# Deep Semantic Learning: Teach machines to understand text, image, and knowledge graph

Xiaodong He

DLTC, Microsoft Research, Redmond, WA, USA

Invited talk at CVPR *DeepVision* workshop, June 11, 2015

# Why should vision people ever care about language?

1. How to *teach* machines to understand images?

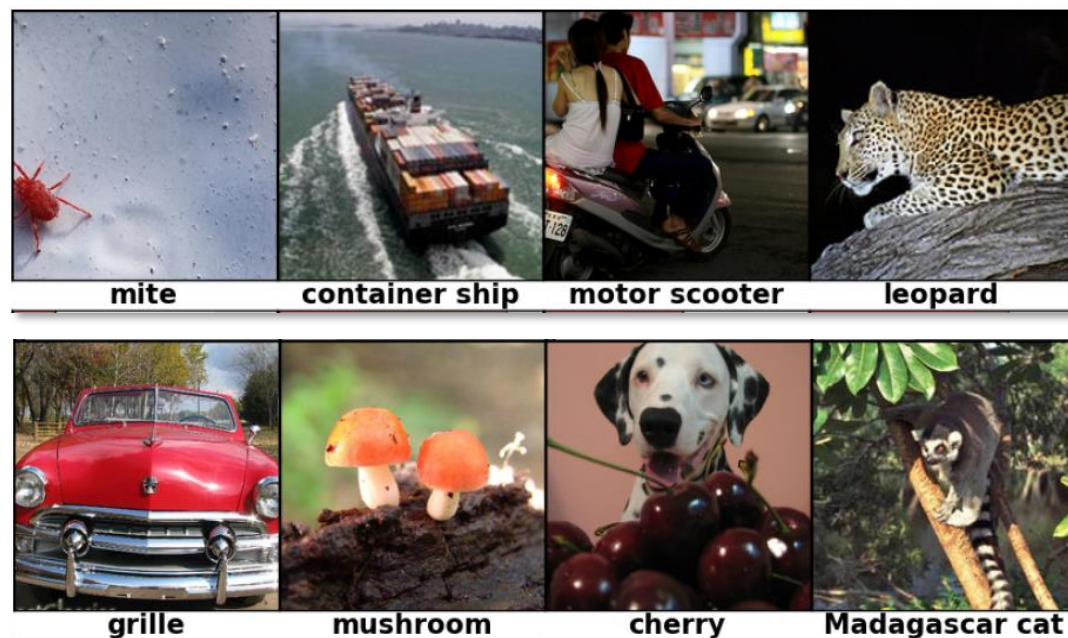2. How to *test* if a machine understands an image or not?

# How to teach machines to understand images?

## For image classification, we can label each image by a category and train the machine to predict

E.g., ImageNet provides hundreds to thousands of images for each category, aka **synset**, in the WordNet.



mite | container ship | motor scooter | leopard
grille | mushroom | cherry | Madagascar cat
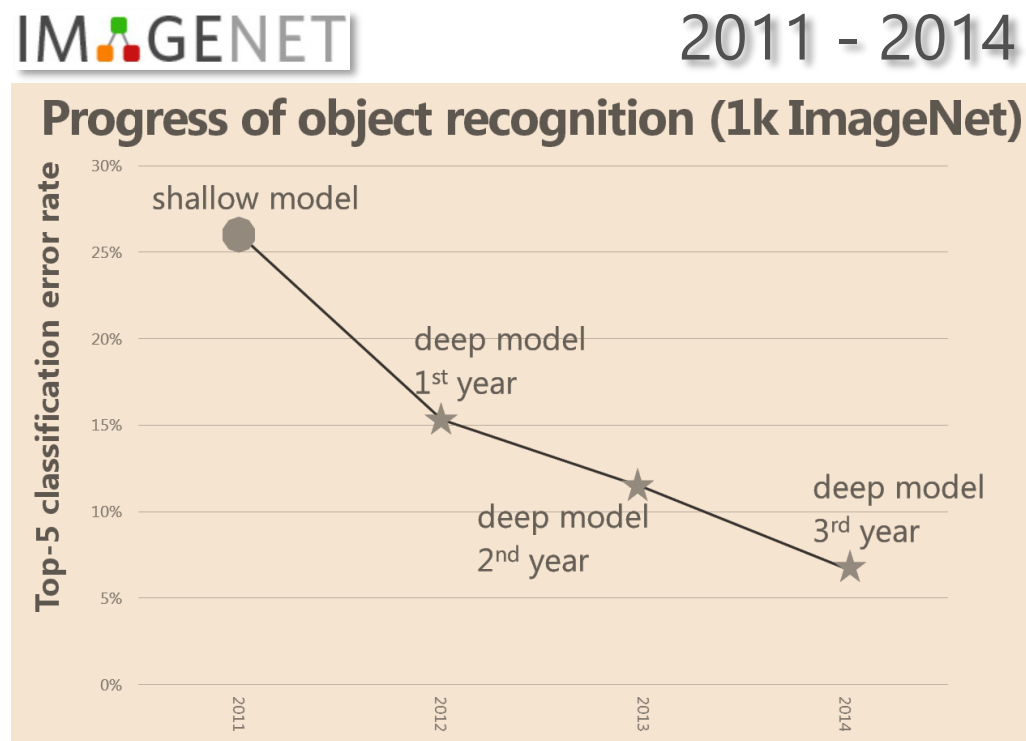
[Russakovsky, Deng, et al., 2014]

# How to test if a machine understand a image?

## For image classification, just check the prediction error rate

Dramatic progress in recent years thanks to deep CNN [LeCun, Bottou, Bengio, Haffner, 1998, Krizhevsky, Sutskever, Hinton, 2012].
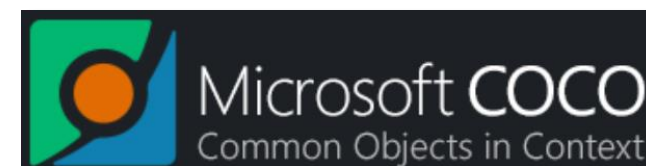
First time surpassed human-level performance (top5 err < 5%) on ImageNet classification in 2015 [He, Zhang, Ren, Sun, 2015]



**IM▲GENET**　　2011 - 2014

**Progress of object recognition (1k ImageNet)**

- Top-5 classification error rate
- shallow model
- deep model 1st year
- deep model 2nd year
- deep model 3rd year

But for complex scenes with a rich context, not possible to define all fine-grained subtle differences by categorization.

The best supervision is a full description in natural language

e.g., MS COCO provides 5 descriptions for each image that has a rich content.



Each description is:

- a coherent story.
- focused on salient info.
- with clear semantic meaning.
- reflecting certain common sense.

Could be a big variety.

- a woman is playing a frisbee with a dog.
- a woman is playing frisbee with her large dog.
- a girl holding a frisbee with a dog coming at her.
- a woman kneeling down holding a frisbee in front of a white dog.
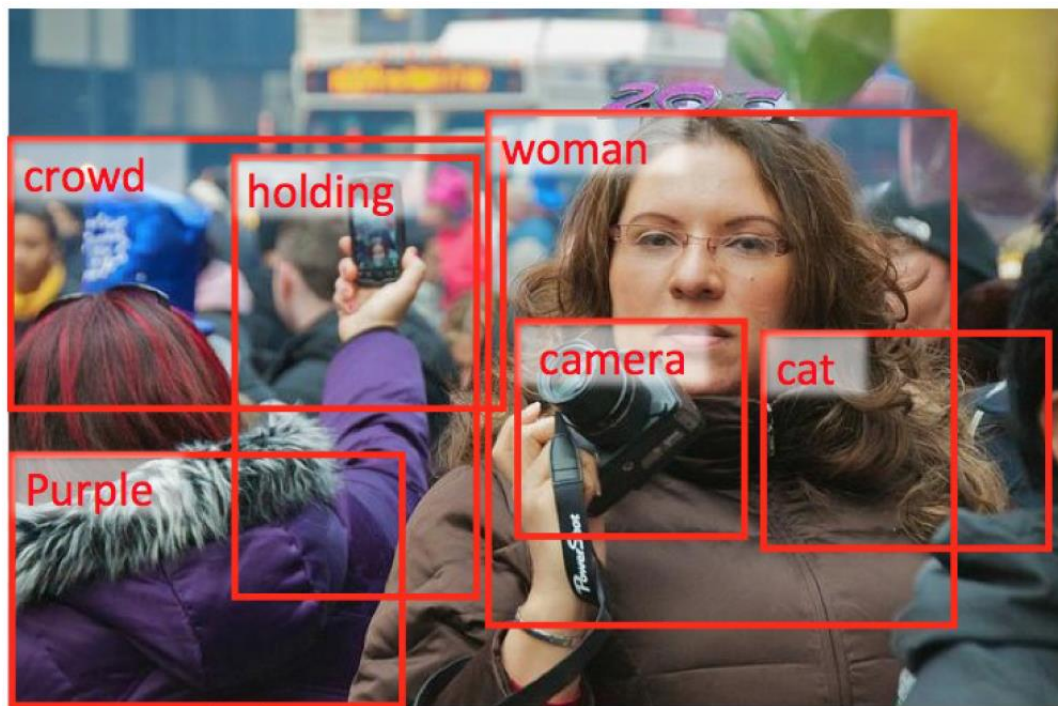- a young lady is playing frisbee with her dog.

[Lin, et al., 2014]

# How we test if a machine understands a complex scene?

## -- let's do a Turing Test!

ask the machine to describe the image in human language and see whether it reads like generated by a human


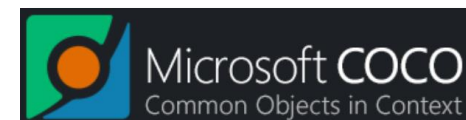
a woman holding a camera in a crowd.

# How much can machines understand complex scenes?

MS COCO Challenge: generate descriptive captions for images

The state-of-the-art at the MS COCO Captioning Challenge 2015

Measure the quality of the captions by human judge. (auto-metrics have big issues, see literature)

Great progress, but still a *big gap* vs. *Human*. (huge room for improvement)

| | | % of captions that pass the Turing Test | |
|---|---|---|---|
| Human | | 67.5% | -- |
| MSR | [Fang+ 15] | 32.2% | 1st(tie) |
| Google | [Vinyals+ 15] | 31.7% | 1st(tie) |
| MSR Captivator | [Devlin+ 15] | 30.1% | 3rd(tie) |
| Montreal/Toronto | [Xu+ 15] | 27.2% | 3rd(tie) |
| Berkeley LRCN | [Donahue+ 15] | 26.8% | 5th |

Understanding language is necessary for building strong vision intelligence

Moreover, knowledge bases, from WordNet to Freebase, are extremely helpful, too.

# Natural Language Understanding

- Build an intelligent system that can interact with human using natural language

- Talk's outline
  - Learning semantic representation of text
  - Knowledge base and question answering
  - Multimodal (image-text) semantic models



http://csunplugged.org/turing-test

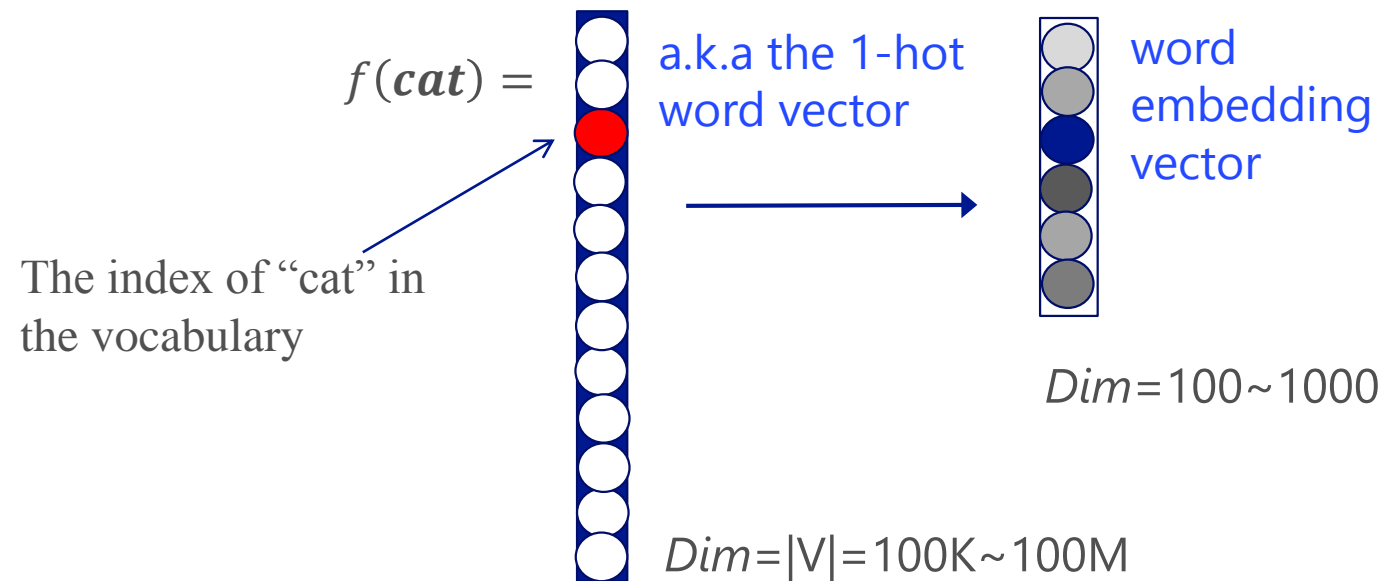# Background

- Word embedding (Word2Vec) [Bengio 03, Mikolov+ 10, 13]
  - representing word meaning in a continuous space

$$f(\textbf{cat}) =$$

a.k.a the 1-hot word vector

The index of "cat" in the vocabulary

$Dim=|V|=100K\sim100M$

word embedding vector

$Dim=100\sim1000$

# Background

- Neural net based language modeling [Bengio+ 03, Schwenk+ 06, Mikolov+ 10]
    - predict the next word given the context, e.g., Cisco issued earnings _?_



Feedforward NN

Recurrent NN

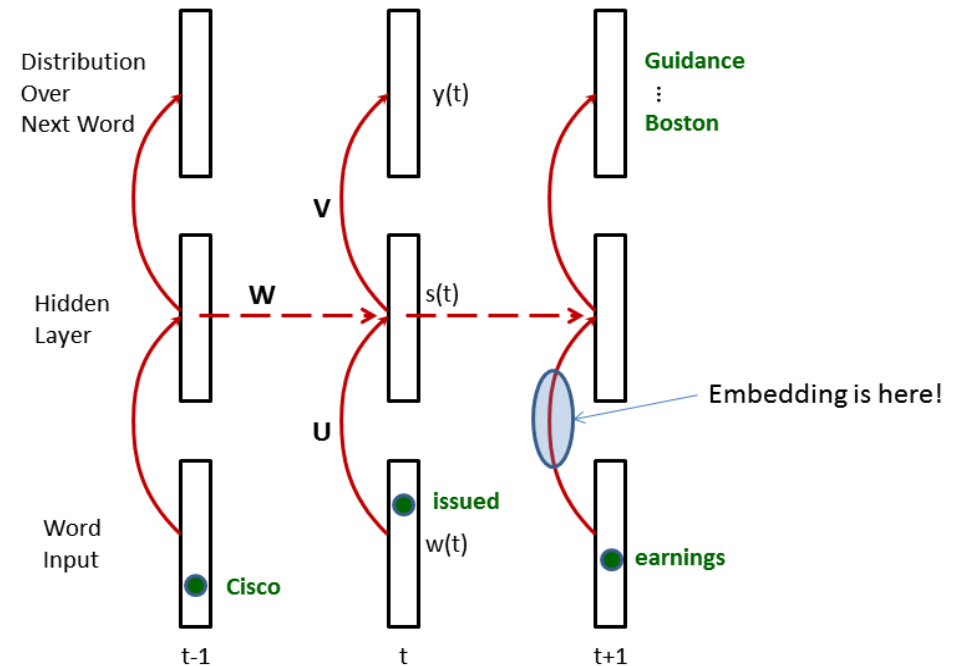Plotting 3K words in 2D

*Deep Learning in Computer Vision 2015*

# Learning semantic representation for a sentence
## e.g., from a raw sentence to an abstract semantic vector (Sent2Vec)

Abstract representation
in the semantic space

$W_4$

H3

$W_3$

each non-linear layer gradually
extracts deeper invariance

H2

$W_2$

H1

$W_1$

Raw text, e.g., a
sequence of words

Input 1

*a woman holding a camera in a crowd*

# The supervision problem:



*a woman holding a camera in a crowd*

However
- the semantic meaning of texts – to be learned – is latent
- no clear target for the model to learn

Fortunately
- we usually know if two texts are "similar" or not.
- That's the signal for semantic representation learning.

# Deep Structured Semantic Model (DSSM)

**Deep Structured Semantic Model/Deep Semantic Similarity Model**
   *Sentence to vector!*

**Built upon sub-word units** for scalability and generalizability
   e.g., letter-trigrams, phones, roots/morphs, instead of *words*

Trained by optimizing an similarity-driven objective
   Using a structure similar to auto-encoder / Siamese net, projecting semantically similar sentences to vectors close to each other

Semi-supervised/weak supervised learning
   semantically-similar text pairs, e.g., user behavior log data, contextual text

[Huang, He, Gao, Deng, Acero, Heck, "Learning deep structured semantic models for web search using clickthrough data," CIKM, October, 2013]

# DSSM: a similarity-driven Sent2Vec model

**Initialization:**

Neural networks are initialized with random weights

Semantic vector $\rightarrow$ $v_s$ $\qquad v_{t^+} \qquad v_{t^-}$

| | $v_s$ | $v_{t^+}$ | $v_{t^-}$ |
|---|---|---|---|
| | d=300 | d=300 | d=300 |
| | $W_{s,4}$ | $W_{t,4}$ | $W_{t,4}$ |
| | d=500 | d=500 | d=500 |
| | $W_{s,3}$ | $W_{t,3}$ | $W_{t,3}$ |
| Letter-trigram embedding matrix | d=500 | d=500 | d=500 |
| $\rightarrow$ | $W_{s,2}$ | $W_{t,2}$ | $W_{t,2}$ |
| Letter-trigram encoding matrix (fixed) | dim = 50K | dim = 50K | dim = 50K |
| $\rightarrow$ | $W_{s,1}$ | $W_{t,1}$ | $W_{t,1}$ |
| Bag-of-words vector | dim = 100M | dim = 100M | dim = 100M |
| Input word/phrase | $s$: "**racing car**" | $t^+$: "**formula one**" | $t^-$: "**racing to me**" |

# DSSM: a similarity-driven Sent2Vec model

**Training:**

Compute Cosine similarity between semantic vectors

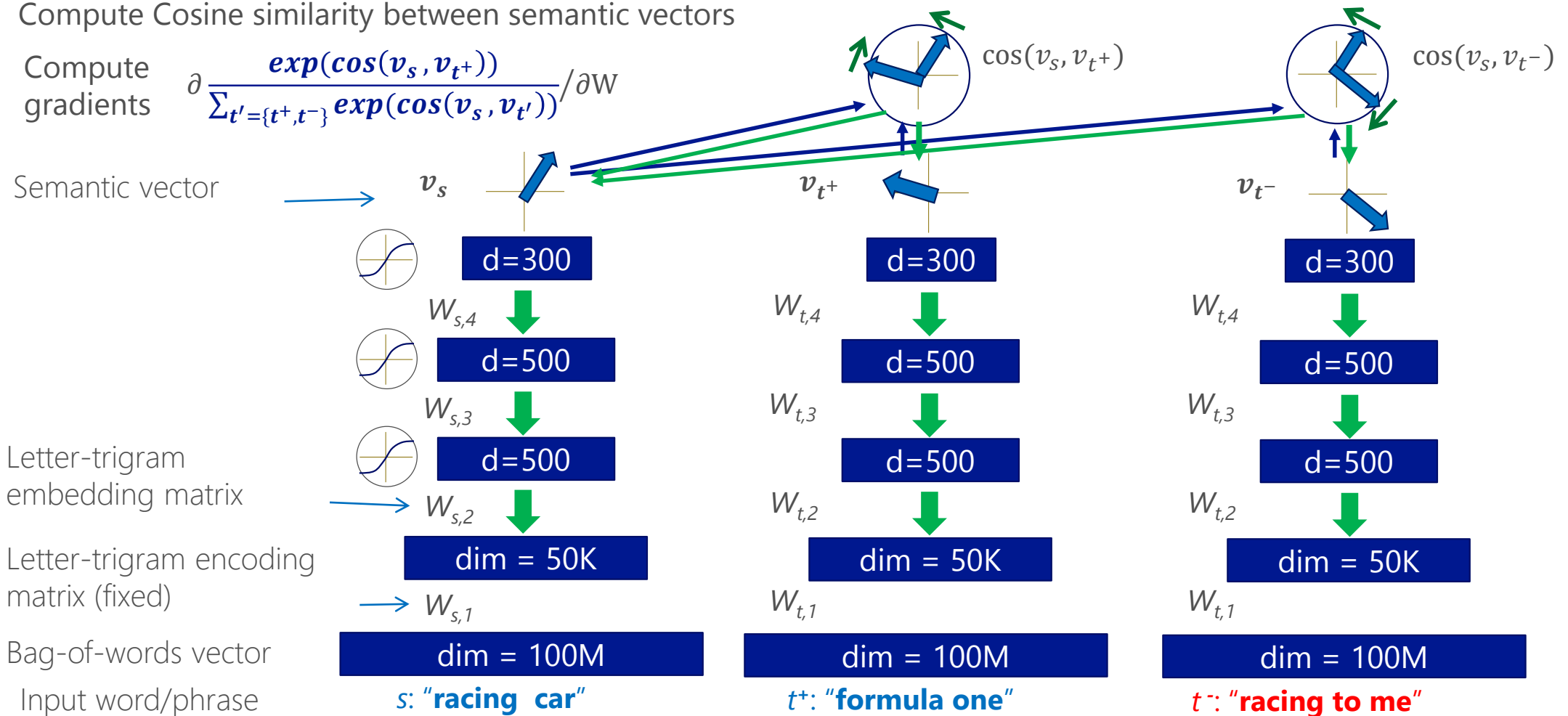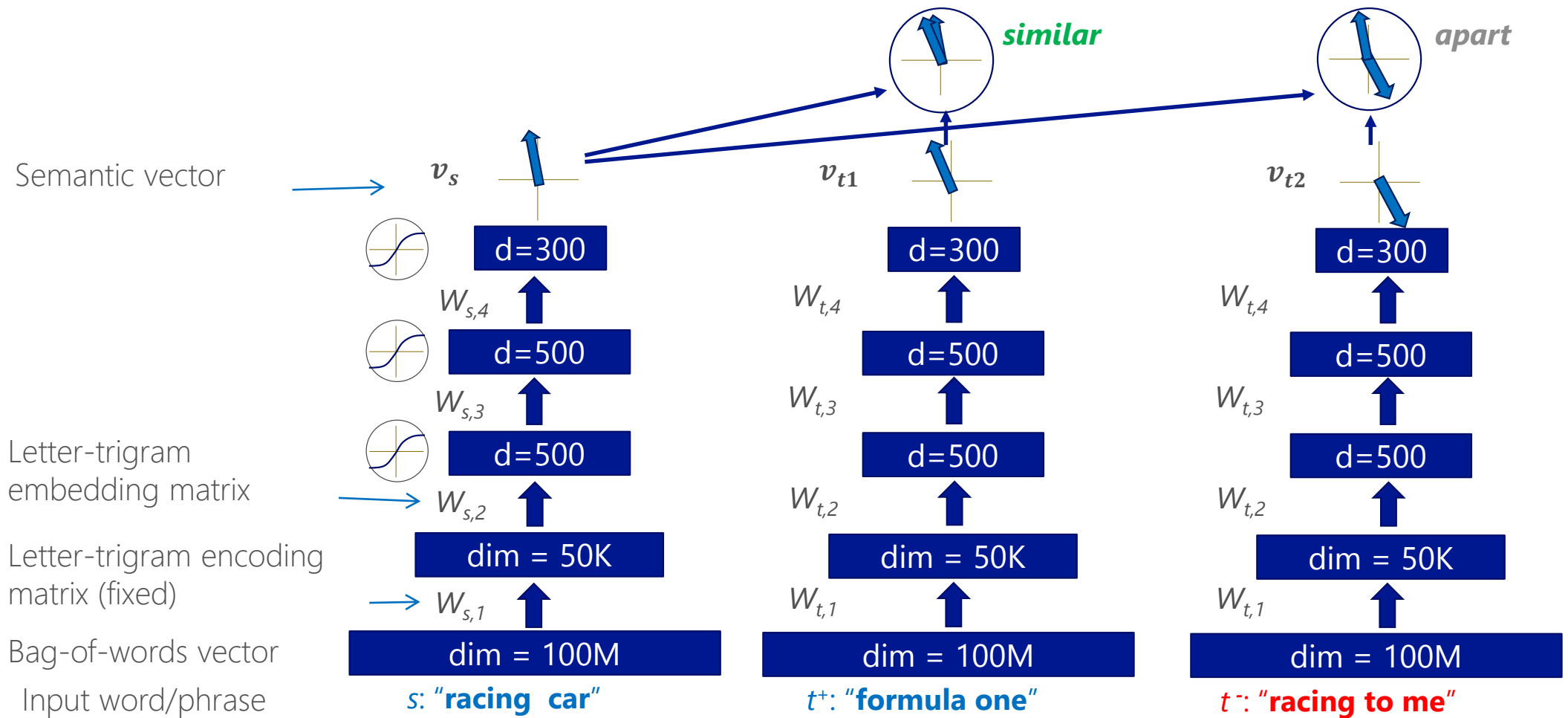Compute gradients $\quad \partial \dfrac{exp(cos(v_s, v_{t^+}))}{\sum_{t'=\{t^+, t^-\}} exp(cos(v_s, v_{t'}))} / \partial w$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad cos(v_s, v_{t^+}) \qquad\qquad\qquad\qquad\qquad cos(v_s, v_{t^-})$

Semantic vector $\qquad\qquad\qquad v_s \qquad\qquad\qquad\qquad\qquad\qquad v_{t^+} \qquad\qquad\qquad\qquad\qquad v_{t^-}$

| | d=300 | | d=300 | | d=300 |
|---|---|---|---|---|---|
| | $W_{s,4}$ | | $W_{t,4}$ | | $W_{t,4}$ |
| | d=500 | | d=500 | | d=500 |
| | $W_{s,3}$ | | $W_{t,3}$ | | $W_{t,3}$ |
| | d=500 | | d=500 | | d=500 |

Letter-trigram
embedding matrix $\qquad\qquad W_{s,2} \qquad\qquad\qquad\qquad W_{t,2} \qquad\qquad\qquad\qquad W_{t,2}$

Letter-trigram encoding
matrix (fixed) $\qquad\qquad\quad$ dim = 50K $\qquad\qquad$ dim = 50K $\qquad\qquad$ dim = 50K

$\qquad\qquad\qquad\qquad W_{s,1} \qquad\qquad\qquad\qquad\qquad W_{t,1} \qquad\qquad\qquad\qquad W_{t,1}$

Bag-of-words vector $\qquad$ dim = 100M $\qquad\qquad$ dim = 100M $\qquad\qquad$ dim = 100M

Input word/phrase $\qquad$ *s:* "**racing car**" $\qquad$ *t⁺:* "**formula one**" $\qquad$ *t⁻:* "**racing to me**"

Xiaodong He $\qquad$ Deep Semantic Learning: Teach machines to understand text, image, and knowledge graph $\qquad$ 17 $\qquad$ CVPR 2015

*Deep Vision: Deep Learning in Computer Vision 2015*

# DSSM: a similarity-driven Sent2Vec model

**Runtime:**



Semantic vector $\quad\quad v_s\quad\quad\quad\quad\quad\quad\quad v_{t1}\quad\quad\quad\quad\quad\quad\quad v_{t2}$

similar

apart

| | $d=300$ | | $d=300$ | | $d=300$ |
| $W_{s,4}$ | | $W_{t,4}$ | | $W_{t,4}$ | |
| | $d=500$ | | $d=500$ | | $d=500$ |
| $W_{s,3}$ | | $W_{t,3}$ | | $W_{t,3}$ | |

Letter-trigram embedding matrix

| | $d=500$ | | $d=500$ | | $d=500$ |
| $W_{s,2}$ | | $W_{t,2}$ | | $W_{t,2}$ | |
| | dim = 50K | | dim = 50K | | dim = 50K |

Letter-trigram encoding matrix (fixed)

$W_{s,1}\quad\quad\quad\quad\quad\quad W_{t,1}\quad\quad\quad\quad\quad\quad W_{t,1}$

Bag-of-words vector

| dim = 100M | dim = 100M | dim = 100M |

Input word/phrase

$s$: "**racing car**"  $\quad\quad\quad$ $t^+$: "**formula one**" $\quad\quad\quad$ $t^-$: "**racing to me**"

# Sent2Vec is crucial in many NLP tasks

| Tasks | Source | Target |
|---|---|---|
| **Web search** | **search query** | **web documents** |
| Ad selection | search query | ad keywords |
| Contextual entity ranking | mention (highlighted) | entities |
| Online recommendation | doc in reading | interesting things / other docs |
| Machine translation | phrases in language S | phrases in language T |
| **Knowledge-base construction** | **entity** | **entity** |
| **Question answering** | **pattern | mention** | **relation | entity** |
| Personalized recommendation | user | app, movie, etc. |
| Image search | query | image |
| **Image captioning** | **image** | **text caption** |
| ... | | |

# DSSM: built on top of sub-word units

Decompose *any* word into sub-word units (SWU), e.g., letter-trigram

embedding vector

$W \rightarrow U \times V$

embedding vector

dim=500

word embedding matrix: $500 \times 100M$

$W$

dim = 100M

Bag-of-words vector

dim=500

SWU embedding matrix: $500 \times 50K$

$U$

dim = 50K

SWU encoding matrix

$V$

dim = 100M

Bag-of-words vector

Could go up to extremely large

Preferable for large scale NL tasks
- Arbitrary size of vocabulary (*scalability*)
- Misspellings, word fragments, new words, etc. (*generalizability*)

[Huang, He, Gao, Deng, Acero, Heck, CIKM2013]

# Sub-word unit encoding

- E.g., letter-trigram based *Word Hashing* of "cat"
  - -> #cat#
  - Tri-letters: #-c-a, c-a-t, a-t-#.

- Compact representation
  - |Voc| (500K) → |Letter-trigram| (30K)

- Generalize to unseen words

- Robust to misspelling, inflection, etc.

What if different words have the same word hashing vector (collision)?

$$x\ (cat) = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

← The index of word *cat* in the vocabulary

$$f(cat) = \begin{bmatrix} \vdots \\ 1 \\ \vdots \\ 1 \\ 1 \\ \vdots \end{bmatrix}$$

Indices of *#-c-a, c-a-t, a-t-#* in the letter-tri-gram list, respectively.

| Vocabulary size | Unique letter-tg observed in voc | Number of Collisions |
|---|---|---|
| 40K | 10306 | 2  (0.005%) |
| 500K | 30621 | 22 (0.004%) |

# Other options of sub-word units (SWU):

- Letters, context-dependent phones
- context-dependent morphs, positioned-roots/morphs

[e.g., Zhang and LeCun, Text Understanding from Scratch, 2015]

CVPR 2015

DeepVision: *Deep Learning in Computer Vision 2015*

# Training objectives

Objective: cosine similarity based loss

Using web search as an example:

- a query $q$ and a list of docs $D = \{d^+, d_1^-, \dots d_K^-\}$
  - $d^+$ positive doc; $d_1^-, \dots d_K^-$ are negative docs to $q$ ( e.g., sampled from not clicked docs)

- Objective: the posterior probability of the clicked doc given the query

$$P(d^+|q) = \frac{\exp\left(\gamma\, cos(v_{\boldsymbol\theta}(q), v_{\boldsymbol\theta}(d^+))\right)}{\sum_{d \in D} \exp\left(\gamma\, cos(v_{\boldsymbol\theta}(q), v_{\boldsymbol\theta}(d))\right)}$$

e.g., $v_{\boldsymbol\theta}(q) = \sigma(W_{s,4} \times \sigma(W_{s,3} \times \sigma(W_{s,2} \times ltg(q))))$

$v_{\boldsymbol\theta}(d) = \sigma(W_{t,4} \times \sigma(W_{t,3} \times \sigma(W_{t,2} \times ltg(d))))$

where $\theta = \{W_{s,2\sim4}, W_{t,2\sim4}\}$, $\sigma()$ is a tanh function.

# Convolutional DSSM

Model local context at the convolutional layer
Model global context at the pooling layer

Semantic layer: $y$

Affine projection matrix: $W_s$

Max pooling layer: $v$

Max pooling operation

Convolutional layer: $h_t$

Convolution matrix: $W_c$

Word hashing layer: $f_t$

Word hashing matrix: $W_f$

Word sequence: $x_t$

**Figure 1: Illustration of the C-DSSM. A convolutional layer with the window size of three is illustrated.**

[Shen, He, Gao, Deng, Mesnil, WWW2014 & CIKM2014]

– ## What does the model learn at the convolutional layer?

Capture the local context dependent word sense

- Learn one embedding vector for each local context-dependent word



$$h_t = W_c \times [f_{t-1}, f_t, f_{t+1}]$$

semantic space

auto **body** repair
car **body** shop   car **body** kits
auto **body** part

wave **body** language
     calculate **body** fat
forcefield **body** armour

The similarity between different "**body**" within contexts

| car **body** shop | cosine similarity |
|---|---|
| car **body** kits | 0.698 |
| auto **body** repair | 0.578 |
| auto **body** parts | 0.555 |
| wave **body** language | 0.301 |
| calculate **body** fat | 0.220 |
| forcefield **body** armour | 0.165 |

**high similarity**

**low similarity**

# CDSSM: What happens at the max-pooling layer?



$$v(i) = \max_{t=1,\dots,T} \{h_t(i)\}$$

where $i = 1,\dots,300$

- Aggregate *local topics* to form the *global intent*
- Identify salient words/phrase at the max-pooling layer

Words that win the most active neurons at the **max-pooling layers:**

auto body repair cost calculator software

Those are salient words containing clear intents/topics

BTW, with the new *attention* model, these info could modeled in a more principled way [Bahdanau, Cho, Bengio, 2014; Xu et al, 2015]

# Mine semantically-similar text pairs from Search Logs

*how to deal with stuffy nose?*

*stuffy nose treatment*

*cold home remedies*

**Best Home Remedies for Cold and Flu**
Wind Heat External Pathogens
By: Catherine Browne, L.Ac., MH, Dipl. Ac.

In Chinese medicine, colds and flu's are delineated
into several different energetic classifications.
Here we will outline the different types of cold
and flu viruses that you will likely encounter, and
then describe the best home remedies for these

| QUERY (Q) | Clicked Doc Title (T) |
|---|---|
| how to deal with stuffy nose | best home remedies for cold and flu |
| stuffy nose treatment | best home remedies for cold and flu |
| cold home remedies | best home remedies for cold and flu |
| ... ... | ... ... |
| go israel | forums goisrael community |
| skate at wholesale at pr | wholesale skates southeastern skate supply |
| breastfeeding nursing blister baby | clogged milk ducts babycenter |
| thank you teacher song | lyrics for teaching educational children s music |
| immigration canada lacolle | cbsa office detailed information |

[Gao, He, Nie, CIKM2010]

CVPR 2015
Deep Vision: *Deep Learning in Computer Vision 2015*

# DSSM for Information Retrieval

- Training Dataset
  - 30 Million (Query, Document) Click Pairs

- Testing Dataset
  - **12,071** English queries
  - around 65 web document associated to each query in average
  - Human gives each <query, doc> pair the label, with range **0 to 4**
  - 0: Bad      1: Fair      2: Good      3: Perfect      4: Excellent

- Evaluation Metric: (higher the better)
  - NDCG

- GPU (NVidia GPU K40)



Dist. of query and doc title length

# Main Experiment Results

CDSSM: Shen et al. 2014

## NDCG@1 Results



| | BM25 | ULM | PLSA | BLTM | WTM | DSSM | CDSSM |
|---|---|---|---|---|---|---|---|
| ■ NDCG@1 | 30.5 | 30.4 | 30.5 | 31.6 | 31.5 | 32.7 | 34.8 |

Chart labels: Lexical Matching Models, Topic Models, Click-Through based Translation Models, Deep Semantic Model, Convolutional Deep Semantic Model

# An example

sarcoidosis is a disease, a symptom is excessive amount of calcium in one's urine and blood. So medicines that increase the absorbing of calcium should be avoid. While **Vitamin d** is closely associated to **calcium absorbing**.

We observed that "sarcoidosis" in the document title and "absorbs" "excessive" and "vitamin (d)" in the query have high activations at neurons 90, 66, 79, indicating that the model knows that **"sarcoidosis" share similar** semantic meaning with "absorbs" "excessive" "vitamin (d)", collectively.

what happens if our body **absorbs** **excessive** amount **vitamin** **d**

| 88 | 90 | 66 | 79 | 102 | 35 | 16 | 94 |

| 88 | 90 | 66 | 79 | 102 | 35 | 16 | 94 |

**calcium** supplements and **vitamin** **d** discussion stop **sarcoidosis**

Most active neurons at the **max-pooling layers** of the query and document nets, respectively

# Recurrent DSSM

- Encode the word one by one in the recurrent hidden layer
- The hidden layer at the last word codes the semantics of the full sentence
- Model is trained by a cosine similarity driven objective

Embedding vector



[Palangi, Deng, Shen, Gao, He, Chen, Song, Ward, 2015]

# Using LSTM cells

## LSTM (long short term memory) uses special cells in RNN (Hochreiter and Schmidhuber, 1997)



Embedding vector

$$\mathbf{y}_g(t) = g(\mathbf{W}_4 \mathbf{l}_1(t) + \mathbf{W}_{rec4} \mathbf{y}(t-1) + \mathbf{b}_4)$$

$$\mathbf{i}(t) = \sigma(\mathbf{W}_3 \mathbf{l}_1(t) + \mathbf{W}_{rec3} \mathbf{y}(t-1) + \mathbf{W}_{p3} \mathbf{c}(t-1) + \mathbf{b}_3)$$

$$\mathbf{f}(t) = \sigma(\mathbf{W}_2 \mathbf{l}_1(t) + \mathbf{W}_{rec2} \mathbf{y}(t-1) + \mathbf{W}_{p2} \mathbf{c}(t-1) + \mathbf{b}_2)$$

$$\mathbf{c}(t) = \mathbf{f}(t) \circ \mathbf{c}(t-1) + \mathbf{i}(t) \circ \mathbf{y}_g(t)$$

$$\mathbf{o}(t) = \sigma(\mathbf{W}_1 \mathbf{l}_1(t) + \mathbf{W}_{rec1} \mathbf{y}(t-1) + \mathbf{W}_{p1} \mathbf{c}(t) + \mathbf{b}_1)$$

$$\mathbf{y}(t) = \mathbf{o}(t) \circ h(\mathbf{c}(t)) \tag{2}$$

where $\circ$ denotes Hadamard (element-wise) product.

Figure 2. The basic LSTM architecture used for sentence embedding

[Palangi, Deng, Shen, Gao, He, Chen, Song, Ward, Deep Sentence Embedding Using the LSTM network: Analysis and Application to IR, 2015]

# Results

| Model | NDCG@1 | NDCG@3 | NDCG@10 |
|---|---|---|---|
| BM25 | 30.5% | 32.8% | 38.8% |
| PLSA (T=500) | 30.8% | 33.7% | 40.2% |
| DSSM (nhid = 288/96), 2 Layers | 31.0% | 34.4% | 41.7% |
| CLSM (nhid = 288/96), 2 Layers | 31.8% | 35.1% | 42.6% |
| RNN (nhid = 288), 1 Layer | 31.7% | 35.0% | 42.3% |
| LSTM-RNN (ncell = 96), 1 Layer | **33.1%** | **36.5%** | **43.6%** |

LSTM learns much faster than regular RNN

LSTM effectively represents the semantic information of a sentence using a vector

# Related work

Embedding vector — Embedding vector

[Palangi et al, 2015]



Source sentence

Target sentence

## Minimize *sentence-level* semantic matching loss

*vs.*

Embedding vector



Source sentence

Target sentence

## Minimize *word-level* cross-entropy loss

[Sutskever, Vinyals, Le, 2014. Sequence to Sequence Learning with Neural Networks]

# Some other related work

**Deep CNN for text input**
Mainly classification tasks in the paper

[Kalchbrenner, Grefenstette, Blunsom, A Convolutional Neural Network for Modelling Sentences, ACL2014]

**Paragraph Vector**
Learn a vector for a paragraph

Quoc Le, Tomas Mikolov, Distributed Representations of Sentences and Documents, in ICML 2014

**Recursive NN (ReNN)**
Tree structure, e.g., for parsing

[Socher, Lin, Ng, Manning, "Parsing natural scenes and natural language with recursive neural networks", 2011]

**Tensor product representation (TPR)**
Tree representation

[Smolensky and Legendre: The Harmonic Mind, From Neural Computation to Optimality-Theoretic Grammar, MIT Press, 2006]

**Tree-structured LSTM Network**
Tree structure LSTM

[Tai, Socher, Manning. 2015. Improved Semantic Representations From Tree-Structured LSTM Networks.]

# From Natural Language to Knowledge Base

- Captures world knowledge by storing properties of millions of entities, as well as relations among them



Freebase
DBpedia
YAGO
NELL
OpenIE/ReVerb

# Neural Knowledge Base Embedding



**Vectors** for entities,
**matrices** for relations

$$S_{(e_1, r, e_2)} = (a^T M_r b)$$

Bordes+ 2012, Collobert &
Weston 2008, Socher+ 2013,
Yang+ 2015

$$G_r\left(y_{e_1}, y_{e_2}\right)$$

$$y_{e_1} = f(W x_{e_1}) \in R^n \qquad y_{e_2} = f(W x_{e_2}) \in R^n$$

$$W \qquad W$$

$$x_{e_1} \in R^m \qquad x_{e_2} \in R^m$$

$$e_1 \qquad r \qquad e_2$$

*Barack Obama*     *Born-in*     *Hawaii*

A neural network framework for multi-relational learning

# Mining Horn-clause Rules

- Can relation embedding capture relation composition?

$$BornInCity(a, b) \wedge CityInCountry(b, c) \Rightarrow Nationality(a, c)$$

- Embedding-based Horn-clause rule extraction
  - For each relation $r$, find a chain of relations $r_1 \cdots r_n$, such that:
  $$dist(M_r, M_1 \circ M_2 \circ \cdots \circ M_n) < \theta$$
  - $r_1(e_1, e_2) \wedge r_2(e_2, e_3) \cdots \wedge r_n(e_n, e_{n+1}) \rightarrow r(e_1, e_{n+1})$

- Advantages vs. Inductive Logic Programming
  - Search the relation space instead of instance space

# Aggregated Precision of Top Length-2 Rules



- AMIE [Galárraga+, WWW-2013] is an association rule-mining approach for large-scale KBs.
- Data: FB15k-401
- Execution time:
  - AMIE: 9 min.
  - EmbedRule: 2 min.

Yang, Yih, He, Gao, Deng, ICLR2015

*Figure 6*: Relation embeddings of DistMult

# Semantic Parsing and Question Answering w/ KB

*Who is Justin Bieber's sister?*

*Jazmyn Bieber*

Knowledge Base

semantic parsing

$\lambda x.\,\text{sister\_of}(\text{justin\_bieber}, x)$

query

matching

$\text{sibling\_of}(\text{justin\_bieber}, x) \wedge \text{gender}(x, \text{female})$

# Key Challenge – Language Mismatch

- Lots of ways to ask the same question
  - *"What was the date that Minnesota became a state?"*
  - *"Minnesota became a state on?"*
  - *"When was the state Minnesota created?"*
  - *"Minnesota's date it entered the union?"*
  - *"When was Minnesota established as a state?"*
  - *"What day did Minnesota officially become a state?"*

- Need to map them to the predicate defined in KB
  - location.dated_location.date_founded

Deep Semantic Learning: Teach machines to understand
text, image, and knowledge graph

CVPR 2015
DeepVision: *Deep Learning in Computer Vision 2015*

# Matching Question and Relation

- Similar text can map to very different relations

  - *Q=Who is the father of King George VI?*
  - *R*=people.person.parents

  - *Q=Who is the father of the Periodic Table?*
  - *R*=law.invention.inventor

# Staged Query Graph Generation

- ## Query graph
  - Resembles subgraphs of the knowledge base
  - Can be directly mapped to a logical form in $\lambda$-calculus
  - Semantic parsing: a search problem that *grows* the graph through actions

- Who first voiced Meg on Family Guy?
- $\lambda x. \exists y. \text{cast}(\text{FamilyGuy}, y) \wedge \text{actor}(y, x) \wedge \text{character}(y, \text{MegGriffin})$



topic entity

constraints

inferential chain

argmin — from — character → Meg Griffin

Family Guy — cast → $y$ — actor → $x$

# Staged Graph Generation

*"Who first voiced Meg on Family Guy?"*

1. Topic Entity Linking [Yang&Chang ACL-15]

$s_0$ $\phi$ → $s_1$ Family Guy

$s_2$ Meg Griffin

*"Who first voiced Meg on **X**?" ∧ **X**=Family Guy* ⇔ Family Guy

or

*"Who first voiced **X** on Family Guy?" ∧ **X** = Meg* ⇔ Meg Griffin

---

Query graph that represents the question:
- Identify possible entities in the question (e.g., Meg, Family Guy)
- Only search relations around these entities in the KB
- Narrow down the search space significantly

# Staged Graph Generation

## 2. Core Inferential Chain (DSSM)

Given an **M**ention/**E**ntity match:

$$X = Family\ Guy \Leftrightarrow \boxed{Family\ Guy}$$

Next, need to match **P** ⇔ **R**

"*Who first voiced Meg on X?*" ⇔ **?R**



DSSM measures the semantic matching between **P**attern and **R**elation:

*who first voiced Meg on X*

And

"cast-actor"

"writer-start"

"genre"



*Who first voiced Meg on X*          cast-actor

Matching (multi-hop) relations: concatenate multiple relations to a long relation on-the-fly, the DSSM takes care the issues of aggregating semantics from individual relations.

# Graph Generation Stages (cont'd)

- **Who first voiced Meg on Family Guy?**

3. Augment constraints



The Freebase

Find the answer

# WEBQUESTIONS Benchmark [Berant+ EMNLP-2013]

- *What character did Natalie Portman play in Star Wars?* ⇒ Padme Amidala
- *What kind of money to take to Bahamas?* ⇒ Bahamian dollar
- *What currency do you use in Costa Rica?* ⇒ Costa Rican colon
- *What did Obama study in school?* ⇒ political science
- *What do Michelle Obama do for a living?* ⇒ writer, lawyer
- *What killed Sammy Davis Jr?* ⇒ throat cancer

[Examples from Berant]

- 5,810 questions crawled from Google Suggest API and answered using Amazon MTurk
  - 3,778 training, 2,032 testing
  - A question may have multiple answers → using Avg. F1 (~accuracy)

## Other work: Subgraph Embedding [Bordes+ EMNLP-2014 ]

# Avg. F1 (Accuracy) on WEBQUESTIONS Test Set

(Benchmark leaderboard on Codalab)



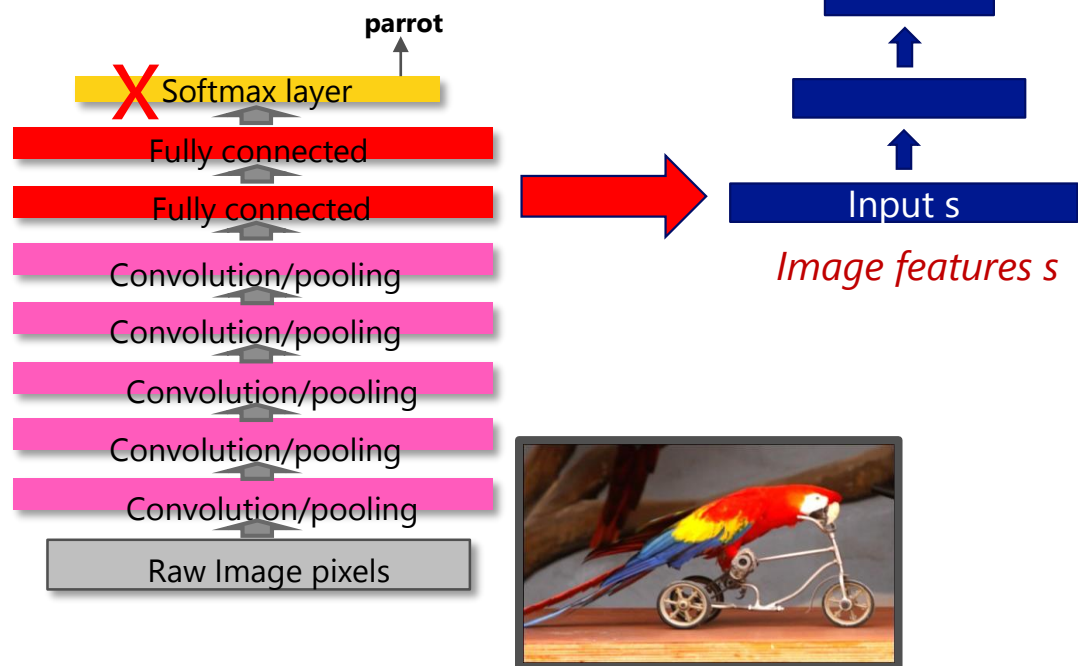| | Avg. F1 |
|---|---|
| Bordes-14a | 29.7 |
| Yao-14 | 33 |
| Berant-13 | 35.7 |
| Bao-14 | 37.5 |
| Bordes-14b | 39.2 |
| Berant-14 | 39.9 |
| Yang-14 | 41.3 |
| Wang-14 | 45.3 |
| Yih-15 | 52.5 |

Yih, Chang, He, and Gao, Semantic Parsing via Staged Query Graph Generation:
Question Answering with Knowledge Base, ACL, July 2015

# Deep Multimodal Similarity Model (DMSM)
## Multimodal DSSM for image-text joint learning

- Recall DSSM for text inputs: *s, t*
- Now: replace text **s** by image **s**
- Pick complete captions affinitize to complete images

Distance(s,t)

**parrot**

X Softmax layer

Fully connected

Fully connected

Convolution/pooling

Convolution/pooling

Convolution/pooling

Convolution/pooling

Convolution/pooling

Raw Image pixels

Input s

*Image features s*

Input t1

Text: *a parrot riding a tricycle*

$Q$ = image, $D$ = caption, $R$ = relevance

**Relevance:** $R(Q, D) = \mathrm{cosine}(y_Q, y_D) = \dfrac{y_Q^T y_D}{\|y_Q\|\|y_D\|}$

**Caption probability:** $P(D|Q) = \dfrac{\exp(\gamma R(Q, D))}{\Sigma_{D' \in \mathbb{D}} \exp(\gamma R(Q, D'))}$

Candidate captions ↗  ↖ Smoothing factor

**Objective:** $L(\Lambda) = -\log \displaystyle\prod_{(Q, D^+)} P(D^+|Q)$

↖ Correct caption

Deep Semantic Learning: Teach machines to
text, image, and knowledge grap

# The convolutional network at the image side

Feed the pre-trained image feature vector into the image side of the DMSM



Dense feature vector for input image

Raw pixels from input box

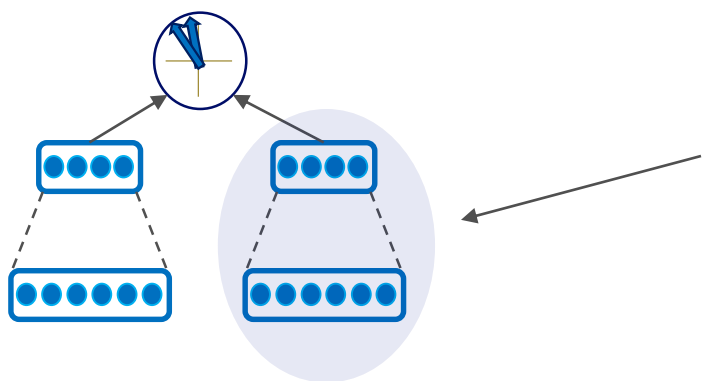Trained to predict object in image

Pretrained from ImageNet [Krizhevsky et al., 2012]

# The convolutional network at the caption side

Models fine-grained structural language information in the caption



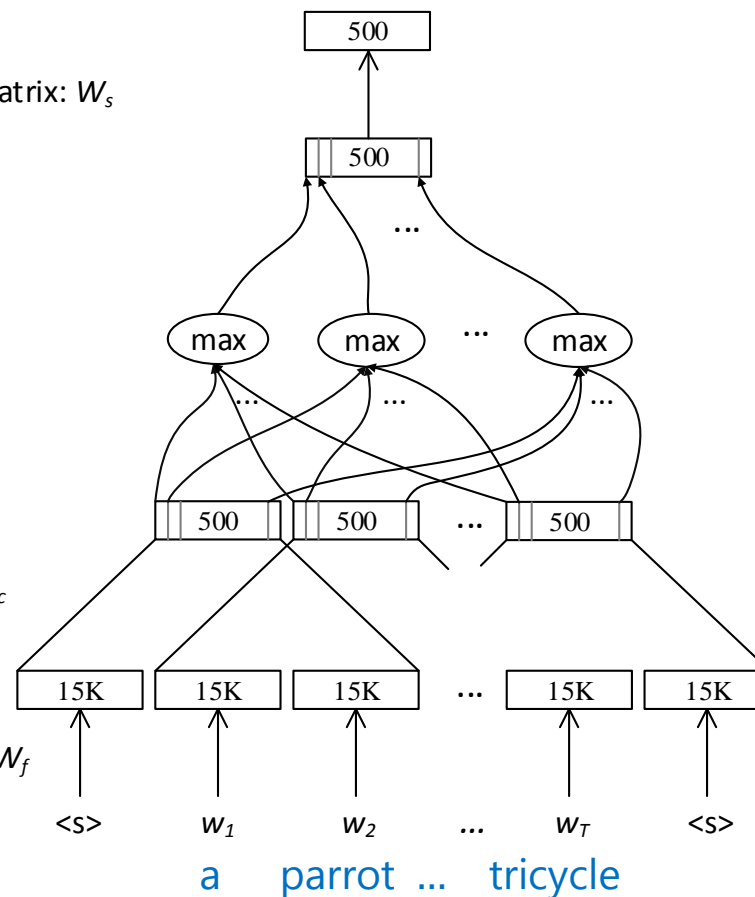Semantic layer: $y$

Semantic projection matrix: $W_s$

Max pooling layer: $v$

Max pooling operation

Convolutional layer: $h_t$

Convolution matrix: $W_c$

Word hashing layer: $f_t$

Word hashing matrix: $W_f$

Word sequence: $x_t$

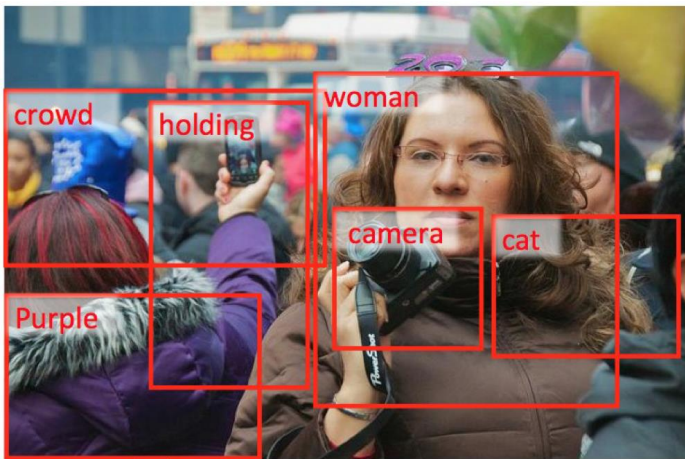Using convolutional neural network for the text caption side

# The task: Image -> Language

- ## Why important?
  For building intelligent machines that understand the semantics in complex scenes

- ## Why difficult?
  Need to capture the salient, coherent semantic information embedded in a picture.



➡ a woman holding a camera in a crowd.

# The MSR system

## Understand the image stage by stage:

## Image word detection

Deep-learned features, applied to likely items in the image, trained to produce words in captions

## Language generation

Maxent language model, trained on caption, conditional on words detected from the image

## Global semantic re-ranking

Hypothetical captions re-ranked by deep-learned multi-modal similarity model looking at the entire image
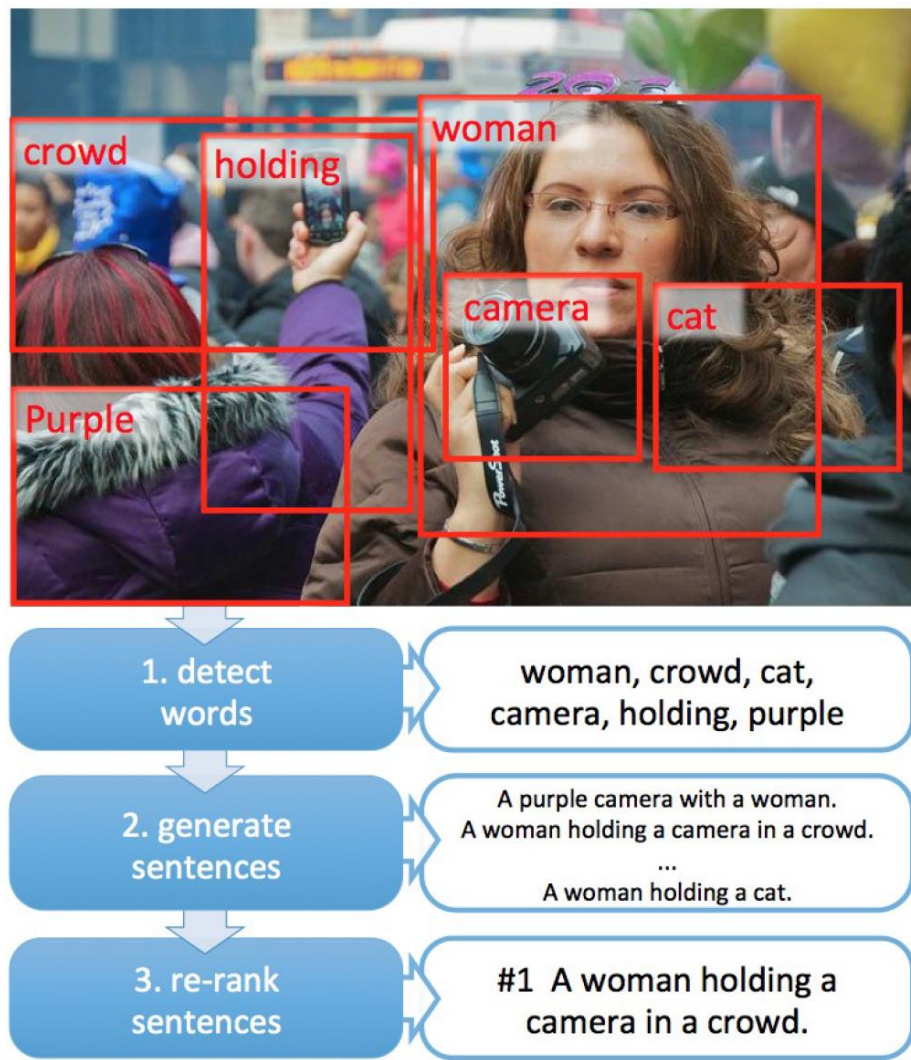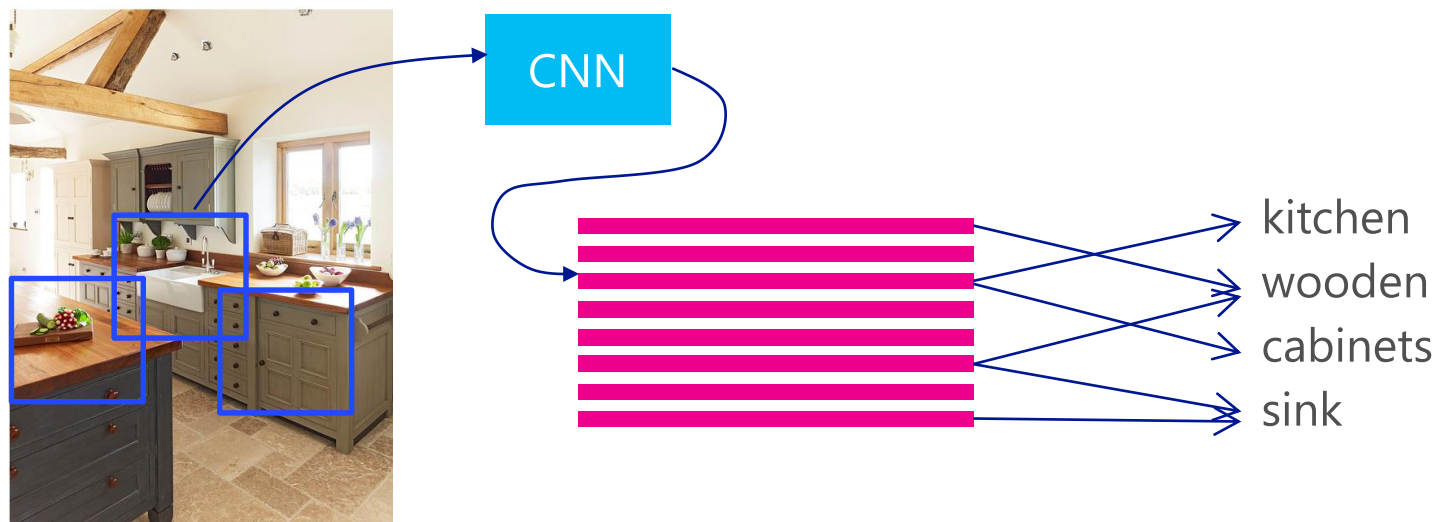


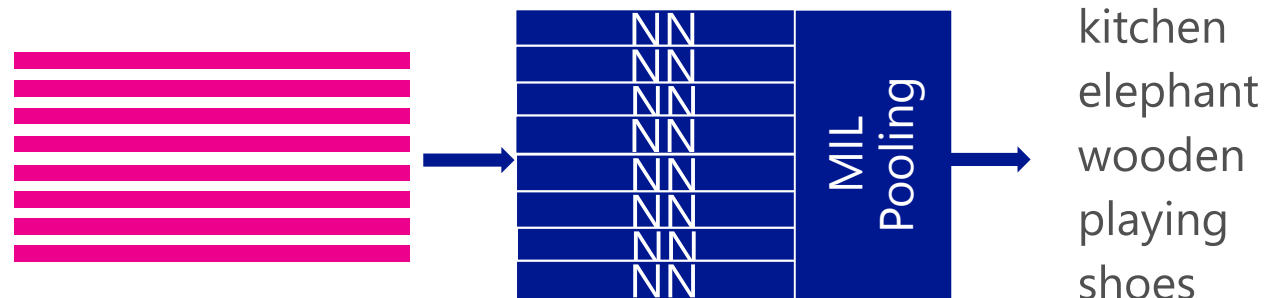Figure 1. An illustrative example of our pipeline.

# Train to predict words in captions



Which words should be detected? Let a neural network figure it out

The prob that the j-th box of the i-th image corresponds to word $w$ is

$$p_{ij}^w = \frac{1}{1 + \exp\left(-(\mathbf{v_w^t}\phi(b_{ij}) + u_w)\right)}$$



Vocabulary = the 1000 most common words in the training captions (92% of data)

# Map features to likely image words

- Train with Multiple Instance Learning (MIL)
  - Use noisy-OR version (Zhang et al., 2005)

- For each word $w$, MIL uses positive and negative bags of bounding boxes
  - For each image $i$:
    - We have the "bag of boxes", $b_i$
    - $b_i$ is **positive** if $w$ in $i$'s description
    - $b_i$ is **negative** if $w$ not in $i$'s description
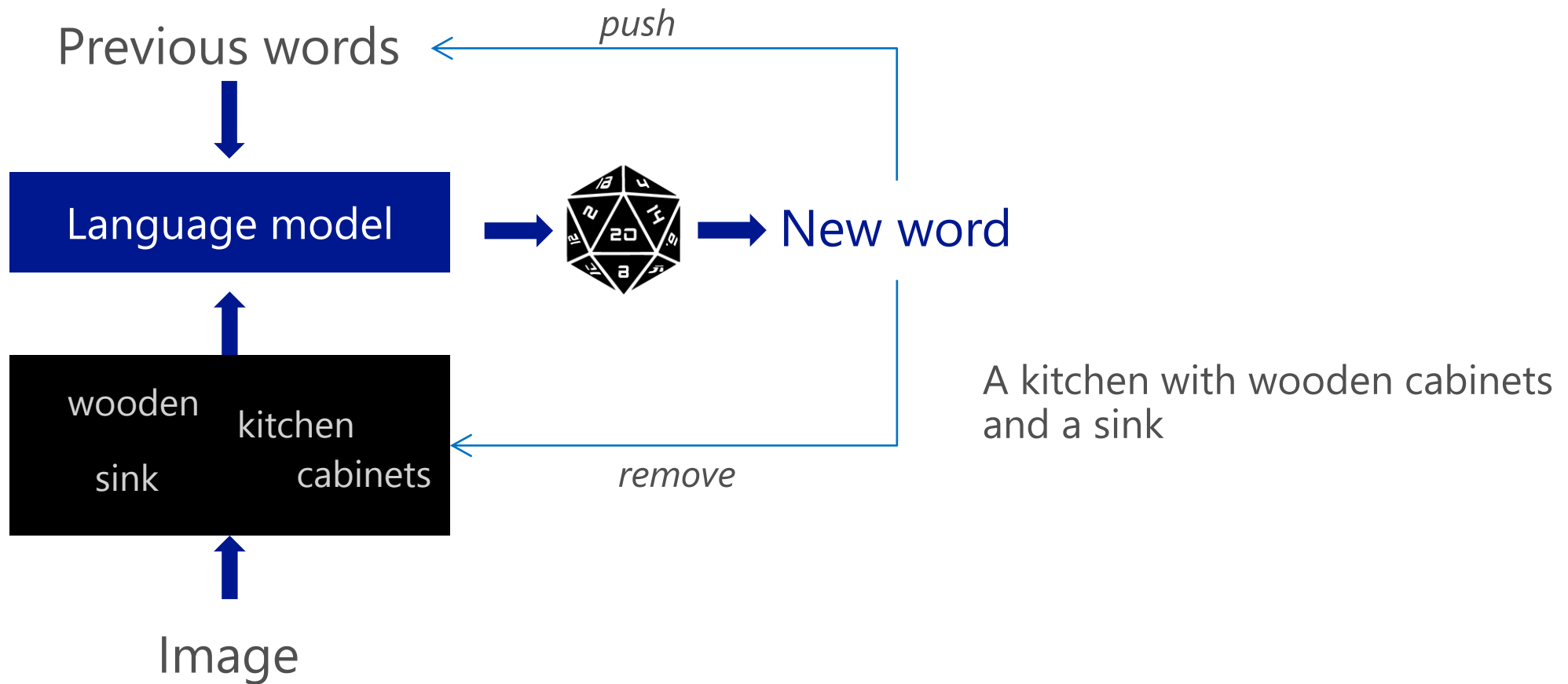  - Probability that image $i$ manifests word $w$, $p_i^w$:

$$p_i^w = 1 - \prod_{j \in b_i} \left(1 - p_{ij}^w\right)$$

Each bounding box in image $\longrightarrow$          Calculated from CNN (last slide)

# Language models with a blackboard

A LM generates 500 caption candidates given detected words

Previous words

*push*

Language model → 🎲 → New word

wooden kitchen sink cabinets

*remove*

A kitchen with wooden cabinets and a sink

Image

# Rerank hypotheses globally using DMSM

Top 500 hypotheses from the language model

A man sitting on a bench

A man sitting on a table

A white bench sitting on top of a table

A man sitting at a table with plates of food

Single hypothesis → Sentence-level features →

Whole image

Similarity (from **DMSM** neural net) →
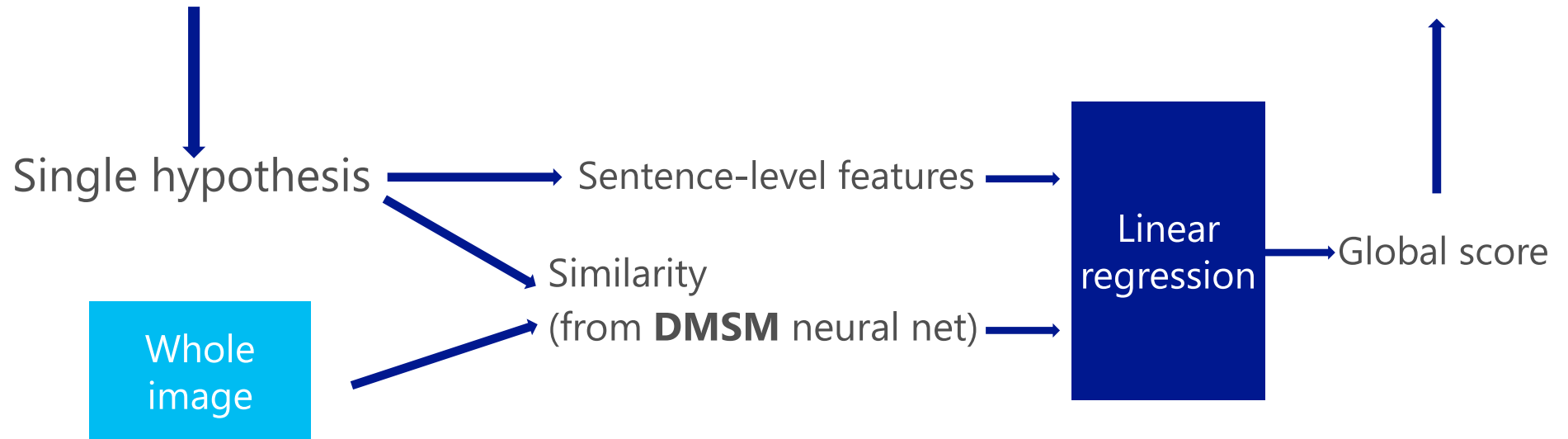
Linear regression → Global score

Return best hypothesis

Image features from AlexNet (Krizhevsky et al., 2012) or VGG (Simonyan and Zisserman, 2014). They are fine-tuned with in-domain image data for DMSM

# The MS COCO Benchmark

**Microsoft COCO**
Common Objects in Context

## What is Microsoft COCO?

Microsoft COCO is a new image recognition, segmentation, and captioning dataset. Microsoft COCO has several features:

- ✔ **Object segmentation**
- ✔ **Recognition in Context**
- ✔ **Multiple objects per image**
- ✔ **More than 300,000 images**
- ✔ **More than 2 Million instances**
- ✔ **80 object categories**
- ✔ **5 captions per image**

## Collaborators

**Tsung-Yi Lin** Cornell Tech

**Michael Maire** TTI Chicago

**Serge Belongie** Cornell Tech

**Lubomir Bourdev** Facebook AI

**Ross Girshick** Microsoft Research

**James Hays** Brown University

**Pietro Perona** Caltech

**Deva Ramanan** UC Irvine

**Larry Zitnick** Microsoft Research

**Piotr Dollár** Facebook AI

CORNELL NYCTECH

Caltech

facebook

Brown University

UCIrvine
University of California, Irvine

Microsoft Research

The man at bat readies to swing at the pitch while the umpire looks on.

A large bus sitting next to a very tall building.

Xiaodong He

Deep Sem

# Results

| System | PPLX | BLEU | METEOR | ≈human | >human | ≥human |
|---|---|---|---|---|---|---|
| 1. Unconditioned | 24.1 | 1.2% | 6.8% | | | |
| 2. Shuffled Human | – | 1.7% | 7.3% | | | |
| 3. Baseline | 20.9 | 16.9% | 18.9% | 9.9% (±1.5%) | 2.4% (±0.8%) | 12.3% (±1.6%) |
| 4. Baseline+Score | 20.2 | 20.1% | 20.5% | 16.9% (±2.0%) | 3.9% (±1.0%) | 20.8% (±2.2%) |
| 5. Baseline+Score+DMSM | 20.2 | 21.1% | 20.7% | 18.7% (±2.1%) | 4.6% (±1.1%) | 23.3% (±2.3%) |
| 6. Baseline+Score+DMSM+ft | 19.2 | 23.3% | 22.2% | – | – | – |
| 7. VGG+Score+ft | 18.1 | 23.6% | 22.8% | – | – | – |
| 8. VGG+Score+DMSM+ft | 18.1 | 25.7% | 23.6% | 26.2% (±2.1%) | 7.8% (±1.3%) | **34.0% (±2.5%)** |
| Human-written captions | – | 19.3% | 24.1% | | | |

\* we use 4 references when measuring BLEU and METEOR, while the official COCO eval server uses 5 references.

DMSM gives additional 2.1 pt BLEU over a strong system (e.g., #8 vs. #7).
Also show significant improvement by human judge (e.g., #5 vs. #4)

# Related work

Use CNN to generate a whole-image feature vector,
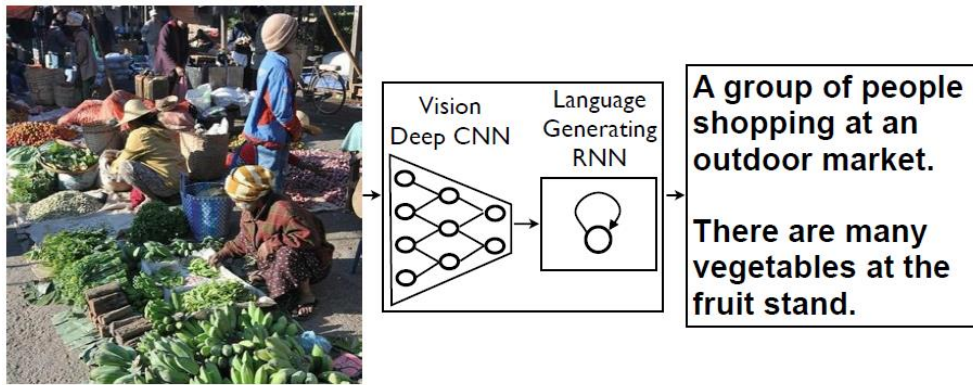then feed it into a LSTM language model to generate the caption.



Figure 1. NIC, our model, is based end-to-end on a neural network consisting of a vision CNN followed by a language generating RNN. It generates complete sentences in natural language from an input image, as shown on the example above.
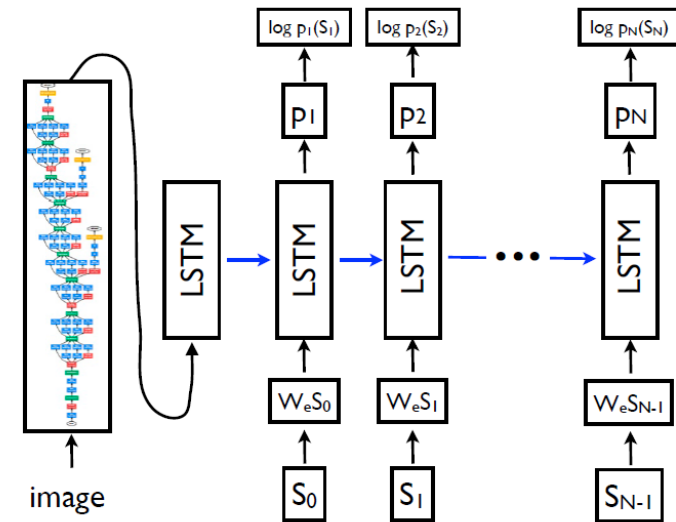


Figure 3. LSTM model combined with a CNN image embedder (as defined in [30]) and word embeddings. The unrolled connections between the LSTM memories are in blue and they correspond to the recurrent connections in Figure 2. All LSTMs share the same parameters.

Vinyals, Toshev, Bengio, Erhan, "Show and Tell: A Neural Image Caption Generator", CVPR 2015

# Some other related work

Andrej and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions". CVPR 2015
Use CNN to generate an image feature vector, then input it, at the 1st step, into a multimodal RNN language model to generate the caption.

Kiros, Salakhutdinov, Zemel, "Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models". TACL 2015
Use LSTM for image-language encoding and decoding

Mao, Xu, Yang, Wang, Huang, Yuille. "Deep Captioning with Multimodal Recurrent Neural Networks (m-RNN)," ICLR 2015
Use CNN to generate a whole-image feature vector, then input it, at every step, into a multimodal RNN language model to generate the caption.

Xu, Ba, Kiros, Cho, Courville, Salakhutdinov, Zemel, Bengio, 2015. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention.
Use CNN to generate a whole-image feature vector, then input it, at every step, into a multimodal RNN language model to generate the caption.

Hill and Korhonen, 2014 Learning Abstract Concept Embeddings from Multi-Modal Data: Since You Probably Can't See What I Mean

CVPR 2015

DeepVision: *Deep Learning in Computer Vision 2015*

# MS COCO Image Captioning Challenge 2015

Measure the quality of the captions by human judge.

| | % of ≥ human | % passing Turing Test | Overall rank |
|---|---|---|---|
| Human | 63.8% | 67.5% | |
| MSR          [Fang+ 15] | 26.8% | 32.2% | 1st (tie) |
| Google       [Vinyals+ 15] | 27.3% | 31.7% | 1st (tie) |
| MSR Captivator   [Devlin+ 15] | 25.0% | 30.1% | 3rd (tie) |
| Montreal/Toronto    [Xu+ 15] | 26.2% | 27.2% | 3rd (tie) |
| Berkeley LRCN [Donahue+ 15] | 24.6% | 26.8% | 5th |

http://mscoco.org/dataset/#leaderboard-cap

# m-DSSM gives the global semantically matching caption for a given image



**Baseline**: a large jetliner sitting on top of a stop sign at an intersection on a city street

**w/ m-DSSM**: a stop light on a city street

**Baseline**: a clock tower in front of a building

**w/ m-DSSM**: a clock tower in the middle of the street

**Baseline**: a red brick building

**w/ m-DSSM**: a living room filled with furniture and a flat screen tv sitting on top of a brick building

**Baseline**: a large jetliner sitting on top of a table

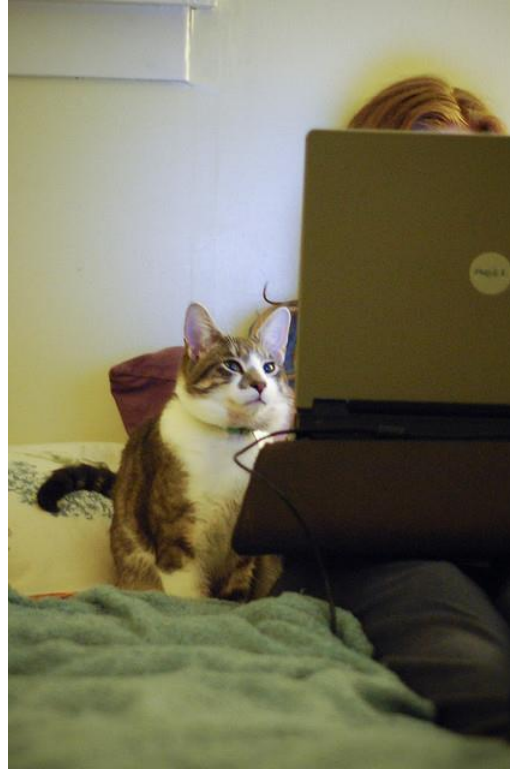**w/ m-DSSM**: a display in a grocery store filled with lots of food on a table

mantic text,

# m-DSSM gives the global semantically matching caption for a given image



**Baseline**: a young man riding a skateboard down a street holding a tennis racquet on a tennis court

**w/ m-DSSM**: a man riding a skateboard down a street



**Baseline**: a cat sitting on a table

**w/ m-DSSM**: a cat sitting on top of a bed



**Baseline**: a group of people standing in a kitchen

**w/ m-DSSM**: a group of people posing for a picture



**Baseline**: two elephants standing next to a baby elephant walking behind a fence

**w/ m-DSSM**: a baby elephant standing next to a fence

CVPR 2015

DeepVision: *Deep Learning in Computer Vision 2015*

# Interpretability



Our system not only generates the caption, but can also interpret it.

# Interpretability



Our system not only generates the caption, but can also interpret it.

# Interpretability



baseball (1.00)

a **baseball**

Our system not only generates the caption, but can also interpret it.

# Interpretability



player (1.00)

a baseball **player**

Our system not only generates the caption, but can also interpret it.
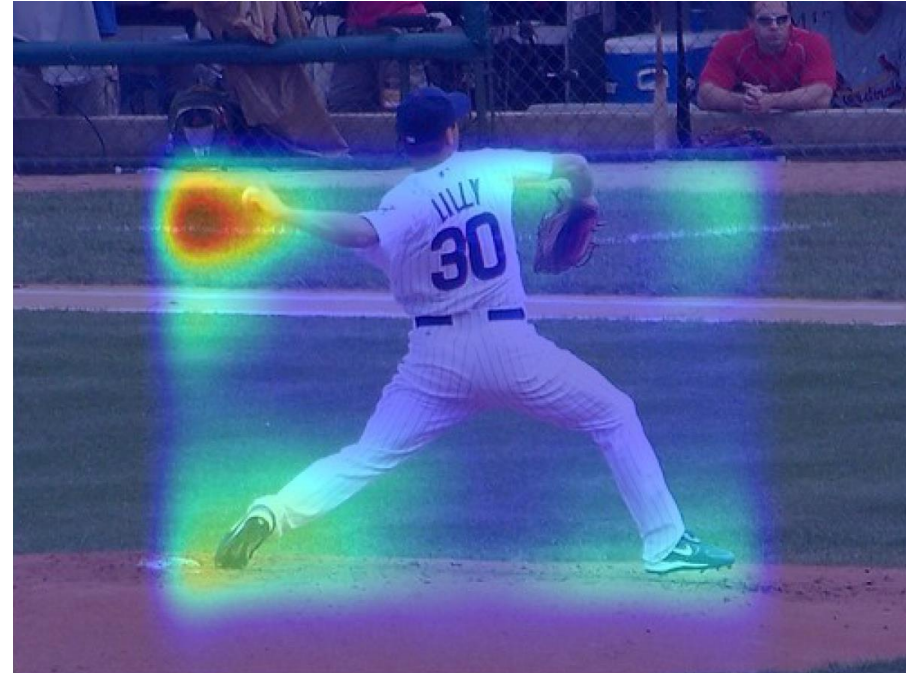
# Interpretability



throwing (0.86)

a baseball player **throwing**

Our system not only generates the caption, but can also interpret it.

# Interpretability



ball (1.00)

a baseball player throwing a **ball**

Our system not only generates the caption, but can also interpret it.
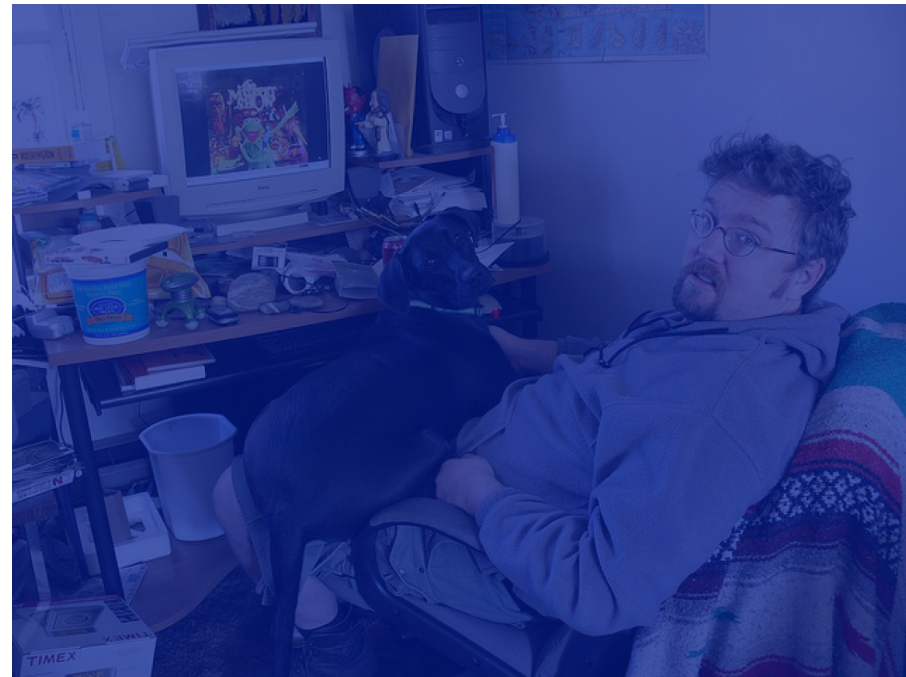
# Interpretability



Our system not only generates the caption, but can also interpret it.

# Interpretability



Our system not only generates the caption, but can also interpret it.

# Interpretability



man (0.93)

a **man**

Our system not only generates the caption, but can also interpret it.

# Interpretability



sitting (0.83)

a man **sitting**

Our system not only generates the caption, but can also interpret it.
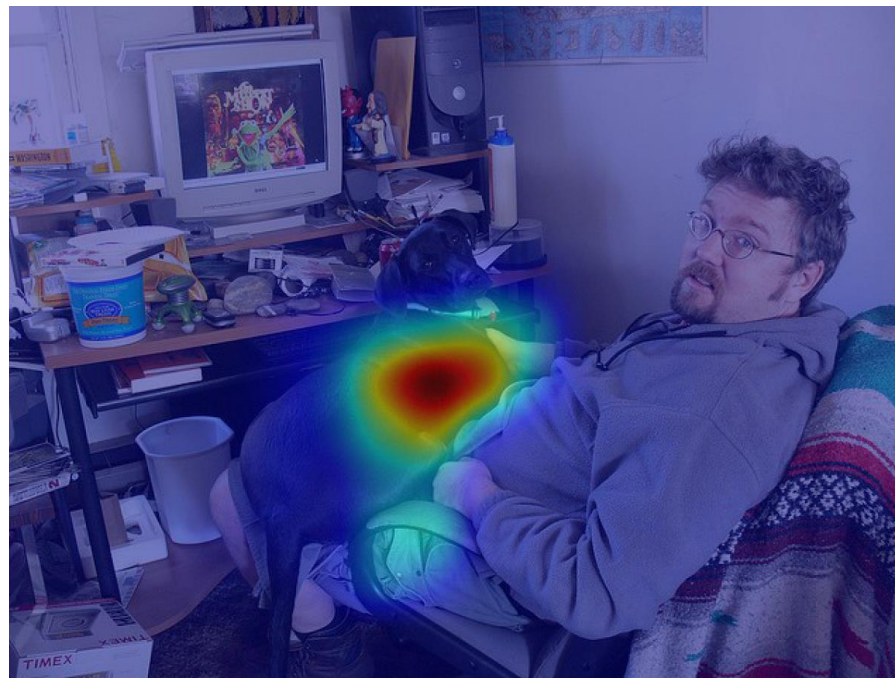
# Interpretability



couch (0.66)

a man sitting in a **couch**

Our system not only generates the caption, but can also interpret it.

# Interpretability



dog (1.00)
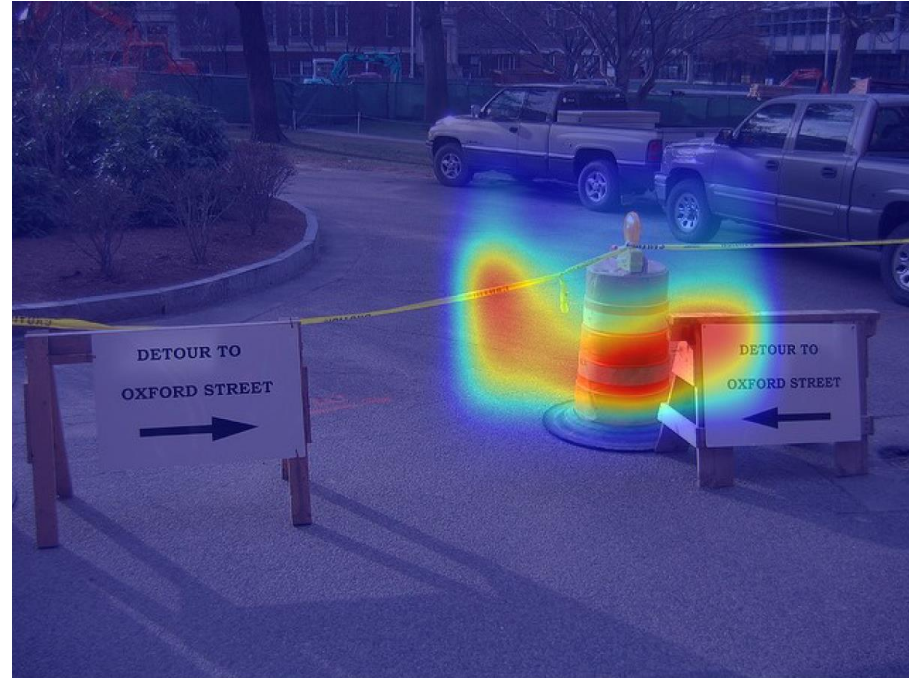
a man sitting in a couch with a **dog**

# Interpretability

# Interpretability
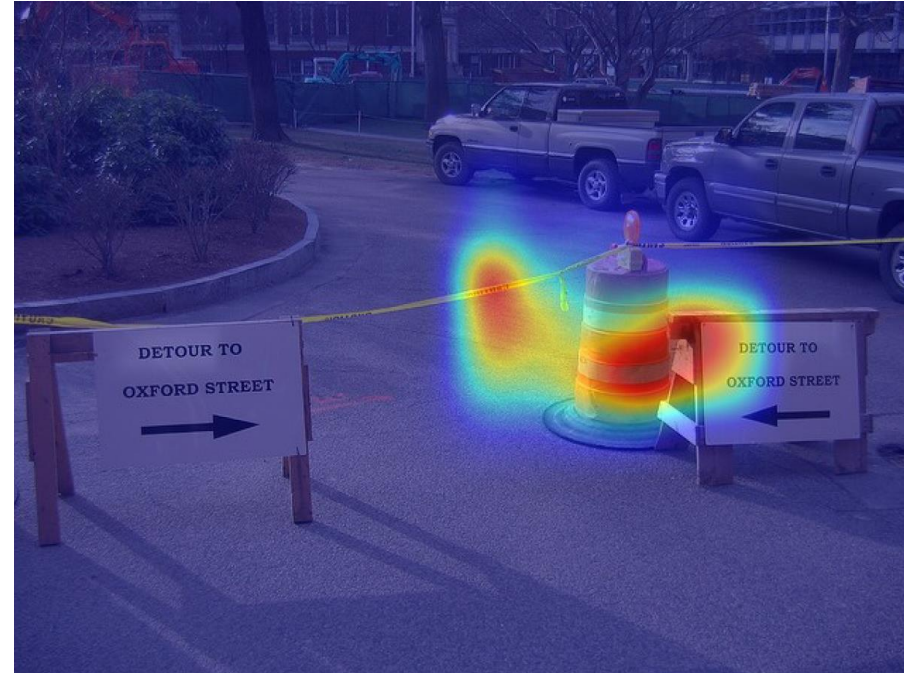
# Interpretability



fire (1.00)

a **fire**

# Interpretability



hydrant (1.00)

a fire **hydrant**

# Interpretability



city (0.69)

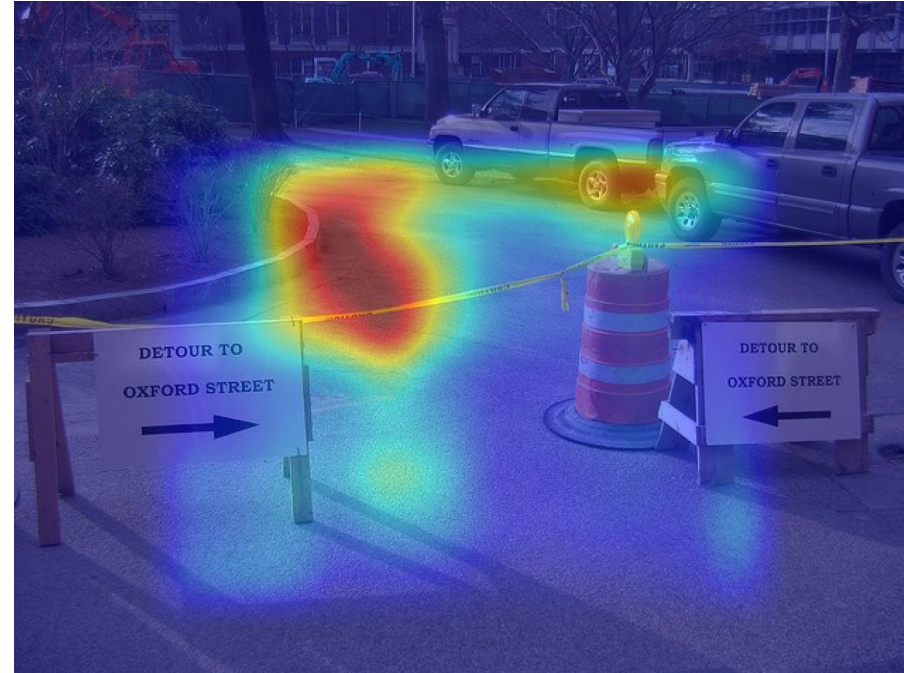a fire hydrant on a **city**

# Interpretability



street (1.00)

a fire hydrant on a city **street**

# Summary

Exciting advances in learning semantic meaning representations

Text, Image, and Knowledge

Sent2Vec Tool kit available:  http://aka.ms/sent2vec/

Looking forward

Building an universal intelligence space

Text, Image, Knowledge, Reasoning,…

From component models to end-to-end solutions

# Acknowledgement:

# Thanks!

## & selected publications from MSR-DLTC

- J. Devlin, H. Cheng, H. Fang, S. Gupta, L. Deng, X. He, G. Zweig, M. Mitchell, Language Models for Image Captioning: The Quirks and What Works, ACL, July 2015

- W. Yih, M. Chang, X. He, and J. Gao, Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base, ACL, July 2015

- H. Fang, S. Gupta, F. Iandola, R. Srivastava, L. Deng, P. Dollar, J. Gao, X. He, M. Mitchell, J. Platt, L. Zitnick, and G. Zweig, From Captions to Visual Concepts and Back, CVPR, June 2015

- X. Liu, J. Gao, X. He, L. Deng, Kevin Duh, and Y. Wang, Representation Learning Using Multi-Task Deep Neural Networks for Semantic Classification and Information Retrieval, NAACL, May 2015

- A. Elkahky, Y. Song, and X. He, A Multi-View Deep Learning Approach for Cross Domain User Modeling in Recommendation Systems, WWW, May 2015

- B. Yang, W. Yih, X. He, J. Gao, and L. Deng, Embedding Entities and Relations for Learning and Inference in Knowledge Bases, ICLR, May 2015

- G. Mesnil, Y. Dauphin, K. Yao, Y. Bengio, L. Deng, D. Hakkani-Tur, X. He, L. Heck, G. Tur, D. Yu, and G. Zweig, Using Recurrent Neural Networks for Slot Filling in Spoken Language Understanding, IEEE/ACM TASLP, March 2015

- H. Palangi, L. Deng, Y. Shen, J. Gao, X. He, J. Chen, X. Song, and R. Ward, Deep Sentence Embedding Using the Long Short Term Memory Network: Analysis and Application to Information Retrieval, *arXiv:1502.06922*, February 2015

- Y. Shen, X. He, J. Gao, L. Deng, and G. Mesnil, A Latent Semantic Model with Convolutional-Pooling Structure for Information Retrieval , CIKM, November 2014

- J. Gao, P. Pantel, M. Gamon, X. He, L. Deng, et al., Modeling Interestingness with Deep Neural Networks, EMNLP, October 2014

- W. Yih, X. He, and C. Meek, Semantic Parsing for Single-Relation Question Answering, ACL, June 2014

- J. Gao, X. He, W. Yih, and L. Deng, Learning Continuous Phrase Representations for Translation Modeling, ACL, June 2014

- P. Huang, X. He, J. Gao, L. Deng, A. Acero, and L. Heck, Learning Deep Structured Semantic Models for Web Search using Clickthrough Data, CIKM, October 2013

- G. Mesnil, X. He, L. Deng, and Y. Bengio, Investigation of Recurrent-Neural-Network Architectures and Learning Methods for Spoken Language Understanding, in *Interspeech*, August 2013

- G. Tur, L. Deng, D. Hakkani-Tur, and X. He, Towards Deeper Understanding Deep Convex Networks for Semantic Utterance Classification, ICASSP, March 2012

# References

- Auli, M., Galley, M., Quirk, C. and Zweig, G., 2013. Joint language and translation modeling with recurrent neural networks. In EMNLP.
- Auli, M., and Gao, J., 2014. Decoder integration and expected bleu training for recurrent neural network language models. In ACL.
- Bahdanau, D., Cho, K., and Bengio, Y. 2015. Neural machine translation by joingly learning to align and translate, in ICLR 2015.
- Bejar, I., Chaffin, R. and Embretson, S. 1991. Cognitive and psychometric analysis of analogical problem solving. Recent research in psychology.
- Bengio, Y., 2009. Learning deep architectures for AI. Foundumental Trends Machine Learning, vol. 2.
- Bengio, Y., Courville, A., and Vincent, P. 2013. Representation learning: A review and new perspectives. IEEE Trans. PAMI, vol. 38, pp. 1798-1828.
- Bengio, Y., Ducharme, R., and Vincent, P., 2000. A Neural Probabilistic Language Model, in NIPS.
- Berant, J., and Liang, P. 2014. Semantic parsing via paraphrasing. In ACL.
- Blei, D., Ng, A., and Jordan M. 2001. Latent dirichlet allocation. In NIPS.
- Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J. and Yakhnenko, O. 2013. Translating Embeddings for Modeling Multi-relational Data. In NIPS.
- Bordes, A., Chopra, S., and Weston, J. 2014. Question answering with subgraph embeddings. In EMNLP.
- Bordes, A., Glorot, X., Weston, J. and Bengio Y. 2012. Joint Learning of Words and Meaning Representations for Open-Text Semantic Parsing. In AISTATS.
- Brown, P., deSouza, P. Mercer, R., Della Pietra, V., and Lai, J. 1992. Class-based n-gram models of natural language. Computational Linguistics 18 (4).
- Chandar, A. P. S., Lauly, S., Larochelle, H., Khapra, M. M., Ravindran, B., Raykar, V., and Saha, A. (2014). An autoencoder approach to learning bilingual word representations. In NIPS.
- Chang, K., Yih, W., and Meek, C. 2013. Multi-Relational Latent Semantic Analysis. In EMNLP.
- Chang, K., Yih, W., Yang, B., and Meek, C. 2014. Typed Tensor Decomposition of Knowledge Bases for Relation Extraction. In EMNLP.
- Collobert, R., and Weston, J. 2008. A Unified Architecture for Natural Language Processing: Deep Neural Networks with Multitask Learning. In ICML.
- Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P., 2011. Natural language processing (almost) from scratch. in JMLR, vol. 12.
- Cui, L., Zhang, D., Liu, S., Chen, Q., Li, M., Zhou, M., and Yang, M. (2014). Learning topic representation for smt with neural networks. In ACL.
- Dahl, G., Yu, D., Deng, L., and Acero, 2012. A. Context-dependent, pre-trained deep neural networks for large vocabulary speech recognition, IEEE Trans. Audio, Speech, & Language Proc., Vol. 20 (1), pp. 30-42.
- Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T., and Harshman, R. 1990. Indexing by latent semantic analysis. J. American Society for Information Science, 41(6): 391-407
- Deng, L., Seltzer, M., Yu, D., Acero, A., Mohamed, A., and Hinton, G., 2010. Binary Coding of Speech Spectrograms Using a Deep Auto-encoder, in Interspeech.
- Deng, L., Tur, G, He, X, and Hakkani-Tur, D. 2012. Use of kernel deep convex networks and end-to-end learning for spoken language understanding, Proc. IEEE Workshop on Spoken Language Technologies.
- Deng, L., Yu, D. and Acero, A. 2006. Structured speech modeling, IEEE Trans. on Audio, Speech and Language Processing, vol. 14, no. 5, pp. 1492-1504.

# References

- Deng, L., Yu, D. and Acero, A. 2006. Structured speech modeling, IEEE Trans. on Audio, Speech and Language Processing, vol. 14, no. 5, pp. 1492-1504.
- Deng, L., Yu, D., and Platt, J. 2012. Scalable stacking and learning for building deep architectures, Proc. ICASSP.
- Deng, L. and Yu, D. 2014. Deeping learning methods and applications. Foundations and Trends in Signal Processing 7:3-4.
- Deoras, A., and Sarikaya, R., 2013. Deep belief network based semantic taggers for spoken language understanding, in INTERSPEECH.
- Devlin, J., Zbib, R., Huang, Z., Lamar, T., Schwartz, R., and Makhoul, J., 2014. Fast and Robust Neural Network Joint Models for Statistical Machine Translation, ACL.
- Duh, K. 2014. Deep learning for natural language processing and machine translation. Tutorial. 2014.
- Duh, K., Neubig, G., Sudoh, K., and Tsukada, H. (2013). Adaptation data selection using neural language models: Experiments in machine translation. In ACL.
- Fader, A., Zettlemoyer, L., and Etzioni, O. 2013. Paraphrase-driven learning for open question answering. In ACL.
- Fang, H., Gupta, S., Iandola, F., Srivastava, R., Deng, L., Dollár, P., Gao, J., He, X., Mitchell, M., Platt, J., Zitnick, L., Zweig, G., "From Captions to Visual Concepts and Back," arXiv:1411.4952
- Faruqui, M. and Dyer, C. (2014). Improving vector space word representations using multilingual correlation. In EACL.
- Firth, J. R. 1957. *Papers in Linguistics 1934–1951*, Oxford University Press, 1957
- Frome, A., Corrado, G., Shlens, J., Bengio, S., Dean, J., Ranzato, M., and Mikolov, T., 2013. DeViSE: A Deep Visual-Semantic Embedding Model, Proc. NIPS.
- Gao, J., He, X., Yih, W-t., and Deng, L. 2014a. Learning continuous phrase representations for translation modeling. In ACL.
- Gao, J., He, X., and Nie, J-Y. 2010. Clickthrough-based translation models for web search: from word models to phrase models. In CIKM.
- Gao, J., Pantel, P., Gamon, M., He, X., Deng, L., and Shen, Y. 2014b. Modeling interestingness with deep neural networks. In EMNLP
- Gao, J., Toutanova, K., Yih., W-T. 2011. Clickthrough-based latent semantic models for web search. In SIGIR.
- Gao, J., Yuan, W., Li, X., Deng, K., and Nie, J-Y. 2009. Smoothing clickthrough data for web search ranking. In SIGIR.
- Gao, J., and He, X. 2013. Training MRF-based translation models using gradient ascent. In NAACL-HLT.
- Getoor, L., and Taskar, B. editors. 2007. Introduction to Statistical Relational Learning. The MIT Press.
- Graves, A., Jaitly, N., and Mohamed, A., 2013a. Hybrid speech recognition with deep bidirectional LSTM, Proc. ASRU.
- Graves, A., Mohamed, A., and Hinton, G., 2013. Speech recognition with deep recurrent neural networks, Proc. ICASSP.
- He, X. and Deng, L., 2013. Speech-Centric Information Processing: An Optimization-Oriented Approach, in Proceedings of the IEEE.
- He, X. and Deng, L., 2012. Maximum Expected BLEU Training of Phrase and Lexicon Translation Models , ACL.
- He, X., Deng, L., and Chou, W., 2008. Discriminative learning in sequential pattern recognition, Sept. IEEE Sig. Proc. Mag.

# References

- Hermann, K. M. and Blunsom, P. (2014). Multilingual models for compositional distributed semantics. In ACL.
- Hinton, G., Deng, L., Yu, D., Dahl, G., Mohamed, A., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T., and Kingsbury, B., 2012. Deep Neural Networks for Acoustic Modeling in Speech Recognition, IEEE Signal Processing Magazine, vol. 29, no. 6, pp. 82-97.
- Hinton, G., Osindero, S., and The, Y-W. 2006. A fast learning algorithm for deep belief nets. *Neural Computation*, 18: 1527-1554.
- Hinton, G., and Salakhutdinov, R., 2010. Discovering binary codes for documents by learning deep generative models. Topics in Cognitive Science.
- Hu, Y., Auli, M., Gao, Q., and Gao, J. 2014. Minimum translation modeling with recurrent neural networks. In EACL.
- Huang, E., Socher, R., Manning, C, and Ng, A. 2012. Improving word representations via global context and multiple word prototypes, Proc. ACL.
- Huang, P., He, X., Gao, J., Deng, L., Acero, A., and Heck, L. 2013. Learning deep structured semantic models for web search using clickthrough data. In CIKM.
- Hutchinson, B., Deng, L., and Yu, D., 2012. A deep architecture with bilinear modeling of hidden representations: Applications to phonetic recognition, Proc. ICASSP.
- Hutchinson, B., Deng, L., and Yu, D., 2013. Tensor deep stacking networks, IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 35, pp. 1944 – 1957.
- Jurgens, D., Mohammad, S., Turney, P. and Holyoak, K. 2012. SemEval-2012 Task 2: Measuring degrees of relational similarity. In SemEval.
- Kalchbrenner, N. and Blunsom, P. (2013). Recurrent continuous translation models., in EMNLLP
- Kiros, R., Zemel, R., and Salakhutdinov, R. 2013. Multimodal Neural Language Models, Proc. NIPS Deep Learning Workshop.
- Klementiev, A., Titov, I., and Bhattarai, B. (2012). Inducing crosslingual distributed representations of words. In COLING.
- Kocisky, T., Hermann, K. M., and Blunsom, P. (2014). Learning bilingual word representations by marginalizing alignments. In ACL.
- Koehn, P. 2009. Statistical Machine Translation. Cambridge University Press.
- Krizhevsky, A., Sutskever, I, and Hinton, G., 2012. ImageNet Classification with Deep Convolutional Neural Networks, NIPS.
- Landauer. T., 2002. On the computational basis of learning and cognition: Arguments from LSA. Psychology of Learning and Motivation, 41:43–84.
- Lao, N., Mitchell, T., and Cohen, W. 2011. Random walk inference and learning in a large scale knowledge base. In EMNLP.
- Lauly, S., Boulanger, A., and Larochelle, H. (2013). Learning multilingual word representations using a bag-of-words autoencoder. In NIPS.
- Le, H-S, Oparin, I., Allauzen, A., Gauvain, J-L., Yvon, F., 2013. Structured output layer neural network language models for speech recognition, IEEE Transactions on Audio, Speech and Language Processing.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. 1998. Gradient-based learning applied to document recognition, Proceedings of the IEEE, Vol. 86, pp. 2278-2324.
- Levy, O., and Goldberg, Y. 2014. Linguistic Regularities in Sparse and Explicit Word Representations. In CoNLL.

# References

- Li, P., Hastie, T., and Church, K.. 2006. Very sparse random projections, in Proc. SIGKDD.
- Li, P., Liu, Y., and Sun, M. (2013). Recursive autoencoders for ITG-based translation. In EMNLP.
- Li, P., Liu, Y., Sun, M., Izuha, T., and Zhang, D. (2014b). A neural reordering model for phrase-based translation. In COLING.
- Liu, S., Yang, N., Li, M., and Zhou, M. (2014). A recursive recurrent neural network for statistical machine translation. In ACL.
- Liu, L., Watanabe, T., Sumita, E., and Zhao, T. (2013). Additive neural networks for statistical machine translation. In ACL.
- Lu, S., Chen, Z., and Xu, B. (2014). Learning new semi-supervised deep auto-encoder features for statistical machine translation. In ACL.
- Maskey, S., and Zhou, B. 2012. Unsupervised deep belief feature for speech translation, in ICASSP.
- Mesnil, G., He, X., Deng, L., and Bengio, Y., 2013. Investigation of Recurrent-Neural-Network Architectures and Learning Methods for Spoken Language Understanding, in Interspeech.
- Mikolov, T. 2012. Statistical Language Models based on Neural Networks, Ph.D. thesis, Brno University of Technology.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. 2013. Efficient estimation of word representations in vector space, Proc. ICLR.
- Mikolov, T., Kombrink,. S., Burget, L., Cernocky, J., Khudanpur, S., 2011. Extensions of Recurrent Neural Network LM. ICASSP.
- Mikolov, T., Yih, W., Zweig, G., 2013. Linguistic Regularities in Continuous Space Word Representations. In NAACL-HLT.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G., and Dean, J. 2013. Distributed Representations of Words and Phrases and their Compositionality. In NIPS.
- Mohamed, A., Yu, D., and Deng, L. 2010. Investigation of full-sequence training of deep belief networks for speech recognition, Proc. Interspeech.
- Mohammad, S., Dorr, Bonnie., and Hirst, G. 2008. Computing word pair antonymy. In EMNLP.
- Ngiam, J., Khosla, A., Kim, M., Nam, J., Lee, H., and Ng, A. 2011. Multimodal deep learning, Proc. ICML.
- Nickel, M., Tresp, V., and Kriegel, H. 2011. A three-way model for collective learning on multi-relational data. In ICML.
- Niehues, J. and Waibel, A. (2013). Continuous space language models using Restricted Boltzmann Machines, in IWLT.
- Reddy, S., Lapata, M., and Steedman, M. 2014. Large-scale semantic parsing without question-answer pairs. Transactions of the Association for Computational Linguistics (TACL).
- Sainath, T., Mohamed, A., Kingsbury, B., and Ramabhadran, B. 2013. Convolutional neural networks for LVCSR, Proc. ICASSP.
- Salakhutdinov R., and Hinton, G., 2007 Semantic hashing. in Proc. SIGIR Workshop Information Retrieval and Applications of Graphical Models
- Salton, G. and McGill, M. 1983. Introduction to Modern Information Retrieval. McGraw Hill.
- Sarikaya, R., Hinton, G., and Ramabhadran, B., 2011. Deep belief nets for natural language call-routing, in Proceedings of the ICASSP.
- Schwenk, H. 2012. Continuous space translation models for phrase-based statistical machine translation, in COLING.

# References

- Schwenk, H., Rousseau, A., and Attik, M., 2012. Large, pruned or continuous space language models on a gpu for statistical machine translation, in NAACL-HLT 2012 Workshop.
- Seide, F., Li, G., and Yu, D. 2011. Conversational speech transcription using context-dependent deep neural networks, Proc. Interspeech
- Shen, Y., He, X., Gao, J., Deng, L., and Mesnil, G. 2014. Learning Semantic Representations Using Convolutional Neural Networks for Web Search, in Proceedings of WWW.
- Shen, Y., He, X., Gao, J., Deng, L., and Mesnil, G. 2014. A convolutional latent semantic model for web search. CIKM
- Socher, R., Chen, D., Manning, C., and Ng, A. 2013. Reasoning With Neural Tensor Networks For Knowledge Base Completion. In NIPS.
- Socher, R., Huval, B., Manning, C., Ng, A., 2012. Semantic compositionality through recursive matrix-vector spaces. In EMNLP.
- Socher, R., Lin, C., Ng, A., and Manning, C. 2011. Learning continuous phrase representations and syntactic parsing with recursive neural networks, Proc. ICML.
- Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C., Ng A., and Potts. C. 2013. Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank, Proc. EMNLP
- Son, L. H., Allauzen, A., and Yvon, F. (2012). Continuous space translation models with neural networks. In NAACL.
- Song, X. He, X., Gao. J., and Deng, L. 2014. Learning Word Embedding Using the DSSM. MSR Tech Report.
- Song, Y., Wang, H., and He, X., 2014. Adapting Deep RankNet for Personalized Search. Proc. WSDM.
- Songyot, T. and Chiang, D. (2014). Improving word alignment using word similarity. In EMNLP.
- Sundermeyer, M., Alkhouli, T., Wuebker, J., and Ney, H. (2014). Translation modeling with bidirectional recurrent neural networks, in EMNLP.
- Tamura, A., Watanabe, T., and Sumita, E. (2014). Recurrent neural networks for word alignment model. In ACL.
- Tran, K. M., Bisazza, A., and Monz, C. (2014). Word translation prediction for morphologically rich languages with bilingual neural networks. In EMNLP.
- Tur, G., Deng, L., Hakkani-Tur, D., and He, X., 2012. Towards Deeper Understanding Deep Convex Networks for Semantic Utterance Classification, in ICASSP.
- Turney P. 2008. A uniform approach to analogies, synonyms, antonyms, and associations. In COLING. Songyot, T. and Chiang, D. (2014). Improving word alignment using word similarity. In EMNLP.
- Vaswani, A., Zhao, Y., Fossum, V., and Chiang, D. 2013. Decoding with large-scale neural language models improves translation, in EMNLP.
- Wright, S., Kanevsky, D., Deng, L., He, X., Heigold, G., and Li, H., 2013. Optimization Algorithms and Applications for Speech and Language Processing, in IEEE Transactions on Audio, Speech, and Language Processing, vol. 21, no. 11.
- Wu, H., Dong, D., Hu, X., Yu, D., He, W., Wu, H., Wang, H., and Liu, T. (2014a). Improve statistical machine translation with context-sensitive bilingual semantic embedding model. In EMNLP.
- Wu, Y., Watanabe, T., and Hori, C. (2014b). Recurrent neural network-based tuple sequence model for machine translation. In COLING.

# References

- Yang, B., Yih, W., He, X., Gao, J., and Deng L. 2014. In NIPS-2014 Workshop Learning Semantics.
- Yang, N., Liu, S., Li, M., Zhou, M., and Yu, N. (2013). Word alignment modeling with context dependent deep neural network. In ACL.
- Yao, K., Zweig, G., Hwang, M-Y. , Shi, Y., Yu, D., 2013. Recurrent neural networks for language understanding, submitted to Interspeech.
- Yann, D., Tur, G., Hakkani-Tur, D., Heck, L., 2014. Zero-Shot Learning and Clustering for Semantic Utterance Classification Using Deep Learning, in ICLR.
- Yih, W., Toutanova, K., Platt, J., and Meek, C. 2011. Learning discriminative projections for text similarity measures. In CoNLL.
- Yih, W., Zweig, G., Platt, J. 2012. Polarity Inducing Latent Semantic Analysis. In EMNLP-CoNLL.
- Yih, W., He, X., Meek, C. 2014. Semantic Parsing for Single-Relation Question Answering, in ACL.
- Yih, W., Chang, M-W., He, X., Gao, J. 2015. Semantic parsing via staged query graph generation: question answering with knowledge base, In ACL.
- Zeiler, M. and Fergus, R. 2013. Visualizing and understanding convolutional networks, arXiv:1311.2901, pp. 1-11.
- Zhang, J., Liu, S., Li, M., Zhou, M., and Zong, C. (2014). Bilingually-constrained phrase embeddings for machine translation. In ACL.
- Zou, W. Y., Socher, R., Cer, D., and Manning, C. D. (2013). Bilingual word embeddings for phrase-based machine translation. In EMNLP.