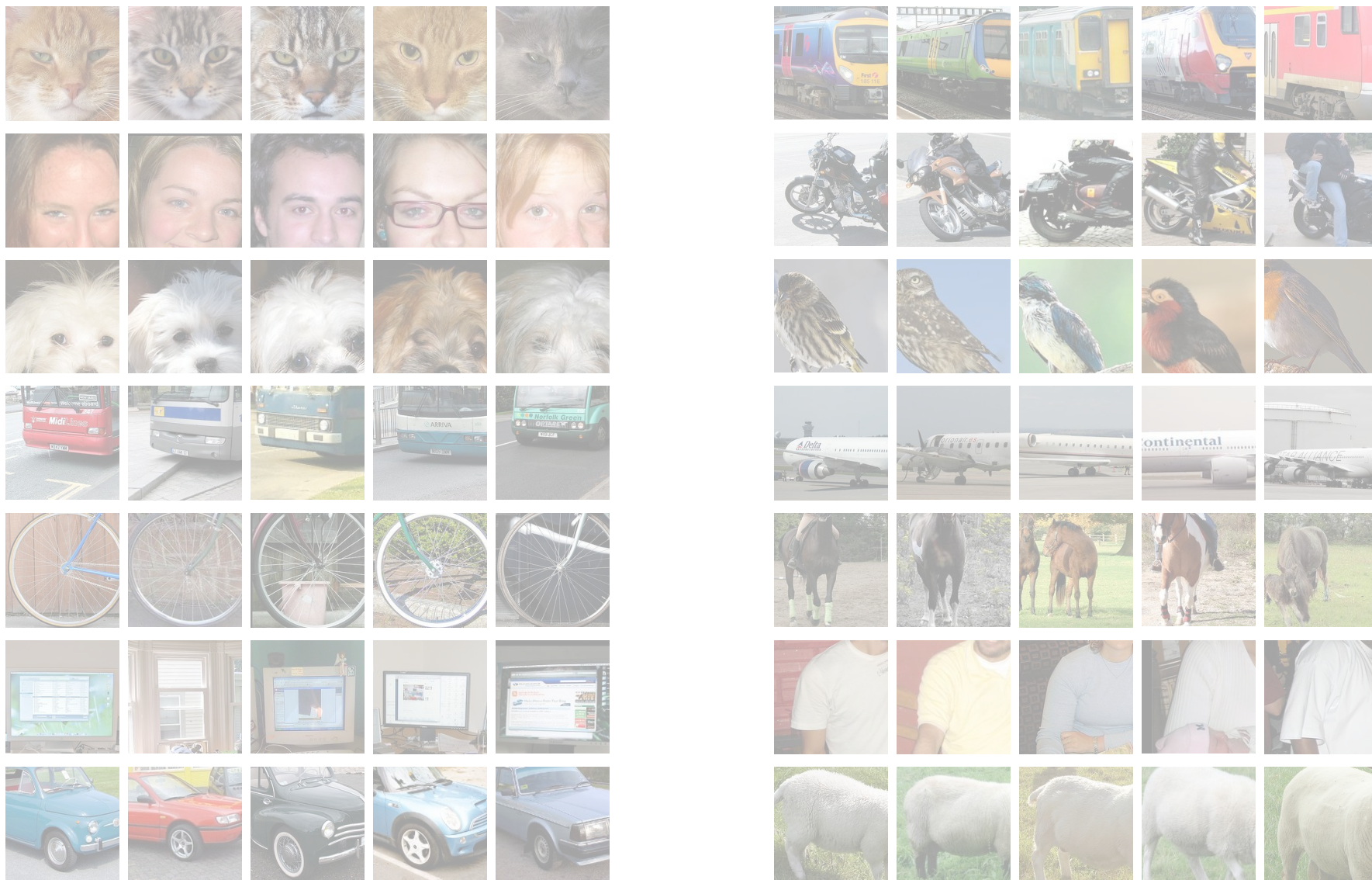
Visual Data Mining



Via Geometric
Verification
Simplified from [Chum et al 2007]



Mined from Pascal VOC2011



Pre-Training for R-CNN

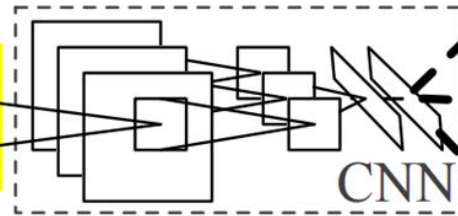


1. Input image



2. Extract region proposals (~2k)

warped region



CNN

3. Compute CNN features

aeroplane? no.

⋮

person? yes.

⋮

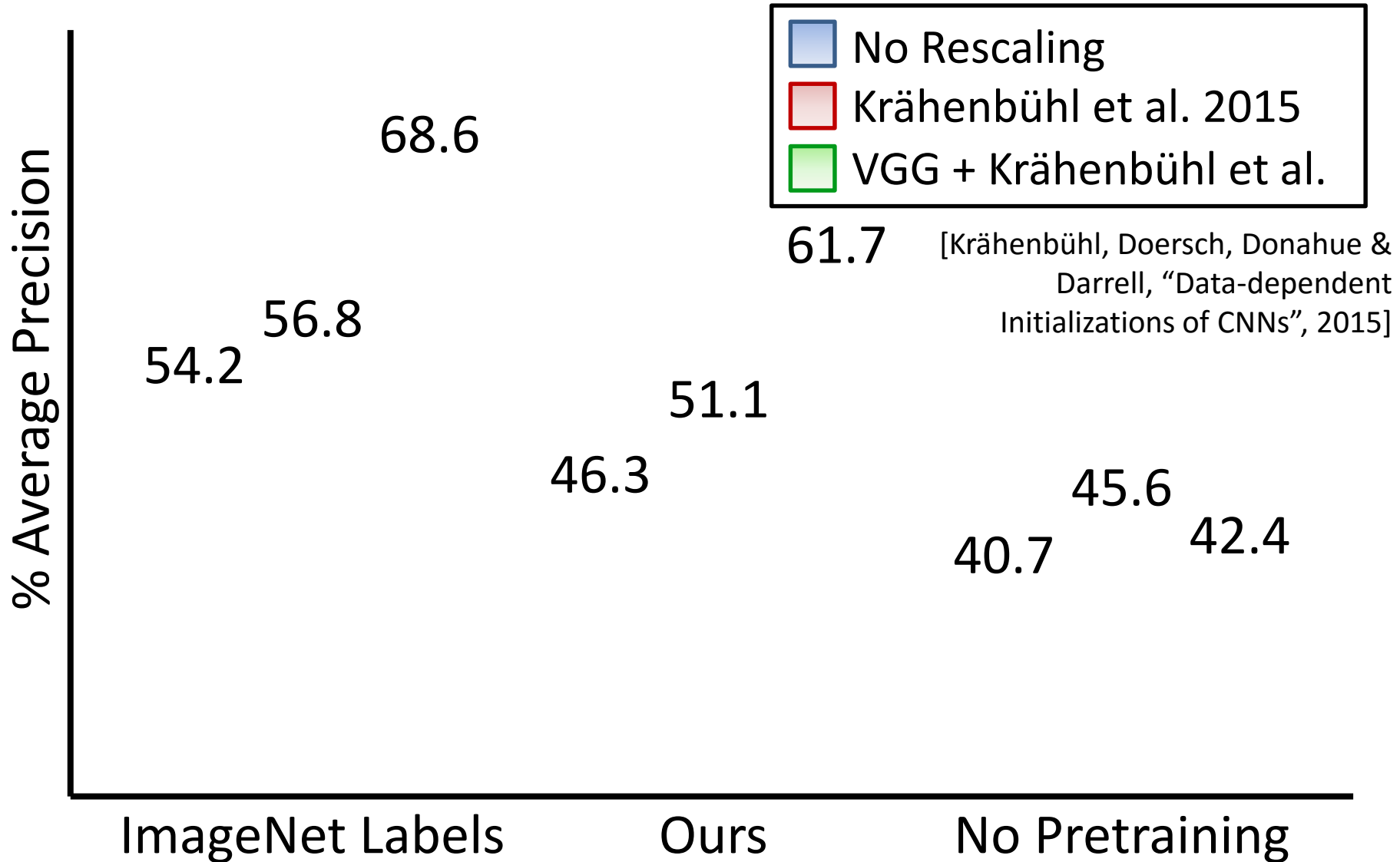
tvmonitor? no.

4. Classify regions

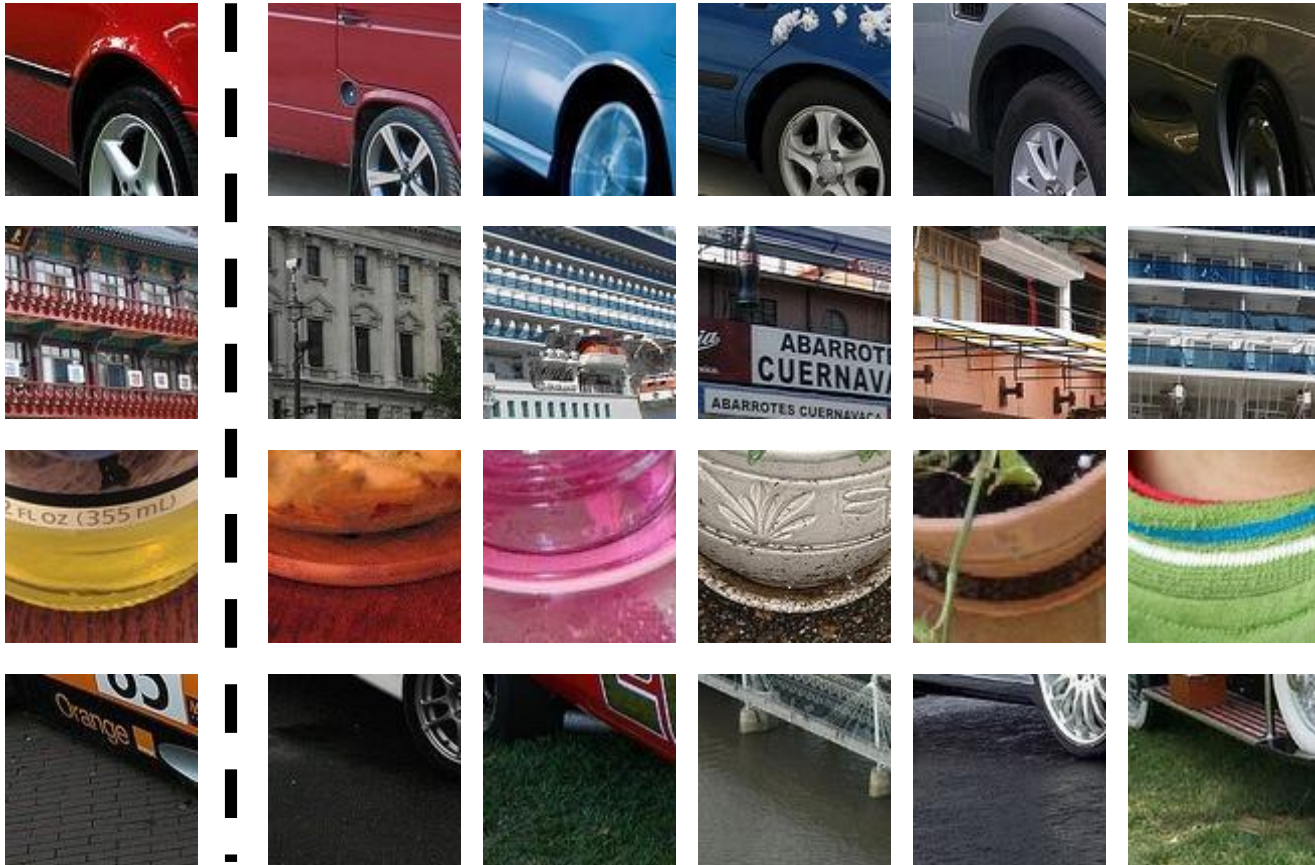
Pre-train on relative-position task, w/o labels

VOC 2007 Performance

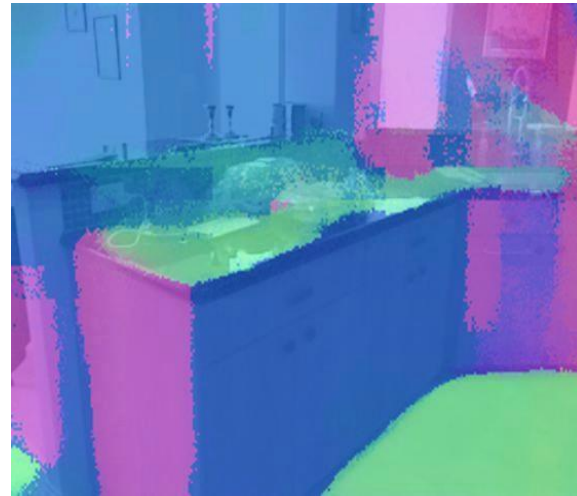
(pretraining for R-CNN)



Capturing Geometry?



Surface-normal Estimation

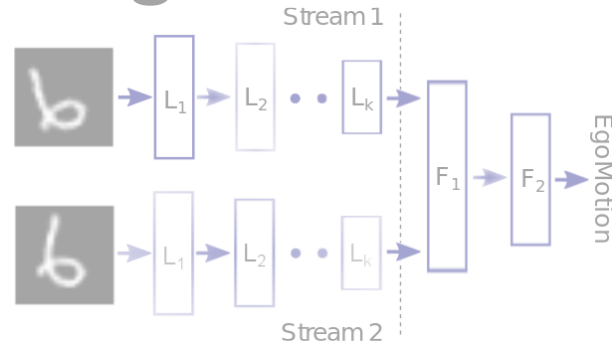


Method	Error (Lower Better)		% Good Pixels (Higher Better)		
	Mean	Median	11.25°	22.5°	30.0°
No Pretraining	38.6	26.5	33.1	46.8	52.5
Ours	33.2	21.3	36.0	51.2	57.8
ImageNet Labels	33.3	20.8	36.7	51.7	58.1

So, do we need semantic labels?

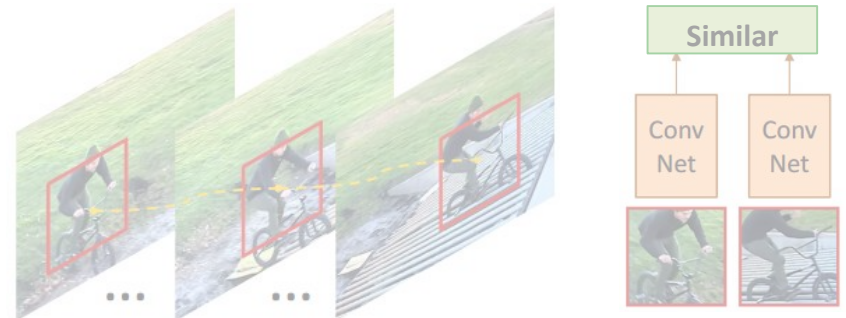
“Self-Supervision” and the Future

Ego-Motion



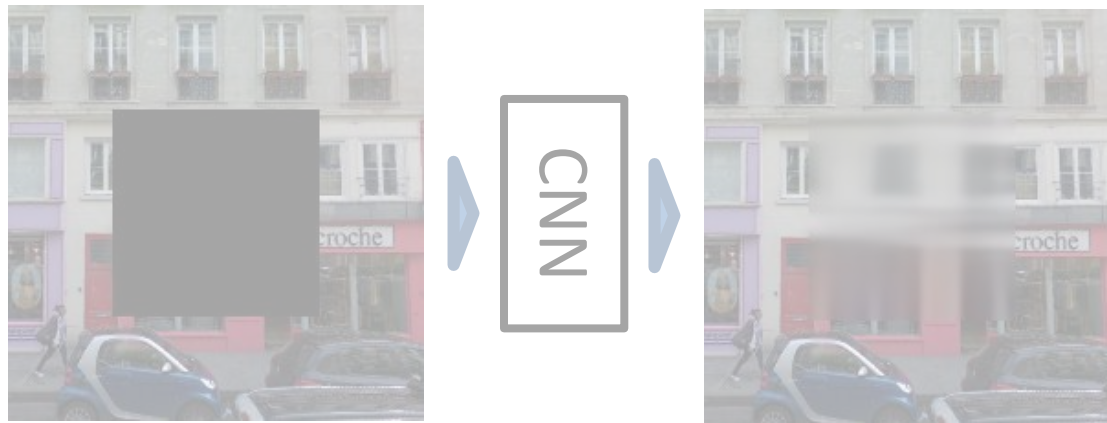
[Agrawal et al. 2015; Jayaraman et al. 2015]

Video



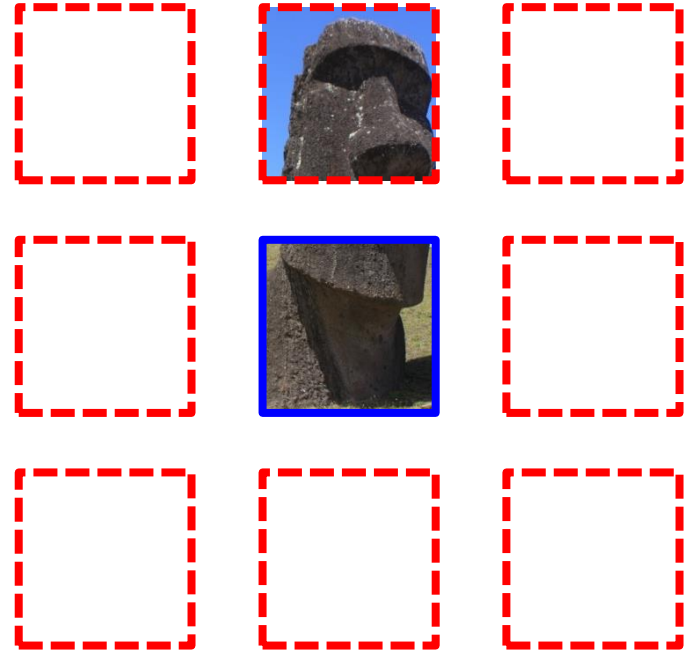
[Wang et al. 2015; Srivastava et al 2015; ...]

Context

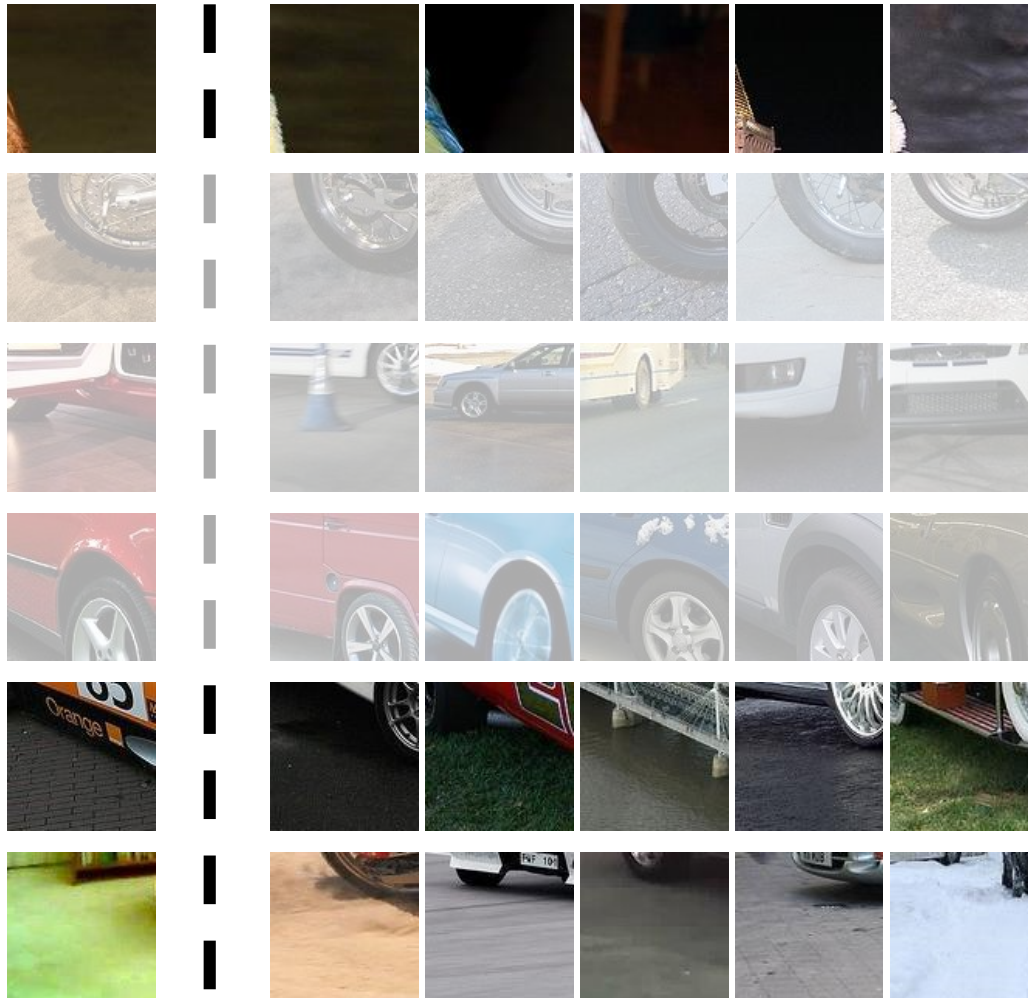


[Doersch et al. 2014; Pathak et al. 2015; Isola et al. 2015]

Thank you!

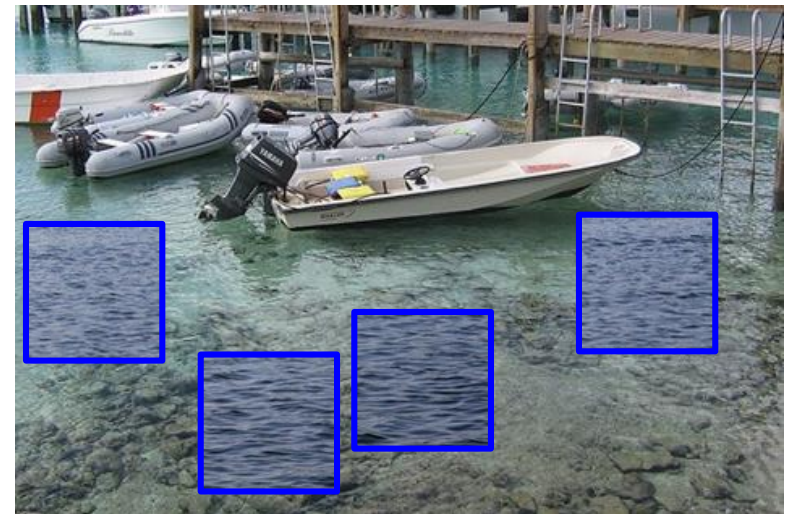
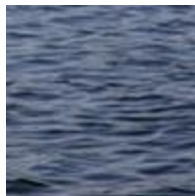
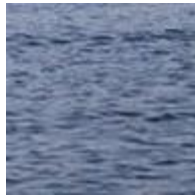
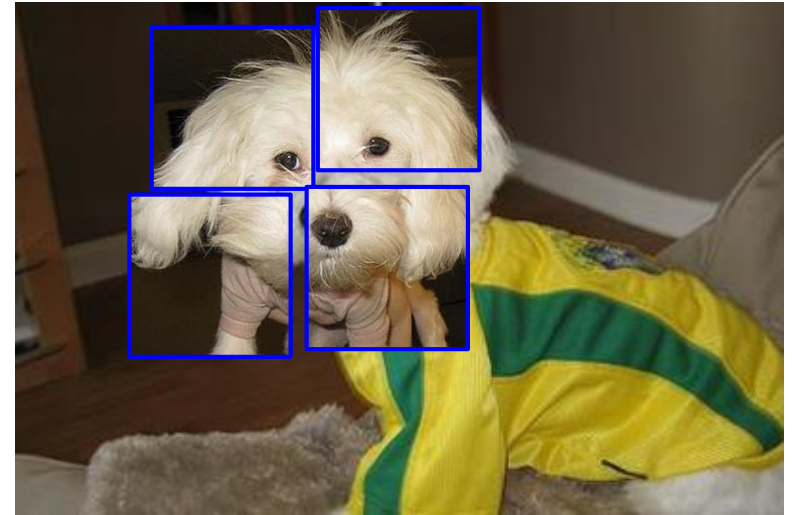


Visual Data Mining?



Geometric Verification

Like [Chum et al. 2007], but simpler

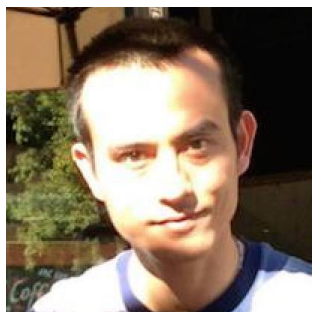


Geometric Verification

Like [Chum et al. 2007], but simpler



Learning Spatiotemporal Features with 3D Convolutional Networks



Du Tran
(1,2)



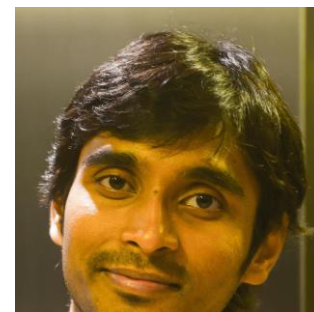
Lubomir Bourdev
(2)



Rob Fergus
(2,3)



Lorenzo Torresani
(1)

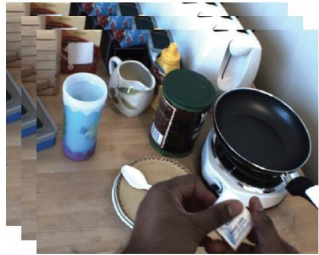


Manohar Paluri
(2)

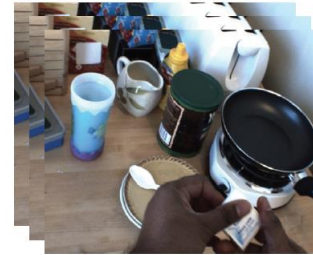
(1) Dartmouth College, (2) Facebook AI Research, (3) New York University

Introduction

Video Understanding problem:



What Objects?



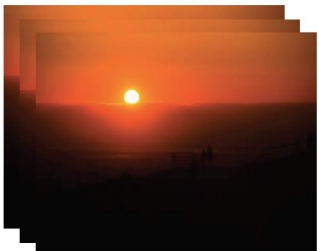
- cup, mug, hands
- indoor
- making coffee



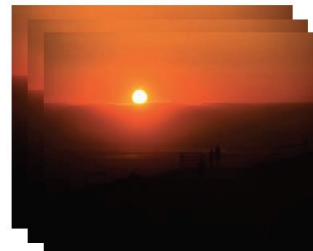
What Scene?



- person, surfboard
- sea
- surfing



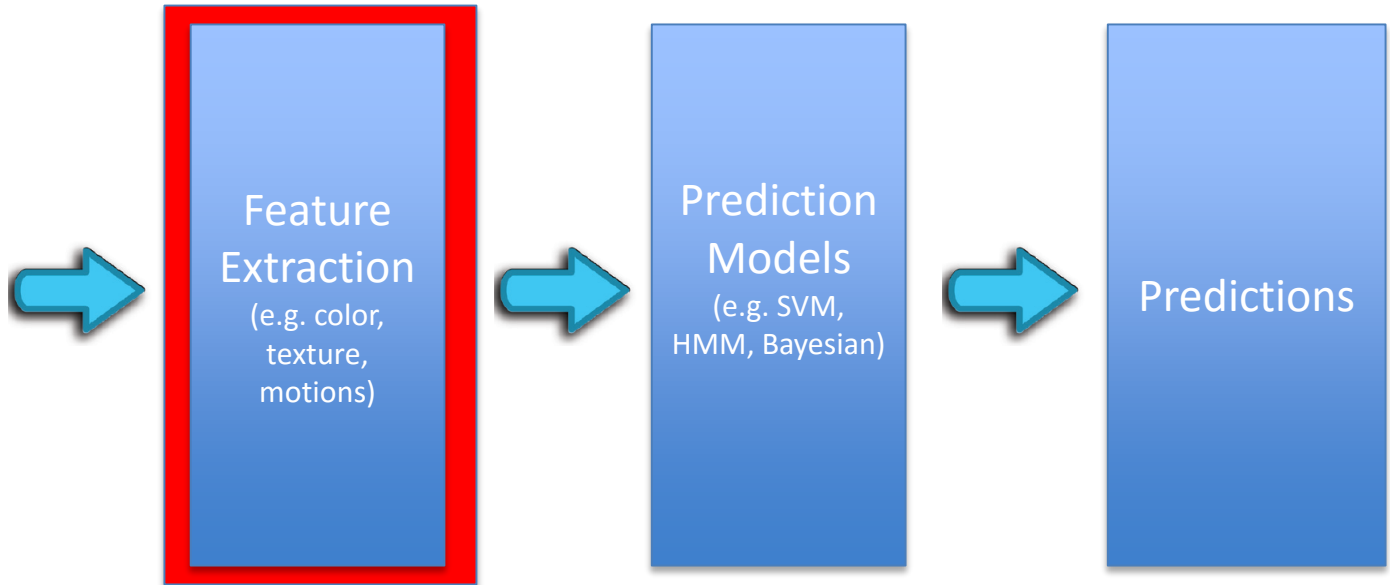
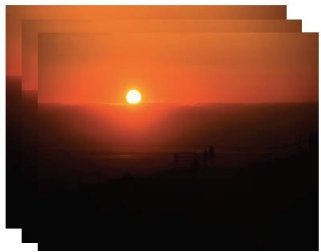
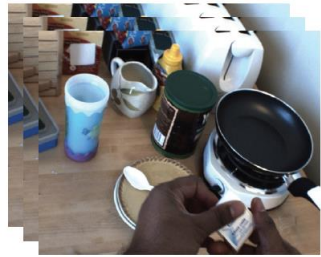
What Actions?



- sun
- mountain
- NA

Many potential applications due to a large-growing number of internet videos

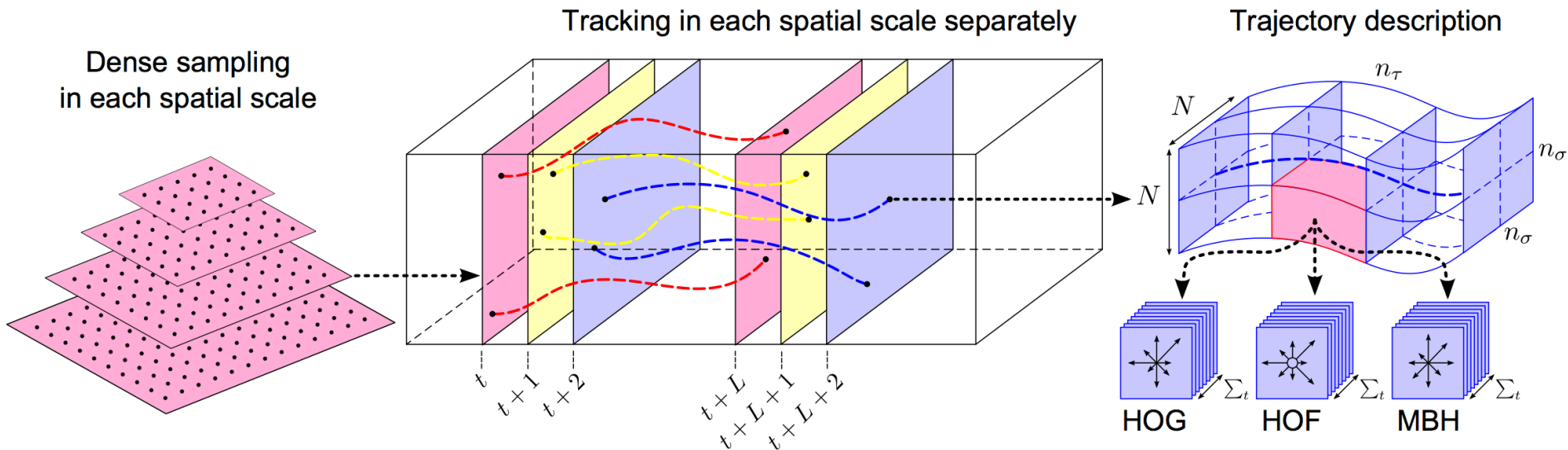
Traditional Computer Vision Pipeline



Focus of this work

Current Best Video Features

- Improved Dense Trajectories (iDT)



Wang et al. IJCV'13

Pros:

- Don't need to learn
- Don't need large-scale training data

Cons:

- Highly hand-crafted
- Computational intensive
- Hard to parallelize

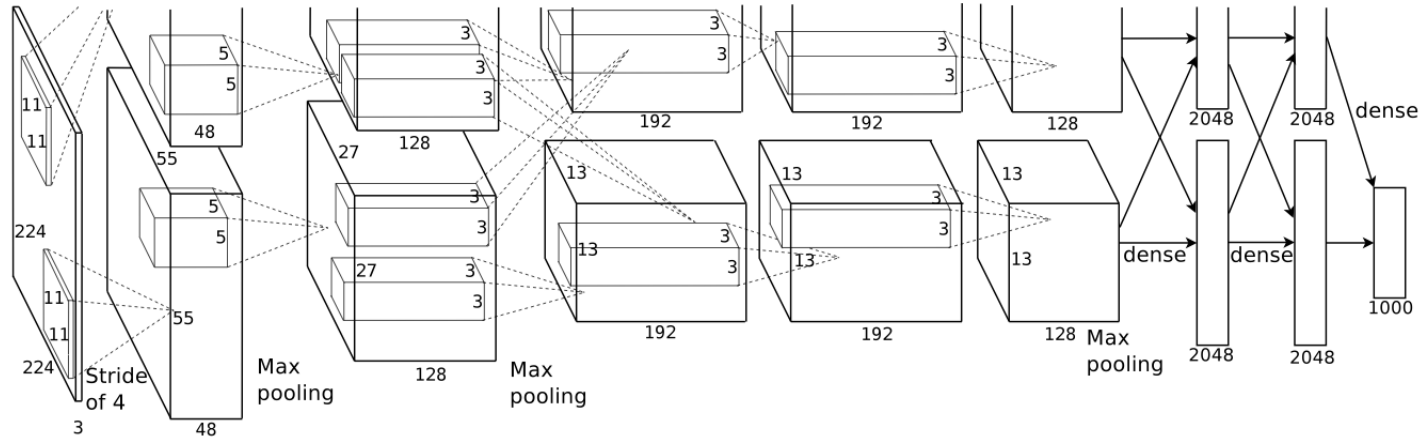
What If We Have Big Data?

- Learn features directly from data (no more human biases)
 - Normally built on deep learning, e.g. deep features
- Does it work?
 - Showed to work well for images [Donahue et al. ICML'14]
- How about videos?

Deep Image-based Features



Russakovsky et al. IJCV'15



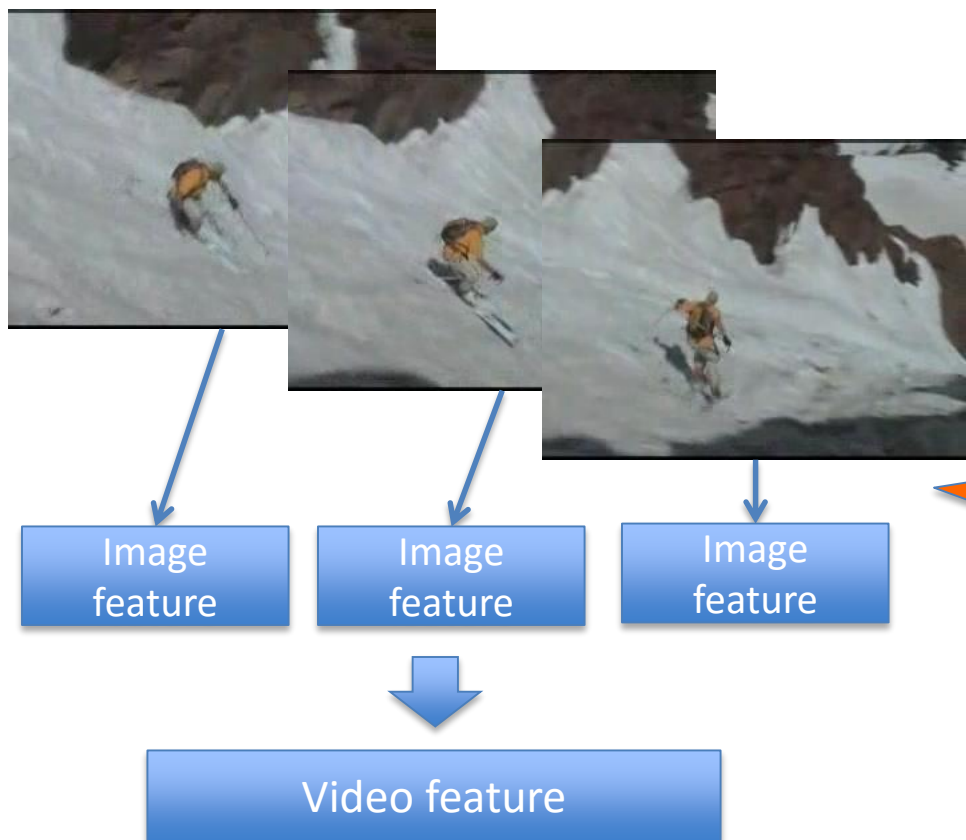
Krizhevsky et al. NIPS'12

The marriage of Big Data and Good Deep Learning Models

- Fully-supervised trained on large-scale dataset.
- Activations are used as features for transferring tasks.

Can Image-based Features Applied to Videos?

- A video is a sequence of images?



Any problems?

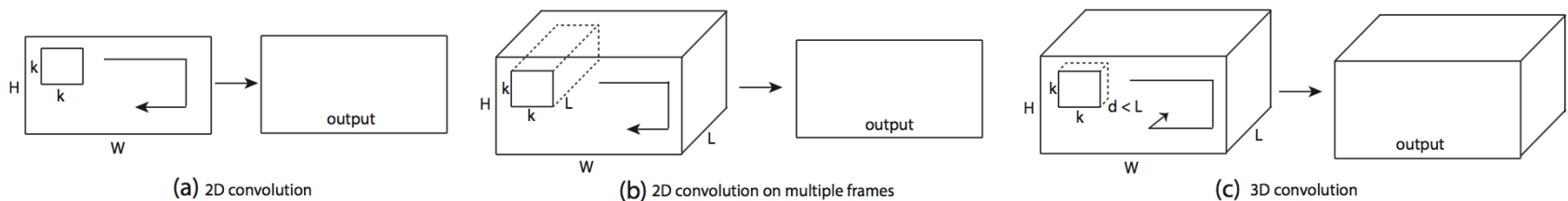
No explicit
motion
modeling

What is right for jointly modeling
appearance & motion?

3D ConvNets instead of 2D
ConvNets

2D ConvNet vs. 3D ConvNet

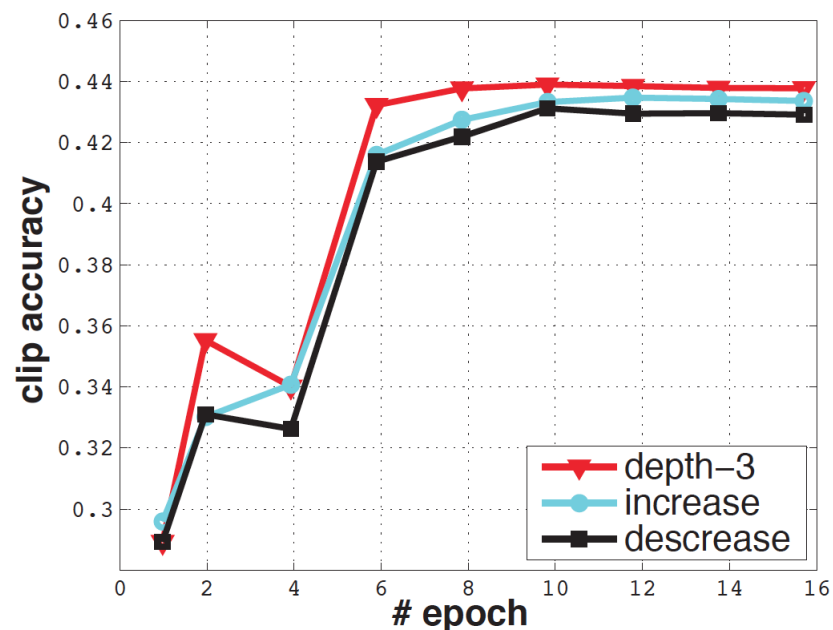
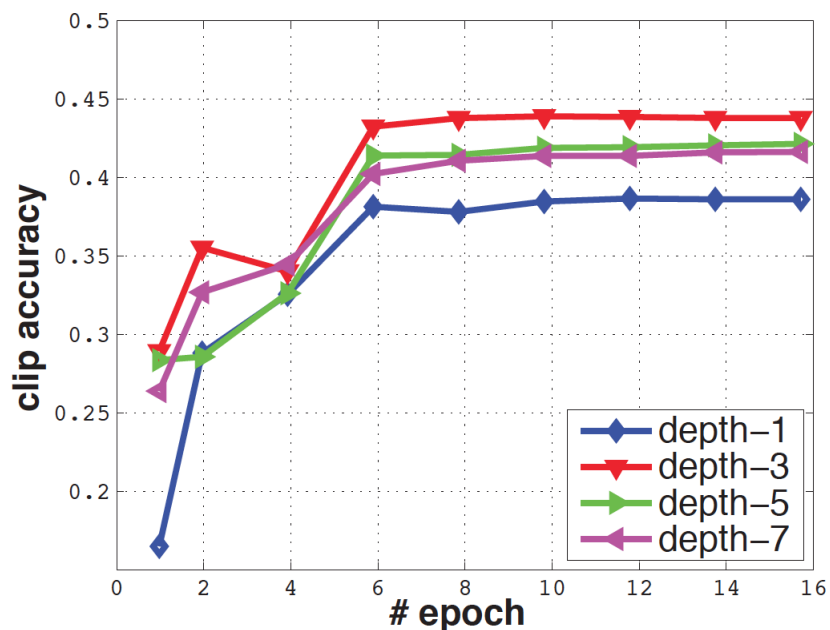
- Basic operations: 2D vs. 3D convolution
- Most of current work
 - Use 2D convolution on images or videos
 - Cannot model temporal information (motions)



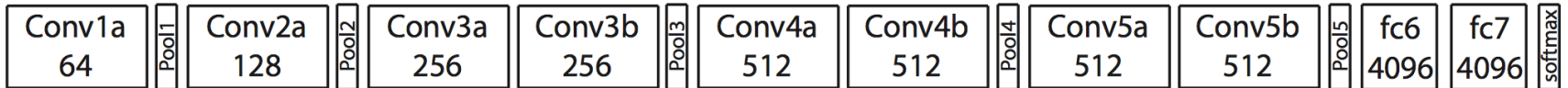
- We propose to use 3D ConvNets for video feature learning

What is a Good Architecture for 3D ConvNets?

- Dataset: UCF101 (13K videos of 101 actions)
- Use similar architecture, varying kernel temporal depth



Learning Video Features with C3D



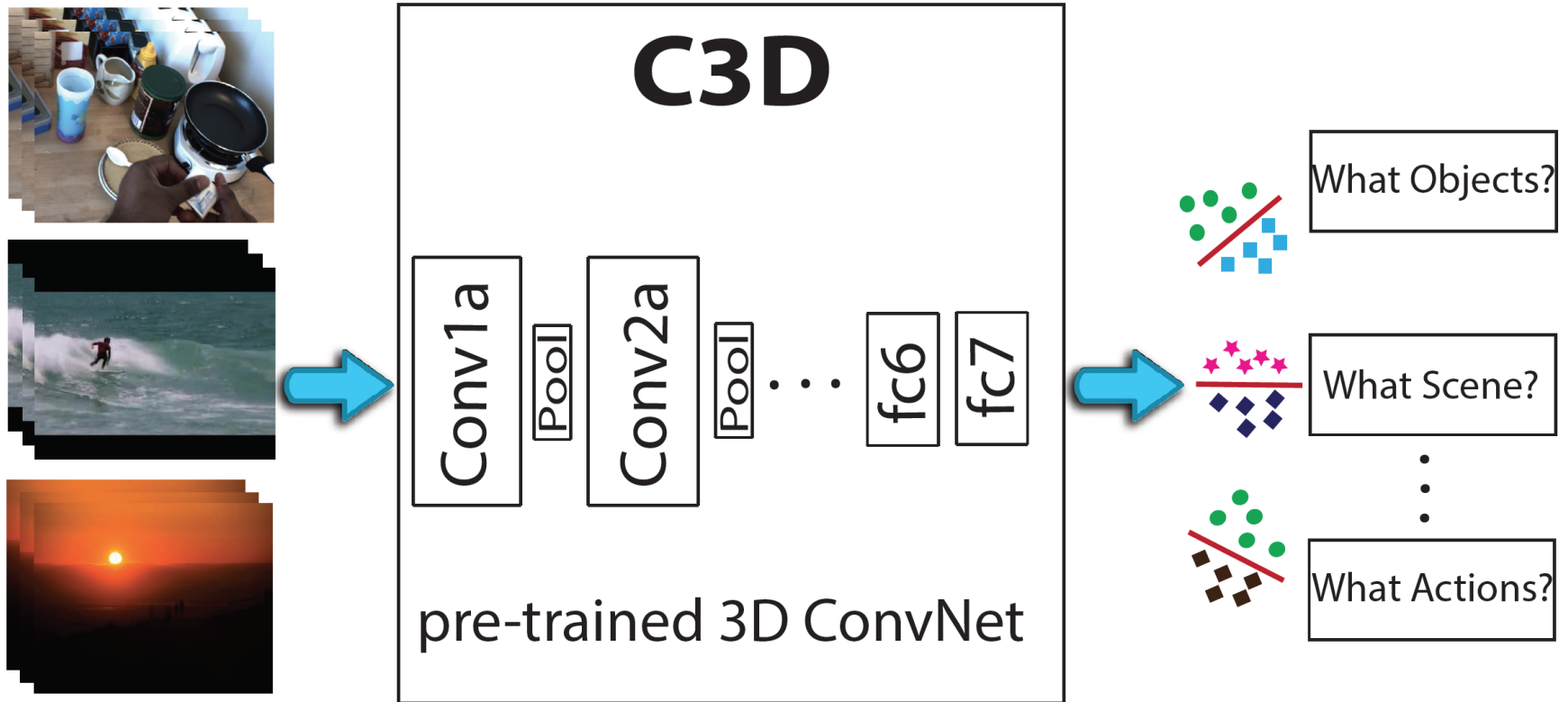
- C3D Architecture
 - 8 convolution, 5 pool, 2 fully-connected layers
 - 3x3x3 convolution kernels
 - 2x2x2 pooling kernels
- Dataset: Sports-1M [Karpathy et al. CVPR'14]
 - 1.1M videos of 487 different sport categories
 - Train/test splits are provided

Sport Classification Results



Method	Number of Nets	Clip hit@1	Video hit@1	Video hit@5
Deep Video's Single-Frame + Multires [19]	3 nets	42.4	60.0	78.5
Deep Video's Slow Fusion [19]	1 net	41.9	60.9	80.2
C3D (trained from scratch)	1 net	44.9	60.0	84.4
C3D (fine-tuned from I380K pre-trained model)	1 net	46.1	61.1	85.2

C3D as Generic Features



Simple recipe: C3D + linear SVM = good performance

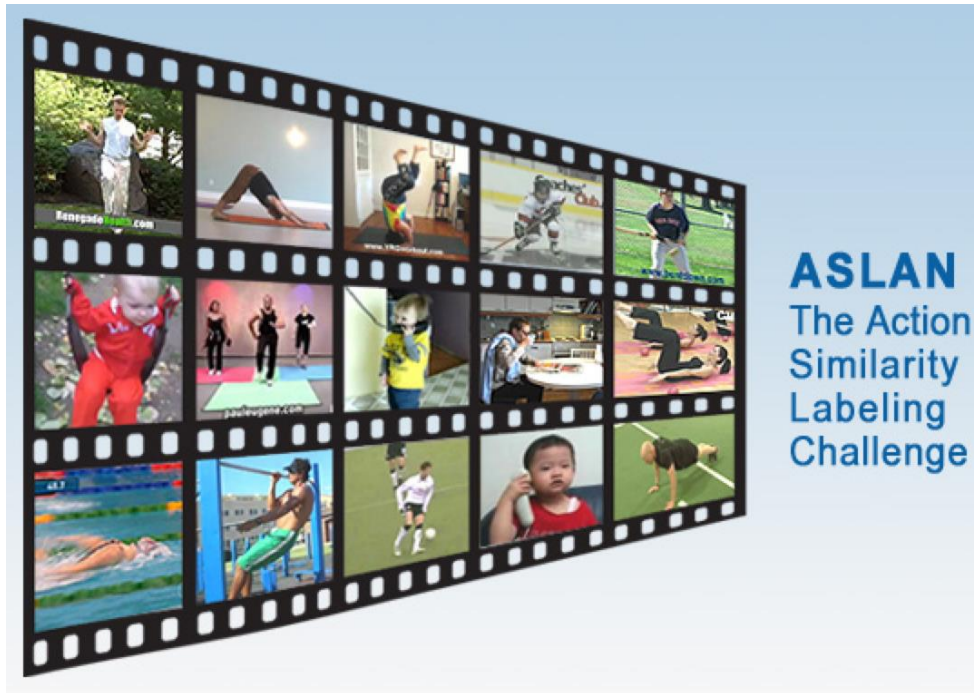
Action Recognition



Action Recognition Results

	Method	Accuracy (%)
Baselines	Imagenet	68.8
	iDT	76.2
Use raw pixel inputs	Deep networks [19]	65.4
	Spatial stream network [36]	72.6
	LRCN [7]	71.1
	LSTM composite model [39]	75.8
	C3D (1 net)	82.3
	C3D (3 nets)	85.2
Use optical flows	iDT with Fisher vector [31]	87.9
	Temporal stream network [36]	83.7
	Two-stream networks [36]	88.0
	LRCN [7]	82.9
	LSTM composite model [39]	84.3
	Multi-skip feature stacking [26]	89.1
	C3D (3 nets) + iDT	90.4

Action Similarity Labeling

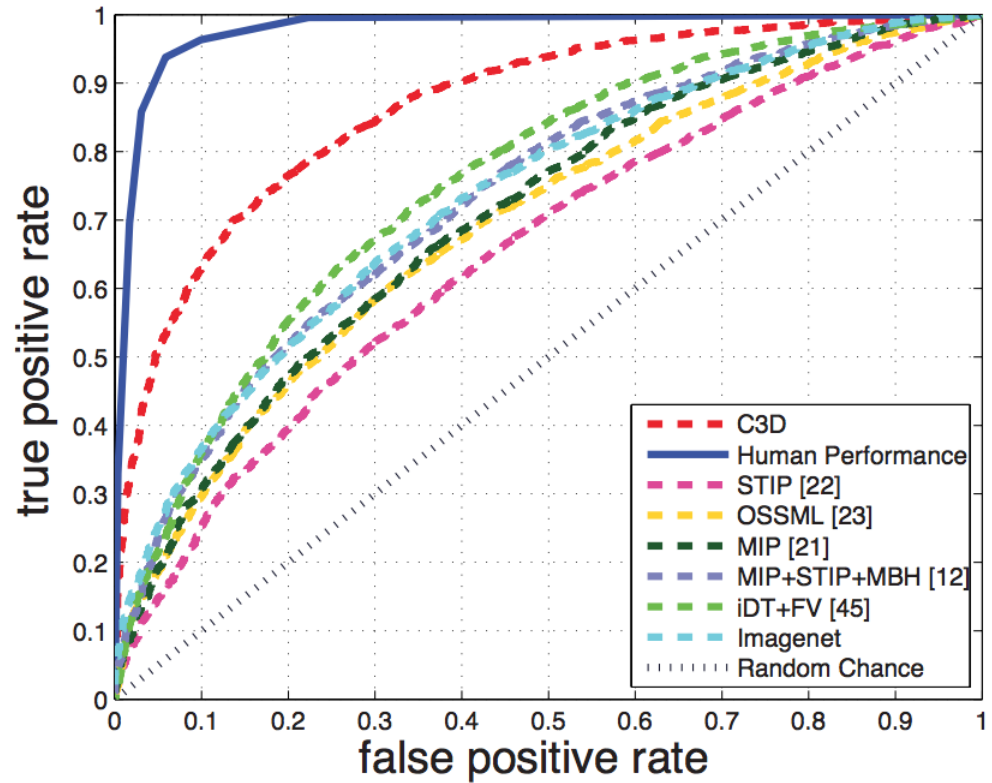


TASK: Given a pair of clips, predict same or different actions

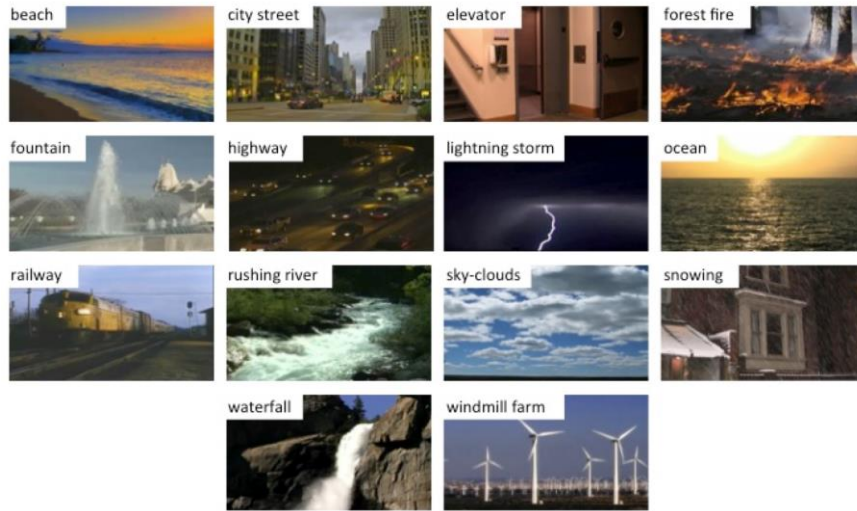
Very challenging evaluation setting: train and test on different categories of actions

ASLAN Results

Method	Features	Model	Acc.	AUC
[22]	STIP	linear	60.9	65.3
[23]	STIP	metric	64.3	69.1
[21]	MIP	metric	65.5	71.9
[12]	MIP+STIP+MBH	metric	66.1	73.2
[45]	iDT+FV	metric	68.7	75.4
Baseline	Imagenet	linear	67.5	73.8
Ours	C3D	linear	78.3	86.5



Dynamic Scene Classification



YUPENN



Maryland

Dataset	[5]	[41]	[9]	[10]	Imagenet	C3D
Maryland	43.1	74.6	67.7	77.7	87.7	87.7
YUPENN	80.7	85.0	86.0	96.2	96.7	98.1

Object Classification



Egocentric object dataset

Dataset Task	Object object recognition
Method	[31]
Result	12.0
C3D	22.3
Δ	10.3

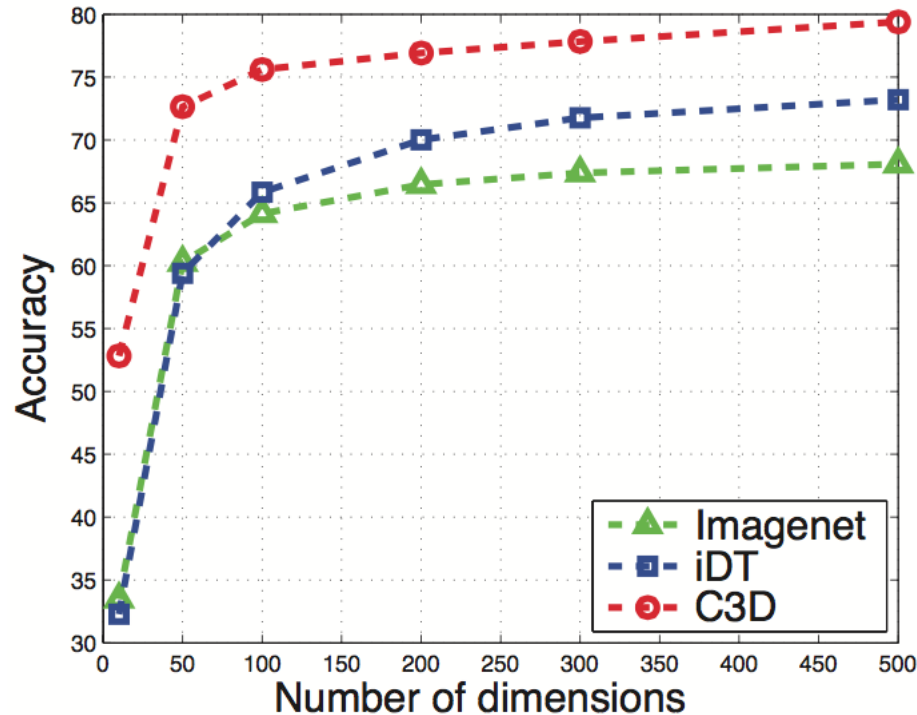
Result Summary

C3D performance compared with current methods

Dataset Task	Sport1M action recognition	UCF101 action recognition	ASLAN action similarity labeling	YUPENN scene classification	UMD scene classification	Object object recognition
Method	[19]	[39]([26])	[31]	[10]	[10]	[32]
Result	80.2	75.8 (89.1)	68.7	96.2	77.7	12.0
C3D	85.2	85.2 (90.4)	78.3	98.1	87.7	22.3
Δ	5.0	9.4 (1.3)	9.6	1.9	10.0	10.3

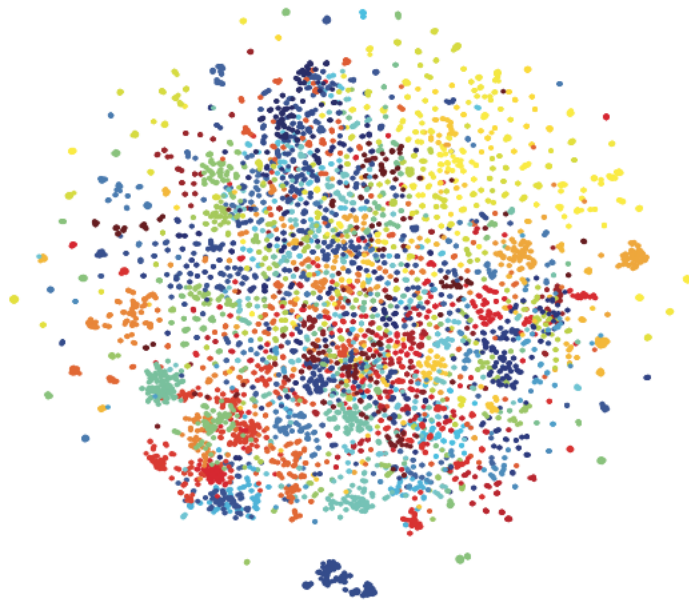
**Consistently outperforms state-of-the-art methods
on 4 different tasks and 6 different datasets**

C3D is Compact

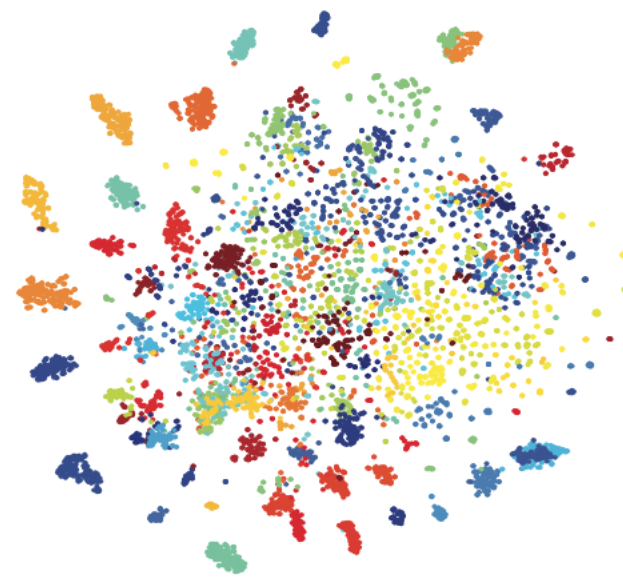


- 10-20% better than Imagenet and iDT at low dimension
- Obtains 52.8% using only 10-dim (random chance is less than 0.96%)

Qualitative Comparison



Imagenet



C3D

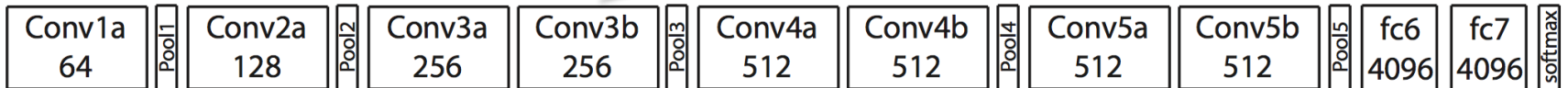
C3D is Efficient

- Extract features on full UCF101
- 91x faster than iDT
- 276x faster than optical-flow-based methods

Why does C3D works so well?

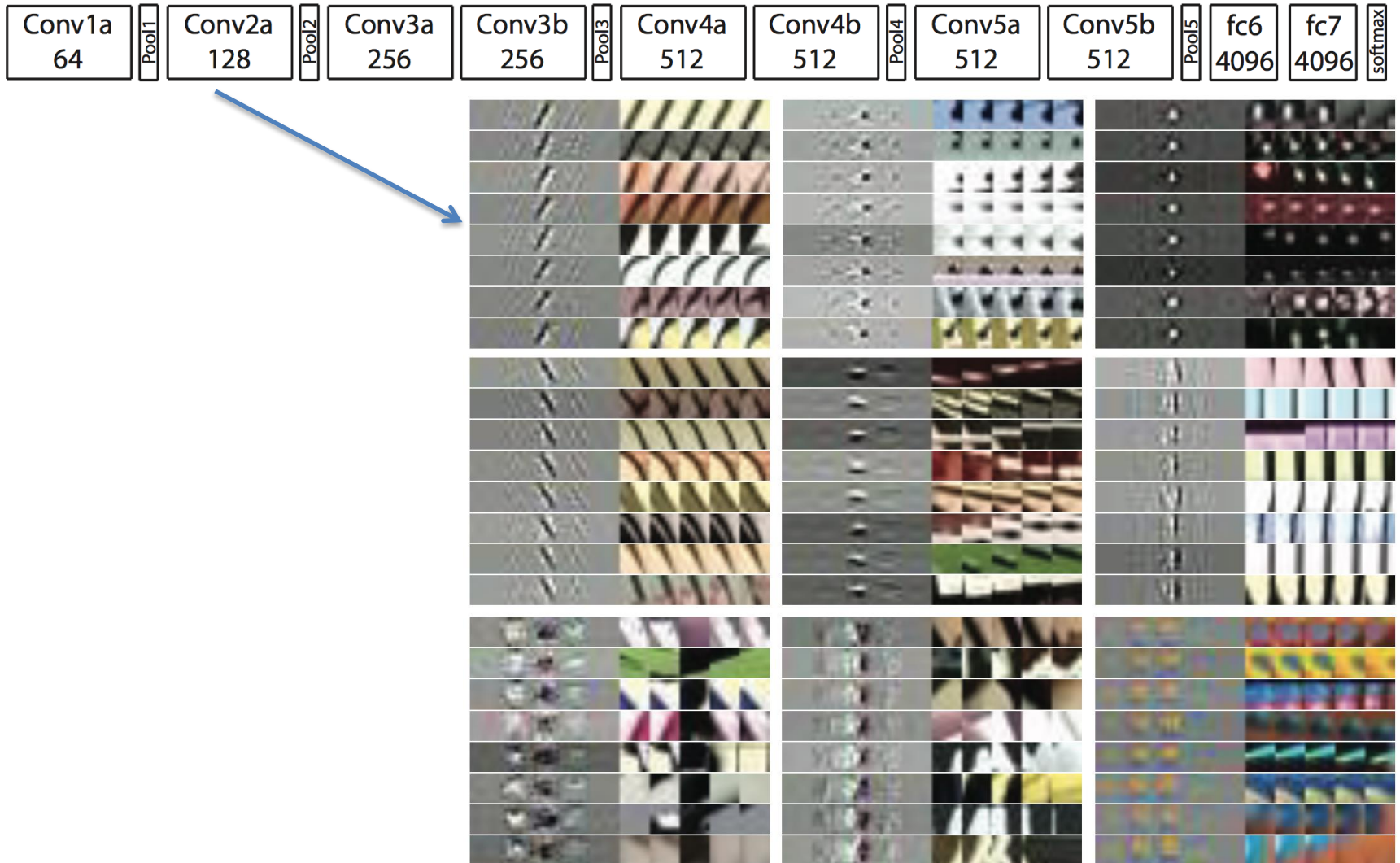
- What does C3D learn at internal layers?

What has actually
been learned
here?

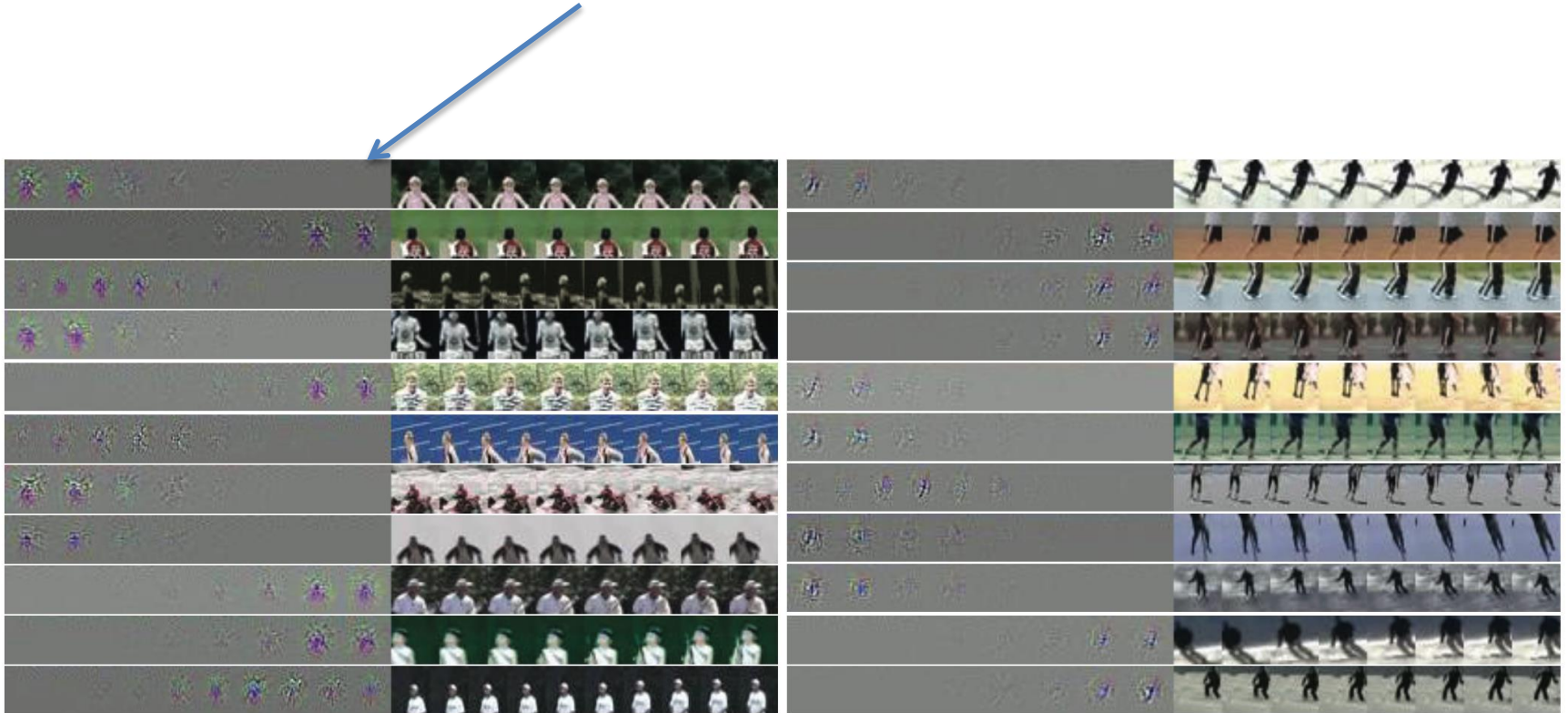
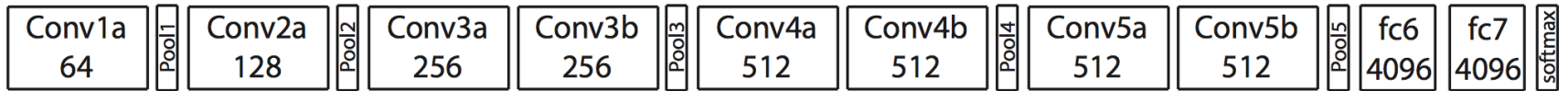


- Use Deconvolution method [Zeiler & Fergus ECCV'14] to visualize C3D learned features of some internal layers.

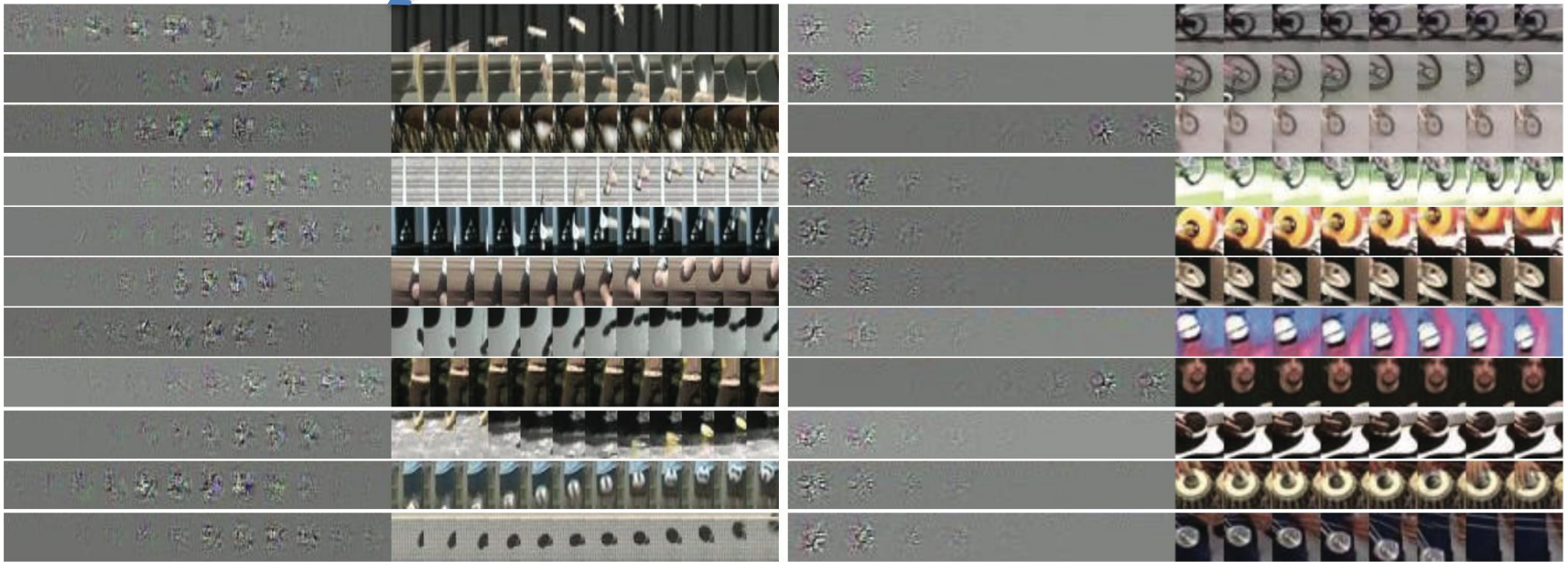
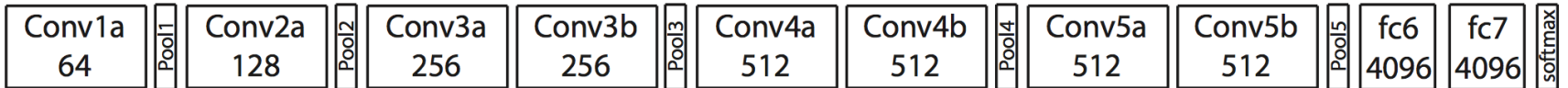
Deconvolutions of conv2a



Deconvolutions of conv3b



Deconvolutions of conv3b



Conclusions

- 3D ConvNet is well-suited for spatiotemporal feature learning.
- C3D is a good architecture for 3D ConvNet
- C3D is a good generic video features
 - Accurate
 - Compact
 - Efficient to compute
 - Easy to use

Source code & models are available at <http://vlg.cs.dartmouth.edu/c3d>

Thank you

- Q&A
- Demo



Words, Pictures, and Common Sense

a.k.a.

Learning by Playing

Devi Parikh

Virginia Tech

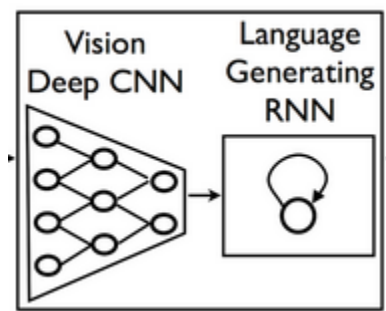


what i think

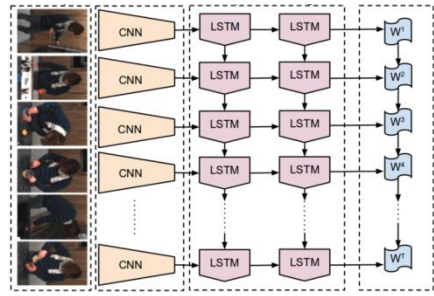


what i say

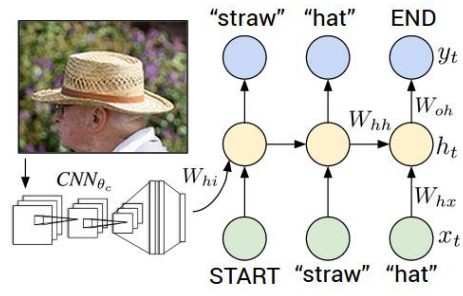
Image captioning is receiving a lot of attention



Vinyals et al., 2015



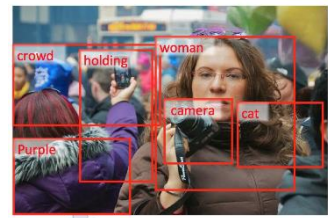
Donahue et al., 2015



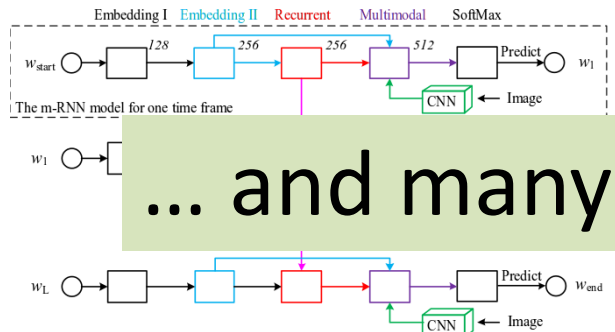
Karpathy and Fei-Fei, 2015



Hodosh et al., 2013

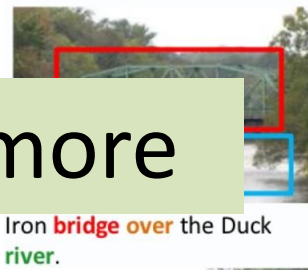


Fang et al., 2015

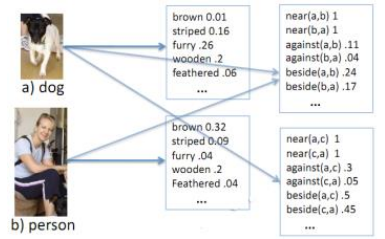


... and many more

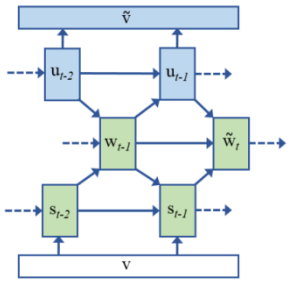
Mao et al., 2015



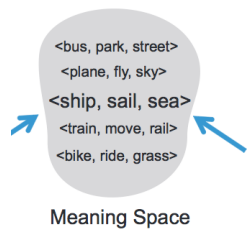
Ordonez et al., 2011



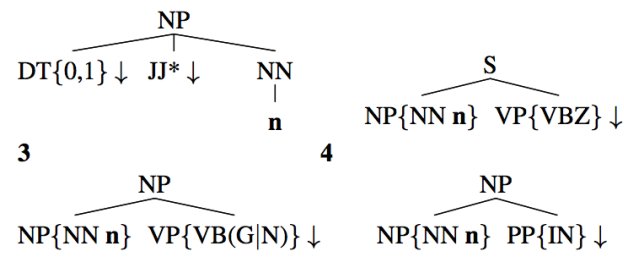
Kulkarni et al., 2011



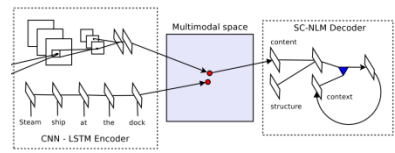
Chen and Zitnick, 2015



Farhadi et al., 2010



Mitchell et al., 2012



Kiros et al., 2015



"man in blue wetsuit is surfing on
wave."

Karpathy and Fei-Fei, CVPR 2015



A group of young people playing a game of frisbee.

Vinyas et al., CVPR 2015



LZ

a car is parked
in the middle
of nowhere .

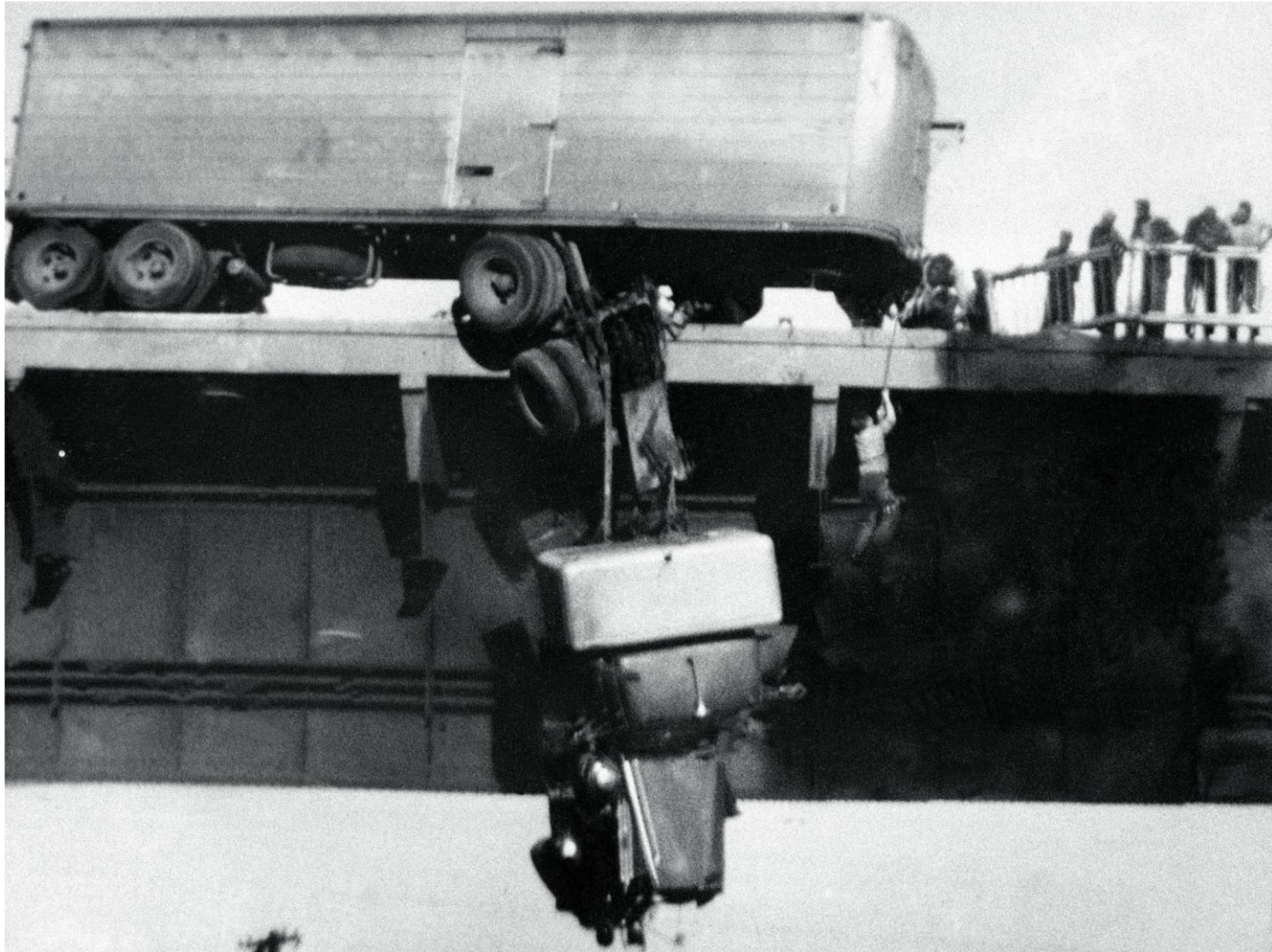
Kiros et al., TACL 2015



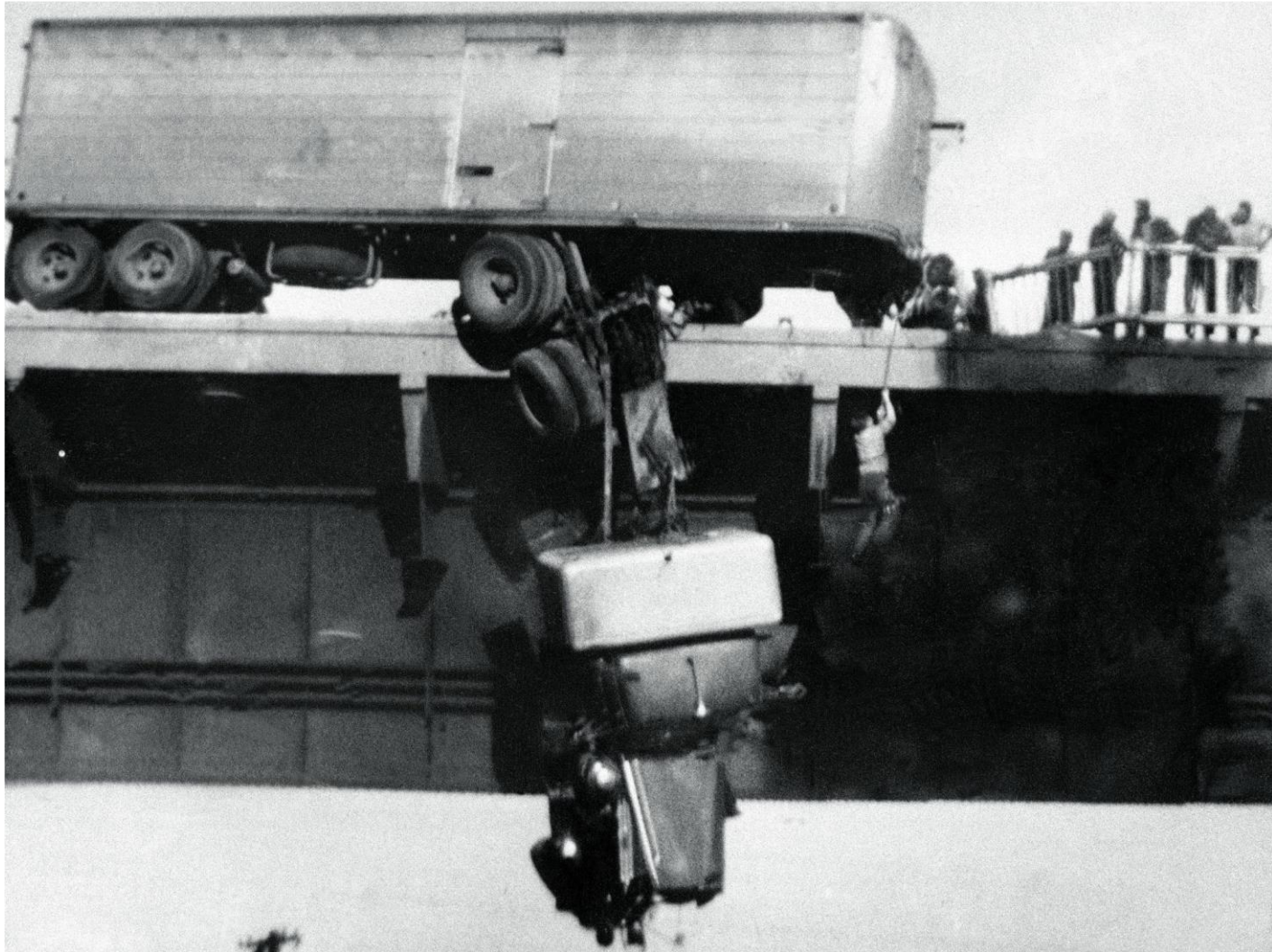
a pot of broccoli on a stove

Fang et al. CVPR 2015

A man is rescued from his truck that is hanging dangerously from a bridge.



A man is *rescued* from his truck that is hanging *dangerously* from a bridge.



Learning Common Sense

- Text
 - Reporting bias

Reporting bias in text

<i>Word</i>	<i>Teraword</i>	<i>Knext</i>	<i>Word</i>	<i>Teraword</i>	<i>Knext</i>
spoke	11,577,917	244,458	hugged	610,040	10,378
laughed	3,904,519	169,347	blinked	390,692	20,624
murdered	2,843,529	11,284	was late	368,922	31,168
inhaled	984,613	4,412	exhaled	168,985	3,490
breathed	725,034	34,912	was punctual	5,045	511

[Gordon et al. 2013]

Reporting bias in text

<i>Word</i>	<i>Teraword</i>	<i>Knext</i>	<i>Word</i>	<i>Teraword</i>	<i>Knext</i>
spoke	11,577,917	244,458	hugged	610,040	10,378
laughed	3,904,519	11,284	was late	368,922	31,168
murdered	2,843,529	11,284	was late	368,922	31,168
inhaled	984,613	4,412	exhaled	168,985	3,490
breathed	725,034	34,912	was punctual	5,045	511

inhale:exhale = 6:1

[Gordon et al. 2013]

Reporting bias in text

<i>Word</i>	<i>Teraword</i>	<i>Knext</i>	<i>Word</i>	<i>Teraword</i>	<i>Knext</i>
spoke	11,577,917	244,458	hugged	610,040	10,378
laughed	3,904,519	169,347	blinked	390,692	20,624
murdered	2,843,529	11,284	was late	368,922	31,168
inhaled	984,613	4,412	exhaled	168,985	3,490
breathed	725,034	34,912	was punctual	5,045	511

murder:exhale = 17:1

[Gordon et al. 2013]

Reporting bias in text

<i>Body Part</i>	<i>Teraword</i>	<i>Knext</i>	<i>Body Part</i>	<i>Teraword</i>	<i>Knext</i>
Head	18,907,427	1,332,154	Liver	246,937	10,474
Eye(s)	18,455,030	1,090,640	Kidney(s)	183,973	5,014
Arm(s)	6,345,039	458,018	Spleen	47,216	1,414
Ear(s)	3,543,711	230,367	Pancreas	24,230	1,140
Brain	3,277,326	260,863	Gallbladder	17,419	1,556

[Gordon et al. 2013]

Reporting bias in text

<i>Body Part</i>	<i>Teraword</i>	<i>Knext</i>	<i>Body Part</i>	<i>Teraword</i>	<i>Knext</i>
Head	18,907,427	1,332,154	Liver	246,937	10,474
Eye(s)	18,455,030	1,090,640	Kidney(s)	183,973	5,014
Arm(s)	14,455,532	452,015	Pituitary	47,216	1,414
Ear(s)	3,543,711	230,367	Pancreas	24,230	1,140
Brain	3,277,326	260,863	Gallbladder	17,419	1,556

People have heads:gallbladders = 1085:1

[Gordon et al. 2013]

Do birds fly?

bird

/bɜːd/

noun

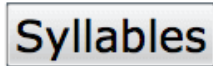
1. any warm-blooded egg-laying vertebrate of the class *Aves*, characterized by a body covering of feathers and forelimbs modified as wings. Birds vary in size between the ostrich and the hummingbird *related adjectives* avian ornithic

Do birds fly?

penguin

[**peng**-gwin, **pen**-]

 Spell

 Syllables

[Examples](#)

[Word Origin](#)

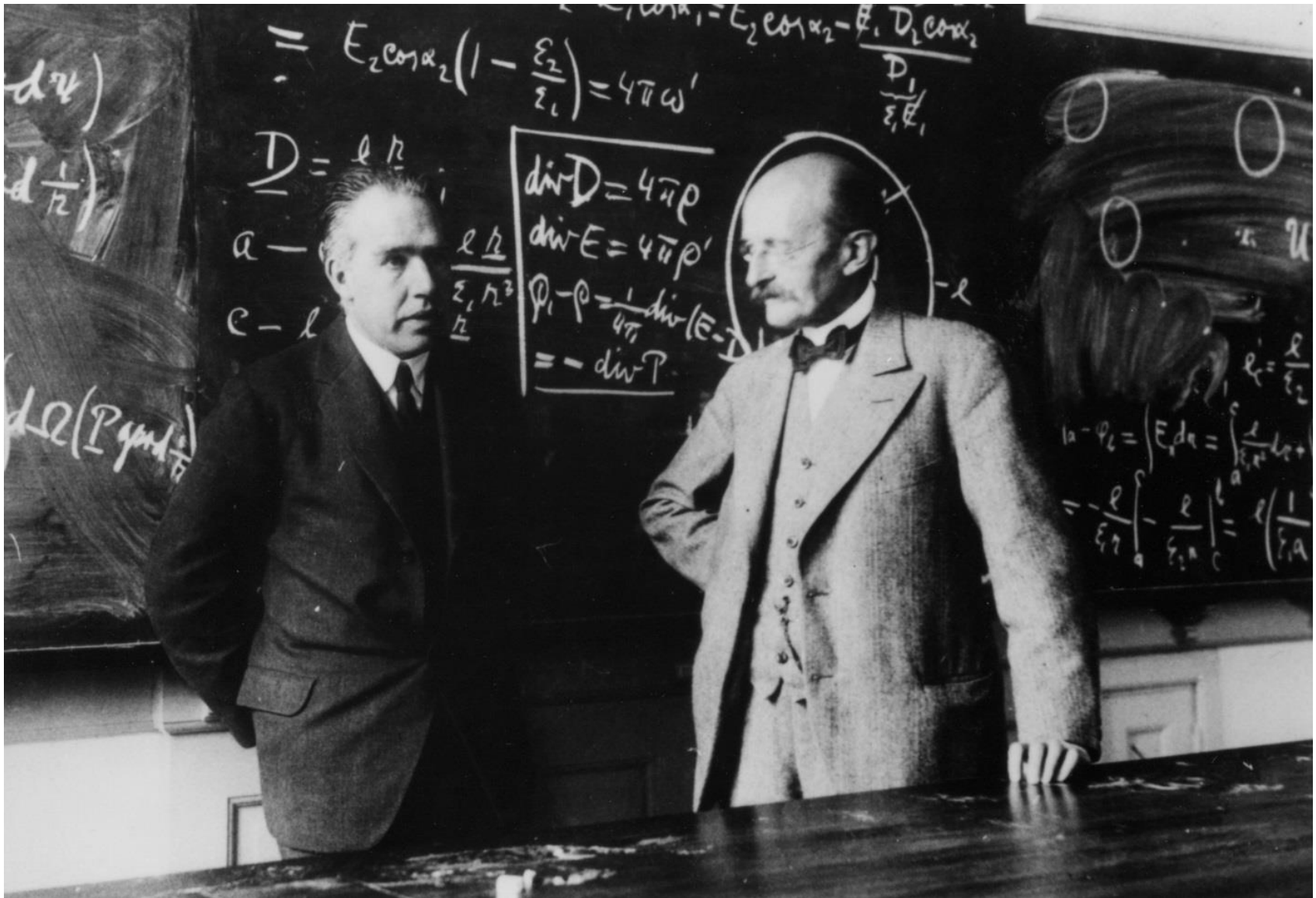
noun, *Ornithology*

1. any of several flightless, aquatic birds of the family Spheniscidae, of the Southern Hemisphere, having webbed feet and wings reduced to flippers.

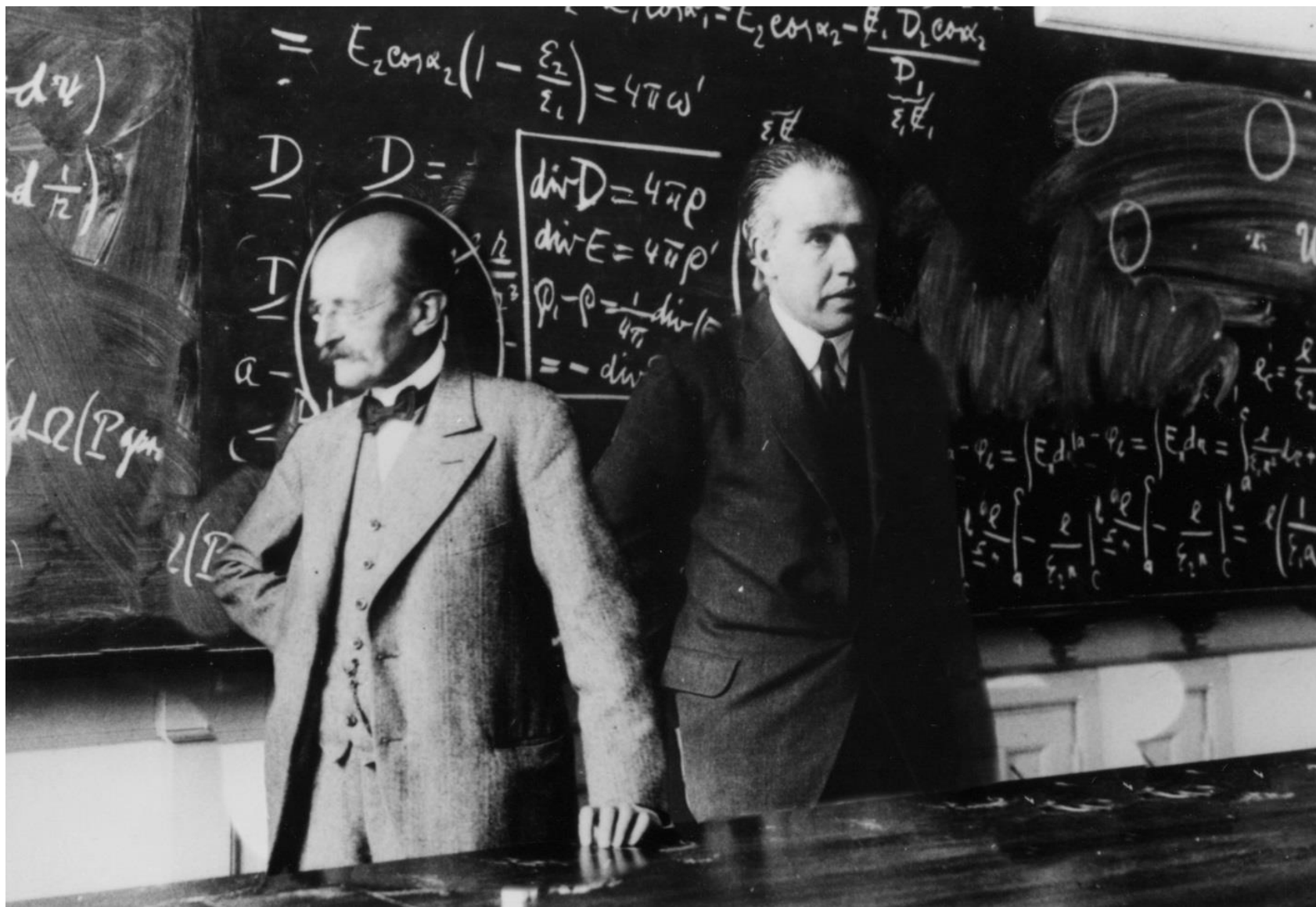
Learning Common Sense

- Text
 - Reporting bias
- From structure in our visual world?

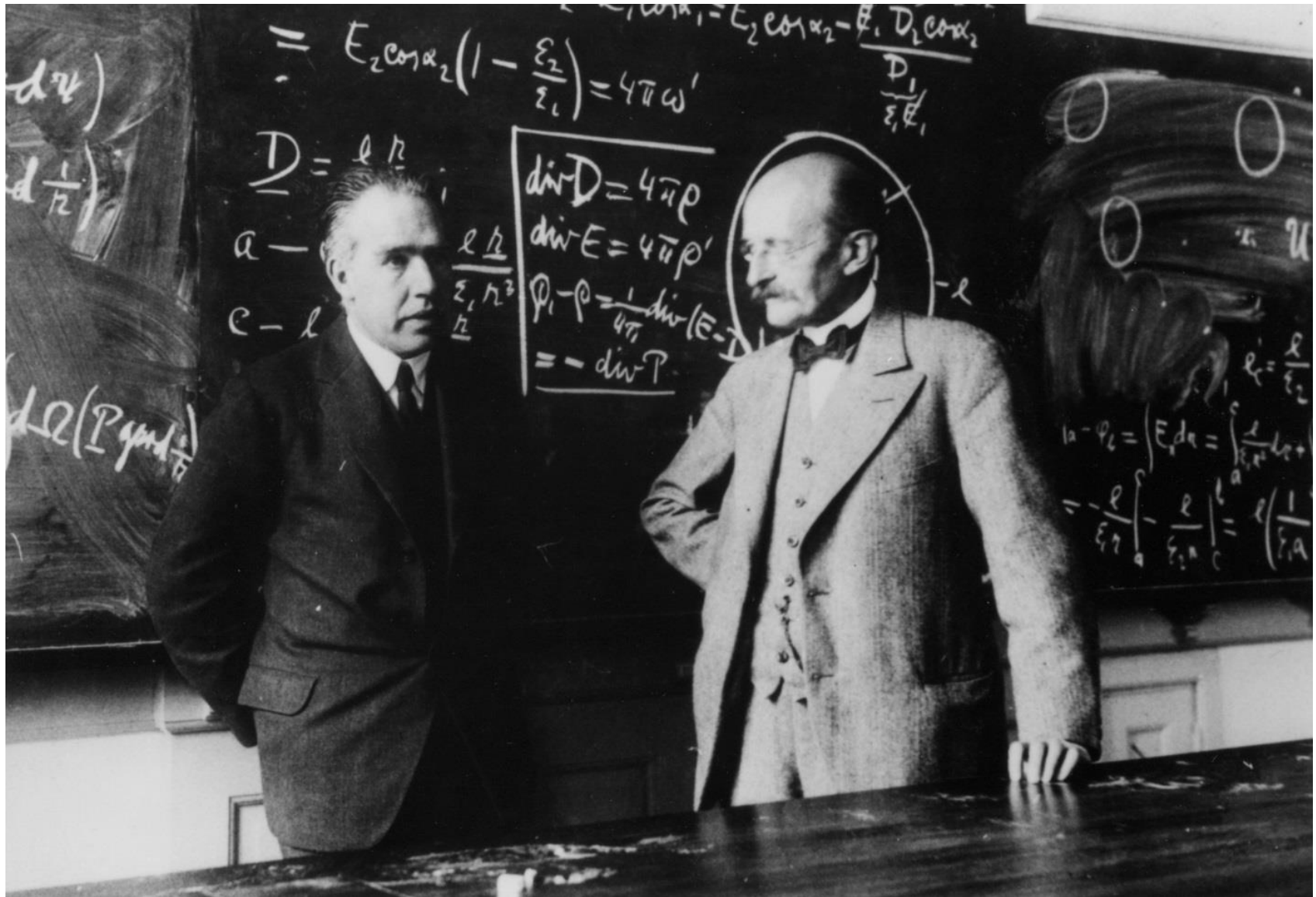
Two professors converse in front of a blackboard.



Two professors stand in front of a blackboard.

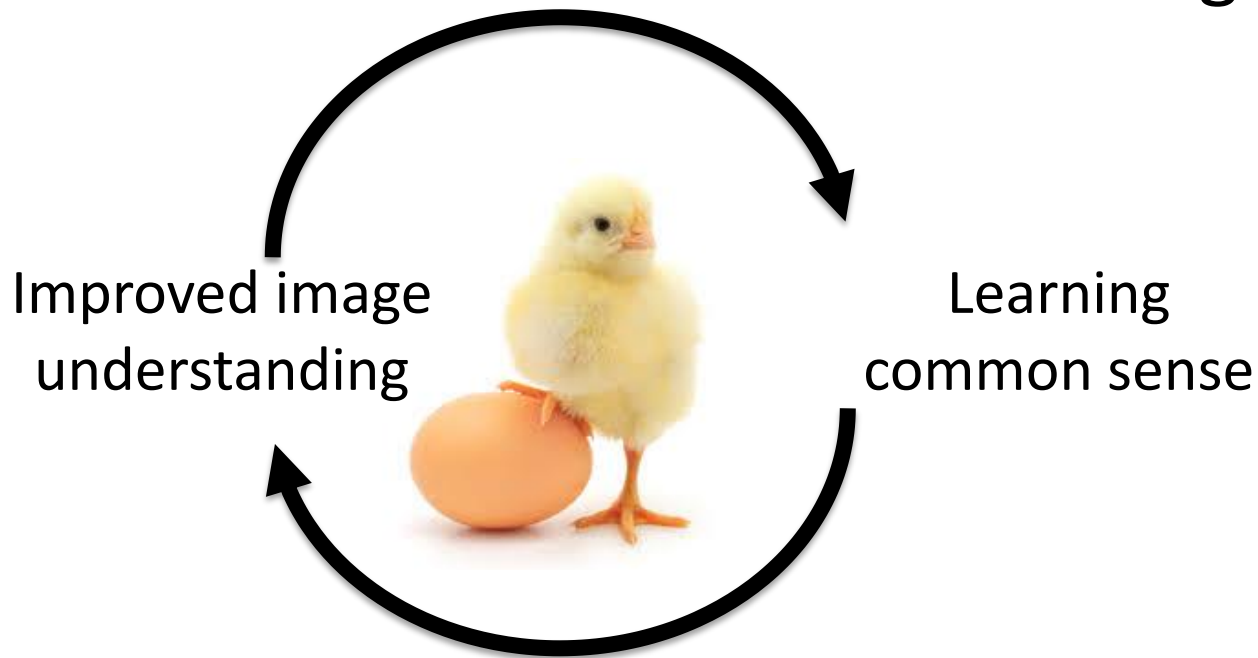


Two professors converse in front of a blackboard.



Challenges

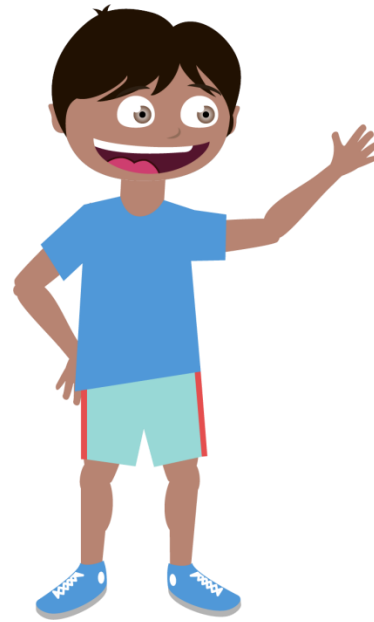
- Lacking visual density
- Annotations are expensive
- Computer vision doesn't work well enough



Is photorealism necessary?

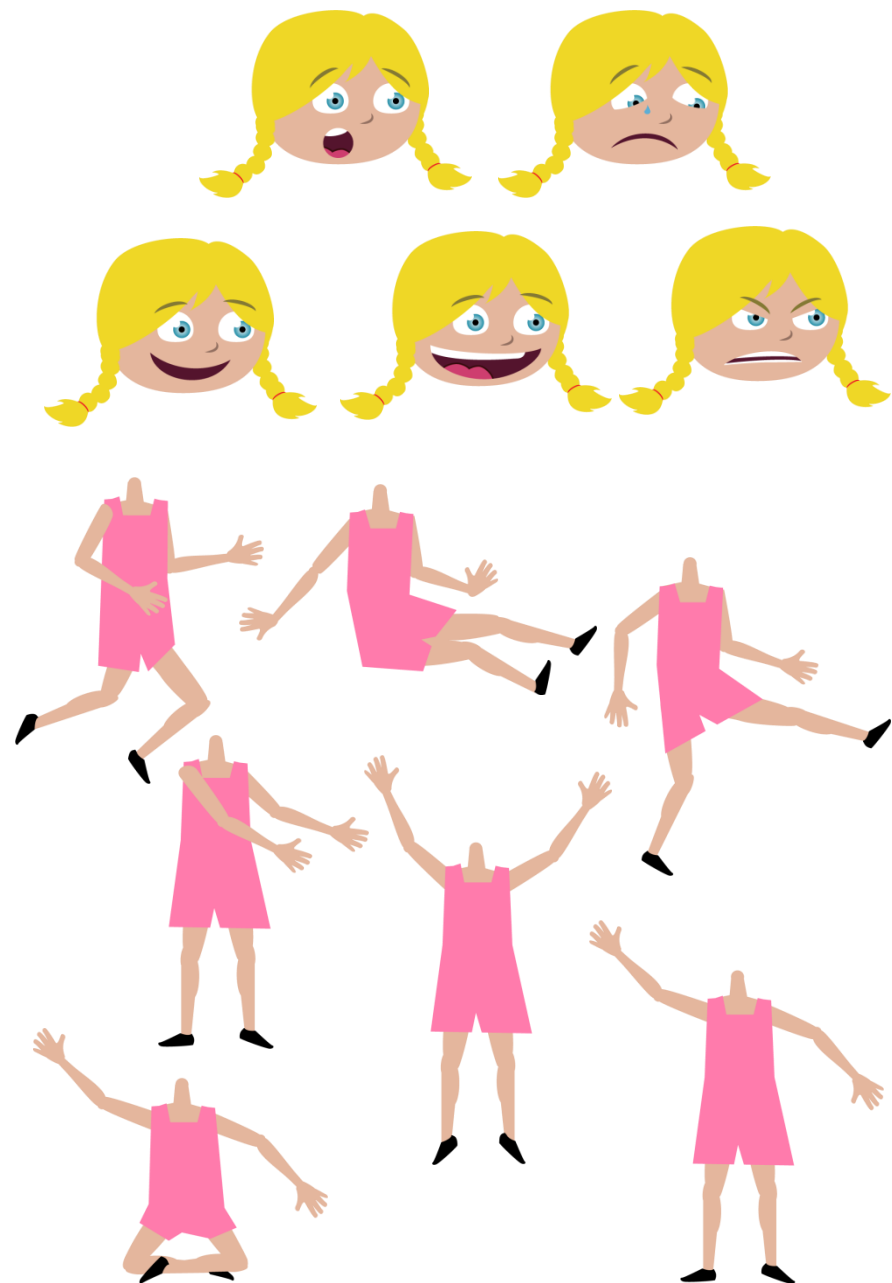
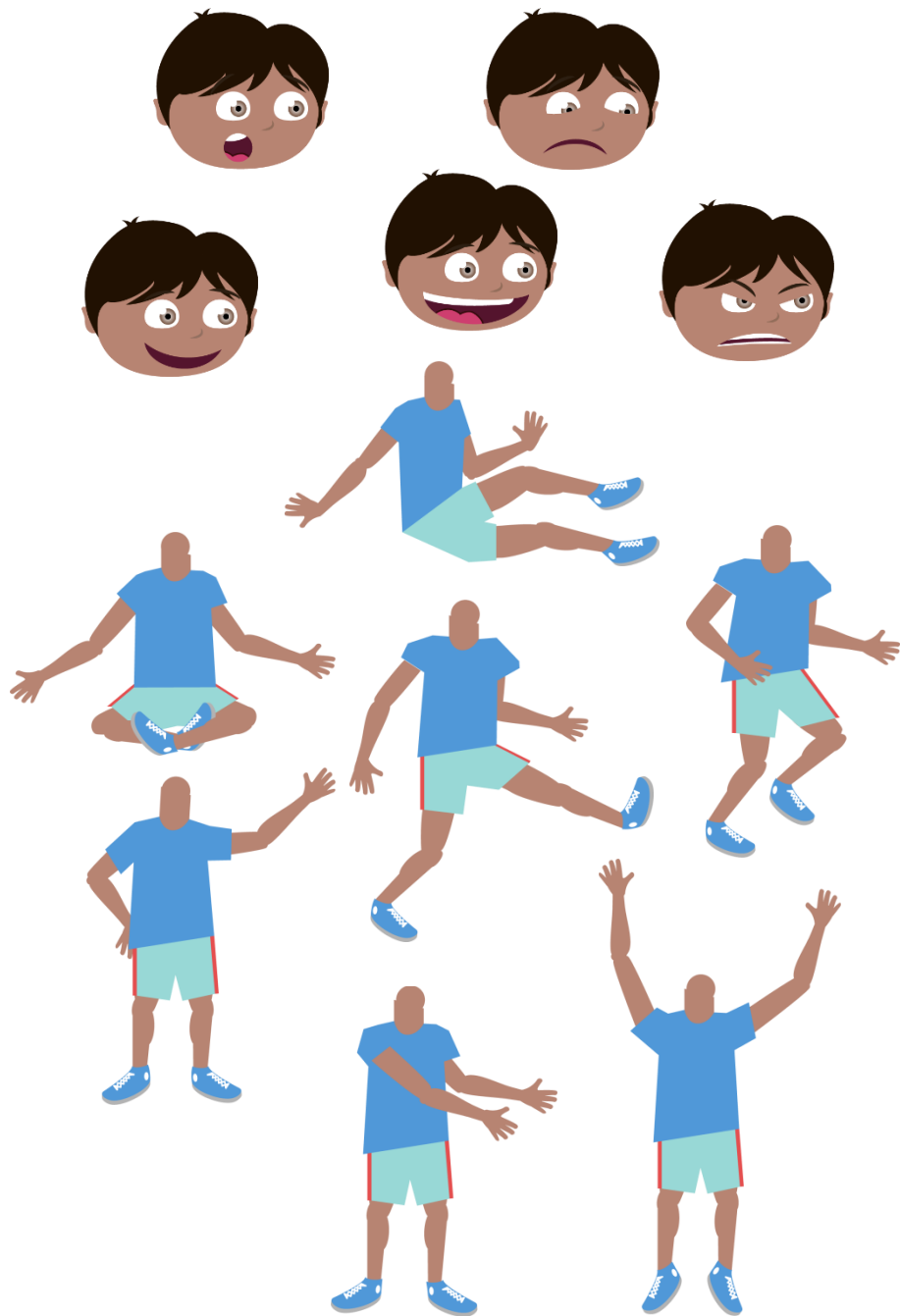


Jenny



Mike









Create a children's illustration!

Please help us create an illustration for a children's story book by creating an outdoor scene from the clipart below. Use your imagination! Clipart may be added by dragging the clipart onto the scene, and removed by dragging it off. The clipart may be resized or flipped, and each clipart may only be added once. Please use at least 6 pieces of clipart in each scene. You will be asked to complete 3 different scenes. Press "Next" when finished with the current scene and "Done" when all are finished. Thanks!

Scene: 1/3

Size  **Flip** 

Clipart

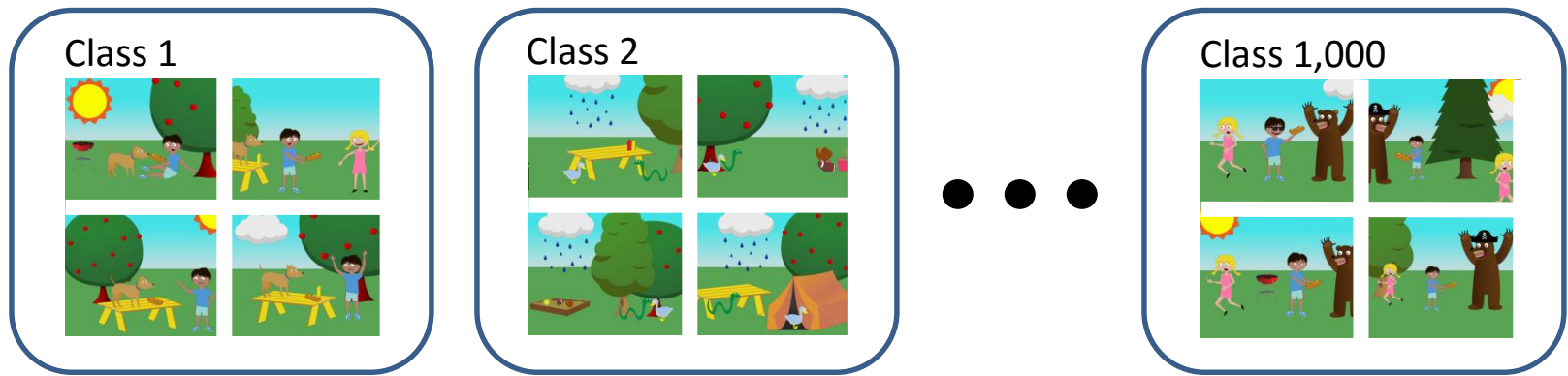


Mike fights off a bear by giving him a hotdog while Jenny runs away.



Dataset

1,000 classes of semantically similar scenes:



1,000 classes x 10 scenes per class = 10,000 scenes

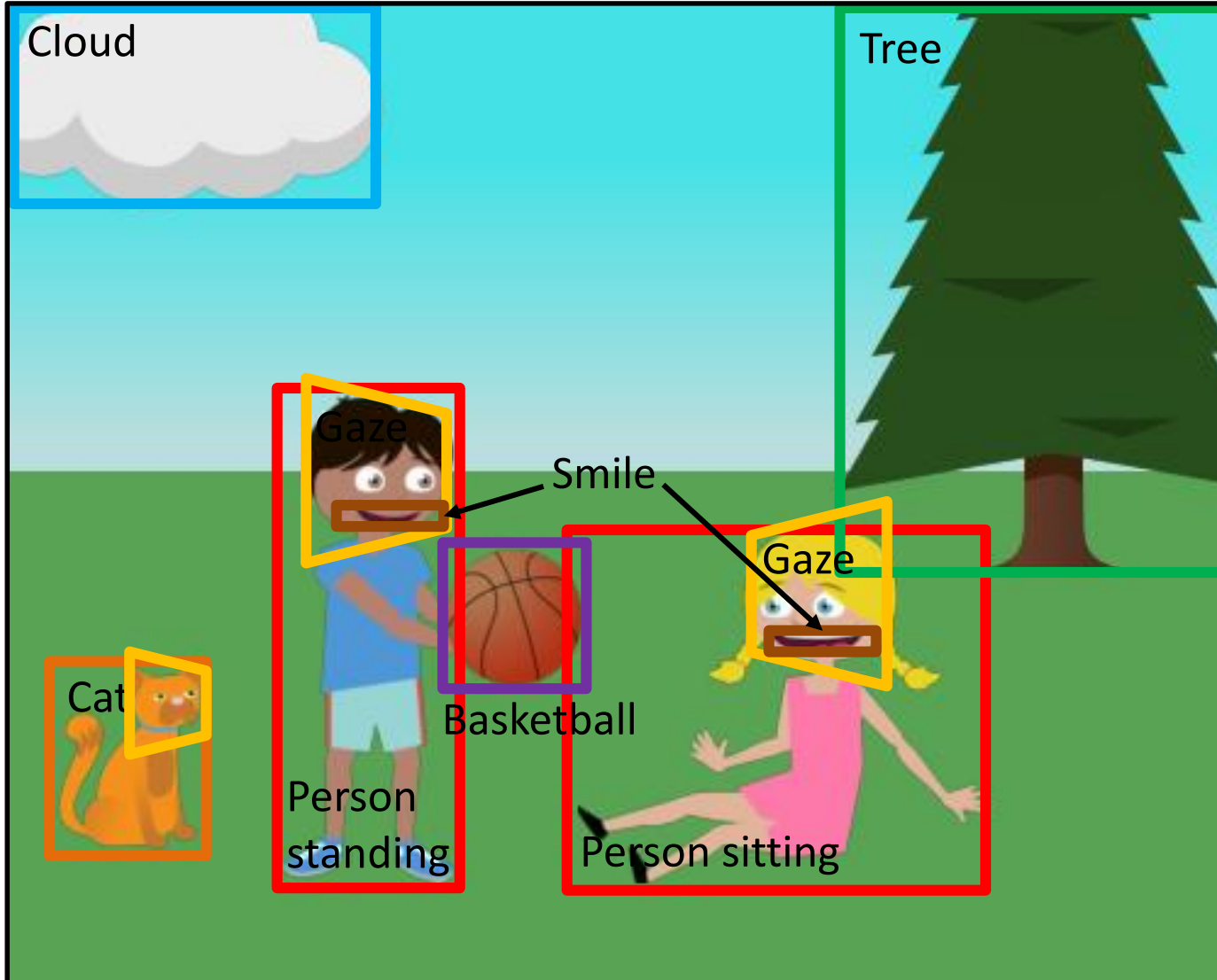
[Zitnick and Parikh, CVPR 2013, Oral]



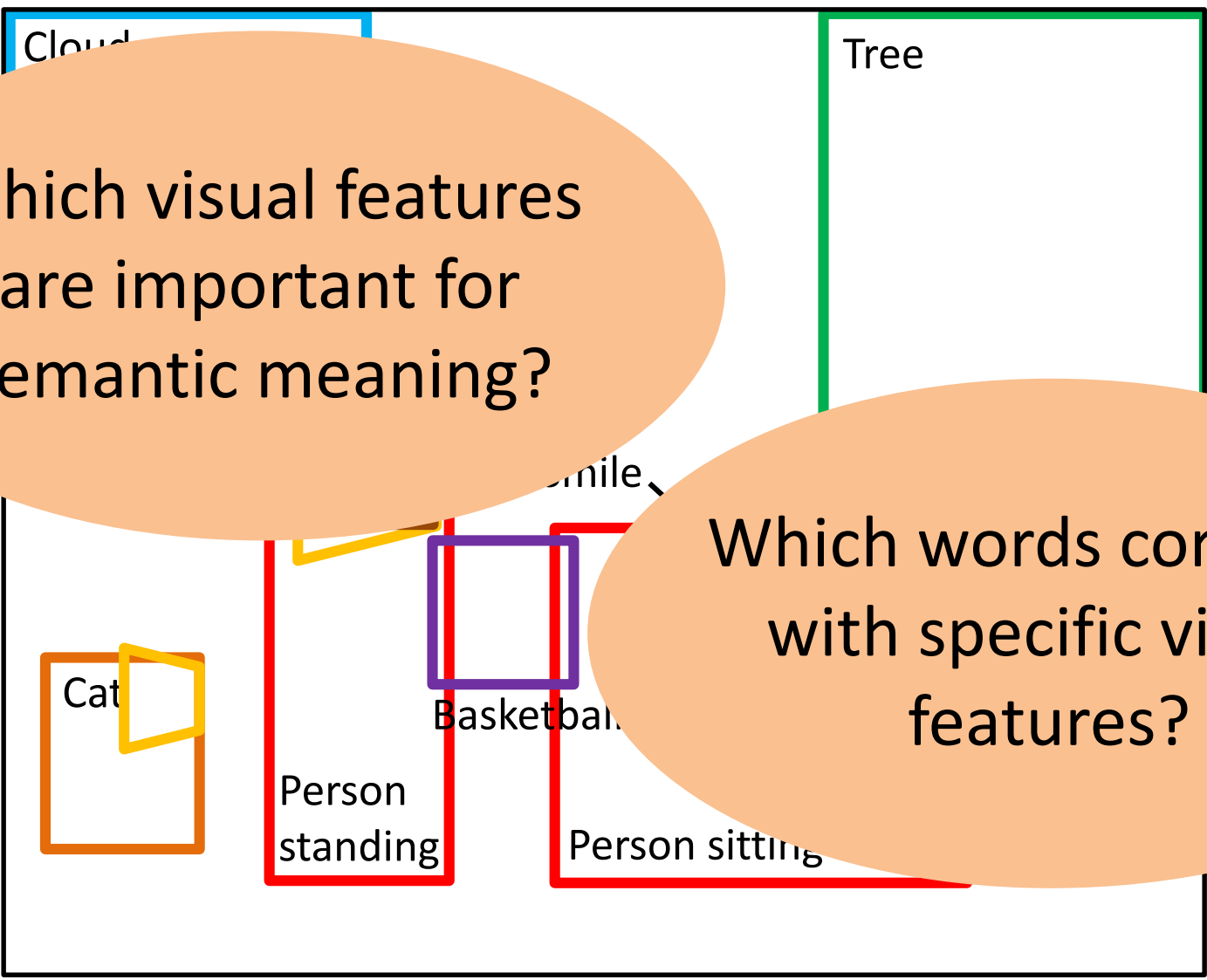
Visual Features



Visual Features



Visual Features



Which visual features are important for semantic meaning?

Which words correlate with specific visual features?

Generate Scenes

Input: Jenny is catching the ball. Mike is kicking the ball. The table is next to the tree.

Tuples: <<Jenny>,<catch>,<ball>> <<Mike>,<kick>,<ball>> <<table>,<be>,<>>



Automatically Generated



Human Generated

[Zitnick, Parikh and Vanderwende, ICCV 2013]

Generate Scenes

Jenny was mad and tried to kick Mike.

<<Jenny>, <be mad>, <>>
 <<Jenny>, <try>, <kick>>
 <<Jenny>, <kick>, <Mike>>

Mike is holding a baseball bat.

<<Mike>, <hold>, <bat>>

Mike is standing in front of the table.

<<Mike>, <stand_in_front_of>, <table>>

Mike is wearing a blue hat with a star.

<<Mike>, <wear>, <hat>>
 <<hat>, <with>, <star>>

Jenny is happy to see Mike.

<<Jenny>, <be happy>, <>>
 <<Jenny>, <see>, <Mike>>

The cat is watching Jenny and Mike.

<<cat>, <watch>, <Jenny>>
 <<cat>, <watch>, <Mike>>

Jenny wants Mike to share the bat.

<<Jenny>, <want>, <share>>
 <<Mike>, <share>, <bat>>

Mike and Jenny are wearing hats.

<<Mike>, <wear>, <hat>>
 <<Jenny>, <wear>, <hat>>

Jenny is sad because she wants the ball.

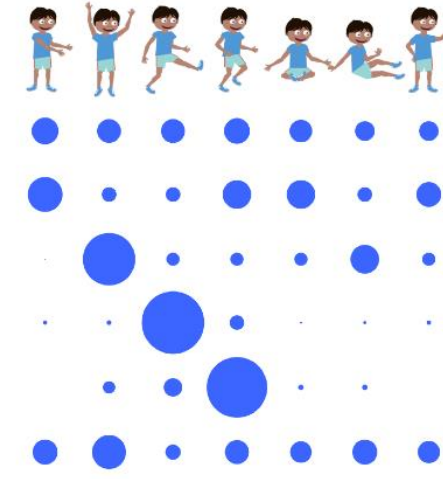
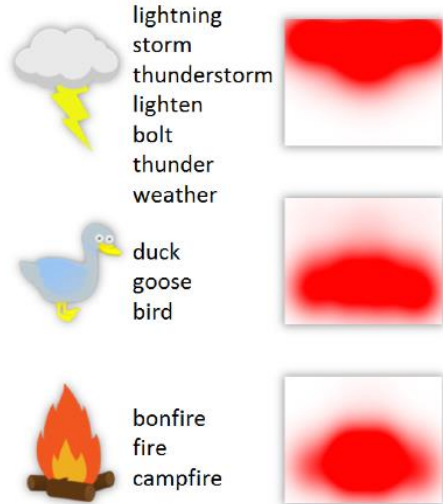
<<Jenny>, <be sad>, <>>
 <<she>, <want>, <ball>>

Mike is eating a burger

<<Mike>, <eat>, <burger>>

Jenny is on the swings.

<<Jenny>, <be>, <>>

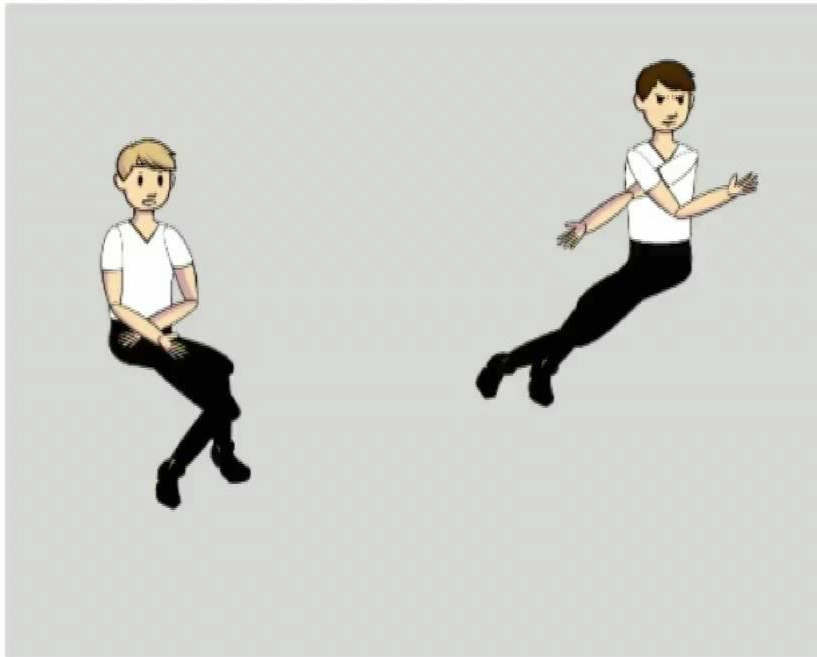


[Zitnick, Parikh and Vanderwende, ICCV 2013]

Learning Fine-grained Interactions

Illustrate this sentence:

Sentence 1/2: Person 1 is dancing with Person 2



Gender

Flip

Expression

A grid of face icons used for fine-grained control. The 'Gender' section shows two columns: the first column has a blonde-haired person and a brown-haired person; the second column has a blonde-haired person and a brown-haired person. The 'Flip' section shows two columns: the first column has a blonde-haired person and a brown-haired person; the second column has a blonde-haired person and a brown-haired person. The 'Expression' section shows two rows of seven faces each, with various expressions like neutral, smiling, and angry. A mouse cursor is pointing at the first face in the first row of the 'Expression' section.

Who is Person 1 in your creation? Blonde-haired person Brown-haired person

Who is Person 2 in your creation? Blonde-haired person Brown-haired person

Next

3x

[Antol, Zitnick and Parikh, ECCV 2014]

Learning Fine-grained Interactions

jumping over



holding hands with

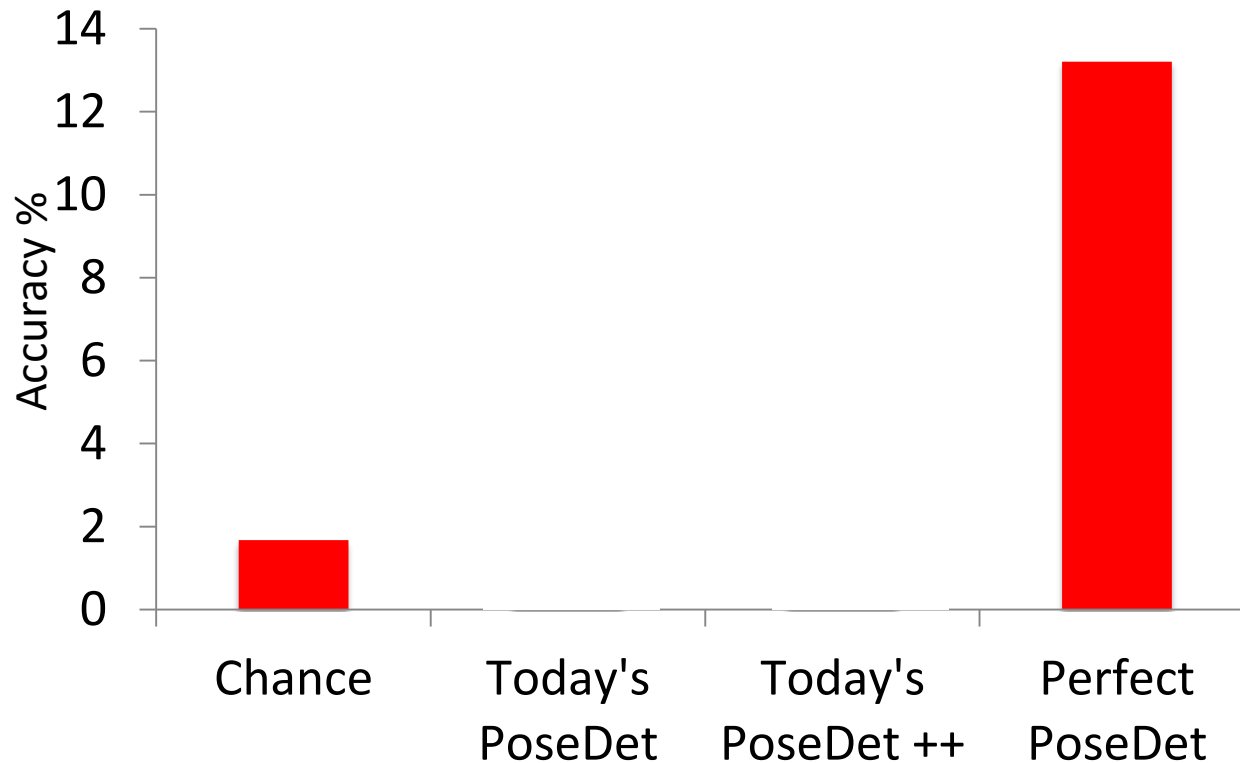


dancing with



Train on clipart, test on real

Results: 60 categories



[Antol, Zitnick and Parikh, ECCV 2014]

Learning Common Sense

- Assess plausibility of relations
 - man holds meal
 - tree grows in table
- Plausibility: similarity to other relations we know are plausible
 - person holds sandwich
 - man eats pizza
 - ...
- Textual and visual similarity

Results

Given *any* tuple, can assess its plausibility

	Average Precision	Rank Correlation
Text alone		
Visual alone		
Text + visual		

[Vedantam, Lin, Batra, Zitnick, and Parikh, ICCV 2015]

Online Demo

Online Demo

Home

Predicting Plausibility of Common Sense Assertions

Based on [Ramakrishna Vedantam*](#), [Xiao Lin*](#), [Tanmay Batra](#), [C. Lawrence Zitnick](#), [Devi Parikh](#), [Learning Common Sense Through Visual Abstraction](#), ICCV 2015. *Equal Contribution

Demo prepared by [Arijit Ray](#).

Enter the Primary, Relation and Secondary Phrases of a Tuple whose plausibility you want to assess:

Examples ▾

Or, submit your text file ▾

Entered Tuple : man eats cake

Predicted Plausibility Score : 0.3096

Show More Details ▾

Fill-in-the-blank:

Mike is having lunch when
he sees a bear.

_____.

- A. Mike orders a pizza.
- B. Mike hugs the bear.
- C. Bears are mammals.
- D. Mike tries to hide.

Fill-in-the-blank

Question

_____. Mike is wearing a blue cap. Mike is telling Jenny to get off the swing

Options and Generated Scenes

A. There is a tree near a table.

B. The brown dog is standing next to Mike.

C. The sun is in the sky.

D. Jenny is standing dangerously on the swing

Visual Paraphrasing

It is a sunny day.

Mike is sitting with a pizza.

Jenny is playing with a soccer ball.

Mike is eating a pizza.

Jenny is playing soccer.

A cat is eating a hot dog.



Same or different?

Results

	Fill-in-the-blanks (FITB) Accuracy (+/- ~0.15)	Visual Paraphrasing (VP) AP (+/- ~0.02)
Random	25.00	33.33

[Lin and Parikh, CVPR 2015]

Visual Abstraction For...

- Studying mappings between
- Zero-shot
- Studying
- Learning common sense knowledge
- Rich annotation modality
 - Ask for descriptions
 - Ask for scene graphs
 - Show scene graphs
 - Perturb a scene and ask for descriptions
 - ...

Study high-level image understanding tasks without waiting for lower-level vision tasks to be solved

Learning by playing!



50k scenes Available online!

Semantic Image Understanding



"Color College Avenue", Blacksburg, VA, May 2012

Semantic Image Understanding

Semantic



Words

Mike is wearing a blue hat with a star.

<<Mike>, <wear>, <hat>>

<<hat>, <with>, <star>>

Jenny is happy to see Mike.

<<Jenny>, <be happy>, <>>

<<Jenny>, <see>, <Mike>>

Image



Pictures

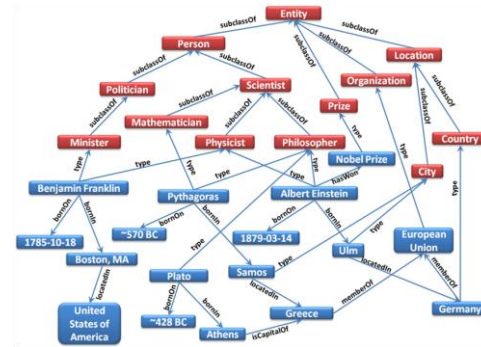


Understanding



Reasoning

(Common Sense, Knowledge Base)



Thank you.

