

Deep Learning for AI

from Machine Perception to Machine Cognition

Li Deng

*Chief Scientist of AI,
Microsoft Applications/Services Group (ASG) &
MSR Deep Learning Technology Center (DLTC)*

A Plenary Presentation at IEEE-ICASSP, March 24, 2016

Thanks go to many colleagues at DLTC & MSR, collaborating universities,
and at Microsoft's engineering groups (ASG+)



WIKIPEDIA
The Free Encyclopedia

Deep learning

From Wikipedia, the free encyclopedia

Definition

Deep learning is a class of machine learning algorithms that[\[1\]](#)(pp199–200)

- use a cascade of **many layers of nonlinear processing** .
- are part of the broader machine learning field of learning representations of data facilitating **end-to-end optimization**.
- learn multiple levels of representations that correspond to **hierarchies of concept abstraction**
- ..., ...



WIKIPEDIA
The Free Encyclopedia

Artificial intelligence

From Wikipedia, the free encyclopedia

Artificial intelligence (AI) is the intelligence exhibited by machines or software. It is also the name of the academic field of study on how to create computers and computer software that are capable of intelligent behavior.

Artificial general intelligence

From Wikipedia, the free encyclopedia

Artificial general intelligence (AGI) is the intelligence of a (hypothetical) machine that could successfully perform any intellectual task that a human being can. It is a primary goal of artificial intelligence research and an important topic for science fiction writers and futurists. Artificial general intelligence is also referred to as "strong AI"...

AI/(A)GI & Deep Learning: the main thesis

AI/GI = machine perception (speech, image, video, gesture, touch...)

+ machine cognition (natural language, *reasoning, attention, memory/learning, knowledge, decision making, action, interaction/conversation, ...*)

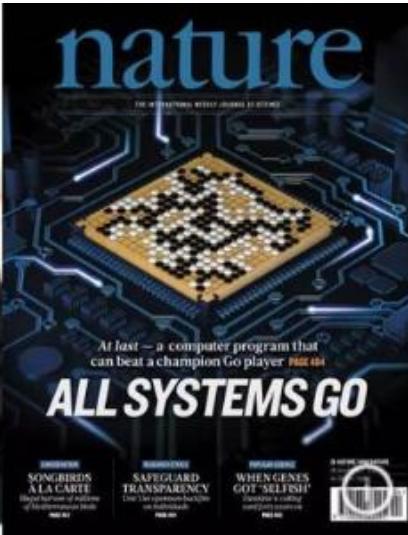
GI: AI that is flexible, general, adaptive, learning from 1st principles

Deep Learning + Reinforcement/Unsupervised Learning
→AI/GI

AI/GI & Deep Learning: how AlphaGo fits

AI/GI = machine perception (**speech, image, video, gesture, touch...**)

+ machine cognition (**natural language, reasoning, attention, memory/learning, knowledge, decision making, action, interaction**)



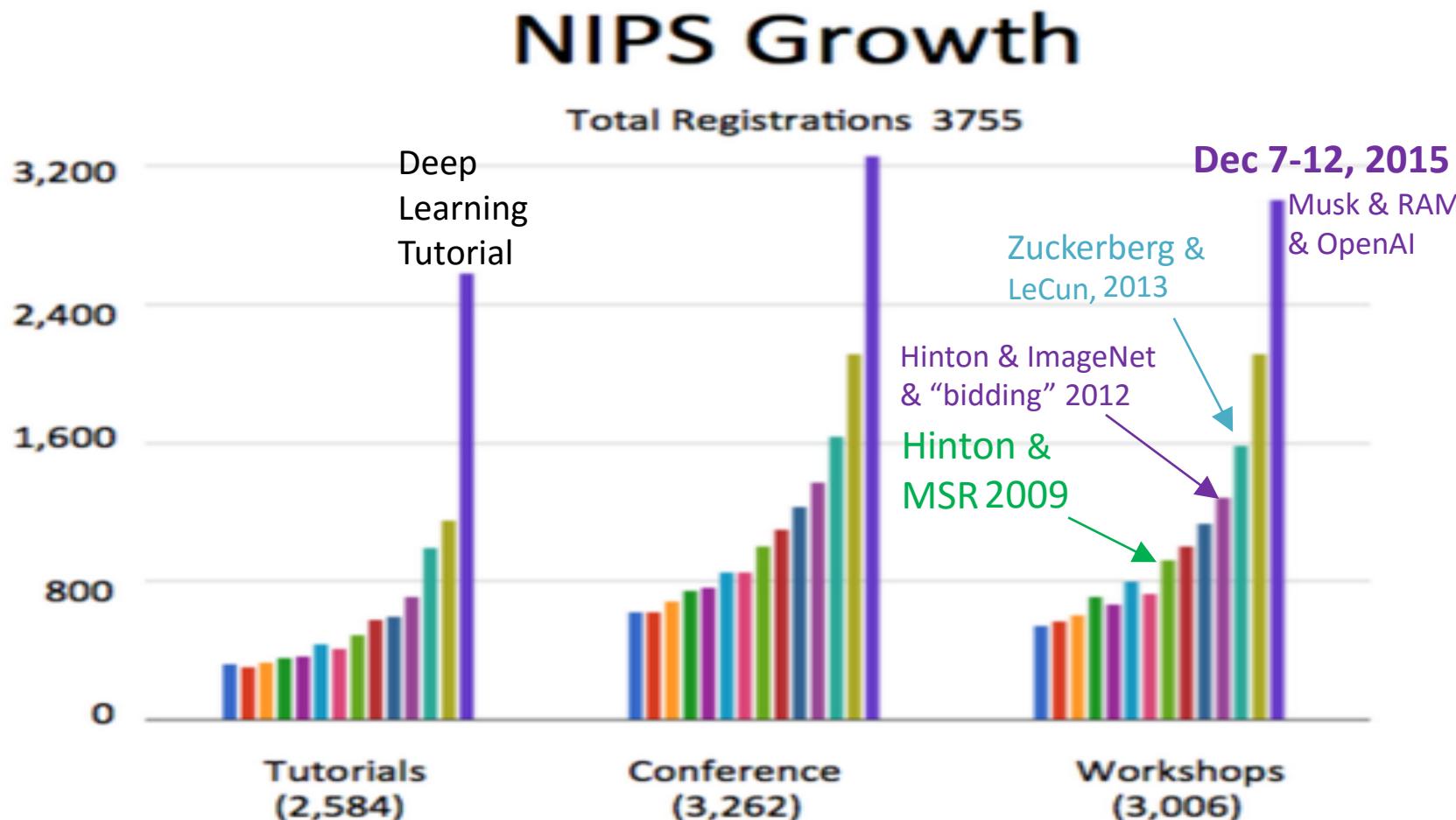
AGI: AI that follows the principles of biological systems
Deep Learning
→AI/AGI

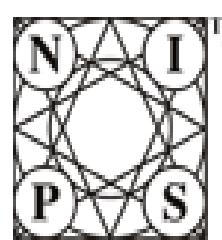
creative, learning from 1st supervised Learning

Outline

- Deep learning for machine **perception**
 - Speech
 - Image
- Deep learning for machine **cognition**
 - Semantic modeling
 - Natural language
 - Multimodality
 - Reasoning, attention, memory (RAM)
 - Knowledge representation/management/exploitation
 - Optimal decision making (by deep reinforcement learning)
- Three hot areas/challenges of deep learning & AI research

Deep learning Research: centered at **NIPS** (Neural Information Processing Systems)





[NIPS Home](#)

[Overview](#)

[Conference Videos](#)

[Workshop Videos](#)

[Program Highlights](#)

[Tutorials](#)

[Conference Sessions](#)

[Workshops](#)

[Publication Models](#)

[Demonstrations](#)

[Mini Symposia](#)

[Accepted Papers](#)

[Dates](#)

[Committees](#)

Li Deng, Dong Yu, Geoffrey Hinton

Microsoft Research; Microsoft Research; University of Toronto

Deep Learning for Speech Recognition and Related Applications

7:30am - 6:30pm Saturday, **December 12, 2009**

Location: Hilton: Cheakamus

Abstract: Over the past 25 years or so, speech recognition technology has been dominated by a "shallow" architecture --- hidden Markov models (HMMs). Significant technological success has been achieved using complex and carefully engineered HMMs. The next generation of the technology requires solutions to remaining challenges under diversified deployment environments. These challenges, not fully addressed in the past, arise from the many types of variability present in the speech generation process. Overcoming these challenges is likely to require "deep" architectures, with efficient learning algorithms. For speech recognition and related sequential recognition applications, some attempts have been made in the past to develop computational architectures that are "deeper" than conventional HMMs, such as

The Universal
Translator ..comes true!



Deep learning
technology enabled
speech-to-speech
translation

The New York Times

Scientists See Promise in Deep-Learning Programs

John Markoff

November 23, 2012

Tianjin, China, October, 25, 2012



A voice recognition program translated a speech given by Richard F. Rashid, Microsoft's top scientist, into Mandarin Chinese.



Microsoft Research



Investigation of full-sequence training of DBNs for speech recognition., Interspeech, Sept 2010

Binary coding of speech spectrograms using a deep auto-encoder, Interspeech, Sept 2010

Roles of Pre-Training & Fine-Tuning in CD-DBN-HMMs for Real-World ASR, NIPS, Dec. 2010

Large Vocabulary Continuous Speech Recognition With CD-DNN-HMMS, ICASSP, April 2011

Conversational Speech Transcription Using Context-Dependent DNN, Interspeech, Aug. 2011



Making deep belief networks effective for LVCSR, ASRU, Dec. 2011

Microsoft

Application of Pretrained DNNs to Large Vocabulary Speech Recognition., ICASSP, 2012

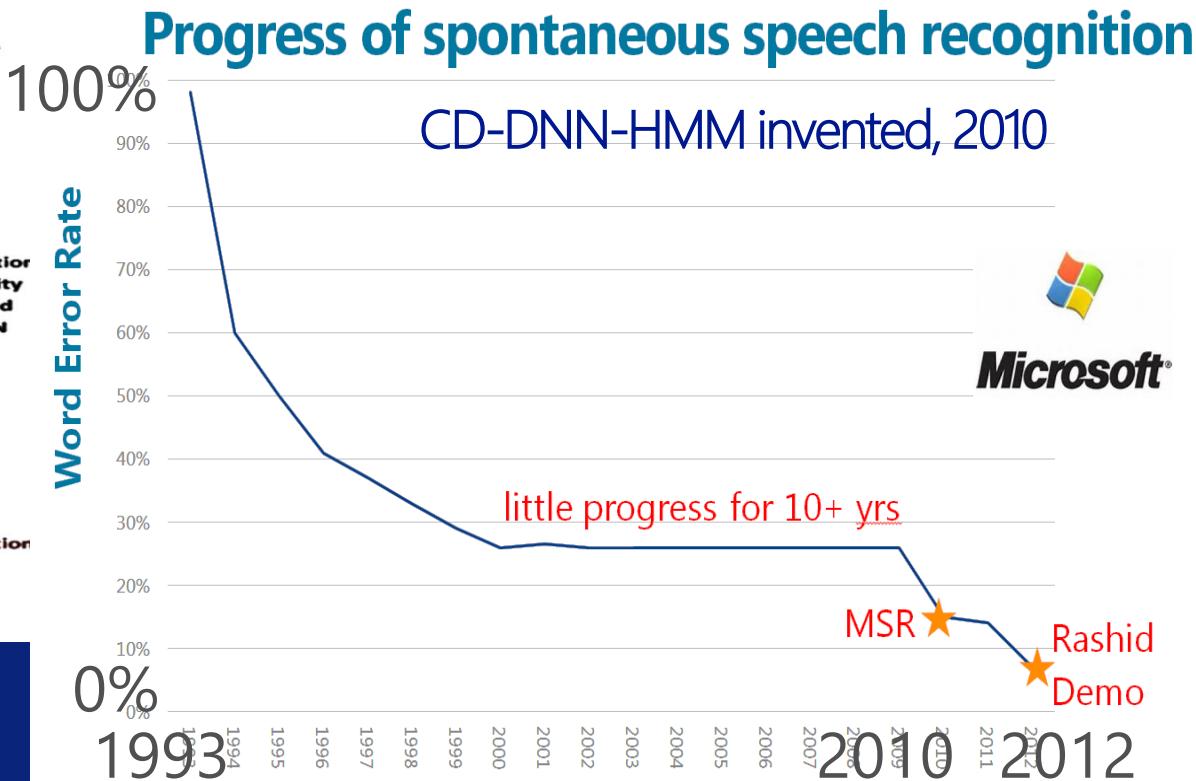
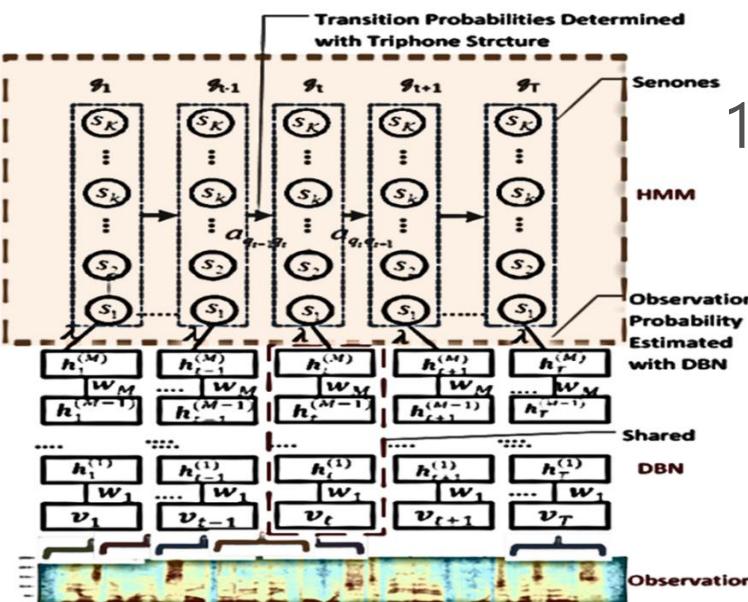
Google

【胡郁】讯飞超脑 2.0 是怎样炼成的？2011, 2015

Later years with rapid progress,

Baidu Research

科大讯飞
iFLYTEK

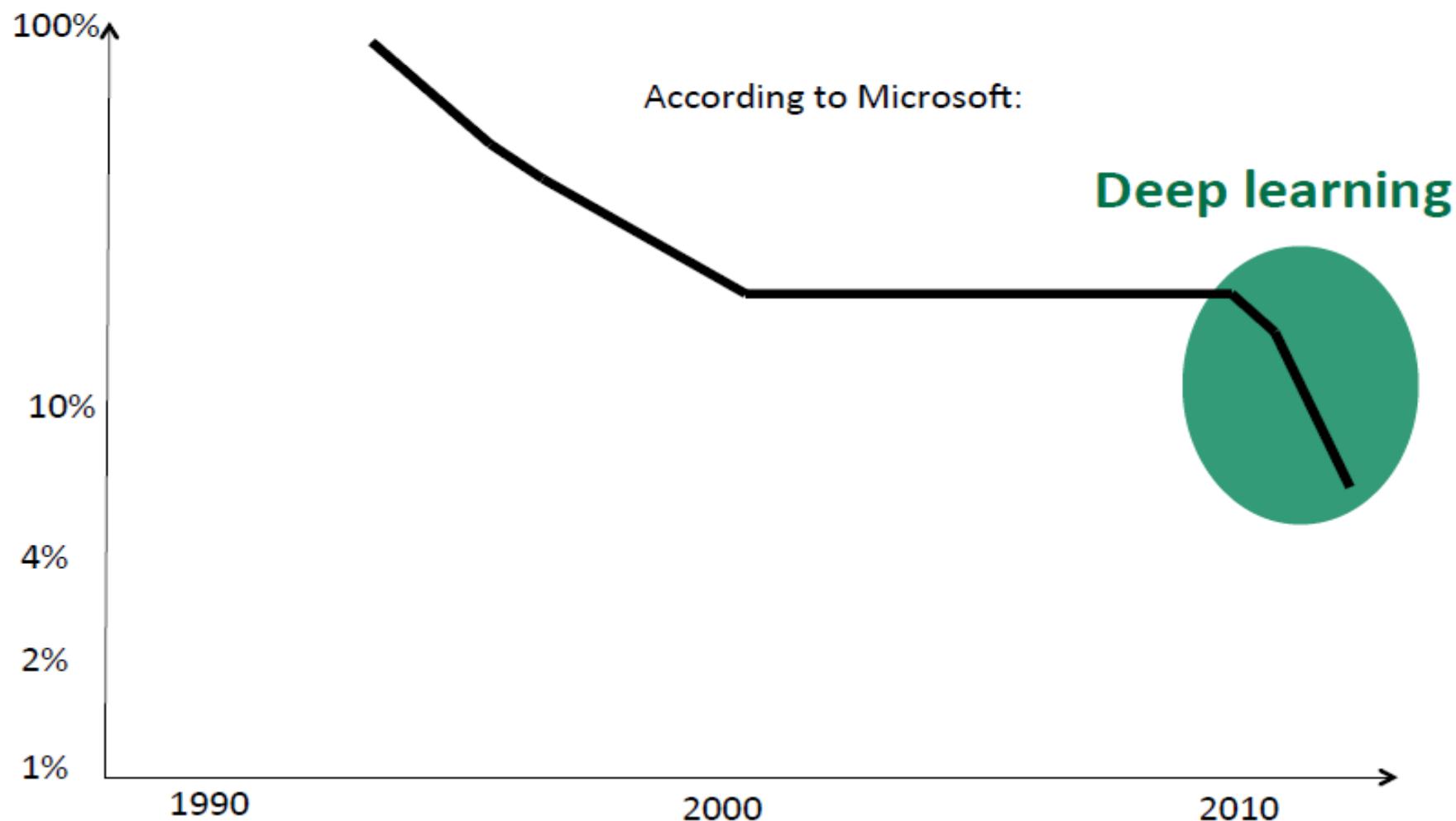


Microsoft



Microsoft Research

2010-2012: Breakthrough in speech
recognition → in Androids by 2012

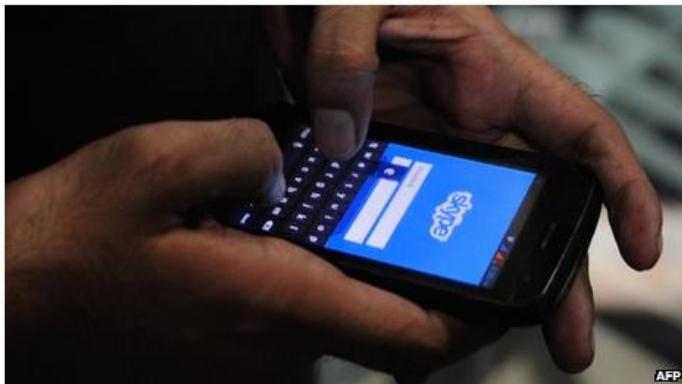


Across-the-Board Deployment of DNN in Speech Industry

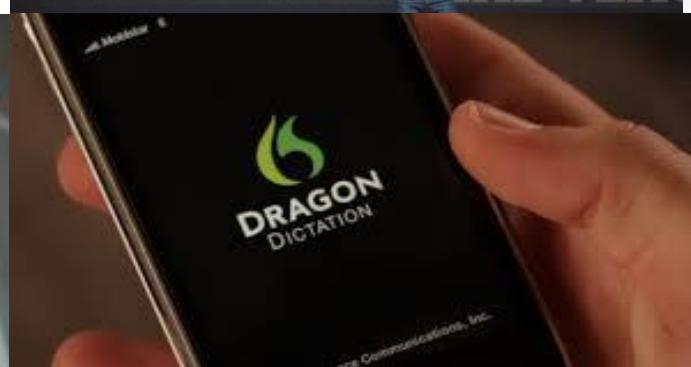
(+ in university labs & DARPA programs)

(2012-2014)

Skype to get 'real-time' translator



Analysts say the translation feature could have wide ranging applications



Enabling Cross-Lingual Conversations in Real Time

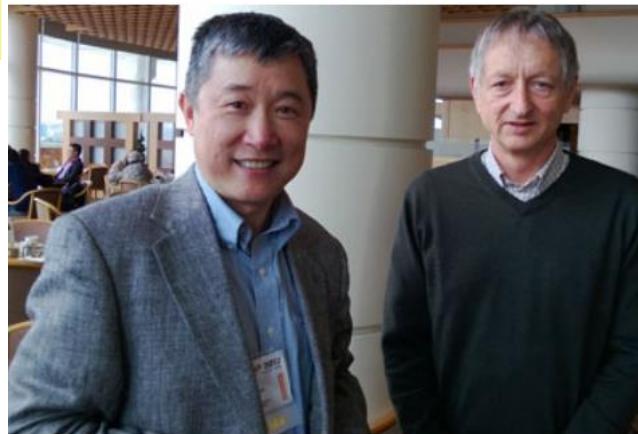
Microsoft Research

May 27, 2014 5:58 PM PT

View milestones
on the path to
Skype Translator
[#speech2speech](#)



HOW SKYPE USED AI TO BUILD ITS AMAZING NEW LANGUAGE TRANSLATOR



Taking a cue from science fiction,
Microsoft demos 'universal translator'

By Jacopo Prisco, for CNN
Updated 12:35 PM ET, Thu October 16, 2014



In the academic world



Deep learning

Yann LeCun, Yoshua Bengio & Geoffrey Hinton

Affiliations | Corresponding author

Nature 521, 436–444 (28 May 2015) | doi:10.1038/nature14539

Received 25 February 2015 | Accepted 01 May 2015 | Published online



[Geoffrey Hinton, Li Deng, Dong Yu, George E. Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N. Sainath, and Brian Kingsbury]

“This joint paper (2012) → from the major speech recognition laboratories details the first major industrial application of deep learning.”

Deep Neural Networks for Acoustic Modeling in Speech Recognition

[The shared views of four research groups]



State-of-the-Art Speech Recognition Today

(& tomorrow --- roles of unsupervised learning)

ASR: Neural Network Architectures at Google

Single Channel:

LSTM acoustic model trained with connectionist temporal classification (CTC)

Results on a 2,000-hr English Voice Search task show an 11% relative improvement

Papers: [H. Sak et al - ICASSP 2015, Interspeech 2015, A. Senior et al - ASRU 2015]

Model	WER
LSTM w/ conventional modeling	14.0
LSTM w/ CTC	12.9%

(Sainath, Senior, Sak, Vinyals)

Multi-Channel:

Multi-channel raw-waveform input for each channel

Initial network layers factored to do spatial and spectral filtering

Output passed to a CLDNN acoustic model, entire network trained jointly

Results on a 2,000-hr English Voice Search task show more than 10% relative improvement

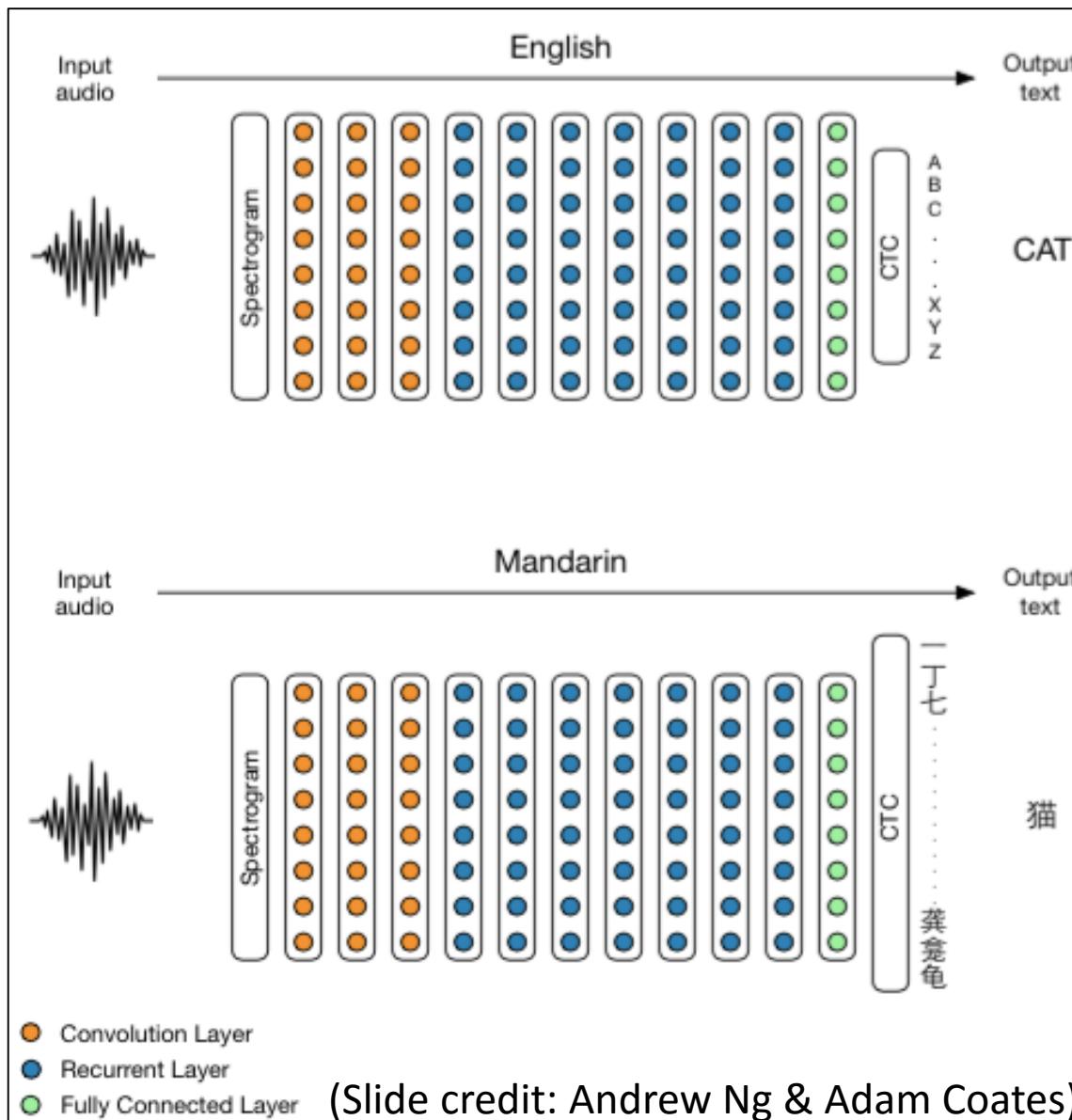
Papers: [T. N. Sainath et al - ASRU 2015, ICASSP 2016]

Model	WER
raw-waveform, 1ch	19.2
delay+sum, 8 channel	18.7
MVDR, 8 channel	18.8
factored raw-waveform, 2ch	17.1

(Slide credit: Tara Sainath & Andrew Senior)

Baidu's Deep Speech 2

End-to-End DL System for Mandarin and English



Paper: bit.ly/deepspeech2

- Human-level Mandarin recognition on short queries:
 - DeepSpeech: 3.7% - 5.7% CER
 - Humans: 4% - 9.7% CER
- Trained on 12,000 hours of conversational, read, mixed speech.
- 9 layer RNN with CTC cost:
 - 2D invariant convolution
 - 7 recurrent layers
 - Fully connected output
- Trained with SGD on heavily-optimized HPC system. “SortaGrad” curriculum learning.
- “Batch Dispatch” framework for low-latency production deployment.



Learning transition probabilities in DNN-HMM ASR

DNN outputs include not only state posterior outputs but also HMM transition probabilities

Real-time reduction of 16%
WER reduction of 10%

State posteriors Transition probs

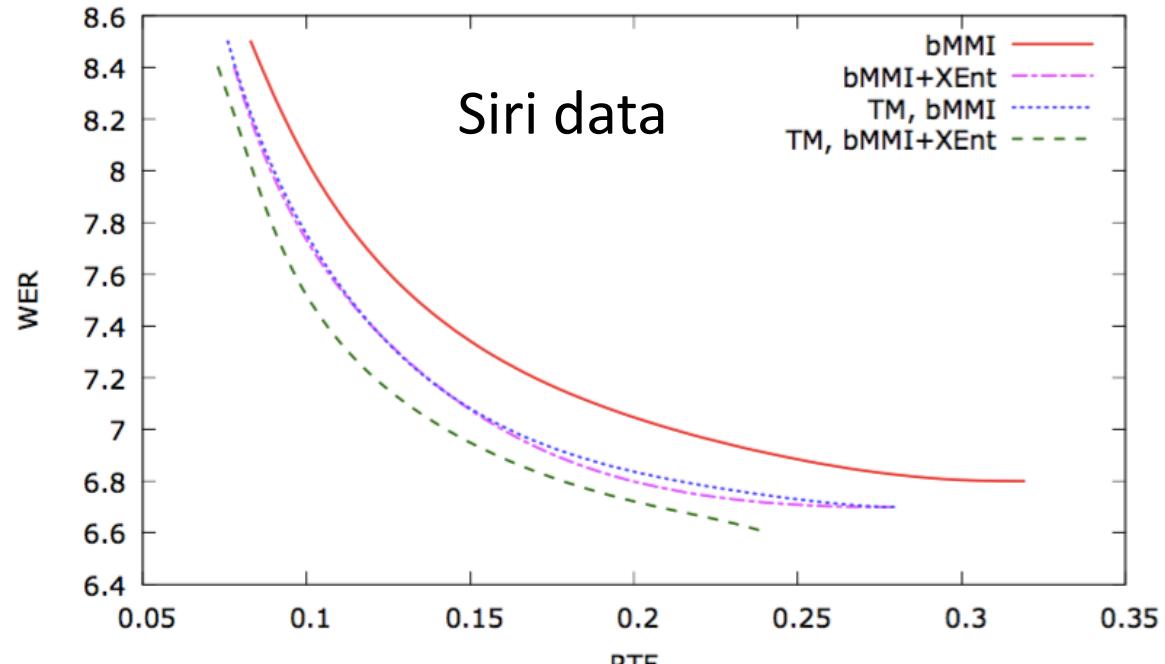
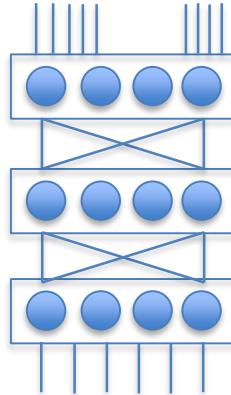


Figure 2: WER vs. RTF (dev set)

Matthias Paulik, "Improvements to the Pruning Behavior of DNN Acoustic Models". Interspeech 2015

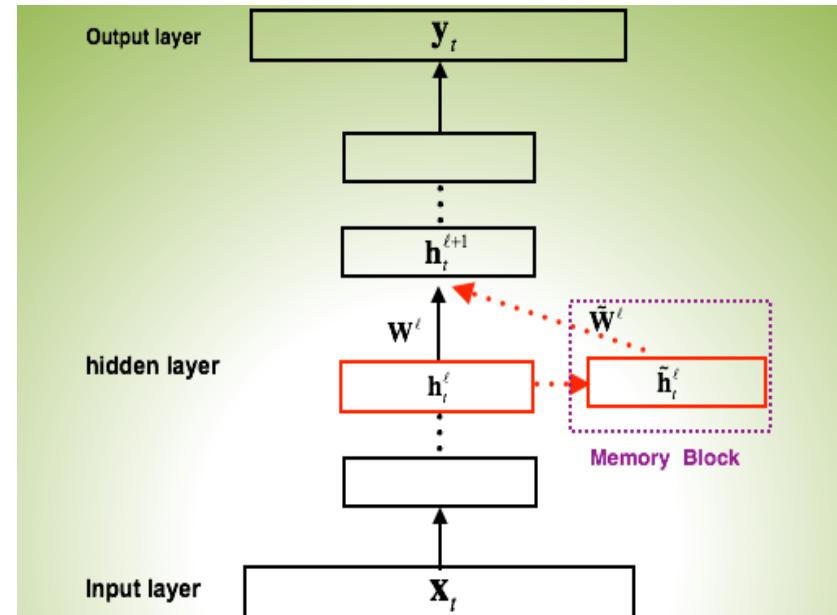
(Slide: Alex Acero)

FSMN-based LVCSR System

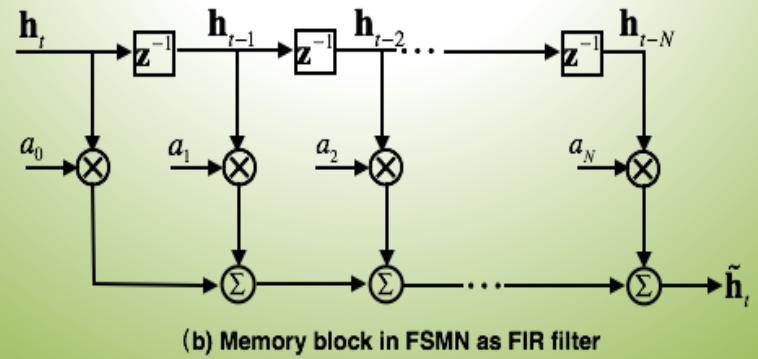


- ❑ Feed-forward Sequential Memory Network(FSMN)
- ❑ Results on 10,000 hours Mandarin short message dictation task
 - 8 hidden layers
 - Memory block with -/+ 15 frames
 - CTC training criteria
- ❑ Comparable results to DBLSTM with smaller model size
- ❑ Training costs only 1 day using 16 GPUs and ASGD algorithm

Model	#Para.(M)	CER (%)
ReLU DNN	40	6.40
LSTM	27.5	5.25
BLSTM	45	4.67
FSMN	19.8	4.61



(a) Feedforward sequential memory neural network (FSMN)



(b) Memory block in F FSMN as FIR filter

Shiliang Zhang, Cong Liu, Hui Jiang, Si Wei, Lirong Dai, Yu Hu. “Feedforward Sequential Memory Networks: A New Structure to Learn Long-term Dependency”. [arXiv:1512.08031](https://arxiv.org/abs/1512.08031), 2015.

(slide credit: Cong Liu & Yu Hu)

English Conversational Telephone Speech Recognition*

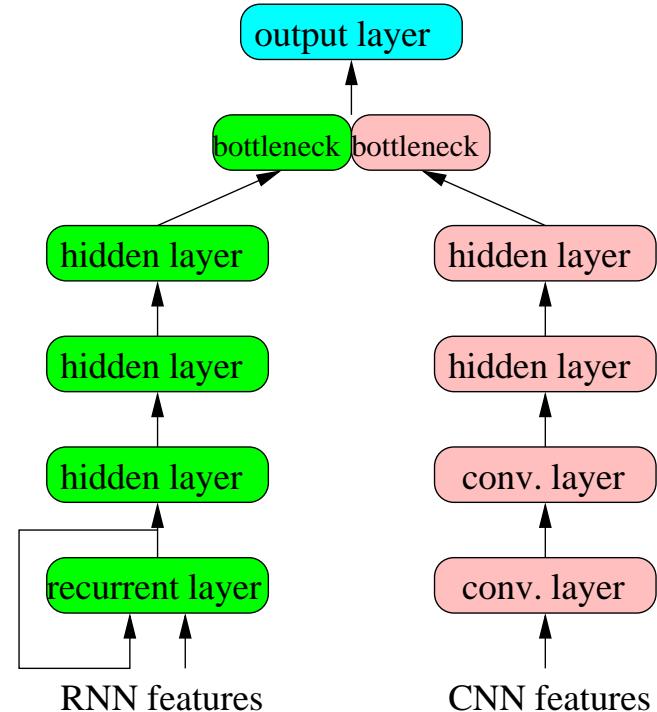


Key ingredients:

- Joint RNN/CNN acoustic model trained on 2000 hours of publicly available audio
- Maxout activations
- Exponential and NN language models

WER Results on Switchboard Hub5-2000:

Model	WER SWB	WER CH
CNN	10.4	17.9
RNN	9.9	16.3
Joint RNN/CNN	9.3	15.6
+ LM rescoring	8.0%	14.1



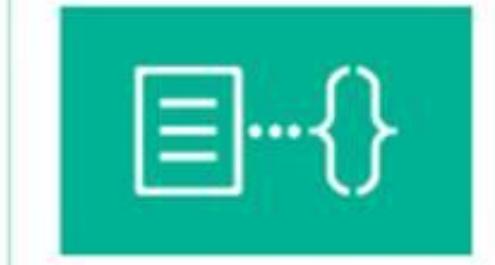
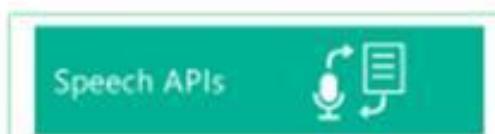
*Saon et al. "The IBM 2015 English Conversational Telephone Speech Recognition System", Interspeech 2015.



Microsoft

MICROSOFT PROJECT OXFORD SERVICES

PROJECT OXFORD: <http://www.projectoxford.ai>



Emotion APIs:
BETA

Understand your users with Emotion Recognition.

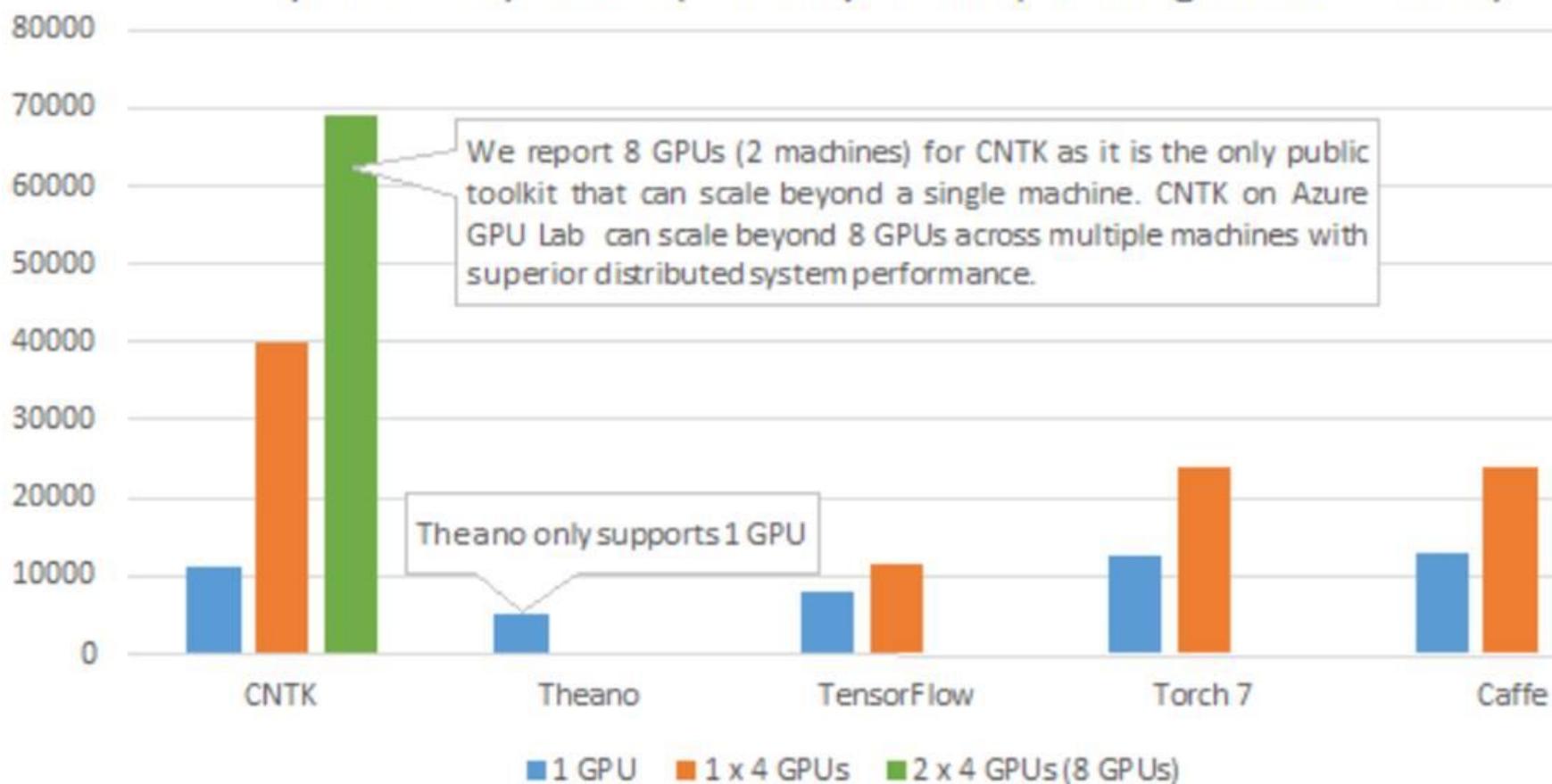
Spell Check APIs
BETA

Detect and correct common and uncommon spelling errors, via the Bing document index.

Language Understanding Intelligent Service (LUIS)
BETA

Understand natural language commands tailored to your application

Speed Comparison (Frames/Second, The Higher the Better)



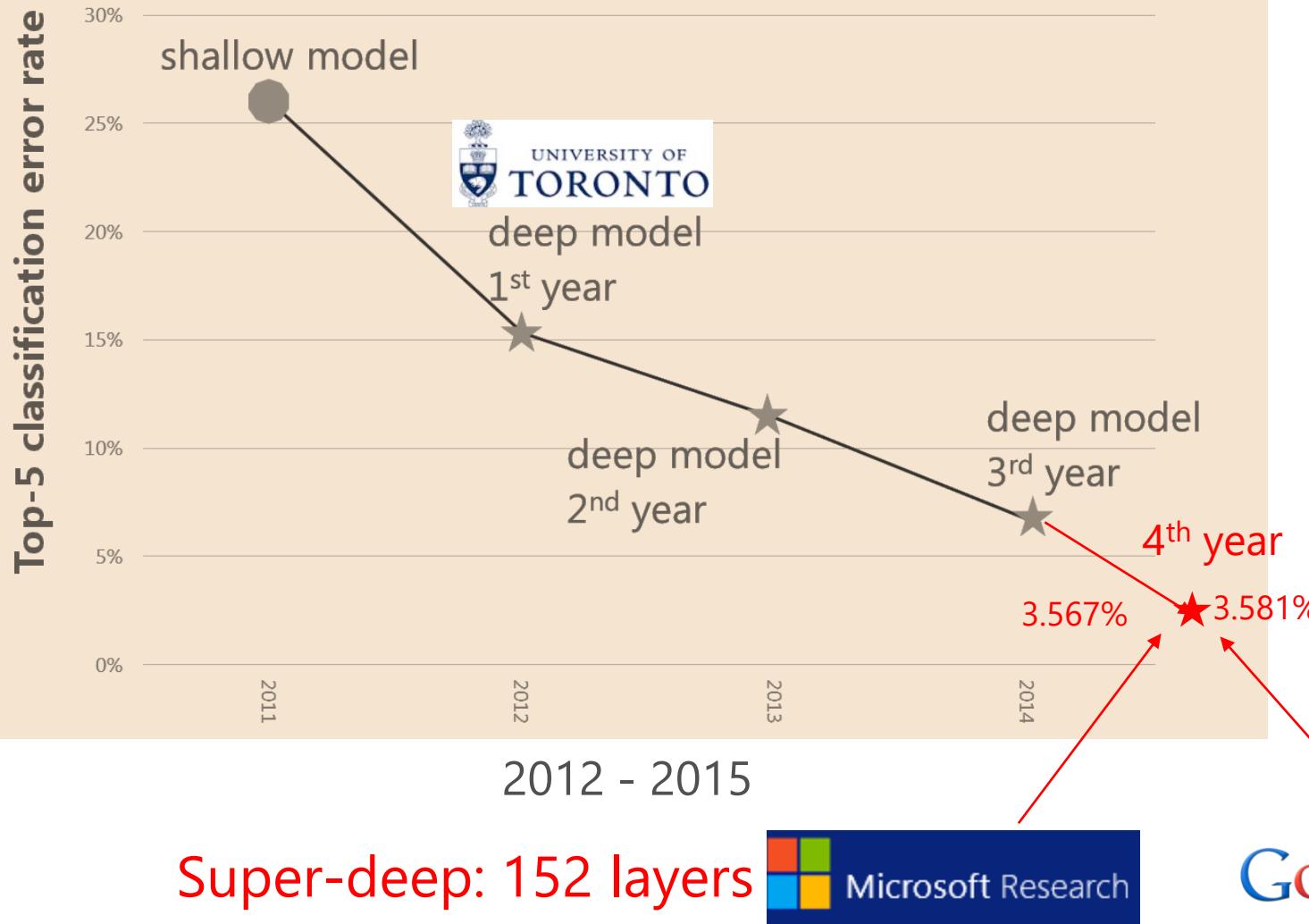
*Google updated that TensorFlow can now scale to support multiple machines recently; comparisons have not been made yet

- Recent Research at MS (ICASSP-2016):
 - “SCALABLE TRAINING OF DEEP LEARNING MACHINES BY INCREMENTAL BLOCK TRAINING WITH INTRA-BLOCK PARALLEL OPTIMIZATION AND BLOCKWISE MODEL-UPDATE FILTERING”
 - “HIGHWAY LSTM RNNs FOR DISTANCE SPEECH RECOGNITION”
 - “SELF-STABILIZED DEEP NEURAL NETWORKS”

Deep Learning also Shattered Image Recognition
(since 2012)

IMAGENET Competition

Progress of object recognition (1k ImageNet)



CADE METZ BUSINESS 01.14.16 7:00 AM

MICROSOFT NEURAL NET SHOWS DEEP LEARNING CAN GET WAY DEEPER

Microsoft beats Google, Intel, Tencent, and Qualcomm in image recognition competition

JORDAN NOVET DECEMBER 10, 2015 4:14 PM

TAGS: ARTIFICIAL INTELLIGENCE, DEEP LEARNING, IBM, IMAGE RECOGNITION, IMAGENET, MICROSOFT, MICROSOFT RESEARCH, NVIDIA, SOFTLAYER

iPod

(trademark) a pocket-sized device used to play music files

1283 pictures

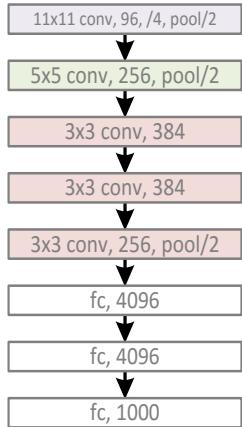
94.2% Popularity Percentile



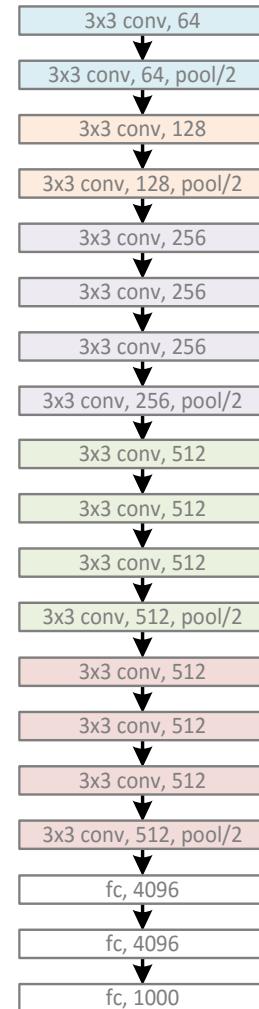
Microsoft Research

Depth is of crucial importance

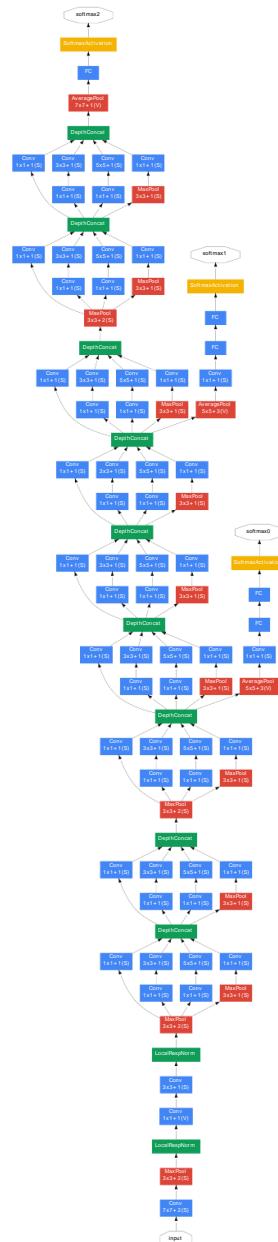
AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)



GoogleNet, 22 layers
(ILSVRC 2014)



ILSVRC (Large Scale Visual Recognition Challenge)

(slide credit: Jian Sun, MSR)

Depth is of crucial importance

AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)



ResNet, **152 layers**
(ILSVRC 2015)

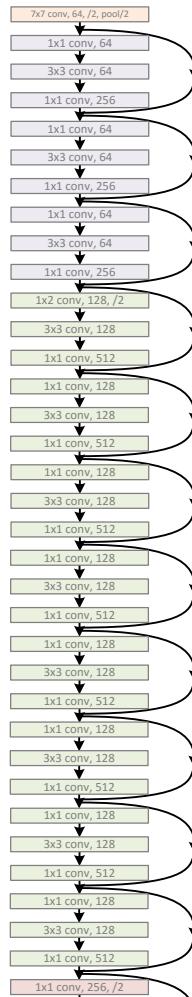


ILSVRC (Large Scale Visual Recognition Challenge)

(slide credit: Jian Sun, MSR)

Depth is of crucial importance

ResNet, 152 layers

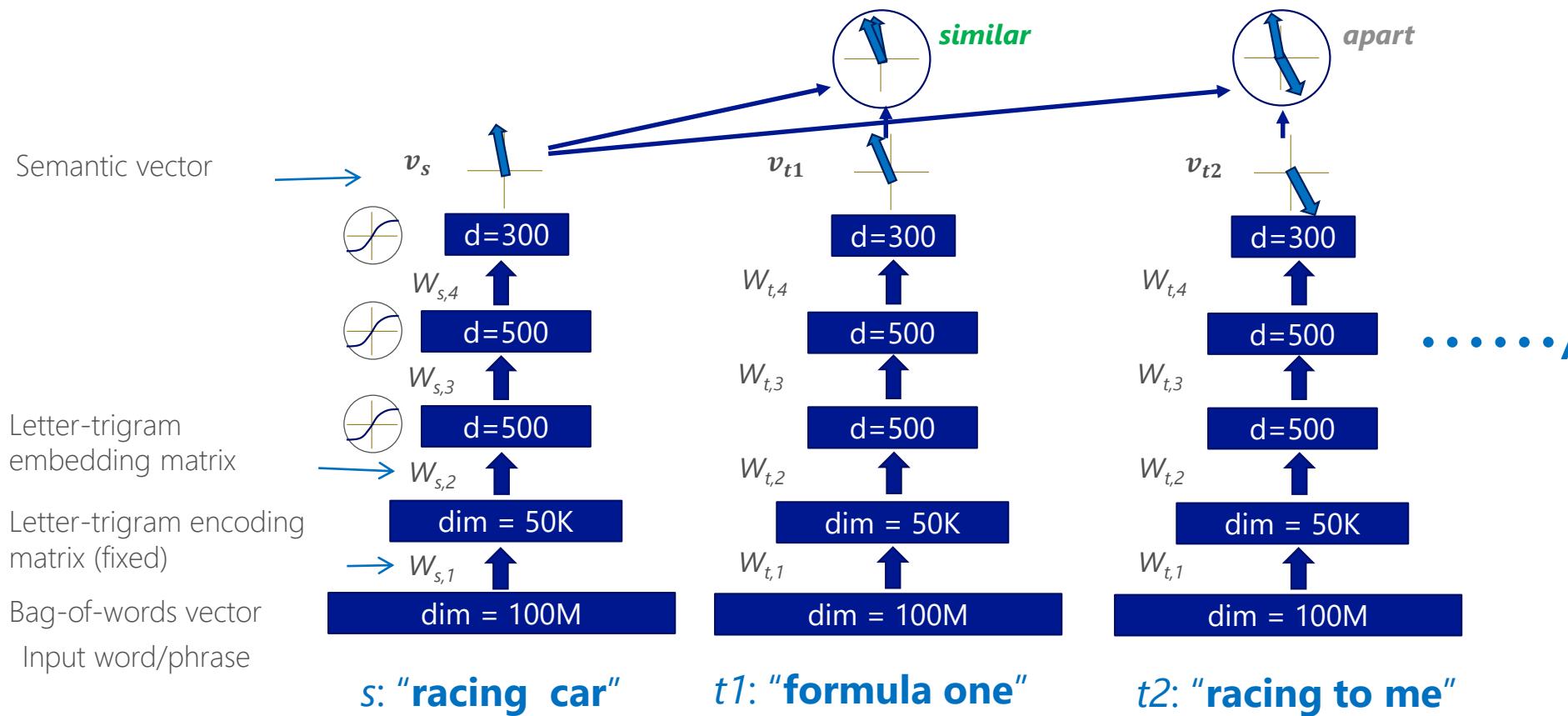


(slide credit: Jian Sun, MSR)

Outline

- Deep learning for machine perception
 - Speech
 - Image
- Deep learning for machine **cognition**
 - Semantic modeling
 - Natural language
 - Multimodality
 - Reasoning, attention, memory (RAM)
 - Knowledge representation/management/exploitation
 - Optimal decision making (by deep reinforcement learning)
- Three hot areas/challenges of deep learning & AI research

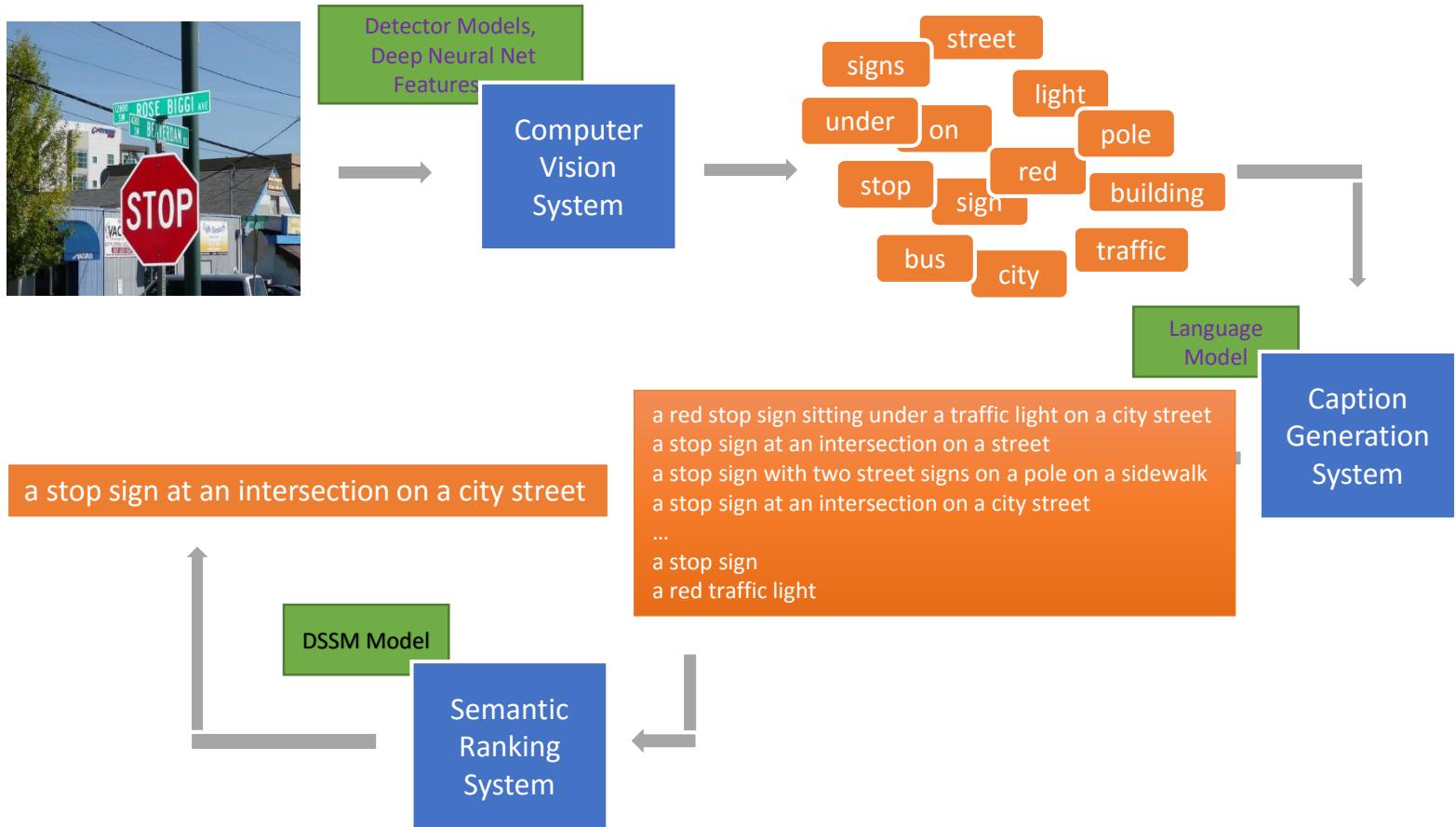
Deep Semantic Model for Symbol Embedding



Many applications of Deep Semantic Modeling: Learning semantic relationship between “Source” and “Target”

Tasks	<i>Source</i>	<i>Target</i>
Word semantic embedding	<i>context</i>	<i>word</i>
Web search	<i>search query</i>	<i>web documents</i>
Query intent detection	<i>Search query</i>	<i>Use intent</i>
Question answering	<i>pattern / mention (in NL)</i>	<i>relation / entity (in knowledge base)</i>
Machine translation	<i>sentence in language a</i>	<i>translated sentences in language b</i>
Query auto-suggestion	<i>Search query</i>	<i>Suggested query</i>
Query auto-completion	<i>Partial search query</i>	<i>Completed query</i>
Apps recommendation	<i>User profile</i>	<i>recommended Apps</i>
Distillation of survey feedbacks	<i>Feedbacks in text</i>	<i>Relevant feedbacks</i>
Automatic image captioning	<i>image</i>	<i>text caption</i>
Image retrieval	<i>text query</i>	<i>images</i>
Natural user interface	<i>command (text / speech / gesture)</i>	<i>actions</i>
Ads selection	<i>search query</i>	<i>ad keywords</i>
Ads click prediction	<i>search query</i>	<i>ad documents</i>
Email analysis: people prediction	<i>Email content</i>	<i>Recipients, senders</i>
Email search	<i>Search query</i>	<i>Email content</i>
Email decluttering	<i>Email contents</i>	<i>Email contents in similar threads</i>
Knowledge-base construction	<i>entity from source</i>	<i>entity fitting desired relationship</i>
Contextual entity search	<i>key phrase / context</i>	<i>entity / its corresponding page</i>
Automatic highlighting	<i>documents in reading</i>	<i>key phrases to be highlighted</i>
Text summarization	<i>long text</i>	<i>summarized short text</i>

Automatic image captioning (MSR system)





A

a woman in a kitchen preparing food

B

woman working on counter near kitchen sink preparing a meal



Machine:

Human:

a woman in a kitchen preparing food

woman working on counter near kitchen sink preparing a meal

COCO Challenge Results (CVPR-2015, Boston)

Tied for
1st prize

	M1	M2	M3	M4	M5
Human ^[5]	0.638	0.675	4.836	3.428	0.352
MSR ^[8]	0.268	0.322	4.137	2.662	0.234
Google ^[4]	0.273	0.317	4.107	2.742	0.233
MSR Captivator ^[9]	0.250	0.301	4.149	2.565	0.233
Montreal/Toronto ^[10]	0.262	0.272	3.932	2.832	0.197
Berkeley LRCN ^[2]	0.246	0.268	3.924	2.786	0.204
Nearest Neighbor ^[11]	0.216	0.255	3.801	2.716	0.196

M1

Percentage of captions that are evaluated as better or equal to human caption.

M2

Percentage of captions that pass the Turing Test.

M3

Average correctness of the captions on a scale 1-5 (incorrect - correct).

M4

Average amount of detail of the captions on a scale 1-5 (lack of details - very detailed).

M5

Percentage of captions that are similar to human description.

Deep Learning for Machine Cognition

- Deep reinforcement learning
- “Optimal” actions: control and business decision making

Human-level control through deep reinforcement learning

Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{1*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fidjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹

The theory of reinforcement learning provides a normative account¹, deeply rooted in psychological² and neuroscientific³ perspectives on animal behaviour, of how agents may optimize their control of an environment. To use reinforcement learning successfully in situations approaching real-world complexity however, agents are confronted

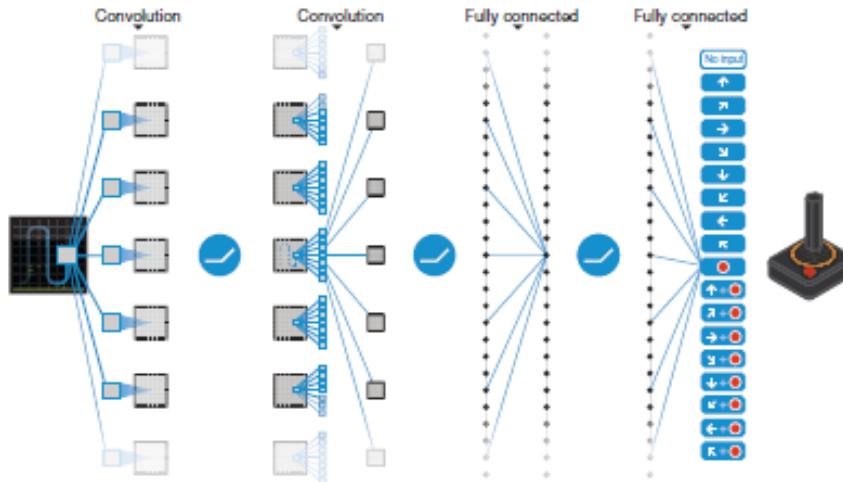
agent is to select actions in a fashion that maximizes cumulative future reward. More formally, we use a deep convolutional neural network to approximate the optimal action-value function

$$Q^*(s,a) = \max_{\pi} \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = s, a_t = a, \pi],$$

Reinforcement learning from “non-working” to “working”, due to Deep Learning (much like DNN for speech)

Deep Q-Network (DQN)

mapping raw screen pixels



to predictions
of final score
for each of 18
joystick actions

- Input layer: image vector of s
- Output layer: a single output Q-value for each action a , $Q(s, a, \theta)$
- DNN parameters: θ

Reinforcement Learning

--- optimizing long-term values

	Short-term	Long-term
Playing the Breakout game		
Optimizing Business Decision Making	<i>Maximize immediate reward</i>	<i>Optimize life-time revenue, service usages, and customer satisfaction</i>



Self play to improve skills

ARTICLE PREVIEW

[view full access options ▶](#)

NATURE | ARTICLE

日本語要約

Mastering the game of Go with deep neural networks and tree search

David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel & Demis Hassabis

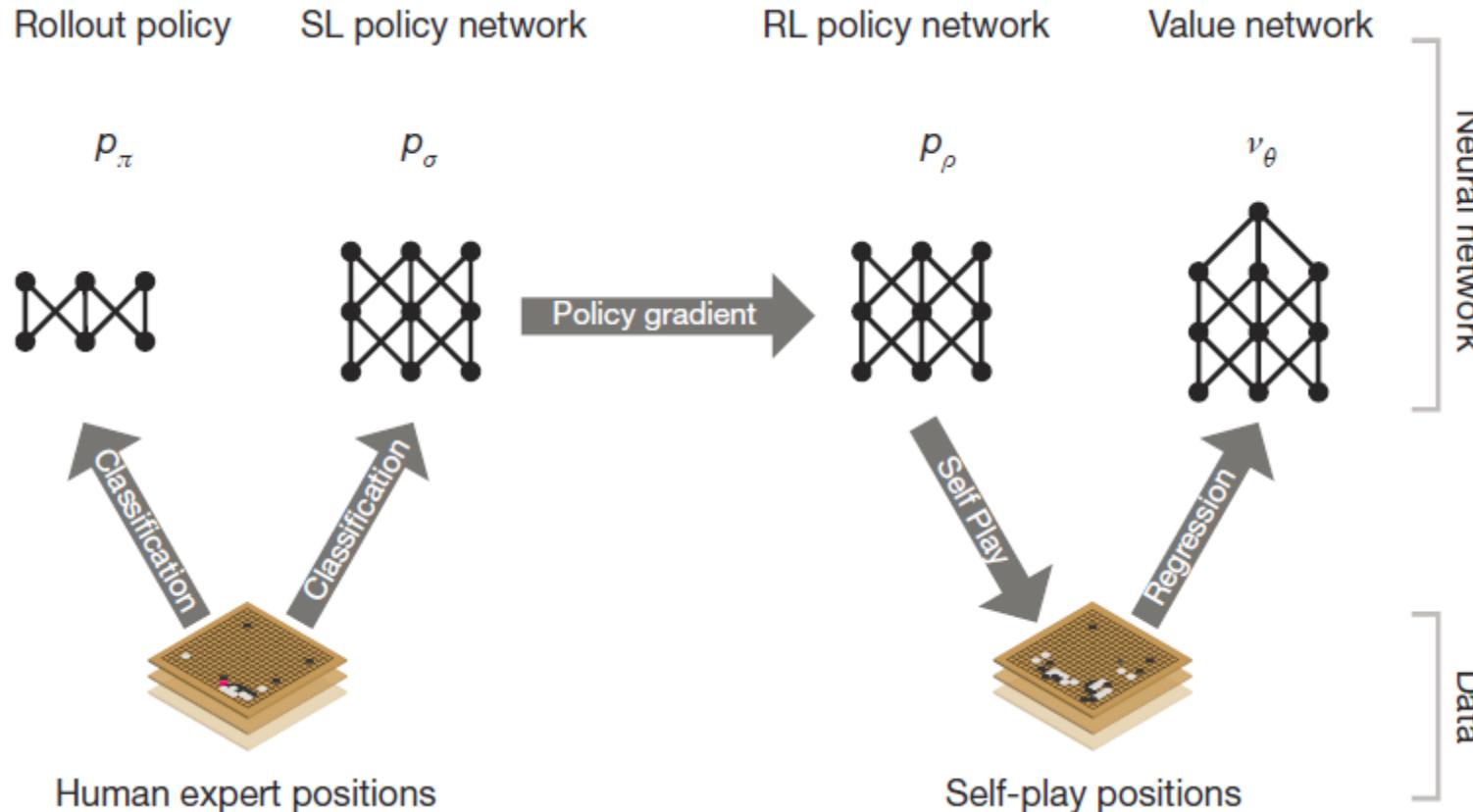
[Affiliations](#) | [Contributions](#) | [Corresponding authors](#)

Nature 529, 484–489 (28 January 2016) | doi:10.1038/nature16961

Received 11 November 2015 | Accepted 05 January 2016 | Published online 27 January 2016



DNN learning pipeline in

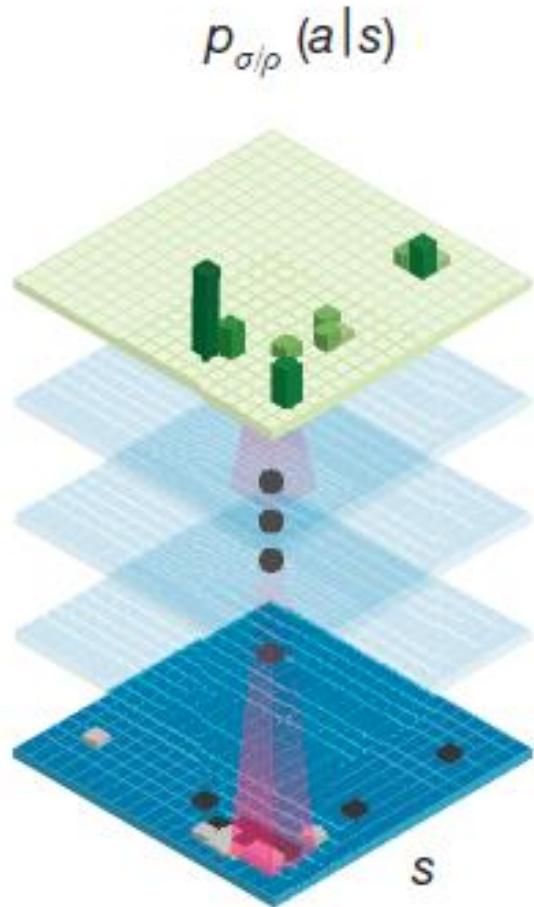


DNN architecture used in

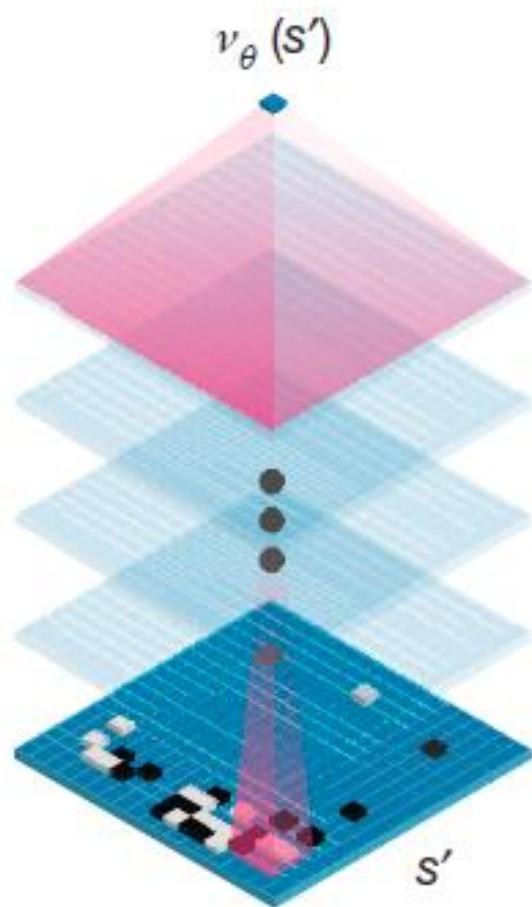


b

Policy network



Value network

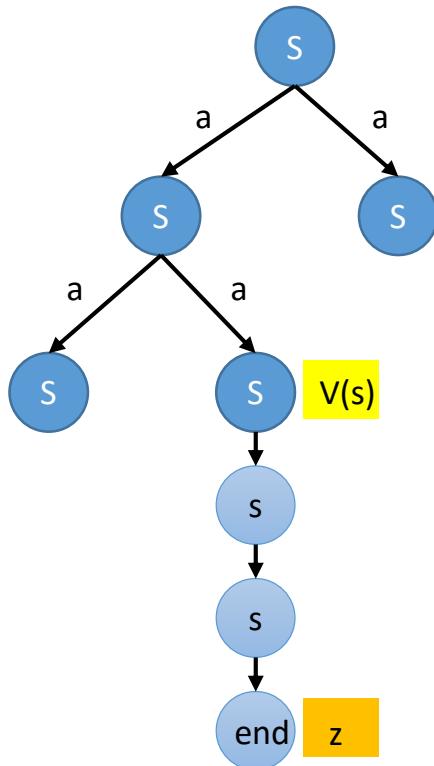


Analysis of four DNNs in



DNNs	Properties	Architecture	Additional details
$\pi_{SL}(a s)$	Slow, accurate stochastic supervised learning policy, trained on 30M (s,a) pairs	13 layer network; alternating ConvNets and rectifier nonlinearities; output dist. over all legal moves	Evaluation time: 3 ms Accuracy vs. corpus: 57% Train time: 3 weeks
$\tilde{\pi}_{SL}(a s)$	Fast, less accurate stochastic SL policy, trained on 30M (s,a) pairs	Linear softmax of small pattern features	Evaluation time: 2 us Accuracy vs. corpus: 24%
$\pi_{RL}(a s)$	Stochastic RL policy, trained by self-play	Same as π_{SL}	Win vs. π_{SL} : 80%
$V(s)$	Value function: % chance of π_{RL} winning by starting in state s	Same as π_{SL} , but with one output (% chance of winning)	15K less computation than evaluating π_{RL} with roll-outs

Monte Carlo Tree Search in



$$\pi(s) = \operatorname{argmax}_a Q(s, a)$$

$$Q(s, a) = Q'(s, a) + u(s, a)$$

Roll-out
estimate

Exploration
bonus

Mixture weight

$$Q'(s, a) = \frac{1}{N(s, a)} \sum_i [(1 - \lambda)V(s_L^i) + \lambda z_L^i]$$

of times action a
taken in state s
Value function
computed in advance
Win/loss result of 1
roll-out with $\tilde{\pi}_{SL}(a|s)$

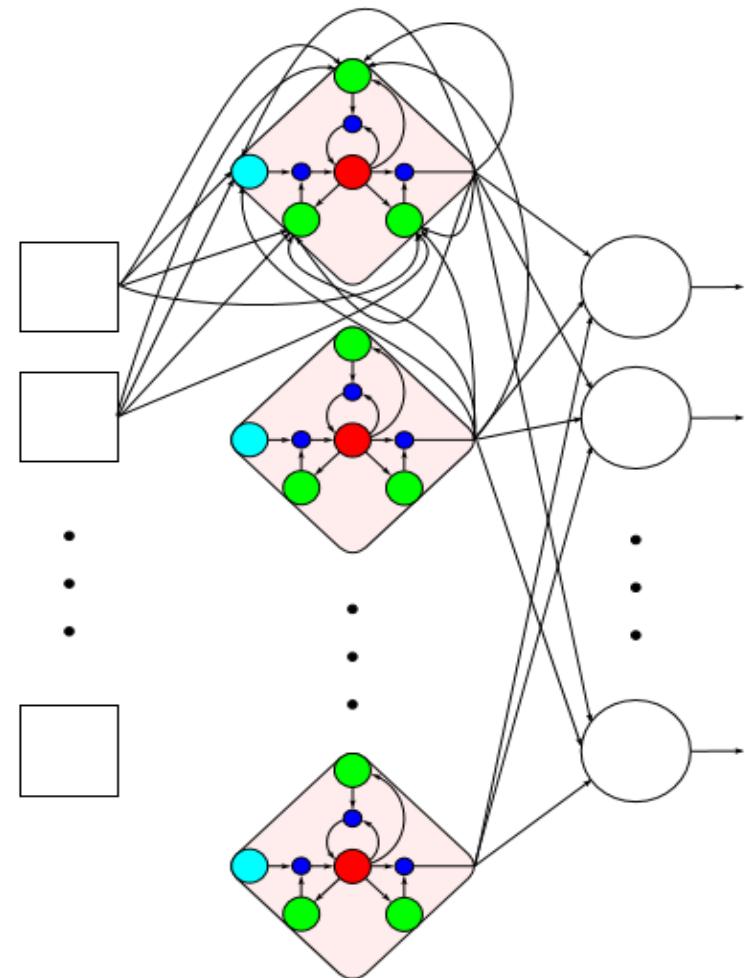
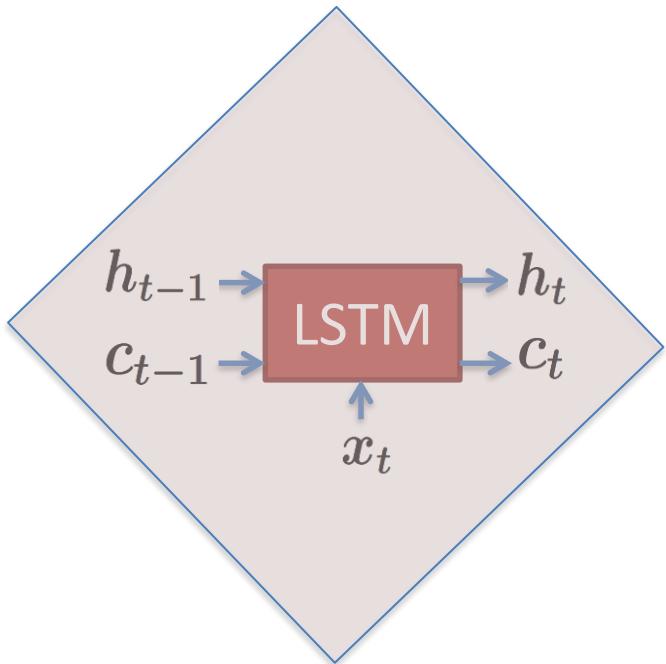
$$u(s, a) = c \cdot \pi_{SL}(a|s) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}$$

- Think of this MCTS component as a highly efficient “decoder”, a concept familiar to ASR
- -> A* search and fast match in speech recognition literature during 80’s-90’s
- This is tree search (GO-specific), not graph search (A*)
- Speech is a relatively simple signal → sequential beam search sufficient, no need for A* or tree
- Key innovation in AlphaGO: “scores” in MCTS computed by DNNs with RL

Deep Learning for Machine Cognition

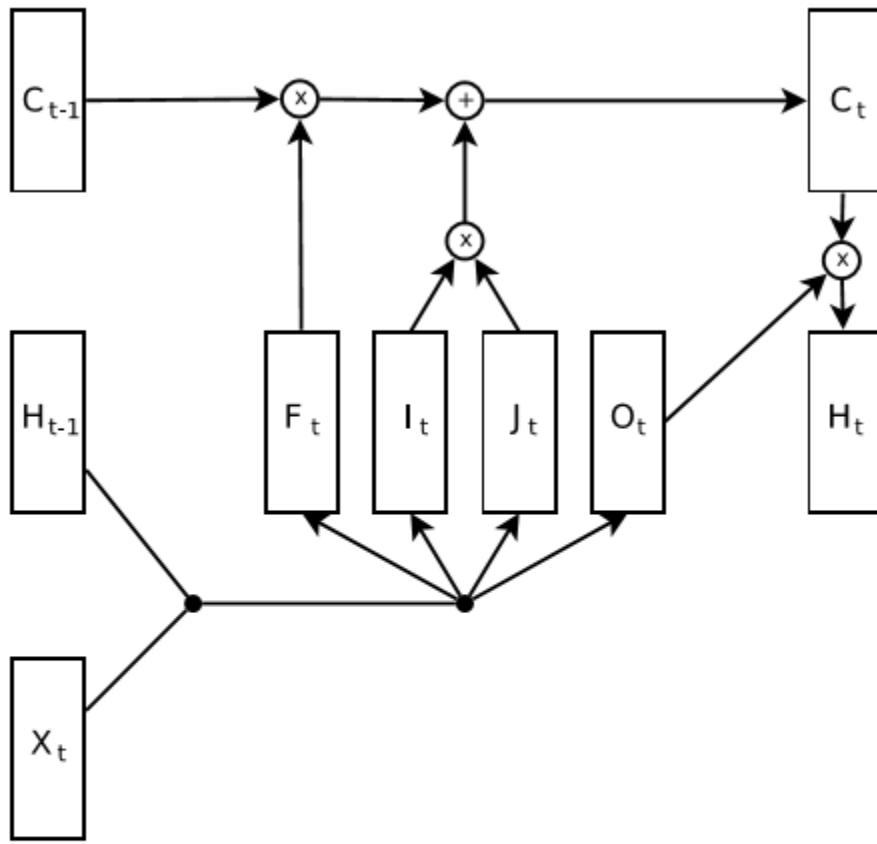
--- Memory & attention (applied to machine translation)

Long Short-Term Memory RNN



(Hochreiter & Schmidhuber, 1997)

LSTM cell unfolding over time



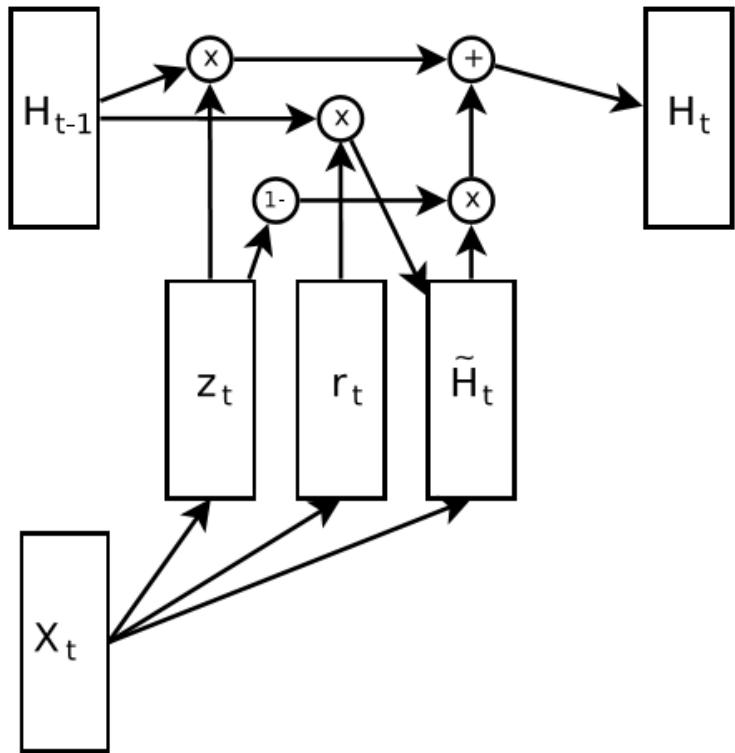
$$\begin{aligned} i_t &= \tanh(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \\ j_t &= \text{sigm}(W_{xj}x_t + W_{hj}h_{t-1} + b_j) \\ f_t &= \text{sigm}(W_{xf}x_t + W_{hf}h_{t-1} + b_f) \\ o_t &= \tanh(W_{xo}x_t + W_{ho}h_{t-1} + b_o) \\ c_t &= c_{t-1} \odot f_t + i_t \odot j_t \\ h_t &= \tanh(c_t) \odot o_t \end{aligned}$$

Figure 1. The LSTM architecture. The value of the cell is increased by $i_t \odot j_t$, where \odot is element-wise product. The LSTM's output is typically taken to be h_t , and c_t is not exposed. The forget gate f_t allows the LSTM to easily reset the value of the cell.

(Jozefowics, Zaremba, Sutskever,
ICML 2015)

Gated Recurrent Unit (GRU)

(simpler than LSTM; no output gates)



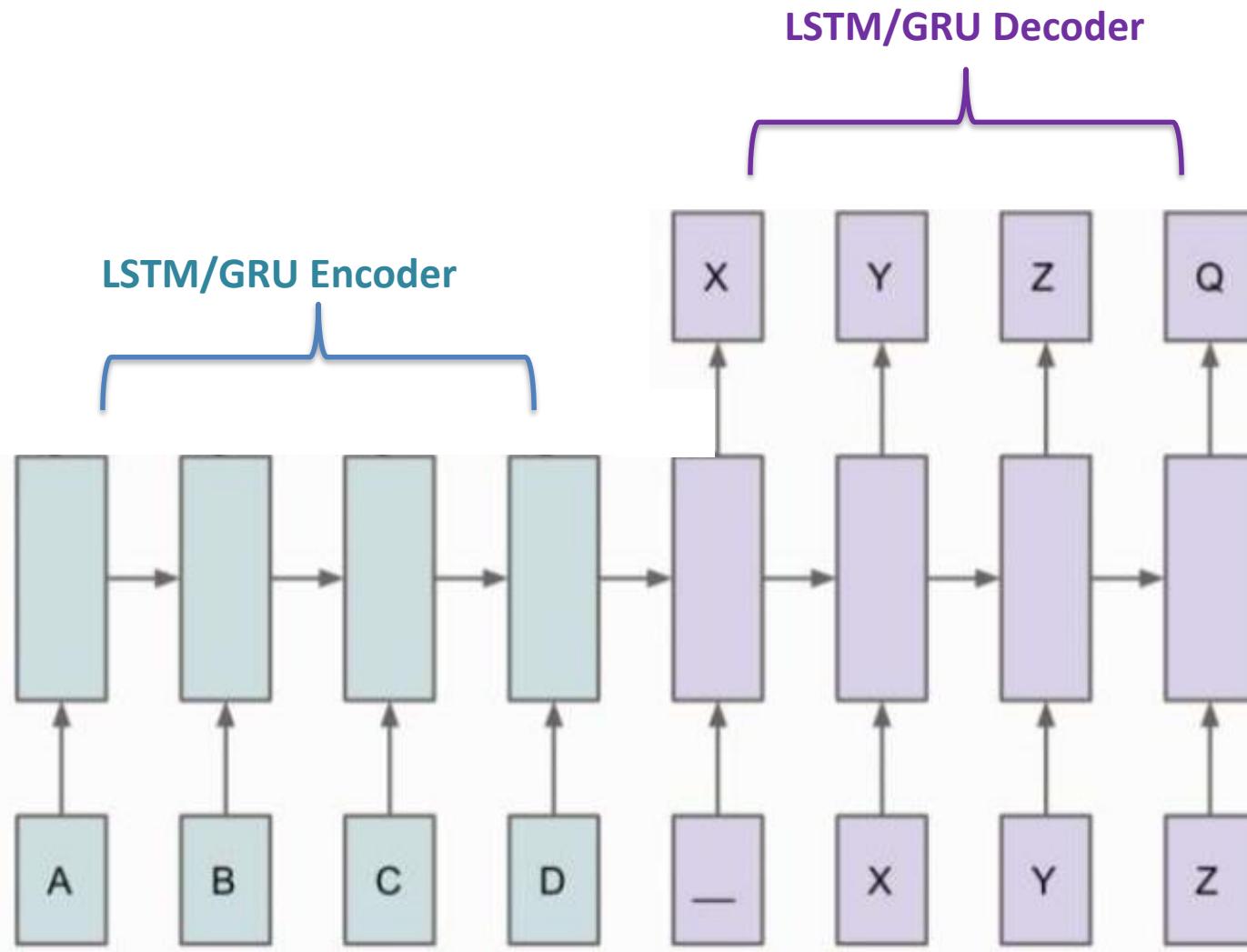
$$\begin{aligned} r_t &= \text{sigm}(W_{xr}x_t + W_{hr}h_{t-1} + b_r) \\ z_t &= \text{sigm}(W_{xz}x_t + W_{hz}h_{t-1} + b_z) \\ \tilde{h}_t &= \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h) \\ h_t &= z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t \end{aligned}$$

Figure 2. The Gated Recurrent Unit. Like the LSTM, it is hard to tell, at a glance, which part of the GRU is essential for its functioning.

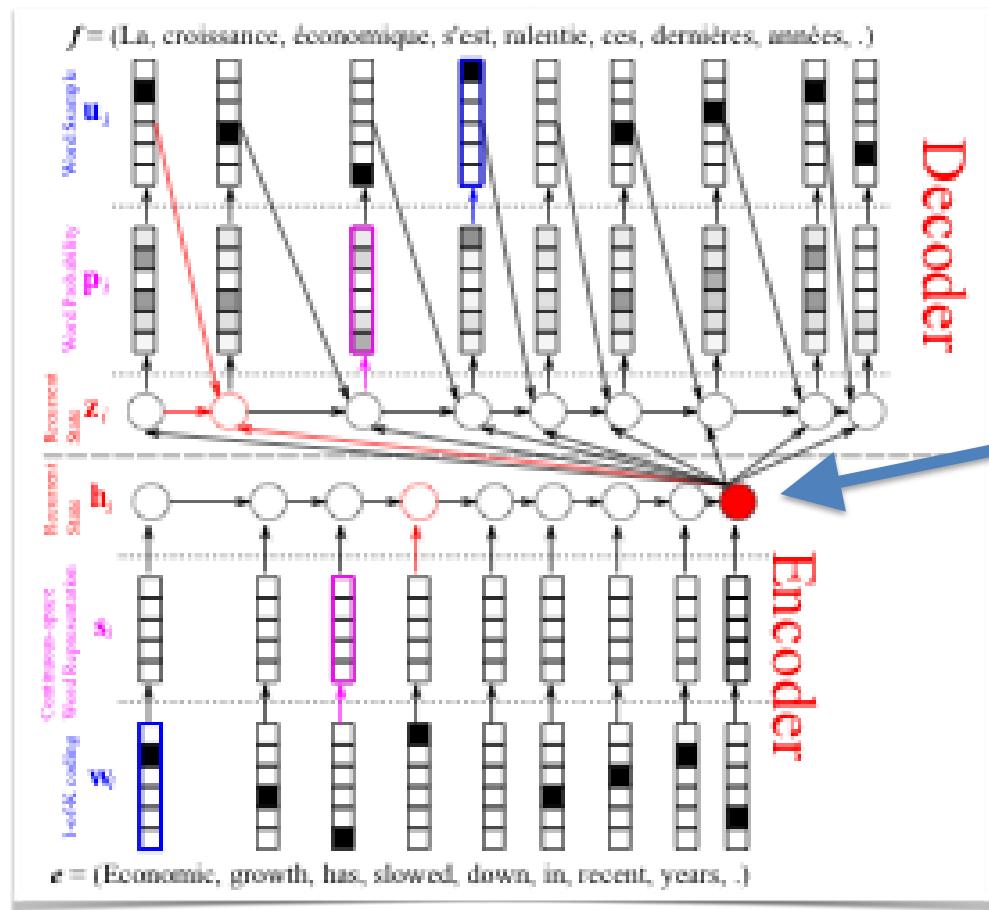
(Jozefowics, Zaremba, Sutskever, ICML 2015; Google Kumar et al., arXiv, July, 2015; Metamind)

Seq-2-Seq Learning (Neural Machine Translation)

[Sutskever, Vinyals, Le, NIPS, 2014]



Neural Machine Translation



“Thought vector”

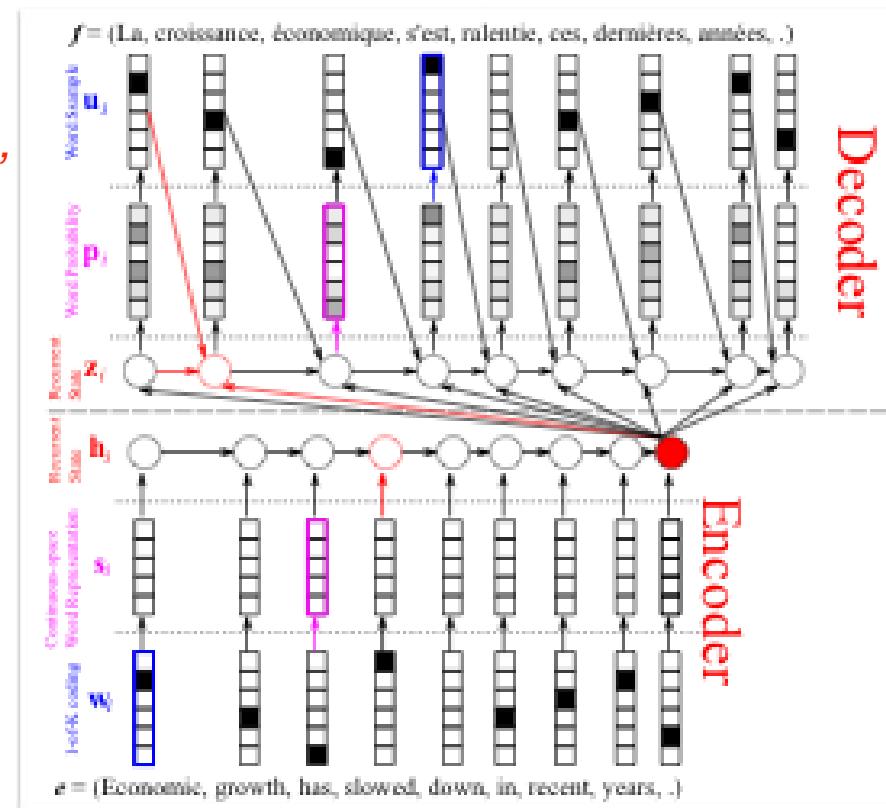
(Forcada&Ñeco, 1997;
Castaño&Casacuberta, 1997;
Kalchbrenner&Blunsom, 2013;
Sutskever et al., 2014;
Cho et al., 2014)

Neural Machine Translation

- This model relying “thought vector” does not perform well
- Especially for long source sentences
- *Because:*

“You can’t cram the meaning of a whole sentence into a single &#!# vector!”*

Ray Mooney



Neural Machine Translation with Attention

Attention-based Model

- Encoder: Bidirectional RNN
 - A set of *annotation* vectors

$$\{h_1, h_2, \dots, h_T\}$$

- Attention-based Decoder

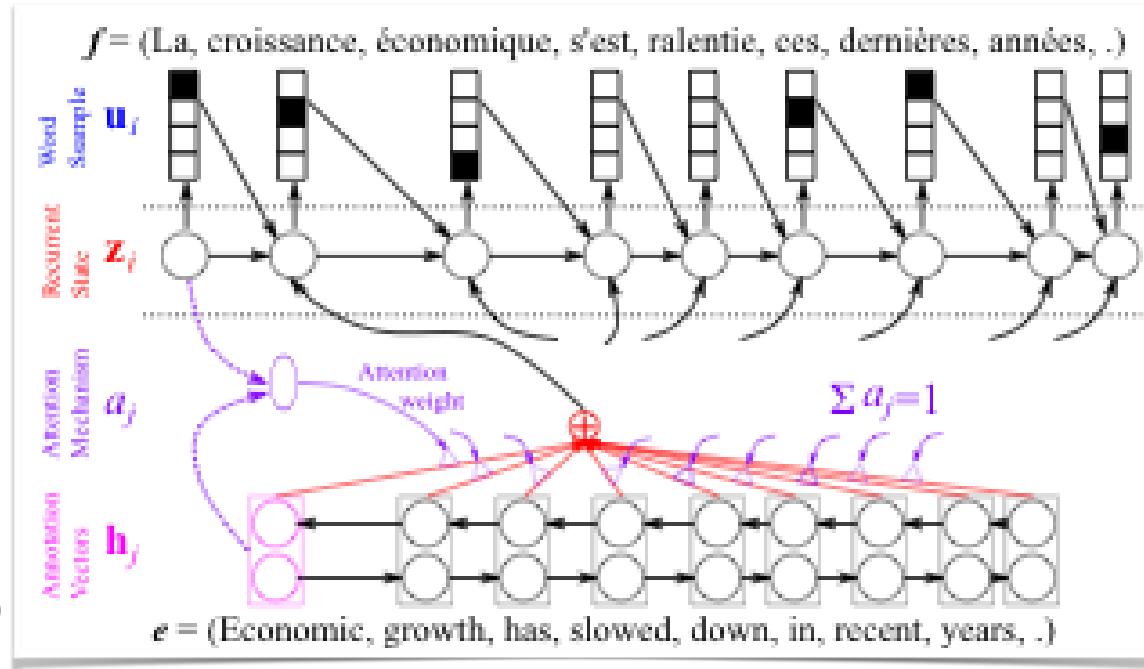
- (1) Compute attention weights

$$\alpha_{t',t} \propto \exp(e(z_{t'-1}, u_{t'-1}, h_t))$$

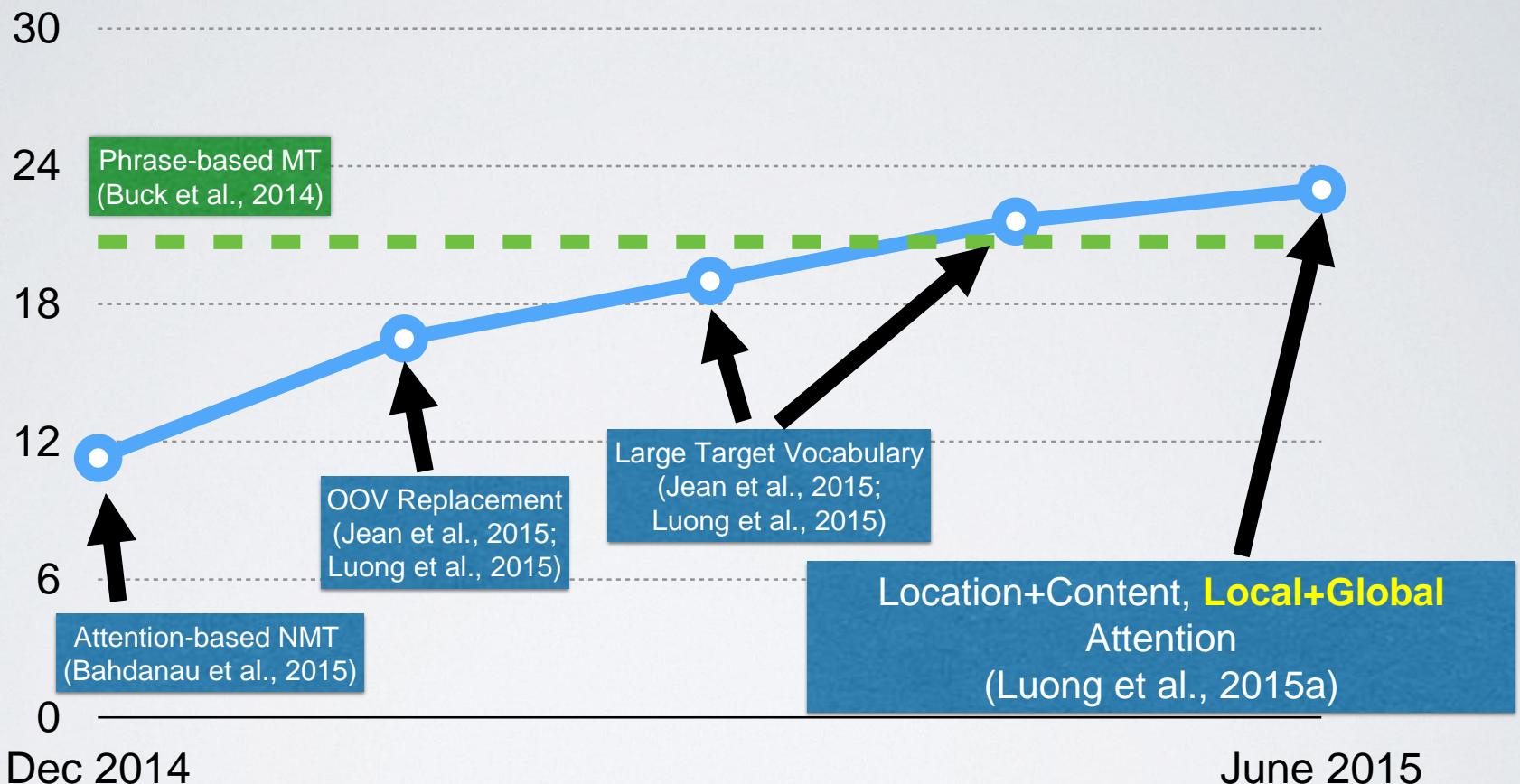
- (2) Weighted-sum of the annotation vectors

$$c_{t'} = \sum_{t=1}^T \alpha_{t',t} h_t$$

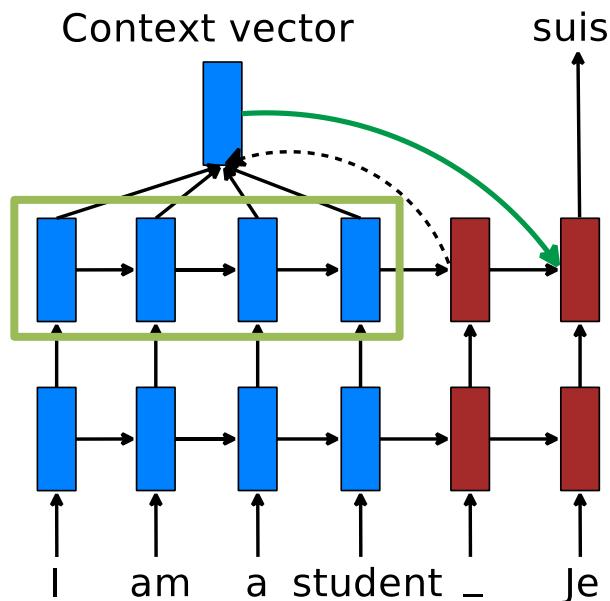
- (3) Use $c_{t'}$ to replace “thought vector” h_T



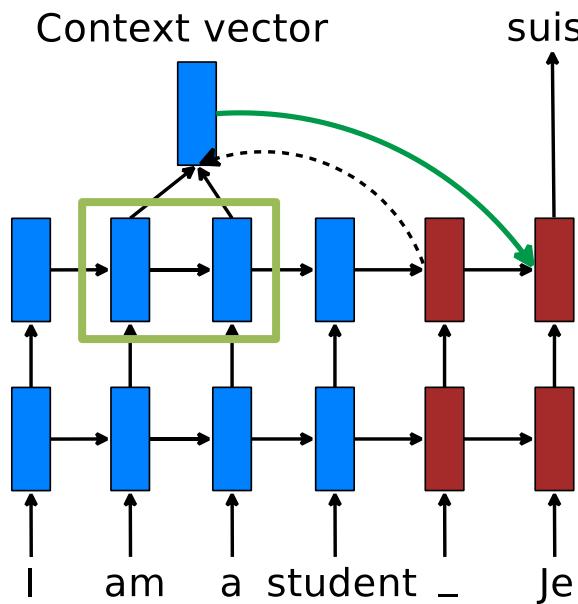
BENCHMARK: WMT'14 EN-DE



Models for Global & Local Attention

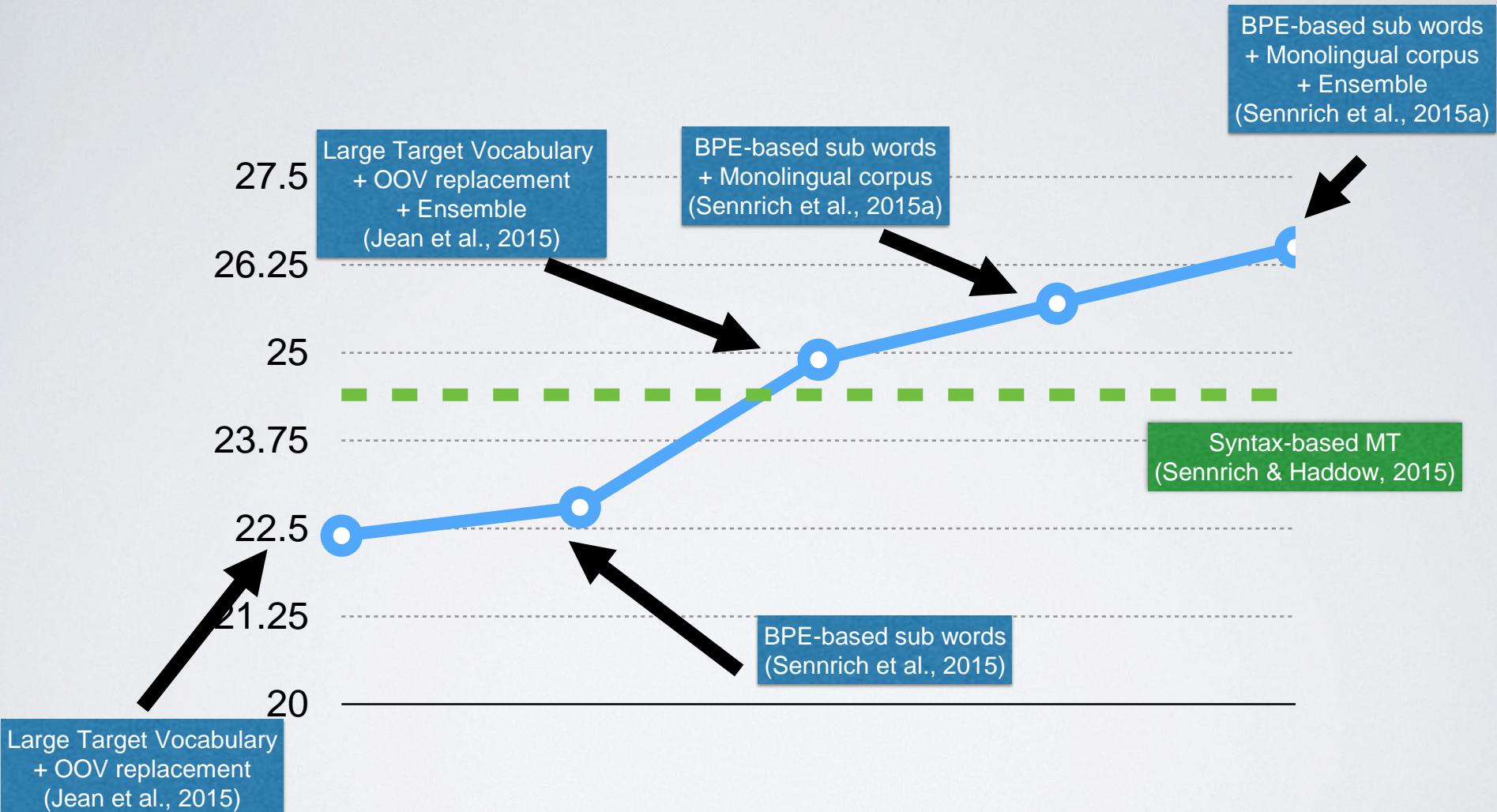


Global: *all* source states.



Local: *subset* of source states.

BENCHMARK: WMT'15 EN-DE



(modified from: Kyunghyun Cho)

Same Attention Model applied to Image Captioning

Topics: Beyond Natural Languages
— Image Caption Generation

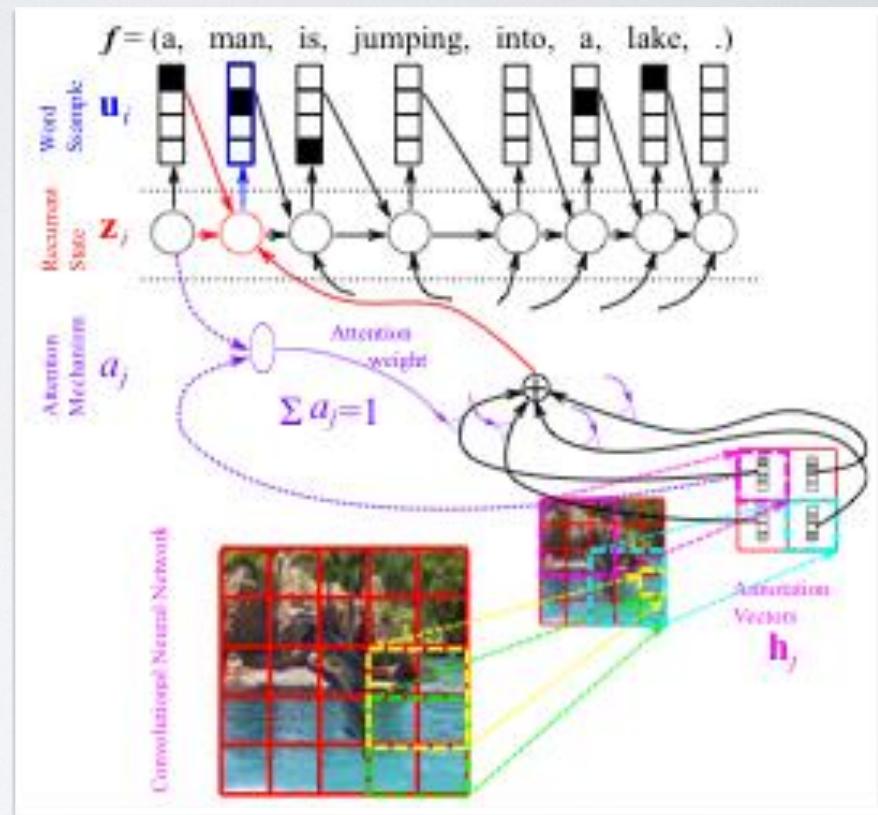
- Conditional language modelling

$p(\text{Two, dolphins, are, diving} |$



) =?

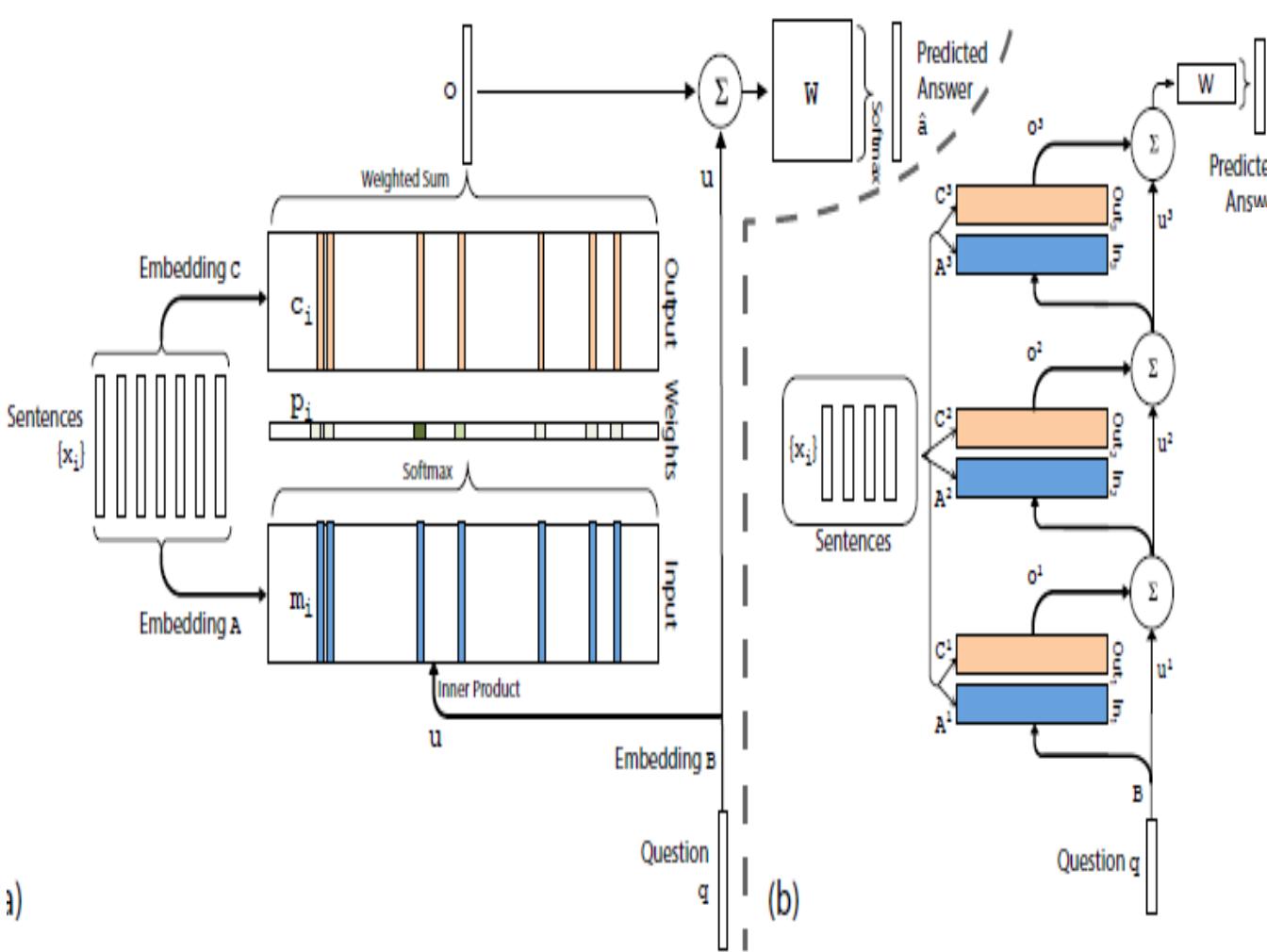
- Encoder: convolutional network
- Pretrained as a classifier or autoencoder
- Decoder: recurrent neural network
- RNN Language model
- With attention mechanism (Xu et al., 2015)



Deep Learning for Machine Cognition

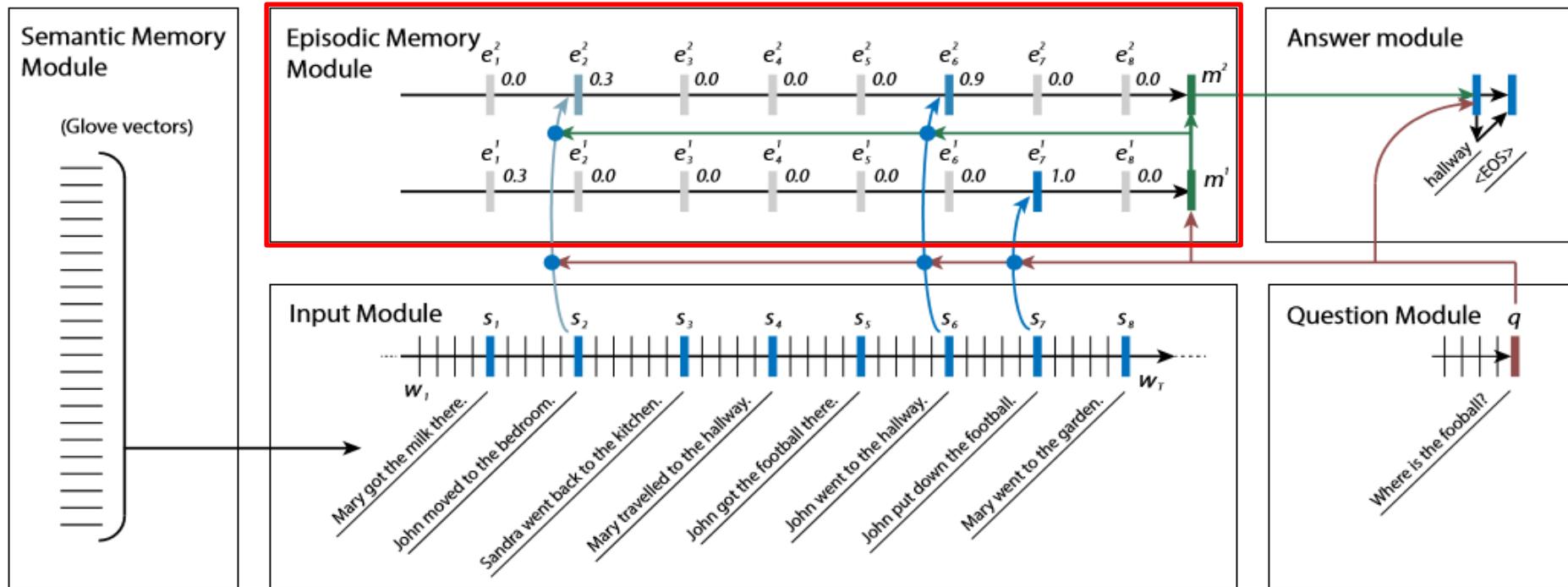
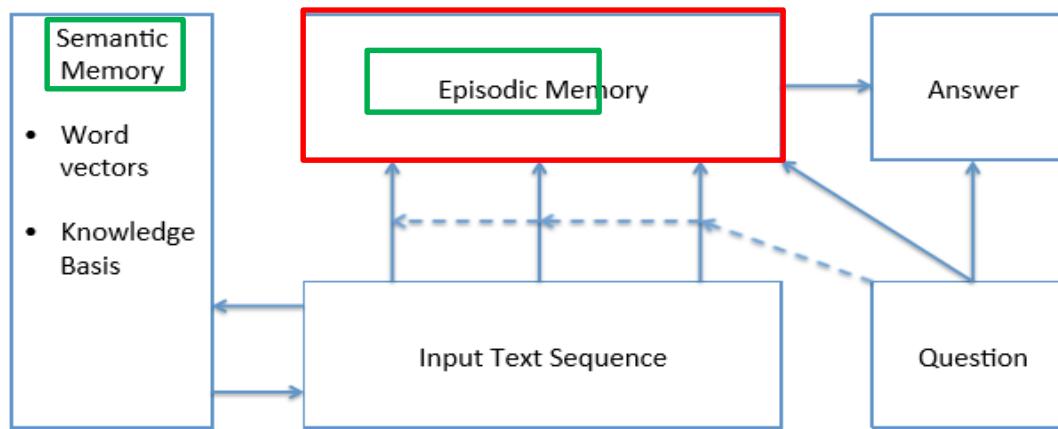
- Neural reasoning: memory network
- Better neural reasoning: Tensor Product Representations (TPR) with structured knowledge representation

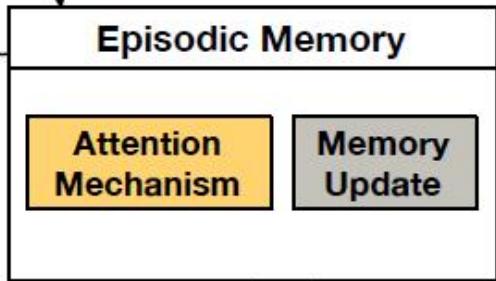
Memory Networks for Reasoning



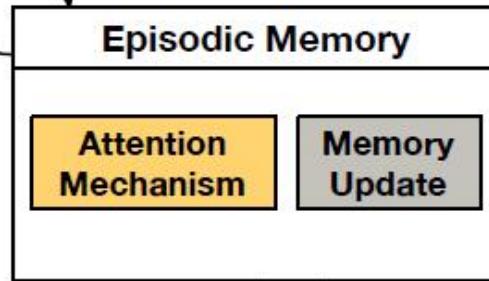
- Rather than placing “attention” to part of a sentence, it can be placed to cognitive space with many sentences
- This allows “reasoning”
- Embedding input
$$m_i = Ax_i$$
$$c_i = Cx_i$$
$$u = Bq$$
- Attention over memories
$$p_i = \text{softmax}(u^T m_i)$$
- Generating the final answer
$$o = \sum_i p_i c_i$$

$$a = \text{softmax}(W(o + u))$$

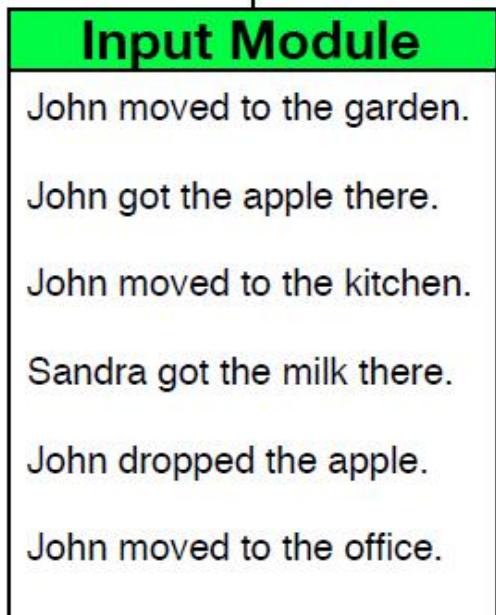




Answer
Kitchen



Answer
Palm



Question
Where is the apple?



(a) Text Question-Answering

(b) Visual Question-Answering

TPR: Neural Representation of Structure

- Structured embedding vectors via tensor-product rep (TPR)



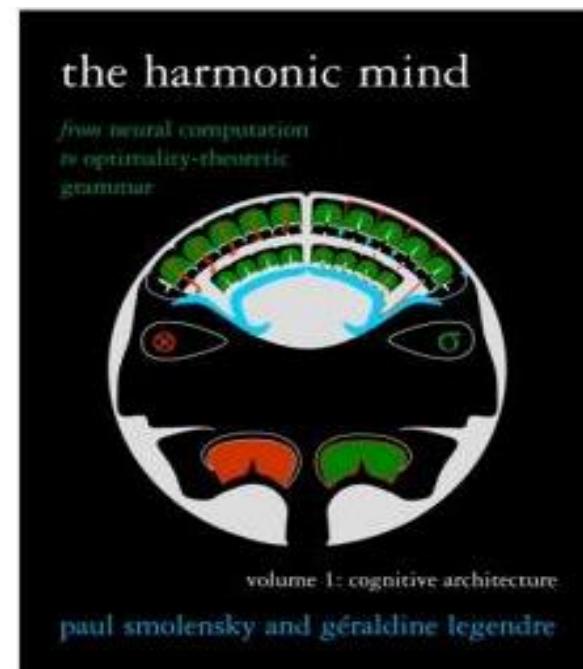
symbolic semantic parse tree (complex relation)

Then, reasoning in symbolic-space (traditional AI) can be beautifully carried out in the continuous-space in human cognitive and neural-net terms

**Paul Smolensky & G. Legendre:
The Harmonic Mind, MIT Press, 2006**

From Neural Computation to Optimality-Theoretic Grammar

Volume I: Cognitive Architecture; Volume 2: Linguistic Implications



Outline

- Deep learning for machine perception
 - Speech
 - Image
- Deep learning for machine cognition
 - Semantic modeling
 - Natural language
 - Multimodality
 - Reasoning, attention, memory (RAM)
 - Knowledge representation/management/exploitation
 - Optimal decision making (by deep reinforcement learning)
- Three hot areas/challenges of deep learning & AI research

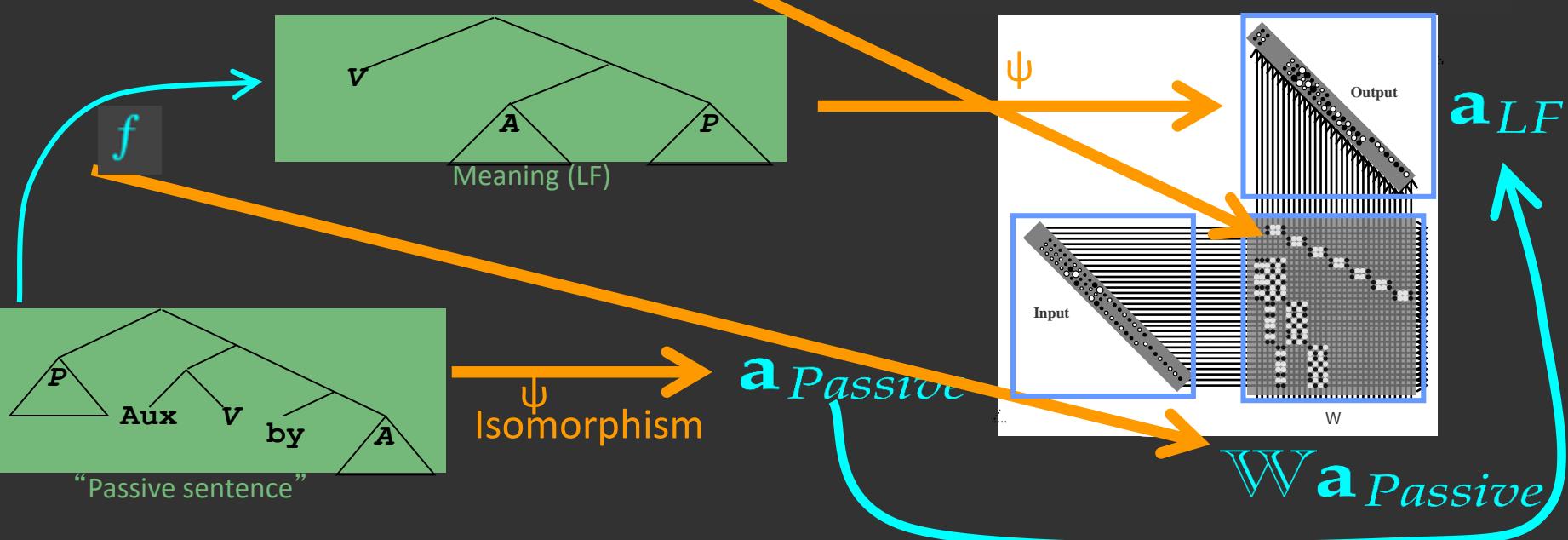
Challenges for Future Research

1. Structured embedding for better reasoning: integrate symbolic/neural representations
2. Integrate deep discriminative & generative/Bayesian models
3. Deep Unsupervised Learning

Few leaders are admired by George Bush \xrightarrow{f} admire(George Bush, few leaders)

$$f(s) = \text{cons}(\text{ex}_1(\text{ex}_0(\text{ex}_1(s))), \text{cons}(\text{ex}_1(\text{ex}_1(\text{ex}_1(s))), \text{ex}_0(s)))$$

$$W = W_{\text{cons}_0}[W_{\text{ex}_1}W_{\text{ex}_0}W_{\text{ex}_1}] + \\ W_{\text{cons}_1}[W_{\text{cons}_0}(W_{\text{ex}_1}W_{\text{ex}_1}W_{\text{ex}_1}) + W_{\text{cons}_1}(W_{\text{ex}_0})]$$



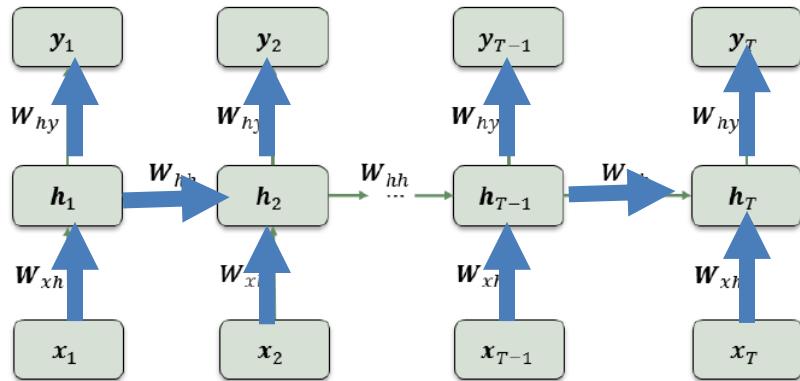
Recurrent NN vs. Dynamic System

$$\mathbf{h}_t = f(\mathbf{h}_{t-1}; \mathbf{W}_{hh}, \mathbf{W}_{xh}, \mathbf{x}_t)$$

$$\mathbf{y}_t = g(\mathbf{h}_t; \mathbf{W}_{hy})$$

Parameterization:

- $\mathbf{W}_{hh}, \mathbf{W}_{hy}, \mathbf{W}_{xh}$: all unstructured regular matrices

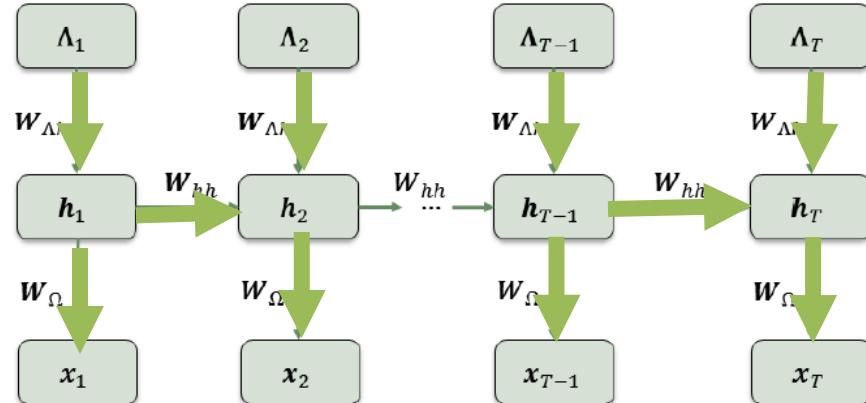


$$\mathbf{h}_t = q(\mathbf{h}_{t-1}; \mathbf{W}_{l_t}, \mathbf{t}_{l_t})$$

$$\mathbf{x}_t = r(\mathbf{h}_t, \Omega_{l_t})$$

Parameterization:

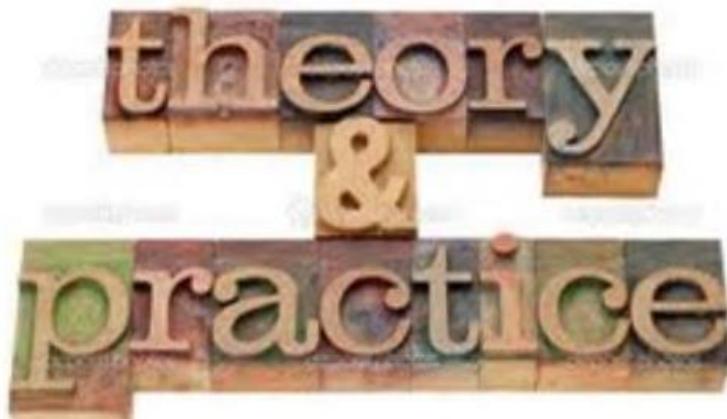
- $\mathbf{W}_{hh} = \mathbf{M}(\gamma_l)$; sparse system matrix
- $\mathbf{W}_{\Omega} = (\Omega_l)$; Gaussian-mix params; MLP
- $\Lambda = \mathbf{t}_l$



	Deep Discriminative NN	Deep Generative (Bayesian)
Structure	Graphical; info flow: bottom-up	Graphical; info flow: top-down
Incorp constraints & domain knowledge	Harder; less fine-grained	Easier; more fine grained
Semi/unsupervised	Hard or impossible	<i>Easier, at least possible</i>
Interpretability	Harder	Easy (generative “story” on data and hidden variables)
Representation	Distributed	Localist (mostly); can be distributed also
Inference/decode	Easy	Harder (but note recent progress)
Scalability/compute	Easier (regular computes/GPU)	Harder (but note recent progress)
Incorp. uncertainty	Hard	Easy
Empirical goal	Classification, feature learning, ...	Classification (via Bayes rule), latent variable inference...
Terminology	Neurons, activation/gate functions, weights ...	Random vars, stochastic “neurons”, potential function, parameters ...
Learning algorithm	A single, unchallenged, algorithm -- BackProp	A major focus of open research, many algorithms, & more to come
Evaluation	On a black-box score – end performance	On almost every intermediate quantity
Implementation	Hard (but increasingly easier)	Standardized but insights needed
Experiments	Massive, real data	Modest, often simulated data
Parameterization	Dense matrices	Sparse (often PDFs); can be dense

Deep Unsupervised Learning

- Unsupervised learning (UL) has recently been a very hot topic in deep learning
- Need to have a task to ground UL --- e.g. help improve prediction
- Examples of speech recognition and image captioning:
 - 3000 hrs of paired acoustics (X) & word label (Y)
 - How can we exploit 300,000+ hrs of speech acoustics with no paired labels?
- 4 sources of knowledge
 - Strong structure prior of “labels” Y (sequences)
 - Strong structure prior of input data X (conventional UL)
 - Dependency of x on y (generative modeling for embedding knowledge)
 - Dependency of y on x (state of the art systems w. supervised learning)



“In theory there is no difference between theory and practice, but in practice there is”

[Jan L. A. van de Snepscheut]

End
(of Chapter 1)

Thank you!
Q/A

深度学习是人工智能领域最近20年里最受瞩目的研究方向，近年来显著推动了语音、图像、自然语言理解、机器翻译，甚至是控制等众多技术方向的发展。本书原著者微软研究院的邓力博士和俞栋博士是语音识别和深度学习方面的先驱之一，对于深度学习的进展有丰富的实践经验和深刻理解。这个学科处于快速进展之际，本书对当前的进展进行全景式系统性的梳理无疑是很有意义的，因为毕竟对于每一位读者，从这几年浩如烟海的论文中准确把握可以沉淀下来的进展是不容易的。谢磊教授受邓力博士之约在百忙之中对这本书进行翻译，对于深度学习在中国的发展具有重大意义。邓力博士和谢磊教授都是我所熟知的学者和好友。我相信，本书作为他们这次合作的成果，对于有志于了解和学习深度学习的中国读者会有极大的帮助。

——余凯

地平线机器人技术 创始人/CEO

前百度研究院常务副院长、深度学习实验室主任

深度学习方法及应用

谢磊
译

深度学习 方法及应用

[美] 邓力 (Li Deng) 俞栋 (Dong Yu) 著

谢磊 译

机械工业出版社

ISBN 978-7-111-52906-4
9 787111 529064
定价：24.00元

地址：北京市百万庄大街22号
邮政编码：100037
电话服务
服务咨询热线：010-88361066
读者购书热线：010-68326294
010-88379203
网络服务
机工官网：www.cmpbook.com
机工微博：weibo.com/cmp1952
金书网：www.golden-book.com
教育服务网：www.cmpedu.com
封面无防伪标均为盗版



机械工业出版社微信公众号
ISBN 978-7-111-52906-4
策划编辑〇王康

机械工业出版社
CHINA MACHINE PRESS

解析深度学习 语音识别实践

语音识别在最近几年里取得了长足的进展，这些进展主要源于在语音识别中引入了深度学习技术。本书作者俞栋博士和邓力博士长期致力于语音识别技术的研究，有丰富的理论和实践经验。他们最先将深度学习技术与传统语音识别技术相结合，成功地大幅降低了大词汇量语音识别系统的错误率。本书译者俞凯博士和钱彦旻博士也是语音识别领域的专家，还是这一进展的积极推动者。他们合作的这本中文译本是第一部系统性地介绍基于深度学习的语音识别技术的专著。在本书中，俞栋博士和邓力博士以第一手资料详细介绍了这一技术发生的背景、发展过程、理论根据、关键技术细节，以及思维方式。本书对所有从事语音识别研究或想了解语音识别技术最新进展和发展方向的读者都是很好的参考书。对想将深度学习技术应用到诸如视觉和文本处理等其他领域的读者，本书也有很强的借鉴意义。

微软公司杰出工程师兼首席语音科学家 黄学东博士



博文视点Broadview



新浪微博
weibo.com
@博文视点Broadview

本书将由电子工业出版社于2016年5月出版，欢迎有兴趣的读者联系我们采购，先睹为快！
联系人：杜老师
联系电话：010-8825 4888
邮箱：duca@phei.com.cn
网址：www.phei.com.cn

策划编辑：刘 娅
责任编辑：李利健
封面设计：吴海燕

Signals and Communication Technology
Automatic Speech Recognition
A Deep Learning Approach

解析深度学习 语音识别实践

【美】俞栋 邓力 著
俞凯 钱彦旻 等译



中国工信出版集团



电子工业出版社
THE HOUSE OF ELECTRONIC INDUSTRY
<http://www.phei.com.cn>

Tensor Product Rep for reasoning

- Facebook's reasoning task

Category 4. Two Argument Relation

- 01: The office is north of the kitchen.
02: The garden is south of the kitchen.
03: What is north of the kitchen? office 1

- 01: The kitchen is west of the garden.
02: The hallway is west of the kitchen.
03: What is the garden east of? kitchen 1

Category 17: Positional Reasoning

- 01: The triangle is above the pink rectangle.
02: The blue square is to the left of the triangle.
03: Is the pink rectangle to the right of the blue square? yes 1 2

- 01: The red sphere is below the yellow square.
02: The red sphere is above the blue square.
03: Is the blue square below the yellow square?
yes 2 1

Category 19: Path Finding

- 01: The bedroom is south of the hallway.
02: The bathroom is east of the office.
03: The kitchen is west of the garden.
04: The garden is south of the office.
05: The office is south of the bedroom.
06: How do you go from the garden to the bedroom? n.n 4 5

arXiv.org > cs > arXiv:1511.06426

Computer Science > Computation and Language

Reasoning in Vector Space: An Exploratory Study of Question Answering

Moontae Lee, Xiaodong He, Wen-tau Yih, Jianfeng Gao, Li Deng, Paul Smolensky

(Submitted on 19 Nov 2015)

Accepted to ICLR, May 2016

Structured Knowledge Representation & Reasoning via TPR

#	Statements/Questions	Relational Translations/Answers	Encodings/Clues
1	Mary went to the kitchen.	<i>Mary belongs to the kitchen (from nowhere).</i>	mk^T $m(k \circ n)^T$
3	Mary got the football there.	<i>The football belongs to Mary.</i>	fm^T fm^T
4	Mary travelled to the garden.	<i>Mary belongs to the garden (from the kitchen).</i>	mg^T $m(g \circ k)^T$
5	Where is the football?	<i>garden</i>	3, 4
9	Mary dropped the football.	<i>The football belongs to where Mary belongs to.</i>	fg^T fg^T
10	Mary journeyed to the kitchen.	<i>Mary belongs to the kitchen (from the garden).</i>	mk^T $m(k \circ g)^T$
11	Where is the football?	<i>garden</i>	9, 4

- Given *containee-container* relationship
- Encode all entities (e.g., actors (*mary*), objects (*football*), and locations (*nowhere*, *kitchen*, *garden*)) by **vectors**
- Encode each statement by a **matrix** via **binding** (tensor product of two vectors), mk^T
- Reasoning (transitivity) by **matrix multiplication**, $(fm^T) \cdot (mg^T) = f(m^T \cdot m)g^T = fg^T$
- Generate answer (e.g., where is the football in #5) via **unbinding** (inner product)
 - Left-multiply by f^T all statements prior to the current time. (Yields $f^T \cdot mk^T, f^T \cdot fg^T$)
 - Pick the most recent container where 2-norms of the multiplications in (a) are approximately 1.0. (Yields g^T .)

TPR Results on FB's bAbI task

Type	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
Accuracy	100%	100%	100%	100%	99.3%	100%	96.9%	96.5%	100%	99%
Model	MNN	MNN	MNN	MNN	DMN	MNN	DMN	DMN	DMN	SSVM
Type	C11	C12	C13	C14	C15	C16	C17	C18	C19	C20
Accuracy	100%	100%	100%	100%	100%	100%	72%	95%	36%	100%
Model	MNN	MNN	MNN	DMN	MNN	MNN	Multitask	MNN	MNN	MNN

Table 4: Best accuracies for each category and the best model which achieved the best accuracy. MNN indicates Strongly-Supervised MemNN trained with the clue numbers, and DMN indicates Dynamic MemNN, and SSVM indicates Structured SVM with the coreference resolution and SRL features. Multitask indicates multitask training.

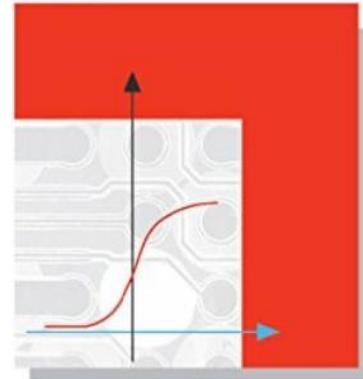
Type	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
Training	100%	100%	100%	100%	99.8%	100%	100%	100%	100%	100%
Test	100%	100%	100%	100%	99.8%	100%	100%	100%	100%	100%
Type	C11	C12	C13	C14	C15	C16	C17	C18	C19	C20
Training	100%	100%	100%	100%	100%	99.4%	100%	100%	100%	100%
Test	100%	100%	100%	100%	100%	99.5%	100%	100%	100%	100%

Table 5: Accuracies on training and test data on our models. We achieve near-perfect accuracies in almost every category including positional reasoning and path finding.

HANDBOOK OF PATTERN RECOGNITION AND COMPUTER VISION

5th Edition

editor
C H Chen



CHAPTER 1.2

DEEP DISCRIMINATIVE AND GENERATIVE MODELS FOR PATTERN RECOGNITION



Li Deng¹ and Navdeep Jaitly²

¹*Microsoft Research, One Microsoft Way, Redmond, WA 98052*

²*Google Research, 1600 Amphitheatre Parkway, Mountain View, CA 94043*

E-mails: deng@microsoft.com; ndjaitly@google.com

In this chapter we describe deep generative and discriminative models as they have been applied to speech recognition and related pattern recognition problems. The former models describe the distribution of data or the joint distribution of data and the corresponding targets, whereas the latter models describe the distribution of targets conditioned on data. Both models are characterized as being ‘deep’ as they use layers of latent or hidden variables. Understanding