# CS 2770: Computer Vision
# Vision, Language, Reasoning

Prof. Adriana Kovashka
University of Pittsburgh
March 31, 2020

# Plan for this lecture

- Learning the relation between images and text
  - Recurrent neural networks
  - Applications: Captioning
  - Transformers
- Visual question answering
  - Graph convolutional networks
  - Incorporating knowledge and reasoning

# Motivation: Descriptive Text for Images



"It was an arresting face, pointed of chin, square of jaw. Her eyes were pale green without a touch of hazel, starred with bristly black lashes and slightly tilted at the ends. Above them, her thick black brows slanted upward, cutting a startling oblique line in her magnolia-white skin–that skin so prized by Southern women and so carefully guarded with bonnets, veils and mittens against hot Georgia suns"

Scarlett O'Hara described in Gone with the Wind

Tamara Berg

# Some pre-RNN good results



This is a picture of one sky, one road and one sheep. The gray sky is over the gray road. The gray sheep is by the gray road.



Here we see one road, one sky and one bicycle. The road is near the blue sky, and near the colorful bicycle. The colorful bicycle is within the blue sky.



This is a picture of two dogs. The first dog is near the second furry dog.

# Some pre-RNN bad results

## Missed detections:



Here we see one potted plant.



This is a picture of one dog.

## False detections:



There are one road and one cat. The furry road is in the furry cat.



This is a picture of one tree, one road and one person. The rusty tree is under the red road. The colorful person is near the rusty tree, and under the red road.

## Incorrect attributes:



This is a photograph of two sheeps and one grass. The first black sheep is by the green grass, and by the second black sheep. The second black sheep is by the green grass.



This is a photograph of two horses and one grass. The first feathered horse is within the green grass, and by the second feathered horse. The second feathered horse is within the green grass.

Kulkarni et al., CVPR 2011

# Results with Recurrent Neural Networks



"man in black shirt is playing guitar."

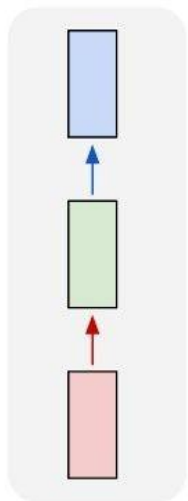"construction worker in orange safety vest is working on road."

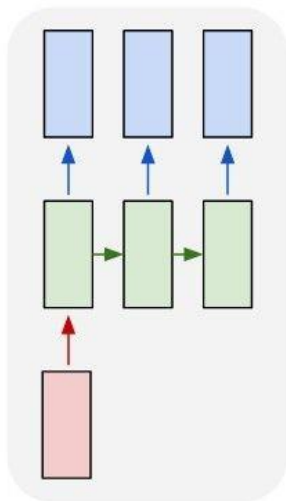"two young girls are playing with lego toy."

"boy is doing backflip on wakeboard."

Karpathy and Fei-Fei, CVPR 2015

# Recurrent Networks offer a lot of flexibility:



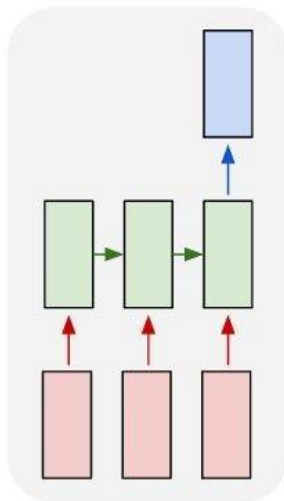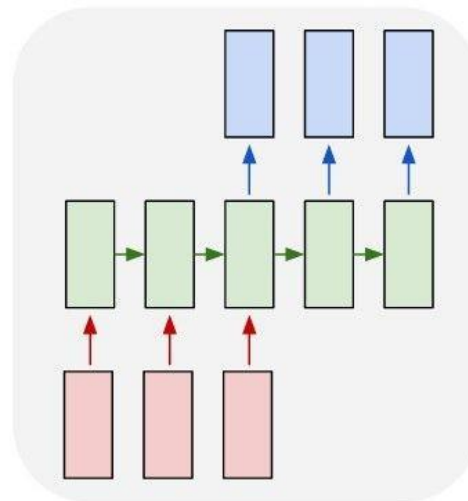one to one    one to many    many to one    many to many    many to many

**vanilla neural networks**

Andrej Karpathy

# Recurrent Networks offer a lot of flexibility:



e.g. **image captioning**
image -> sequence of words

# Recurrent Networks offer a lot of flexibility:



e.g. **sentiment classification**
sequence of words -> sentiment

# Recurrent Networks offer a lot of flexibility:



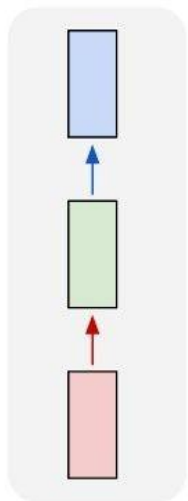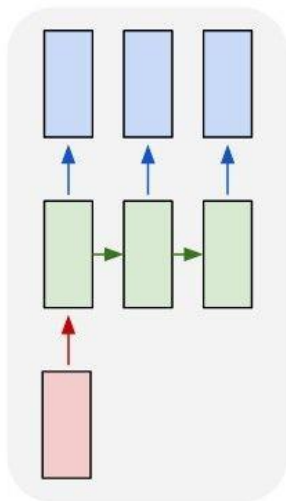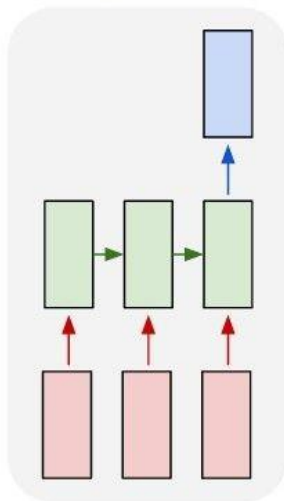one to one   one to many   many to one   many to many   many to many
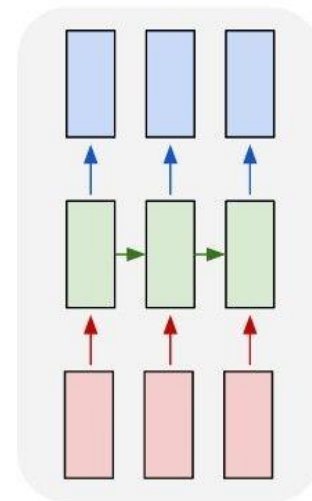
e.g. **machine translation**
seq of words -> seq of words

# Recurrent Networks offer a lot of flexibility:



e.g. **video classification on frame level**

# Recurrent Neural Network

# Recurrent Neural Network



usually want to output a prediction at some time steps

# Recurrent Neural Network

We can process a sequence of vectors **x** by
applying a recurrence formula at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state

some function
with parameters W

old state

input vector at
some time step

# Recurrent Neural Network

We can process a sequence of vectors **x** by
applying a recurrence formula at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

Notice: the same function and the same set
of parameters are used at every time step.

# (Vanilla) Recurrent Neural Network

The state consists of a single *"hidden"* vector **h**:

$$h_t = f_W(h_{t-1}, x_t)$$

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

# Example

**Character-level language model example**

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**

# Example

**Character-level language model example**

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**

# Example

**Character-level language model example**

Vocabulary: [h,e,l,o]

Example training sequence: **"hello"**

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

# Example

**Character-level language model example**

Vocabulary: [h,e,l,o]

Example training sequence: **"hello"**

# The vanishing gradient problem

- The error at a time step ideally can tell a previous time step from many steps away to change during backprop
- But we're multiplying together many values between 0 and 1

# The vanishing gradient problem

- Total error is the sum of each error at time steps t

$$\frac{\partial E}{\partial W} = \sum_{t=1}^{T} \frac{\partial E_t}{\partial W}$$

- Chain rule:

$$\frac{\partial E_t}{\partial W} = \sum_{k=1}^{t} \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W}$$

- More chain rule:

$$\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^{t} \frac{\partial h_j}{\partial h_{j-1}}$$

- Derivative of vector wrt vector is a Jacobian matrix of partial derivatives; norm of this matrix can become very small or very large quickly [Bengio et al 1994, Pascanu et al 2013], leading to vanishing/exploding gradient

Adapted from Richard Socher

# The vanishing gradient problem for language models

- In the case of language modeling or question answering words from time steps far away are not taken into consideration when training to predict the next word

- Example:

  Jane walked into the room. John walked in too. It was late in the day. Jane said hi to _____

Richard Socher

# Gated Recurrent Units (GRUs)

- More complex hidden unit computation in recurrence!

- Introduced by Cho et al. 2014

- Main ideas:

    - keep around memories to capture long distance dependencies

    - allow error messages to flow at different strengths depending on the inputs

Richard Socher

# Gated Recurrent Units (GRUs)

- Standard RNN computes hidden layer at next time step directly:  $h_t = f\left(W^{(hh)}h_{t-1} + W^{(hx)}x_t\right)$

- GRU first computes an update **gate** (another layer) based on current input word vector and hidden state

$$z_t = \sigma\left(W^{(z)}x_t + U^{(z)}h_{t-1}\right)$$

- Compute reset gate similarly but with different weights

$$r_t = \sigma\left(W^{(r)}x_t + U^{(r)}h_{t-1}\right)$$

Richard Socher

# Gated Recurrent Units (GRUs)

- Update gate
$$z_t = \sigma\left(W^{(z)}x_t + U^{(z)}h_{t-1}\right)$$

- Reset gate
$$r_t = \sigma\left(W^{(r)}x_t + U^{(r)}h_{t-1}\right)$$

- New memory content: $\tilde{h}_t = \tanh\left(Wx_t + r_t \circ Uh_{t-1}\right)$
  If reset gate unit is ~0, then this ignores previous memory and only stores the new word information

- Final memory at time step combines current and previous time steps: $h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$

# Gated Recurrent Units (GRUs)



Final memory — $h_{t-1}$ ... $h_t$

Memory (reset) — $\tilde{h}_{t-1}$ ... $\tilde{h}_t$

Update gate — $z_{t-1}$ ... $z_t$

Reset gate — $r_{t-1}$ ... $r_t$

Input: — $x_{t-1}$ ... $x_t$

$$z_t = \sigma \left( W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$

$$r_t = \sigma \left( W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$

$$\tilde{h}_t = \tanh \left( W x_t + r_t \circ U h_{t-1} \right)$$

$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

Richard Socher

# Gated Recurrent Units (GRUs)

- If reset is close to 0, ignore previous hidden state: Allows model to drop information that is irrelevant in the future

$$z_t = \sigma\left(W^{(z)}x_t + U^{(z)}h_{t-1}\right)$$

$$r_t = \sigma\left(W^{(r)}x_t + U^{(r)}h_{t-1}\right)$$

$$\tilde{h}_t = \tanh\left(Wx_t + r_t \circ Uh_{t-1}\right)$$

$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

- Update gate z controls how much of past state should matter now

  - If z close to 1, then we can copy information in that unit through many time steps! **Less vanishing gradient**!

- Units with short-term dependencies often have reset gates (r) very active; ones with long-term dependencies have active update gates (z)

Richard Socher

# Long-short-term-memories (LSTMs)

- Proposed by Hochreiter and Schmidhuber in 1997

- We can make the units even more complex

- Allow each time step to modify

  - Input gate (current cell matters) $\quad i_t = \sigma\left(W^{(i)} x_t + U^{(i)} h_{t-1}\right)$

  - Forget (gate 0, forget past) $\quad f_t = \sigma\left(W^{(f)} x_t + U^{(f)} h_{t-1}\right)$

  - Output (how much cell is exposed) $o_t = \sigma\left(W^{(o)} x_t + U^{(o)} h_{t-1}\right)$

  - New memory cell $\quad \tilde{c}_t = \tanh\left(W^{(c)} x_t + U^{(c)} h_{t-1}\right)$

- Final memory cell: $\quad c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$

- Final hidden state: $\quad h_t = o_t \circ \tanh(c_t)$

Adapted from Richard Socher

# Long-short-term-memories (LSTMs)



Intuition: memory cells can keep information intact, unless inputs makes them forget it or overwrite it with new input

Cell can decide to output this information or just store it

# Plan for this lecture

- Learning the relation between images and text
  - Recurrent neural networks
  - Applications: Captioning
  - Transformers
- Visual question answering
  - Graph convolutional networks
  - Incorporating knowledge and reasoning

# Generating poetry with RNNs

# Generating poetry with RNNs

at first:

```
tyntd-iafhatawiaoihrdemot  lytdws  e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tklrgd t o idoe ns,smtt   h ne etie h,hregtrs nigtike,aoaenns lng
```

↓ train more

```
"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."
```

↓ train more

```
Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.
```

↓ train more

```
"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.
```

More info: http://karpathy.github.io/2015/05/21/rnn-effectiveness/

Andrej Karpathy

# Generating poetry with RNNs

```
PANDARUS:
Alas, I think he shall be come approached and the day
When little srain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:
They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:
Well, your wit is in the care of side and that.

Second Lord:
They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.

Clown:
Come, sir, I will make did behold your worship.

VIOLA:
I'll drink it.
```

```
VIOLA:
Why, Salisbury must find his flesh and thought
That which I am not aps, not a man and in fire,
To show the reining of the raven and the wars
To grace my hand reproach within, and not a fair are hand,
That Caesar and my goodly father's world;
When I was heaven of presence and our fleets,
We spare with hours, but cut thy council I am great,
Murdered and by thy master's ready there
My power to give thee but so much as hell:
Some service in the noble bondman here,
Would show him to her wine.

KING LEAR:
O, if you were a feeble sight, the courtesy of your law,
Your sight and several breath, will wear the gods
With his heads, and my hands are wonder'd at the deeds,
So drop upon your lordship's head, and your opinion
Shall be against your honour.
```

# Generating textbooks with RNNs

## open source textbook on algebraic geometry



Latex source

# Generating textbooks with RNNs

For $\bigoplus_{n=1,\ldots,m}$ where $\mathcal{L}_{m_\bullet} = 0$, hence we can find a closed subset $\mathcal{H}$ in $\mathcal{H}$ and any sets $\mathcal{F}$ on $X$, $U$ is a closed immersion of $S$, then $U \to T$ is a separated algebraic space.

*Proof.* Proof of (1). It also start we get

$$S = \mathrm{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \to V$. Consider the maps $M$ along the set of points $Sch_{fppf}$ and $U \to U$ is the fibre category of $S$ in $U$ in Section, **??** and the fact that any $U$ affine, see Morphisms, Lemma **??**. Hence we obtain a scheme $S$ and any open subset $W \subset U$ in $Sh(G)$ such that $\mathrm{Spec}(R') \to S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that $f_i$ is of finite presentation over $S$. We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x', s'' \in S'$ such that $\mathcal{O}_{X,x'} \to \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma **??** we can define a map of complexes $\mathrm{GL}_{S'}(x'/S'')$ and we win. $\square$

To prove study we see that $\mathcal{F}|_U$ is a covering of $\mathcal{X}'$, and $\mathcal{T}_i$ is an object of $\mathcal{F}_{X/S}$ for $i > 0$ and $\mathcal{F}_p$ exists and let $\mathcal{F}_i$ be a presheaf of $\mathcal{O}_X$-modules on $\mathcal{C}$ as a $\mathcal{F}$-module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widetilde{M}^\bullet = \mathcal{I}^\bullet \otimes_{\mathrm{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1}\mathcal{F})$$

is a unique morphism of algebraic stacks. Note that

$$\mathrm{Arrows} = (Sch/S)_{fppf}^{opp}, (Sch/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \longmapsto (U, \mathrm{Spec}(A))$$

is an open subset of $X$. Thus $U$ is affine. This is a continuous map of $X$ is the inverse, the groupoid scheme $S$.

*Proof.* See discussion of sheaves of sets. $\square$

The result for prove any open covering follows from the less of Example **??**. It may replace $S$ by $X_{spaces,\acute{e}tale}$ which gives an open subspace of $X$ and $T$ equal to $S_{Zar}$, see Descent, Lemma **??**. Namely, by Lemma **??** we see that $R$ is geometrically regular over $S$.

---

**Lemma 0.1.** *Assume (3) and (3) by the construction in the description.*

*Suppose $X = \lim |X|$ (by the formal open covering $X$ and a single map $\underline{\mathrm{Proj}}_X(\mathcal{A}) = \mathrm{Spec}(B)$ over $U$ compatible with the complex*

$$\mathrm{Set}(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X,\mathcal{O}_X}).$$

*When in this case of to show that $\mathcal{Q} \to \mathcal{C}_{Z/X}$ is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition **??** (without element is when the closed subschemes are catenary. If $T$ is surjective we may assume that $T$ is connected with residue fields of $S$. Moreover there exists a closed subspace $Z \subset X$ of $X$ where $U$ in $X'$ is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem*

(1) *$f$ is locally of finite type. Since $S = \mathrm{Spec}(R)$ and $Y = \mathrm{Spec}(R)$.*

*Proof.* This is form all sheaves of sheaves on $X$. But given a scheme $U$ and a surjective étale morphism $U \to X$. Let $U \cap U = \coprod_{i=1,\ldots,n} U_i$ be the scheme $X$ over $S$ at the schemes $X_i \to X$ and $U = \lim_i X_i$. $\square$

The following lemma surjective restrocomposes of this implies that $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{\mathcal{X},\ldots,0}$.

**Lemma 0.2.** *Let $X$ be a locally Noetherian scheme over $S$, $E = \mathcal{F}_{X/S}$. Set $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}'_n$. Since $\mathcal{I}^n \subset \mathcal{I}^n$ are nonzero over $i_0 \leq \mathfrak{p}$ is a subset of $\mathcal{J}_{n,0} \circ \overline{A}_2$ works.*

**Lemma 0.3.** *In Situation **??**. Hence we may assume $\mathfrak{q}' = 0$.*

*Proof.* We will use the property we see that $\mathfrak{p}$ is the mext functor (**??**). On the other hand, by Lemma **??** we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where $K$ is an $F$-algebra where $\delta_{n+1}$ is a scheme over $S$. $\square$

Andrej Karpathy

# Generating textbooks with RNNs

**Proof.** Omitted. □

**Lemma 0.1.** *Let $\mathcal{C}$ be a set of the construction.*
*Let $\mathcal{C}$ be a gerber covering. Let $\mathcal{F}$ be a quasi-coherent sheaves of $\mathcal{O}$-modules. We have to show that*

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

.

**Proof.** This is an algebraic space with the composition of sheaves $\mathcal{F}$ on $X_{\acute{e}tale}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where $\mathcal{G}$ defines an isomorphism $\mathcal{F} \to \mathcal{F}$ of $\mathcal{O}$-modules. □

**Lemma 0.2.** *This is an integer $\mathcal{Z}$ is injective.*

**Proof.** See Spaces, Lemma ??. □

**Lemma 0.3.** *Let $S$ be a scheme. Let $X$ be a scheme and $X$ is an affine open covering. Let $\mathcal{U} \subset \mathcal{X}$ be a canonical and locally of finite type. Let $X$ be a scheme. Let $X$ be a scheme which is equal to the formal complex.*

*The following to the construction of the lemma follows.*

*Let $X$ be a scheme. Let $X$ be a scheme covering. Let*

$$b : X \to Y' \to Y \to Y \to Y' \times_X Y \to X.$$

*be a morphism of algebraic spaces over $S$ and $Y$.*

**Proof.** Let $X$ be a nonzero scheme of $X$. Let $X$ be an algebraic space. Let $\mathcal{F}$ be a quasi-coherent sheaf of $\mathcal{O}_X$-modules. The following are equivalent

(1) $\mathcal{F}$ is an algebraic space over $S$.
(2) If $X$ is an affine open covering.

Consider a common structure on $X$ and $X$ the functor $\mathcal{O}_X(U)$ which is locally of finite type. □

---

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram



is a limit. Then $\mathcal{G}$ is a finite type and assume $S$ is a flat and $\mathcal{F}$ and $\mathcal{G}$ is a finite type $f_*$. This is of finite type diagrams, and

- the composition of $\mathcal{G}$ is a regular sequence,
- $\mathcal{O}_{X'}$ is a sheaf of rings. □

**Proof.** We have see that $X = \mathrm{Spec}(R)$ and $\mathcal{F}$ is a finite type representable by algebraic space. The property $\mathcal{F}$ is a finite morphism of algebraic stacks. Then the cohomology of $X$ is an open neighbourhood of $U$. □

**Proof.** This is clear that $\mathcal{G}$ is a finite presentation, see Lemmas ??.
A *reduced above* we conclude that $U$ is an open covering of $\mathcal{C}$. The functor $\mathcal{F}$ is a "field

$$\mathcal{O}_{X,x} \longrightarrow \mathcal{F}_{\overline{x}} \quad -1(\mathcal{O}_{X_{\acute{e}tale}}) \longrightarrow \mathcal{O}_{X_i}^{-1}\mathcal{O}_{X_\lambda}(\mathcal{O}_{X_\eta}^{\overline{v}})$$

is an isomorphism of covering of $\mathcal{O}_{X_i}$. If $\mathcal{F}$ is the unique element of $\mathcal{F}$ such that $X$ is an isomorphism.
The property $\mathcal{F}$ is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme $\mathcal{O}_X$-algebra with $\mathcal{F}$ are opens of finite type over $S$.
If $\mathcal{F}$ is a scheme theoretic image points. □

If $\mathcal{F}$ is a finite direct sum $\mathcal{O}_{X_\lambda}$ is a closed immersion, see Lemma ??. This is a sequence of $\mathcal{F}$ is a similar morphism.

Andrej Karpathy

# Generating code with RNNs

```c
static void do_command(struct seq_file *m, void *v)
{
  int column = 32 << (cmd[2] & 0x80);
  if (state)
    cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
  else
    seq = 1;
  for (i = 0; i < 16; i++) {
    if (k & (1 << 1))
      pipe = (in_use & UMXTHREAD_UNCCA) +
        ((count & 0x00000000ffffff8) & 0x000000f) << 8;
    if (count == 0)
      sub(pid, ppc_md.kexec_handle, 0x20000000);
    pipe_set_bytes(i, 0);
  }
  /* Free our user pages pointer to place camera if all dash */
  subsystem_info = &of_changes[PAGE_SIZE];
  rek_controls(offset, idx, &soffset);
  /* Now we want to deliberately put it to device */
  control_check_polarity(&context, val, 0);
  for (i = 0; i < COUNTER; i++)
    seq_puts(s, "policy ");
}
```

Generated
C code

# Image Captioning



CVPR 2015:
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei
Show and Tell: A Neural Image Caption Generator, Vinyals et al.
Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.
Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

# Image Captioning

# Image Captioning


test image

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096
FC-1000
softmax

test image

Andrej Karpathy

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

X

test image

Andrej Karpathy

# Image Captioning



test image



| image |
|---|
| conv-64 |
| conv-64 |
| maxpool |
| conv-128 |
| conv-128 |
| maxpool |
| conv-256 |
| conv-256 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| FC-4096 |
| FC-4096 |

x0
<START>

<START>

# Image Captioning



test image

| |
|---|
| image |
| conv-64 |
| conv-64 |
| maxpool |
| conv-128 |
| conv-128 |
| maxpool |
| conv-256 |
| conv-256 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| FC-4096 |
| FC-4096 |

im

**Wih**

y0

h0

x0
<START>

<START>

**before:**

$h = \tanh(W_{xh} * x + W_{hh} * h)$

**now:**

$h = \tanh(W_{xh} * x + W_{hh} * h + W_{ih} * im)$

Andrej Karpathy

# Image Captioning



test image

sample!

image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096

y0

h0

x0
<START>

straw

<START>

# Image Captioning



test image

<START>

Andrej Karpathy

# Image Captioning



test image

sample!

<START>

Andrej Karpathy

# Image Captioning



test image

<START>

Andrej Karpathy

# Image Captioning



test image

Caption generated:
"straw hat"

sample
<END> token
=> finish.

image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096

y0    y1    y2

h0 → h1 → h2

x0
<START>    straw    hat

<START>

# Image Captioning



"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."

"a young boy is holding a baseball bat."

"a cat is sitting on a couch with a remote control."

"a woman holding a teddy bear in front of a mirror."

"a horse is standing in the middle of a road."

# Video Captioning

Generate descriptions for events depicted in video clips



A monkey pulls a dog's tail and is chased by the dog.

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning



| English Sentence | → | RNN encoder | → ◯◯◯ → | RNN decoder | → | French Sentence | [Sutskever et al. NIPS'14] |

| 📷 | → | Encode | → ◯◯◯ → | RNN decoder | → | Sentence | [Donahue et al. CVPR'15] [Vinyals et al. CVPR'15] |

| ▶ | → | Encode | → ◯◯◯ → | RNN decoder | → | Sentence | [Venugopalan et. al. NAACL'15] (this work) |

Key Insight:

Generate feature representation of the video and "decode" it to a sentence

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning



Input Video → Sample frames @1/10

Forward propagate
Output: "fc7" features
(activations before classification layer)

CNN

$\frac{1}{n}\Sigma$

fc7: 4096 dimension "feature vector"

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning



Input Video     Convolutional Net     Recurrent Net     Output

$$\frac{1}{n}\sum$$

Mean across all frames

Output: A boy is playing golf <EOS>

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning

**Annotated video data is scarce.**

Key Insight:
Use supervised pre-training on data-rich
auxiliary tasks and transfer.

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning

## CNN pre-training



fc7: 4096 dimension "feature vector"

CNN

- Caffe Reference Net - variation of Alexnet [Krizhevsky et al. NIPS'12]
- 1.2M+ images from ImageNet ILSVRC-12 [Russakovsky et al.]
- Initialize weights of our network.

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning

**Image-Caption pre-training**



Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning

## Fine-tuning



1. Video dataset
2. Mean pooled feature
3. Lower learning rate

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning





- A man appears to be plowing a rice field with a plow being pulled by two oxen.
- A man is plowing a mud field.
- Domesticated livestock are helping a man plow.
- A man leads a team of oxen down a muddy path.
- A man is plowing with some oxen.
- A man is tilling his land with an ox pulled plow.
- Bulls are pulling an object.
- Two oxen are plowing a field.
- The farmer is tilling the soil.
- A man in ploughing the field.

- A man is walking on a rope.
- A man is walking across a rope.
- A man is balancing on a rope.
- A man is balancing on a rope at the beach.
- A man walks on a tightrope at the beach.
- A man is balancing on a volleyball net.
- A man is walking on a rope held by poles
- A man balanced on a wire.
- The man is balancing on the wire.
- A man is walking on a rope.
- A man is standing in the sea shore.

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning

MT metrics (BLEU, METEOR) to compare the system generated sentences against (all) ground truth references.

| Model | BLEU | METEOR |
|---|---|---|
| **Best Prior Work**<br>[Thomason et al. COLING'14] | 13.68 | 23.90 |
| **Only Images** | 12.66 | 20.96 |
| **Only Video** | 31.19 | 26.87 |
| **Images+Video** | **33.29** | **29.07** |

Pre-training only, no fine-tuning

No pre-training

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning



FGM: A person is dancing with the person on the stage.

YT: A group of men are riding the forest.

I+V: **A group of people are dancing.**

GT: Many men and women are dancing in the street.



FGM: A person is cutting a potato in the kitchen.

YT: A man is slicing a tomato.

I+V: **A man is slicing a carrot.**

GT: A man is slicing carrots.



FGM: A person is walking with a person in the forest.

YT: A monkey is walking.

I+V: **A bear is eating a tree.**

GT: Two bear cubs are digging into dirt and plant matter at the base of a tree.



FGM: A person is riding a horse on the stage.

YT: A group of playing are playing in the ball.

I+V: **A basketball player is playing**.

GT: Dwayne wade does a fancy layup in an allstar game.

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning



| | | |
|---|---|---|
| English Sentence → | RNN encoder → ○○○ → RNN decoder → | French Sentence | [Sutskever et al. NIPS'14] |

[Donahue et al. CVPR'15]
[Vinyals et al. CVPR'15]

[Venugopalan et. al. NAACL'15]

[Venugopalan et. al. ICCV'15] (this work)

3

Venugopalan et al., "Sequence to Sequence - Video to Text", ICCV 2015

# Video Captioning



S2VT Overview

Now decode it to a sentence!

Encoding stage

Decoding stage

A    man    is    talking    ...

Venugopalan et al., "Sequence to Sequence - Video to Text", ICCV 2015

# Sequence-to-sequence: the bottleneck problem

Encoding of the
source sentence.

Target sentence (output)

he    hit    me    with    a    pie    <END>

Encoder RNN

Decoder RNN

il    a    m'    entarté    <START>    he    hit    me    with    a    pie

Source sentence (input)

Problems with this architecture?

Abigail See

# Sequence-to-sequence: the bottleneck problem

Encoding of the source sentence. This needs to capture *all information* about the source sentence. Information bottleneck!

Target sentence (output)

he    hit    me    with    a    pie    <END>

Encoder RNN

Decoder RNN

il    a    m'    entarté

<START>    he    hit    me    with    a    pie

Source sentence (input)

Abigail See

# Attention

- **Attention** provides a solution to the bottleneck problem.

- Core idea: on each step of the decoder, use *direct connection to the encoder* to *focus on a particular part* of the source sequence



- First we will show via diagram (no equations), then we will show with equations

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

*il*    *a*    *m'*    *entarté*    *<START>*

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

il      a      m'      entarté      <START>

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

il     a     m'     entarté          <START>

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

il   a   m'   entarté   <START>

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



On this decoder timestep, we're mostly focusing on the first encoder hidden state ("he")

Take softmax to turn the scores into a probability distribution

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

il    a    m'    entarté    <START>

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



Use the attention distribution to take a **weighted sum** of the encoder hidden states.

The attention output mostly contains information from the hidden states that received high attention.

Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

*il     a     m'     entarté          <START>*

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

*he*

$y_1$

Concatenate attention output with decoder hidden state, then use to compute $y_1$ as before

*il*    *a*    *m'*    *entarté*      *<START>*

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



Attention output

hit

$y_2$

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

il    a    m'    entarté        <START>    he

Source sentence (input)

Sometimes we take the **attention output** from the previous step, and also feed it into the decoder (along with the usual decoder input).

Abigail See

# Attention: in equations

- We have encoder hidden states $h_1, \ldots, h_N \in \mathbb{R}^h$

- On timestep $t$, we have decoder hidden state $s_t \in \mathbb{R}^h$

- We get the attention scores $e^t$ for this step:

$$e^t = [s_t^T h_1, \ldots, s_t^T h_N] \in \mathbb{R}^N$$

- We take softmax to get the attention distribution $\alpha^t$ for this step (this is a probability distribution and sums to 1)

$$\alpha^t = \text{softmax}(e^t) \in \mathbb{R}^N$$

- We use $\alpha^t$ to take a weighted sum of the encoder hidden states to get the attention output $a_t$

$$a_t = \sum_{i=1}^{N} \alpha_i^t h_i \in \mathbb{R}^h$$

- Finally we concatenate the attention output $a_t$ with the decoder hidden state $s_t$ and proceed as in the non-attention seq2seq model

$$[a_t; s_t] \in \mathbb{R}^{2h}$$

# Visual Description



Berkeley LRCN [Donahue et al. CVPR'15]:
A brown bear standing on top of a lush green field.

MSR CaptionBot [http://captionbot.ai/]:
A large brown bear walking through a forest.

**MSCOCO**
**80 classes**

Venugopalan et al., "Captioning Images With Diverse Objects", CVPR 2017

# Object Recognition

Can identify hundreds of categories of objects.



IMAGENET  14M images, 22K classes [Deng et al. CVPR'09]

mammal → placental → carnivore → canine → dog → working dog → husky

vehicle → craft → watercraft → sailing vessel → sailboat → trimaran

Venugopalan et al., "Captioning Images With Diverse Objects", CVPR 2017

# Novel Object Captioner (NOC)

We present Novel Object Captioner which can compose descriptions of 100s of objects in context.



NOC (ours): Describe novel objects without paired image-caption data.

IMAGENET + MSCOCO +

An **okapi** standing in the middle of a field.

Visual Classifiers. IMAGENET

**okapi**

Existing captioners. IMAGENET init + train MSCOCO

A horse standing in the dirt.

Venugopalan et al., "Captioning Images With Diverse Objects", CVPR 2017

# Insights

1. Need to recognize and describe objects outside of image-caption datasets.

         okapi

Venugopalan et al., "Captioning Images With Diverse Objects", CVPR 2017

# Insights

2. Describe unseen objects that are similar
   to objects seen in image-caption datasets.



okapi ⟷ zebra

Venugopalan et al., "Captioning Images With Diverse Objects", CVPR 2017

# Insight 2: Capture semantic similarity of words

Venugopalan et al., "Captioning Images With Diverse Objects", CVPR 2017

# Insight 3: Jointly train on multiple sources



Venugopalan et al., "Captioning Images With Diverse Objects", CVPR 2017

# Qualitative Evaluation: ImageNet



**Instruments**

A man holding a **banjo** in a park.

A large **chime** hanging on a metal pole

**Vehicles**

A **snowplow** truck driving down a snowy road.

A group of people standing around a large white **warship**.

**Land Animals**

A **okapi** is in the grass with a **okapi**.

A small brown and white **jackal** is standing in a field.

**Household**

A large metal **candelabra** next to a wall.

A black and white photo of a **corkscrew** and a **corkscrew**.

Venugopalan et al., "Captioning Images With Diverse Objects", CVPR 2017

# Qualitative Evaluation: ImageNet



**Birds**

A small **pheasant** is standing in a field.

A **osprey** flying over a large grassy area.

**Outdoors**

A large **glacier** with a mountain in the background.

A group of people are sitting in a **baobab**.

**Water Animals**

A **humpback** is flying over a large body of water.

A man is standing on a beach holding a **snapper**.

**Misc**

A table with a **cauldron** in the dark.

A woman is posing for a picture with a **chiffon** dress.

Venugopalan et al., "Captioning Images With Diverse Objects", CVPR 2017

# Plan for this lecture

- Learning the relation between images and text
  - Recurrent neural networks
  - Applications: Captioning
  - Transformers
- Visual question answering
  - Graph convolutional networks
  - Incorporating knowledge and reasoning

# Transformers: Motivation

- We want **parallelization** but RNNs are inherently sequential

- Despite GRUs and LSTMs, RNNs still need attention mechanism  to deal with long range dependencies – **path length** between  states grows with sequence otherwise

- But if **attention** gives us access to any state… maybe we can just  use attention and don't need the RNN?

# Transformer Overview

Attention is all you need. 2017.  Aswani, Shazeer, Parmar, Uszkoreit,  Jones, Gomez, Kaiser, Polosukhin
https://arxiv.org/pdf/1706.03762.pdf

- Non-recurrent sequence-to-sequence encoder-decoder model

- Task: machine translation with parallel corpus

- Predict each translated word

- Final cost/error function is standard cross-entropy error on top of a softmax classifier

This and related figures from paper ⇑



Christopher Manning

# Dot-Product Attention (Extending our previous def.)

- Inputs: a query q and a set of key-value (k-v) pairs to an output
- Query, keys, values, and output are all vectors

- Output is weighted sum of values, where
- Weight of each value is computed by an inner product of query and corresponding key
- Queries and keys have same dimensionality $d_k$ value have $d_v$

$$A(q, K, V) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$$

# Attention visualization in layer 5

- Words start to pay attention to other words in sensible ways

# BERT: Devlin, Chang, Lee, Toutanova (2018)

- Mask out $k$% of the input words, and then predict the masked words
  - They always use $k = 15$%

<p style="text-align:center">store      gallon</p>

<p style="text-align:center">↑       ↑</p>

<p style="text-align:center">the man went to the [MASK] to buy a [MASK] of milk</p>

- Too little masking: Too expensive to train
- Too much masking: Not enough context

# Additional task: Next sentence prediction

- To learn *relationships* between sentences, predict whether Sentence B is actual sentence that proceeds Sentence A, or a random sentence
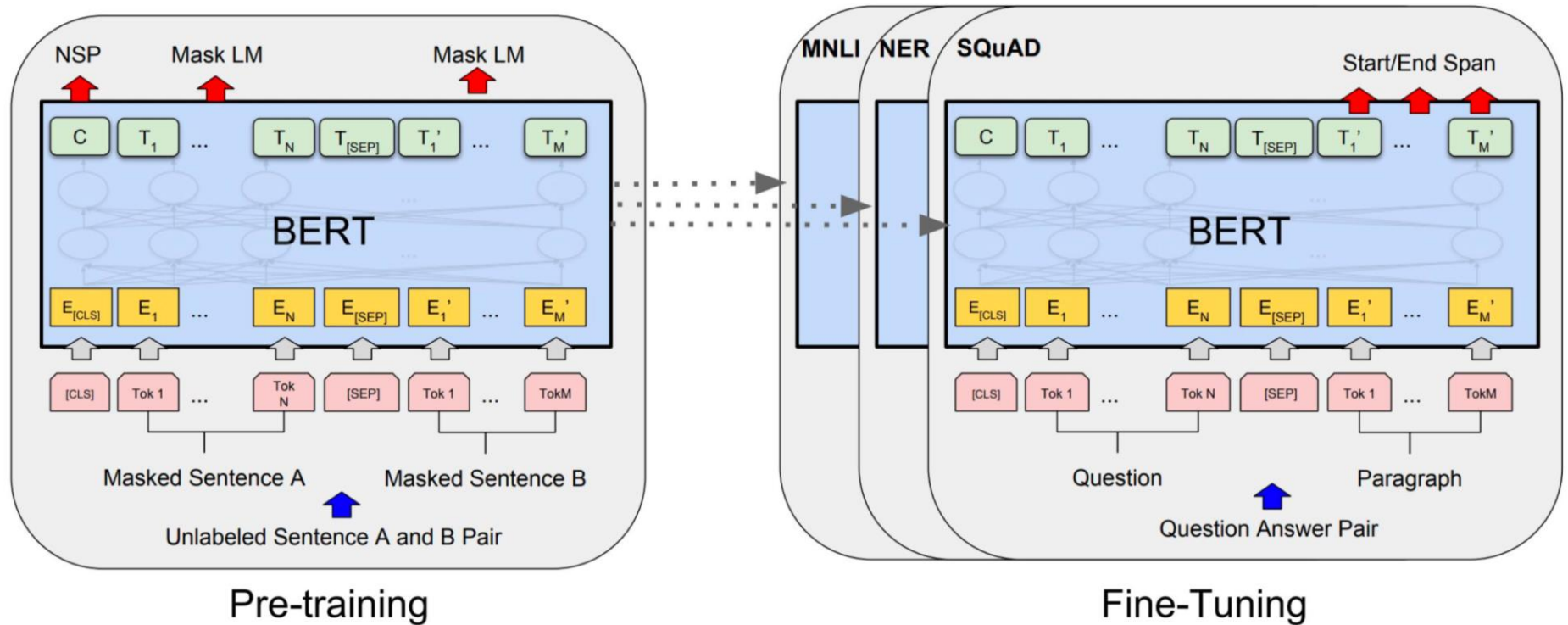
**Sentence A** = The man went to the store.
**Sentence B** = He bought a gallon of milk.
**Label** = IsNextSentence

**Sentence A** = The man went to the store.
**Sentence B** = Penguins are flightless.
**Label** = NotNextSentence

# BERT model fine tuning

- Simply learn a classifier built on the top layer for each task that you fine tune for

# SQuAD 2.0 leaderboard, 2019-02-07

| Rank | Model | EM | F1 |
|------|-------|-----|-----|
| | Human Performance<br>*Stanford University*<br>(Rajpurkar & Jia et al. '18) | 86.831 | 89.452 |
| 1<br>Jan 15, 2019 | BERT + MMFT + ADA (ensemble)<br>*Microsoft Research Asia* | **85.082** | **87.615** |
| 2<br>Jan 10, 2019 | BERT + Synthetic Self-Training (ensemble)<br>*Google AI Language*<br>https://github.com/google-research/bert | 84.292 | 86.967 |
| 3<br>Dec 13, 2018 | BERT finetune baseline (ensemble)<br>*Anonymous* | 83.536 | 86.096 |
| 4<br>Dec 16, 2018 | Lunet + Verifier + BERT (ensemble)<br>*Layer 6 AI NLP Team* | 83.469 | 86.043 |
| 4<br>Dec 21, 2018 | PAML+BERT (ensemble model)<br>*PINGAN GammaLab* | 83.457 | 86.122 |
| 5<br>Dec 15, 2018 | Lunet + Verifier + BERT (single model)<br>*Layer 6 AI NLP Team* | 82.995 | 86.035 |

Christopher Manning

# Cross-modal transformers



Figure 1: Our ViLBERT model consists of two parallel streams for visual (green) and linguistic (purple) processing that interact through novel co-attentional transformer layers. This structure allows for variable depths for each modality and enables sparse interaction through co-attention. Dashed boxes with multiplier subscripts denote repeated blocks of layers.

Lu et al., "ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks", NeurIPS 2019

# Cross-modal transformers



Figure 1: The LXMERT model for learning vision-and-language cross-modality representations. 'Self' and 'Cross' are abbreviations for self-attention sub-layers and cross-attention sub-layers, respectively. 'FF' denotes a feed-forward sub-layer.

Tan and Bansal, "LXMERT: Learning Cross-Modality Encoder Representationsfrom Transformers", EMNLP 2019

# Cross-modal transformers



Figure 1: Overview of the proposed UNITER model (best viewed in color), consisting of an Image Embedder, a Text Embedder and a multi-layer self-attention Transformer, learned through three pre-training tasks.

Chen et al., "UNITER: Learning UNiversal Image-TExt Representations", arxiv 2019

# Visual Commonsense Reasoning Leaderboard



Why is [person4 👤] pointing at [person1 👤]?
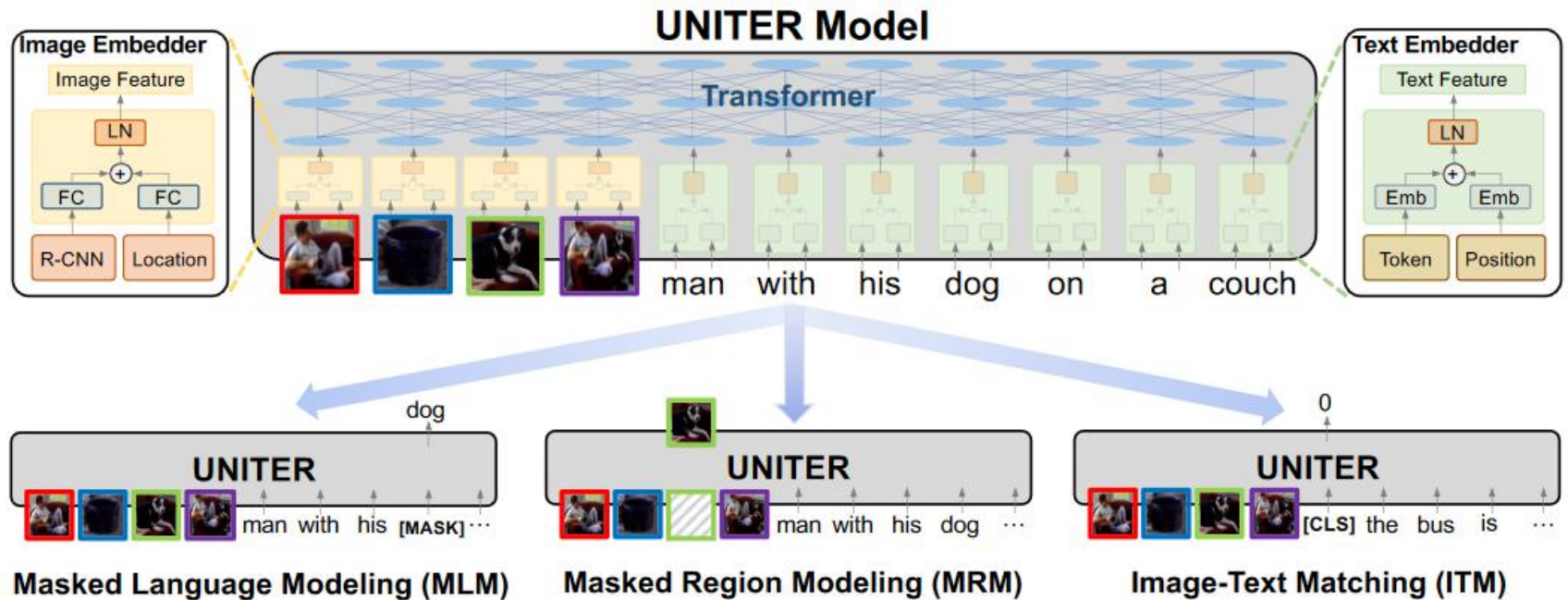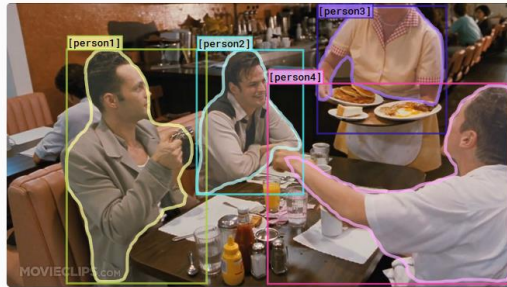
a) He is telling [person3 👤] that [person1 👤] ordered the pancakes.
b) He just told a joke.
c) He is feeling accusatory towards [person1 👤].
d) He is giving [person1 👤] directions.

Rationale: I think so because...

a) [person1 👤] has the pancakes in front of him.
b) [person4 👤] is taking everyone's order and asked for clarification.
c) [person3 👤] is looking at the pancakes both she and [person2 👤] are smiling slightly.
d) [person3 👤] is delivering food to the table, and she might not know whose order is whose.

| Rank | Model | Q->A | QA->R | Q->AR |
|---|---|---|---|---|
| | **Human Performance** *University of Washington* (Zellers et al. '18) | 91.0 | 93.0 | 85.0 |
| 📷 September 30, 2019 | **UNITER-large (ensemble)** *MS D365 AI* https://arxiv.org /abs/1909.11740 | **79.8** | **83.4** | **66.8** |
| 2 September 23, 2019 | UNITER-large (single model) *MS D365 AI* https://arxiv.org /abs/1909.11740 | 77.3 | 80.8 | 62.8 |
| 3 August 9,2019 | ViLBERT (ensemble of 10 models) *Georgia Tech & Facebook AI Research* https://arxiv.org /abs/1908.02265 | 76.4 | 78.0 | 59.8 |
| 4 September 23,2019 | VL-BERT (single model) *MSRA & USTC* https://arxiv.org /abs/1908.08530 | 75.8 | 78.4 | 59.7 |
| 5 August 9,2019 | ViLBERT (ensemble of 5 models) *Georgia Tech & Facebook AI Research* https://arxiv.org /abs/1908.02265 | 75.7 | 77.5 | 58.8 |

https://visualcommonsense.com/leaderboard/

# Plan for this lecture

- Learning the relation between images and text
  - Recurrent neural networks
  - Applications: Captioning
  - Transformers
- Visual question answering
  - Graph convolutional networks
  - Incorporating knowledge and reasoning

# Visual Question Answering (VQA)

Task: Given an image and a natural language open-ended question, generate a natural language answer.



What color are her eyes?
What is the mustache made of?

How many slices of pizza are there?
Is this a vegetarian pizza?

Is this person expecting company?
What is just under the tree?

Does it appear to be rainy?
Does this person have 20/20 vision?

Agrawal et al., "VQA: Visual Question Answering", ICCV 2015

# Visual Question Answering (VQA)



Image Embedding

Neural Network
Softmax
over top K answers

4096-dim

Convolution Layer + Non-Linearity  Pooling Layer  Convolution Layer + Non-Linearity  Pooling Layer  Fully-Connected

$h_1^{(2)}$  $h_2^{(2)}$  $h_3^{(2)}$  +1

$P(y = 0 \mid x)$
$P(y = 1 \mid x)$
$P(y = 2 \mid x)$

Input (Features II)  Softmax classifier

Question Embedding

"How many horses are in this image?"

1024-dim

LSTM

Agrawal et al., "VQA: Visual Question Answering", ICCV 2015

# Visual Question Answering (VQA)



Figure 2. Our proposed framework: given an image, a CNN is first applied to produce the attribute-based representation $V_{att}(I)$. The internal textual representation is made up of image captions generated based on the image-attributes. The hidden state of the caption-LSTM after it has generated the last word in each caption is used as its vector representation. These vectors are then aggregated as $V_{cap}(I)$ with average-pooling. The exter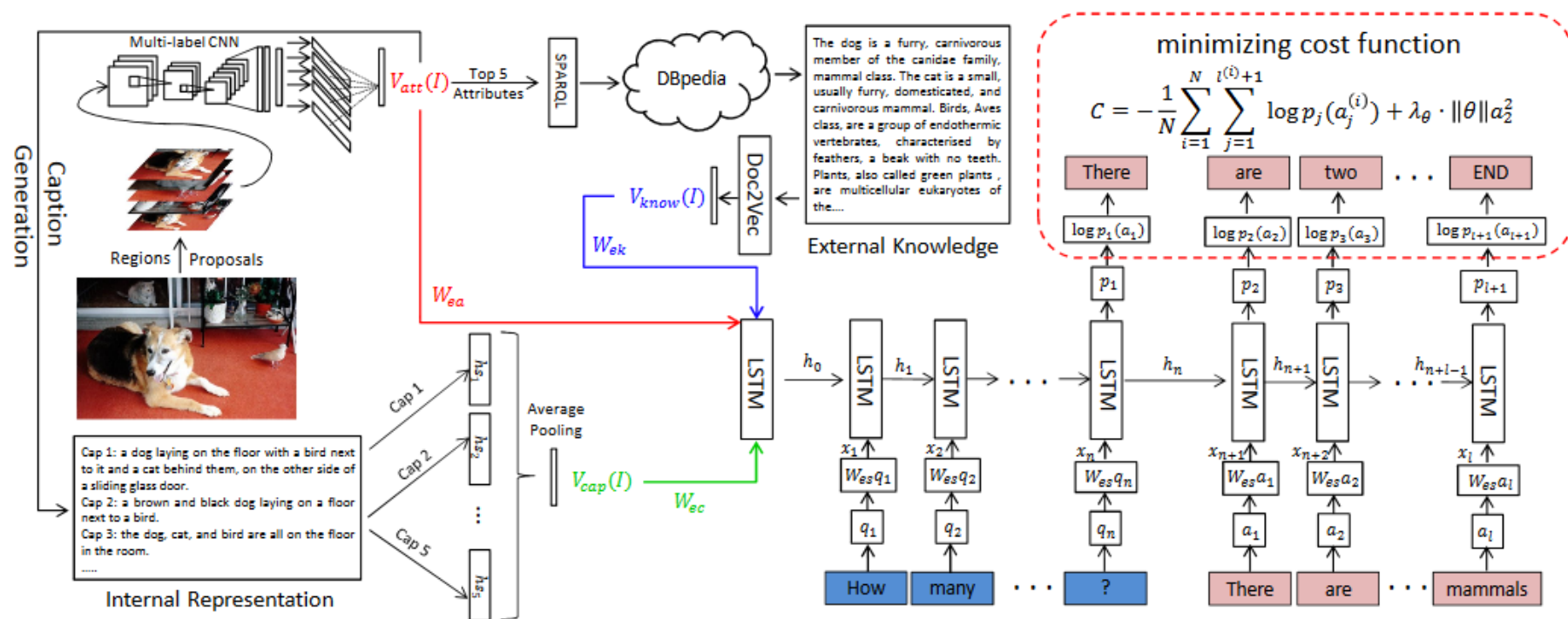nal knowledge is mined from the KB (in this case DBpedia) and the responses encoded by Doc2Vec, which produces a vector $V_{know}(I)$. The 3 vectors V are combined into a single representation of scene content, which is input to the VQA LSTM model which interprets the question and generates an answer.

Wu et al., "Ask Me Anything: Free-Form Visual Question Answering Based on Knowledge From External Sources", CVPR 2016

# Visual Question Answering (VQA)



Agrawal et al., "VQA: Visual Question Answering", ICCV 2015

# Reasoning for VQA



Is there a red shape above a circle?

yes

red

exists ↦ true

above

Andreas et al., "Neural Module Networks", CVPR 2016

# Reasoning for VQA



**Question**: *Are there more cubes than yellow things?*   **Answer**: *Yes*

Johnson et al., "Inferring and Executing Programs for Visual Reasoning", ICCV 2017

# Reasoning for VQA



Question: Which object in the image can be used to eat with?
Relation: UsedFor
Associated Fact: (Fork, UsedFor, Eat)
Answer Source: Image
Answer: Fork

Question: What do the animals in the image eat?
Relation: RelatedTo
Associated Fact: (Sheep, RelatedTo, Grass Eater)
Answer Source: Knowledge Base
Answer: Grass

Question: Which equipment in this image is used to hit baseball?
Relation: CapableOf
Associated Fact: (Baseball bat, CapableOf, Hit a baseball)
Answer Source: Image
Answer: Baseball bat

**Fig. 1.** The FVQA dataset expects methods to answer questions about images utilizing information from the image, as well as fact-based knowledge bases. Our method makes use of the image, and question text features, as well as high-level visual concepts extracted from the image in combination with a learned fact-ranking neural network. Our method is able to answer both visually grounded as well as fact based questions.

Narasimhan and Schwing, "Straight to the Facts: Learning Knowledge Base Retrieval for Factual Visual Question Answering", ECCV 2018
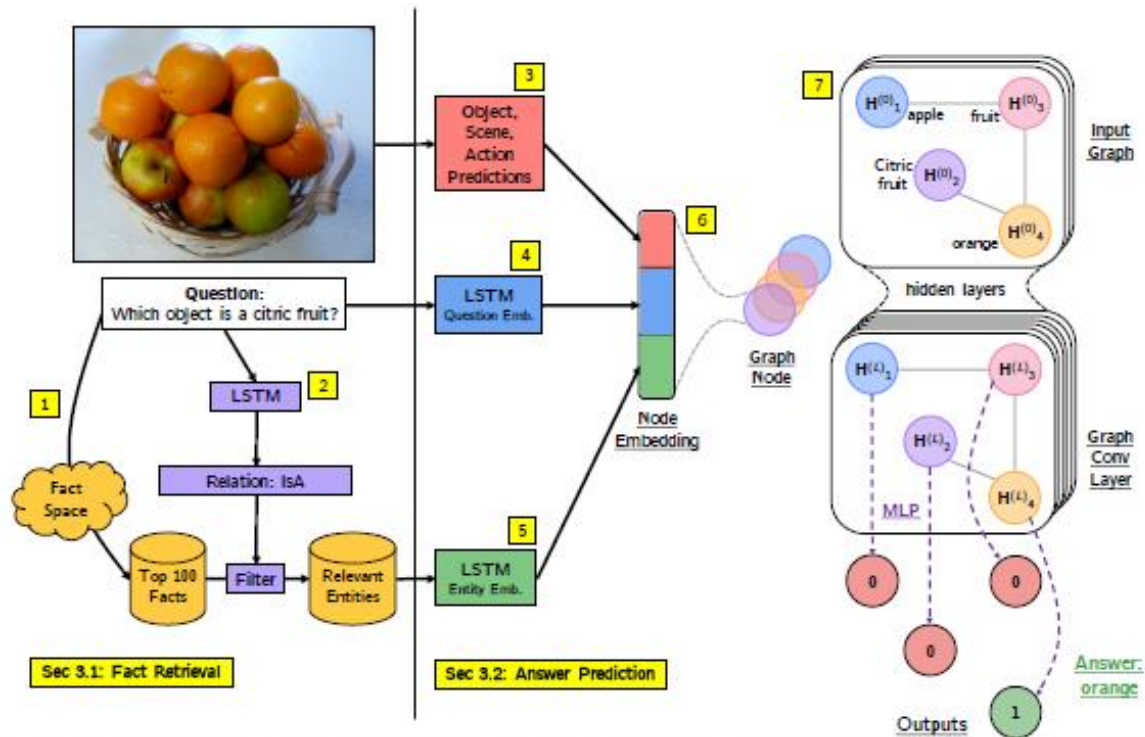
# Reasoning for VQA



Figure 2: Outline of the proposed approach: Given an image and a question, we use a similarity scoring technique (1) to obtain relevant facts from the fact space. An LSTM (2) predicts the relation from the question to further reduce the set of relevant facts and its entities. An entity embedding is obtained by concatenating the visual concepts embedding of the image (3), the LSTM embedding of the question (4), and the LSTM embedding of the entity (5). Each entity forms a single node in the graph and the relations constitute the edges (6). A GCN followed by an MLP performs joint assessment (7) to predict the answer. Our approach is trained end-to-end.

Narasimhan and Schwing, "Out of the Box: Reasoning with Graph Convolution Nets for Factual Visual Question Answering", NeurIPS 2018

# Graph convolutional networks

Recall: Single CNN layer with 3x3 filter:



(Animation by Vincent Dumoulin)



**Update for a single pixel:**

- Transform messages individually $\mathbf{W}_i \mathbf{h}_i$

- Add everything up $\sum_i \mathbf{W}_i \mathbf{h}_i$

Full update: $\mathbf{h}_4^{(l+1)} = \sigma \left( \mathbf{W}_0^{(l)} \mathbf{h}_0^{(l)} + \mathbf{W}_1^{(l)} \mathbf{h}_1^{(l)} + \cdots + \mathbf{W}_8^{(l)} \mathbf{h}_8^{(l)} \right)$

Kipf and Welling, "Semi-supervised learning with deep generative models", ICLR 2017 (slides by Thomas Kipf)

# Graph convolutional networks

**What if our data looks like this?**

or this:



**Real-world examples:**

- Social networks
- World-wide-web
- Protein-interaction networks

- Telecommunication networks
- Knowledge graphs
- …

# Graph convolutional networks

**Graph**: $G = (\mathcal{V}, \mathcal{E})$

**Adjacency matrix:  A**



$$
\begin{array}{c}
 & \begin{array}{ccccc} A & B & C & D & E \end{array} \\
\begin{array}{c} A \\ B \\ C \\ D \\ E \end{array} &
\left(\begin{array}{ccccc}
0 & 1 & 1 & 1 & 0 \\
1 & 0 & 0 & 1 & 1 \\
1 & 0 & 0 & 1 & 0 \\
1 & 1 & 1 & 0 & 1 \\
0 & 1 & 0 & 1 & 0
\end{array}\right)
\end{array}
$$

Kipf and Welling, "Semi-supervised learning with deep generative models", ICLR 2017 (slides by Thomas Kipf)

# Graph convolutional networks

Consider this undirected graph:

Calculate update for node in red:



**Update rule:**

$$\mathbf{h}_i^{(l+1)} = \sigma \left( \mathbf{h}_i^{(l)} \mathbf{W}_0^{(l)} + \sum_{j \in \mathcal{N}_i} \frac{1}{c_{ij}} \mathbf{h}_j^{(l)} \mathbf{W}_1^{(l)} \right)$$

$\mathcal{N}_i$ : neighbor indices

$c_{ij}$ : norm. constant (per edge)

Note: We could also choose simpler or more general functions over the neighborhood

Kipf and Welling, "Semi-supervised learning with deep generative models", ICLR 2017 (slides by Thomas Kipf)

# Graph convolutional networks

**Input**: Feature matrix $\mathbf{X} \in \mathbb{R}^{N \times E}$ , preprocessed adjacency matrix $\hat{\mathbf{A}}$



$$\mathbf{H}^{(l+1)} = \sigma\left(\hat{\mathbf{A}}\mathbf{H}^{(l)}\mathbf{W}^{(l)}\right)$$

Kipf and Welling, "Semi-supervised learning with deep generative models", ICLR 2017 (slides by Thomas Kipf)

# Graph convolutional networks

## Semi-supervised classification on graphs

**Setting:**

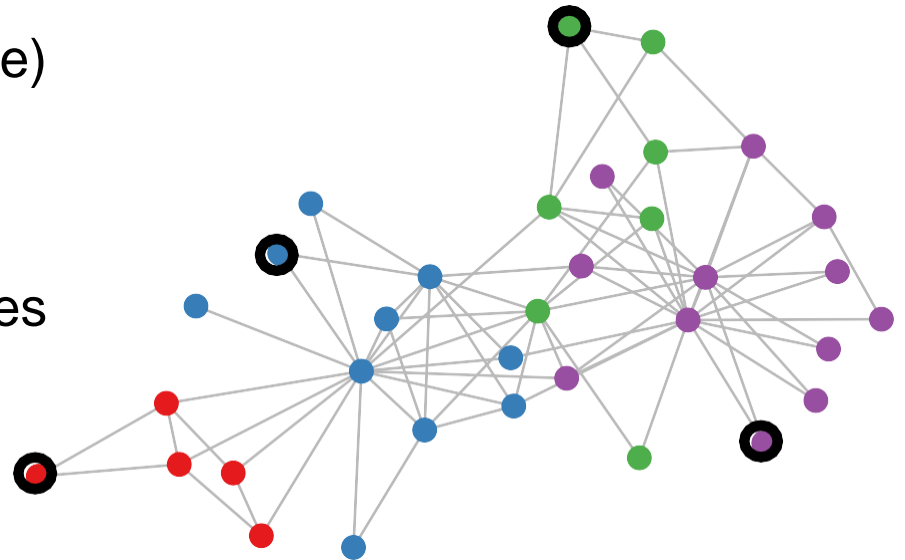Some nodes are labeled (black circle)
All other nodes are unlabeled

**Task:**

Predict node label of unlabeled nodes



Kipf and Welling, "Semi-supervised learning with deep generative models", ICLR 2017 (slides by Thomas Kipf)

# Decoding image advertisements

- What message does the ad convey (*action*), and what arguments does it provide for taking the suggested action (*reason*)?
- Multiple-choice task: Given *k* options for action-reason statements, pick one that matches the image



- I should drink evian because it helps you recover
- I should drink Evian because it will keep me like a baby
- I should buy Evian because it keeps us young

*Hussain, Zhang, Zhang, Ye, Thomas, Agha, Ong and Kovashka, CVPR 2017*

# Ads are effective



- Campaign ran for 25 years, and contained over 1,500 ads
- Absolut's share of the US vodka market went from 2.5% to about 25%

# Retrieve the best action-reason statement



Ye et al., TPAMI 2019

# Experimental results (image features only)

- We outperform prior art by a large margin, for both statement ranking and classification

| Method | Rank (Lower ↓ is better) | | Recall@3 (Higher ↑ is better) | |
|---|---|---|---|---|
| | PSA | Product | PSA | Product |
| 2-WAY NETS | 4.836 (± 0.090) | 4.170 (± 0.023) | 0.923 (± 0.016) | 1.212 (± 0.004) |
| VSE | 4.155 (± 0.091) | 3.202 (± 0.019) | 1.146 (± 0.017) | 1.447 (± 0.004) |
| VSE++ | 4.139 (± 0.094) | 3.110 (± 0.019) | 1.197 (± 0.017) | 1.510 (± 0.004) |
| HUSSAIN-RANKING | 3.854 (± 0.088) | 3.093 (± 0.019) | 1.258 (± 0.017) | 1.515 (± 0.004) |
| ADVISE (ours) | **3.013** (± 0.075) | **2.469** (± 0.015) | **1.509** (± 0.017) | **1.725** (± 0.004) |

- Our methods accurately capture the rhetoric, even in deliberately confusing ads



**VSE++ on Ads:** I should wear Revlon makeup because it will make me more attractive"

**ADVISE (ours):** "I should stop smoking because it doesn't make me pretty"

*Ye and Kovashka, ECCV 2018*

# Incorporating external knowledge



- Expand image representation using external knowledge (from DBPedia); represent regions, slogans, KB nuggets in a graph

- To prevent overfitting and break non-generalizable shortcuts, we randomly mask parts of training samples (e.g. slogan, words in KB nugget)

*Ye, Zhang and Kovashka*

# Incorporating external knowledge

- Training via metric learning: match image to human-annotated action-reason statements

- Image representation is a graph

- Slogan node updates:

$$\mathbf{t}_i^{(1)} = \underbrace{\alpha_{i,0}\mathbf{t}_i^{(0)}}_{\text{original meaning}} + \underbrace{\sum_{j=1}^{|\phi(t_i)|} \alpha_{i,j}\mathbf{k}_{i,j}}_{\text{descriptions from extra knowledge}}$$

- Global node update:

$$\mathbf{h} = \underbrace{\sum_{i=1}^{|V|} \beta_i \mathbf{v}_i}_{\text{messages from proposals}} + \underbrace{\sum_{i=|V|+1}^{|V|+|T|} \beta_i \mathbf{t}_i^{(1)}}_{\text{messages from slogans}}$$

- Edge weights α, β allow model to choose what knowledge to use

*Ye, Zhang and Kovashka*

# Incorporating external knowledge

- We stochastically mask aspects of training data, to prevent model from relying too much on word-matching or object-matching

- Three strategies; can also learn how to mask:
    - $M_t$ randomly drops a detected textual (T) slogan, with a probability of 0.5
    - $M_s$ randomly sets the KB query words (e.g. "WWF" or "Nike") in the human-annotated statements (S) to the out-of-vocabulary token, with probability 0.5
    - $M_k$ replaces the DBpedia queries in the retrieved knowledge contents with the out-of-vocabulary token

*Ye, Zhang and Kovashka*

# Incorporating external knowledge

- Outperform prior state of the art

| Methods | Accuracy (%) |
|---|---|
| VSE [31] | 62.0 |
| ADNET [6] | 65.0 |
| ADVISE [31] | 69.0 |
| CYBERAGENT [18] | 82.0 |
| RHETORIC [32] | 83.3 |
| OURS | **87.3** |

- Using external knowledge helps when data masked

| Method | P@1 | P@3 | P@5 | P@10 | R@1 | R@3 | R@5 | R@10 | Min Rank | Avg Rank | Med Rank |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Results on the Challenge-15 task | | | | | | | | | | | |
| V,T | **87.3** | 76.6 | 55.1 | 30.6 | **28.4** | 74.2 | 87.9 | 97.5 | 1.26 | 3.02 | 2.77 |
| V,T+K | **87.3** | 76.6 | 55.1 | 30.6 | **28.4** | 74.3 | 87.9 | 97.6 | 1.25 | 3.02 | 2.77 |
| V,T+K($M_t$,$M_s$,$M_k$) | **87.3** | **77.5** | **55.9** | **30.8** | **28.4** | **75.2** | **89.2** | **98.2** | **1.23** | **2.91** | **2.69** |
| Results on the Sampled-100 task | | | | | | | | | | | |
| V,T | 79.8 | 66.5 | 46.9 | 26.2 | 26.0 | 64.4 | 74.9 | 83.5 | 2.38 | 7.52 | 5.86 |
| V,T+K | 80.0 | 67.0 | 47.0 | 26.1 | 26.0 | 64.9 | 75.1 | 83.4 | 2.29 | 7.49 | 5.81 |
| V,T+K($M_t$,$M_s$,$M_k$) | **80.2** | **67.9** | **47.9** | **26.8** | **26.1** | **65.8** | **76.6** | **85.4** | **2.14** | **6.56** | **5.19** |
| Results on the Sampled-500 task | | | | | | | | | | | |
| V,T | **65.5** | 52.3 | 37.8 | 21.7 | **21.3** | 50.5 | 60.4 | 69.0 | 8.18 | 30.1 | 21.6 |
| V,T+K | 65.4 | 52.3 | 38.0 | 21.9 | **21.3** | 50.6 | 60.7 | 69.6 | 7.60 | 30.0 | 21.4 |
| V,T+K($M_t$,$M_s$,$M_k$) | 64.8 | **52.4** | **38.3** | **22.1** | 21.1 | **50.7** | **61.1** | **70.6** | **6.89** | **25.1** | **18.2** |

*Ye, Zhang and Kovashka*

# Incorporating external knowledge

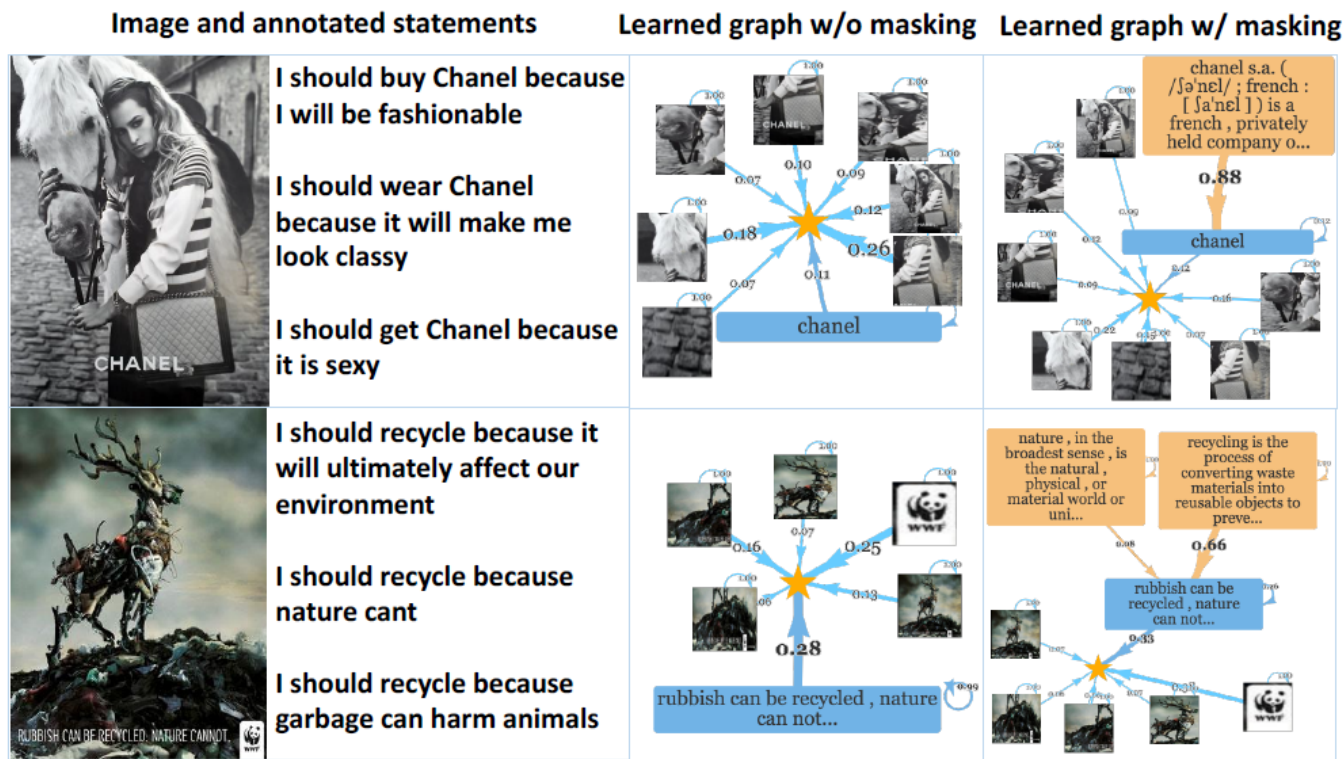Quantitatively: Without masking we retrieve relevant info with accuracy 25%, vs 54% with masking.



Fig. 4: **Examples of the learned graphs (best with zoom).** We show the ad image and annotated action-reason statements on the left, the graph learned without masking in the middle, and that learned with masking (our approach) on the right. We show slogans in blue, DBpedia comments in orange, and the global node as a star. **Arrow thickness is correlated with learned weights $\alpha, \beta$.** For visualization we removed all edges with small weights (threshold=0.05). We see our method more effectively leverages external information.

*Ye, Zhang and Kovashka*

# What's next?