



Language
Technologies
Institute

Carnegie
Mellon
University

Multimodal Machine Learning

Lecture 5.1: Multimodal alignment

Louis-Philippe Morency

* Original version co-developed with Tadas Baltrusaitis

Administrative Stuff



Piazza Live Q&A – Reminder

The screenshot displays the Piazza web interface for a class. The browser address bar shows the URL `piazza.com/class/kcncr11wq24q6z7?cid=43`. The page header includes the Piazza logo, course ID `11777-A`, and navigation tabs for `Q & A`, `Resources`, `Statistics`, and `Manage Class`. The user profile for `Louis-Philippe Morency` is visible in the top right.

The left sidebar contains a navigation menu with a `LIVE Q&A` folder highlighted in red. Below it, a `New Post` button is highlighted in orange. The sidebar also lists several posts, including a question about lecture start times, a pinned post about a project preference form, and a course website announcement.

The main content area shows a question titled `question @44` with the text `When is the lecture starting?`. The question is tagged `live_q&a` and has `0` views. It was updated just now by Louis-Philippe Morency. Below the question, the `the instructors' answer` is displayed, stating `At 3:20pm EST`. The answer also has `0` views and was updated just now by Louis-Philippe Morency.

Upcoming Schedule

First project assignment:

- Proposal presentations (Friday 10/9)
- First project reports (Sunday 10/11)

Midterm project assignment

- Midterm presentations (Friday 11/12)
- Midterm reports (Sunday 11/14)

Final project assignment

- Final presentations (Friday 12/11)
- Final reports (Sunday 12/13)

Unimodal Representation Analyses

Main goals:

- Get familiar with unimodal representations
 - Learn about tools based on CNNs, word2vec, ...
- Understand the structure in your unimodal data
 - Perform some visualization of the unimodal data
- Explore qualitatively the unimodal data
 - How does it relate to your labels? Look at specific examples

Examples of unimodal analyses:

- What are the different verbs used in the VQA questions?
- What objects do not get detected? Are they important?
- Visualize face embeddings with respect of emotion labels



Language
Technologies
Institute

Carnegie
Mellon
University

Multimodal Machine Learning

Lecture 5.1: Multimodal alignment

Louis-Philippe Morency

* Original version co-developed with Tadas Baltrusaitis

Lecture objectives

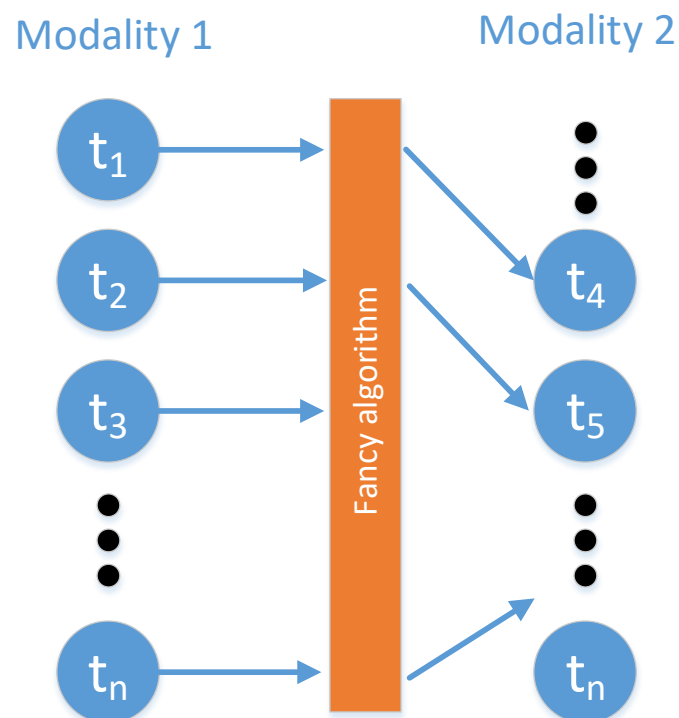
- Multimodal alignment
 - Implicit
 - Explicit
- Explicit signal alignment
 - Dynamic Time Warping
 - Canonical Time Warping
- Attention models in deep learning (implicit and explicit alignment)
 - Soft attention
 - Hard attention
 - Spatial Transformer Networks

Multimodal alignment



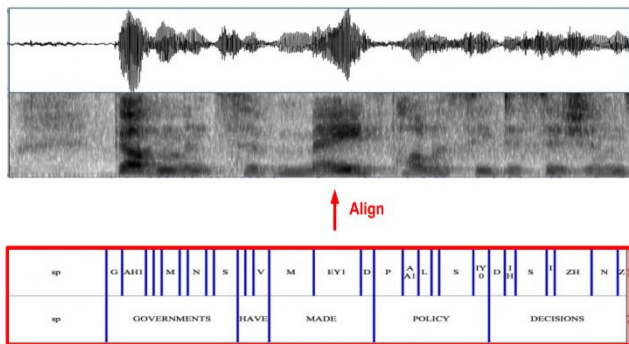
Multimodal-alignment

- Multimodal alignment – finding relationships and correspondences between two or more modalities
- Two types
 - **Explicit** – alignment is the task in itself
 - **Implicit / Latent** – alignment helps when solving a different task (for example “Attention” models)
- Examples ?
 - Images with captions
 - Recipe steps with a how-to video
 - Phrases/words of translated sentences



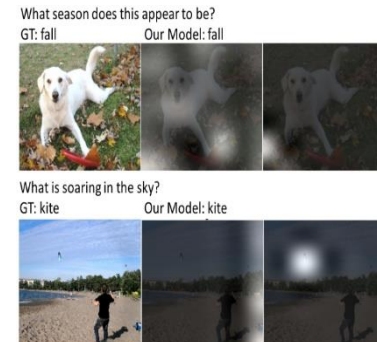
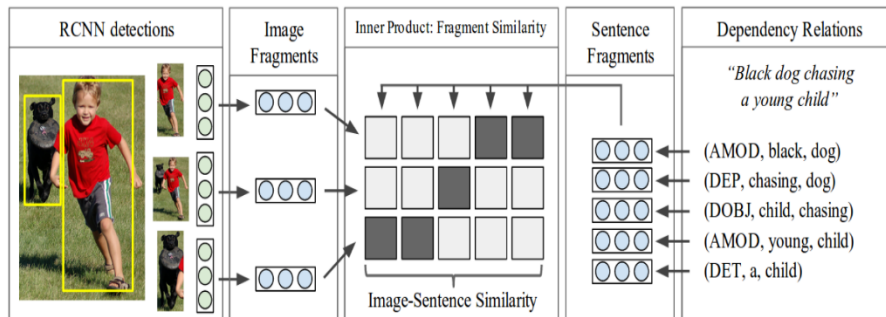
Explicit multimodal-alignment

- Explicit alignment - goal is to find correspondences between modalities
 - Aligning speech signal to a transcript
 - Aligning two out-of sync sequences
 - Co-referring expressions



Implicit multimodal-alignment

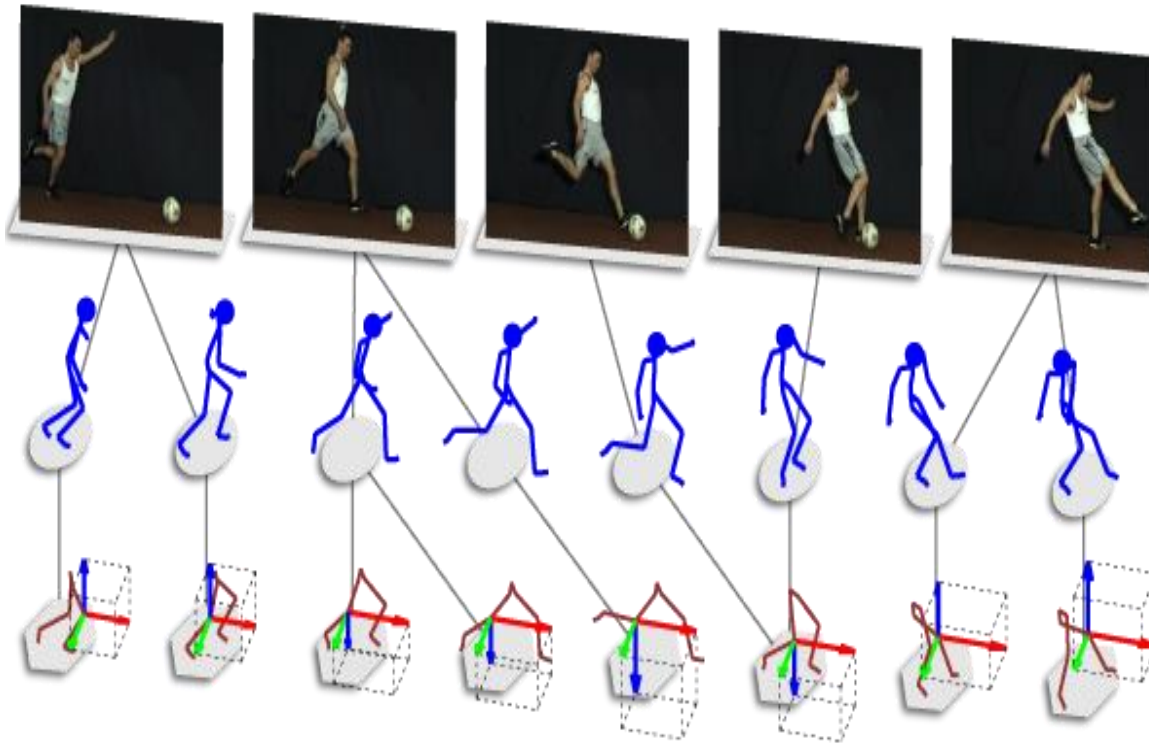
- Implicit alignment - uses internal latent alignment of modalities in order to better solve various problems
 - Machine Translation
 - Cross-modal retrieval
 - Image & Video Captioning
 - Visual Question Answering



Explicit alignment



Temporal sequence alignment



Applications:

- Re-aligning asynchronous data
- Finding similar data across modalities (we can estimate the aligned cost)
- Event reconstruction from multiple sources

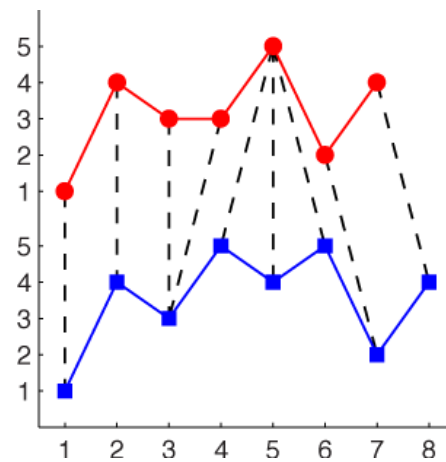
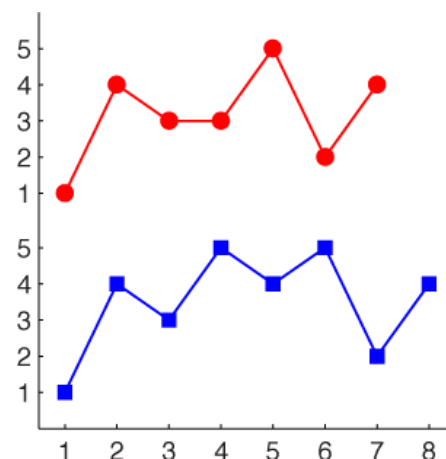


Let's start unimodal – Dynamic Time Warping

- We have two unaligned temporal unimodal signals
 - $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n_x}] \in \mathbb{R}^{d \times n_x}$
 - $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{n_y}] \in \mathbb{R}^{d \times n_y}$
- Find set of indices to minimize the alignment difference:

$$L(\mathbf{p}^x, \mathbf{p}^y) = \sum_{t=1}^l \left\| \mathbf{x}_{p_t^x} - \mathbf{y}_{p_t^y} \right\|_2^2$$

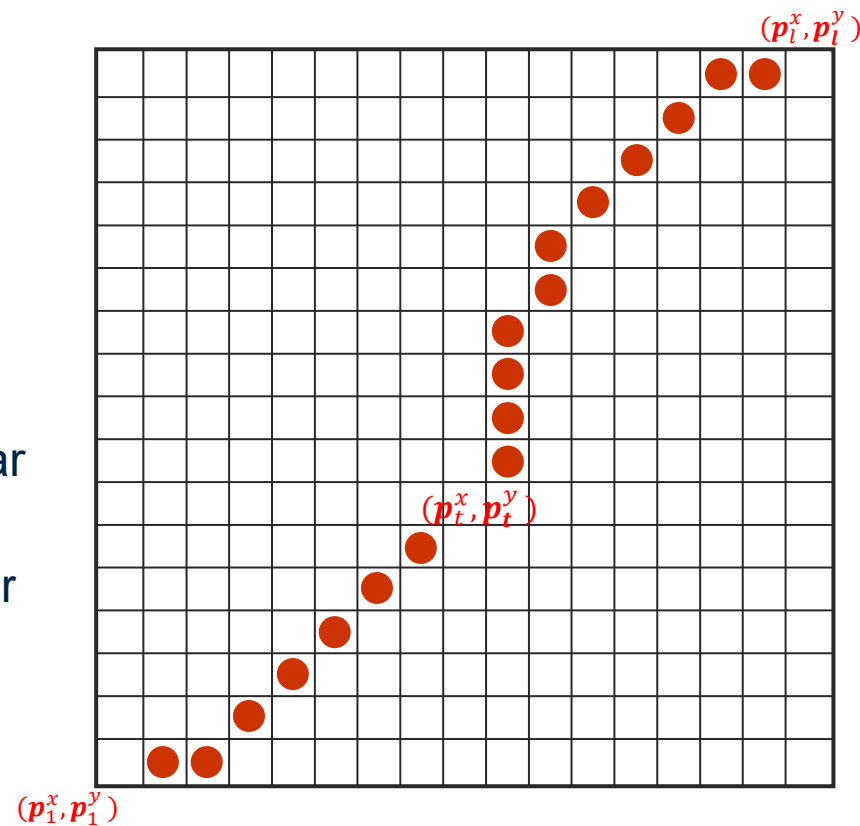
- Where \mathbf{p}^x and \mathbf{p}^y are index vectors of same length
- Dynamic Time Warping is designed to find these index vectors



Dynamic Time Warping continued

Lowest cost path in a cost matrix

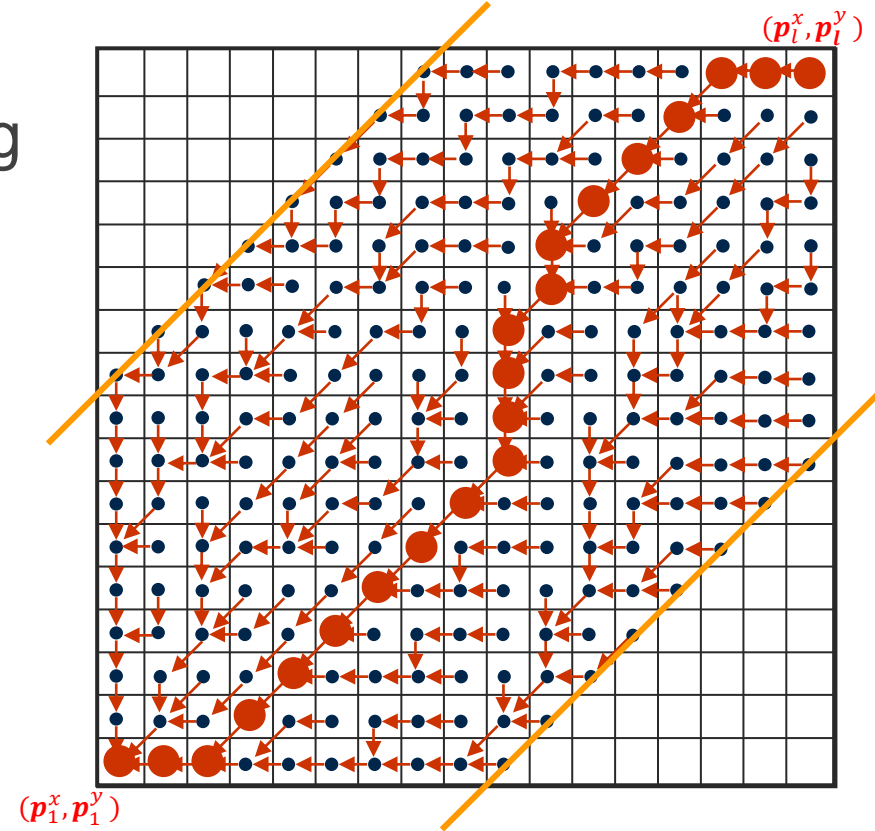
- Restrictions?
 - Monotonicity – no going back in time
 - Continuity - no gaps
 - Boundary conditions - start and end at the same points
 - Warping window - don't get too far from diagonal
 - Slope constraint – do not insert or skip too much



Dynamic Time Warping continued

Lowest cost path in a cost matrix

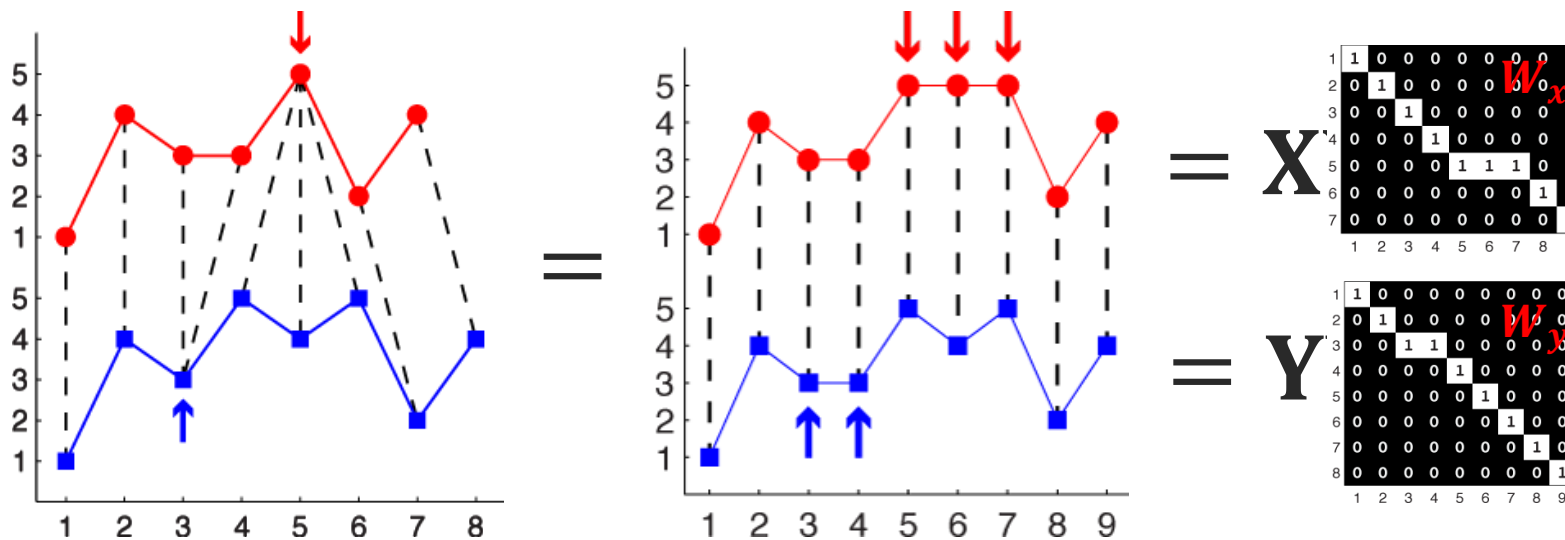
- Solved using dynamic programming while respecting the restrictions



DTW alternative formulation

$$L(\mathbf{p}^x, \mathbf{p}^y) = \sum_{t=1}^l \left\| \mathbf{x}_{p_t^x} - \mathbf{y}_{p_t^y} \right\|_2^2$$

Replication doesn't change the objective!



Alternative objective:

$$L(\mathbf{W}_x, \mathbf{W}_y) = \left\| \mathbf{X}\mathbf{W}_x - \mathbf{Y}\mathbf{W}_y \right\|_F^2$$

\mathbf{X}, \mathbf{Y} – original signals (same #rows, possibly different #columns)

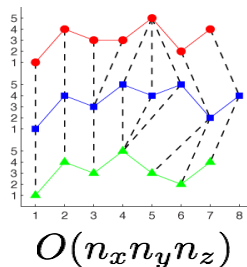
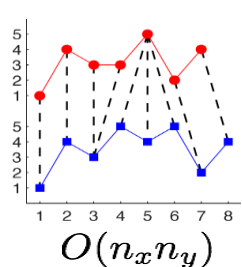
$\mathbf{W}_x, \mathbf{W}_y$ - alignment matrices

Frobenius norm $\|\mathbf{A}\|_F^2 = \sum_i \sum_j |a_{i,j}|^2$



DTW – Some Limitations

- Computationally complex



m sequences

$$O\left(\prod_{i=1}^m n_i\right)$$

- Sensitive to outliers

- Unimodal!



Canonical Correlation Analysis reminder

maximize: $tr(U^T \Sigma_{XY} V)$

subject to: $U^T \Sigma_{YY} U = V^T \Sigma_{YY} V = I$, $u_{(j)}^T \Sigma_{XY} v_{(i)} = 0$ for $i \neq j$

1

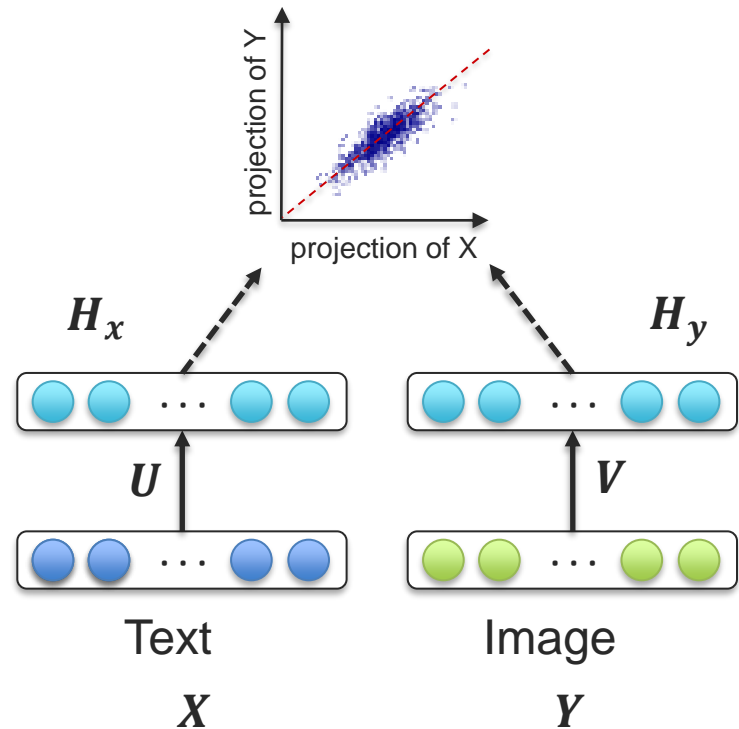
Linear projections maximizing correlation

2

Orthogonal projections

3

Unit variance of the projection vectors

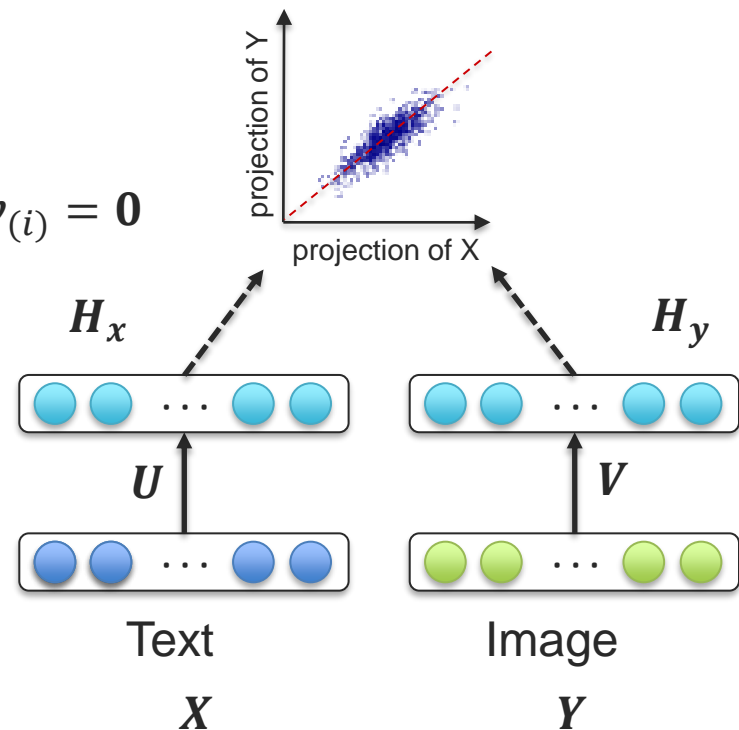


Canonical Correlation Analysis reminder

- When data is normalized it is actually equivalent to smallest RMSE reconstruction
- CCA loss can also be re-written as:

$$L(\mathbf{U}, \mathbf{V}) = \|\mathbf{U}^T \mathbf{X} - \mathbf{V}^T \mathbf{Y}\|_F^2$$

subject to: $\mathbf{U}^T \Sigma_{YY} \mathbf{U} = \mathbf{V}^T \Sigma_{YY} \mathbf{V} = \mathbf{I}, \mathbf{u}_{(j)}^T \Sigma_{XY} \mathbf{v}_{(i)} = 0$



Canonical Time Warping

- Dynamic Time Warping + Canonical Correlation Analysis = Canonical Time Warping

$$L(\mathbf{U}, \mathbf{V}, \mathbf{W}_x, \mathbf{W}_y) = \|\mathbf{U}^T \mathbf{X} \mathbf{W}_x - \mathbf{V}^T \mathbf{Y} \mathbf{W}_y\|_F^2$$

- Allows to align multi-modal or multi-view (same modality but from a different point of view)
- $\mathbf{W}_x, \mathbf{W}_y$ – temporal alignment
- \mathbf{U}, \mathbf{V} – cross-modal (spatial) alignment

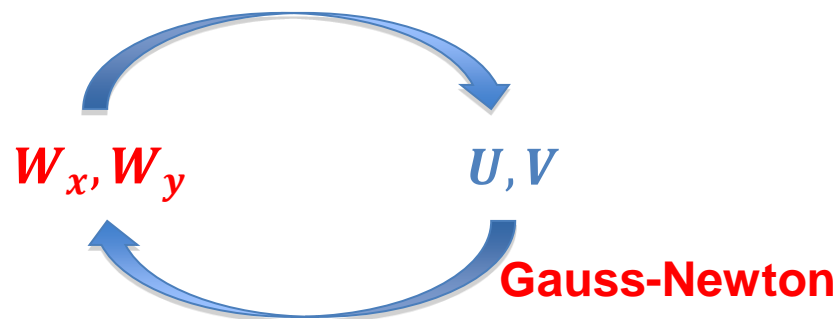
[Canonical Time Warping for Alignment of Human Behavior, Zhou and De la Torre, 2009]

Canonical Time Warping

$$L(\mathbf{U}, \mathbf{V}, \mathbf{W}_x, \mathbf{W}_y) = \|\mathbf{U}^T \mathbf{X} \mathbf{W}_x - \mathbf{V}^T \mathbf{Y} \mathbf{W}_y\|_F^2$$

Optimized by Coordinate-descent – fix one set of parameters, optimize another

Generalized Eigen-decomposition



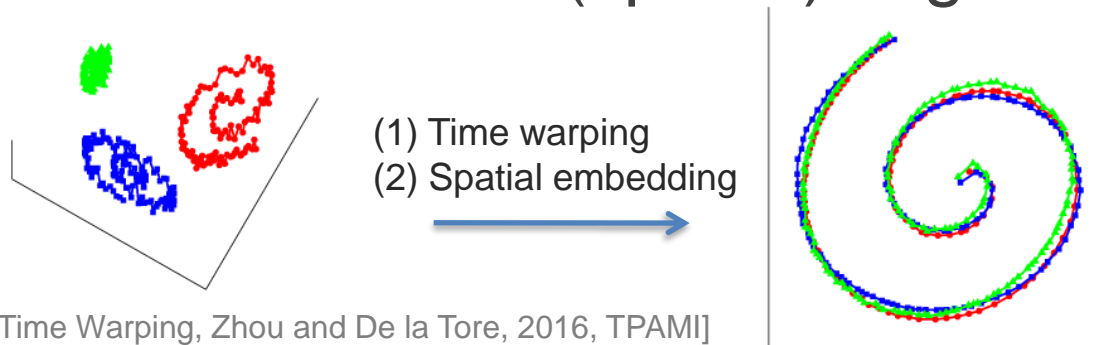
[Canonical Time Warping for Alignment of Human Behavior, Zhou and De la Tore, 2009, NIPS]

Generalized Time warping

- Generalize to multiple sequences all of different modality

$$L(\mathbf{U}_i, \mathbf{W}_i) = \sum_{i=1} \sum_{j=1} \|\mathbf{U}_i^T \mathbf{x}_i \mathbf{W}_i - \mathbf{U}_j^T \mathbf{x}_j \mathbf{W}_j\|_F^2$$

- \mathbf{W}_i – set of temporal alignments
- \mathbf{U}_i – set of cross-modal (spatial) alignments



[Generalized Canonical Time Warping, Zhou and De la Torre, 2016, TPAMI]

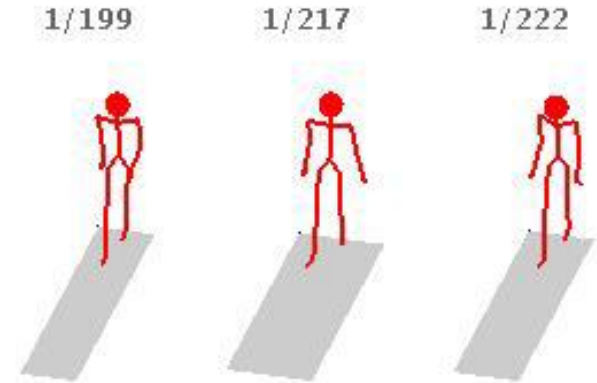
Alignment examples (unimodal)

CMU Motion Capture

Subject 1: 199 frames

Subject 2: 217 frames

Subject 3: 222 frames

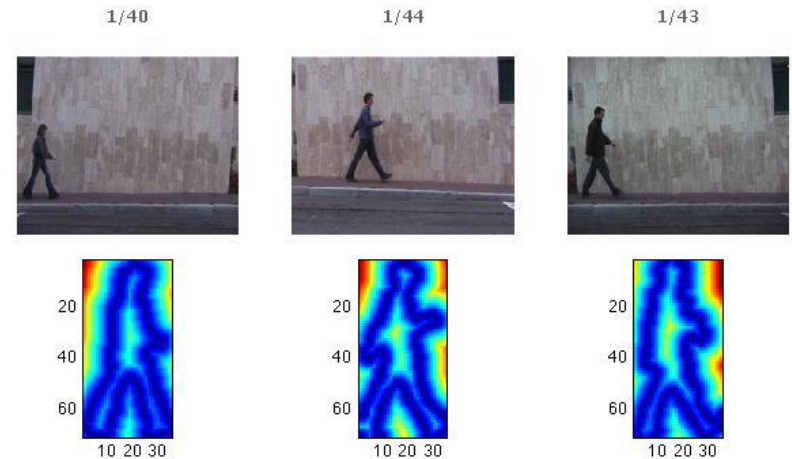


Weizmann

Subject 1: 40 frames

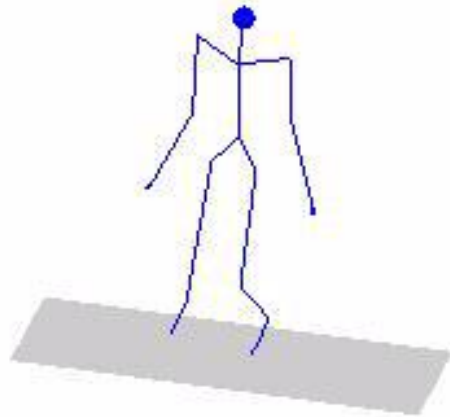
Subject 2: 44 frames

Subject 3: 43 frames



Alignment examples (multimodal)

1/273



1/51



1/127



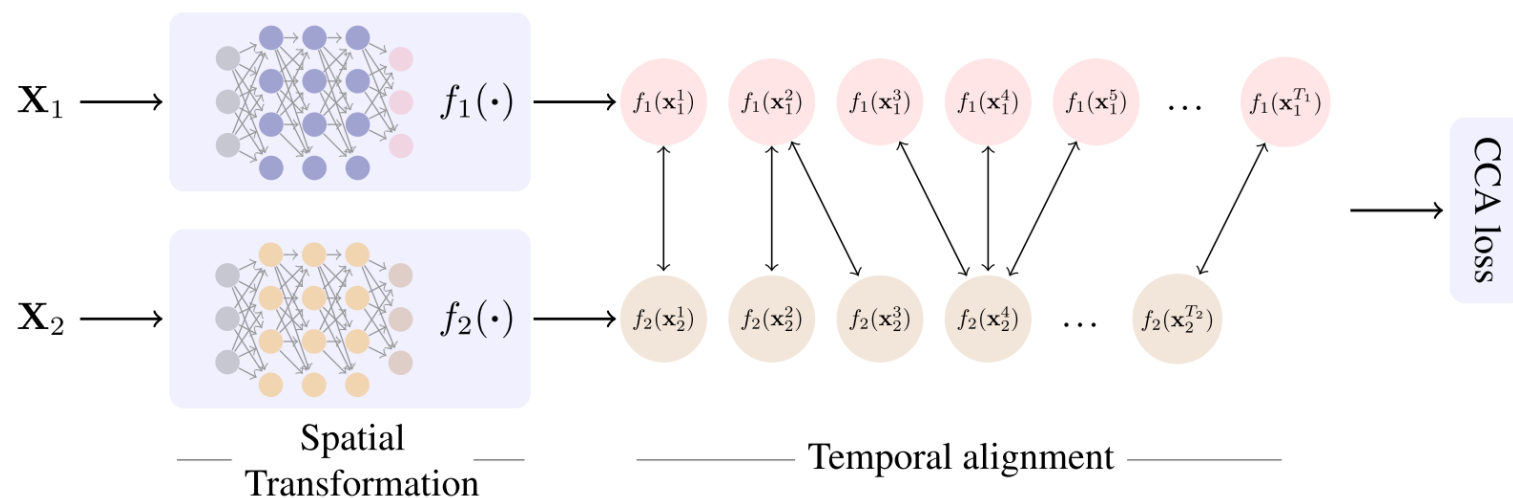
Canonical time warping - limitations

- Linear transform between modalities
- How to address this?

Deep Canonical Time Warping

$$L(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \mathbf{W}_x, \mathbf{W}_y) = \left\| f_{\boldsymbol{\theta}_1}(\mathbf{X})\mathbf{W}_x - f_{\boldsymbol{\theta}_1}(\mathbf{Y})\mathbf{W}_y \right\|_F^2$$

- Could be seen as generalization of DCCA and GTW



[Deep Canonical Time Warping, Trigeorgis et al., 2016, CVPR]

Deep Canonical Time Warping

$$L(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \mathbf{W}_x, \mathbf{W}_y) = \left\| f_{\boldsymbol{\theta}_1}(\mathbf{X})\mathbf{W}_x - f_{\boldsymbol{\theta}_1}(\mathbf{Y})\mathbf{W}_y \right\|_F^2$$

- The projections are orthogonal (like in DCCA)
- Optimization is again iterative:
 - Solve for alignment $(\mathbf{W}_x, \mathbf{W}_y)$ with fixed projections $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$
 - Eigen decomposition
 - Solve for projections $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$ with fixed alignment $(\mathbf{W}_x, \mathbf{W}_y)$
 - Gradient descent
- Repeat till convergence

[Deep Canonical Time Warping, Trigeorgis et al., 2016, CVPR]

Implicit alignment



Implicit alignment

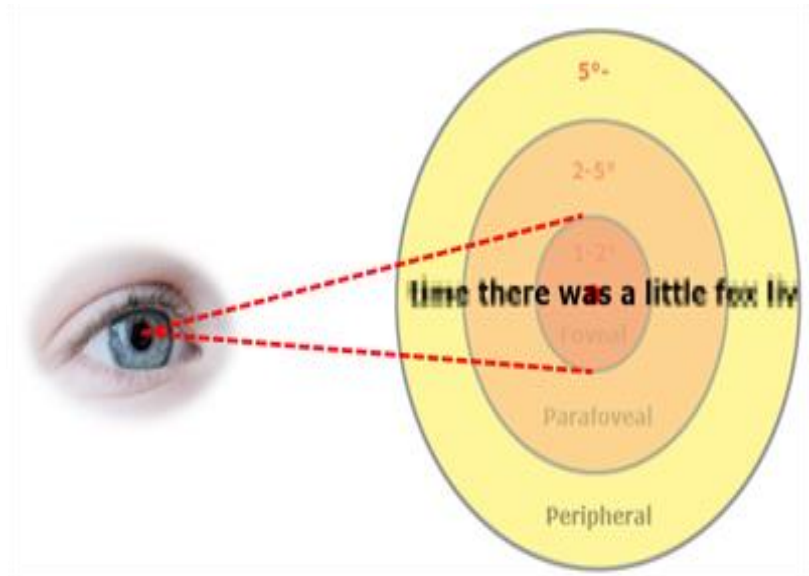
- We looked how to explicitly align temporal data
- Could use that as an internal (hidden) step in our models?
- Can we instead encourage the model to align data when solving a different problem?
- Yes!
 - Graphical models
 - Neural attention models (focus of today's lecture)

Attention models



Attention in humans

- Foveal vision – we only see in “high resolution” in 2 degrees of vision
- We focus our attention selectively to certain words (for example our names)
- We attend to relevant speech in a noisy room



Attention models in deep learning

- Many examples of attention models in recent years!
- Why:
 - Allows for implicit data alignment
 - Good results empirically
 - In some cases faster (don't need to focus on all the image)
 - Better Interpretability

Types of Attention Models

- Recent attention models can be roughly split into three major categories
 1. Soft attention
 - Acts like a gate function. Deterministic inference.
 2. Transform network
 - Warp the input to better align with canonical view
 3. Hard attention
 - Includes stochastic processes. Related to reinforcement learning.

Soft attention



Machine Translation

Given a sentence in one language translate it to another

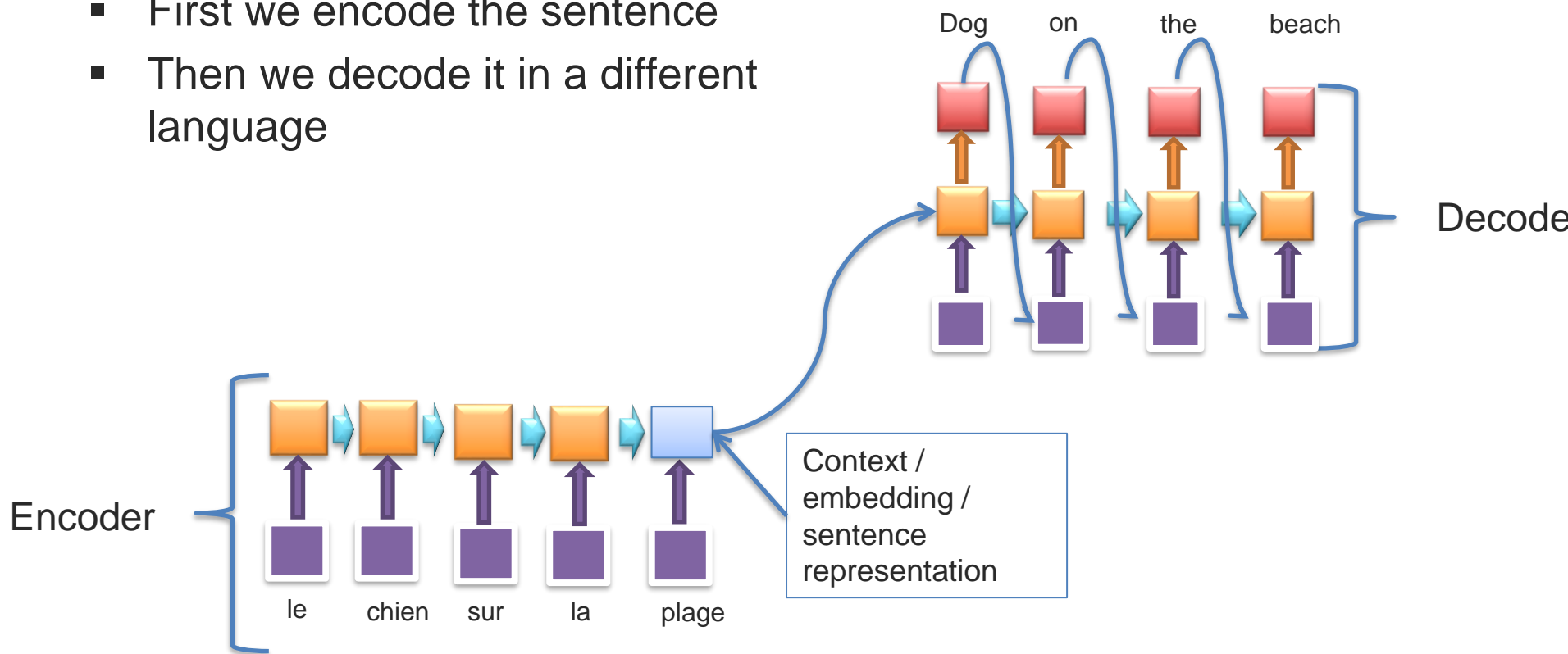
Dog on the beach → le chien sur la plage

- Not exactly multimodal task – but a good start! Each language can be seen almost as a modality.

Machine Translation with RNNs

A quick reminder about encoder decoder frameworks

- First we encode the sentence
- Then we decode it in a different language



Machine Translation with RNNs

What is the problem with this?

What happens when the sentences are very long?

- We expect the encoders hidden state to capture everything in a sentence, a very complex state in a single vector, such as

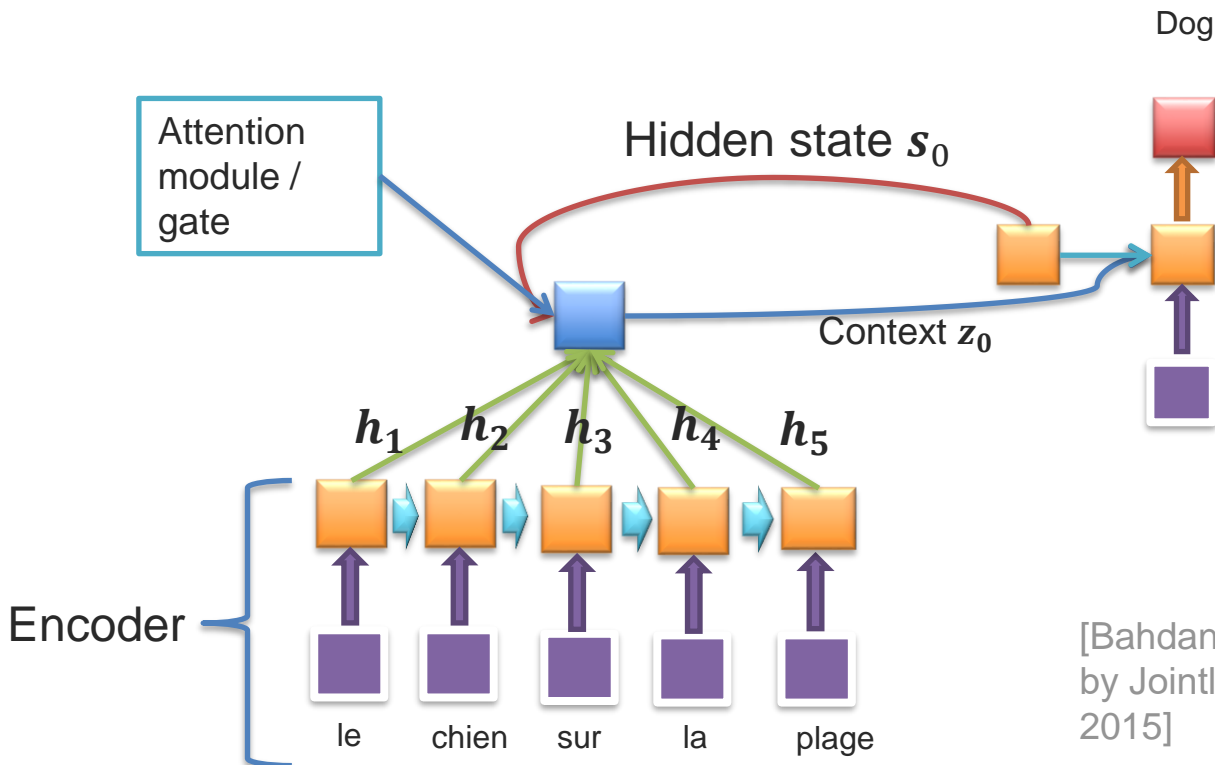
The agreement on the European Economic Area was signed in August 1992.



L' accord sur la zone économique européenne a été signé en août 1992.

Decoder – attention model

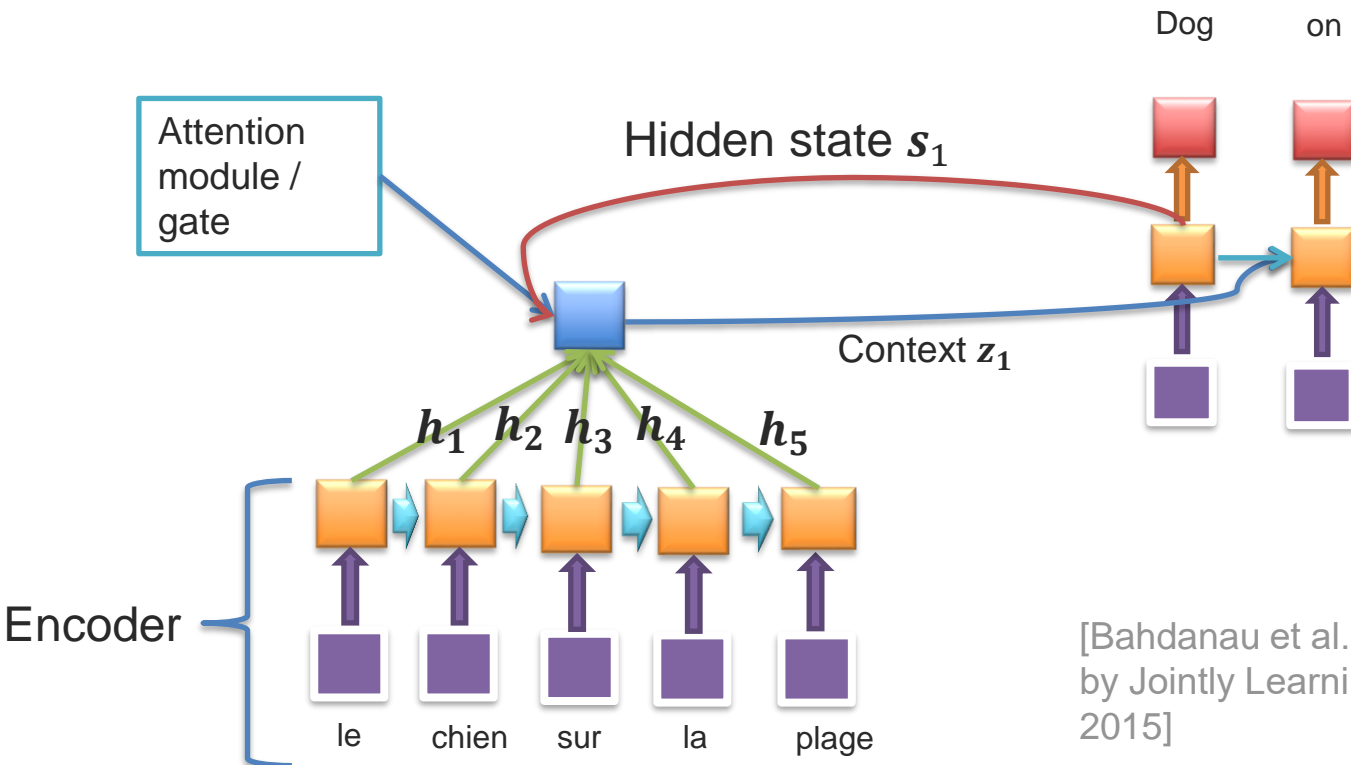
- Before encoder would just take the final hidden state, now we actually care about the intermediate hidden states



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

Decoder – attention model

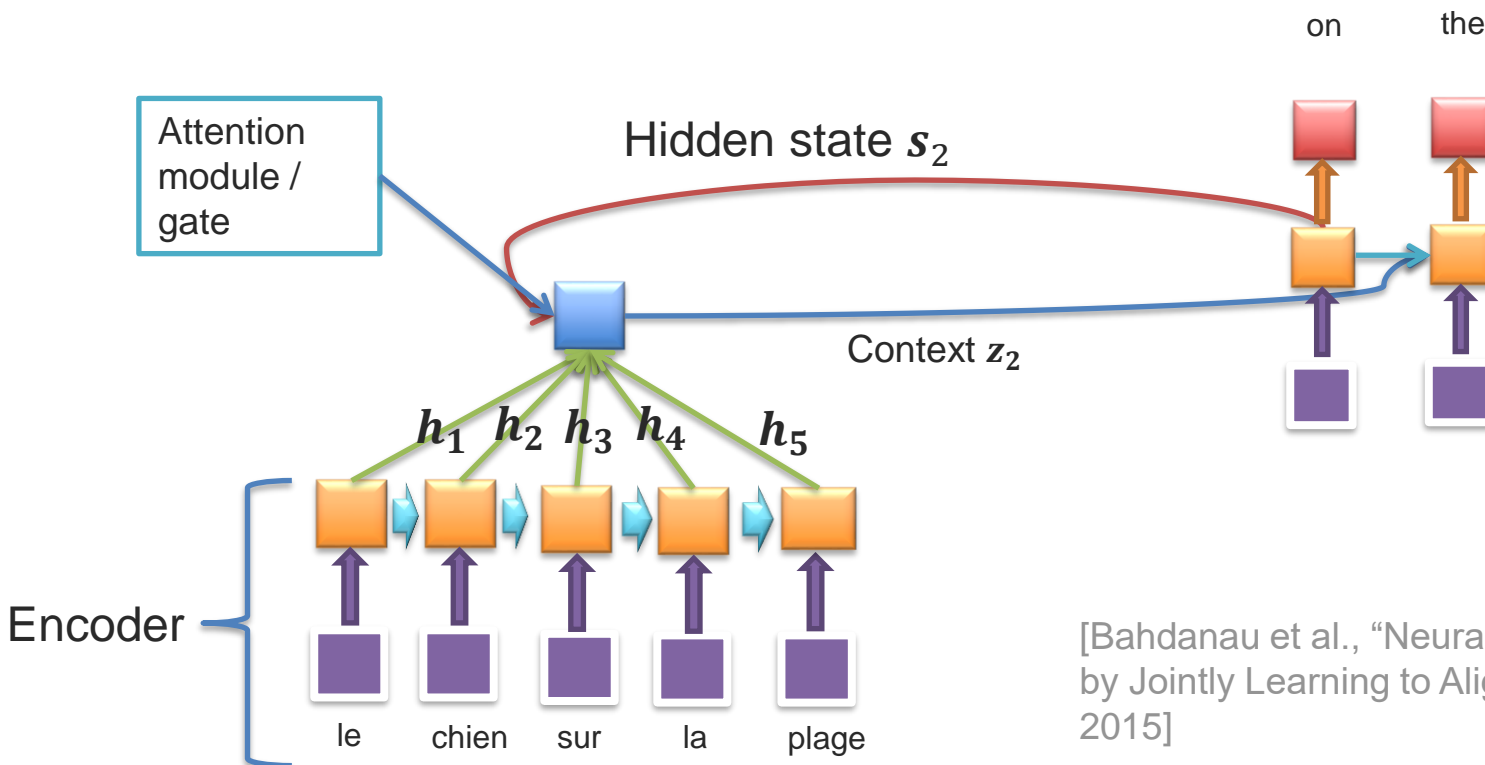
- Before encoder would just take the final hidden state, now we actually care about the intermediate hidden states



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

Decoder – attention model

- Before encoder would just take the final hidden state, now we actually care about the intermediate hidden states



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

How do we encode attention?

Before:

$p(y_i | y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, \mathbf{s}_i, \mathbf{z})$, where $\mathbf{z} = \mathbf{h}_T$, and \mathbf{s}_i - the current state of the decoder

Now:

$$p(y_i | y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, \mathbf{s}_i, \mathbf{z}_i)$$

- Have an attention “gate”
 - A different context \mathbf{z}_i used at each time step!
 - $\mathbf{z}_i = \sum_{j=i}^{T_x} \alpha_{ij} \mathbf{h}_j$

α_{ij} is the (scalar) attention for word j at generation step i

MT with attention

So how do we determine α_{ij} ,

- $\alpha_{i,j} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$ - softmax, making sure they sum to 1

where:

- $e_{ij} = \mathbf{v}^T \sigma(W \mathbf{s}_{i-1} + U \mathbf{h}_j)$

a feedforward network that can tell us given the current state of decoder how important the current encoding is now

\mathbf{v} , W , U – learnable weights

$$\mathbf{z}_i = \sum_{j=1}^{T_x} \alpha_{ij} \mathbf{h}_j$$

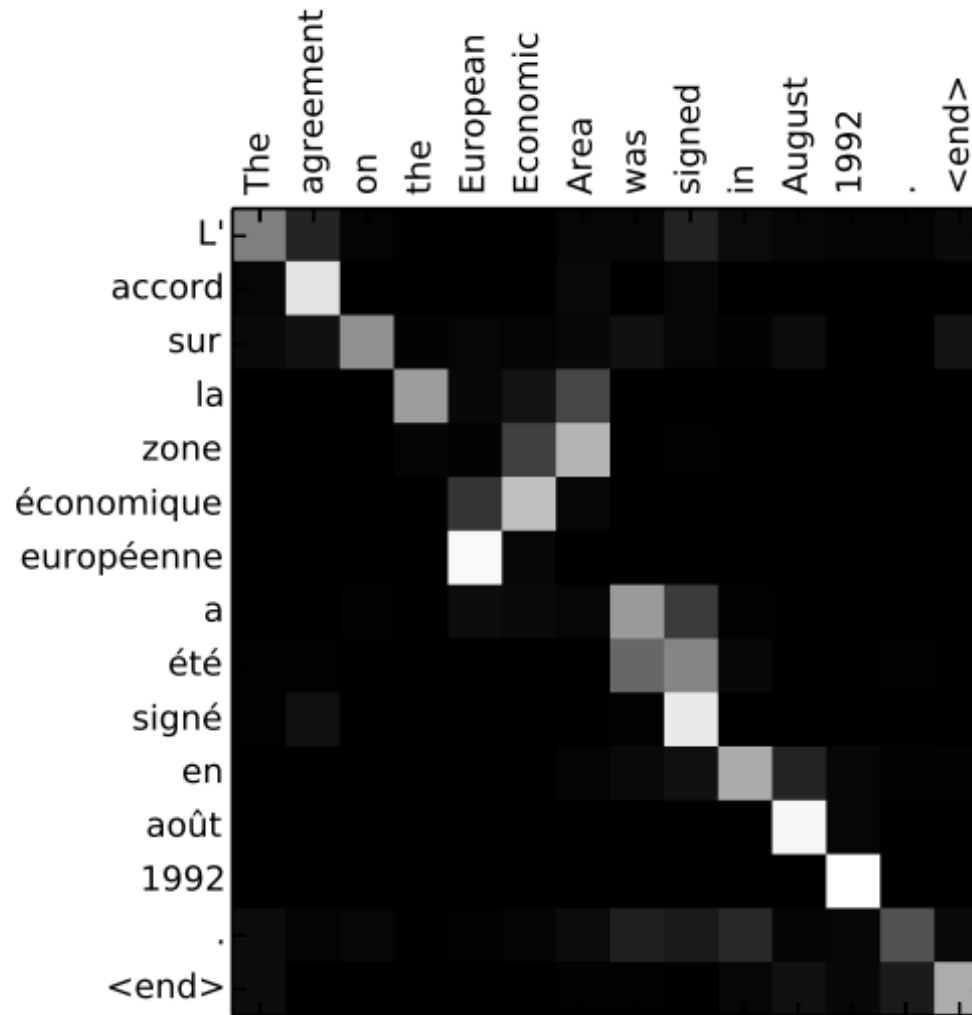
← expectation of the context (a fancy way to say it's a weighted average)

MT with attention

Basically we are using a neural network to tell us where a neural network should be looking!

- We can use with RNN, LSTM or GRU
- Encoder being used is the same structure as before
 - Can use uni-directional
 - Can use bi-directional
- Model can be trained using our regular back-propagation through time, all of the modules are differentiable

Does it work?



MT with attention recap

- It gives good translation results (especially for long sentences)
- Also get a (soft) alignment of sentences in different languages
 - Extra interpretability of method functioning

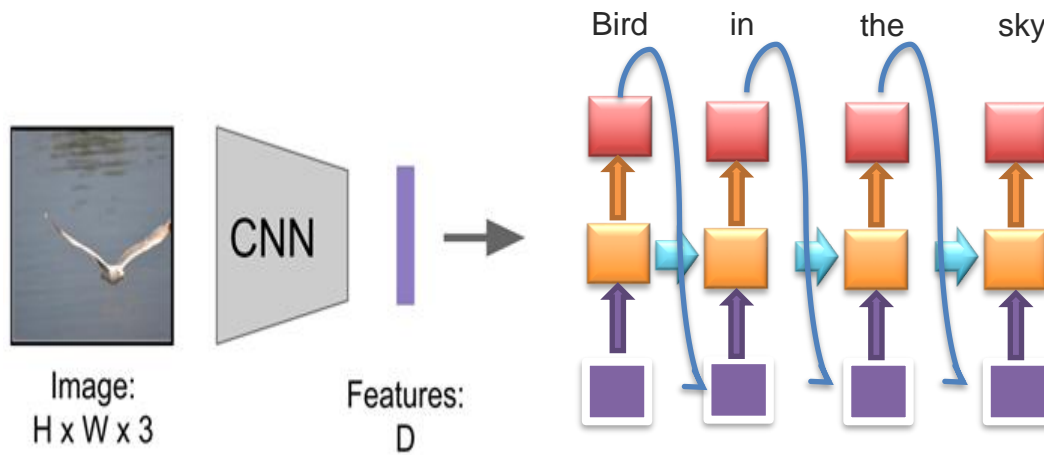
How do we move to multimodal?

Visual captioning with soft attention



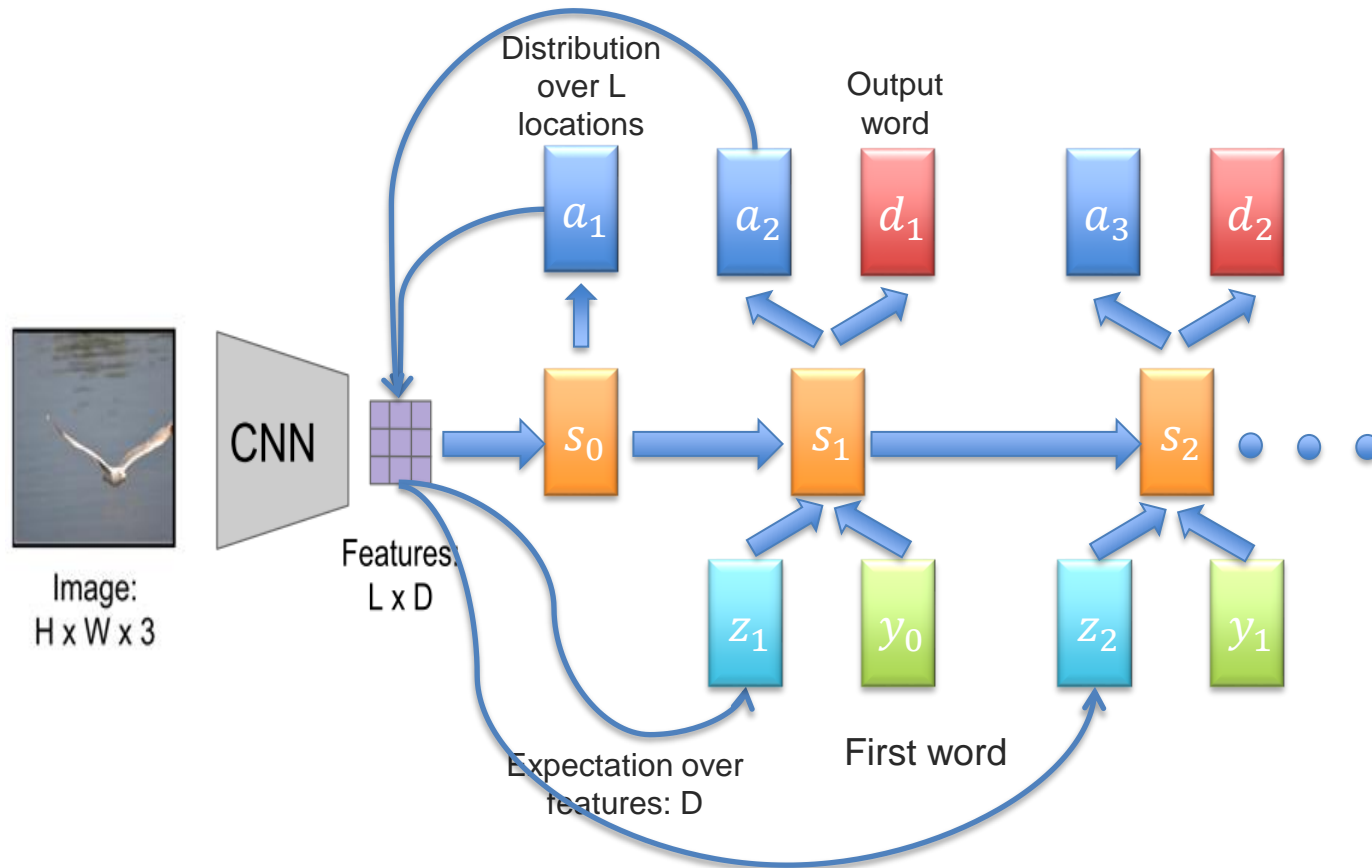
[Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, Xu et al., 2015]

Recap RNN for Captioning



Why not using final layer of the CNN?

Looking at more fine grained features



Soft attention

- Allows for latent data alignment
- Allows us to get an idea of what the network “sees”
- Can be optimized using back propagation

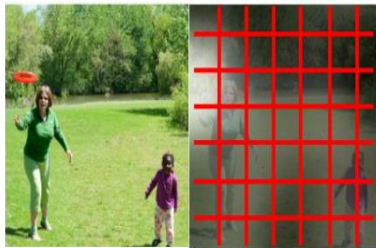
- Good at paper naming!
 - Show, Attend and Tell (extension of Show and Tell)
 - Listen, Attend and Walk
 - Listen, Attend and Spell
 - Ask, Attend and Answer

Spatial Transformer networks



Some limitations of grid-based attention

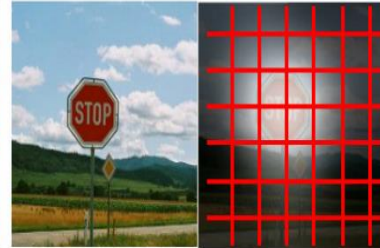
Can we fixate on small parts of image but still have easy end-to-end training?



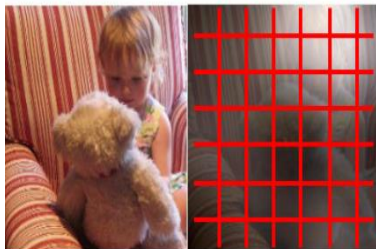
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



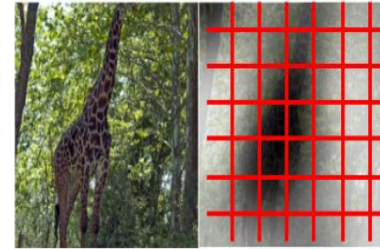
A stop sign is on a road with a mountain in the background.



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

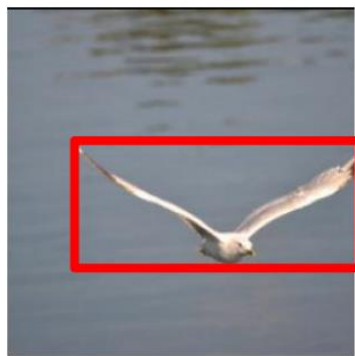


A group of people sitting on a boat in the water.



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

Spatial Transformer Networks



Input image:
 $H \times W \times 3$

Box Coordinates:
 (x_c, y_c, w, h)

Can we make this
function differentiable?

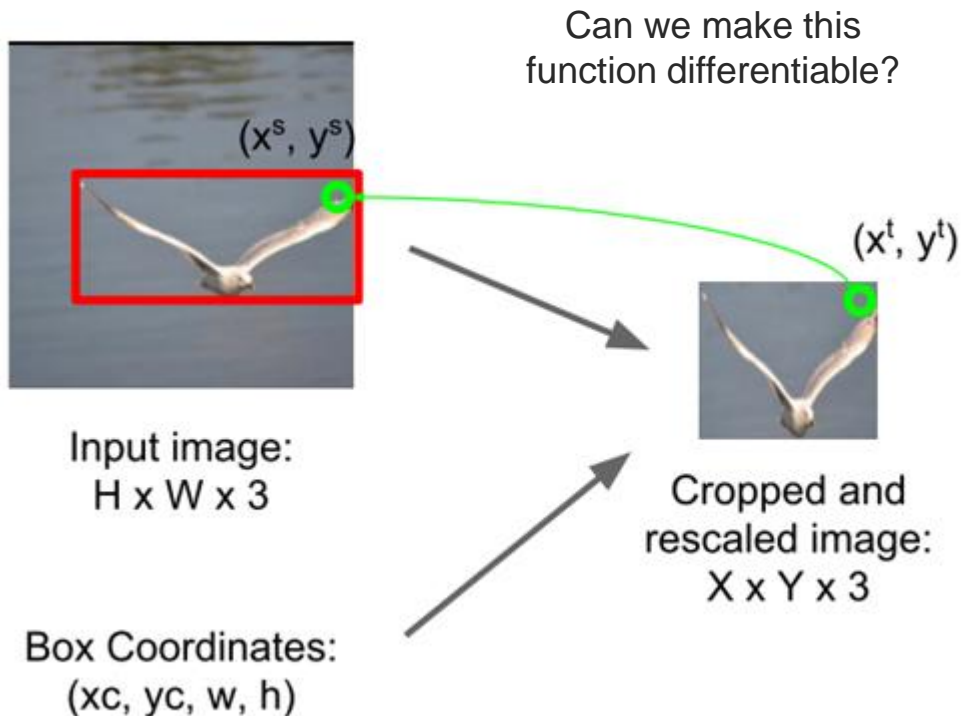


Cropped and
rescaled image:
 $X \times Y \times 3$



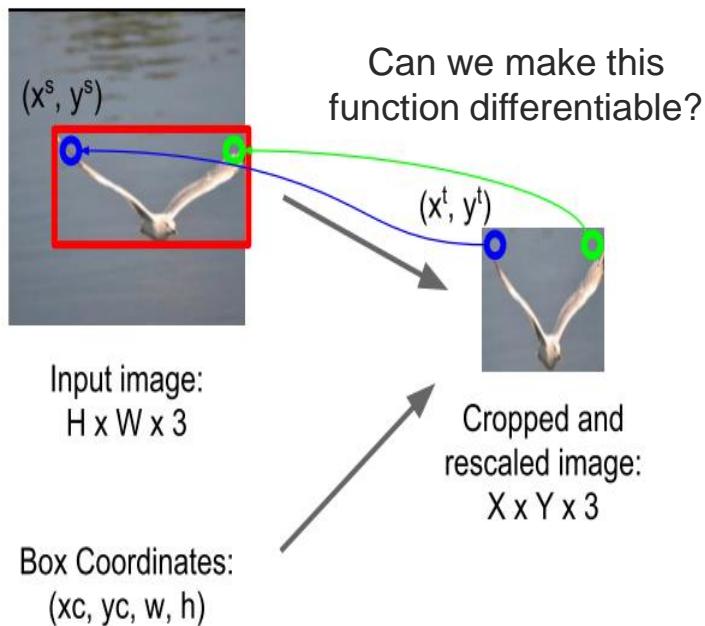
Spatial Transformer Networks

Idea: Function mapping pixel coordinates (x^t, y^t) of output to pixel coordinates (x^s, y^s) of input



$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{1,1} & \theta_{1,2} & \theta_{1,3} \\ \theta_{2,1} & \theta_{2,2} & \theta_{2,3} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$

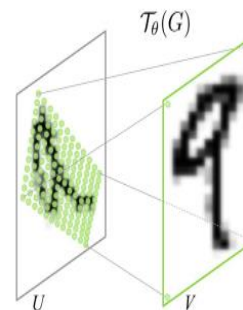
Spatial Transformer Networks



Idea: Function mapping pixel coordinates (x^t, y^t) of output to pixel coordinates (x^s, y^s) of input

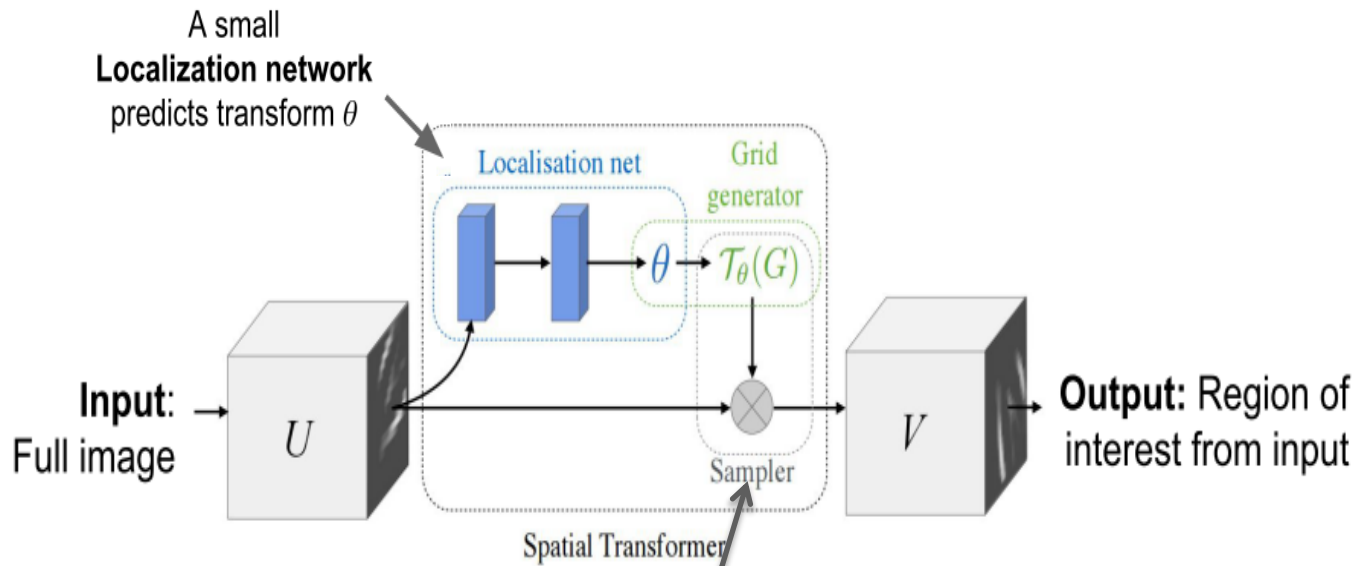
$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{1,1} & \theta_{1,2} & \theta_{1,3} \\ \theta_{2,1} & \theta_{2,2} & \theta_{2,3} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$

Network "attends" to input by predicting θ



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

Spatial Transformer Networks



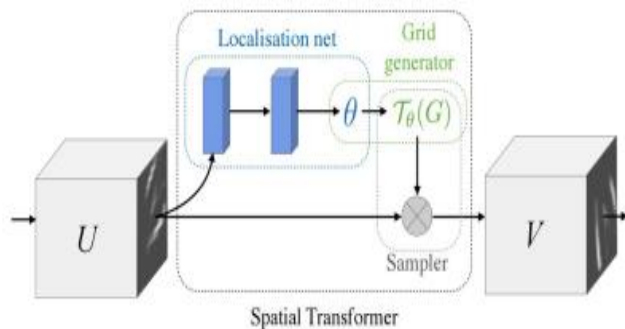
Sampler uses bilinear interpolation to produce output

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

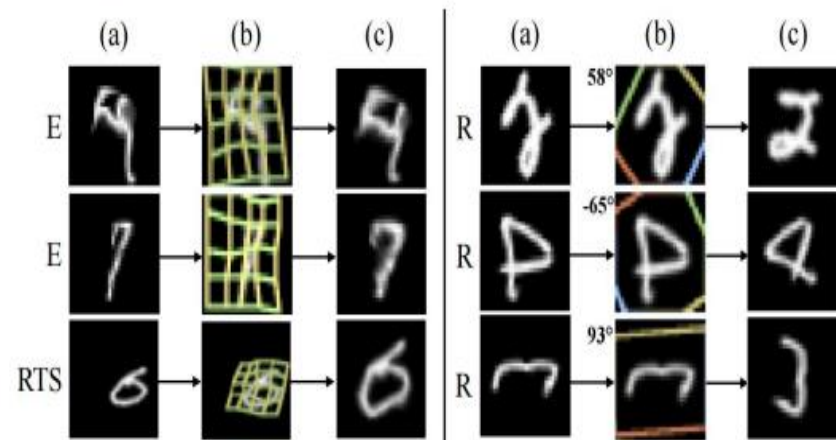


Spatial Transformer Networks

Differentiable “attention / transformation” module



Insert spatial transformers into a classification network and it learns to attend and transform the input



Examples on real world data

Results on traffic sign recognition



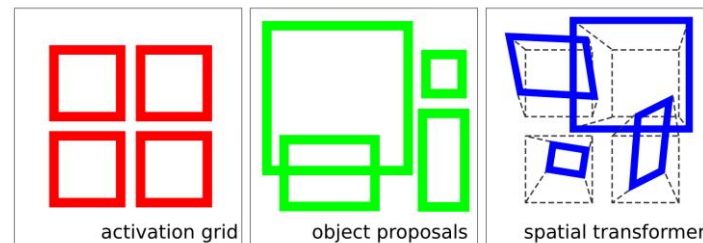
Code available http://torch.ch/blog/2015/09/07/spatial_transformers.html

Recap on Spatial Transformer Networks

- Differentiable so we can just use back-prop for training end-to-end
- Can be used with complex transformations to focus on an image
 - Affine and Piece-Wise Affine, Perspective, Thin Plate Splines
- We can use it instead of grid based soft and hard attention for multi-modal tasks



A **man** is flying a **kite** on a sandy **beach**.



Glimpse Network (Hard Attention)

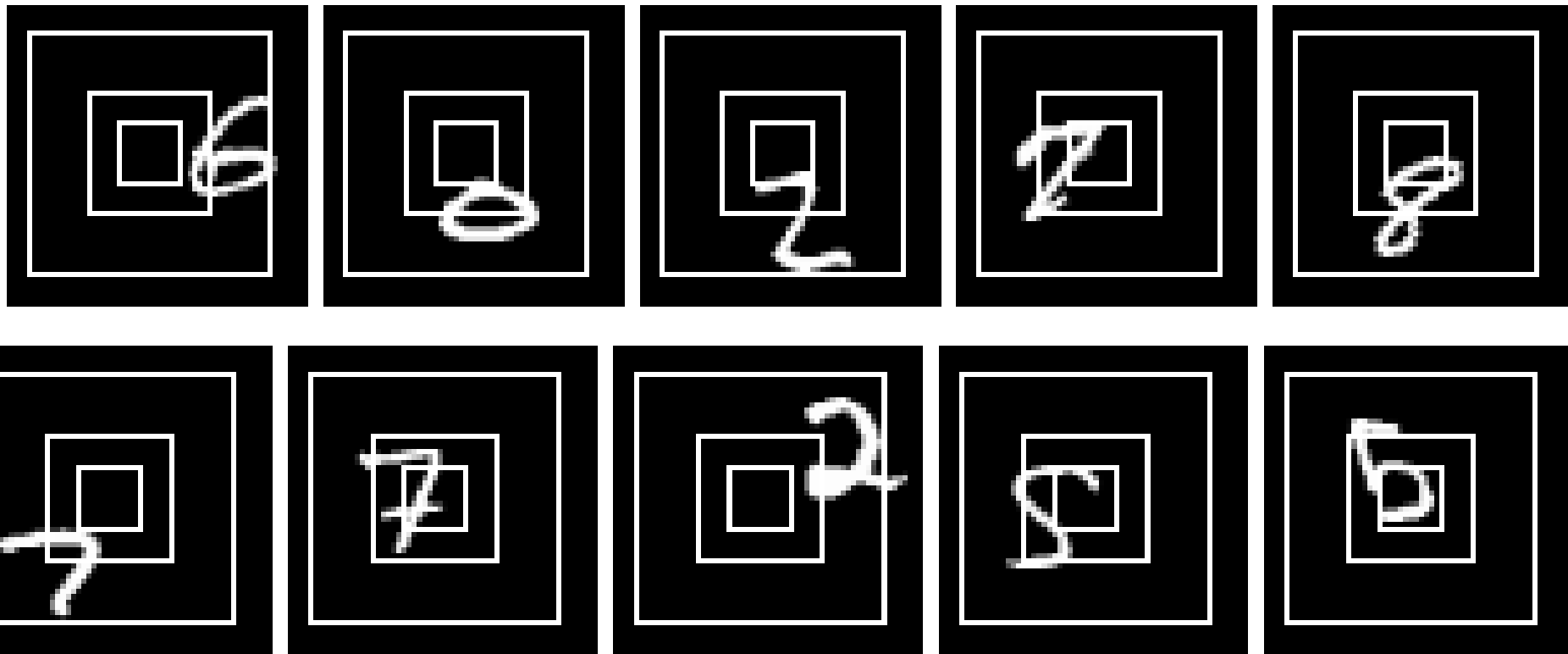


Hard attention

- Soft attention requires computing a representation for the whole image or sentence
- Hard attention on the other hand forces looking only at one part
- Main motivation was reduced computational cost rather than improved accuracy (although that happens a bit as well)
- **Saccade followed by a glimpse – how human visual system works**

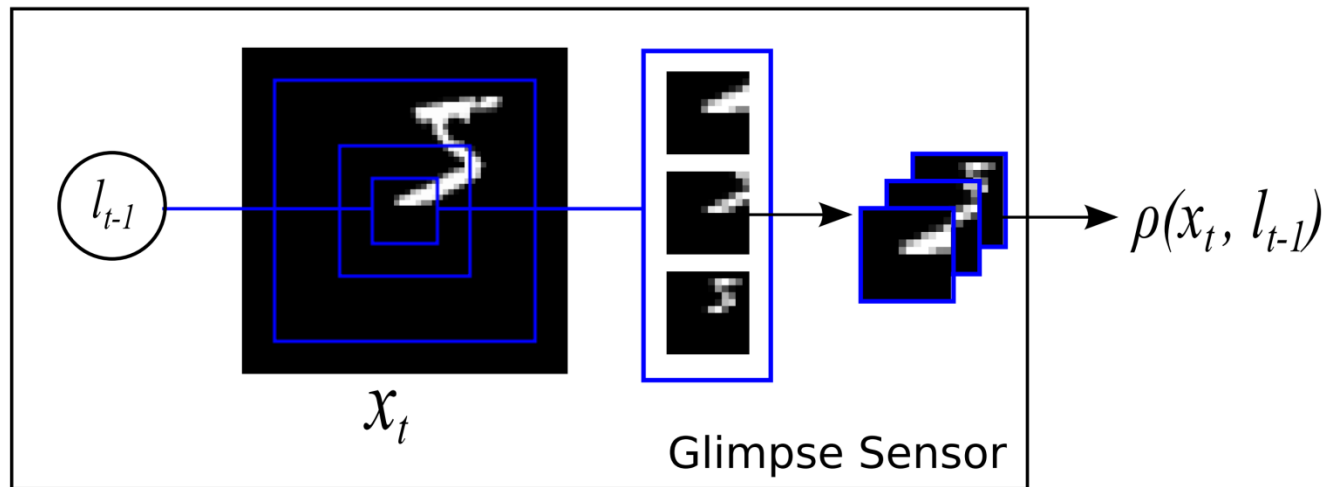
[Recurrent Models of Visual Attention, Mnih, 2014]
[Multiple Object Recognition with Visual Attention, Ba, 2015]

Hard attention examples



Glimpse Sensor

Looking at a part of an image at different scales

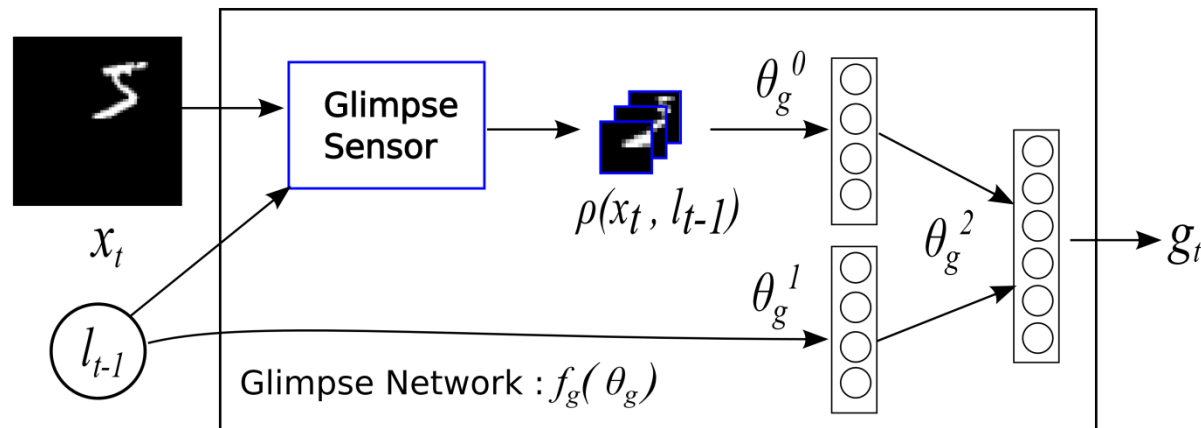


- At a number of different scales combined to a single multichannel image (human retina like representation)
- Given a location l_t output an image summary at that location

[Recurrent Models of Visual Attention, Mnih, 2014]

Glimpse network

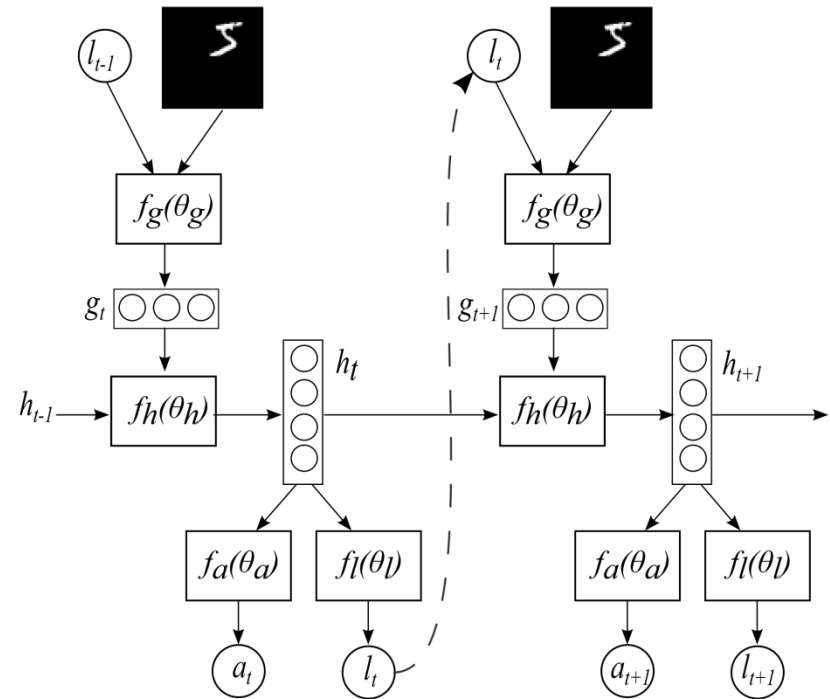
- Combining the Glimpse and the location of the glimpse into a joint network



- The glimpse is followed by a feedforward network (CNN or a DNN)
- The exact formulation of how the location and appearance are combined varies, the important thing is combining **what** and **where**
- Differentiable with respect to glimpse parameters but not the location

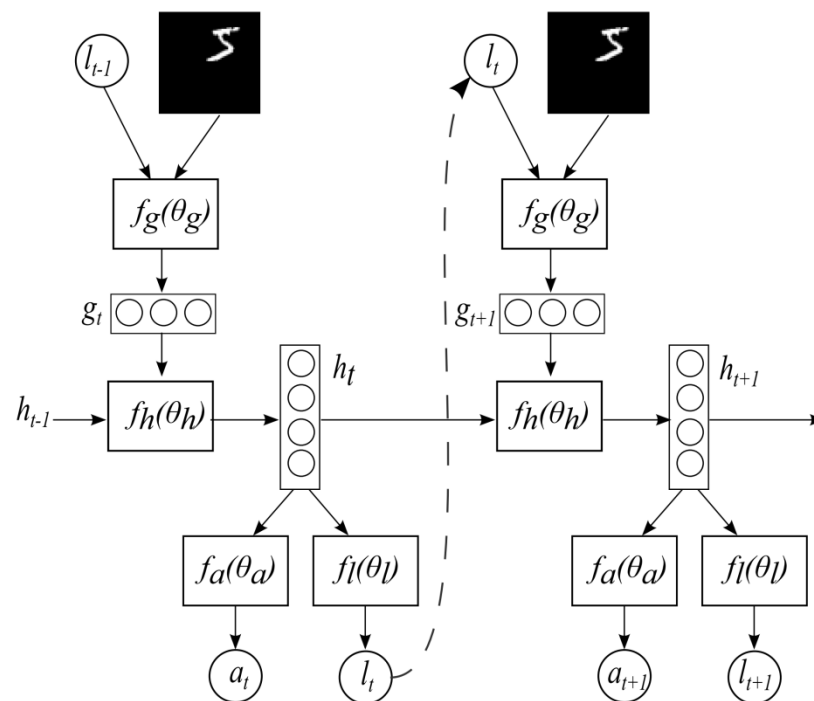
Overall Architecture - Emission network

- Given an image a glimpse location l_t , and optionally an action a_t
- Action can be:
 - Some action in a dynamic system – press a button etc.
 - Classification of an object
 - Word output
- This is an RNN with two output gates and a slightly more complex input gate!



Recurrent model of Visual Attention (RAM)

- Sample locations of glimpses leading to updates in the network
- Use gradient descent to update the weights (the glimpse network weights are differentiable)
- The emission network is an RNN
- Not as simple as backprop but doable
- Turns out this is very similar and in some cases equivalent to reinforcement learning using the REINFORCE learning rule [Williams, 1992]



Multi-modal alignment recap

Multimodal-alignment recap

- Explicit alignment - aligns two or more modalities (or views) as an actual task. The goal is to find correspondences between modalities
 - Dynamic Time Warping
 - Canonical Time Warping
 - Deep Canonical Time Warping
- Implicit alignment - uses internal latent alignment of modalities in order to better solve various problems
 - Attention models
 - Soft attention
 - Spatial transformer networks
 - Hard attention