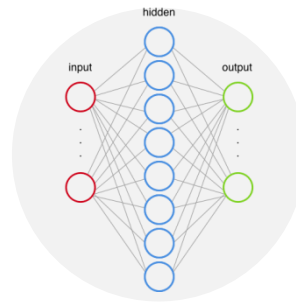


DEEP LEARNING

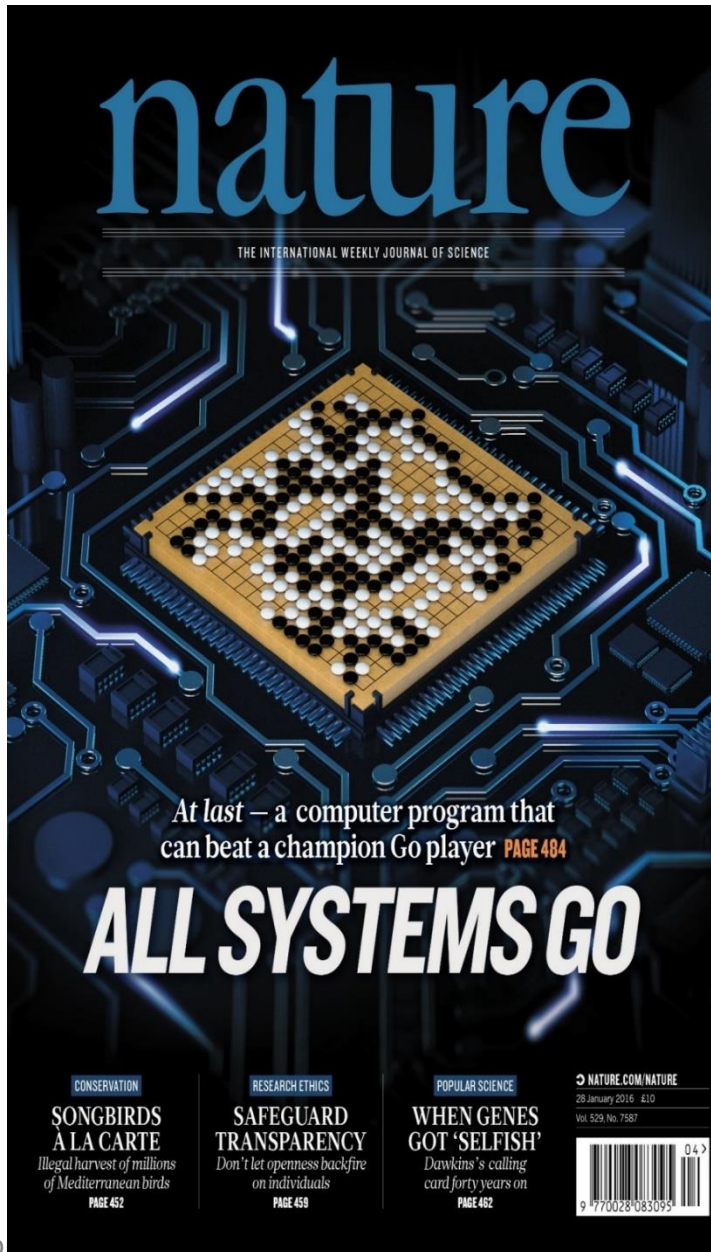
May 21, 2019



1. Technology Landscape
2. Overview
3. Customer References
4. Key Service Offerings
5. Human Resources
6. Case Studies



DEEP LEARNING - INTRODUCTION



In “Nature” 27 January 2016:

- “DeepMind’s program AlphaGo beat Fan Hui, the European Go champion, five times out of five in tournament conditions...”
- “AlphaGo was not preprogrammed to play Go: rather, it learned using a general-purpose algorithm that allowed it to interpret the game’s patterns.”
- “...AlphaGo program applied **deep learning** in neural networks (convolutional NN) — brain-inspired programs in which connections between layers of simulated neurons are strengthened through examples and experience.”

Insufficient depth can hurt

- With shallow architecture (SVM, NB, KNN, etc.), the required number of nodes in the graph (i.e. computations, and also number of parameters, when we try to learn the function) may grow very large.
- Many functions that can be represented efficiently with a deep architecture cannot be represented efficiently with a shallow one.

The brain has a deep architecture

- The visual cortex shows a sequence of areas each of which contains a representation of the input, and signals flow from one to the next.
- Note that representations in the brain are in between dense distributed and purely local: they are **sparse**: about 1% of neurons are active simultaneously in the brain.

Cognitive processes seem deep

- Humans organize their ideas and concepts hierarchically.
- Humans first learn simpler concepts and then compose them to represent more abstract ones.
- Engineers break-up solutions into multiple levels of abstraction and processing

a challenge for ML, CV, AI, Neuroscience, Cognitive Science...

■ **How do we learn representations of the perceptual world?**

- ▶ How can a perceptual system build itself by looking at the world?
- ▶ How much prior structure is necessary

■ **ML/AI: how do we learn features or feature hierarchies?**

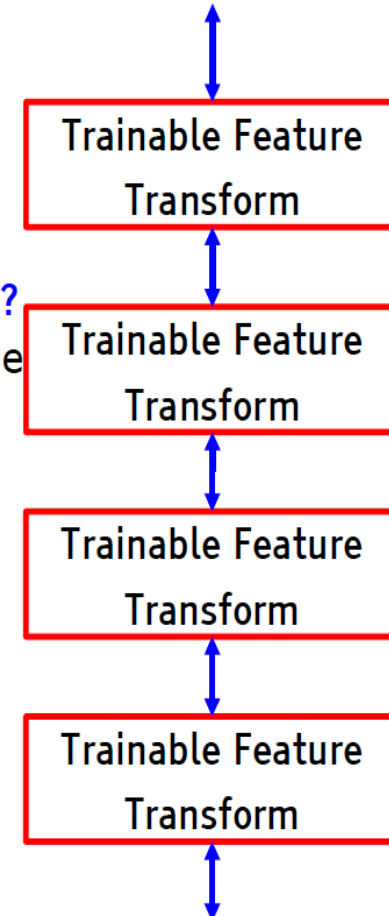
- ▶ What is the fundamental principle? What is the learning algorithm? What is the architecture?

■ **Neuroscience: how does the cortex learn perception?**

- ▶ Does the cortex “run” a single, general learning algorithm? (or a small number of them)

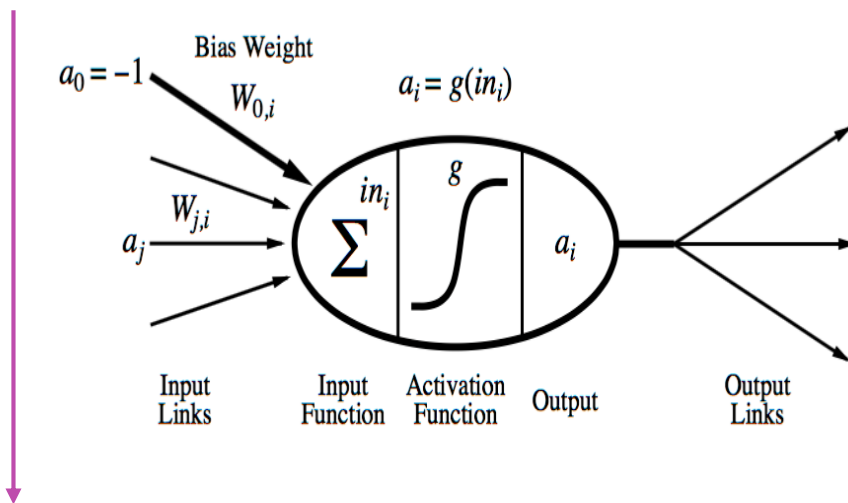
■ **CogSci: how does the mind learn abstract concepts on top of less abstract ones?**

■ **Deep Learning addresses the problem of learning hierarchical representations with a single algorithm**

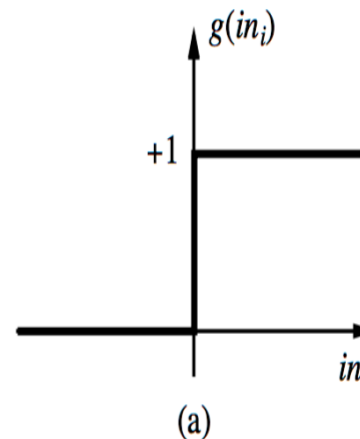


- A neuron is a computational unit in the neural network that exchanges messages with each other.

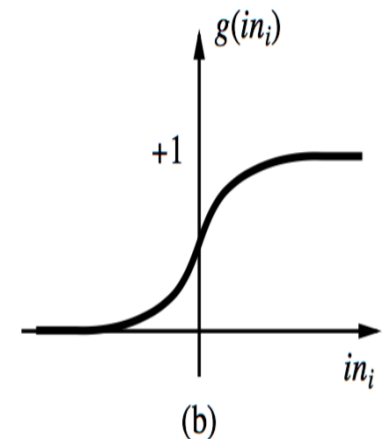
$$a_i \leftarrow g(in_i) = g(\sum_j W_{j,i} a_j)$$



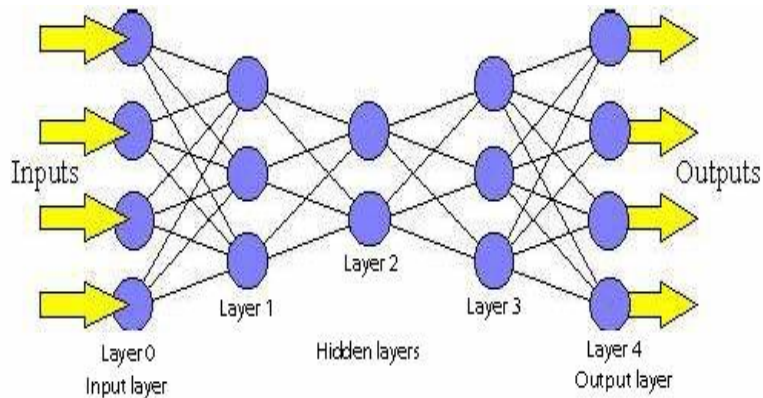
Possible activation functions:



Step function/threshold function

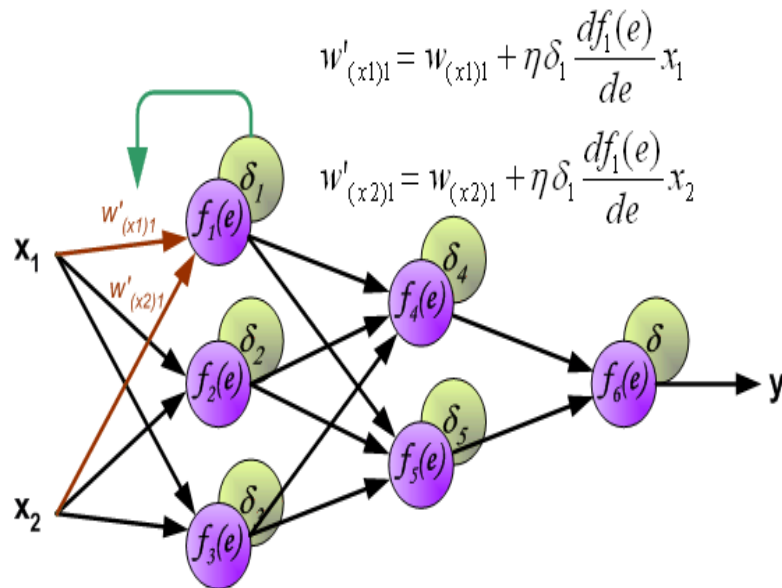


Sigmoid function (a.k.a, logistic function)



Feed forward algorithm

- Activate the neurons from the bottom to the top.



Backpropagation

- Randomly initialize the parameters
- Calculate total error at the top, $f_6(e)$
- Then calculate contributions to error, δ_n , at each step going backwards.

Random initialization + densely connected networks lead to:

■ High cost

- Each neuron in the neural network can be considered as a logistic regression.
- Training the entire neural network is to train all the interconnected logistic regressions.

■ Difficult to train as the number of hidden layers increases

- Recall that logistic regression is trained by gradient descent.
- In backpropagation, gradient is progressively getting more dilute. That is, below top layers, the correction signal δ_n is minimal.

■ Stuck in local optima

- The objective function of the neural network is usually not convex.
- The random initialization does not guarantee starting from the proximity of global optima.

■ Solutions

- Deep Learning/Learning multiple levels of representation

■ Deep Belief Networks & Autoencoders

We will cover two major deep learning models:

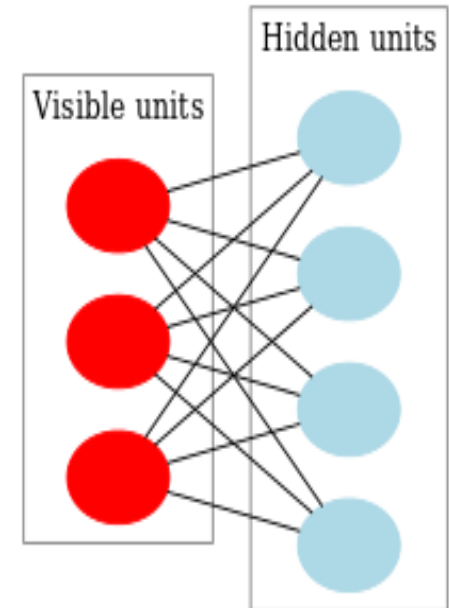
- Employs layer-wise unsupervised learning to initialize each layer and capture multiple levels of representation simultaneously.
 - Hinton, G. E, Osindero, S., and Teh, Y. W. (2006). A fast learning algorithm for deep belief nets. *Neural Computation*, 18:1527-1554.
 - Bengio, Y., Lamblin, P., Popovici, P., Larochelle, H. (2007). Greedy Layer-Wise Training of Deep Networks, *Advances in Neural Information Processing Systems* 19

■ Convolutional Neural Network

- Organizes neurons based on animal's visual cortex system, which allows for learning patterns at both local level and global level.
 - Y. LeCun, L. Bottou, Y. Bengio and P. Haffner: Gradient-Based Learning Applied to Document Recognition, *Proceedings of the IEEE*, 86(11):2278-2324, November 1998

- **A *deep belief network* (DBN)** is a probabilistic, generative model made up of multiple layers of hidden units.
 - A composition of simple learning modules that make up each layer
- A DBN can be used to generatively pre-train a DNN by using the learned DBN weights as the initial DNN weights.
 - Back-propagation or other discriminative algorithms can then be applied for fine-tuning of these weights.
- Advantages:
 - Particularly helpful when limited training data are available
 - These pre-trained weights are closer to the optimal weights than are randomly chosen initial weights.

- A DBN can be efficiently trained in an unsupervised, layer-by-layer manner, where the layers are typically made of *restricted Boltzmann machines* (RBM).
- RBM: undirected, generative energy-based model with a "visible" input layer and a hidden layer, and connections between the layers but not within layers.
- The training method for RBMs proposed by Geoffrey Hinton for use with training "Product of Expert" models is called *contrastive divergence* (CD).



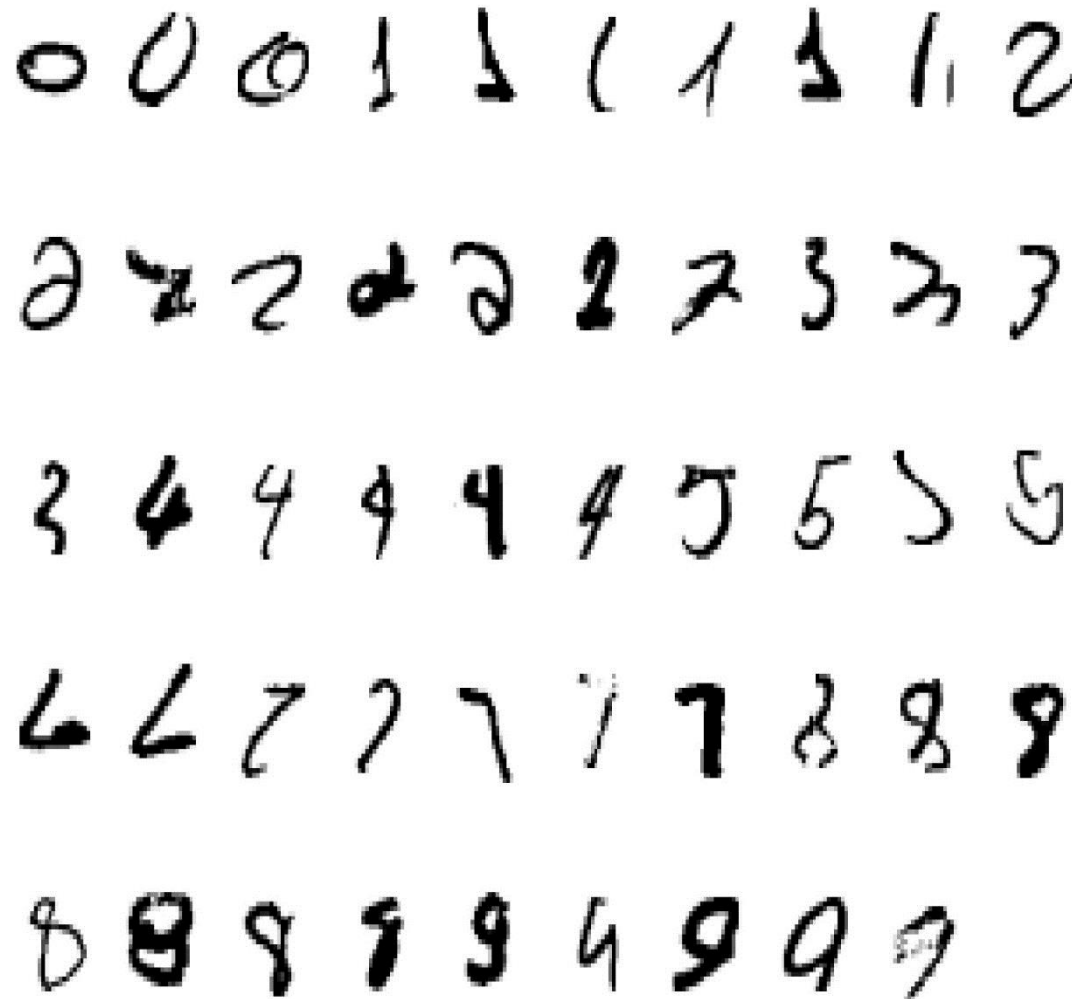
An RBM with fully connected visible and hidden units. Note there are no hidden-hidden or visible-visible connections.

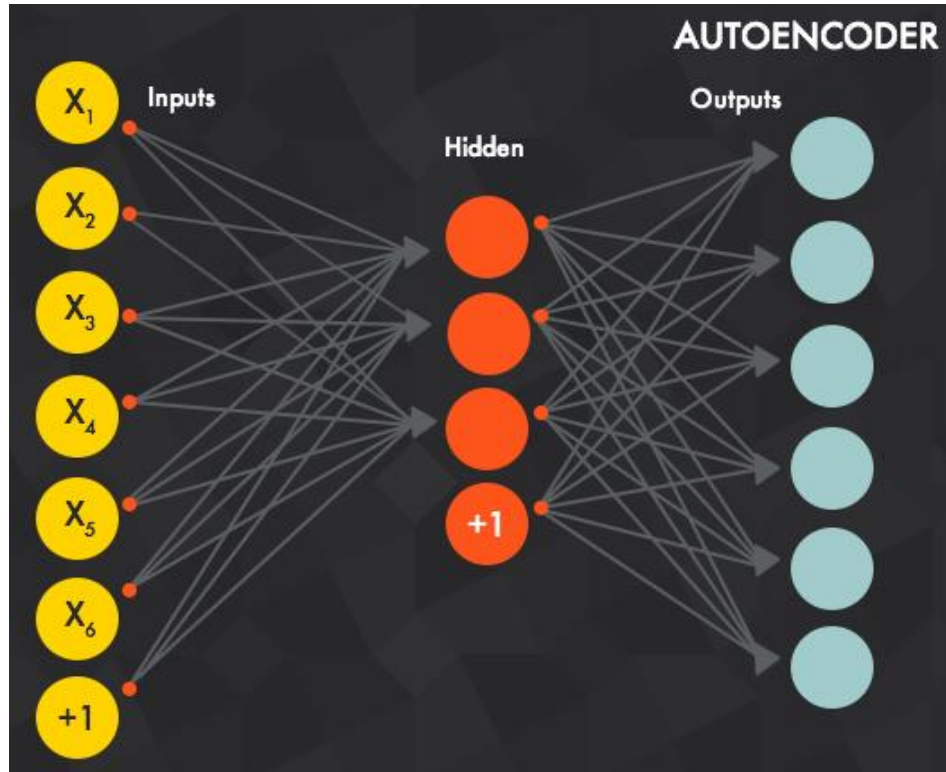
- CD provides an approximation to the maximum likelihood method that would ideally be applied for learning the weights of the RBM with gradient ascent:

$$\Delta w_{ij}(t+1) = w_{ij}(t) + \eta \frac{\partial \log(p(v))}{\partial w_{ij}}$$

- $p(v)$: the probability of a visible vector, given by $p(v) = \frac{\sum_h \exp(-E(v,h))}{Z}$.
- $E(v, h)$: the energy function assigned to the state of the network,
 - Given by $E(v, h) = -v^T \mathbf{W} h$
 - A lower energy indicates the network is in a more “desirable” configuration.
- $\frac{\partial \log(p(v))}{\partial w_{ij}}$ has the simple form $\langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{model}$.
 - $\langle \dots \rangle_p$ represent averages with respect to distribution p .
 - Intuition: fitting the data (reducing error)
- Challenge: sampling $\langle v_i h_j \rangle_{model}$ requires alternating Gibbs sampling for a long time.
- CD replaces this step by running alternating Gibbs sampling for n steps. After n steps, the data are sampled and that sample is used in place of $\langle v_i h_j \rangle_{model}$.

Input

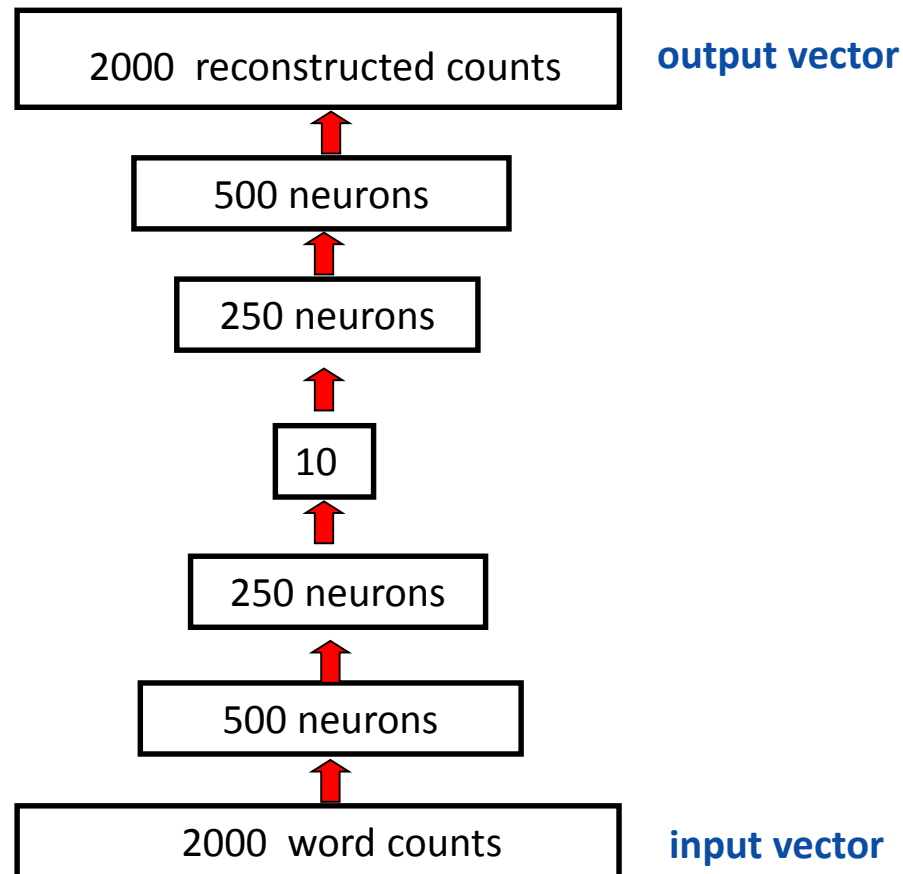




- The *auto encoder* idea is motivated by the concept of a good representation.
 - For example, for a classifier, a good representation can be defined as one that will yield a better performing classifier.
- An *encoder* is a deterministic mapping f_{θ} that transforms an input vector x into hidden representation y
 - $\theta = \{\mathbf{W}, b\}$, where \mathbf{W} is the weight matrix and b is bias (an offset vector)
- A *decoder* maps back the hidden representation y to the reconstructed input z via g_{θ} .
- Auto encoding: compare the reconstructed input z to the original input x and try to minimize this error to make z as close as possible to x .

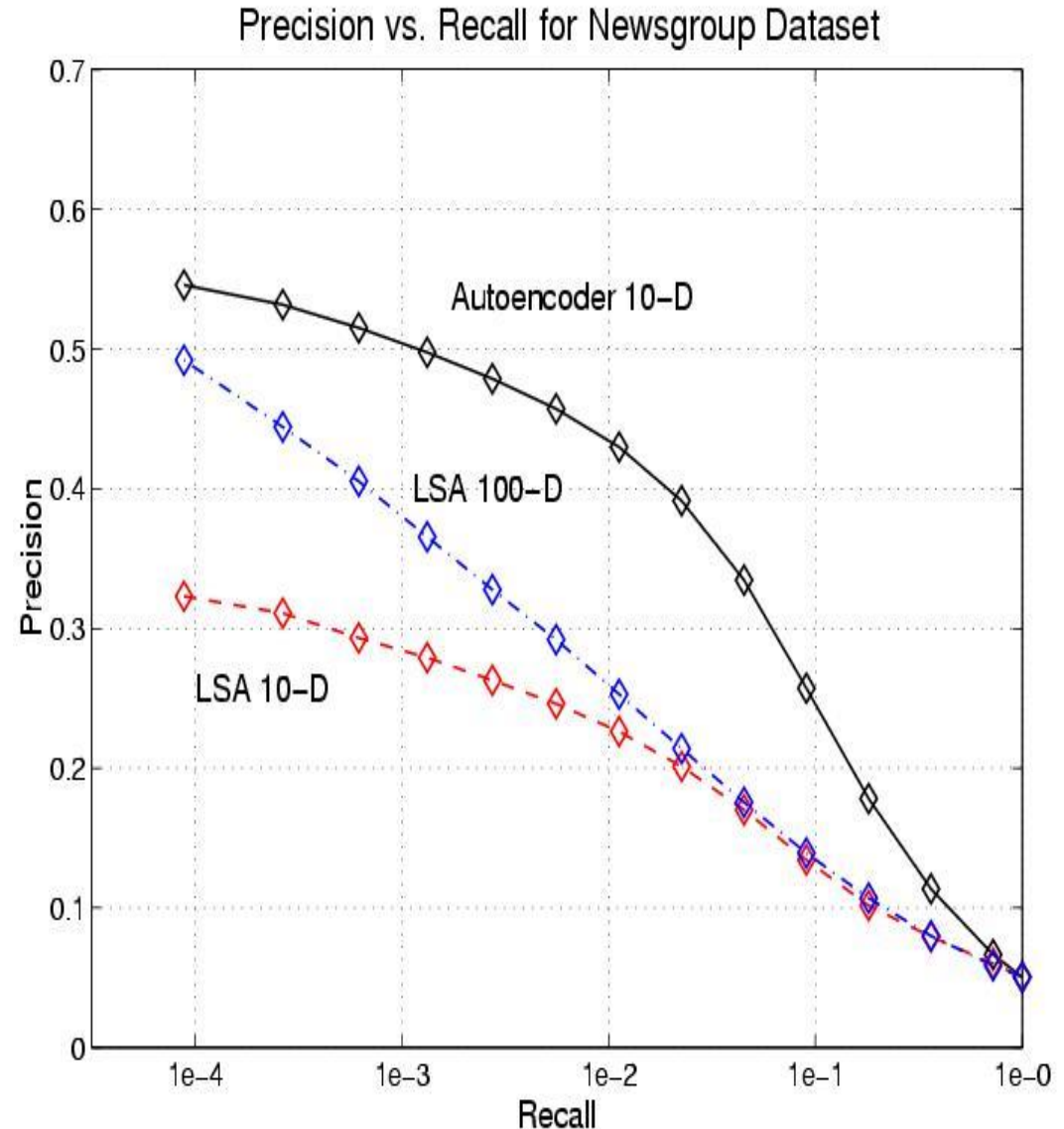
Example: Information Retrieval

- We can use an auto encoder to find low-dimensional codes for documents that allow fast and accurate retrieval of similar documents from a large set.
- We start by converting each document into a “bag of words”. This a 2000 dimensional vector that contains the counts for each of the 2000 commonest words.
- We train the neural network to reproduce its input vector as its output
- This forces it to compress as much information as possible into the 10 numbers in the central bottleneck.
- These 10 numbers are then a good way to compare documents.



- Train on bags of 2000 words for 400,000 training cases of business documents.
 - First train a stack of RBM's. Then fine-tune with backprop.
- Test on a separate 400,000 documents.
 - Pick one test document as a query. Rank order all the other test documents by using the cosine of the angle between codes.
 - Repeat this using each of the 400,000 test documents as the query (requires 0.16 trillion comparisons).
- Plot the number of retrieved documents against the proportion that are in the same hand-labeled class as the query document.

- The shortlist found using binary codes actually improves the precision-recall curves of TF-IDF.
 - Locality sensitive hashing (the fastest other method) is 50 times slower and has worse precision-recall curves.
 - Vs. Latent Semantic Analysis (LSA)

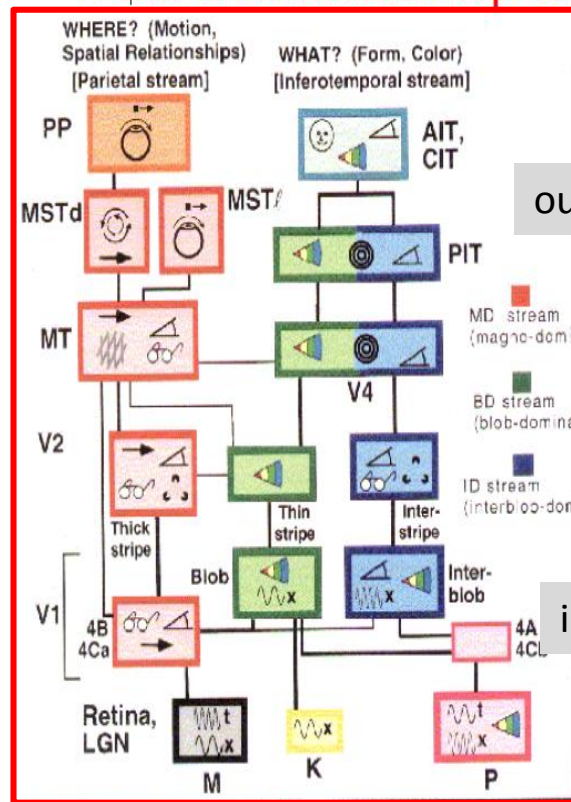
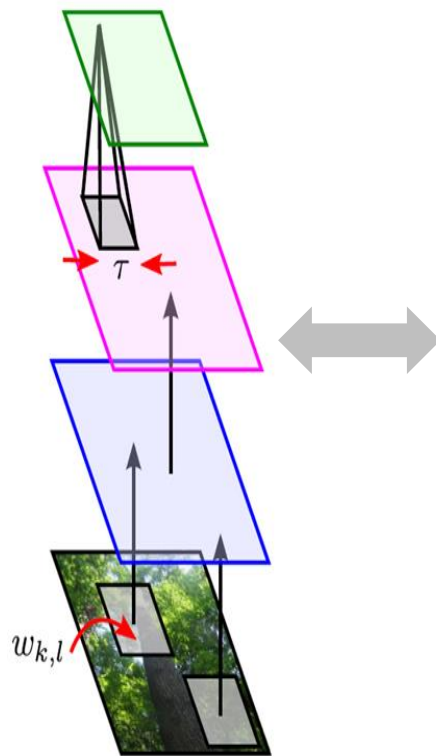


- Convolutional Neural Networks are inspired by mammalian visual cortex.
 - The visual cortex contains a complex arrangement of cells, which are sensitive to small sub-regions of the visual field, called a receptive field. These cells act as local filters over the input space and are well-suited to exploit the strong spatially local correlation present in natural images.
 - Two basic cell types:
 - Simple cells respond maximally to specific edge-like patterns within their receptive field.
 - Complex cells have larger receptive fields and are locally invariant to the exact position of the pattern.

The Mammalian Visual Cortex Inspires CNN

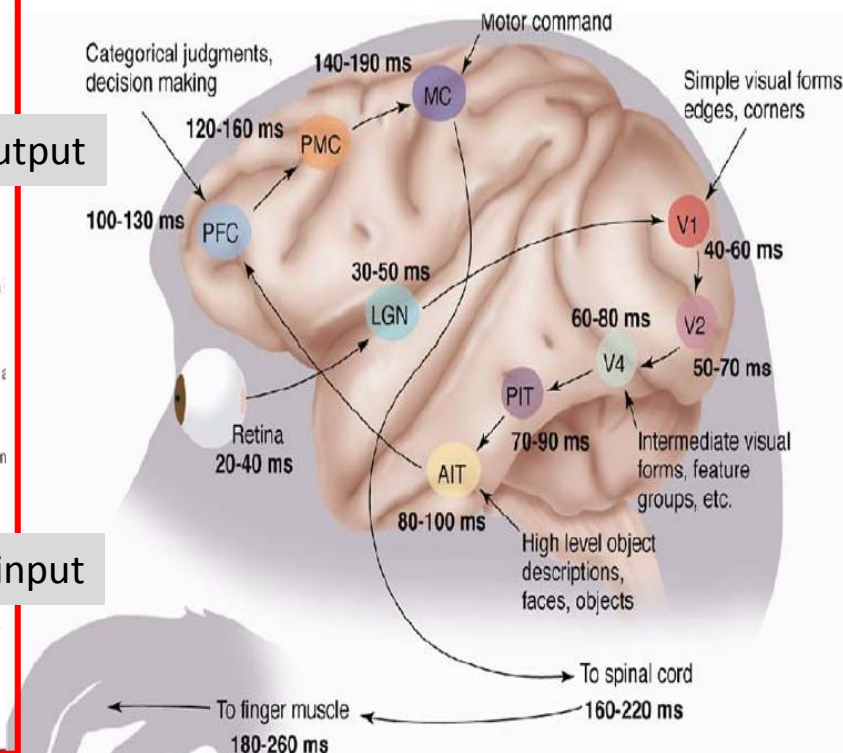
Convolutional Neural Net

- The ventral (recognition) pathway in the visual cortex has multiple stages
- Retina - LGN - V1 - V2 - V4 - PIT - AIT
- Lots of intermediate representations



output

input

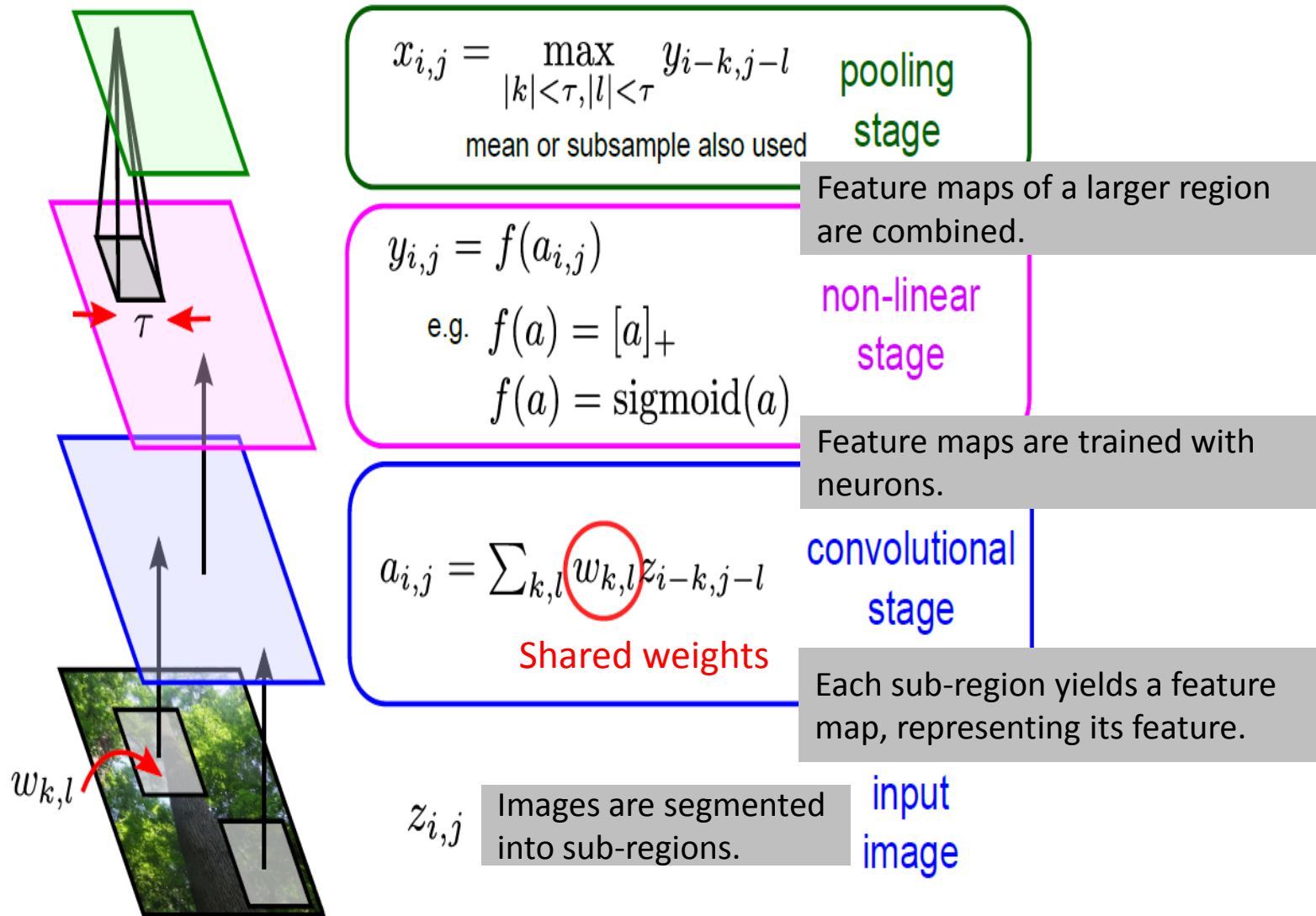


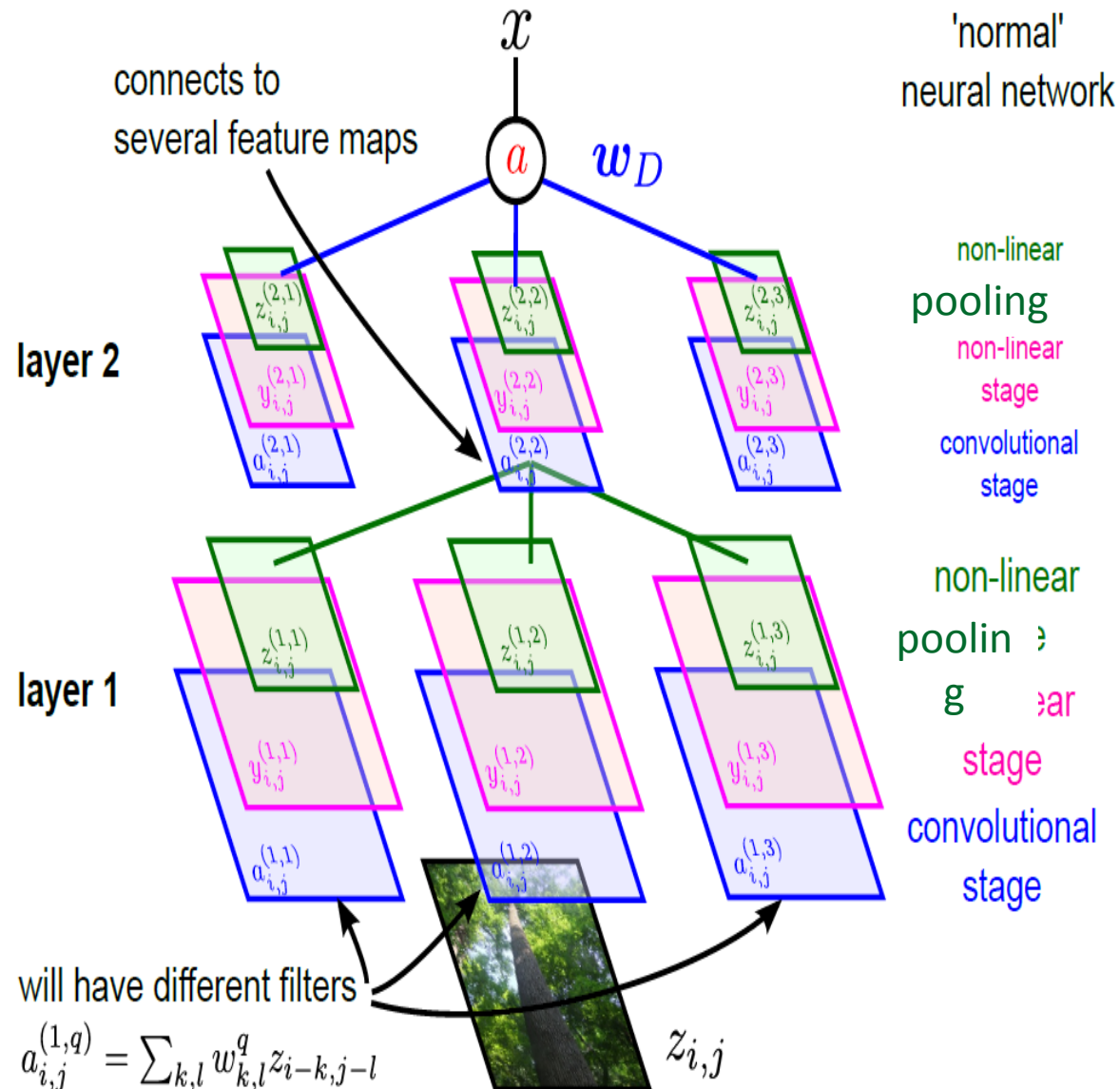
[picture from Simon Thorpe]

[Gallant & Van Essen]

- Intuition: Neural network with specialized connectivity structure,
 - Stacking multiple layers of feature extractors
 - Low-level layers extract local features.
 - High-level layers extract learn global patterns.
- A CNN is a list of layers that transform the input data into an output class/prediction.
- There are a few distinct types of layers:
 - Convolutional layer
 - Non-linear layer
 - Pooling layer

Building-blocks for CNN's





Building deep learning models on textual data requires:

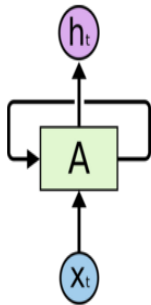
- Representation of the basic text unit, word.
- Neural network structure that can hierarchically capture the sequential nature of text.

Deep learning models for text mining use:

- Vector representation of words (i.e., word embedding)
- Neural network structures
 - *Recurrent* Neural Network
 - *Recursive* Neural Network

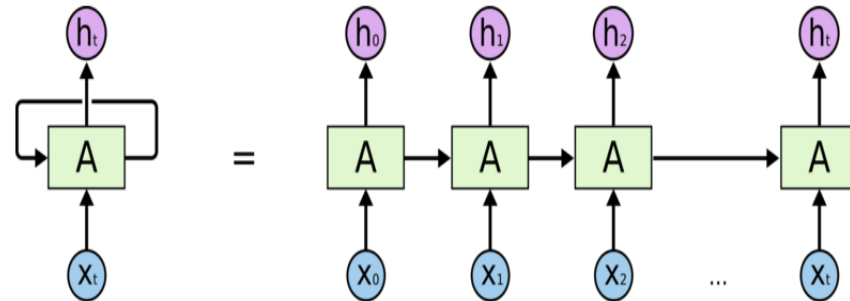
- The applications of standard Neural Networks (and also Convolutional Networks) are limited due to:
 - They only accepted a fixed-size vector as input (e.g., an image) and produce a fixed-size vector as output (e.g., probabilities of different classes).
 - These models use a fixed amount of computational steps (e.g. the number of layers in the model).
- Recurrent Neural Networks are unique as they allow us to operate over sequences of vectors.
 - Sequences in the input, the output, or in the most general case both

- Recurrent Neural Networks are networks with loops in them, allowing information to persist.



Recurrent Neural Networks have loops.

In the above diagram, a chunk of neural network, **A**, looks at some input x_t and outputs a value h_t . A loop allows information to be passed from one step of the network to the next.

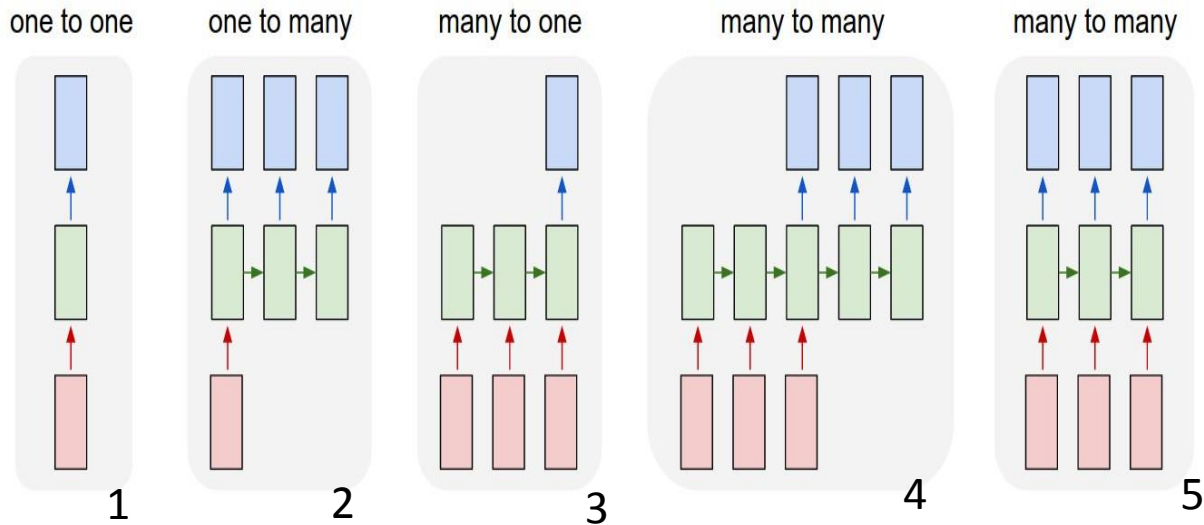


An unrolled recurrent neural network.

A recurrent neural network can be thought of as multiple copies of the same network, each passing a message to a successor. The diagram above shows what happens if we unroll the loop.

- Intuition of Recurrent Neural Networks
 - Human thoughts have persistence; humans don't start their thinking from scratch every second.
 - As you read this sentence, you understand each word based on your understanding of previous words.
 - One of the appeals of RNNs is the idea that they are able to connect previous information to the present task
 - E.g., using previous video frames to inform the understanding of the present frame.
 - E.g., a language model tries to predict the next word based on the previous ones.

- Examples of Recurrent Neural Networks



- Each rectangle is a vector and arrows represent functions (e.g. matrix multiply).
- Input vectors are in red, output vectors are in blue and green vectors hold the RNN's state

(1) Standard mode of processing without RNN, from fixed-sized input to fixed-sized output (e.g. image classification).

(2) Sequence output (e.g. image captioning takes an image and outputs a sentence of words).

(3) Sequence input (e.g. sentiment analysis where a given sentence is classified as expressing positive or negative sentiment).

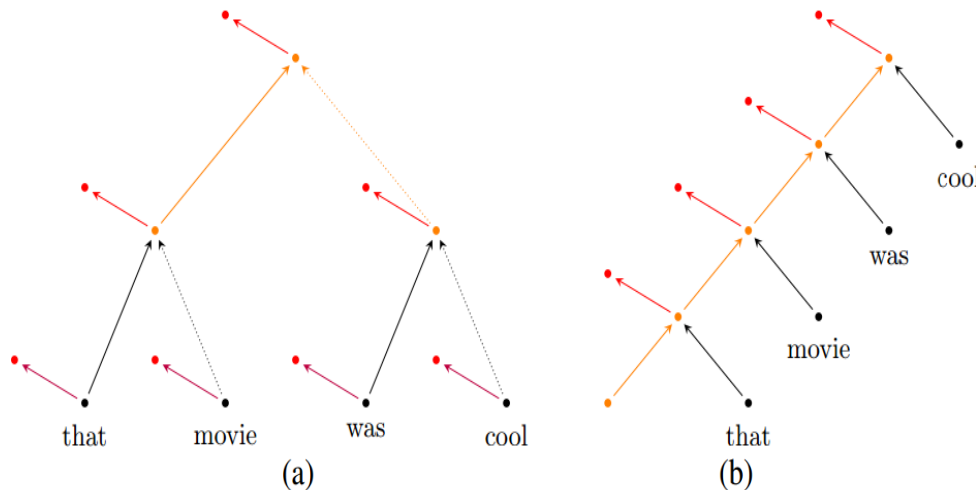
(4) Sequence input and sequence output (e.g. Machine Translation: an RNN reads a sentence in English and then outputs a sentence in French).

(5) Synced sequence input and output (e.g. video classification where we wish to label each frame of the video).

- Incredible success applying RNNs to language modeling and sequence learning problems

Task	Input Sequence	Output Sequence
Machine translation (Sutskever et al. 2014)	English	French
Question answering (Bordes et al. 2014)	Question	Answer
Speech recognition (Graves et al. 2013)	Voice	Text
Handwriting prediction (Graves 2013)	Handwriting	Text
Opinion mining (Irsoy et al. 2014)	Text	Opinion expression

- A recursive NN can be seen as a generalization of the recurrent NN
 - A recurrent neural network is in fact a recursive neural network with the structure of a linear chain
 - Recursive NNs operate on hierarchical structure, Recurrent NN operate on progression of time

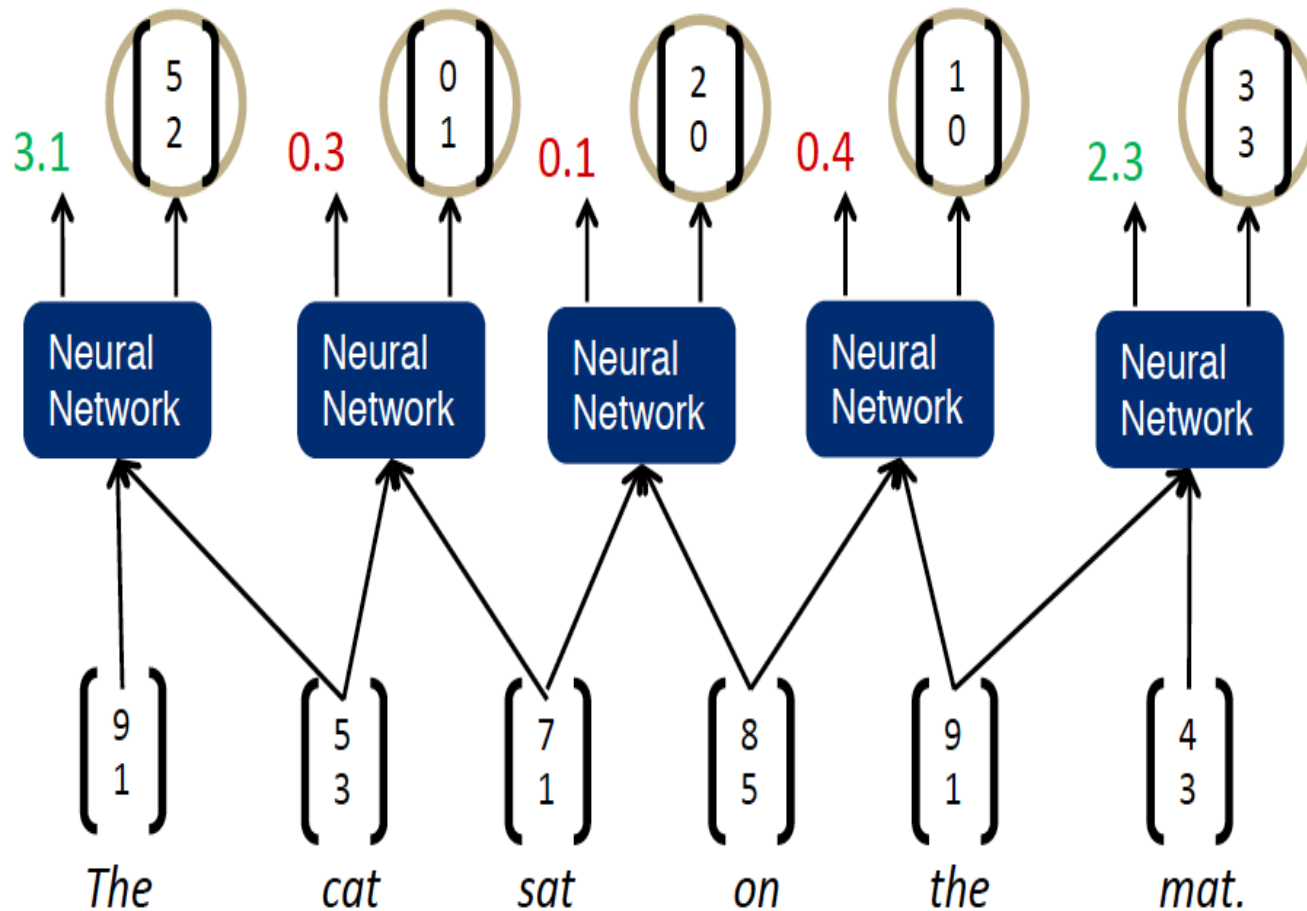


- Operation of a recursive net (a), and a recurrent net (b) on an example sentence. Note the linear chain in (b)
- Black, orange and red dots represent input, hidden and output layers, respectively.
- Directed edges having the same color-style combination denote shared connections.

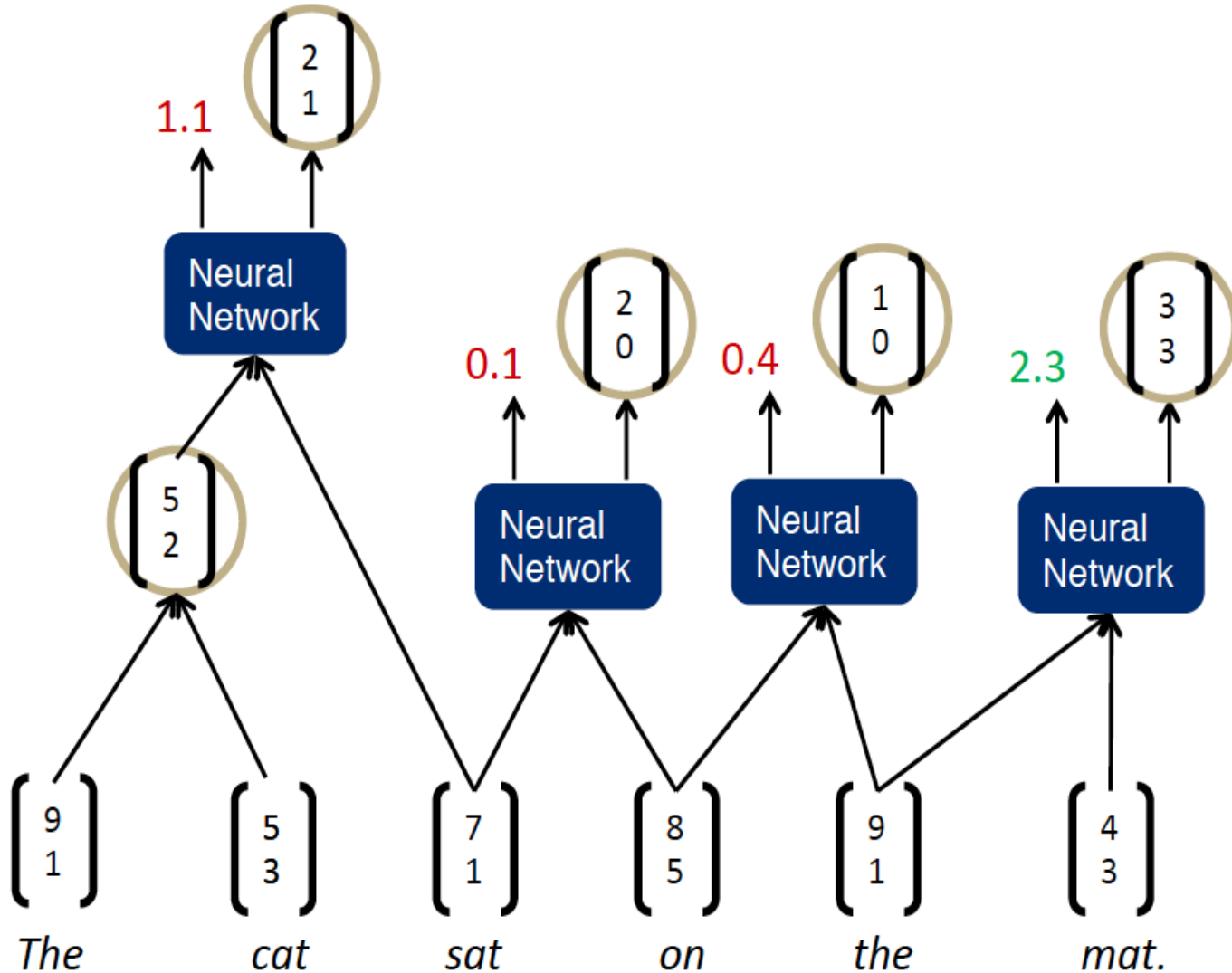
- Applied to parsing, sentence-level sentiment analysis, and paraphrase detection.

- Given the structural representation of a sentence, e.g. a parse tree:
 - Recursive Neural Networks recursively generate parent representations in a bottom-up fashion,
 - By combining tokens to produce representations for phrases, eventually producing the whole sentence.
 - The sentence-level representation (or, alternatively, its phrases) can then be used to make a final classification for a given input sentence — e.g. whether it conveys a positive or a negative sentiment.

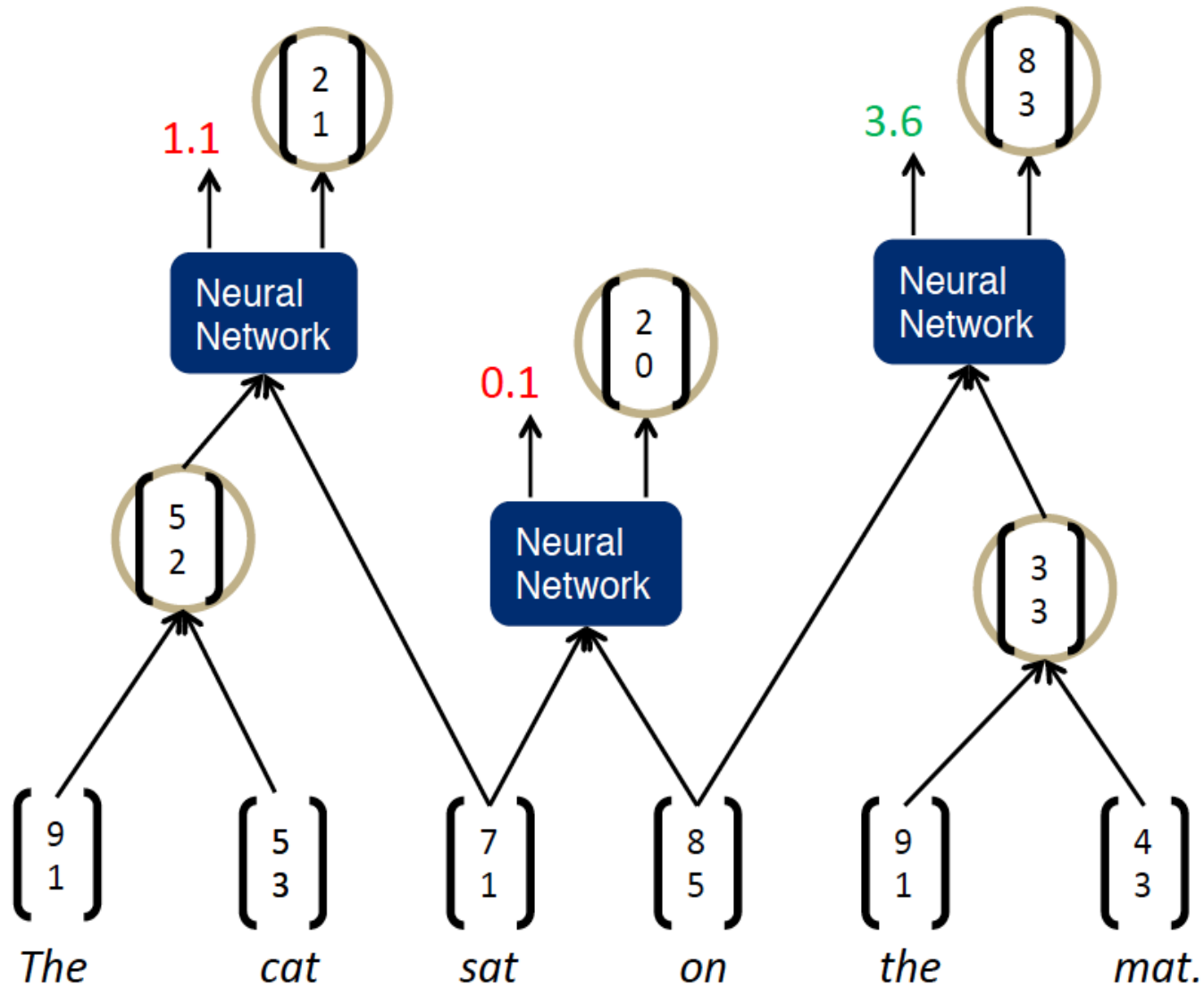
Recursive Neural Network Illustration (Animation)



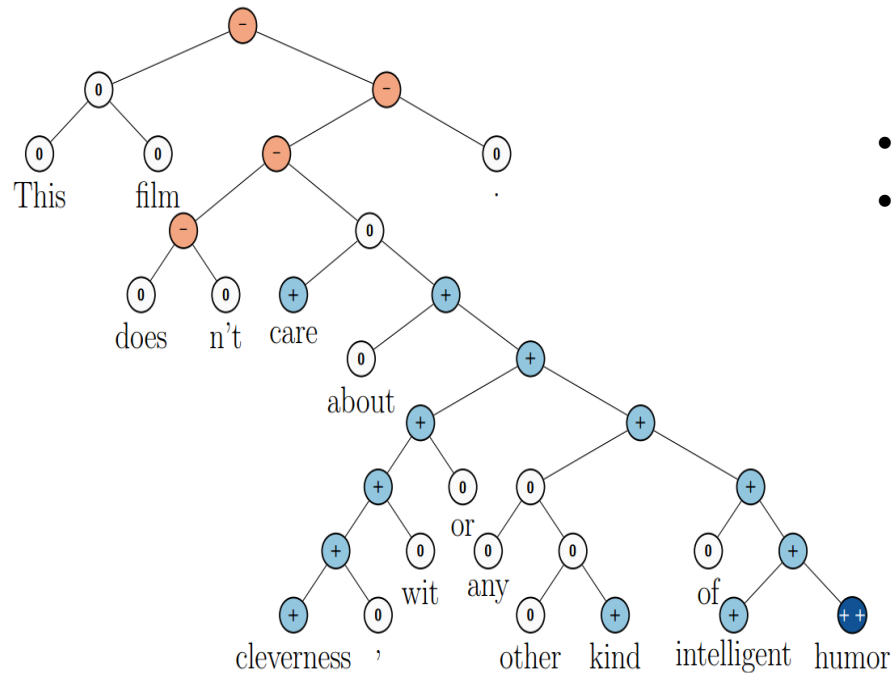
Recursive Neural Network Illustration (Animation)



Recursive Neural Network Illustration (Animation)



- Recursive neural networks have had significant successes in a number of NLP tasks.
 - Socher *et al.* (2013) uses a recursive neural network to predict sentence sentiment:



- Fine-grained sentiment analysis
- Allocate positive and negative sentiment from words to phrases and to the entire sentence

- SemEval 2014: Predicting semantic relatedness of two sentences (Tai et al. 2015)

Ranking by mean word vector cosine similarity	Score	Ranking by Dependency Tree LSTM model	Score
a woman is slicing potatoes		a woman is slicing potatoes	
a woman is cutting potatoes	0.96	a woman is cutting potatoes	4.82
a woman is slicing herbs	0.92	potatoes are being sliced by a woman	4.70
a woman is slicing tofu	0.92	tofu is being sliced by a woman	4.39
a boy is waving at some young runners from the ocean		a boy is waving at some young runners from the ocean	
a man and a boy are standing at the bottom of some stairs , which are outdoors	0.92	a group of men is playing with a ball on the beach	3.79
a group of children in uniforms is standing at a gate and one is kissing the mother	0.90	a young boy wearing a red swimsuit is jumping out of a blue kiddies pool	3.37
a group of children in uniforms is standing at a gate and there is no one kissing the mother	0.90	the man is tossing a kid into the swimming pool that is near the ocean	3.19
two men are playing guitar		two men are playing guitar	
some men are playing rugby	0.88	the man is singing and playing the guitar	4.08
two men are talking	0.87	the man is opening the guitar for donations and plays with the case	4.01
two dogs are playing with each other	0.87	two men are dancing and singing in front of a crowd	4.00



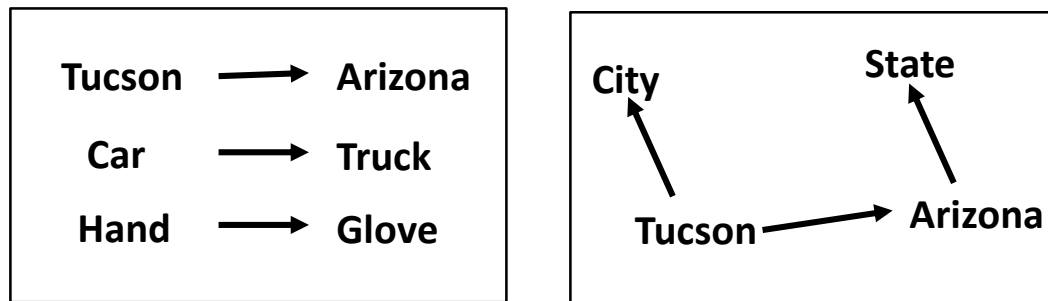
Recursive NN based model can pick up more subtle semantic relatedness in sentences compared to word vector model, e.g., ocean and beach.

Table 4: Most similar sentences from a 1000-sentence sample drawn from the SICK test set. The Tree-LSTM model is able to pick up on more subtle relationships, such as that between “beach” and “ocean” in the second example.

- Vector space models (VSMs) represent (embed) words in a continuous vector space
 - Theoretical foundation in Linguistics: Distributional Hypothesis
 - Words with similar meanings will occur with similar neighbors if enough text material is available (Rubenstein et al. 1967).
- Approaches that leverage VSMs can be divided into two categories

Approach	Example	Description
Count-based methods	Latent semantic analysis	Compute how often some word co-occurs with its neighbor words in a large text corpus, and then map these count-statistics down to a small, dense vector for each word
Predictive methods	Neural probabilistic language model	Directly predict a word from its neighbors in terms of learned small, dense embedding vectors (considered parameters of the model)

- Word2vec: computationally-efficient, 2-layer predictive NN for learning word embedding from raw text
 - Considered deep for its ability to digest expansive data sets quickly
- Can be used for unsupervised learning of words
 - Relationships between different words
 - Ability to abstract higher meaning between words (e.g., Tucson is a city in the state of Arizona)



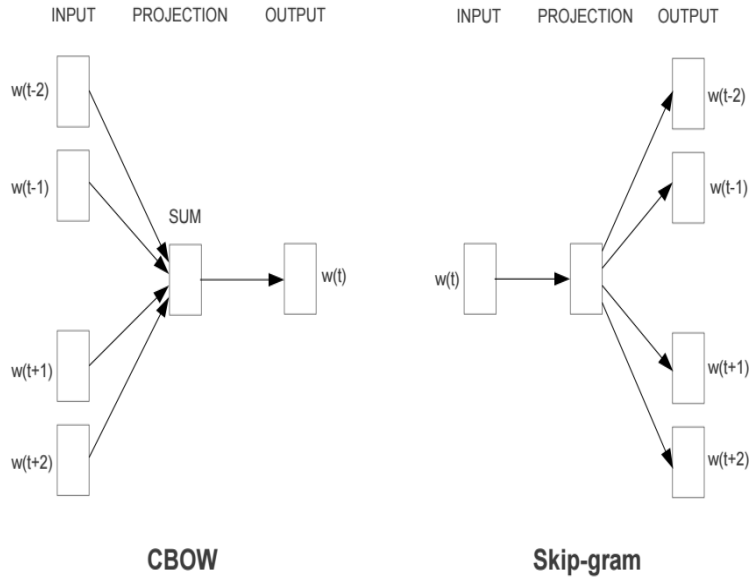
- Useful for language modeling, sentiment analysis, and more

- Its input is a text corpus and its output is a set of vectors or “embeddings” (feature vectors for words in that corpus)
 - Similarity between two embeddings represents conceptual similarity of words
- Example results: words associated with *Sweden*, in order of proximity:

Word	Cosine distance

norway	0.760124
denmark	0.715460
finland	0.620022
switzerland	0.588132
belgium	0.585835
netherlands	0.574631
iceland	0.562368
estonia	0.547621
slovenia	0.531408

- Word2vec comes with two models



Model	Approach	Speed and Performance	Use case
Continuous Bag-of-Words model (CBOW)	The CBOW predicts the current word based on the context.	Faster to train than the skip-gram model	Predicts frequent words better
Skip-Gram model	Skip-gram predicts surrounding words given the current word.	Usually performs better than CBOW	Predicts rare words better

- Skip-gram learning:
 - Given w_0 , predict w_{-2} , w_{-1} , w_1 , and w_2

w_{-2}	w_{-1}	w_0	w_1	w_2
?	?	Network	?	?



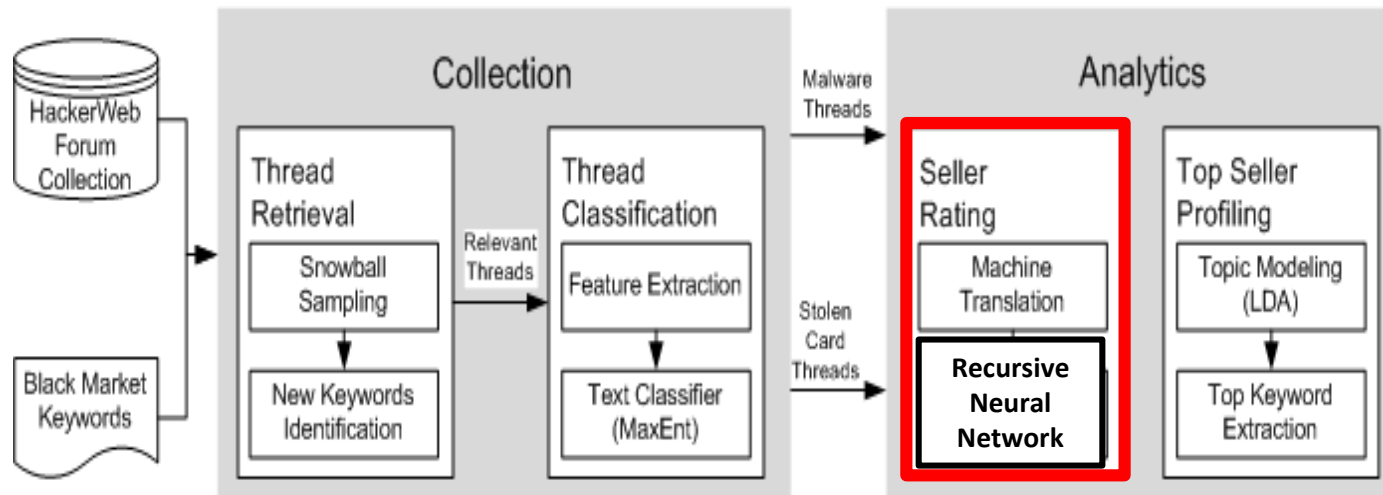
- Conversely, CBOW tries to predict w_0 when given w_{-2} , w_{-1} , w_1 , and w_2

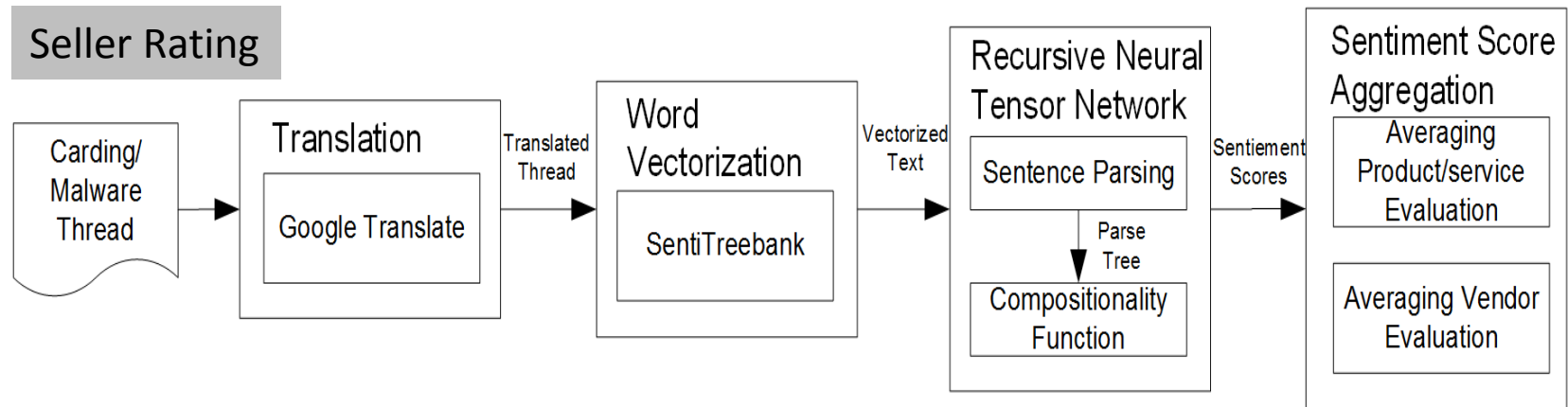
w_{-2}	w_{-1}	w_0	w_1	w_2
Recurrent	Neural		Language	Model

- The carding community rely heavily on the online black market for exchanging malwares and stolen data.
- Online black market for carders:
 - *Sellers* advertise their malwares and/or stolen data by posting a thread in the carding community
 - *Buyers* leave feedback commenting the their evaluation of the seller.
- **Objective:** to identify key carding sellers through evaluating buyers feedback.

Recursive Neural Network Example – Identifying Key Carding Sellers

- Input: Original threads from hacker forums
- Preprocessing:
 - Thread Retrieval: Identifying threads related to the underground economy by conducting snowball sampling-based keywords search
 - Thread Classification: Identifying advertisement threads using MaxEnt classifier
 - Focusing on malware advertisements and stolen card advertisement
 - Can be generalized to other advertisements.

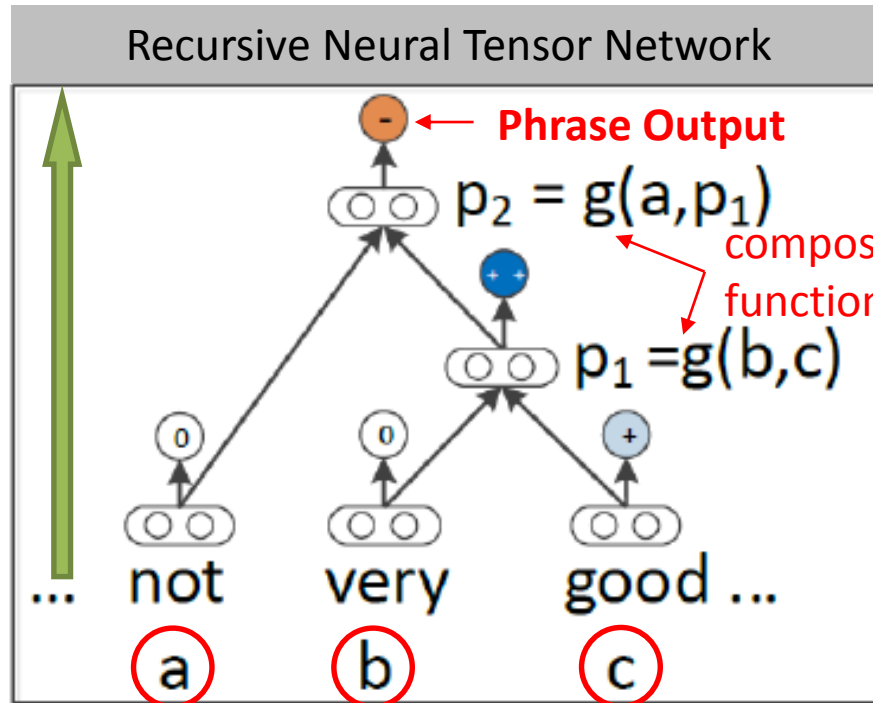




- *Translation* translates the thread content from the original language to English.
- *Word Vectorization* vectorizes each word in the text into a five dimensional vector using SentiTreeBank, a dictionary for customer feedback(Socher et al., 2013).
 - *SentiTreebank*: corpus with fully labeled parse trees with sentiment on each level
- *Recursive Neural Tensor Network* parses the sentence into a binary tree and aggregate semantics from constituents recursively.
- *Sentiment Score Aggregation* averages the sentiment measures for each advertisement and seller.

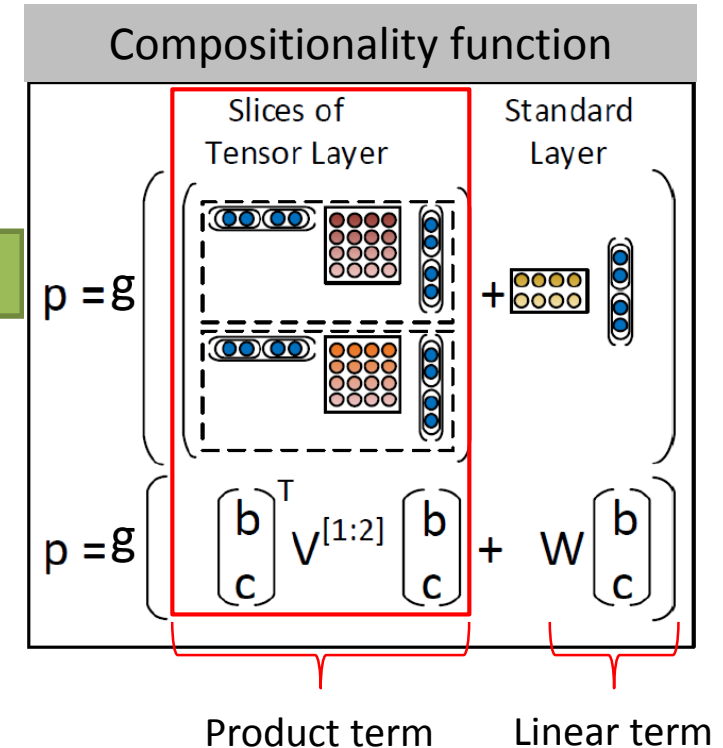
Recursive Neural Network Example – Identifying Key Carding Sellers

Example feedback: “*The dump is not very good*”



Word embeddings trained by SentiTreeBank

Output of the example feedback: **Negative**



The additional *tensor layer* makes the compositional function more powerful by introducing the product of the constituent word embeddings.

Top 3 Best/Worst Malware and Stolen Data Sellers for Each Forum								
	Top 3 Best				Top 3 Worst			
	Malware		Stolen Data		Malware		Stolen Data	
Rank	User	Score	User	Score	User	Score	User	Score
	Antichat							
1	L**G	5	i**o	3.6	N**g	1.8	I**s	2.3
2	V**U	4.5	a**s	3.5	k**a	2	P**A	2.3
3	g**l	4	D**R	3.4	D**i	2	s**8	2.4
	CrdPro							
1	H**l	4	R**r	4.4	N**0	2	f**4	1.3
2	b**t	4	F**x	4	1**4	2	s**3	1.3
3	S**r	4	R**c	4	M**D	2	l**u	1.3
	Zloy							
1	P**t	5	B**r	4	r**t	1.5	r**y	1
2	D**n	4	B**1	4	w**0	1.6	m**n	1
3	D**f	4	s**c	4	g**n	2	j**a	2

COMPUTER VISION

May 21,
2019

