



Language
Technologies
Institute

Carnegie
Mellon
University

Tutorial on **Multimodal Machine Learning**

Louis-Philippe (LP) Morency
Tadas Baltrusaitis

**CMU Multimodal Communication and
Machine Learning Laboratory** [MultiComp Lab]

Your Instructors



Louis-Philippe Morency
morency@cs.cmu.edu



Tadas Baltrusaitis
tbaltrus@cs.cmu.edu

CMU Course 11-777: Multimodal Machine Learning

piazza 11-777 ▾ Q & A Resources Statistics Manage Class Louis-Philippe Morency

Carnegie Mellon University - Spring 2016

11-777: Advanced Multimodal Machine Learning

Syllabus

Course Information Staff Resources Groups

Description

Multimodal machine learning (MMML) is a vibrant multi-disciplinary research field which addresses some of the original goals of artificial intelligence by integrating and modeling multiple communicative modalities, including linguistic, acoustic and visual messages. With the initial research on audio-visual speech recognition and more recently with language & vision projects such as image and video captioning, this research field brings some unique challenges for multimodal researchers given the heterogeneity of the data and the contingency often found between modalities. This course will teach fundamental mathematical concepts related to MMML including multimodal alignment and fusion, heterogeneous representation learning and multi-stream temporal modeling. We will also review recent papers describing state-of-the-art probabilistic models and computational algorithms for MMML and discuss the current and upcoming challenges.

The main technical topics are: (1) multimodal representation learning, including multimodal auto-encoder and deep learning, (2) multimodal component analysis and fusion, including deep canonical correlation analysis and multi-kernel learning, (3) multimodal alignment and multi-stream modeling, including multi-instance learning and multimodal recurrent neural networks, and (4) multi-sensor computational modeling, including nonparametric Bayesian networks

Announcements [show all](#)

Room assignments for paper discussion

(4/21/2016)
4/21/16 3:41 PM

The randomized room assignment for the discussion tomorrow Thursday 4/21 at 4:30pm is shown below. Be sure to be there on time as the discussion will be shorter due to 6 presentations at the end of it.

Room WEH 4220	
Bagher Zadeh	Amirali
Bharadwaj	Akash
Correia	Joana
Jang	Hyeju
Jo	Yohan



Tutorial Schedule

■ Introduction

- What is Multimodal?
 - Historical view, multimodal vs multimedia
- Why multimodal
 - Multimodal applications: image captioning, video description, AVSR,...
- Core technical challenges
 - Representation learning, translation, alignment, fusion and co-learning



Tutorial Schedule

- **Basic concepts – Part 1**
 - Linear models
 - Score and loss functions, regularization
 - Neural networks
 - Activation functions, multi-layer perceptron
 - Optimization
 - Stochastic gradient descent, backpropagation



Tutorial Schedule

- **Unimodal representations**
 - Visual representations
 - Convolutional neural networks
 - Acoustic representations
 - Spectrograms, autoencoders



Tutorial Schedule

- **Multimodal representations**
 - Joint representations
 - Visual semantic spaces, multimodal autoencoder
 - Tensor fusion representation
 - Coordinated representations
 - Similarity metrics, canonical correlation analysis
- Coffee break [20 mins]



Tutorial Schedule

- **Basic concepts – Part 2**
 - Recurrent neural networks
 - Long Short-Term Memory models
 - Optimization
 - Backpropagation through time



Tutorial Schedule

- **Translation and alignment**
 - Translation applications
 - Machine translation, image captioning
 - Explicit alignment
 - Dynamic time warping, deep canonical time warping
 - Implicit alignment
 - Attention models, multi instance learning
 - Temporal attention-gated model



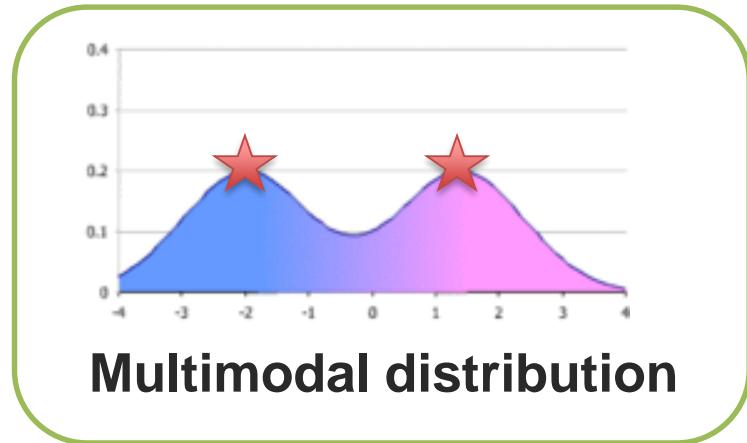
Tutorial Schedule

- **Multimodal fusion**
 - Model free approaches
 - Early and late fusion, hybrid models
 - Kernel-based fusion
 - Multiple kernel learning
 - Multimodal graphical models
 - Factorial HMM, Multi-view Hidden CRF
 - Multi-view LSTM model



What is Multimodal?

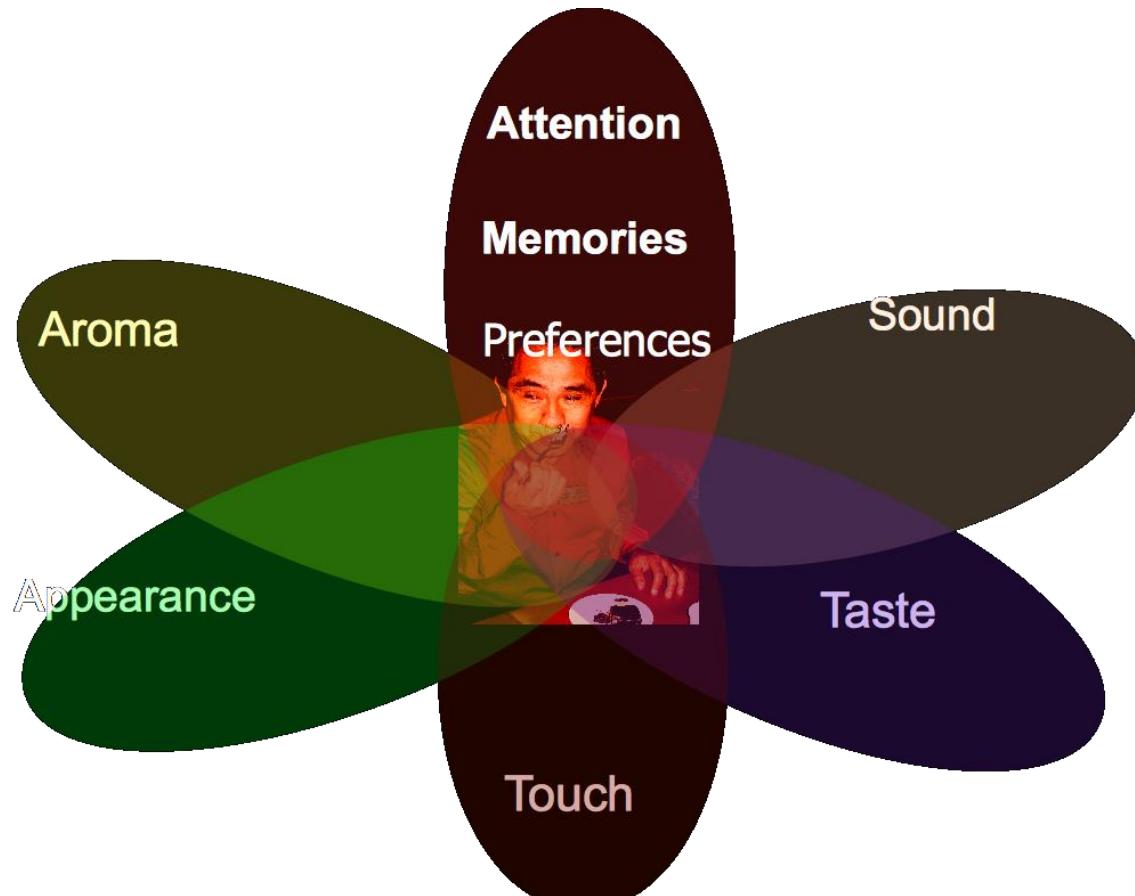
What is Multimodal?



- Multiple modes, i.e., distinct “peaks” (local maxima) in the probability density function



What is Multimodal?



Sensory Modalities



Language Technologies Institute

Carnegie Mellon University

What is Multimodal?

Modality

The way in which something happens or is experienced.

- Modality refers to a certain type of information and/or the representation format in which information is stored.
- Sensory modality: one of the primary forms of sensation, as vision or touch; channel of communication.

Medium (“middle”)

A means or instrumentality for storing or communicating information; system of communication/transmission.

- Medium is the means whereby this information is delivered to the senses of the interpreter.



Examples of Modalities

- Natural language (both spoken or written)
- Visual (from images or videos)
- Auditory (including voice, sounds and music)
- Haptics / touch
- Smell, taste and self-motion
- Physiological signals
 - Electrocardiogram (ECG), skin conductance
- Other modalities
 - Infrared images, depth images, fMRI



Multimodal Communicative Behaviors

Verbal

Lexicon

Words

Syntax

Part-of-speech
Dependencies

Pragmatics

Discourse acts

Vocal

Prosody

Intonation
Voice quality

Vocal expressions

Laughter, moans

Visual

Gestures

Head gestures
Eye gestures
Arm gestures

Body language

Body posture
Proxemics

Eye contact

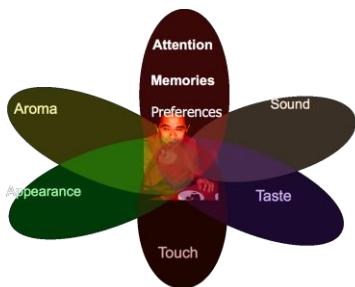
Head gaze
Eye gaze

Facial expressions

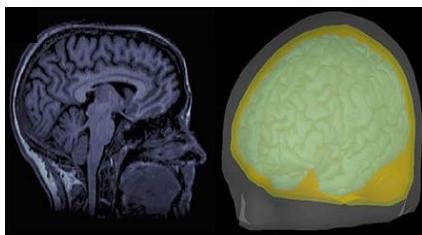
FACS action units
Smile, frowning



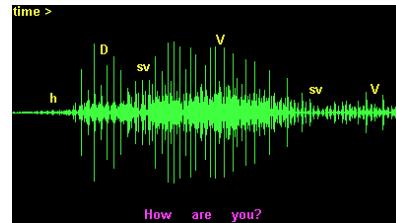
Multiple Communities and Modalities



Psychology



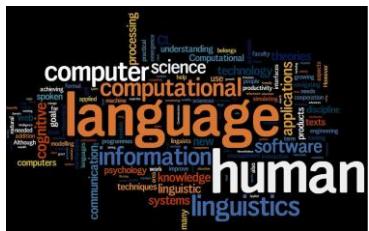
Medical



Speech



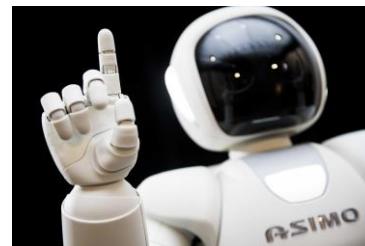
Vision



Language



Multimedia



Robotics

$$\text{ca} \quad a, \sigma^2(S_1) = \frac{\lambda - a}{\sigma^2} \int f_{a,\sigma}(\xi_1) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{|\xi_1 - a|}{\sigma^2}\right)$$
$$\int T(x) \cdot \frac{\partial}{\partial \theta} f(x, \theta) dx = M \left[T(\xi) \cdot \frac{\partial}{\partial \theta} \ln L(\xi, \theta) \right] \int_{R_n}$$
$$\int T(x) \cdot \left(\frac{\partial}{\partial \theta} \ln L(x, \theta) \right) \cdot f(x, \theta) dx = \int T(\xi) \cdot \left(\frac{\partial}{\partial \theta} \ln L(\xi, \theta) \right) \cdot f_{a,\sigma}(\xi) d\xi$$
$$\frac{\partial}{\partial \theta} M T(\xi) = \frac{\partial}{\partial \theta} \int_{R_n} T(x) f(x, \theta) dx = \int_{R_n} \frac{\partial}{\partial \theta} T(x) f(x, \theta) dx$$
$$= \int_{R_n} (\xi - a)^2 \frac{\partial}{\partial \theta} f_{a,\sigma}(\xi) d\xi$$

Learning



A Historical View

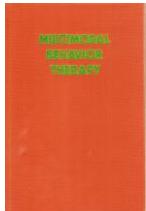
Prior Research on “Multimodal”

Four eras of multimodal research

- The “behavioral” era (1970s until late 1980s)
- The “computational” era (late 1980s until 2000)
- The “interaction” era (2000 - 2010)
- The “deep learning” era (2010s until ...)
 - ❖ Main focus of this tutorial



The “Behavioral” Era (1970s until late 1980s)



Multimodal Behavior Therapy by Arnold Lazarus [1973]

- 7 dimensions of personality (or modalities)

Multi-sensory integration (in psychology):

- Multimodal signal detection: Independent decisions vs. integration [1980]
- Infants' perception of substance and temporal synchrony in multimodal events [1983]
- A multimodal assessment of behavioral and cognitive deficits in abused and neglected preschoolers [1984]

□ TRIVIA: Geoffrey Hinton received his B.A. in Psychology ☺



Language and Gestures



David McNeill

University of Chicago

Center for Gesture and Speech Research

*“For McNeill, gestures are *in effect* the speaker’s thought in action, and integral components of speech, not merely accompaniments or additions.”*



1970

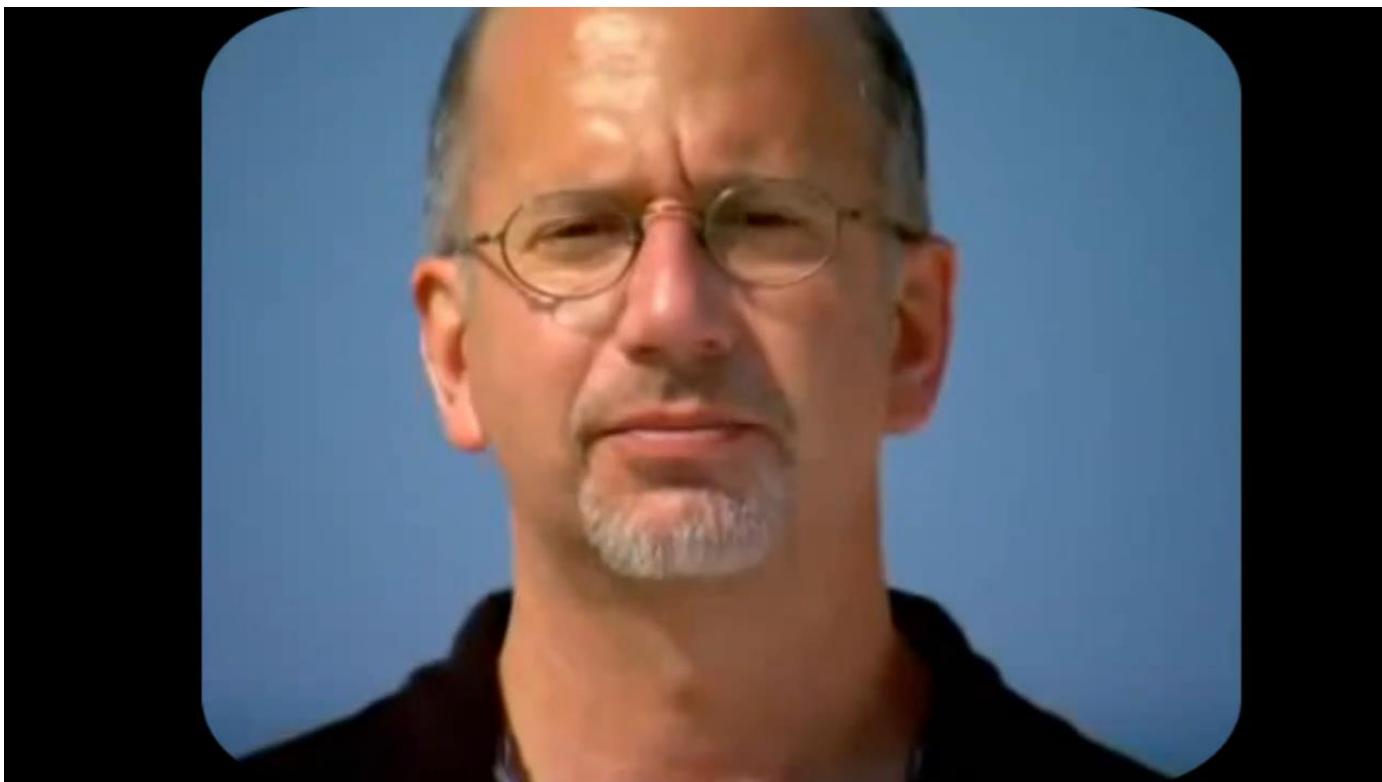
1980

1990

2000

2010

The McGurk Effect (1976)



Hearing lips and seeing voices – Nature



1970

1980

1990

2000

2010

The McGurk Effect (1976)



Hearing lips and seeing voices – Nature



1970

1980

1990

2000

2010

➤ The “Computational” Era(Late 1980s until 2000)

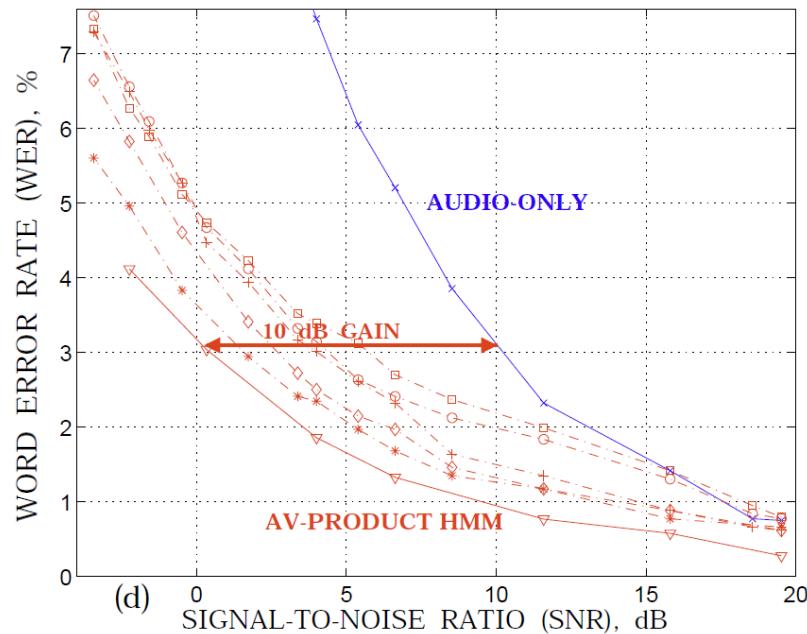
1) Audio-Visual Speech Recognition (AVSR)

- Motivated by the McGurk effect
 - First AVSR System in 1986
“Automatic lipreading to enhance speech recognition”
 - Good survey paper [2002]
“Recent Advances in the Automatic Recognition of Audio-Visual Speech”
- TRIVIA: The first multimodal deep learning paper was about audio-visual speech recognition [ICML 2011]



➤ The “Computational” Era(Late 1980s until 2000)

1) Audio-Visual Speech Recognition (AVSR)



1970

1980

1990

2000

2010

➤ The “Computational” Era (Late 1980s until 2000)

2) Multimodal/multisensory interfaces

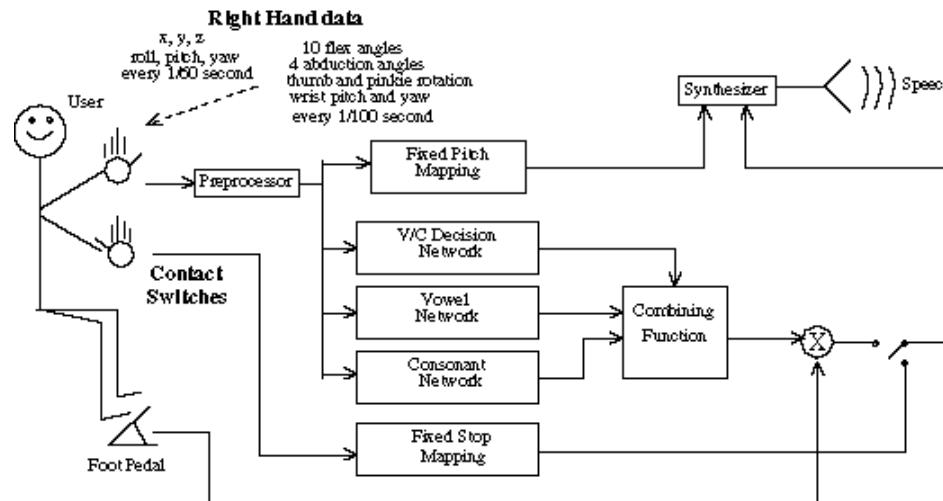
- Multimodal Human-Computer Interaction (HCI)

“Study of how to design and evaluate new computer systems where human interact through multiple modalities, including both input and output modalities.”



➤ The “Computational” Era (Late 1980s until 2000)

2) Multimodal/multisensory interfaces



Glove-talk: A neural network interface between a data-glove and a speech synthesizer By Sidney Fels & Geoffrey Hinton [CHI'95]



➤ The “Computational” Era (Late 1980s until 2000)

2) Multimodal/multisensory interfaces



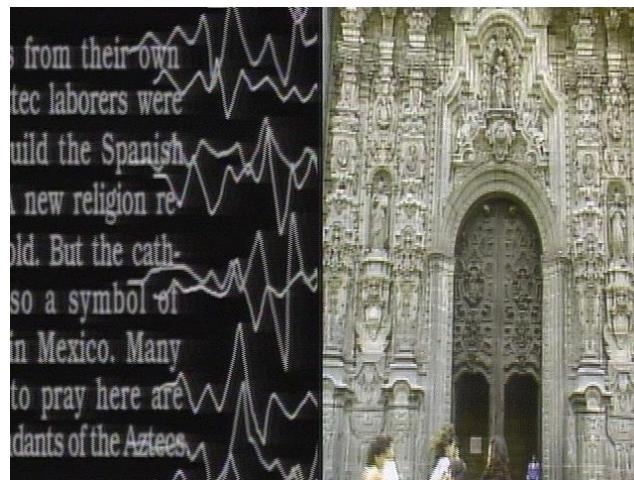
Rosalind Picard

Affective Computing is computing that relates to, arises from, or deliberately influences emotion or other affective phenomena.



➤ The “Computational” Era (Late 1980s until 2000)

3) Multimedia Computing



Carnegie
Mellon
University
informed media
arts of the
digital video understanding

[1994-2010]

“The Informed Media Digital Video Library Project automatically combines speech, image and natural language understanding to create a full-content searchable digital video library.”



Language Technologies Institute

➤ The “Computational” Era (Late 1980s until 2000)

3) Multimedia Computing

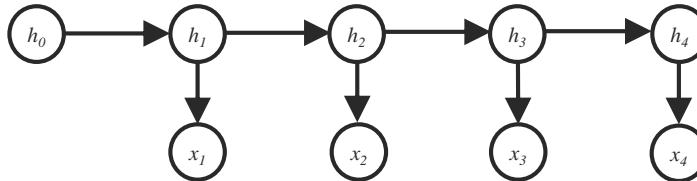
Multimedia content analysis

- **Shot-boundary detection (1991 -)**
 - Parsing a video into continuous camera shots
- **Still and dynamic video abstracts (1992 -)**
 - Making video browsable via representative frames (keyframes)
 - Generating short clips carrying the essence of the video content
- **High-level parsing (1997 -)**
 - Parsing a video into semantically meaningful segments
- **Automatic annotation (indexing) (1999 -)**
 - Detecting prespecified events/scenes/objects in video

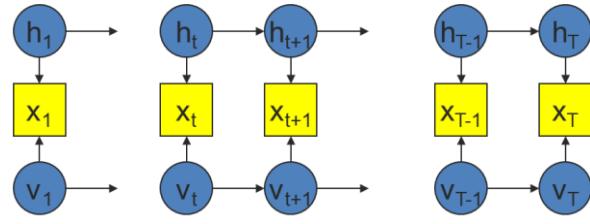


Multimodal Computation Models

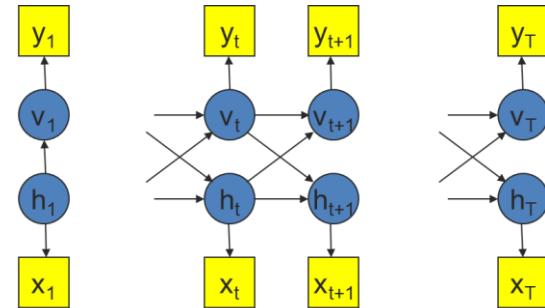
- Hidden Markov Models [1960s]



- Factorial Hidden Markov Models [1996]



- Coupled Hidden Markov Models [1997]



1970

1980

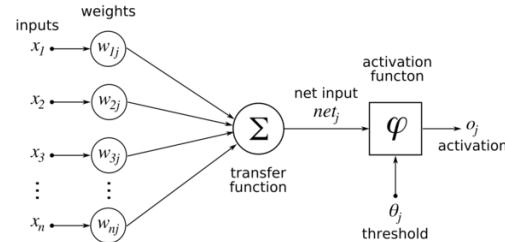
1990

2000

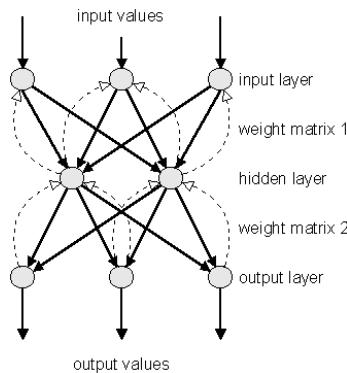
2010

Multimodal Computation Models

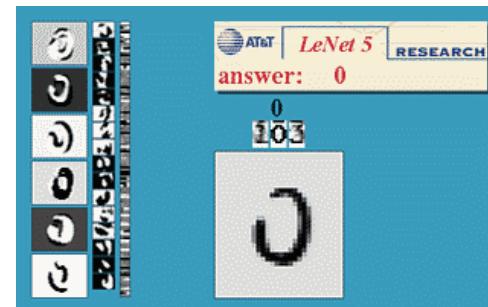
- Artificial Neural Networks [1940s]



- Backpropagation [1975]



- Convolutional neural networks [1980s]



1970

1980

1990

2000

2010

➤ The “Interaction” Era (2000s)

1) Modeling Human Multimodal Interaction



AMI Project [2001-2006, IDIAP]

- 100+ hours of meeting recordings
- Fully synchronized audio-video
- Transcribed and annotated



CHIL Project [Alex Waibel]

- Computers in the Human Interaction Loop
- Multi-sensor multimodal processing
- Face-to-face interactions

□ **TRIVIA:** Samy Bengio started at IDIAP working on AMI project



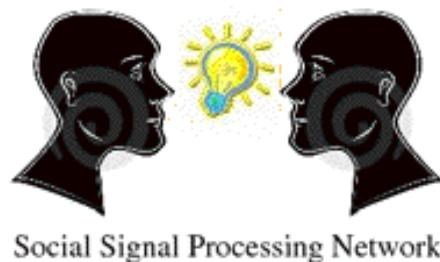
➤ The “Interaction” Era (2000s)

1) Modeling Human Multimodal Interaction



CALO Project [2003-2008, SRI]

- Cognitive Assistant that Learns and Organizes
- Personalized Assistant that Learns (PAL)
- Siri was a spinoff from this project



SSP Project [2008-2011, IDIAP]

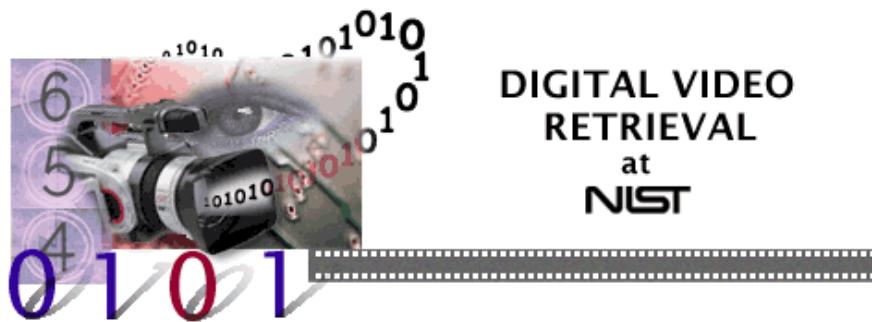
- Social Signal Processing
- First coined by Sandy Pentland in 2007
- Great dataset repository: <http://sspnet.eu/>

□ TRIVIA: LP's PhD research was partially funded by CALO ☺



➤ The “Interaction” Era (2000s)

2) Multimedia Information Retrieval

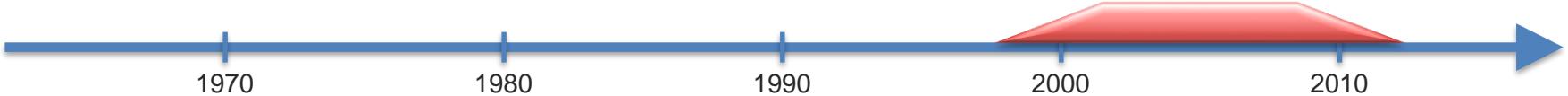


“Yearly competition to promote progress in content-based retrieval from digital video via open, metrics-based evaluation”

[Hosted by NIST, 2001-2016]

Research tasks and challenges:

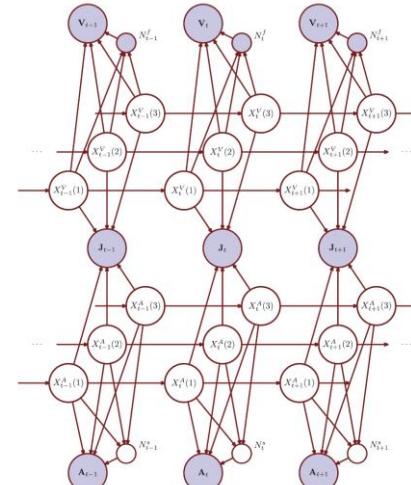
- Shot boundary, story segmentation, search
- “High-level feature extraction”: semantic event detection
- Introduced in 2008: copy detection and surveillance events
- Introduced in 2010: Multimedia event detection (MED)



Multimodal Computational Models

- Dynamic Bayesian Networks
 - Kevin Murphy's PhD thesis and Matlab toolbox
 - Asynchronous HMM for multimodal [Samy Bengio, 2007]

Audio-visual
speech
segmentation



1970

1980

1990

2000

2010

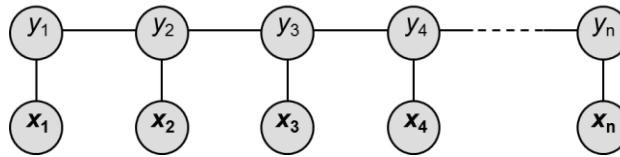


Language Technologies Institute

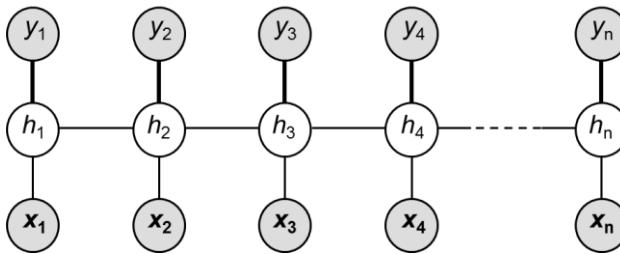
Carnegie Mellon University

Multimodal Computational Models

- Discriminative sequential models
 - Conditional random fields [Lafferty et al., 2001]



- Latent-dynamic CRF [Morency et al., 2007]



➤ The “deep learning” era (2010s until ...)

Representation learning (a.k.a. deep learning)

- Multimodal deep learning [ICML 2011]
- Multimodal Learning with Deep Boltzmann Machines [NIPS 2012]
- Visual attention: Show, Attend and Tell: Neural Image Caption Generation with Visual Attention [ICML 2015]

Key enablers for multimodal research:

- New large-scale multimodal datasets
- Faster computer and GPUS
- High-level visual features
- “Dimensional” linguistic features

Our tutorial focuses on this era!



➤ The “deep learning” era (2010s until ...)

Many new challenges and multimodal corpora !!

Audio-Visual Emotion Challenge (AVEC, 2011-)



- Emotional dimension estimation
- Standardized training and test sets
- Based on the SEMAINE dataset

Emotion Recognition in the Wild Challenge (EmotiW 2013-)



- Emotional dimension estimation
- Standardized training and test sets
- Based on the SEMAINE dataset



➤ The “deep learning” era (2010s until ...)

Renew of multimedia content analysis !

- Image captioning

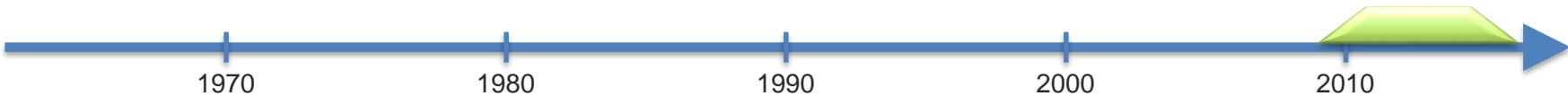


The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.

- Video description
- Visual Question-Answer



Real-World Tasks Tackled by Multimodal Research

- Affect recognition
 - Emotion
 - Persuasion
 - Personality traits
- Media description
 - Image captioning
 - Video captioning
 - Visual Question Answering
- Event recognition
 - Action recognition
 - Segmentation
- Multimedia information retrieval
 - Content based/Cross-media



Core Technical Challenges

Core Challenges in “Deep” Multimodal ML

Representation

Alignment

Fusion

Translation

Co-Learning

Multimodal Machine Learning: A Survey and Taxonomy

By Tadas Baltrusaitis, Chaitanya Ahuja,
and Louis-Philippe Morency

<https://arxiv.org/abs/1705.09406>

- 5 core challenges
- 37 taxonomic classes
- 253 referenced citations

These challenges are non-exclusive.



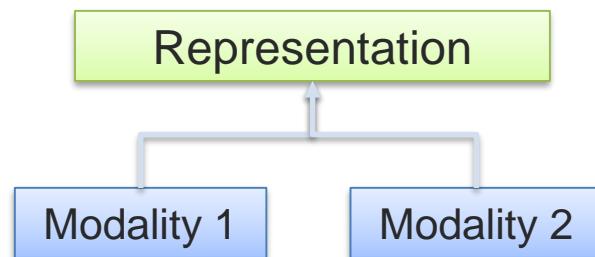
Language Technologies Institute

Carnegie Mellon University

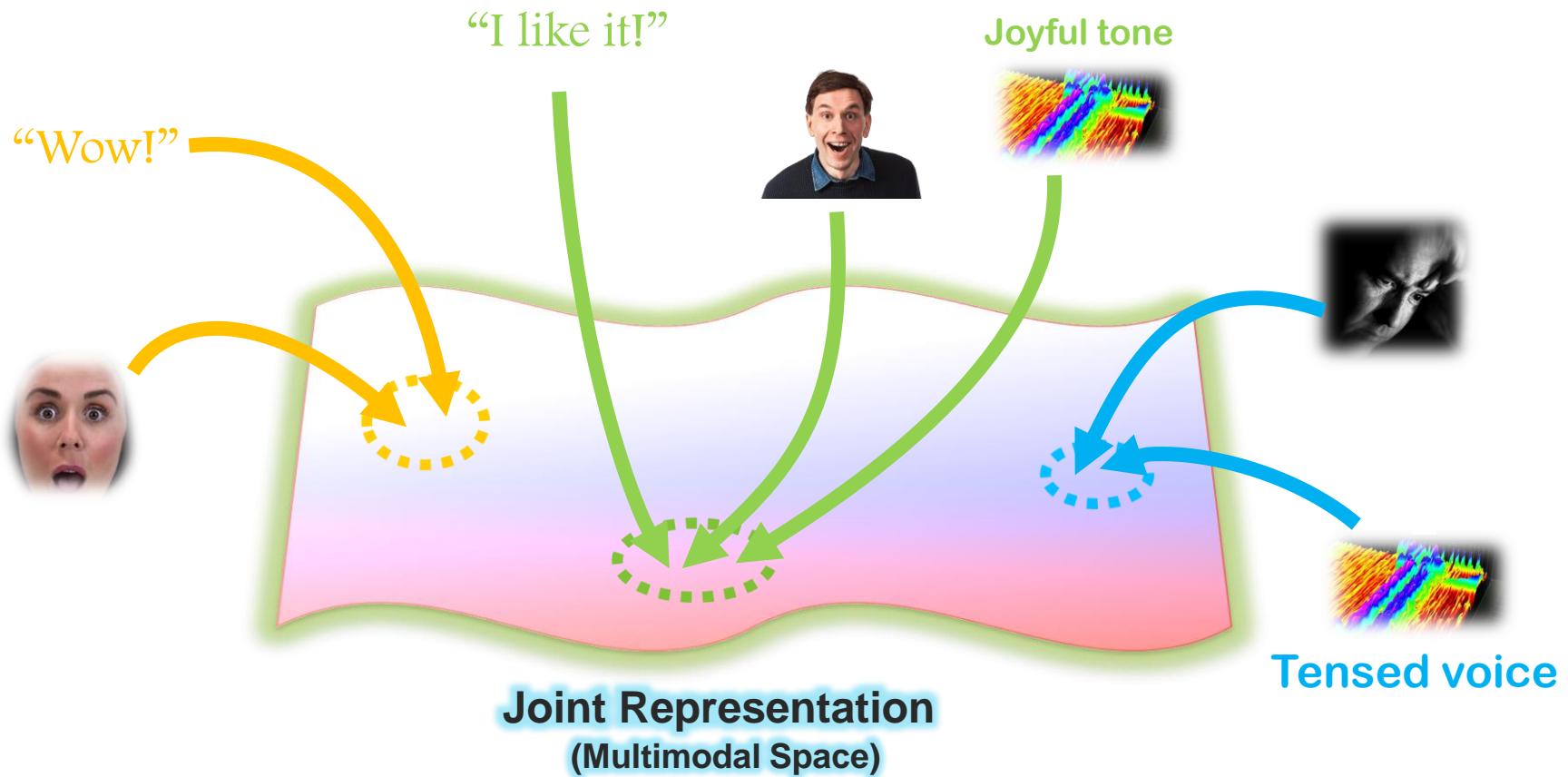
Core Challenge 1: Representation

Definition: Learning how to represent and summarize multimodal data in a way that exploits the complementarity and redundancy.

A Joint representations:



Joint Multimodal Representation



Joint Multimodal Representations

Audio-visual speech recognition

[Ngiam et al., ICML 2011]

- Bimodal Deep Belief Network

Image captioning

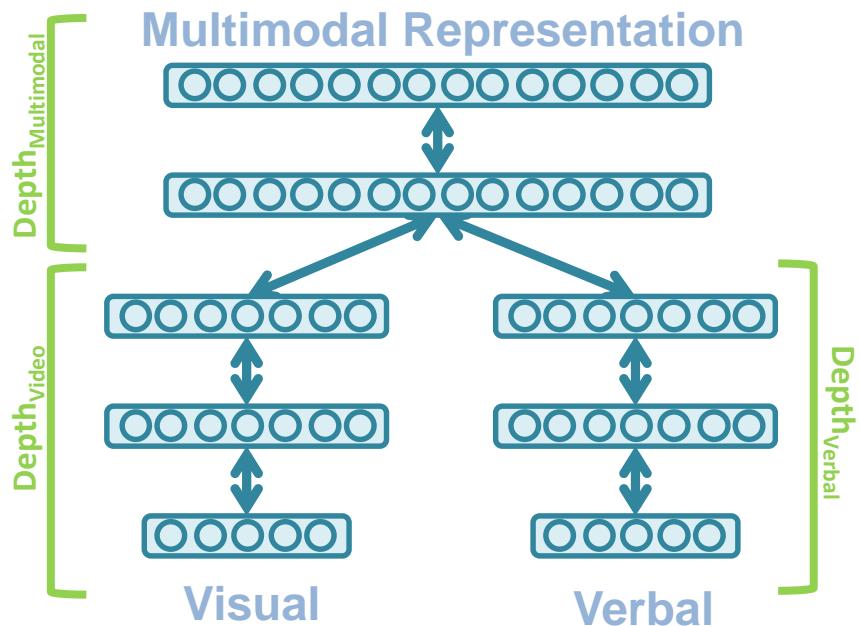
[Srivastava and Salakhutdinov, NIPS 2012]

- Multimodal Deep Boltzmann Machine

Audio-visual emotion recognition

[Kim et al., ICASSP 2013]

- Deep Boltzmann Machine



Multimodal Vector Space Arithmetic

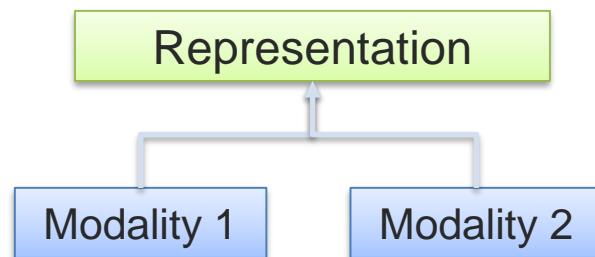


[Kiros et al., Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models, 2014]

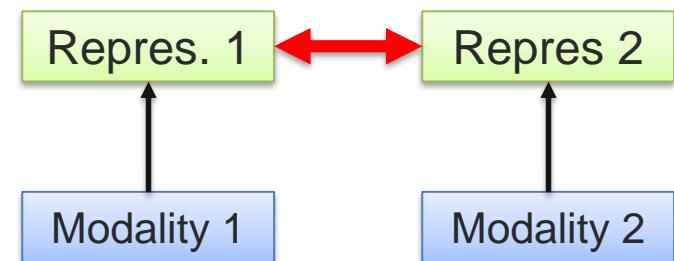
Core Challenge 1: Representation

Definition: Learning how to represent and summarize multimodal data in a way that exploits the complementarity and redundancy.

A Joint representations:



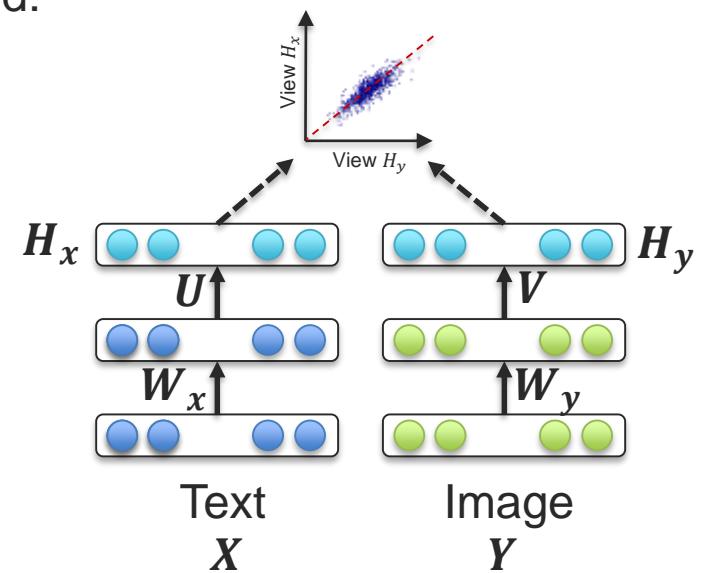
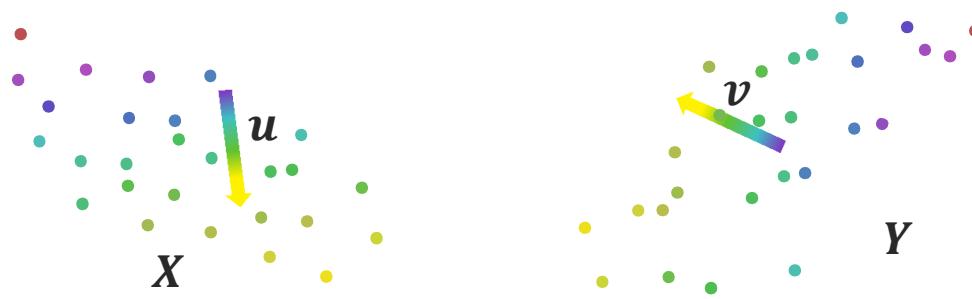
B Coordinated representations:



Coordinated Representation: Deep CCA

Learn linear projections that are maximally correlated:

$$(\mathbf{u}^*, \mathbf{v}^*) = \underset{\mathbf{u}, \mathbf{v}}{\operatorname{argmax}} \operatorname{corr}(\mathbf{u}^T \mathbf{X}, \mathbf{v}^T \mathbf{Y})$$

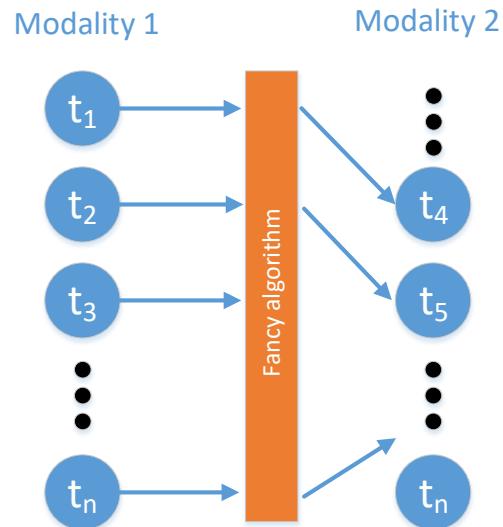


Andrew et al., ICML 2013



Core Challenge 2: Alignment

Definition: Identify the direct relations between (sub)elements from two or more different modalities.



A Explicit Alignment

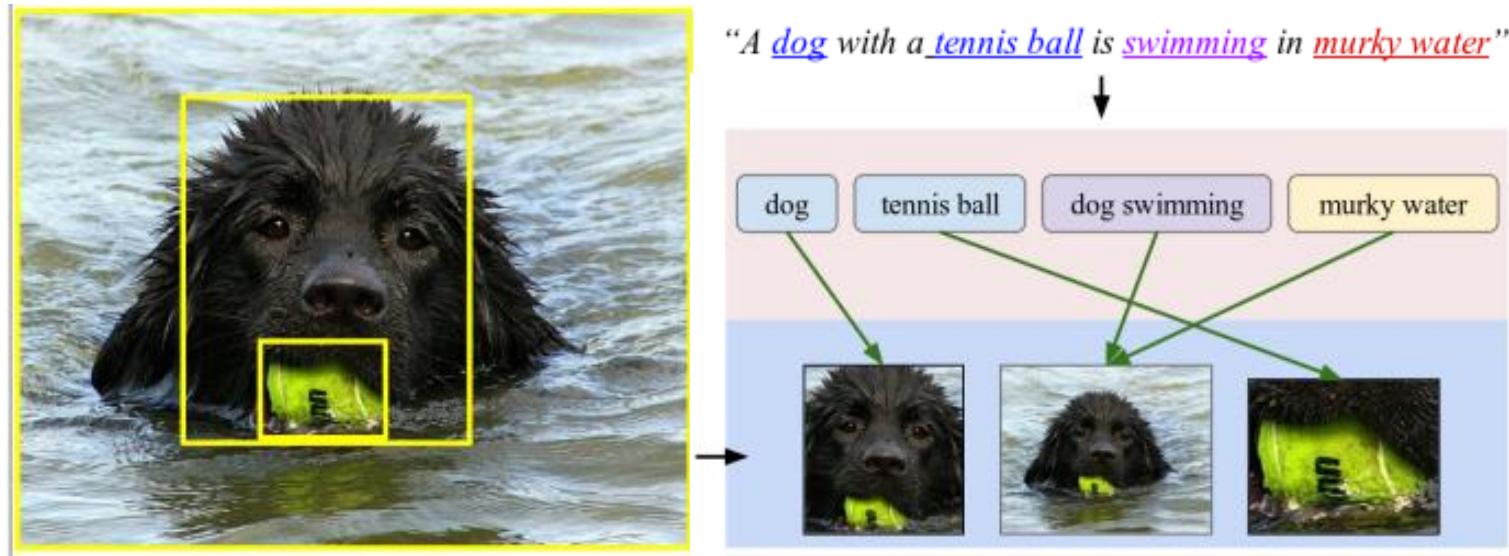
The goal is to directly find correspondences between elements of different modalities

B Implicit Alignment

Uses internally latent alignment of modalities in order to better solve a different problem

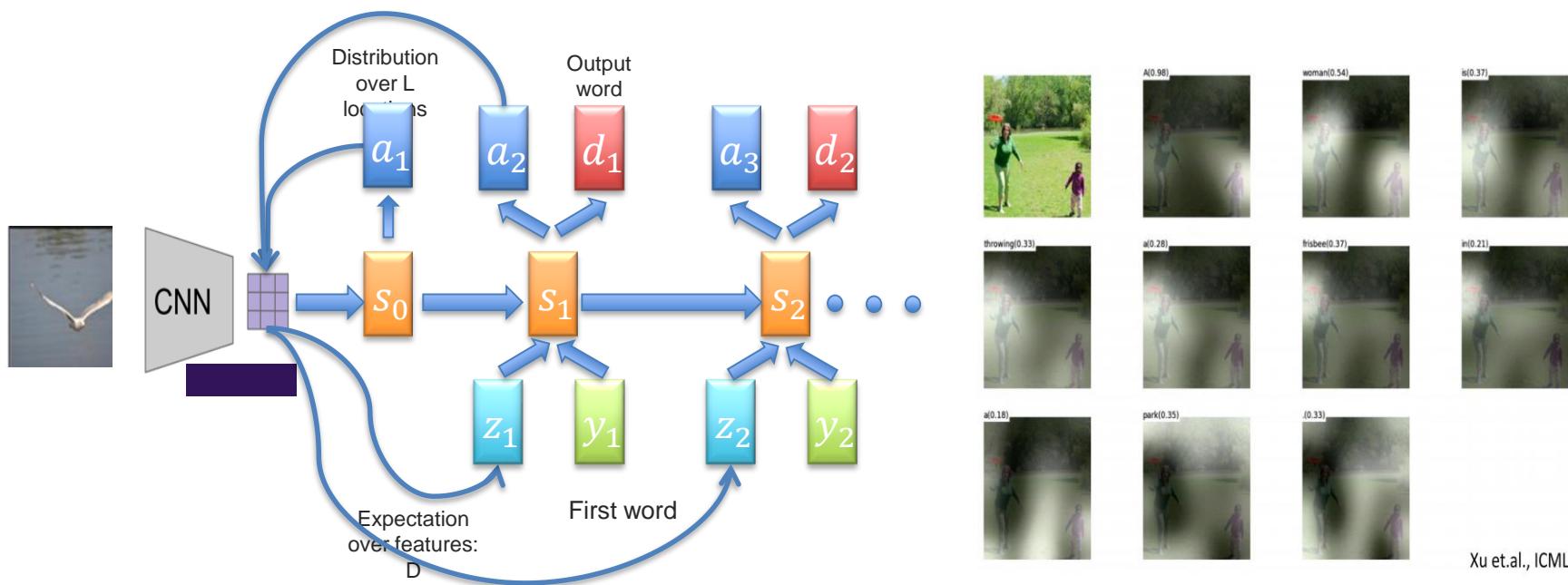


Implicit Alignment



Karpathy et al., Deep Fragment Embeddings for Bidirectional Image Sentence Mapping,
<https://arxiv.org/pdf/1406.5679.pdf>

Attention Models for Image Captioning

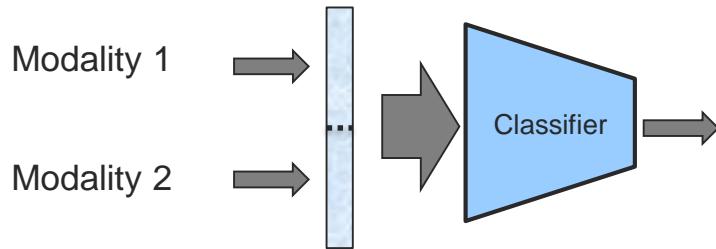


Core Challenge 3: Fusion

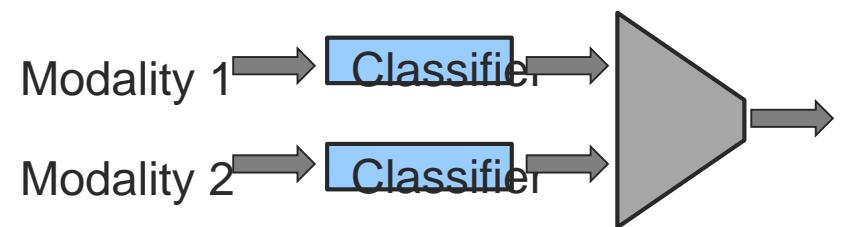
Definition: To join information from two or more modalities to perform a prediction task.

A Model-Agnostic Approaches

1) Early Fusion



2) Late Fusion

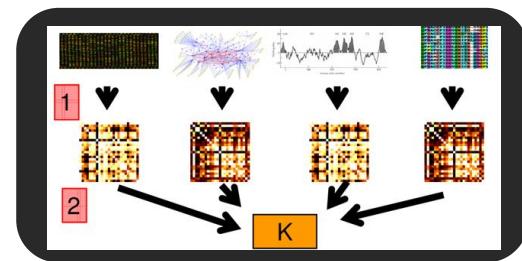


Core Challenge 3: Fusion

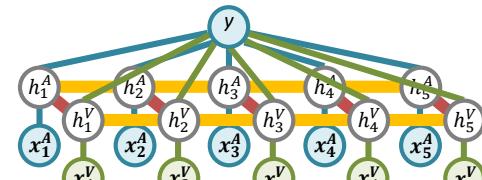
Definition: To join information from two or more modalities to perform a prediction task.

B Model-Based (Intermediate) Approaches

- 1) Deep neural networks
- 2) Kernel-based methods
- 3) Graphical models



Multiple kernel learning



Multi-View Hidden CRF

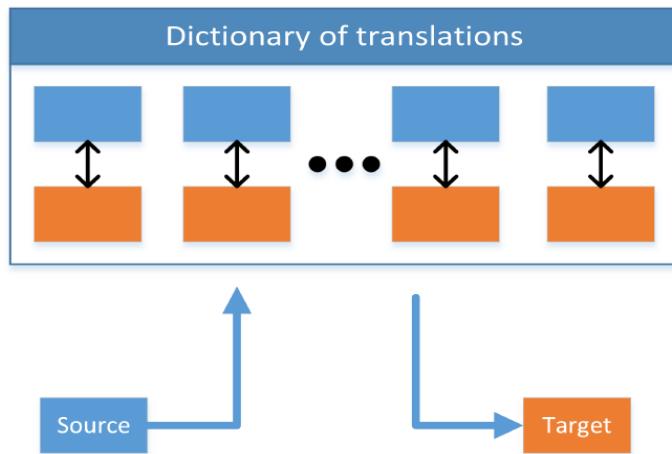


Core Challenge 4: Translation

Definition: Process of changing data from one modality to another, where the translation relationship can often be open-ended or subjective.

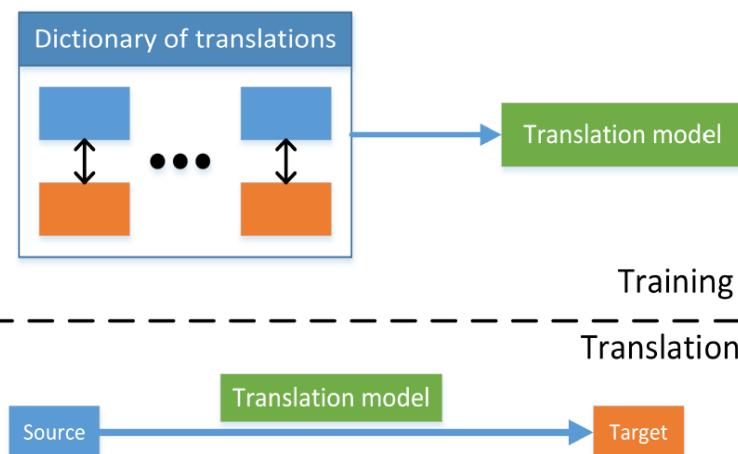
A

Example-based



B

Model-driven



Core Challenge 4 – Translation



Visual gestures
(both speaker and
listener gestures)

Transcriptions
+
Audio streams

Marsella et al., Virtual character performance from speech, SIGGRAPH/Eurographics Symposium on Computer Animation, 2013

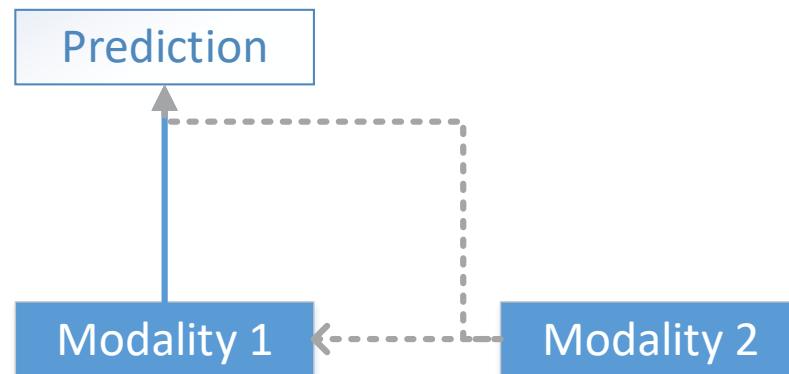


Language Technologies Institute

Carnegie Mellon University

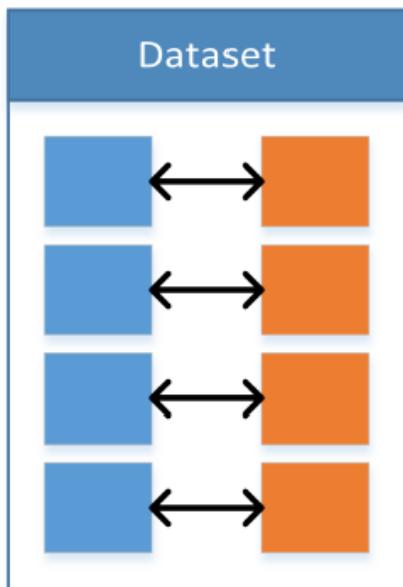
Core Challenge 5: Co-Learning

Definition: Transfer knowledge between modalities, including their representations and predictive models.

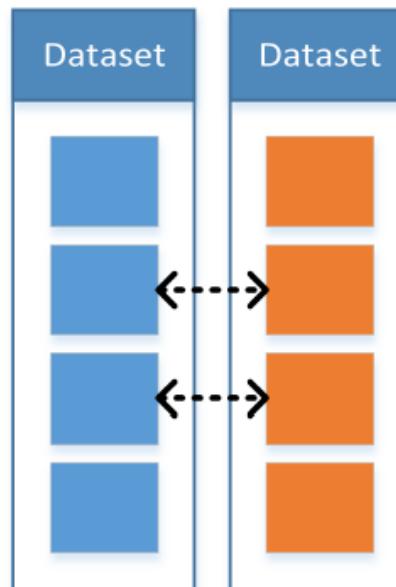


Core Challenge 5: Co-Learning

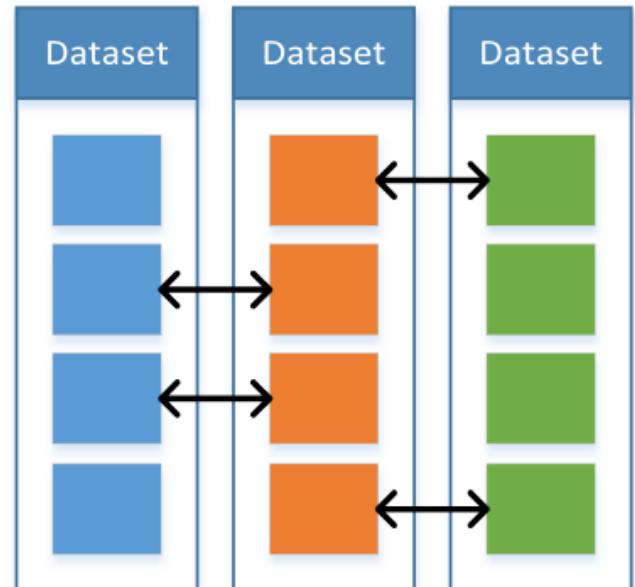
(A) Parallel



(B) Non-Parallel



(C) Hybrid



Taxonomy of Multimodal Research

[<https://arxiv.org/abs/1705.09406>]

Representation

- Joint
 - Neural networks
 - Graphical models
 - Sequential
- Coordinated
 - Similarity
 - Structured

Translation

- Example-based
 - Retrieval
 - Combination
- Model-based
 - Grammar-based

- Encoder-decoder
- Online prediction

Alignment

- Explicit
 - Unsupervised
 - Supervised
- Implicit
 - Graphical models
 - Neural networks

Fusion

- Model agnostic
 - Early fusion
 - Late fusion
 - Hybrid fusion

- Model-based
 - Kernel-based
 - Graphical models
 - Neural networks

Co-learning

- Parallel data
 - Co-training
 - Transfer learning
- Non-parallel data
 - Zero-shot learning
 - Concept grounding
 - Transfer learning
- Hybrid data
 - Bridging

Tadas Baltrusaitis, Chaitanya Ahuja, and Louis-Philippe Morency, Multimodal Machine Learning: A Survey and Taxonomy



Language Technologies Institute

Carnegie Mellon University

Multimodal Applications

[<https://arxiv.org/abs/1705.09406>]

APPLICATIONS	CHALLENGES				
	REPRESENTATION	TRANSLATION	FUSION	ALIGNMENT	CO-LEARNING
Speech Recognition and Synthesis Audio-visual Speech Recognition (Visual) Speech Synthesis	✓ ✓	✓	✓	✓	✓
Event Detection Action Classification Multimedia Event Detection	✓ ✓		✓ ✓		✓ ✓
Emotion and Affect Recognition Synthesis	✓ ✓	✓	✓	✓	✓
Media Description Image Description Video Description Visual Question-Answering Media Summarization	✓ ✓ ✓ ✓	✓ ✓ ✓ ✓	✓ ✓ ✓ ✓	✓ ✓ ✓	✓ ✓ ✓
Multimedia Retrieval Cross Modal retrieval Cross Modal hashing	✓ ✓	✓		✓	✓

Tadas Baltrusaitis, Chaitanya Ahuja, and Louis-Philippe Morency, Multimodal Machine Learning: A Survey and Taxonomy



Language Technologies Institute

Carnegie Mellon University

Basic Concepts: Score and Loss Functions

Linear Classification (e.g., neural network)

Image



(Size: 32*32*3)

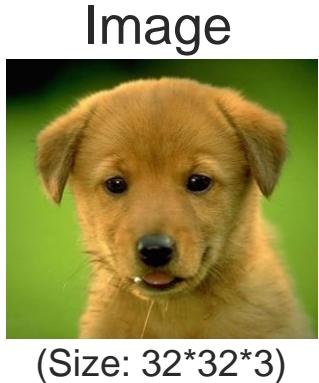


?

- 1. Define a (linear) score function**
- 2. Define the loss function (possibly nonlinear)**
- 3. Optimization**



1) Score Function



Duck ?
Cat ?
Dog ?
Pig ?
Bird ?

**What should be
the prediction
score for each
label class?**

For linear classifier:

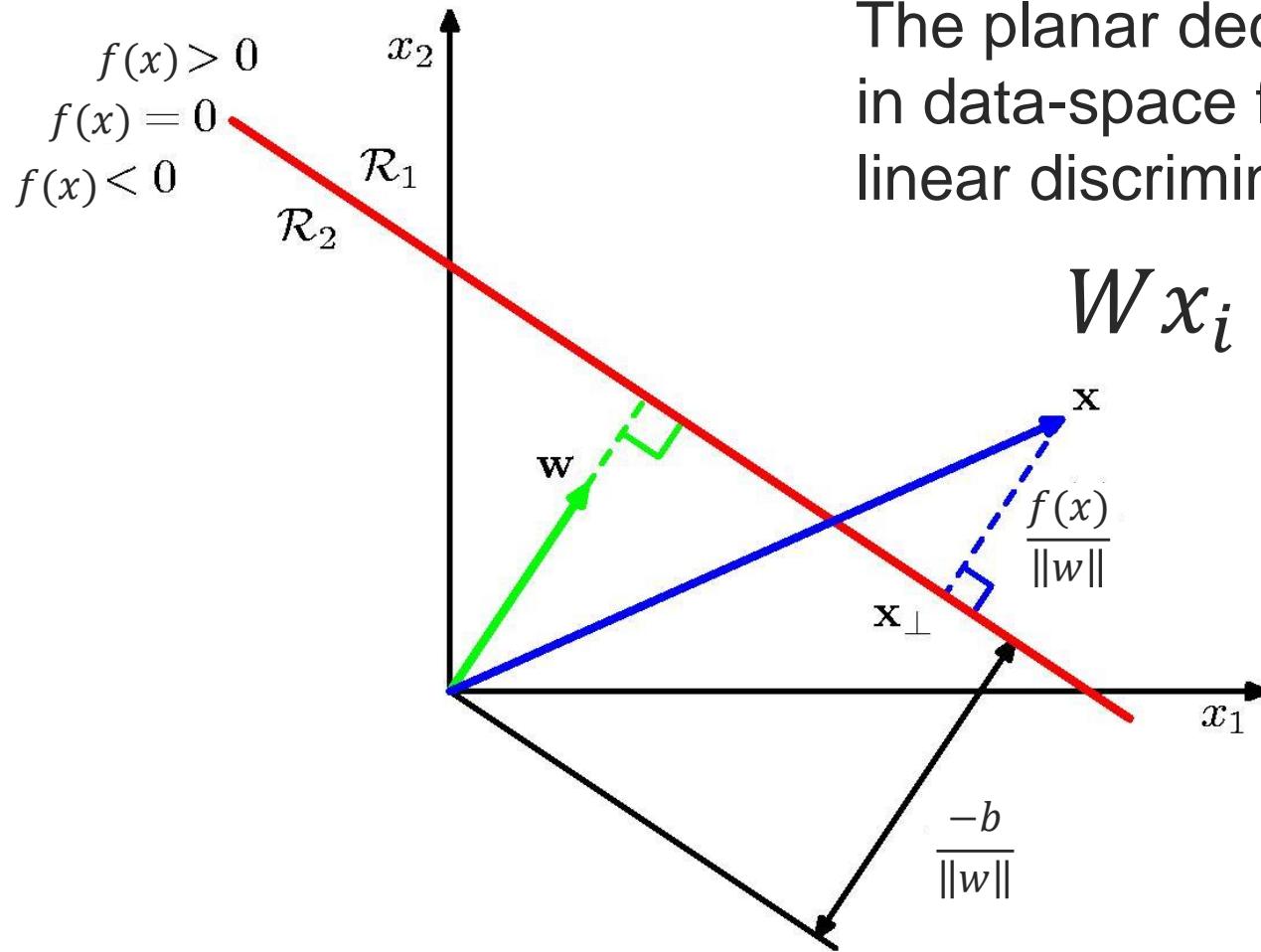
$$f(x_i; W, b) = Wx_i + b$$

Class score [10x1] Input observation (i^{th} element of the dataset) [3072x1]

Weights [10x3072] Bias vector [10x1]

Parameters [10x3073]

Interpreting a Linear Classifier



Some Notation Tricks – Multi-Label Classification

$$W = [W_1 \quad W_2 \quad \dots \quad W_N]$$

$$f(x_i; W, b) = Wx_i + b \quad \longrightarrow \quad f(x_i; W) = Wx_i$$

Weights \times Input + Bias
[10x3072] [3072x1] [10x1]

Weights \times Input
[10x3073] [3073x1]

The bias vector will
be the last column of
the weight matrix

Add a “1” at the
end of the input
observation vector



Some Notation Tricks

General formulation of linear classifier: $f(x_i; W, b)$

“dog” linear classifier:

$$f(x_i; W_{\text{dog}}, b_{\text{dog}}) \quad \text{or}$$

$$f(x_i; W, b)_{\text{dog}} \quad \text{or} \quad f_{\text{dog}}$$

Linear classifier for label j:

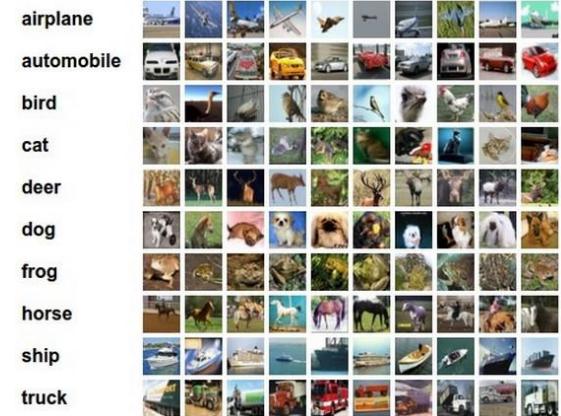
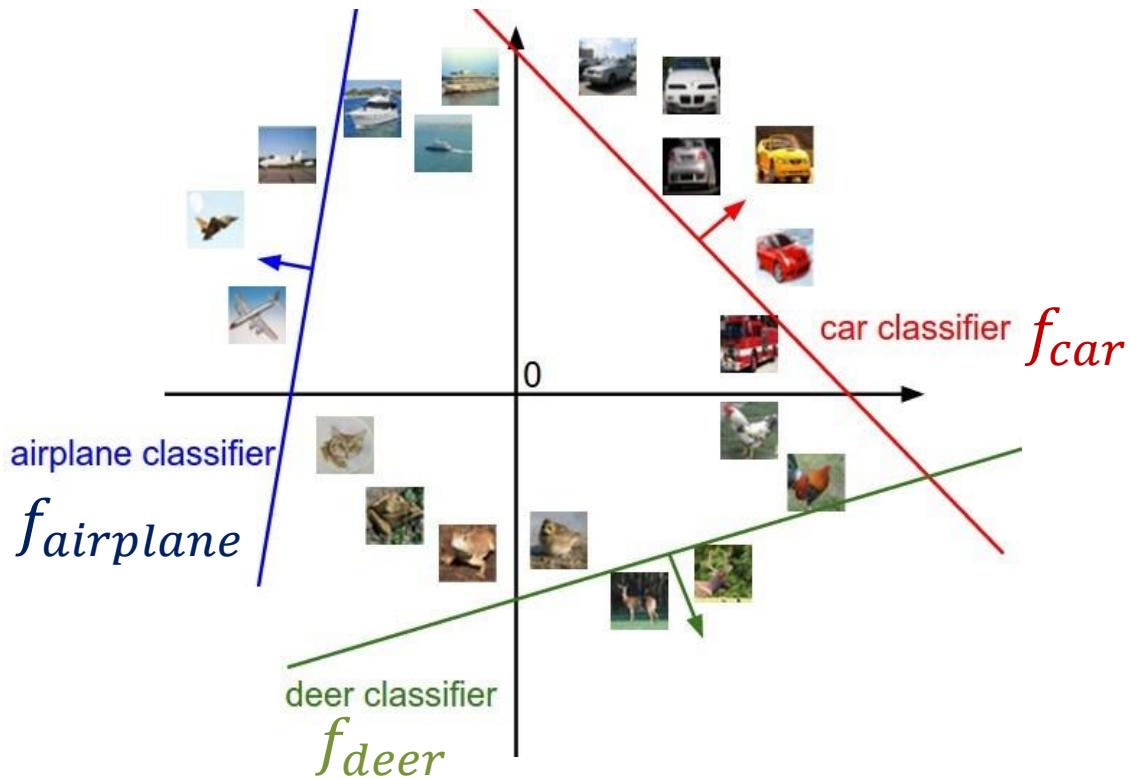
$$f(x_i; W_j, b_j) \quad \text{or}$$

$$f(x_i; W, b)_j \quad \text{or} \quad f_j$$



Interpreting Multiple Linear Classifiers

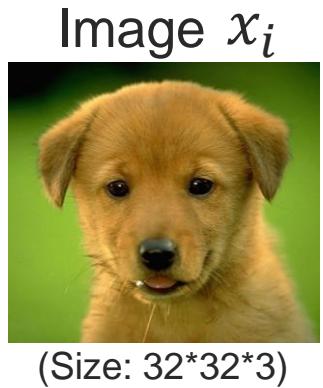
$$f(x_i; W_j, b_j) = W_j x_i + b_j$$



CIFAR-10 object
recognition dataset

Linear Classification: 2) Loss Function

(or cost function or objective)



Multi-class problem

Scores	Label	→ Loss
$f(x_i; W)$	$y_i = 2$ (dog)	$L_i = ?$
0 (duck) ?	-12.3	
1 (cat) ?	45.6	
2 (dog) ?	98.7	
3 (pig) ?	12.2	
4 (bird) ?	-45.3	

How to assign
only one number
representing
how “unhappy”
we are about
these scores?

The loss function quantifies the amount by which
the prediction scores deviate from the actual values.

A first challenge: how to normalize the scores?



Language

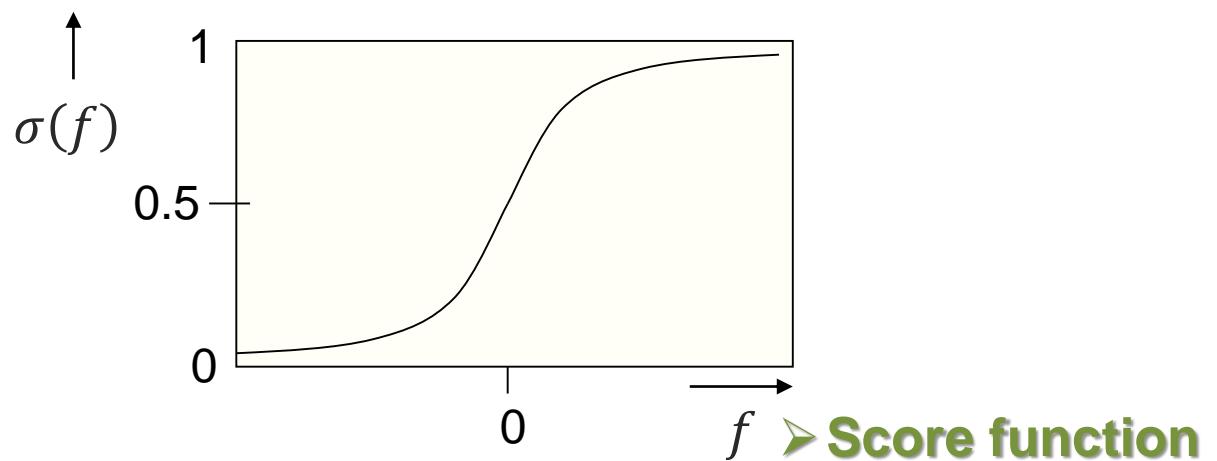
University

First Loss Function: Cross-Entropy Loss

(or logistic loss)

Logistic function:

$$\sigma(f) = \frac{1}{1 + e^{-f}}$$



First Loss Function: Cross-Entropy Loss

(or logistic loss)

Logistic function:

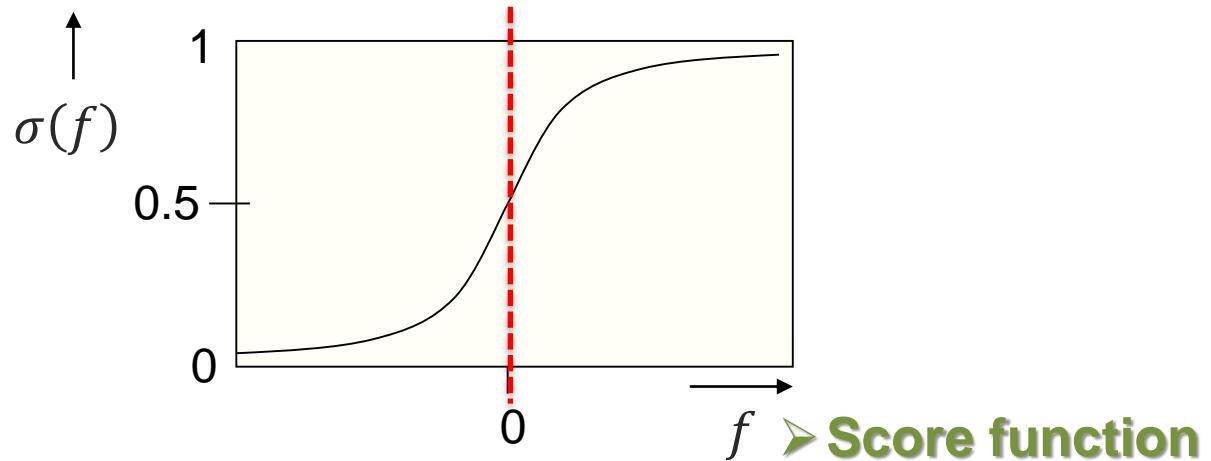
$$\sigma(f) = \frac{1}{1 + e^{-f}}$$

Logistic regression:
(two classes)

$$p(y_i = "dog" | x_i; w) = \sigma(w^T x_i)$$

= true

for two-class problem



First Loss Function: Cross-Entropy Loss

(or logistic loss)

Logistic function:

$$\sigma(f) = \frac{1}{1 + e^{-f}}$$

Logistic regression:
(two classes)

$$p(y_i = "dog" | x_i; w) = \sigma(w^T x_i)$$

= true
for two-class problem

Softmax function:
(multiple classes)

$$p(y_i | x_i; W) = \frac{e^{f_{y_i}}}{\sum_j e^{f_j}}$$



First Loss Function: Cross-Entropy Loss

(or logistic loss)

Cross-entropy loss:

$$L_i = -\log \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right)$$

Softmax function

Minimizing the negative log likelihood.



Second Loss Function: Hinge Loss

(or max-margin loss or Multi-class SVM loss)

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

↑
loss due to
example i

↑
sum over all
incorrect labels

difference between the correct class
score and incorrect class score



Second Loss Function: Hinge Loss

(or max-margin loss or Multi-class SVM loss)

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

↑
e.g. 10

Example: $f(x_i, W) = [13, -7, 11]$

$$y_i = 0$$

$$L_i = \max(0, -7 - 13 + 10) + \max(0, 11 - 13 + 10)$$

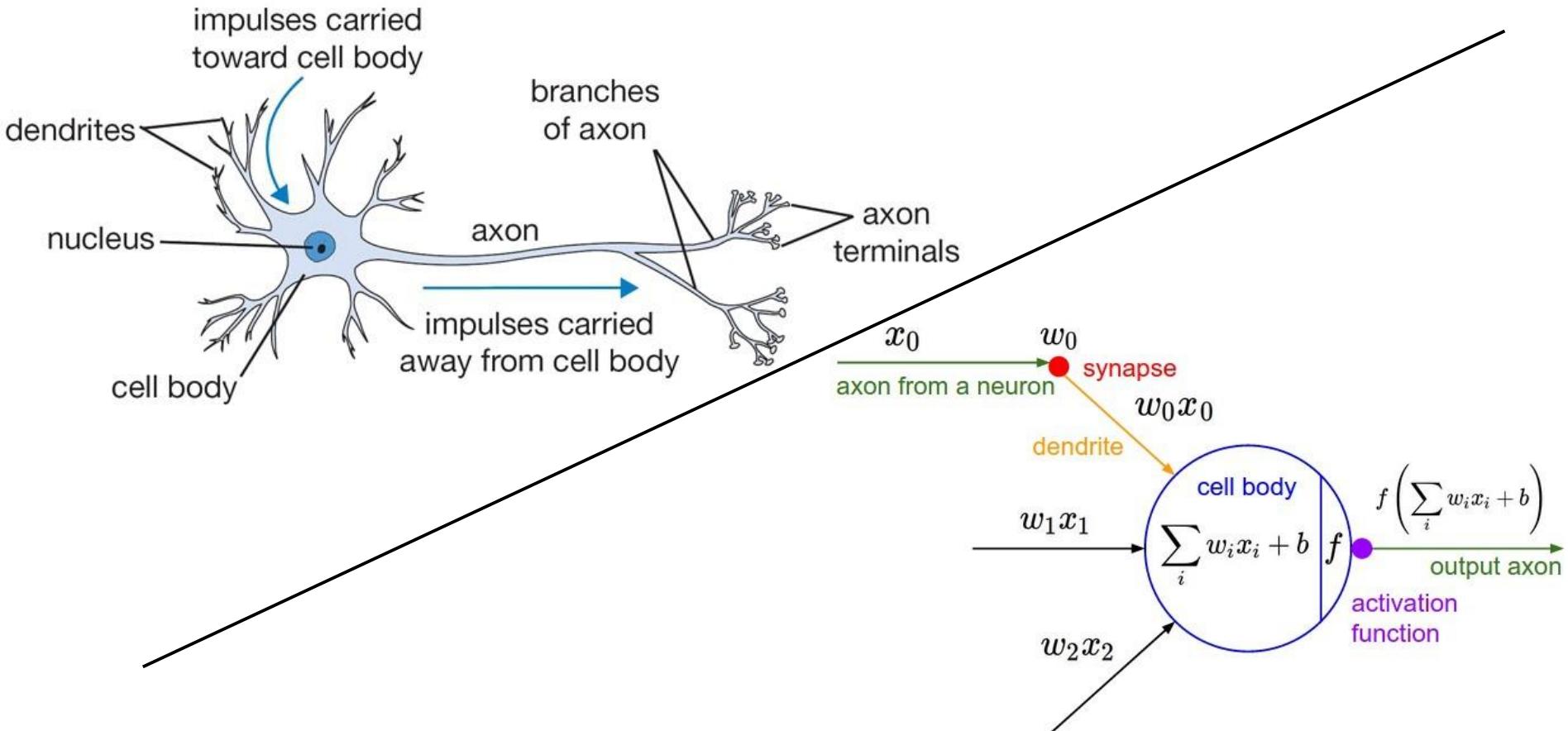
How to find the optimal W ?



Basic Concepts: Neural Networks

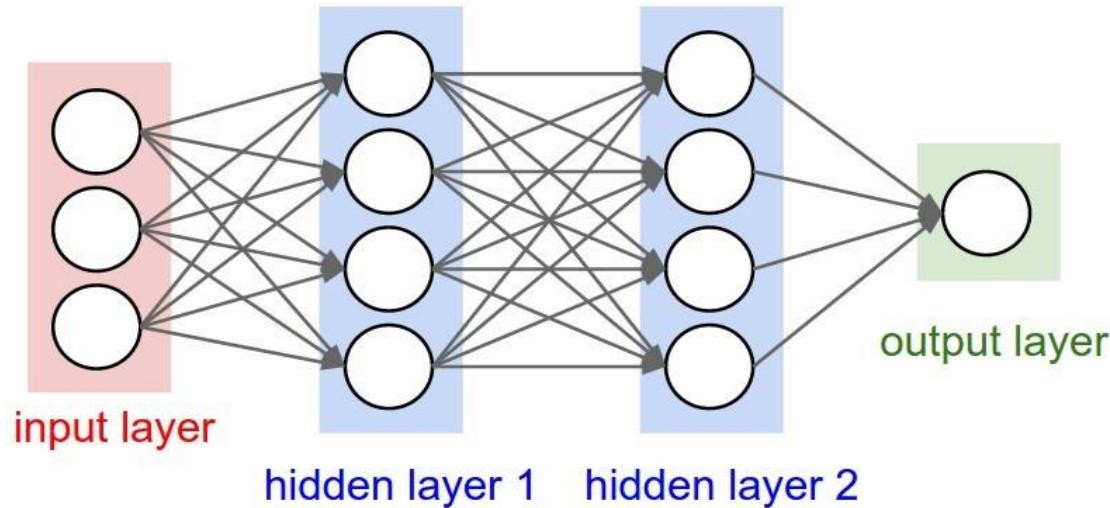
Neural Networks – inspiration

- Made up of artificial neurons



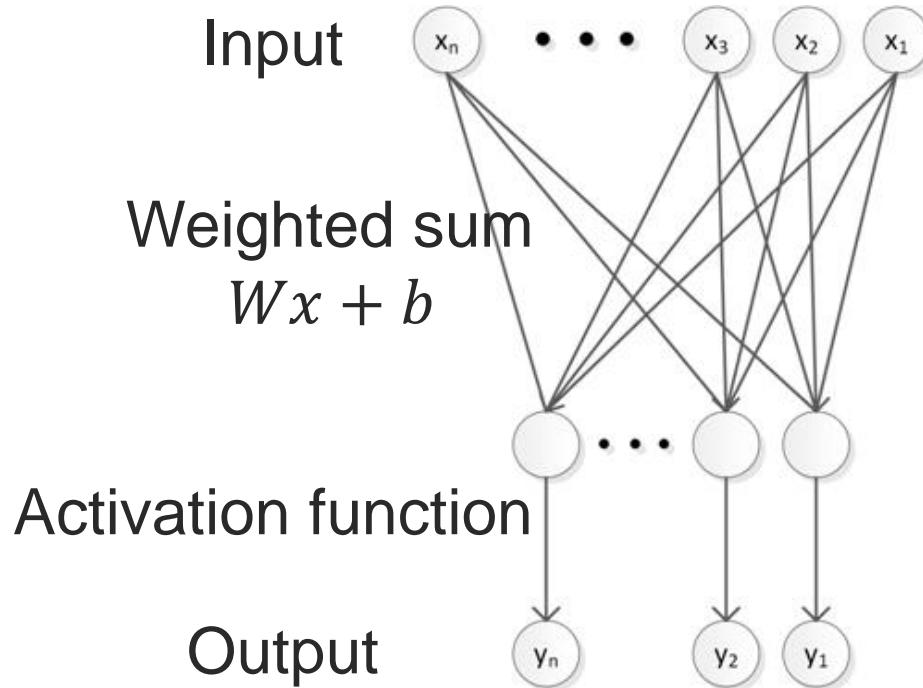
Neural Networks – score function

- Made up of artificial neurons
 - Linear function (dot product) followed by a nonlinear activation function
- Example a Multi Layer Perceptron



Basic NN building block

- Weighted sum followed by an activation function

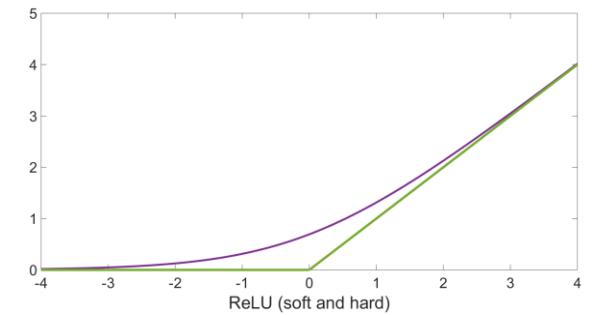
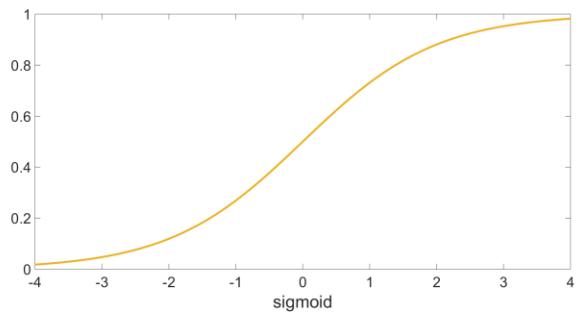
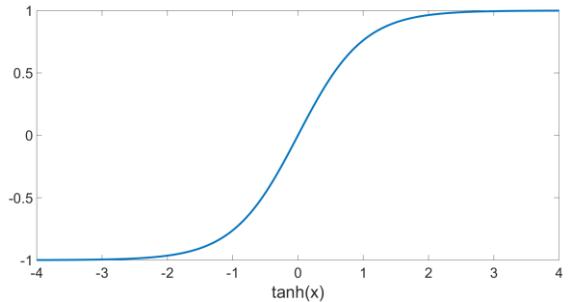


$$y = f(Wx + b)$$



Neural Networks – activation function

- $f(x) = \tanh(x)$
- Sigmoid - $f(x) = (1 + e^{-x})^{-1}$
- Linear – $f(x) = ax + b$
- ReLU $f(x) = \max(0, x) \sim \log(1 + \exp(x))$
 - Rectifier Linear Units
 - Faster training - no gradient vanishing
 - Induces sparsity



Neural Networks – loss function

- Already discussed it – cross-entropy, Euclidean loss, cosine similarity, etc.
- Combine it with the score function to have an end-to-end training objective
- As example use Euclidean loss for data-point i

$$L_i = (f(x_i) - y_i)^2 = (f_{3;W_3}(f_{2;W_2}(f_{1;W_1}(x_i))))^2$$

- Full loss is computed across all training samples

$$L = \sum_i (f(x_i) - y_i)^2$$



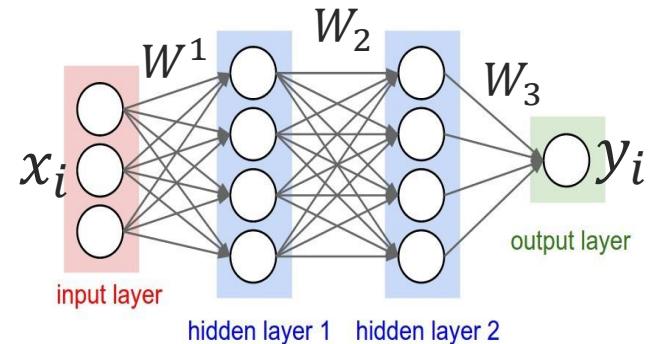
Multi-Layer Feedforward Network

Activation functions (individual layers)

$$f_{1;W_1}(x) = \sigma(W_1 x + b_1)$$

$$f_{2;W_2}(x) = \sigma(W_2 x + b_2)$$

$$f_{3;W_3}(x) = \sigma(W_3 x + b_3)$$



Score function

$$y_i = f(x_i) = f_{3;W_3}(f_{2;W_2}(f_{1;W_1}(x_i)))$$

Loss function (e.g., Euclidean loss)

$$L_i = (f(x_i) - y_i)^2 = (f_{3;W_3}(f_{2;W_2}(f_{1;W_1}(x_i))))^2$$



Basic Concepts: Optimization

Optimizing a generic function

- We want to find a minimum (or maximum) of a generic function
- How do we do that?
 - Searching everywhere (global optimum) is computationally infeasible
 - We could search randomly from our starting point (mostly picked at random) – impractical and not accurate
 - Instead we can follow the gradient



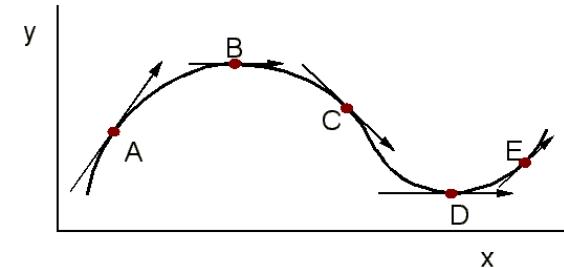
What is a gradient?

- Geometrically

- Points in the direction of the greatest rate of increase of the function and its magnitude is the slope of the graph in that direction

- More formally in 1D

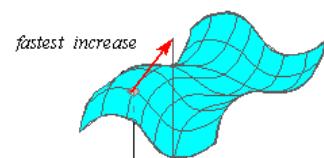
$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$



- In higher dimensions

$$\frac{\partial f}{\partial x_i}(a_1, \dots, a_n) = \lim_{h \rightarrow 0} \frac{f(a_1, \dots, a_i + h, \dots, a_n) - f(a_1, \dots, a_i, \dots, a_n)}{h}$$

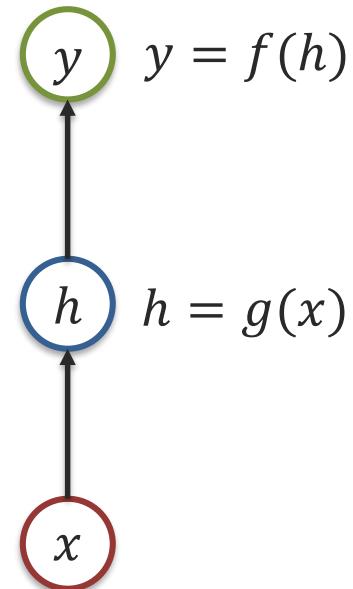
- In multiple dimension, the **gradient** is the vector of (partial derivatives) and is called a **Jacobian**.



Gradient Computation

Chain rule:

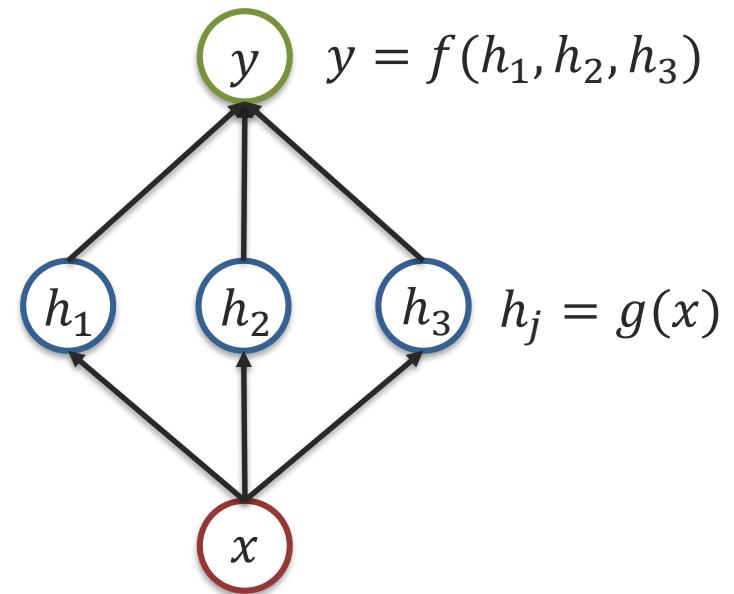
$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial h} \frac{\partial h}{\partial x}$$



Optimization: Gradient Computation

Multiple-path chain rule:

$$\frac{\partial y}{\partial x} = \sum_j \frac{\partial y}{\partial h_j} \frac{\partial h_j}{\partial x}$$



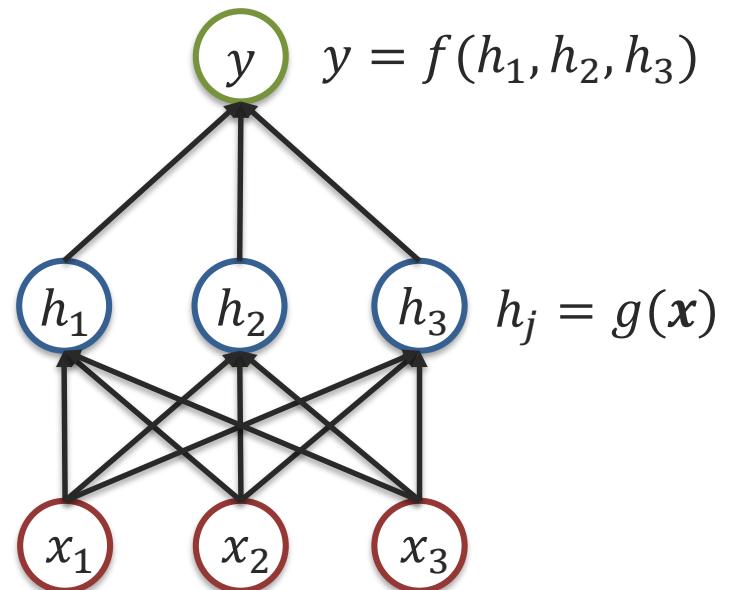
Optimization: Gradient Computation

Multiple-path chain rule:

$$\frac{\partial y}{\partial x_1} = \sum_j \frac{\partial y}{\partial h_j} \frac{\partial h_j}{\partial x_1}$$

$$\frac{\partial y}{\partial x_2} = \sum_j \frac{\partial y}{\partial h_j} \frac{\partial h_j}{\partial x_1}$$

$$\frac{\partial y}{\partial x_3} = \sum_j \frac{\partial y}{\partial h_j} \frac{\partial h_j}{\partial x_1}$$



Optimization: Gradient Computation

Vector representation:

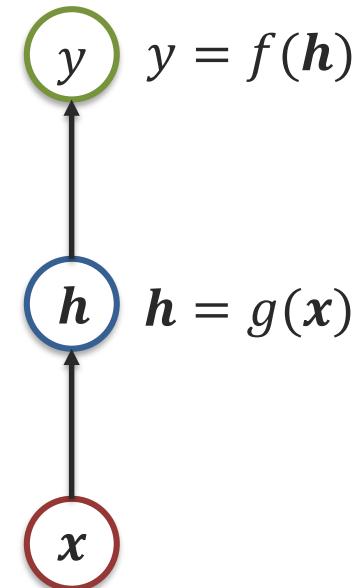
$$\nabla_x y = \left[\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \frac{\partial y}{\partial x_3} \right]$$

Gradient

$$\nabla_x y = \left(\frac{\partial h}{\partial x} \right)^T \nabla_h y$$

“local” Jacobian
(matrix of size $|h| \times |x|$ computed
using partial derivatives)

“backprop” Gradient



Backpropagation Algorithm (efficient gradient)

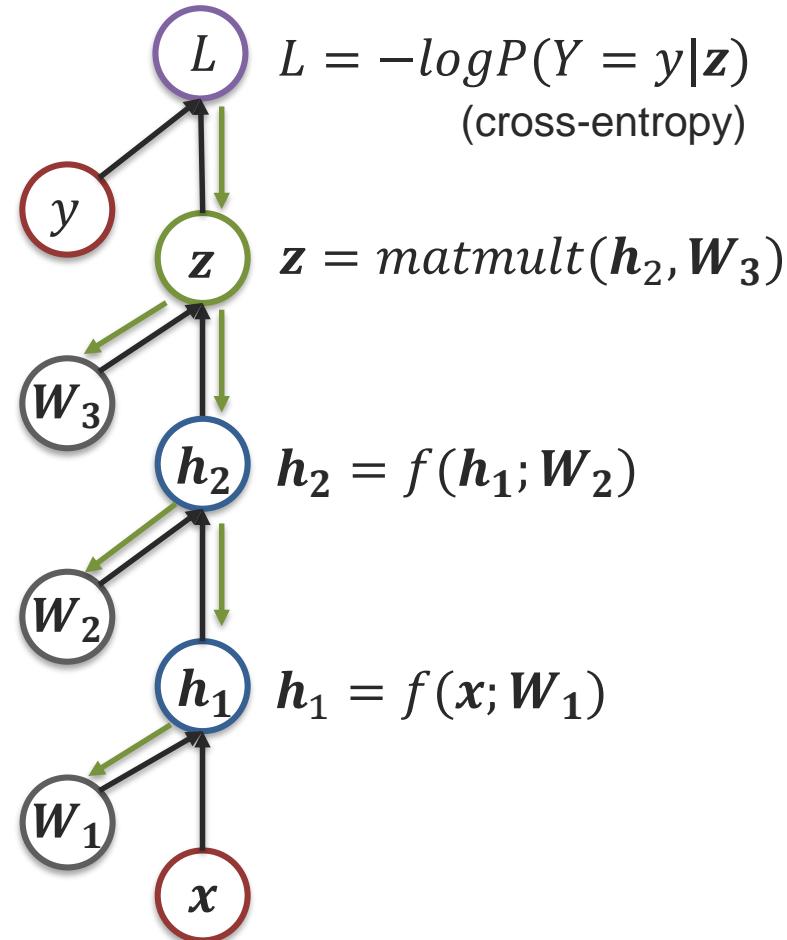
Forward pass

- Following the graph topology, compute value of each unit

Backpropagation pass

- Initialize output gradient = 1
- Compute “local” Jacobian matrix using values from forward pass
- Use the chain rule:

Gradient = “local” Jacobian \times
“backprop” gradient



How to follow the gradient

- Many methods for optimization
 - **Gradient Descent (actually the “simplest” one)**
 - Newton methods (use Hessian – second derivative)
 - Quasi-Newton (use approximate Hessian)
 - BFGS
 - LBFGS
 - Don’t require learning rates (fewer hyperparameters)
 - But, do not work with stochastic and batch methods so rarely used to train modern Neural Networks
- **All of them look at the gradient**
 - Very few non gradient based optimization methods



Parameter Update Strategies

Gradient descent:

$$\theta^{(t+1)} = \theta^t - \epsilon_k \nabla_{\theta} L$$

Annotations for the top equation:

- New model parameters: arrow pointing to $\theta^{(t+1)}$
- Previous parameters: arrow pointing to θ^t
- Learning rate at iteration k: arrow pointing to ϵ_k
- Gradient of our loss function: red box around $\nabla_{\theta} L$

$$\epsilon_k = (1 - \alpha)\epsilon_0 + \alpha \epsilon_{\tau}$$

Annotations for the bottom equation:

- Learning rate at iteration k: green arrow pointing to ϵ_k
- Decay: purple arrow pointing to $(1 - \alpha)$
- Initial learning rate: purple arrow pointing to ϵ_0
- Decay learning rate linearly until iteration τ : red box around ϵ_{τ}

- Extensions:
- Stochastic (“batch”)
 - with momentum
 - AdaGrad
 - RMSProp

Unimodal representations: Language Modality

Unimodal Classification – Language Modality

Written language

A row of five yellow five-pointed star icons, used to represent a rating or review.

Masterful!

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in
disguises who likes to see the subject
tackled in a **humourous** manner.

MARTHA (CON'T)
Look around you. Look at all the great things you've done and the people you've helped.

CLARK
But you've only put up the good things they say about me.

MARTHA

Clark, honey. If I were to use the bad things they say I could cover the barn, the house and the outhouse.

“one-hot” vector

$|x_i|$ = number of words in dictionary

Word-level classification

Part-of-speech ?

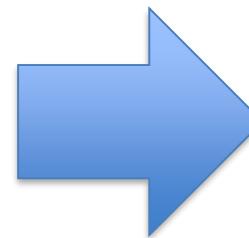
(noun, verb,...)

Sentiment ?

(positive or negative)

Named entity ?

(names of person,...)



Unimodal Classification – Language Modality

Written language

A row of five yellow five-pointed star icons, used to represent a rating or review.

Masterful!

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in
disguises who likes to see the subject
tackled in a humourous manner.

0 of 4 people found this review helpful

Spoken language

MARTHA (CON ' T)

Look around you. Look at all the great things you've done and the people you've helped.

CLARK

But you've only put up the good things they say about me.

MARTHA

Clark, honey. If I were to use the bad things they say I could cover the barn, the house and the outhouse.

“bag-of-word” vector

$|x_i|$ = number of words in dictionary

Document-level classification

Sentiment ?

(positive or negative)

How to Learn (Better) Language Representations?

- **Distribution hypothesis:** Approximate the word meaning by its surrounding words
- Words used in a similar context will lie close together

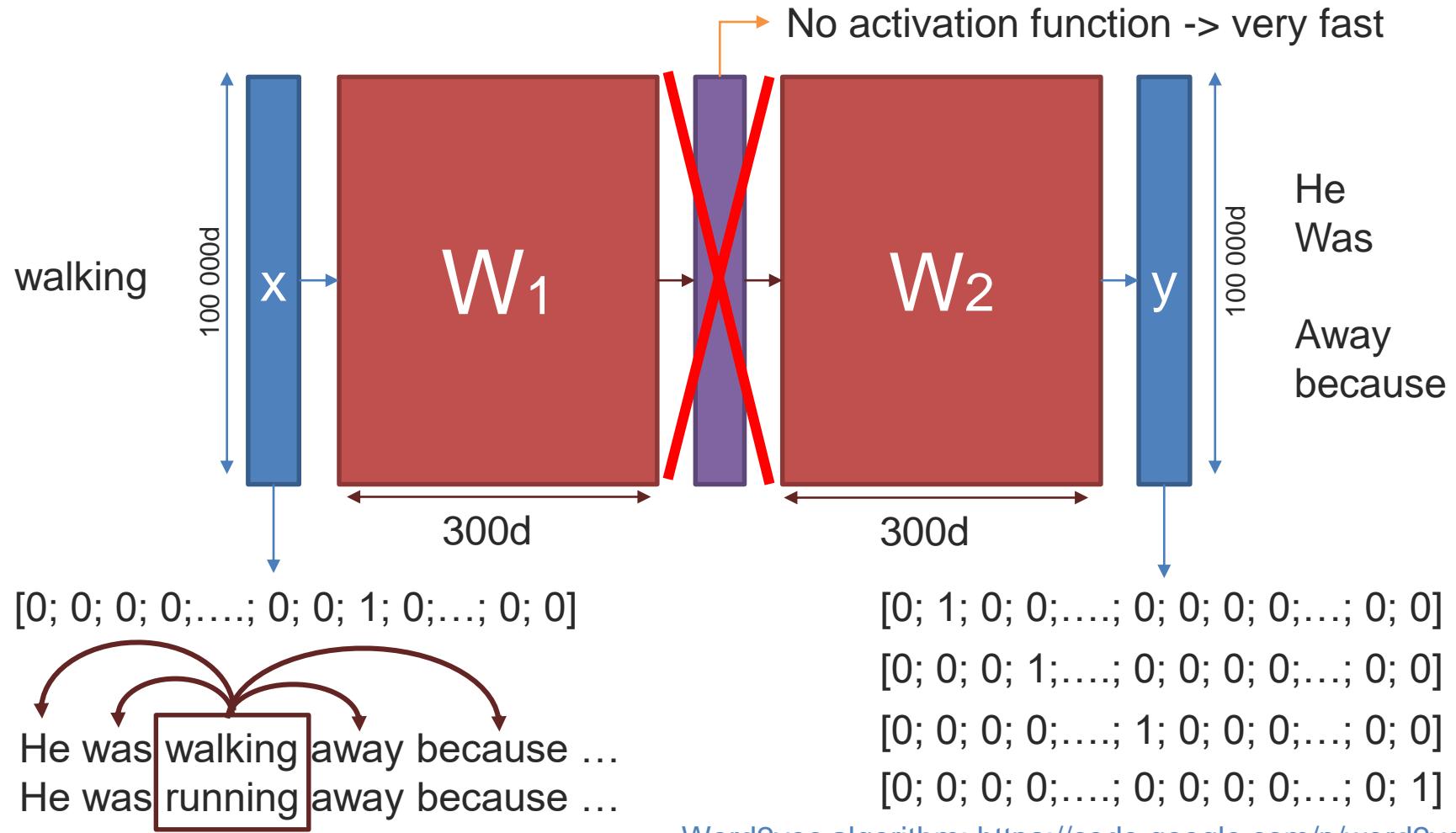
He was walking away because ...
He was running away because ...

- Instead of capturing co-occurrence counts directly, predict surrounding words of every word

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$



How to Learn (Better) Language Representations?



Word2vec algorithm: <https://code.google.com/p/word2vec/>



How to use these word representations

If we would have a vocabulary of 100 000 words:

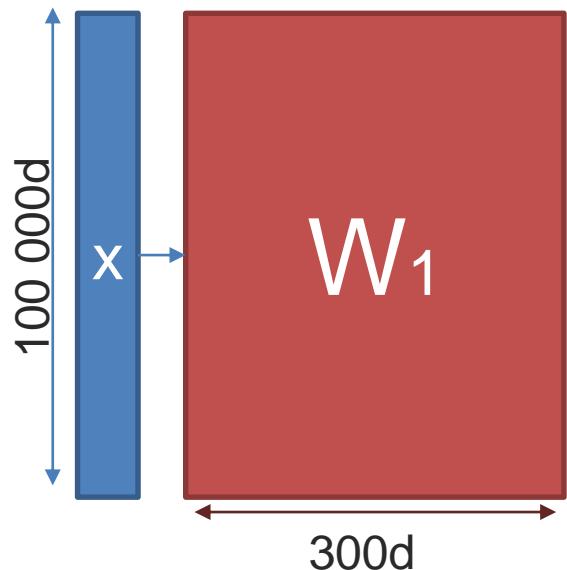
Classic NLP: $\xleftarrow{100\,000 \text{ dimensional vector}}$

Walking: [0; 0; 0; 0; ...; 0; 0; 1; 0; ...; 0; 0]

Running: [0; 0; 0; 0; ...; 0; 0; 0; 0; ...; 1; 0]

→ Similarity = 0.0

↓ Transform: $x' = x^*W$



Goal: $\xleftarrow{300 \text{ dimensional vector}}$

Walking: [0,1; 0,0003; 0; ...; 0,02; 0,08; 0,05]

Running: [0,1; 0,0004; 0; ...; 0,01; 0,09; 0,05]

→ Similarity = 0.9



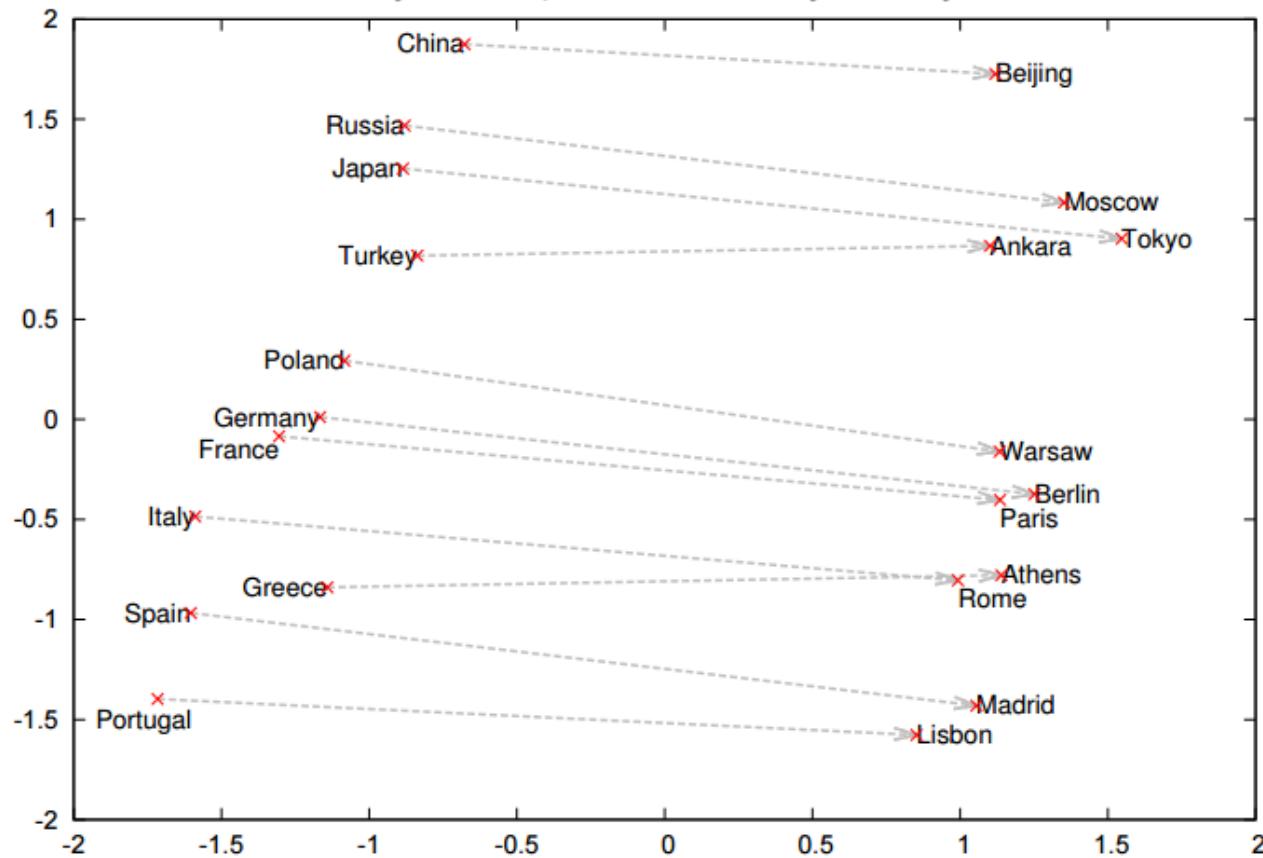
Vector space models of words

- While learning these word representations, we are actually building a vector space in which all words reside with certain relationships between them
- Encodes both syntactic and semantic relationships
- This vector space allows for algebraic operations:

$$\text{Vec(king)} - \text{vec(man)} + \text{vec(woman)} \approx \text{vec(queen)}$$



Vector space models of words: semantic relationships



Trained on the Google news corpus with over 300 billion words

Mikolov et al., "Distributed Representations of Words and Phrases and their Compositionality", NIPS 2013

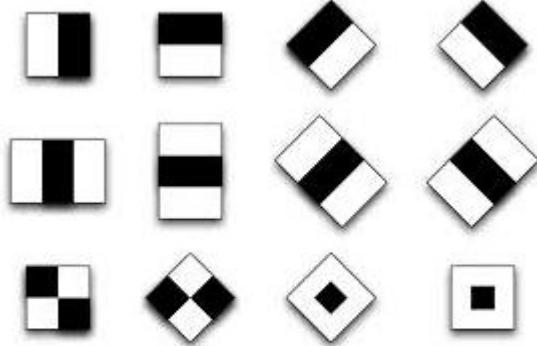
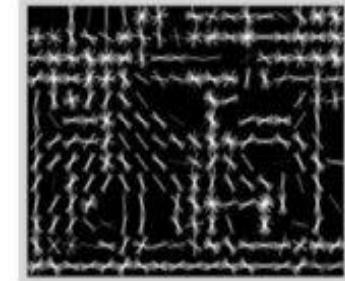
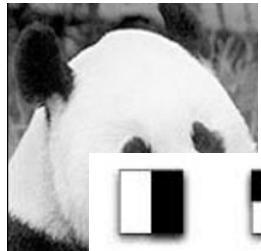


Language Technologies Institute

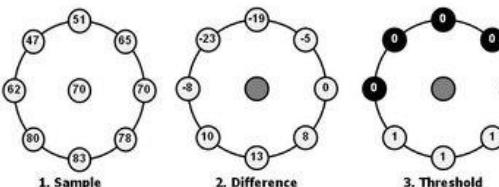
Carnegie Mellon University

Unimodal representations: Visual Modality

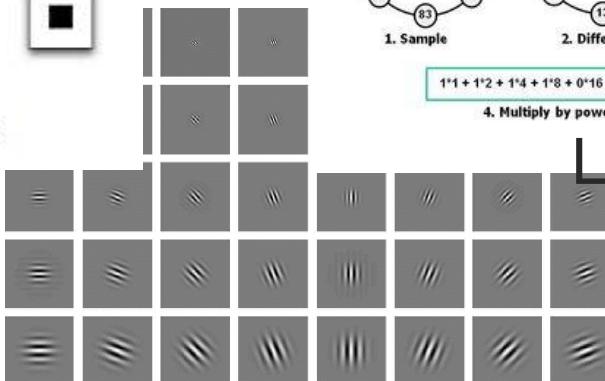
Visual Descriptors



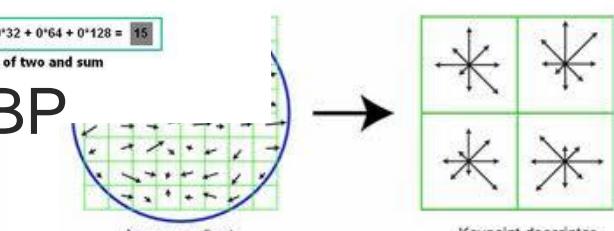
Edge



Haar Wavelets



Optical Flow



SIFT descriptors

Gabor Jets

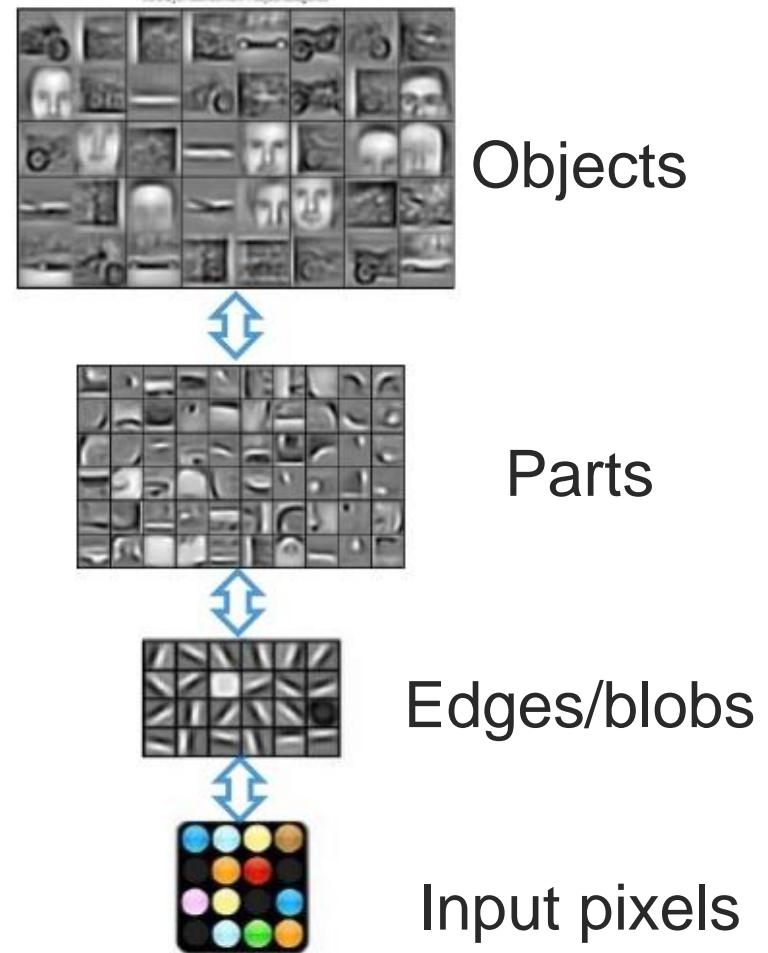


Language Technologies Institute

Carnegie Mellon University

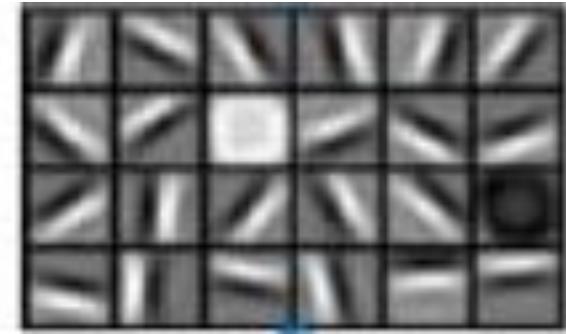
Why use Convolutional Neural Networks

- Using basic Multi Layer Perceptrons does not work well for images
- Intention to build more abstract representation as we go up every layer



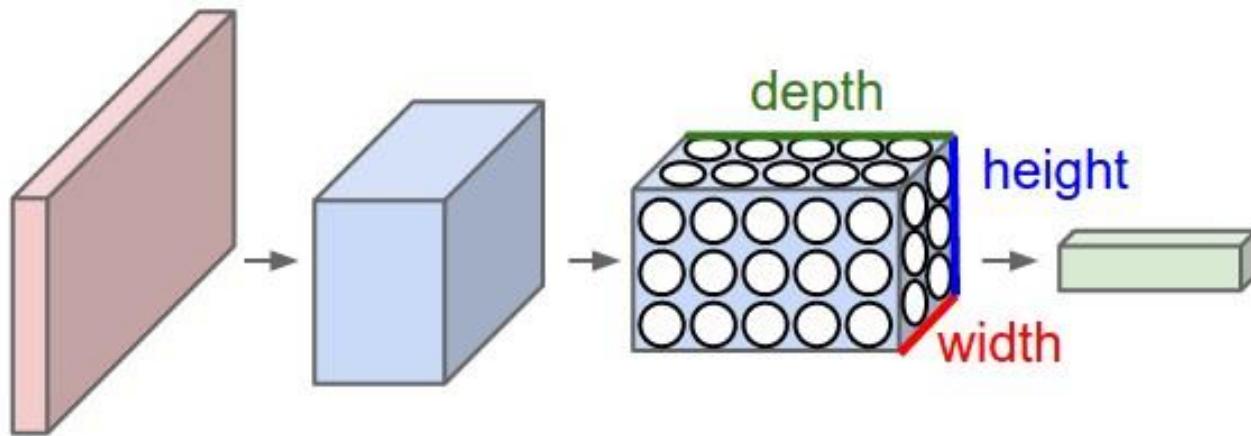
Why not just use an MLP for images (1)?

- MLP connects each pixel in an image to each neuron
- Does not exploit redundancy in image structure
 - Detecting edges, blobs
 - Don't need to treat the top left of image differently from the center
- Too many parameters
 - For a small 200×200 pixel RGB image the first matrix would have $120000 \times n$ parameters for the first layer alone
- MLP does not exploit translation invariance
- MLP does not necessarily encourage visual abstraction

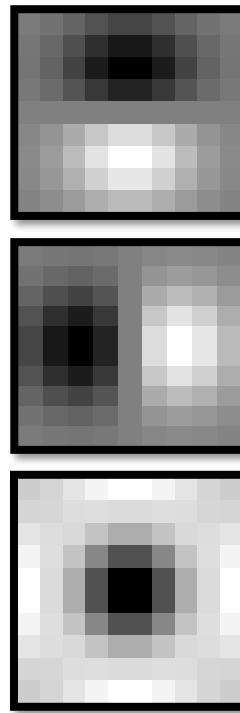
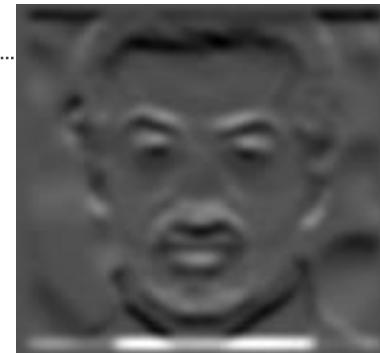


Main differences of CNN from MLP

- Addition of:
 - Convolution layer
 - Pooling layer
- Everything else is the same (loss, score and optimization)
- MLP layer is called Fully Connected layer

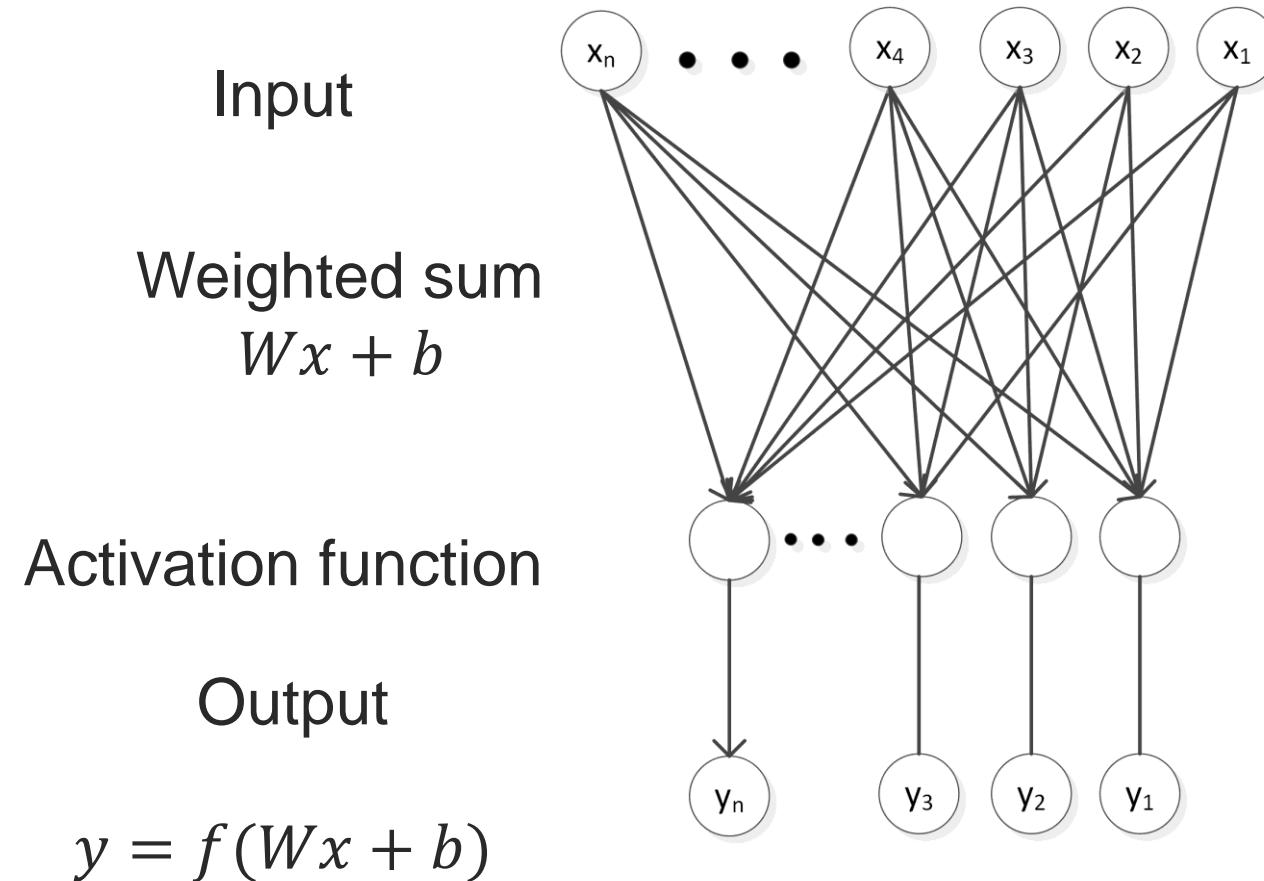


Convolution in 2D

 $*$  $=$ 

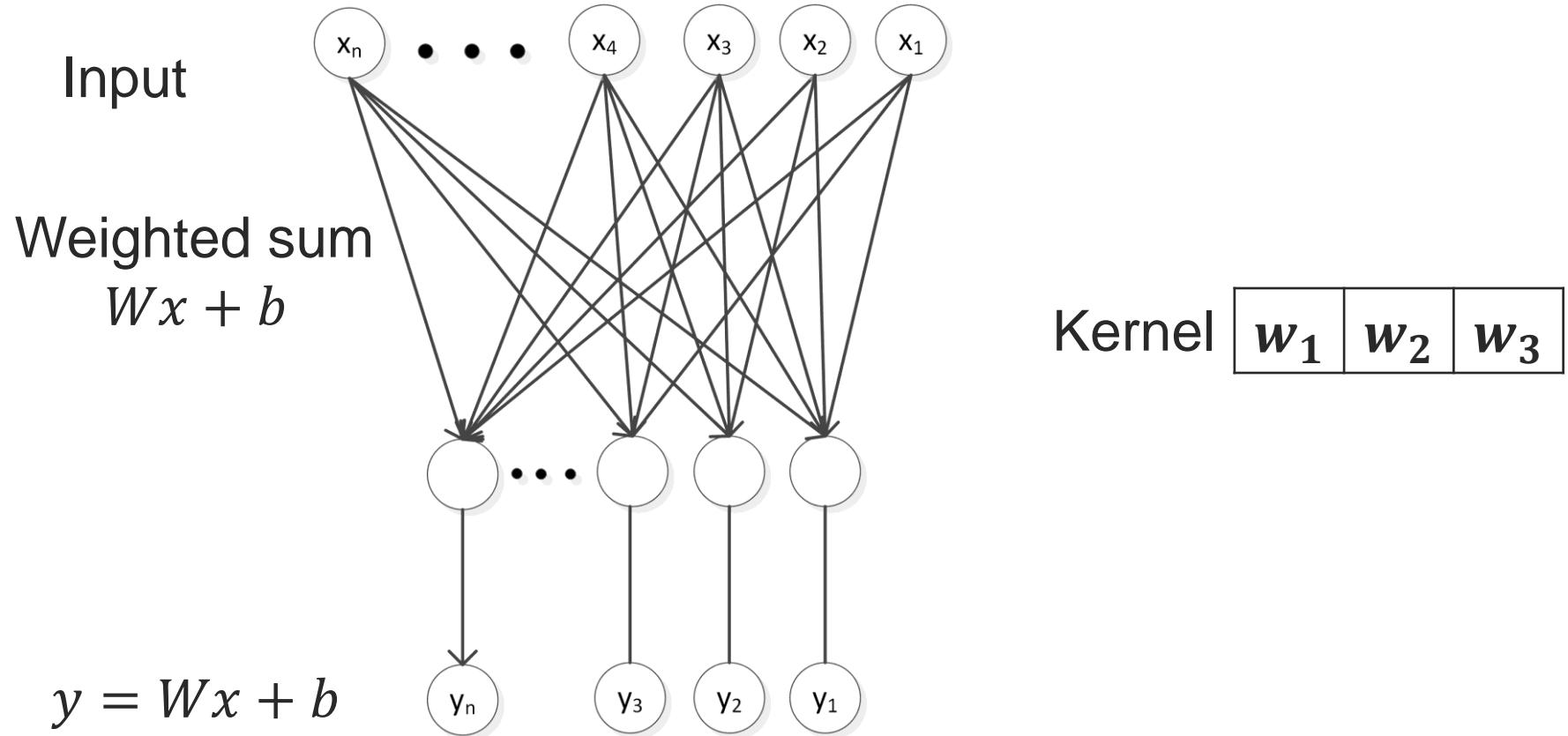
Fully connected layer

- Weighted sum followed by an activation function



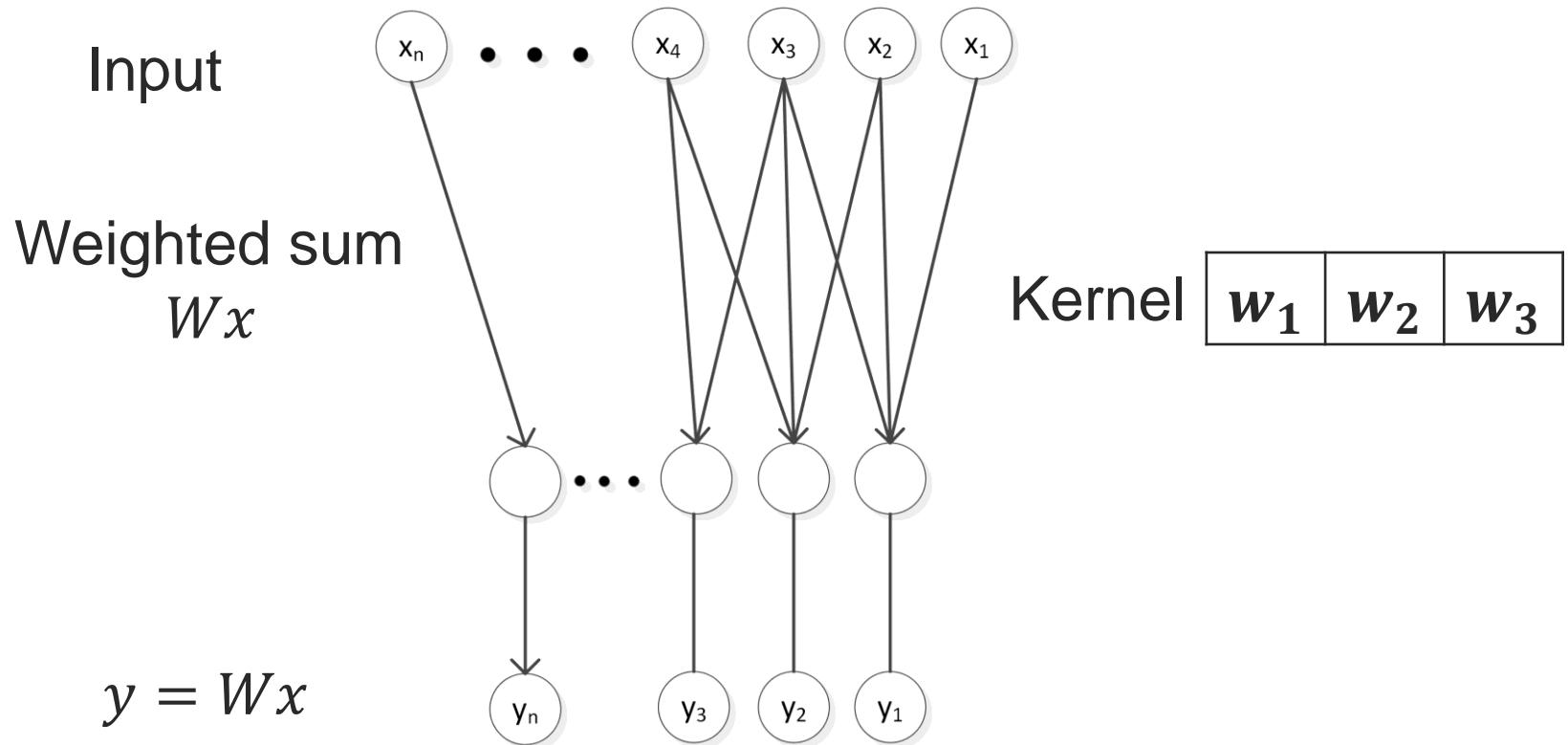
Convolution as MLP (1)

- Remove activation



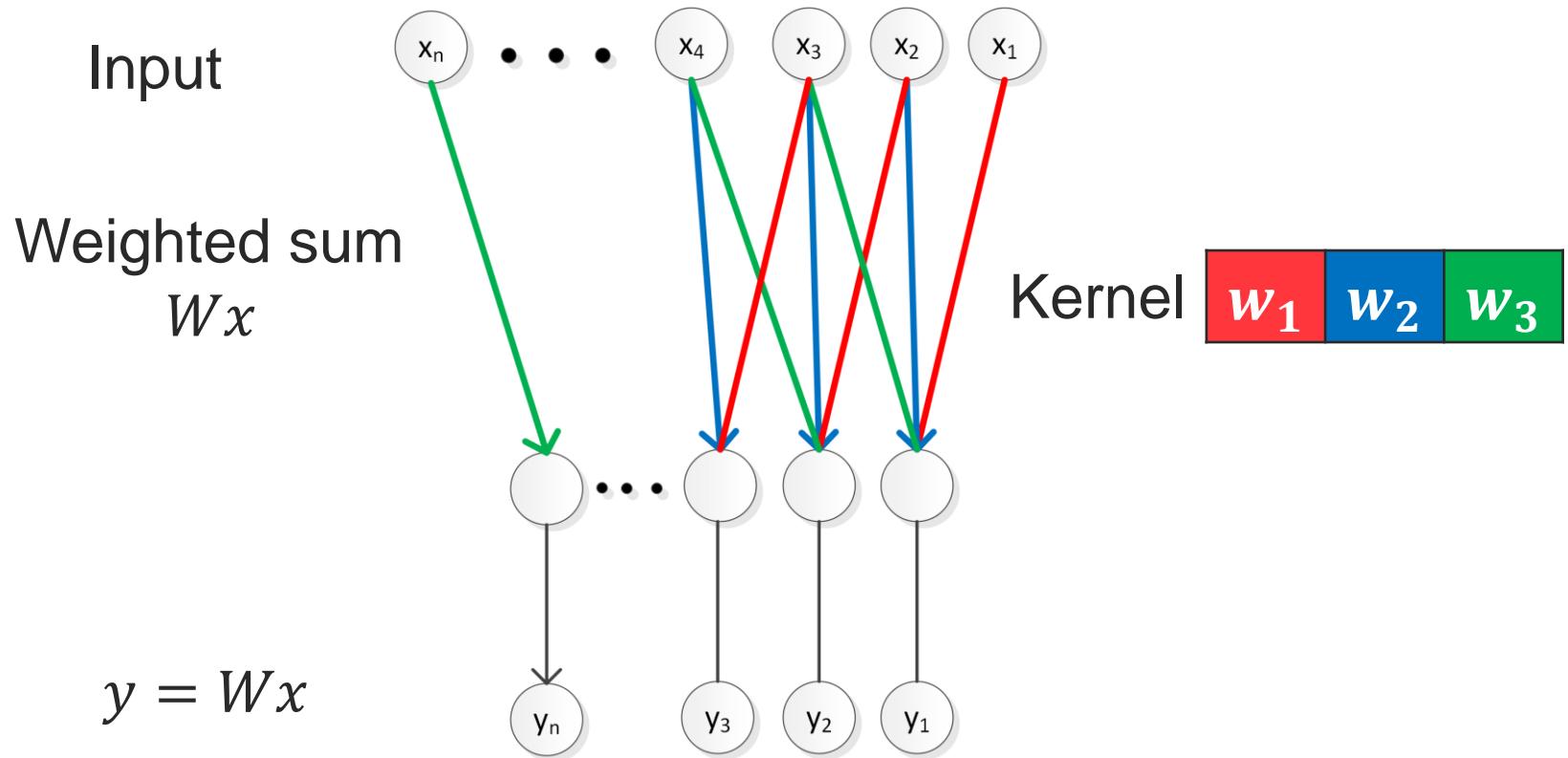
Convolution as MLP (2)

- Remove redundant links making the matrix W sparse
(optionally remove the bias term)



Convolution as MLP (3)

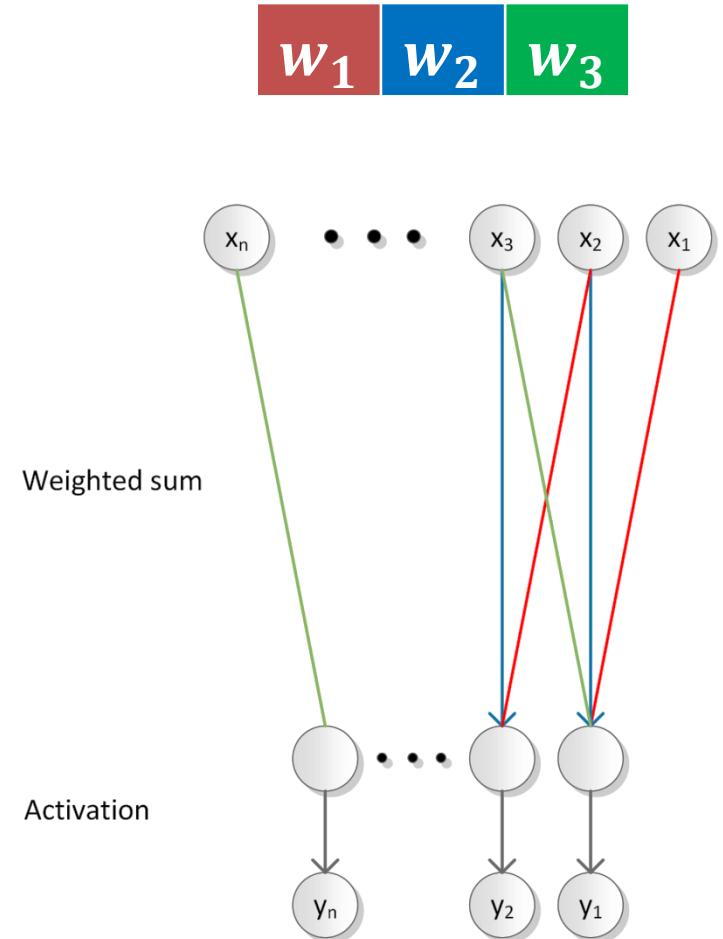
- We can also share the weights in matrix W not to do redundant computation



How do we do convolution in MLP recap

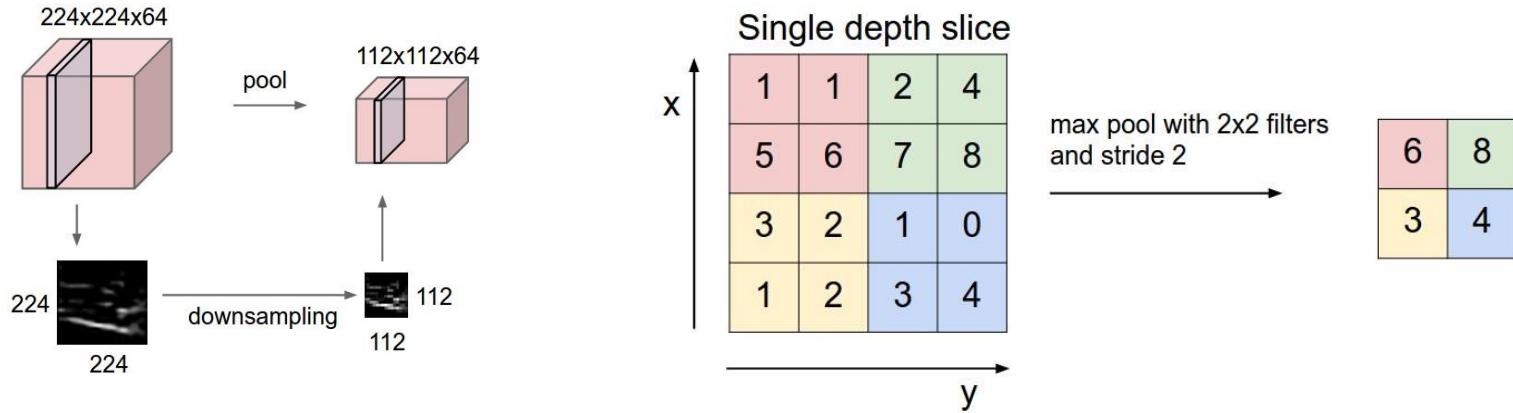
- Not a fully connected layer anymore
- Shared weights
 - Same colour indicates same (shared) weight

$$W = \begin{pmatrix} w_1 & w_2 & w_3 & & 0 & 0 & 0 \\ 0 & w_1 & w_2 & \dots & 0 & 0 & 0 \\ 0 & 0 & w_1 & & 0 & 0 & 0 \\ \vdots & & \ddots & & \vdots & & \\ 0 & 0 & 0 & & w_3 & 0 & 0 \\ 0 & 0 & 0 & \dots & w_2 & w_3 & 0 \\ 0 & 0 & 0 & & w_1 & w_2 & w_3 \end{pmatrix}$$



Pooling layer

- Used for sub-sampling



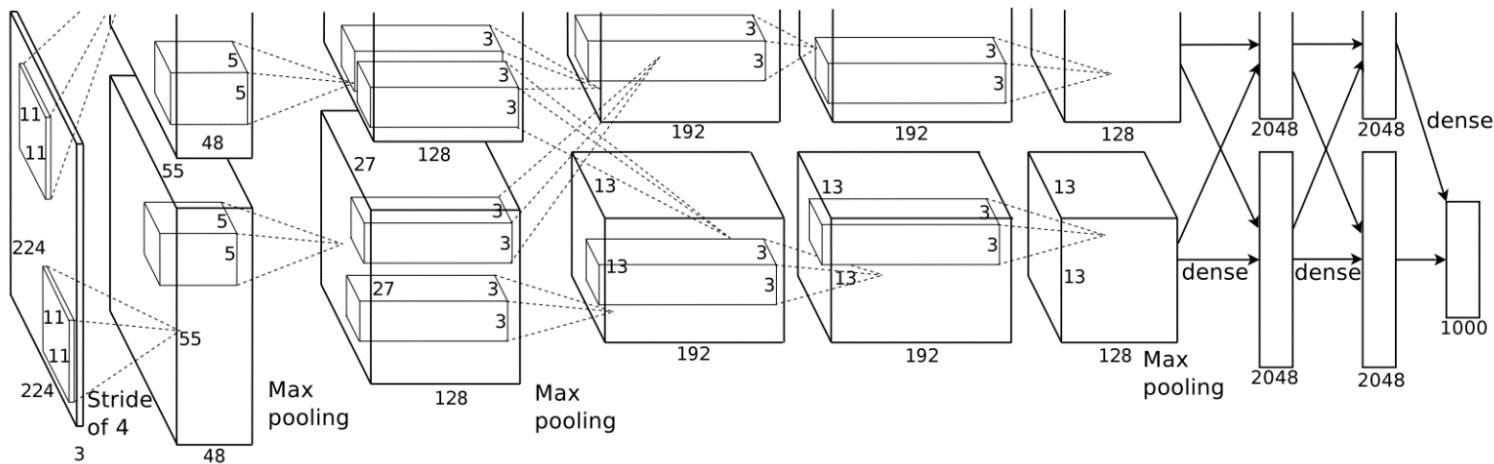
Pick the maximum value from input using a smooth and differentiable approximation

$$y = \frac{\sum_{i=1}^n x_i e^{\alpha x_i}}{\sum_{i=1}^n e^{\alpha x_i}}$$



Example: AlexNet Model

- Used for object classification task
 - 1000 way classification task – pick one



Unimodal representations: Acoustic Modality

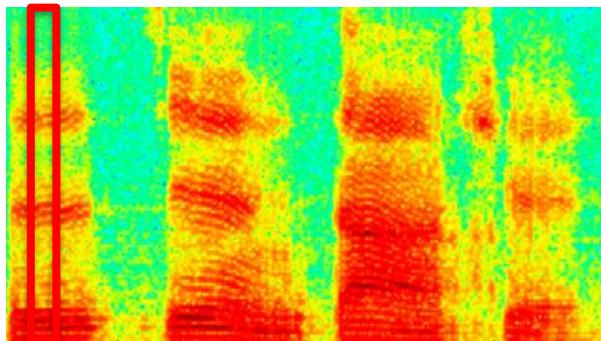
Unimodal Classification – Acoustic Modality

Digitalized acoustic signal



- Sampling rates: 8~96kHz
- Bit depth: 8, 16 or 24 bits
- Time window size: 20ms
 - Offset: 10ms

Input observation x_i
0.21
0.14
0.56
0.45
0.9
0.98
0.75
0.34
0.24
0.11
0.02



Spectrogram

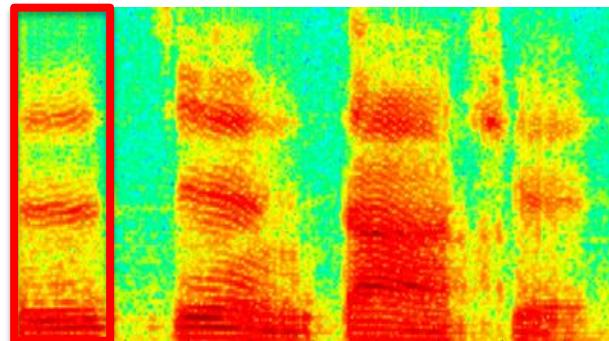


Unimodal Classification – Acoustic Modality

Digitalized acoustic signal

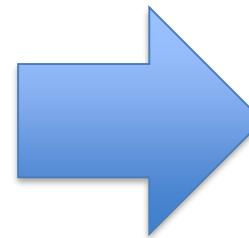


- Sampling rates: 8~96kHz
- Bit depth: 8, 16 or 24 bits
- Time window size: 20ms
 - Offset: 10ms



Spectrogram

Input observation x_i
0.21
0.14
0.56
0.45
0.9
0.98
0.75
0.34
0.24
0.11
0.02
0.24
0.26
0.58
0.9
0.99
0.79
0.45
0.34
0.24
⋮



Emotion ?

Spoken word ?

Voice quality ?



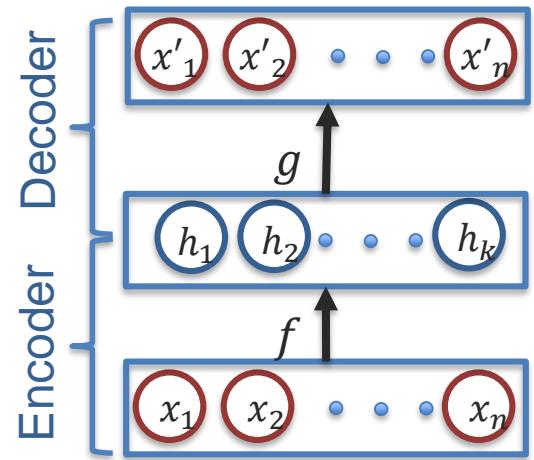
Audio representation for speech recognition

- Speech recognition systems historically much more complex than vision systems – language models, vocabularies etc.
- Large breakthrough of using representation learning instead of hand-crafted features
 - [Hinton et al., Deep Neural Networks for Acoustic Modeling in Speech Recognition: The Shared Views of Four Research Groups, 2012]
- A huge boost in performance (up to 30% on some datasets)



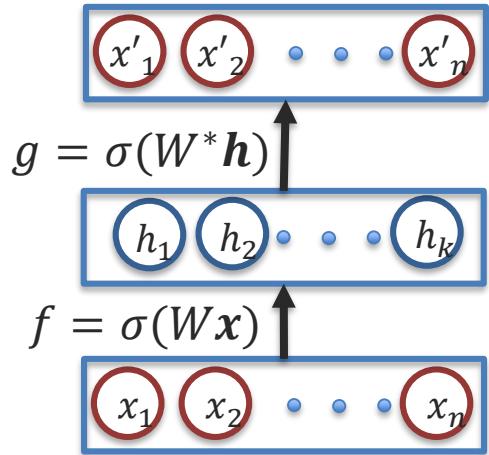
Autoencoders

- What does auto mean?
 - Greek for self – self encoding
- Feed forward network intended to reproduce the input
- Two parts encoder/decoder
 - $x' = f(g(x))$ – score function
 - g - encoder
 - f - decoder



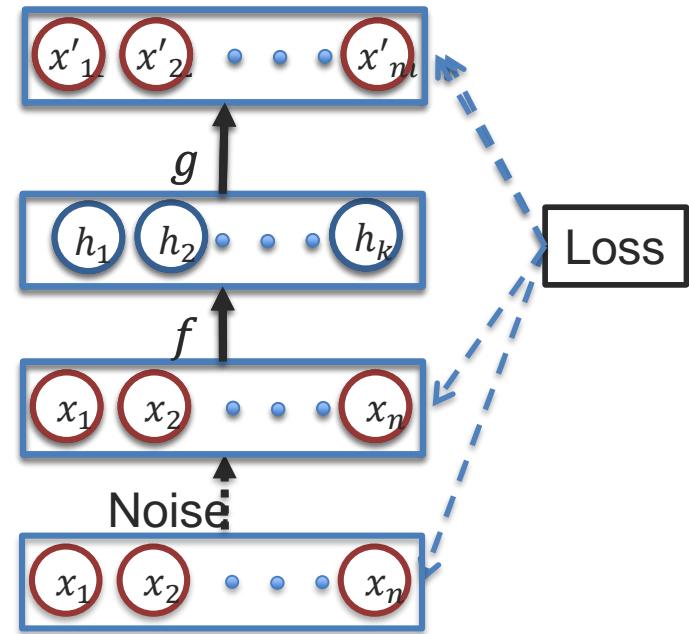
Autoencoders

- Mostly follows Neural Network structure
 - A matrix multiplication followed by a sigmoid
- Activation will depend on type of x
 - Sigmoid for binary
 - Linear for real valued
- Often we use tied weights to force the sharing of weights in encoder/decoder
 - $W^* = W^T$
- word2vec is actually a bit similar to autoencoder (except for the auto part)



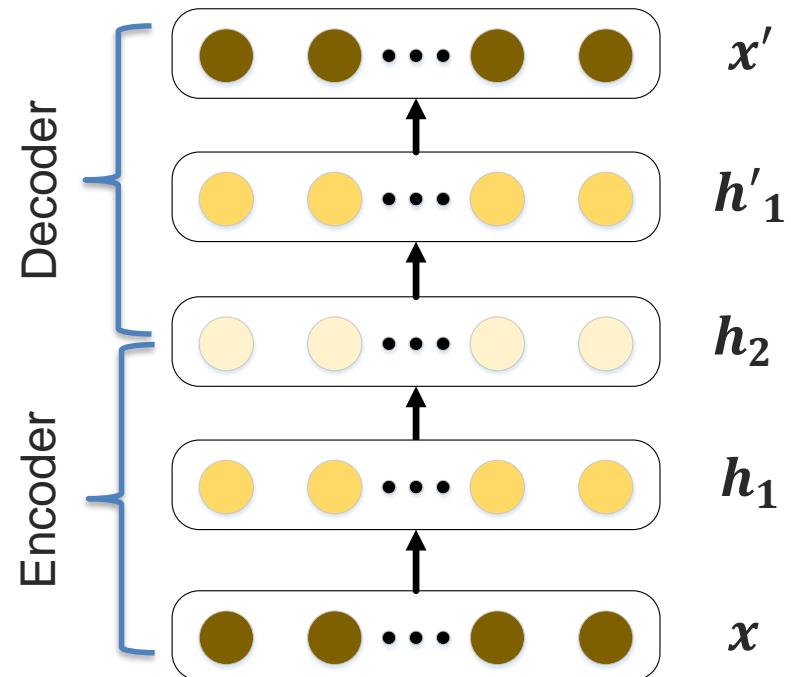
Denoising autoencoder

- Simple idea
 - Add noise to input x but learn to reconstruct original
- Leads to a more robust representation and prevents copying
- Learns what the relationship is to represent a certain x
- Different noise added during each epoch



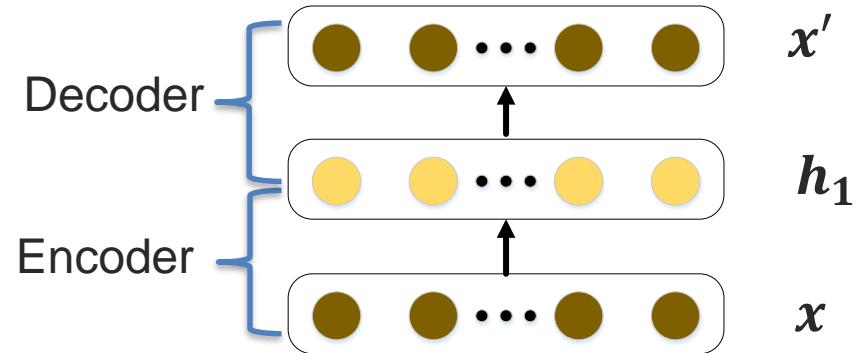
Stacked autoencoders

- Can stack autoencoders as well
- Each encoding unit has a corresponding decoder
- Inference as before is feed forward structure, but now with more hidden layers



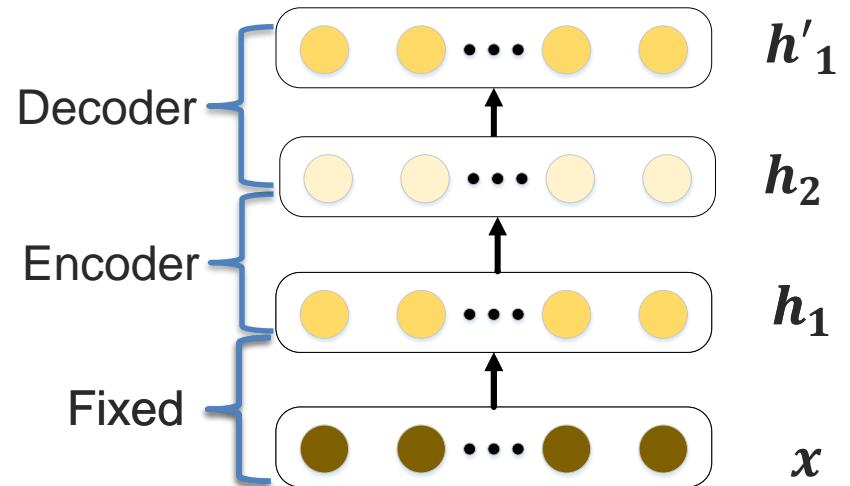
Stacked autoencoders

- Greedy layer-wise training
- Start with training first layer
 - Learn to encode x to h_1 and to decode x from h_1
 - Use backpropagation



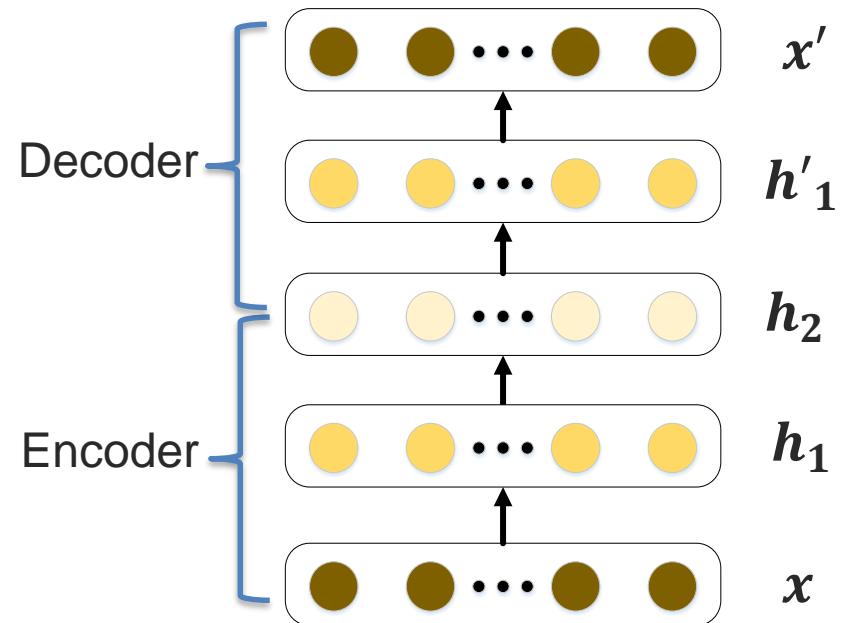
Stacked autoencoders

- Map from all x 's to h_1 's
 - Discard decoder for now
- Train the second layer
 - Learn to encode h_1 to h_2 and to decode h_2 from h_1
 - Repeat for as many layers



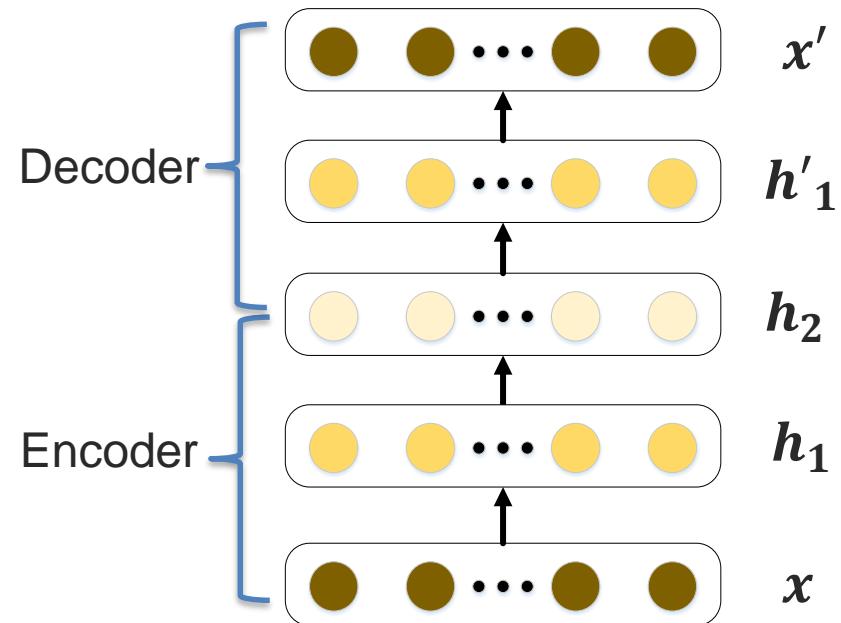
Stacked autoencoders

- Reconstruct using previously learned decoders mappings
- Fine-tune the full network end-to-end



Stacked denoising autoencoders

- Can extend this to a denoising model
- Add noise when training each of the layers
 - Often with increasing amount of noise per layer
 - 0.1 for first, 0.2 for second, 0.3 for third

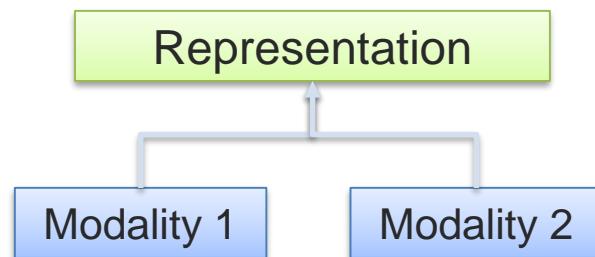


Multimodal Representations

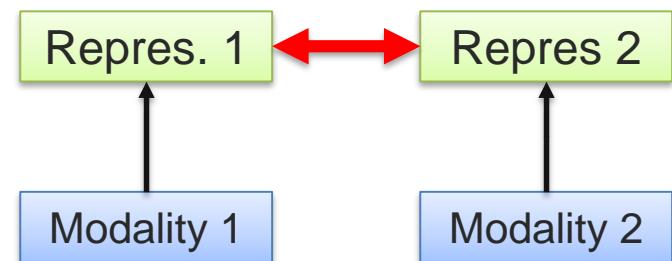
Core Challenge: Representation

Definition: Learning how to represent and summarize multimodal data in a way that exploits the complementarity and redundancy.

A Joint representations:

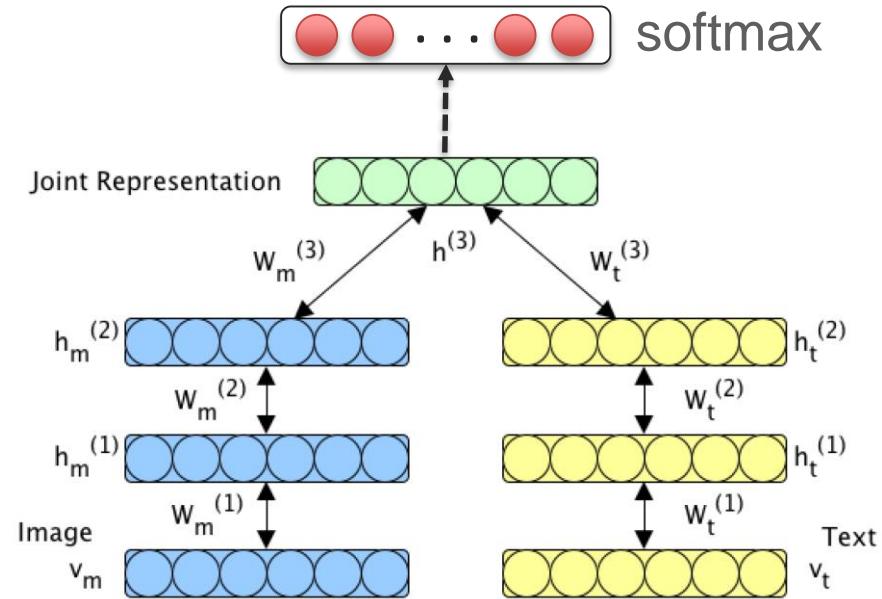


B Coordinated representations:



Deep Multimodal Boltzmann machines

- Generative model
- Individual modalities trained like a DBN
- Multimodal representation trained using Variational approaches
- Used for image tagging and cross-media retrieval
- Reconstruction of one modality from another is a bit more “natural” than in autoencoder representation
- Can actually sample text and images



[Srivastava and Salakhutdinov, Multimodal Learning with Deep Boltzmann Machines, 2012, 2014]



Deep Multimodal Boltzmann machines

Image	Given Tags	Generated Tags	Input Text	2 nearest neighbours to generated image features
	pentax, k10d, kangarooisland, southaustralia, sa, australia, australiassealion, 300mm	beach, sea, surf, strand, shore, wave, seascape, sand, ocean, waves	nature, hill scenery, green clouds	 
	<no text>	night, lights, christmas, nightshot, nacht, nuit, notte, longexposure, noche, nocturna	flower, nature, green, flowers, petal, petals, bud	 
	aheram, 0505 sarahc, moo	portrait, bw, blackandwhite, woman, people, faces, girl, blackwhite, person, man	blue, red, art, artwork, painted, paint, artistic surreal, gallery bleu	 
	unseulpixel, naturey crap	fall, autumn, trees, leaves, foliage, forest, woods, branches, path	bw, blackandwhite, noiretblanc, biancoenero blancoynegro	 

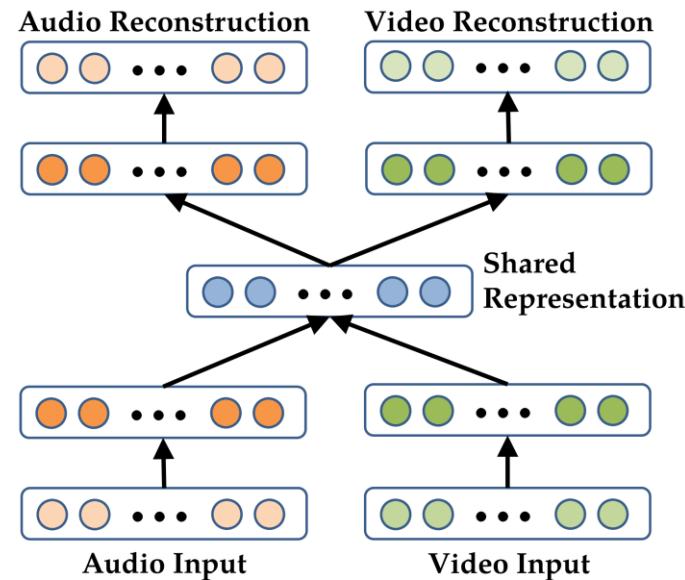
Model	MAP	Prec@50
Random	0.124	0.124
SVM (Huiskes et al., 2010)	0.475	0.758
LDA (Huiskes et al., 2010)	0.492	0.754
DBM	0.526 ± 0.007	0.791 ± 0.008
DBM (using unlabelled data)	0.585 ± 0.004	0.836 ± 0.004

Srivastava and Salakhutdinov, “Multimodal Learning with Deep Boltzmann Machines”, NIPS 2012



Deep Multimodal autoencoders

- A deep representation learning approach
- A bimodal auto-encoder
 - Used for Audio-visual speech recognition



[Ngiam et al., Multimodal Deep Learning, 2011]

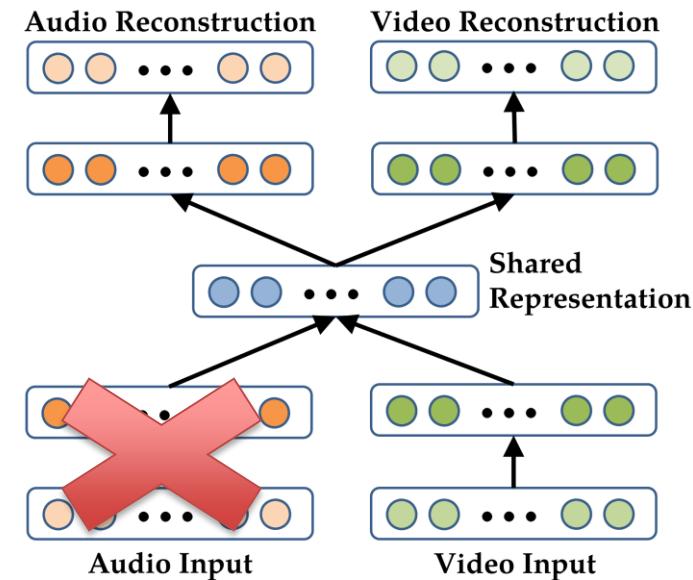


Language Technologies Institute

Carnegie Mellon University

Deep Multimodal autoencoders - training

- Individual modalities can be pretrained
 - RBMs
 - Denoising Autoencoders
- To train the model to reconstruct the other modality
 - Use both
 - Remove audio



[Ngiam et al., Multimodal Deep Learning, 2011]

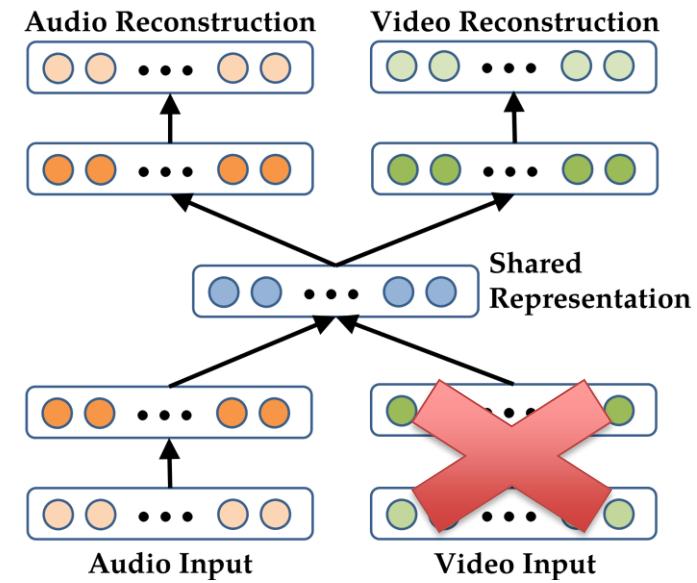


Language Technologies Institute

Carnegie Mellon University

Deep Multimodal autoencoders - training

- Individual modalities can be pretrained
 - RBMs
 - Denoising Autoencoders
- To train the model to reconstruct the other modality
 - Use both
 - Remove audio
 - Remove video



[Ngiam et al., Multimodal Deep Learning, 2011]

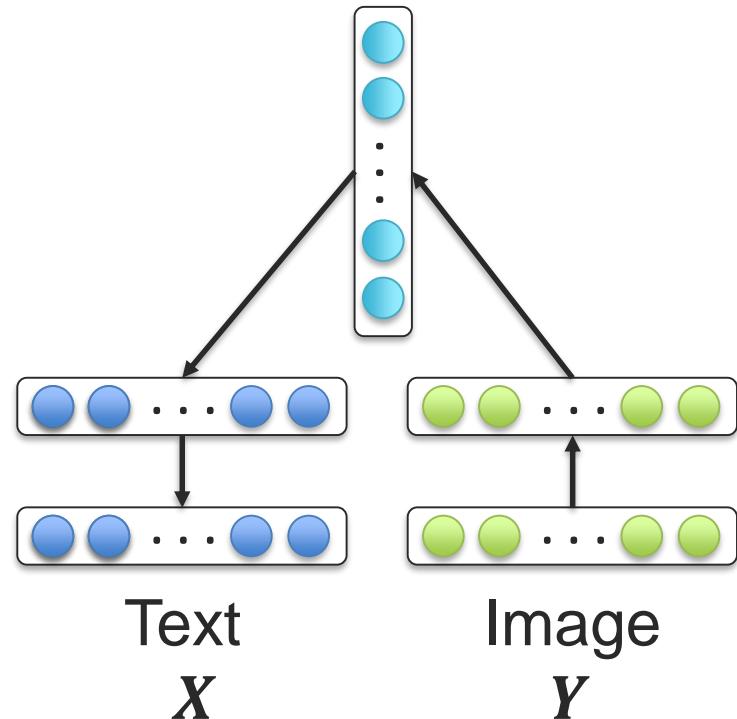


Language Technologies Institute

Carnegie Mellon University

Multimodal Encoder-Decoder

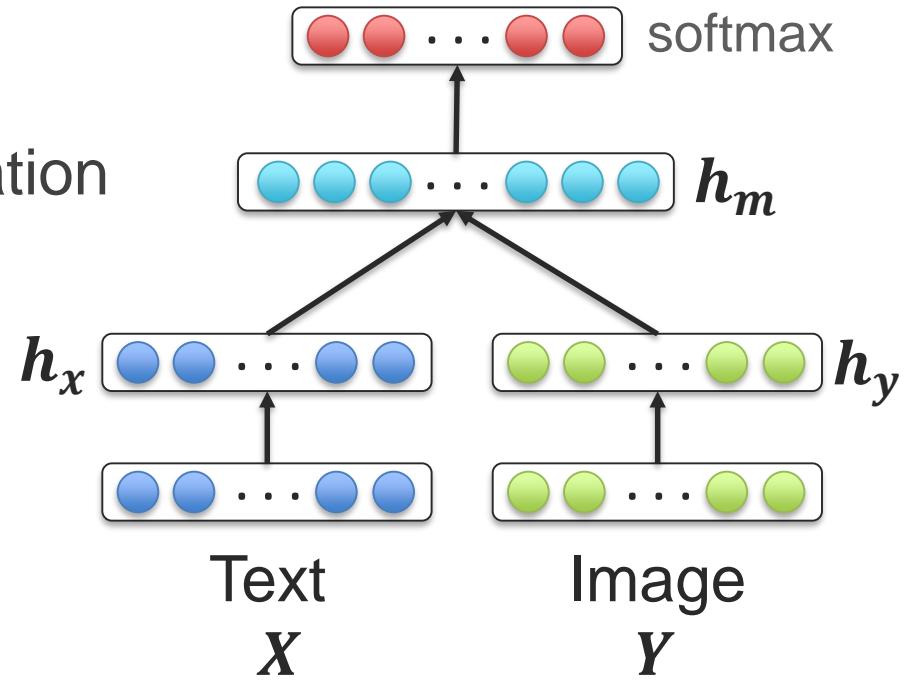
- Visual modality often encoded using CNN
- Language modality will be decoded using LSTM
 - A simple multilayer perceptron will be used to translate from visual (CNN) to language (LSTM)



Multimodal Joint Representation

- For supervised learning tasks
- Joining the unimodal representations:
 - Simple concatenation
 - Element-wise multiplication or summation
 - Multilayer perceptron
- How to explicitly model both unimodal and bimodal interactions?

e.g. Sentiment



Multimodal Sentiment Analysis

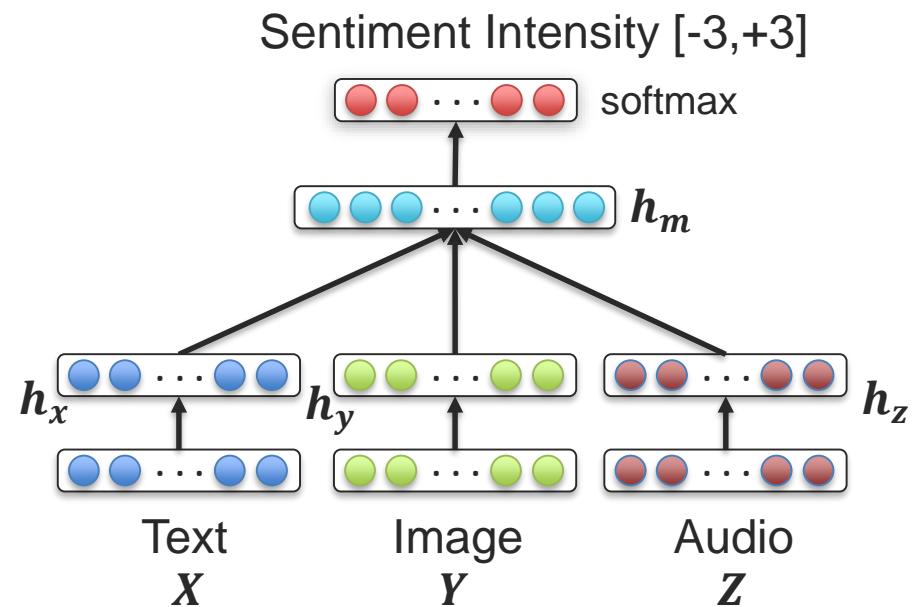
MOSI dataset (Zadeh et al, 2016)



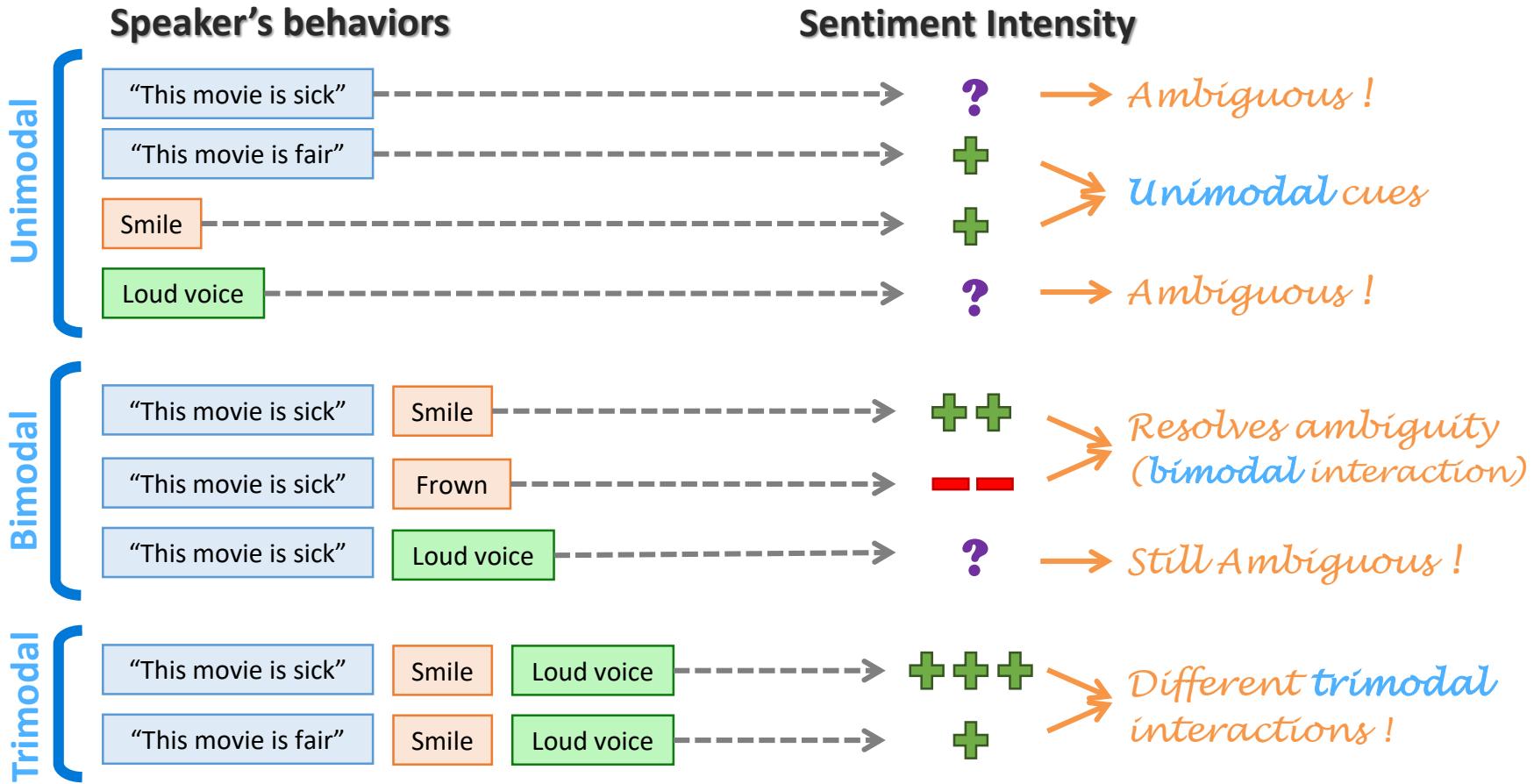
- 2199 subjective video segments
- Sentiment intensity annotations
- 3 modalities: text, video, audio

Multimodal joint representation:

$$h_m = f(W \cdot [h_x, h_y, h_z])$$



Unimodal, Bimodal and Trimodal Interactions



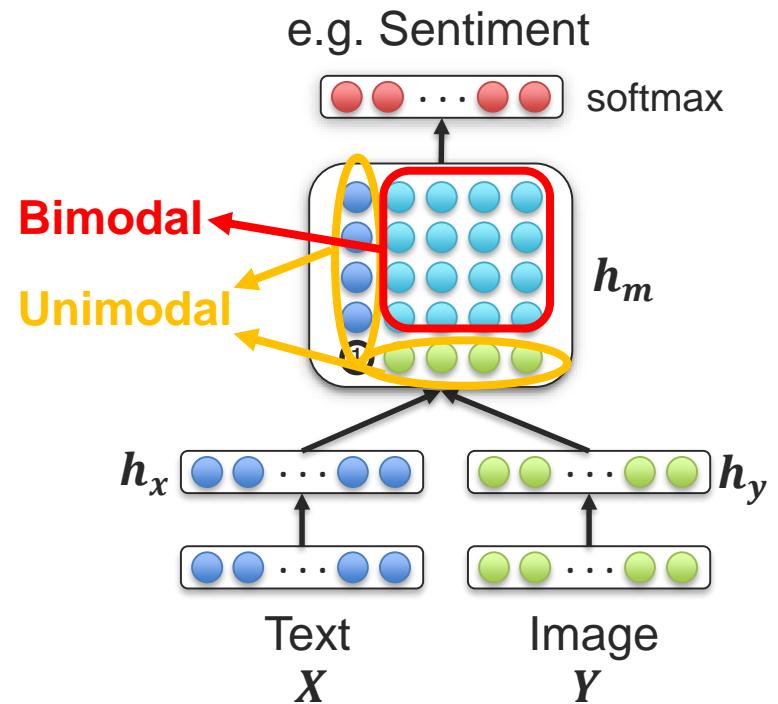
Multimodal Tensor Fusion Network (TFN)

Models both unimodal and bimodal interactions:

$$h_m = [h_x] \otimes [h_y] = \begin{bmatrix} h_x \\ 1 \end{bmatrix} \otimes \begin{bmatrix} h_y \\ 1 \end{bmatrix} = \begin{bmatrix} h_x & h_x \otimes h_y \\ 1 & h_y \end{bmatrix}$$

Important !

[Zadeh, Jones and Morency, EMNLP 2017]



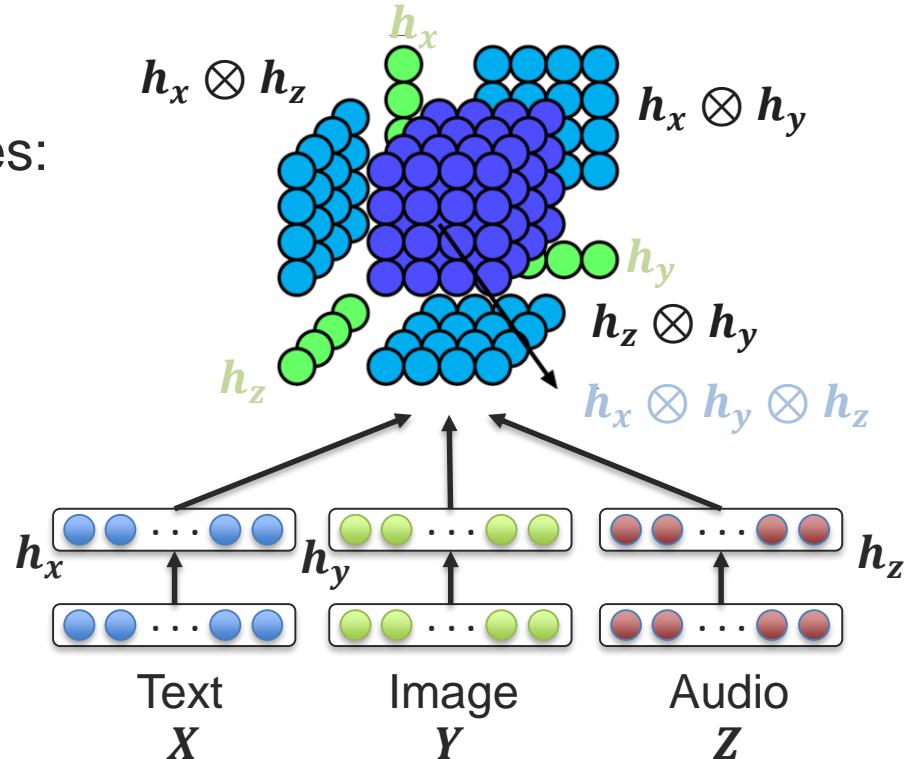
Multimodal Tensor Fusion Network (TFN)

Can be extended to three modalities:

$$\mathbf{h}_m = \begin{bmatrix} \mathbf{h}_x \\ 1 \end{bmatrix} \otimes \begin{bmatrix} \mathbf{h}_y \\ 1 \end{bmatrix} \otimes \begin{bmatrix} \mathbf{h}_z \\ 1 \end{bmatrix}$$

Explicitly models unimodal, bimodal and trimodal interactions !

[Zadeh, Jones and Morency, EMNLP 2017]



Experimental Results – MOSI Dataset

Multimodal Baseline	Binary		5-class		Regression	
	Acc(%)	F1	Acc(%)	MAE	r	
Random	50.2	48.7	23.9	1.88	-	
C-MKL	73.1	75.2	35.3	-	-	
SAL-CNN	73.0	-	-	-	-	
SVM-MD	71.6	72.3	32.0	1.10	0.53	
RF	71.4	72.1	31.9	1.11	0.51	
TFN	77.1	77.9	42.0	0.87	0.70	
Human	85.7	87.5	53.9	0.71	0.82	
Δ^{SOTA}	↑ 4.0	↑ 2.7	↑ 6.7	↓ 0.23	↑ 0.17	

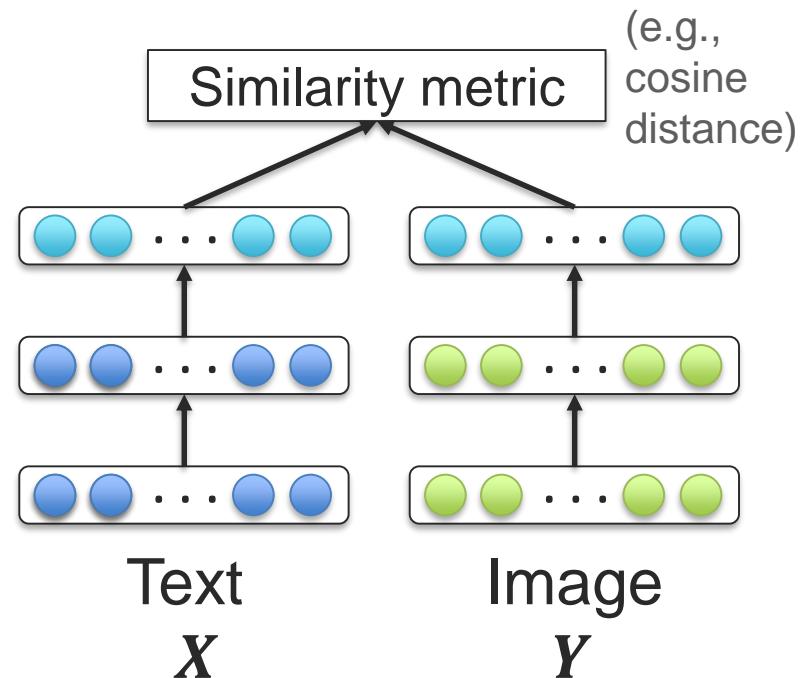
Improvement over State-Of-The-Art

Baseline	Binary		5-class		Regression	
	Acc(%)	F1	Acc(%)	MAE	r	
TFN _{language}	74.8	75.6	38.5	0.99	0.61	
TFN _{visual}	66.8	70.4	30.4	1.13	0.48	
TFN _{acoustic}	65.1	67.3	27.5	1.23	0.36	
TFN _{bimodal}	75.2	76.0	39.6	0.92	0.65	
TFN _{trimodal}	74.5	75.0	38.9	0.93	0.65	
TFN _{not trimodal}	75.3	76.2	39.7	0.919	0.66	
TFN	77.1	77.9	42.0	0.87	0.70	
TFN _{early}	75.2	76.2	39.0	0.96	0.63	

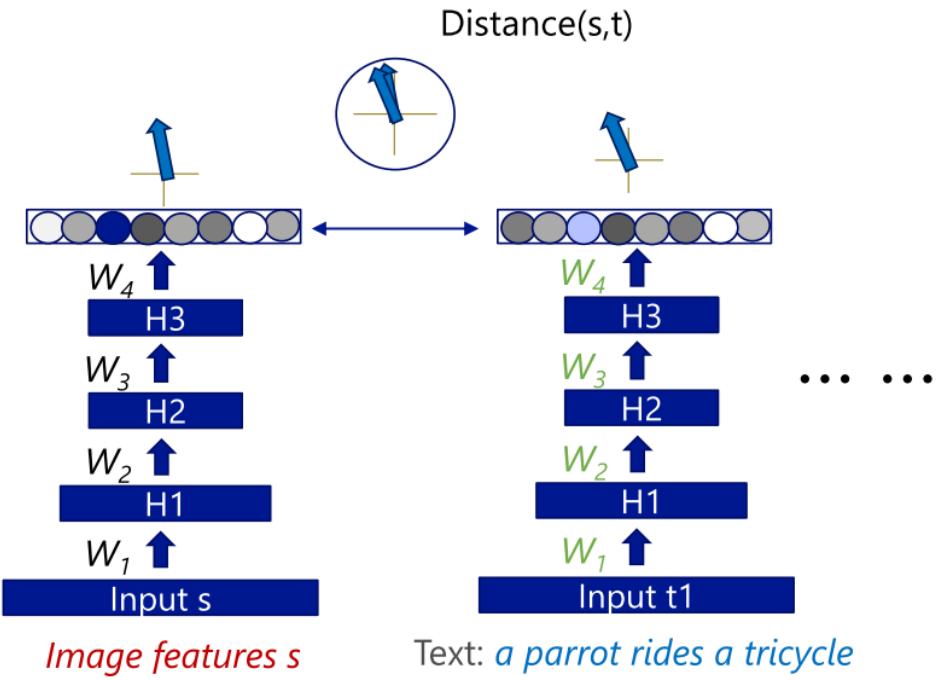
Coordinated Multimodal Representations

Coordinated Multimodal Representations

Learn (unsupervised) two or more coordinated representations from multiple modalities. A loss function is defined to bring closer these multiple representations.



Coordinated Multimodal Embeddings



[Huang et al., Learning Deep Structured Semantic Models for Web Search using Clickthrough Data, 2013]



Language Technologies Institute

Carnegie Mellon University

Multimodal Vector Space Arithmetic



- blue + red =

- blue + yellow =

- yellow + red =

- white + red =

Nearest images



[Kiros et al., Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models, 2014]



Language Technologies Institute

Carnegie Mellon University

Multimodal Vector Space Arithmetic



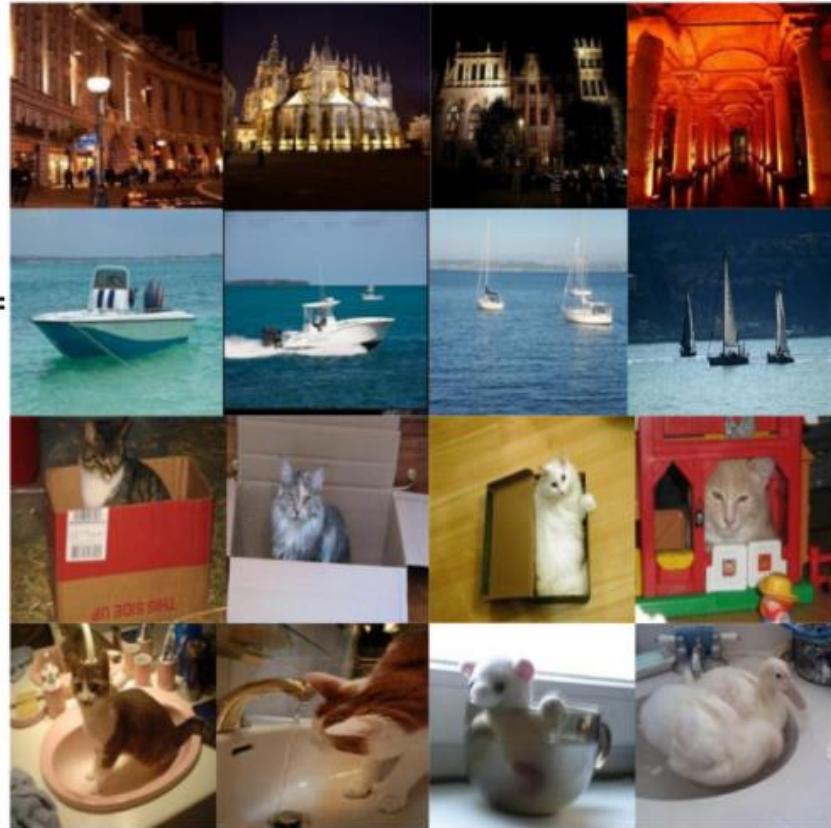
- day + night =

- flying + sailing =

- bowl + box =

- box + bowl =

Nearest images



[Kiros et al., Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models, 2014]



Language Technologies Institute

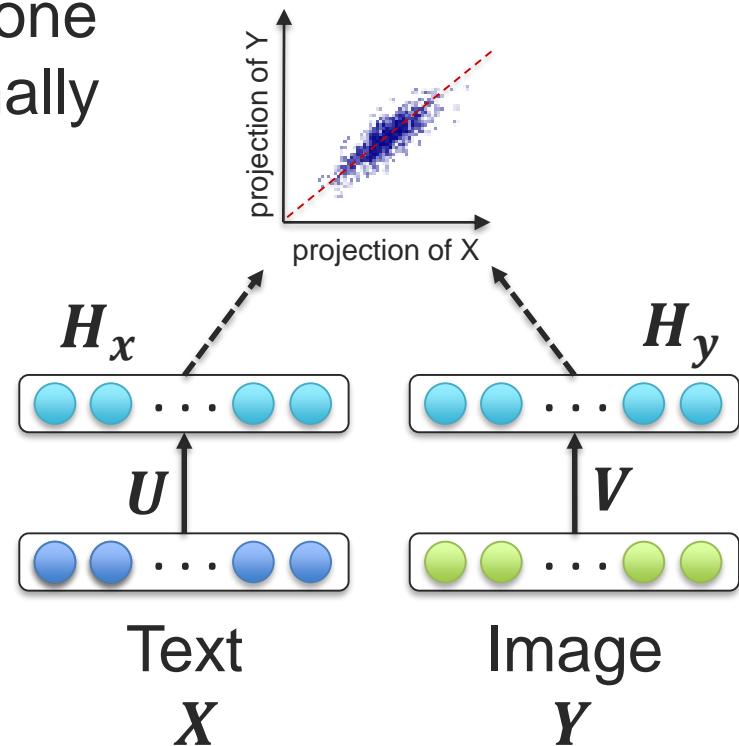
Carnegie Mellon University

Canonical Correlation Analysis

“canonical”: reduced to the simplest or clearest schema possible

- 1 Learn two linear projections, one for each view, that are maximally correlated:

$$\begin{aligned}(\boldsymbol{u}^*, \boldsymbol{v}^*) &= \operatorname{argmax}_{\boldsymbol{u}, \boldsymbol{v}} \operatorname{corr}(\boldsymbol{H}_x, \boldsymbol{H}_y) \\ &= \operatorname{argmax}_{\boldsymbol{u}, \boldsymbol{v}} \operatorname{corr}(\boldsymbol{u}^T \boldsymbol{X}, \boldsymbol{v}^T \boldsymbol{Y})\end{aligned}$$



Correlated Projection

- 1 Learn two linear projections, one for each view, that are maximally correlated:

$$(\mathbf{u}^*, \mathbf{v}^*) = \operatorname{argmax}_{\mathbf{u}, \mathbf{v}} \operatorname{corr}(\mathbf{u}^T \mathbf{X}, \mathbf{v}^T \mathbf{Y})$$



Two views X, Y where same instances have the same color



Canonical Correlation Analysis

We want to learn multiple projection pairs $(\mathbf{u}_{(i)}X, \mathbf{v}_{(i)}Y)$:

$$(\mathbf{u}_{(i)}^*, \mathbf{v}_{(i)}^*) = \underset{\mathbf{u}_{(i)}, \mathbf{v}_{(i)}}{\operatorname{argmax}} \operatorname{corr}(\mathbf{u}_{(i)}^T X, \mathbf{v}_{(i)}^T Y) \approx \mathbf{u}_{(i)}^T \Sigma_{XY} \mathbf{v}_{(i)}$$

- 2 We want these multiple projection pairs to be orthogonal (“canonical”) to each other:

$$\mathbf{u}_{(i)}^T \Sigma_{XY} \mathbf{v}_{(j)} = \mathbf{u}_{(j)}^T \Sigma_{XY} \mathbf{v}_{(i)} = \mathbf{0} \quad \text{for } i \neq j$$

$$\mathbf{U} \Sigma_{XY} \mathbf{V} = \operatorname{tr}(\mathbf{U} \Sigma_{XY} \mathbf{V}) \quad \text{where } \mathbf{U} = [\mathbf{u}_{(1)}, \mathbf{u}_{(2)}, \dots, \mathbf{u}_{(k)}] \\ \text{and } \mathbf{V} = [\mathbf{v}_{(1)}, \mathbf{v}_{(2)}, \dots, \mathbf{v}_{(k)}]$$



Canonical Correlation Analysis

- ③ Since this objective function is invariant to scaling, we can constraint the projections to have unit variance:

$$\mathbf{U}^T \boldsymbol{\Sigma}_{XX} \mathbf{U} = I \quad \mathbf{V}^T \boldsymbol{\Sigma}_{YY} \mathbf{V} = I$$

Canonical Correlation Analysis:

maximize: $\text{tr}(\mathbf{U}^T \boldsymbol{\Sigma}_{XY} \mathbf{V})$

subject to: $\mathbf{U}^T \boldsymbol{\Sigma}_{YY} \mathbf{U} = \mathbf{V}^T \boldsymbol{\Sigma}_{YY} \mathbf{V} = I$



Canonical Correlation Analysis

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

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

$$\Sigma = \begin{bmatrix} \Sigma_{XX} & \Sigma_{YX} \\ \hline \Sigma_{XY} & \Sigma_{YY} \end{bmatrix} \xrightarrow{\mathbf{U}, \mathbf{V}} \begin{bmatrix} 1 & 0 & 0 & \lambda_1 & 0 & 0 \\ 0 & 1 & 0 & 0 & \lambda_2 & 0 \\ 0 & 0 & 1 & 0 & 0 & \lambda_3 \\ \hline \lambda_1 & 0 & 0 & 1 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 & 1 & 0 \\ 0 & 0 & \lambda_3 & 0 & 0 & 1 \end{bmatrix}$$

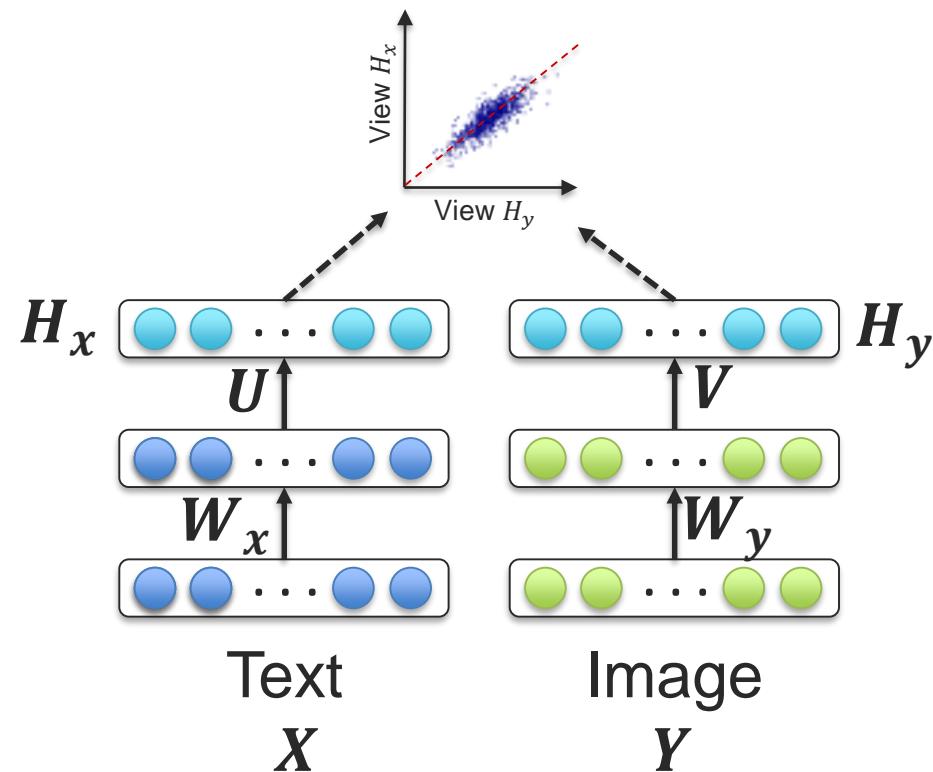


Deep Canonical Correlation Analysis

Same objective function as CCA:

$$\underset{V, U, W_x, W_y}{\operatorname{argmax}} \operatorname{corr}(H_x, H_y)$$

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

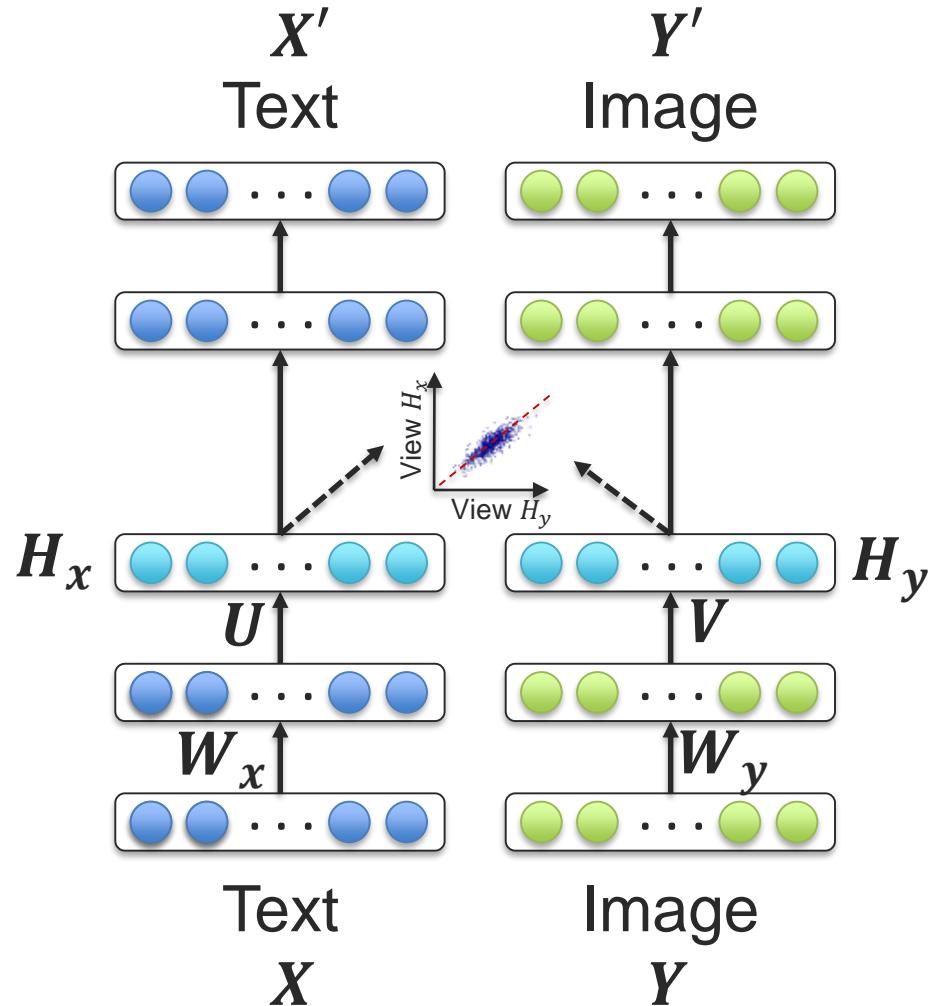


Andrew et al., ICML 2013

Deep Canonically Correlated Autoencoders (DCCAE)

Jointly optimize for DCCA and autoencoders loss functions

- A trade-off between multi-view correlation and reconstruction error from individual views

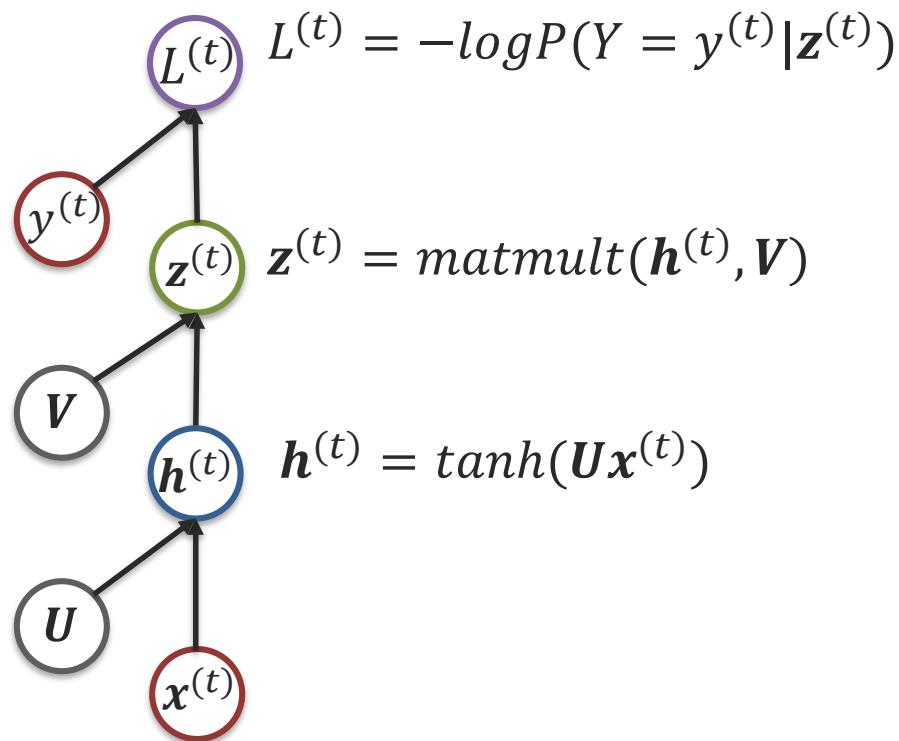


Wang et al., ICML 2015

Basic Concepts: Recurrent Neural Networks

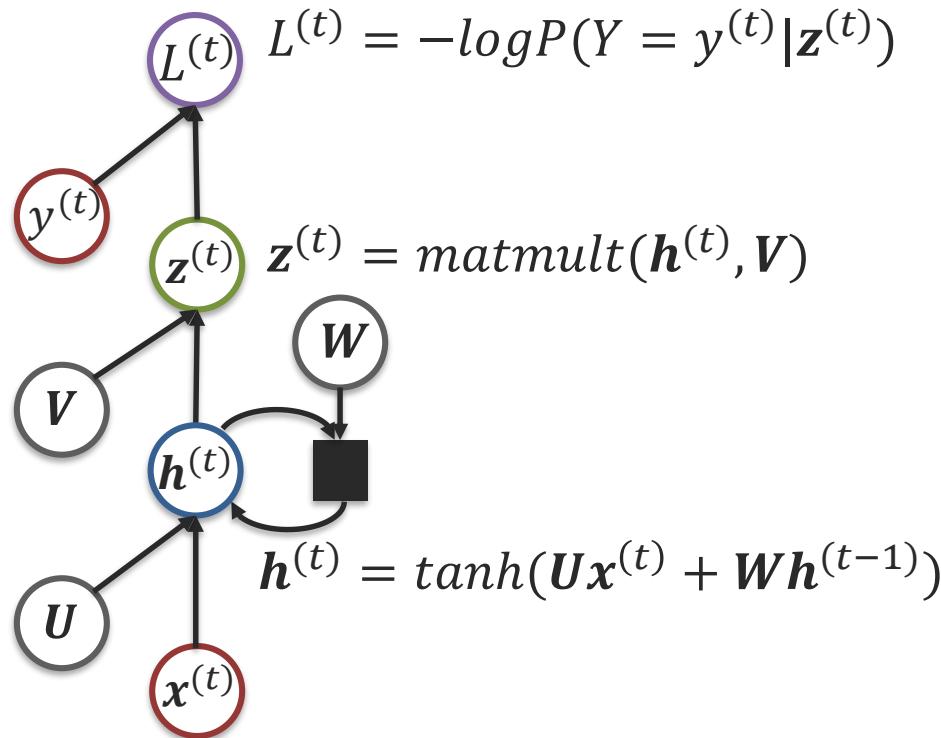


Feedforward Neural Network



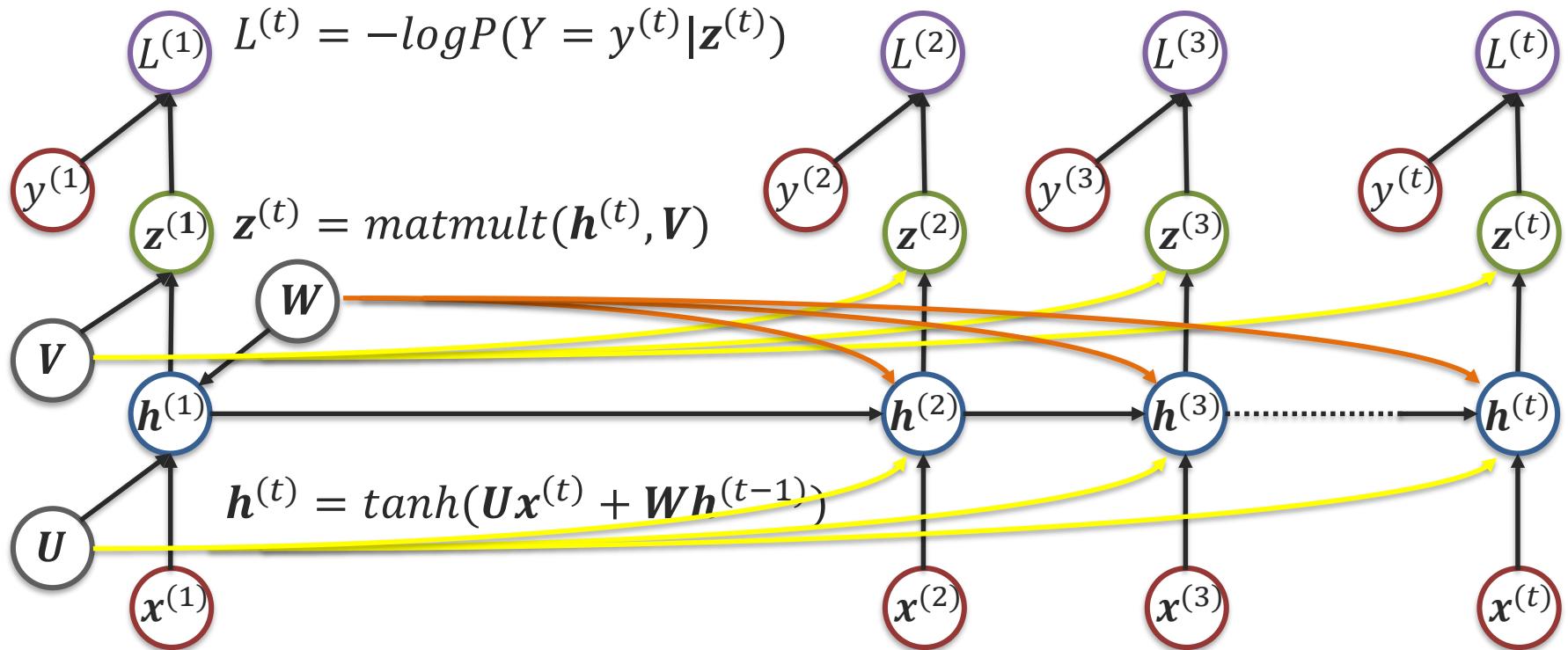
Recurrent Neural Networks

$$L = \sum_t L^{(t)}$$



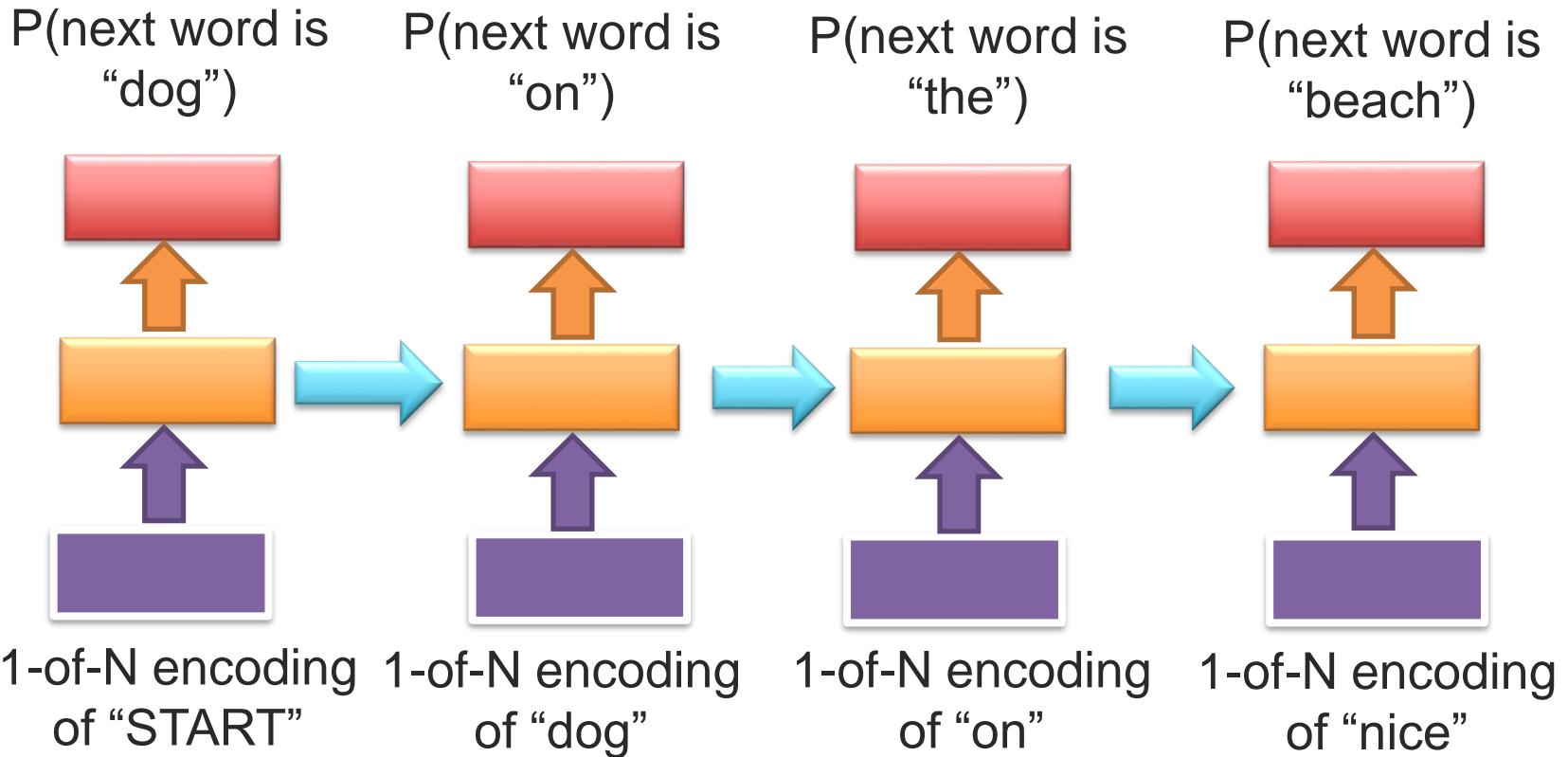
Recurrent Neural Networks - Unrolling

$$L = \sum_t L^{(t)}$$



Same model parameters are used for all time parts.

Recurrent Neural Networks – Language models

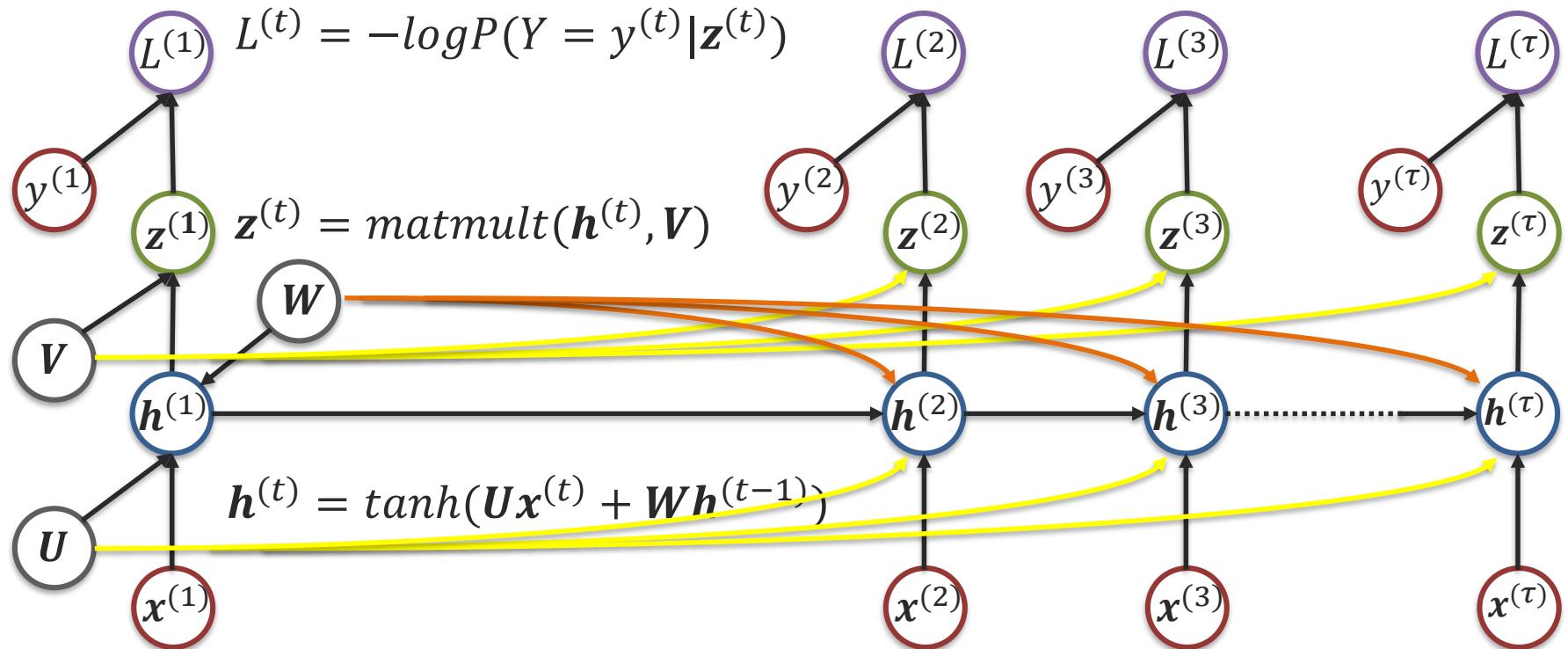


➤ Model long-term information



Recurrent Neural Networks

$$L = \sum_t L^{(t)}$$



Backpropagation Through Time

$$L = \sum_t L^{(t)} = -\sum_t \log P(Y = y^{(t)} | \mathbf{z}^{(t)})$$

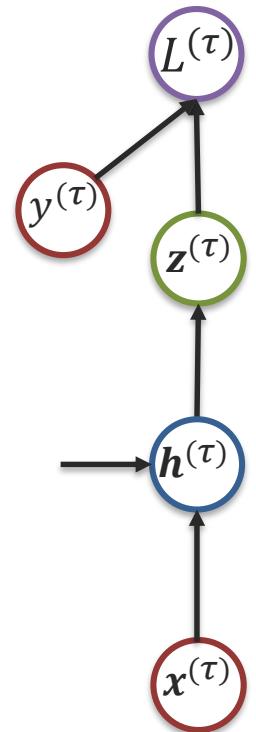
$L^{(\tau)}$ or $L^{(t)}$ $\frac{\partial L}{\partial L^{(t)}} = 1$

Gradient = “backprop” gradient
x “local” Jacobian

$\mathbf{z}^{(\tau)}$ or $\mathbf{z}^{(t)}$ $(\nabla_{\mathbf{z}^{(t)}} L)_i = \frac{\partial L}{\partial z_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial z_i^{(t)}} = \text{sigmoid}(z_i^t) - \mathbf{1}_{i,y^{(t)}}$

$\mathbf{h}^{(\tau)}$ $\nabla_{\mathbf{h}^{(\tau)}} L = \nabla_{\mathbf{z}^{(\tau)}} L \frac{\partial z^{(\tau)}}{\partial \mathbf{h}^{(\tau)}} = \nabla_{\mathbf{z}^{(\tau)}} L \mathbf{V}$

$\mathbf{h}^{(t)} \rightarrow \mathbf{h}^{(t+1)}$ $\nabla_{\mathbf{h}^{(t)}} L = \nabla_{\mathbf{z}^{(t)}} L \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} + \nabla_{\mathbf{z}^{(t+1)}} L \frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{h}^{(t)}}$



Backpropagation Through Time

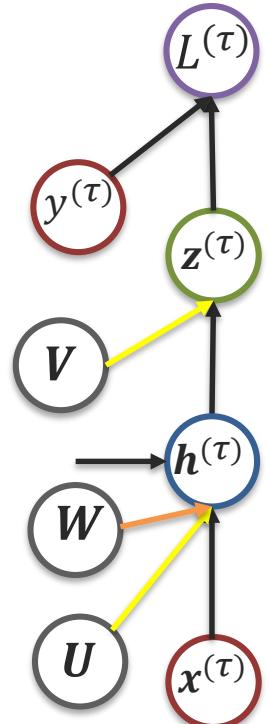
$$L = \sum_t L^{(t)} = -\sum_t \log P(Y = y^{(t)} | \mathbf{z}^{(t)})$$

Gradient = “backprop” gradient
x “local” Jacobian

(V) $\nabla_V L = \sum_t (\nabla_{\mathbf{z}^{(t)}} L) \frac{\partial \mathbf{z}^{(t)}}{\partial V}$

(W) $\nabla_W L = \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \frac{\partial \mathbf{h}^{(t)}}{\partial W}$

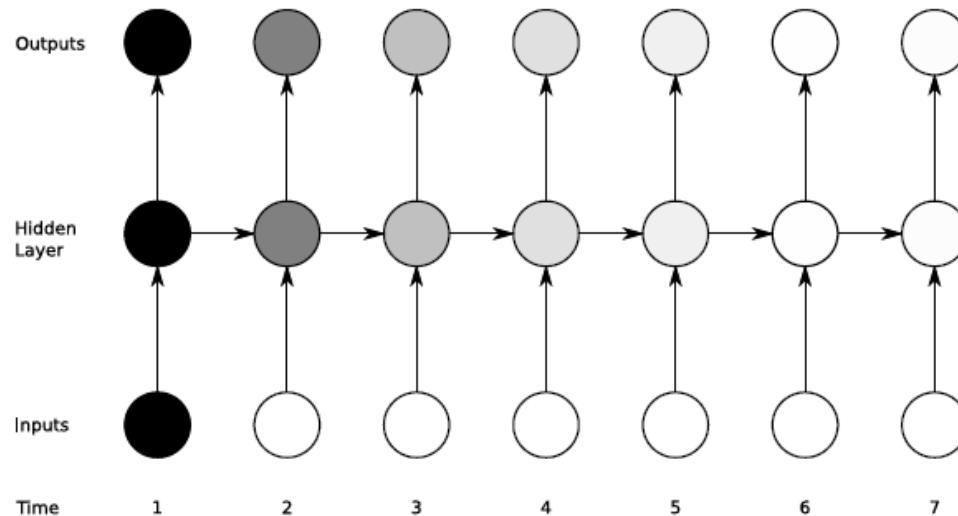
(U) $\nabla_U L = \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \frac{\partial \mathbf{h}^{(t)}}{\partial U}$



Long-term Dependencies

Vanishing gradient problem for RNNs:

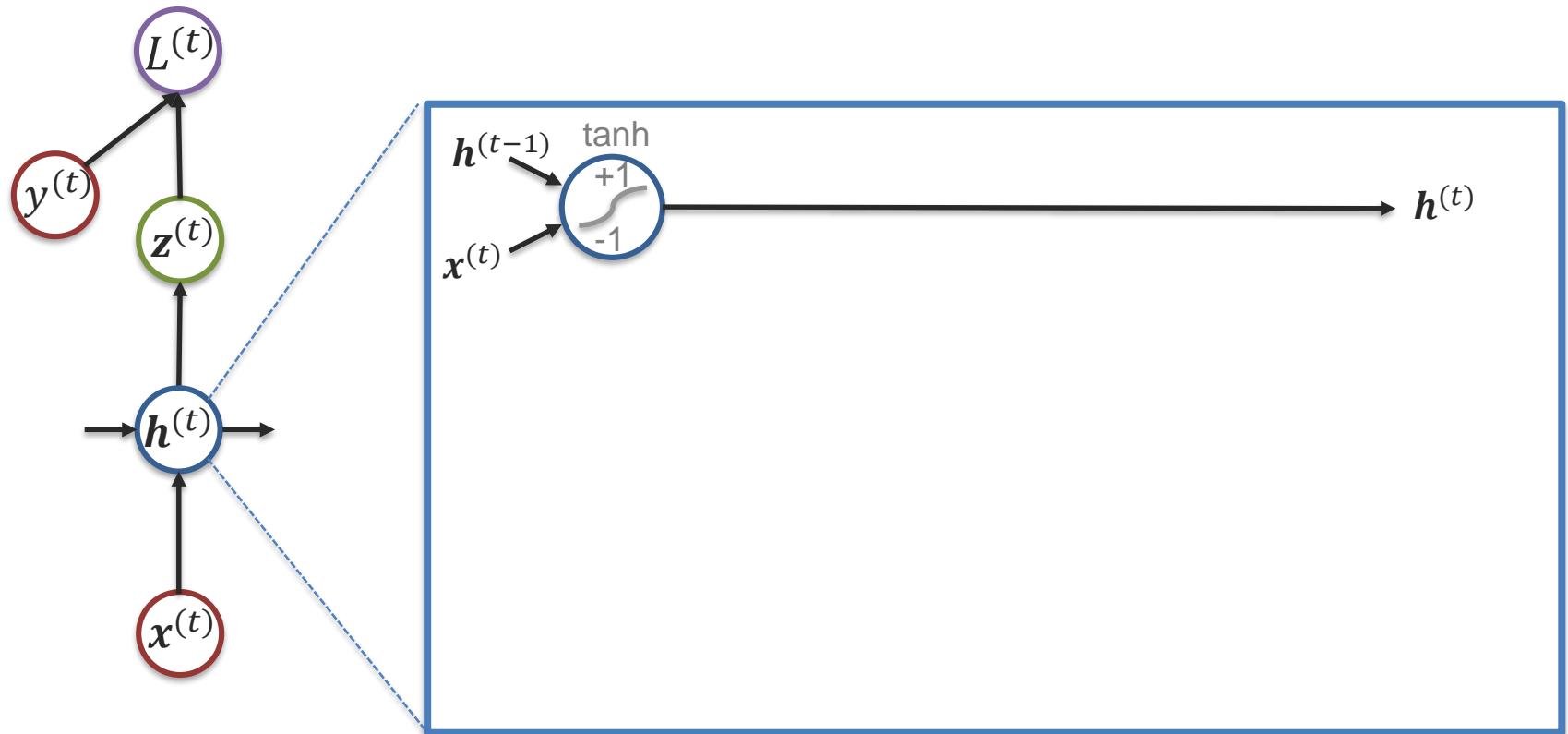
$$\mathbf{h}^{(t)} \sim \tanh(\mathbf{W}\mathbf{h}^{(t-1)})$$



- The influence of a given input on the hidden layer, and therefore on the network output, either decays or blows up exponentially as it cycles around the network's recurrent connections.



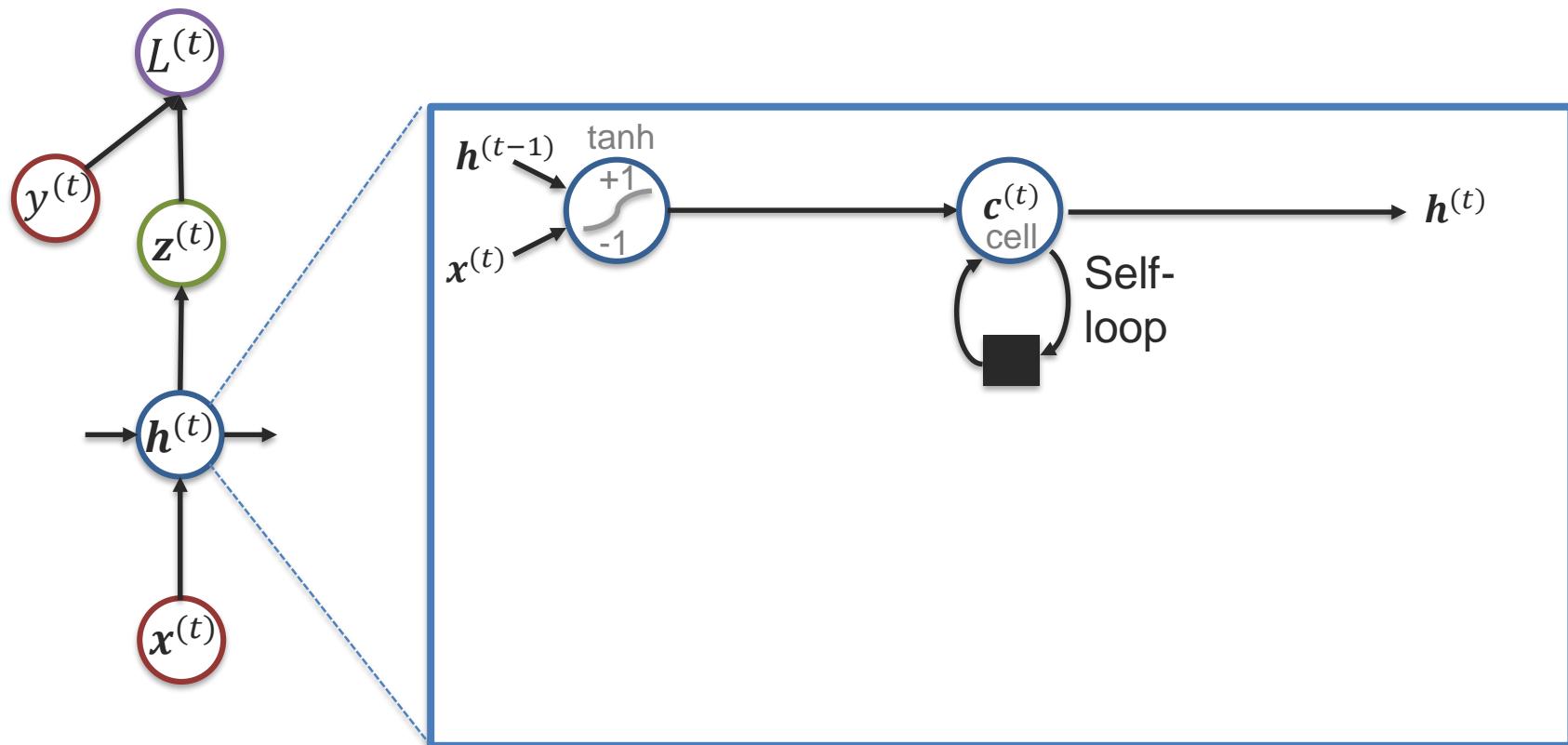
Recurrent Neural Networks



LSTM ideas: (1) “Memory” Cell and Self Loop

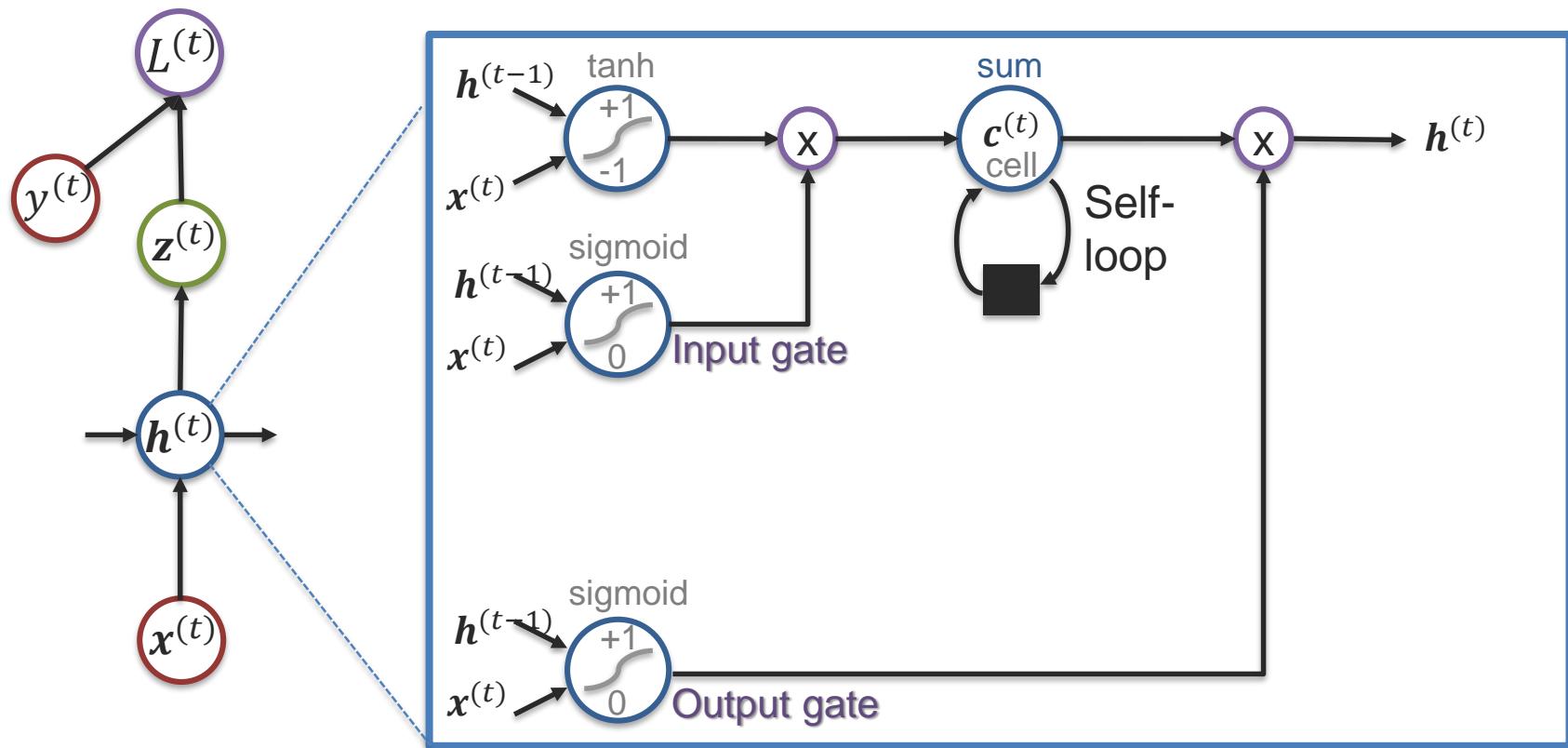
[Hochreiter and Schmidhuber, 1997]

Long Short-Term Memory (LSTM)



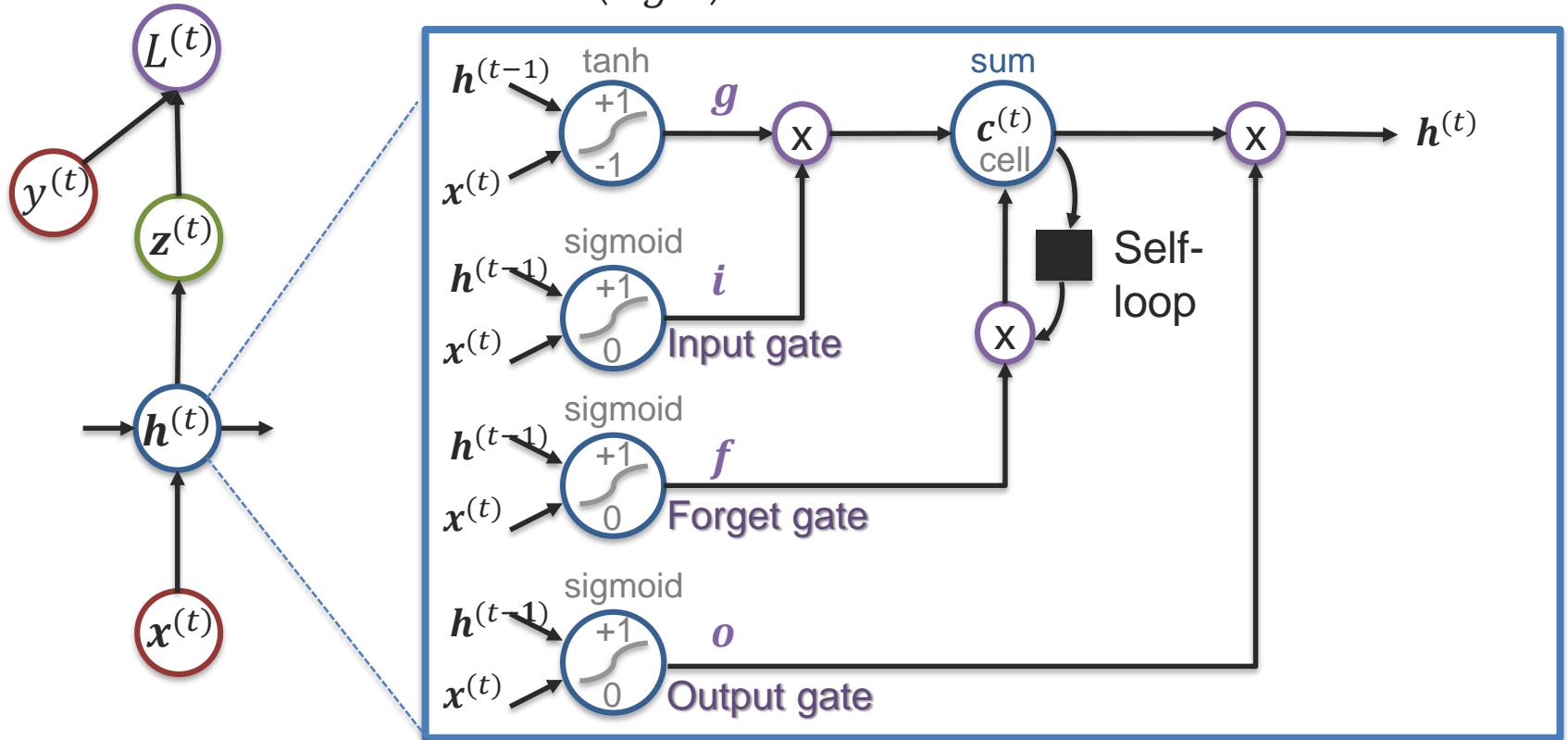
LSTM Ideas: (2) Input and Output Gates

[Hochreiter and Schmidhuber, 1997]

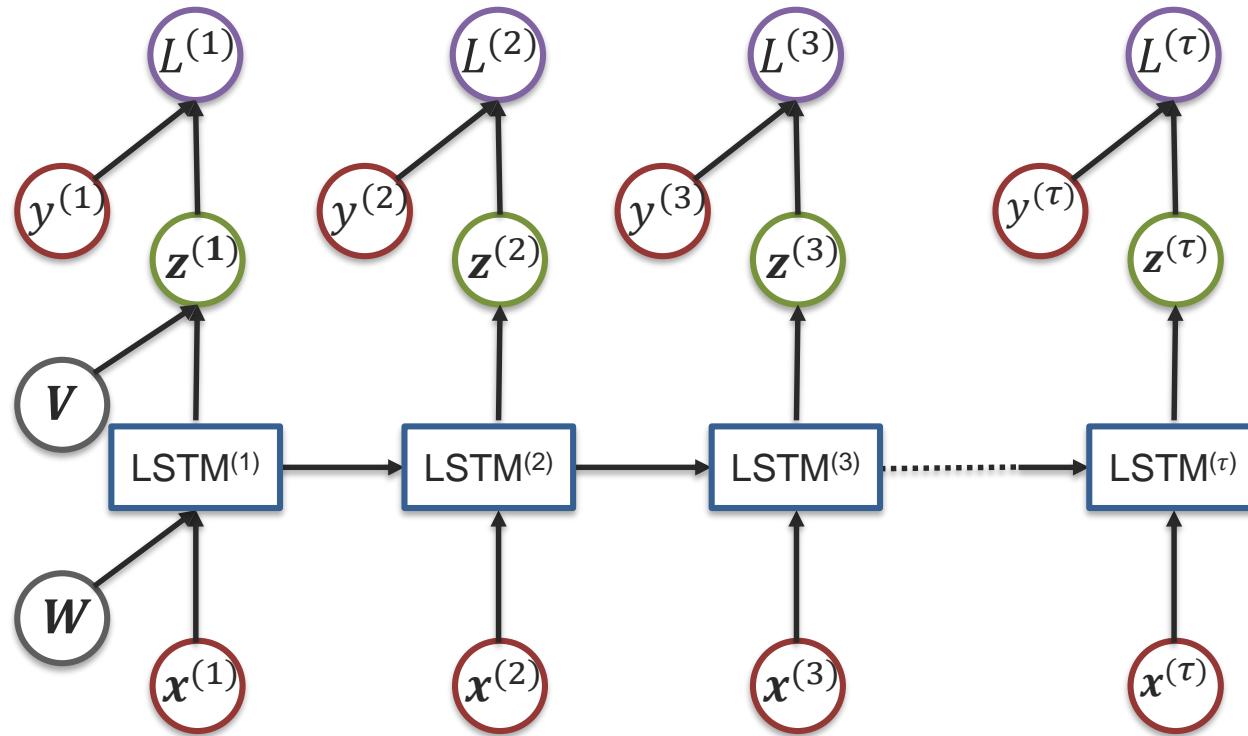


LSTM Ideas: (3) Forget Gate [Gers et al., 2000]

$$\begin{pmatrix} g \\ i \\ f \\ o \end{pmatrix} = \begin{pmatrix} \tanh \\ \text{sigm} \\ \text{sigm} \\ \text{sigm} \end{pmatrix} W \begin{pmatrix} h^{(t-1)} \\ x^{(t)} \end{pmatrix}$$
$$c^{(t)} = f \odot c^{(t-1)} + i \odot g$$
$$h^{(t)} = o \odot \tanh(c^{(t)})$$



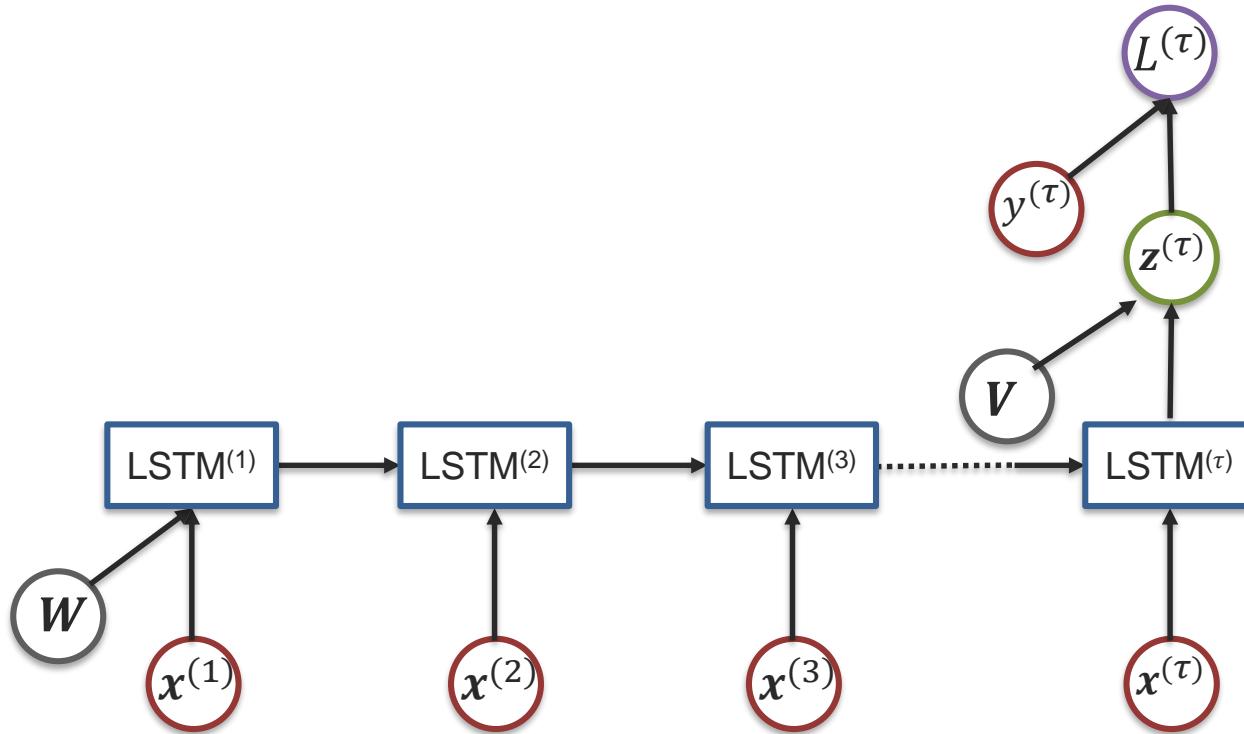
Recurrent Neural Network using LSTM Units



Gradient can still be computer using backpropagation!



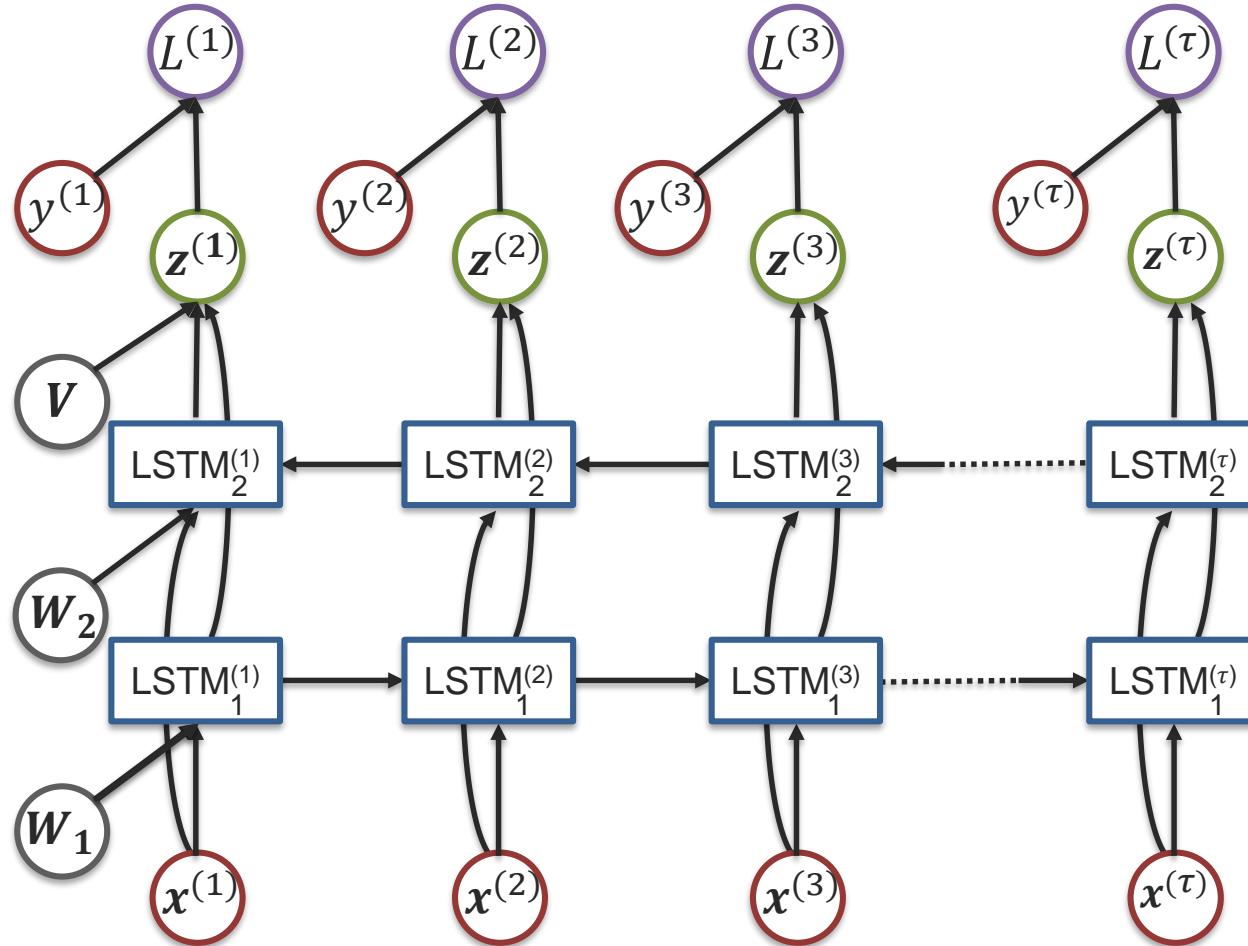
Recurrent Neural Network using LSTM Units



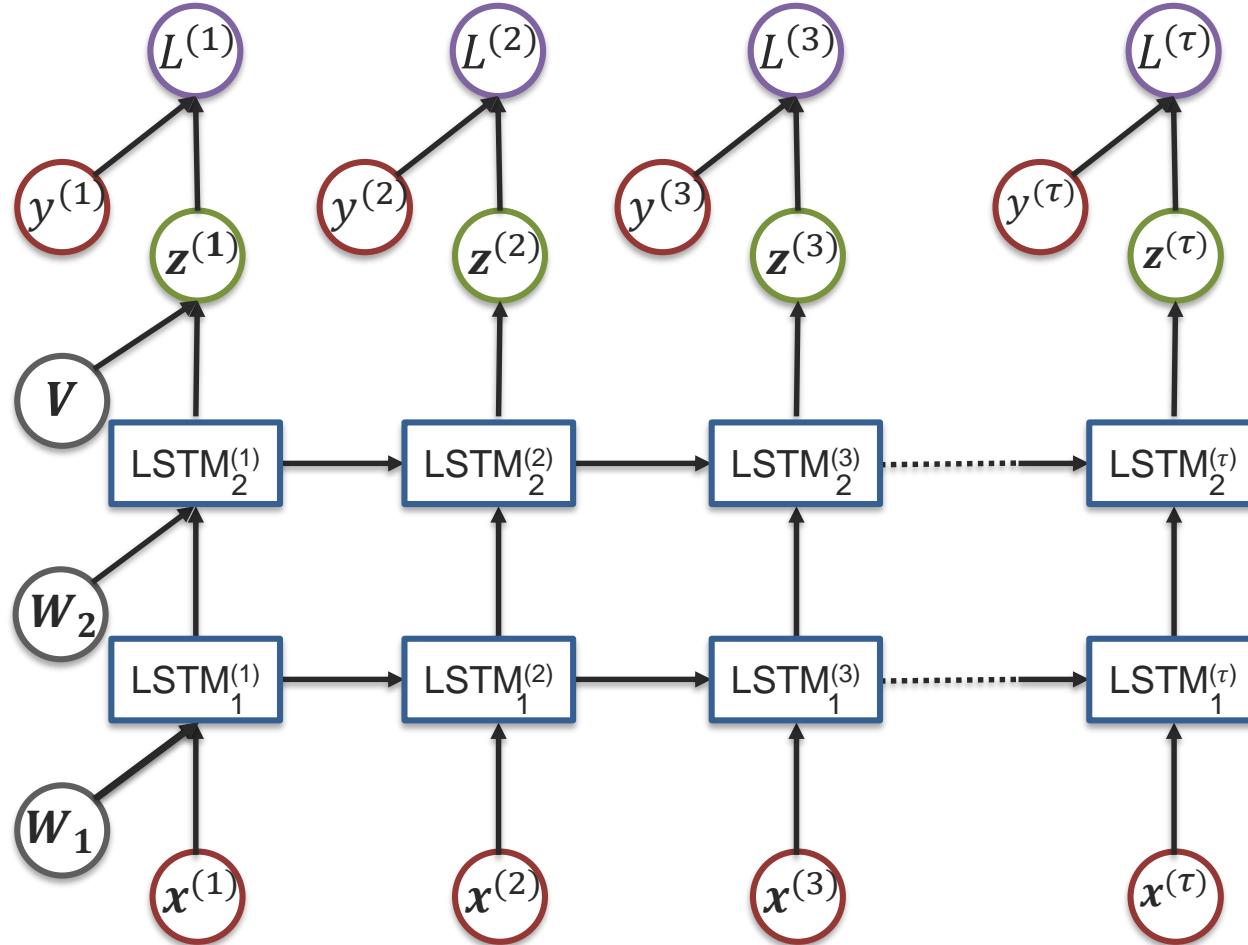
Gradient can still be computer using backpropagation!



Bi-directional LSTM Network



Deep LSTM Network



Translation and Alignment

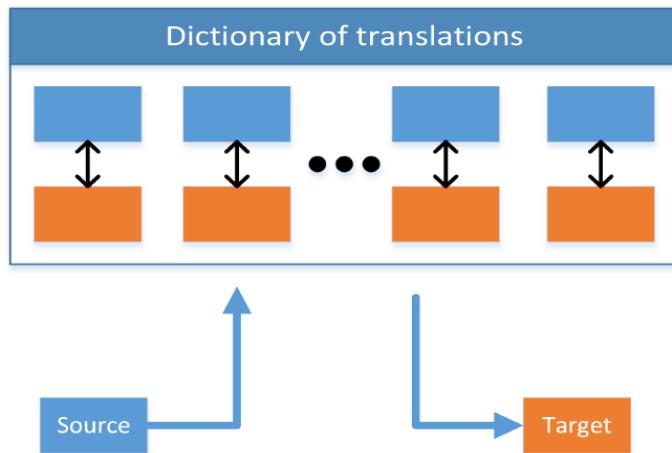


Core Challenge 4: Translation

Definition: Process of changing data from one modality to another, where the translation relationship can often be open-ended or subjective.

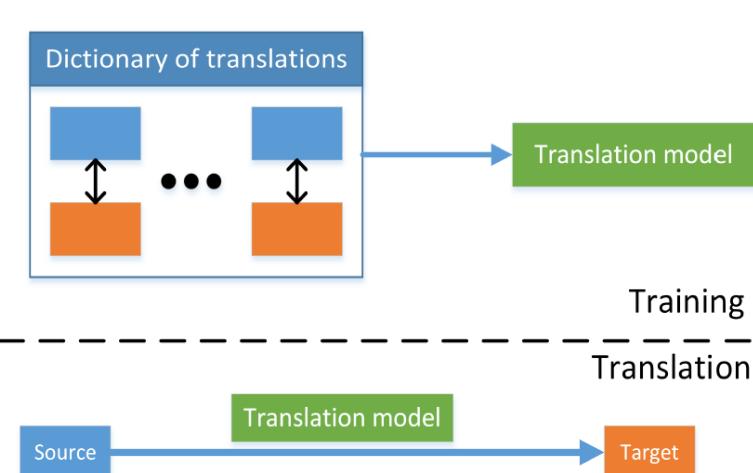
A

Example-based



B

Model-driven



Translation

➤ Visual animations



➤ Image captioning



➤ Speech synthesis



Challenges:

- I. Different representations
- II. Multiple source modalities
- III. Open ended translations
- IV. Subjective evaluation
- V. Repetitive processes



Example-based translation

- Cross-media retrieval – bounded task
- Multimodal representation plays a key role here



... 'Iniesta is really impressing me.' said Zinedine Nods of approval could be seen across the continent: Andres Iniesta was named the best player of Euro 2012. In six Spain games in Poland and Ukraine, Iniesta did not score once but appreciation for the 28-year-old extends well beyond goals, it is now as broad as Europe. Iniesta has not quite gained the inevitability of gravity but the reliability of his talent is unquestionable ...

[Wei et al. 2015]

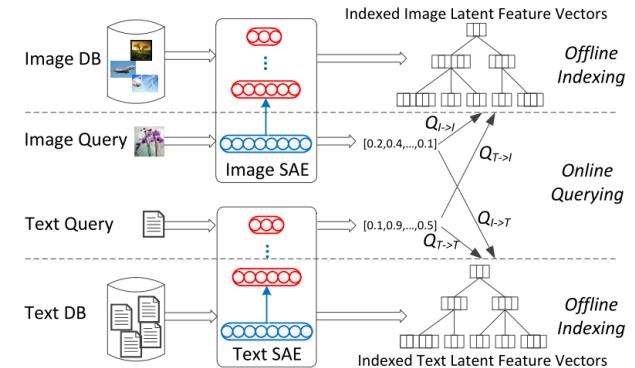
Kobe Bryant said, "To be really frank with you, I really do not look at it as that, for the simple fact that Michael Jordan has really taught me a lot. Really taught me a lot. The trainer of his, Tim Grover, he's passed on to me and I work with him a great deal, and he's shown me a lot. So I can't sit there and say, well, I'm trying to catch Michael Jordan at six, I want to pass him after six."



...

Example-based translation

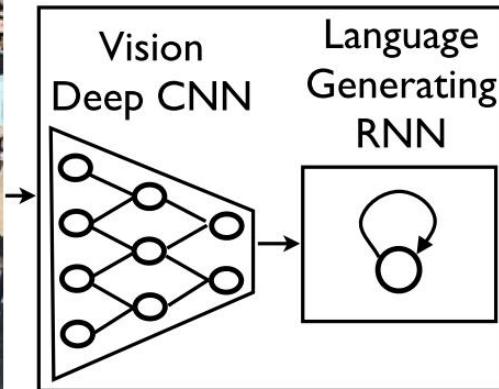
- Need a way to measure similarity between the modalities
- Remember multimodal representations
 - CCA
 - Coordinated
 - Joint
 - Hashing
- Can use pairs of instances to train them and retrieve closest ones during retrieval stage
- Objective and bounded task



[Wang et al. 2014]



Model-based Image captioning with Encoder-Decoder



A group of people shopping at an outdoor market.

There are many vegetables at the fruit stand.

[Vinyals et al., "Show and Tell: A Neural Image Caption Generator", CVPR 2015]

Visual Question Answering

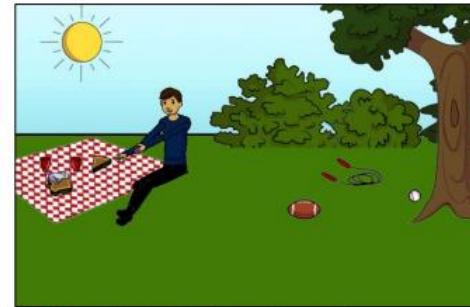
- A very new and exciting task created in part to address evaluation problems with the above task
- Task - Given an image and a question answer the question (<http://www.visualqa.org/>)



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?



Is this person expecting company?
What is just under the tree?



Does it appear to be rainy?
Does this person have 20/20 vision?



Language Technologies Institute

Carnegie Mellon University

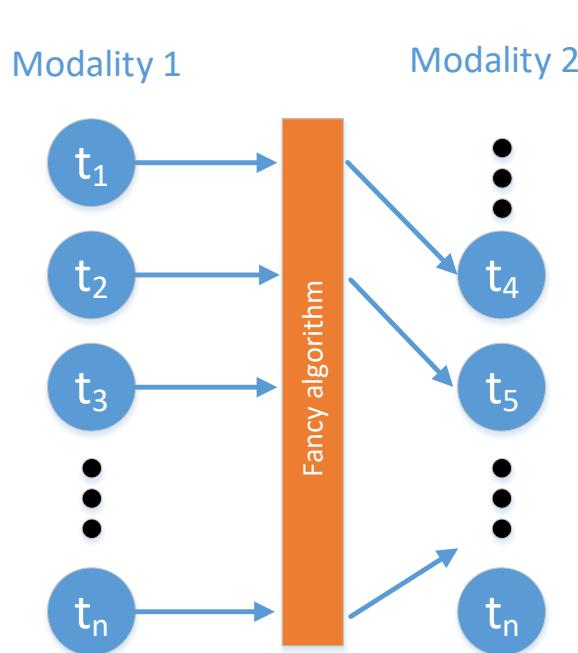
Evaluation on “Unbounded” Translations

- Tricky to do automatically!
- Ideally want humans to evaluate
 - What do you ask?
 - Can't use human evaluation for validating models – too slow and expensive
- Using standard machine translation metrics instead
 - BLEU, ROUGE CIDEER, Meteor



Core Challenge: Alignment

Definition: Identify the direct relations between (sub)elements from two or more different modalities.



A Explicit Alignment

The goal is to directly find correspondences between elements of different modalities

B Implicit Alignment

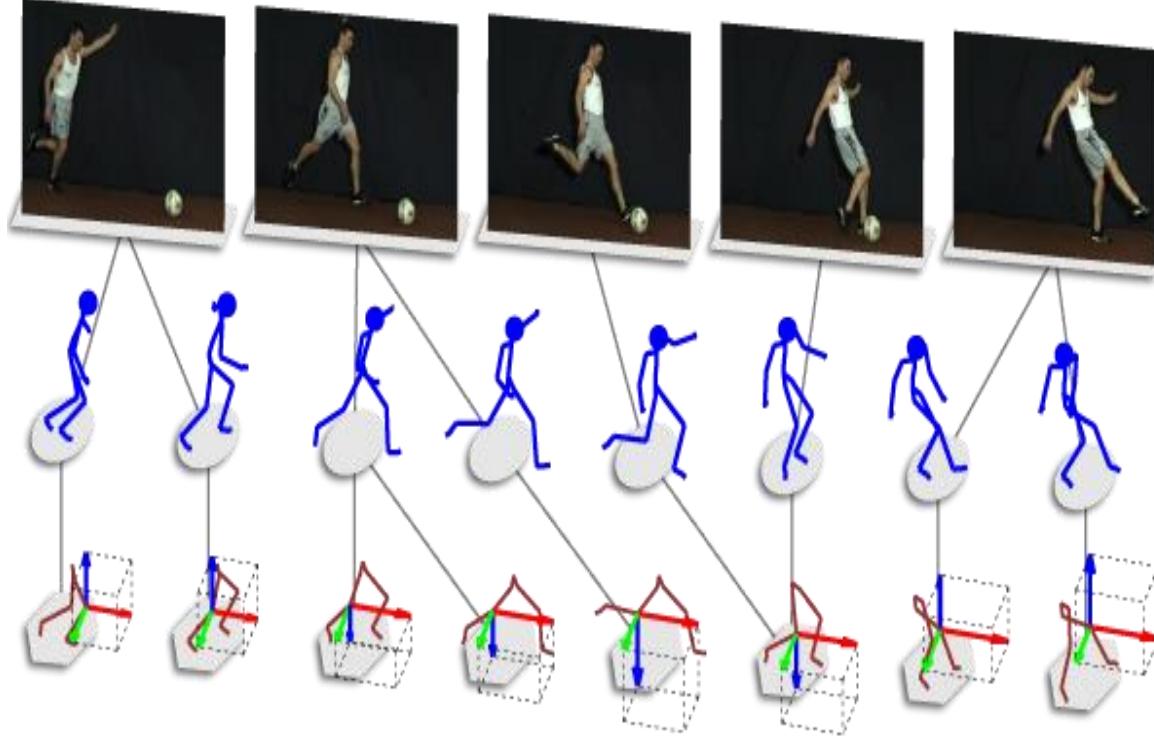
Uses internally latent alignment of modalities in order to better solve a different problem



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

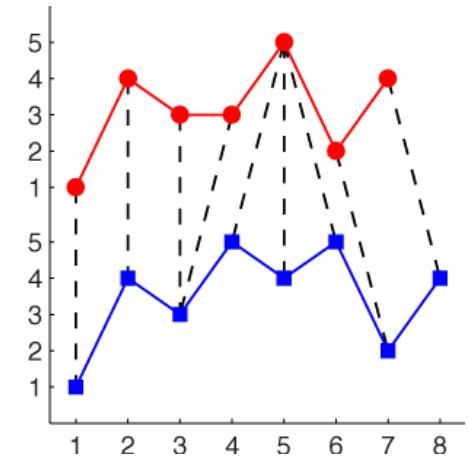
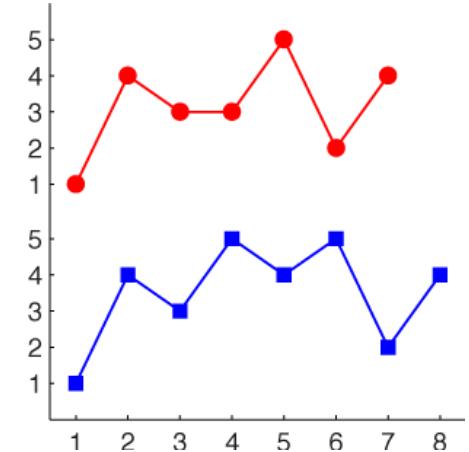


Let's start unimodal – Dynamic Time Warping

- We have two unaligned temporal unimodal signals
 - $\mathbf{X} = [x_1, x_2, \dots, x_{n_x}] \in \mathbb{R}^{d \times n_x}$
 - $\mathbf{Y} = [y_1, y_2, \dots, y_{n_y}] \in \mathbb{R}^{d \times n_y}$
- Find set of indices to minimize the alignment difference:

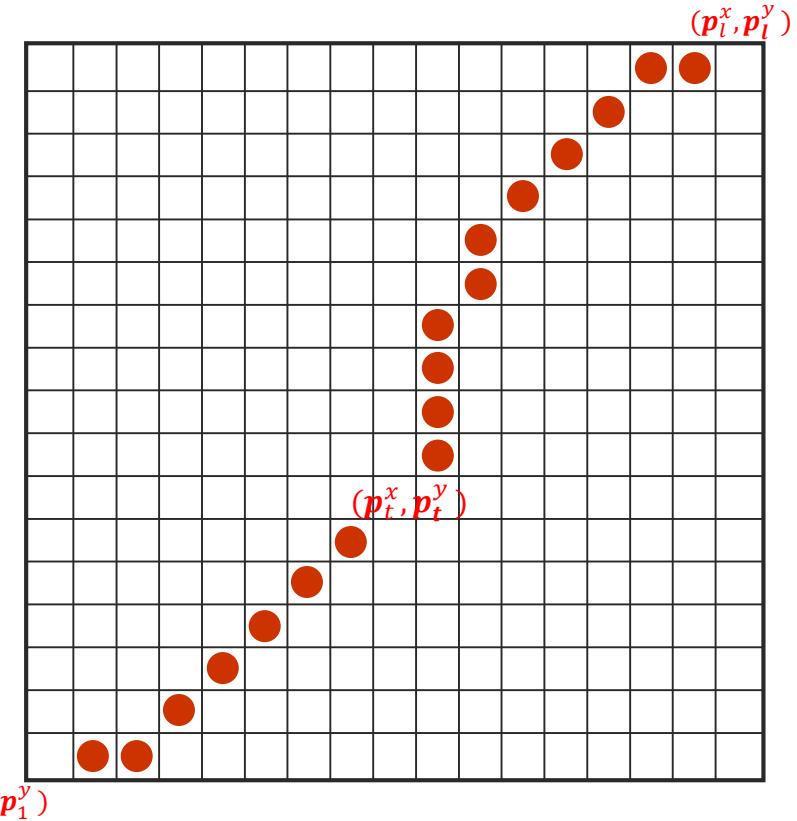
$$L(\mathbf{p}_t^x, \mathbf{p}_t^y) = \sum_{t=1}^l \|x_{\mathbf{p}_t^x} - y_{\mathbf{p}_t^y}\|_2^2$$

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



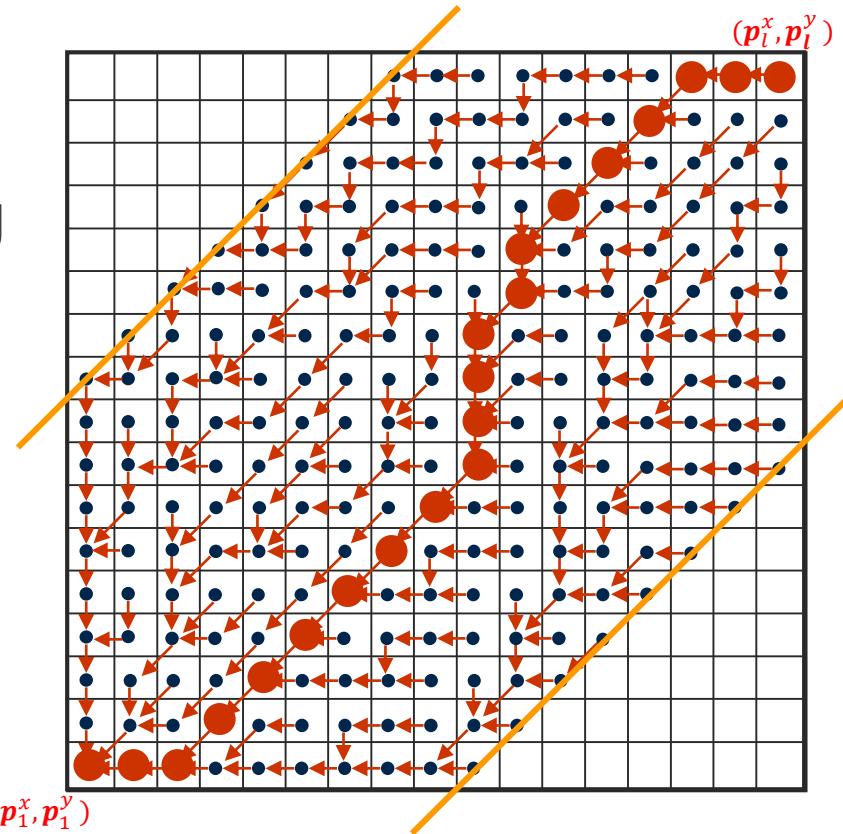
Dynamic Time Warping continued

- Lowest cost path in a cost matrix
- 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

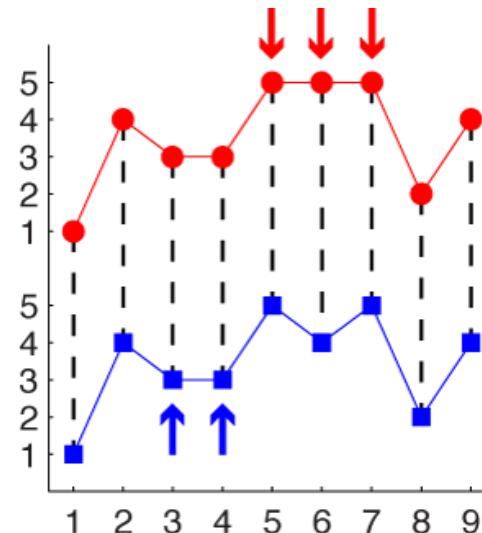
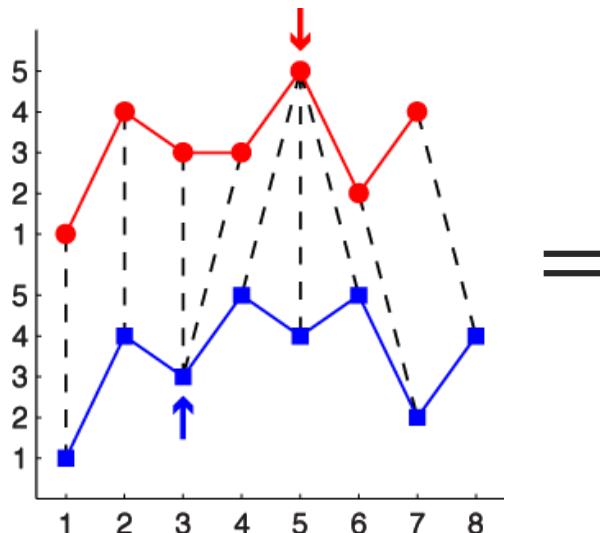
- Lowest cost path in a cost matrix
- Solved using dynamic programming whilst respecting the restrictions



DTW alternative formulation

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

Replication doesn't change the objective!



$$= \mathbf{X}$$

1	1	0	0	0	0	0	0	0	0
2	0	1	0	0	0	0	0	0	0
3	0	0	1	0	0	0	0	0	0
4	0	0	0	1	0	0	0	0	0
5	0	0	0	0	1	1	1	0	0
6	0	0	0	0	0	0	0	1	0
7	0	0	0	0	0	0	0	0	1
8	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0

$$= \mathbf{Y}$$

1	1	0	0	0	0	0	0	0	0
2	0	1	0	0	0	0	0	0	0
3	0	0	1	1	0	0	0	0	0
4	0	0	0	0	1	0	0	0	0
5	0	0	0	0	0	1	0	0	0
6	0	0	0	0	0	0	1	0	0
7	0	0	0	0	0	0	0	1	0
8	0	0	0	0	0	0	0	0	1
9	0	0	0	0	0	0	0	0	0

Alternative objective:

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



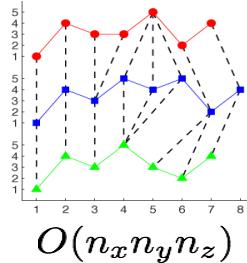
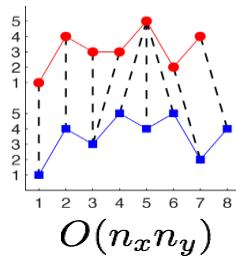
$$\text{Frobenius norm } \|A\|_F^2 = \sum_i \sum_j |a_{i,j}|^2$$

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

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

DTW - limitations

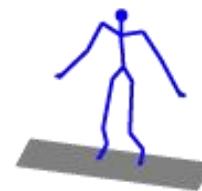
- Computationally complex



m sequences

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

- Sensitive to outliers
- Unimodal!



Canonical Correlation Analysis reminder

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

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

1

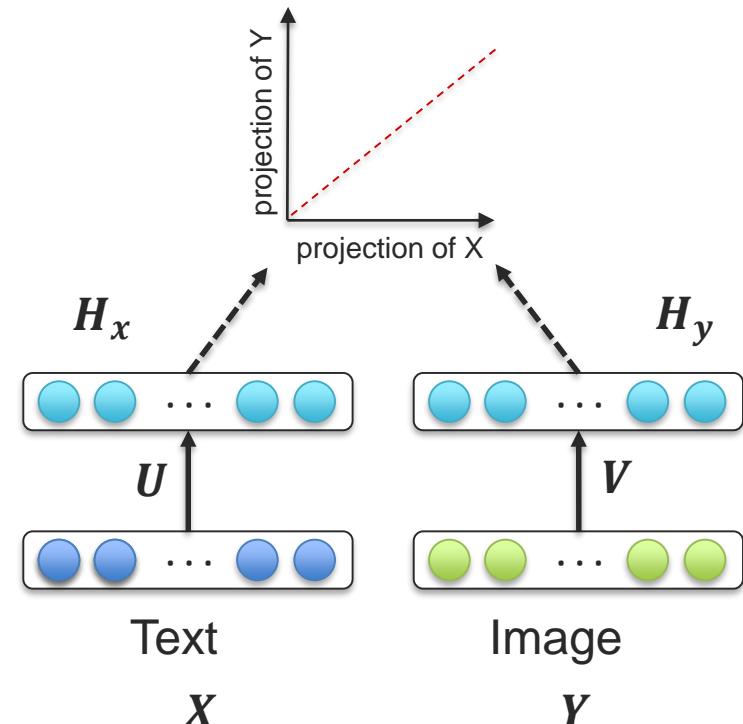
Linear projections maximizing correlation

2

Orthogonal projections

3

Unit variance of the projection vectors

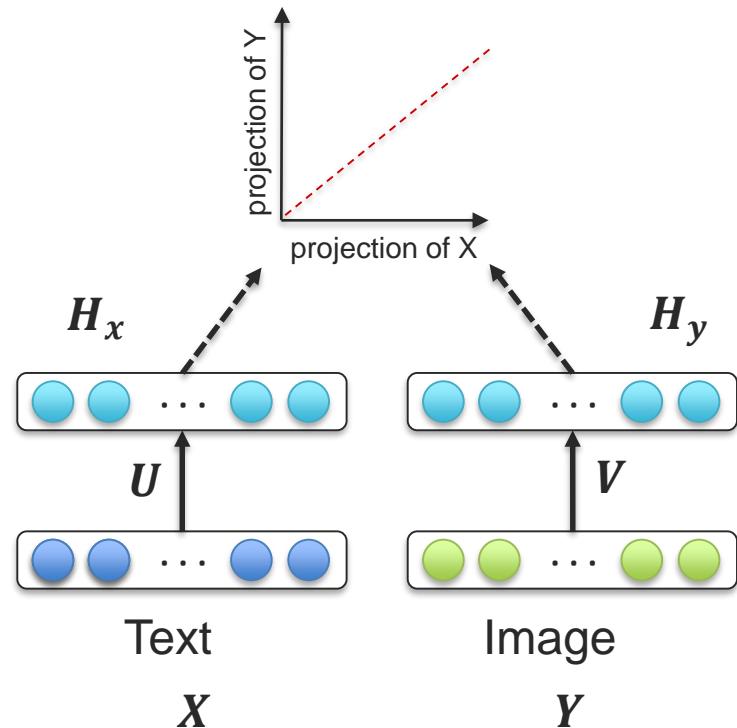


Canonical Correlation Analysis reminder

- 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: $\mathbf{U}^T \Sigma_{YY} \mathbf{U} = \mathbf{V}^T \Sigma_{YY} \mathbf{V} = I$



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)
- $\mathbf{W}_x, \mathbf{W}_y$ – temporal alignment
- \mathbf{U}, \mathbf{V} – cross-modal (spatial) alignment

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

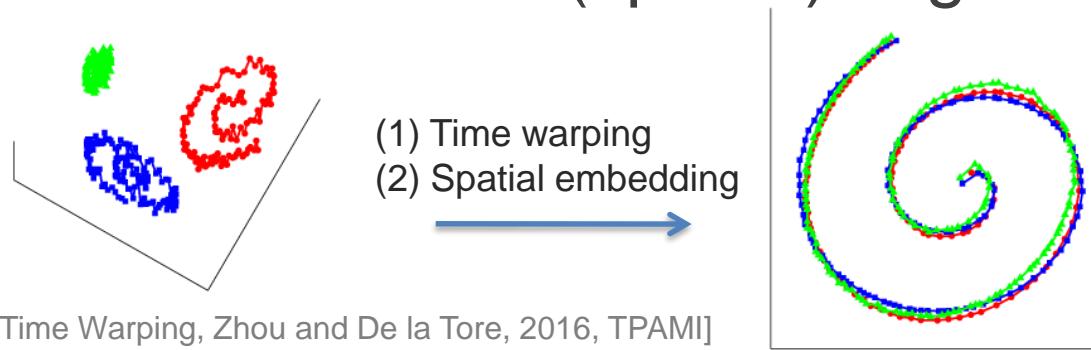


Generalized Time warping

- Generalize to multiple sequences all of different modality

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

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



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



Language Technologies Institute

Carnegie Mellon University

Alignment examples (unimodal)

CMU Motion Capture

Subject 1: 199 frames

Subject 2: 217 frames

Subject 3: 222 frames

1/199



1/217



1/222



Weizmann

Subject 1: 40 frames

Subject 2: 44 frames

Subject 3: 43 frames

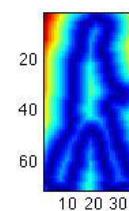
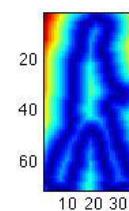
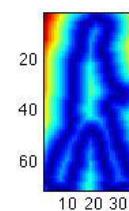
1/40



1/44

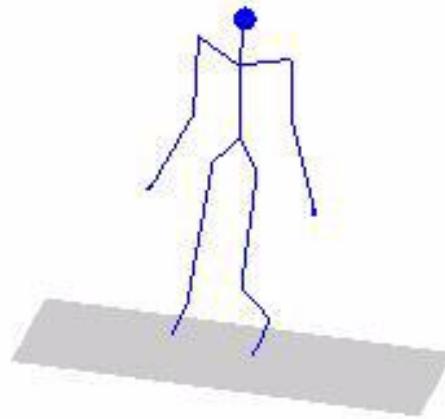


1/43



Alignment examples (multimodal)

1/273



1/51



1/127

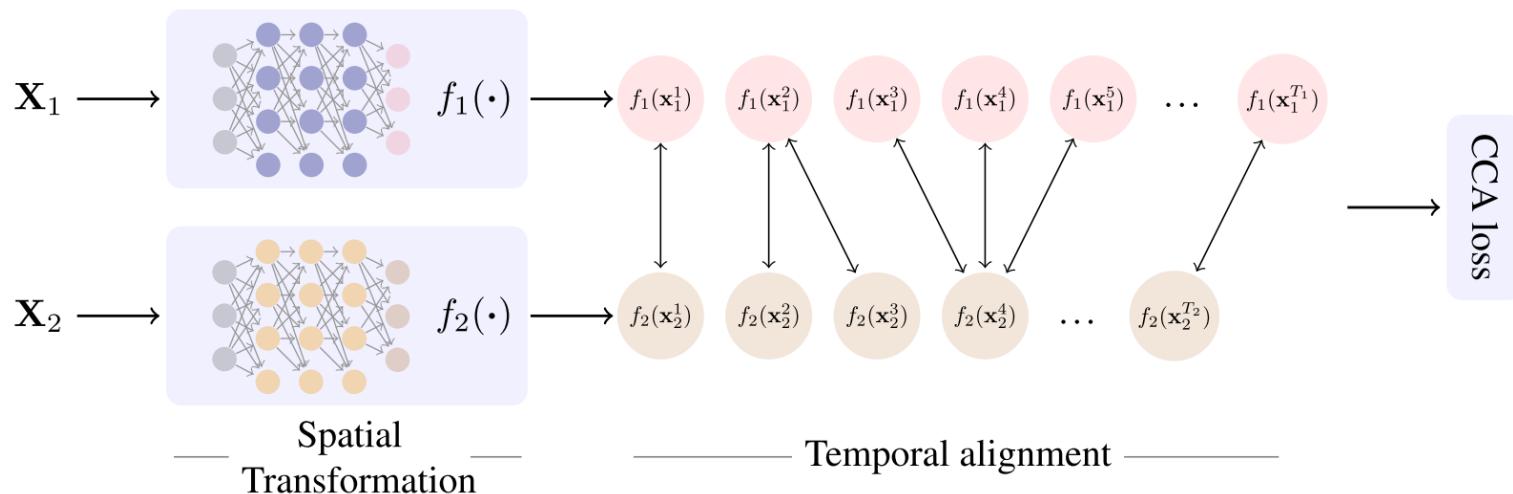


But how to model non-linear alignment functions?

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



Machine Translation

- Given a sentence in one language translate it to another

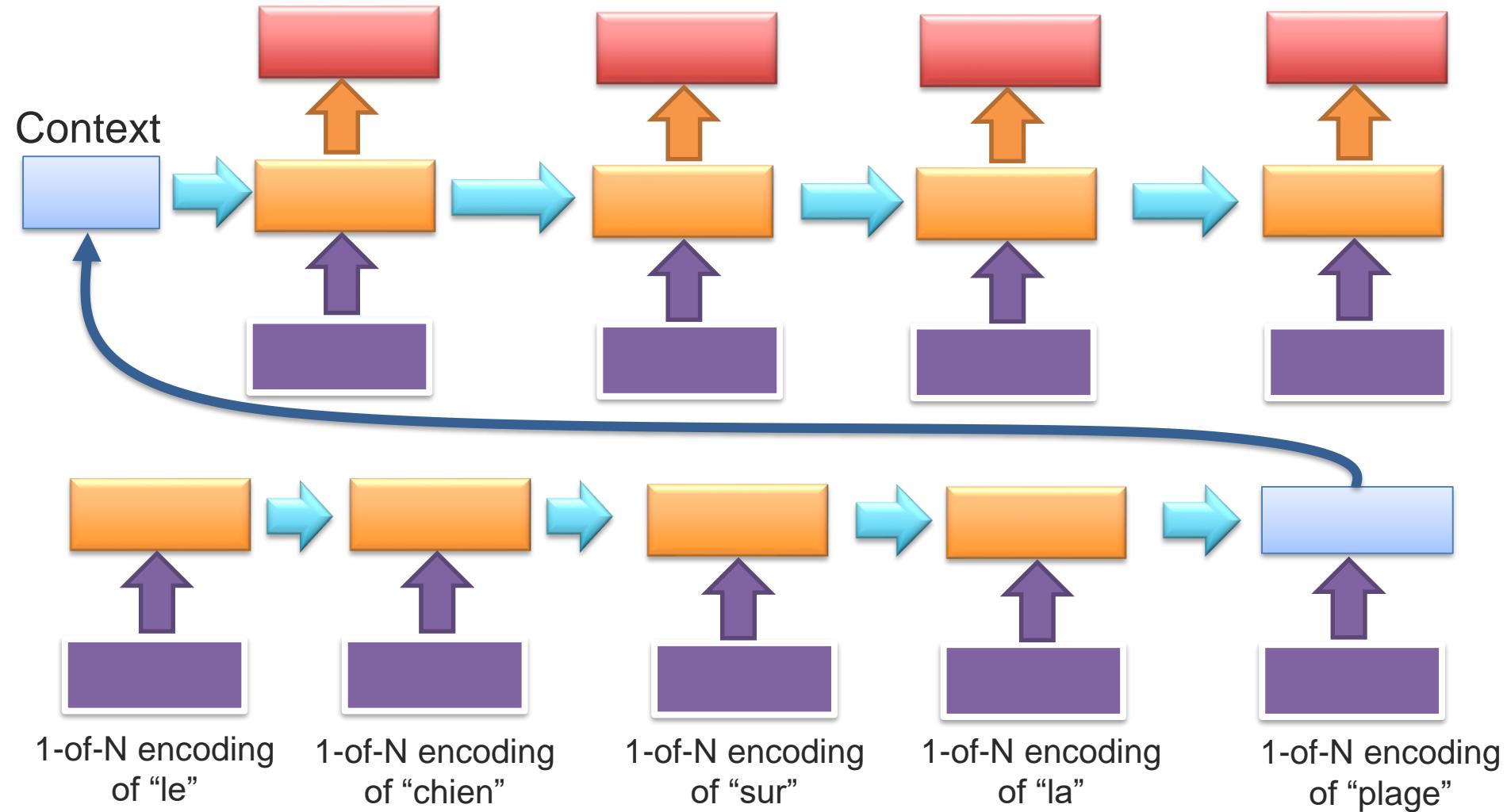
Dog on the beach → le chien sur la plage

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



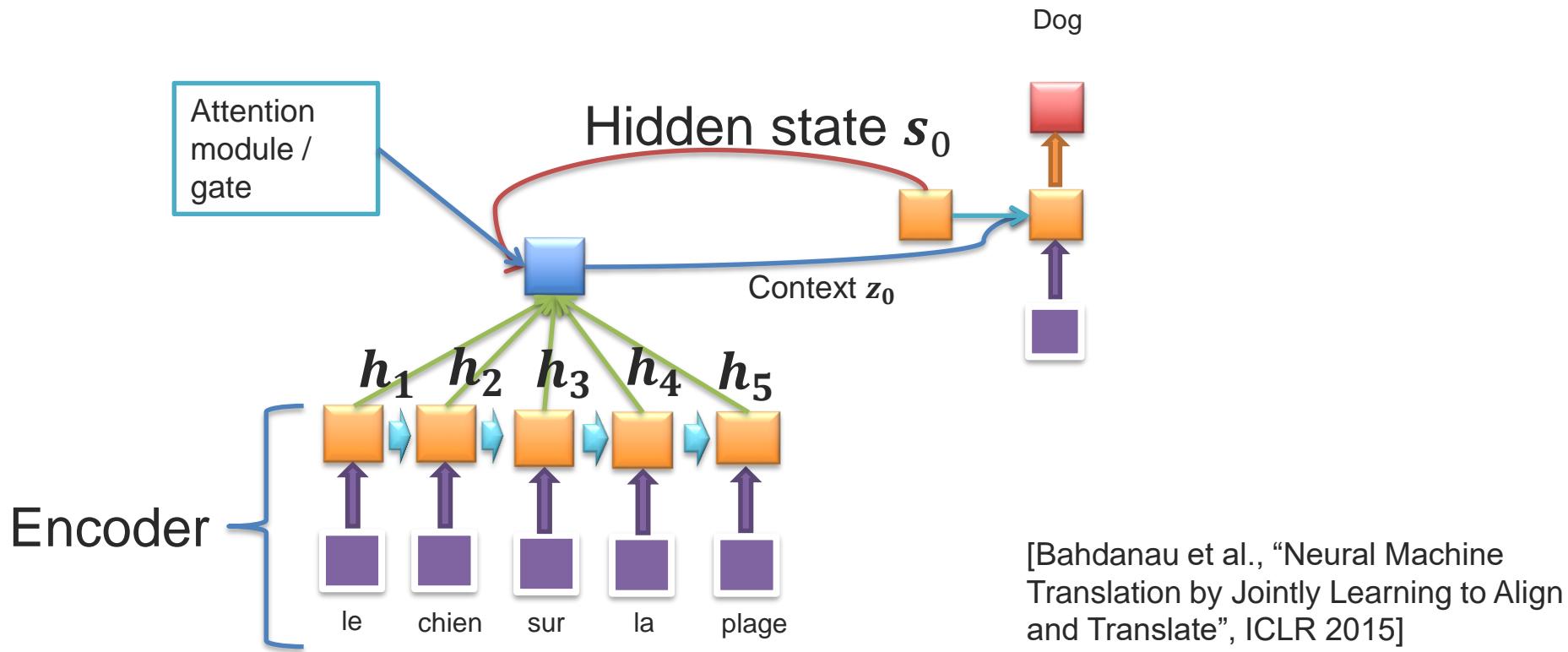
Encoder-Decoder Architecture for Machine Translation

[Cho et al., "Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation", EMNLP 2014]



Attention Model for Machine Translation

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

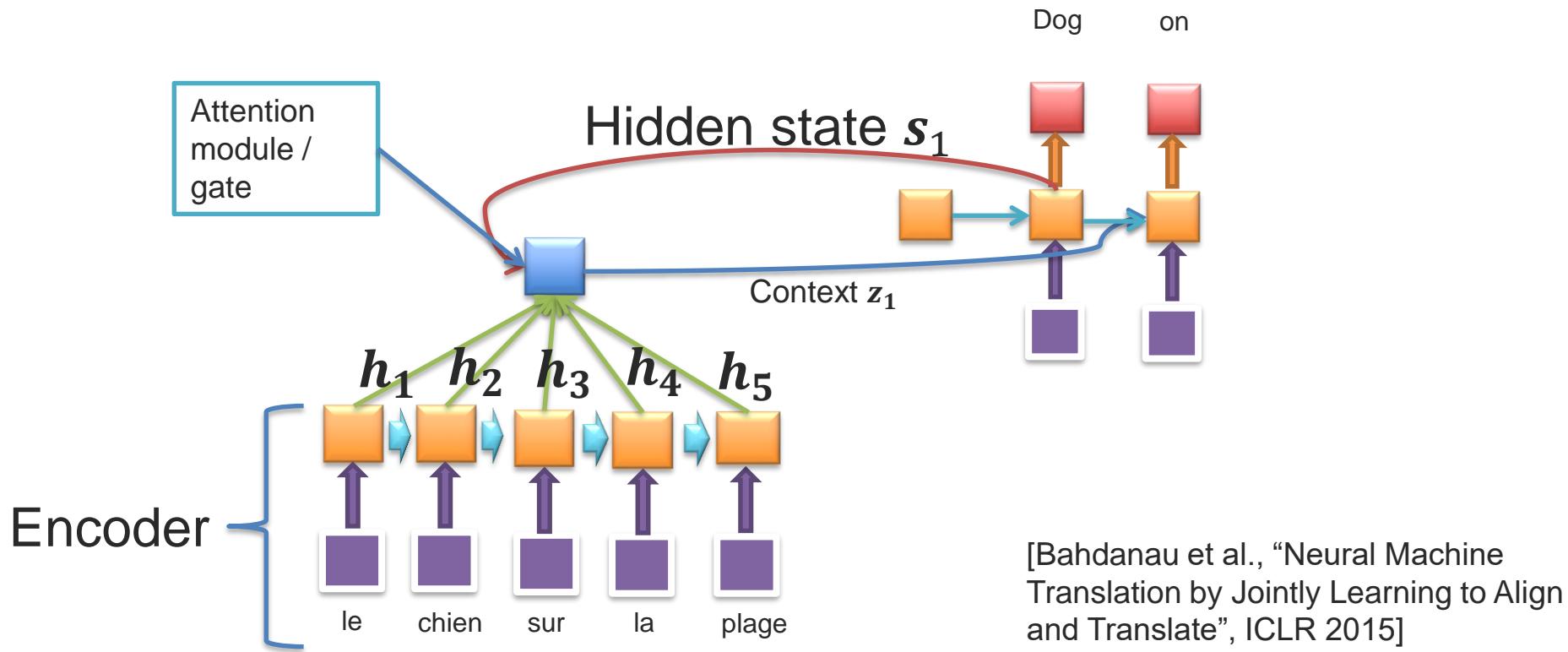


[Bahdanau et al., “Neural Machine Translation by Jointly Learning to Align and Translate”, ICLR 2015]



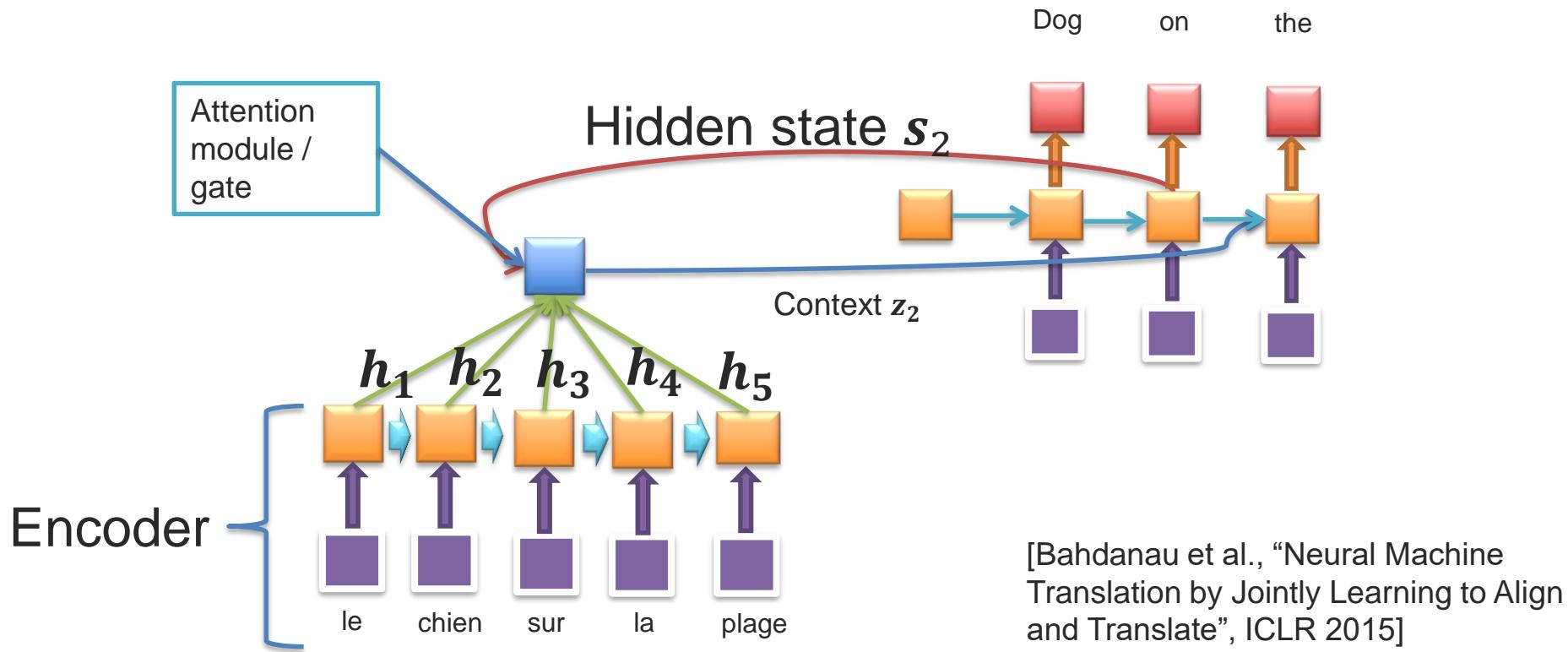
Attention Model for Machine Translation

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

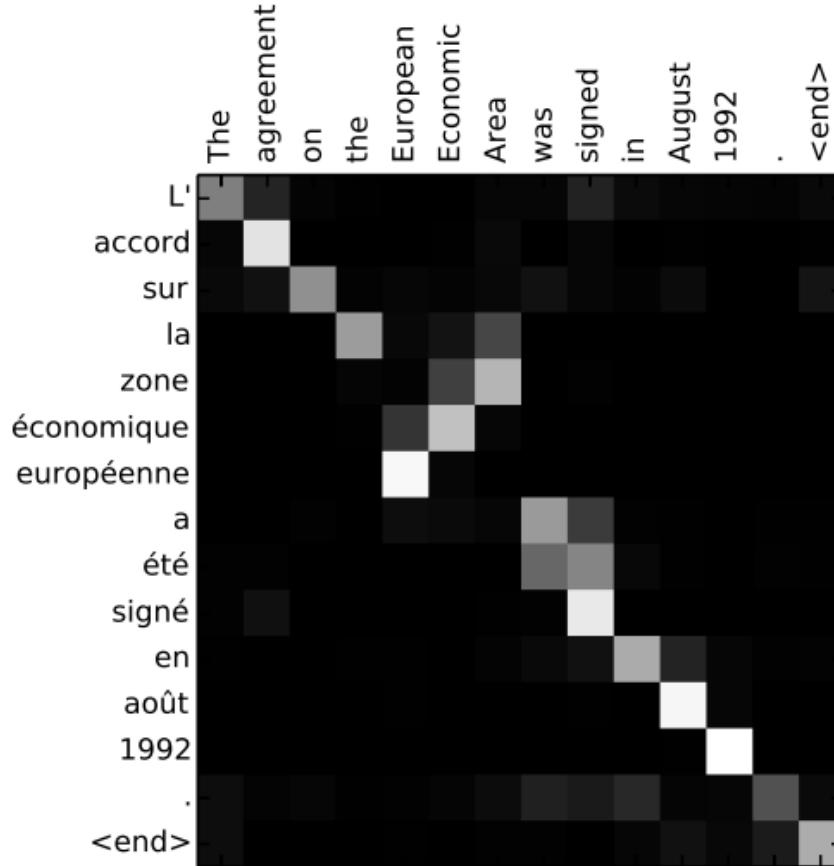


Attention Model for Machine Translation

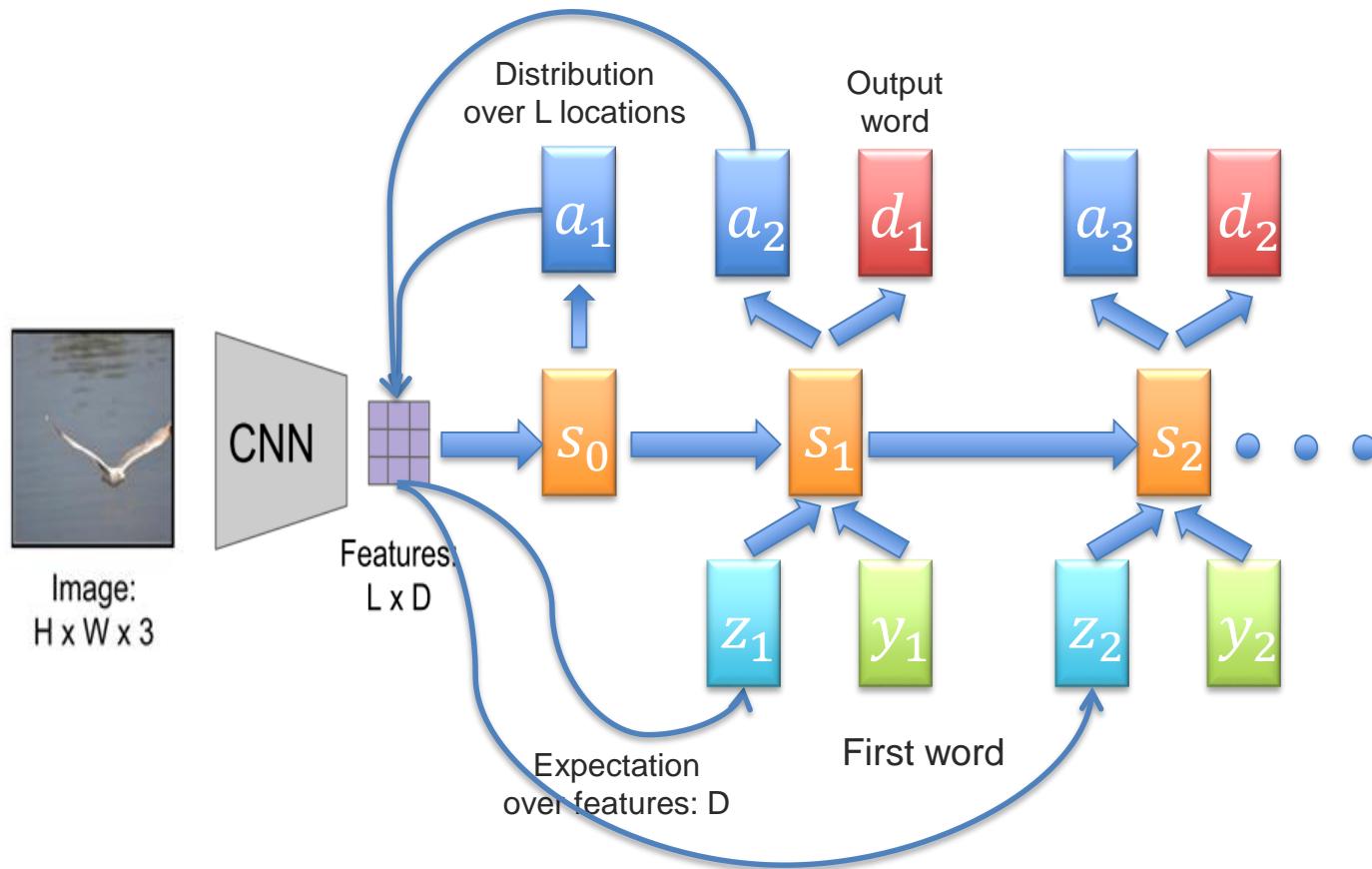
- Before encoder would just take the final hidden state, now we actually care about the intermediate hidden states



Attention Model for Machine Translation



Attention Model for Image Captioning



Attention Model for Image Captioning



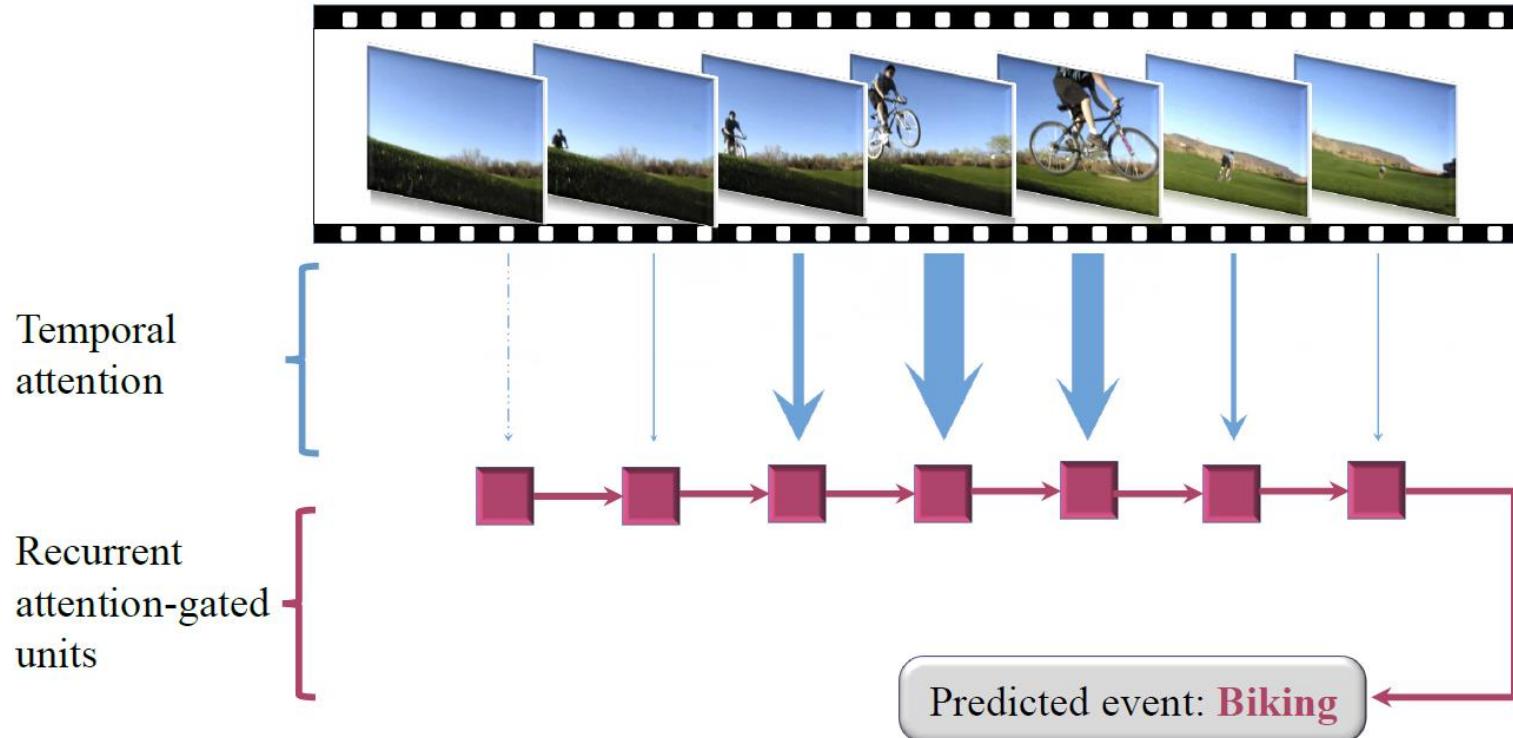
Xu et.al., ICML 2015



Language Technologies Institute

Carnegie Mellon University

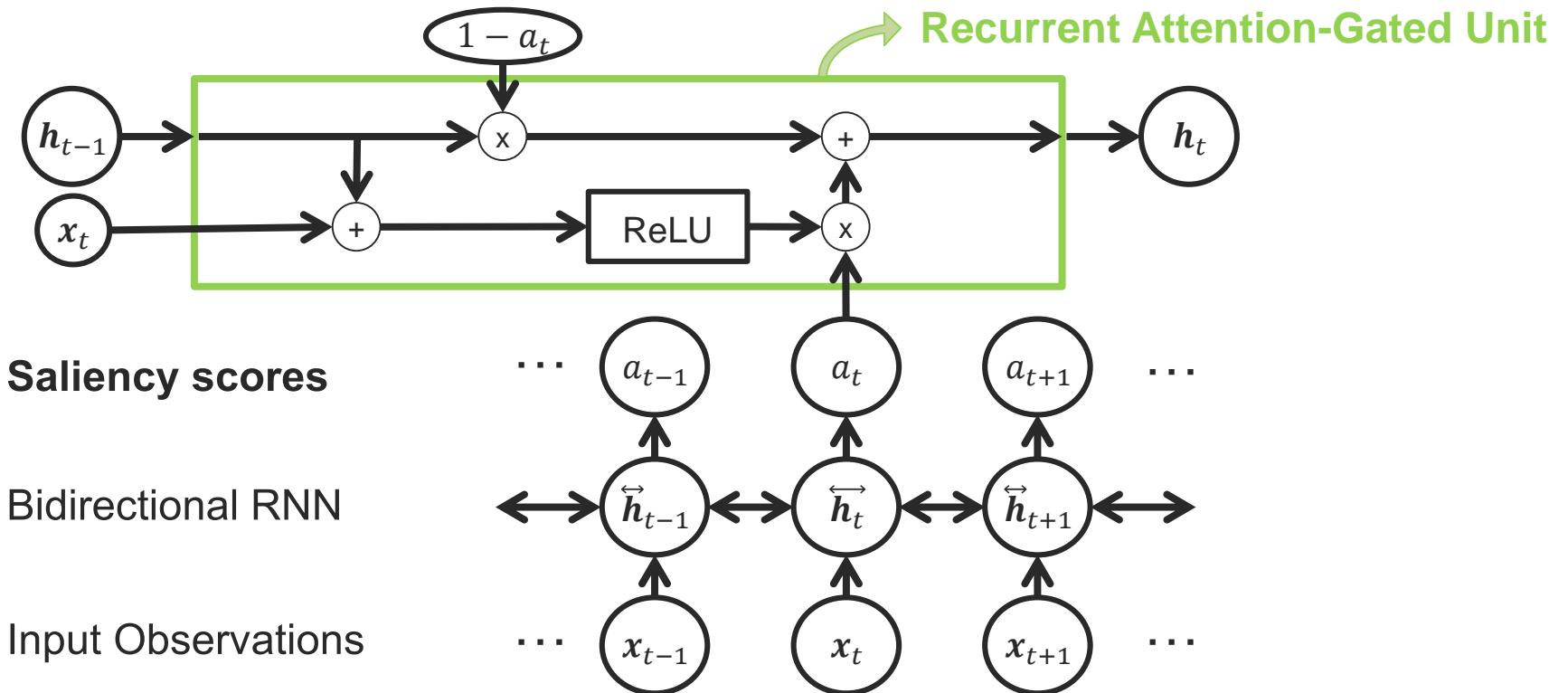
Attention Model for Video Sequences



[Pei, Baltrušaitis, Tax and Morency. Temporal Attention-Gated Model for Robust Sequence Classification, **CVPR**, 2017]



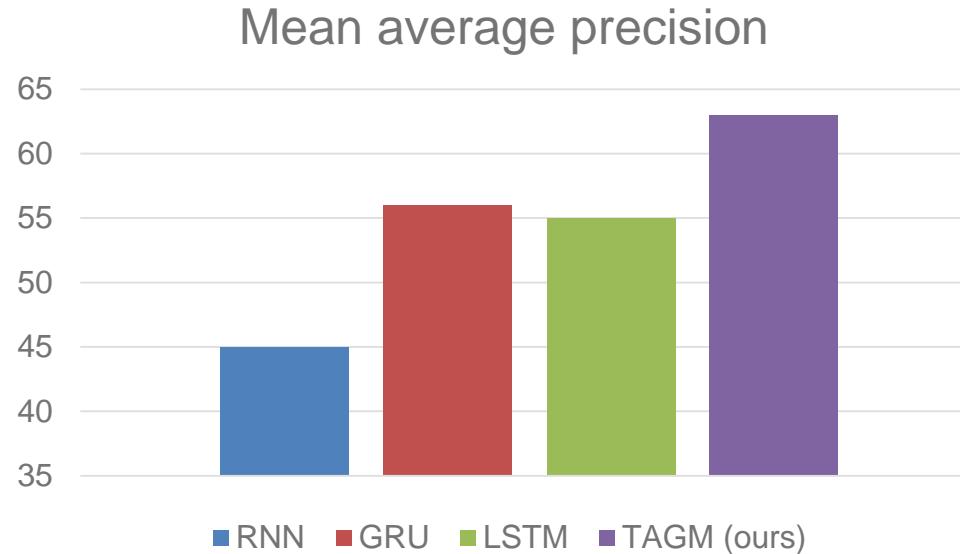
Temporal Attention-Gated Model (TAGM)



Temporal Attention Gated Model (TAGM)

CCV dataset

- 20 video categories
- Biking, birthday, wedding etc.



[Pei, Baltrušaitis, Tax and Morency. Temporal Attention-Gated Model for Robust Sequence Classification, CVPR, 2017]



Temporal Attention Gated Model (TAGM)



Biking

[Pei, Baltrušaitis, Tax and Morency. Temporal Attention-Gated Model for Robust Sequence Classification, **CVPR**, 2017]

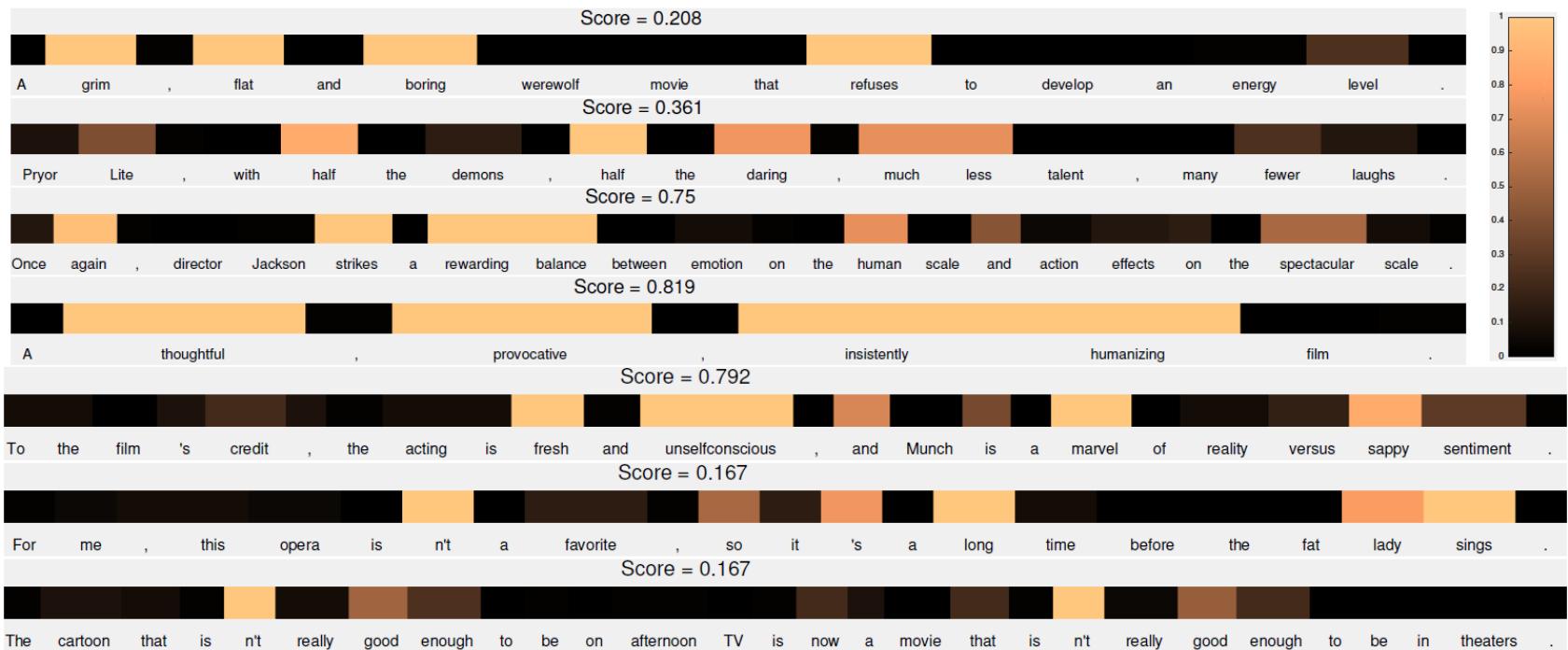


Language Technologies Institute

Carnegie Mellon University

Temporal Attention Gated Model (TAGM)

Text-based Sentiment Analysis



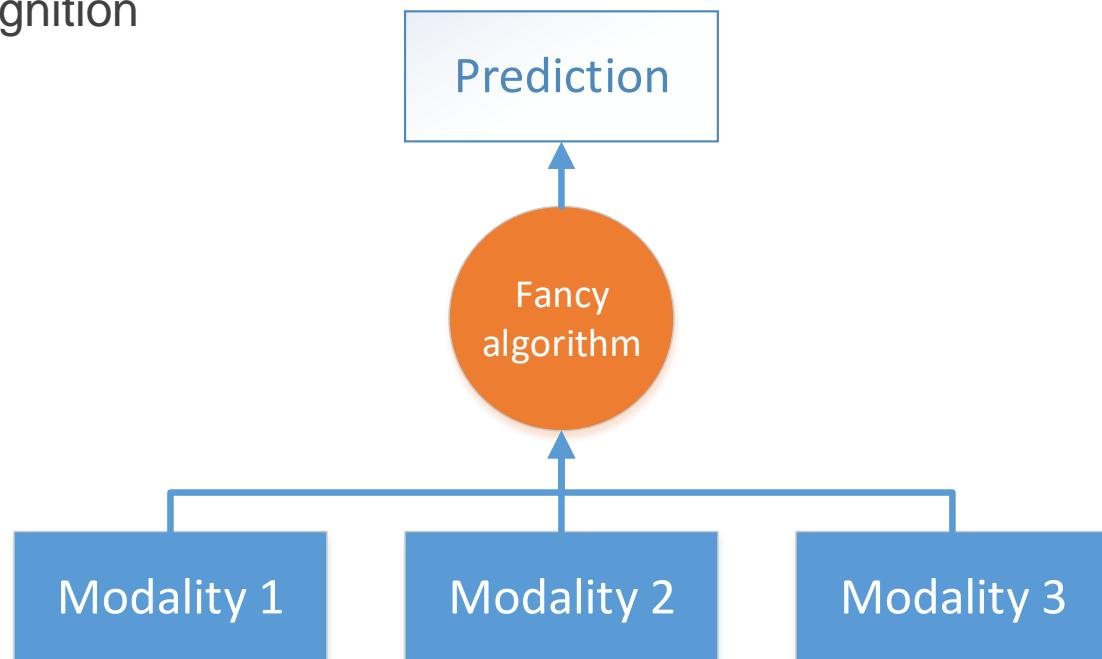
[Pei, Baltrušaitis, Tax and Morency. Temporal Attention-Gated Model for Robust Sequence Classification, CVPR, 2017]



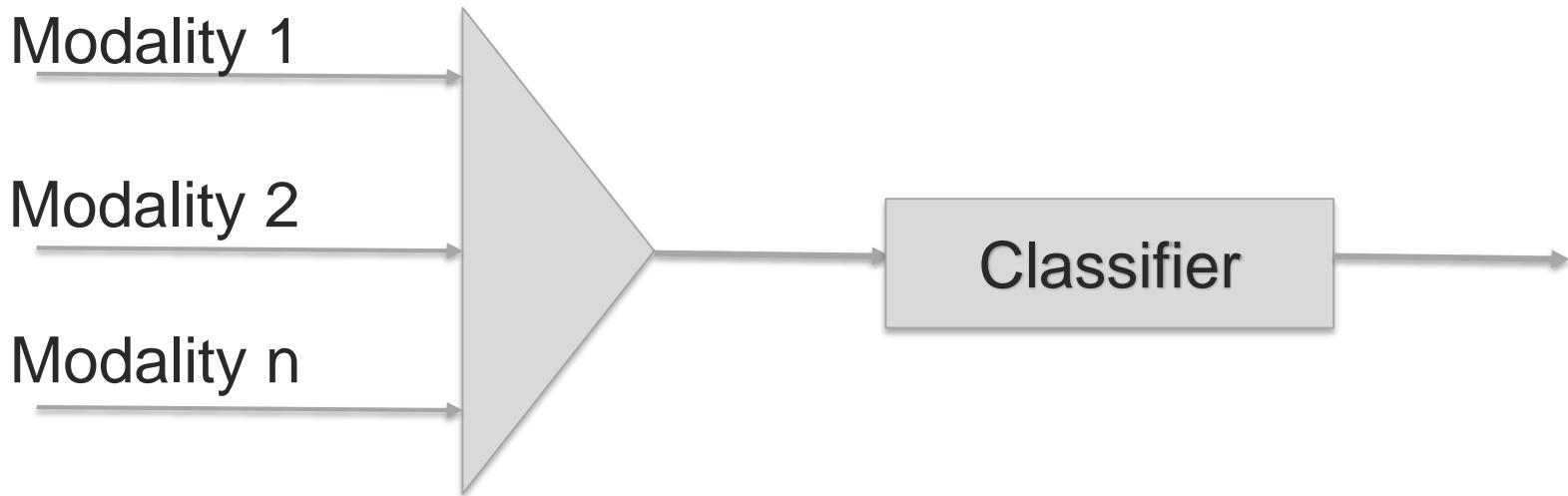
Multimodal Fusion

Multimodal Fusion

- Process of joining information from two or more modalities to perform a prediction
 - One of the earlier and more established problems
 - e.g. audio-visual speech recognition, multimedia event detection, multimodal emotion recognition
- Two major types
- Model Free
 - Early, late, hybrid
- Model Based
 - Kernel Methods
 - Graphical models
 - Neural networks



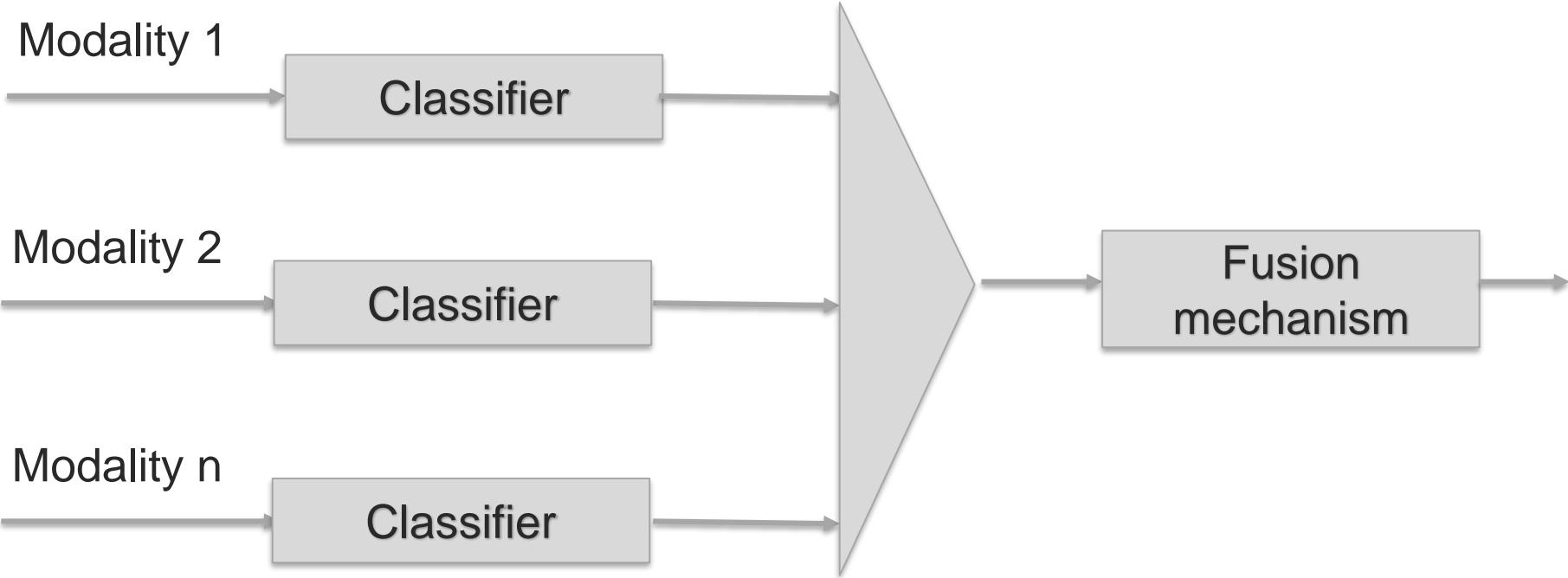
Model free approaches – early fusion



- Easy to implement – just concatenate the features
- Exploit dependencies between features
- Can end up very high dimensional
- More difficult to use if features have different framerates



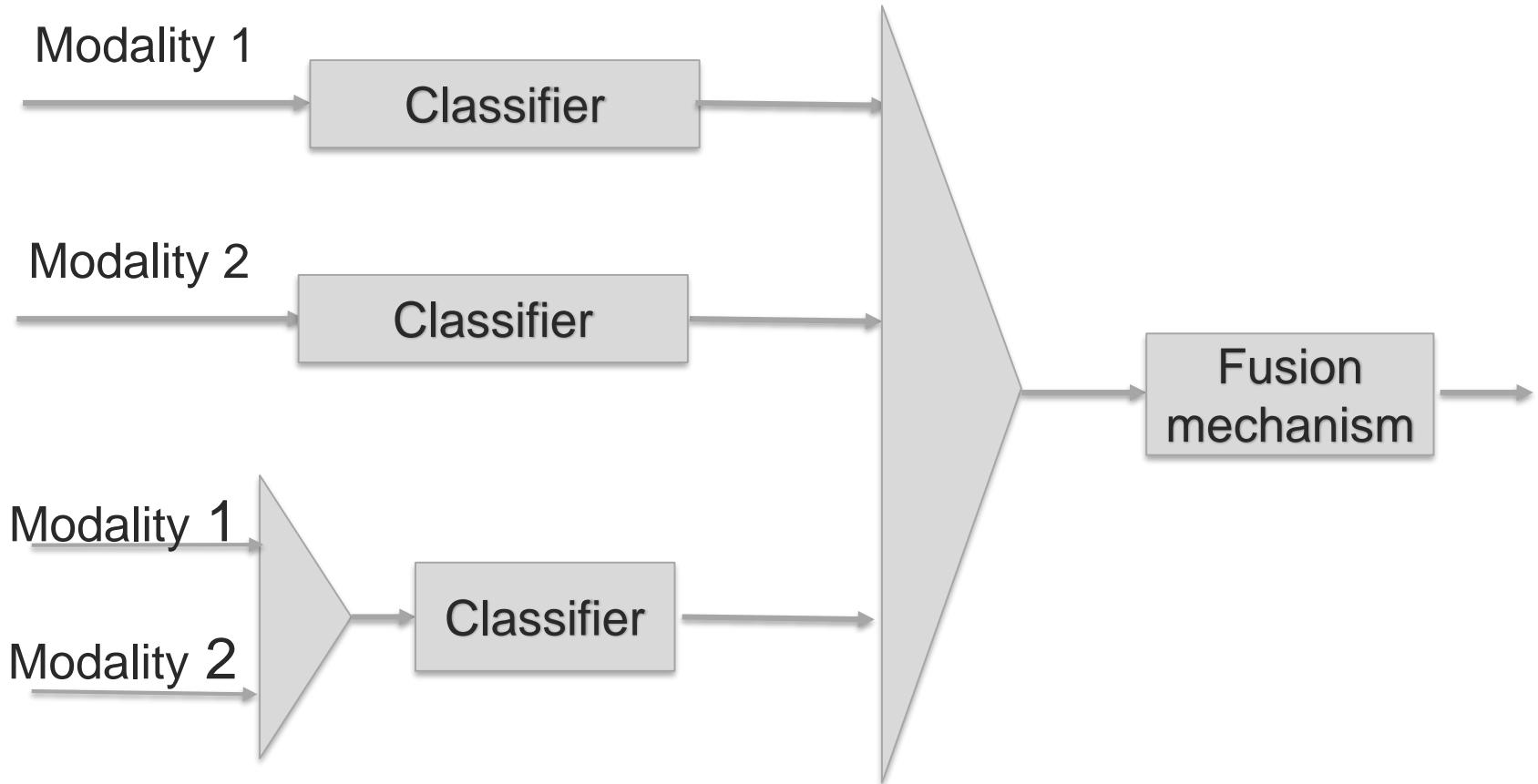
Model free approaches – late fusion



- Train a unimodal predictor and a multimodal fusion one
- Requires multiple training stages
- Do not model low level interactions between modalities
- Fusion mechanism can be voting, weighted sum or an ML approach



Model free approaches – hybrid fusion

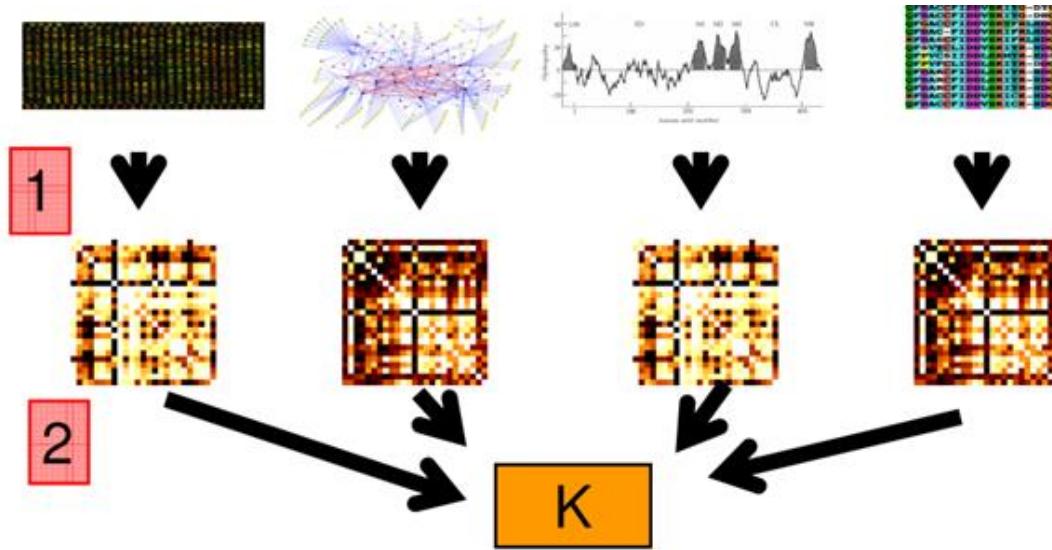


- Combine benefits of both early and late fusion mechanisms



Multiple Kernel Learning

- Pick a family of kernels for each modality and learn which kernels are important for the classification case
- Generalizes the idea of Support Vector Machines
- Works as well for unimodal and multimodal data, very little adaptation is needed



[Lanckriet 2004]



Multimodal Fusion for Sequential Data

Modality-*private* structure

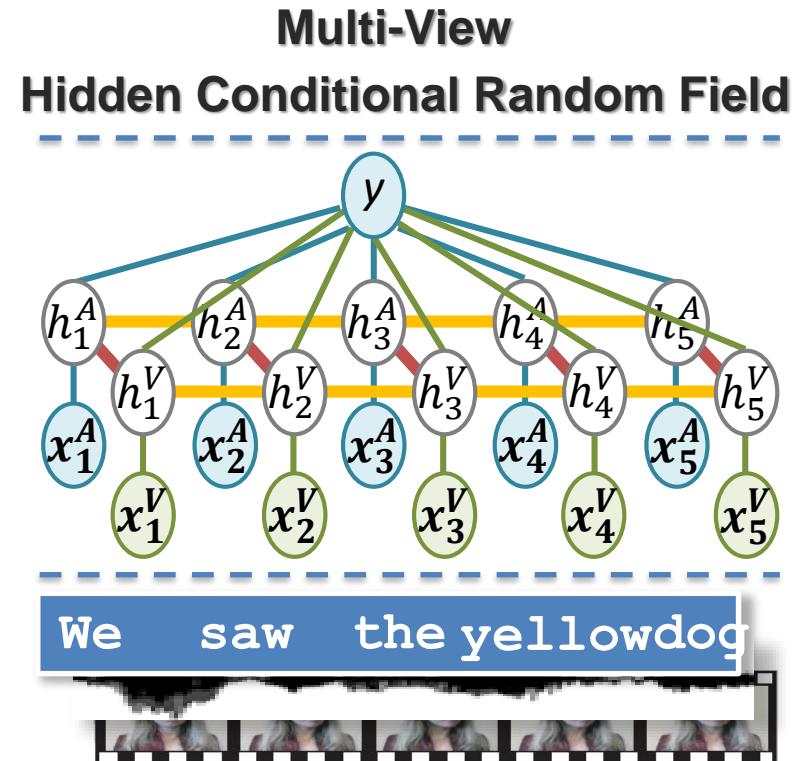
- Internal grouping of observations

Modality-*shared* structure

- Interaction and synchrony

$$p(y|x^A, x^V; \theta) = \sum_{h^A, h^V} p(y, h^A, h^V | x^A, x^V; \theta)$$

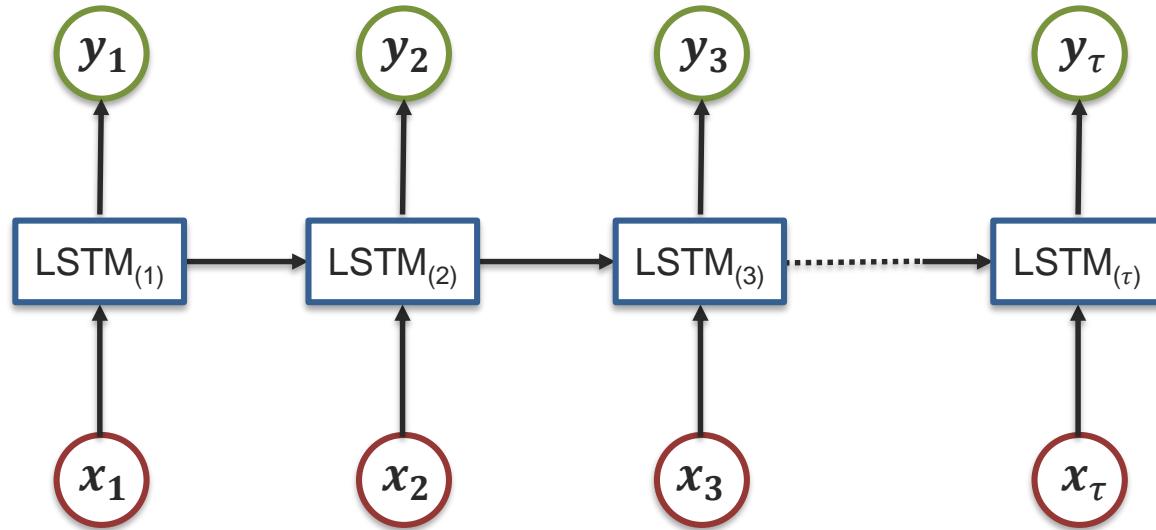
- Approximate inference using loopy-belief



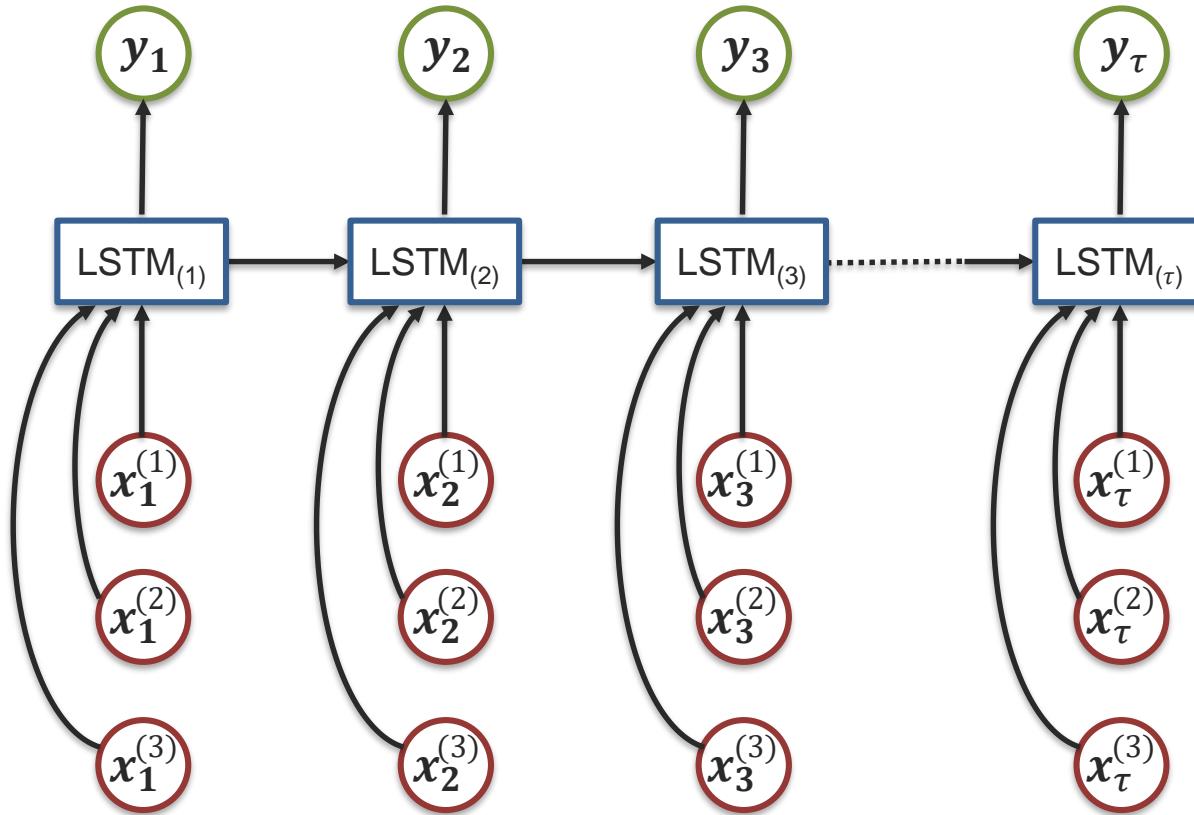
[Song, Morency and Davis, CVPR 2012]



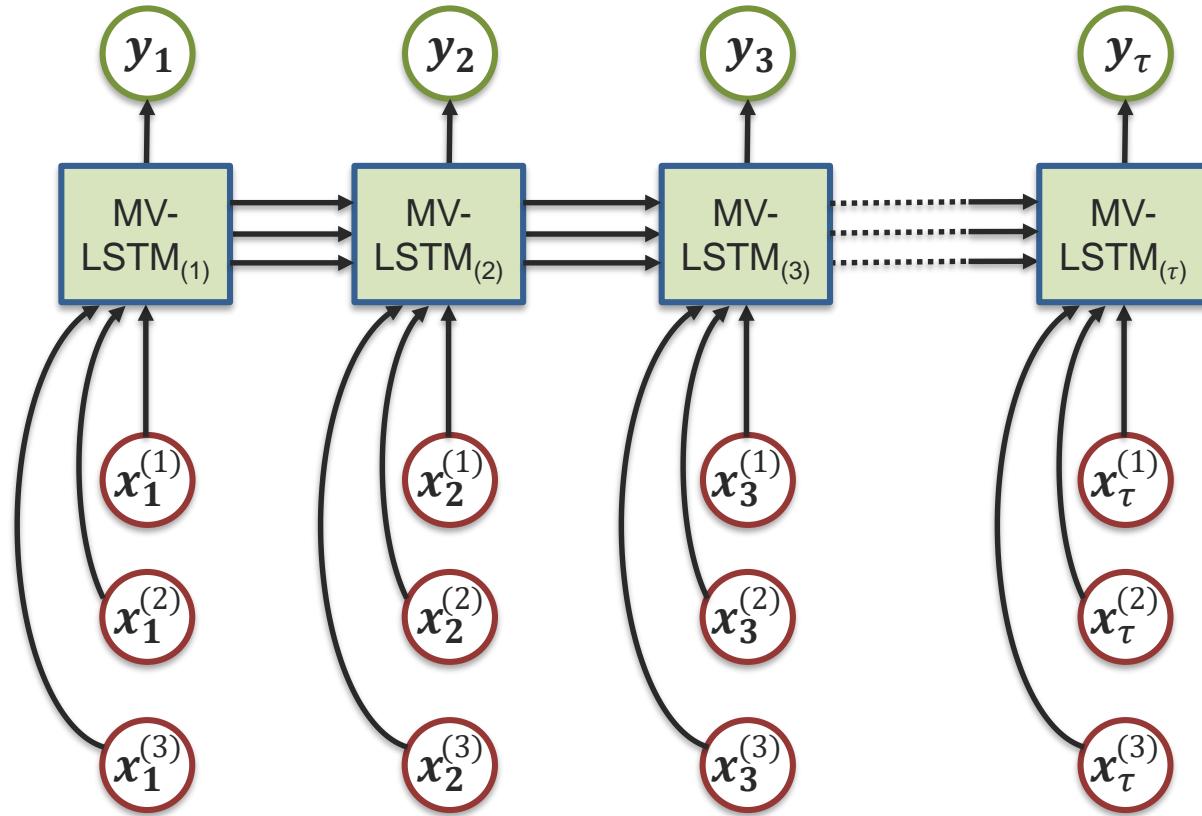
Sequence Modeling with LSTM



Multimodal Sequence Modeling – Early Fusion



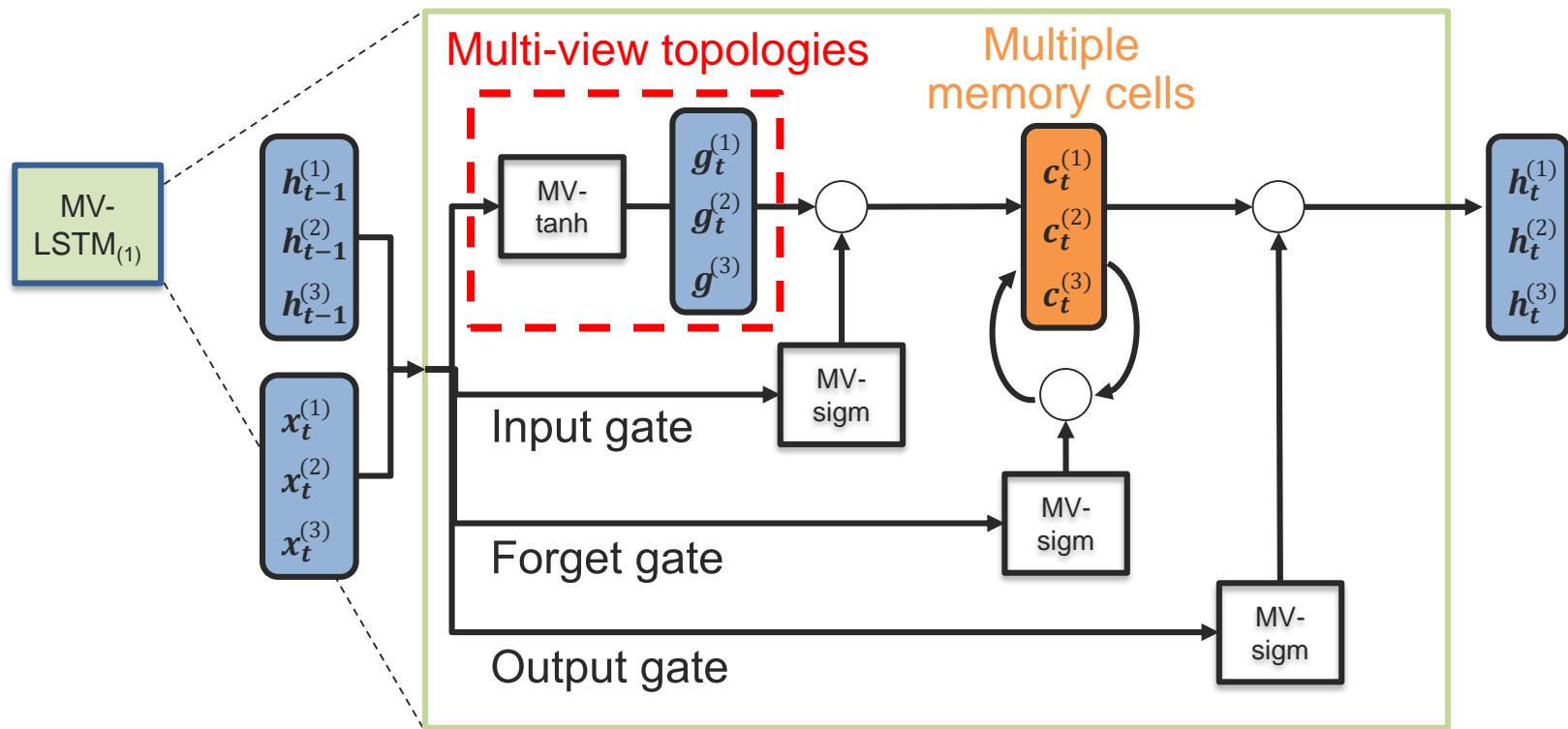
Multi-View Long Short-Term Memory (MV-LSTM)



[Shyam, Morency, et al. Extending Long Short-Term Memory for Multi-View Structured Learning, ECCV, 2016]



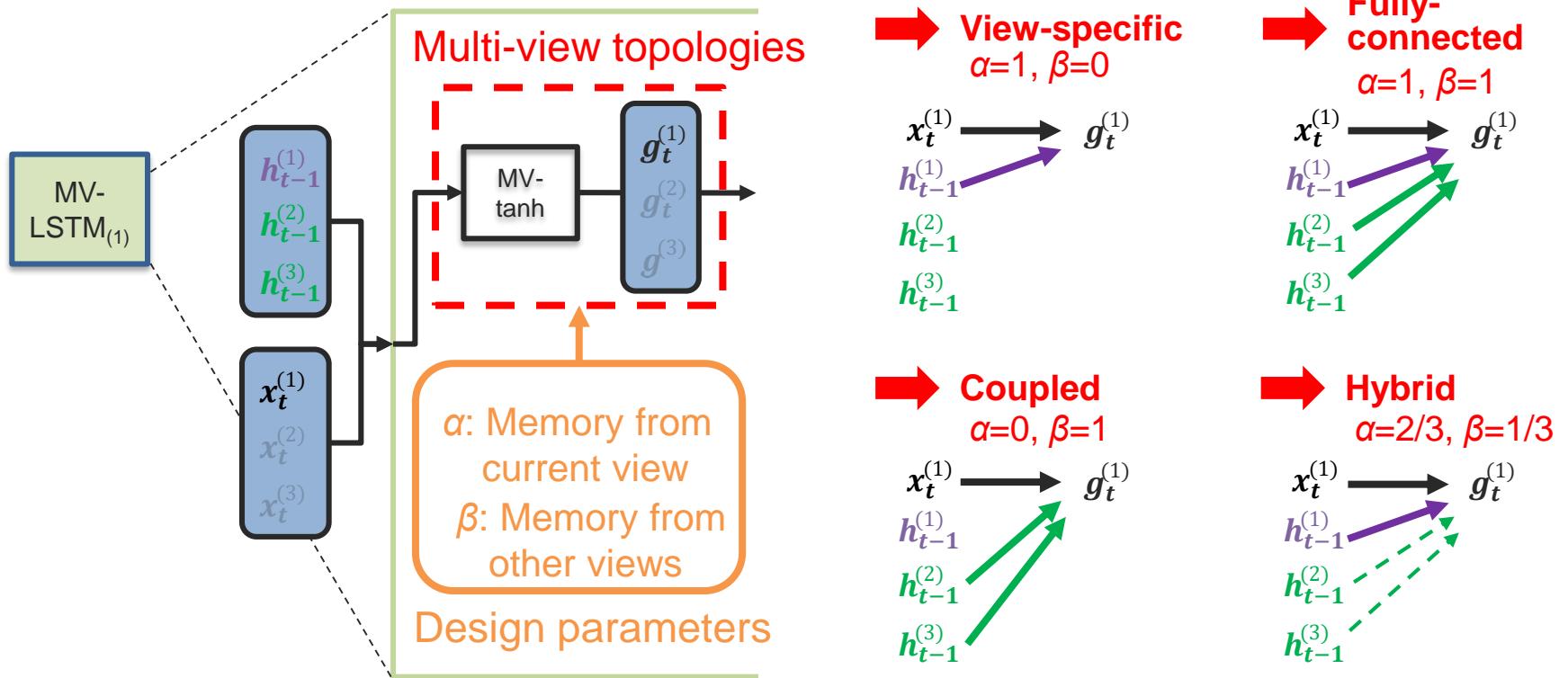
Multi-View Long Short-Term Memory



[Shyam, Morency, et al. Extending Long Short-Term Memory for Multi-View Structured Learning, ECCV, 2016]



Topologies for Multi-View LSTM



[Shyam, Morency, et al. Extending Long Short-Term Memory for Multi-View Structured Learning, ECCV, 2016]



Multi-View Long Short-Term Memory (MV-LSTM)

Multimodal prediction of children engagement

Class labels	Model	Precision	Recall	F1
Easy to engage	LSTM (Early fusion)	0.75	0.81	0.78
	MV-LSTM Full	0.81	0.81	0.81
	MV-LSTM Coupled	0.79	0.81	0.80
	MV-LSTM Hybrid	0.80	0.86	0.83
Difficult to engage	LSTM (Early fusion)	0.63	0.55	0.59
	MV-LSTM Full	0.68	0.68	0.68
	MV-LSTM Coupled	0.67	0.64	0.65
	MV-LSTM Hybrid	0.74	0.64	0.68

[Shyam, Morency, et al. Extending Long Short-Term Memory for Multi-View Structured Learning, ECCV, 2016]

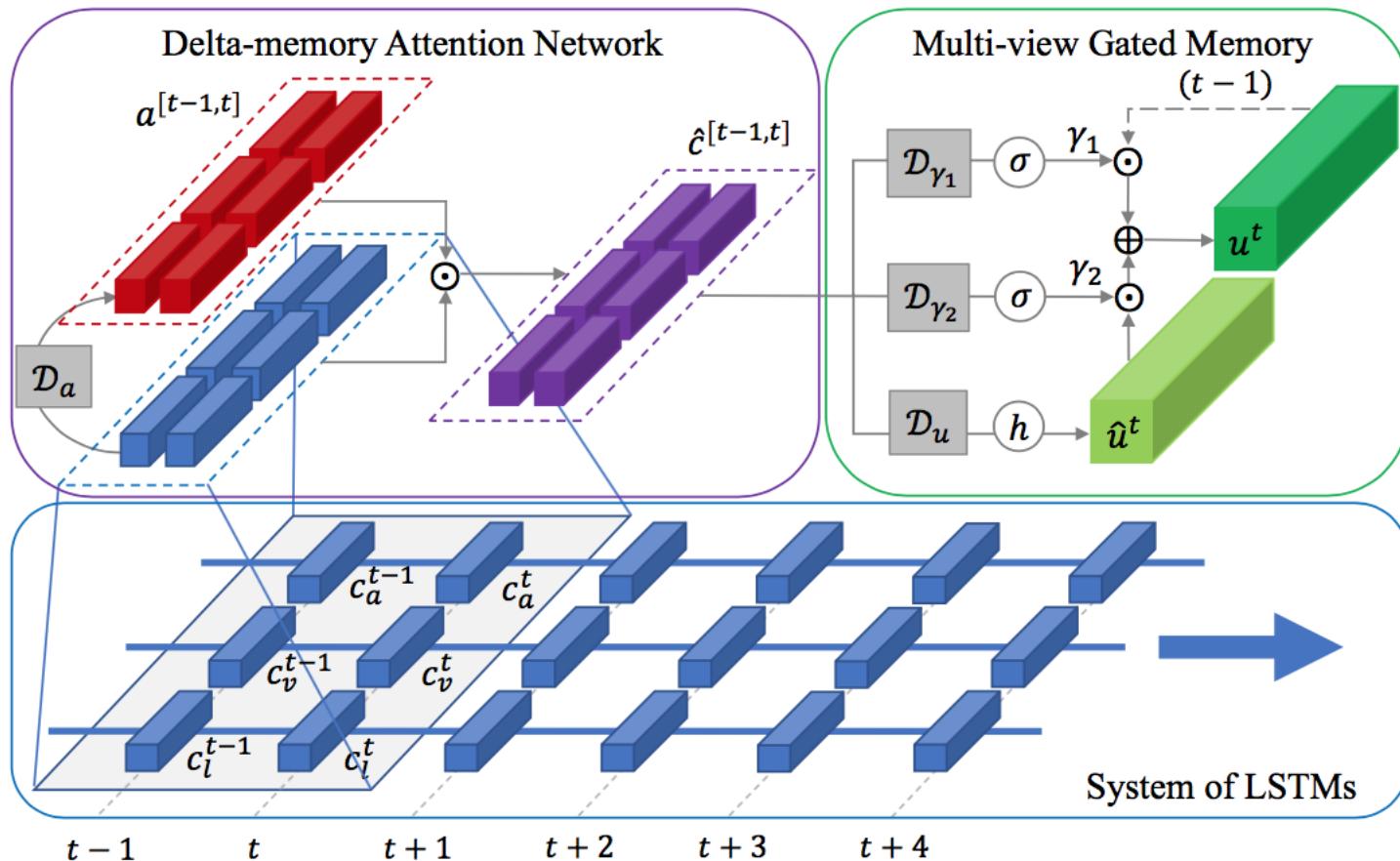


Memory Based

- A memory accumulates multimodal information over time.
- From the representations throughout a source network.
- No need to modify the structure of the source network, only attached the memory.



Memory Based



[Zadeh et al., Memory Fusion Network for Multi-view Sequential Learning, AAAI 2018]

Multimodal Machine Learning

Representation

Alignment

Fusion

Translation

Co-Learning

Multimodal Machine Learning: A Survey and Taxonomy

By Tadas Baltrusaitis, Chaitanya Ahuja,
and Louis-Philippe Morency

<https://arxiv.org/abs/1705.09406>

- 5 core challenges
- 37 taxonomic classes
- 253 referenced citations

