# Attention

M. Soleymani

Sharif University of Technology

Fall 2017

Most slides have been adopted from Fei Fei Li and colleagues lectures, cs231n, Stanford 2016 and some from John Canny, cs294-129, Berkeley, 2016.

# Attention

Focusing on a subset of the given information.

# 2014: Neural Translation Breakthroughs

- Devlin et al, ACL'2014

- Cho et al EMNLP'2014

- Bahdanau, Cho & Bengio, arXiv sept. 2014

- Jean, Cho, Memisevic & Bengio, arXiv dec. 2014

- Sutskever et al NIPS'2014

# Other Applications

- Ba et al 2014, **Visual attention for recognition**

- Mnih et al 2014, **Visual attention for recognition**

- Chorowski et al, 2014, **Speech recognition**

- Graves et al 2014, **Neural Turing machines**

- Yao et al 2015, **Video description generation**

- Vinyals et al, 2015, **Conversational Agents**

- Xu et al 2015, **Image caption generation**

- Xu et al 2015, **Visual Question Answering**

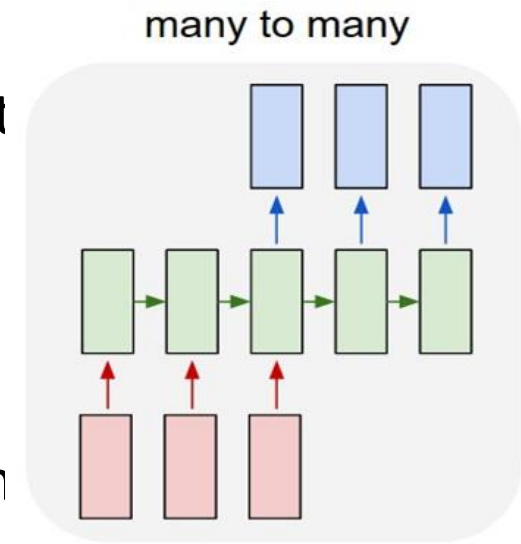# Soft vs Hard Attention Models

**Hard attention:**

- Attend to a single input location among the set of locations.
- Can't use gradient descent.
- Need **reinforcement learning.**

**Soft attention:**

- Compute a weighted combination (attention) over some inputs using an attention network.
- Can use backpropagation to train end-to-end.

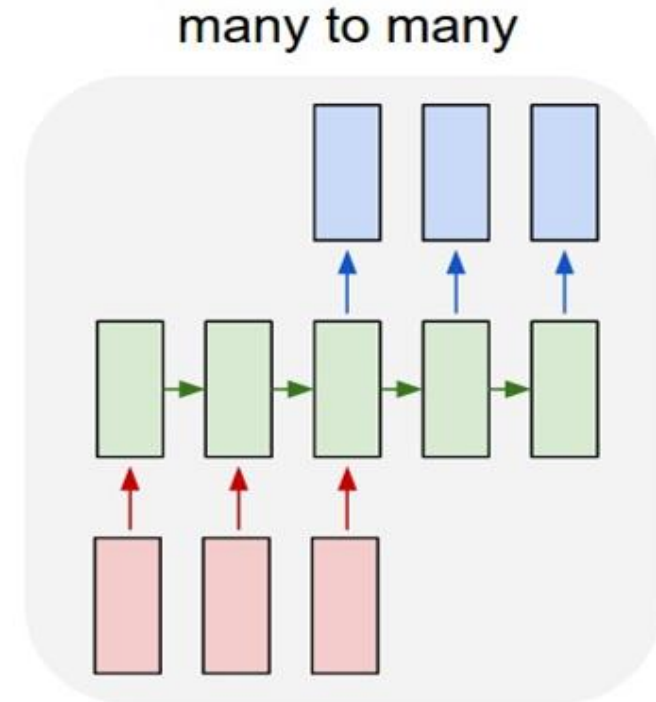# Attention instead of simple encoder-decoder

- Encoder-decoder models
  - needs to be able to compress all the necessary information of a source sentence into a fixed-length vector
  - performance deteriorates rapidly as the length of an input sentence increases.

- Attention avoids this by:
  - allowing the RNN generating the output to focus on hidden states (generated by the first RNN) as they become relevant.

many to many

Bahdanau et al, "Neural Machine Translation by Jointly Learning to Align and Translate", ICLR 2015
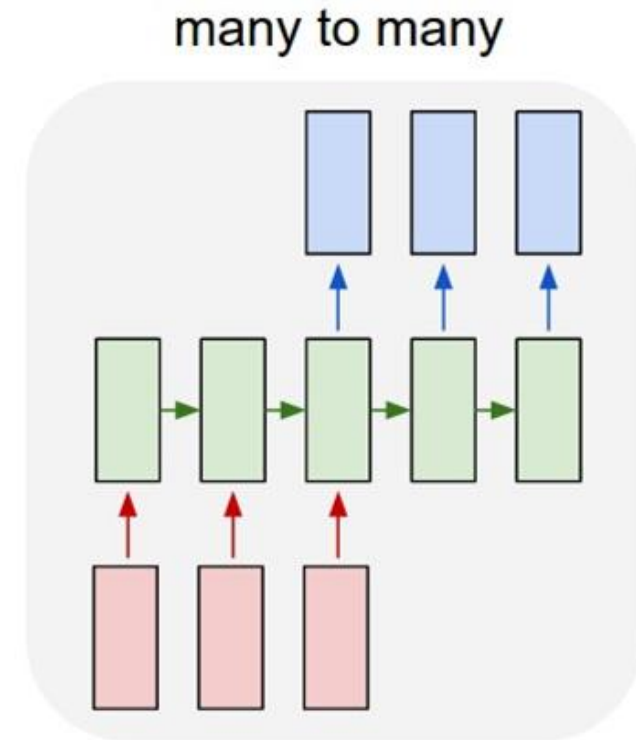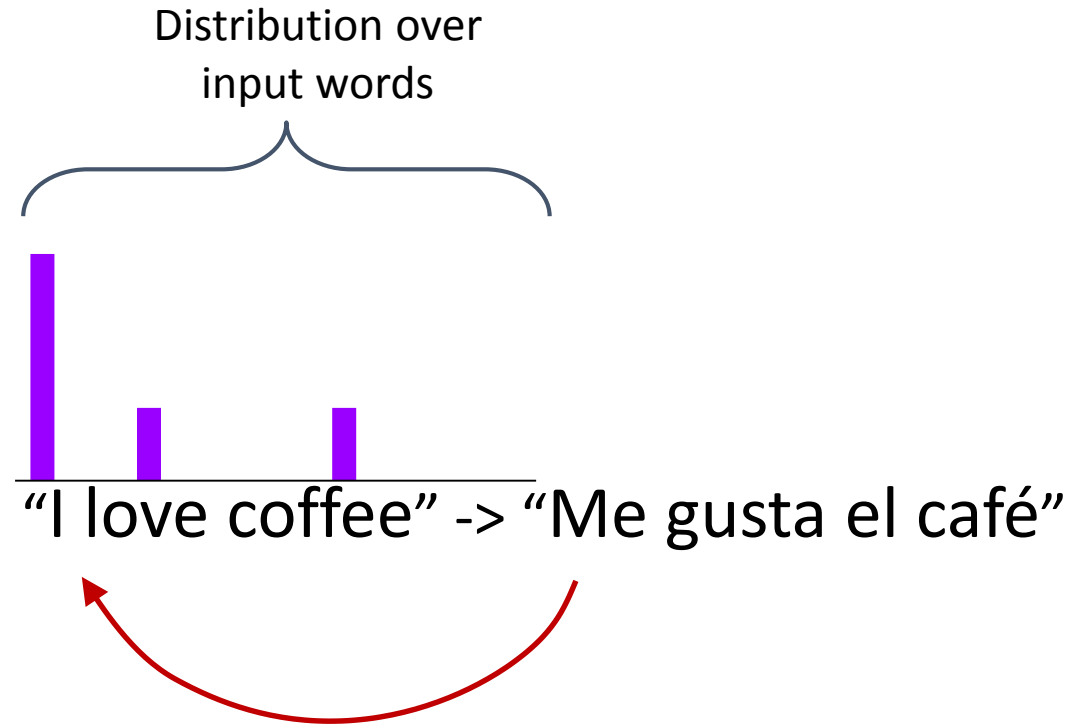
# Soft Attention for Translation

An RNN can attend over the output of another RNN. At every time step, it focuses on different positions in the other RNN.
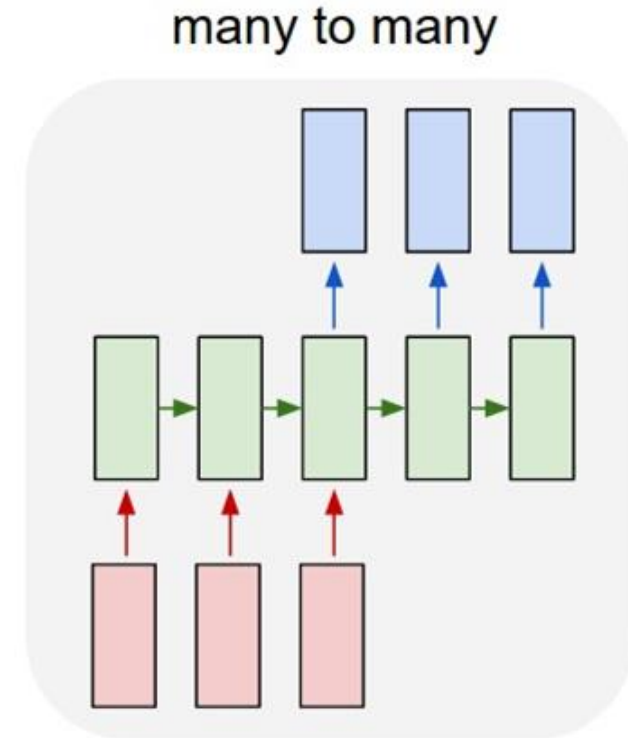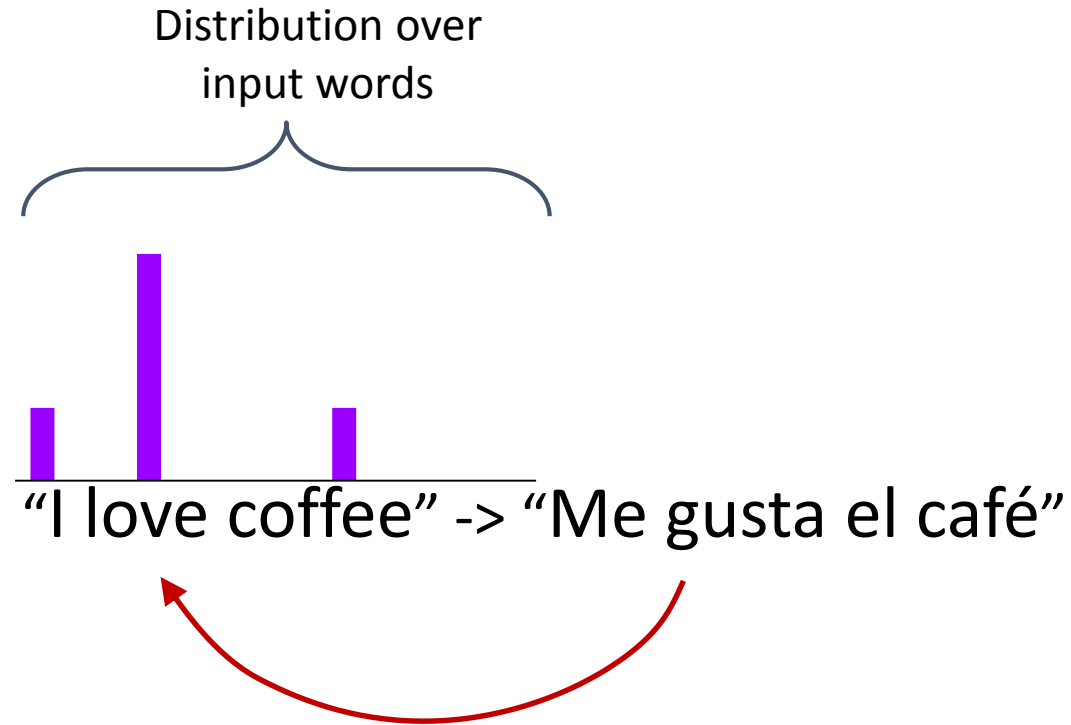
"I love coffee" -> "Me gusta el café"



many to many

Bahdanau et al, "Neural Machine Translation by Jointly Learning to Align and Translate", ICLR 2015

# Soft Attention for Translation

Distribution over
input words



"I love coffee" -> "Me gusta el café"

many to many

Bahdanau et al, "Neural Machine Translation by Jointly
Learning to Align and Translate", ICLR 2015

# Soft Attention for Translation

Distribution over
input words



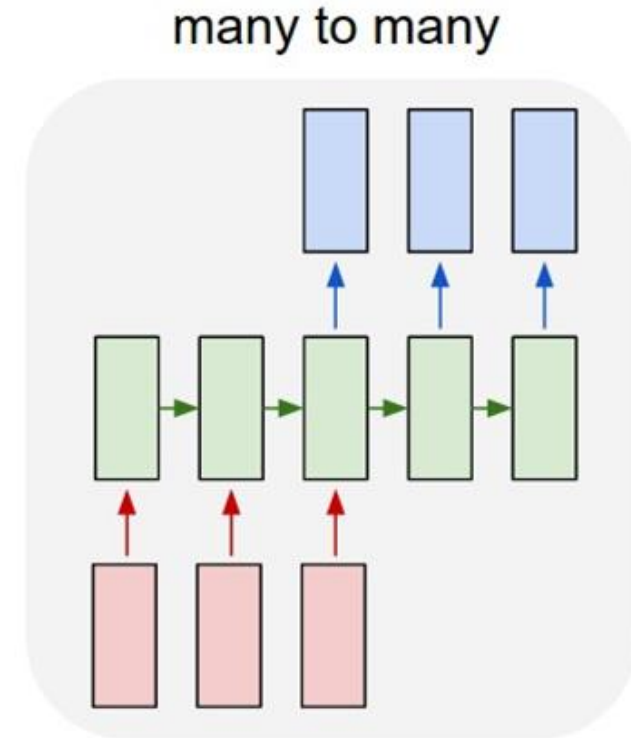"I love coffee" -> "Me gusta el café"

many to many

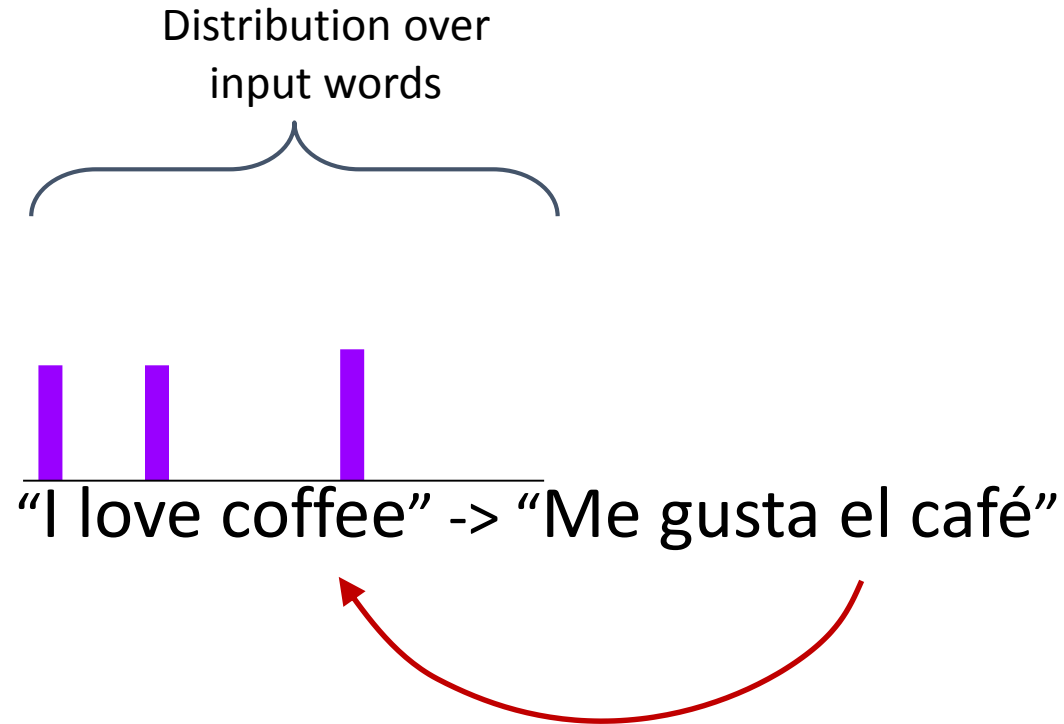Bahdanau et al, "Neural Machine Translation by Jointly
Learning to Align and Translate", ICLR 2015

# Soft Attention for Translation

Distribution over
input words



"I love coffee" -> "Me gusta el café"

many to many

Bahdanau et al, "Neural Machine Translation by Jointly
Learning to Align and Translate", ICLR 2015

# Soft Attention for Translation

Distribution over
input words



"I love coffee" -> "Me gusta el café"


many to many

Bahdanau et al, "Neural Machine Translation by Jointly
Learning to Align and Translate", ICLR 2015

# Soft Attention for Translation



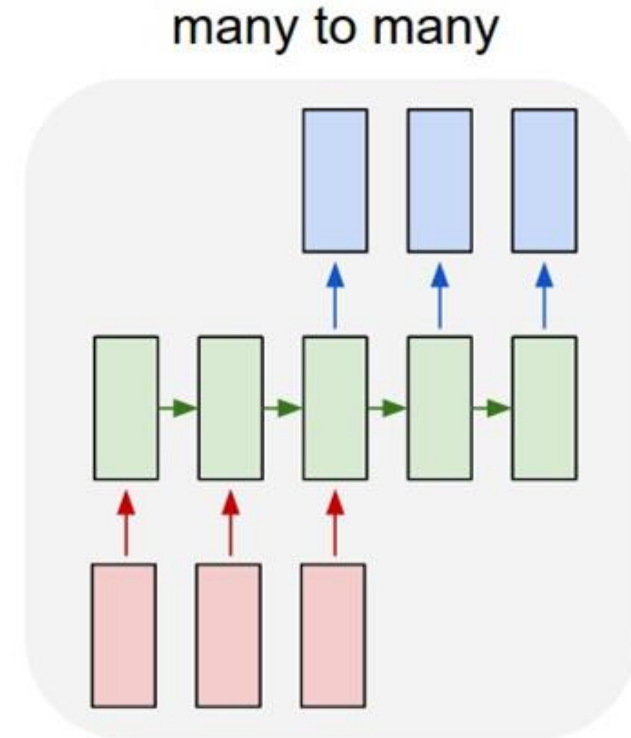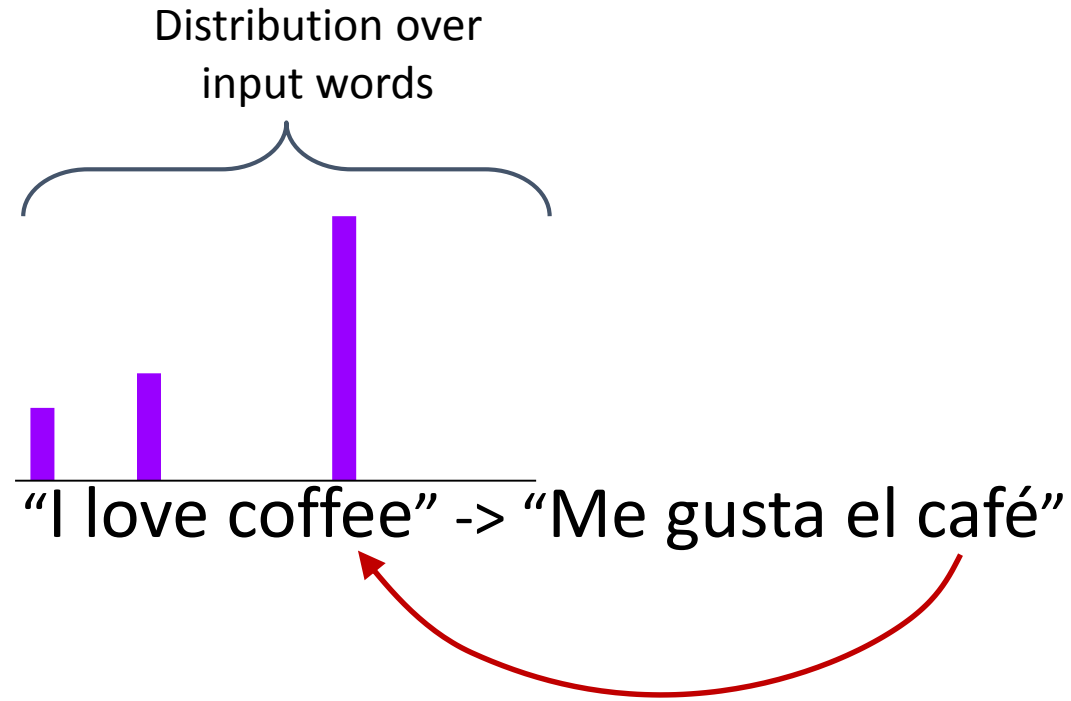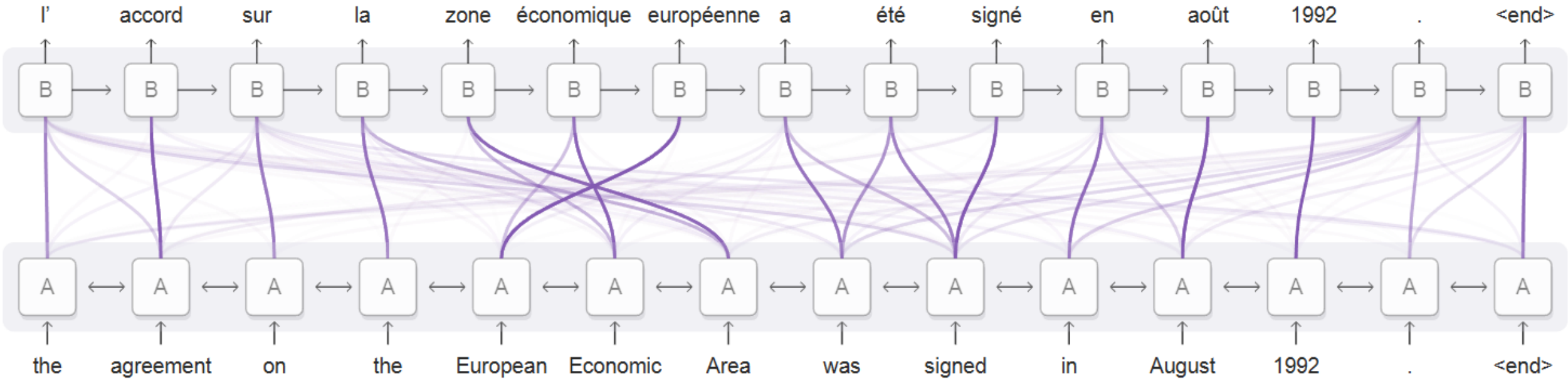Diagram derived from Fig. 3 of Bahdanau, *et al.* 2014

Source: https://distill.pub/2016/augmented-rnns/

# Soft Attention for voice recognition



Figure derived from Chan, et al. 2015

Source: https://distill.pub/2016/augmented-rnns/

# Simple soft attention mechanism



The attending RNN generates a query describing what it wants to focus on.

Each item is dot producted with the query to produce a score, describing how well it matches the query. The scores are fed into a softmax to create the attention distribution.

Source: https://distill.pub/2016/augmented-rnns/

# Soft Attention for Translation

Context vector (input to decoder):

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

Mixture weights:

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

Alignment score (how well do input words near j match output words at position i):

$$e_{ij} = a(s_{i-1}, h_j)$$

$$s_i = f(s_{i-1}, y_{i-1}, c_i)$$



Bahdanau et al, "Neural Machine Translation by Jointly Learning to Align and Translate", ICLR 2015

alignment model: s a feedforward neural network which is jointly trained with all the other components of the proposed system

# Alleviate fixed length encoding

- The decoder decides parts of the source sentence to pay attention to.

- By letting the decoder have an attention mechanism, we relieve the encoder from the burden of source sentence into a fixed length vector
  - the information can be spread throughout the sequence
  - and can be selectively retrieved by the decoder accordingly.

# Soft Attention for Translation



From Y. Bengio CVPR 2015 Tutorial

# Soft Attention for Translation



From Y. Bengio CVPR 2015 Tutorial

# Soft Attention for Translation



From Y. Bengio CVPR 2015 Tutorial

# Soft Attention for Translation



From Y. Bengio CVPR 2015 Tutorial

# Soft Attention for Translation



(a)

(b)

Bahdanau et al, "Neural Machine Translation by Jointly
Learning to Align and Translate", ICLR 2015

Stacked LSTM (c.f. bidirectional flat encoder in Bahdanau et al):



Effective Approaches to Attention-based Neural Machine Translation
Minh-Thang Luong, Hieu Pham, Christopher D. Manning, EMNLP 15

# Global Attention Model

Global attention model is similar but simpler than Badanau's:

Different word matching functions were used



$$a_t(s) = \text{align}(h_t, \bar{h}_s)$$

$$= \frac{\exp\left(\text{score}(h_t, \bar{h}_s)\right)}{\sum_{s'} \exp\left(\text{score}(h_t, \bar{h}_{s'})\right)}$$

$$\text{score}(h_t, \bar{h}_s) = \begin{cases} h_t^\top \bar{h}_s & \text{dot} \\ h_t^\top W_a \bar{h}_s & \text{general} \\ v_a^\top \tanh\left(W_a[h_t; \bar{h}_s]\right) & \text{concat} \end{cases}$$

Effective Approaches to Attention-based Neural Machine Translation
Minh-Thang Luong Hieu Pham Christopher D. Manning, EMNLP 15

# Local Attention Model

- Compute a best aligned position $p_t$ first

- Then compute a context vector centered at that position



$$p_t = S \cdot \text{sigmoid}(\boldsymbol{v}_p^\top \tanh(\boldsymbol{W_p h_t})),$$

$$\boldsymbol{a}_t(s) = \text{align}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_s) \exp\left(-\frac{(s - p_t)^2}{2\sigma^2}\right)$$

Effective Approaches to Attention-based Neural Machine Translation
Minh-Thang Luong Hieu Pham Christopher D. Manning, EMNLP 15

# Results

| System | Ppl | BLEU |
|---|---|---|
| Winning WMT'14 system – *phrase-based* + *large LM* (Buck et al., 2014) | | 20.7 |
| *Existing NMT systems* | | |
| RNNsearch (Jean et al., 2015) | | 16.5 |
| RNNsearch + unk replace (Jean et al., 2015) | | 19.0 |
| RNNsearch + unk replace + large vocab + *ensemble* 8 models (Jean et al., 2015) | | **21.6** |
| *Our NMT systems* | | |
| Base | 10.6 | 11.3 |
| Base + reverse | 9.9 | 12.6 (+*1.3*) |
| Base + reverse + dropout | 8.1 | 14.0 (+*1.4*) |
| Base + reverse + dropout + global attention (*location*) | 7.3 | 16.8 (+*2.8*) |
| Base + reverse + dropout + global attention (*location*) + feed input | 6.4 | 18.1 (+*1.3*) |
| Base + reverse + dropout + local-p attention (*general*) + feed input | 5.9 | 19.0 (+*0.9*) |
| Base + reverse + dropout + local-p attention (*general*) + feed input + unk replace | | 20.9 (+*1.9*) |
| *Ensemble* 8 models + unk replace | | **23.0** (+*2.1*) |

Local *and* global models

Effective Approaches to Attention-based Neural Machine Translation
Minh-Thang Luong Hieu Pham Christopher D. Manning, EMNLP 15

# Image Captioning with Attention



RNN focuses its attention at a different spatial location when generating each word

1. Input Image
2. Convolutional Feature Extraction
14x14 Feature Map
3. RNN with attention over the image
LSTM
4. Word by word generation

A bird flying over a body of water

Image:
H x W x 3

# Recall: RNN for Captioning



Image:
H x W x 3

CNN

Features:
D

Image:
H x W x 3

Features:
D

Hidden
state: H

CNN

h0

# Recall: RNN for Captioning



Distribution over vocab

d1

CNN

h0 → h1 → d1

Image:
H x W x 3

Features:
D

Hidden
state: H

y1

First
word

# Recall: RNN for Captioning

# Recall: RNN for Captioning



Distribution over vocab

RNN only looks at whole image, once

Image: H x W x 3

CNN

Features: D

Hidden state: H

First word

Second word

# Recall: RNN for Captioning



Distribution over vocab

RNN only looks at whole image, once

Image: H x W x 3

Features: D

Hidden state: H

First word

Second word

What if the RNN looks at different parts of the image at each timestep?

# Soft Attention for Captioning



Image:
H x W x 3

CNN

Features:
L x D

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Soft Attention for Captioning



Image:
H x W x 3

CNN

Features:
L x D

h0

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Soft Attention for Captioning



Distribution over L locations

a1

CNN

h0

Image:
H x W x 3

Features:
L x D

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Soft Attention for Captioning



Distribution over L locations

a1

CNN

h0

Features: L x D

Image: H x W x 3

Weighted features: D

z1

Weighted combination of features

# Soft Attention for Captioning

# Soft Attention for Captioning

# Soft Attention for Captioning

# Soft Attention for Captioning

# Soft Attention for Captioning

# Soft vs Hard Attention



Image:
H x W x 3

CNN

Grid of features
(Each D-dimensional)

$$\begin{array}{|c|c|} \hline a & b \\ \hline c & d \\ \hline \end{array}$$

From
RNN:

$$\begin{array}{|c|c|} \hline p_a & p_b \\ \hline p_c & p_d \\ \hline \end{array}$$

Distribution over
grid locations
$p_a + p_b + p_c + p_c = 1$

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Soft vs Hard Attention



Image:
H x W x 3

CNN

Grid of features
(Each D-dimensional)

Context vector z
(D-dimensional)

From
RNN:

$p_a$ | $p_b$
$p_c$ | $p_d$

Distribution over
grid locations
$p_a + p_b + p_c + p_c = 1$

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Soft vs Hard Attention



Image:
H x W x 3

Grid of features
(Each D-dimensional)

CNN

From
RNN:

Distribution over
grid locations
$p_a + p_b + p_c + p_c = 1$

Context vector z
(D-dimensional)

**Soft attention:**
Summarize ALL locations
$z = p_a a + p_b b + p_c c + p_d d$

Derivative dz/dp is nice!
Train with gradient descent

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Soft vs Hard Attention



Image:
H x W x 3

CNN

Grid of features
(Each D-dimensional)

From
RNN:

Distribution over
grid locations
$p_a + p_b + p_c + p_c = 1$

Context vector z
(D-dimensional)

**Soft attention:**
Summarize ALL locations
$z = p_a a + p_b b + p_c c + p_d d$

Derivative dz/dp is nice!
Train with gradient descent

**Hard attention**:
Sample ONE location
according to p, z = that vector

With argmax, dz/dp is zero
almost everywhere …
Can't use gradient descent;
need reinforcement learning

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Soft Attention for Captioning



Model want to attend to salient part of an image while generating its caption

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Soft Attention for Captioning



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.

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Visual Question Answering



**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 ¾ Rd.
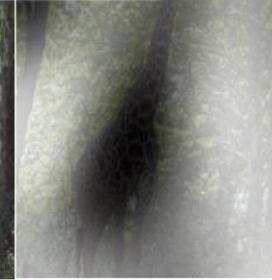**A:** Onto 25 ¾ Rd.
**A:** Onto 23 ¾ Rd.
**A:** Onto Main Street.

**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.

Agrawal et al, "VQA: Visual Question Answering", ICCV 2015
Zhu et al, "Visual 7W: Grounded Question Answering in Images", CVPR 2016

Zhu et al, "Visual 7W: Grounded Question Answering in Images", CVPR 2016
Figures from Zhu et al, copyright IEEE 2016. Reproduced for educational purposes.

# Soft Attention for Video

"Describing Videos by Exploiting Temporal Structure," Li Yao et al, arXiv 2015.

# Soft Attention for Video

The attention model:



Features-Extraction     Soft-Attention     Caption Generation

"Describing Videos by Exploiting Temporal Structure," Li Yao et al, arXiv 2015.

# Soft Attention for Video

Table 1. Performance of different variants of the model on the Youtube2Text and DVS datasets.

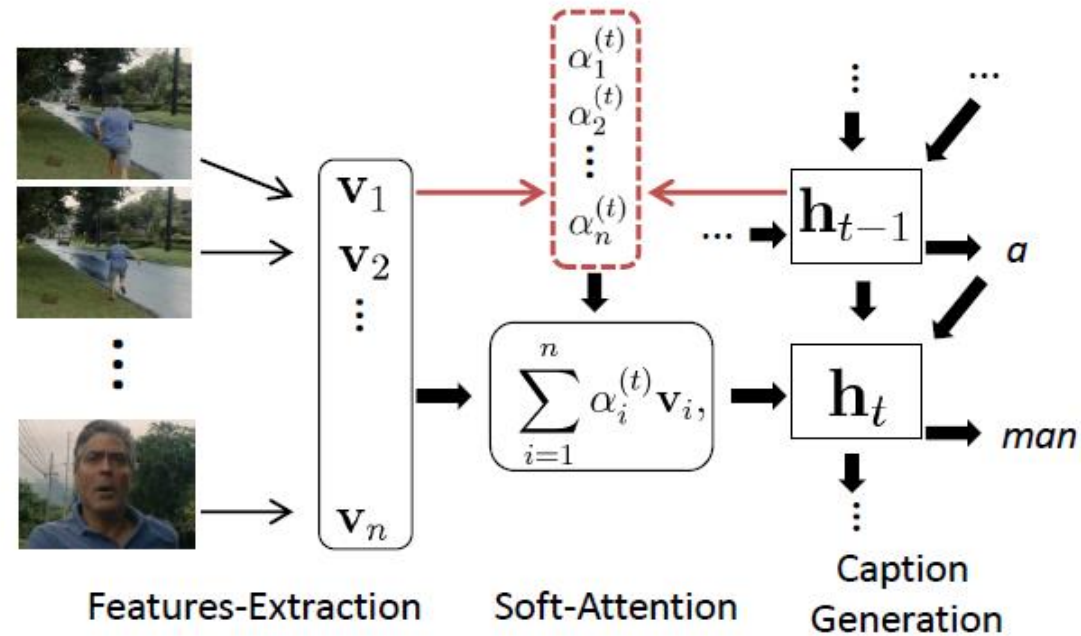| Model | Youtube2Text | | | | DVS | | | |
|---|---|---|---|---|---|---|---|---|
| | BLEU | METEOR | CIDEr | Perplexity | BLEU | METEOR | CIDEr | Perplexity |
| Enc-Dec (Basic) | 0.3869 | 0.2868 | 0.4478 | 33.09 | 0.003 | 0.044 | 0.044 | 88.28 |
| + Local (3-D CNN) | 0.3875 | 0.2832 | 0.5087 | 33.42 | 0.004 | 0.051 | 0.050 | 84.41 |
| + Global (Temporal Attention) | 0.4028 | 0.2900 | 0.4801 | 27.89 | 0.003 | 0.040 | 0.047 | 66.63 |
| + Local + Global | **0.4192** | **0.2960** | **0.5167** | **27.55** | **0.007** | **0.057** | **0.061** | **65.44** |
| Venugopalan *et al.* [41] | 0.3119 | 0.2687 | - | - | - | - | - | - |
| + Extra Data (Flickr30k, COCO) | 0.3329 | 0.2907 | - | - | - | - | - | - |
| Thomason *et al.* [37] | 0.1368 | 0.2390 | - | - | - | - | - | - |

"Describing Videos by Exploiting Temporal Structure," Li Yao et al, arXiv 2015.

# Soft Attention for Captioning

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.

Image:
H x W x 3

Features:
L x D

A woman is throwing a frisbee in a park.

Attention mechanism from Show, Attend, and Tell only lets us softly attend to fixed grid positions ... can we do better?

# Spatial Transformer Networks



Input
image:
H x W x 3

Box
Coordinates:
(xc, yc, w, h)

Cropped and
rescaled image:
X x Y x 3

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

# Spatial Transformer Networks



Can we make this function differentiable?

Input image: H x W x 3

Box Coordinates: (xc, yc, w, h)

Cropped and rescaled image: X x Y x 3

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

# Spatial Transformer Networks

Can we make this function differentiable?

Input image: H x W x 3

Box Coordinates: (xc, yc, w, h)

Cropped and rescaled image: X x Y x 3

**Idea**: Function mapping *pixel coordinates* (xt, yt) of output to *pixel coordinates* (xs, ys) of input

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

# Spatial Transformer Networks

(x$^s$, y$^s$)

Can we make this function differentiable?

(x$^t$, y$^t$)

Input image: H x W x 3

Box Coordinates: (xc, yc, w, h)

Cropped and rescaled image: X x Y x 3

**Idea**: Function mapping *pixel coordinates* (xt, yt) of output to *pixel coordinates* (xs, ys) of input

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

# Spatial Transformer Networks

(x^s, y^s)

Can we make this function differentiable?
(x^t, y^t)

Input image:
H x W x 3

Cropped and rescaled image:
X x Y x 3

Box Coordinates:
(xc, yc, w, h)

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

**Idea**: Function mapping *pixel coordinates* (xt, yt) of output to *pixel coordinates* (xs, ys) of input

$$
\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}
$$

# Spatial Transformer Networks

Input
image:
H x W x 3

Box
Coordinates:
(xc, yc, w, h)
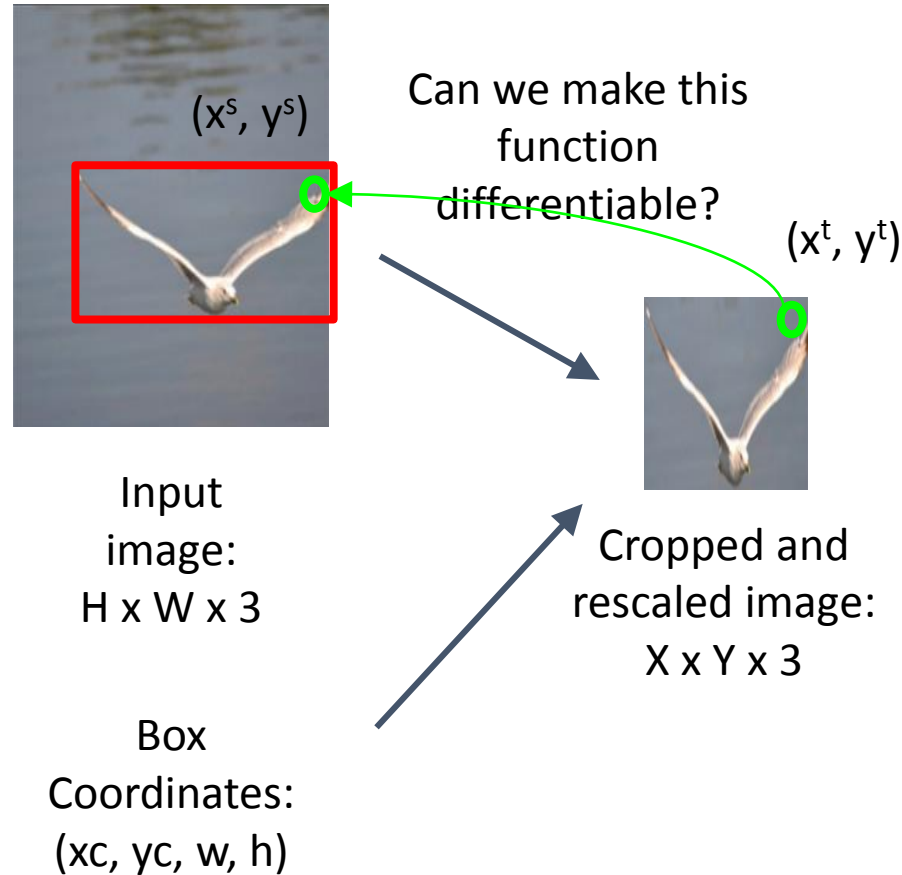
Can we make this
function
differentiable?

Cropped and
rescaled image:
X x Y x 3

**Idea**: Function mapping *pixel
coordinates* (xt, yt) of output
to *pixel coordinates* (xs, ys) of
input

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$

Repeat for all
pixels in *output*
to get a **sampling
grid**

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

# Spatial Transformer Networks

Can we make this function differentiable?

Input image: H x W x 3

Box Coordinates: (xc, yc, w, h)

Cropped and rescaled image: X x Y x 3

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

**Idea**: Function mapping *pixel coordinates* (xt, yt) of output to *pixel coordinates* (xs, ys) of input

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$
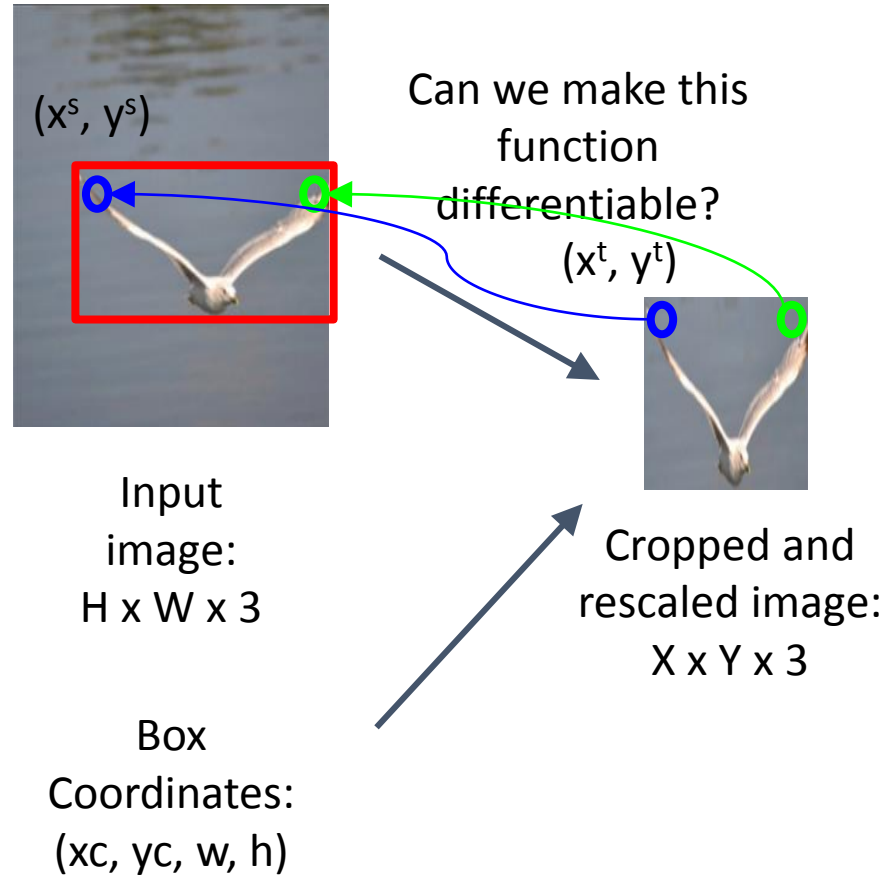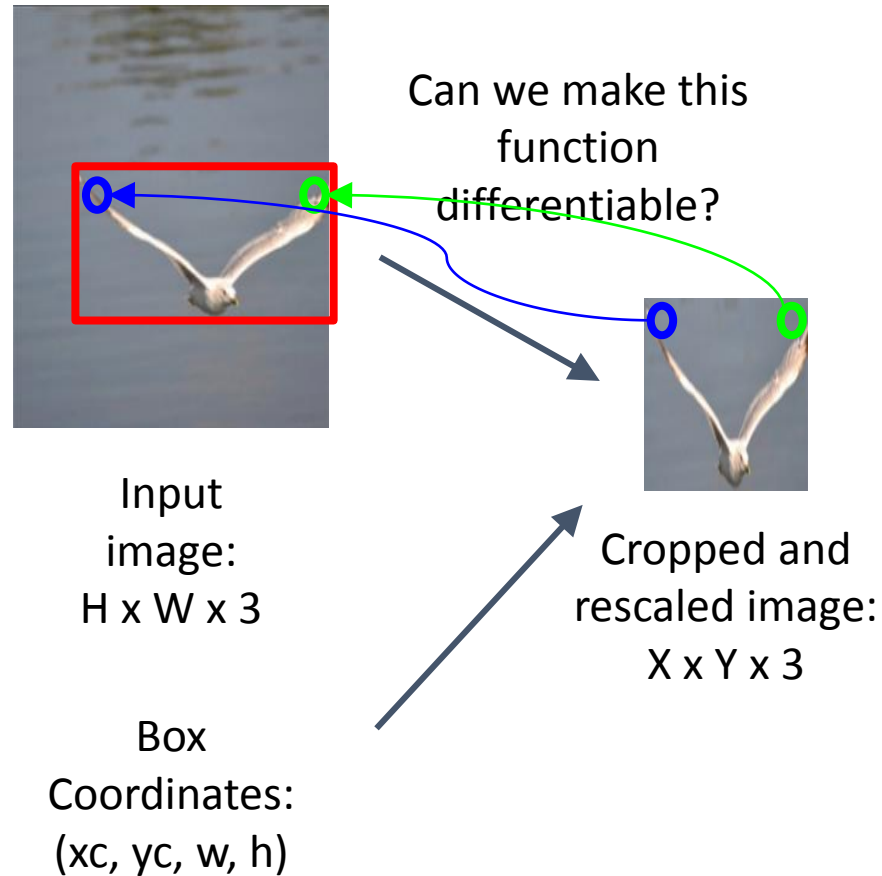
$\mathcal{T}_\theta(G)$

Repeat for all pixels in *output* to get a **sampling grid**

Then use **bilinear interpolation** to compute output

$U$    $V$

# Spatial Transformer Networks

Can we make this function differentiable?

**Idea**: Function mapping *pixel coordinates* (xt, yt) of output to *pixel coordinates* (xs, ys) of input

Network attends to input by predicting $\theta$

Input image: H x W x 3

Box Coordinates: (xc, yc, w, h)

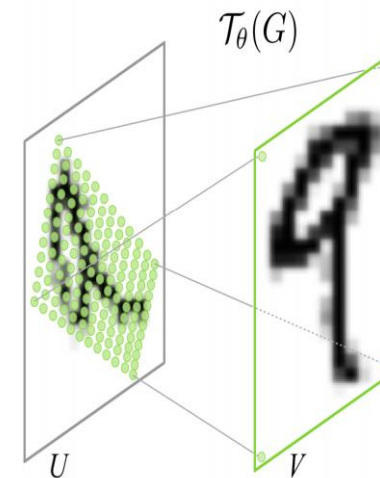Cropped and rescaled image: X x Y x 3

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$
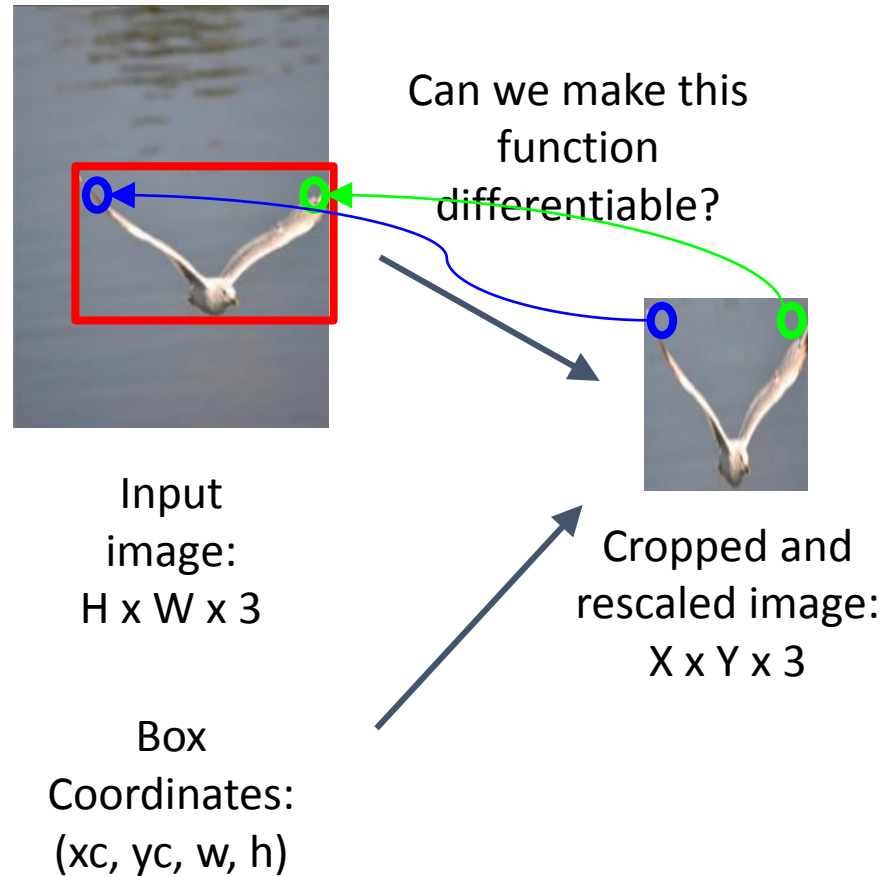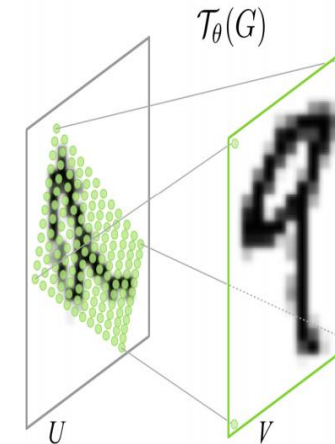
$\mathcal{T}_\theta(G)$

$U$

$V$

Repeat for all pixels in *output* to get a **sampling grid**

Then use **bilinear interpolation** to compute output

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

# Spatial Transformer Networks



**Input:** Full image

**Output:** Region of interest from input

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

# Spatial Transformer Networks



A small **Localization network** predicts transform $\theta$

**Input**: Full image

**Output:** Region of interest from input

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

# Spatial Transformer Networks

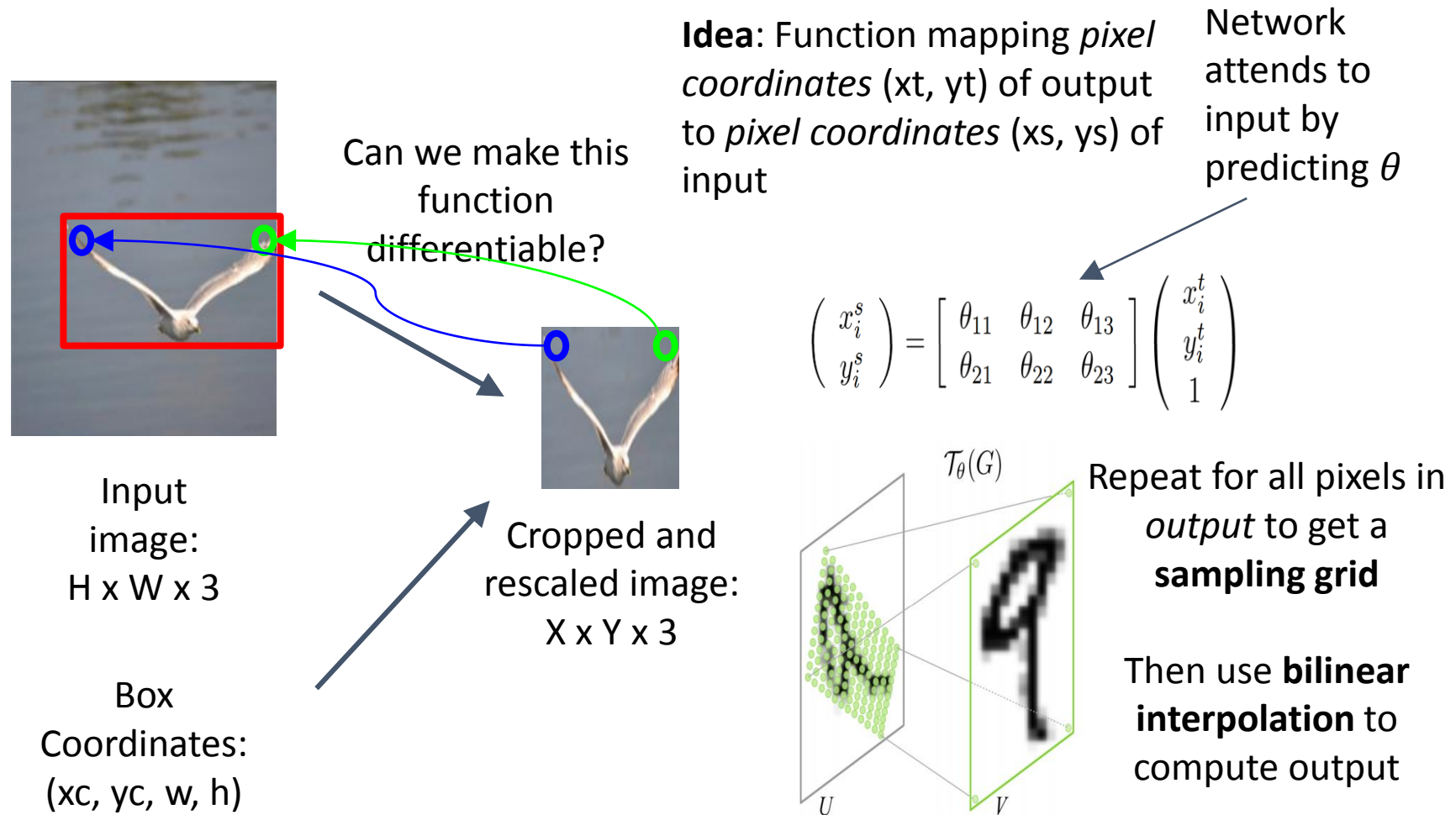**Grid generator** uses $\theta$ to compute sampling grid

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$
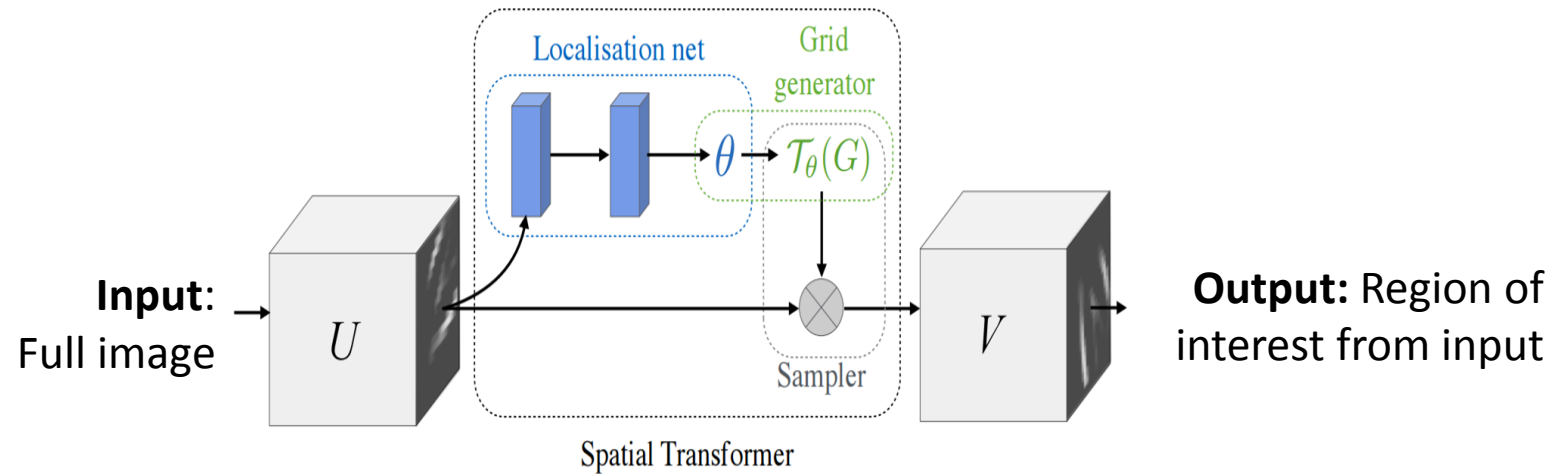
$\mathcal{T}_\theta(G)$

A small **Localization network** predicts transform $\theta$

Localisation net

Grid generator

$\theta \rightarrow \mathcal{T}_\theta(G)$

**Input**: Full image

$U$

Sampler

$V$

**Output:** Region of interest from input

Spatial Transformer

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

$U$

$V$

# Spatial Transformer Networks

**Grid generator** uses $\theta$ to compute sampling grid

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$

A small **Localization network** predicts transform $\theta$

**Input:** Full image

**Output:** Region of interest from input

**Sampler** uses bilinear interpolation to produce output

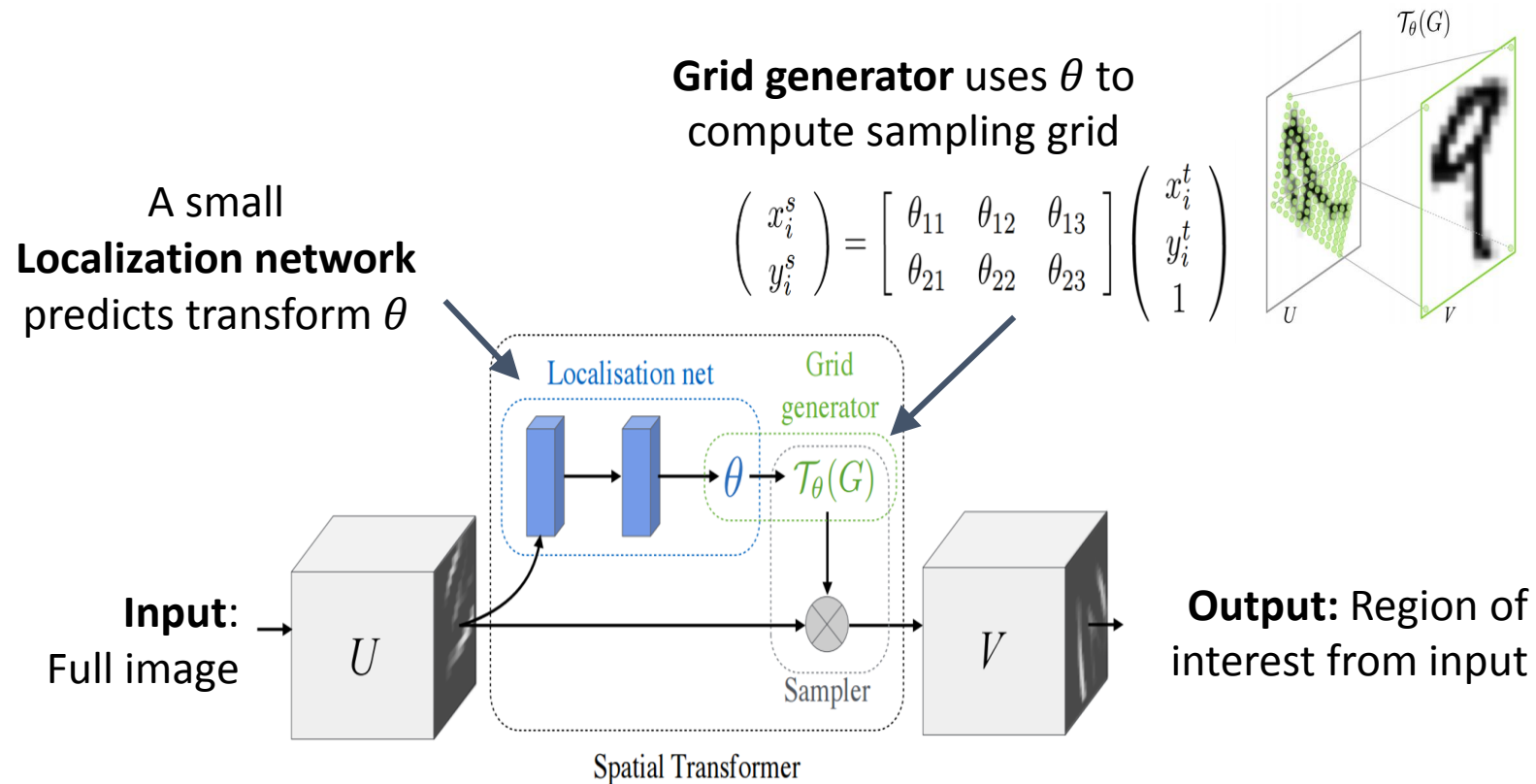$$V_i^c = \sum_n^H \sum_m^W U_{nm}^c \max(0, 1 - |x_i^s - m|) \max(0, 1 - |y_i^s - n|)$$

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

# Spatial Transformer Networks

Insert spatial transformers into a classification network and it learns to attend and transform the input
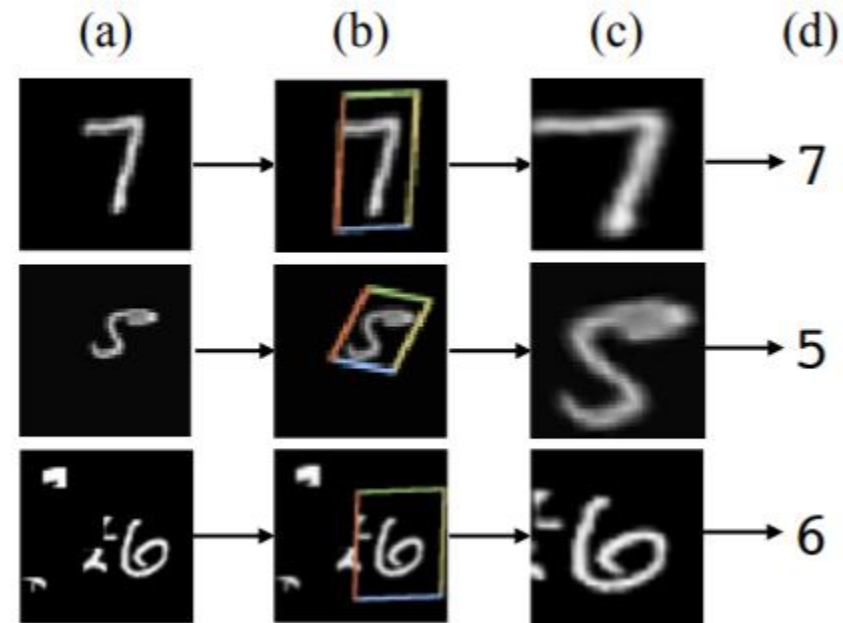
Differentiable "attention / transformation" module



Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

# Attention Takeaways

**Performance:**

- Attention models can ***improve accuracy*** and ***reduce computation*** at the same time.

**Complexity:**

- There are many design choices.
- Those choices have a big effect on performance.
- Ensembling has unusually large benefits.
- Simplify where possible!

# Attention Takeaways

**Explainability:**

- Attention models encode explanations.
- Both locus and trajectory help understand what's going on.

**Hard vs. Soft:**

- Soft models are easier to train, hard models require reinforcement learning.
- They can be combined, as in Luong et al.