

Deep Learning for Computer Vision

Fall 2022

<https://cool.ntu.edu.tw/courses/189345> (NTU COOL)

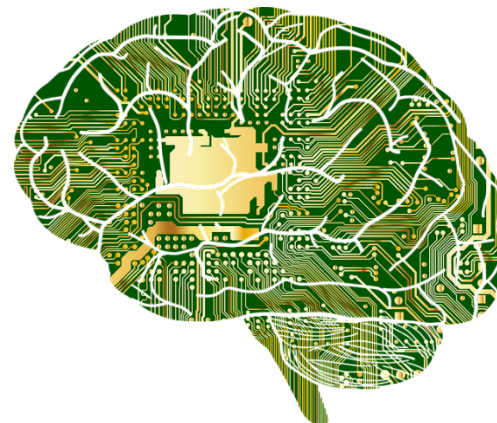
<http://vllab.ee.ntu.edu.tw/dlcv.html> (Public website)

Yu-Chiang Frank Wang 王鈺強, Professor

Dept. Electrical Engineering, National Taiwan University

What to Be Covered Today...

- Transfer Learning
 - Visual Classification – Domain Adaptation
 - Visual Synthesis – Style Transfer
- Recurrent Neural Networks
 - From RNN to LSTM & GRU
 - Selected Models for Sequence-to-Sequence Learning
 - Attention in RNN



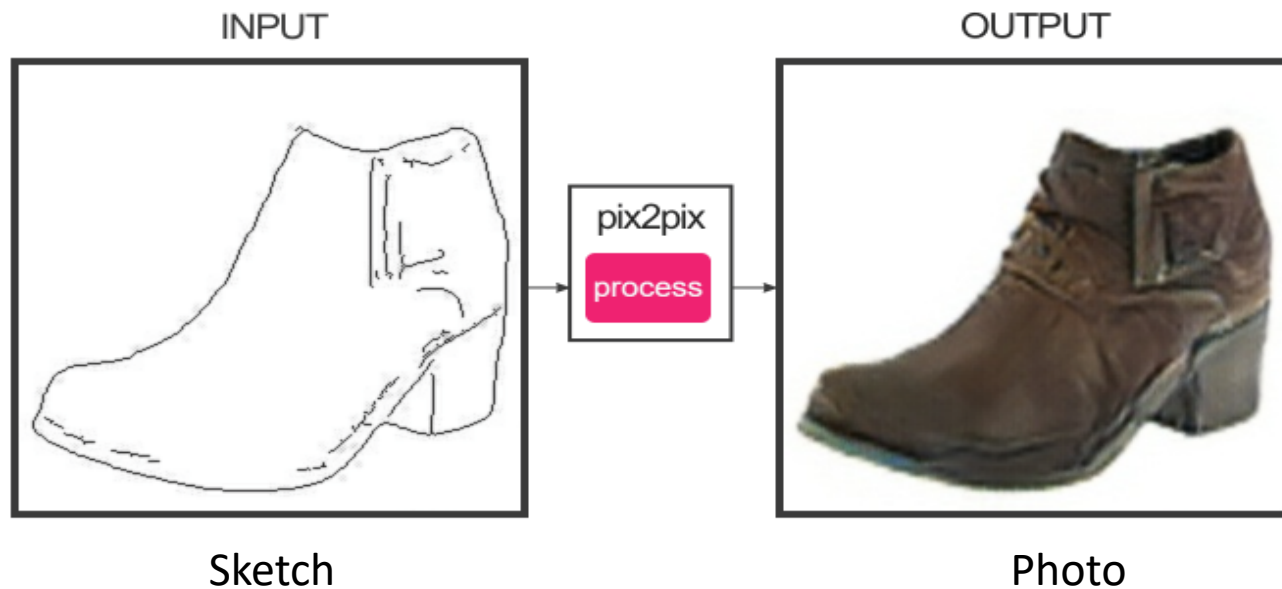
Transfer Learning for Image Synthesis

- Cross-Domain Image Translation
 - Pix2pix: Pairwise cross-domain training data
 - CycleGAN/DualGAN/DiscoGAN: Unpaired cross-domain training data
 - UNIT: Learning cross-domain image representation (with unpaired training data)
 - AdaIN: Single-image arbitrary style transfer in real-time
 - Beyond image translation



Pix2pix

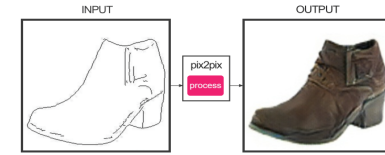
- Image-to-image translation with conditional adversarial networks (CVPR'17)
 - Can be viewed as image style transfer



Pix2pix

$$\begin{aligned} \{\tilde{x}, x_{in}\} &\rightarrow [D] \rightarrow T/F \\ \{x_{real}, x_{in}\} &\rightarrow [D] \rightarrow T/F \end{aligned}$$

Testing Phase



- **Goal / Problem Setting**

- Image translation across two distinct domains (e.g., sketch v.s. photo)

- ✓ • **Pairwise** training data

- **Method: Conditional GAN**

- Example: Sketch to Photo

- **Generator**

Input: Sketch

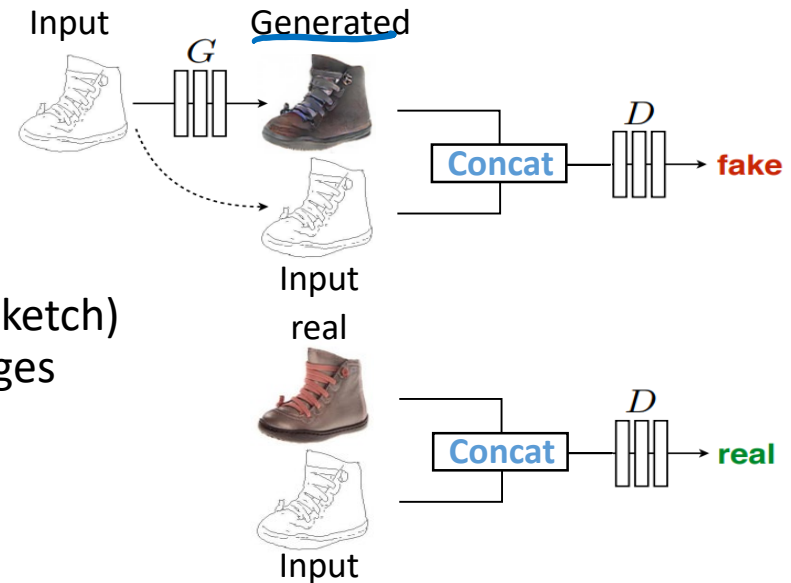
Output: Photo

- **Discriminator**

Input: Concatenation of Input(Sketch) & Synthesized/Real(Photo) images

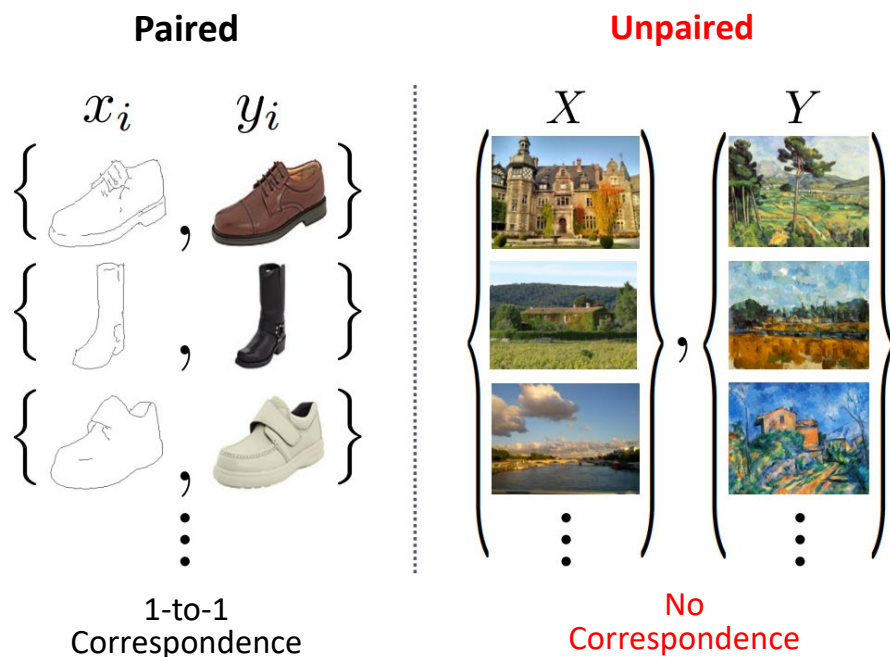
Output: Real or Fake

Training Phase



CycleGAN/DiscoGAN/DualGAN

- CycleGAN
 - Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks -to-image translation with conditional adversarial networks



- Easier to collect training data
- More practical

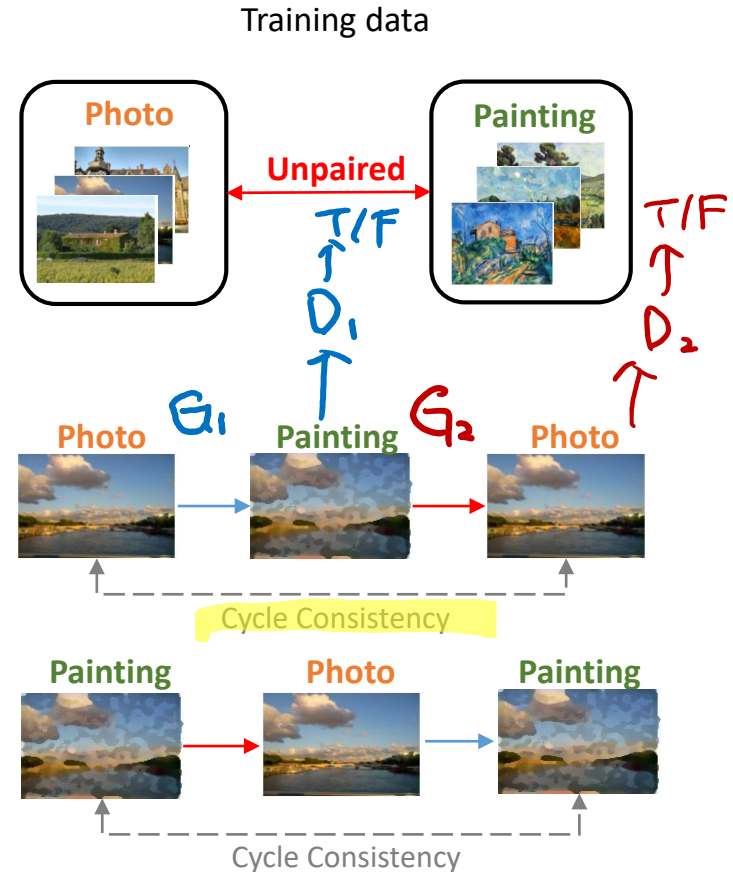
CycleGAN

- **Goal / Problem Setting**

- Image translation across two distinct domains
- **Unpaired** training data

- **Idea**

- Autoencoding-like image translation
- **Cycle consistency** between two domains

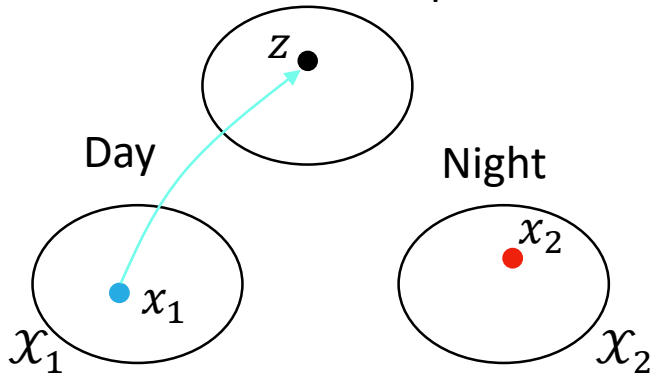


UNIT

- Unsupervised Image-to-Image Translation Networks (NIPS'17)
 - Image translation via learning cross-domain joint representation

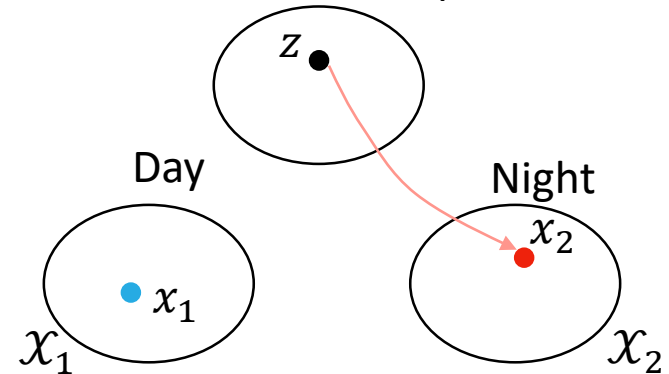
Stage1: Encode to the joint space

\mathcal{Z} : Joint latent space



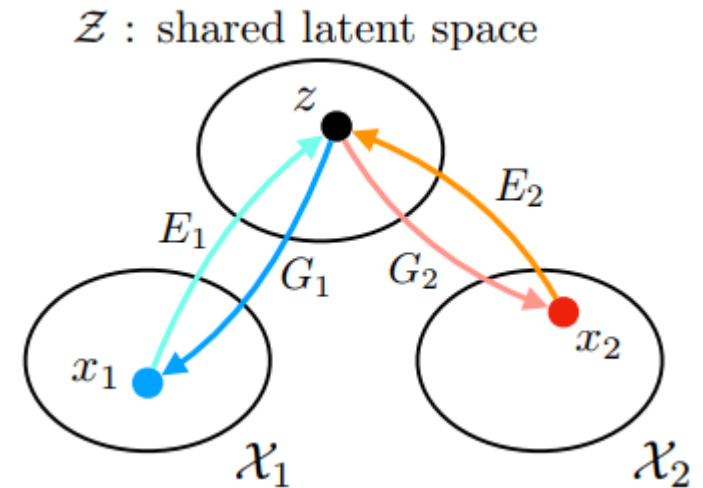
Stage2: Generate cross-domain images

\mathcal{Z} : Joint latent space



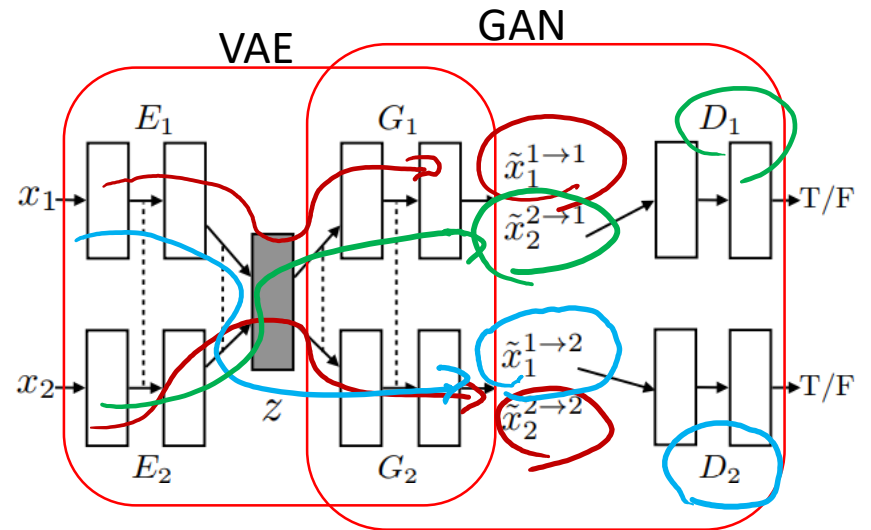
UNIT

- Goal/Problem Setting
 - Image translation across two distinct domains
 - Unpaired training image data
- Idea
 - Based on two parallel VAE-GAN models



source

tgt



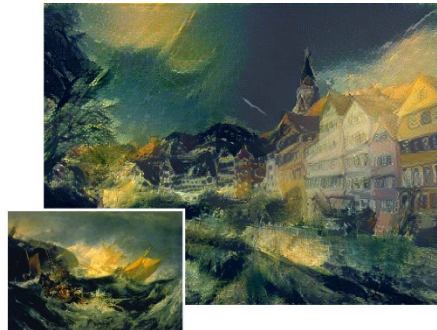
Transfer Learning for Image Synthesis

- Cross-Domain Image Translation
 - Pix2pix: Pairwise cross-domain training data
 - CycleGAN/DualGAN/DiscoGAN: Unpaired cross-domain training data
 - UNIT: Learning cross-domain image representation (with unpaired training data)
 - AdaIN: Single-image arbitrary style transfer in real-time
 - Beyond image translation

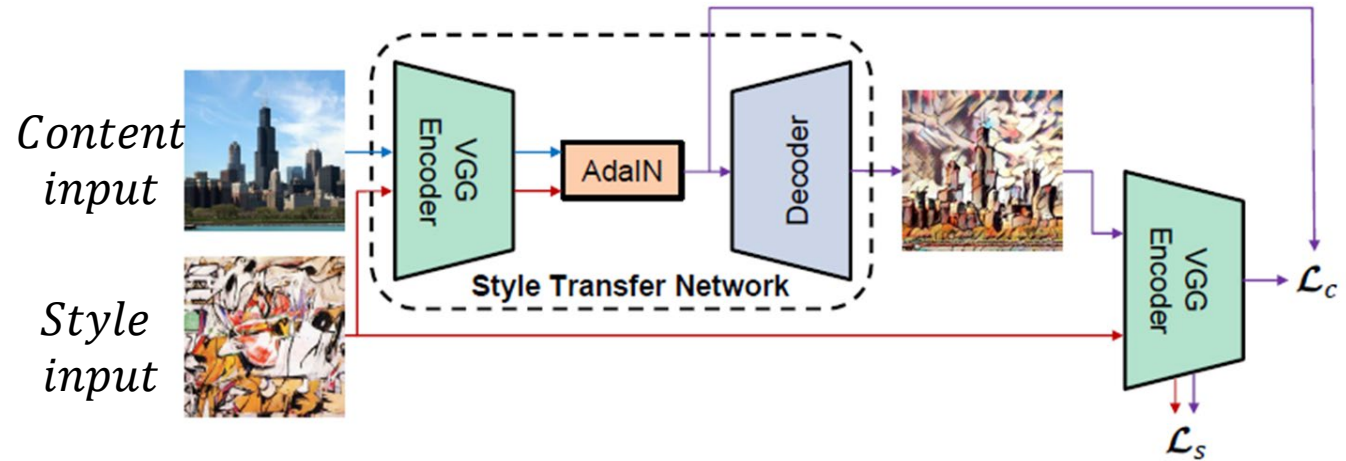


AdaIN

- We've talked about style transfer methods like Pix2Pix or CycleGAN.
- Arbitrary Style Transfer in Real-time with Adaptive Instance Normalization (ICCV'17)
 - Single-image arbitrary style transfer in real-time



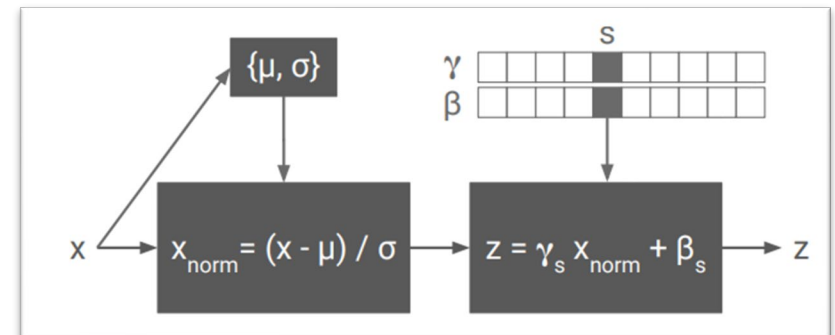
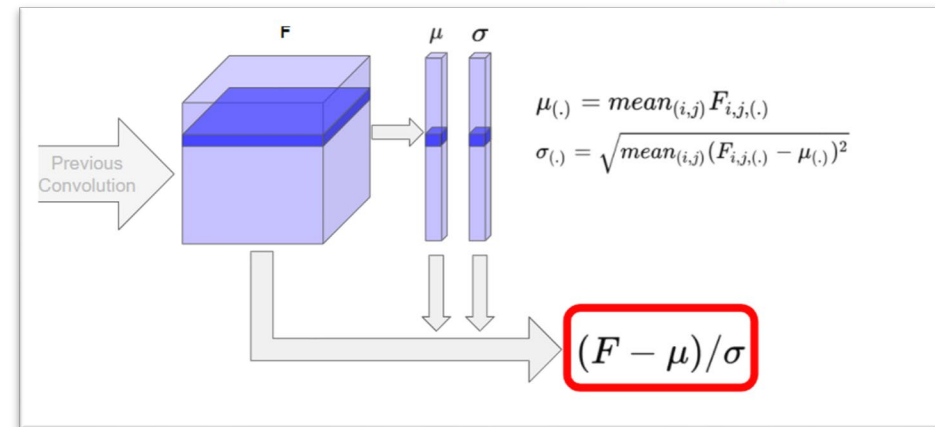
AdaIN



- Adaptive Instance Normalization

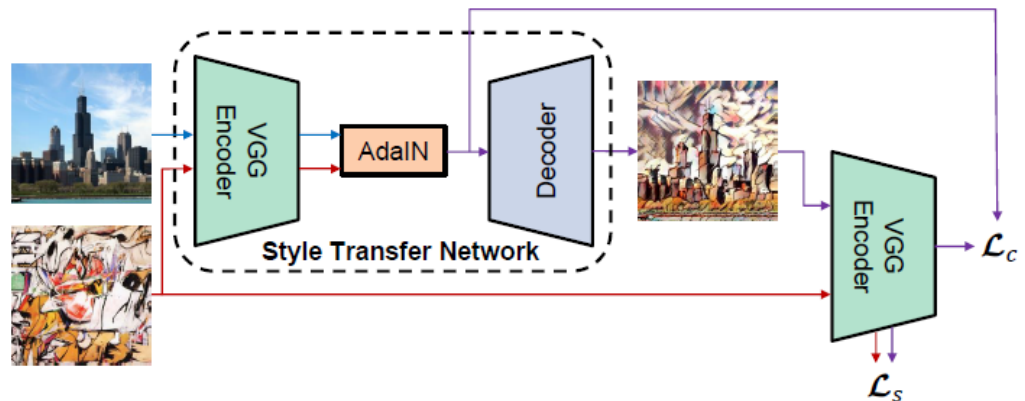
$$\text{AdaIN}(x, y) = \sigma(y) \left(\frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y)$$

- x : content input, y : style input
- No learnable affine parameters
- Perform style transfer in the feature space



AdaIN (cont'd)

- f : Encoder, g : Decoder



Content loss (via content/perceptual consistency):

$$\mathcal{L}_{content}(\vec{p}, \vec{x}, l) = \frac{1}{2} \sum_{i,j} (F_{ij}^l - P_{ij}^l)^2 .$$

Style loss (via Gram matrix loss):

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l \quad \Rightarrow \quad E_l = \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2$$

$$\mathcal{L}_{style}(\vec{a}, \vec{x}) = \sum_l^L w_l E_l$$

$$\Rightarrow \mathcal{L}_{total}(\vec{p}, \vec{a}, \vec{x}) = \alpha \mathcal{L}_{content}(\vec{p}, \vec{x}) + \beta \mathcal{L}_{style}(\vec{a}, \vec{x})$$

AdaIN

- Qualitative results



Style

Content

Ours

Chen and Schmidt

Ulyanov *et al.*

Gatys *et al.*

Transfer Learning for Image Synthesis

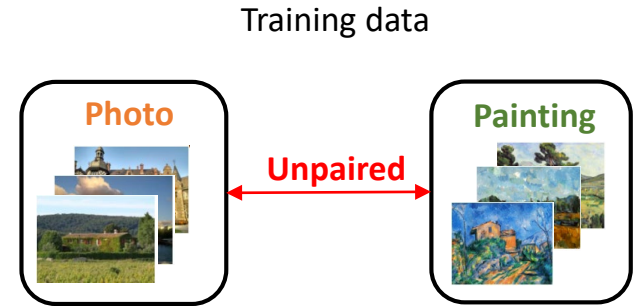
- Cross-Domain Image Translation
 - Pix2pix: Pairwise cross-domain training data
 - CycleGAN/DualGAN/DiscoGAN: Unpaired cross-domain training data
 - UNIT: Learning cross-domain image representation (with unpaired training data)
 - AdaIN
 - Beyond image translation



Revisit: CycleGAN

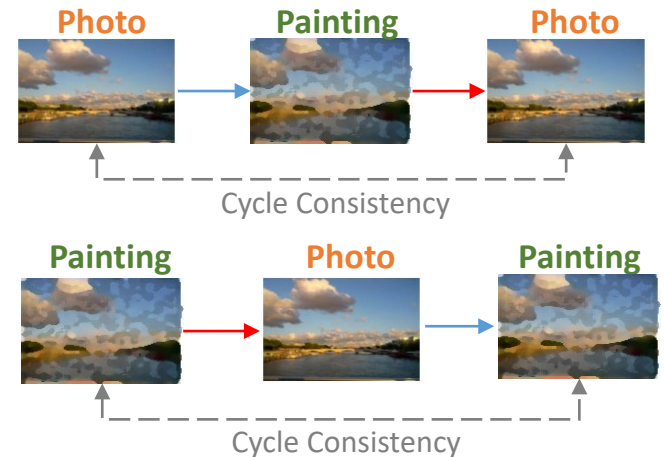
- **Goal / Problem Setting**

- Image translation across two distinct domains
- **Unpaired** training data

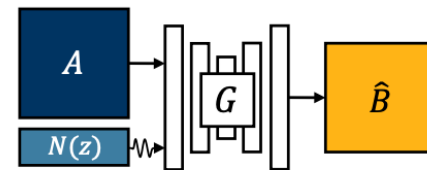


- **Idea**

- Autoencoding-like image translation
- **Cycle consistency** between two domains

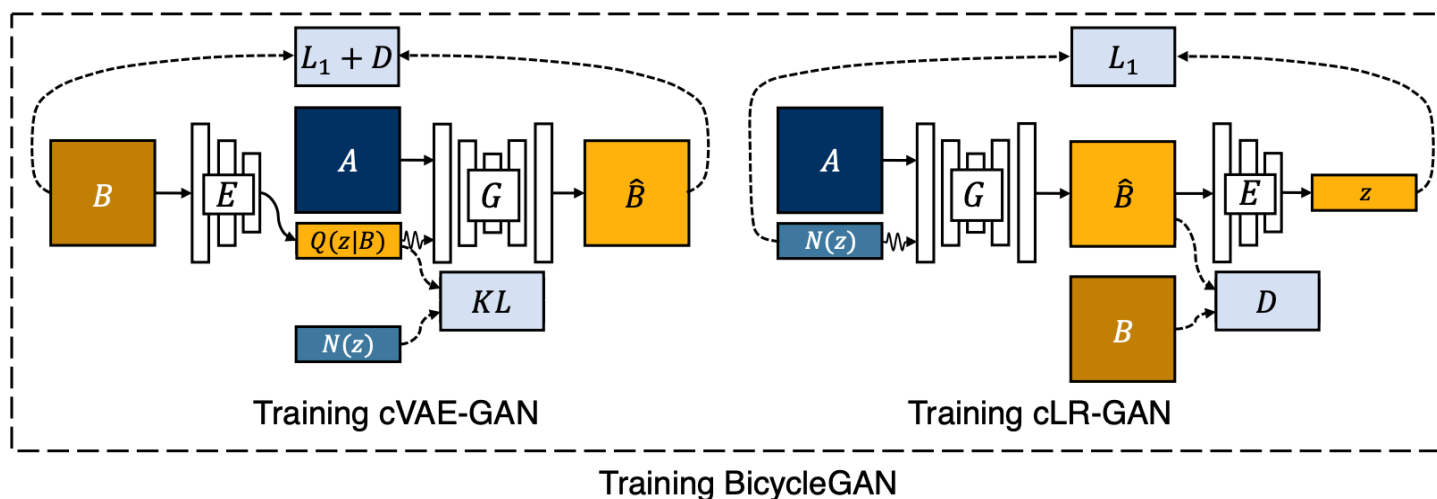


BicycleGAN

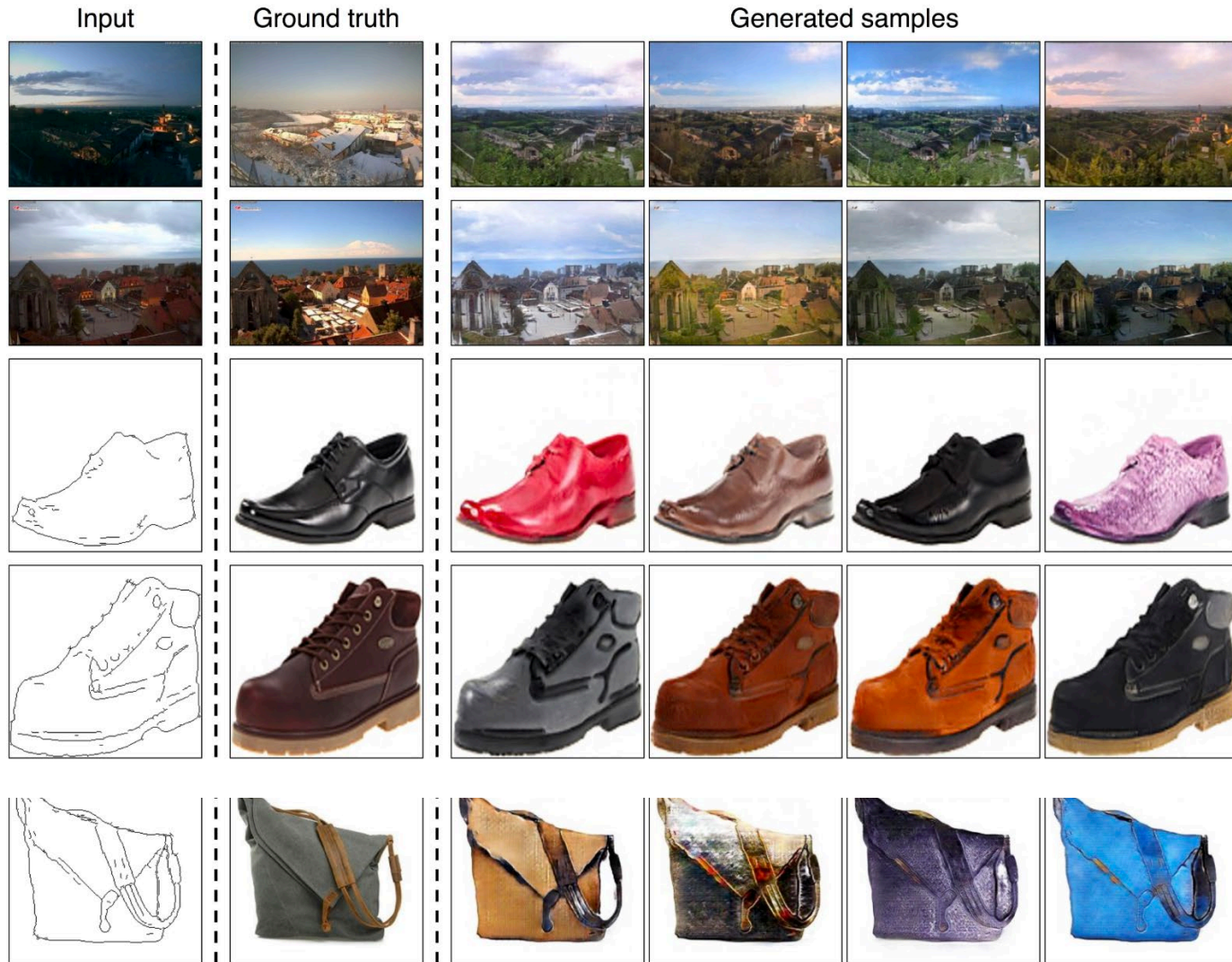


(a) Testing Usage for all models

- Toward Multimodal Image-to-Image Translation (NIPS'17)
- Goal / Problem Setting
 - Producing **diverse** images across two distinct domains.
 - **Pairwise** training data
- Idea
 - Combine conditional VAE-GAN and conditional Latent Regressor GAN.



BicycleGAN - Experiment



DRIT

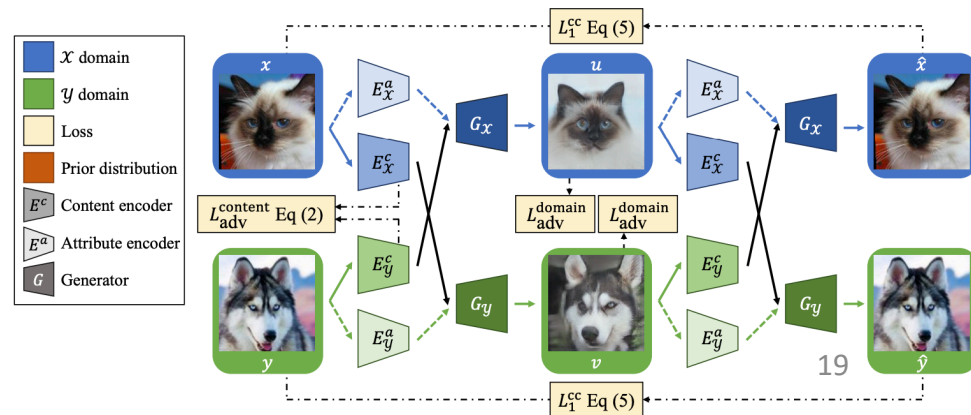
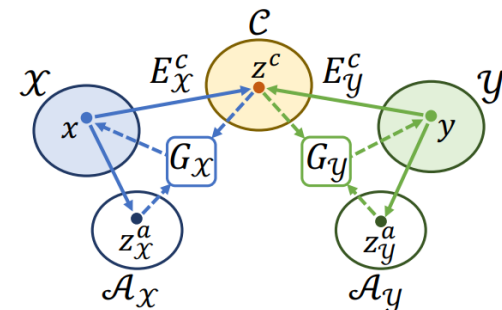
- **Diverse Image-to-Image Translation via Disentangled Representations**
(ECCV'18 oral)

- **Goal / Problem Setting**

- Producing diverse images across two distinct domains.
 - **Unpaired** training data

- **Idea**

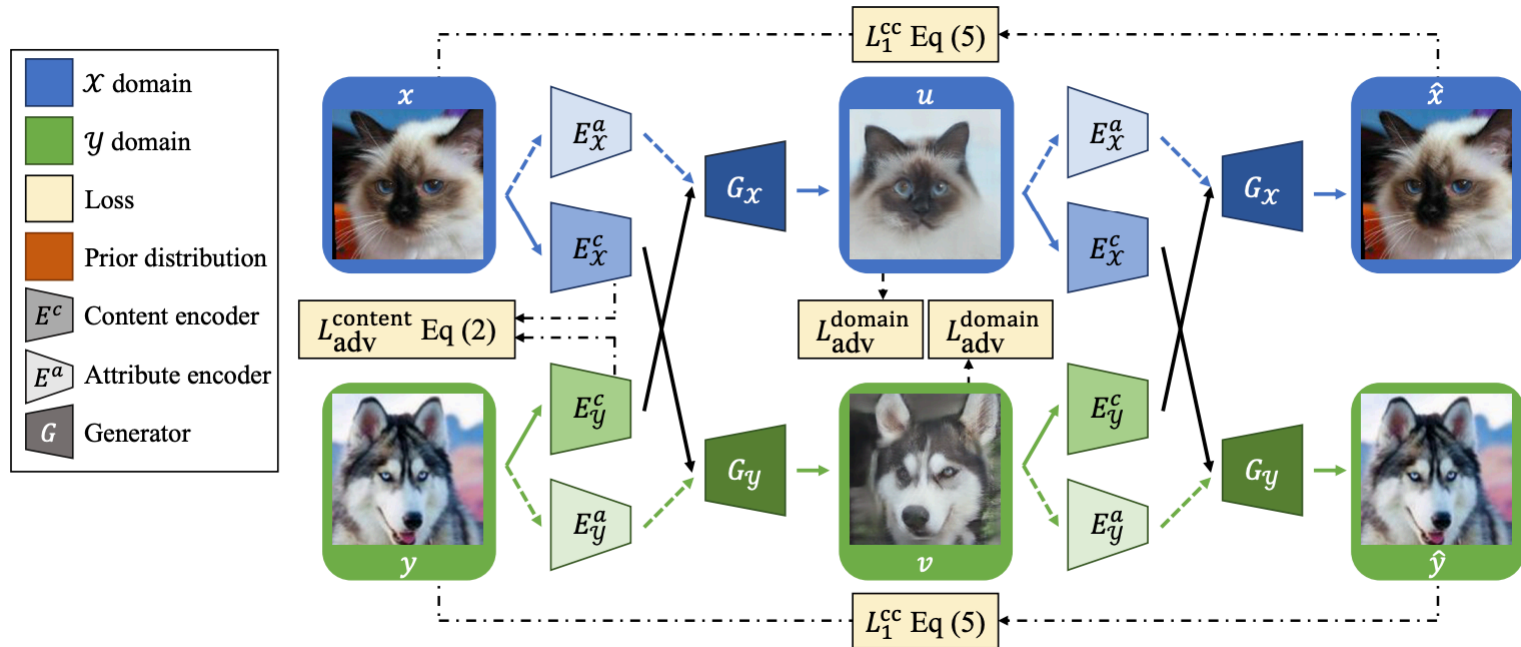
- Disentangle latent representation into **domain-invariant** and **domain-specific** features
 - Generate cross-domain images by swapping the latent feature from each domain.
 - Applied cross-cycle consistency



Method – Main Framework

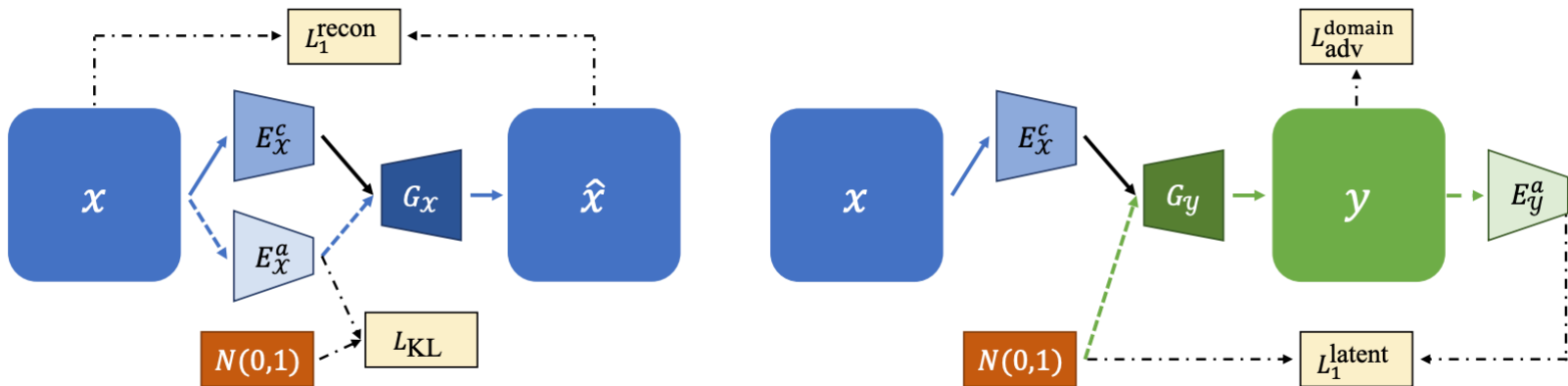
Attribute: species

Content: pose (style)

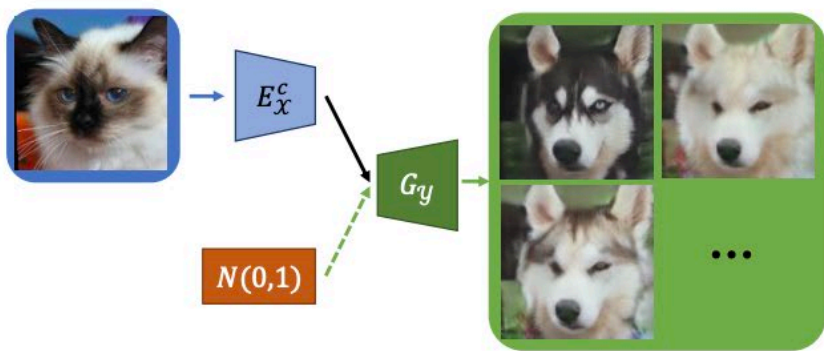


Method – For Attribute Features

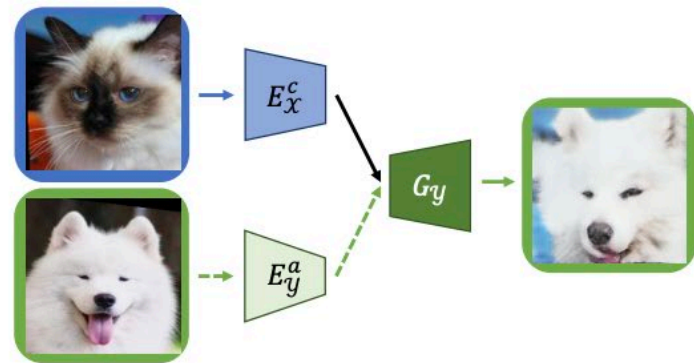
- **KL loss:**
perform stochastic sampling at test time.
- **Latent regression loss:**
encourage invertible mapping btw image and latent representations



Method – Inference phase

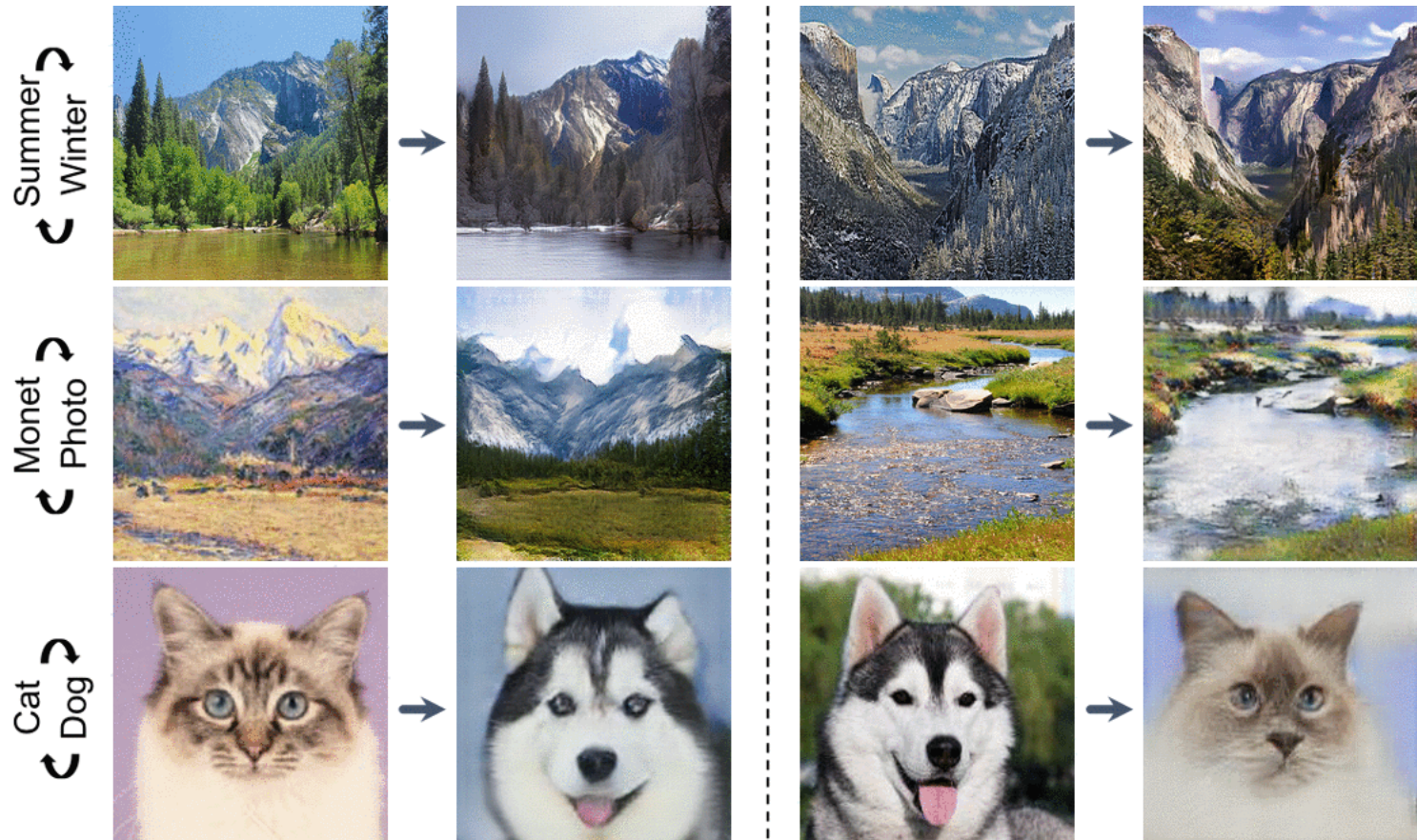


(b) Testing with random attributes



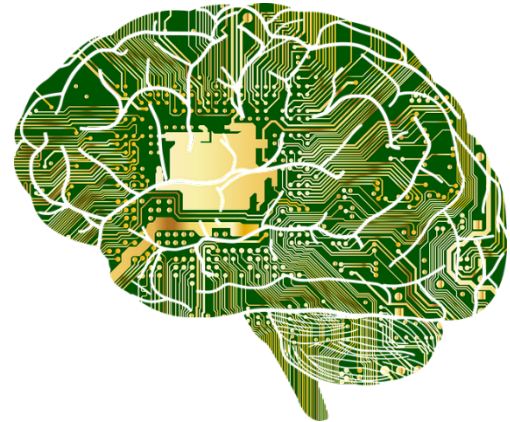
(c) Testing with a given attribute

Example Results



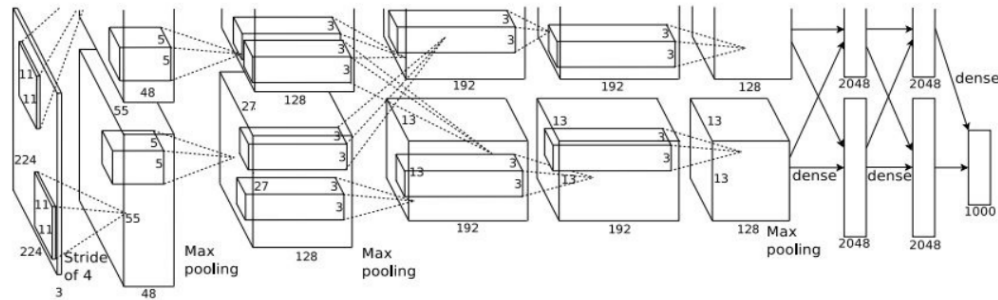
What to Be Covered Today...

- Transfer Learning
 - Visual Classification – Domain Adaptation
 - Visual Synthesis – Style Transfer
- Recurrent Neural Networks
 - From RNN to LSTM & GRU
 - Selected Models for Sequence-to-Sequence Learning
 - Attention in RNN



What Are The Limitations of CNN?

- Deal with image data
 - Both input and output are images/vectors
- Simply feed-forward processing



DOG
CAT
MONKEY

Example of (Visual) Sequential Data



<https://quickdraw.withgoogle.com/#>

More Applications in Vision

Image Captioning

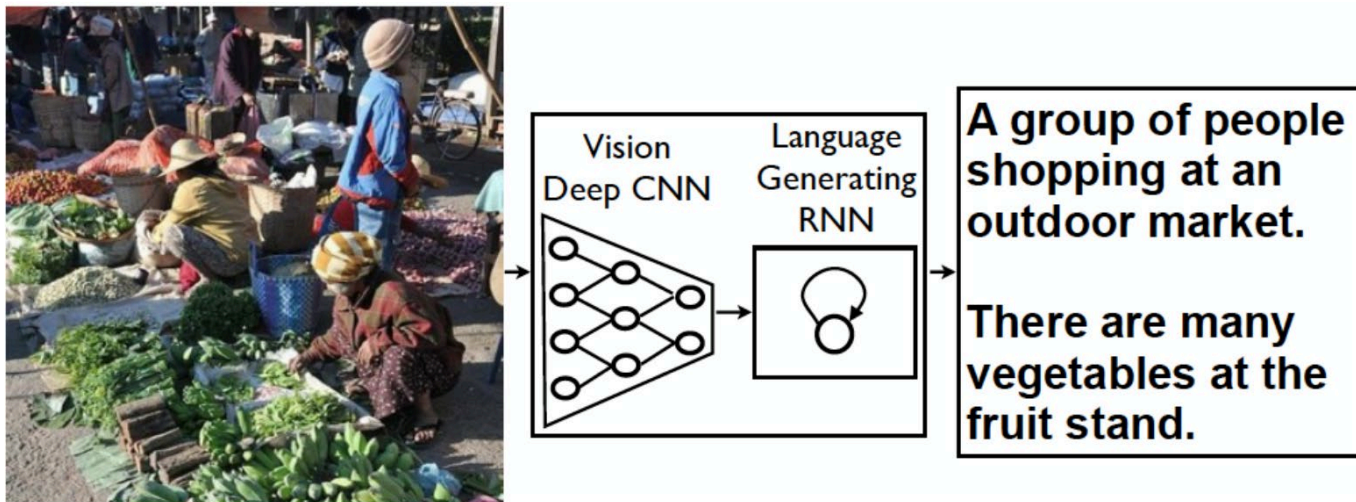


Figure from Vinyals et al, "Show and tell: A neural image caption generator", CVPR 2015

More Applications in Vision

Visual Question Answering (VQA)




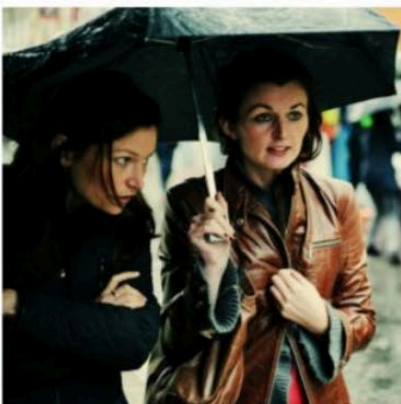
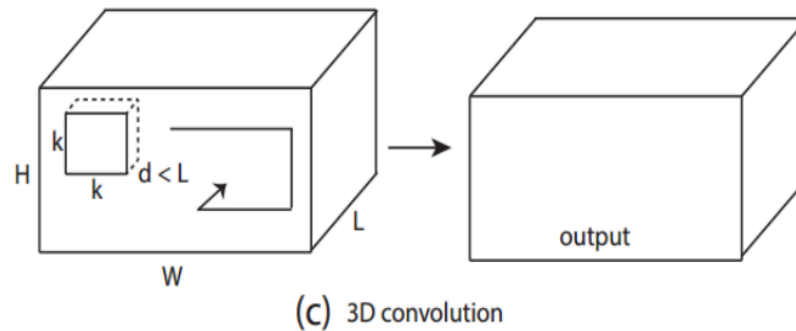
Input				
	<p>Q: What endangered animal is featured on the truck?</p> <p>A: A bald eagle. A: A sparrow. A: A humming bird. A: A raven.</p>	<p>Q: Where will the driver go if turning right?</p> <p>A: Onto 24 1/4 Rd. A: Onto 25 1/4 Rd. A: Onto 23 1/4 Rd. A: Onto Main Street.</p>	<p>Q: When was the picture taken?</p> <p>A: During a wedding. A: During a bar mitzvah. A: During a funeral. A: During a Sunday church service</p>	<p>Q: Who is under the umbrella?</p> <p>A: Two women. A: A child. A: An old man. A: A husband and a wife.</p>
output				

Figure from Zhu et al, "Visual 7W: Grounded Question Answering in Images", CVPR 2016

How to Model Sequential Data?

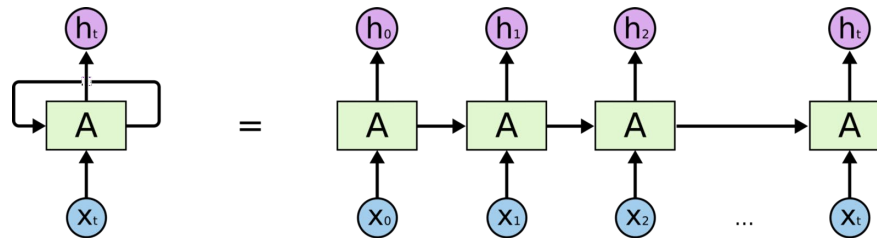
- Deep learning for sequential data
 - 3-dimensional convolution neural networks



3D convolution

How to Model Sequential Data?

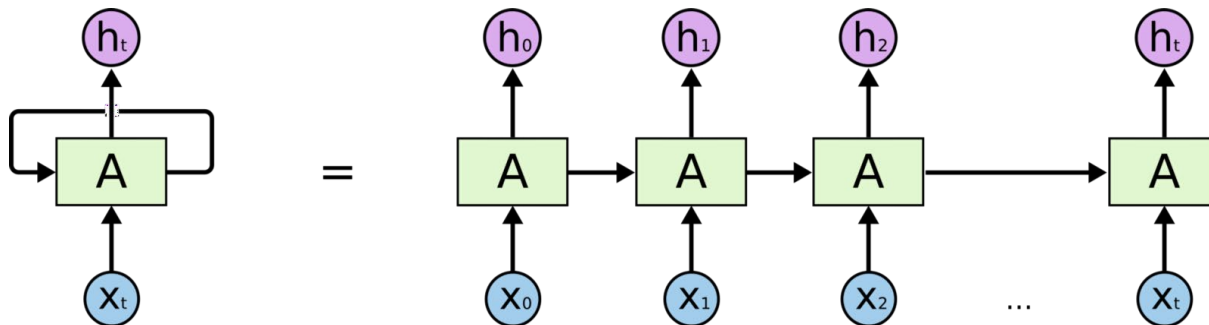
- Deep learning for sequential data
 - Recurrent neural networks (RNN)



RNN

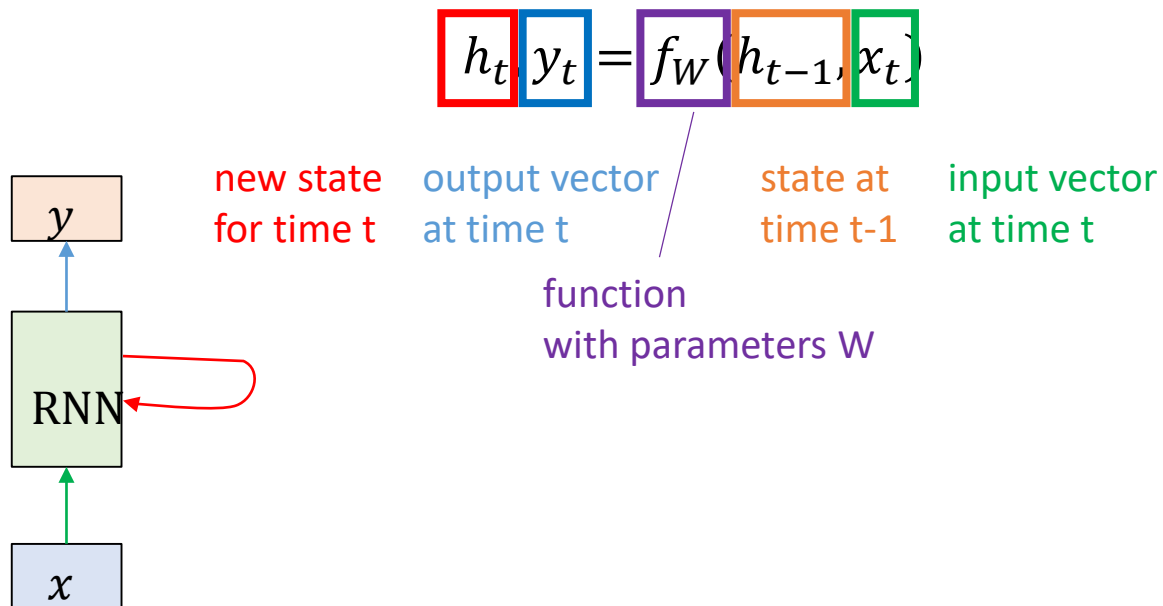
Recurrent Neural Networks

- Parameter sharing + unrolling
 - Keeps the number of parameters fixed
 - Allows sequential data with varying lengths
- Memory ability
 - Capture and preserve information which has been extracted



Recurrence Formula

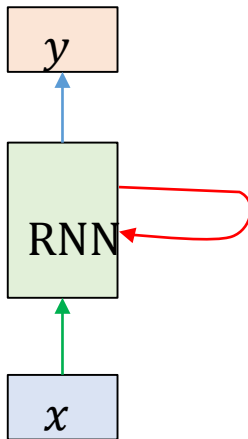
- Same function and parameters used at every time step:



Recurrence Formula

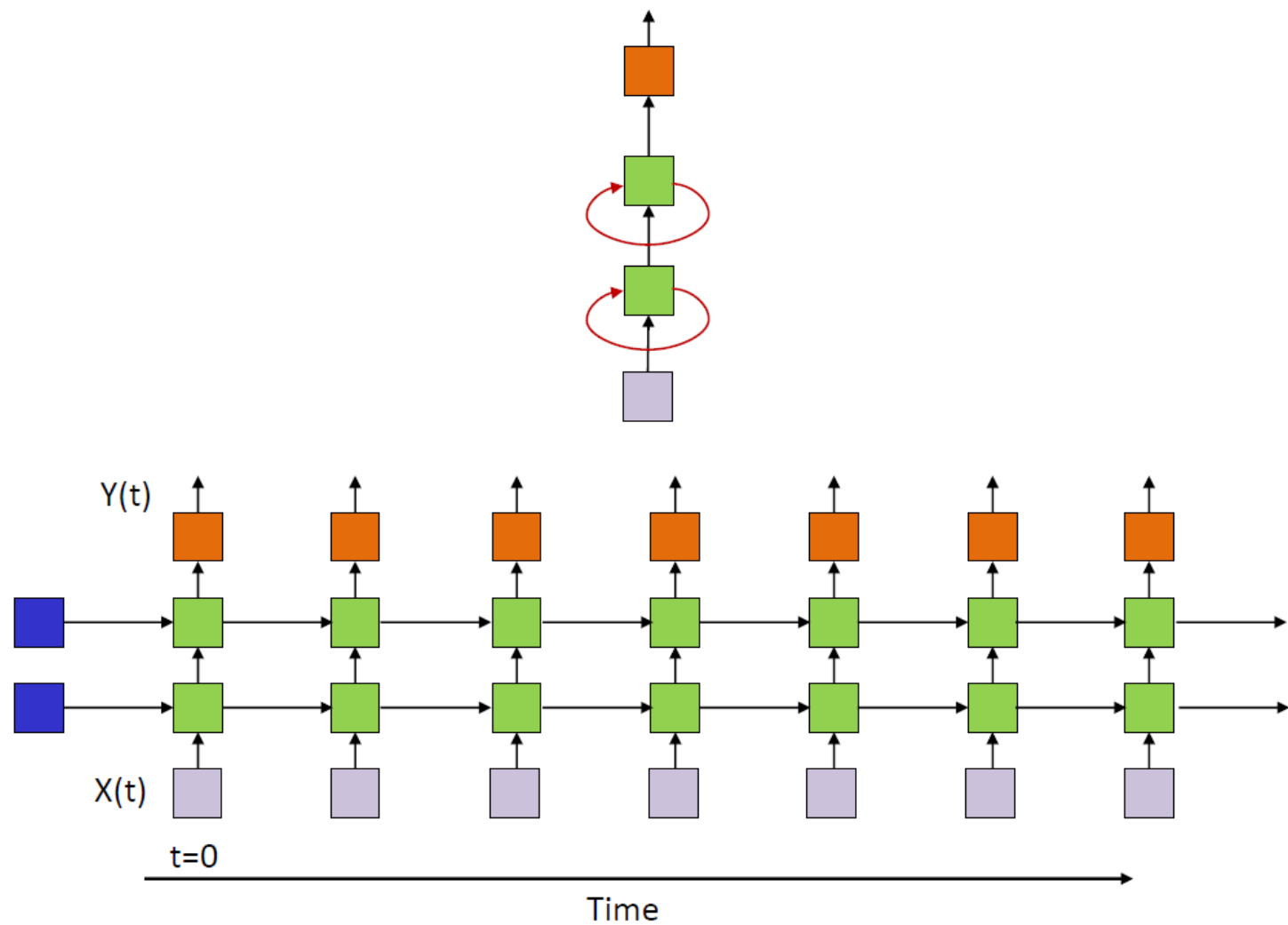
- Same function and parameters used at every time step:

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

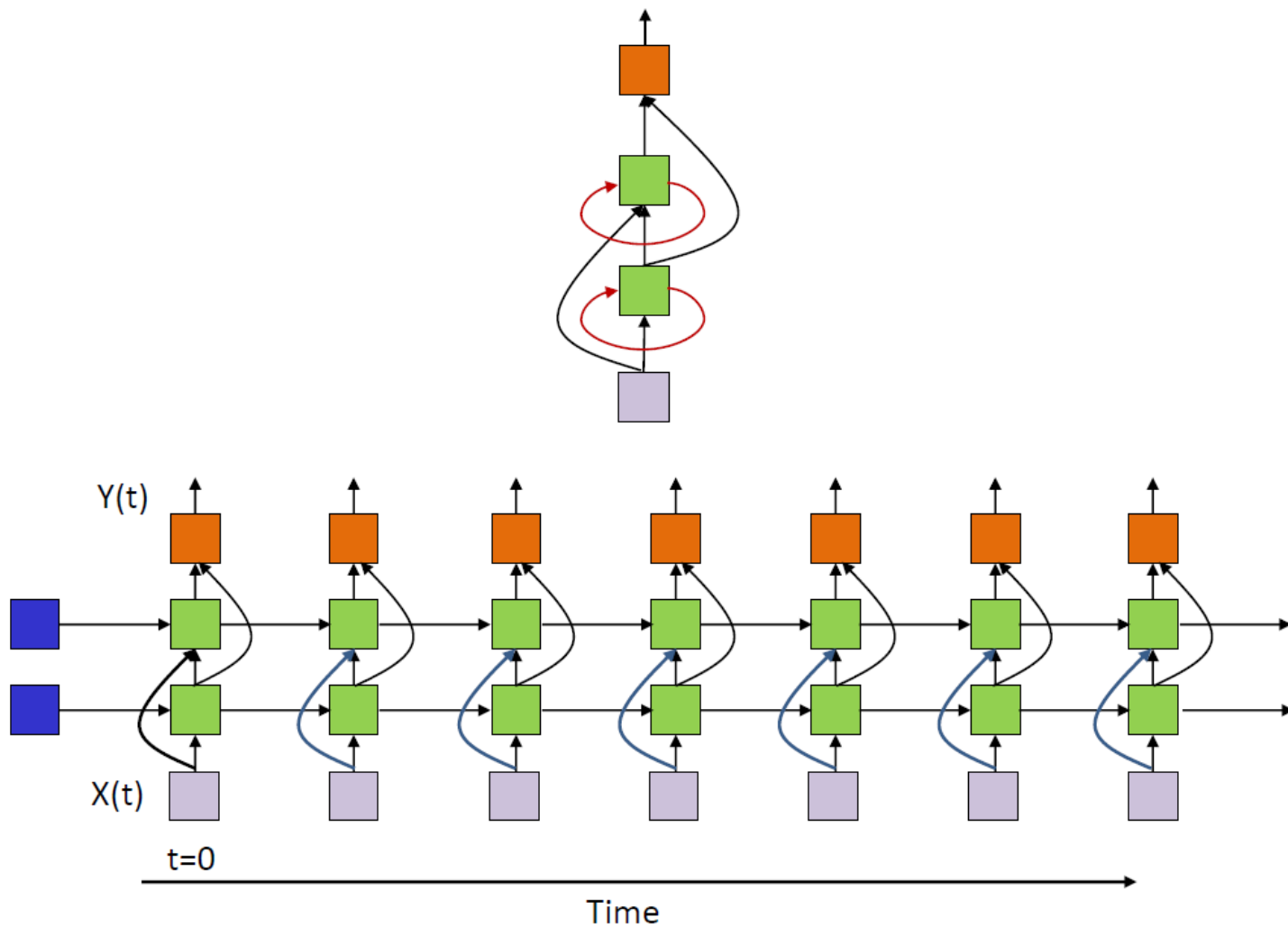


$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$
$$y_t = W_{hy}h_t$$

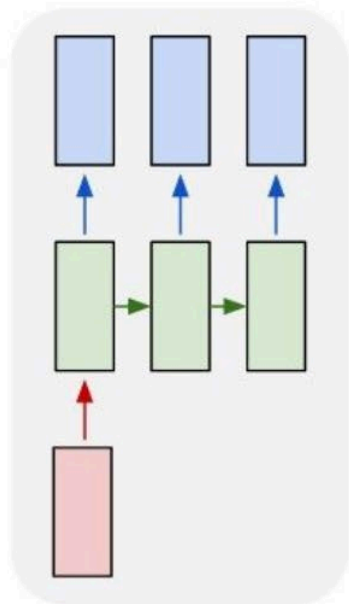
Multiple Recurrent Layers



Multiple Recurrent Layers

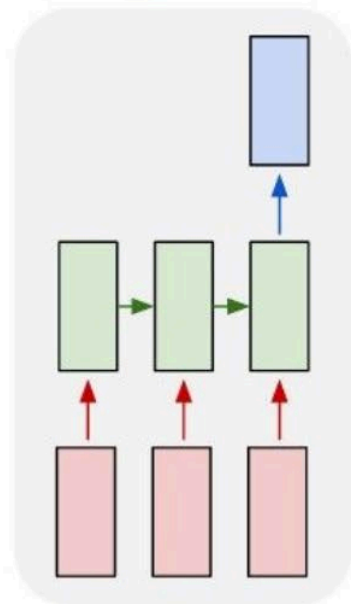


one to many



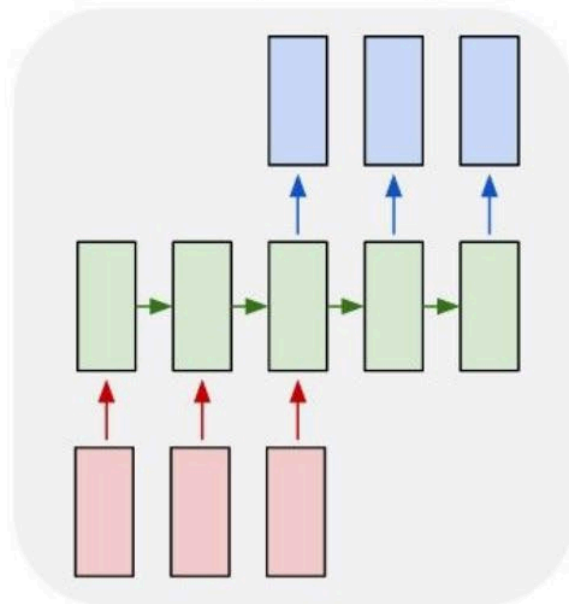
e.g., image caption

many to one



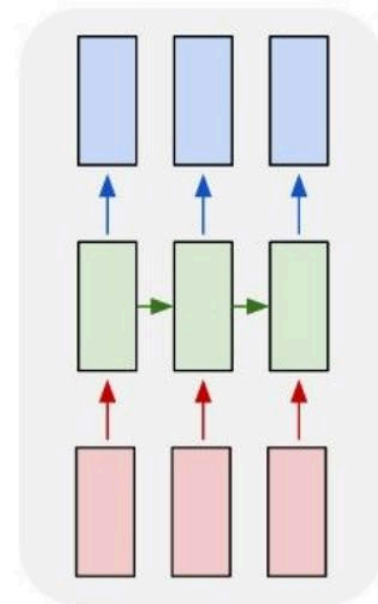
e.g., action recognition

many to many



e.g., video prediction

many to many



e.g., video indexing

Example: Image Captioning

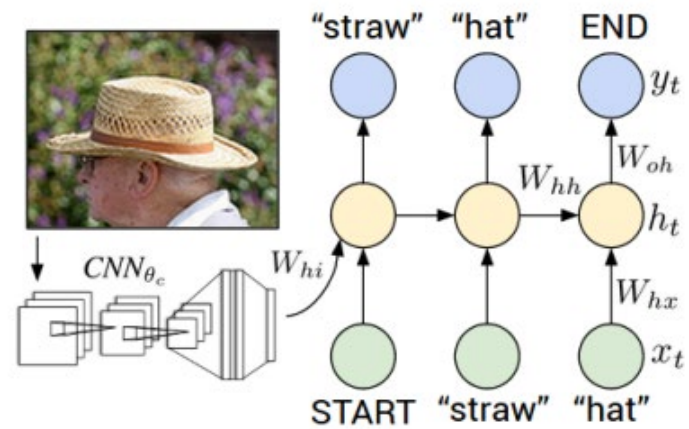
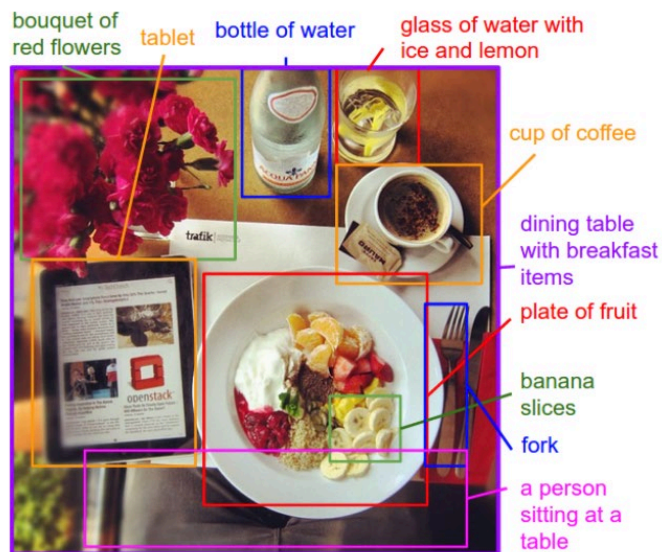


Figure from Karpathy et al, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015

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

CNN

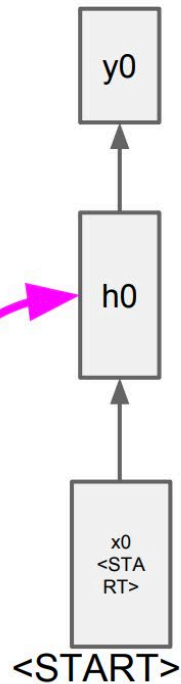


test image



V

Wih



test image

before:

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

now:

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

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



test image

y0

h0

x0
<STA
RT>

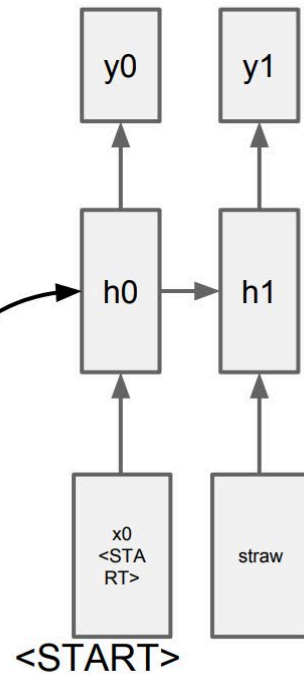
<START>

straw

sample!

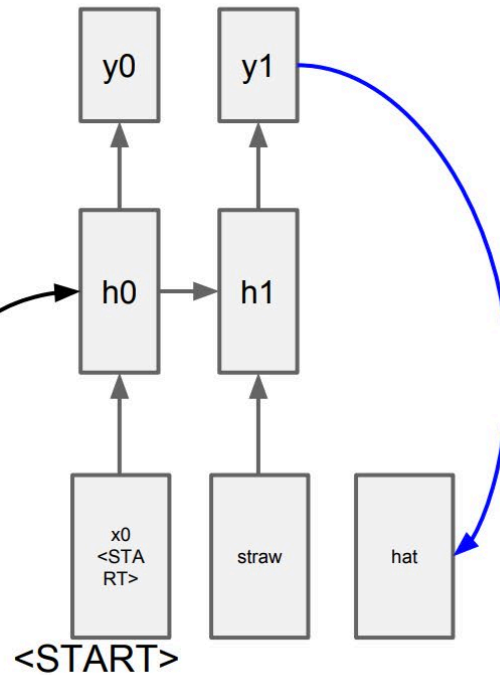


test image





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



test image

y0

y1

y2

h0

h1

h2

x0
<START>

straw

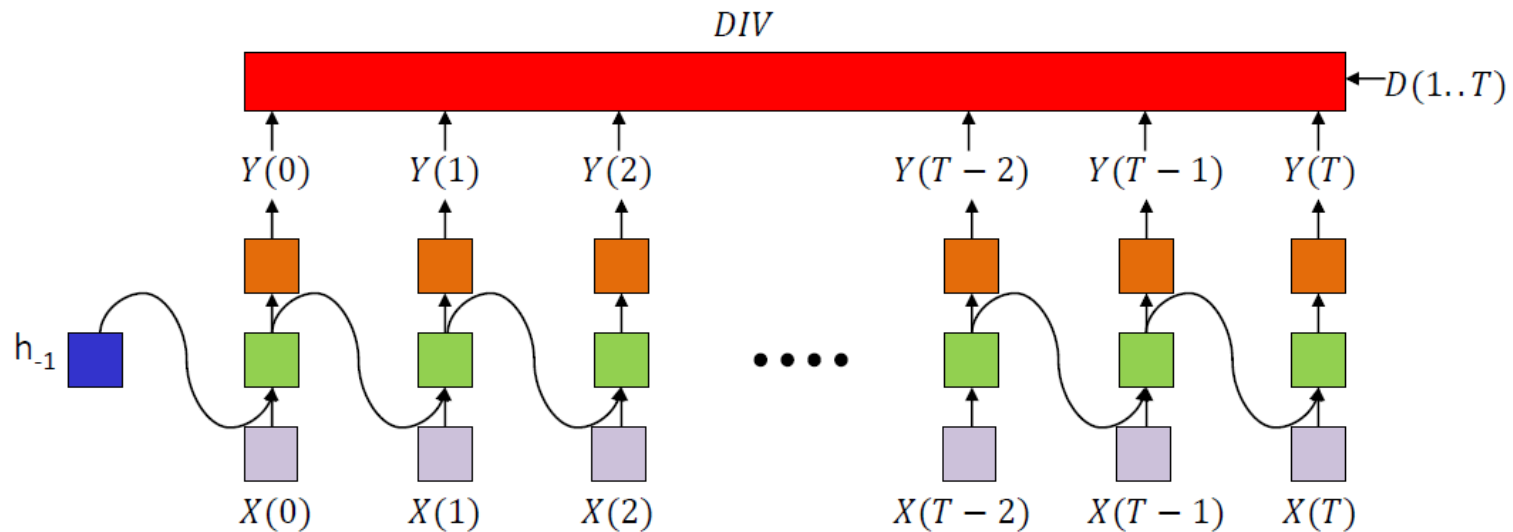
hat

sample
<END> token
=> finish.

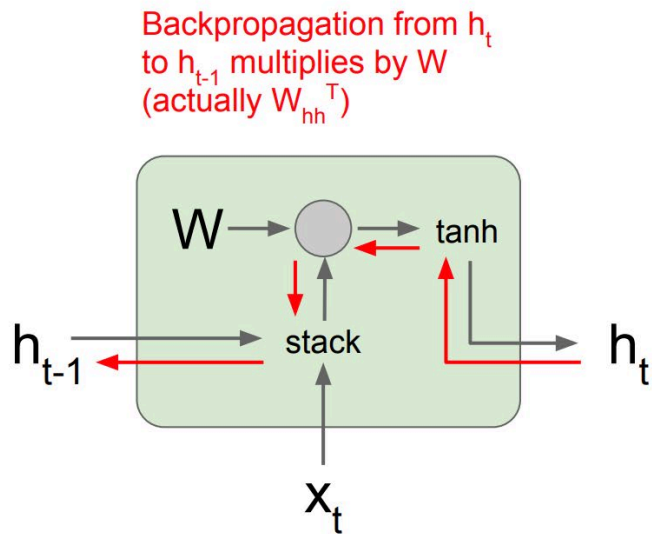
<START>

Training RNNs: Back Propagation Through Time

- Let's focus on one training instance.
- The divergence to be computed is between the sequence of outputs by the network and the desired output sequence.
- Generally, this is not just the sum of the divergences at individual times.

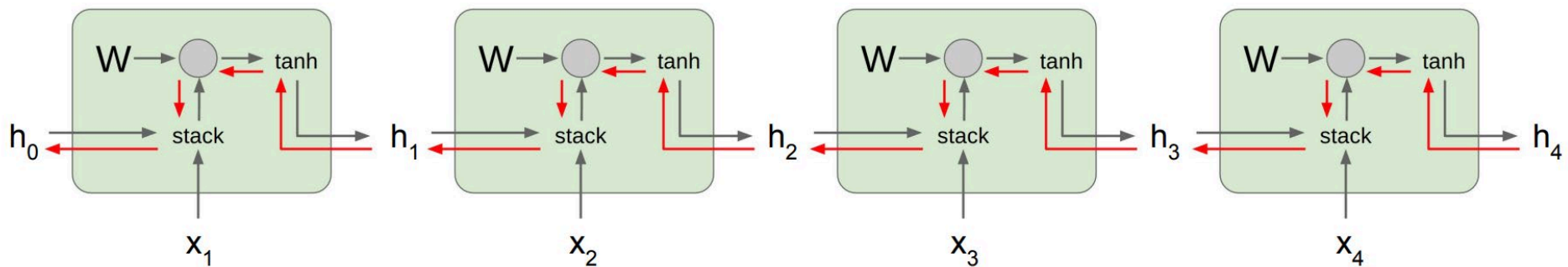


Back Propagation Through Time (BPTT)



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{hx}x_t) \\ &= \tanh\left((W_{hh} \quad W_{hx}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \end{aligned}$$

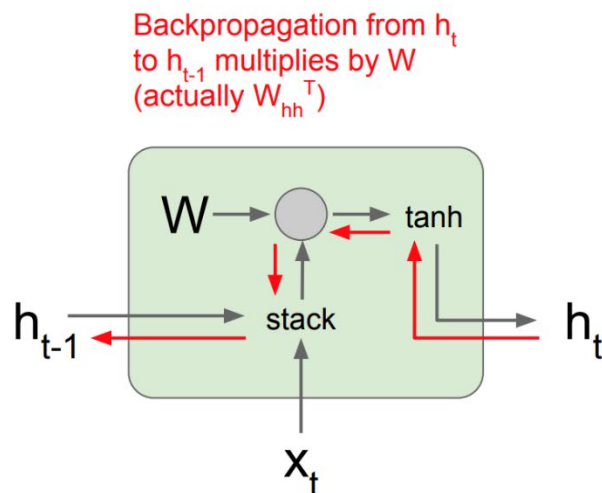
Back Propagation Through Time (BPTT)



Computing gradient of h_0 involves many factors of W (and repeated \tanh)

Gradient Vanishing & Exploding

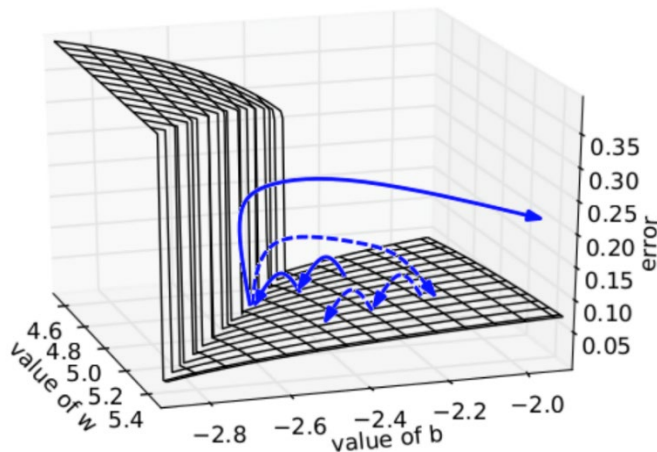
- Computing gradient involves many factors of W
 - Exploding gradients : Largest singular value > 1
 - Vanishing gradients : Largest singular value < 1



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{hx}x_t) \\ &= \tanh\left((W_{hh} \quad W_{hx}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \end{aligned}$$

Solutions...

- Gradients clipping : rescale gradients if too large



→ standard gradient descent trajectories

---→ gradient clipping to fix problem

```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```

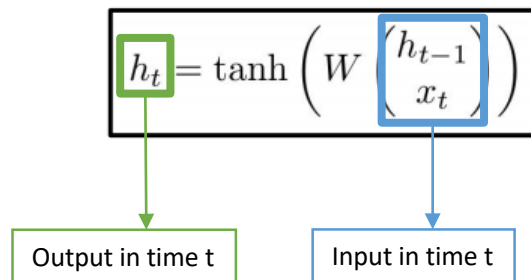
- How about vanishing gradients?
 - Change RNN architecture!

Variants of RNN

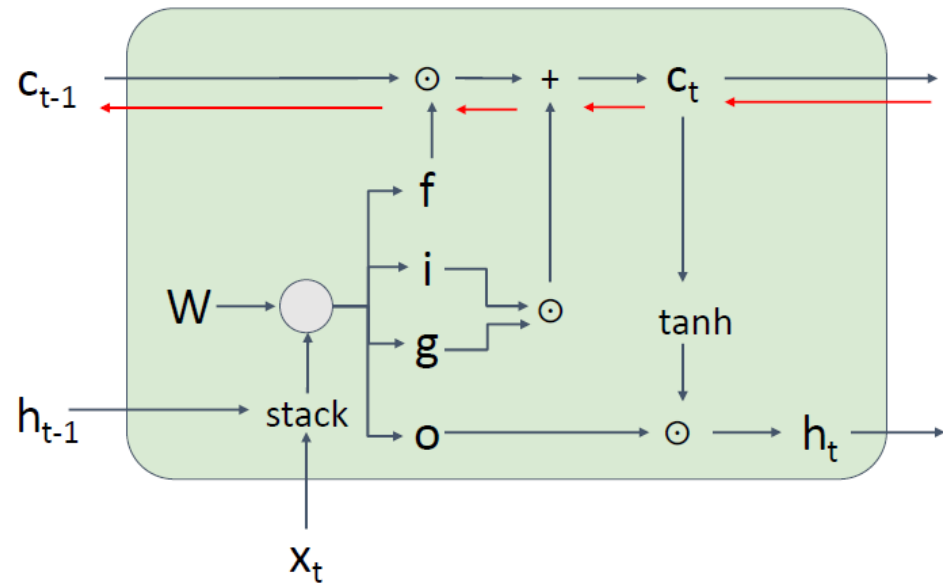
- Long Short-term Memory (LSTM) [Hochreiter et al., 1997]
 - Additional memory cell
 - Input/Forget/Output Gates
 - Handle **gradient vanishing**
 - Learn long-term dependencies
- Gated Recurrent Unit (GRU) [Cho et al., EMNLP 2014]
 - Similar to LSTM
 - No additional memory cell
 - Reset / Update Gates
 - Fewer parameters than LSTM
 - Comparable performance to LSTM [Chung et al., NIPS Workshop 2014]

Vanilla RNN vs. LSTM

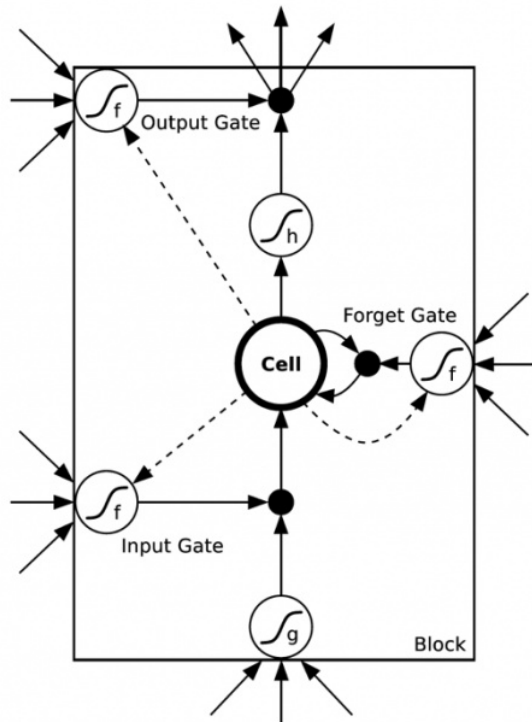
Vanilla RNN



LSTM



Long Short-Term Memory (LSTM)



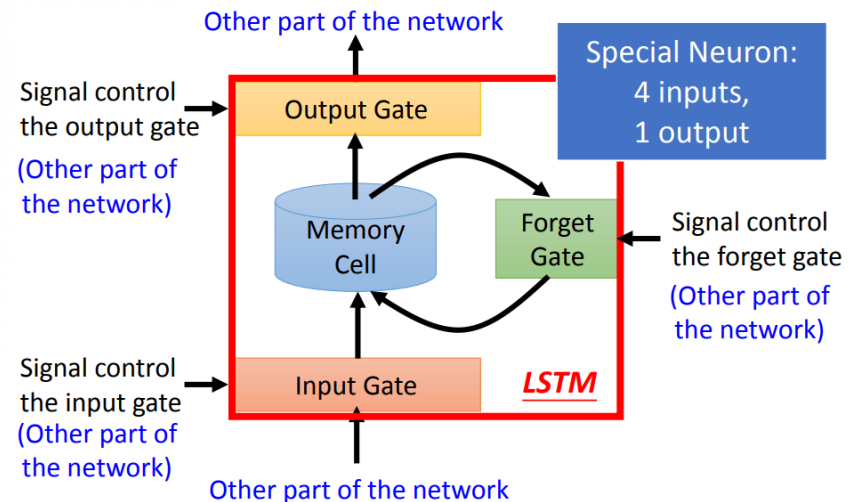
Forget gate
Input gate
Input activation
Output gate
Cell state
Hidden state

$$\begin{aligned}
 f &= \sigma(W_f [h_{t-1}, x_t] + b_f) \\
 i &= \sigma(W_i [h_{t-1}, x_t] + b_i) \\
 \tilde{h} &= \tanh(W_h [h_{t-1}, x_t] + b_h) \\
 o &= \sigma(W_o [h_{t-1}, x_t] + b_o) \\
 c_t &= f \odot c_{t-1} + i \odot \tilde{h} \\
 h_t &= o \odot \tanh(c_t)
 \end{aligned}$$

Memory Cell

Output in time t

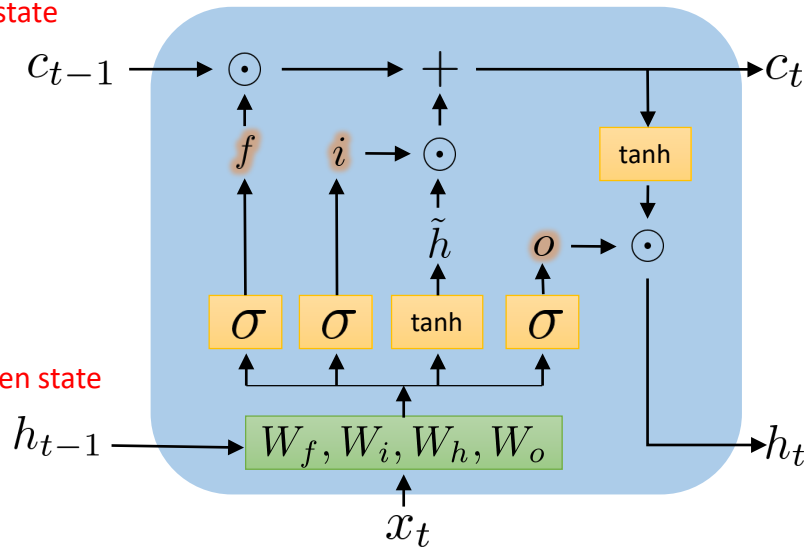
Input in time t



Long Short-Term Memory (LSTM)

Cell state

Hidden state



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

Input gate $i = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$

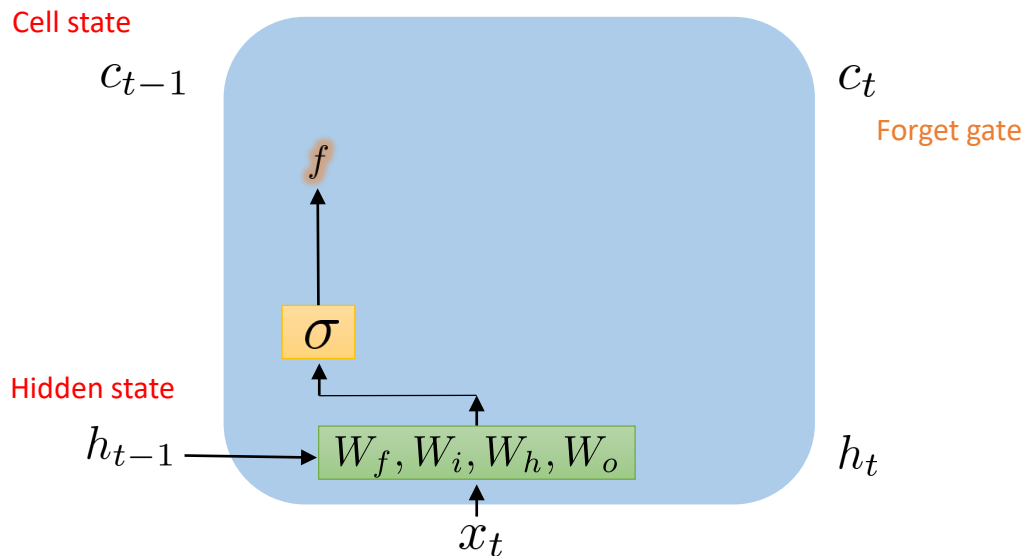
Input activation $\tilde{h} = \tanh (W_h \cdot [h_{t-1}, x_t] + b_h)$

Output gate $o = \sigma (W_o \cdot [h_{t-1}, x_t] + b_o)$

Cell state $c_t = f \odot c_{t-1} + i \odot \tilde{h}$

Hidden state $h_t = o \odot \tanh(c_t)$

Long Short-Term Memory (LSTM)



Calculate forget gate f :
whether to erase cell

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

$$i = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

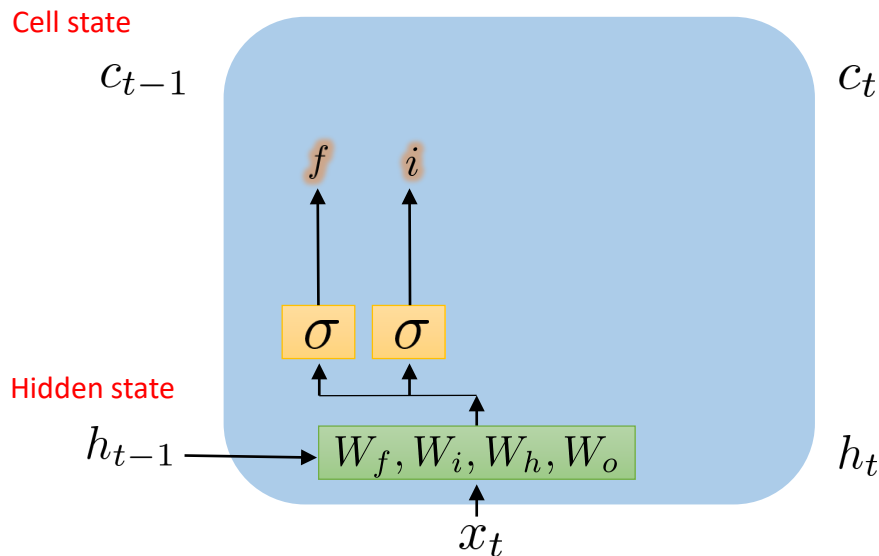
$$\tilde{h} = \tanh (W_h \cdot [h_{t-1}, x_t] + b_h)$$

$$o = \sigma (W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$c_t = f \odot c_{t-1} + i \odot \tilde{h}$$

$$h_t = o \odot \tanh(c_t)$$

Long Short-Term Memory (LSTM)



Calculate input gate i :
whether to write to cell

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

$$i = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

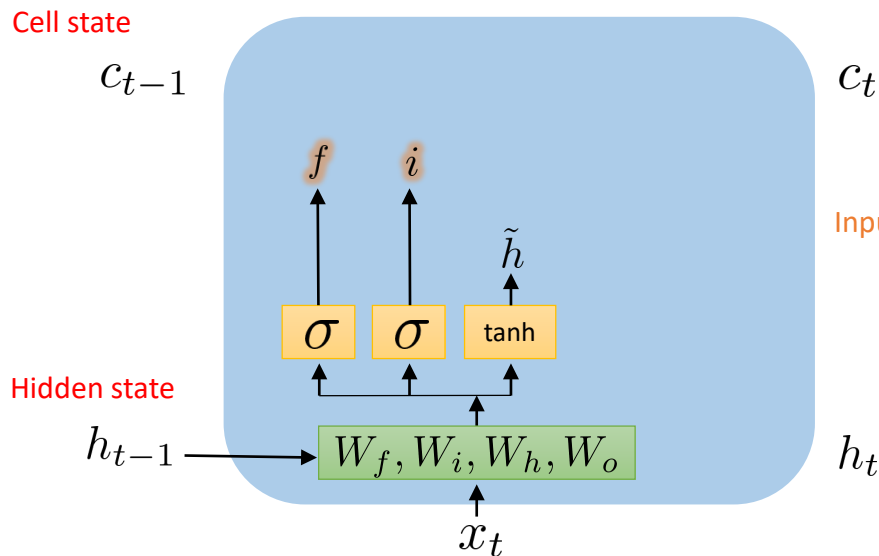
$$\tilde{h} = \tanh (W_h \cdot [h_{t-1}, x_t] + b_h)$$

$$o = \sigma (W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$c_t = f \odot c_{t-1} + i \odot \tilde{h}$$

$$h_t = o \odot \tanh(c_t)$$

Long Short-Term Memory (LSTM)



Calculate input activation:
how much to write to cell

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

$$i = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

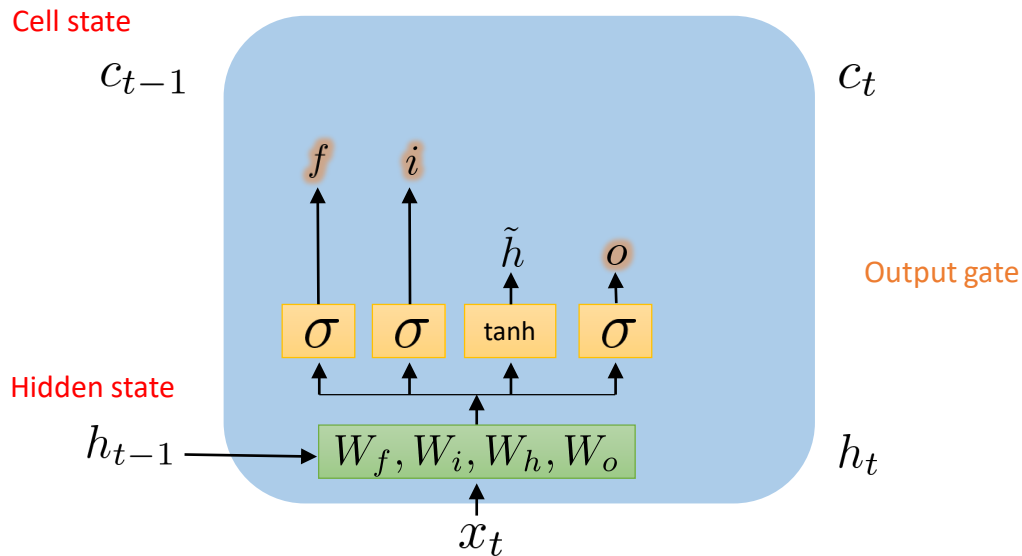
Input activation $\tilde{h} = \tanh (W_h \cdot [h_{t-1}, x_t] + b_h)$

$$o = \sigma (W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$c_t = f \odot c_{t-1} + i \odot \tilde{h}$$

$$h_t = o \odot \tanh(c_t)$$

Long Short-Term Memory (LSTM)



Calculate output gate o :
how much to reveal cell

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

$$i = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

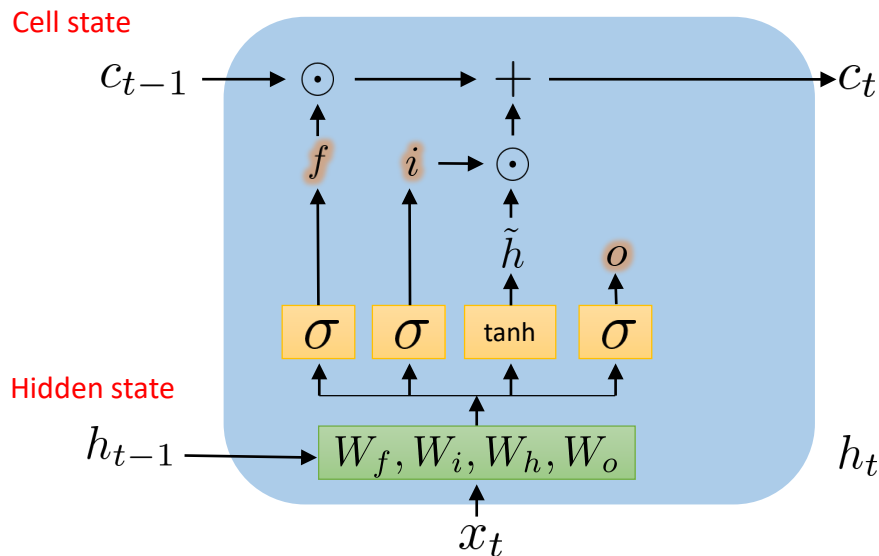
$$\tilde{h} = \tanh (W_h \cdot [h_{t-1}, x_t] + b_h)$$

$$o = \sigma (W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$c_t = f \odot c_{t-1} + i \odot \tilde{h}$$

$$h_t = o \odot \tanh(c_t)$$

Long Short-Term Memory (LSTM)



Update memory cell

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

$$i = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{h} = \tanh (W_h \cdot [h_{t-1}, x_t] + b_h)$$

$$o = \sigma (W_o \cdot [h_{t-1}, x_t] + b_o)$$

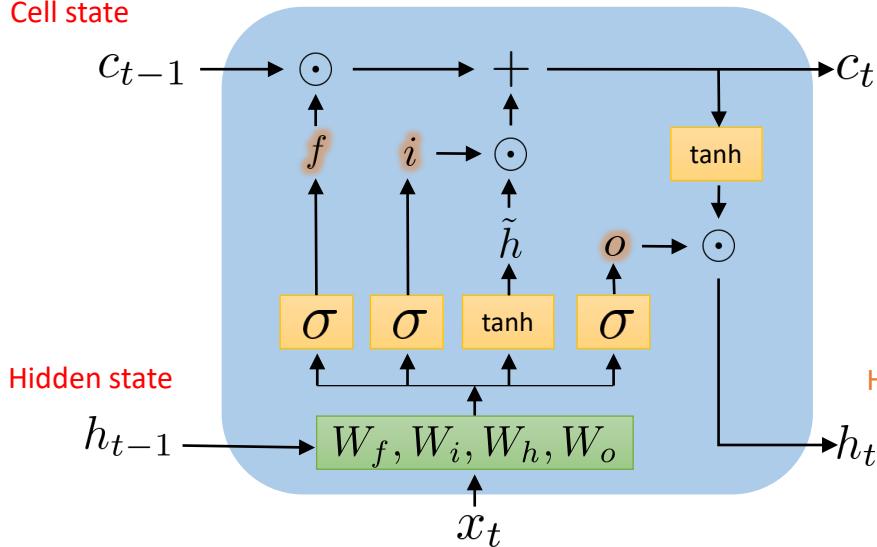
Cell state

$$c_t = f \odot c_{t-1} + i \odot \tilde{h}$$

$$h_t = o \odot \tanh(c_t)$$

Long Short-Term Memory (LSTM)

Cell state



Calculate output h_t

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

$$i = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

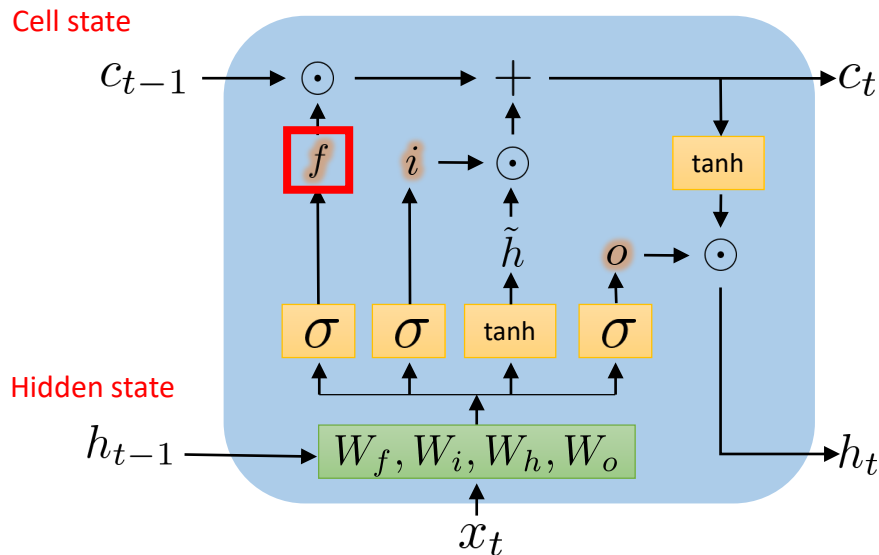
$$\tilde{h} = \tanh(W_h \cdot [h_{t-1}, x_t] + b_h)$$

$$o = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$c_t = f \odot c_{t-1} + i \odot \tilde{h}$$

$$h_t = o \odot \tanh(c_t)$$

Long Short-Term Memory (LSTM)



Prevent gradient vanishing
if forget gate f is open (>0).

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

$$i = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

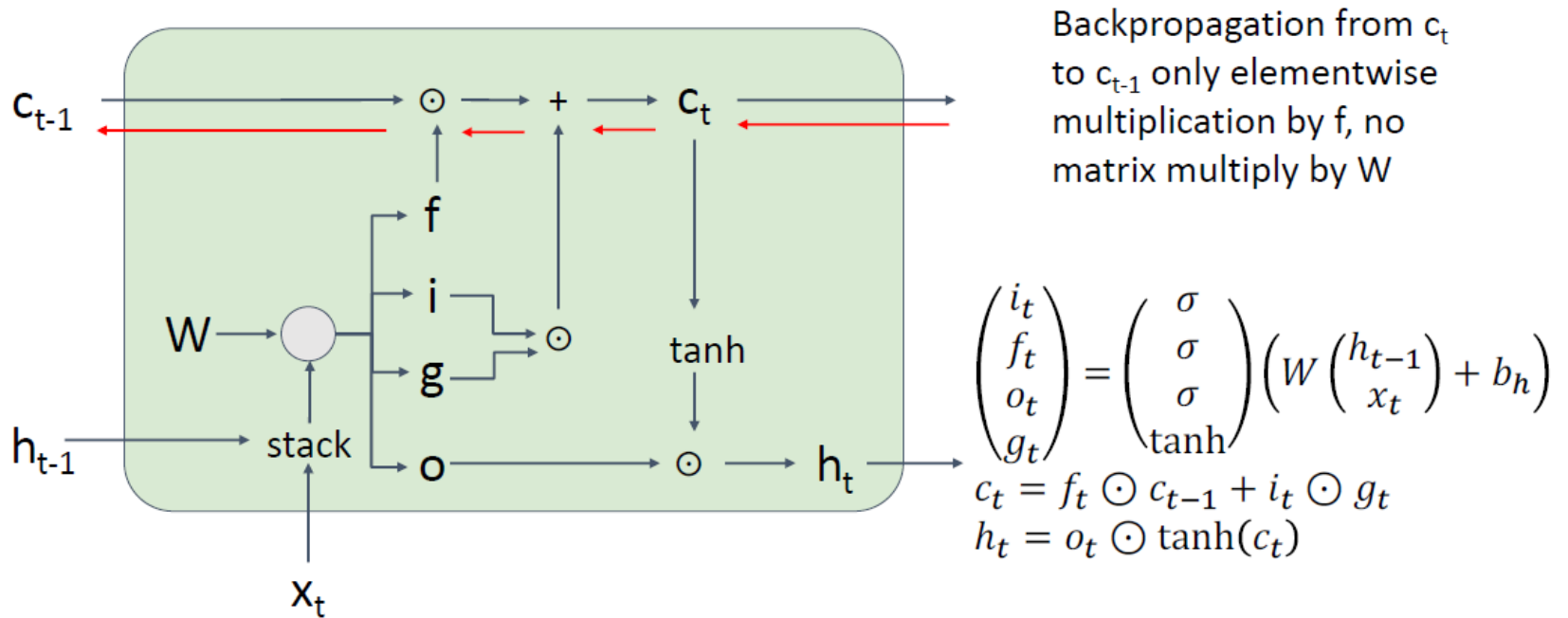
$$\tilde{h} = \tanh(W_h \cdot [h_{t-1}, x_t] + b_h)$$

$$o = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$c_t = f \odot c_{t-1} + i \odot \tilde{h}$$

$$h_t = o \odot \tanh(c_t)$$

LSTM: Gradient Flow

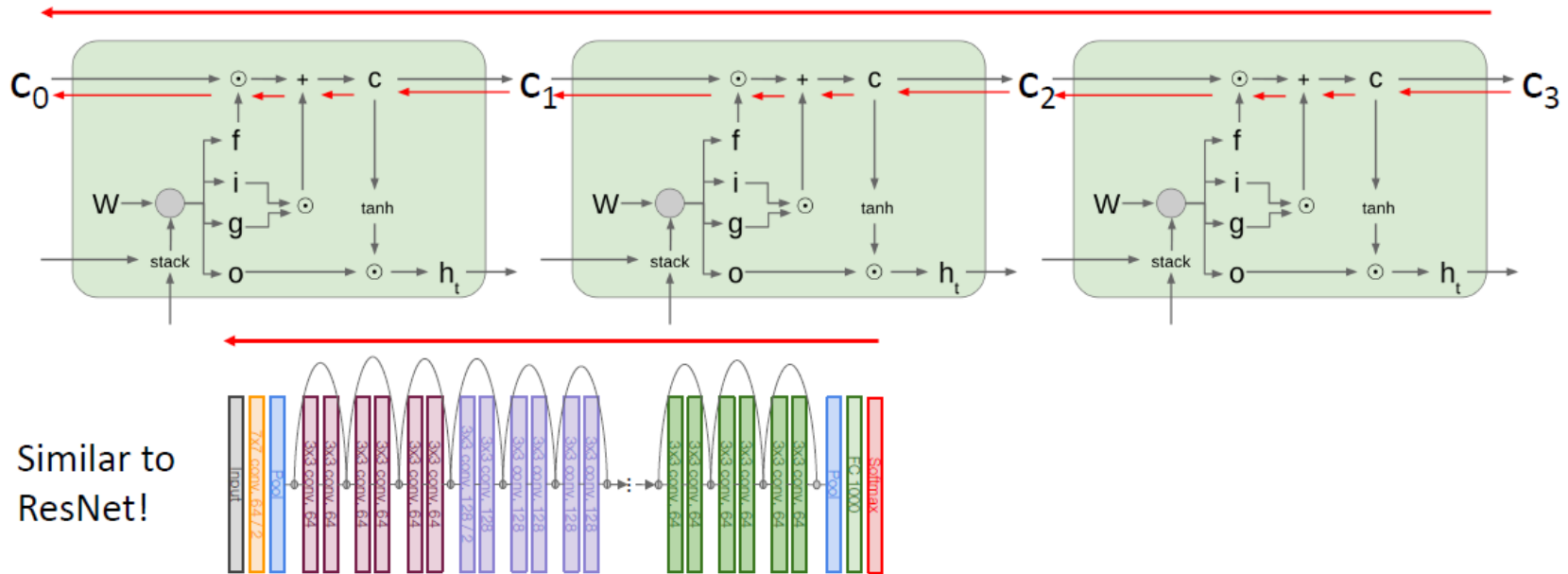


Remarks:

- Linear relationship between c_t and c_{t-1} instead of multiplication relationship of h_t and h_{t-1} in vanilla RNN

LSTM: Gradient Flow

Uninterrupted gradient flow!



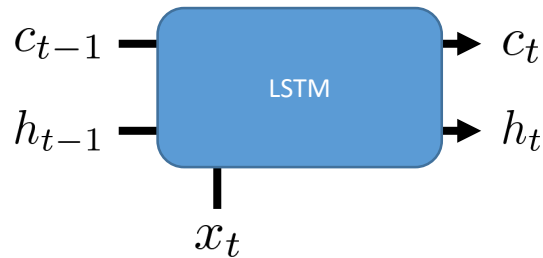
Remarks:

- Linear relationship between c_t and c_{t-1} instead of multiplication relationship of h_t and h_{t-1} in vanilla RNN
- Forget gate f provides shortcut connection across time steps

Variants of RNN

- Long Short-term Memory (LSTM) [*Hochreiter et al., 1997*]
 - Additional memory cell
 - Input / Forget / Output Gates
 - Handle gradient vanishing
 - Learn long-term dependencies
- Gated Recurrent Unit (GRU) [*Cho et al., EMNLP 2014*]
 - Similar to LSTM
 - No additional memory cell
 - Reset/Update Gates
 - Fewer parameters than LSTM
 - Comparable performance to LSTM [*Chung et al., NIPS Workshop 2014*]

LSTM vs. GRU



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

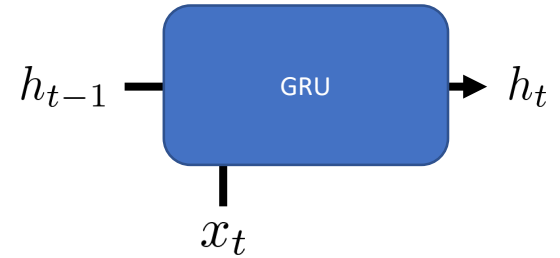
Input gate $i = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$

Input activation $\tilde{h} = \tanh (W_h \cdot [h_{t-1}, x_t] + b_h)$

Output gate $o = \sigma (W_o \cdot [h_{t-1}, x_t] + b_o)$

Cell state $c_t = f \odot c_{t-1} + i \odot \tilde{h}$

Hidden state $h_t = o \odot \tanh(c_t)$



Reset gate $r = \sigma (W_r \cdot [h_{t-1}, x_t] + b_r)$

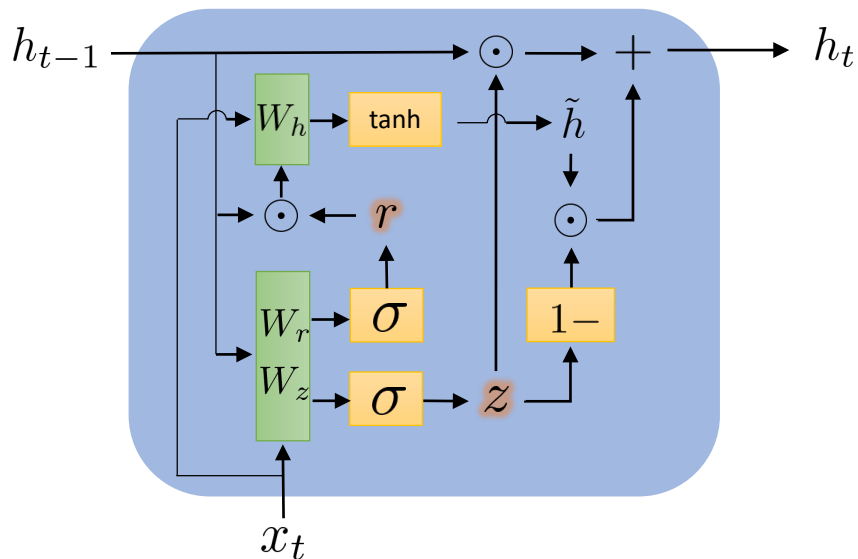
Update gate $z = \sigma (W_z \cdot [h_{t-1}, x_t] + b_z)$

Input activation $\tilde{h} = \tanh (W_h \cdot [(r \odot h_{t-1}), x_t] + b_h)$

Hidden state $h_t = z \odot h_{t-1} + (1 - z) \odot \tilde{h}$

Gated Recurrent Unit (GRU)

Hidden state

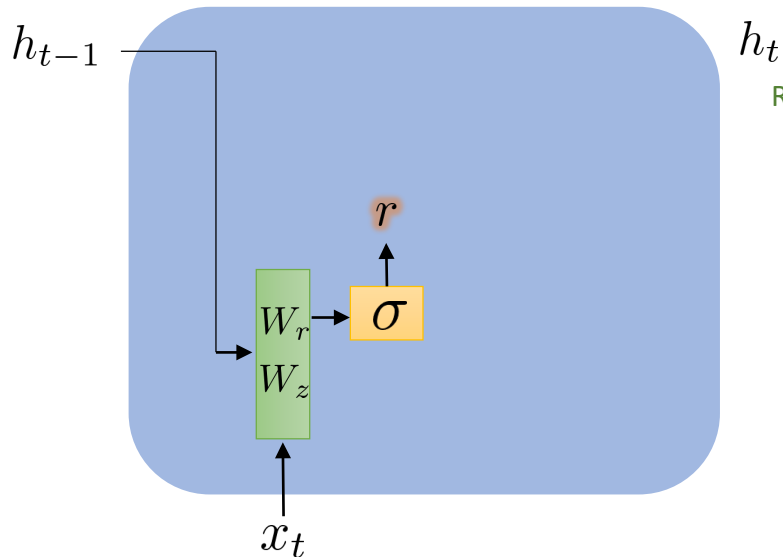


$$\begin{aligned} r &= \sigma(W_r \cdot [h_{t-1}, x_t] + b_r) \\ z &= \sigma(W_z \cdot [h_{t-1}, x_t] + b_z) \\ \tilde{h} &= \tanh(W_h \cdot [(r \odot h_{t-1}), x_t] + b_h) \end{aligned}$$

$$h_t = z \odot h_{t-1} + (1 - z) \odot \tilde{h}$$

Gated Recurrent Unit (GRU)

Hidden state



Reset gate

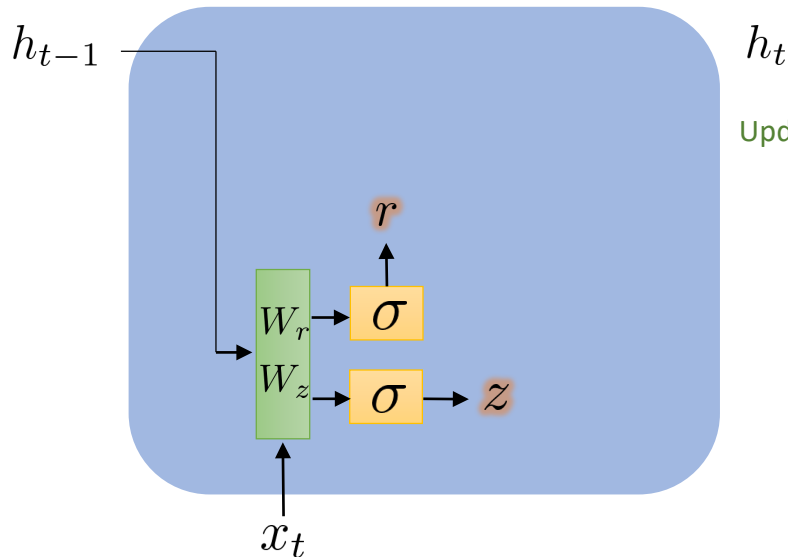
Calculate reset gate r

$$\begin{aligned} r &= \sigma (W_r \cdot [h_{t-1}, x_t] + b_r) \\ z &= \sigma (W_z \cdot [h_{t-1}, x_t] + b_z) \\ \tilde{h} &= \tanh (W_h \cdot [(r \odot h_{t-1}), x_t] + b_h) \end{aligned}$$

$$h_t = z \odot h_{t-1} + (1 - z) \odot \tilde{h}$$

Gated Recurrent Unit (GRU)

Hidden state



Calculate update gate z

$$r = \sigma (W_r \cdot [h_{t-1}, x_t] + b_r)$$

Update gate

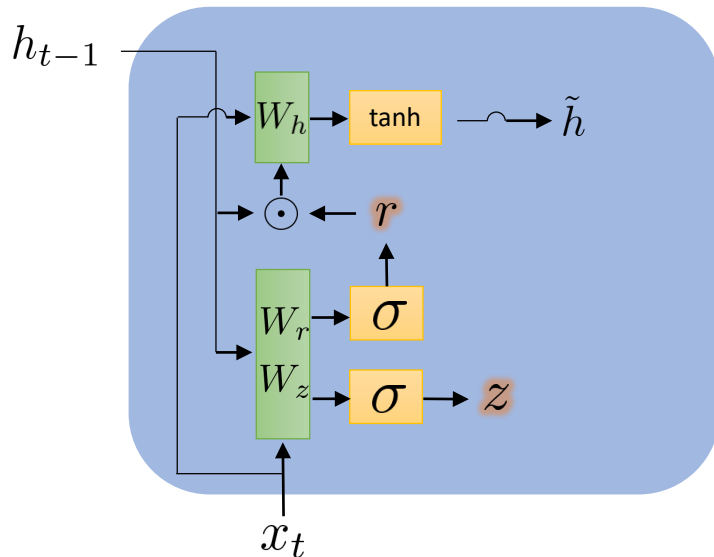
$$z = \sigma (W_z \cdot [h_{t-1}, x_t] + b_z)$$

$$\tilde{h} = \tanh (W_h \cdot [(r \odot h_{t-1}), x_t] + b_h)$$

$$h_t = z \odot h_{t-1} + (1 - z) \odot \tilde{h}$$

Gated Recurrent Unit (GRU)

Hidden state



h_t

Calculate input activation

$$r = \sigma(W_r \cdot [h_{t-1}, x_t] + b_r)$$

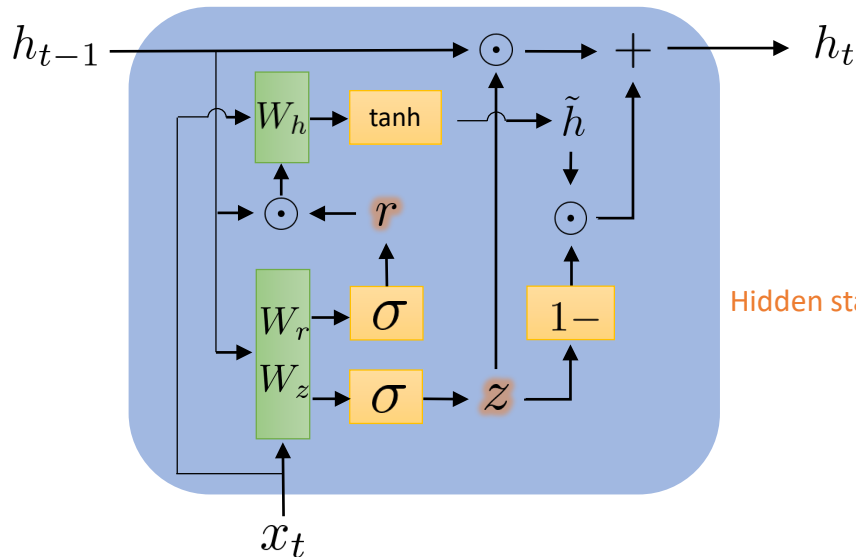
$$z = \sigma(W_z \cdot [h_{t-1}, x_t] + b_z)$$

Input activation $\tilde{h} = \tanh(W_h \cdot [(r \odot h_{t-1}), x_t] + b_h)$

$$h_t = z \odot h_{t-1} + (1 - z) \odot \tilde{h}$$

Gated Recurrent Unit (GRU)

Hidden state



Calculate output

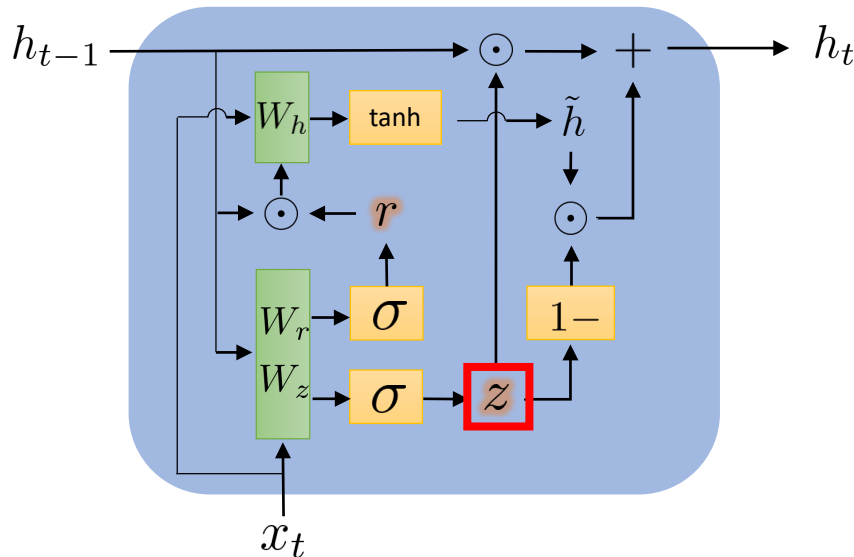
$$\begin{aligned} r &= \sigma(W_r \cdot [h_{t-1}, x_t] + b_r) \\ z &= \sigma(W_z \cdot [h_{t-1}, x_t] + b_z) \\ \tilde{h} &= \tanh(W_h \cdot [(r \odot h_{t-1}), x_t] + b_h) \end{aligned}$$

Hidden state

$$h_t = z \odot h_{t-1} + (1 - z) \odot \tilde{h}$$

Gated Recurrent Unit (GRU)

Hidden state



Prevent gradient vanishing if update gate z is open!

$$r = \sigma(W_r \cdot [h_{t-1}, x_t] + b_r)$$

$$z = \sigma(W_z \cdot [h_{t-1}, x_t] + b_z)$$

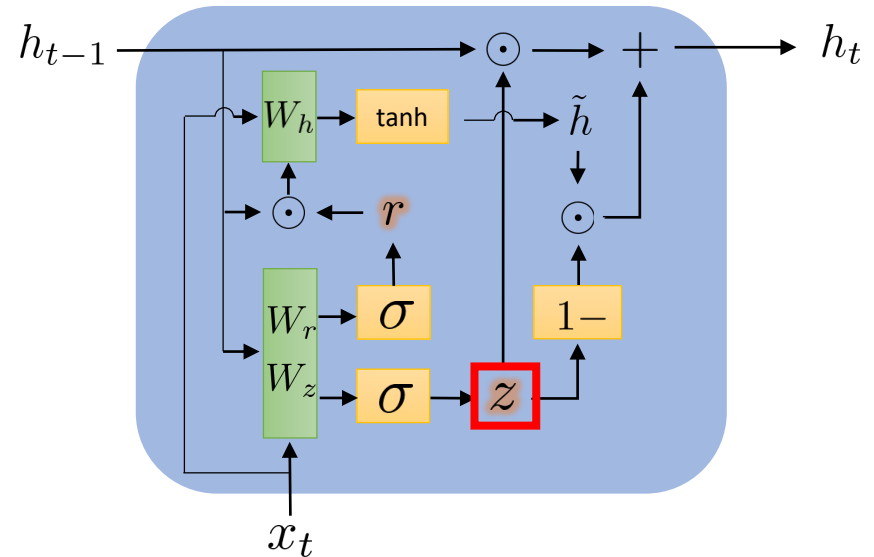
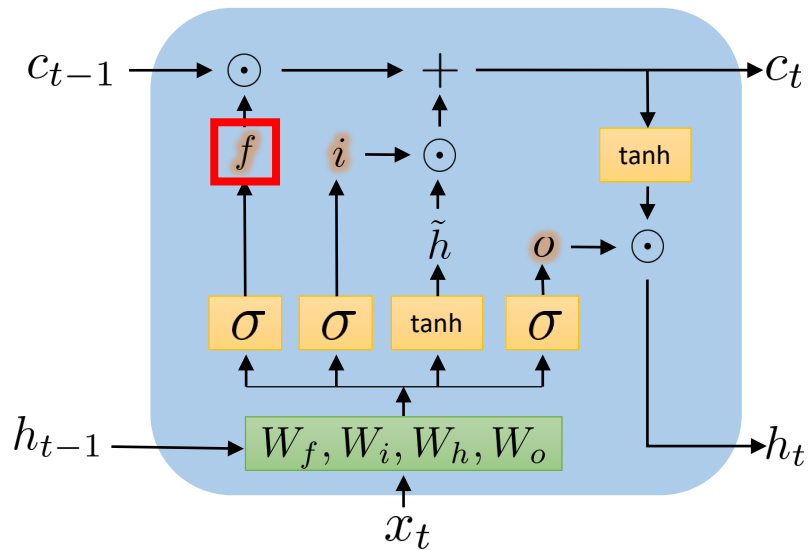
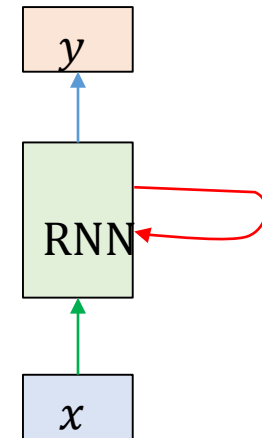
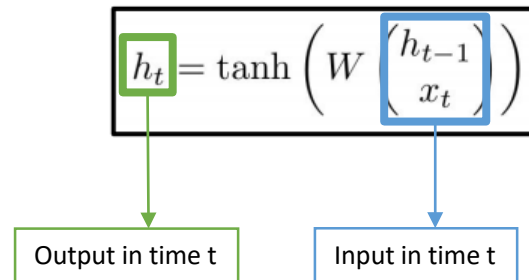
$$\tilde{h} = \tanh(W_h \cdot [(r \odot h_{t-1}), x_t] + b_h)$$

$$h_t = z \odot h_{t-1} + (1 - z) \odot \tilde{h}$$




Remarks:

- Update gate z provides shortcut connection across time steps
- Linear relationship between h_t and h_{t-1} instead of multiplication relationship of h_t and h_{t-1} in vanilla RNN

Vanilla RNN, LSTM, & GRU

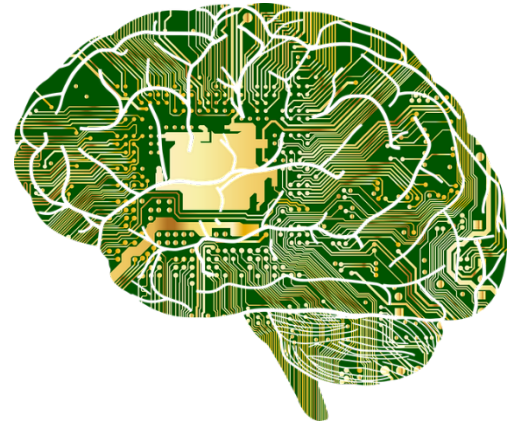


LSTM vs. GRU

	Vanilla RNN	LSTM	GRU
Cell state	X	O	O
Number of Gates	N/A	3	2
Parameters	Least	Most	Fewer
Gradient Vanishing / Exploding			

What to Be Covered Today...

- Transfer Learning
 - Visual Classification – Domain Adaptation
 - Visual Synthesis – Style Transfer
- Recurrent Neural Networks
 - From RNN to LSTM & GRU
 - Selected Models for Sequence-to-Sequence Learning
 - Attention in RNN

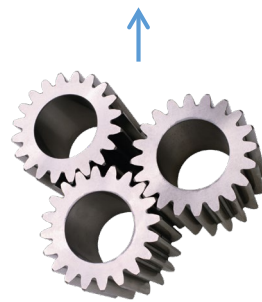


Sequence-to-Sequence Modeling

- Setting
 - An input sequence $\mathbf{X}_1, \dots, \mathbf{X}_N$
 - An output sequence $\mathbf{Y}_1, \dots, \mathbf{Y}_M$
 - Generally $N \neq M$, i.e., **no synchrony** between \mathbf{X} and \mathbf{Y}
- Examples
 - Speech recognition: speech goes in, and a word sequence comes out
 - Machine translation: word sequence goes in, and another comes out
 - Video captioning: video frames goes in, word sequence comes out



深度學習好棒棒!

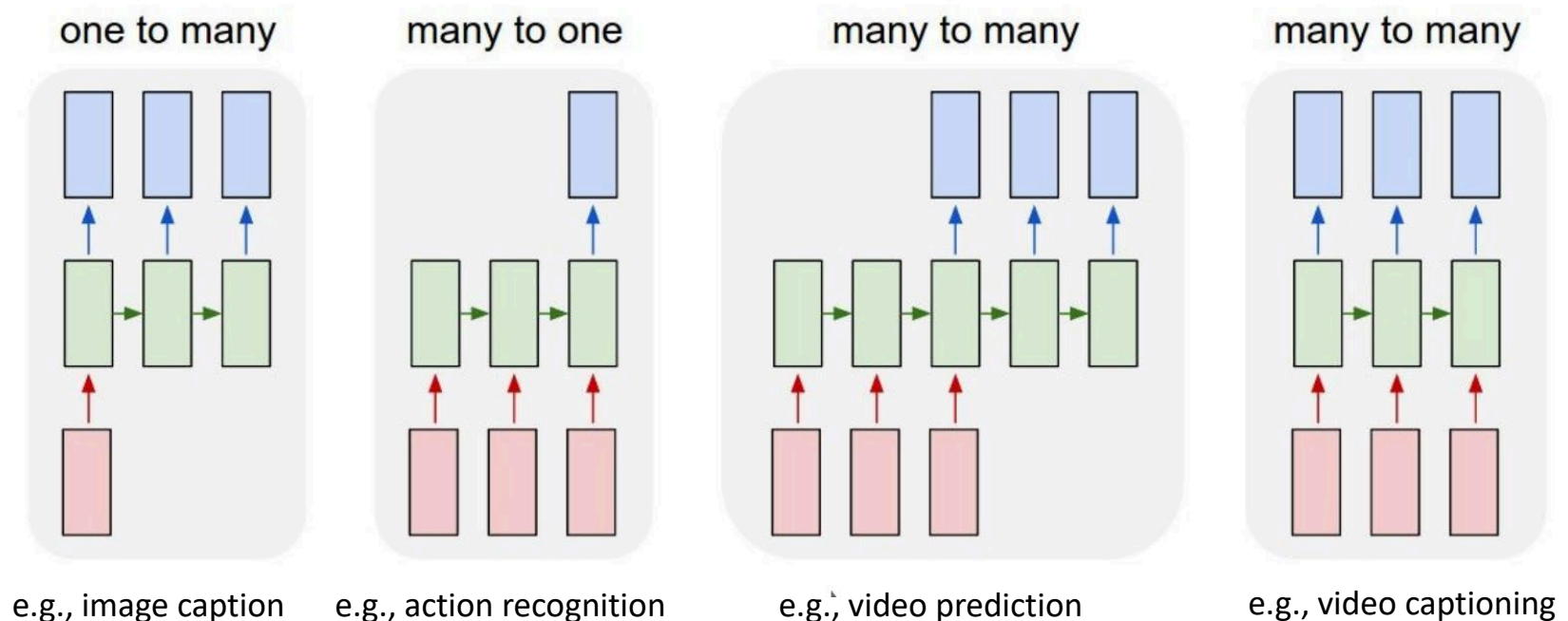


Deep learning rocks!



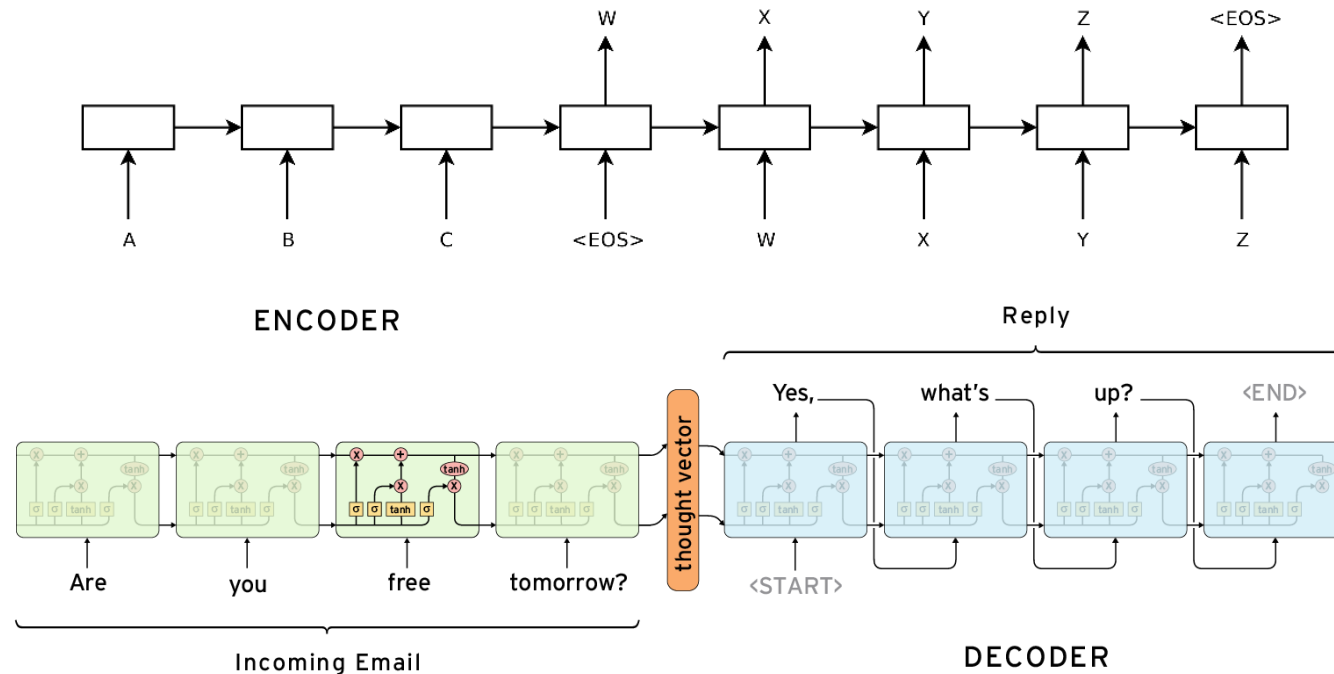
S-to-S Models with Alignment

- The input and output sequences happen in the same order
 - The input/output sequences may be asynchronous.
 - E.g., speech recognition or video captioning, in which the input sequence corresponds to the phoneme/caption sequence out.
 - Recall that...



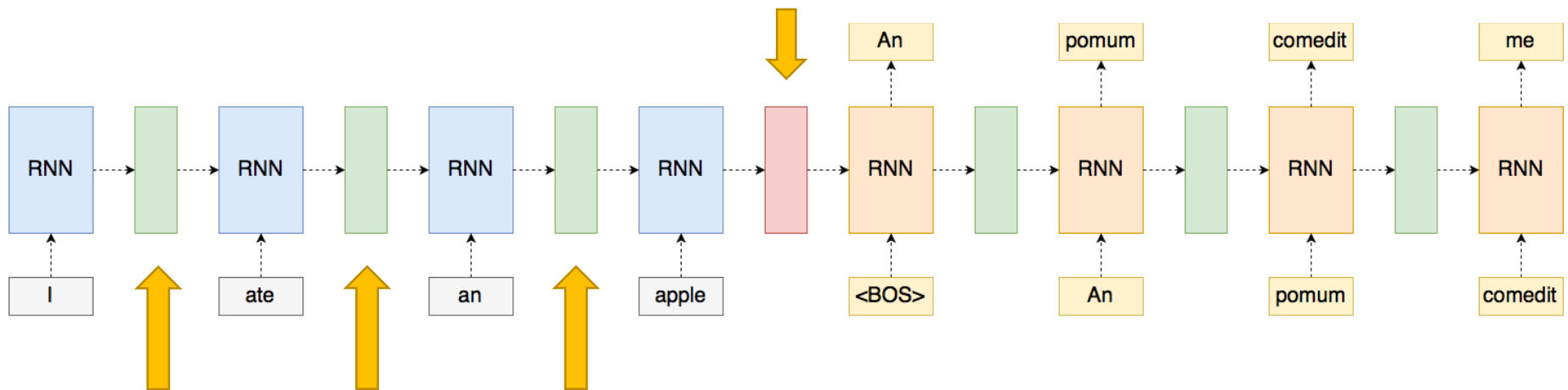
Sequence-to-Sequence Modeling (cont'd)

- Original model proposed in NIPS 2014
 - An encoder-decoder model



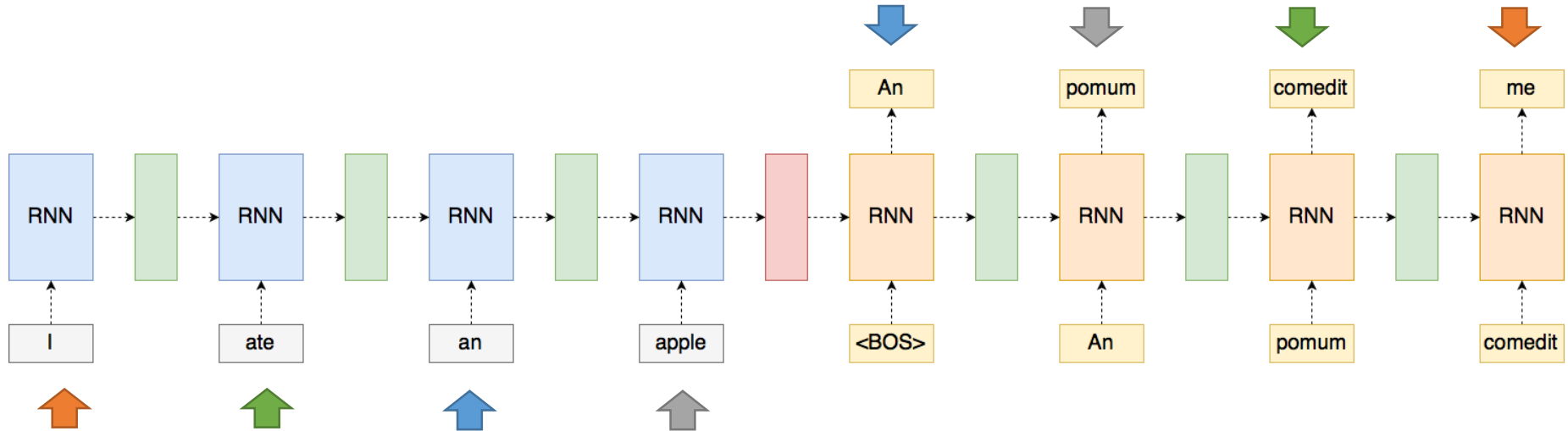
What's the Potential Problem?

- Each hidden state vector extracts/carries information across time steps (some might be diluted downstream).
- However, information of the entire input sequence is embedded into a **single hidden state** vector.



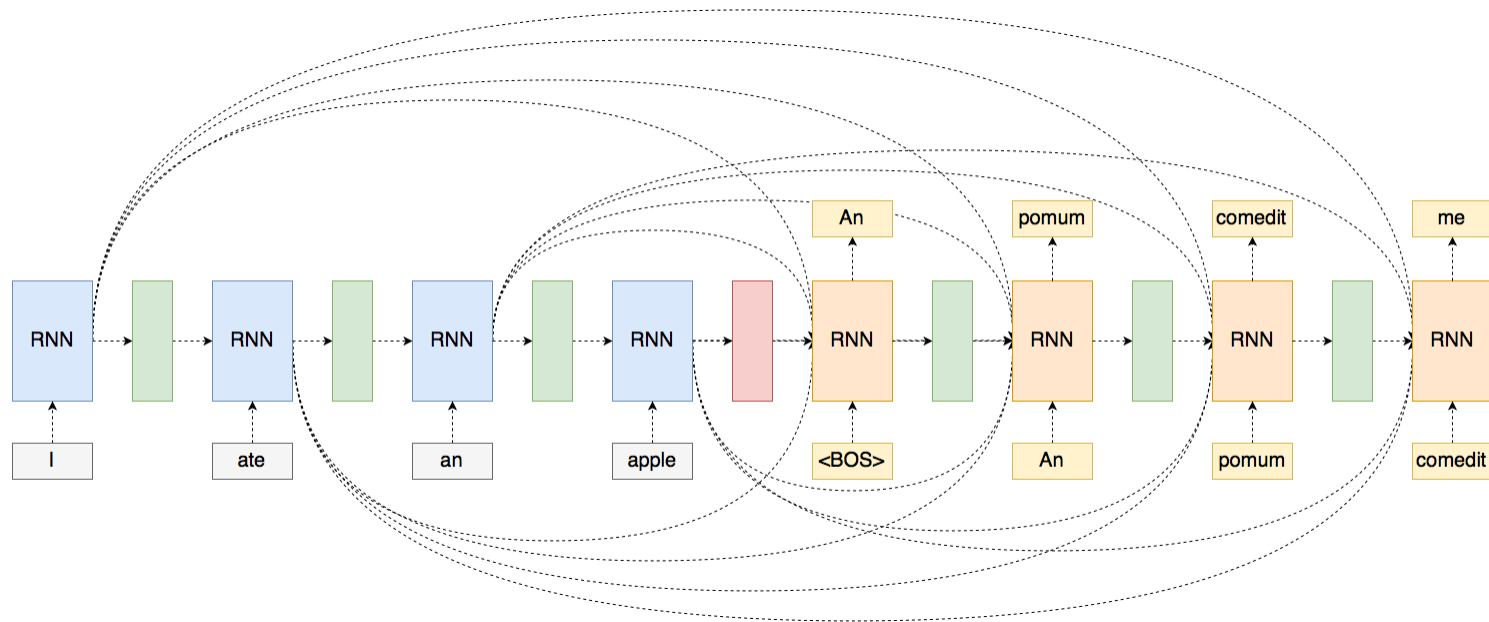
What's the Potential Problem? (cont'd)

- Outputs at different time steps have particular meanings.
- However, **synchrony** between **input** and **output seqs** is **not** required.



What's the Potential Problem? (cont'd)

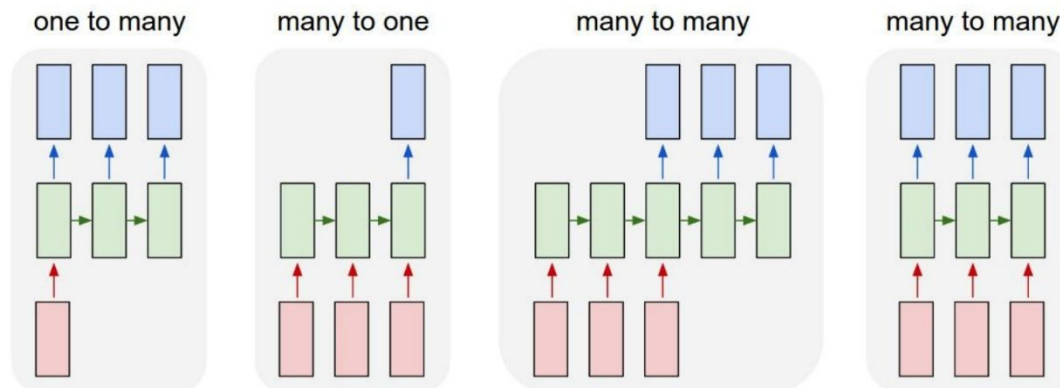
- Connecting every hidden state between encoder and decoder?



- Infeasible!
 - Both inputs and outputs are with varying sizes.
 - Overparameterized
 - Possible solution: attention (will cover next week)

Recent Advances of GAN-based Models for Video-Based Applications

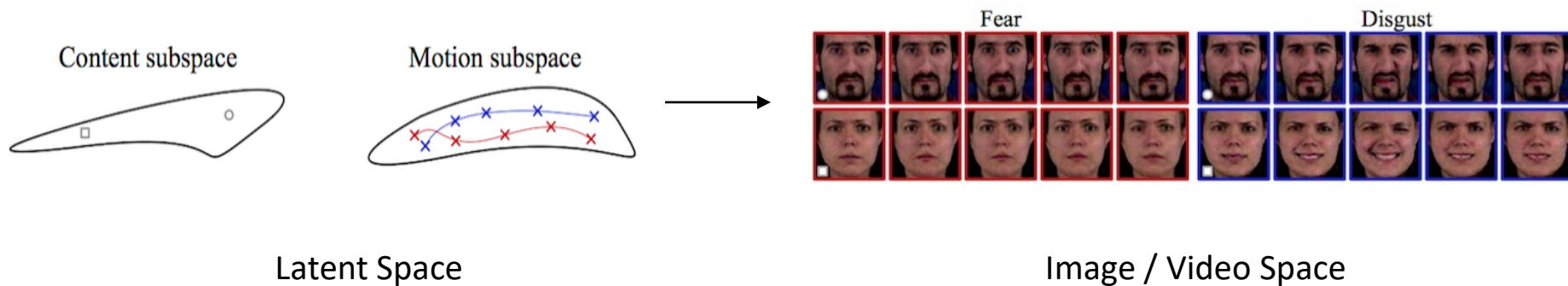
- Video Generation
 - MoCoGAN: Decomposing Motion and Content for Video Generation (CVPR'18)
- Video Prediction
 - Deterministic
 - Unsupervised Learning of Video Representations using LSTMs (ICML'15)
 - Decomposing Motion and Content for Natural Video Sequence Prediction (ICLR'17)
 - Learning to generate long-term future via hierarchical prediction (ICML'17)
 - Stochastic (if time permits)
 - Stochastic Video Generation with a Learned Prior (Denton et al., ArXiv'18)



Video Generation

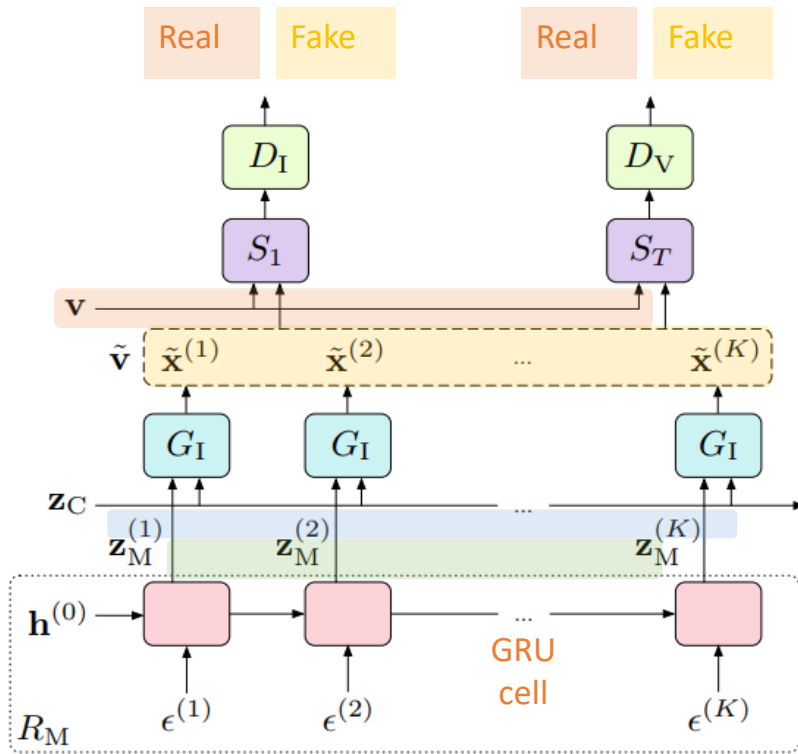
- Learning a latent space to describe image/video data
- Input: latent representation
- Output: sequence of images/frames (i.e., video)

E.g., Latent space for content-motion decomposition



Video Generation

- MoCoGAN: Decomposing Motion and Content for Video Generation (Tulyakov et al, CVPR'18)



Generated Video: $\tilde{\mathbf{v}} = \left[G_I \left(\begin{bmatrix} z_C \\ R_M(1) \end{bmatrix} \right), \dots, G_I \left(\begin{bmatrix} z_C \\ R_M(\kappa) \end{bmatrix} \right) \right]$

Real Video:

\mathbf{v}

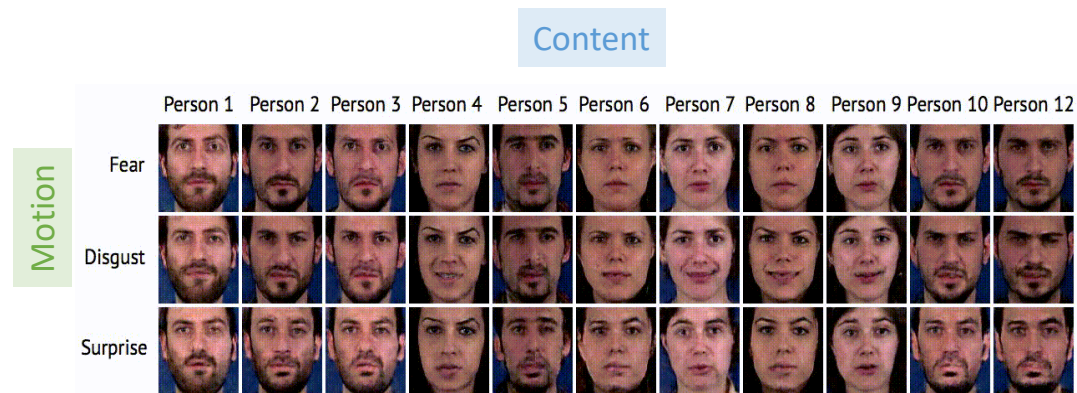
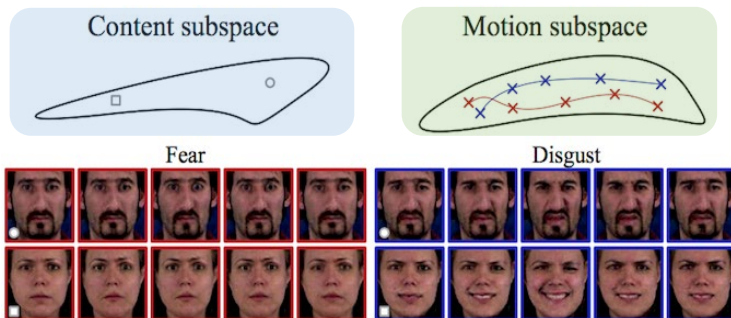
Objective Function:

$$\begin{aligned} & \max_{G_I, R_M} \min_{D_I, D_V} \mathcal{F}_V(D_I, D_V, G_I, R_M) \\ &= \mathbb{E}_{\mathbf{v}}[-\log D_I(S_1(\mathbf{v}))] + \mathbb{E}_{\tilde{\mathbf{v}}}[-\log(1 - D_I(S_1(\tilde{\mathbf{v}})))] + \\ & \quad \mathbb{E}_{\mathbf{v}}[-\log D_V(S_T(\mathbf{v}))] + \mathbb{E}_{\tilde{\mathbf{v}}}[-\log(1 - D_V(S_T(\tilde{\mathbf{v}})))] \end{aligned}$$

Example Results

- MoCoGAN: Decomposing Motion and Content for Video Generation (Tulyakov et al, CVPR'18)

Results:

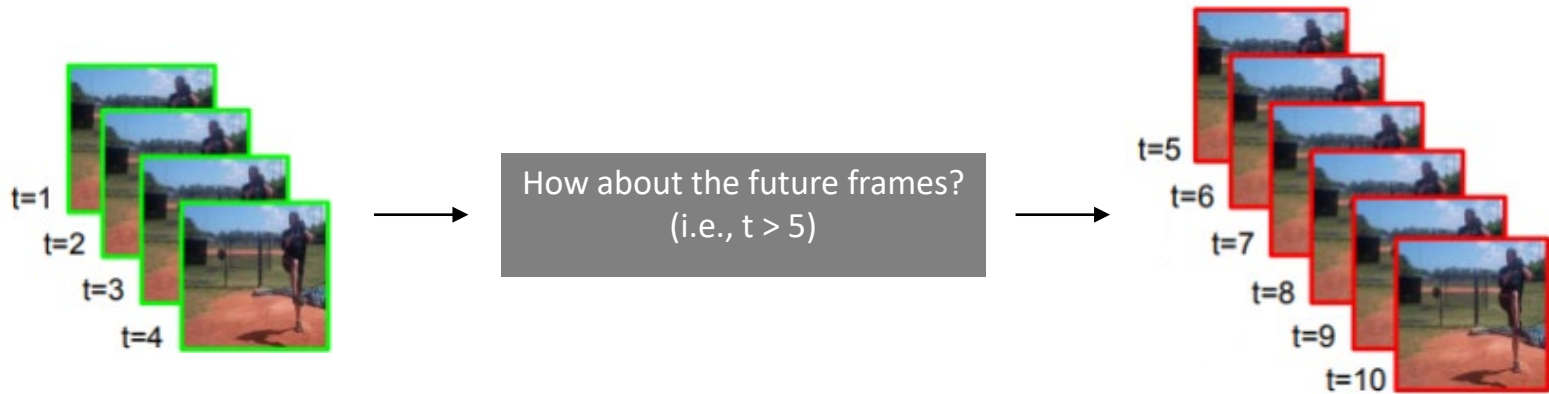


Recent Advances of Attention models in Video-Based Applications

- Video Generation
 - MoCoGAN: Decomposing Motion and Content for Video Generation (CVPR'18)
- Video Prediction
 - Deterministic
 - Unsupervised Learning of Video Representations using LSTMs (ICML'15)
 - Decomposing Motion and Content for Natural Video Sequence Prediction (ICLR'17)
 - Learning to generate long-term future via hierarchical prediction (ICML'17)
 - Stochastic
 - Stochastic Video Generation with a Learned Prior (Denton et al., ArXiv'18)

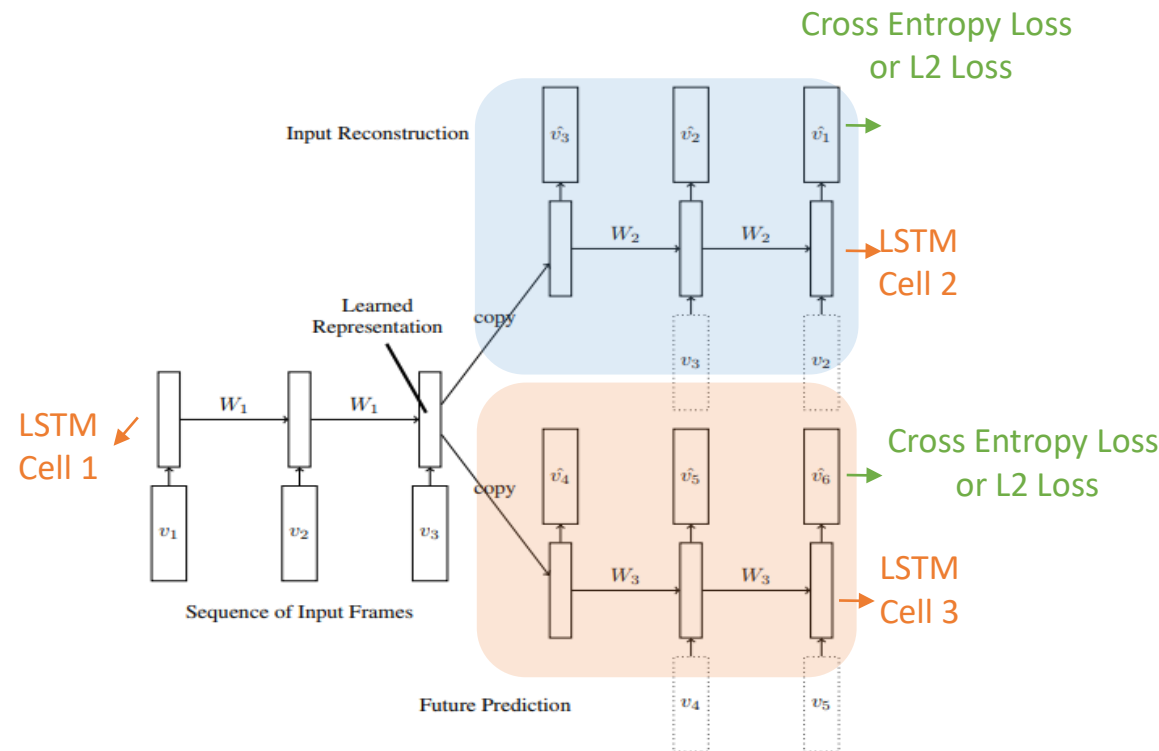
Video Prediction



- Input: A few **known** frames
- Output: **Unknown** future frames



Video Prediction - Deterministic

- Unsupervised Learning of Video Representations using LSTMs
(Srivastava et al., ICML'15)

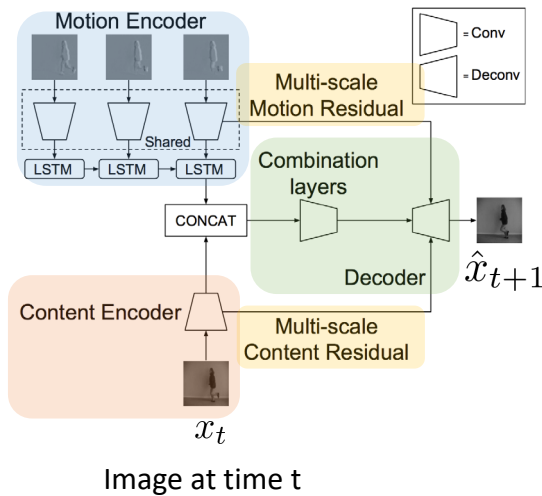


- LSTM Encoder-Decoder model
- Two tasks: Reconstruction & Prediction
- Results (L: ref video, R: output video)
 - Bouncing (Moving) MNIST
 - 
 - Video patches of UCF-101
 - 

Video Prediction - Deterministic

- Decomposing Motion and Content for Natural Video Sequence Prediction (Villegas et al., ICLR'17)

History of image differences



Hidden state
(motion)

Motion Encoder

$$[\mathbf{d}_t, \mathbf{c}_t] = f^{\text{dyn}}(\mathbf{x}_t - \mathbf{x}_{t-1}, \mathbf{d}_{t-1}, \mathbf{c}_{t-1})$$

Content Encoder

$$\mathbf{s}_t = f^{\text{cont}}(\mathbf{x}_t)$$

Multi-scale Motion & Content Residual

$$\mathbf{r}_t^l = f^{\text{res}}([\mathbf{s}_t^l, \mathbf{d}_t^l])^l$$

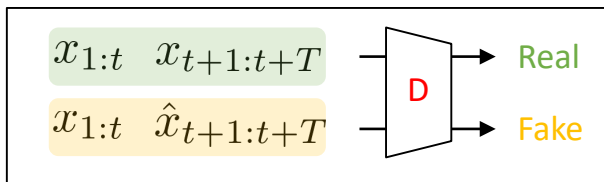
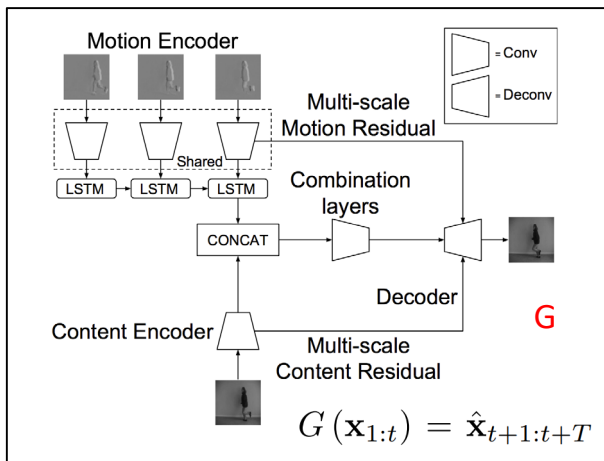
Combination Layers and Decoder

$$\mathbf{f}_t = g^{\text{comb}}([\mathbf{d}_t, \mathbf{s}_t]) \quad \hat{\mathbf{x}}_{t+1} = g^{\text{dec}}(\mathbf{f}_t, \mathbf{r}_t)$$

→ Recurrently generate $\hat{x}_{t+2}, \hat{x}_{t+3}, \dots, \hat{x}_{t+T}$

→ $G(\mathbf{x}_{1:t}) = \hat{\mathbf{x}}_{t+1:t+T}$

- Decomposing Motion and Content for Natural Video Sequence Prediction (Villegas et al., ICLR'17)



Objective Function: Adversarial Training -> alternate btw min L & L_{disc}

Update Image Generation Network (G)

$$\mathcal{L} = \alpha \mathcal{L}_{\text{img}} + \beta \mathcal{L}_{\text{GAN}}$$

Pixel Value similarity ↑ Gradient of Pixel Value similarity ↑

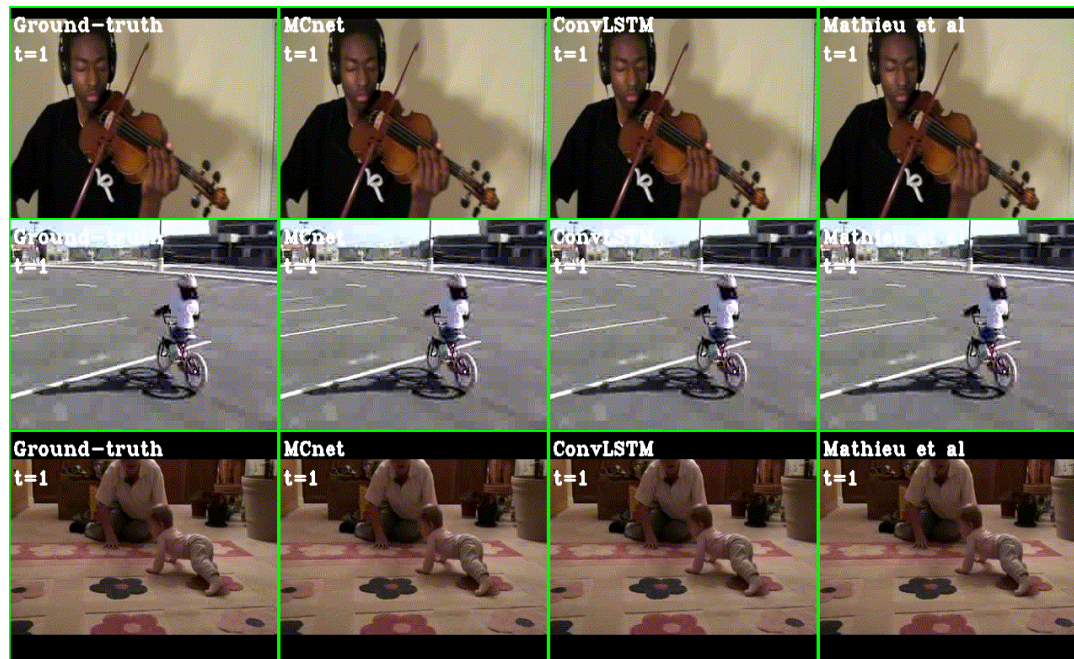
$$\mathcal{L}_{\text{img}} = \mathcal{L}_p(\mathbf{x}_{t+k}, \hat{\mathbf{x}}_{t+k}) + \mathcal{L}_{gdl}(\mathbf{x}_{t+k}, \hat{\mathbf{x}}_{t+k})$$

$$\mathcal{L}_{\text{GAN}} = -\log D([\mathbf{x}_{1:t}, G(\mathbf{x}_{1:t})])$$

Update Discriminator (D)

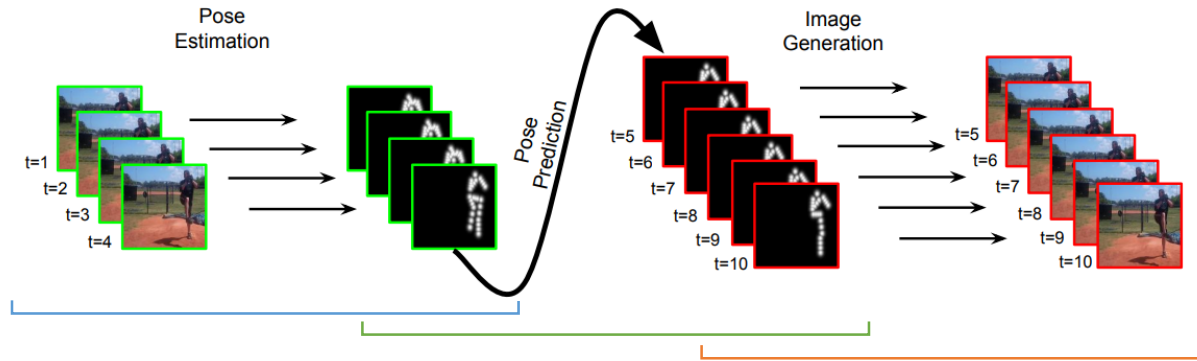
$$\mathcal{L}_{\text{disc}} = -\log D([\mathbf{x}_{1:t}, \mathbf{x}_{t+1:t+T}]) - \log(1 - D([\mathbf{x}_{1:t}, G(\mathbf{x}_{1:t})]))$$

- Example Results



Video Prediction - Deterministic

- Learning to generate long-term future via hierarchical prediction (Villegas et al., ICML'17)



Stage 1:

Pose Estimation

Hourglass network (Newell et al., ECCV'16)

Stage 2:

Future Pose Prediction

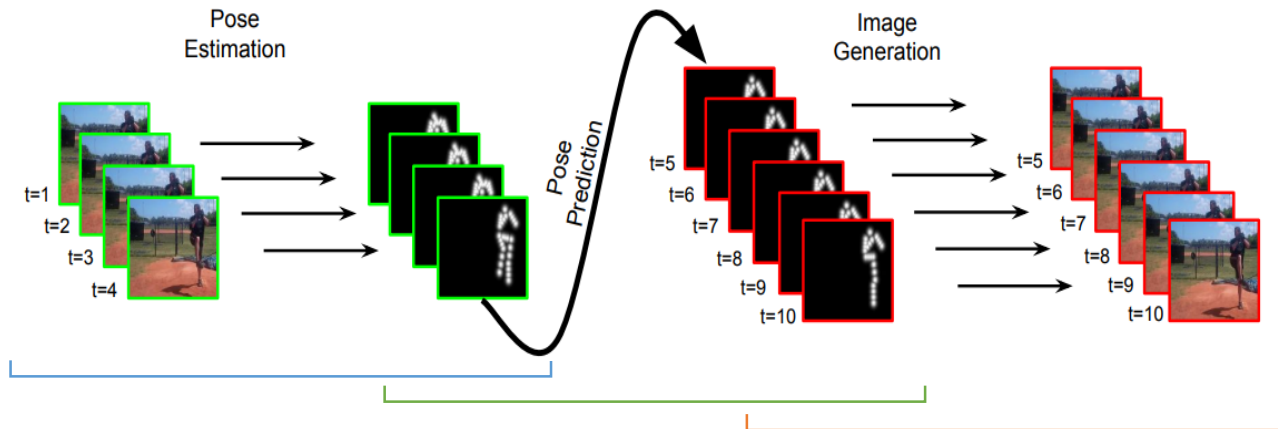
Sequence-to-Sequence model on high-level structure

Stage 3:

Image Generation

Visual-Structure Analogy

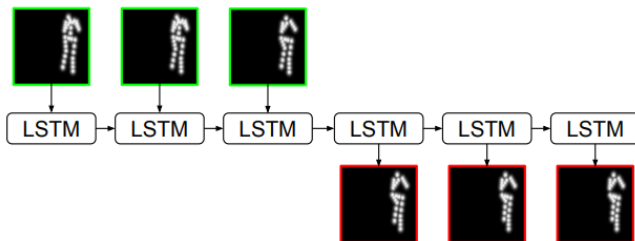
- Learning to generate long-term future via hierarchical prediction (Villegas et al., ICML'17)



Step 2:

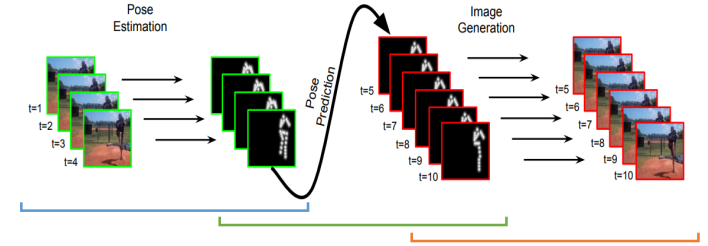
Future Pose Prediction

Sequence-to-Sequence model on high-level structure



Objective Function:

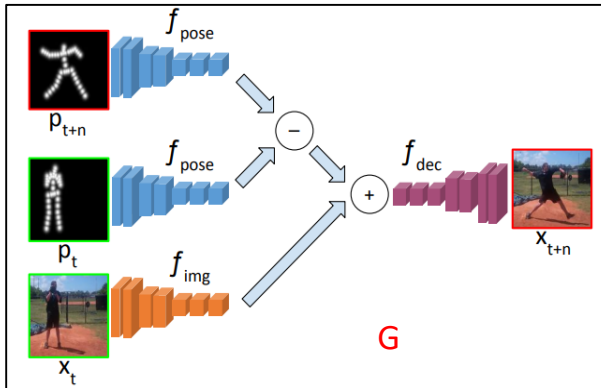
$$\mathcal{L}_{\text{pose}} = \frac{1}{TL} \sum_{t=1}^T \sum_{l=1}^L \mathbb{1}_{\{m_{k+t}^l=1\}} \|\hat{\mathbf{p}}_{k+t}^l - \mathbf{p}_{k+t}^l\|_2^2$$



- Learning to generate long-term future via hierarchical prediction (Villegas et al., ICML'17)

Step 3:

Image Generation



Visual-Structure Analogy

Objective Function:

Adversarial Training -> alternately minimize L & L_{Disc}

Update Image Generation Network (G)

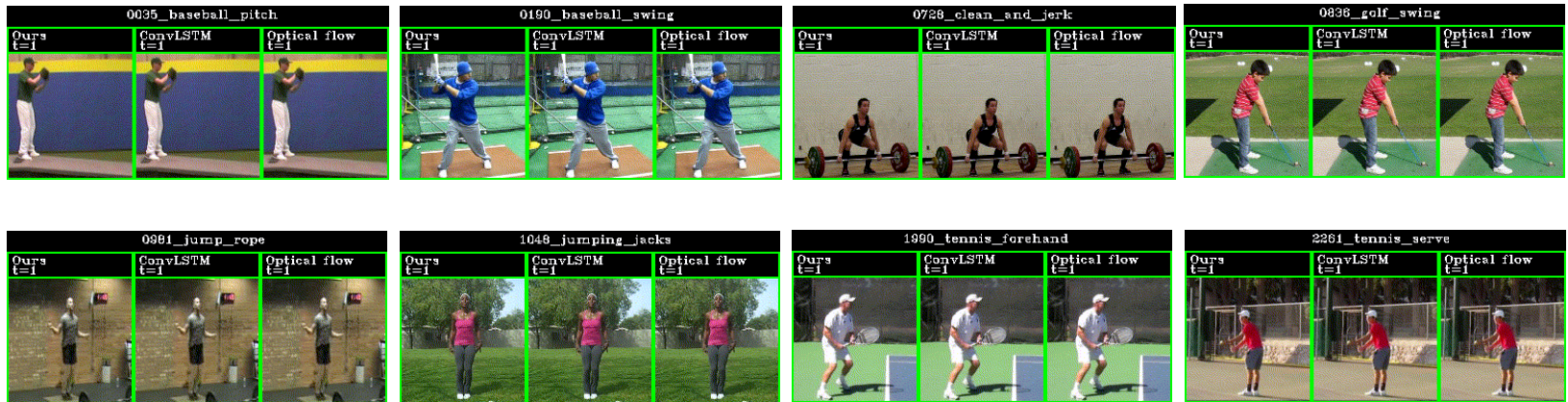
$$\begin{aligned}\mathcal{L} &= \mathcal{L}_{\text{img}} + \mathcal{L}_{\text{feat}} + \mathcal{L}_{\text{Gen}} \\ \mathcal{L}_{\text{img}} &= \|\mathbf{x}_{t+n} - \hat{\mathbf{x}}_{t+n}\|_2^2 \\ \mathcal{L}_{\text{feat}} &= \|C_1(\mathbf{x}_{t+n}) - C_1(\hat{\mathbf{x}}_{t+n})\|_2^2 \\ &\quad + \|C_2(\mathbf{x}_{t+n}) - C_2(\hat{\mathbf{x}}_{t+n})\|_2^2 \\ \mathcal{L}_{\text{Gen}} &= -\log D([\mathbf{p}_{t+n}, \hat{\mathbf{x}}_{t+n}])\end{aligned}$$

Update Discriminator (D)

$$\begin{aligned}\mathcal{L}_{\text{Disc}} &= -\log D([\mathbf{p}_{t+n}, \mathbf{x}_{t+n}]) \\ &\quad - 0.5 \log(1 - D([\mathbf{p}_{t+n}, \hat{\mathbf{x}}_{t+n}])) \\ &\quad - 0.5 \log(1 - D([\mathbf{p}_{t+n}, \mathbf{x}_t])),\end{aligned}$$

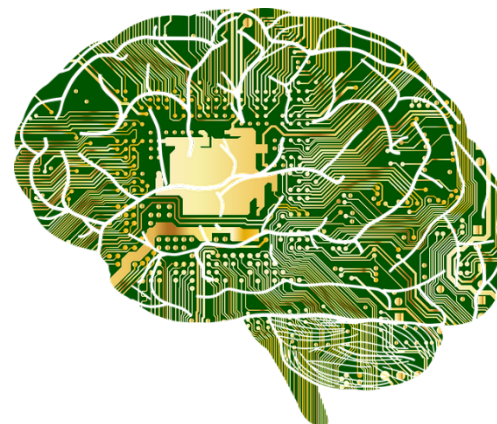
- Example Results

Results on Penn Action Dataset:



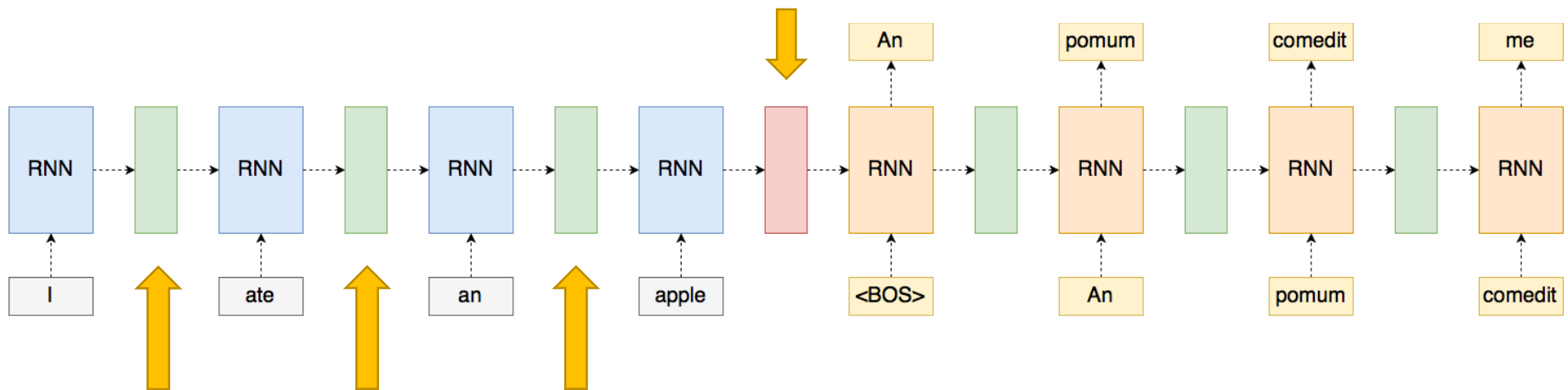
What to Be Covered Today...

- Transfer Learning
 - Visual Classification – Domain Adaptation
 - Visual Synthesis – Style Transfer
- Recurrent Neural Networks
 - From RNN to LSTM & GRU
 - Selected Models for Sequence-to-Sequence Learning
 - Attention in RNN



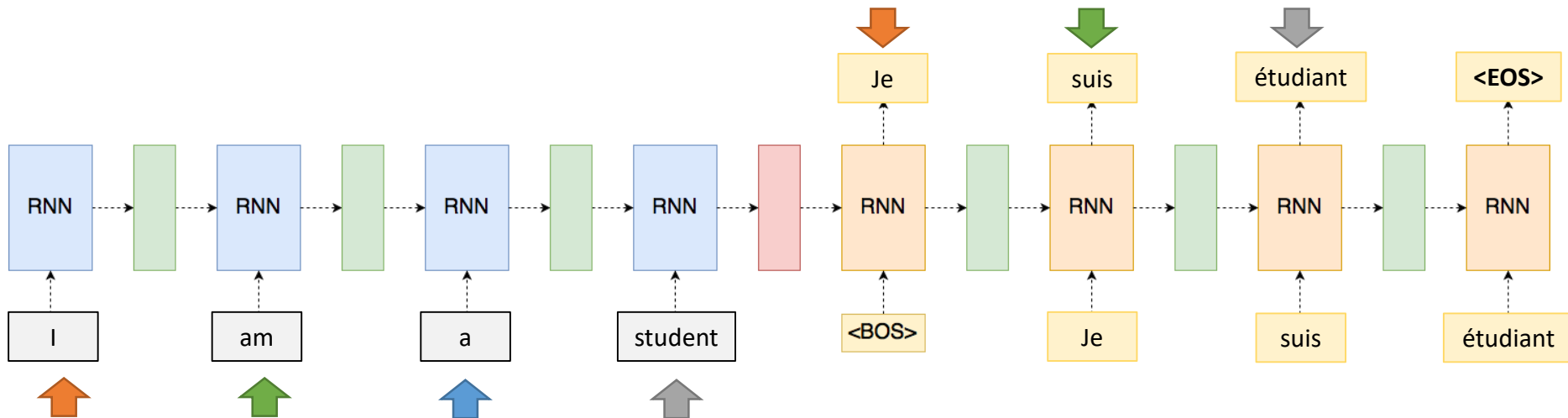
What's the Potential Problem of RNN?

- Each hidden state vector extracts/carries information across time steps (some might be diluted downstream).
- However, information of the entire input sequence is embedded into a **single hidden state** vector.



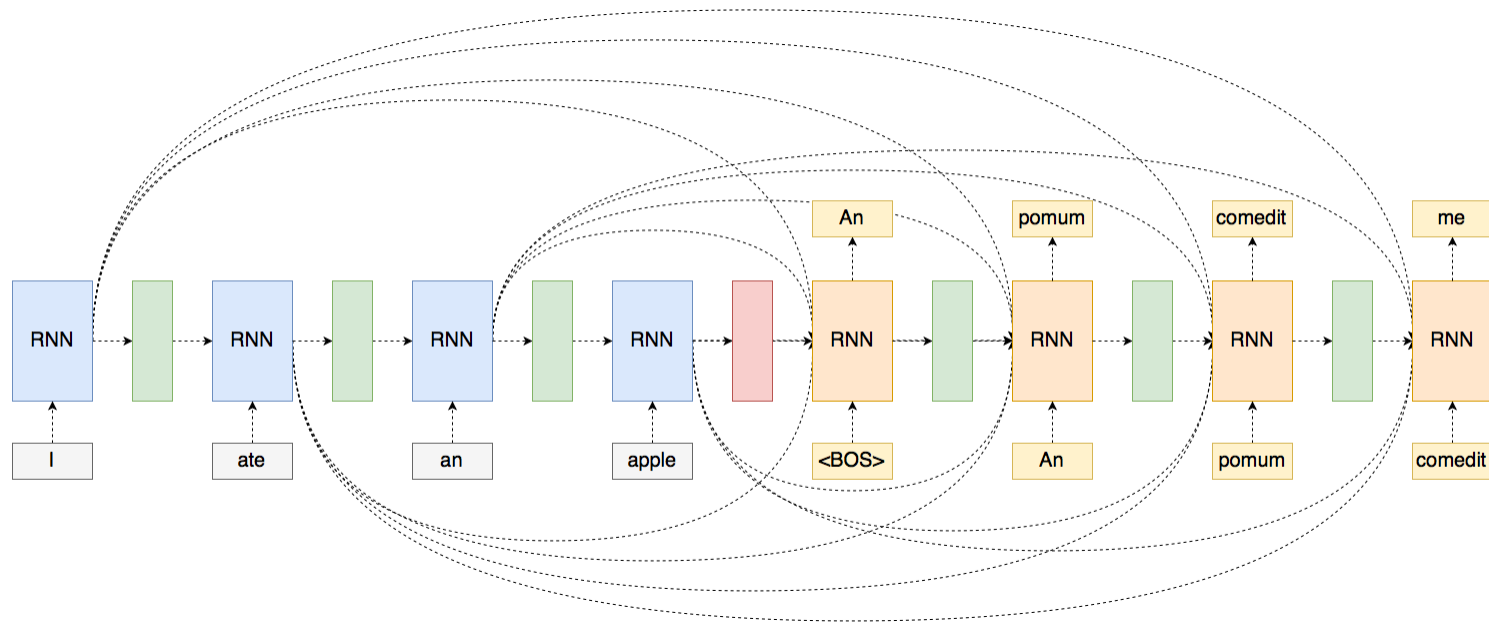
What's the Potential Problem? (cont'd)

- Outputs at different time steps have particular meanings.
- However, **synchrony** between **input** and **output seqs** is **not** required.



RNN with Attention is Good, But..

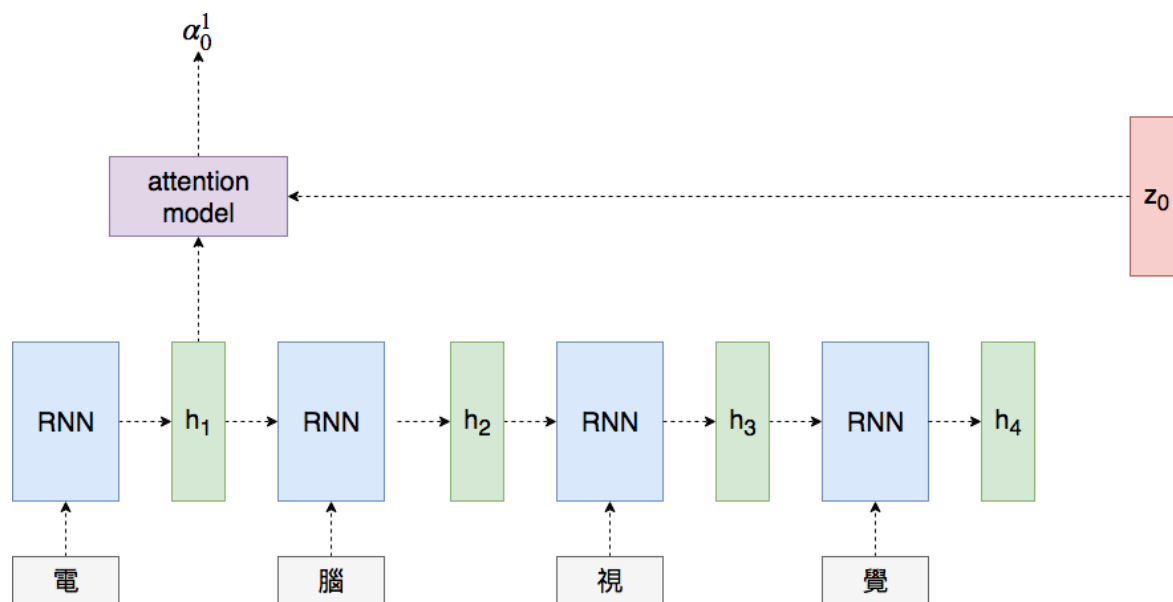
- Connecting every hidden state between encoder and decoder?



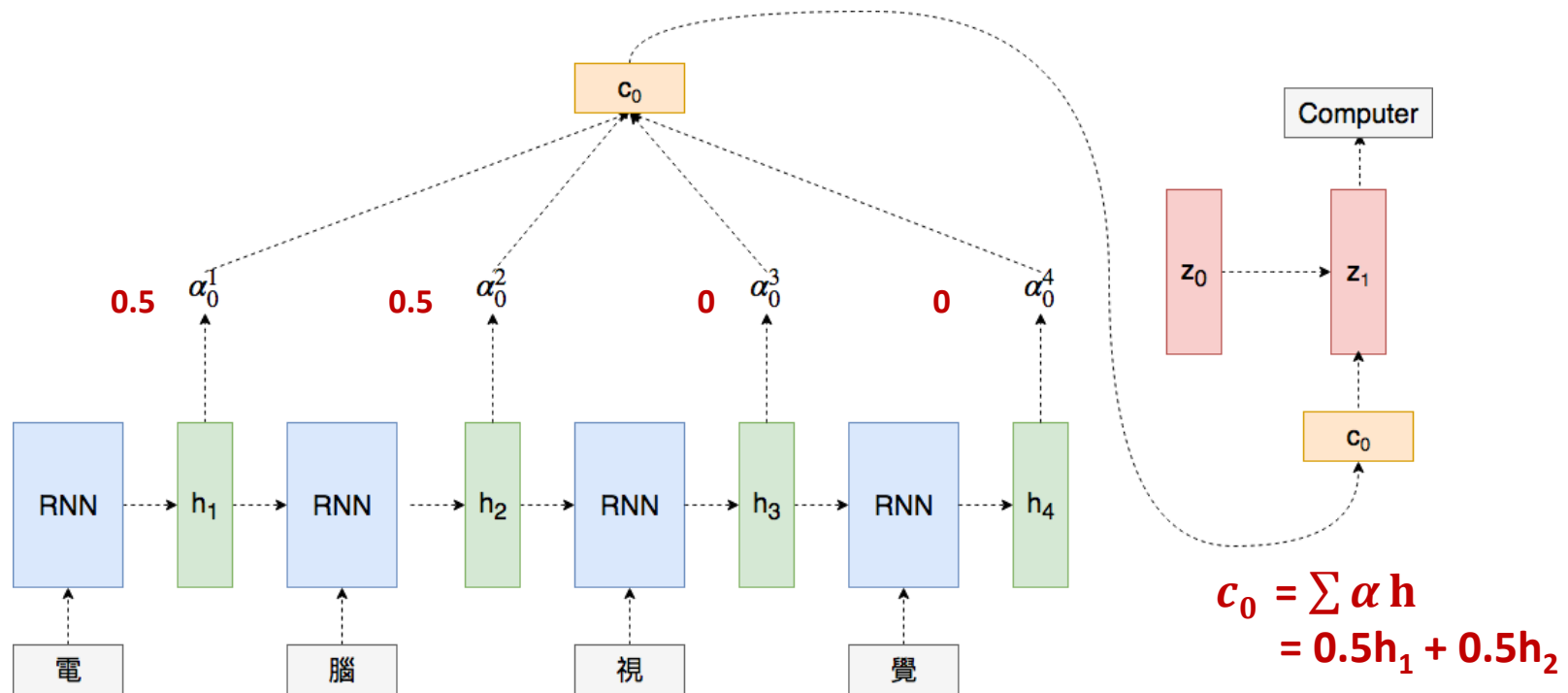
- Infeasible!
 - Both inputs and outputs are with varying sizes.
 - Overparameterized

Solution Ver. 1: Attention Model

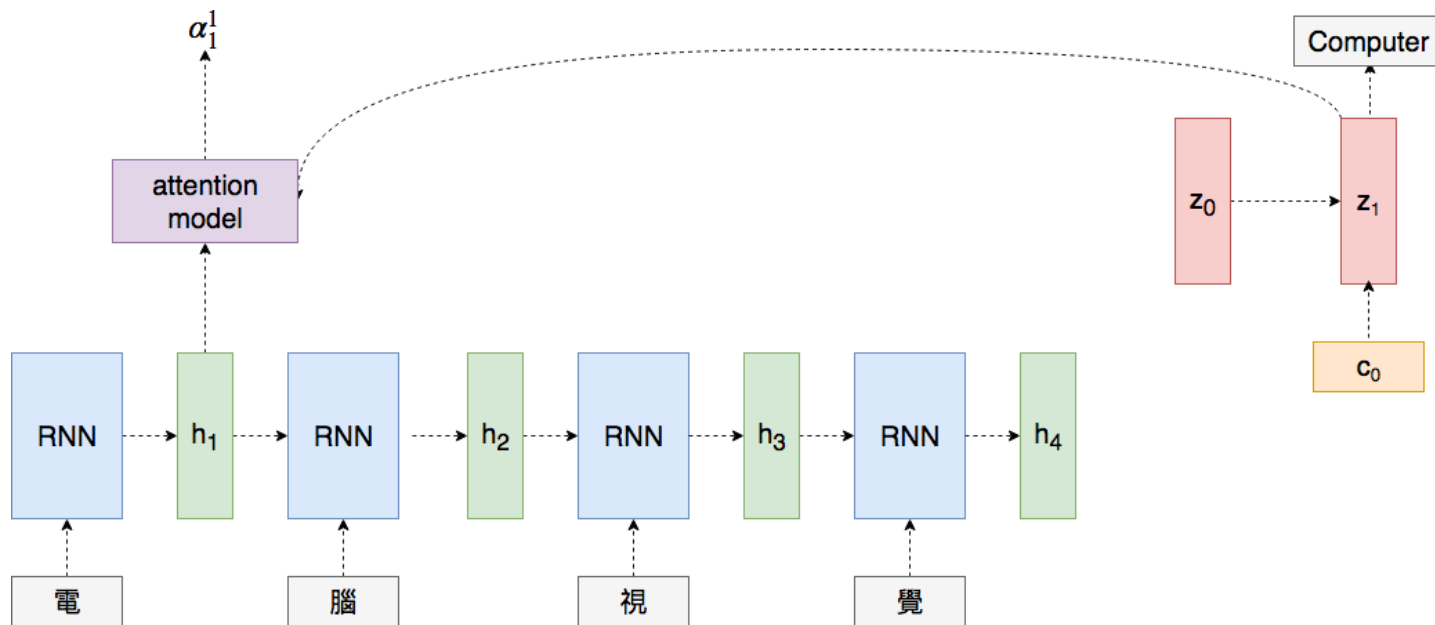
- What should the attention model be?
 - A NN whose inputs are z and h while output is a scalar, indicating the **similarity** between z and h .
- Most attention models are jointly learned or trained with other parts of network (e.g., recognition, etc.)



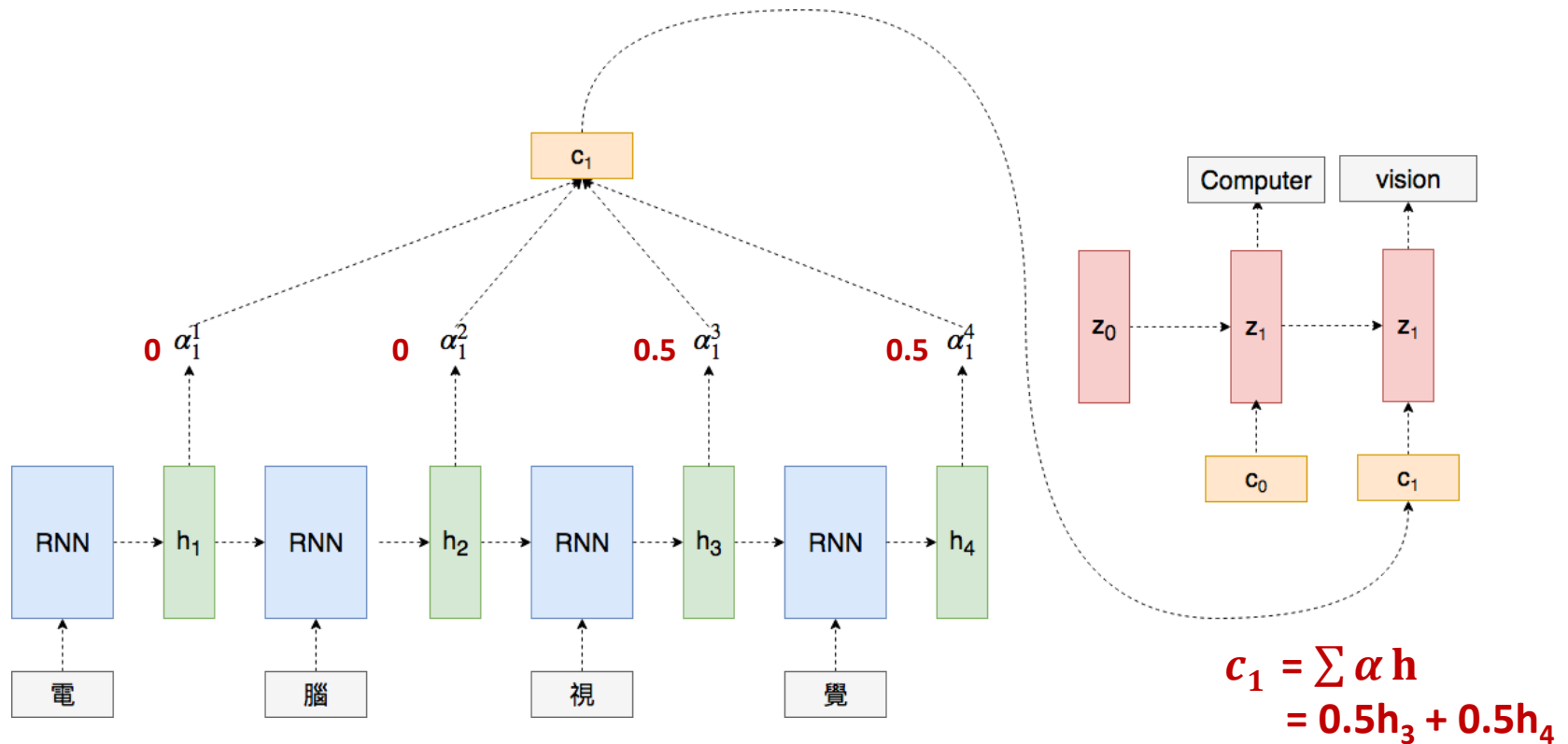
Solution: Attention Model



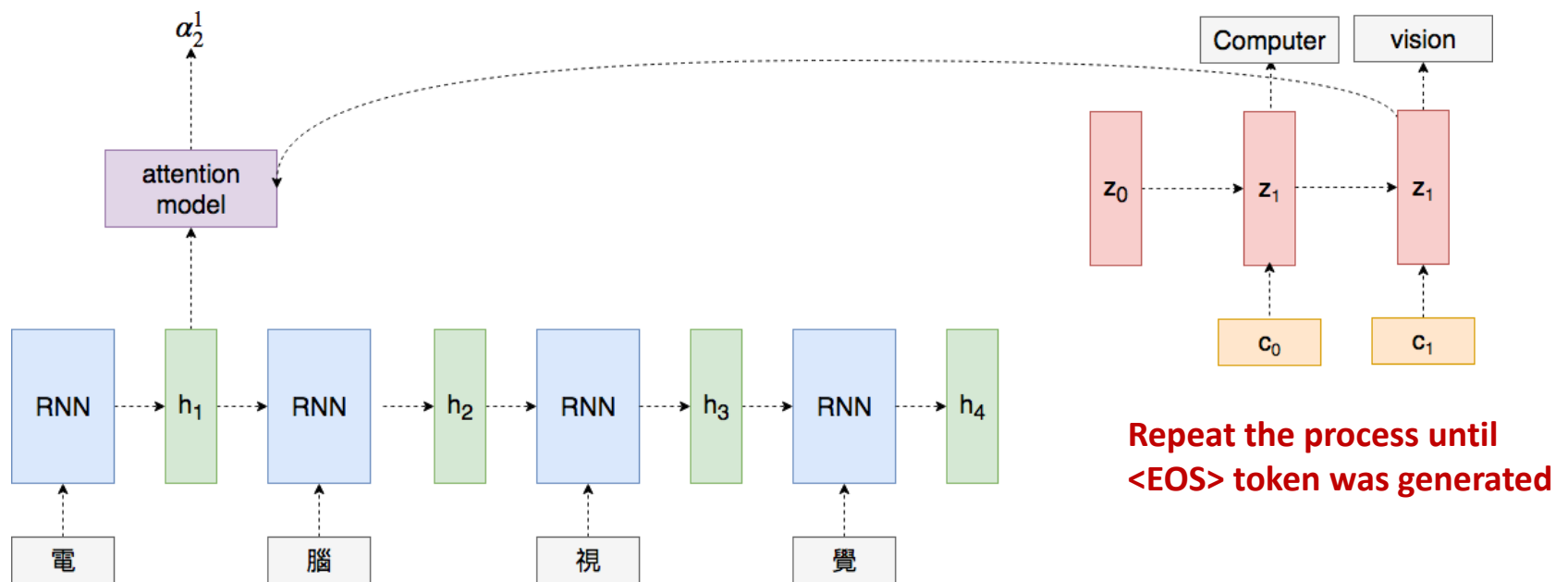
Solution: Attention Model



Solution: Attention Model



Solution: Attention Model



Selected Attention Models for Image-Based Applications

- Image Captioning
 - Xu et al, “Show, Attend and Tell: Neural Image Caption Generation with Visual Attention”, ICML ’15
- Visual Question Answering
 - Zhu et al, “Visual7W: Grounded Question Answering in Images”, CVPR ’16
- Image Classification
 - Mnih et al, “Recurrent Models of Visual Attention”, NIPS ’14
- Image Generation
 - Gregor et al, “DRAW: A Recurrent Neural Network For Image Generation”, ICML ’15

Image Captioning with Attention

- RNN focuses visual attention at different spatial locations when generating corresponding words during captioning.

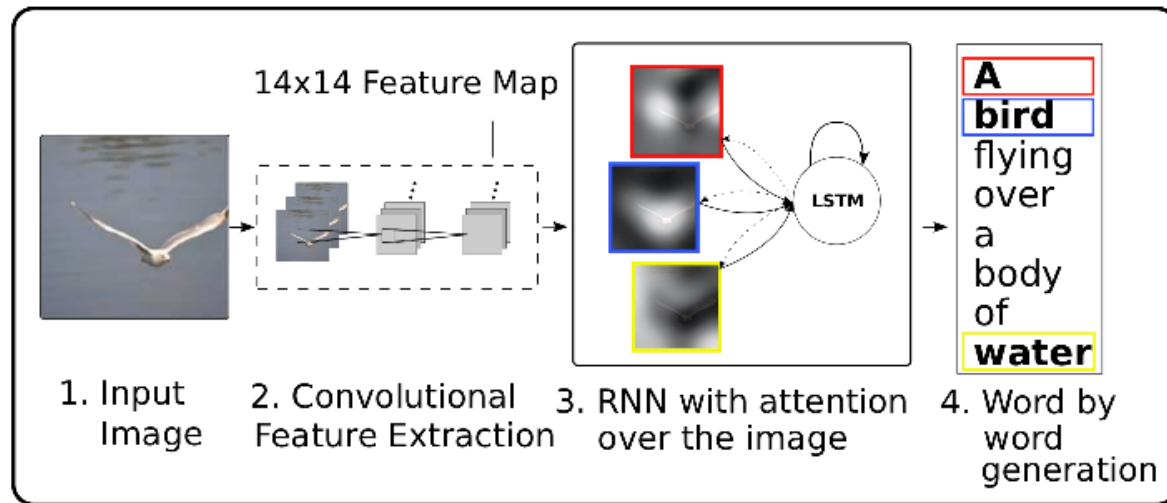


Image Captioning with Attention

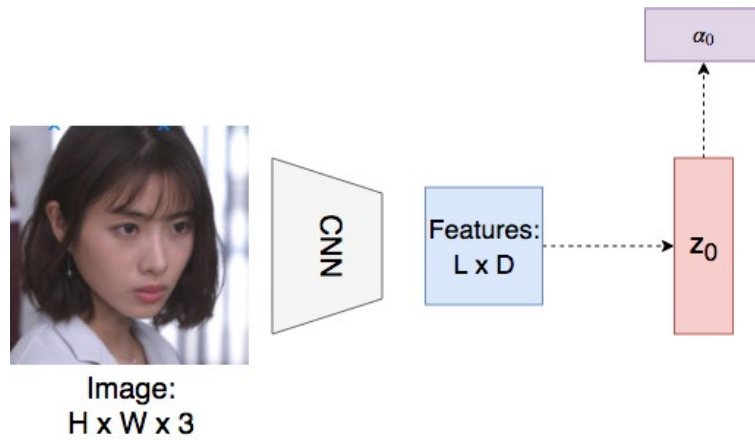


Image Captioning with Attention

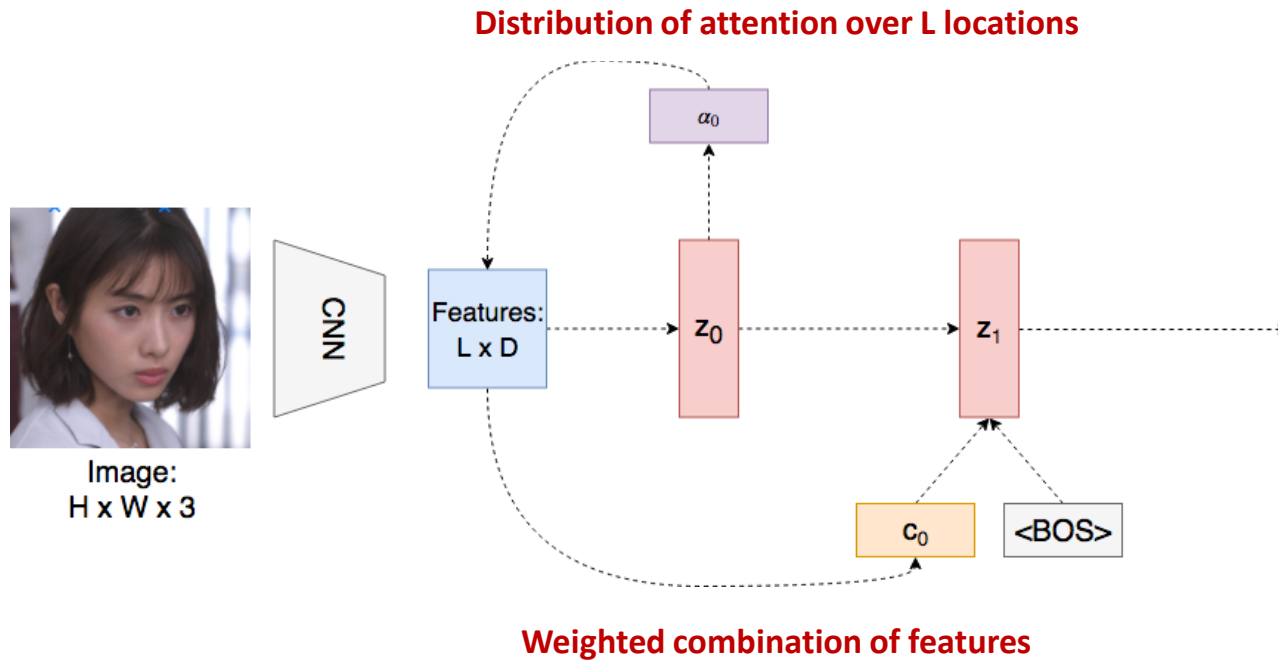


Image Captioning with Attention

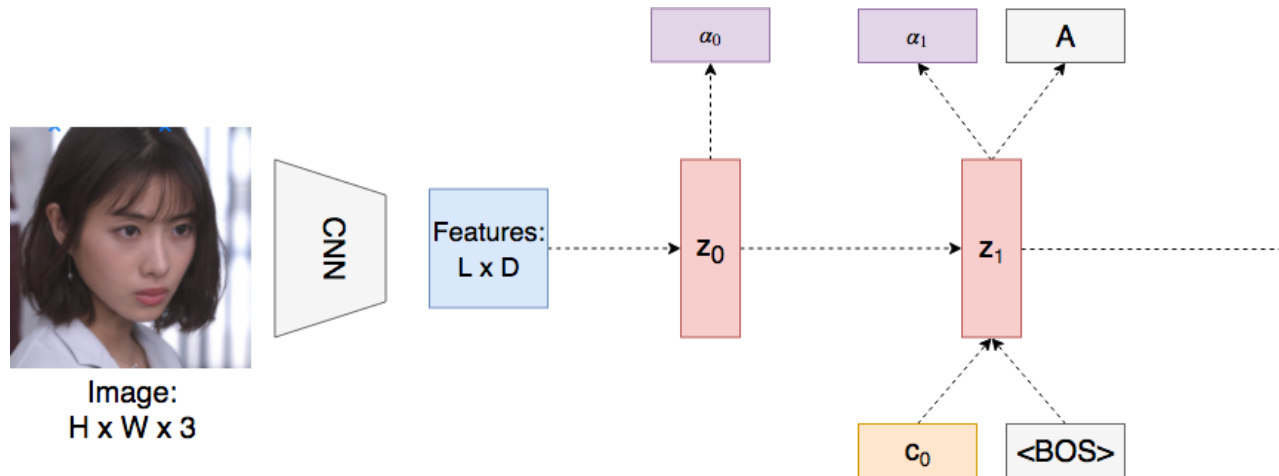


Image Captioning with Attention

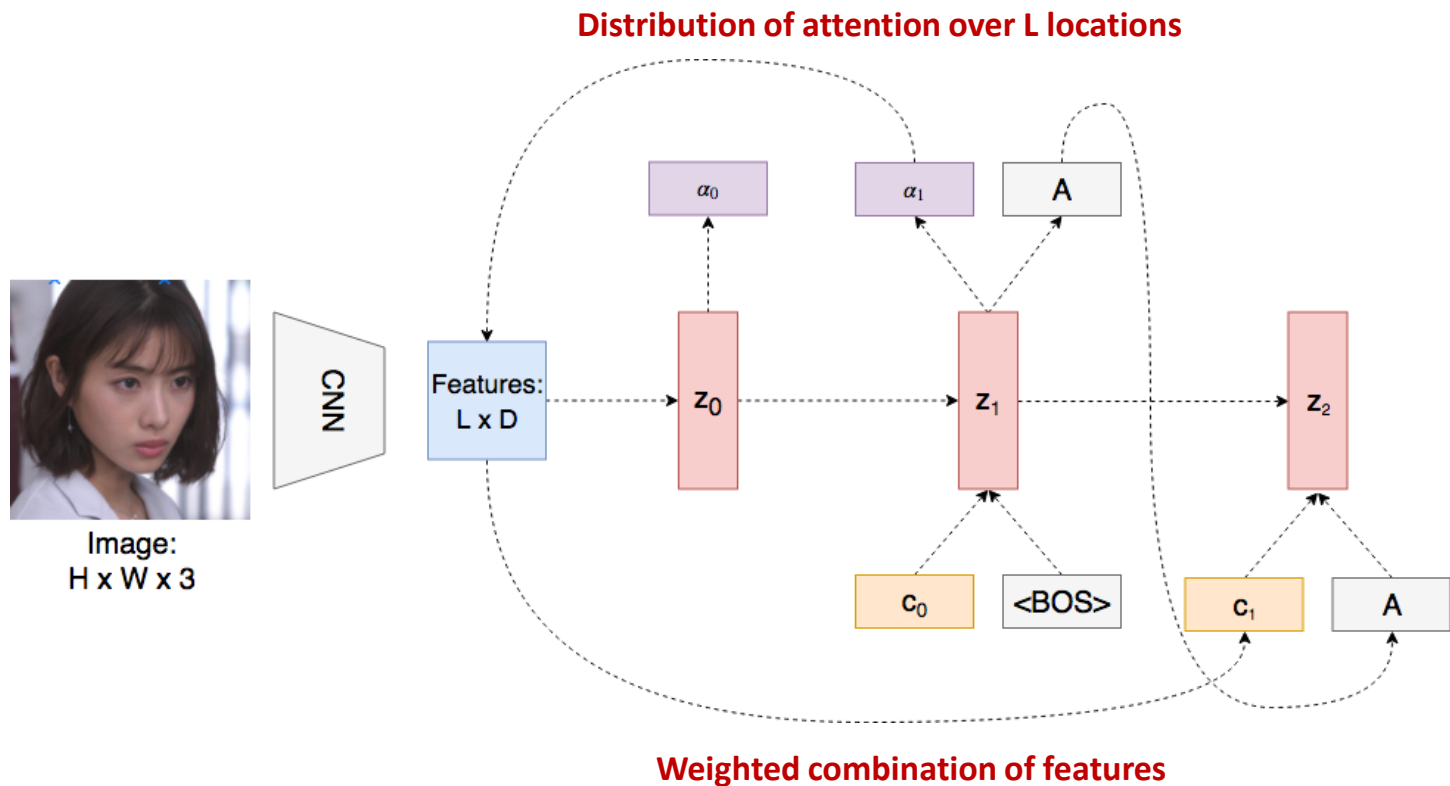
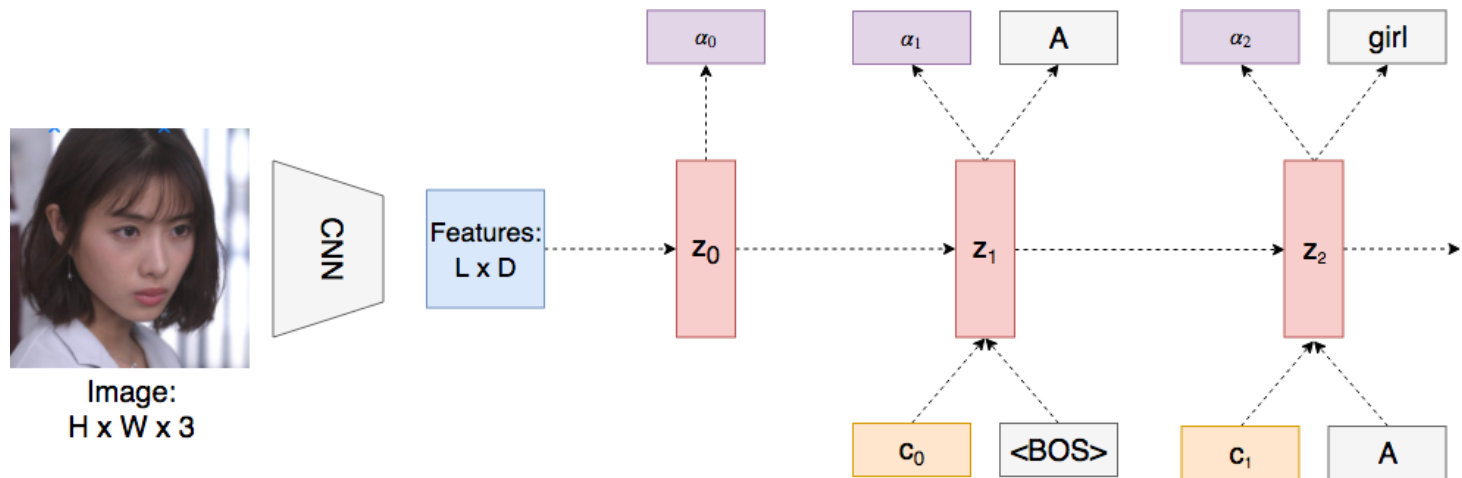


Image Captioning with Attention



Repeat the process until
 $\langle \text{EOS} \rangle$ token was generated

Image Captioning with Attention



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



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



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



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



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

Selected Attention Models for Image-Based Applications

- Image Captioning
 - Xu et al, “Show, Attend and Tell: Neural Image Caption Generation with Visual Attention”, ICML ’15
- Visual Question Answering
 - Zhu et al, “Visual7W: Grounded Question Answering in Images”, CVPR ’16
- Image Classification
 - Mnih et al, “Recurrent Models of Visual Attention”, NIPS ’14
- Image Generation
 - Gregor et al, “DRAW: A Recurrent Neural Network For Image Generation”, ICML ’15

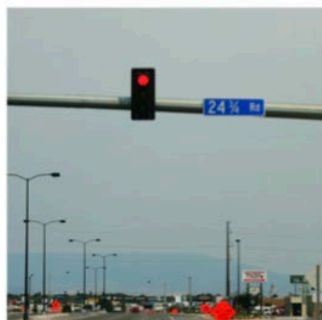
Visual Question Answering

- Examples of multiple-choice QA & pointing QA



Q: What endangered animal is featured on the truck?

- A: **A bald eagle.**
 A: A sparrow.
 A: A humming bird.
 A: A raven.



Q: Where will the driver go if turning right?

- A: **Onto 24 1/2 Rd.**
 A: Onto 25 3/4 Rd.
 A: Onto 23 3/4 Rd.
 A: Onto Main Street.



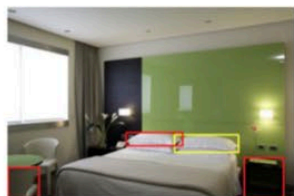
Q: When was the picture taken?

- A: **During a wedding.**
 A: During a bar mitzvah.
 A: During a funeral.
 A: During a Sunday church service.



Q: Who is under the umbrella?

- A: **Two women.**
 A: A child.
 A: An old man.
 A: A husband and a wife.



Q: Which pillow is farther from the window?



Q: Which step leads to the tub?



Q: Which is the small computer in the corner?



Q: Which item is used to cut items?



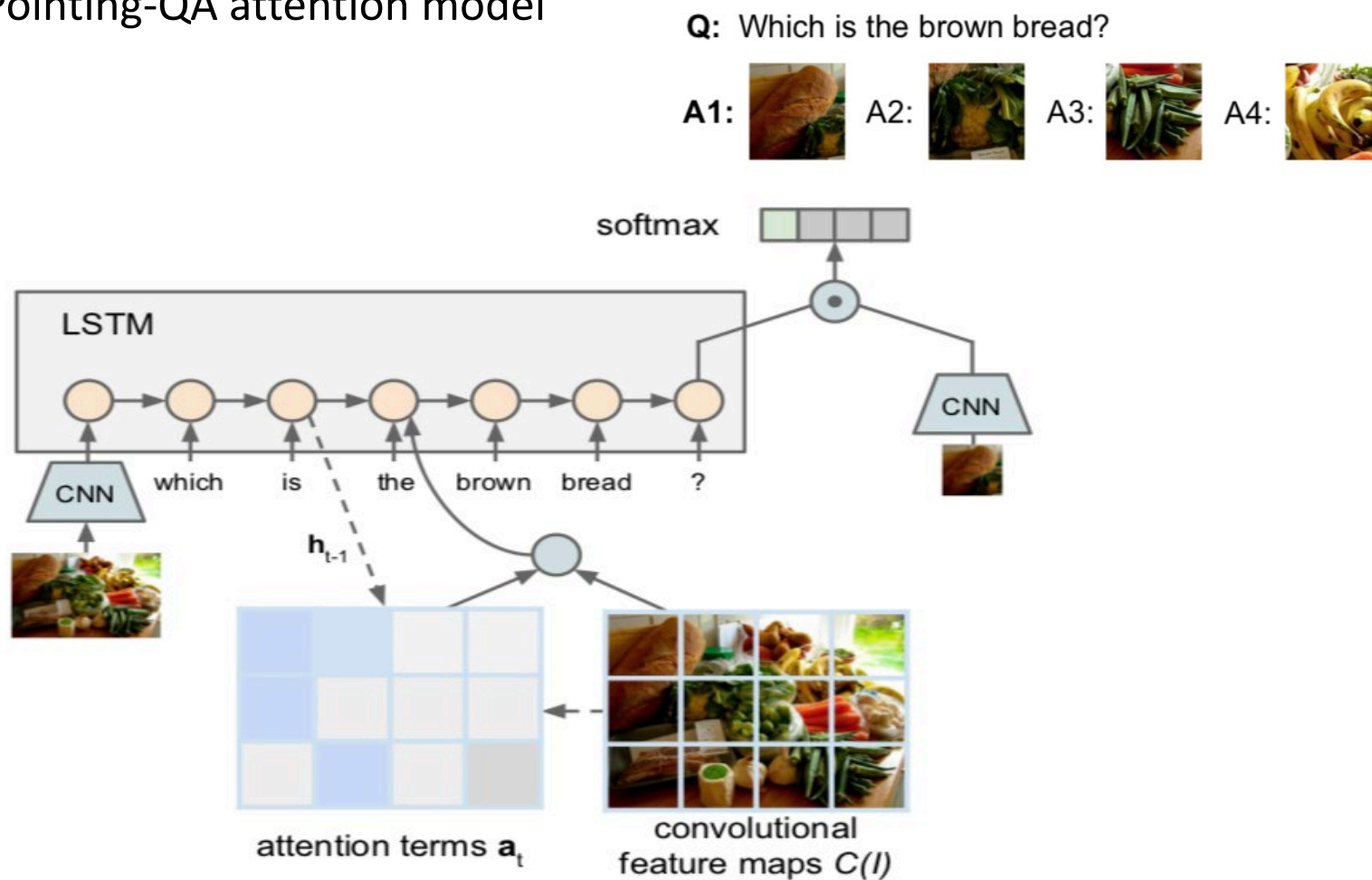
Q: Which doughnut has multicolored sprinkles?



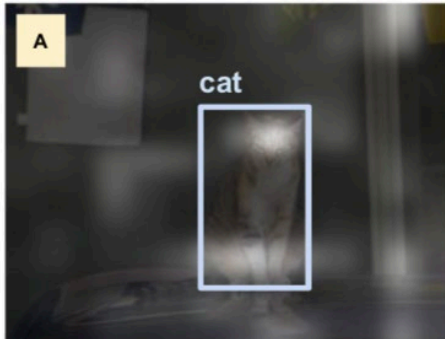
Q: Which man is wearing the red tie?

Visual Question Answering with Attention

- Pointing-QA attention model



Visual Question Answering with Attention



What kind of animal is in the photo?
A **cat**.



Why is the person holding a knife?
To cut the **cake** with.



Where are the carrots?
At the top.



How many people are there?
Three.

The peaks of the attention maps reside in the bounding boxes of the target objects.

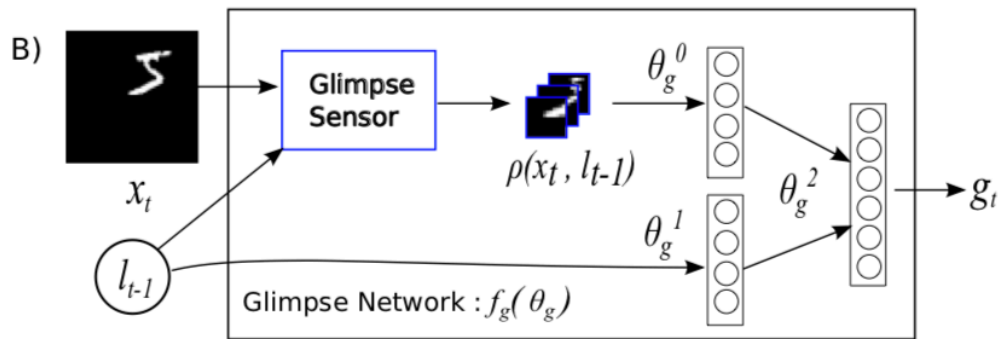
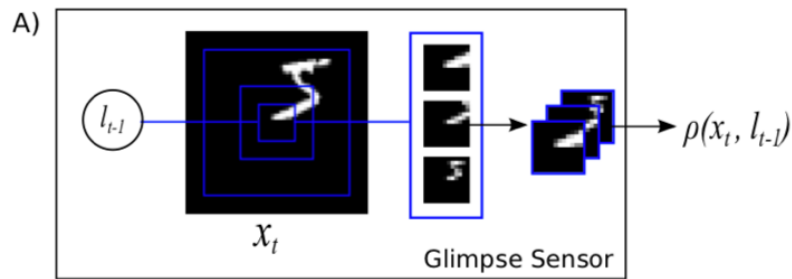
The bottom two examples show QA pairs with answers not explicitly containing objects. The attention heat maps are scattered around the image grids.

Selected Attention Models for Image-Based Applications

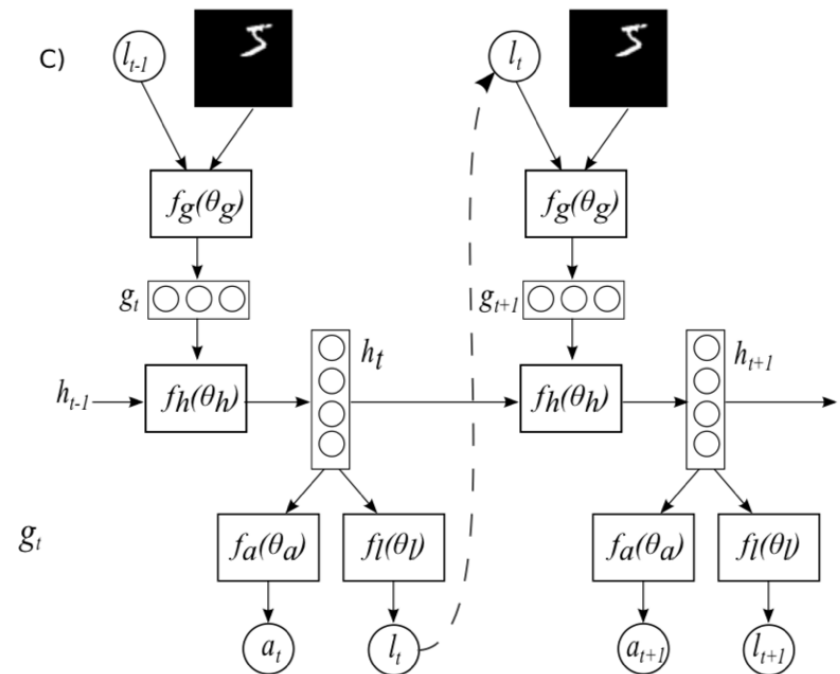
- Image Captioning
 - Xu et al, “Show, Attend and Tell: Neural Image Caption Generation with Visual Attention”, ICML ’15
- Visual Question Answering
 - Zhu et al, “Visual7W: Grounded Question Answering in Images”, CVPR ’16
- Image Classification
 - Mnih et al, “Recurrent Models of Visual Attention”, NIPS ’14
- Image Generation
 - Gregor et al, “DRAW: A Recurrent Neural Network For Image Generation”, ICML ’15

Glimpse Sensor & Glimpse Network

Glimpse sensor: extracts a retina-like representation centered at l_{t-1} that contains multiple resolution patches.

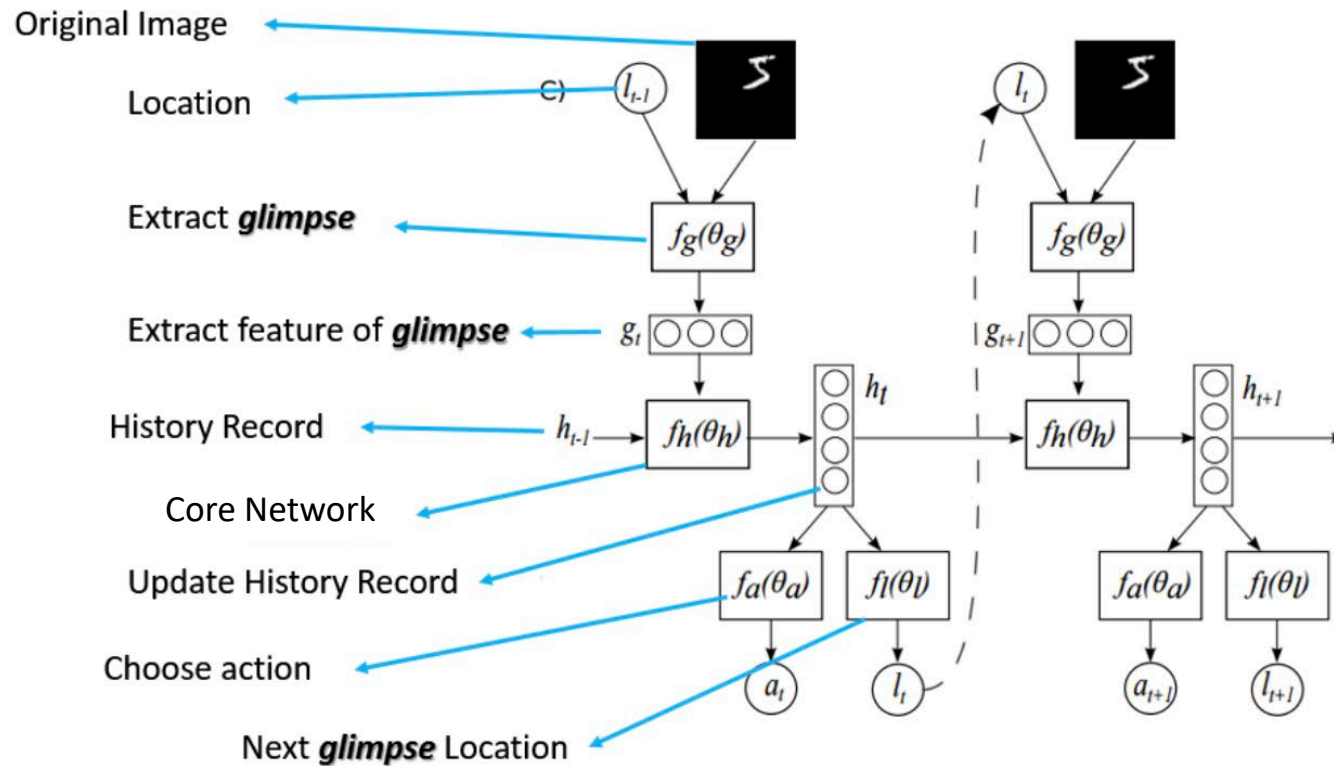


Glimpse network: given location l_{t-1} and image x_t , use the glimpse sensor to extract retina representation, which is mapped into a joint hidden space.



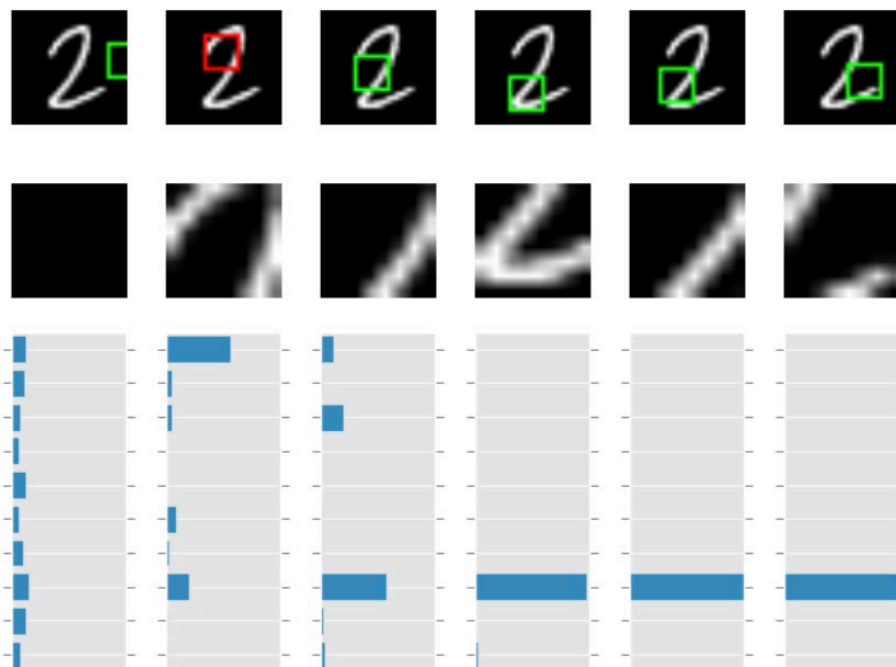
RNN-based model architecture: the core network takes the glimpse representation as input with the hidden state vector from the previous step, and outputs the new hidden state resulting in **location** and **action** networks to predict the next location to attend and the associated action.

Architecture: RNN with Attention Models

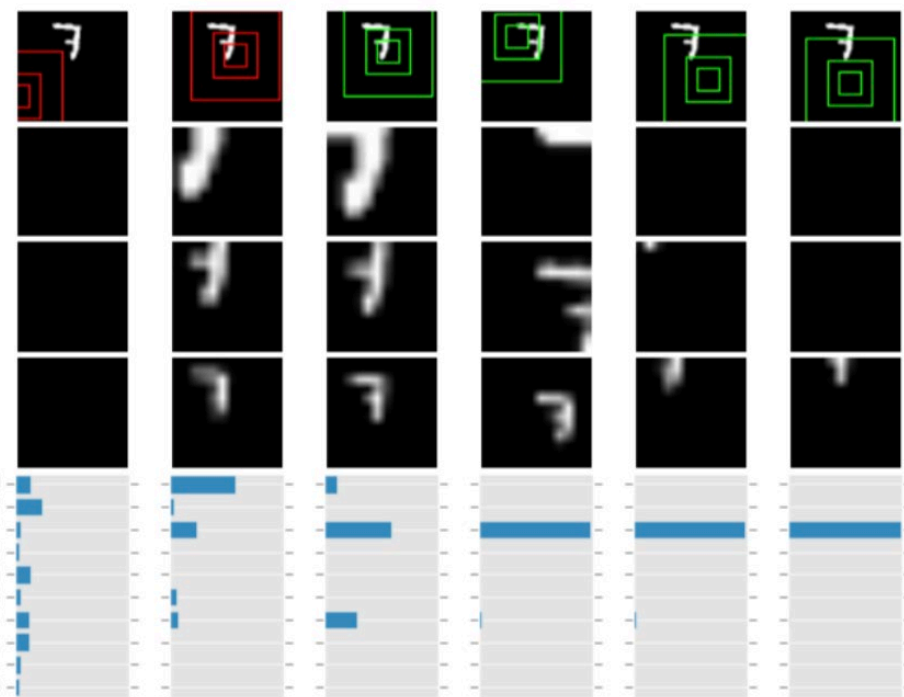


Example Results

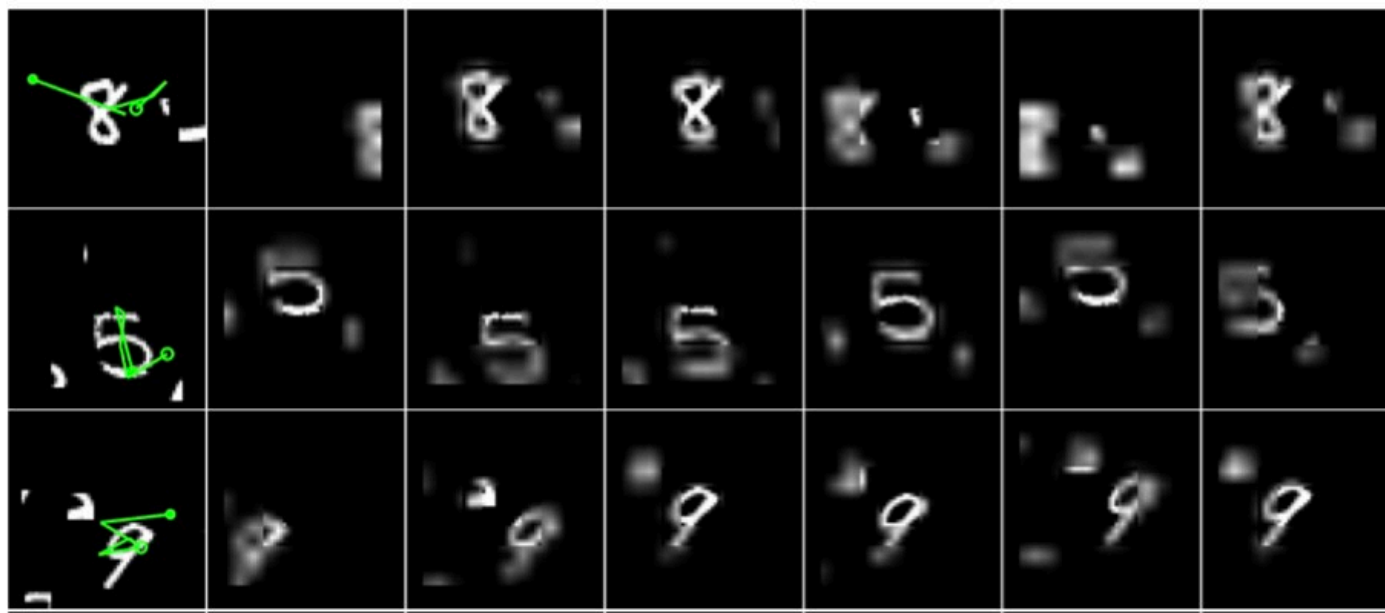
• Original MNIST



• Translated MNIST



Example: Actual Glimpse Path



What We've Covered Today...

- Transfer Learning
 - Visual Synthesis – Style Transfer
- Recurrent Neural Networks
 - From RNN to LSTM & GRU
 - Selected Models for Sequence-to-Sequence Learning
 - Attention in RNN
- Next week: Transformer

