





Advanced Multimodal Machine Learning

Lecture 9.1: Attention models and alignment

Louis-Philippe Morency Tadas Baltrušaitis

Upcoming Schedule

- First project assignment:
 - Proposal presentation (2/21 and 2/22)
 - First project report (3/5)
- Midterm project assignment
 - Midterm presentations (Tuesday 4/4 & Thursday 4/6)
 - Midterm report (Sunday 4/9)
- Final project assignment
 - Final presentation (5/2 & 5/4)
 - Final report (5/7)

Midterm Presentation Instructions

- 8-9 minute presentations (12-18 slides)
 - +2-3 minutes for written feedback and notes
- All team members should be involved during the presentation.
- The ordering of the presentations (Tuesday vs. Thursday) will be inverted based on proposal presentations.
- The presentations will be from 4:30pm 6:15 to give everyone more time
 - Let us know if you need to leave before then

Midterm Presentation Instructions

- Motivation and general definition of your research problem (2-3 slides)
- Mathematical formalization of the research problem, including definition of the main variables and overall objective function (2-4 slides)
- Explain at least one multimodal baseline model for your research problem (2-4 slides)
- Present current results of this baseline model on your dataset. You should study the failure cases of the baseline model (2-4 slides)
- Describe the research directions you are planning to explore.
 Discuss how they will address some of the shortcoming of your baseline model. (2-3 slides)

Midterm Project Report Instructions

PART 1

- Research problem: describe and motivate the research problem you are planning to work on. Explain why this problem is important for the research community and, if possible, the society in general. Define in generic terms the main computational challenges involved in this research problem.
- Related Work: Present an overview of the work happening in this research area. This section should include about 12-15 citations of prior work, grouped in similar topics. Also, you should present in more details the 3-4 research papers most related to your proposed work. The related work section should end emphasizing how your proposed approach differ from previous work.
- Dataset and Input Modalities: Describe the dataset(s) you are planning to use for this project. If many options exist, please motivate your choice of dataset for this research problem. Describe the input modalities and annotations available in this dataset. Specify which subset of these modalities and annotations you are planning to use.

Midterm Project Report Instructions

PART 2

- Problem statement: formalize mathematically your research problem. This should include the mathematical definition of the variables involved in your problem.
- Multimodal baseline models: Describe mathematically at least one multimodal baseline model for your research problem.
- Experimental methodology: Describe your experimental methodology for evaluating the multimodal baseline model(s).
- Results and Discussion: Present in tables and/or figures your experimental results. This section should include more than re-running existing baseline models.
- Proposed approach: Describe what models you are planning to test for the final report experiments. Whenever possible, you should write down the loss function of these models, following the same mathematical formulation previously used.

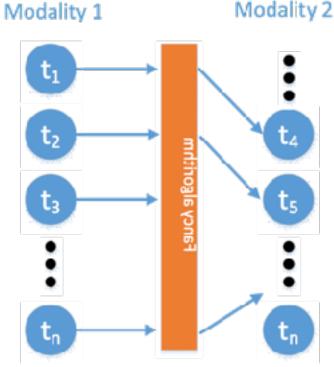
Multi-modal alignment

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

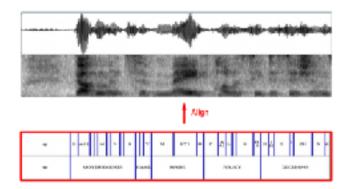
Multimodal-alignment

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



Multimodal-alignment

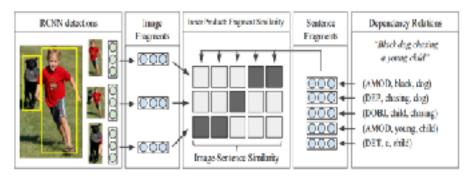
- Explicit alignment goal is to find correspondences between modalities
 - Aligning speech signal to a transcript
 - Aligning two out-of sync sequences
 - Co-referring expressions





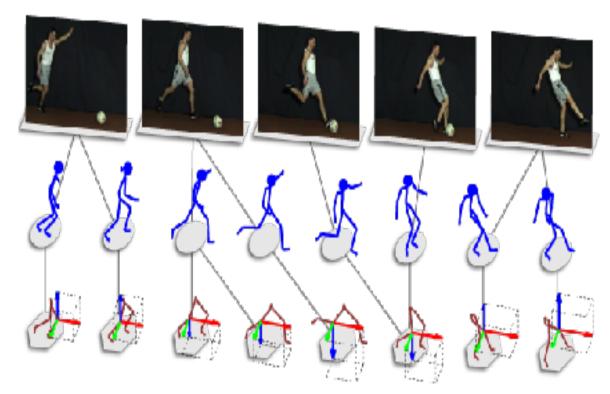
Multimodal-alignment

- 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

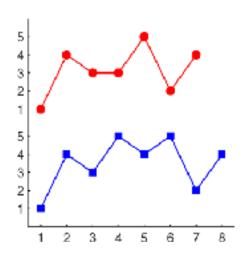
$$\mathbf{X} = [x_1, x_2, \dots, x_{n_x}] \in \mathbb{R}^{d \times n_x}$$

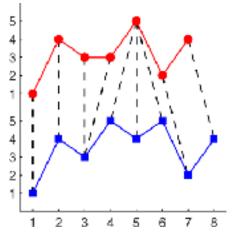
•
$$\mathbf{Y} = \begin{bmatrix} \mathbf{y}_1, \mathbf{y}_2, ..., \mathbf{y}_{n_y} \end{bmatrix} \in \mathbb{R}^{d \times n_y}$$

 Find set of indices to minimize the alignment difference:

$$L(\boldsymbol{p_t^x}, \boldsymbol{p_t^y}) = \sum_{t=1}^{l} \left\| \boldsymbol{x_{p_t^x}} - \boldsymbol{y_{p_t^y}} \right\|_{2}^{2}$$

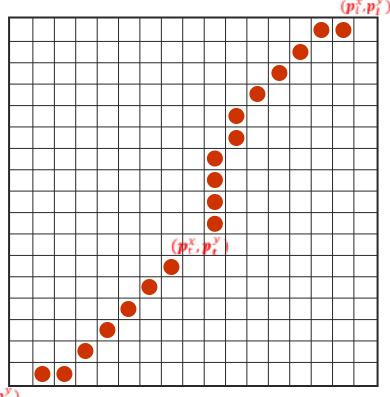
- Where p^x and p^y are index vectors of same length
- Finding these indices is called Dynamic Time Warping





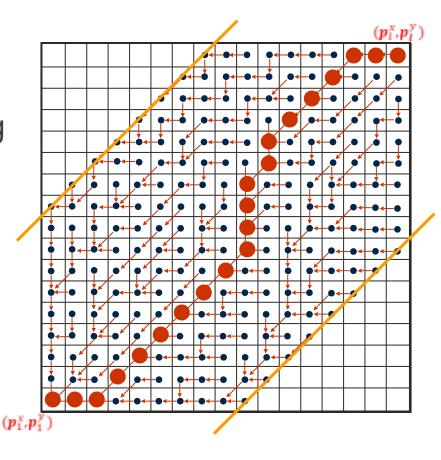
Dynamic Time Warping continued

- Lowest cost path in a cost matrix (expensive to compute)
- 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

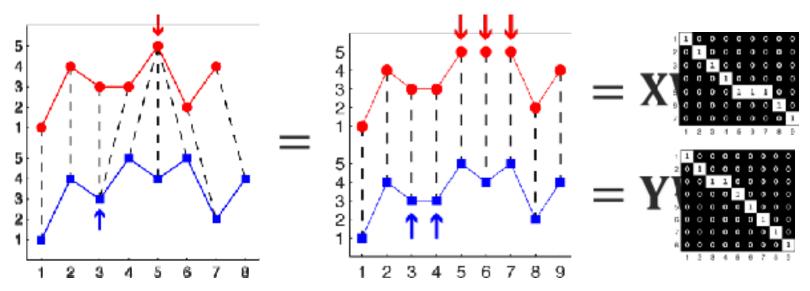
- Lowest cost path in a cost matrix (expensive to compute)
- Solved using dynamic programming whilst respecting the restrictions



DTW alternative formulation

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

Replication doesn't change the objective!



Alternative objective:

$$L(\boldsymbol{W_x}, \boldsymbol{W_y}) = \|\boldsymbol{X}\boldsymbol{W_x} - \boldsymbol{Y}\boldsymbol{W_y}\|_F^2$$

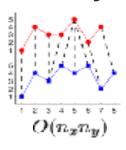
X, Y — original signals (same #rows, possibly different #columns)

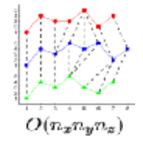
 $oldsymbol{W}_{\chi}$, $oldsymbol{W}_{y}$ - alignment matrices



DTW - limitations

Computationally complex





m sequences

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

Sensitive to outliers

Unimodal!



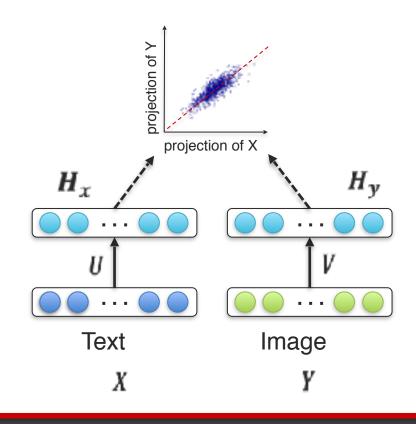


Canonical Correlation Analysis reminder

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

subject to: $U^T \Sigma_{YY} U = V^T \Sigma_{YY} V = I$

- Linear projections maximizing correlation
- 2 Orthogonal projections
- Unit variance of the projection vectors

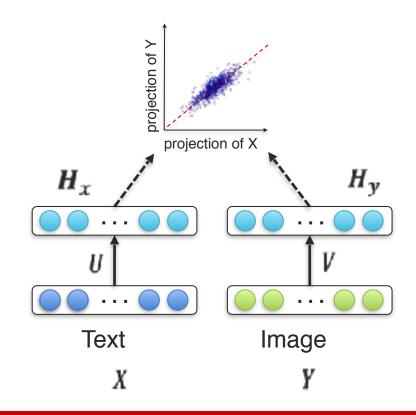


Canonical Correlation Analysis reminder

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

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

subject to: $U^T \Sigma_{yy} U = V^T \Sigma_{yy} V = I$



Canonical Time Warping

Dynamic Time Warping + Canonical Correlation Analysis = Canonical Time Warping

$$L(\mathbf{U}, \mathbf{V}, \mathbf{W}_{x}, \mathbf{W}_{y}) = \left\| \mathbf{U}^{T} \mathbf{X} \mathbf{W}_{x} - \mathbf{V}^{T} \mathbf{Y} \mathbf{W}_{y} \right\|_{F}^{2}$$

- Allows to align multi-modal or multi-view (same modality but from a different point of view)
- W_x, W_y temporal alignment
- U, V − cross-modal (spatial) alignment

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

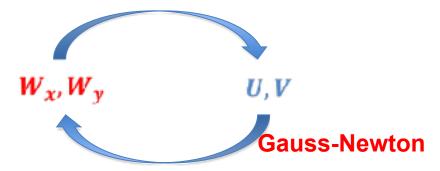


Canonical Time Warping

$$L(\mathbf{U}, \mathbf{V}, \mathbf{W}_{x}, \mathbf{W}_{y}) = \left\| \mathbf{U}^{T} \mathbf{X} \mathbf{W}_{x} - \mathbf{V}^{T} \mathbf{Y} \mathbf{W}_{y} \right\|_{F}^{2}$$

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]

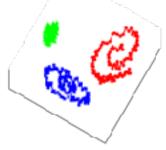
Generalized Time warping

Generalize to multiple sequences all of different modality

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

W_i – set of temporal alignments

U_i – set of cross-modal (spatial) alignments



- (1) Time warping(2) Spatial embedding

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

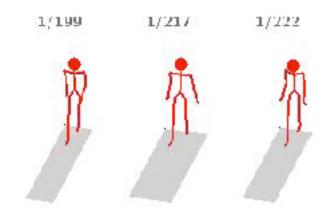
Alignment examples (unimodal)

CMU Motion Capture

Subject 1: 199 frames

Subject 2: 217 frames

Subject 3: 222 frames

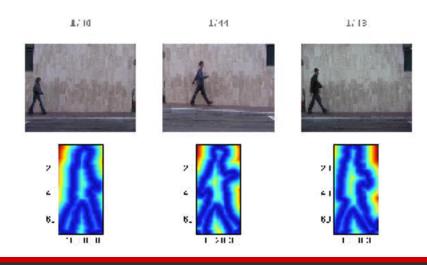


Weizmann

Subject 1: 40 frames

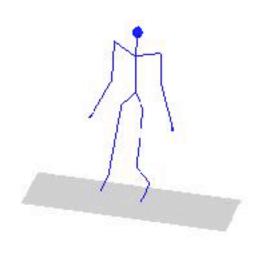
Subject 2: 44 frames

Subject 3: 43 frames



Alignment examples (multimodal)

1/273 1/51 1/127







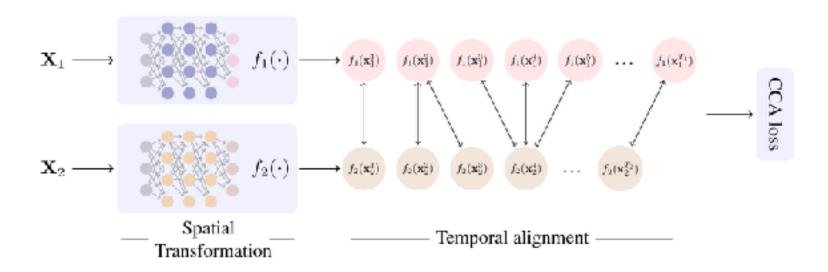
Canonical time warping - limitations

- Linear transform between modalities
- How to address this?

Deep Canonical Time Warping

$$L(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{W}_{\boldsymbol{x}}, \boldsymbol{W}_{\boldsymbol{y}}) = \|f_{\boldsymbol{\theta}_1}(\mathbf{X})\mathbf{W}_{\boldsymbol{x}} - f_{\boldsymbol{\theta}_1}(\mathbf{Y})\mathbf{W}_{\boldsymbol{y}}\|_F^2$$

Could be seen as generalization of DCCA and GTW



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

Deep Canonical Time Warping

$$L(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{W}_{\boldsymbol{x}}, \boldsymbol{W}_{\boldsymbol{y}}) = \|f_{\boldsymbol{\theta}_1}(\boldsymbol{X})\boldsymbol{W}_{\boldsymbol{x}} - f_{\boldsymbol{\theta}_1}(\boldsymbol{Y})\boldsymbol{W}_{\boldsymbol{y}}\|_F^2$$

- Optimization is again iterative:
 - Solve for alignment (W_x, W_y) with fixed projections (θ₁, θ₂)
 - Eigen decomposition
 - Solve for projections (θ_1, θ_2) with fixed alignment (W_x, W_y)
 - Gradient descent
 - Repeat till convergence

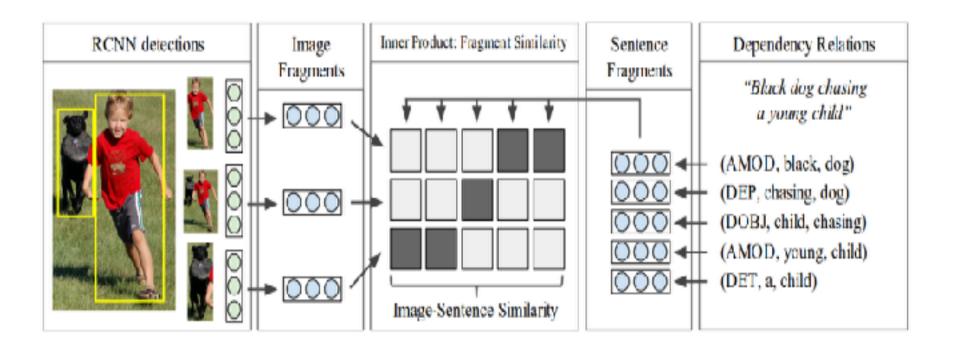
[Deep Canonical Time Warping, Trigeorgis et al., 2016, CVPR]

Implicit alignment

Implicit alignment

- We looked how to explicitly align data
- Could use that as a pre-processing step in our pipelines
- Can we instead allow/encourage the model to align data when solving a particular problem?

Deep fragment embeddings



Visual Question Answering

Xu and Saenko, 2015

What season does this appear to be? GT: fall Our Model: fall



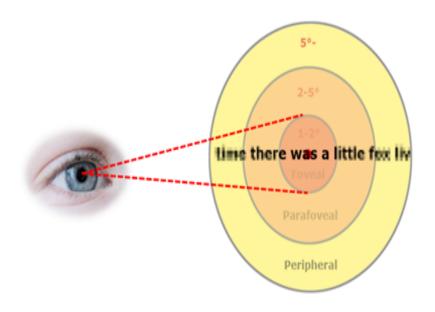
What is soaring in the sky?
GT: kite Our Mo



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

- Lots of attention
- 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 models

- Recent attention models can be roughly split into two major categories
- Soft attention
 - Acts more like a gate function
 - Deterministic
- Hard attention
 - Is more reminiscent of reinforcement learning
 - Stochastic

Soft attention

Machine Translation

Given a sentence in one language translate it to another

Dog on the beach -> le chien sur la plage

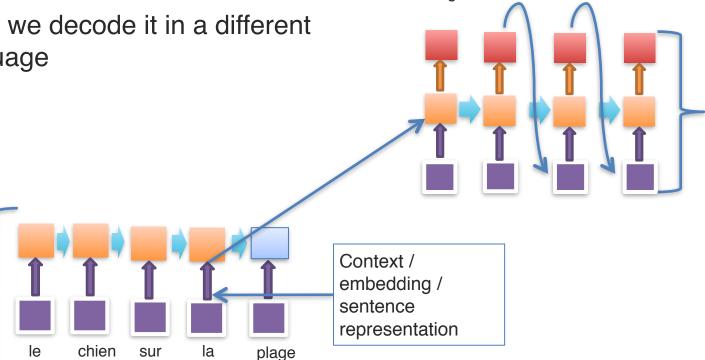
 Not exactly multimodal task – but a good start! Each language can be seen almost as a modality.

Machine Translation with RNNs

A quick reminder about encoder decoder frameworks



Then we decode it in a different language



Dog

on

the

Encoder

beach

Decode

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

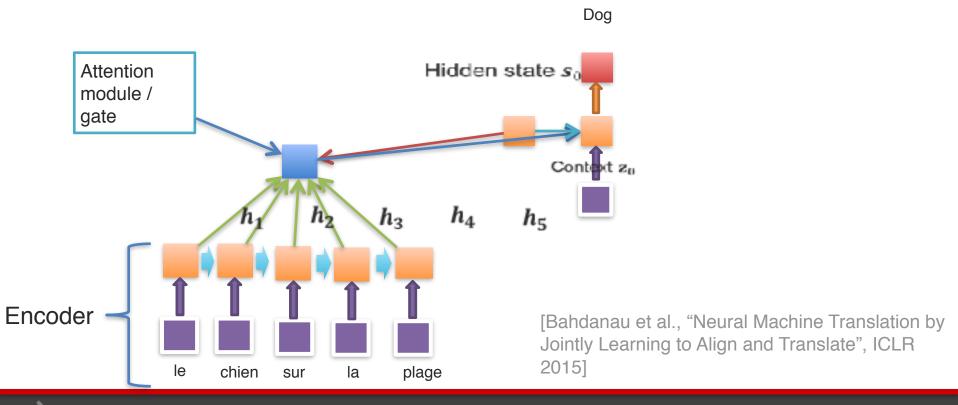
The agreement on the European Economic Area was signed in August 1992.



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

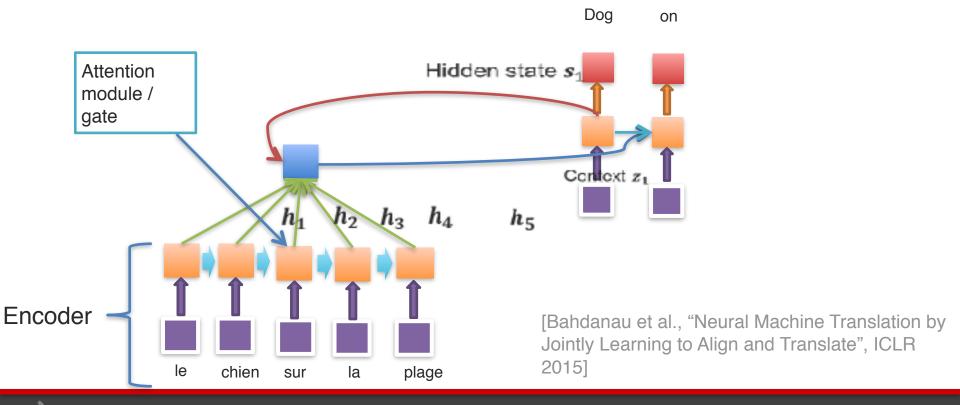
Decoder – attention model

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



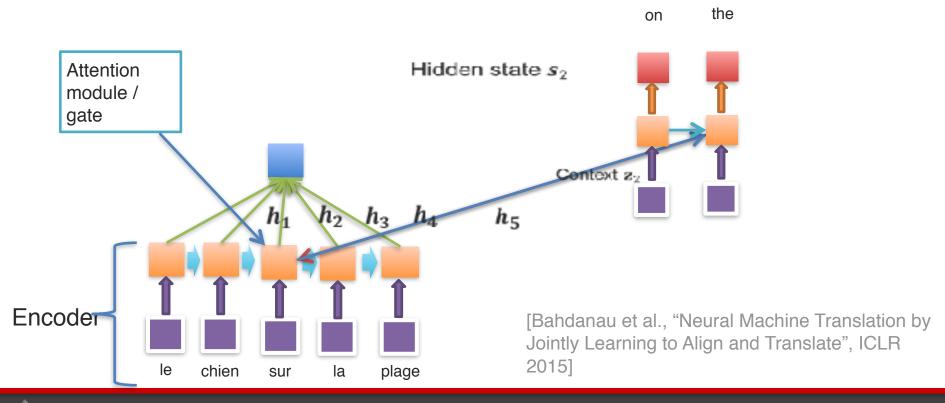
Decoder – attention model

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



Decoder – attention model

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





How do we encode attention

- Before:
 - $p(y_i|y_1,...,y_{i-1},x) = g(y_{i-1},s_i,z)$, where $z = h_T$, and s_i the current state of the decoder
 - Now:
 - $p(y_i|y_1,...,y_{i-1},x) = g(y_{i-1},s_i,z_i)$
 - Have an attention "gate"
 - A different context z_i used at each time step!
 - $\mathbf{z}_i = \sum_{j=i}^{T_{\chi}} \alpha_{ij} \mathbf{h}_j$

 $lpha_{ij}$ - 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:
 - $\bullet \quad e_{ij} = \mathbf{v}^T \sigma \big(W \mathbf{s_{i-1}} + U \mathbf{h_j} \big)$
 - a feedforward network that can tell us given the current state of decoder how important the current encoding is now
 - v, W, U− learnable weights

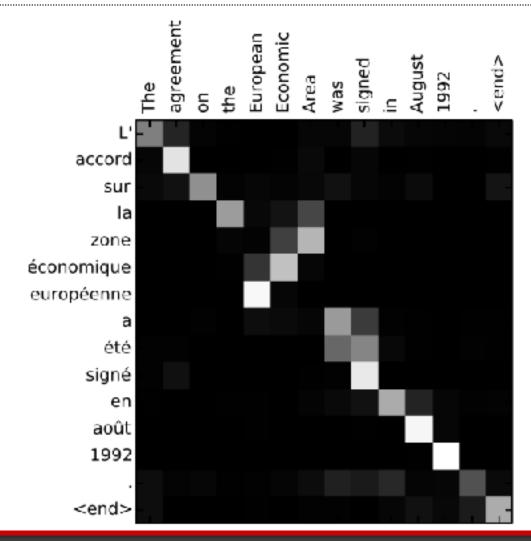
$$z_i = \sum_{j=i}^{T_x} \alpha_{ij} h_j$$

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

MT with attention

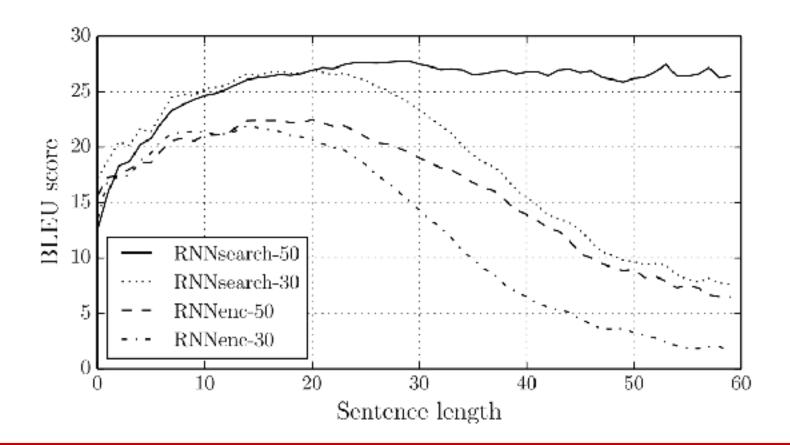
- Basically we are using a neural network to tell us where a neural network should be looking!
- Can use with RNN, LSTM or GRU
- Encoder being used is the same structure as before
 - Can use uni-directional
 - Can use bi-directional
- Model can be trained using our regular back-propagation through time, all of the modules are differentiable

Does it work?



Does it work

Especially good for long sentences



MT with attention recap

- Get 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

- A similar model but with a visual modality over which we pay attention
- This week's reading group paper















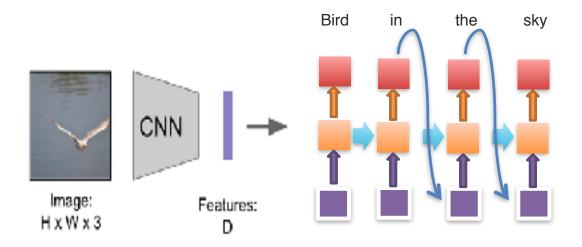






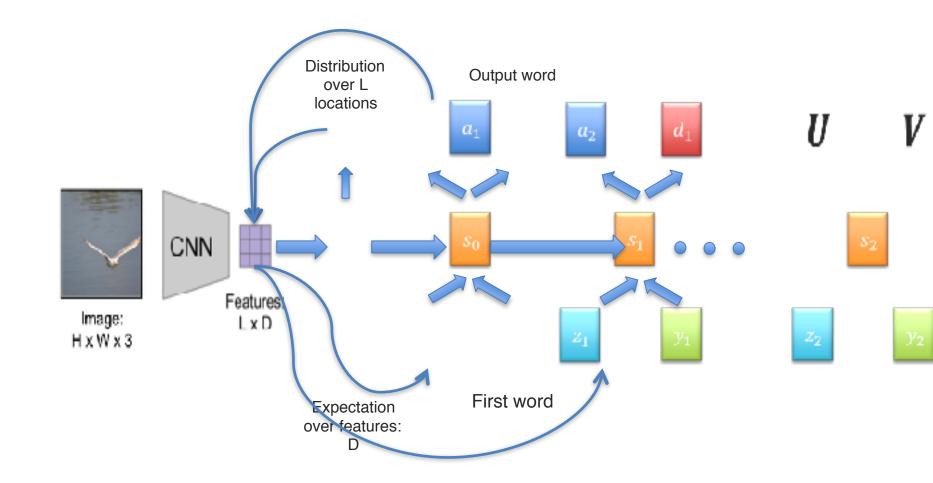
[Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, Xu et al., 2015]

Recap RNN for Captioning



Why might we not want to focus on the final layer?

Looking at more fine grained features



Visual captioning with soft attention

- Works pretty well outperforms a number of baselines
- Allows us to get an idea of what the network "sees"
- A very similar model to that used for translation
- As well can be optimized using back propagation
 - Code is available https://github.com/kelvinxu/arctic-captions
- Also explored hard attention (our next topic)

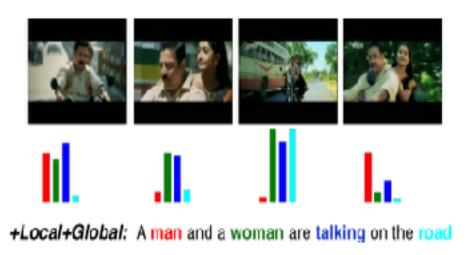
Other attention work

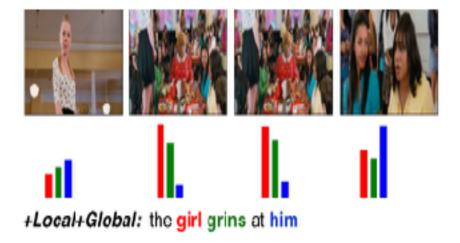
Good at paper naming!

- Show, Attend and Tell (extension of Show and Tell)
- Listen, Attend and Walk
- Listen, Attend and Spell
- Ask, Attend and Answer

Video Descriptions

Yao et al. 2015



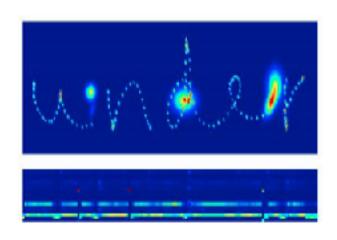


Soft attention model

Speech

- Speech transcription Listen, Attend and Spell [Chan et al.]
- Speech recognition Attention-Based Models for Speech Recognition [Chorowski et al.]

Generative models with attention



Groves, "Generating Secuences with Recurrent Neural Networks", prXiv 2013

more of national temperatures

more of national temperatures

more of material temperatures

more of material temperatures

more of material temperatures

more of material temperatures

DRAW paper from DeepMind

Hard attention

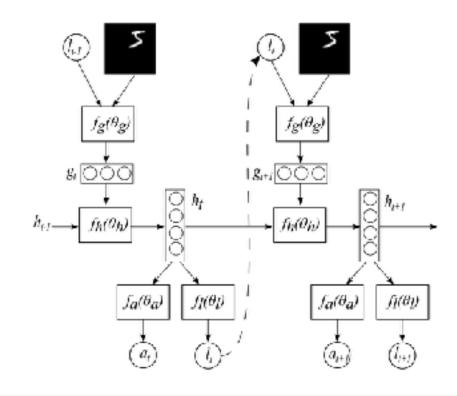
Hard attention

- 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]

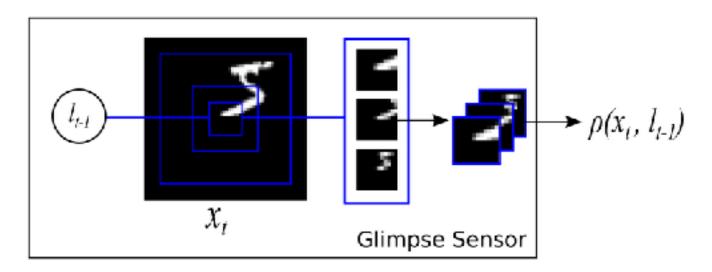
Emission network

- Given an image and a glimpse/attention location the model produces next 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!



Glimpse

Looking at a part of an image at different scales

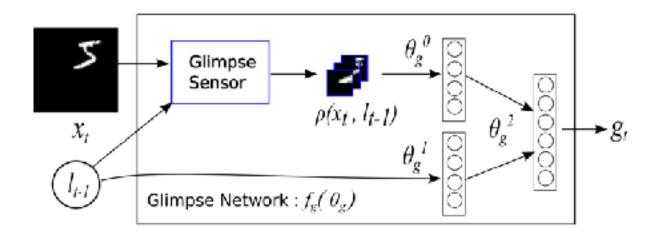


- At a number of different scales combined to a single multichannel image (human retina like representation)
- Given a location l_t output an image summary at that location

[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

Recurrent model of Visual Attention (RAM)

Model definition

$$L = \log[p(y|X,W)] = \log\sum_{l} p(l|X,W)p(y|l,X,W)$$

$$W - \text{set of parameters of the RNN } (\theta_g,\theta_l,\theta_a)$$

$$X \text{ the input (image, frame etc.)}$$

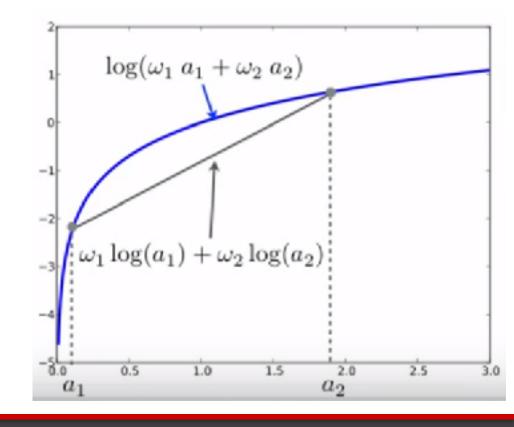
$$l - \text{the set of actions and locations}$$

$$y - \text{correct output (digit classification), correct word prediction etc.}$$

- How do we sum across all of them?
 - It seems summing is the biggest problem in ML
 - Inability to sum was also the main issue in RBM and DBN training

Jensen's Inequality

- For further steps we need and aside
- Weighted average of concave functions
- $\log(\sum_i w_i a_i) \ge \sum_i w_i \log(a_i)$
- If $\sum_i w_i = 1$ and $w_i \ge 0$



Variational bound

Going back to our formulation of our loss:

$$L = \log[p(y|X, W)] = \log \sum_{l} p(l|X, W)p(y|l, X, W)$$

- F is called variational free energy lower bound
- By maximizing it we are indirectly maximizing the true likelihood

Optimization

$$\frac{\partial F}{\partial W} = \sum_{l} p(l|X, W) \left[\frac{\partial \log p(y|l, X, W)}{\partial W} + \log p(y|l, X, W) \frac{\partial \log p(l|X, W)}{\partial W} \right]$$

Still have to sum across all possible glimpses, so stochastic version:

$$\frac{\partial F}{\partial W} = \frac{1}{M} \sum_{m=1}^{M} \left[\frac{\partial \log p(y | \tilde{l}^m, X, W)}{\partial W} + \log p(y | \tilde{l}^m, X, W) \frac{\partial \log p(\tilde{l}^m | X, W)}{\partial W} \right]$$

Drawing M samples from some prior (Gaussian, Bernoulli), for example taking a glimpse or by attending to a particular region for caption generation

Optimization

$$\frac{\partial F}{\partial W} = \sum_{l} p(l|X, W) \left[\frac{\partial \log p(y|l, X, W)}{\partial W} + \log p(y|l, X, W) \frac{\partial \log p(l|X, W)}{\partial W} \right]$$

Still have to sum across all possible glimpses, so use a stochastic version:

$$\frac{\partial F}{\partial W} = \frac{1}{M} \sum_{m=1}^{M} \left[\frac{\partial \log p(y|\tilde{l}^{m}, X, W)}{\partial W} + \log p(y|\tilde{l}^{m}, X, W) \frac{\partial \log p(\tilde{l}^{m}|X, W)}{\partial W} \right]$$

Unbounded as can have very low negative values

Often replace it with an indicator function instead (Deep RAM and RAM)

Or use a moving average (Thursday's reading)

Optimization

$$\frac{\partial F}{\partial W} = \sum_{l} p(l|X, W) \left[\frac{\partial \log p(y|l, X, W)}{\partial W} + \log p(y|l, X, W) \frac{\partial \log p(l|X, W)}{\partial W} \right]$$

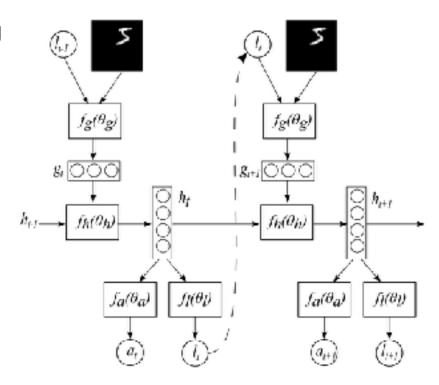
Still have to sum across all possible glimpses, so stochastic version:

$$\frac{\partial F}{\partial W} = \frac{1}{M} \sum_{m=1}^{M} \left[\frac{\partial \log p(y | \tilde{l}^m, X, W)}{\partial W} + \log p(y | \tilde{l}^m, X, W) \frac{\partial \log p(\tilde{l}^m | X, W)}{\partial W} \right]$$

To optimize our model sample stochastically for each iteration, then update the weights. Rinse and repeat.

Putting everything together

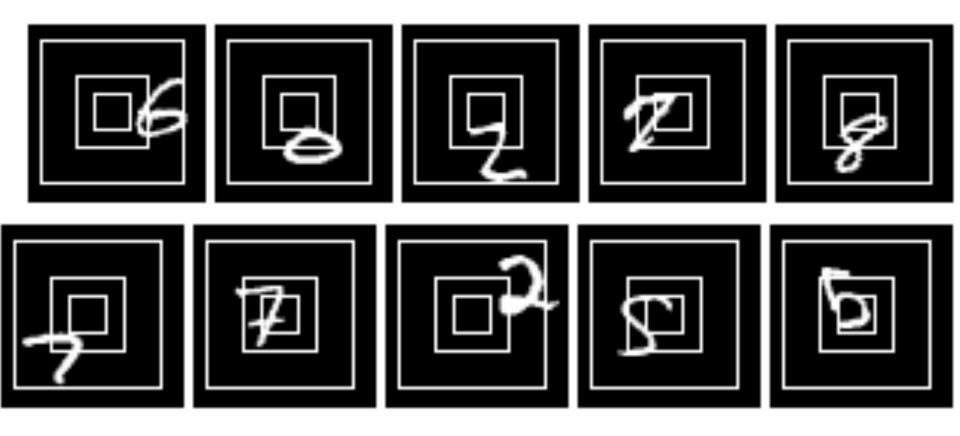
- Sample locations of glimpses leading to updates in the network
- Use gradient descent to update the weights (the glimpse network weights are differentiable)
- The emission network is an RNN
- Not as simple as backprop but doable
- Some similarities to Restricted Boltzmann Machine training



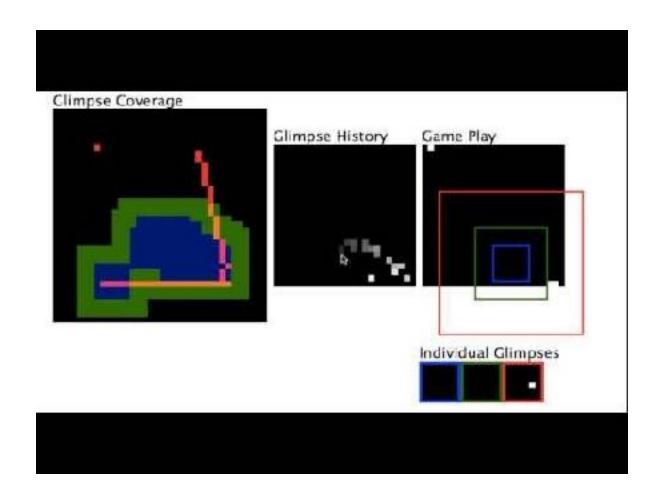
Reinforcement learning

- Turns out this is very similar and in some cases equivalent to reinforcement learning using the REINFORCE learning rule [Williams, 1992]
- Intuition
 - We want to perform a set of actions leading to a reward
 - The reward is correct object classification
 - What actions do I need to take to get there
 - How do I summarize previous actions RNN

Hard attention examples

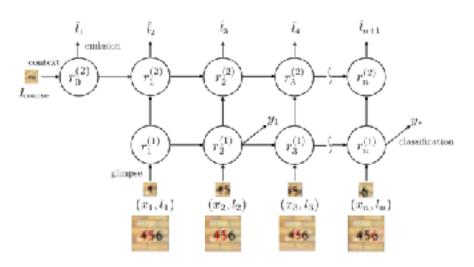


Some results



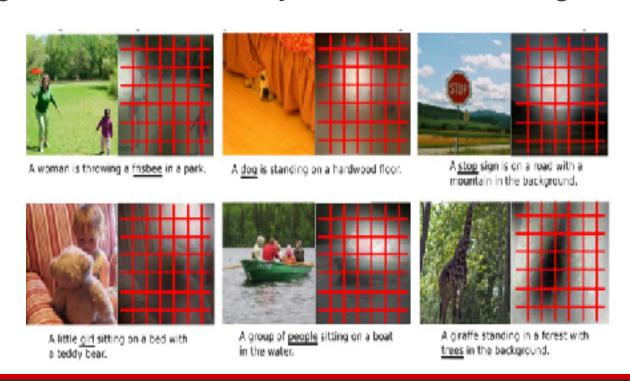
DRAM - a slight extension to RAM model

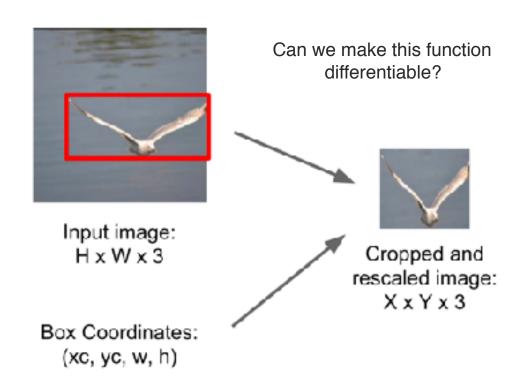
- Coarse image input initially
- Interestingly location output from top LSTM
- But the bottom one does classification
 - Why?
- This prevents the coarse image to be used for classification, thus shortcutting attention

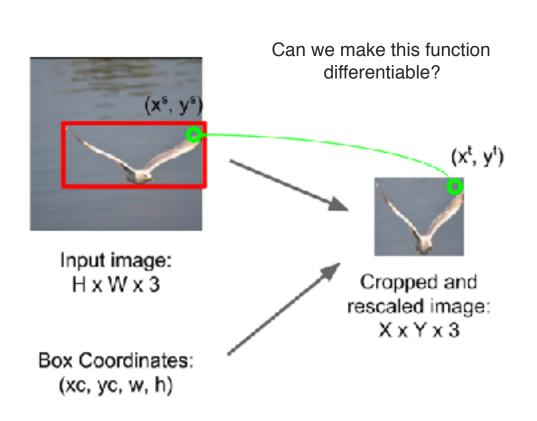


Some limitations of grid based attention

Limited to a fixed grid analysis, glimpses solve this a bit but it's difficult to train, can we fixate on small parts of image but still have easy end-to-end training?

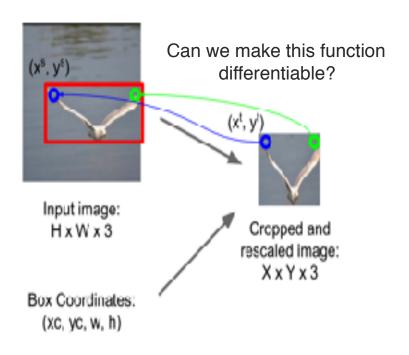






Idea: Function mapping pixel coordinates (x^t, y^t) of output to pixel coordinates (x^s, y^s) of input

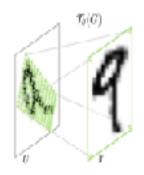
$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{1,1} & \theta_{1,2} & \theta_{1,3} \\ \theta_{2,1} & \theta_{2,2} & \theta_{2,3} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$



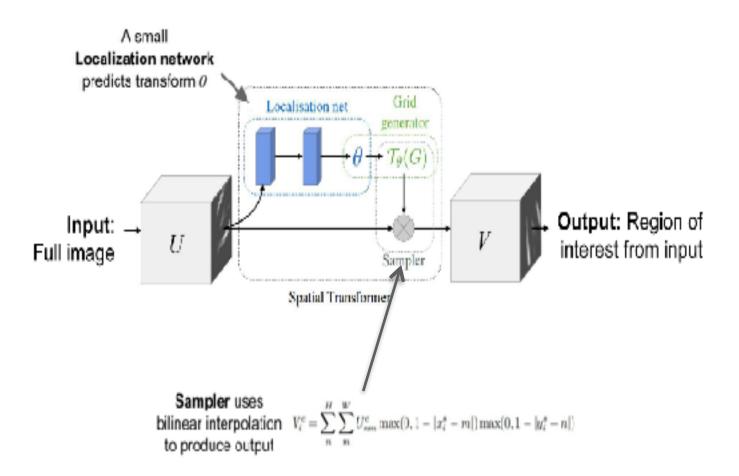
Idea: Function mapping pixel coordinates (x^t, y^t) of output to pixel coordinates (x^s, y^s) of input

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

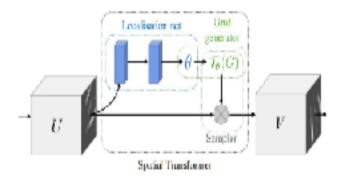
Network "attends" to input by predicting θ



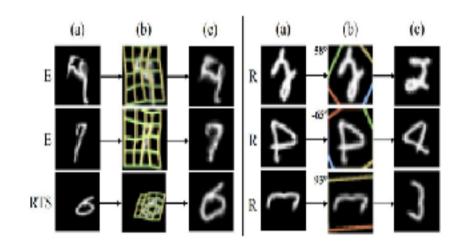
Repeat for all pixels in output to get a sampling grid

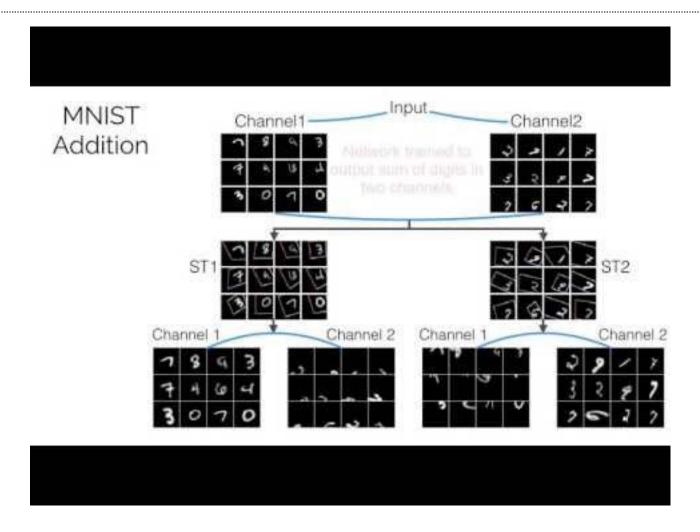


Differentiable "attention / transformation" module



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

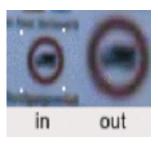




Examples on real world data

State-of-the-art 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 use complex models for focusing on an image
 - Affine and Piece-Wise Affine
 - Perspective
 - This Plate Splines
- Can use to focus on certain parts of an image
- Seems to outperform soft and hard attention models on a number of visual tasks
- We can use it instead of (or in addition to) soft and hard attention for multi-modal tasks

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
 - Hard attention
 - Spatial transformer networks