

Deep Learning for Natural Language Processing

Dr. Minlie Huang (黃民烈)

aihuang@tsinghua.edu.cn

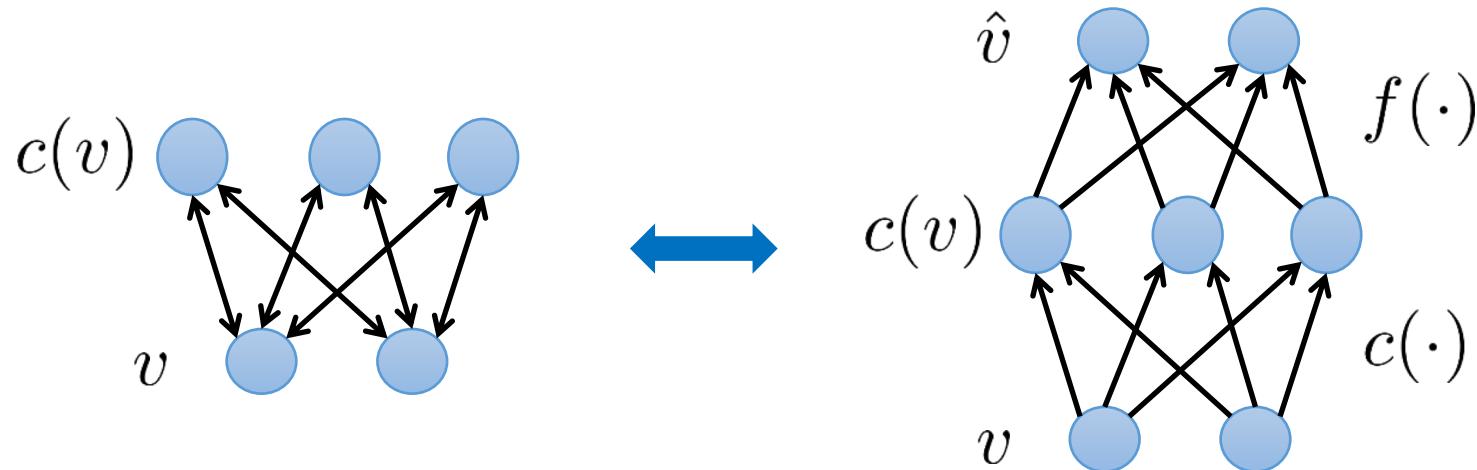
Computer Science Department

Tsinghua University

Homepage: <http://coai.cs.Tsinghua.edu.cn/hml>

Recursive Networks for NLP

Auto-encoder



- Encode the input v into some representation $c(v)$ so that the input can be reconstructed from that representation
 - Encoding function $c(v)$
 - Decoding function $f(c(v))$

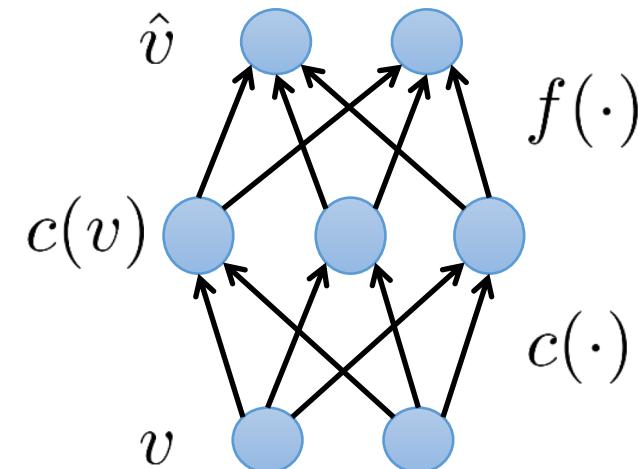
Encoding and decoding functions

- Nonlinear function

$$c(v) = \text{sigmoid}(W_1 v + \theta)$$

$$f(c) = \text{sigmoid}(W_2 c + \eta)$$

- If c and f are binary, then the functions can be used as probabilities



Loss function

- Minimize the reconstruction error or **the negative data log-likelihood**

$$RE = -\langle \ln P(v|c(v)) \rangle \quad \langle \cdot \rangle: \text{average over samples}$$

- Gaussian probability (v is real)

$$P(v|c(v)) \propto \exp\left(\frac{-\|v - f(c(v))\|^2}{2\sigma^2}\right)$$

then

$$RE = \langle \|v - f(c(v))\|^2 \rangle$$

- Binomial probability (v is binary)

$$P(v|c(v)) \propto \prod_i f_i(c(v))^{v_i} (1 - f_i(c(v)))^{1-v_i}$$

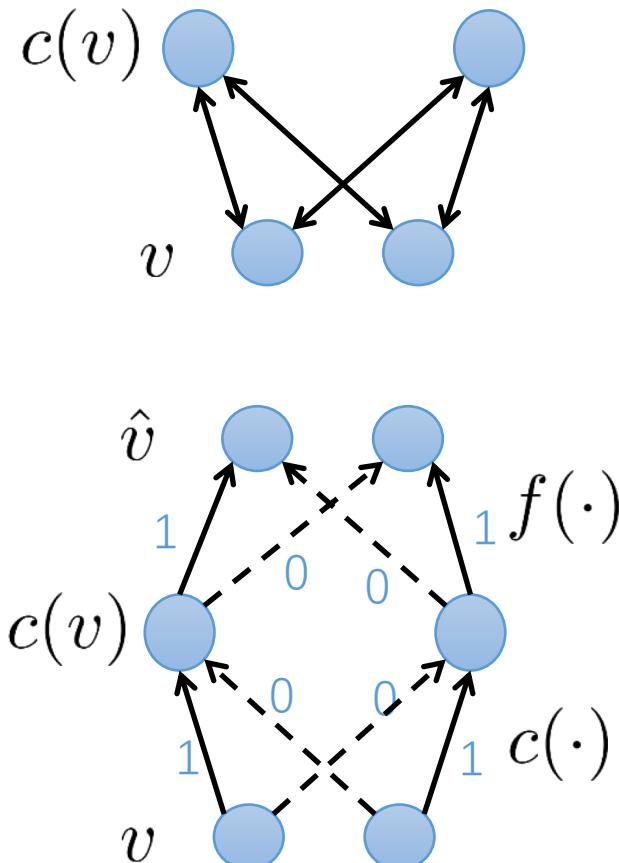
then

$$RE = -\langle \sum_i (v_i \ln f_i(c(v)) + (1 - v_i) \ln(1 - f_i(c(v)))) \rangle$$

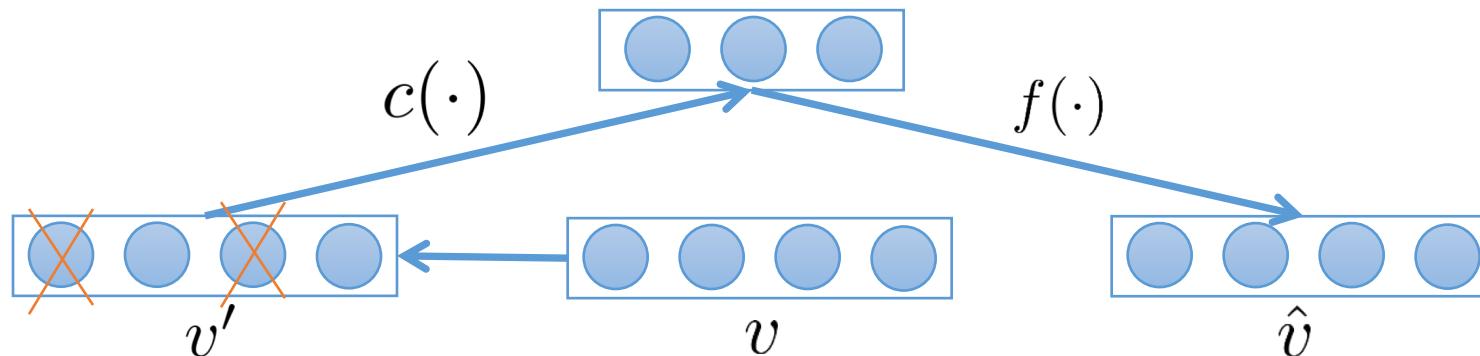
A trivial solution

- In the following case
 - Binary units
 - The number of hidden units is equal to the number of visible units
- There is a trivial solution $W_1 = W_2 = I, \eta = \theta = -0.5$

0	-0.5	0	-0.5	0
1	0.5	1	0.5	1



Denoising auto-encoder

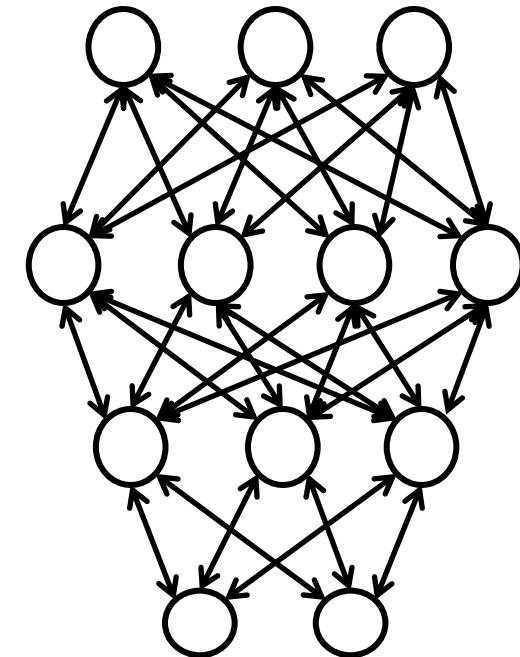


- Corrupt v to v' by randomly setting some elements of v to zero.
- Use v' as input and try to reconstruct v .
 - Ideally, \hat{v} is the clean version of v

Vincent, et al., ICML, 2008

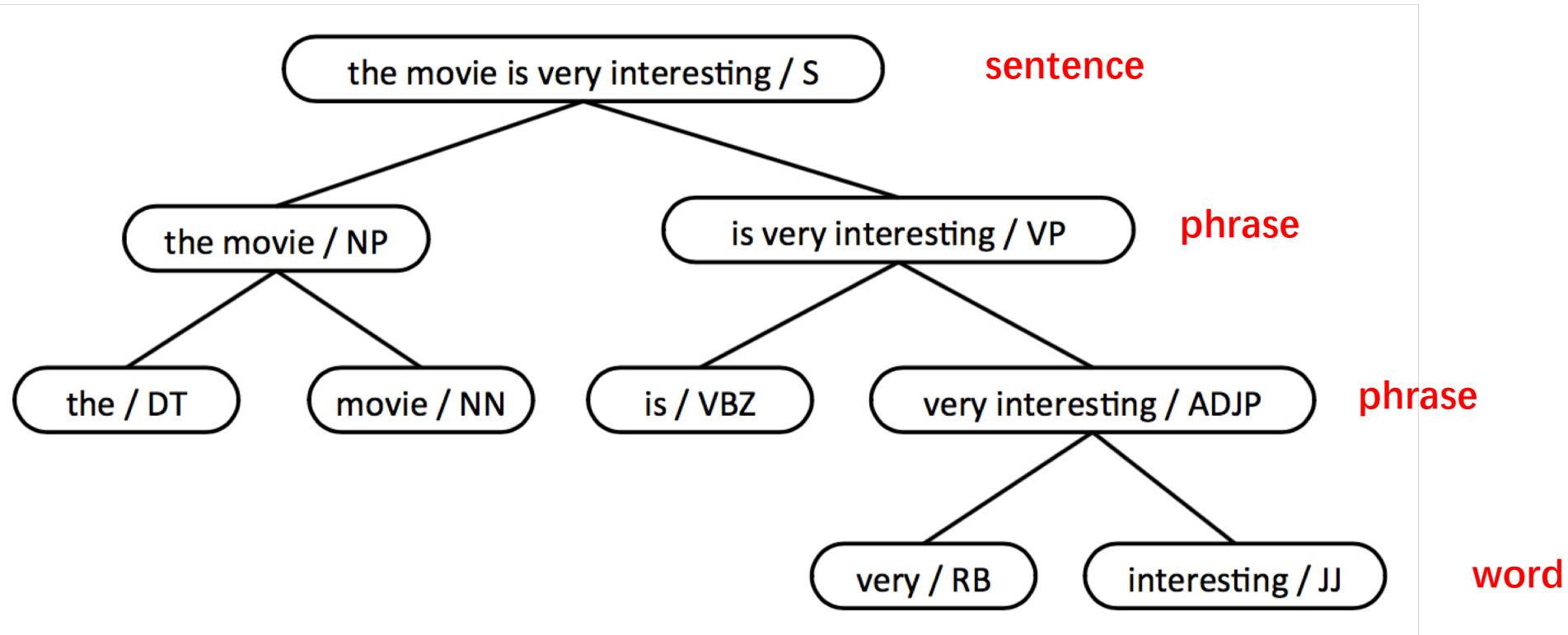
Deep Auto-encoder

- Stack auto-encoders on top of each other
- Train layers one by one
- Fine tune with BP
- Sparsity or other regularizations can be used



Recursive Autoencoders

Sentences/phrases have composition structures!

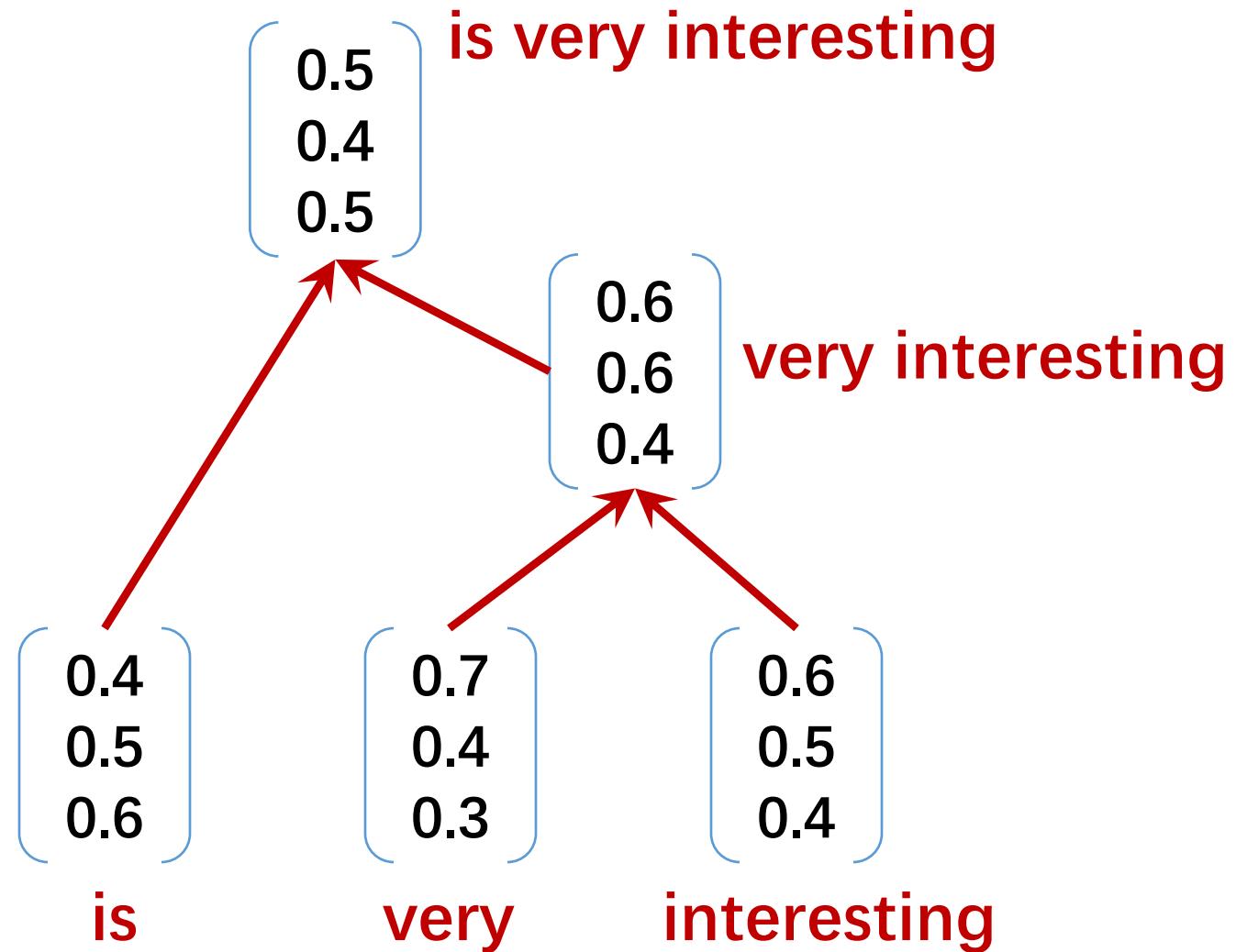


Recursive Autoencoders

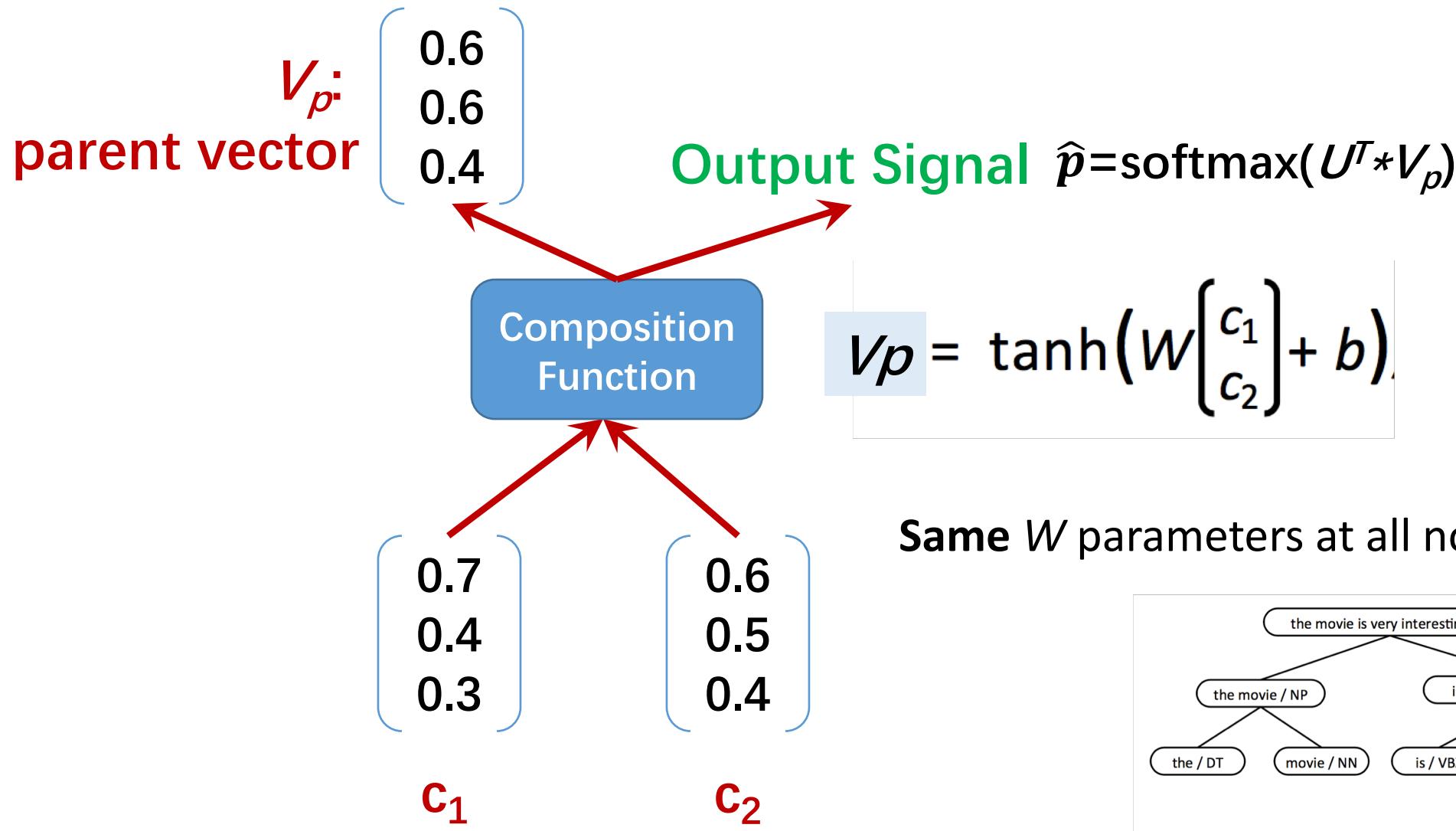
Rules of Compositionality

The meaning (vector) of a sentence is determined by

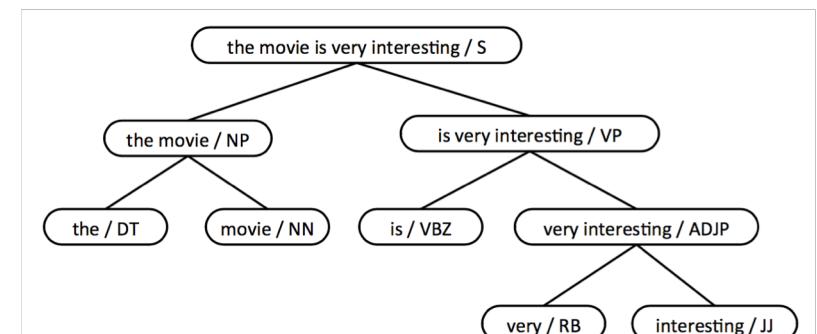
- (1) The meanings of its words
- (2) The rules that combine them



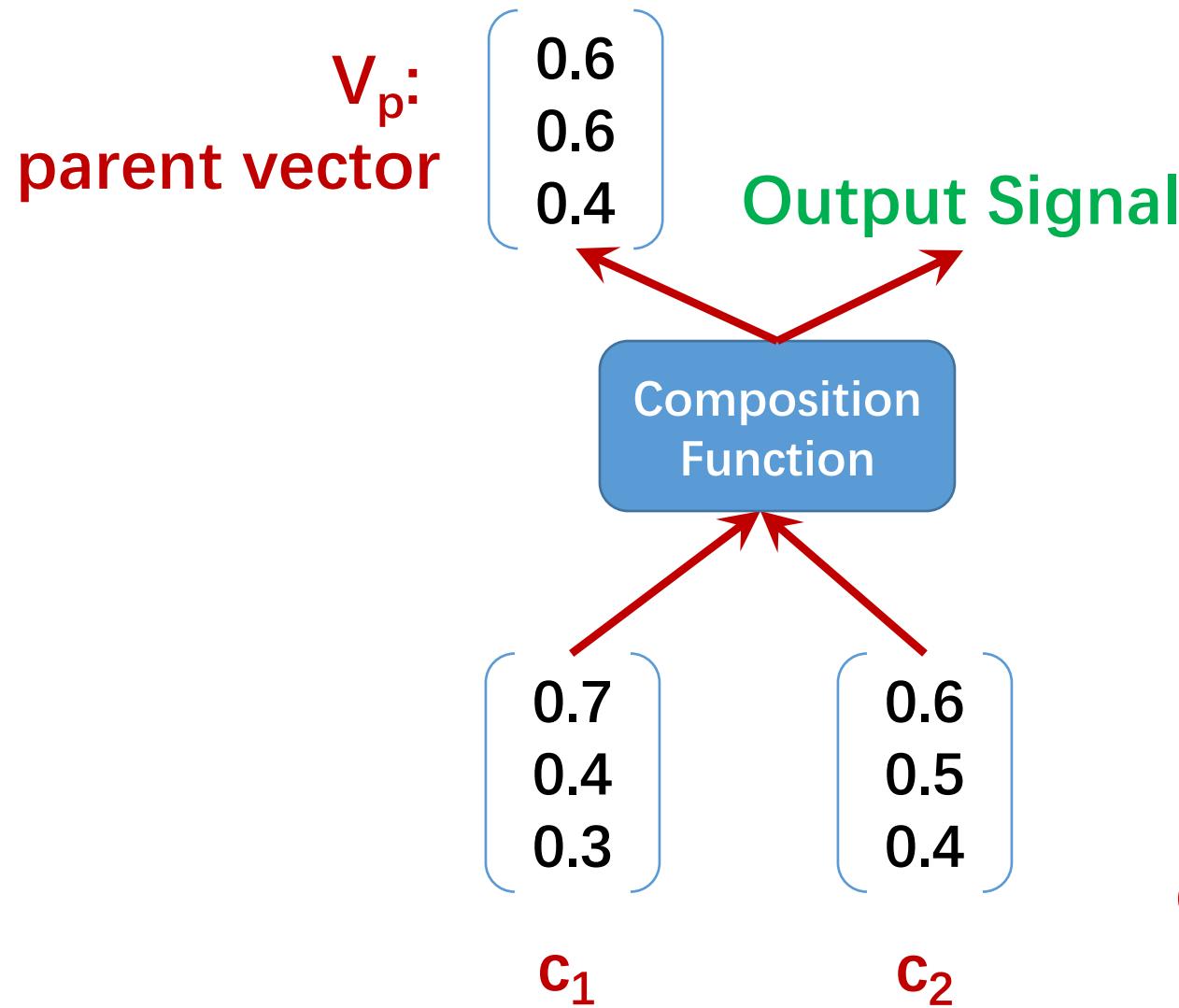
Recursive Autoencoders



Same W parameters at all nodes of the tree



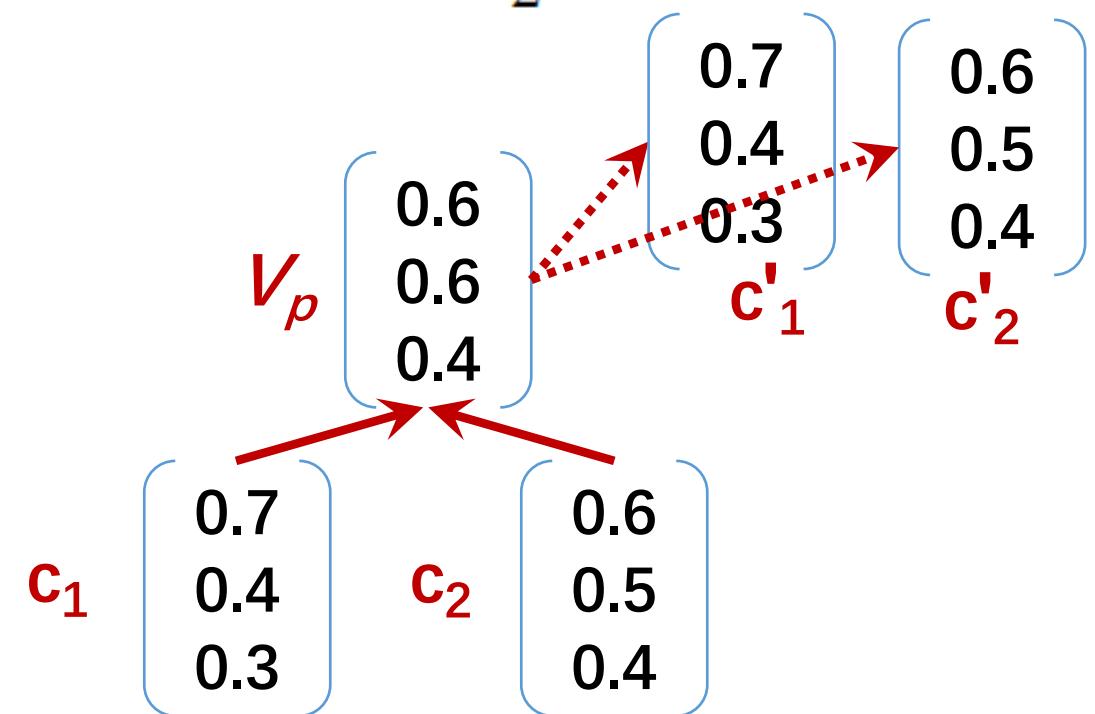
Recursive Autoencoders



Train the model

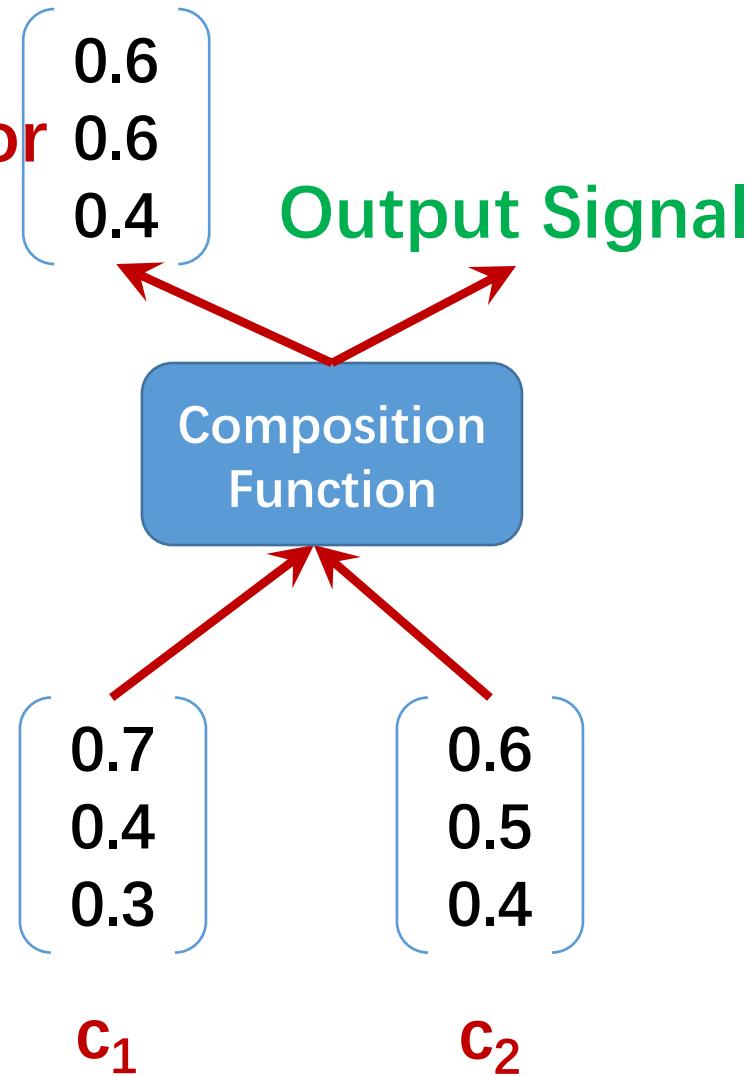
- ✓ No supervision: minimizing reconstruction error

$$E_{rec}([c_1; c_2]) = \frac{1}{2} \|[c_1; c_2] - [c'_1; c'_2]\|^2$$



Recursive Autoencoders

V_p : parent vector



Train the model

- ✓ With supervision:
minimizing cross entropy error

$$E = \sum_{k=1}^K -p(k) \log \hat{p}(k)$$

The gold distribution over class labels k

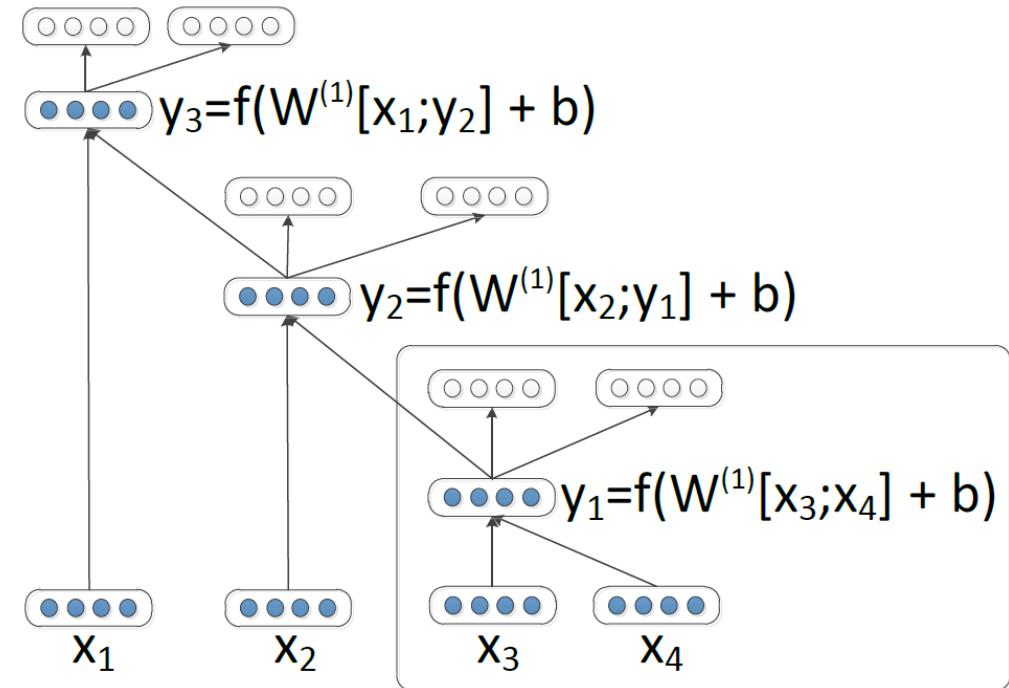
The predicted distribution based on
the parent vector V_p

$$\hat{p} = \text{softmax}(U^T * V_p)$$

Recursive Autoencoder

- Given the tree structure, we can compute all the node vectors from bottom to top.
- Train by minimizing reconstruction error

$$E_{rec}([c_1; c_2]) = \frac{1}{2} \|[c_1; c_2] - [c'_1; c'_2]\|^2$$



Semi-Supervised Recursive Autoencoders

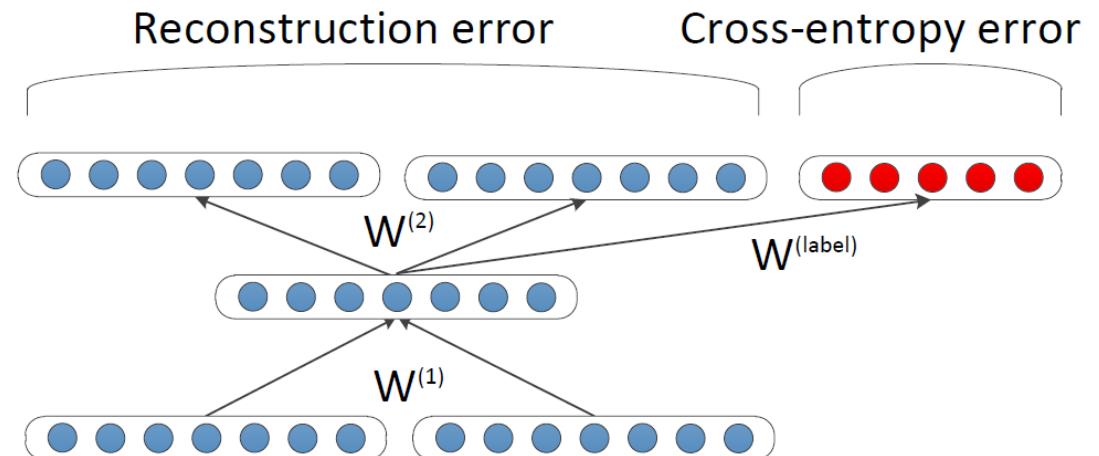
- The task label can be introduced in all nodes of the tree.

$$d(p; \theta) = \text{softmax}(W^{label} p)$$

$$E_{cE}(p, t; \theta) = - \sum_{k=1}^K t_k \log d_k(p; \theta)$$

$$E([c_1; c_2]_s, p_s, t, \theta) =$$

$$\alpha E_{rec}([c_1; c_2]_s; \theta) + (1 - \alpha) E_{cE}(p_s, t; \theta)$$



Comments on Recursive Autoencoders

- Dependent on a tree structure
 - Parser is required
- Deep structures ($h=\log N$)
 - More supervision is required (at the internal nodes)

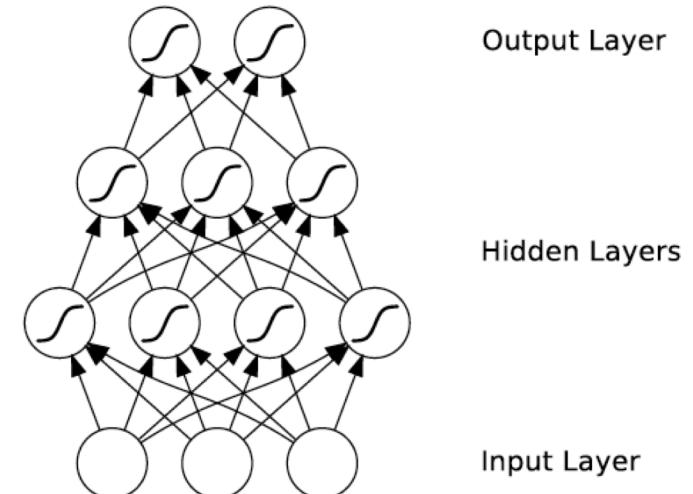
Recurrent Network for NLP

Recurrent Neural Network

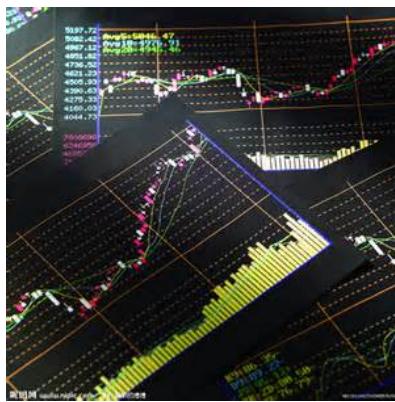
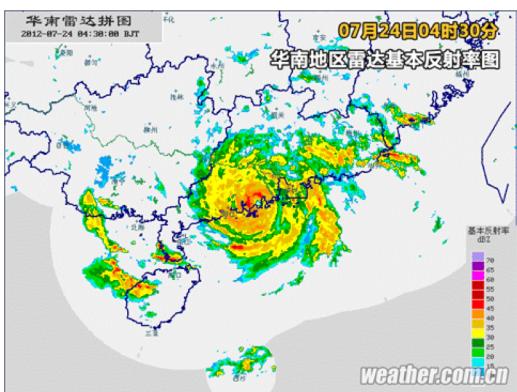
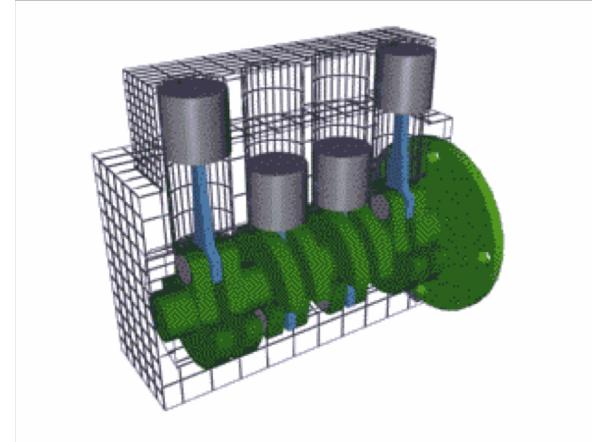
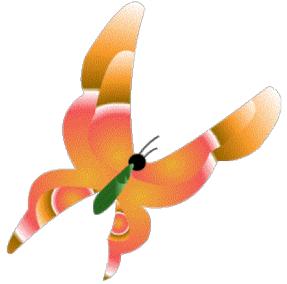
- Standard recurrent neural network (RNN)
- Gated recurrent network (GRU)
- Long Short-term Memory (LSTM)

Issues with MLPs?

- Problem 1: Can not model sequences
 - Fixed-sized Inputs & Outputs
 - No temporal structure, thus no dynamics
- Problem 2: Pure feed-forward processing
 - No **memory**, no **feedback**



Dynamic systems



Sequences are everywhere…

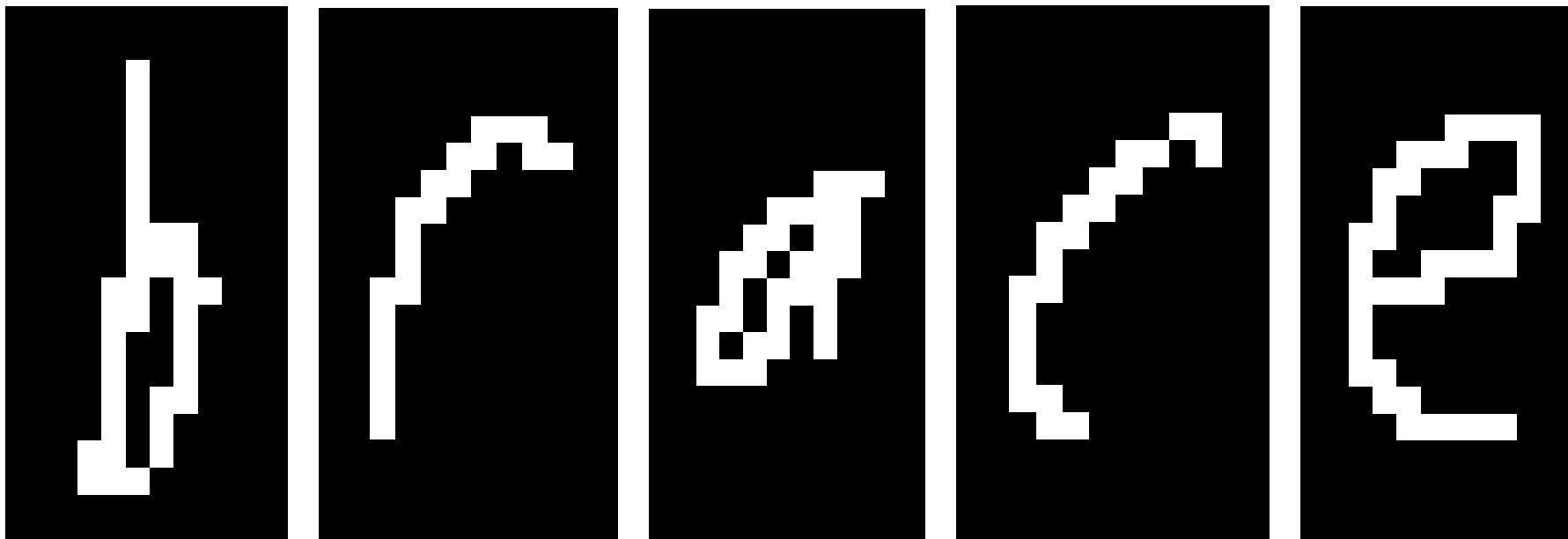
Foreign Minister. → FOREIGN MINISTER.

→ THE SOUND OF

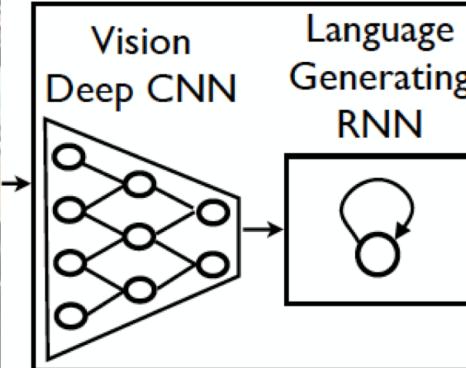
$\mathbf{x} = \begin{matrix} a_1=2 & a_2=0 & a_3=1 & a_4=3 & a_5=4 & a_6=2 & a_7=5 \\ \text{bringen} & \text{sie} & \text{bitte} & \text{das} & \text{auto} & \text{zurück} & . \end{matrix}$

$\mathbf{y} = \begin{matrix} \text{please} & \text{return} & \text{the} & \text{car} & . \end{matrix}$

Why model sequences?

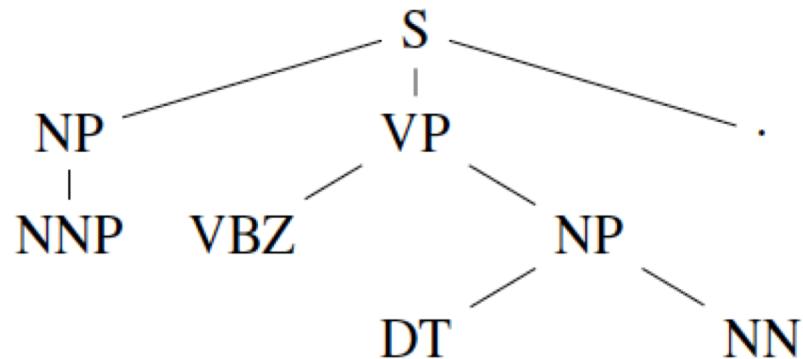


Even where you might not expect a sequence…



A group of people shopping at an outdoor market.
There are many vegetables at the fruit stand.

John has a dog . →



John has a dog . →

(S (NP NNP)_{NP} (VP VBZ (NP DT NN)_{NP})_{VP} .)_S

To be continued