



MIT CSAIL

6.869: Advances in Computer Vision

MIT
COMPUTER
VISION

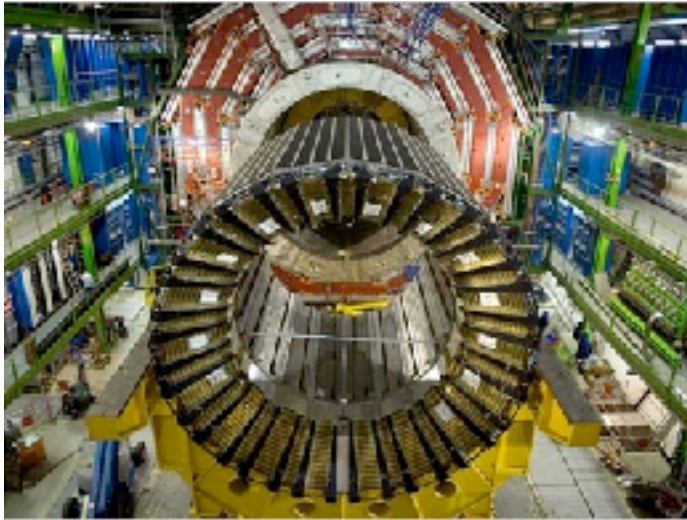
Lecture 20

Words and pictures

Crowdsourcing



The value of data



The Large Hadron Collider
 $\$ 10^{10}$



Amazon Mechanical Turk
 $\$ 10^2 - 10^4$

But can humans collect good
data?

Google

bedroom



Search

About 299,000,000 results (0.19 seconds)

Everything

Related searches: [bedroom designs](#) [master bedroom](#) [modern bedroom](#) [simple bedroom](#) [small bedroom](#)

Images

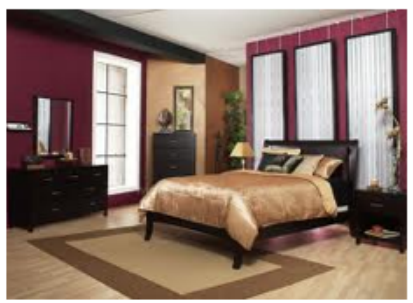
Maps

Videos

News

Shopping

More



Any time

Past 24 hours

Past week

Custom range...



All results

By subject

Personal



Any size

Large

Medium

Icon

Larger than...

Exactly...



Search

About 66,700,000 results (0.15 seconds)



Everything

Images

Maps

Videos

News

Shopping

More

Any time

Past 24 hours

Past week

Custom range...

All results

By subject

Personal

Any size

Large

Medium

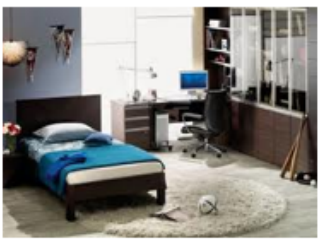
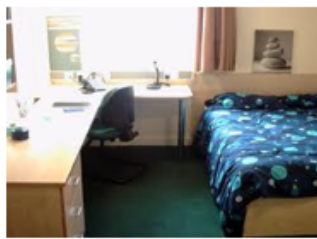
Icon

Larger than...

Exactly...

Any color

Full color





www.bigstock.com - 7067629



Google

mug



Getting more humans in the annotation loop

Labeling to get a Ph.D.



Labeling for fun

Luis Von Ahn and Laura Dabbish 2004



Labeling for money
(Sorokin, Forsyth, 2008)



Labeling because it
gives you added value



Visipedia
(Belongie, Perona, et al)

Just for labeling



Beware of the human in your loop

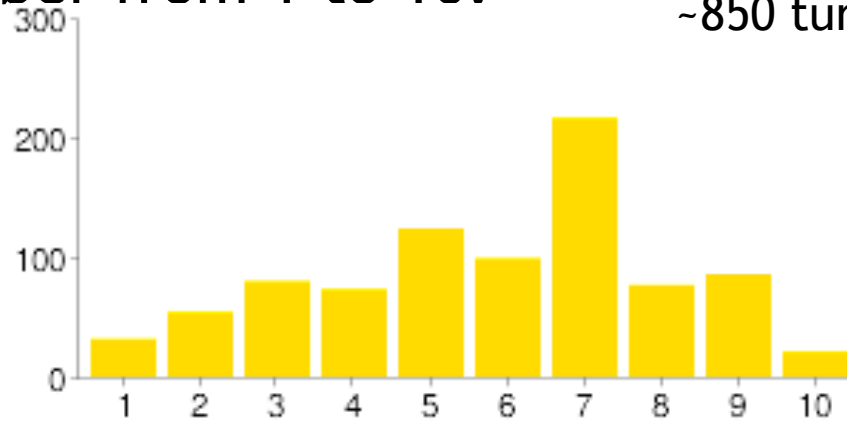
- What do you know about them?
- Will they do the work you pay for?

Let's check a few simple experiments

People has biases...

Turkers were offered 1 cent to pick a number from 1 to 10.

~850 turkers



Do humans have consistent biases?

Choose Item

Requester: SimpleSphere

Reward: \$0.01 per HIT

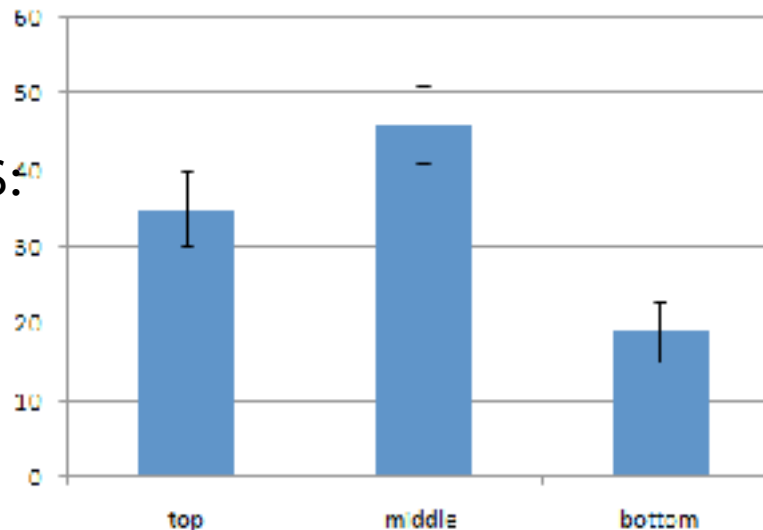
HITs Available: 1

Duration: 50 minutes

Qualifications Required: None

Please choose one of the following:

Results form 100 HITS.

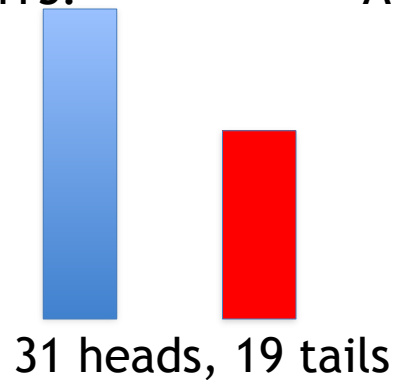


Do humans do what you ask for?

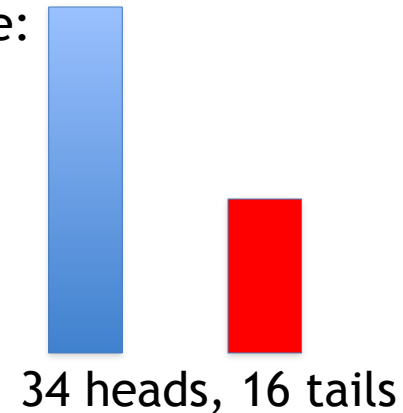
Find a coin
Requester: ROB MILLER Reward: \$0.01 per HIT HITS Available: 1 Duration: 5 minutes
Qualifications Required: None

Please flip an actual coin and type either H or T below.

After 50 HITS:



And 50 more:



Are humans reliable even in simple tasks?

Choose the given item.

Requester: SimpleSolvers

Reward: \$0.01 per HIT

HITs Available: 1

Duration: 50 minutes

Qualifications Required: None

Please click button B:

B

C

A

Results of 100 HITS

A: 2

B: 96

C: 2



Please [contact us](#) if you find any bugs or have any suggestions.



[Show me another image](#)

[Sign In](#) ([why?](#))

With your help, there are **91348** labelled objects in the database ([more stats](#))

Instructions ([Get more help](#))

Use your mouse to click around the boundary of some objects in this image. You will then be asked to enter the name of the object (examples: car, window).



Labeling tools



Polygons in this image [view](#)

- [door](#)
- [door](#)
- [road](#)
- [train](#)
- [window](#)
- [window](#)
- [sidewalk](#)
- [building region](#)
- [house](#)
- [window](#)
- [window](#)
- [window](#)

Label as many objects and regions as you can in this image



Tool went online July 1st, 2005

Labelme.csail.mit.edu

All HITs | HITs Available To You | HITs Assigned To You

Search for HITs containing

that pay at least \$ 0.00 for which you are qualified

Timer: 00:00:13 of 60 minutes

Finished with this HIT?

Submit HIT

Let someone else do it!

Retire HIT

Total Earned: \$9.01
Total HITs Submitted: 12

Automatically accept the next HIT

Labeler: Local objects in this image

Requester: Bryan C. Russell

Qualifications Required: None

Reward: \$0.01 per HIT

HITs Available: 25

Duration: 60 minutes

Please label as many objects as you want in this image. Scroll down to see the entire image.

Search HIT

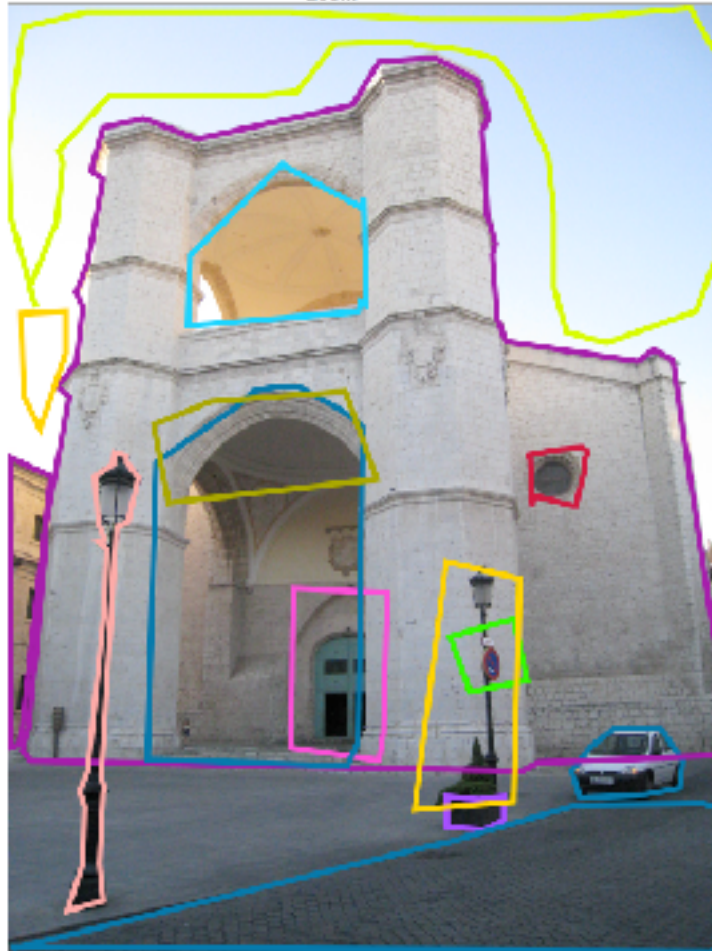


With Bryan Russell

1 cent

Task: Label one object in this image





car sky

car

building

car sky

car sky

car sky

car sky

car sky

car sky

car

car sky

car sky

car

car

car

car sky

LabelMe iterations

- 1) Label as many objects as you can
- 2) Delete any wrong polygon
- 3) Go to 1



Label some objects



Delete any wrong polygons



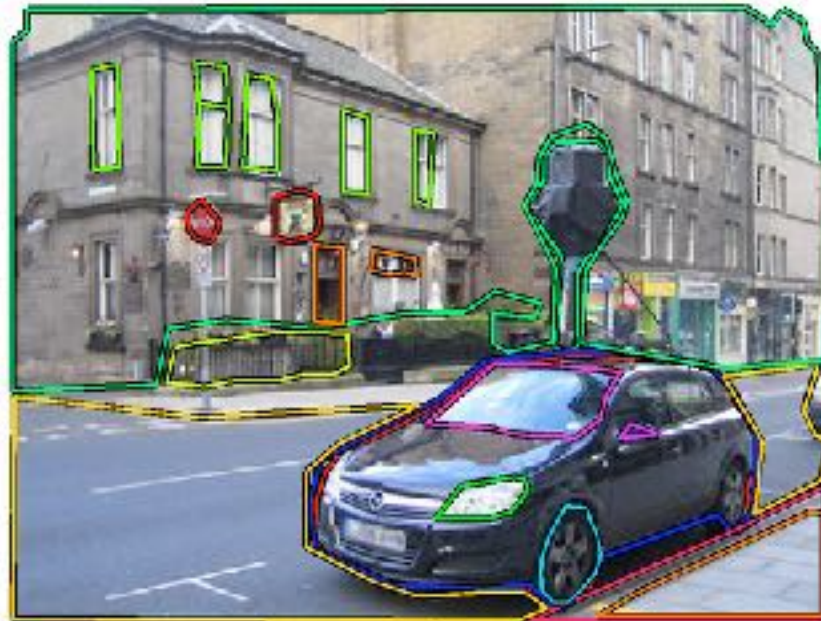
Label some objects



Delete any wrong polygons



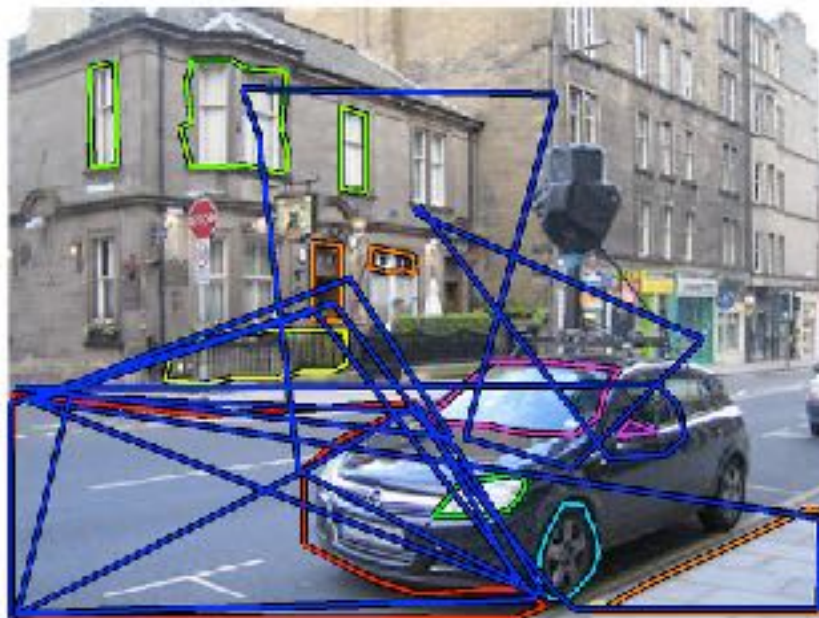
Label some objects



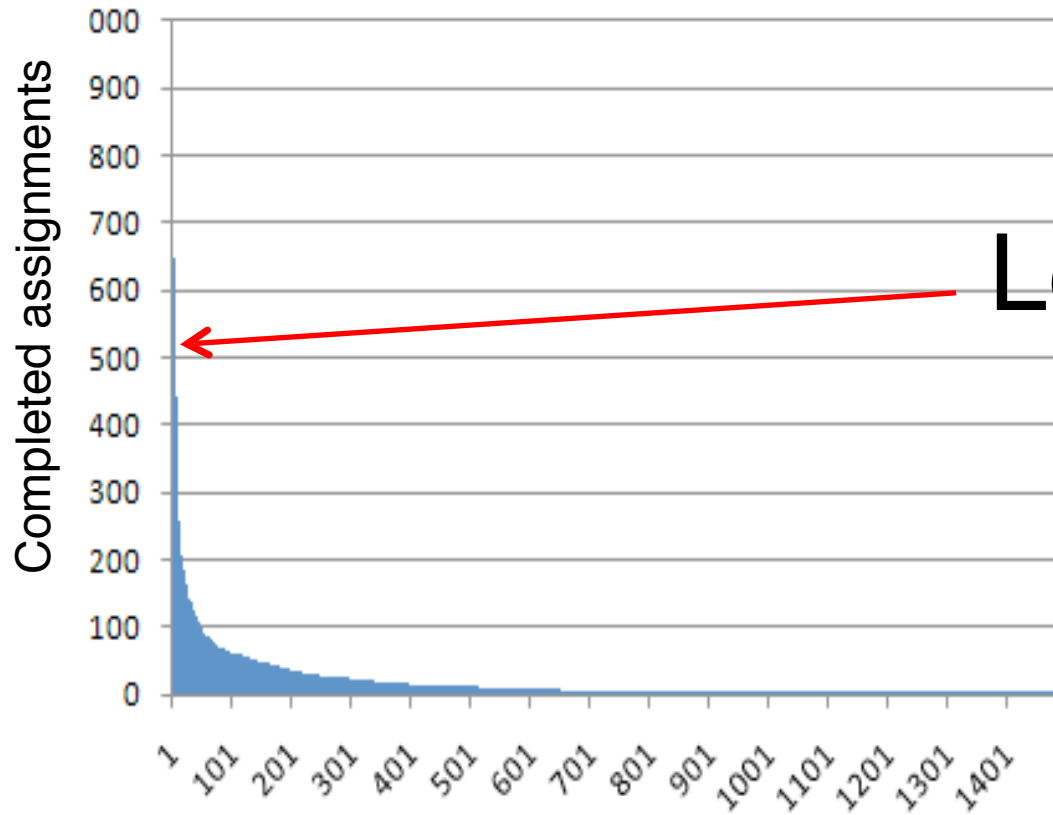
Delete any wrong polygons



Label some objects



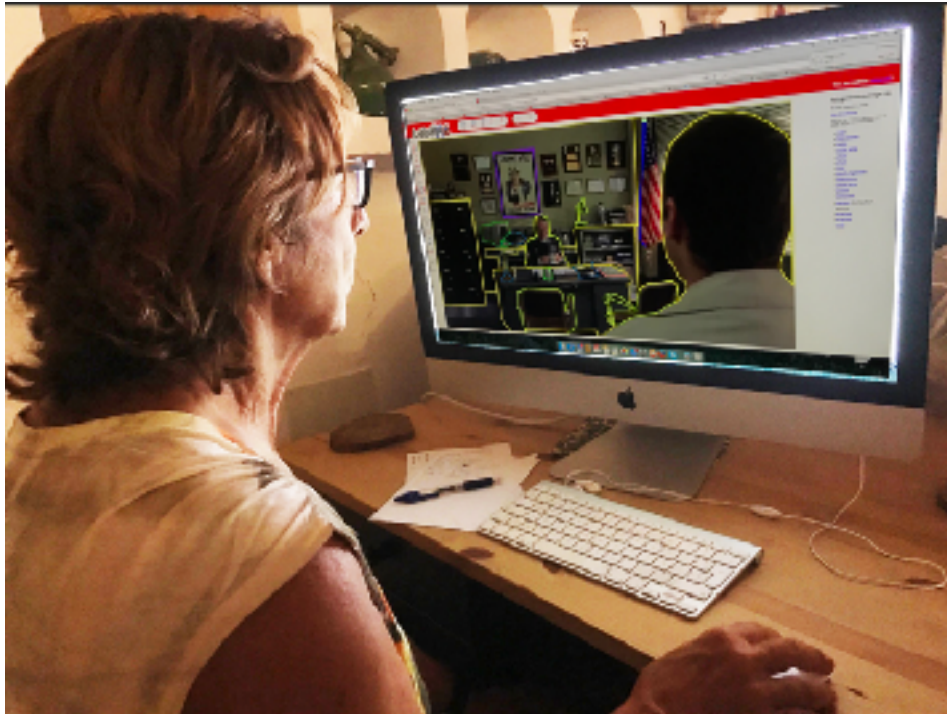
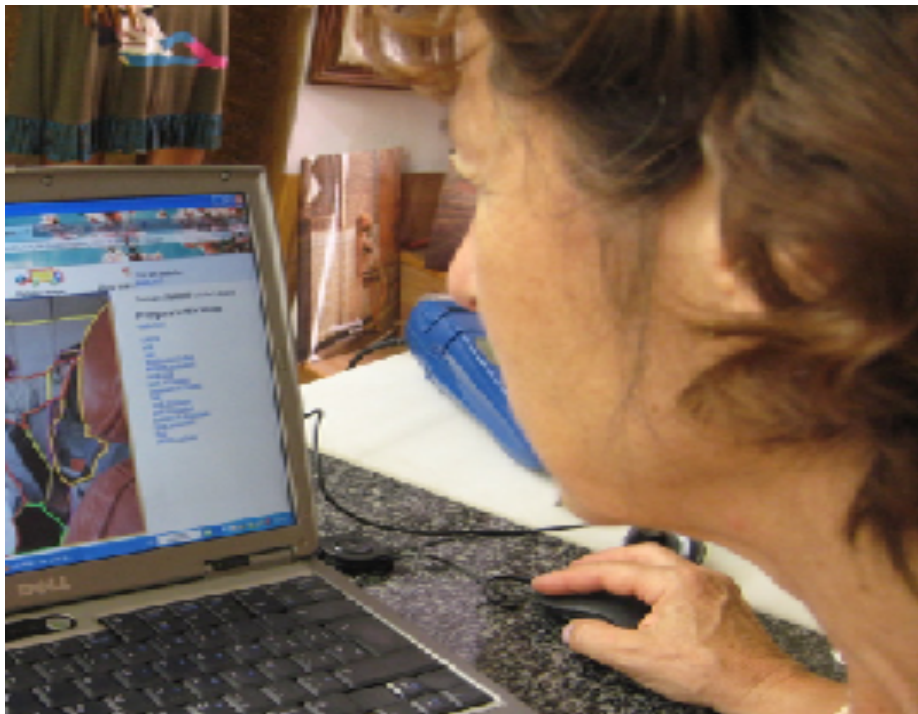
Who does the work?



Let's hire that

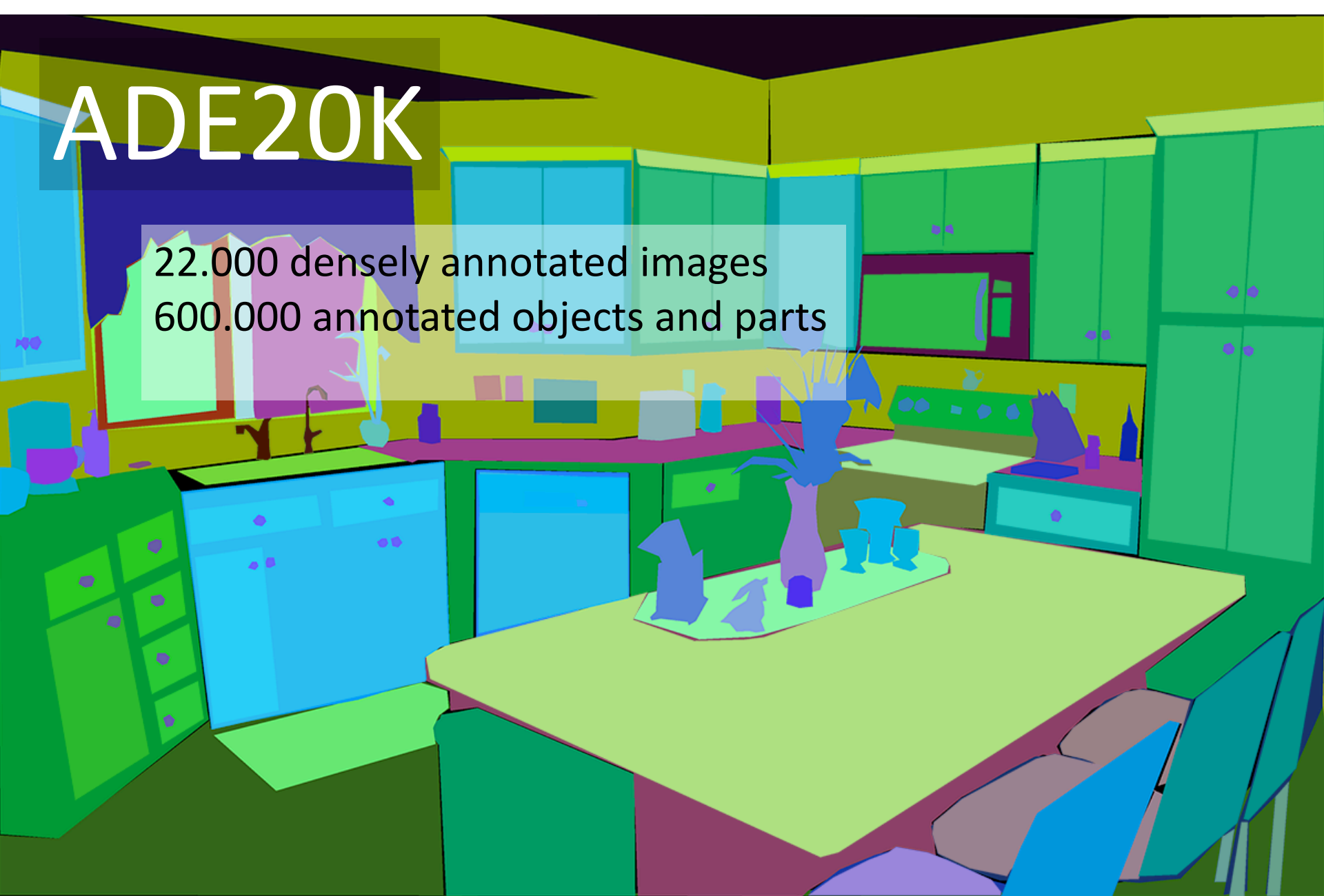
Workers sorted by contribution

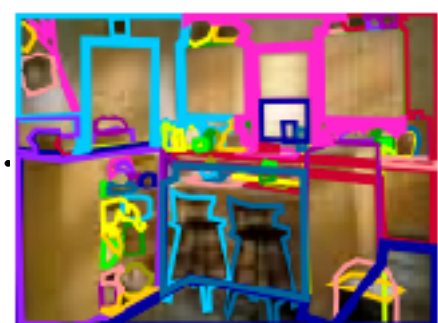
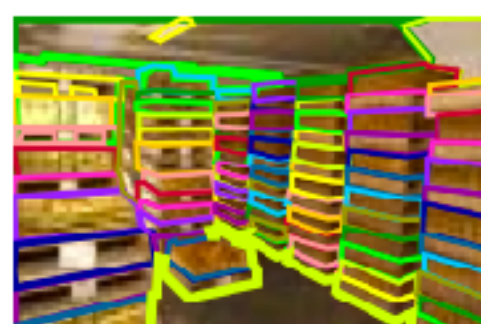
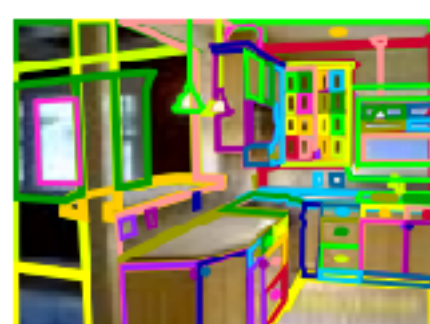
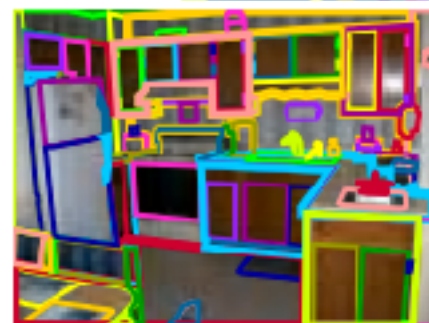
From <http://groups.csail.mit.edu/uid/deneme/>



ADE20K

22.000 densely annotated images
600.000 annotated objects and parts





COCO



ADE20K



	Images	Obj. inst.	Obj. classes	Part inst.	Part classes	Obj. classes per image
COCO	123,287	886,284	91	0	0	3.5
ImageNet*	476,688	534,309	200	0	0	1.7
NYU Depth V2	1,419	34,064	894	0	0	14.1
Cityscapes	25,000	N/A	30	0	0	N/A
SUN	16,873	313,884	4,479	0	0	9.8
OpenSurfaces	22,214	71,460	160	0	0	N/A
PascalContext	10,103	~104,398**	540	181,770	40	5.1
ADE20K	22,000	415,099	2,944	171,148	354	10.5

* has only bounding boxes (no pixel-level segmentation). Sparse annotations.

** PascalContext dataset does not have instance segmentation. In order to estimate the number of instances, we find connected components (having at least 150pixels) for each class label.



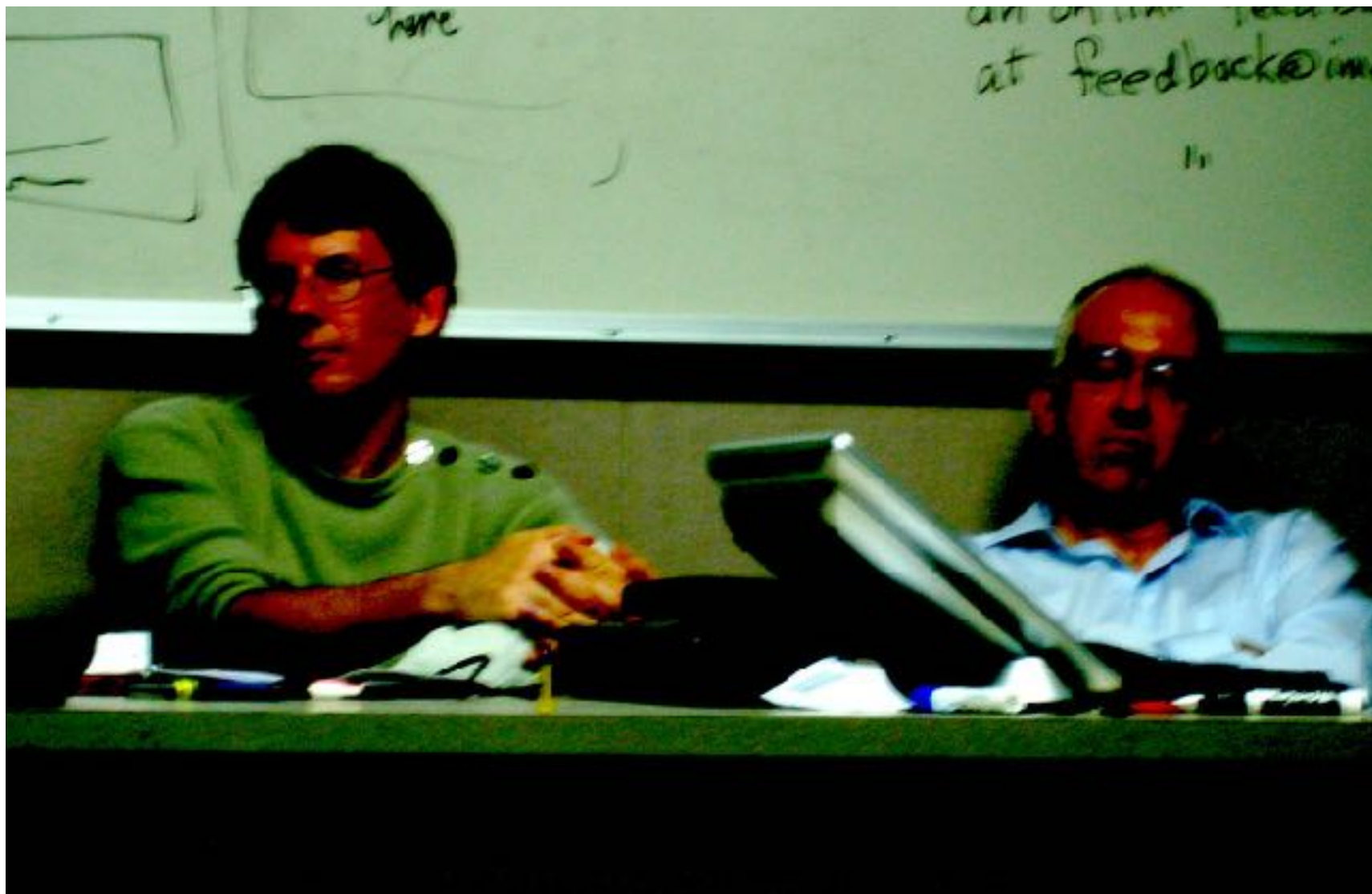
<https://www.youtube.com/watch?v=AlEeakeXV>



Cross modal learning text and images



Two man sitting behind a long table.



Q: Is everyone of these two holding a wine glass? A: No

Q: How many people are there? A: 2

Q: How many are awake? A: 1

Inspired from <http://visualqa.org/index.ht>

Story-like Description

Pietro had a long day of talks at the workshop. At the end of the session, Pietro was invited to participate in a panel. As he is a mature and confident professor, he decided to take a short nap during the discussion. The chair was comfortable. Nobody dared to wake him up as there were other less confident professors at the panel that could answer the questions. ...



“Pictures and words”

- Barnard, Duygulu, de Freitas, Forsyth, Blei, Jordan, Matching words and pictures, JMLR, 2003
- Duygulu, Barnard, de Freitas, Forsyth, Object Recognition as Machine Translation: Learning a lexicon for a fixed image vocabulary , ECCV, 2003
- Blei & Jordan, Modeling annotated data, ACM SIGIR, 2003
- Chang, Goh, Sychay, & Wu, *Soft* annotation using Bayes point machines, IEEE Transactions on Circuits and Systems for Video Technology, 2003
- Goh, Chang, & Cheng, Ensemble of SVM-based classifiers for annotation, 2003
-

Object Recognition as Machine Translation: Learning a Lexicon for a Fixed Image Vocabulary

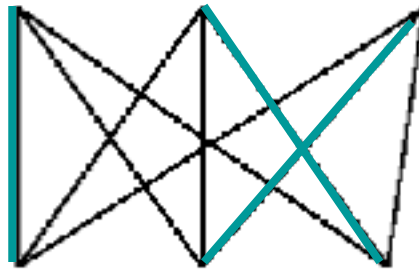
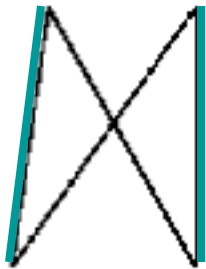
P. Duygulu¹ , K. Barnard¹ , J.F.G. de Freitas² and D.A. Forsyth¹

Computer Science Division, U.C. Berkeley, Berkeley, CA 94720
Department of Computer Science, University of British Columbia, Vancouver
{duygulu, kobus, daf}@cs.berkeley.edu, nando@cs.ubc.ca

Statistical Machine Translation

- Statistically link words in one language to words in another
- Requires aligned bitext
 - eg. Hansard for Canadian parliament

... la maison ... la maison bleue... la fleur ...



... the house ... the blue house ... the flower ...

Multimedia Translation

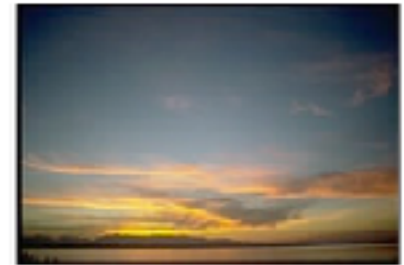
- Data:



116011
WATER HARBOR
SKY CLOUDS



TIGER CAT WATER GRASS



1090
SUN CLOUDS
WATER SKY

– Words are associated with images, but correspondences are unknown



sun sea sky



sun sea sky

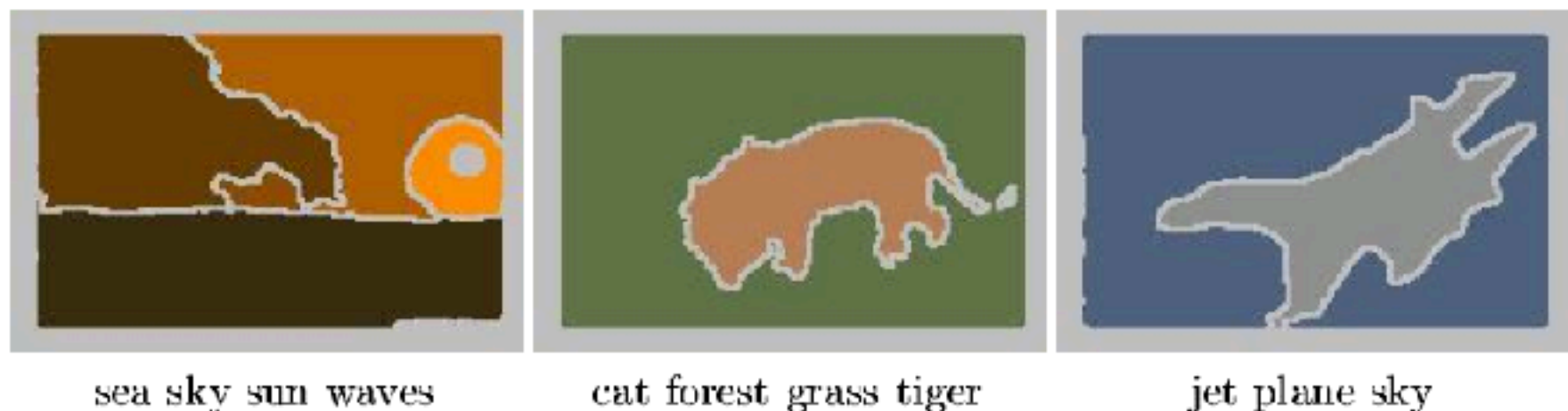


Fig. 1. *Examples from the Corel data set. We have associated keywords and segments for each image, but we don't know which word corresponds to which segment. The number of words and segments can be different; even when they are same, we may have more than one segment for a single word, or more than one word for a single blob. We try to align the words and segments, so that for example an orange stripy blob will correspond to the word tiger.*

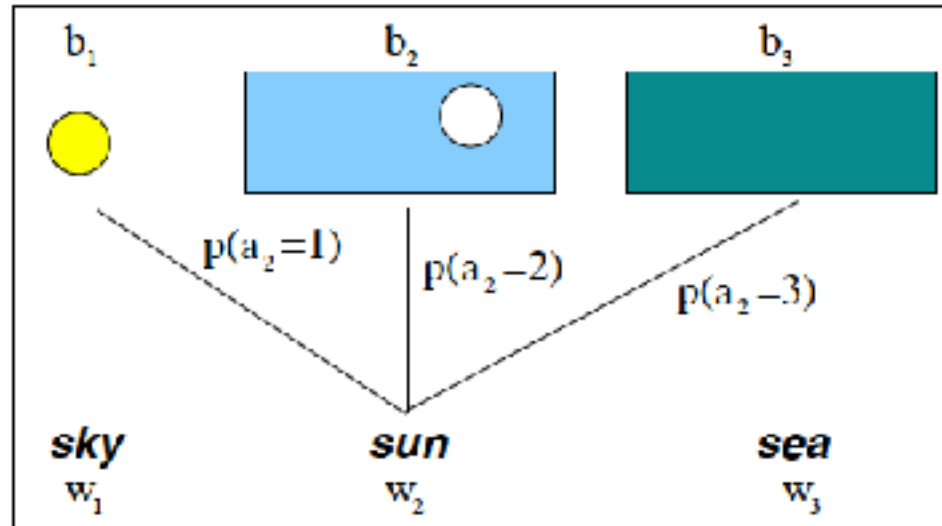


Fig. 3. Example : Each word is predicted with some probability by each blob, meaning that we have a mixture model for each word. The association probabilities provide the correspondences (assignments) between each word and the various image segments. Assume that these assignments are known; then computing the mixture model is a matter of counting. Similarly, assume that the association probabilities are known; then the correspondences can be predicted. This means that EM is an appropriate estimation algorithm.



Fig. 8. Some examples of the labelling results. The words overlaid on the images are the words predicted with top probability for corresponding blob. We are very successful in predicting words like sky, tree and grass which have high recall. Sometimes, the words are correct but not in the right place like tree and buildings in the center image.

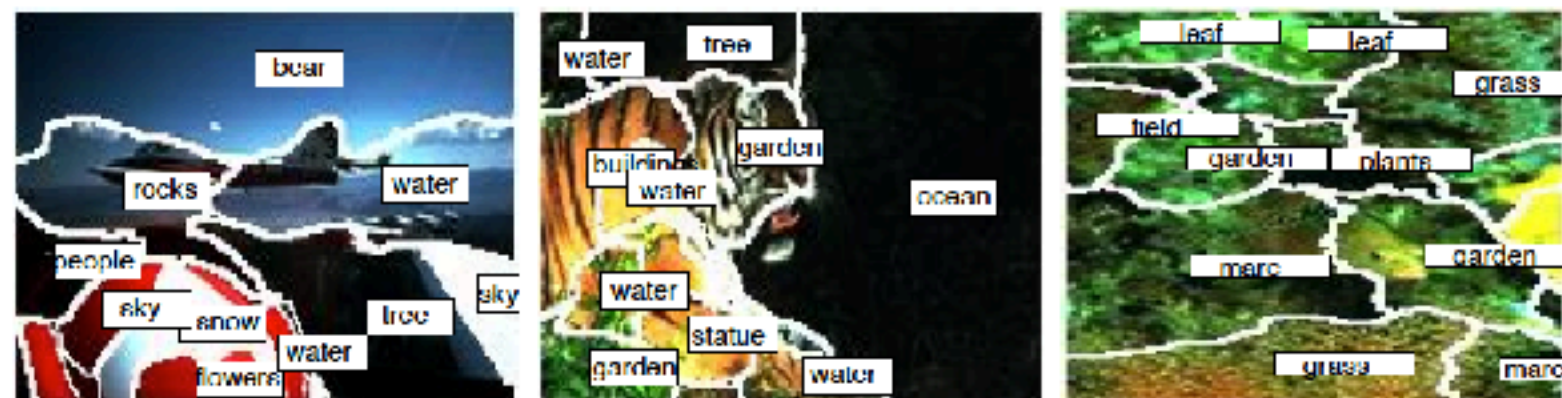
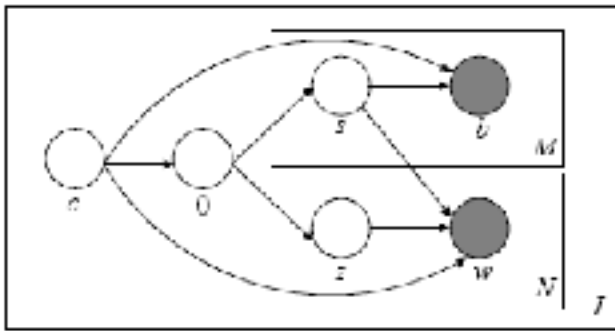
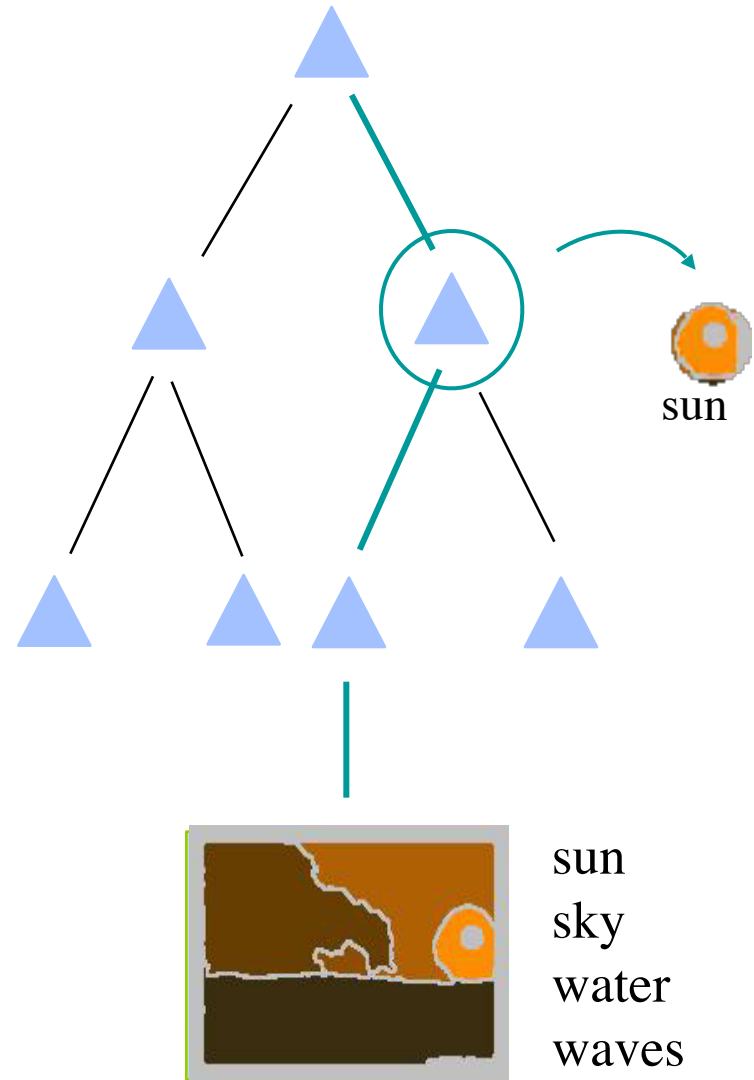


Fig. 9. Some test results which are not satisfactory. Words that are wrongly predicted are the ones with very low recall values. The problem mostly seen in the third image is since green blobs occur mostly with grass, plants or leaf rather than the under water plants.



- A generative model for assembling image data sets from multimodal clusters
 - Chose an image cluster by $p(c)$
 - Chose multimodal concept clusters using $p(s|c)$
 - From each multimodal cluster, sample a Gaussian for blob features, $p(b|s)$, and a multinomial for words, $p(w|s)$
 - (Skip with some probability to account for mismatched numbers of words and blobs)
 - For a given correspondence*

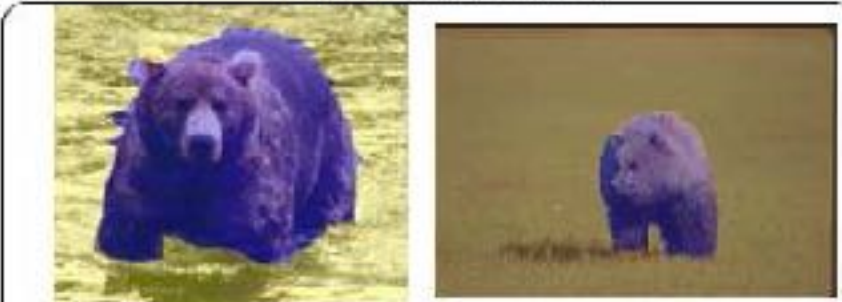
$$p(\{w \Leftrightarrow b\}) = \sum_c p(c) \prod_{\{w \Leftrightarrow b\}} \left(\sum_l p(w|l)p(b|l)p(l|c) \right)$$



“Beyond nouns”



Car is on the Street



Bear in water

Bear is on the field

Co-occurrence:

Red Car/Street

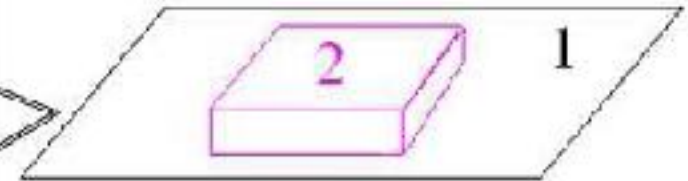
Blue Car/Street

Green Bear

Orange Field

Yellow Water

Our Approach:



2 is on 1

Red Street

Blue Car

Green Bear

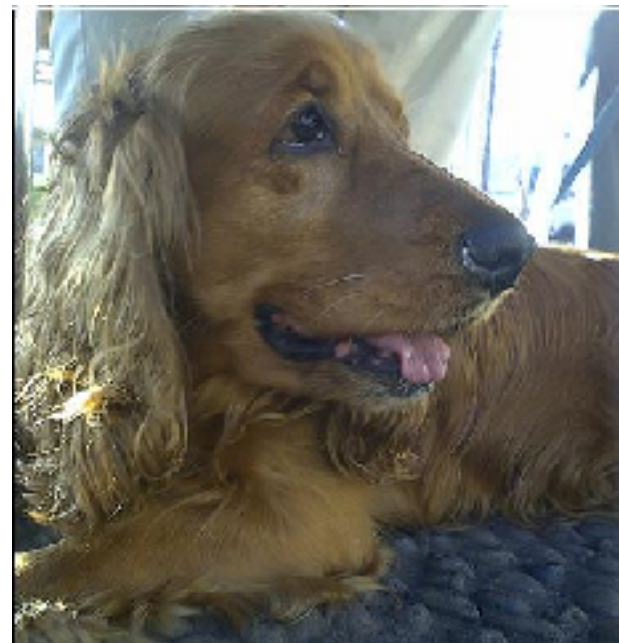
Orange Field

Yellow Water

What, where and who? Classifying events by scene and object recognition



Attribute Examples



Shape:

Part: Head, Ear, Nose,
Mouth, Hair, Face,
Torso, Hand, Arm

Material: Skin, Cloth

Shape:

Part: Head, Ear, Snout, Eye
Material: Furry

Shape:

Part: Head, Ear, Snout,
Eye, Torso, Leg
Material: Furry



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



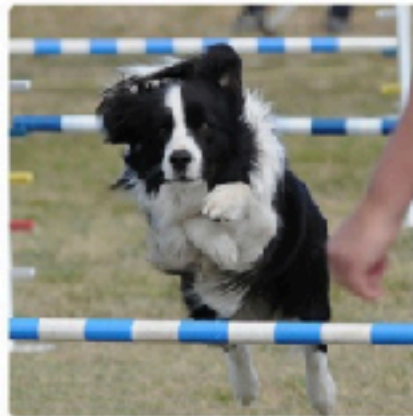
"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"girl in pink dress is jumping in air."



"black and white dog jumps over bar."



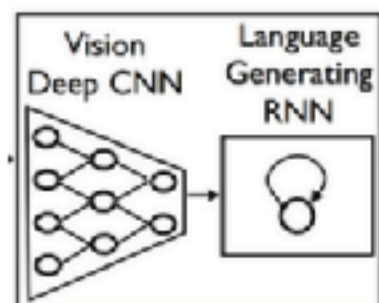
"young girl in pink shirt is swinging on swing."



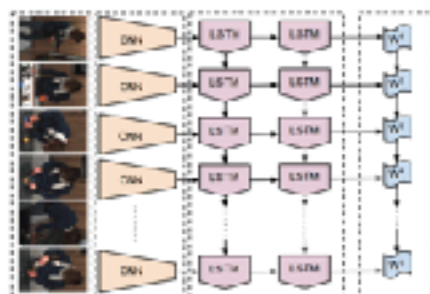
"man in blue wetsuit is surfing on wave."

Slides Credits: Andrej Karpathy, FeiFei Li

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



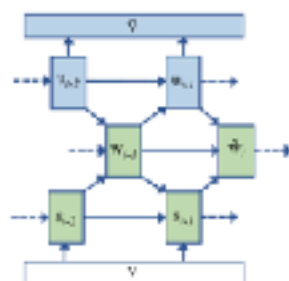
Mao et al., 2015



Ordonez et al., 2011



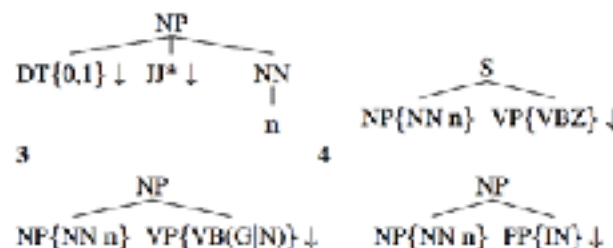
Kulkarni et al., 2011



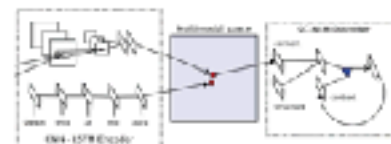
Chen and Zitnick, 2015



Farhadi et al., 2010

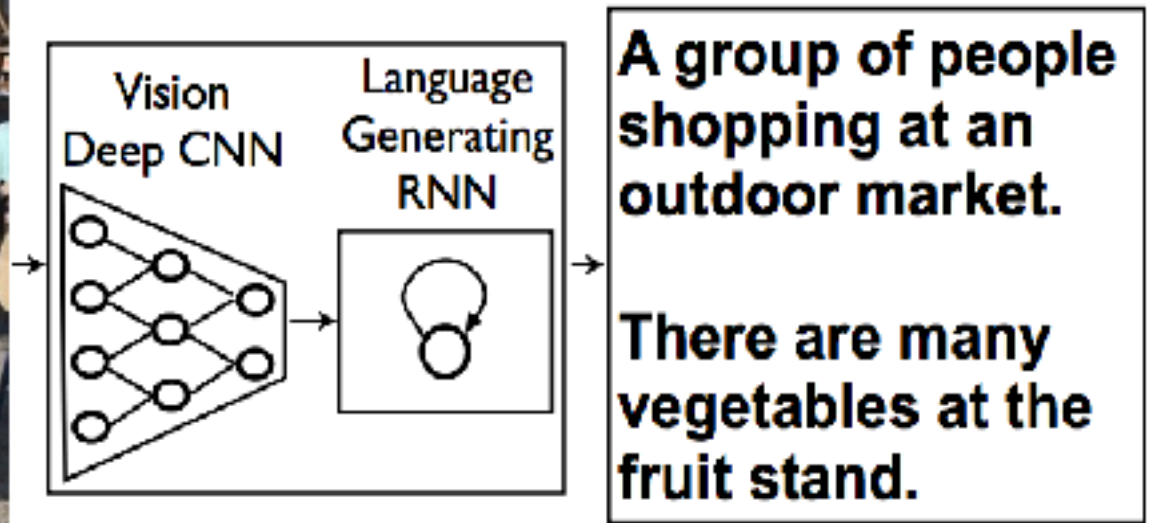


Mitchell et al., 2012



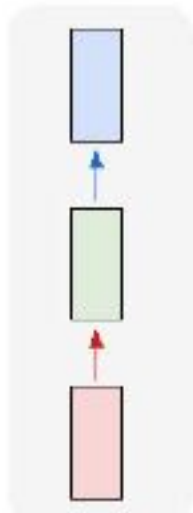
Kiros et al., 2015

Neural Image Caption (NIC) (CVPR 2015)



How do we model sequences?

one to one

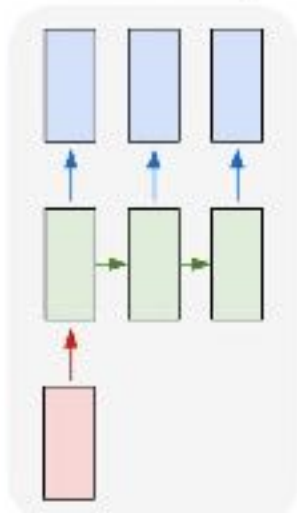


Input: No sequence

Output: No sequence

Example: "standard" classification / regression problems

one to many

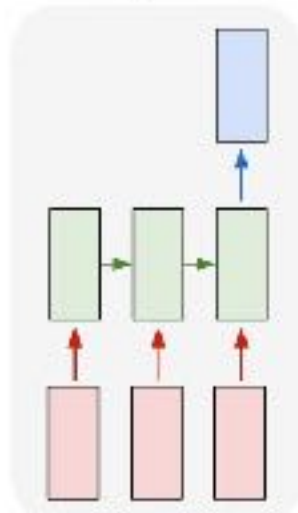


Input: No sequence

Output: Sequence

Example: Im2Caption

many to one

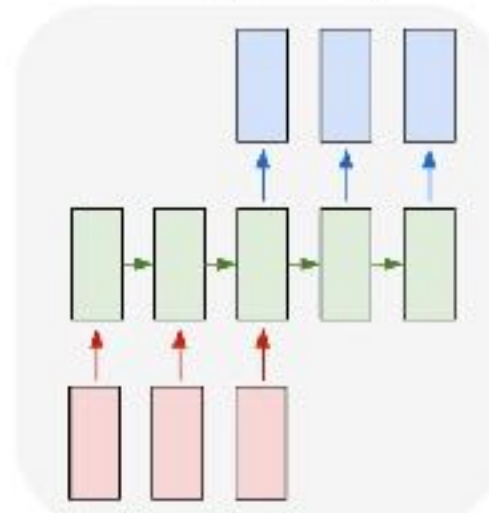


Input: Sequence

Output: No sequence

Example: sentence classification, multiple-choice question answering

many to many

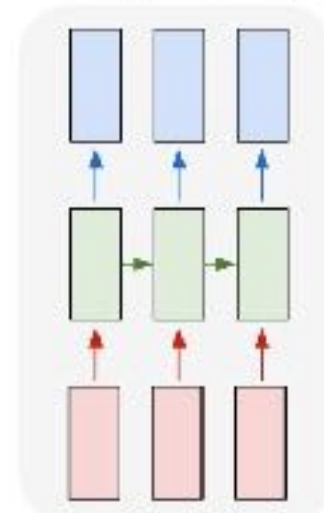


Input: Sequence

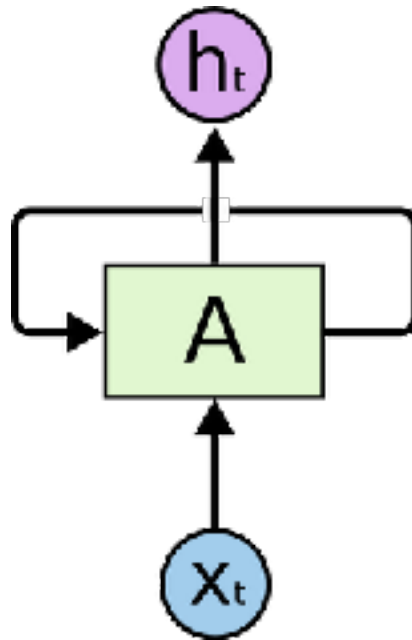
Output: Sequence

Example: machine translation, video captioning, open-ended question answering, video question answering

many to many

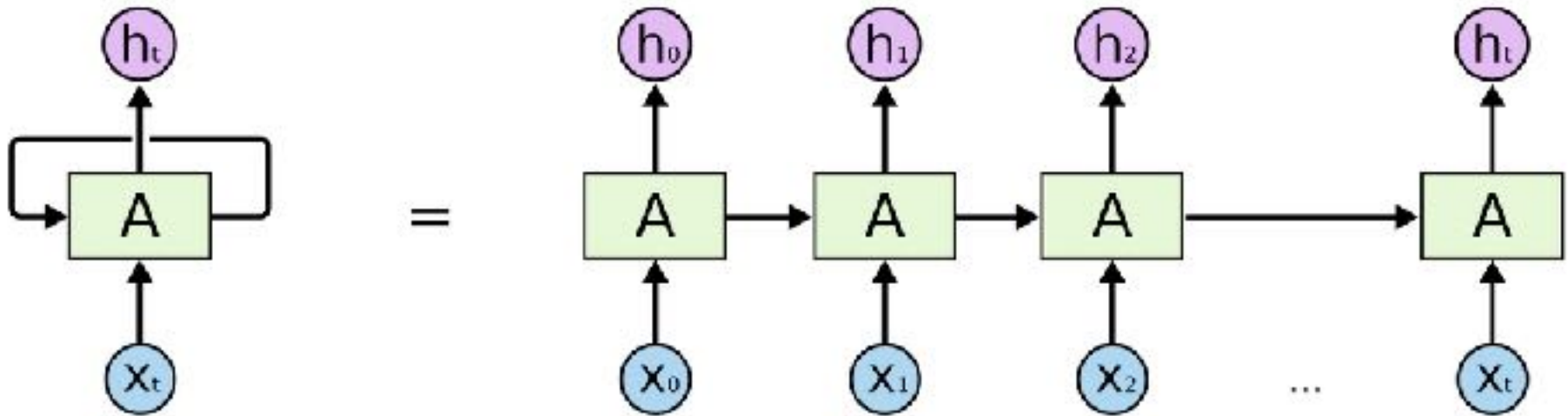


Recurrent Neural Networks (RNNs)



In the above diagram, a chunk of neural network, A , looks at some input x_t and outputs a value h_t . A loop allows information to be passed from one step of the network to the next.

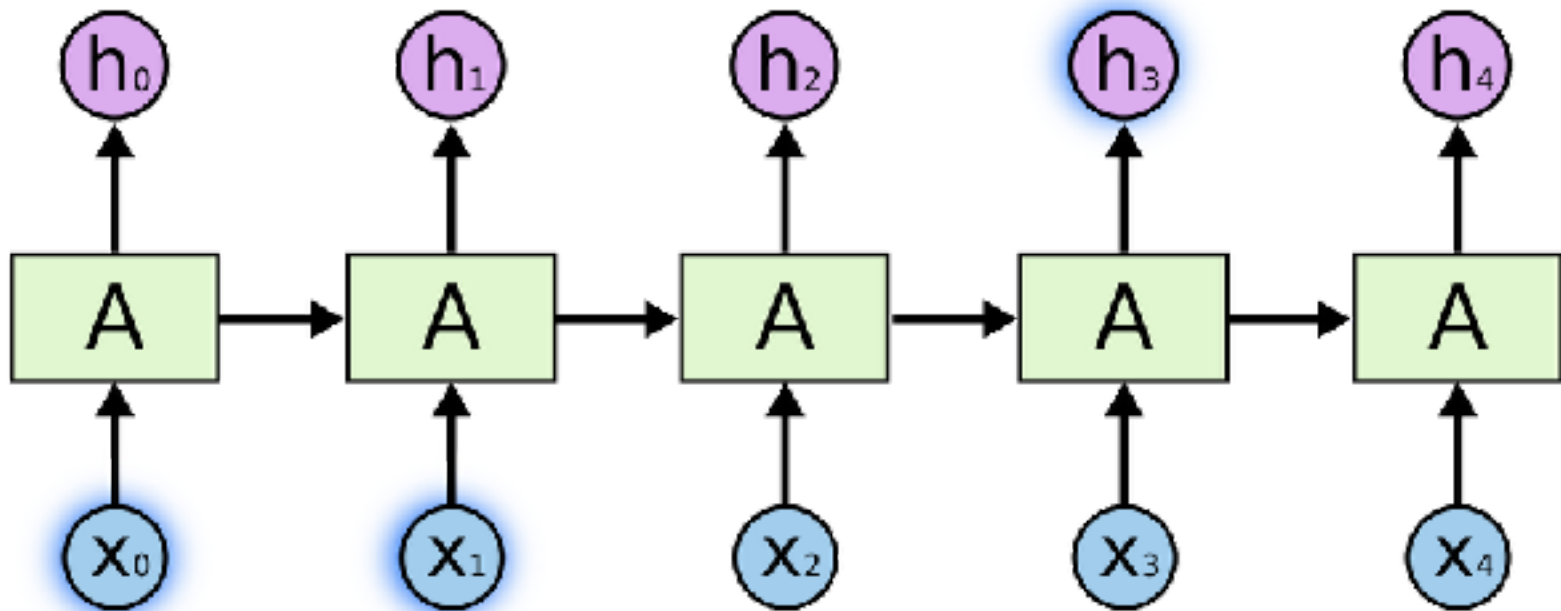
Recurrent Neural Networks (RNNs)



An unrolled recurrent neural network.

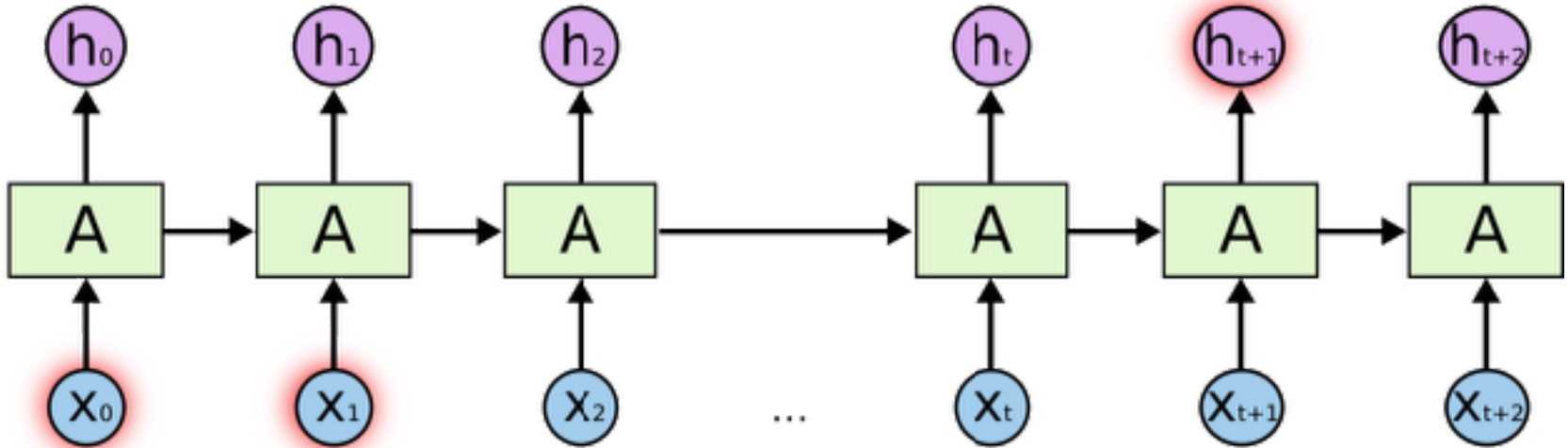
A recurrent neural network can be thought of as multiple copies of the same network, each passing a message to a successor

Recurrent Neural Networks (RNNs)



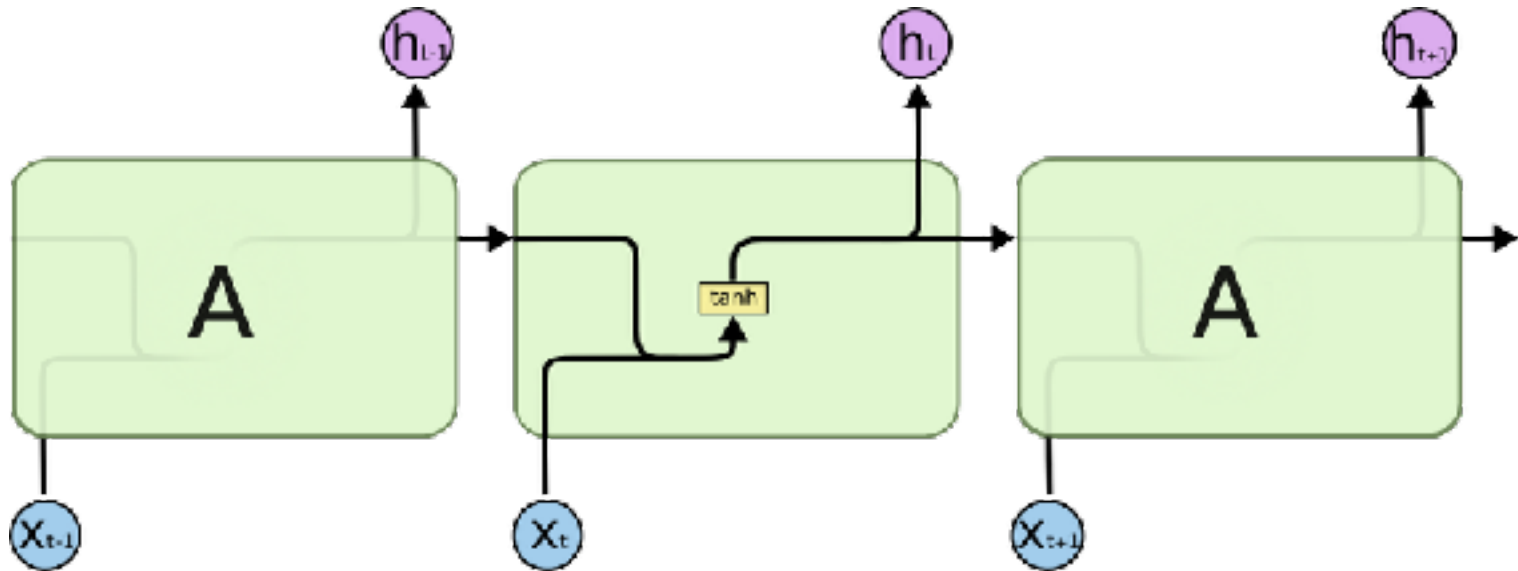
When the gap between the relevant information and the place that it's needed is small, RNNs can learn to use the past information

Long-term dependencies - hard to model!



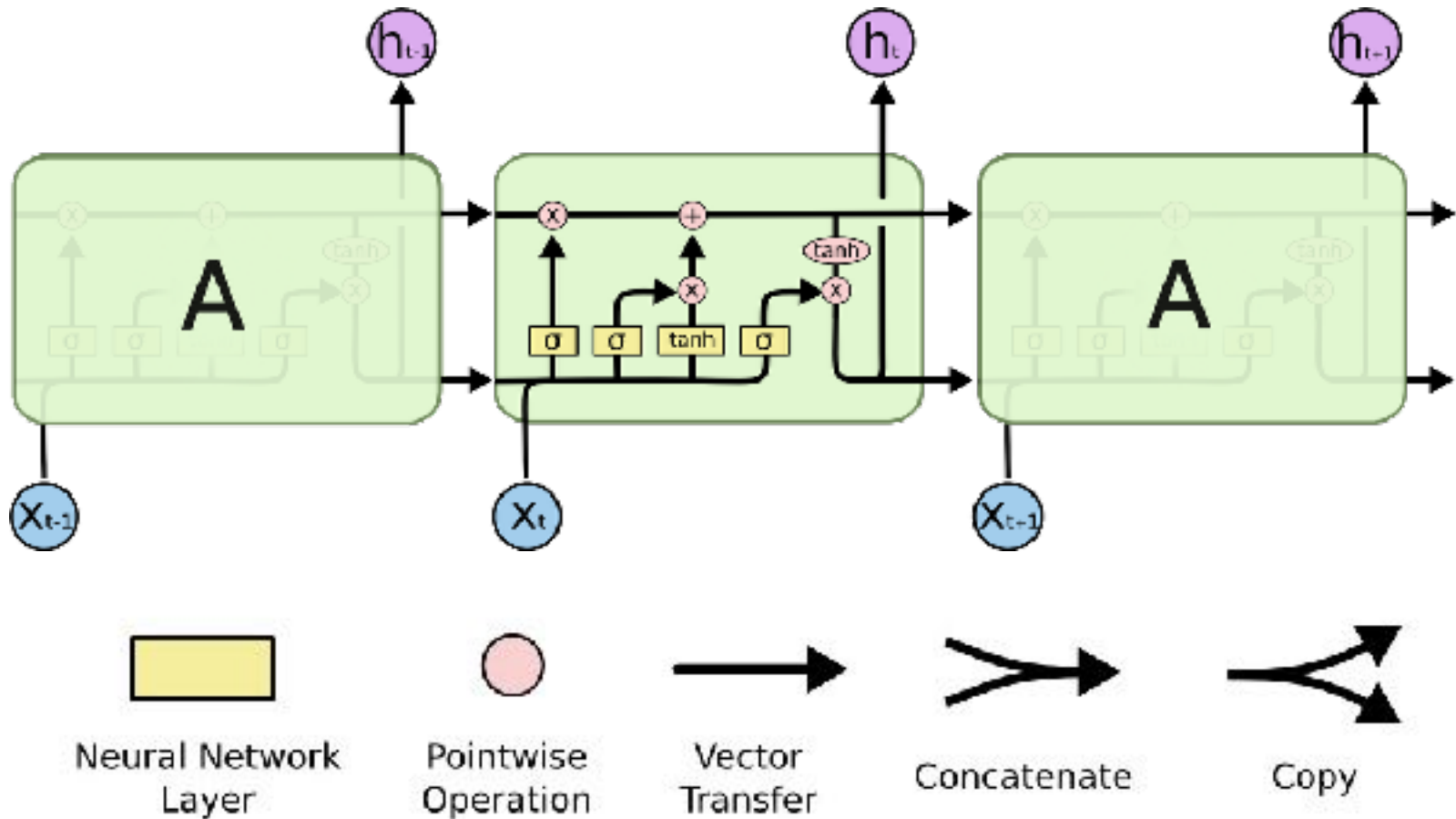
But there are also cases where we need more context.

From plain RNNs to LSTMs



(LSTM: Long Short Term Memory Networks)

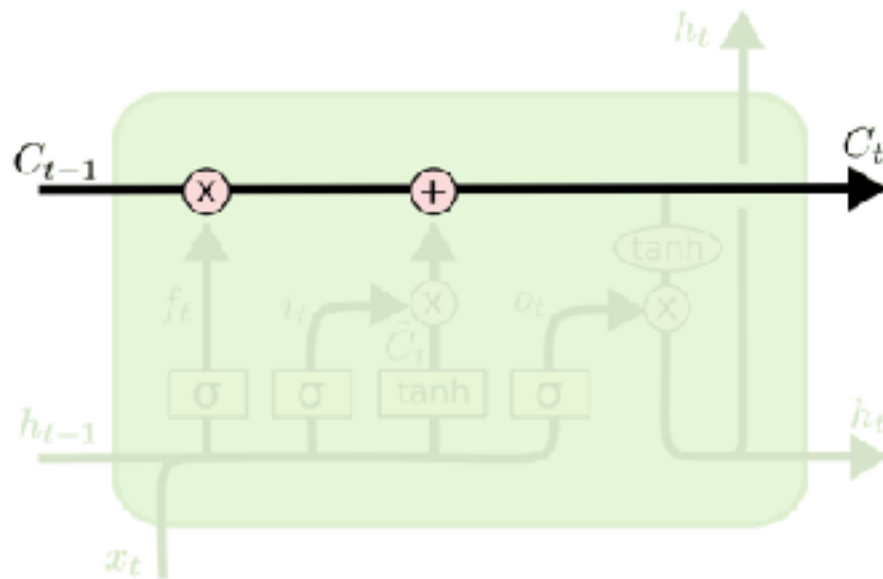
From plain RNNs to LSTMs



(LSTM: Long Short Term Memory Networks)

LSTMs Step by Step: Memory

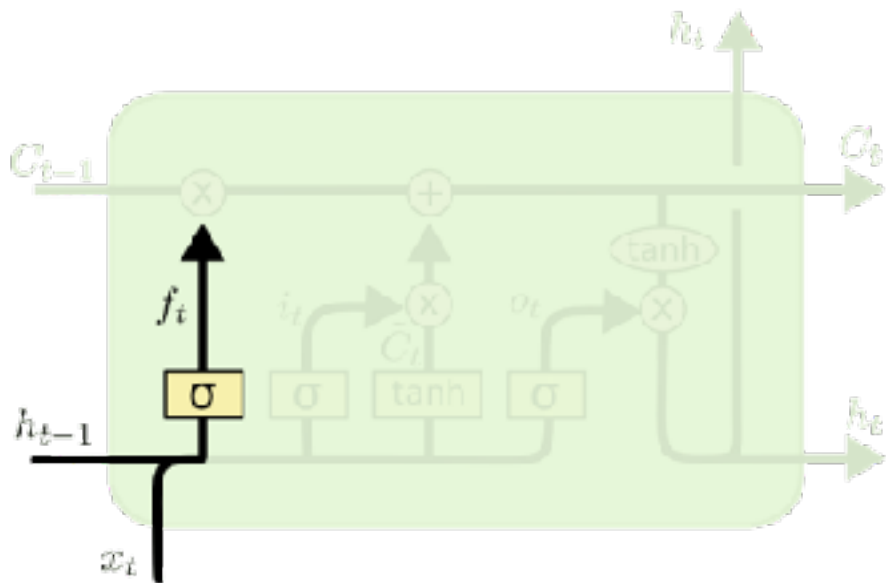
Cell State / Memory



The LSTM does have the ability to remove or add information to the cell state, carefully regulated by structures called gates

LSTMs Step by Step: Forget Gate

Should we continue to remember this “bit” of information or not?

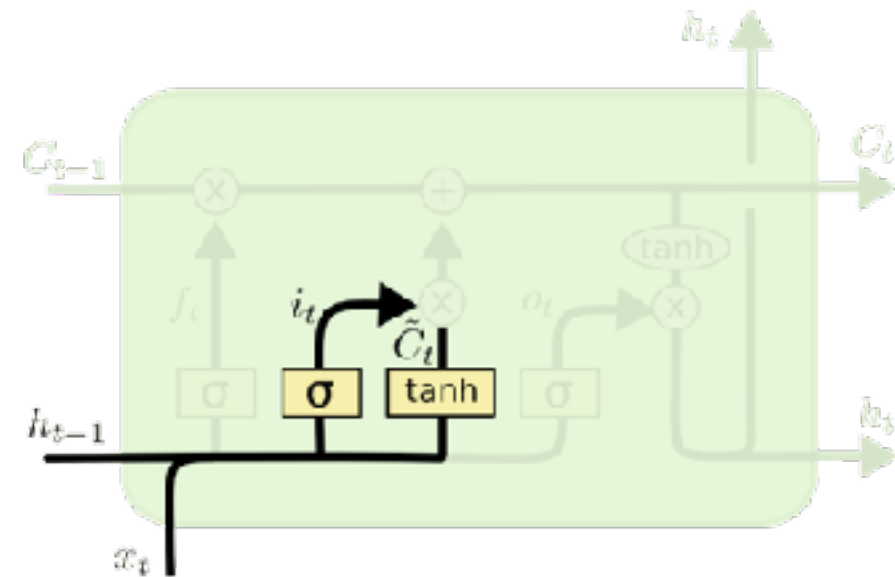


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

The first step in our LSTM is to decide what information we’re going to throw away from the cell state. This decision is made by a sigmoid layer called the “forget gate layer.”

LSTMs Step by Step: Input Gate

Should we update this “bit” of information or not? If so, with what?



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

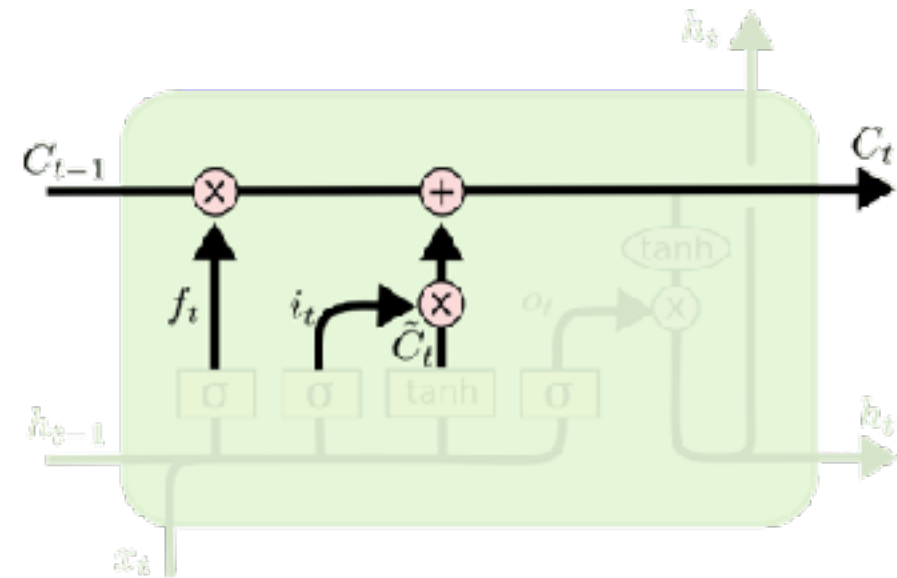
The next step is to decide what new information we’re going to store in the cell state. This has two parts. First, a sigmoid layer called the “input gate layer” decides which values we’ll update. Next, a tanh layer creates a vector of new candidate values, \tilde{C}_t , that could be

added to the state.

Credit: Christopher Olah

LSTMs Step by Step: Memory Update

Decide what will be kept in the cell state/memory

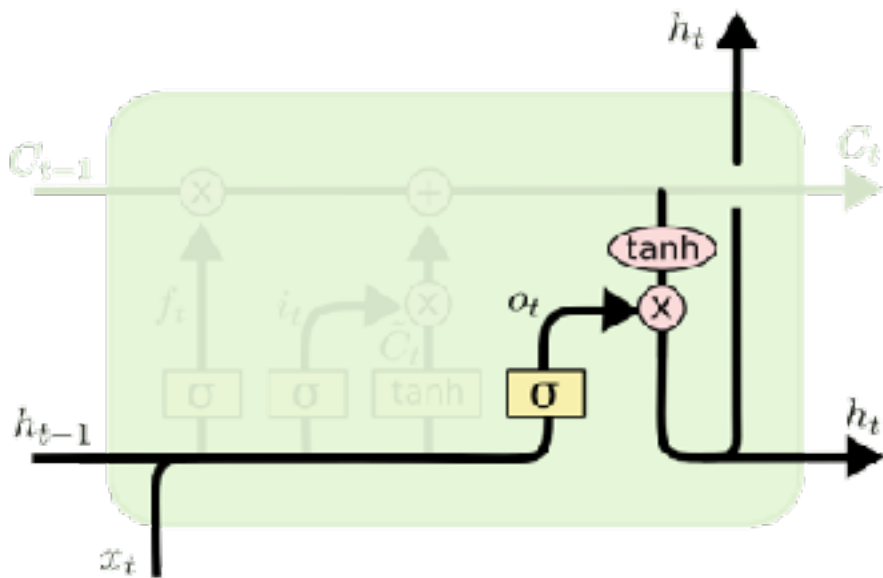


Forget that **Memorize this**

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

LSTMs Step by Step: Output Gate

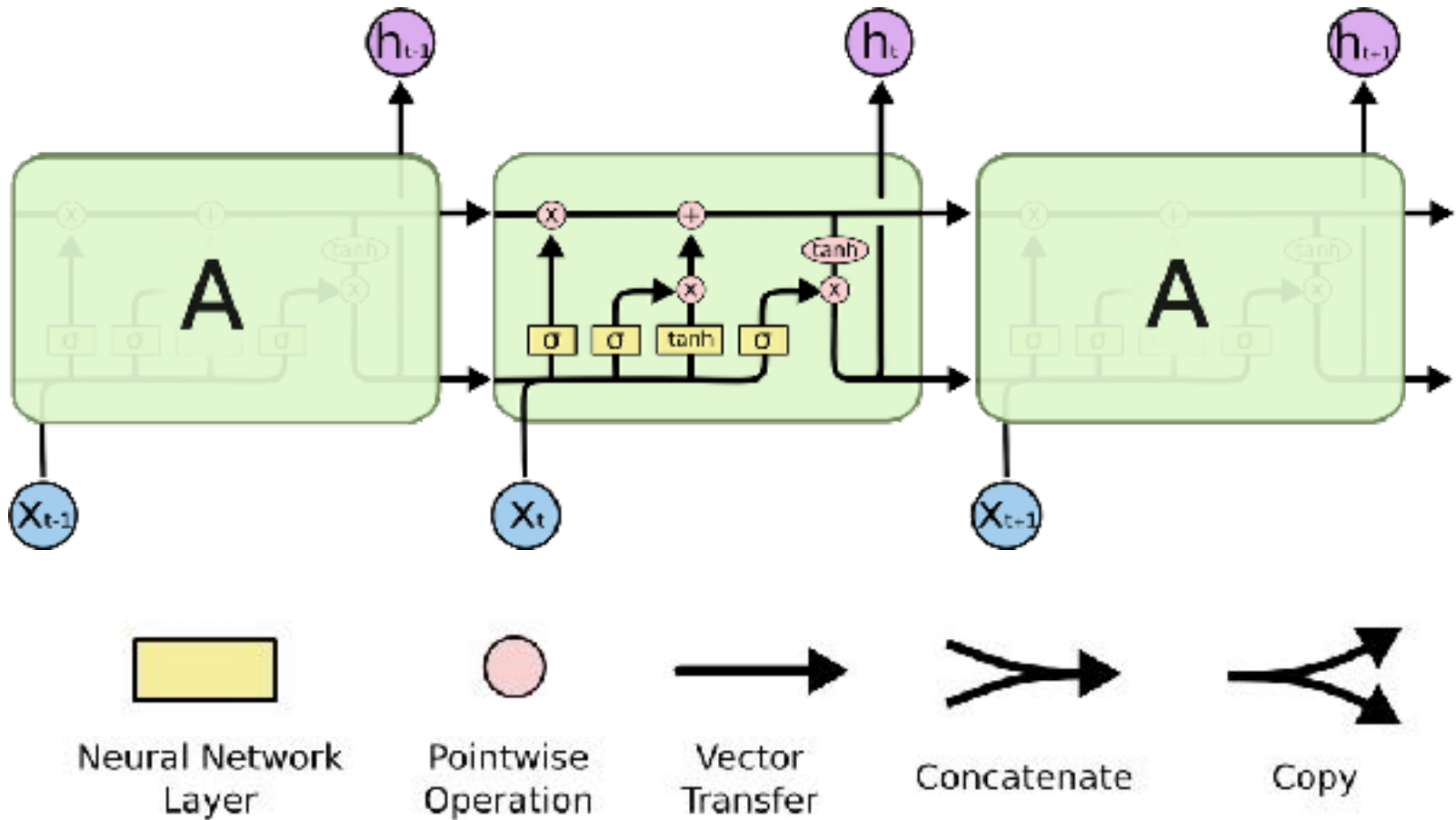
Should we output this “bit” of information?



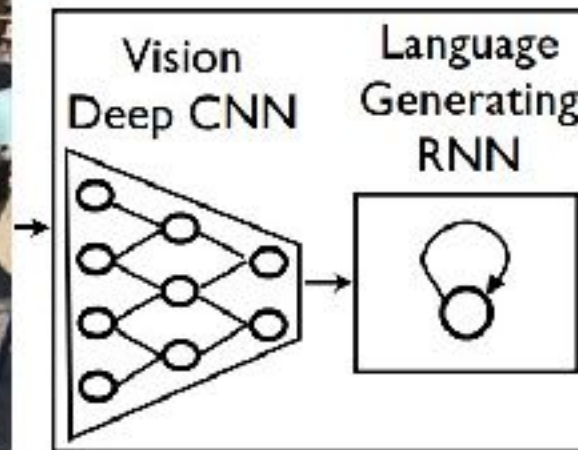
$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$
$$h_t = o_t * \tanh (C_t)$$

This output will be based on our cell state, but will be a filtered version. First, we run a sigmoid layer which decides what parts of the cell state we’re going to output. Then, we put the cell state through tanh (to push the values to be between -1 and 1) and multiply it by the output of the sigmoid gate, so that we only output the parts we decided to.

Complete LSTM - A pretty sophisticated cell



Show and Tell: A Neural Image Caption Generator



A group of people shopping at an outdoor market.

There are many vegetables at the fruit stand.

Show and Tell: A Neural Image Caption Generator

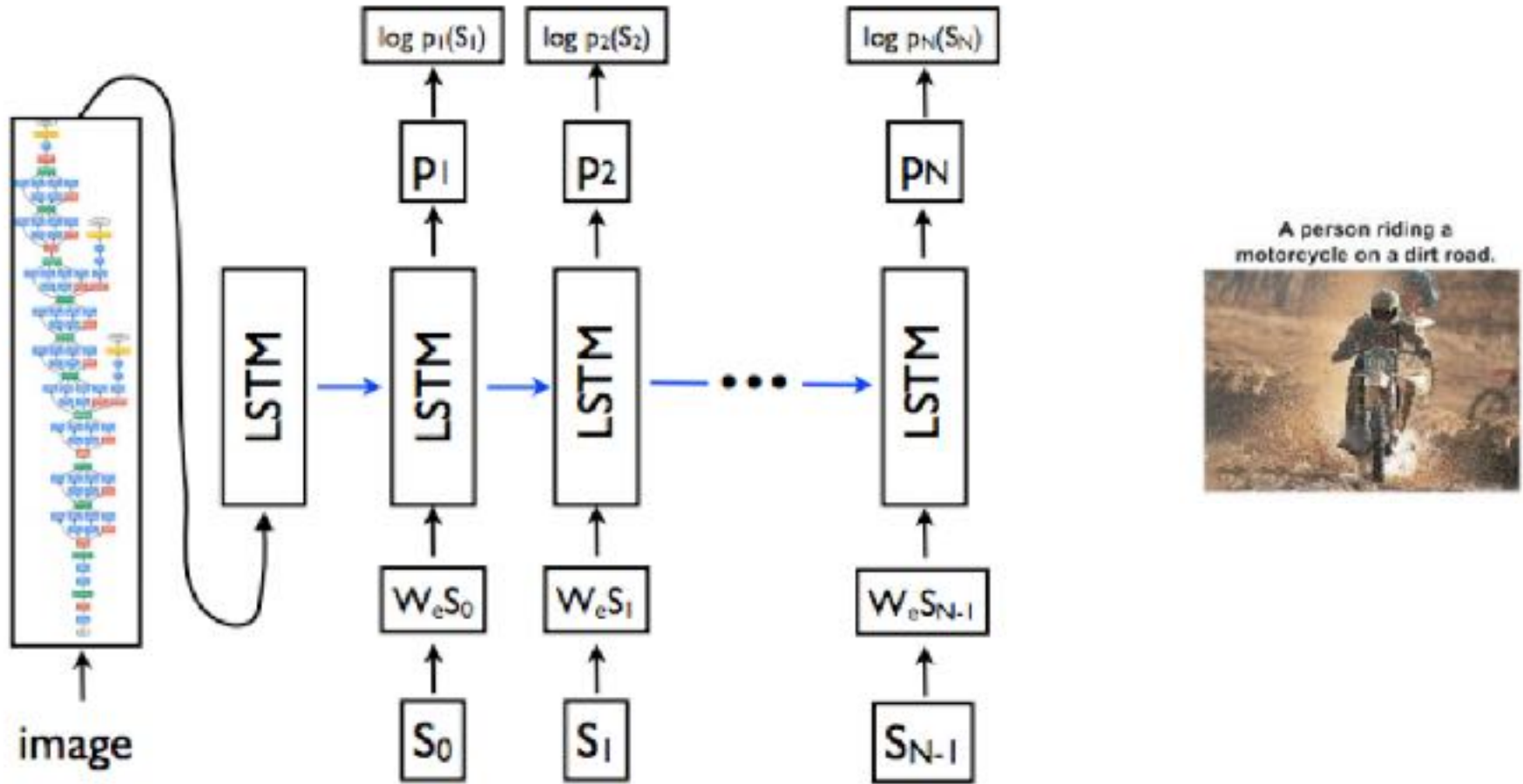


Image Caption Generator Results

A person riding a motorcycle on a dirt road.



Two dogs play in the grass.



A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A refrigerator filled with lots of food and drinks.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



A red motorcycle parked on the side of the road.



A yellow school bus parked in a parking lot.



Describes without errors

Describes with minor errors

Somewhat related to the image

Unrelated to the image

Cross-modal learning

Description (eg, Wikipedia article)

Snares penguin

From Wikipedia, the free encyclopedia

The **Snares penguin** (*Eudyptes robustus*), also known as the **Snares crested penguin** and the **Snares Islands penguin**, is a [penguin](#) from [New Zealand](#). The species breeds on [The Snares](#), a group of islands off the southern coast of the [South Island](#). This is a medium-small, yellow-crested penguin, at a size of 50–70 cm (19.5–27.5 in) and a weight of 2.5–4 kg (5.5–8.8 lb). It has dark blue-black upperparts and white underparts. It has a bright yellow eyebrow-stripe which extends over the eye to form a drooping, bushy crest. It has bare pink skin at the base of its large red-brown bill.

- Lots of descriptions/entries in Wikipedia available

Images



Zero-shot Learning

Description (eg, Wikipedia article)

Cardinal (bird)

From Wikipedia, the free encyclopedia

This article is about the bird family. For other uses, see [Cardinal](#).

Cardinals, in the family **Cardinalidae**, are [passerine birds](#) found in [North](#) and [South America](#). They are also known as cardinal-grosbeaks and cardinal-buntlings. The South American cardinals in the [genus *Paroaria*](#) are placed in another family, the [Thraupidae](#) (previously placed in [Emberizidae](#)).

Can we predict an image classifier from a description alone?

Assume:

- In training we have access to wiki articles and labeled images
- For test classes we only have wiki articles
- We want to classify a new image (it can belong to any class)

Zero-shot Learning

- Goal: learn to predict an image classifier from a description
- Linear binary 1-vs-all classifier:

$$y_c = w_c^T x$$

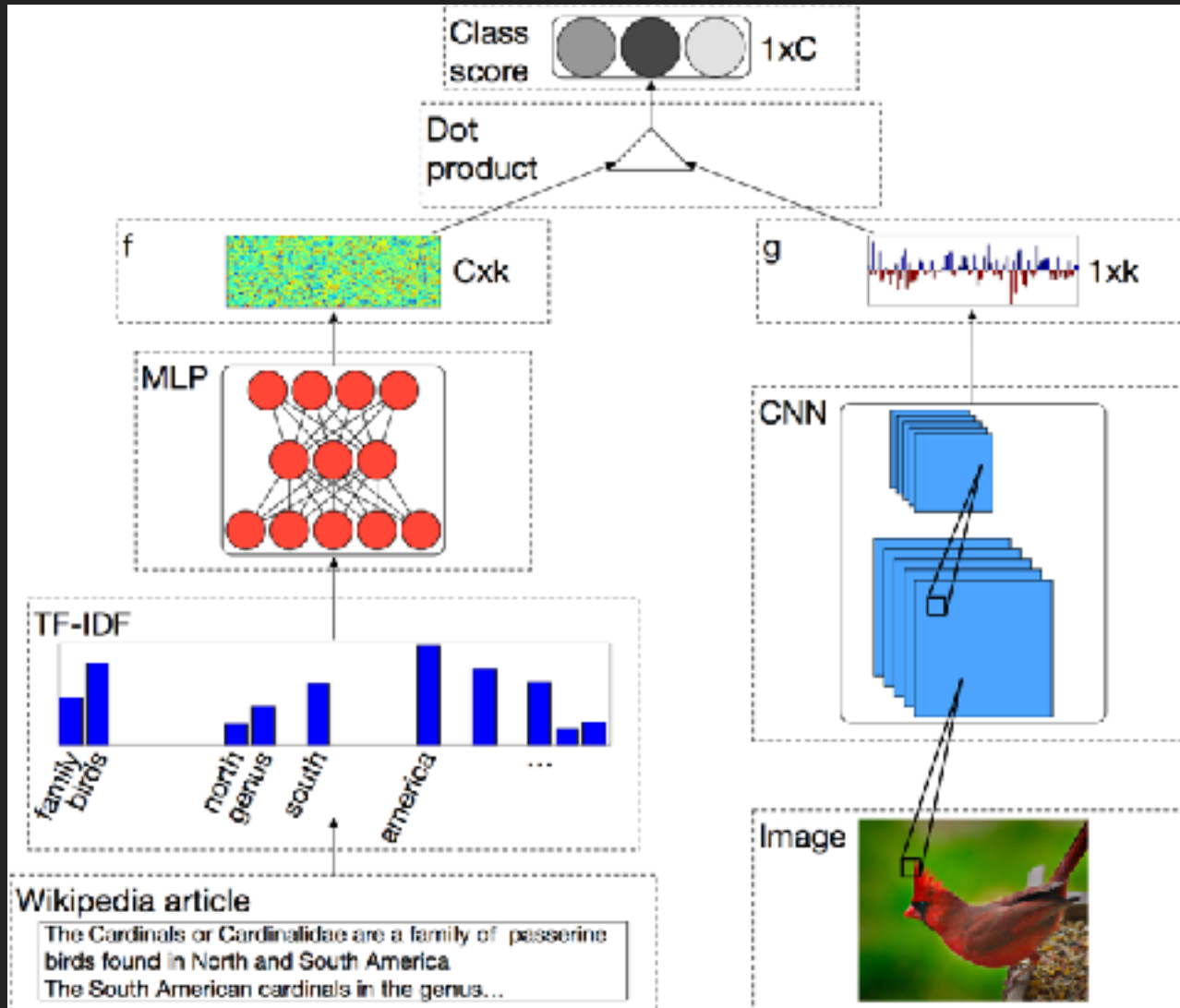
- x ... image feature vector
- w_c ... classifier weight vector for class c
- We are also given t_c , a vector representing a textual description about class c
- We want:

$$w_c = f_t(t_c)$$

- f_c ... a mapping $\mathbb{R}^p \rightarrow \mathbb{R}^d$ that transforms text features to the visual image feature space

Zero-shot Learning

- f_t can be a neural network



g used to compress x to a $k \ll d$ dim

x

Red faced Cormorant

The Red-faced Cormorant, Red-faced Shag or Violet Shag, *Phalacrocorax urile*, is a species of cormorant that is found in the far north of the Pacific Ocean and Bering Sea, from the eastern tip of Hokkaidō in Japan, via the Kuril Islands, the southern tip of the Kamchatka Peninsula and the Aleutian Islands to the Alaska Peninsula and Gulf of Alaska. The Red-faced Cormorant is closely related to the Pelagic Cormorant *P. pelagicus*, which has a similar range, and like the Pelagic Cormorant is placed by some authors (e.g. Johnsgaard) in a genus *Leucocarbo*. Where it nests alongside the Pelagic Cormorant, the Red-faced Cormorant generally breeds the more successfully of the two species, and it is currently increasing in numbers, at least in the easterly parts of its range. It is however listed as being of conservation concern{Verify source|date=September 2009}, partly because relatively little is so far known about it.

The adult bird has glossy plumage that is a deep greenish blue in colour, becoming purplish or bronze on the back and sides. In breeding condition it has a double crest,

.....



Red faced Cormorant

The Red-faced Cormorant, Red-faced Shag or Violet Shag, *Phalacrocorax urile*, is a species of cormorant that is found in the far north of the Pacific Ocean and Bering Sea, from the eastern tip of Hokkaidō in Japan, via the Kuril Islands, the southern tip of the Kamchatka Peninsula and the Aleutian Islands to the Alaska Peninsula and Gulf of Alaska. The Red-faced Cormorant is closely related to the Pelagic Cormorant *P. pelagicus*, which has a similar range, and like the Pelagic Cormorant is placed by some authors (e.g. Johnsgaard) in a genus *Leucocarbo*. Where it nests alongside the Pelagic Cormorant, the Red-faced Cormorant generally breeds the more successfully of the two species, and it is currently increasing in numbers, at least in the easterly parts of its range. It is however listed as being of conservation concern (Verify source|date=September 2009), partly because relatively little is so far known about it.

The adult bird has glossy plumage that is a deep greenish blue in colour, becoming purplish or bronze on the back and sides. In breeding condition it has a double crest,

.....

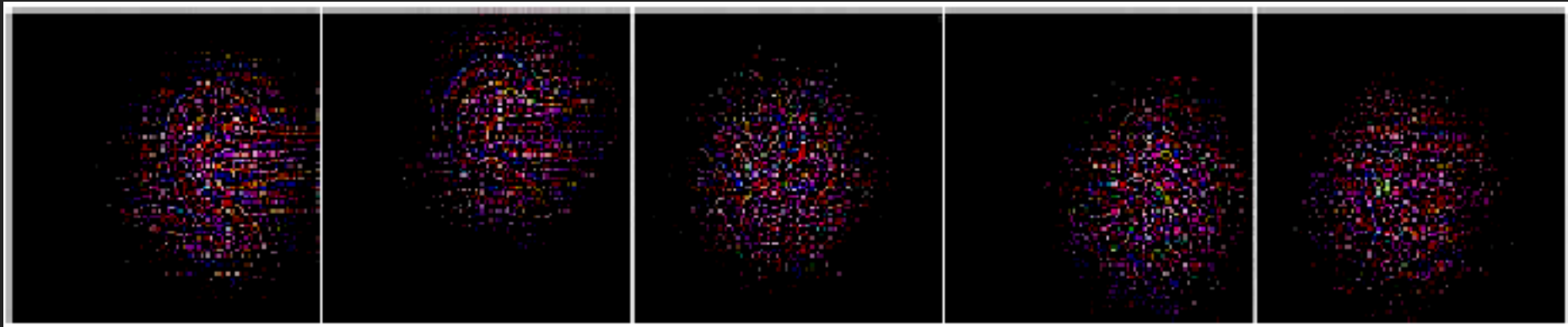


Red faced Cormorant

The Red-faced Cormorant, Red-faced Shag or Violet Shag, *Phalacrocorax urile*, is a species of cormorant that is found in the far north of the Pacific Ocean and Bering Sea, from the eastern tip of Hokkaidō in Japan, via the Kuril Islands, the southern tip of the Kamchatka Peninsula and the Aleutian Islands to the Alaska Peninsula and Gulf of Alaska. The Red-faced Cormorant is closely related to the Pelagic Cormorant *P. pelagicus*, which has a similar range, and like the Pelagic Cormorant is placed by some authors (e.g. Johnsgaard) in a genus *Leucocarbo*. Where it nests alongside the Pelagic Cormorant, the Red-faced Cormorant generally breeds the more successfully of the two species, and it is currently increasing in numbers, at least in the easterly parts of its range. It is however listed as being of conservation concern (Verify source|date=September 2009), partly because relatively little is so far known about it.

The adult bird has glossy plumage that is a deep greenish blue in colour, becoming purplish or bronze on the back and sides. In breeding condition it has a double crest,

.....

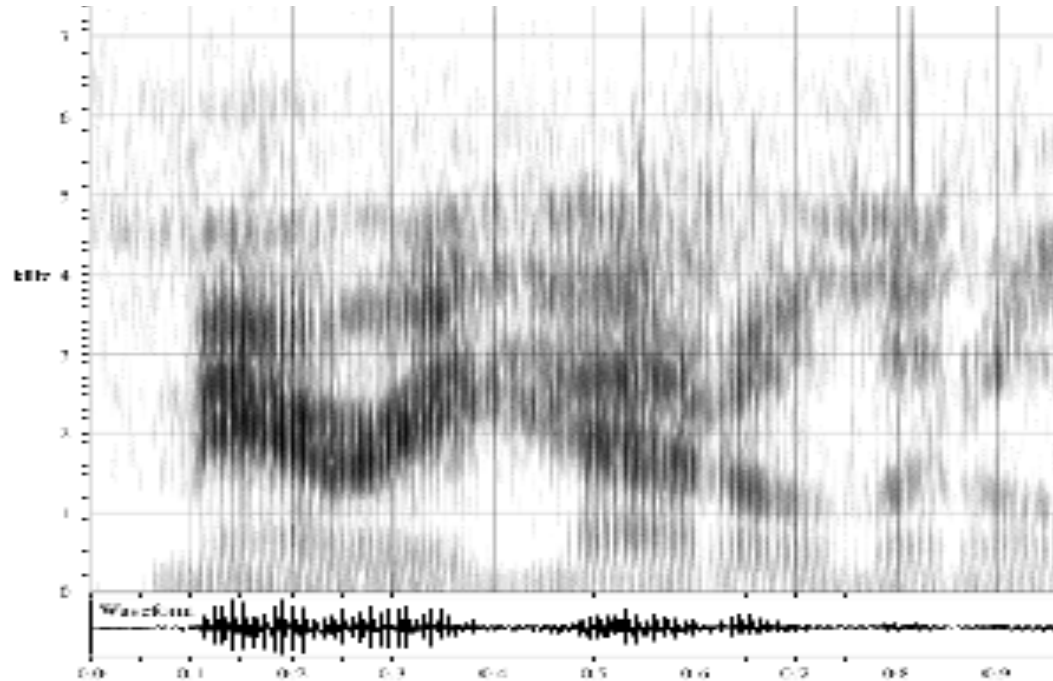




Raw pixels



Raw audio speech



- [Unsupervised Learning of Spoken Language with Visual Context](#). David Harwath, Antonio Torralba, James Glass. Advances in Neural Information Processing Systems (NIPS), 2016.

Unsupervised Learning of Spoken Language with Visual Context.

David Harwath



Jim Glass



[Unsupervised Learning of Spoken Language with Visual Context.](#) David Harwath, Antonio Torralba, James Glass.
Advances in Neural Information Processing Systems (NIPS), 2016.

Crowdsourcing Audio-Visual Data

Instructions

This HIT is a TURBIT audio-recording task. You should see either the HIT is visible and you require no permission. The media of the recording is limited to audio, video, and photos. You may not use any other content. You may not use any other content. You may not use any other content. You may not use any other content.

To complete this task, you must be:


- using a computer equipped with a microphone
- using the Chrome web browser
- in a relatively quiet environment

If your microphone is on and working, the volume meter at the right should move as you speak (after you grant permission for the site to use your microphone). Underneath the microphone volume meter you can see whether you are connected to server for recording. If you become disconnected, please continue recording after a connection is reestablished.

You will be presented with 4 image scenes. For each image, please:

- Press the **Record** button next to the image and then describe the image as if you were describing it to a blind person. During recording, the record button will be replaced with a stop button, and the recording by pressing the **Stop** button next to the image.
- After you record a copy, we will process the recording. If it is acceptable, it will be marked as **Good**; otherwise, the sentence will be marked with a **Bad** and you must redo the recording of that sentence to complete the task.
- After all 3 descriptions have been accepted, the submit button at the bottom of the page will be enabled.

Here's an example of the level of detail we're looking for:



"A man and a woman riding a a bench on top of an elephant. The woman is wearing a pink shirt and a hat. The elephant is standing on a dirt road in front of an old stone structure."

We're looking for a couple of sentences per image. You can talk about specific objects, locations, shapes, colors, etc. in the image.

Poor quality work will be rejected and you will be blocked from completing any more of our HITs.

Please record a description of each image below.

Record



Record



Record



Record



Crowdsourcing Audio-Visual Data



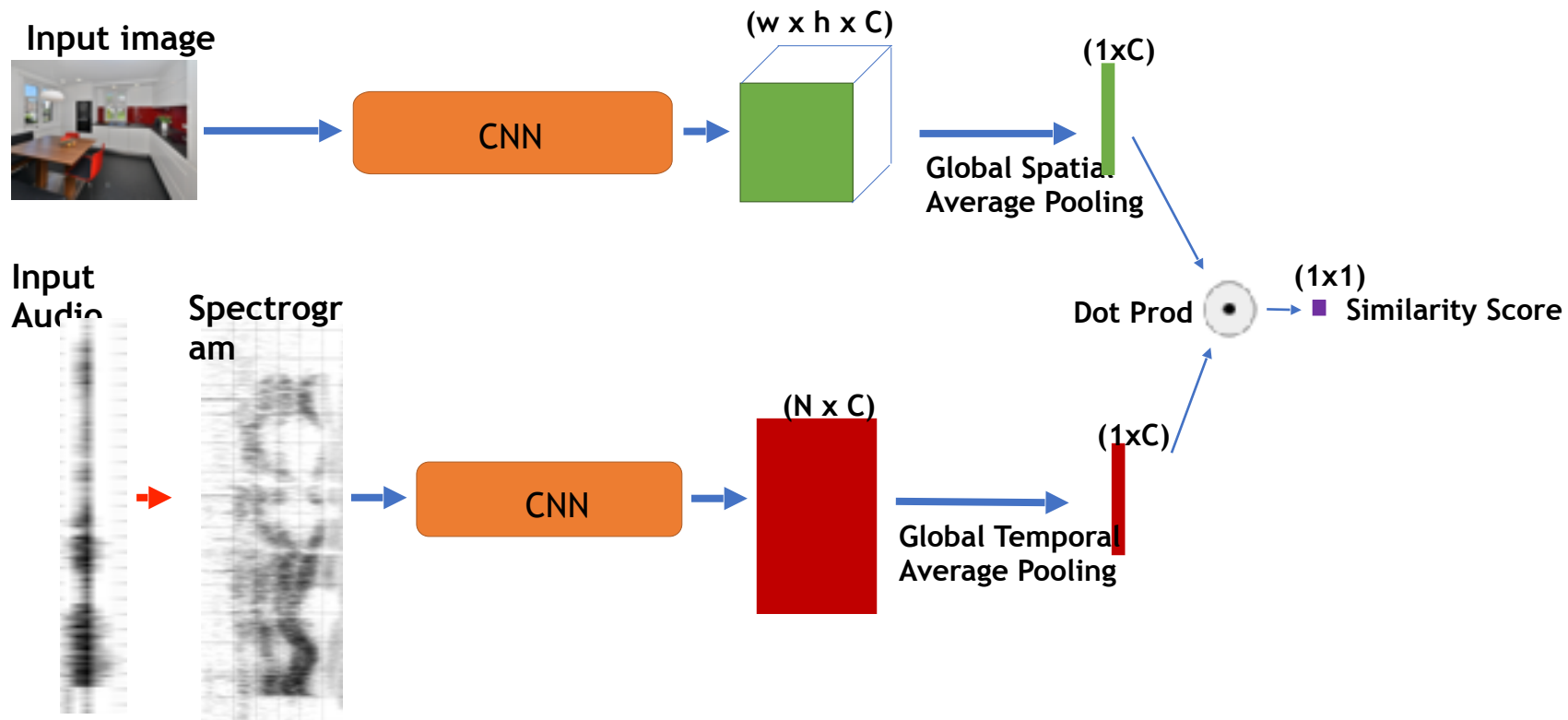
382.060 Speech descriptions on Images from Places dataset.

Crowdsourcing Audio-Visual Data

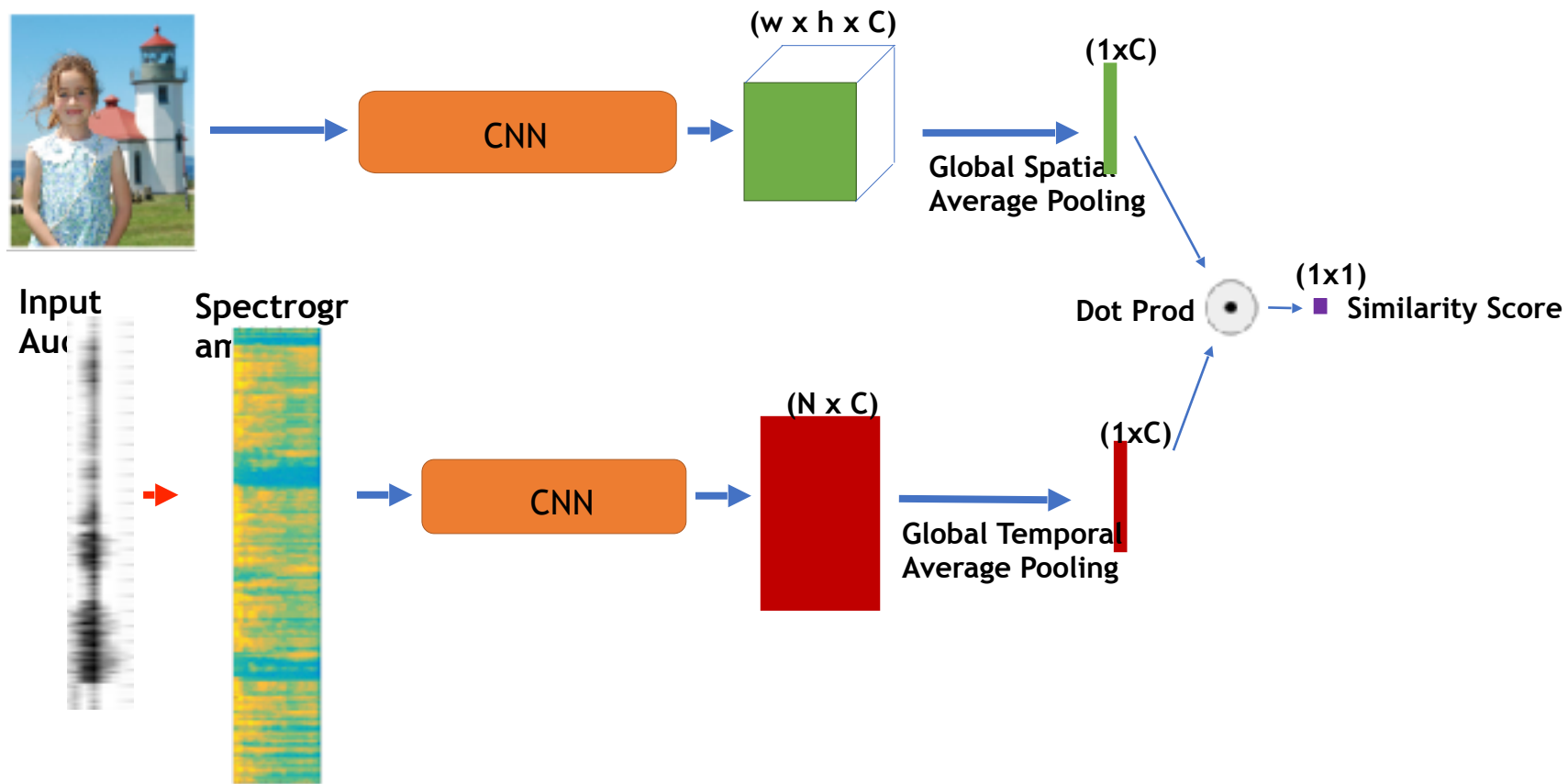


82.060 Speech descriptions on
images from Places dataset.

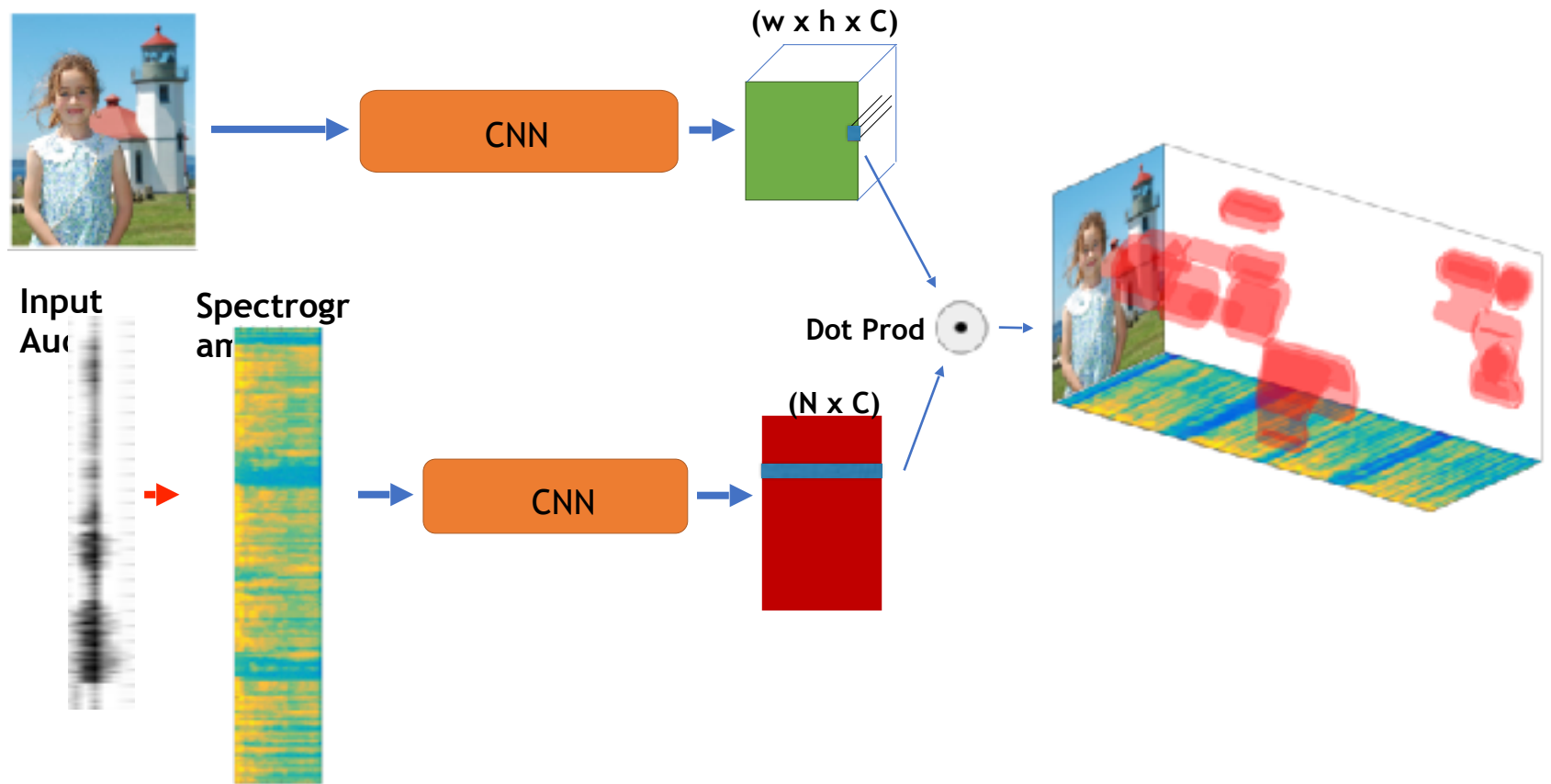
Joint Audio-Visual Architecture



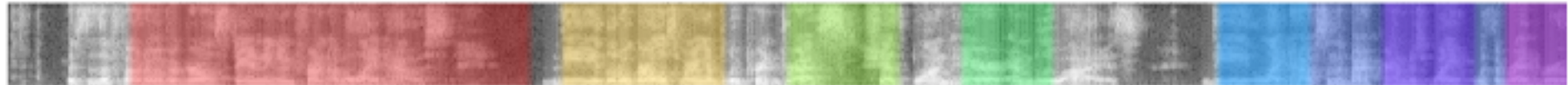
Joint Audio-Visual Architecture



Joint Audio-Visual Architecture



Co-segmentation of speech and image



THE IS A PHOTO OF A GIRL STANDING IN FRONT OF A LIGHTHOUSE THE LITTLE GIRL WEARS A BLUE DRESS SHE HAS BLOND HAIR AND BLUE EYES THE LIGHTHOUSE IN THE BACKGROUND IS WHITE WITH A RED ROOF



Here in this pic there some skier

going up a
mountain















Words

- Need ways to compare words

Next to the 'sofa' is a desk, and a 'person' is sitting behind it.

'armchair'	'man'
'bench'	'woman'
'chair'	'child'
'deck chair'	'teenager'
'ottoman'	'girl'
'seat'	'boy'
'stool'	'baby'
'swivel chair'	'daughter'
'loveseat'	'son'
...	...

Encoding words into vectors

- Need ways to compare words



So that if two words i and j are similar then w_i and w_j are close

Encoding words into vectors

- Need ways to compare words

One-hot representations

'sofa' \longrightarrow $[1, 0, 0, 0, \dots 0]_V$

'person' \longrightarrow $[0, 1, 0, 0, \dots 0]_V$

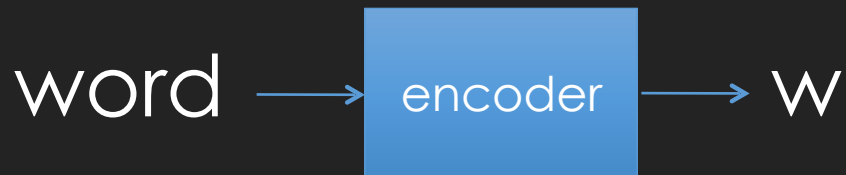
'car' \longrightarrow $[0, 0, 1, 0, \dots 0]_V$

'tree' \longrightarrow $[0, 0, 0, 1, \dots 0]_V$

1-of- V coding, where V is size of the vocabulary

word2vec

- Find better vector encodings



'sofa' → w1

'person' → w2

'car' → w3

'tree' → w4

So that if two words i and j are similar then w_i and w_j are close

But we do not have word similarities...

How do we learn the vectors?

We will use a different task, and hope that similarity will emerge...

We will train a classifier to predict the words surrounding each word.

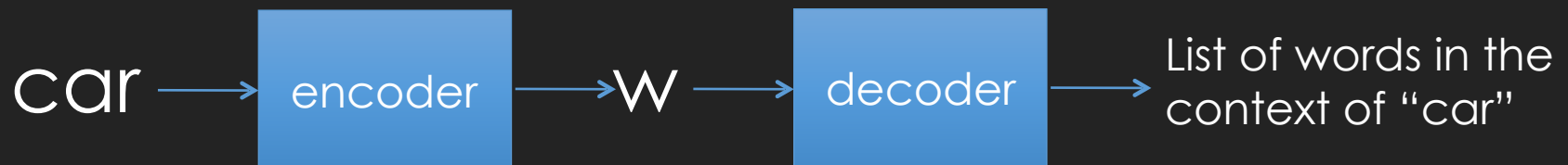
word2vec

I parked the **car** in a nearby street. It is a red **car** with two doors, ...

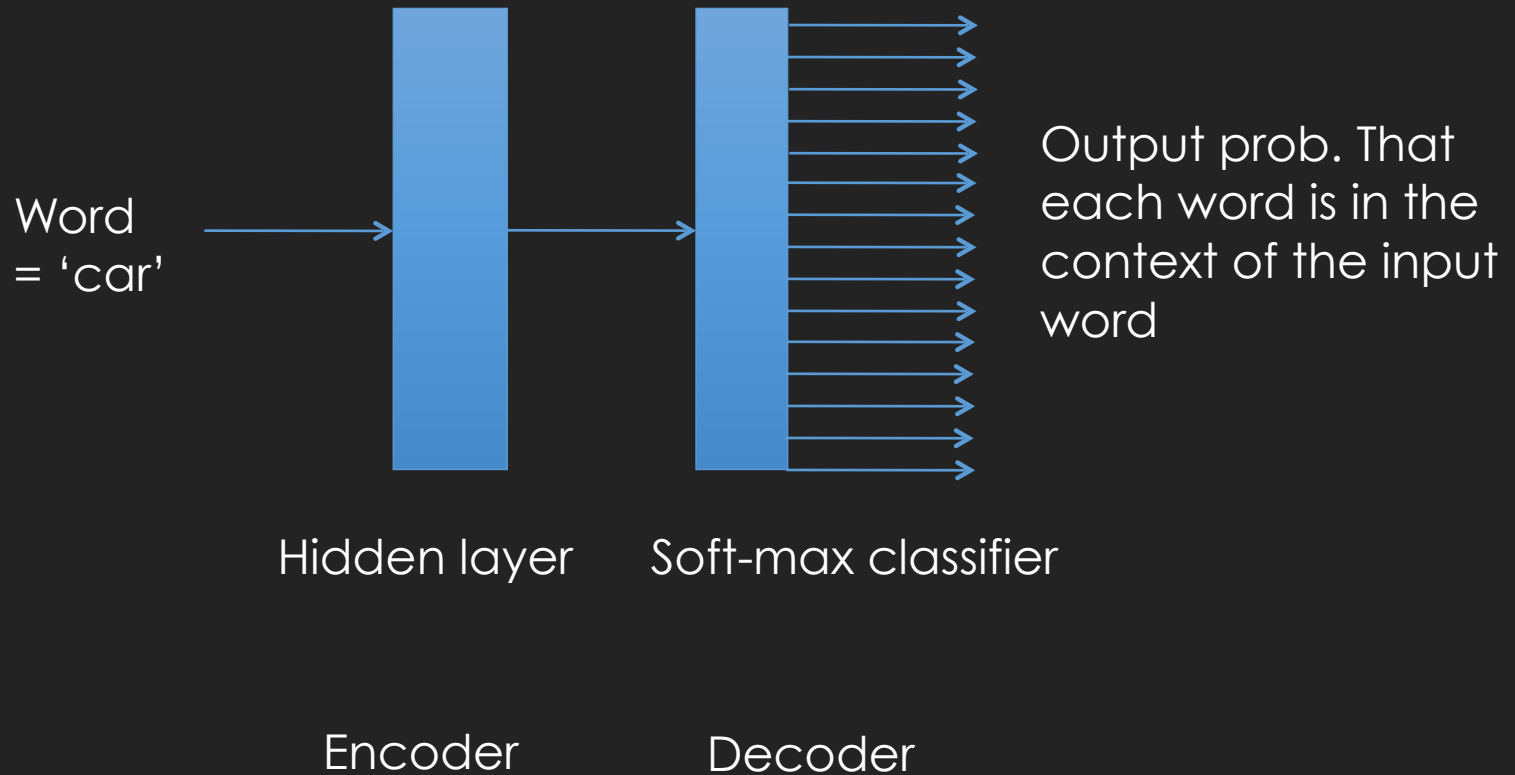
I parked the **vehicle** in a nearby street...

word2vec

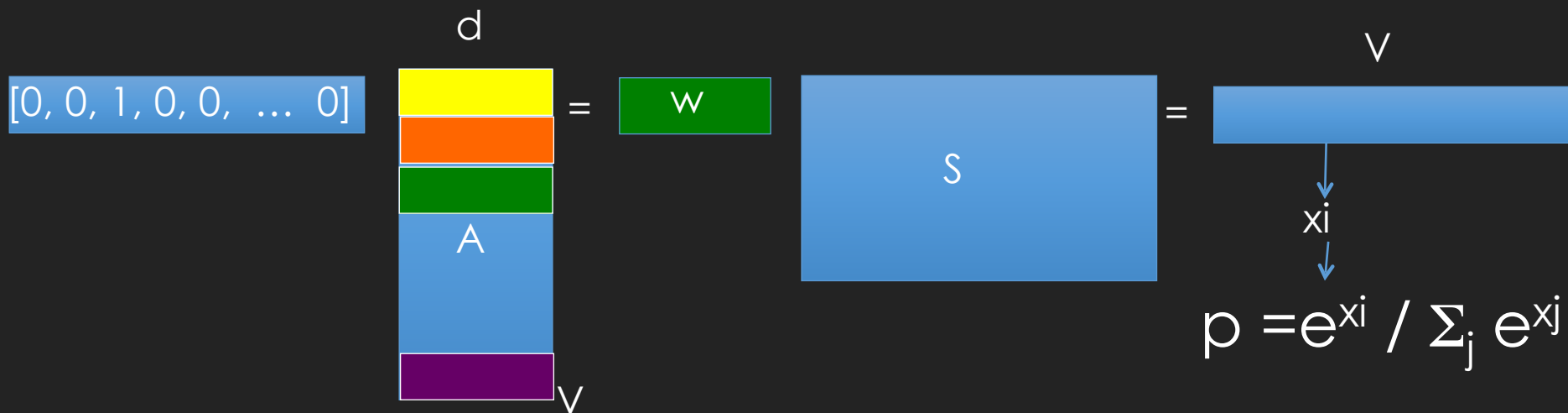
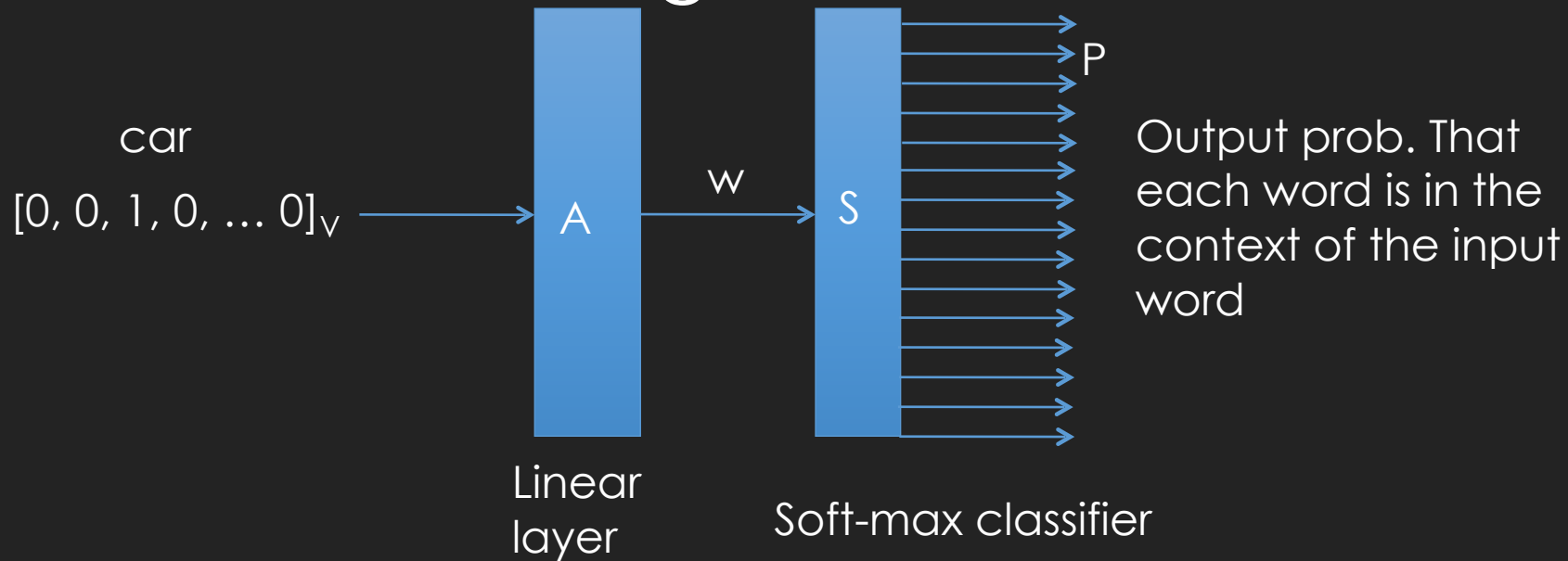
I parked the **car** in a nearby street. It is a red **car** with two doors, ...



word2vec



word2vec, training

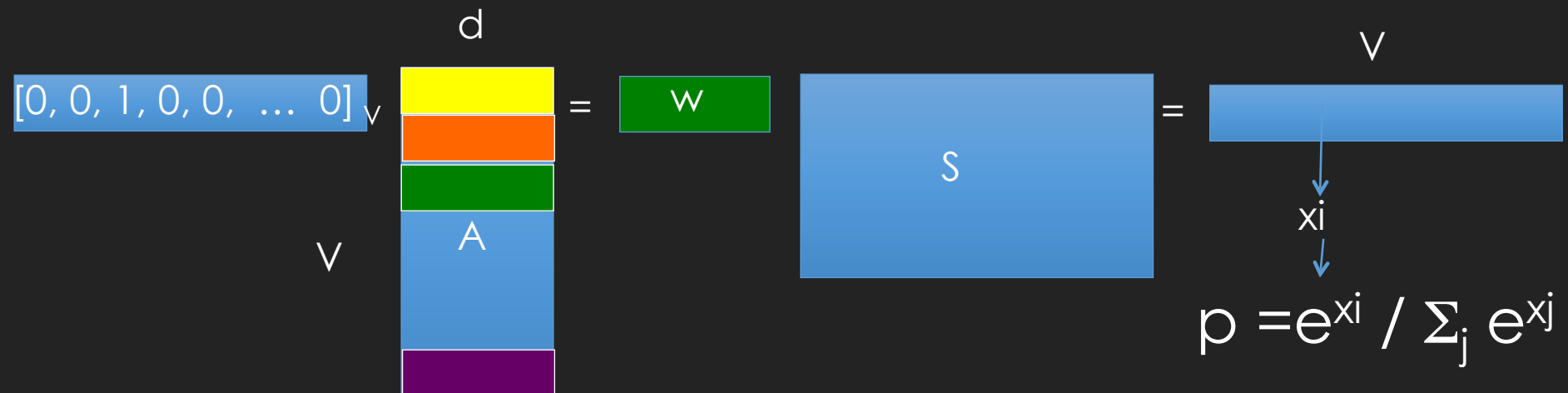


word2vec, training

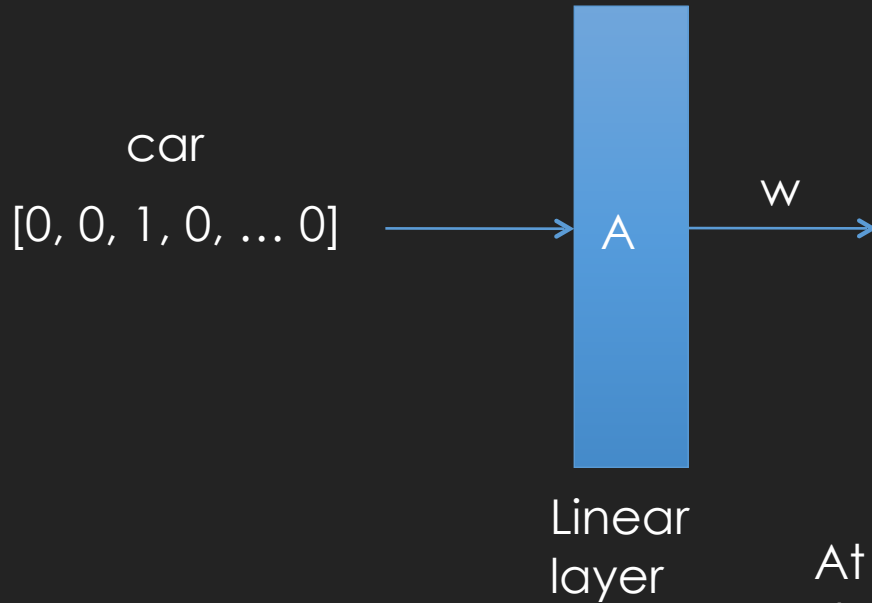
- In training maximize log-likelihood over the training set:

$$\sum_{t=1}^T \sum_{i=-c}^c \log p(w_{t+i} | w_t)$$

T ... training set size
c ... context window size



word2vec, test time



At test time, w is our word embedding.
The encoding is just a look up table.



Algebraic operations with the vector representation of words

$$X = \text{Vector}(\text{"Paris"}) - \text{vector}(\text{"France"}) + \text{vector}(\text{"Italy"})$$

Closest nearest neighbor to X is $\text{vector}(\text{"Rome"})$

Remember: Subtitle – Sentence Similarity



01:00:58 --> 01:01:03

I'm telling you, it's spooky. She knows more about you than you do.

01:01:03 --> 01:01:05

Who doesn't?

Subtitle – Sentence Similarity



She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

Sentences

Skip-Thought Vectors

training corpus: 11K books

.....

They called from outside.

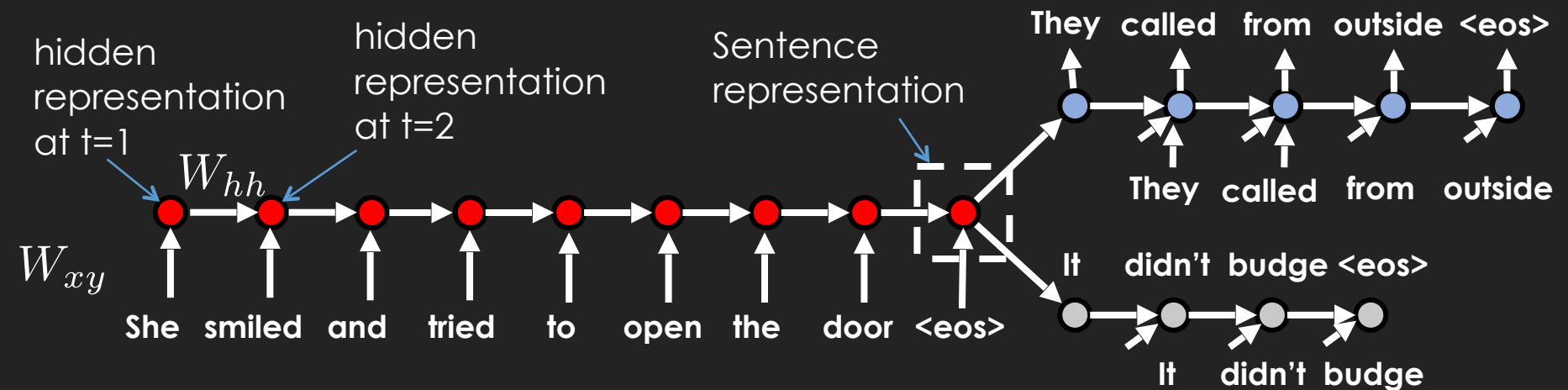
She smiled and tried to open the door.

It didn't budge.

.....

Sentences

Skip-Thought Vectors

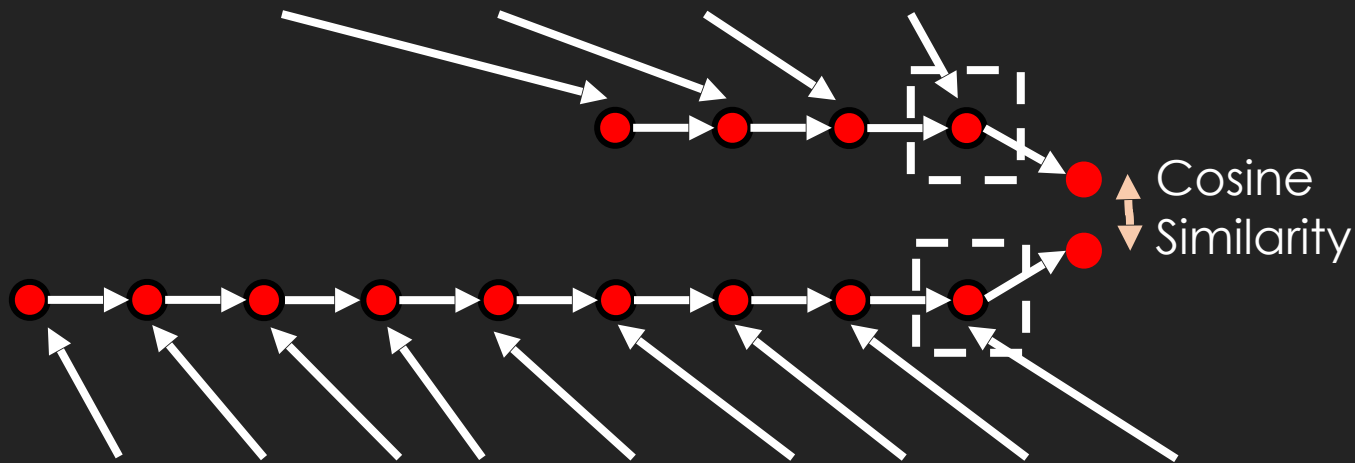


Sentences

On test time:

Sentence 1:

Take the staircase, quick!



Sentence 2:

They sped up a staircase to the third floor.

A photograph of a library or bookstore. The shelves are filled with books of various colors and sizes. In the foreground, there are several tall stacks of books, some with their spines visible. The background shows more shelves extending into the distance, creating a sense of a vast collection of literature.

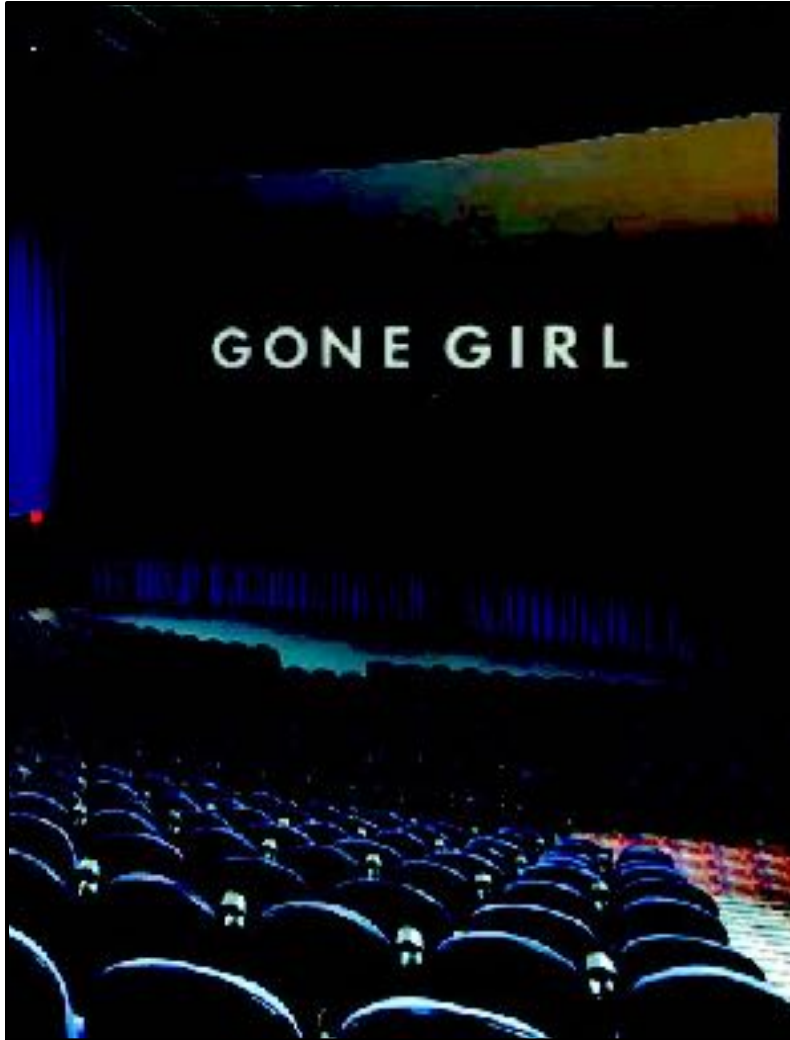
Books Contain Rich Descriptions

But lack rich visual content

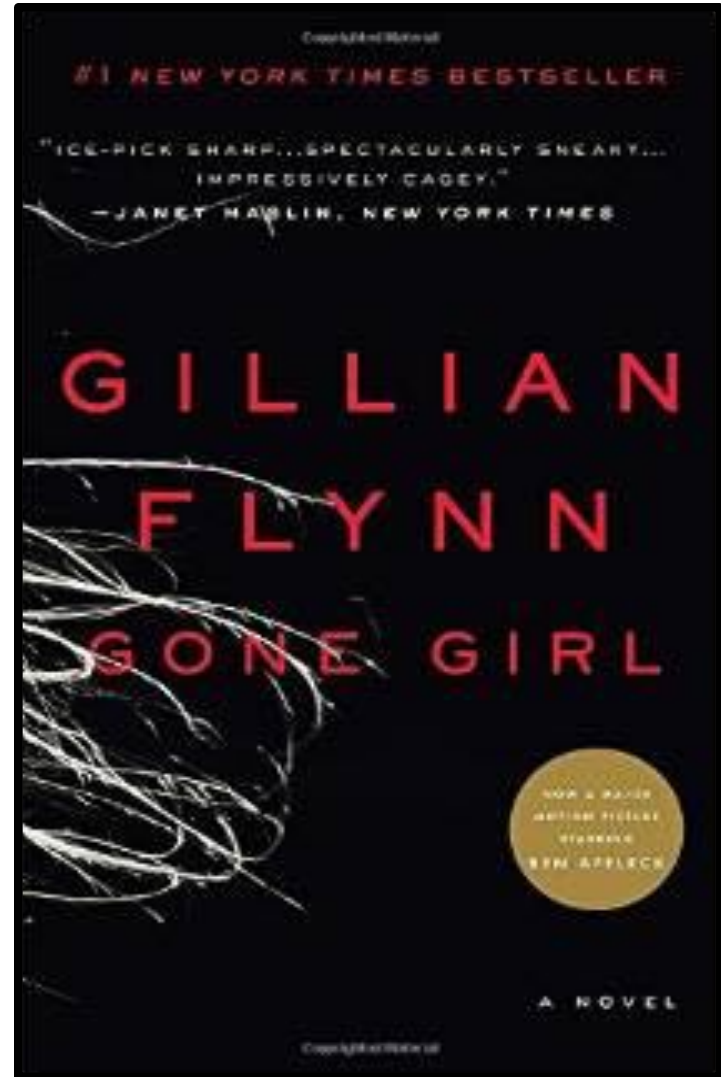
Movies contain rich visual content



Lots of paired books and movies



movie



book

Book: Tells A Story

As I walked toward the bar across the concrete parking lot, I looked straight down the road and saw the river. Moving apace with the river was a long single line of men, eyes aimed at their feet, walking steadfastly nowhere. I felt an immediate need to get inside.



Movie: Visualizes A Story



As I walked toward the bar across the concrete parking lot, I looked straight down the road and saw the river. Moving apace with the river was a long single line of men, eyes aimed at their feet, walking steadfastly nowhere. I felt an immediate need to get inside.

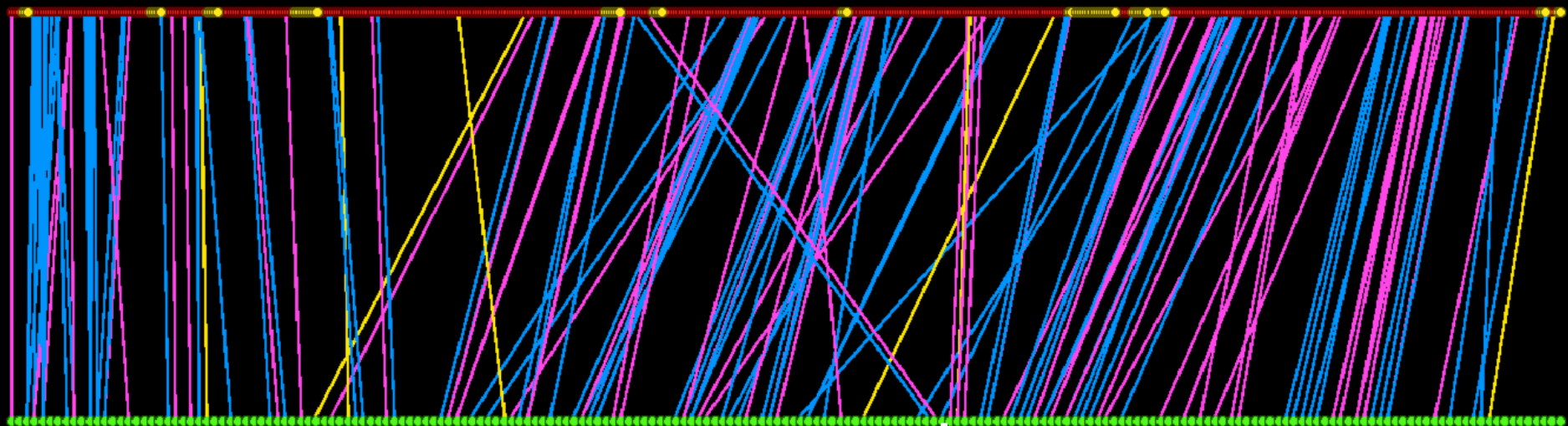


As I walked toward the bar across the concrete parking lot, I looked straight down the road and saw the river. Moving apace with the river was a long single line of men, eyes aimed at their feet, walking steadfastly nowhere. I felt an immediate need to get inside.

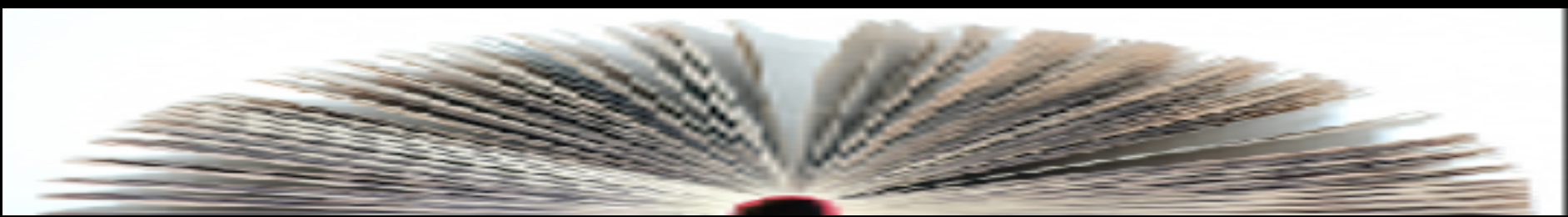




time



paragraphs



— visual match (106)
— dialog match (76)

— different (7)
● not in movie

Sentence Query Example

Query:

- He drove down the street off into the distance.

Sentence Query Example

Query:

- He drove down the street off into the distance.

Top Retrieved Sentences using Skip-Thoughts:

- He started the car, left the parking lot and merged onto the highway a few miles down the road.
- She watched the lights flicker through the trees as the men drove toward the road.

Skip-Thought Vectors

Near state-of-the-art results on standard NLP tasks:

- Semantic relatedness
- Paraphrase detection
- Image-sentence ranking
- Movie review sentiment prediction
- Question type classification

Cross-modal learning

Description (eg, Wikipedia article)

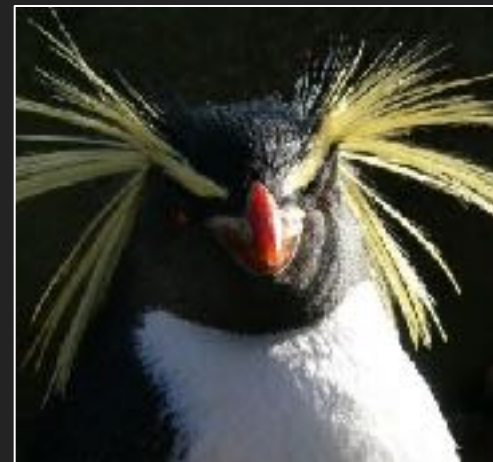
Snares penguin

From Wikipedia, the free encyclopedia

The **Snares penguin** (*Eudyptes robustus*), also known as the **Snares crested penguin** and the **Snares Islands penguin**, is a penguin from **New Zealand**. The species breeds on **The Snares**, a group of islands off the southern coast of the **South Island**. This is a medium-small, yellow-crested penguin, at a size of 50–70 cm (19.5–27.5 in) and a weight of 2.5–4 kg (5.5–8.8 lb). It has dark blue-black upperparts and white underparts. It has a bright yellow eyebrow-stripe which extends over the eye to form a drooping, bushy crest. It has bare pink skin at the base of its large red-brown bill.

- Lots of descriptions/entries in Wikipedia available

Images



Aligning Movies with Books



01:00:58 --> 01:01:03

I'm telling you, it's spooky. She knows more about you than you do.

01:01:03 --> 01:01:05

Who doesn't?

01:01:06 --> 01:01:08

What's happening?

01:01:08 --> 01:01:13

The staircases change, remember?



... As night fell, the promised storm blew up around them. Spray from the high waves splattered the walls of the hut and a fierce wind rattled the filthy windows.

Aunt Petunia found a few moldy blankets in the second room and made up a bed for Dudley on the moth-eaten sofa.

She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

The storm raged more and more ferociously as the night went on, and Harry couldn't sleep.

Aligning Movies with Books



01:00:58 --> 01:01:03

I'm telling you, it's spooky. She knows more about you than you do.



01:01:03 --> 01:01:05

Who doesn't?



01:01:06 --> 01:01:08

What's happening?



01:01:08 --> 01:01:13

The staircases change, remember?



... As night fell, the promised storm blew up around them. Spray from the high waves splattered the walls of the hut and a fierce wind rattled the filthy windows.

Aunt Petunia found a few moldy blankets in the second room and made up a bed for Dudley on the moth-eaten sofa.

She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

The storm raged more and more ferociously as the night went on, and Harry couldn't sleep.

Aligning Movies with Books



01:00:58 --> 01:01:03

I'm telling you, it's spooky. She knows more about you than you do.



01:01:03 --> 01:01:05

Who doesn't?



01:01:06 --> 01:01:08

What's happening?



01:01:08 --> 01:01:13

The staircases change, remember?



... As night fell, the promised storm blew up around them. Spray from the high waves splattered the walls of the hut and a fierce wind rattled the filthy windows.

Aunt Petunia found a few moldy blankets in the second room and made up a bed for Dudley on the moth-eaten sofa.

She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

The storm raged more and more ferociously as the night went on, and Harry couldn't sleep.

Aligning Movies with Books



01:00:58 --> 01:01:03

I'm telling you, it's spooky. She knows more about you than you do.

01:01:03 --> 01:01:05


Who doesn't?

Shot – Sentence Similarity



She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

Aligning Movies with Books

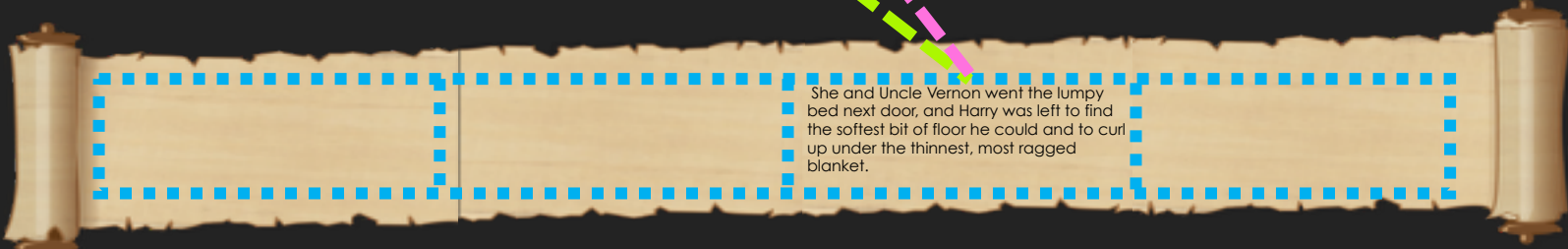


01:00:58 --> 01:01:03
I'm telling you, it's spooky. She knows more about you than you do.

01:01:03 --> 01:01:05
Who doesn't?

Video– Sentence Similarity

Subtitle – Sentence Similarity



She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

Our Method



01:00:58 --> 01:01:03

I'm telling you, it's spooky. She knows more about you than you do.



01:01:03 --> 01:01:05

Who doesn't?



01:01:06 --> 01:01:08

What's happening?

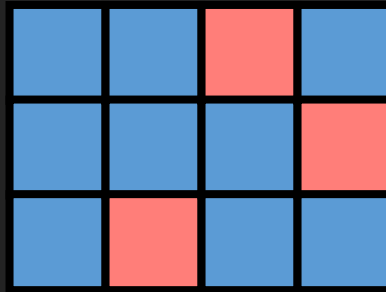


01:01:08 --> 01:01:13


The staircases change, remember?

Shot – Sentence Similarity Matrix

shots in movie



sentences in book



... As night fell, the promised storm blew up around them. Spray from the high waves splattered the walls of the hut and a fierce wind rattled the filthy windows.

Aunt Petunia found a few moldy blankets in the second room and made up a bed for Dudley on the moth-eaten sofa.

She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

The storm raged more and more ferociously as the night went on, and Harry couldn't sleep.

Subtitle – Sentence Similarity



01:00:58 --> 01:01:03

I'm telling you, it's spooky. She knows more about you than you do.

01:01:03 --> 01:01:05

Who doesn't?

- BLEU scores
- Longest subsequence matching + TF-IDF
- Sentence embedding (Skip-Thoughts)

Subtitle – Sentence Similarity



She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

Subtitle – Sentence Similarity



01:00:58 --> 01:01:03

I'm telling you, it's spooky. She knows more about you than you do.

01:01:03 --> 01:01:05

Who doesn't?

Subtitle – Sentence Similarity



She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

Video – Sentence Similarity



01:00:58 --> 01:01:03
I'm telling you, it's spooky. She knows more about you than you do.

01:01:03 --> 01:01:05
Who doesn't?

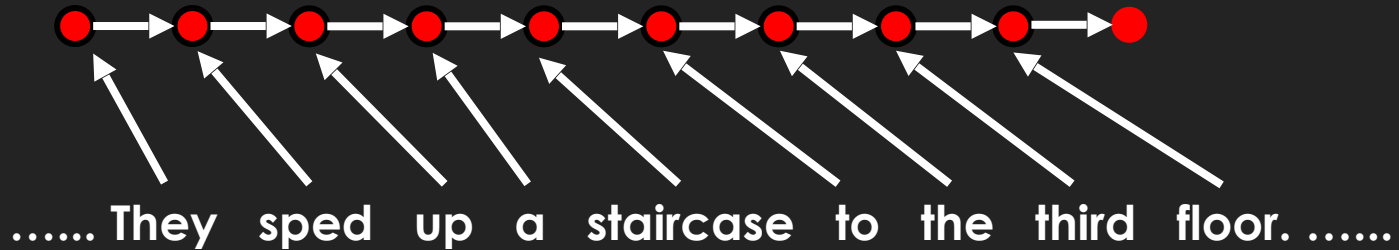
- Visual Semantic Embedding

Video– Sentence Similarity



She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

Video – Sentence Similarity



... As night fell, the promised storm blew up around them. Spray from the high waves splattered the walls of the hut and a fierce wind rattled the filthy windows.

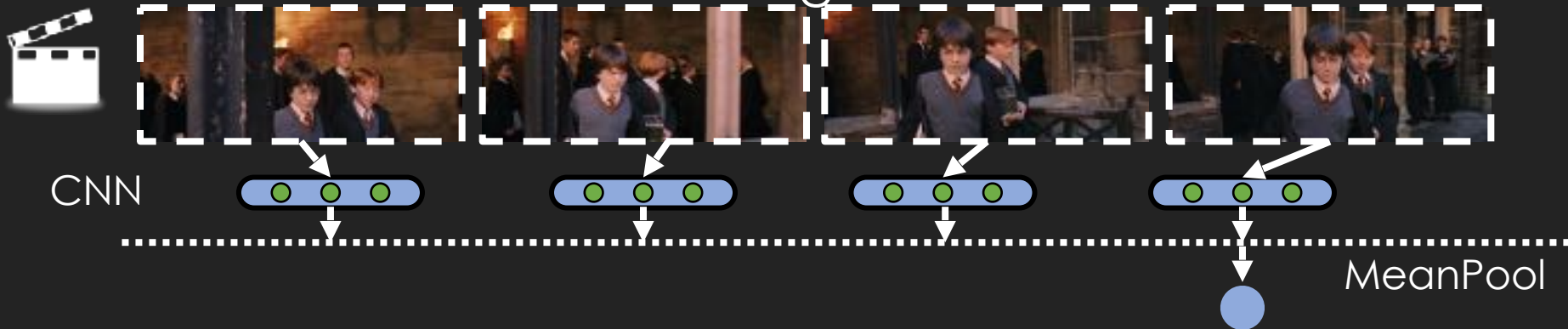
Aunt Petunia found a few moldy blankets in the second room and made up a bed for Dudley on the moth-eaten sofa.

She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

The storm raged more and more ferociously as the night went on, and Harry couldn't sleep.

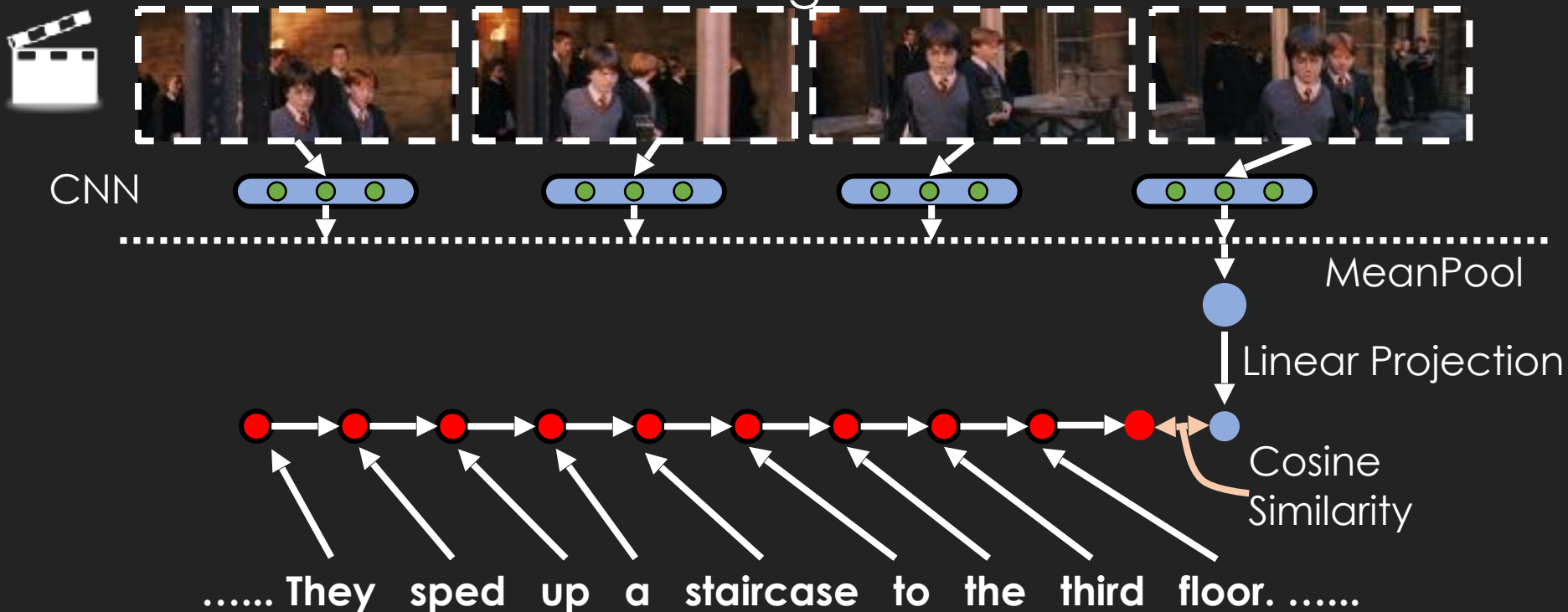
Video – Sentence Similarity

Visual Semantic Embedding



Video – Sentence Similarity

Visual Semantic Embedding



... As night fell, the promised storm blew up around them. Spray from the high waves splattered the walls of the hut and a fierce wind rattled the filthy windows.

Aunt Petunia found a few moldy blankets in the second room and made up a bed for Dudley on the moth-eaten sofa.

She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

The storm raged more and more ferociously as the night went on, and Harry couldn't sleep.

Aligning Movies with Books



01:00:58 --> 01:01:03

I'm telling you, it's spooky. She knows more about you than you do.

01:01:03 --> 01:01:05

Who doesn't?

Video – Sentence Similarity

Subtitle – Sentence Similarity



She and Uncle Vernon went the lumpy bed next door, and Harry was left to find the softest bit of floor he could and to curl up under the thinnest, most ragged blanket.

Dataset

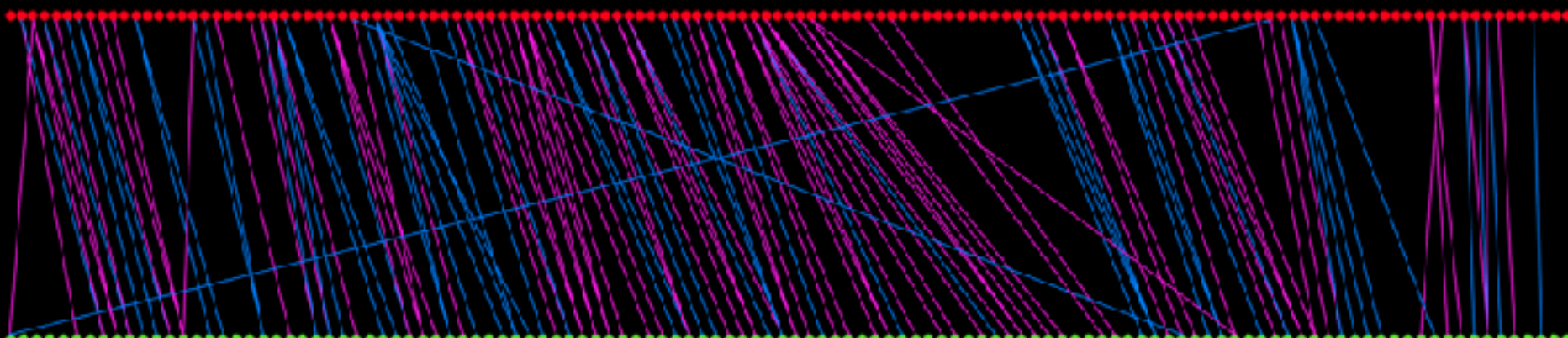


11 movie-book pairs annotated with 2,070 correspondences

Qualitative Results

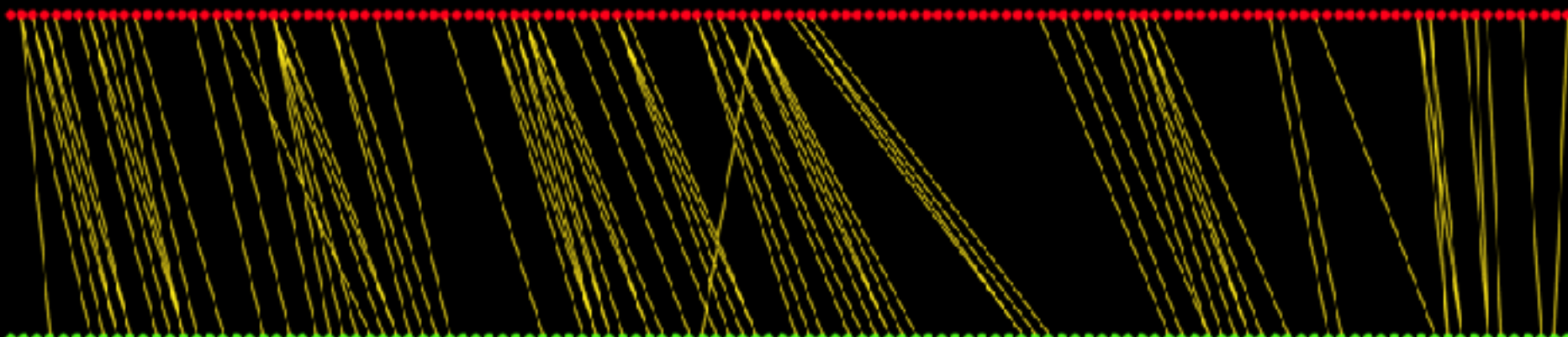
The Green Mile

book (paragraph)



movie (shot)

book (paragraph)



movie (shot)

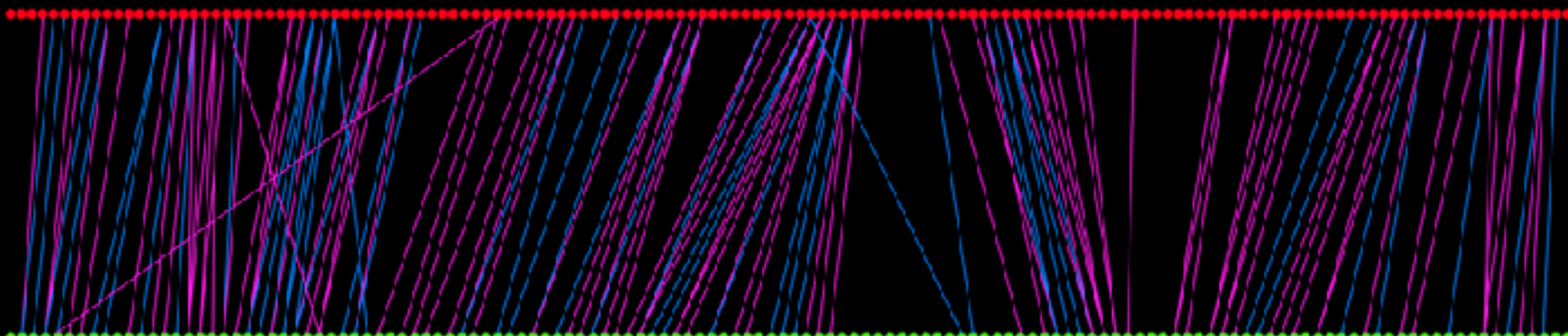
— visual match (113)
— dialog match (214)

— predict (317)

Qualitative Results

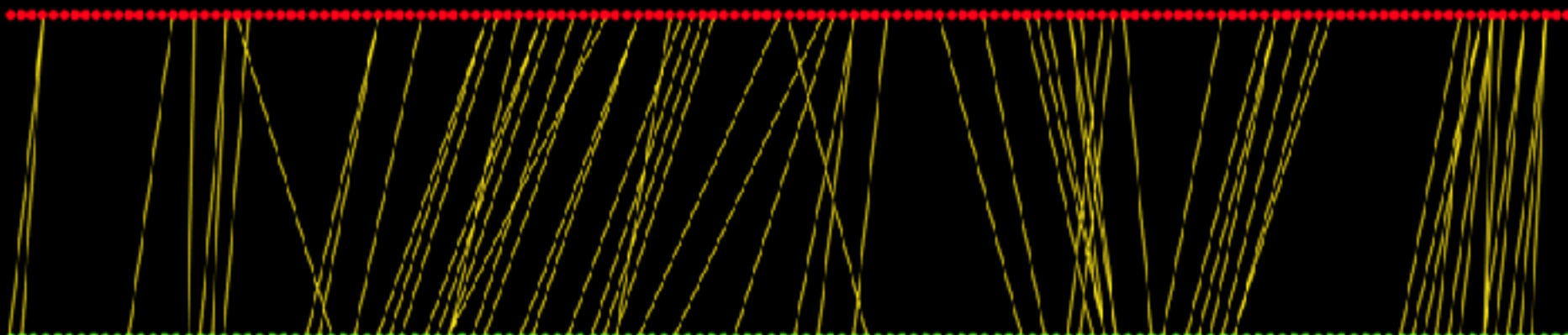
Harry Potter and the Sorcerers Stone

book (paragraph)



movie (shot)

book (paragraph)



movie (shot)

— visual match (77)

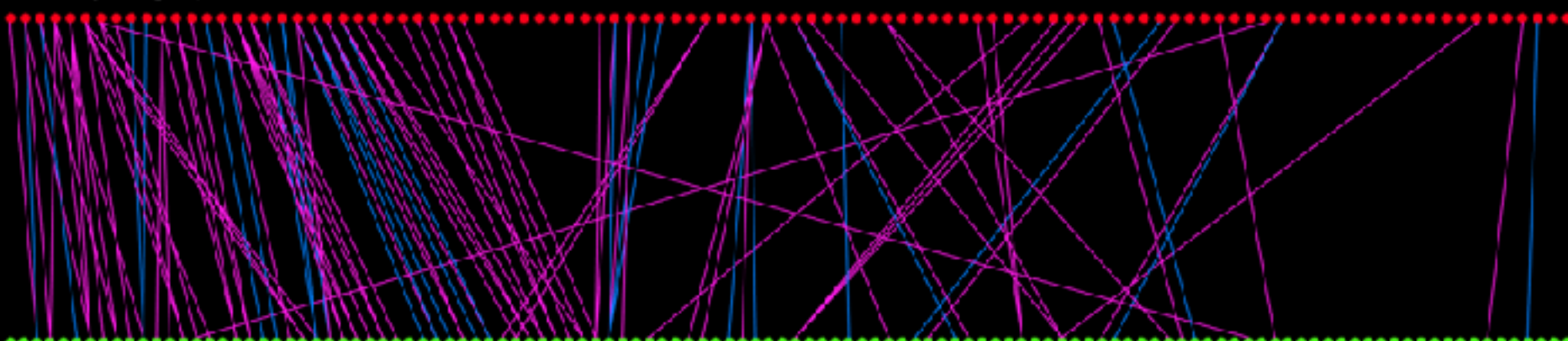
— dialog match (164)

— predict (239)

Qualitative Results

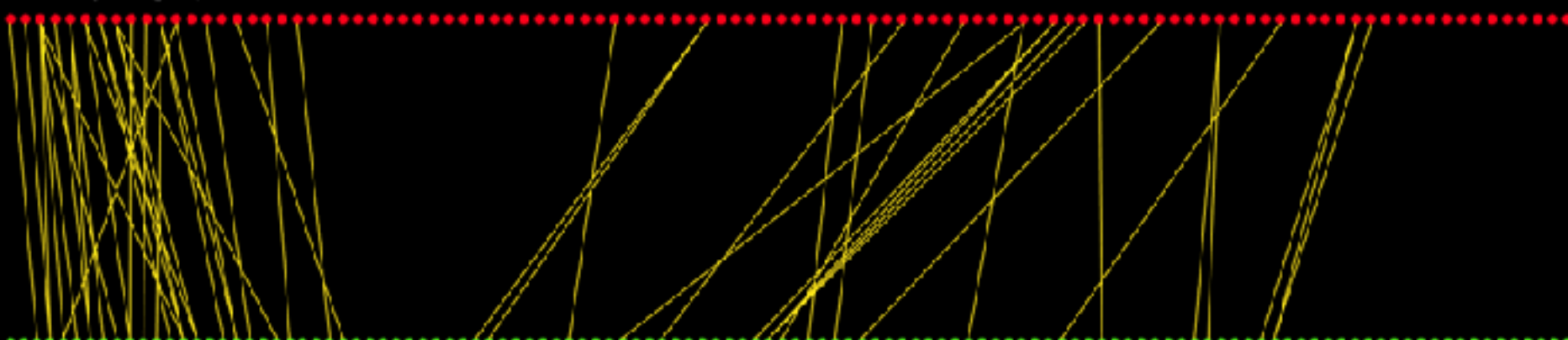
Fight Club

book (paragraph)



movie (shot)

book (paragraph)



movie (shot)

— visual match (46)
— dialog match (172)
— predict (212)

Qualitative Results on Alignment



We narrate the matched paragraph from the Harry Potter book.

Qualitative Results on Alignment

