# Recurrent neural networks

Image source
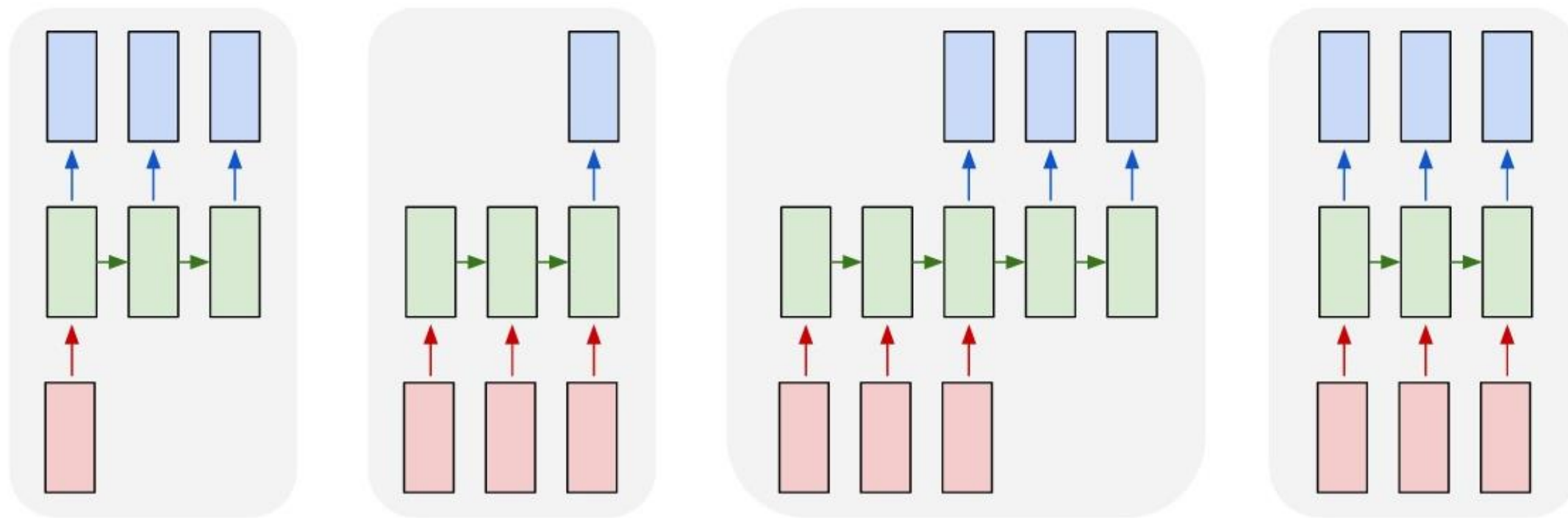
# Outline

- Examples of sequential prediction tasks
- Common recurrent units
    - Vanilla RNN unit (and how to train it)
    - Long Short-Term Memory (LSTM)
    - Gated Recurrent Unit (GRU)
- Recurrent network architectures
- Applications in (a bit) more detail
    - Sequence classification
    - Language modeling
    - Image captioning
    - Machine translation
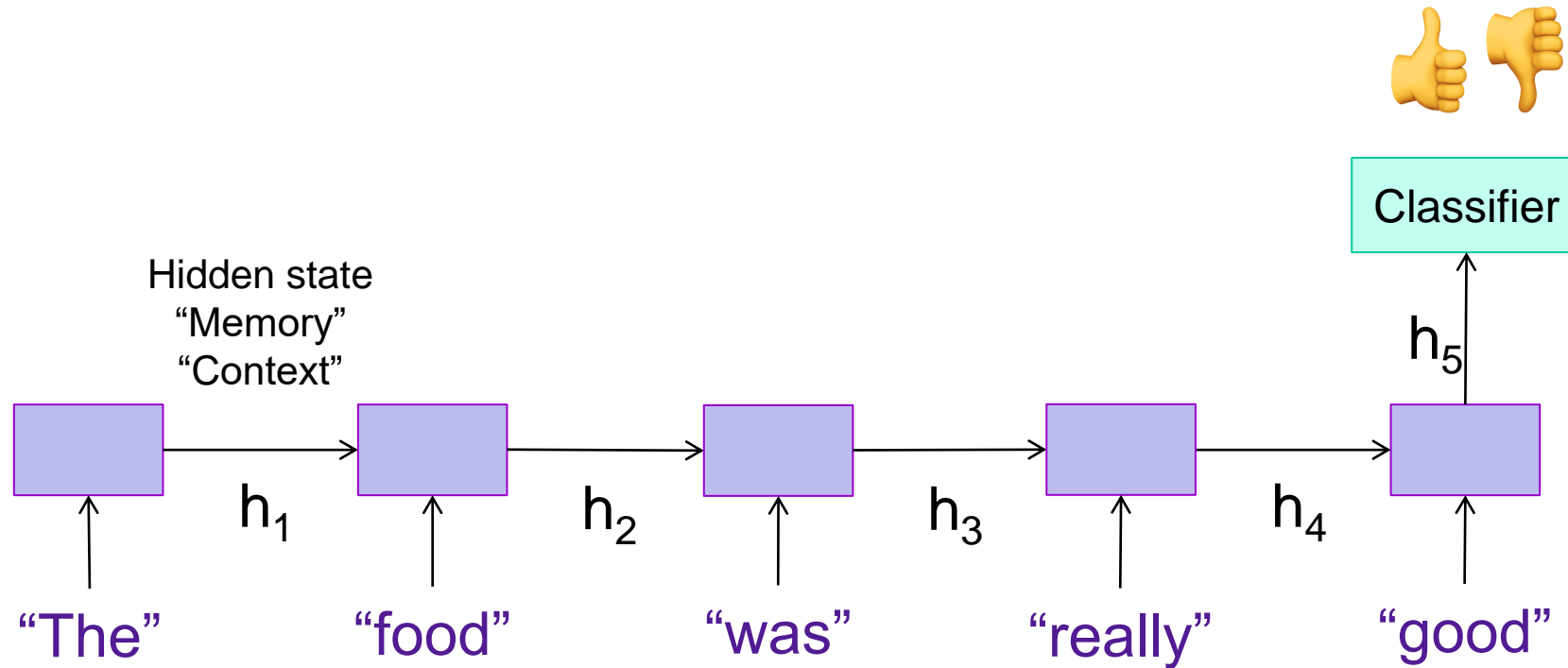
# Sequential prediction tasks

- So far, we focused mainly on prediction problems with fixed-size inputs and outputs

- But what if the input and/or output is a variable-length sequence?

# Example 1: Sentiment classification

- Goal: classify a text sequence (e.g., restaurant, movie or product review, Tweet) as having positive or negative sentiment

  - "The food was really good"
  - "The vacuum cleaner broke within two weeks"
  - "The movie had slow parts, but overall was worth watching"

- What makes this problem challenging?

- What feature representation or predictor structure can we use for this problem?

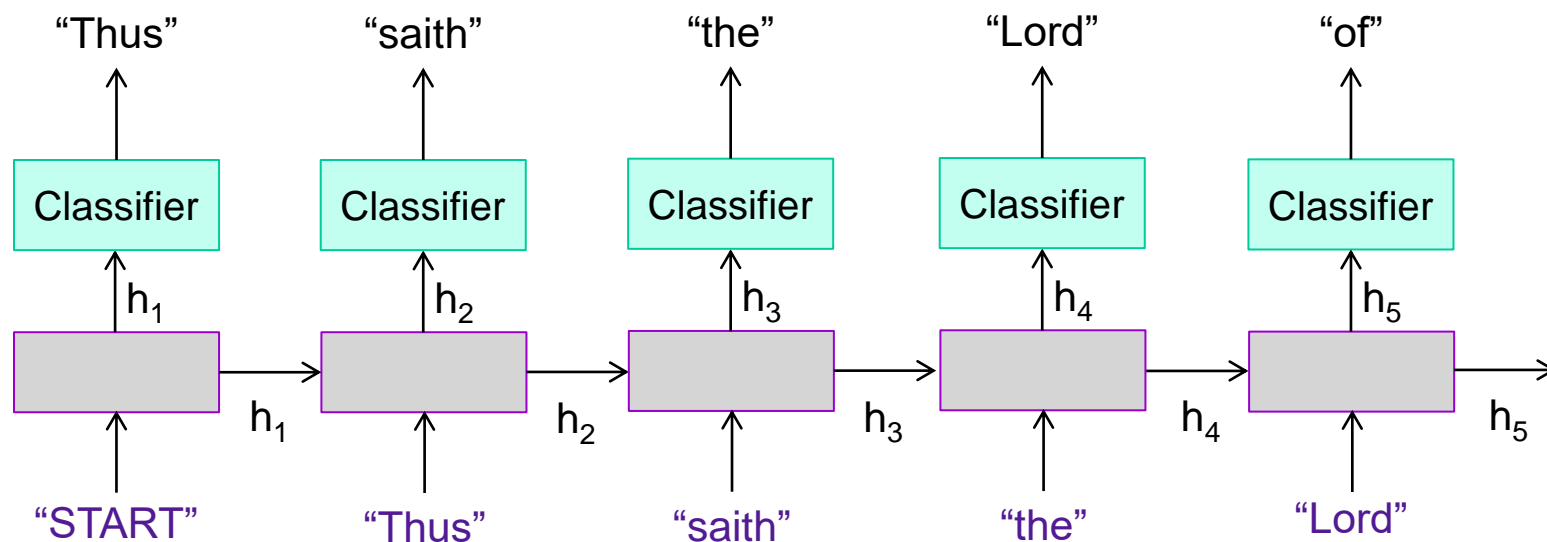# Example 1: Sentiment classification

- Recurrent model:

# Example 2: Text generation

- Sample from the distribution of a given text corpus (also known as language modeling)



**DeepDrumpf**
@DeepDrumpf

I'm a Neural Network trained on Trump's transcripts. Priming text in [ ]s. Donate (gofundme.com/deepdrumpf) to interact! Created by @hayesbh.

📅 Joined March 2016

**7** Following   **24.6K** Followers

| Tweets | Tweets & replies | Media | Likes |

**DeepDrumpf** @DeepDrumpf · May 31, 2017
[Despite the negative press #covfefe] look at what's going on. They shoot media. Usually that's a bad sign of things to come.
💬 6      🔁 38      ♡ 124

**DeepDrumpf** @DeepDrumpf · Apr 7, 2017
When I have to build a hotel, we're bombing the hell out of them. Lots of money. To those suffering, I say vote for Donald. #SyriaStrikes
💬 1      🔁 71      ♡ 173

**DeepDrumpf** @DeepDrumpf · Mar 20, 2017
Replying to @Thomas1774Paine
There will be no amnesty. It is going to pass because the people are going to be gone. I'm giving a mandate. #ComeyHearing @Thomas1774Paine

# Example 2: Text generation

- Sample from the distribution of a given text corpus (also known as language modeling)

- Can be done one character or one word at a time:

# Example 3: Image caption generation



*A cat sitting on a suitcase on the floor*



*A cat is sitting on a tree branch*



*A dog is running in the grass with a frisbee*



*A white teddy bear sitting in the grass*



*Two people walking on the beach with surfboards*
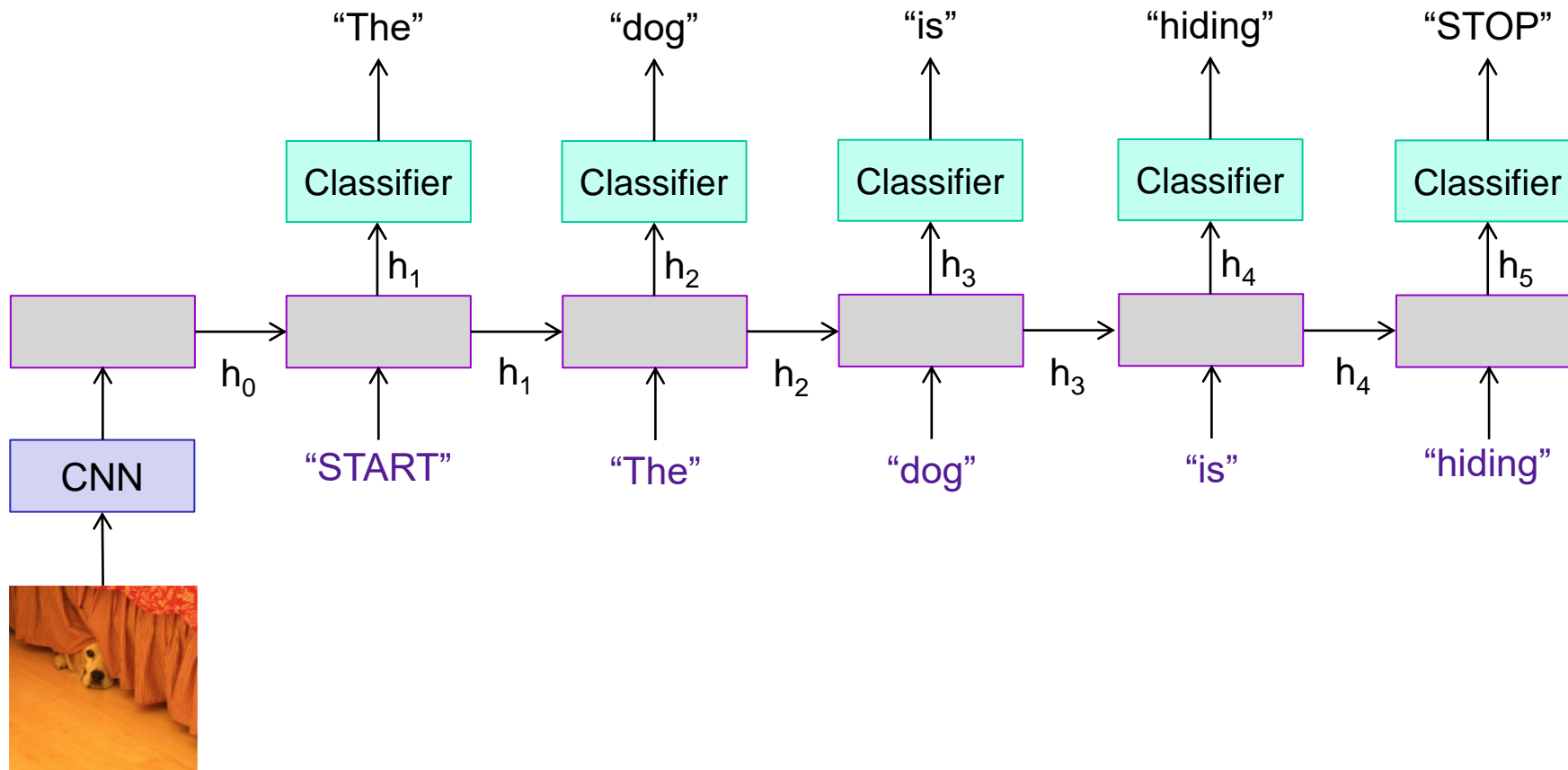


*A tennis player in action on the court*



*Two giraffes standing in a grassy field*



*A man riding a dirt bike on a dirt track*

# Example 3: Image caption generation

# Example 4: Machine translation



https://translate.google.com/

# Example 4: Machine translation

- Multiple input – multiple output (or sequence to sequence) scenario:

# Summary: Input-output scenarios

| | | |
|---|---|---|
| Single - Single | | Feed-forward Network |
| Multiple - Single | | Sequence Classification |
| Single - Multiple | | Sequence generation, captioning |
| Multiple - Multiple | | Sequence generation, captioning |
| Multiple - Multiple | | Translation |

# Outline

- Examples of sequential prediction tasks
- Common recurrent units
  - Vanilla RNN unit
  - Long Short-Term Memory (LSTM)
  - Gated Recurrent Unit (GRU)

# Recurrent unit

# Recurrent unit

Output at time $t$     $y_t$

Classifier

Hidden
representation
at time $t$     $h_t$

Hidden layer

Input at time $t$     $x_t$

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

new state    function of $W$    input at time $t$    old state

# Vanilla RNN cell



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

$$= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

J. Elman, Finding structure in time, Cognitive science 14(2), pp. 179–211, 1990

# Vanilla RNN cell

$h_t$



$W$

$h_{t-1}$    $x_t$

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

$= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$



$\sigma(a)$

$\tanh(a)$

$\tanh(a) = \dfrac{e^a - e^{-a}}{e^a + e^{-a}}$

$= 2\sigma(2a) - 1$

# Vanilla RNN cell



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

$$= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$\frac{d}{da} \tanh(a) = 1 - \tanh^2(a)$$

# Vanilla RNN cell



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

$$= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$= \tanh(W_x x_t + W_h h_{t-1})$$

# RNN forward pass



$$e_t = -\log(y_t(GT_t))$$

$$y_t = \text{softmax}(W_y h_t)$$

$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

------ shared weights

# RNN forward pass: Computation graph



$$e_t = -\log(y_t(GT_t))$$

$$y_t = \mathrm{softmax}(W_y h_t)$$

$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

# Training: Backpropagation through time (BPTT)

- The unfolded network (used during forward pass) is treated as one big feed-forward network that accepts the whole time series as input

- The weight updates are computed for each copy in the unfolded network, then summed (or averaged) and applied to the RNN weights

# Backpropagation through time



Forward through entire sequence to compute loss, then backward to compute gradient

# Backpropagation through time



Loss

Problem: Takes a lot of memory for long sequences!

# Training: Backpropagation through time (BPTT)

- The unfolded network (used during forward pass) is treated as one big feed-forward network that accepts the whole time series as input

- The weight updates are computed for each copy in the unfolded network, then summed (or averaged) and applied to the RNN weights

- In practice, *truncated* BPTT is used: run the RNN forward $k_1$ time steps, propagate backward for $k_2$ time steps

https://machinelearningmastery.com/gentle-introduction-backpropagation-time/
http://www.cs.utoronto.ca/~ilya/pubs/ilya_sutskever_phd_thesis.pdf

# Truncated backpropagation through time



Run forward and backward through chunks of the sequence instead of whole sequence

# Truncated backpropagation through time



Loss

Carry hidden states forward in time further, but only backpropagate for some smaller number of steps

# Truncated backpropagation through time

# RNN backward pass



$$h_t = \tanh(W_x x_t + W_h h_{t-1})$$

$$\frac{\partial e}{\partial W_h} = \frac{\partial e}{\partial h_t} \odot \big(1 - \tanh^2(W_x x_t + W_h h_{t-1})\big) h_{t-1}^T$$

$$\frac{\partial e}{\partial W_x} = \frac{\partial e}{\partial h_t} \odot \big(1 - \tanh^2(W_x x_t + W_h h_{t-1})\big) x_t^T$$

$$\frac{\partial e}{\partial h_{t-1}} = W_h^T \big(1 - \tanh^2(W_x x_t + W_h h_{t-1})\big) \odot \frac{\partial e}{\partial h_t}$$

# Vanishing and exploding gradients



$$\frac{\partial e}{\partial h_{t-1}} = W_h^T \big(1 - \tanh^2(W_x x_t + W_h h_{t-1})\big) \odot \frac{\partial e}{\partial h_t}$$

Computing gradient for $h_0$ involves many multiplications by $W_h^T$ (and rescalings between 0 and 1)

Gradients will *vanish* if largest singular value of $W_h$ is less than 1 and *explode* if it's greater than 1

# Outline

- Examples of sequential prediction tasks
- Common recurrent units
  - Vanilla RNN unit (and how to train it)
  - Long Short-Term Memory (LSTM)
  - Gated Recurrent Unit (GRU)

# Long short-term memory (LSTM)

- Add a *memory cell* that is not subject to matrix multiplication or squashing, thereby avoiding gradient decay



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

# The LSTM cell



$x_t$

$W_g$

$h_{t-1}$

$g_t = \tanh W_g \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$

Cell

$c_t$

$c_t = c_{t-1} + g_t$

$h_t$

$h_t = \tanh c_t$

# The LSTM cell

# The LSTM cell



Input Gate

$W_i$

$i_t$

$$i_t = \sigma\left(W_i\begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_i\right)$$

$x_t$

$W_g$

$x_t$

$h_{t-1}$

$$g_t = \tanh W_g\begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$x_t$

$h_{t-1}$

Cell

$c_t$

$h_t$

$$c_t = c_{t-1} + i_t \odot g_t$$

# The LSTM cell



Input Gate

$W_i$

$i_t$

$i_t = \sigma \left( W_i \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_i \right)$

Output Gate

$W_o$

$o_t$

$o_t = \sigma \left( W_o \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_o \right)$

$x_t$

$h_{t-1}$

$W_g$

$g_t = \tanh W_g \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$

Cell

$c_t$

$h_t$

$h_t = o_t \odot \tanh c_t$

$c_t = c_{t-1} + i_t \odot g_t$

# The LSTM cell



$x_t$   $h_{t-1}$                                $x_t$   $h_{t-1}$

Input Gate   $i_t$   $W_i$                Output Gate   $o_t$   $W_o$

$$i_t = \sigma\left(W_i \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_i\right)$$

$$o_t = \sigma\left(W_o \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_o\right)$$

$x_t$   $W_g$                Cell

$c_t$                $h_t$   $h_t = o_t \odot \tanh c_t$

$$g_t = \tanh W_g \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$h_{t-1}$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

Forget Gate

$W_f$   $f_t$   $$f_t = \sigma\left(W_f \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_f\right)$$

$x_t$   $h_{t-1}$

# LSTM forward pass summary



$$\begin{pmatrix} g_t \\ i_t \\ f_t \\ o_t \end{pmatrix} = \begin{pmatrix} \tanh \\ \sigma \\ \sigma \\ \sigma \end{pmatrix} \begin{pmatrix} W_g \\ W_i \\ W_f \\ W_o \end{pmatrix} \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$
$$h_t = o_t \odot \tanh c_t$$

# LSTM backward pass



Gradient flow from $c_t$ to $c_{t-1}$ only involves back-propagating through addition and elementwise multiplication, not matrix multiplication or tanh

For complete details: [Illustrated LSTM Forward and Backward Pass](#)

# Gated recurrent unit (GRU)



- Get rid of separate cell state

- Merge "forget" and "output" gates into "update" gate

K. Cho et al., Learning phrase representations using RNN encoder-decoder for statistical machine translation, ACL 2014

# Gated recurrent unit (GRU)



$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

# Gated recurrent unit (GRU)



$$r_t = \sigma\left(W_r \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_t\right)$$

$$h'_t = \tanh W \begin{pmatrix} x_t \\ r_t \odot h_{t-1} \end{pmatrix}$$

# Gated recurrent unit (GRU)



$$r_t = \sigma\left(W_r \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_t\right)$$

$$h'_t = \tanh W \begin{pmatrix} x_t \\ r_t \odot h_{t-1} \end{pmatrix}$$

$$z_t = \sigma\left(W_z \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_z\right)$$

# Gated recurrent unit (GRU)



$$r_t = \sigma \left( W_r \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_t \right)$$

$$h'_t = \tanh W \begin{pmatrix} x_t \\ r_t \odot h_{t-1} \end{pmatrix}$$

$$z_t = \sigma \left( W_z \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_z \right)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot h'_t$$

# Outline

- Examples of sequential prediction tasks
- Common recurrent units
  - Vanilla RNN unit (and how to train it)
  - Long Short-Term Memory (LSTM)
  - Gated Recurrent Unit (GRU)
- **Recurrent network architectures**

# Summary: Input-output scenarios

Multiple - Single

Sequence Classification

Single - Multiple

Sequence generation, captioning

Multiple - Multiple

Sequence generation, captioning

Multiple - Multiple

Translation

# Multi-layer RNNs

- We can of course design RNNs with multiple hidden layers



- Anything goes: skip connections across layers, across time, …

# Bi-directional RNNs

- RNNs can process the input sequence in forward and in the reverse direction (common in speech recognition)

# Google Neural Machine Translation (GNMT)



Y. Wu et al., Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation, arXiv 2016

# Outline

- Examples of sequential prediction tasks
- Common recurrent units
  - Vanilla RNN unit
  - Long Short-Term Memory (LSTM)
  - Gated Recurrent Unit (GRU)
- Recurrent network architectures
- **Applications in (a bit) more detail**
  - Sequence classification
  - Language modeling
  - Image captioning
  - Machine translation

# Sequence classification

# Sequence classification

# Sequence classification

# Language modeling: Character RNN



Output symbol $y_i$

Output layer (linear + softmax)

Hidden state $h_i$

One-hot encoding $x_i$

Input symbol

$$p(y_1, y_2, \ldots, y_n)$$

$$= \prod_{i=1}^{n} p(y_i | y_1, \ldots, y_{i-1})$$

$$\approx \prod_{i=1}^{n} P_W(y_i | h_i)$$

# Language modeling: Character RNN

```
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

300th iteration

```
"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

700th iteration

```
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

2000th iteration

```
"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.
```

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

# Searching for interpretable hidden units



quote detection cell

A. Karpathy, J. Johnson, and L. Fei-Fei, Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Searching for interpretable hidden units



line position tracking cell

A. Karpathy, J. Johnson, and L. Fei-Fei, Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Searching for interpretable hidden units



if statement cell

A. Karpathy, J. Johnson, and L. Fei-Fei, Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Searching for interpretable hidden units



quote/comment cell

A. Karpathy, J. Johnson, and L. Fei-Fei, Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Searching for interpretable hidden units

```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

code depth cell

A. Karpathy, J. Johnson, and L. Fei-Fei, Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# RNNs: Outline

- Examples of sequential prediction tasks
- Common recurrent units
  - Vanilla RNN unit
  - Long Short-Term Memory (LSTM)
  - Gated Recurrent Unit (GRU)
- Recurrent network architectures
  - Multilayer, bidirectional, skip connections
- Applications in (a bit) more detail
  - Sequence classification
  - Language modeling
  - Image captioning
  - Machine translation

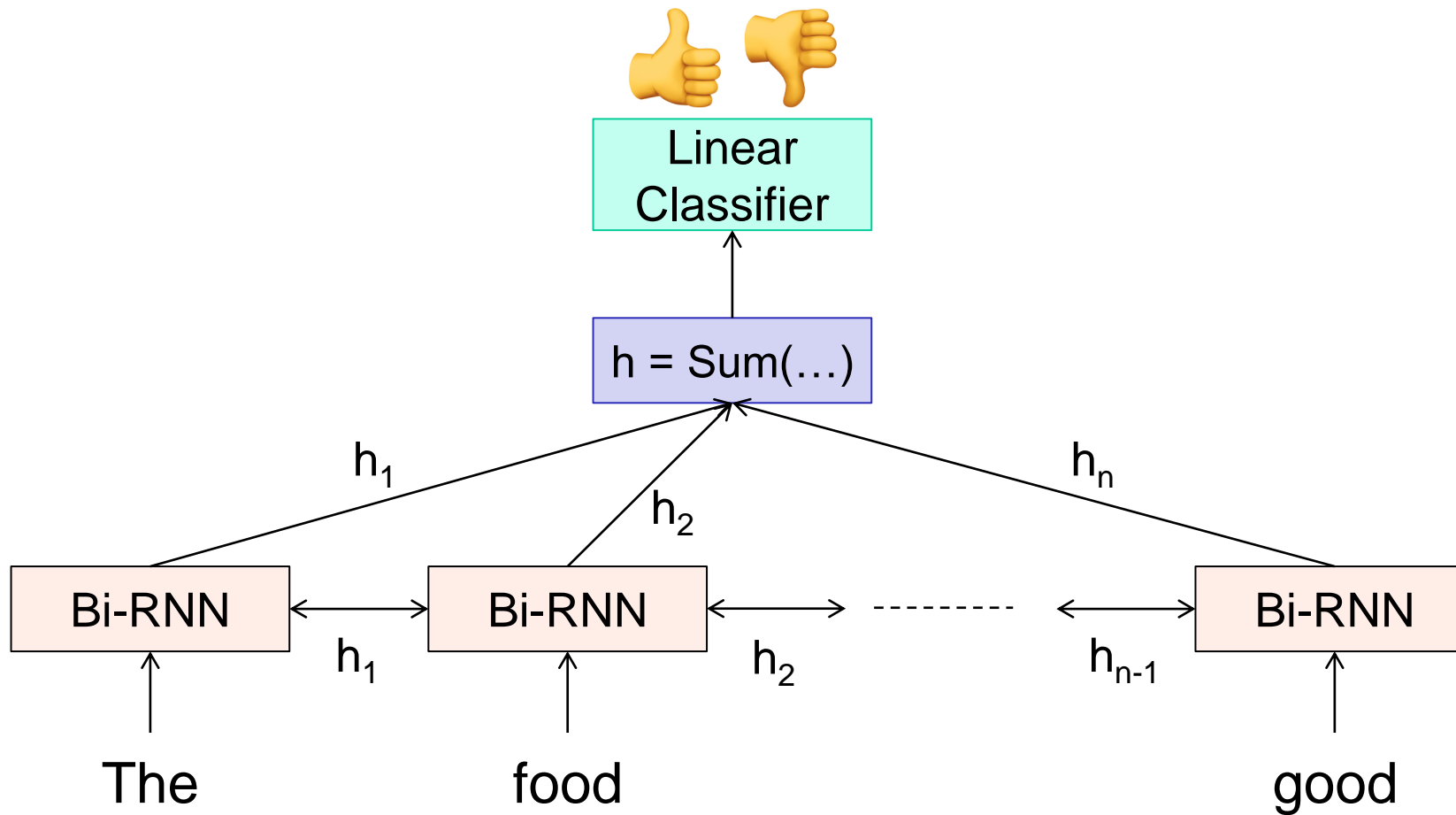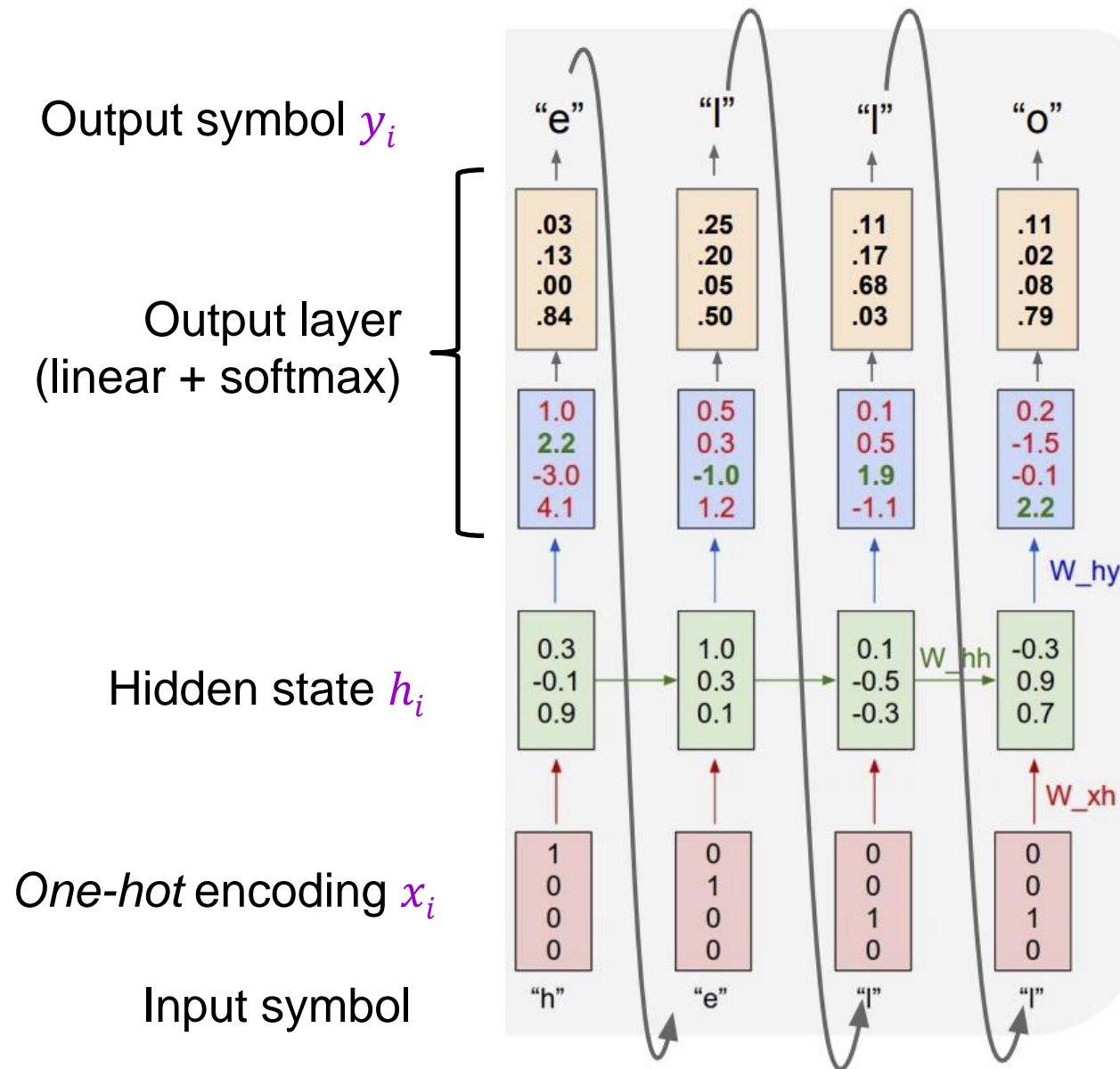# Image caption generation



**Training time**

- Maximize likelihood of reference captions

Log-likelihood of next reference word

Softmax probability over next word

Word embedding

Words of reference caption (one-hot encoding)

O. Vinyals, A. Toshev, S. Bengio, D. Erhan, Show and Tell: A Neural Image Caption Generator, CVPR 2015

# Image caption generation: Training time

- Minimize negative log-likelihood of the ground truth caption $Y^* = (Y_1^*, \dots, Y_N^*)$ given image $I$:

$$L(I, Y^*) = -\sum_{t=1}^{N} \log P_W(Y_i^* | Y_1^*, \dots, Y_{i-1}^*, I)$$

# Image caption generation: Test time

- Sample next word according to posterior distribution of classifier
  - Sentences quickly become incoherent
- Always choose the highest-likelihood word
  - Does this necessarily maximize the likelihood of the overall sentence?

# Image caption generation: Test time

- Beam search:
  - Maintain $k$ (*beam width*) top-scoring candidate sentences according to sum of per-word log-likelihoods (or some other score)
  - At each step, generate all their successors and keep the best $k$

# Image caption generation: Beam search

A person riding a motorcycle on a dirt road.

Two dogs play in the grass.

A skateboarder does a trick on a ramp.

A dog is jumping to catch a frisbee.

A group of young people playing a game of frisbee.

Two hockey players are fighting over the puck.

A little girl in a pink hat is blowing bubbles.

A refrigerator filled with lots of food and drinks.

A herd of elephants walking across a dry grass field.

A close up of a cat laying on a couch.

A red motorcycle parked on the side of the road.

A yellow school bus parked in a parking lot.

Describes without errors | Describes with minor errors | Somewhat related to the image | Unrelated to the image

# How to evaluate image captioning?

Reference sentences (written by human annotators):



- "A dog hides underneath a bed with its face peeking out of the bed skirt"
- "The small white dog is peeking out from under the bed"
- "A dog is peeking its head out from underneath a bed skirt"
- "A dog peeking out from under a bed"
- "A dog that is under a bed on the floor"

Generated sentence:

- "A dog is hiding"

# BLEU: Bilingual Evaluation Understudy

- **N-gram precision**: count the number of n-gram matches between candidate and reference translation, divide by total number of n-grams in candidate translation
  - Clip counts by the maximum number of times an n-gram occurs in any reference translation
  - Multiply by *brevity penalty* to penalize short translations

- Most commonly used measure for image captioning and machine translation despite multiple shortcomings

K. Papineni, S. Roukos, T. Ward, W.-J. Zhu, BLEU: a Method for Automatic Evaluation of Machine Translation, ACL 2002

ⓘ Overview    ⚑ Challenges ▾    ⊙ Download    ▬ Evaluate ▾    ☰ Leaderboard ▾

**Table-C5**   Table-C40   2015 Captioning Challenge          Last update: June 8, 2015. Visit CodaLab for the latest results.

| | CIDEr-D | Meteor | ROUGE-L | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 |
|---|---|---|---|---|---|---|---|
| m-RNN (Baidu/ UCLA)[16] | 0.886 | 0.238 | 0.524 | 0.72 | 0.553 | 0.41 | 0.302 |
| m-RNN[15] | 0.917 | 0.242 | 0.521 | 0.716 | 0.545 | 0.404 | 0.299 |
| MSR Captiva | | | | | | | 0.308 |
| Google[4] | | | | | | | 0.309 |
| Berkeley LRC | | | | | | | 0.277 |
| Nearest Neig | | | | | | | 0.28 |
| MSR[8] | | | | | | | 0.291 |
| Montreal/Toronto[10] | 0.85 | 0.243 | 0.513 | 0.689 | 0.515 | 0.372 | 0.268 |
| PicSOM[13] | 0.833 | 0.231 | 0.505 | 0.683 | 0.51 | 0.377 | 0.281 |
| Tsinghua Bigeye[14] | 0.673 | 0.207 | 0.49 | 0.671 | 0.494 | 0.35 | 0.241 |
| MLBL[7] | 0.74 | 0.219 | 0.499 | 0.666 | 0.498 | 0.362 | 0.26 |
| Human[5] | 0.854 | 0.252 | 0.484 | 0.663 | 0.469 | 0.321 | 0.217 |

Metrics

| | |
|---|---|
| CIDEr-D | CIDEr: Consensus-based Image Description Evaluation |
| METEOR | Meteor Universal: Language Specific Translation Evaluation for Any Target Language |
| Rouge-L | ROUGE: A Package for Automatic Evaluation of Summaries |
| BLEU | BLEU: a Method for Automatic Evaluation of Machine Translation |

**http://mscoco.org/dataset/#captions-leaderboard**

**Overview**    **Challenges**    **Download**    **Evaluate**    **Leaderboard**

Table-C5    Table-C40    **2015 Captioning Challenge**    Last update: June 8, 2015. Visit CodaLab for the latest results.

| | M1 | M2 | M3 | M4 | M5 |
|---|---|---|---|---|---|
| Human[5] | 0.638 | 0.675 | 4.836 | 3.428 | 0.352 |
| Google[4] | | | | | |
| MSR[8] | | | | | |
| Montreal | | | | | |
| MSR Ca | | | | | |
| Berkeley | | | | | |
| m-RNN[1] | | | | | |
| Nearest Neighbor[11] | 0.216 | 0.255 | 3.801 | 2.716 | 0.196 |
| PicSOM[13] | 0.202 | 0.250 | 3.965 | 2.552 | 0.182 |
| Brno University[3] | 0.194 | 0.213 | 3.079 | 3.482 | 0.154 |
| m-RNN (Baidu/ UCLA)[16] | 0.190 | 0.241 | 3.831 | 2.548 | 0.195 |
| MIL[6] | 0.168 | 0.197 | 3.349 | 2.915 | 0.159 |
| MLBL[7] | 0.167 | 0.196 | 3.659 | 2.420 | 0.156 |

| M1 | Percentage of captions that are evaluated as better or equal to human caption. |
|---|---|
| M2 | Percentage of captions that pass the Turing Test. |
| M3 | Average correctness of the captions on a scale 1-5 (incorrect - correct). |
| M4 | Average amount of detail of the captions on a scale 1-5 (lack of details - very detailed). |
| M5 | Percentage of captions that are similar to human description. |

# Generative model for diverse captioning

- We would like to sample diverse captions given an image to accurately reflect intrinsic open-endedness of the task



**LSTM + beam search output lacks diversity**

a close up of a plate of food with a sandwich on a table
a close up of a sandwich on a plate
a close up of a plate of food on a table
a close up of a plate of food with a sandwich on it
a close up of a plate of food on a white plate

L. Wang, A. Schwing, and S. Lazebnik, Diverse and Accurate Image Description Using a Variational Auto-Encoder with an Additive Gaussian Encoding Space, NeurIPS 2017

# Generative model for diverse captioning

- We would like to sample diverse captions given an image to accurately reflect intrinsic open-endedness of the task



**LSTM + beam search output lacks diversity**

a close up of a plate of food with a sandwich on a table
a close up of a sandwich on a plate
a close up of a plate of food on a table
a close up of a plate of food with a sandwich on it
a close up of a plate of food on a white plate

**Conditional variational auto-encoder with additive Gaussian space (AG-CVAE)**

a close up of a plate of food on a table
a table with a plate of food on it
a plate of food with a sandwich on it
a white plate topped with a plate of food
a plate of food on a table next to a cup of coffee

L. Wang, A. Schwing, and S. Lazebnik, Diverse and Accurate Image Description Using a Variational Auto-Encoder with an Additive Gaussian Encoding Space, NeurIPS 2017

# CVAE for captioning

L. Wang, A. Schwing, and S. Lazebnik, Diverse and Accurate Image Description Using a Variational Auto-Encoder with an Additive Gaussian Encoding Space, NeurIPS 2017

# CVAE for captioning

Standard CVAE objective:

Sentence

Latent vector

$$\log p_\theta(x|c) \geq \mathbb{E}_{q_\phi(z|x,c)}[\log p_\theta(x|z,c)] - D_{\mathrm{KL}}[q_\phi(z|x,c), p(z|c)]$$

Image content

Decoder distribution

Encoder distribution

Zero-mean Gaussian prior

D. Kingma and M. Welling, Auto-encoding variational Bayes, ICLR 2014

# CVAE with additive Gaussian prior

Proposed objective: shift prior mean based on image content

$$\max_{\theta,\phi} \sum_{i=1}^{N} \log p_\theta(x^i|z^i,c^i) - D_{\mathrm{KL}}[q_\phi(z|x,c), \boxed{p(z|c)}], \quad \text{s.t.} \; \forall i \; z^i \sim q_\phi(z|x,c).$$

$$p(z|c) = \mathcal{N}\left(z \,\Big|\, \sum_{k=1}^{K} c_k \mu_k, \; \sigma^2 \mathrm{I}\right)$$



"dining table"

"teddy bear"

$C_z$

"cup"

L. Wang, A. Schwing, and S. Lazebnik, Diverse and Accurate Image Description Using a Variational Auto-Encoder with an Additive Gaussian Encoding Space, NeurIPS 2017

# Results

- More controllable captions: changing the conditioning vector of object labels changes the caption in a reasonable way



Predicted Object Labels:
'person' 'cup' 'donut' 'dining table'

AG-CVAE:

a woman sitting at a table with a cup of coffee
a person sitting at a table with a cup of coffee
a table with two plates of donuts and a cup of coffee
a woman sitting at a table with a plate of coffee
a man sitting at a table with a plate of food

LSTM Baseline:

a close up of a table with two plates of coffee
a close up of a table with a plate of food
a close up of a plate of food on a table
a close up of a table with two plates of food
a close up of a table with plates of food

L. Wang, A. Schwing, and S. Lazebnik, Diverse and Accurate Image Description Using a Variational Auto-Encoder with an Additive Gaussian Encoding Space, NeurIPS 2017

# Results

- More controllable captions: changing the conditioning vector of object labels changes the caption in a reasonable way



Object Labels: 'person'
AG-CVAE sentences:
a man and a woman standing in a room
a man and a woman are playing a game
a man standing next to a woman in a room
a man standing next to a woman in a field
a man standing next to a woman in a suit

Object Labels: 'person', 'remote'
AG-CVAE sentences:
a man and a woman playing a video game
a man and a woman are playing a video game
a man and woman are playing a video game
a man and a woman playing a game with a remote
a woman holding a nintendo wii game controller

L. Wang, A. Schwing, and S. Lazebnik, Diverse and Accurate Image Description Using a Variational Auto-Encoder with an Additive Gaussian Encoding Space, NeurIPS 2017

# Acknowledgement

Thanks to the following courses and corresponding researchers for making their teaching/research material online

- Deep Learning, Stanford University

- Introduction to Deep Learning, University of Illinois at Urbana-Champaign

- Introduction to Deep Learning, Carnegie Mellon University

- Convolutional Neural Networks for Visual Recognition, Stanford University

- Natural Language Processing with Deep Learning, Stanford University

- And Many More ......