

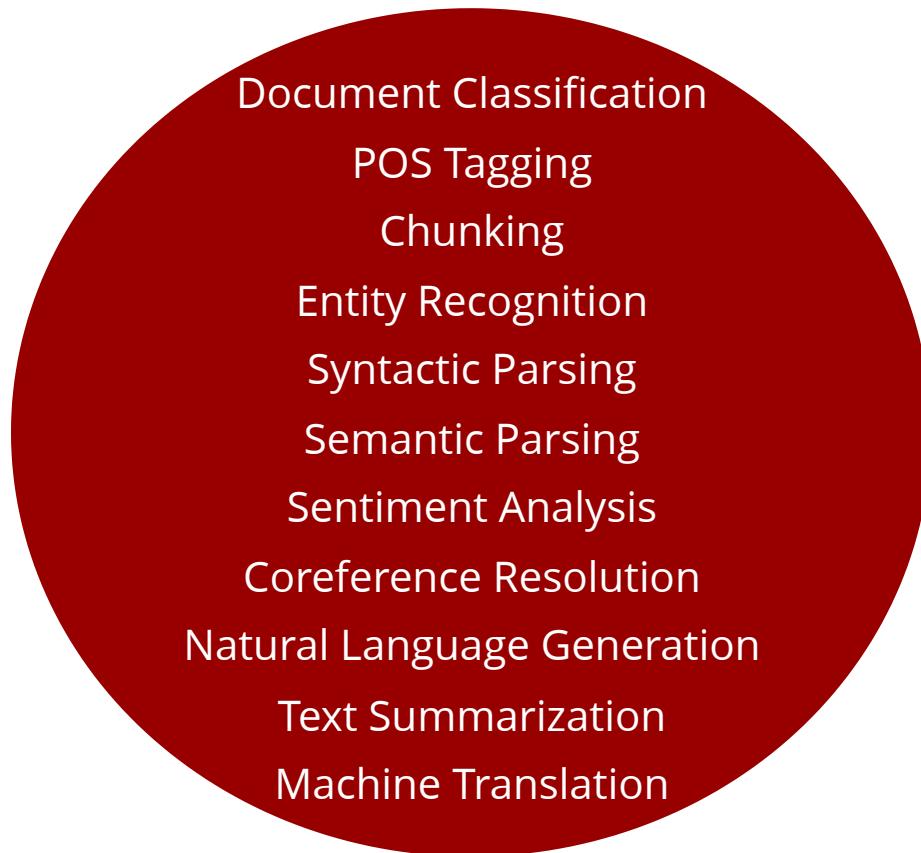
50.040

# Natural Language Processing

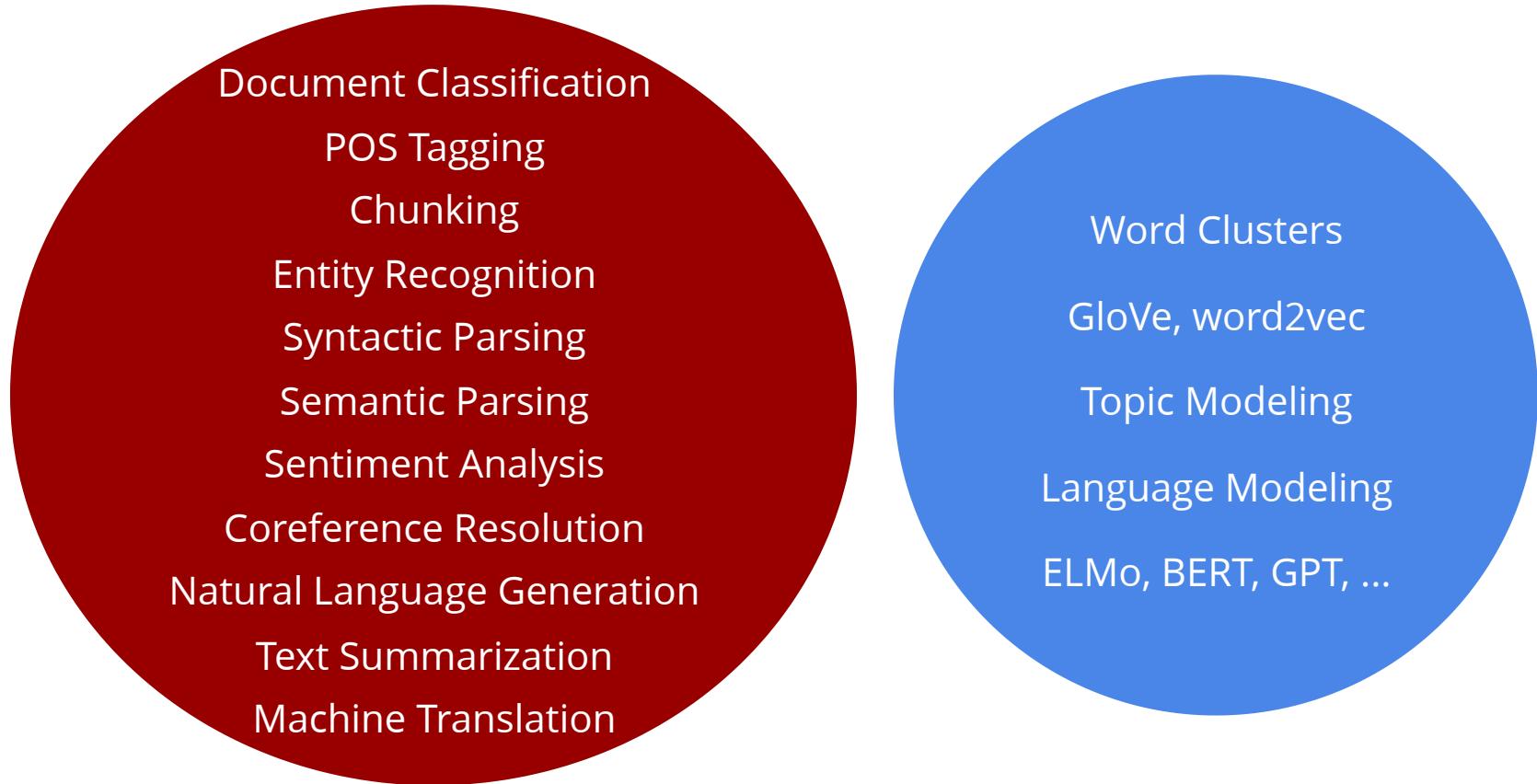
Lu, Wei



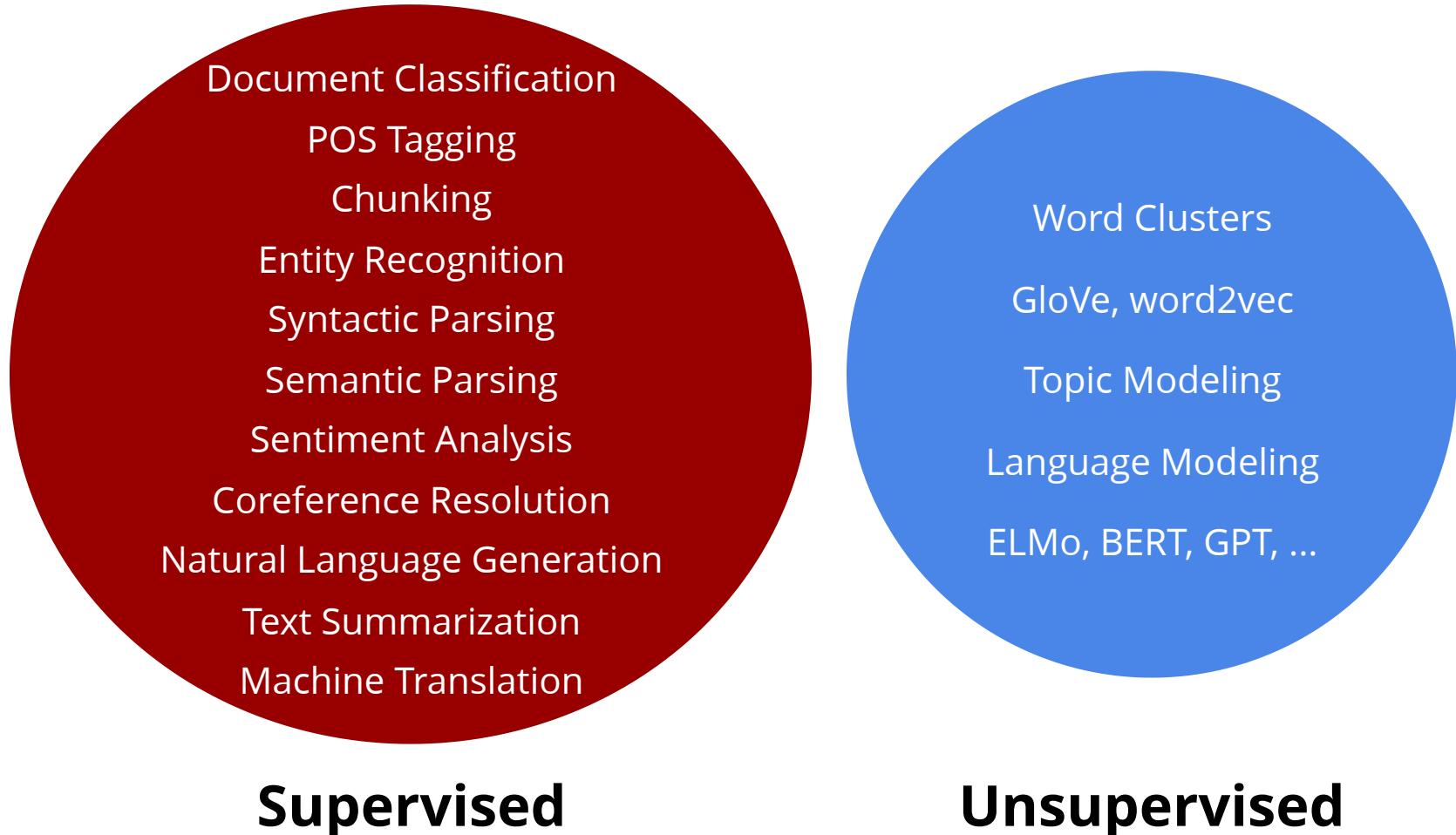
# Tasks in NLP



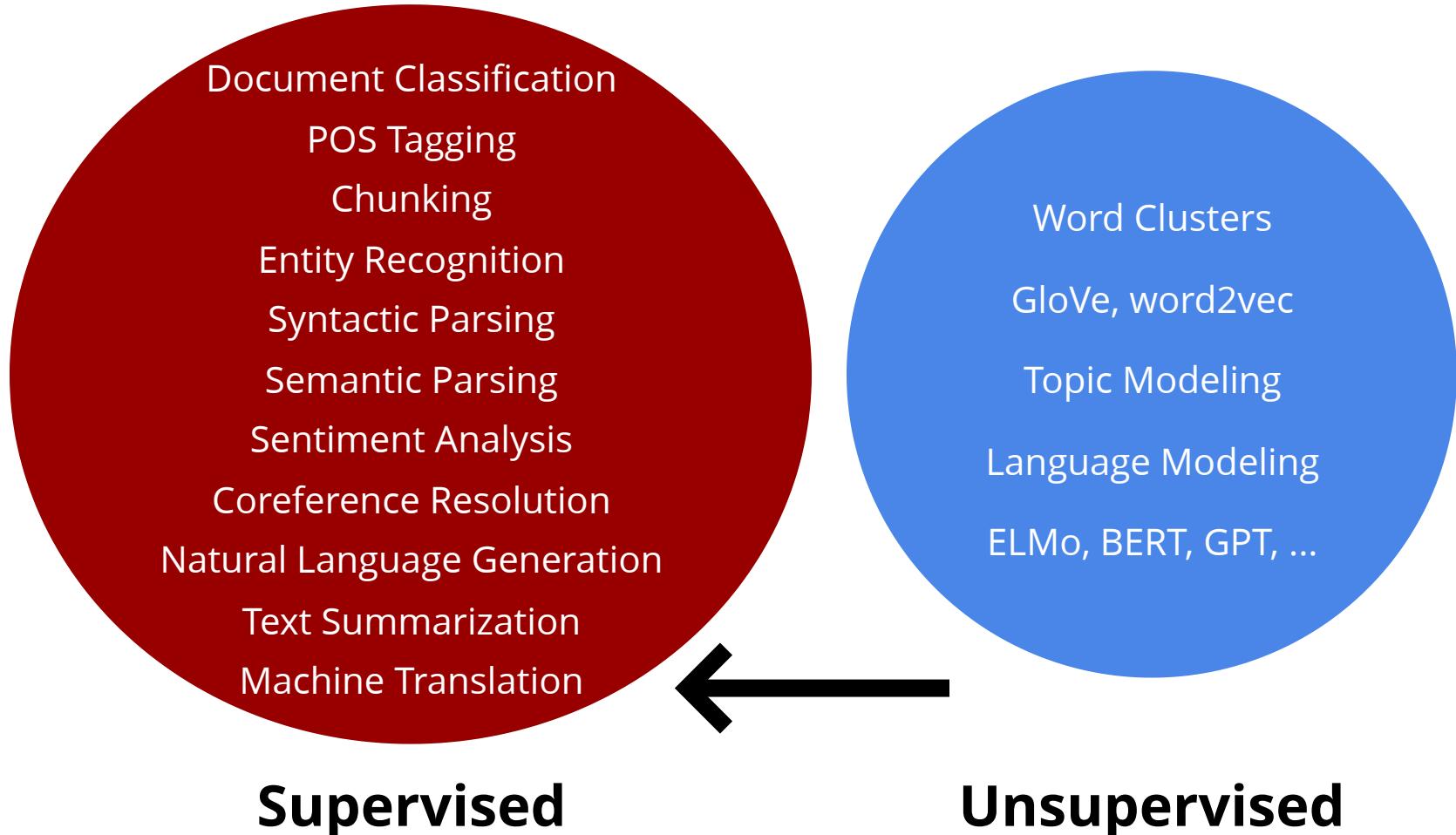
# Tasks in NLP



# Tasks in NLP



# Tasks in NLP

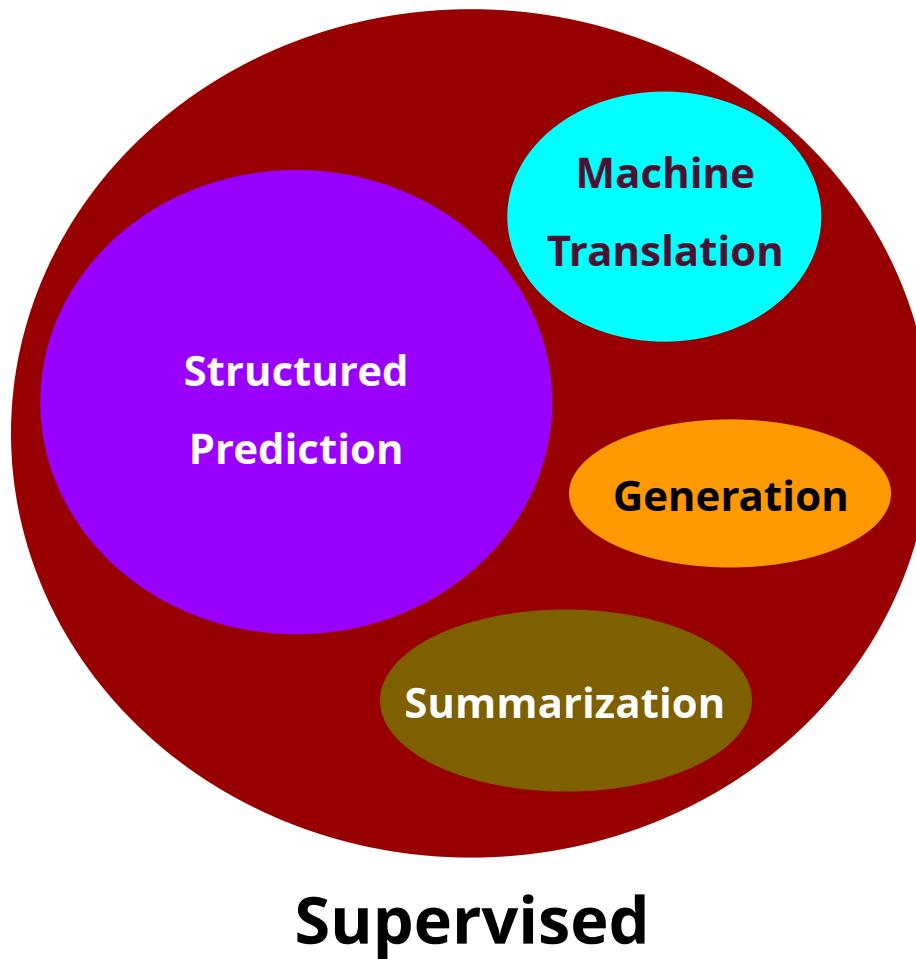


# Tasks in NLP



**Supervised**

# Tasks in NLP



# Part-of-Speech Tagging

A	N	V	D	N
<i>Fruit</i>	<i>flies</i>	<i>like</i>	<i>a</i>	<i>banana</i>

# Noun-Phrase Chunking

NP

*Fruit*

NP

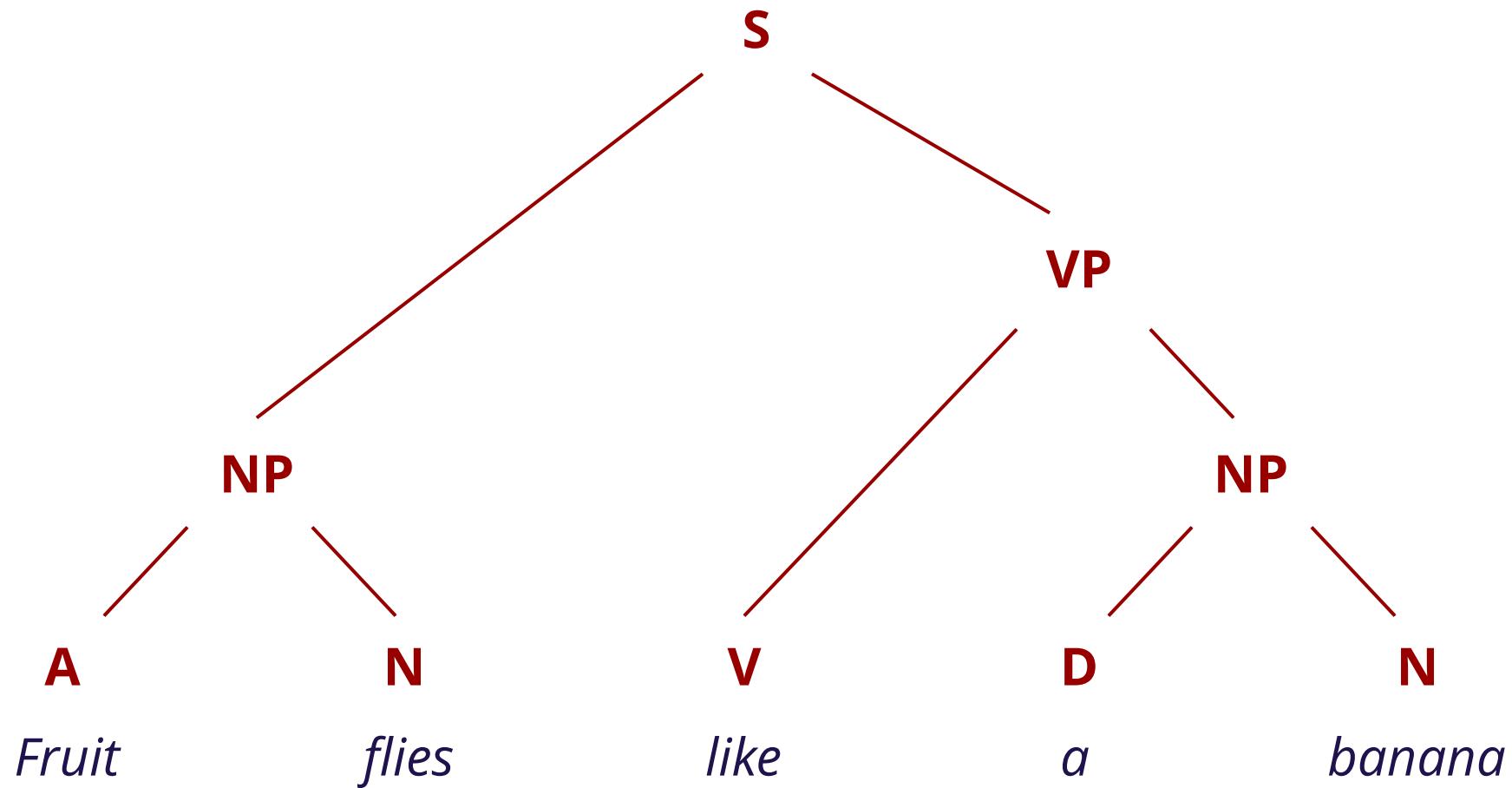
*banana*

*flies*

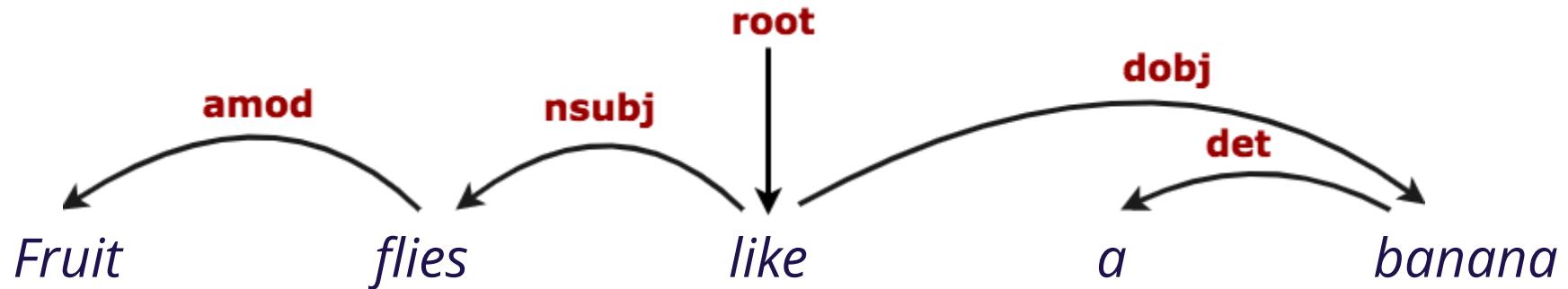
*like*

*a*

# Constituency Parsing



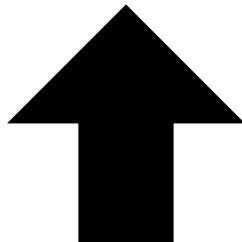
# Dependency Parsing



# Semantic Parsing

logical form

LIKE(F102, B87)



*Fruit*

*flies*

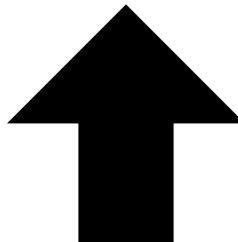
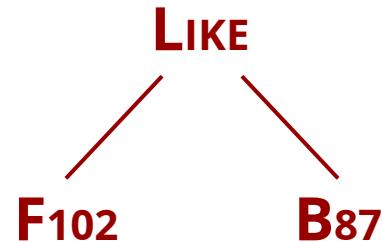
*like*

*a*

*banana*

# Semantic Parsing

logical form



*Fruit*

*flies*

*like*

*a*

*banana*

# Sentiment Analysis

(      **neutral**      )

*Fruit*

*flies*

*like*

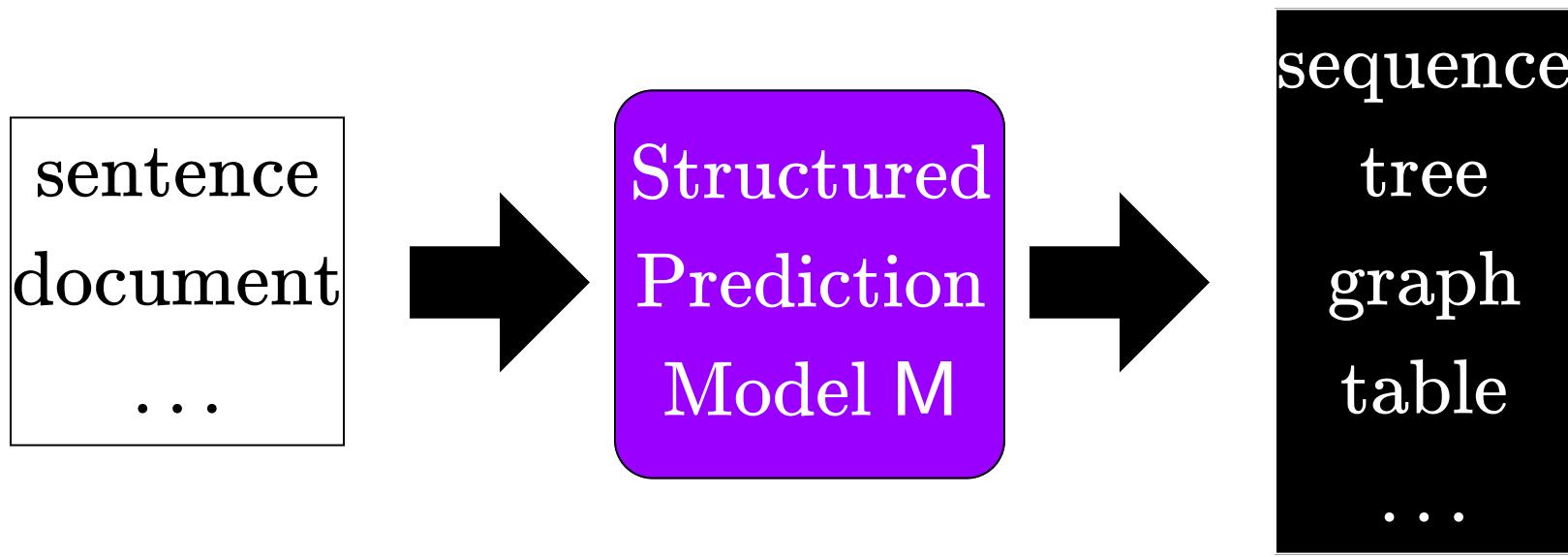
*a*

*banana*

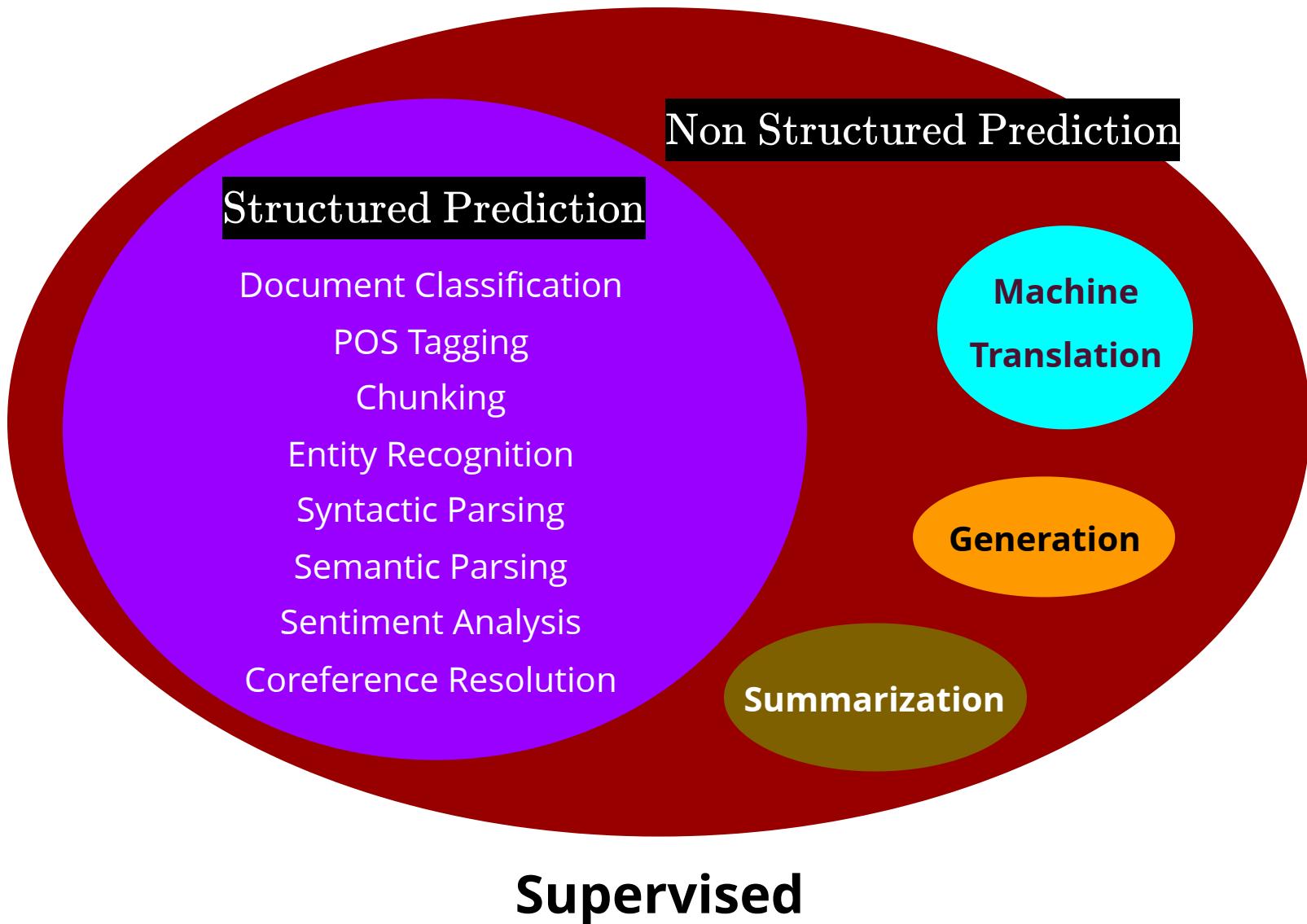
(      **positive**      )

# Structured Prediction

These supervised problems  
are "structured prediction"  
problems



# Tasks in NLP

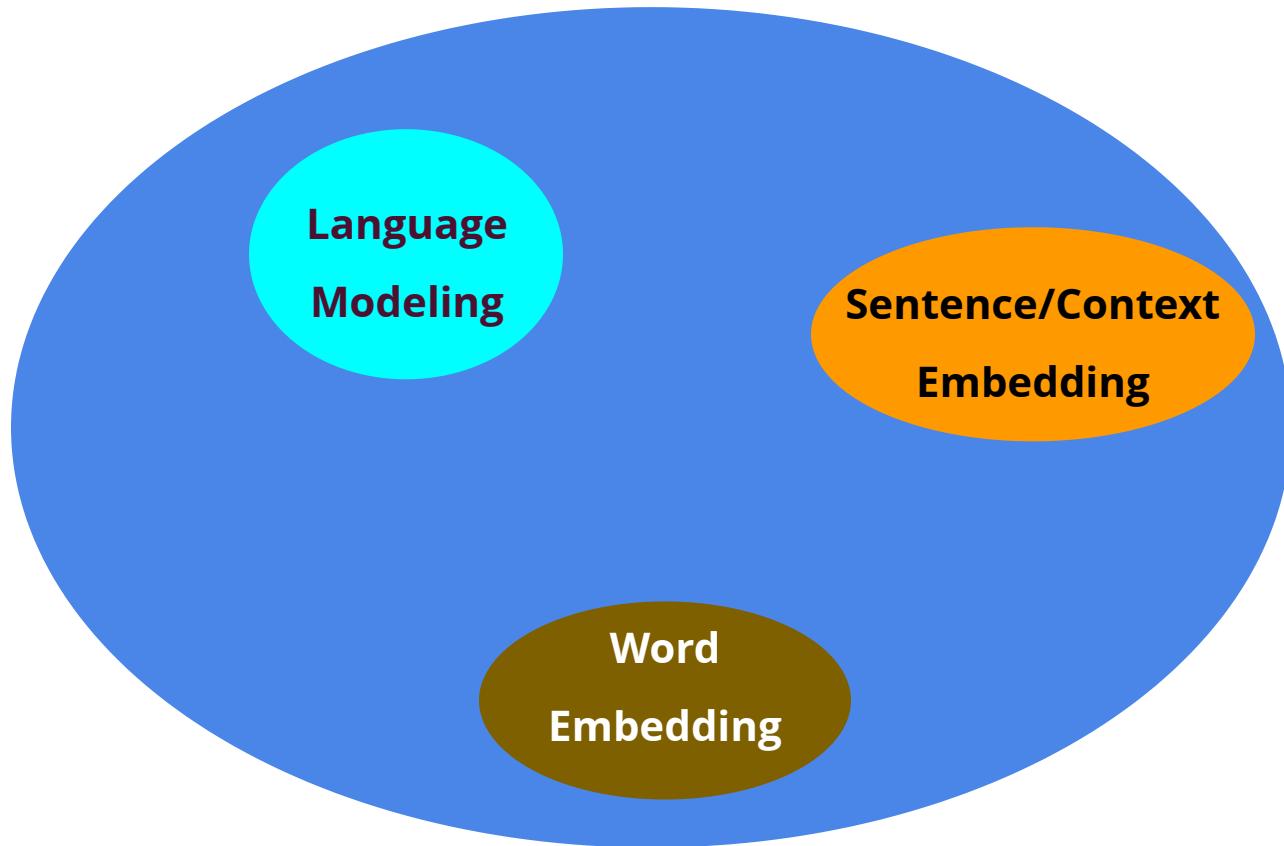


# Tasks in NLP



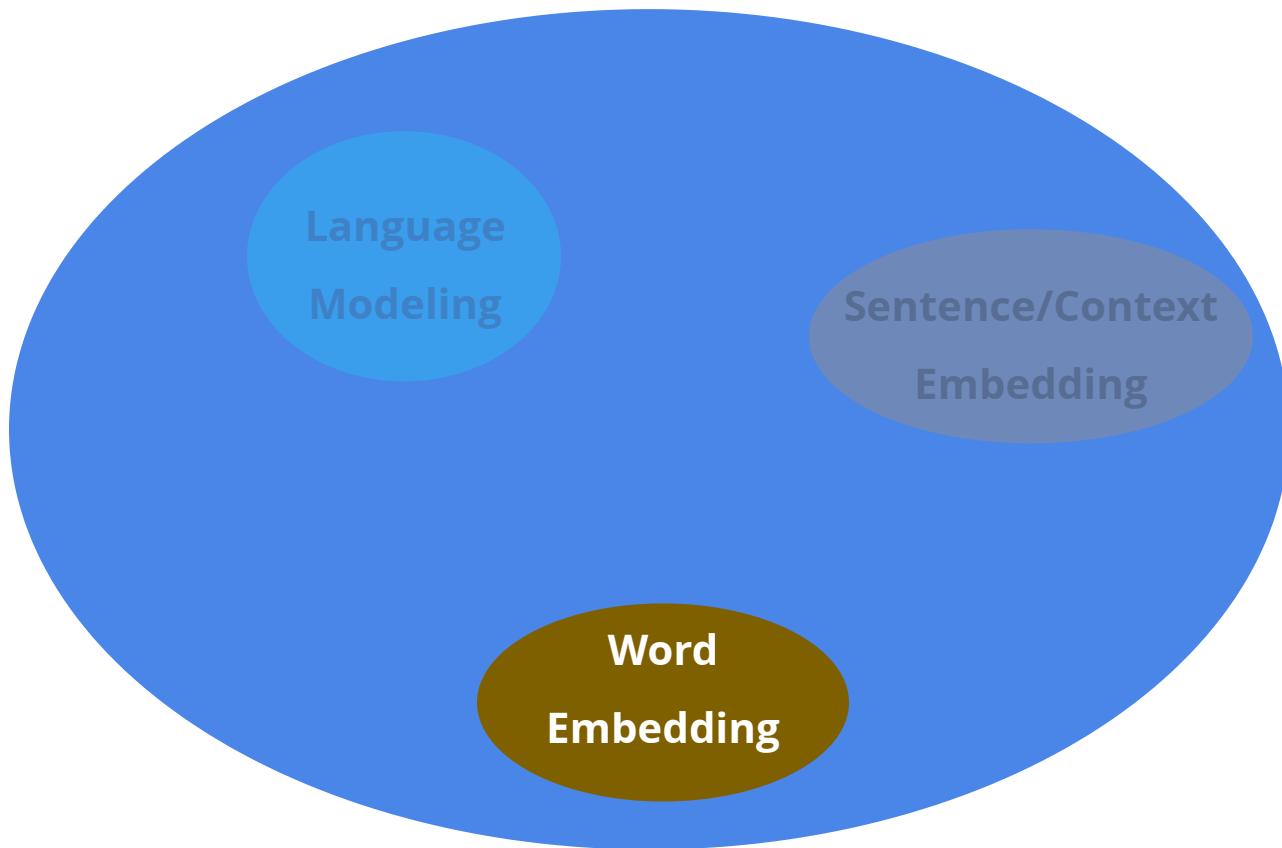
**Unsupervised**

# Tasks in NLP



**Unsupervised**

# Tasks in NLP



**Unsupervised**

# Fundamental Question

How do we capture the  
"meaning" of a word?

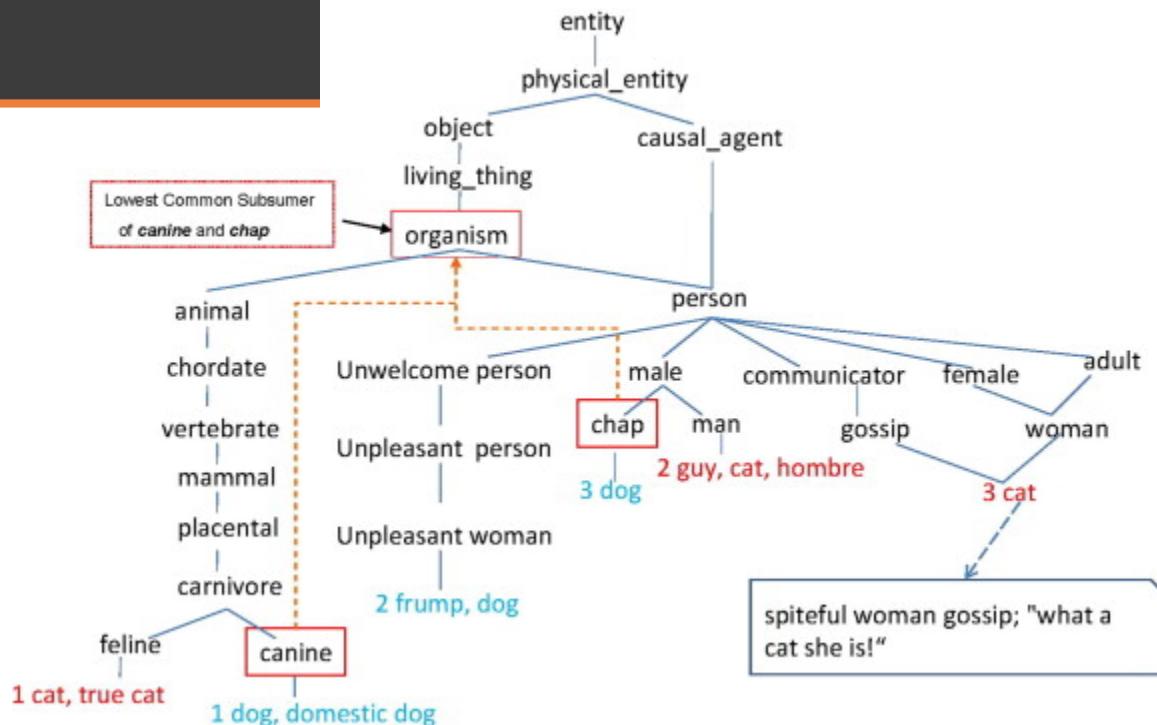
# WordNet

<https://wordnet.princeton.edu>



## WordNet

A Lexical Database for English

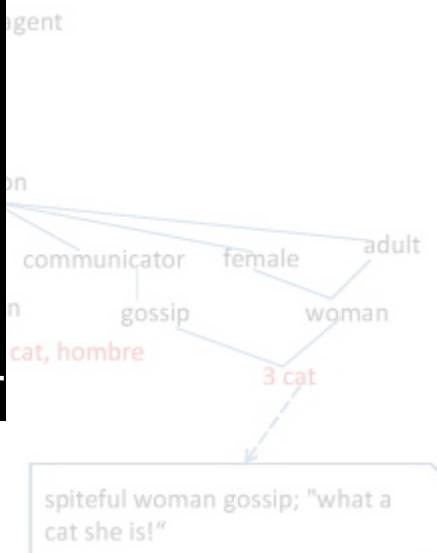


# WordNet

<https://wordnet.princeton.edu>



Limitations:  
Only for English  
Incomplete  
Subjective  
Expensive to Maintain



# WordNet

<https://wordnet.princeton.edu>



Limitations:  
Only for English  
Incomplete  
Subjective  
Expensive to Maintain



Possible to capture word semantics *automatically?*

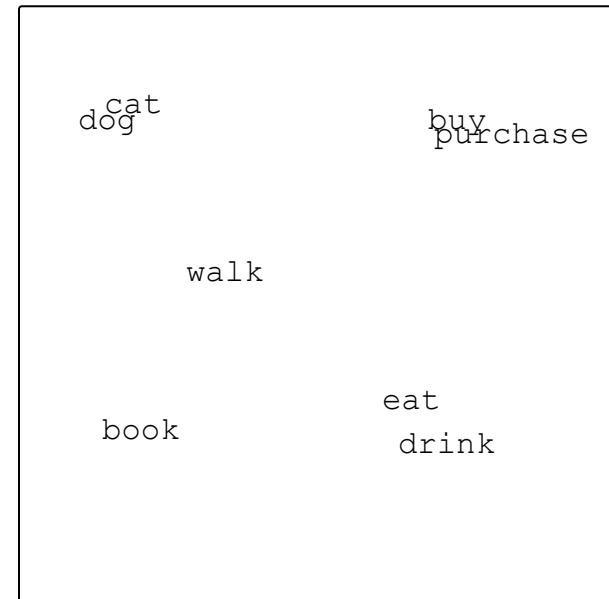
# How to Represent Words

## Traditionally ...

cat	eat	dog
$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \cdot \\ \cdot \\ \cdot \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \cdot \\ \cdot \\ \cdot \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \cdot \\ \cdot \\ \cdot \\ 0 \\ 0 \end{pmatrix}$

# How to Represent Words Ideally ...

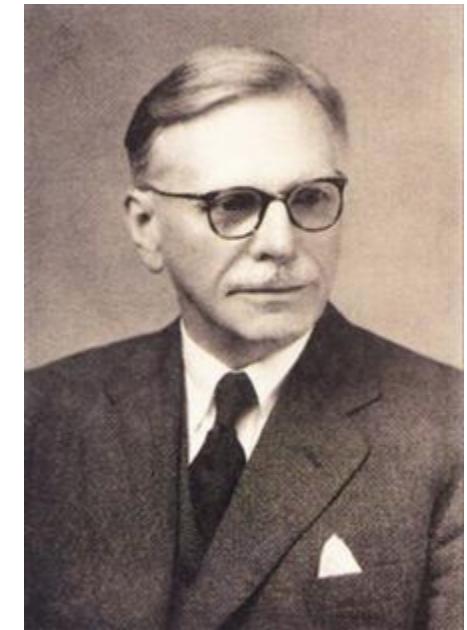
$$\begin{array}{c} \text{cat} \\ \left( \begin{array}{c} 0 \\ 1 \\ 0.232 \\ -0.903 \\ 0.213 \\ 0.239 \\ -0.679 \\ 0.923 \\ 0 \\ 0 \end{array} \right) \end{array} \quad \begin{array}{c} \text{eat} \\ \left( \begin{array}{c} 0 \\ 0 \\ -0.932 \\ -0.903 \\ 0.213 \\ 0.590 \\ 0.609 \\ 0.388 \\ 0 \\ 0 \end{array} \right) \end{array} \quad \begin{array}{c} \text{dog} \\ \left( \begin{array}{c} 0 \\ 0 \\ 0.190 \\ -0.887 \\ 0.193 \\ 0.257 \\ -0.708 \\ 0.906 \\ 0 \\ 0 \end{array} \right) \end{array}$$



# Distributional Semantics

*You shall know a word by the company it keeps*

- Firth, J. R. 1957



# Co-Occurrence Matrix

	cat		eat		dog		
cat	0	2	1	0	0	0	0
eat	2	0	0	1	0	1	0
eat	1	0	0	0	0	0	1
dog	0	1	0	0	1	0	0
dog	0	0	0	1	0	0	1
dog	0	1	0	0	0	0	1
dog	0	0	1	0	0	0	1
dog	0	0	0	0	1	1	1

# Co-Occurrence Matrix

	cat			eat		dog	
cat	0	2	1	0	0	0	0
	2	0	0	1	0	1	0
	1	0	0	0	0	0	1
	0	1	0	0	1	0	0
eat	0	0	0	1	0	0	0
	0	0	0	1	0	0	1
	0	1	0	0	0	0	1
dog	0	0	1	0	0	0	1
	0	0	0	0	1	1	0

# Co-Occurrence Matrix

	cat			eat		dog	
cat	0	2	1	0	0	0	0
	2	0	0	1	0	1	0
	1	0	0	0	0	0	1
	0	1	0	0	1	0	0
eat	0	0	0	1	0	0	1
	0	1	0	0	0	0	1
dog	0	0	1	0	0	0	1
	0	0	0	0	1	1	0

Apply SVD to compress the vectors!

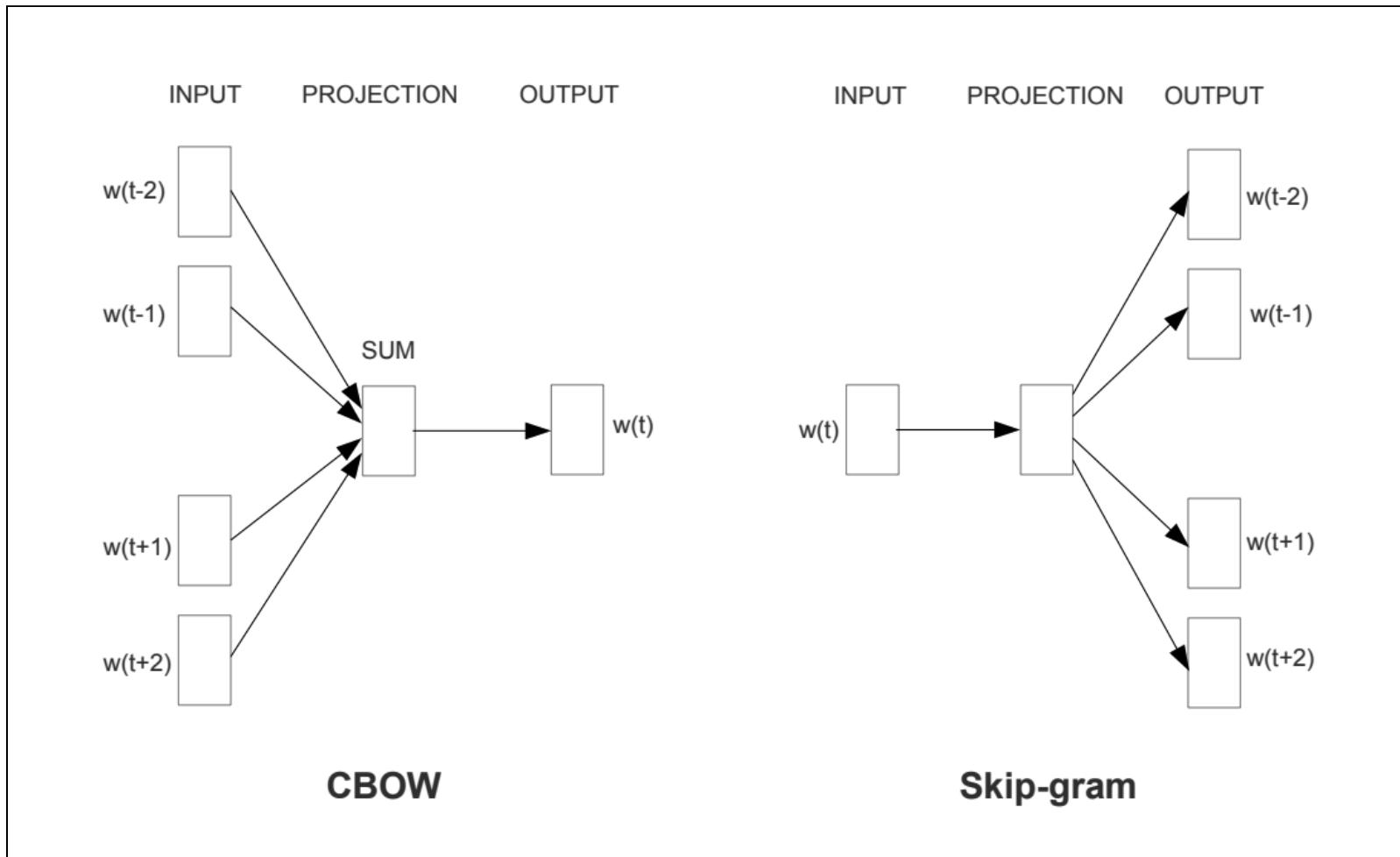
# Co-Occurrence Matrix

$$\begin{bmatrix} & \text{cat} & & \text{eat} & & \text{dog} \\ \text{cat} & \left[ \begin{array}{ccccc} 0 & 2 & 1 & 0 & 0 \\ 2 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0.232 & -0.903 & -0.932 & 0.190 \\ 0 & -0.903 & 0.213 & -0.903 & -0.887 \\ 0 & 0.213 & 0.239 & 0.213 & 0.193 \\ 0 & 0.239 & -0.679 & 0.590 & 0.257 \\ 0 & -0.679 & 0.923 & 0.609 & -0.708 \\ 0 & 0.923 & 0 & 0.388 & 0.906 \end{array} \right] & & \left[ \begin{array}{ccccc} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{array} \right] & & \left[ \begin{array}{ccccc} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{array} \right] \\ \text{eat} & & & & & \\ \text{dog} & & & & & \end{bmatrix}$$

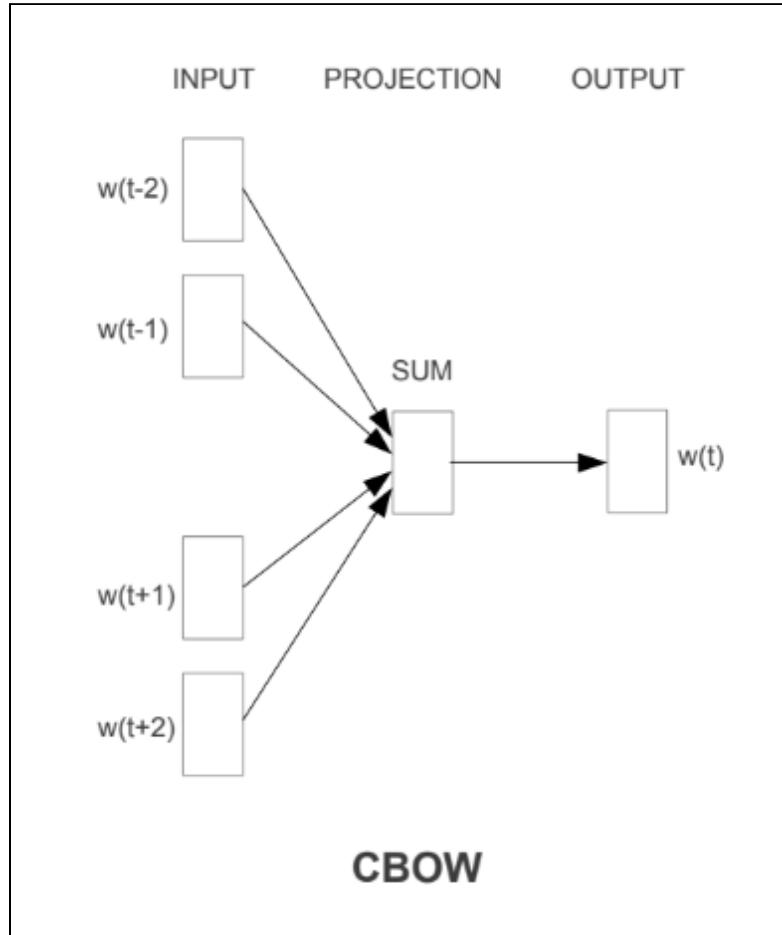
We call the resulting vectors **word embeddings**.

# Word2Vec

## (Mikolov et al. 2013)

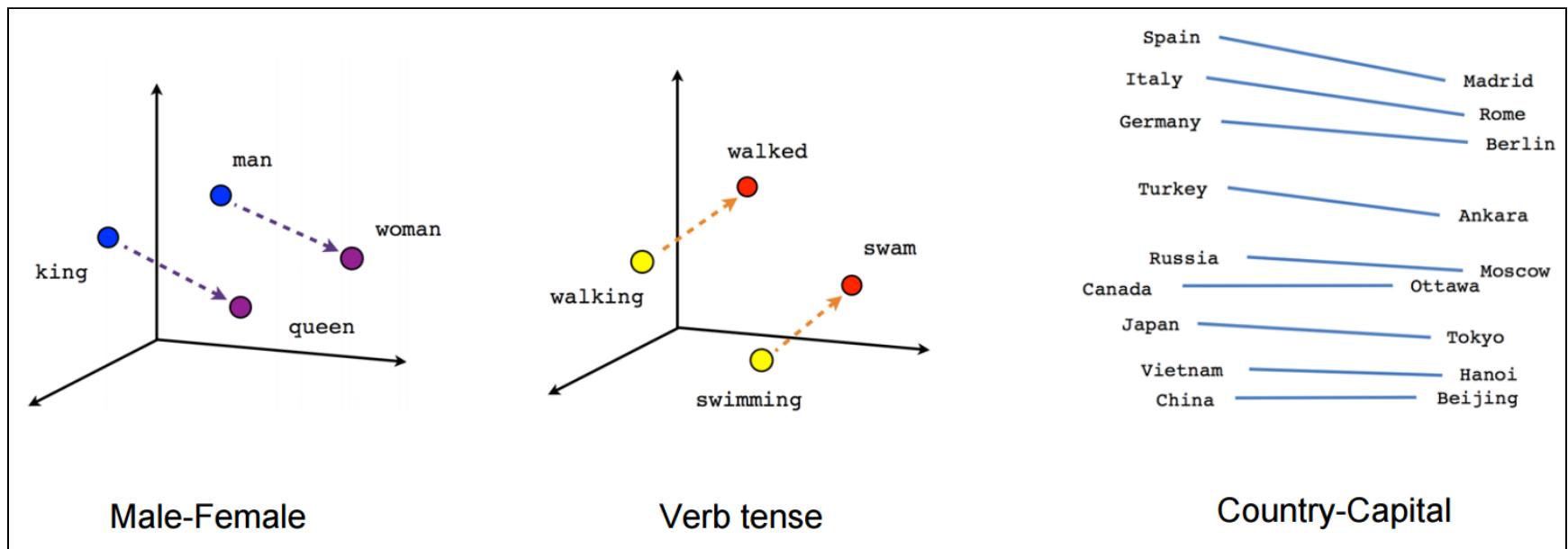


# Word2Vec

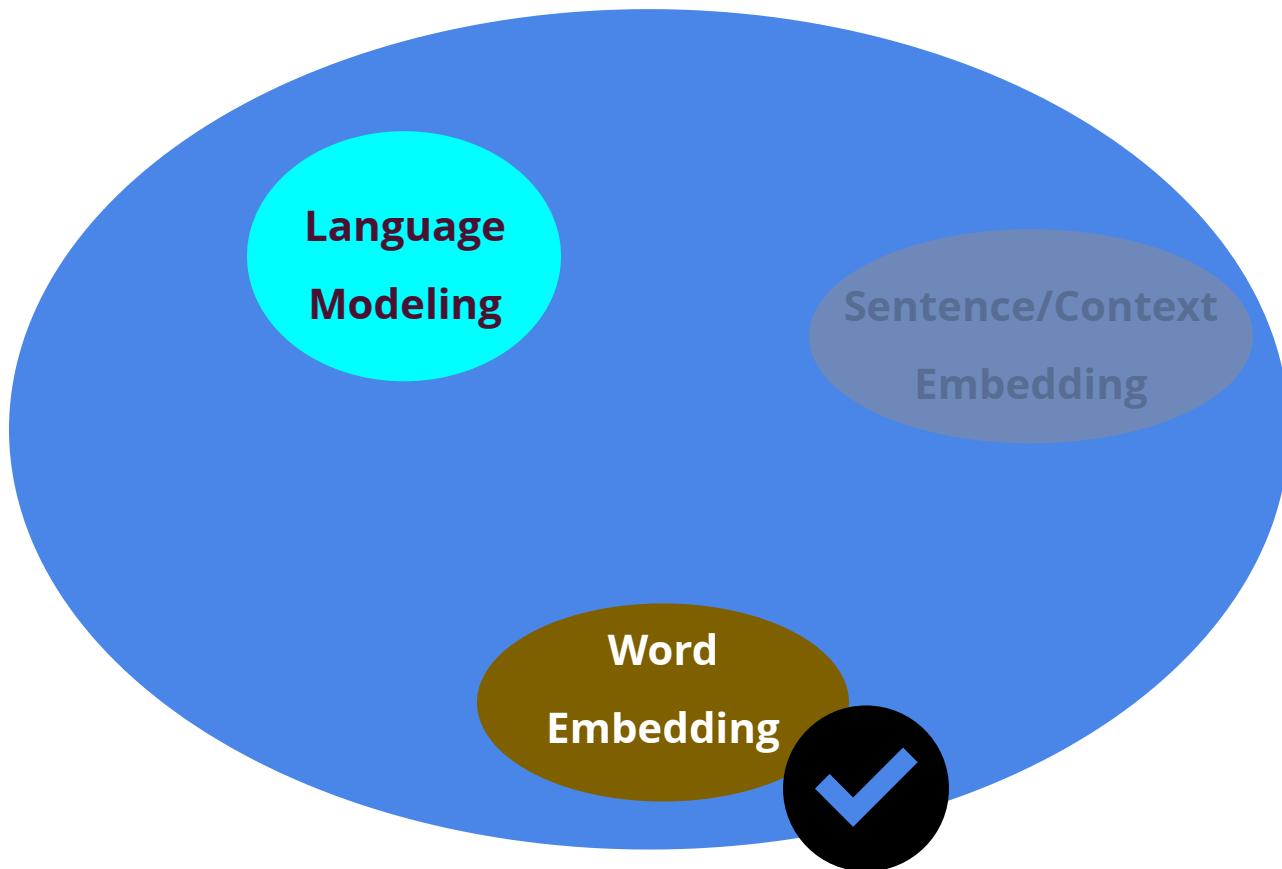


1. Project each context word from  $|V|$  dimensional space to its (current) embedding space
2. Sum these vectors
3. Project it to another space of dimension  $|V|$  (we expect this vector to be similar to the one-hot vector of the target word)

# Word2Vec



# Tasks in NLP



**Unsupervised**

# Fundamental Question

How likely can we ever see  
each English sentence?

# Language Modeling

How likely can we see these sentences in English?

*Fruit flies like a banana*



0.00078

*I love NLP*



0.00428

*flies Fruit*



0.00000001

# Language Model

Each possible sentence should be assigned a probability score (ideally, a "good" sentence should have a higher probability).

The sum of all such probabilities should be 1.

# $n$ -Gram Language Model

$$p(x_1, x_2, \dots, x_m) = \prod_{i=1, \dots, m} p(x_i | x_{i-n+1}, \dots, x_{i-1})$$

## Bigram LM

$$n = 2$$

$$p(x_1, x_2, \dots, x_m) = \prod_{i=1, \dots, m} p(x_i | x_{i-1})$$

# $n$ -Gram Language Model

$$p(x_1, x_2, \dots, x_m) = \prod_{i=1, \dots, m} p(x_i | x_{i-n+1}, \dots, x_{i-1})$$

How many model parameters do we have to store for a n-gram model?

# Curse of Dimensionality

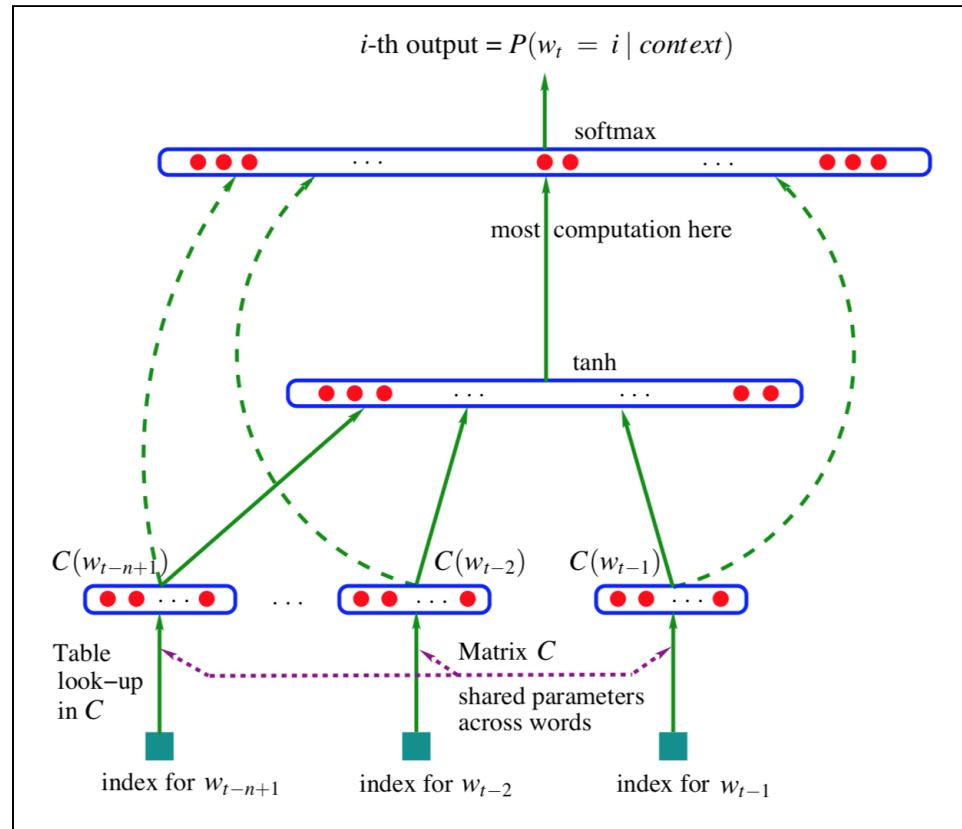
in the order of  $|V|^n$

Too many model parameters!

One Solution:

Learn language model with word embeddings!

# Neural Language Model (Bengio et al. 2003)

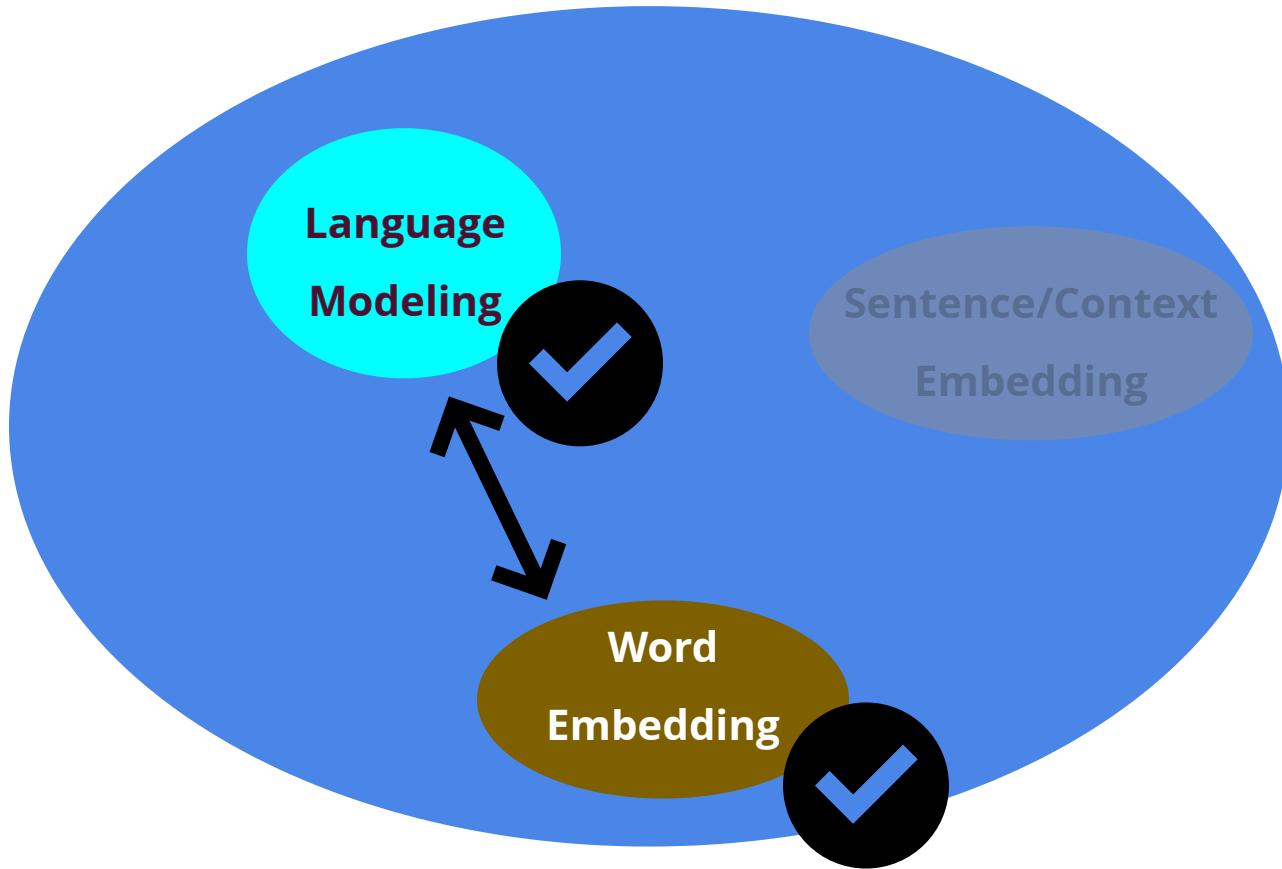


Learn simultaneously the *word embeddings* and the parameters of the language model *probability function*.

# Neural Language Model

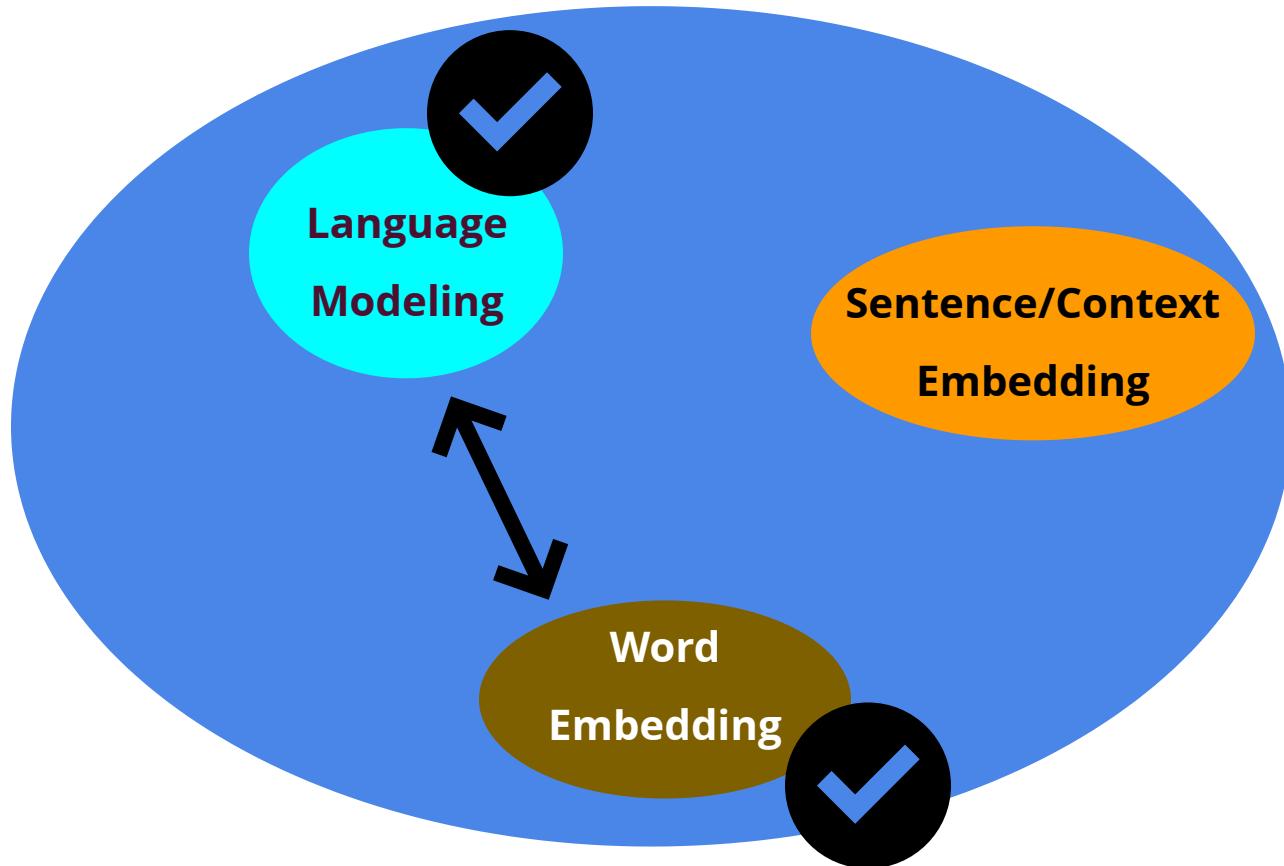
1. associate with each word in the vocabulary a distributed *word embedding*,
2. express the joint probability function of word sequences in terms of the embeddings of these words in the sequence, and
3. learn simultaneously the *word embeddings* and the parameters of that *probability function*.

# Tasks in NLP



**Unsupervised**

# Tasks in NLP



**Unsupervised**

# Word Senses

A word may have multiple senses.

bank<sup>1</sup>  
/baŋk/  
**noun**  
noun: bank; plural noun: banks

1. the land alongside or sloping down to a river or lake.  
"willows lined the bank of the stream"  
synonyms: [edge](#), [side](#), [embankment](#), [levee](#), [border](#), [verge](#), [boundary](#), [margin](#), [rim](#), [fringe](#), [fringes](#), [flank](#), [brick](#), [perimeter](#), [circumference](#), [extremity](#), [periphery](#), [limit](#), [outer limit](#), [limits](#), [bound](#), [bounds](#); [More](#)

bank<sup>2</sup>  
/baŋk/  
**noun**  
noun: bank; plural noun: banks

2. a long, high mass or mound of earth or stones.  
"a grassy bank"  
synonyms: [slope](#), [rise](#), [incline](#)  
• an elevation in the sea  
• a transverse slope giving a boat a sharp entry into the water  
• the sideways tilt of an aircraft  
"a rather steep angle of bank"  
3. a set of similar things, especially of the same kind.  
"the DJ had big banks of speakers."  
synonyms: [array](#), [row](#), [line](#)  
• a tier of oars.  
"the early ships had oars in banks"  
4. the cushion of a pool table.  
"a bank shot"

**verb**  
verb: bank; 3rd person present: banking; past participle: banked;

1. heap (a substance) into a bank.  
"the rain banked the soil up"  
synonyms: pile (up), heap  
More  
• form into a mass or mound.  
"purple clouds banked up in the sky"  
• heave up (a fire) with fuel.  
"she banked up the fire"  
synonyms: damp (down), extinguish  
• edge or surround with.  
"steps banked with pink roses"

2. (with reference to an aircraft) bank (a road, railway, cornering).  
"the track was banked to allow a train to take curves faster while maintaining passenger comfort"

3. BRITISH (of a locomotive) provide additional power for (a train) in ascending an incline.

4. (of an angler) succeed in landing (a fish).  
"it was the biggest rainbow trout that had ever been banked"

5. NORTH AMERICAN (in pool) play (a ball) so that it rebounds off a surface such as a cushion.  
"I banked the eight ball off two cushions"

## WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for: bank

Display Options:  Select option to change  Change

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations  
Display options for sense: (gloss) "an example sentence"

### Noun

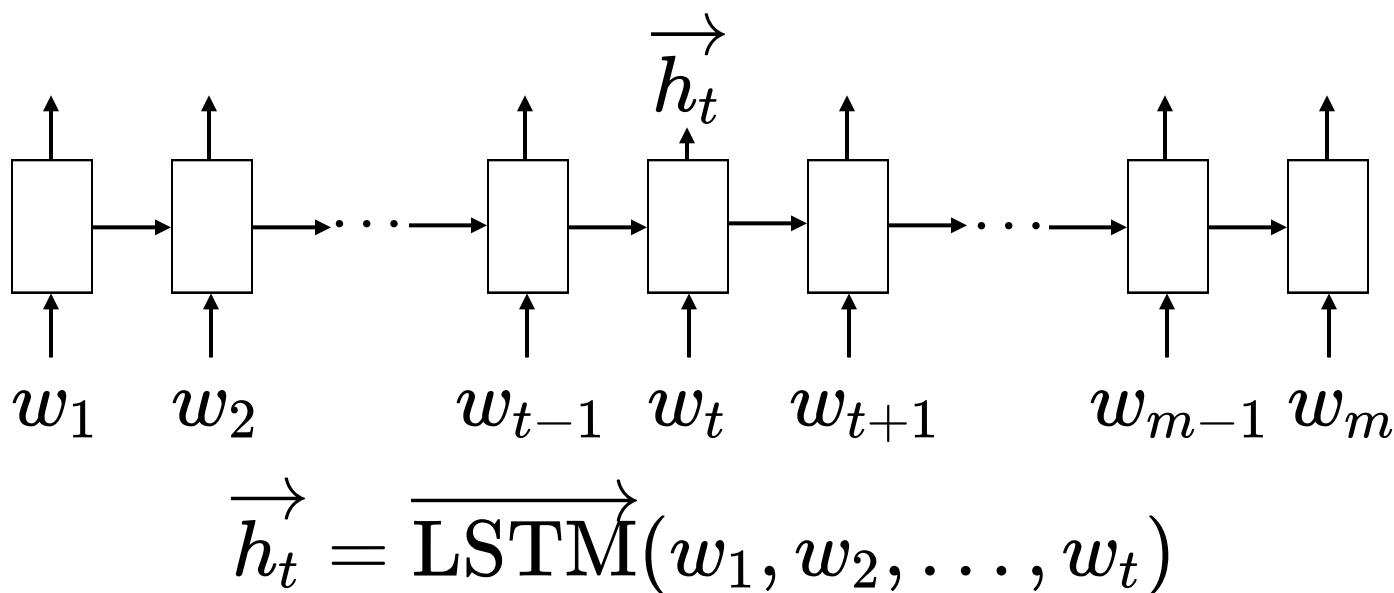
- [S: \(n\) bank](#) (sloping land (especially the slope beside a body of water)) "they pulled the canoe up on the bank"; "he sat on the bank of the river and watched the currents"
- [S: \(n\) depository financial institution](#), [bank](#), [banking concern](#), [banking company](#) (a financial institution that accepts deposits and channels the money into lending activities) "he cashed a check at the bank"; "that bank holds the mortgage on my home"
- [S: \(n\) bank](#) (a long ridge or pile) "a huge bank of earth"
- [S: \(n\) bank](#) (an arrangement of similar objects in a row or in tiers) "he operated a bank of switches"
- [S: \(n\) bank](#) (a supply or stock held in reserve for future use (especially in emergencies))
- [S: \(n\) bank](#) (the funds held by a gambling house or the dealer in some gambling games) "he tried to break the bank at Monte Carlo"
- [S: \(n\) bank](#), [cant](#), [camber](#) (a slope in the turn of a road or track; the outside is higher than the inside in order to reduce the effects of centrifugal force)
- [S: \(n\) savings bank](#), [coin bank](#), [money box](#), [bank](#) (a container (usually with a slot in the top) for keeping money at home) "the coin bank was empty"
- [S: \(n\) bank](#), [bank building](#) (a building in which the business of banking transacted) "the bank is on the corner of Nassau and Witherspoon"
- [S: \(n\) bank](#) (a flight maneuver; aircraft tips laterally about its longitudinal axis (especially in turning)) "the plane went into a steep bank"

### Verb

- [S: \(v\) bank](#) (tip laterally) "the pilot had to bank the aircraft"
- [S: \(v\) bank](#) (enclose with a bank) "bank roads"
- [S: \(v\) bank](#) (do business with a bank or keep an account at a bank) "Where do you bank in this town?"
- [S: \(v\) bank](#) (act as the banker in a game or in gambling)
- [S: \(v\) bank](#) (be in the banking business)
- [S: \(v\) deposit](#), [bank](#) (put into a bank account) "She deposits her paycheck every month"
- [S: \(v\) bank](#) (cover with ashes so to control the rate of burning) "bank a fire"
- [S: \(v\) count](#), [bet](#), [depend](#), [swear](#), [rely](#), [bank](#), [look](#), [calculate](#), [reckon](#) (have faith or confidence in) "you can count on me to help you any time"; "Look to your friends for support"; "You can bet on that!"; "Depend on your family in times of crisis"

# Long Short-Term Memory

(Hochreiter & Schmidhuber 1997)



# LSTM

$c_{t-1}$

The long-term  
memory

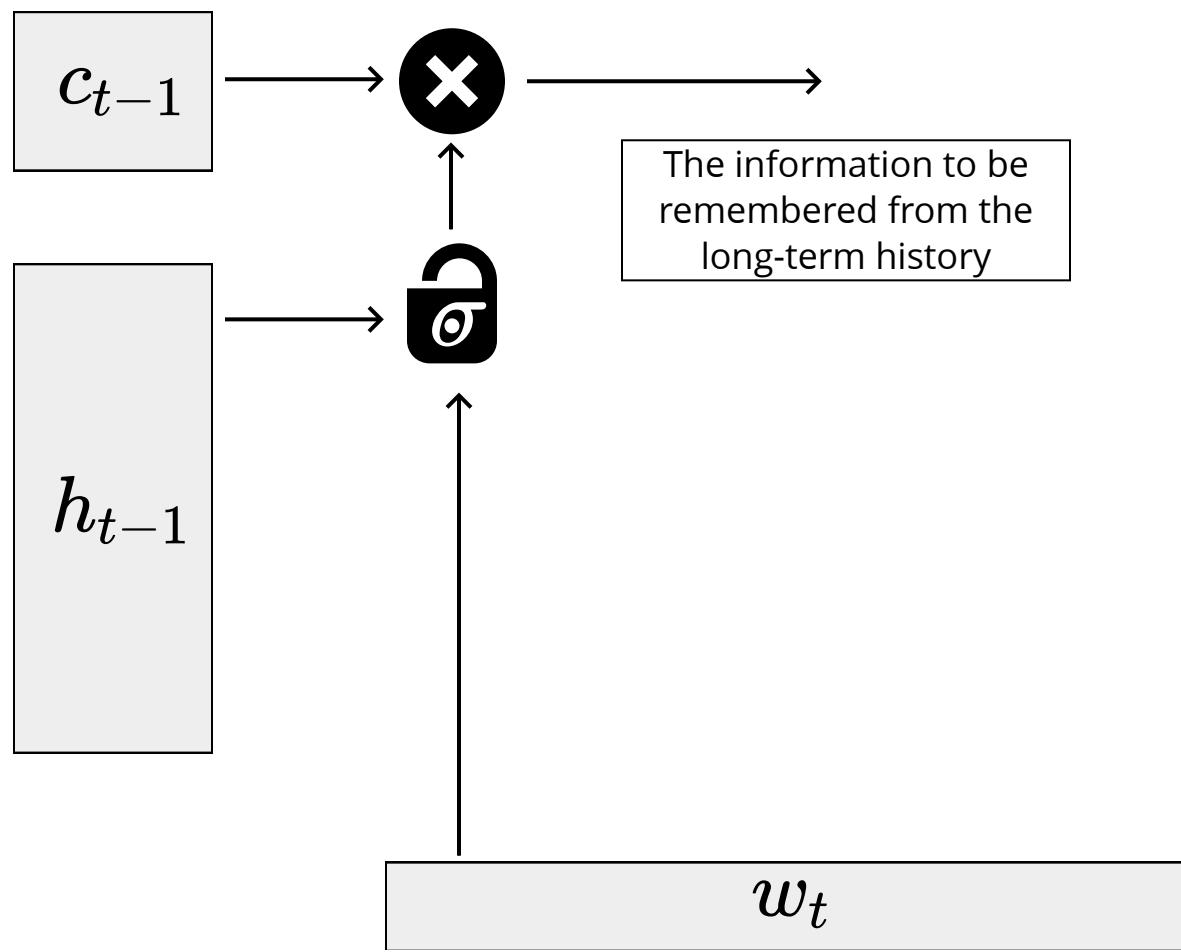
$h_{t-1}$

The working  
memory

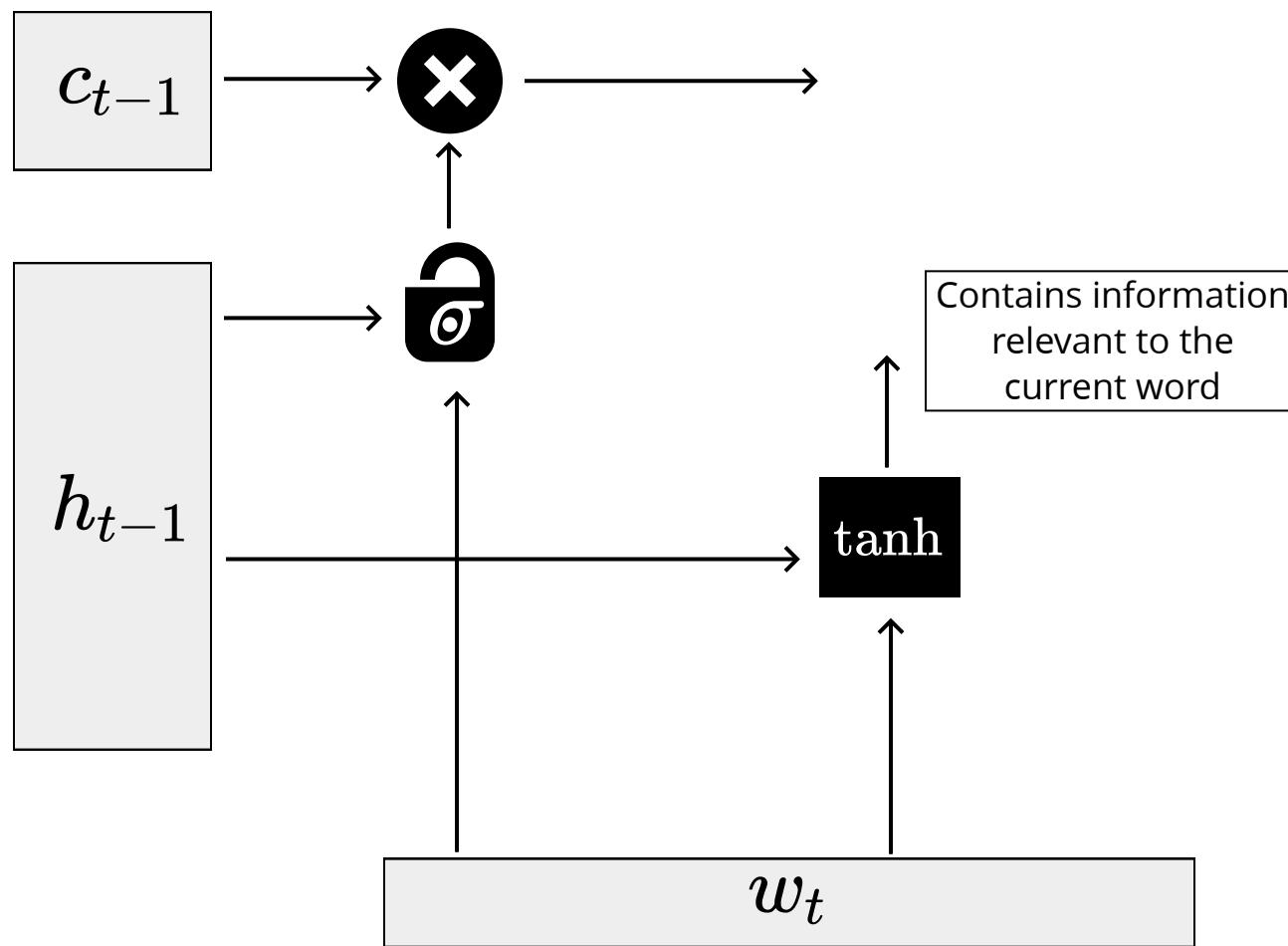
The current word

$w_t$

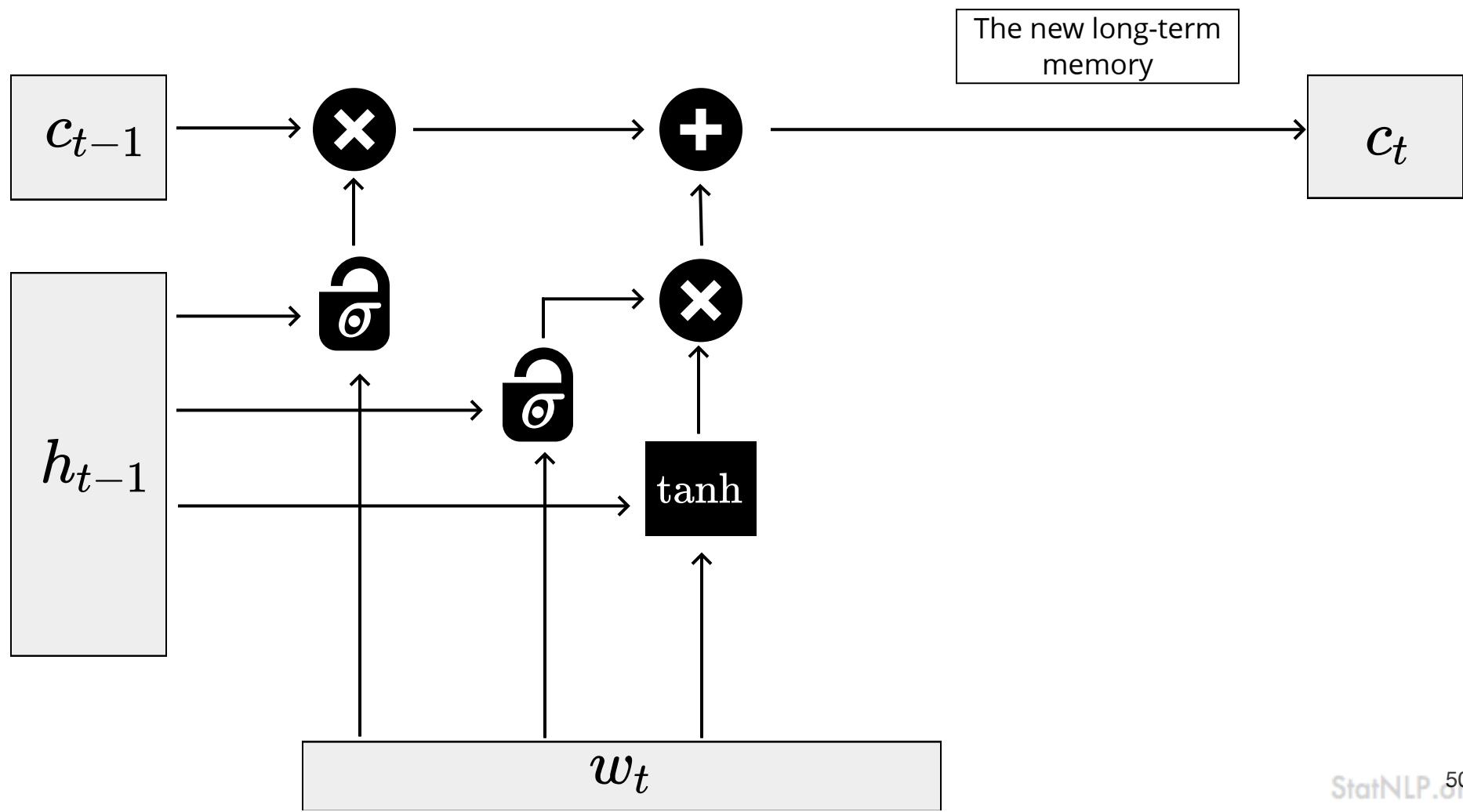
# LSTM



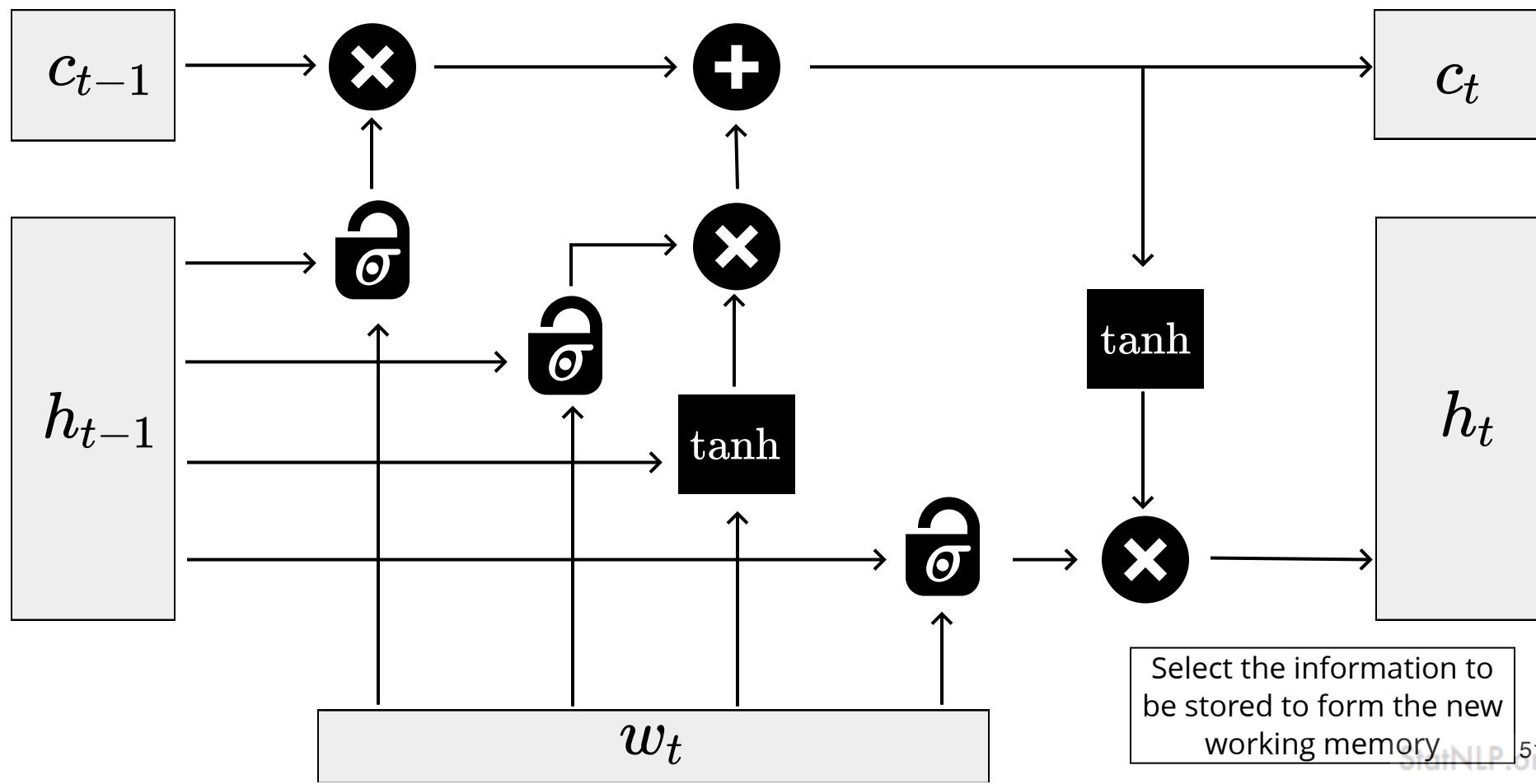
# LSTM



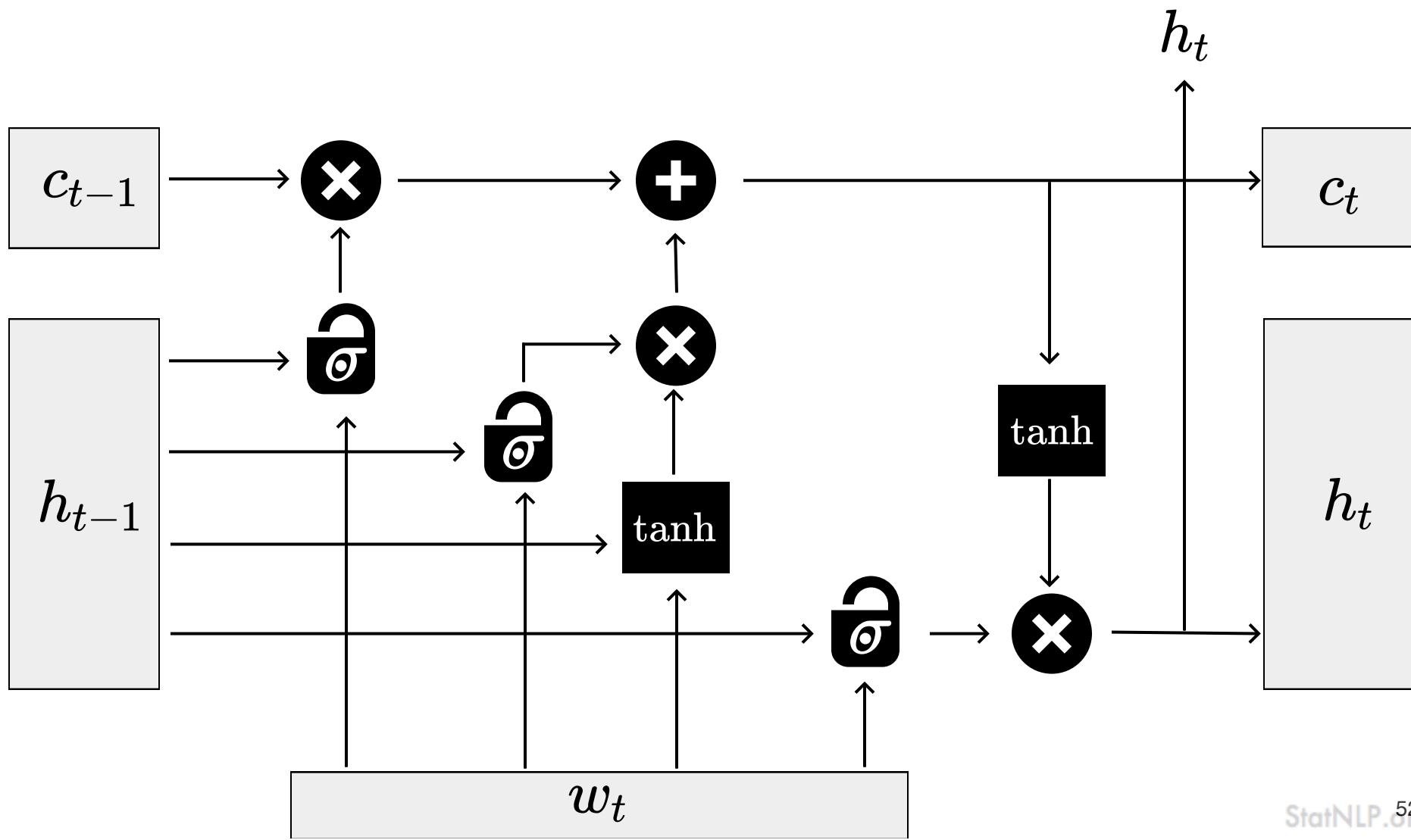
# LSTM



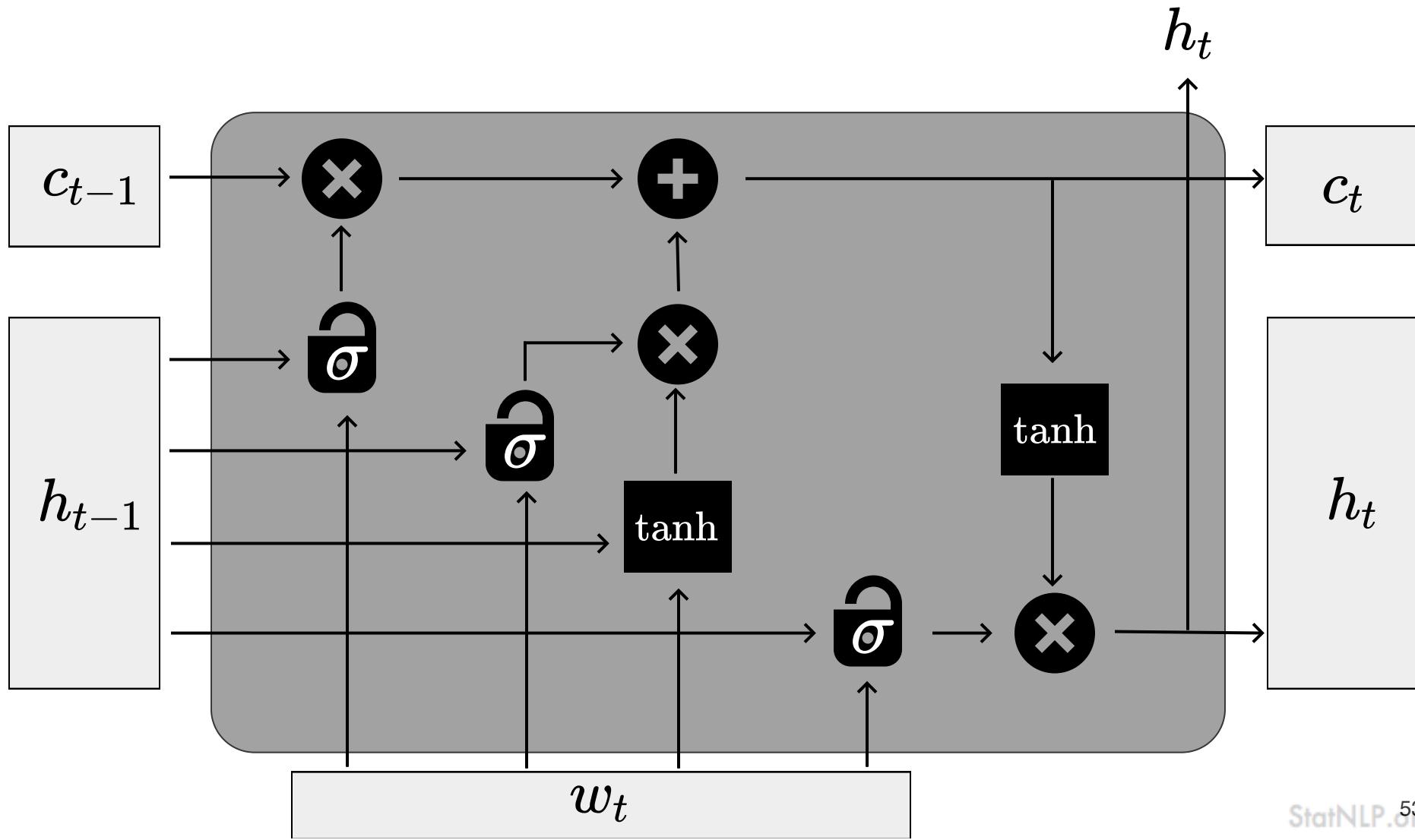
# LSTM



# LSTM

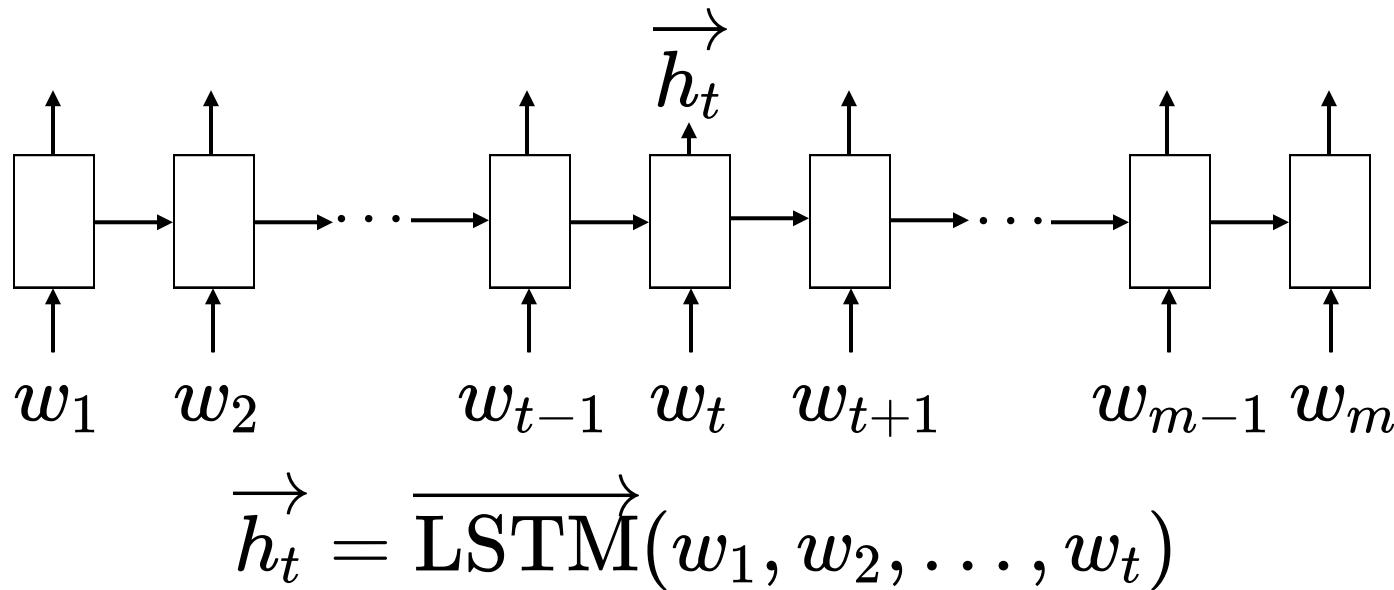


# LSTM



# LSTM Language Model

(Sundermeyer et al. 2012)

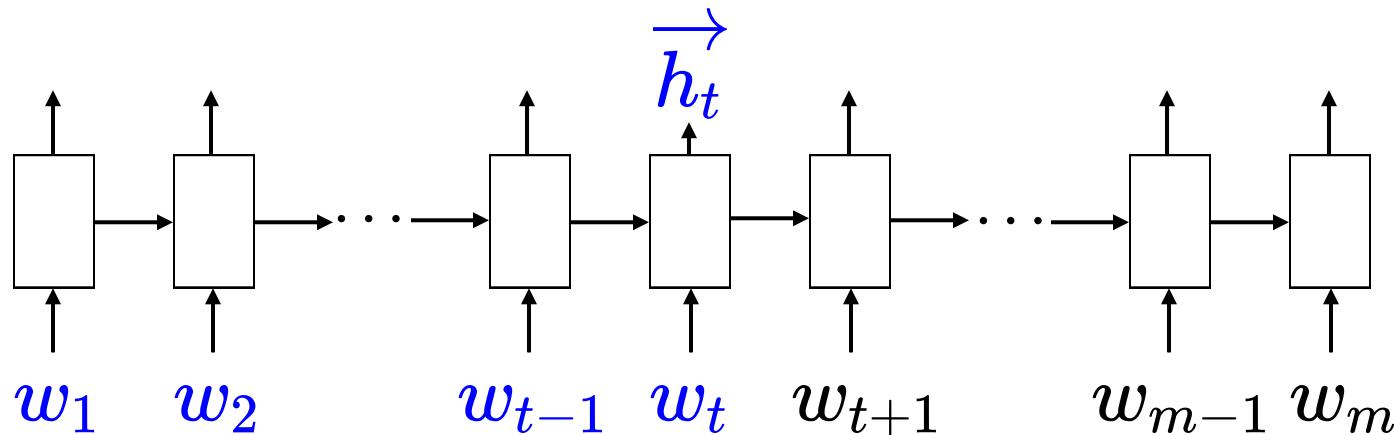


Individual probability at each position:

1. project  $\vec{h}_t$  to a  $|V|$  dimensional space
2. apply softmax on top of the vector
3. get the probability of generating the desired word  $w_{t+1}$

# LSTM Language Model

(Sundermeyer et al. 2012)



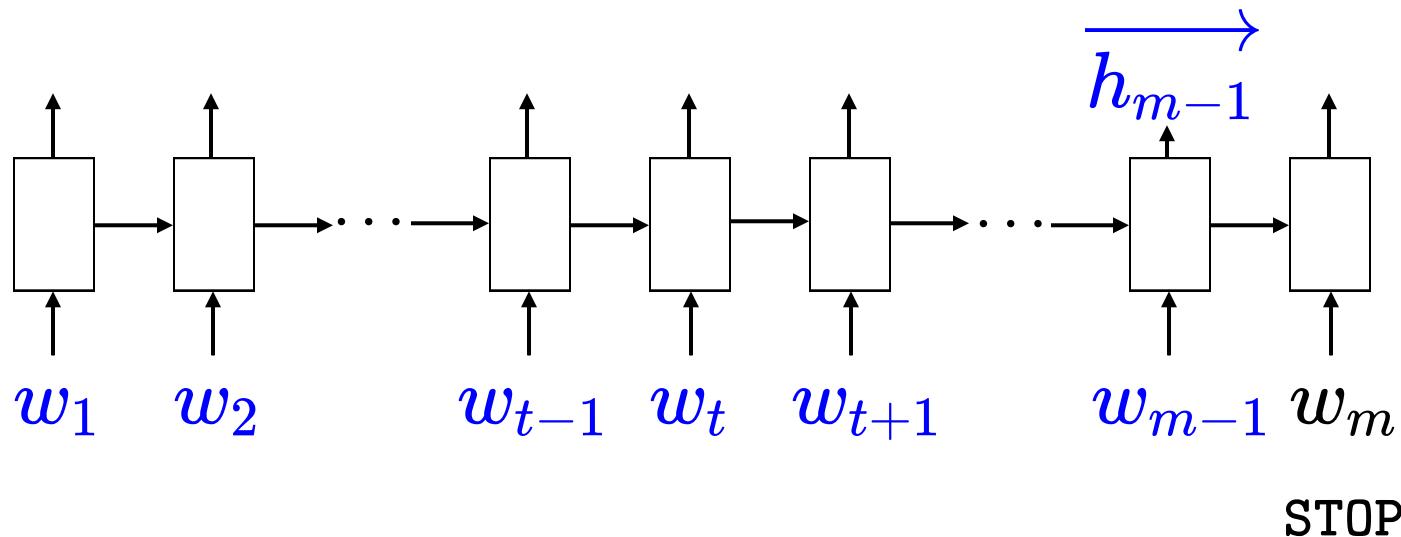
$$\vec{h}_t = \overrightarrow{\text{LSTM}}(w_1, w_2, \dots, w_t)$$

Context Embedding!

A function over a sequence of word embeddings.

# LSTM Language Model

(Sundermeyer et al. 2012)



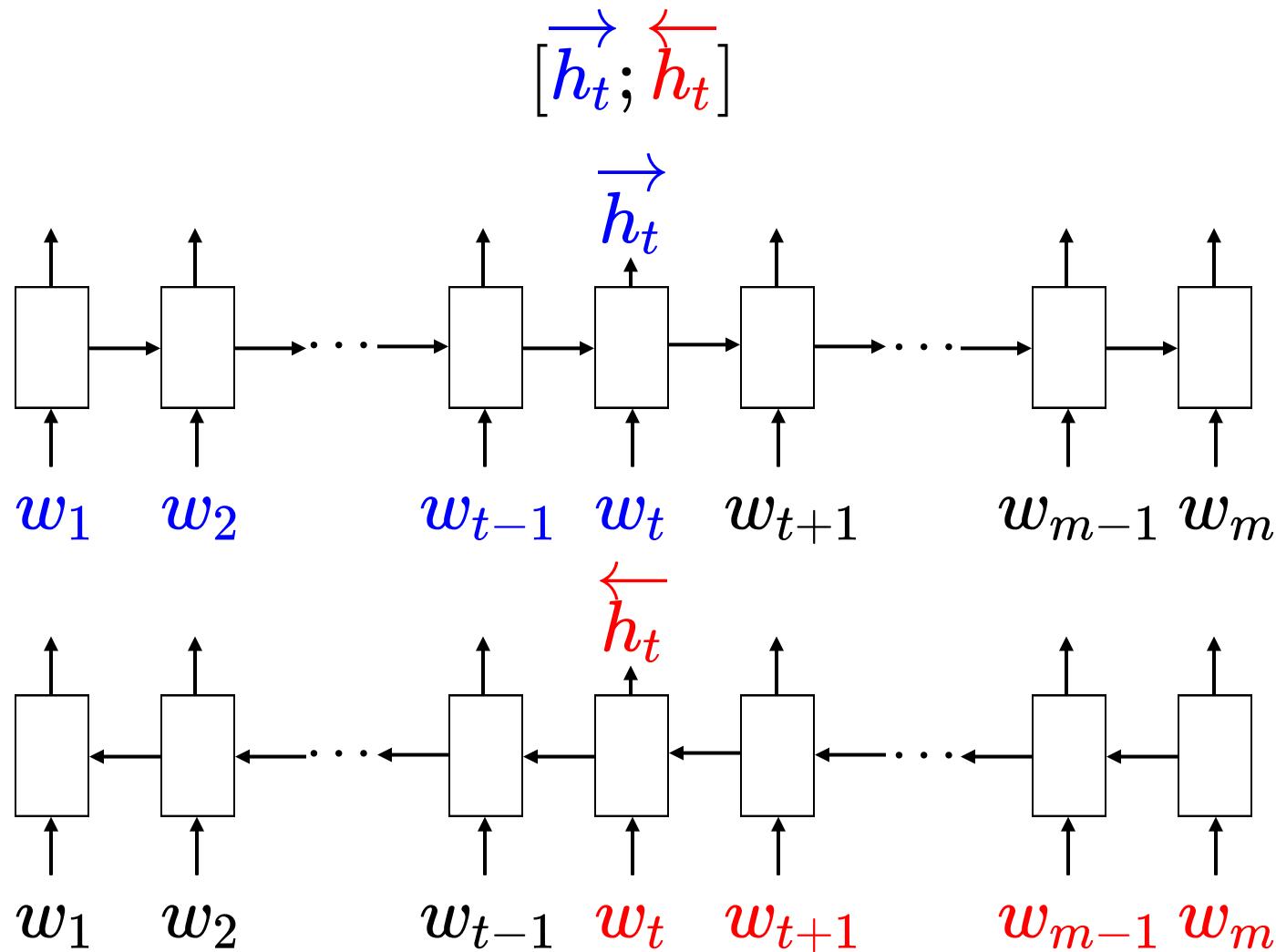
$$\overrightarrow{h_{m-1}} = \overrightarrow{\text{LSTM}}(w_1, w_2, \dots, w_{m-1})$$

Sentence Embedding!

It is essentially a special context embedding.

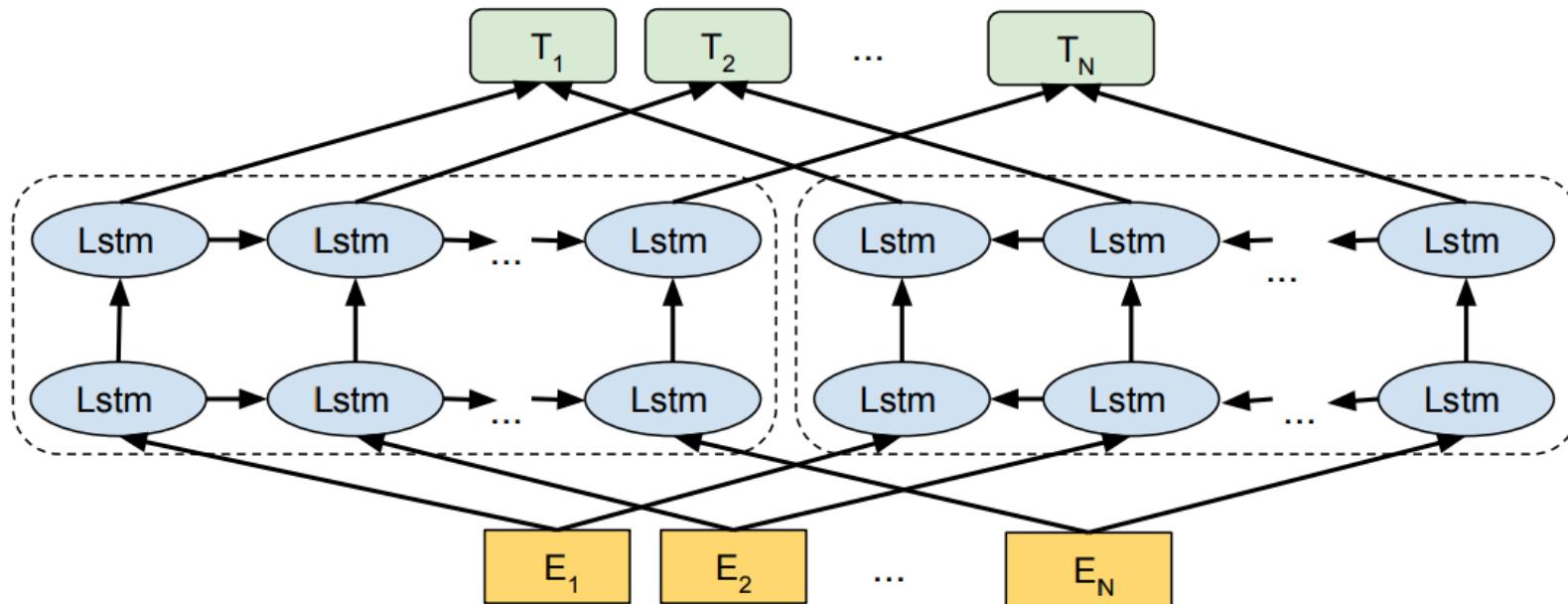
# Bidirectional LSTM

Context Embedding



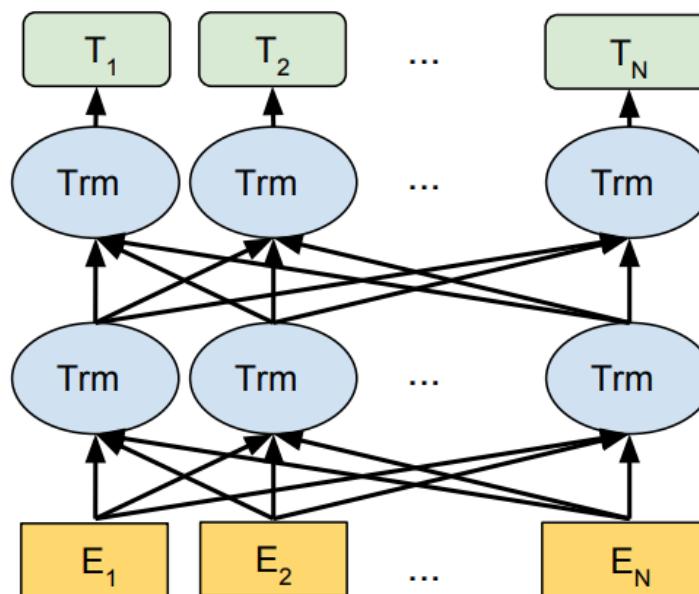
# ELMo

Embeddings from Language Models  
Trained with Two LSTM LMs  
(Peters et al. 2018)



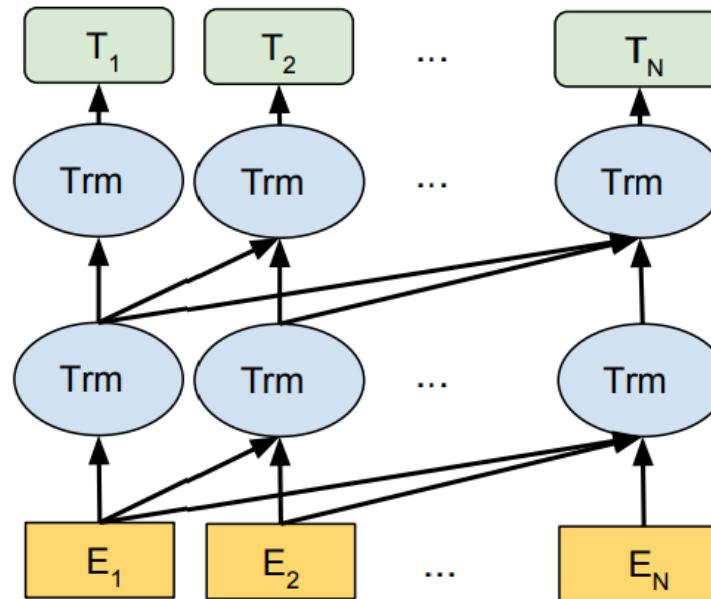
# BERT

Bidirectional Transformer Encoder  
Trained with Masked LM  
(Devlin et al. 2019)

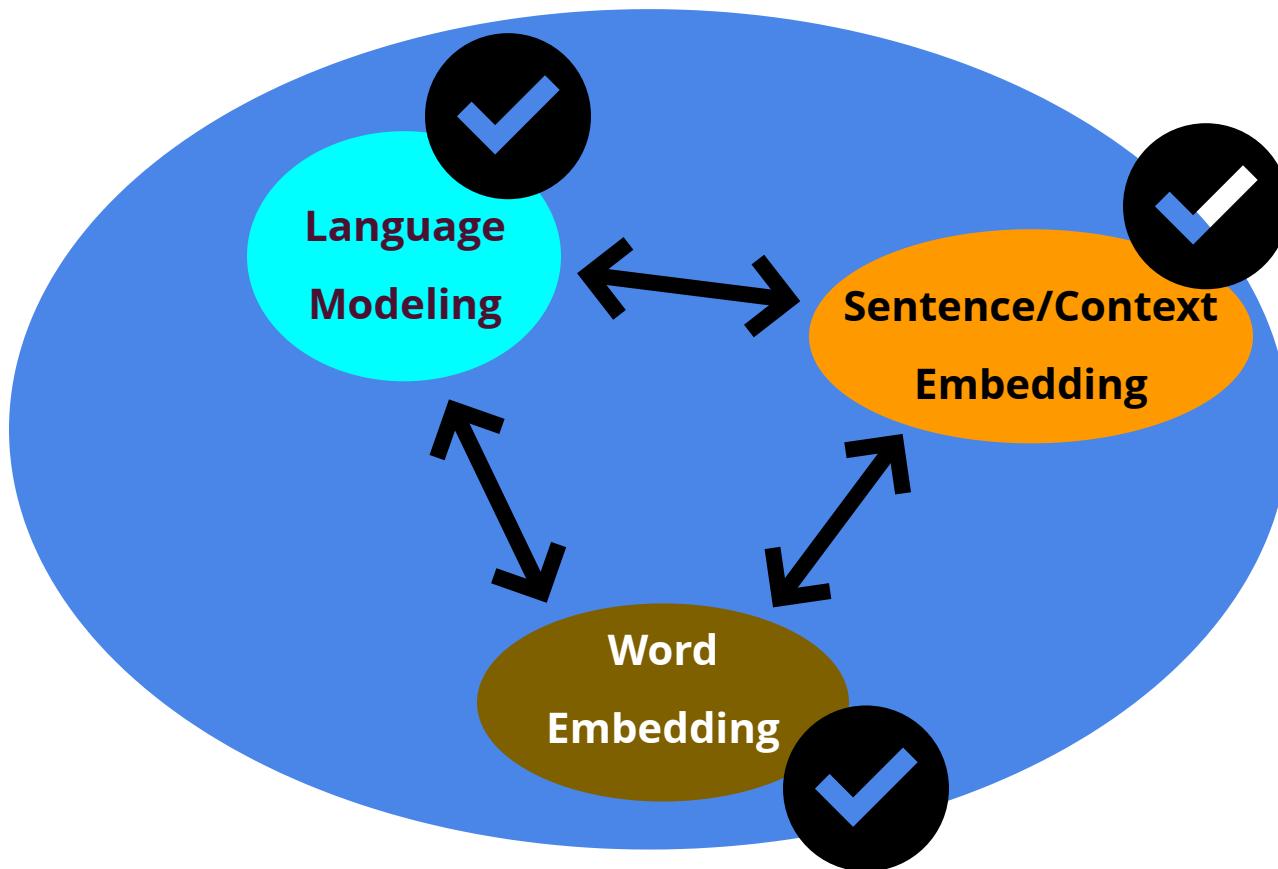


# GPT-3

## (Directional) Transformer Encoder Trained with LM (Brown et al. 2020)

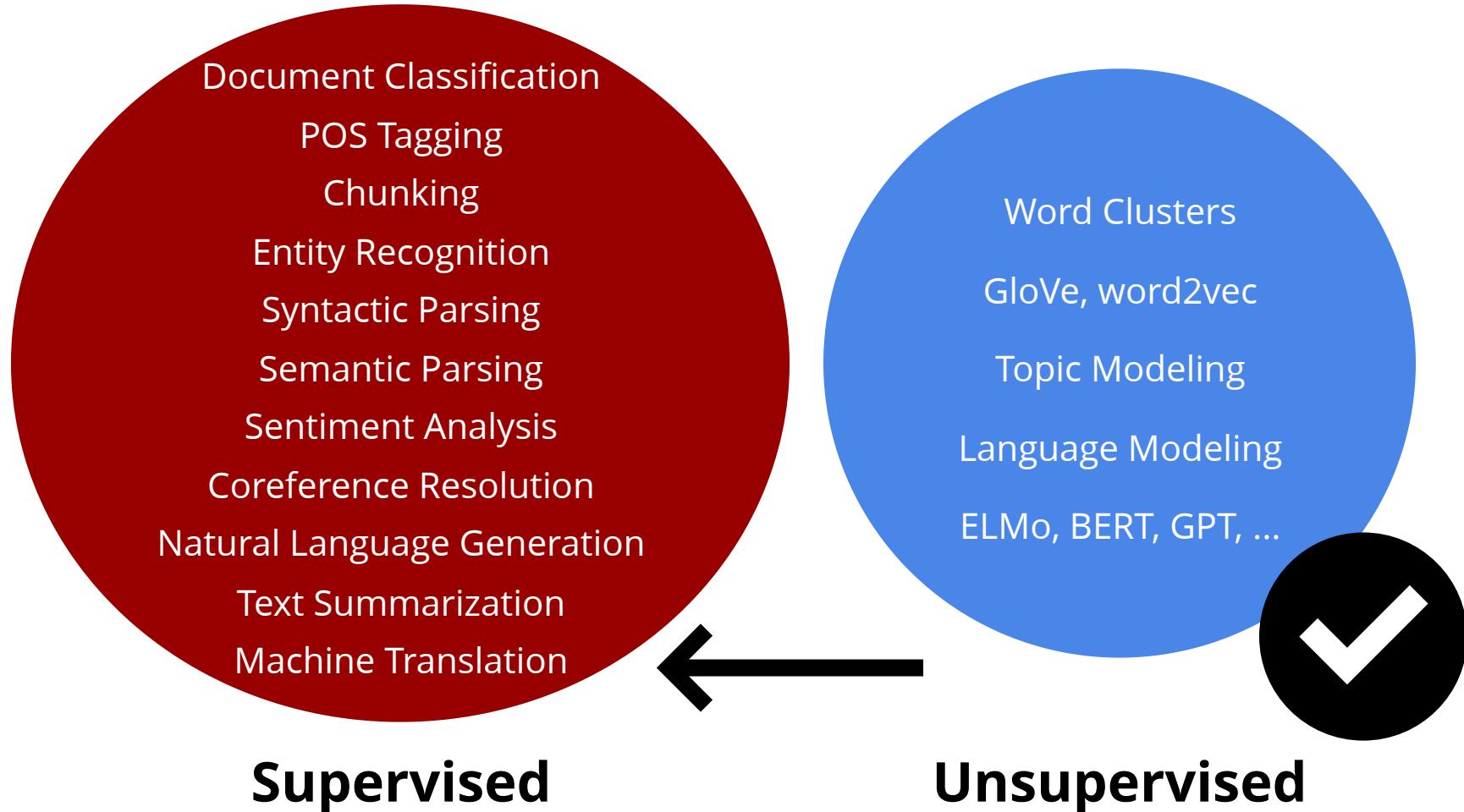


# Tasks in NLP

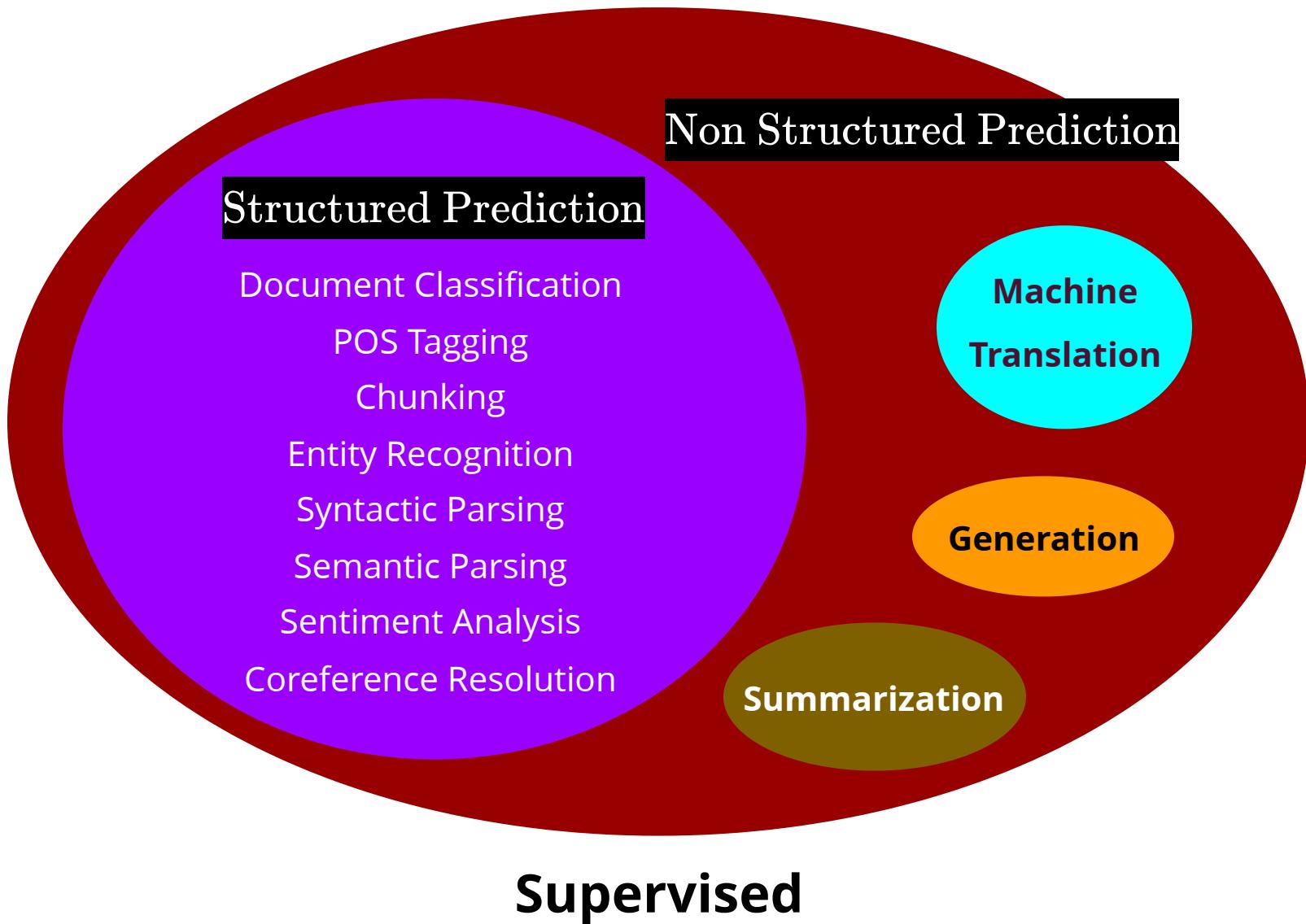


Unsupervised

# Tasks in NLP

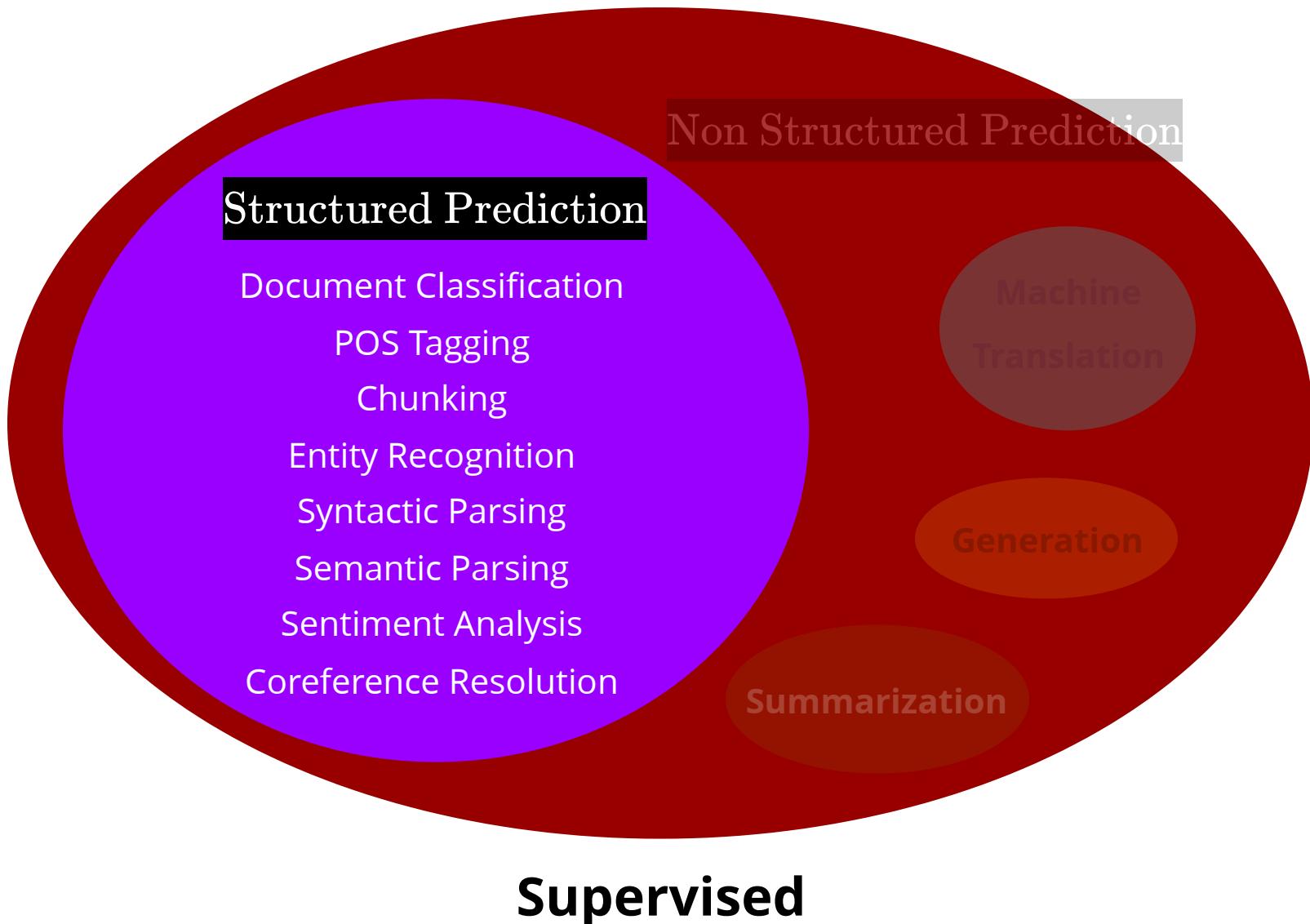


# Tasks in NLP



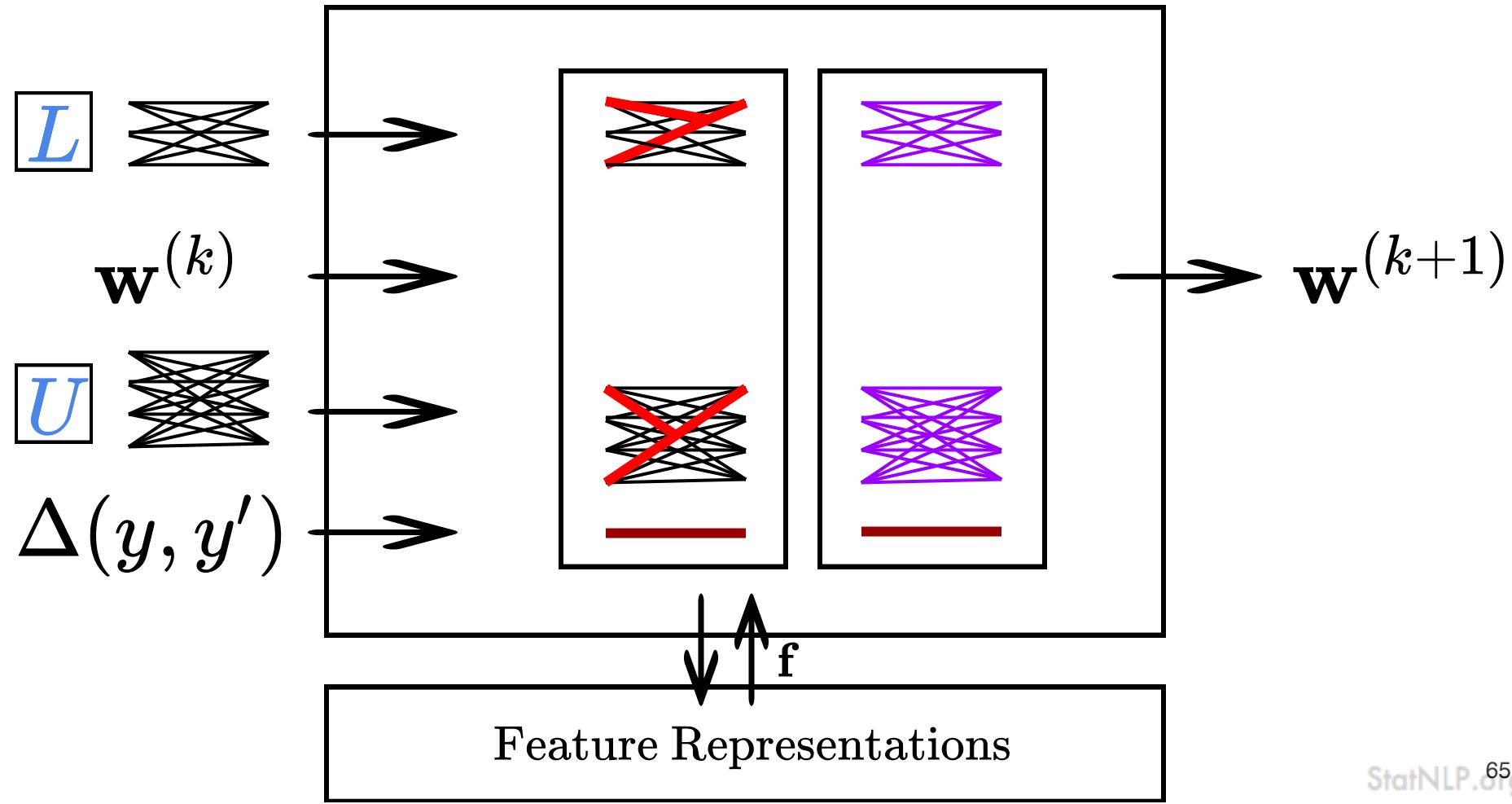
**Supervised**

# Tasks in NLP



**Supervised**

# A Unified Framework for Structured Prediction



# Structured Prediction

## One Assumption

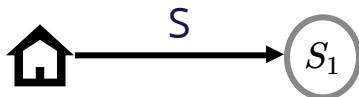
Structures are constructed  
by following a collection of  
discrete actions.

# States, Actions

S: shift

L: left-arc

R:right-arc



↓  
*Fruit*

*flies*

*like*

*a*

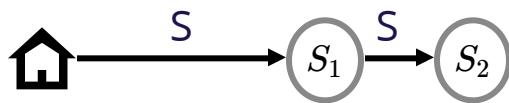
*banana*

# States, Actions

S: shift

L: left-arc

R:right-arc

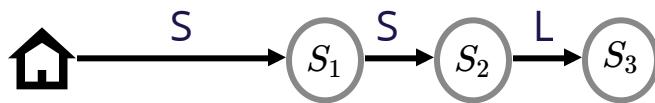


# States, Actions

S: shift

L: left-arc

R:right-arc

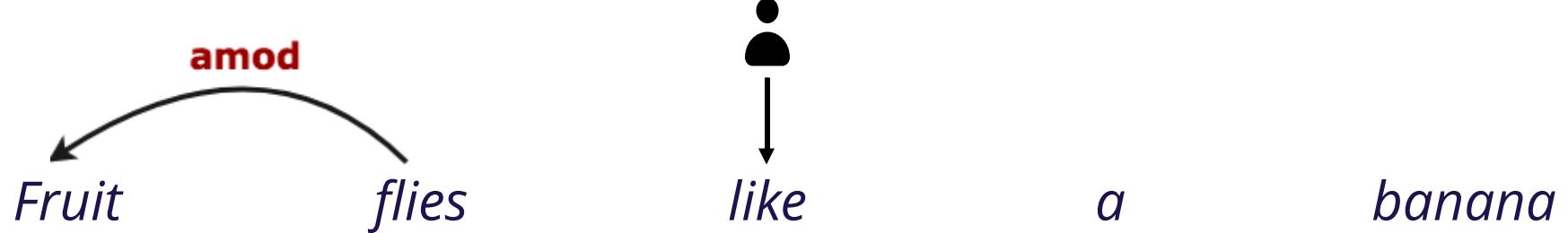
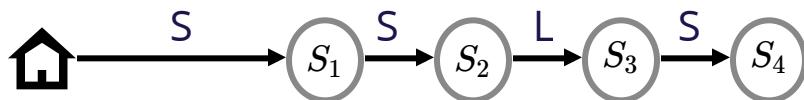


# States, Actions

S: shift

L: left-arc

R:right-arc

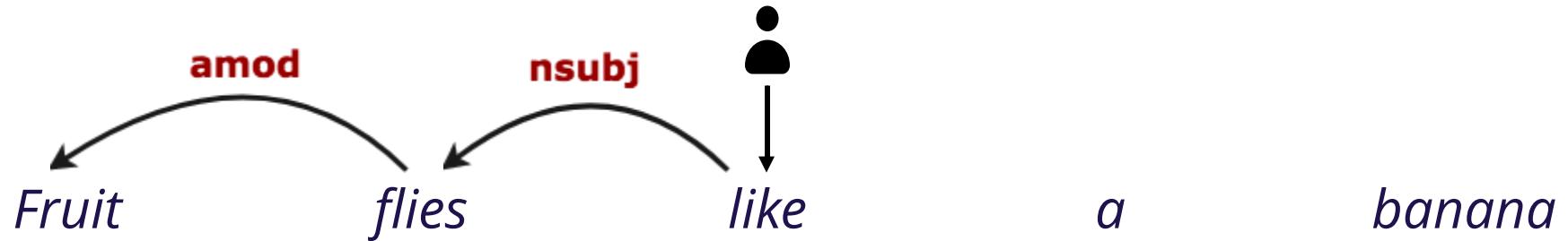
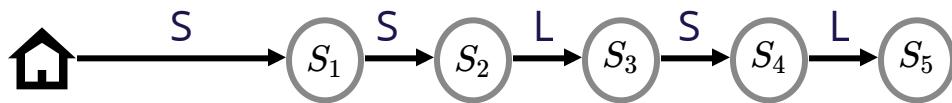


# States, Actions

S: shift

L: left-arc

R:right-arc

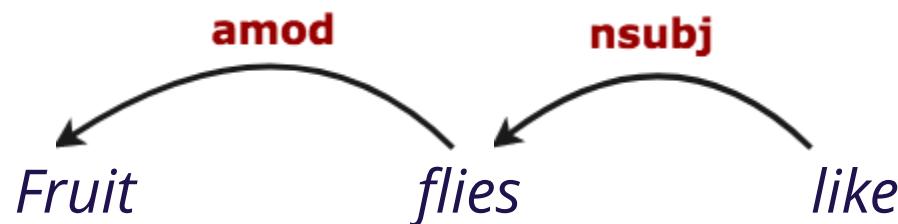
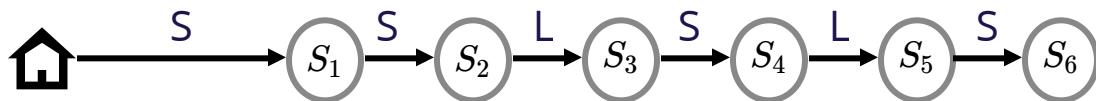


# States, Actions

S: shift

L: left-arc

R:right-arc



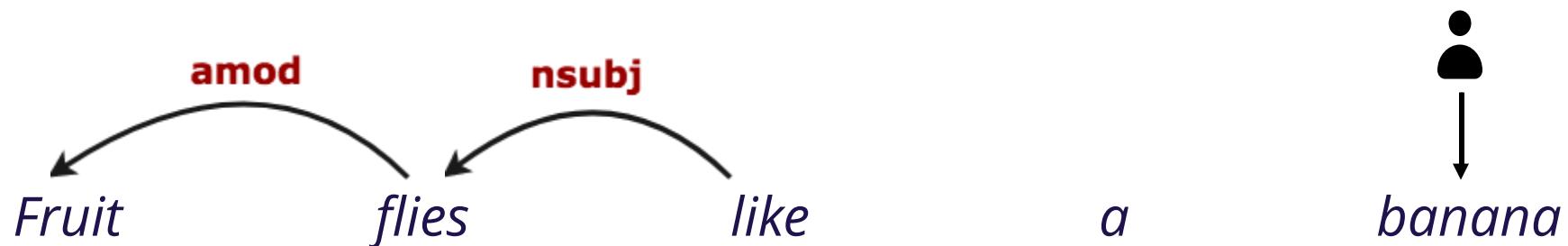
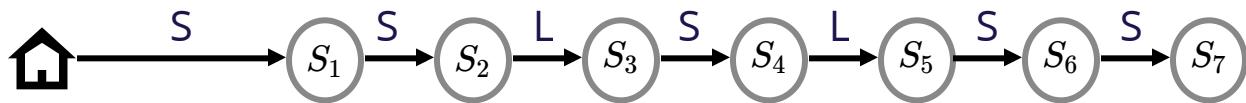
banana

# States, Actions

S: shift

L: left-arc

R:right-arc

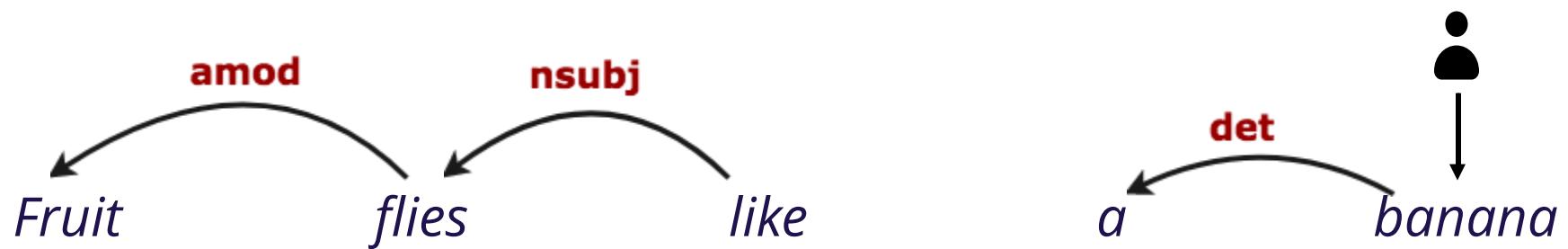
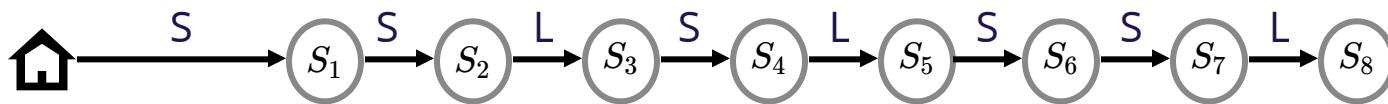


# States, Actions

S: shift

L: left-arc

R:right-arc

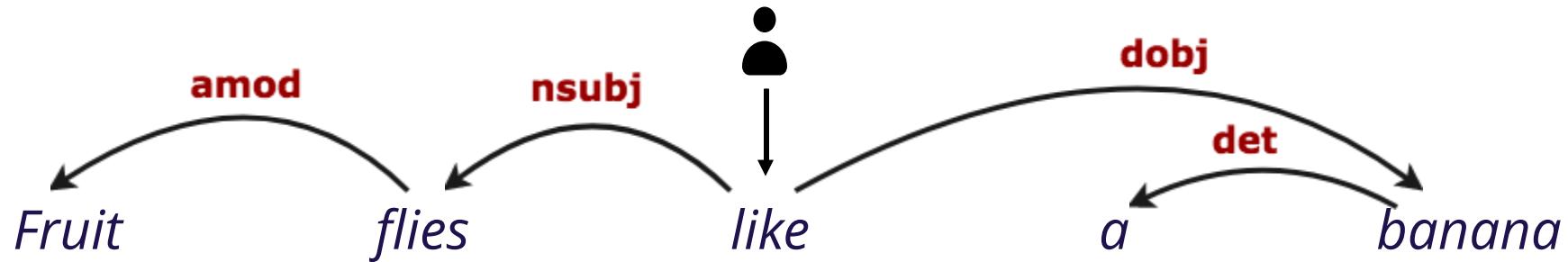
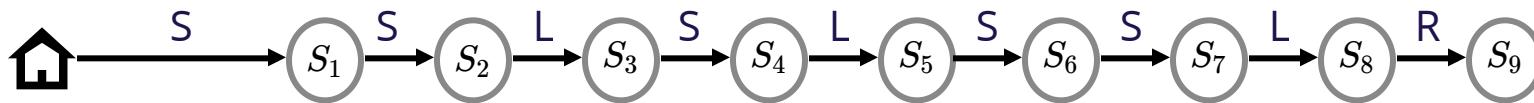


# States, Actions

S: shift

L: left-arc

R:right-arc

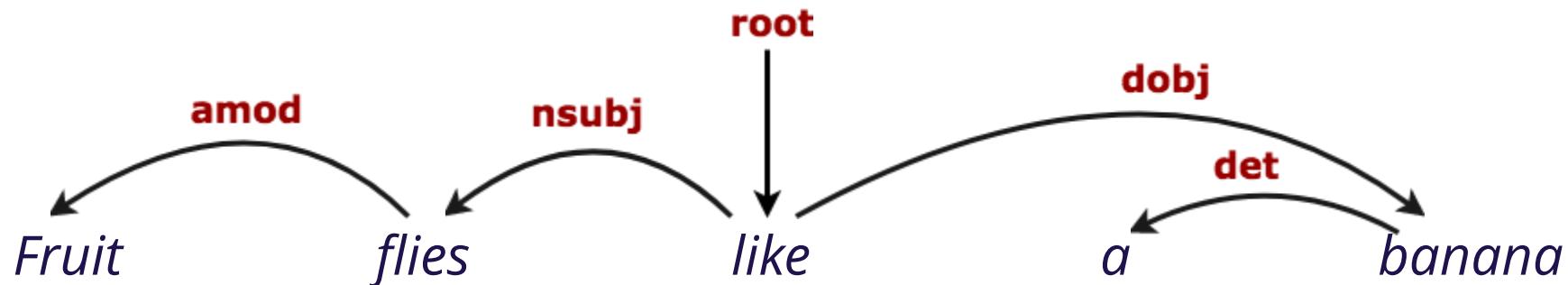
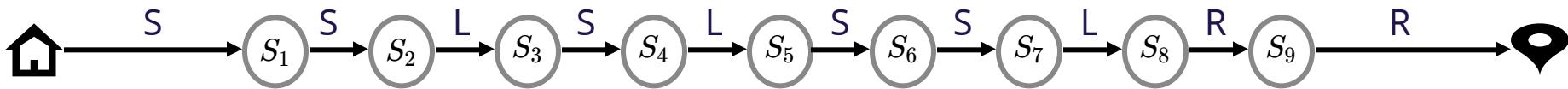


# States, Actions

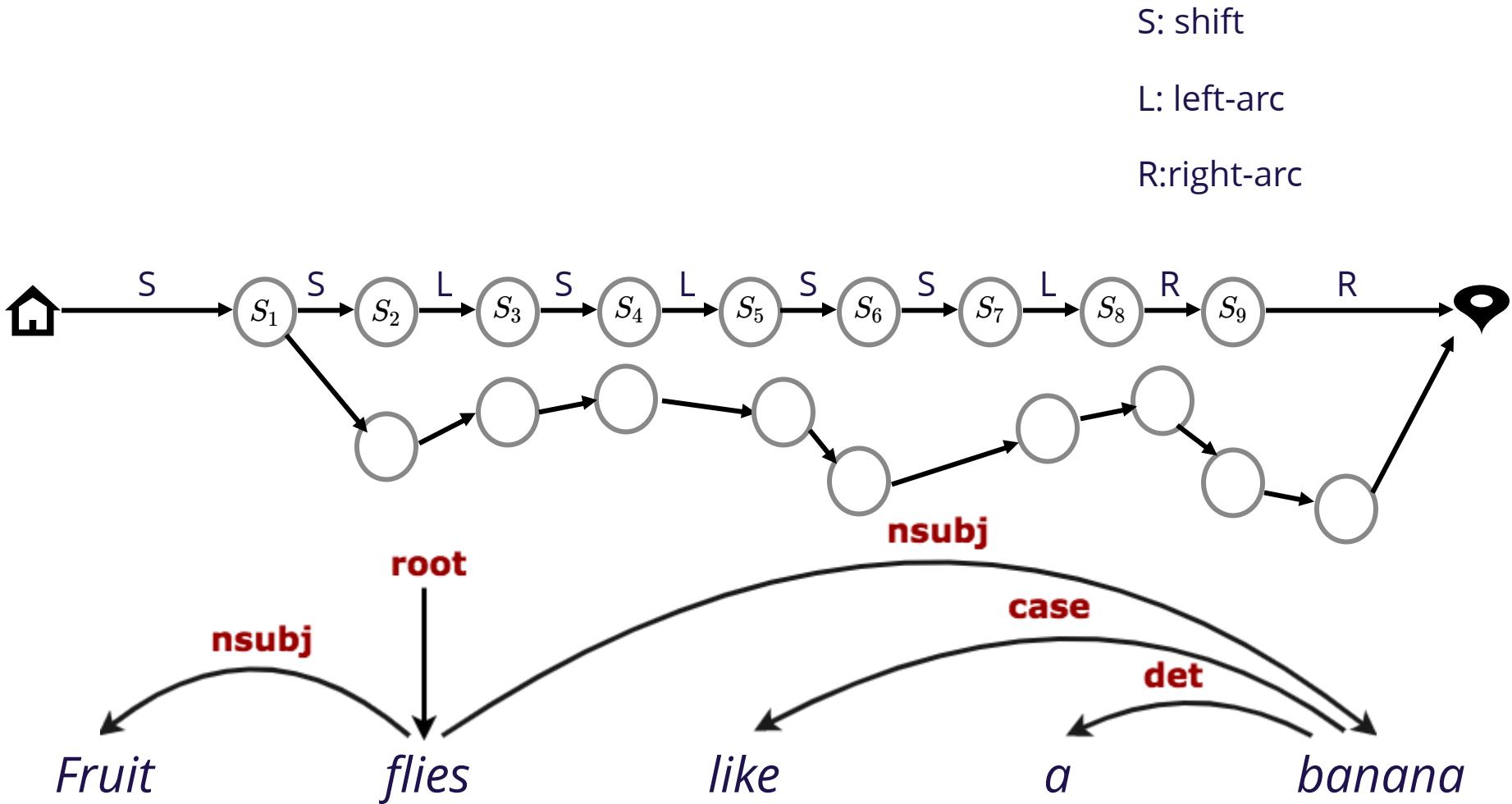
S: shift

L: left-arc

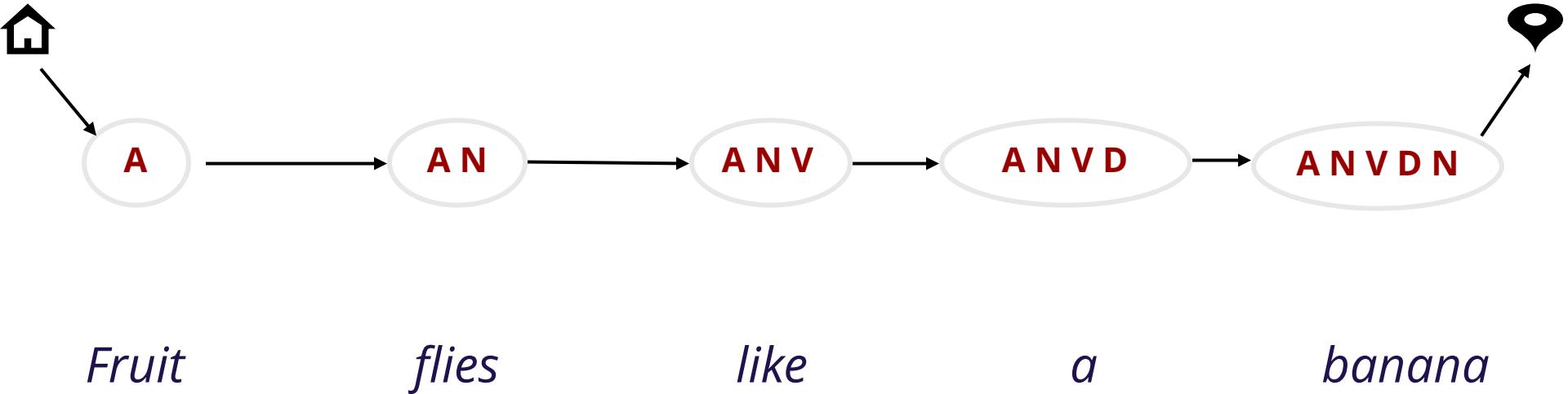
R:right-arc



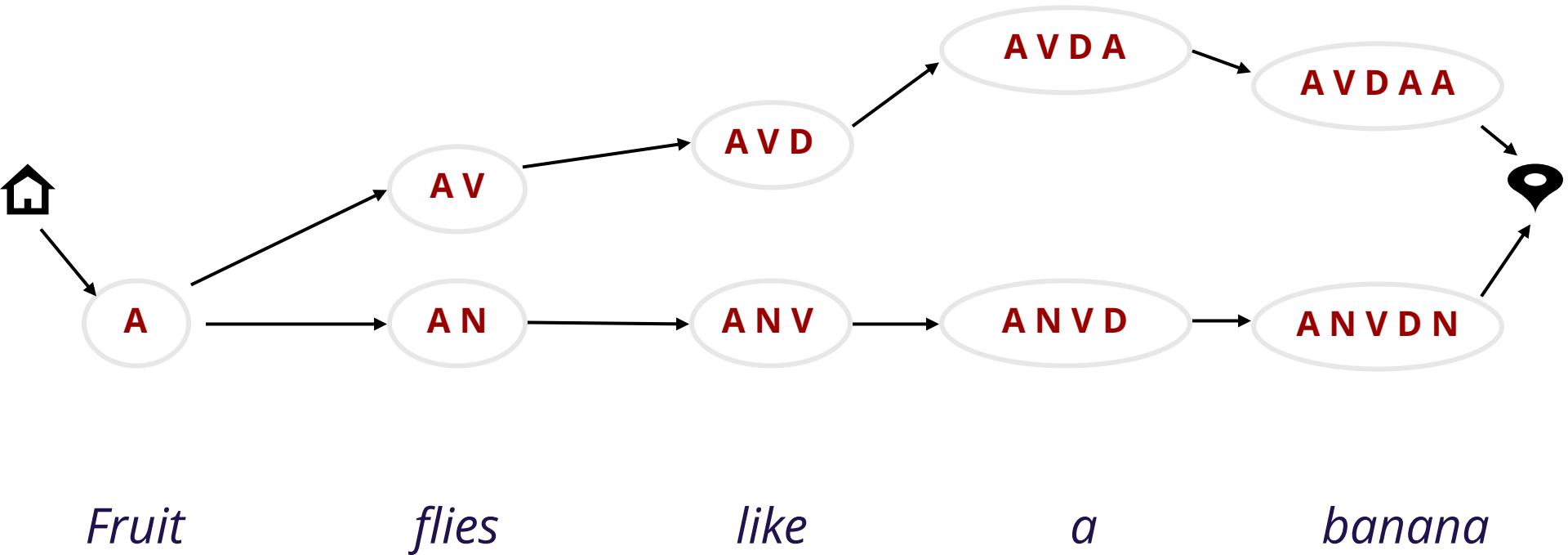
# States, Actions



# States, Actions, Paths



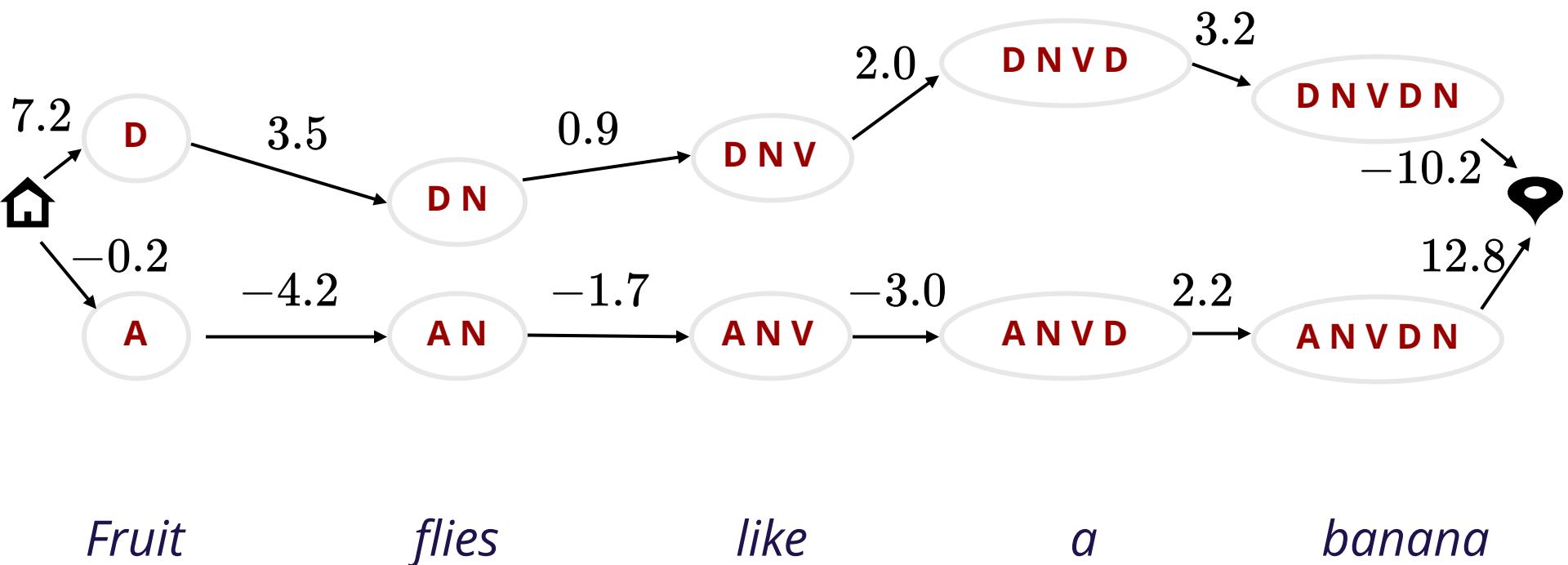
# States, Actions, Paths



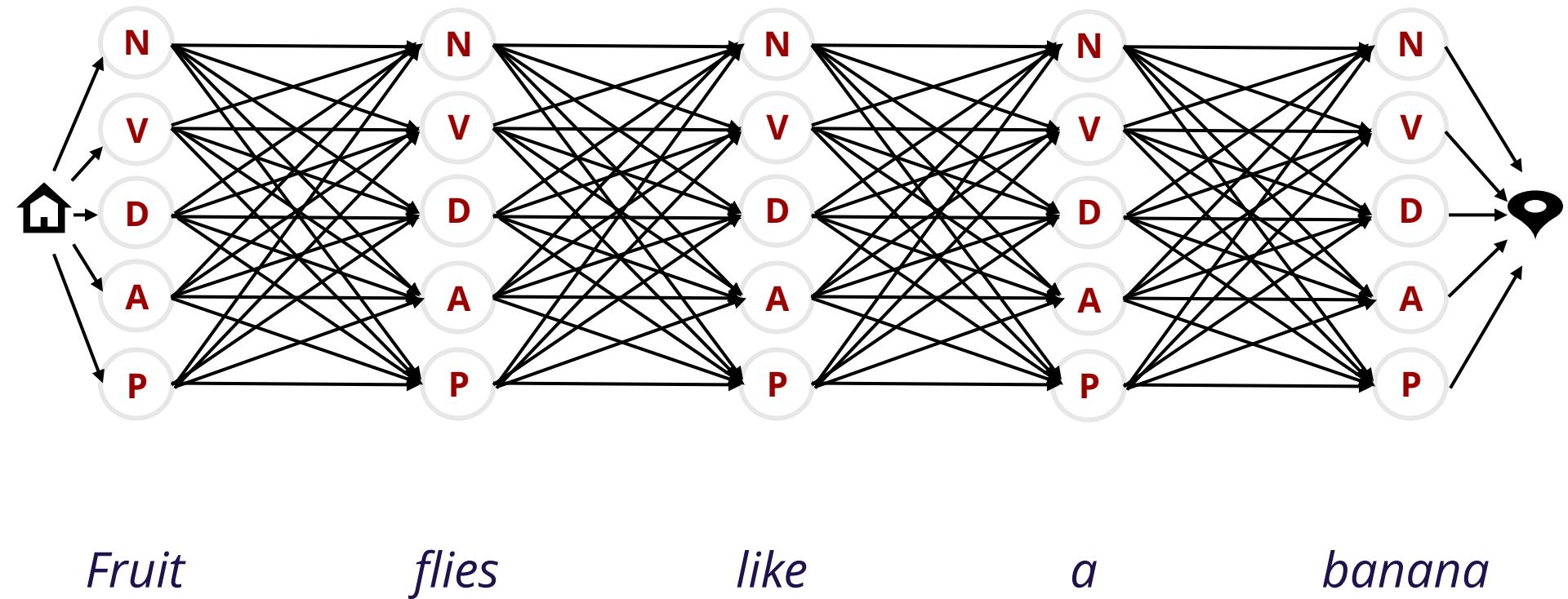
# Score of a Path

$$S_{\mathbf{w}}(p) = \sum_{e \in p} s_{\mathbf{w}}(e)$$

$$s_{\mathbf{w}}(e) = \mathbf{w} \cdot \mathbf{f}(e)$$

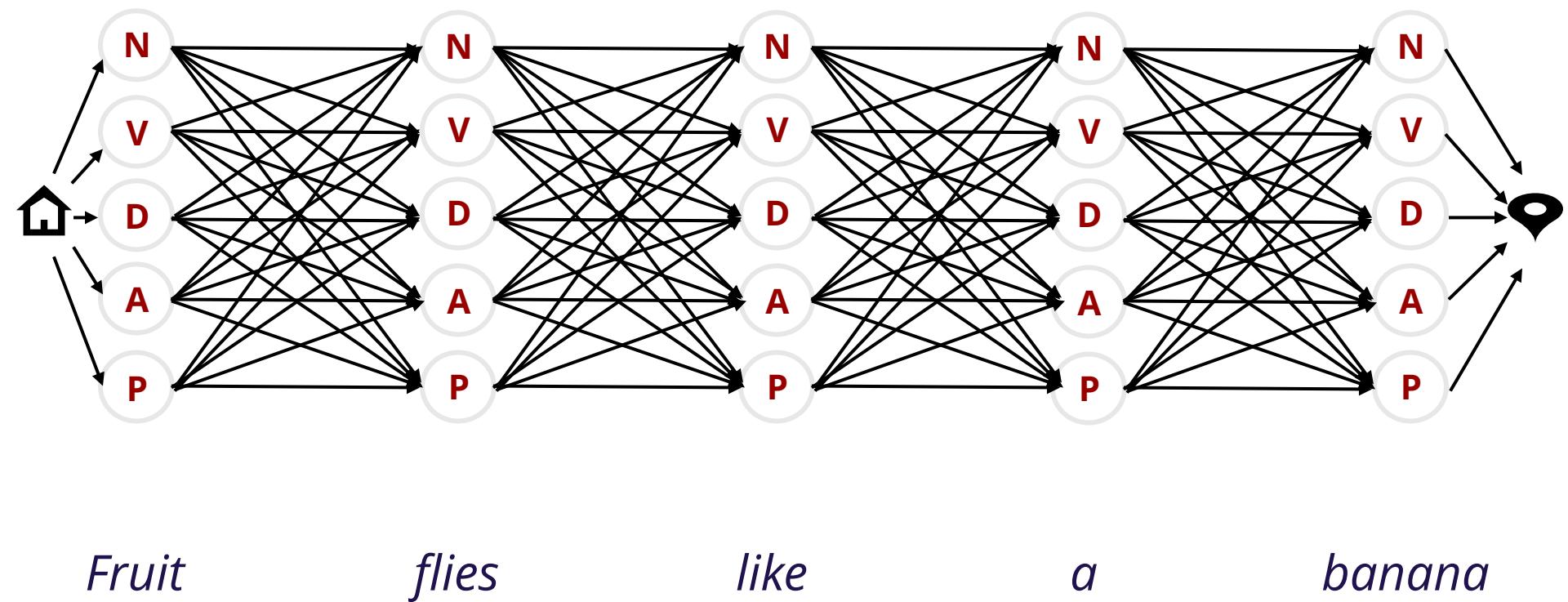


# Compact Search Graph



# Structured Prediction

Given	Find
$x, \mathbf{w}$	$y$



*Fruit*

*flies*

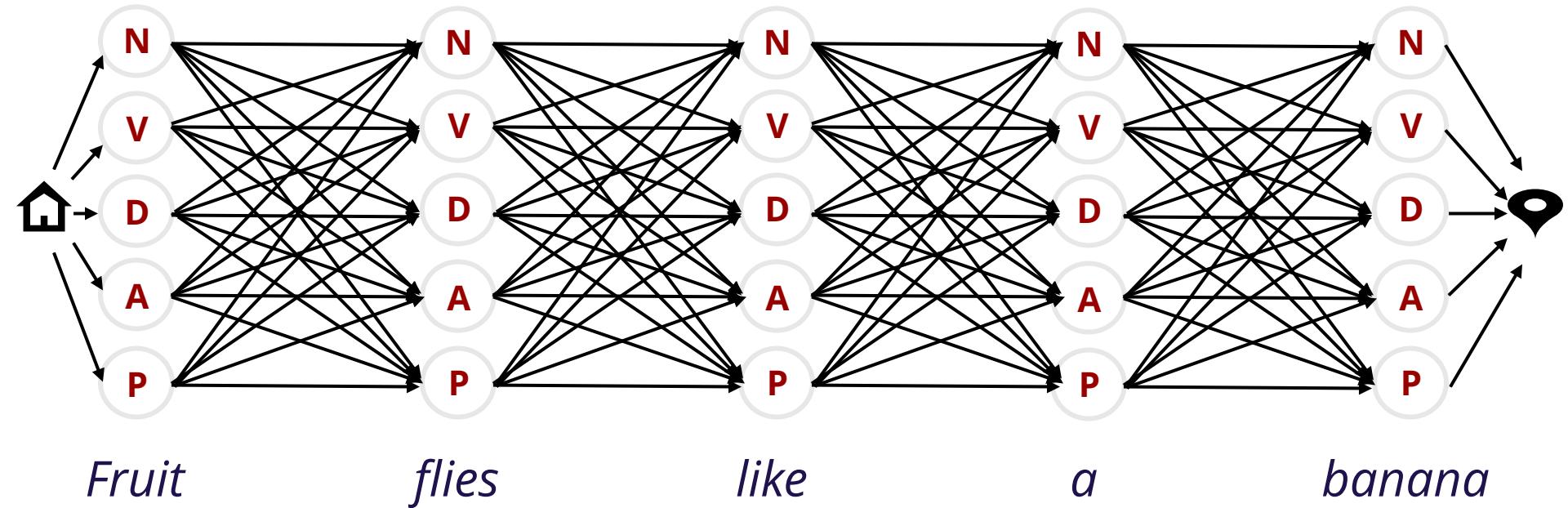
*like*

*a*

*banana*

# The Search Problem

$$\arg \max_y \mathbf{w} \cdot \mathbf{f}(x, y)$$



*Fruit*

*flies*

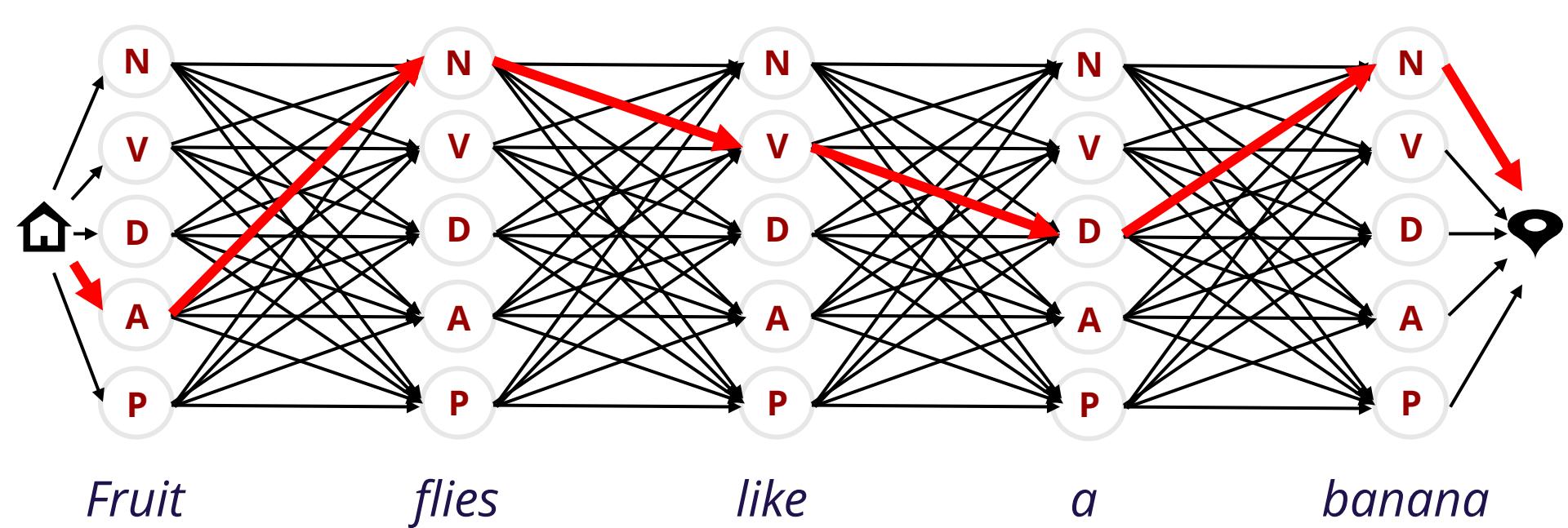
*like*

*a*

*banana*

# MAP Inference

$$\arg \max_y \mathbf{w} \cdot \mathbf{f}(x, y)$$



*Fruit*

*flies*

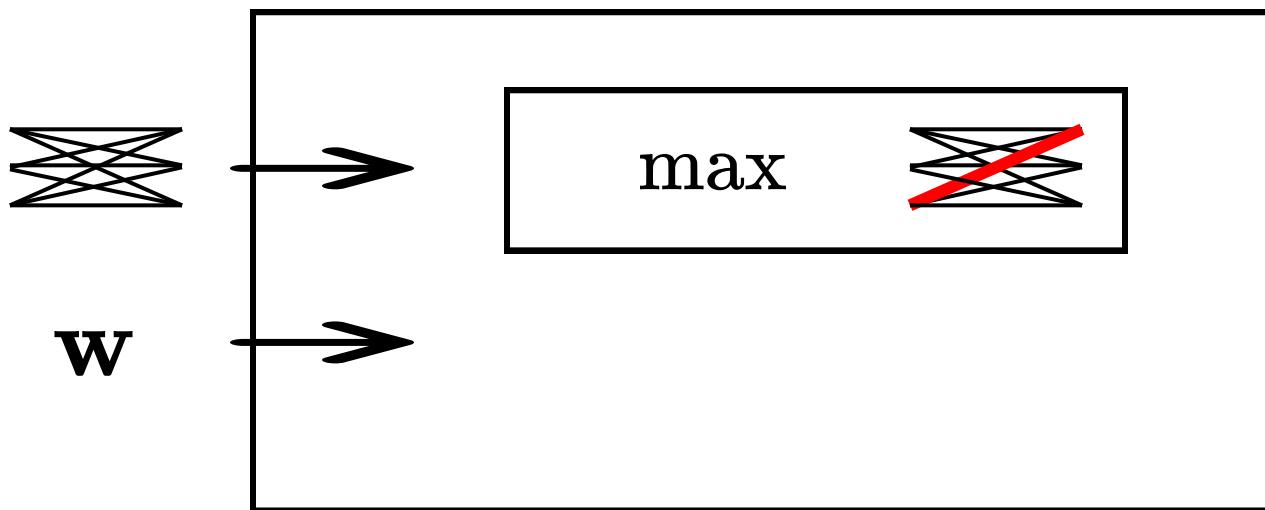
*like*

*a*

*banana*

# Inference

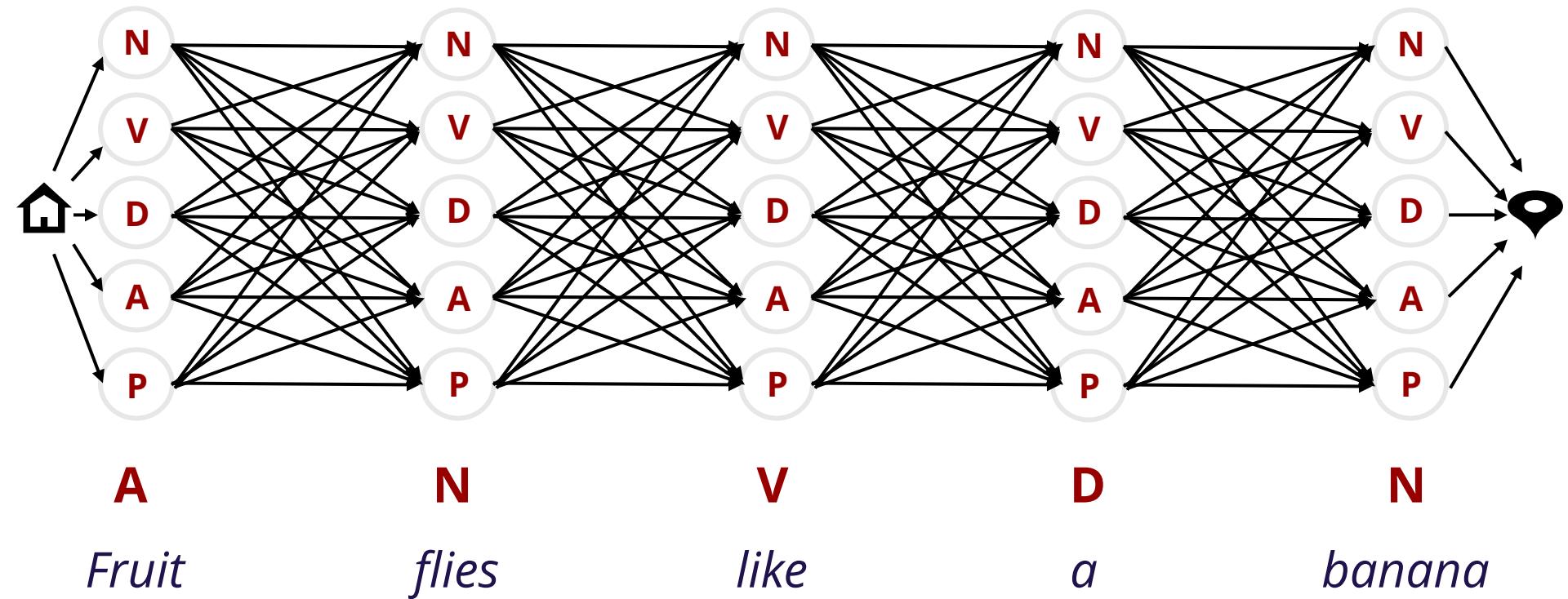
MAP



To find the optimal structure, just construct the graph, and tell the inference module the current weights, and specify the inference algorithm (in this case MAP)

# Learning

Given	Find
$(x_i, y_i)$	$w$



*Fruit*

*flies*

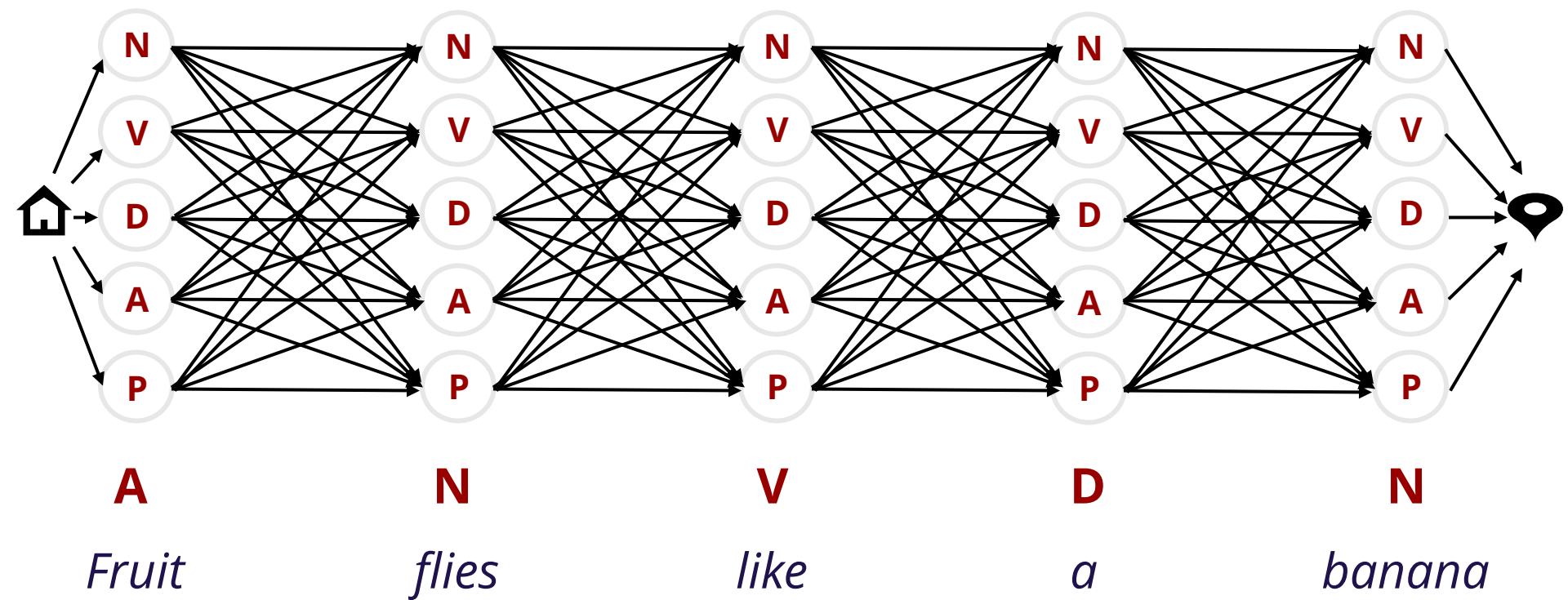
*like*

*a*

*banana*

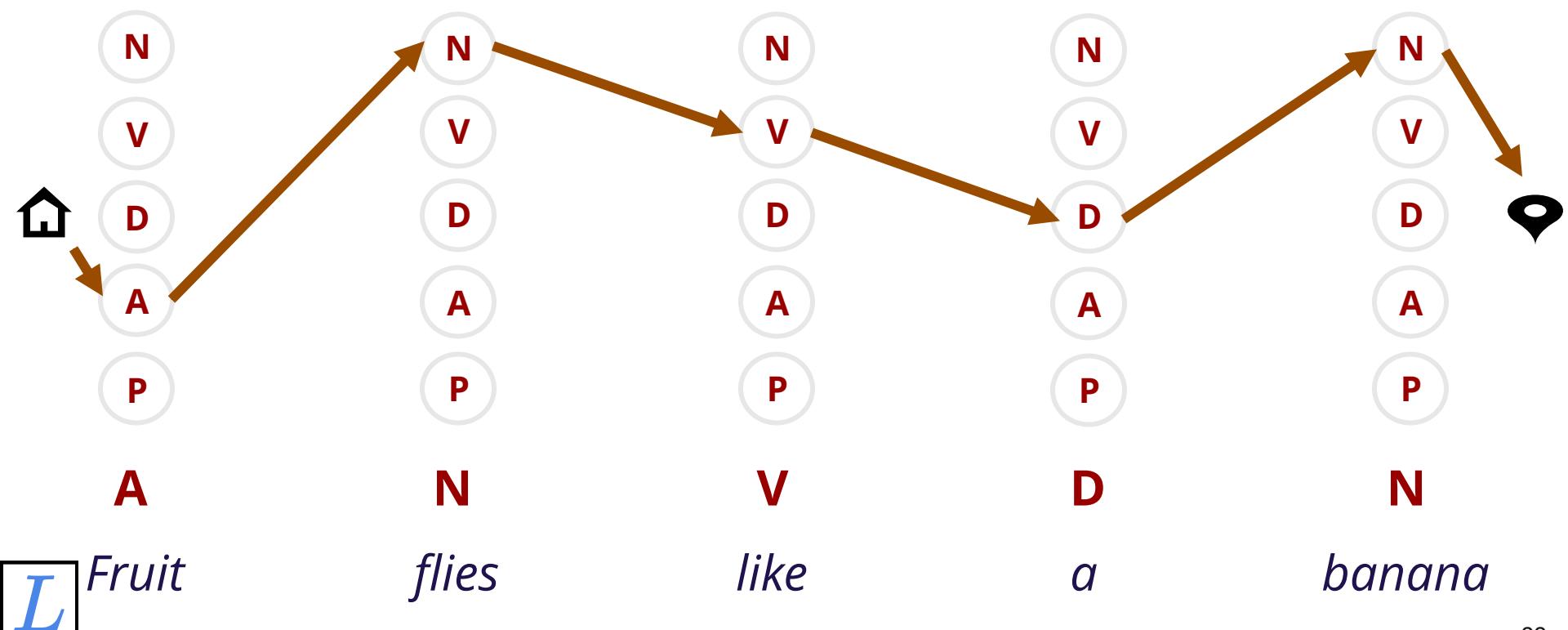
# Structured Perceptron

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \max_y (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$



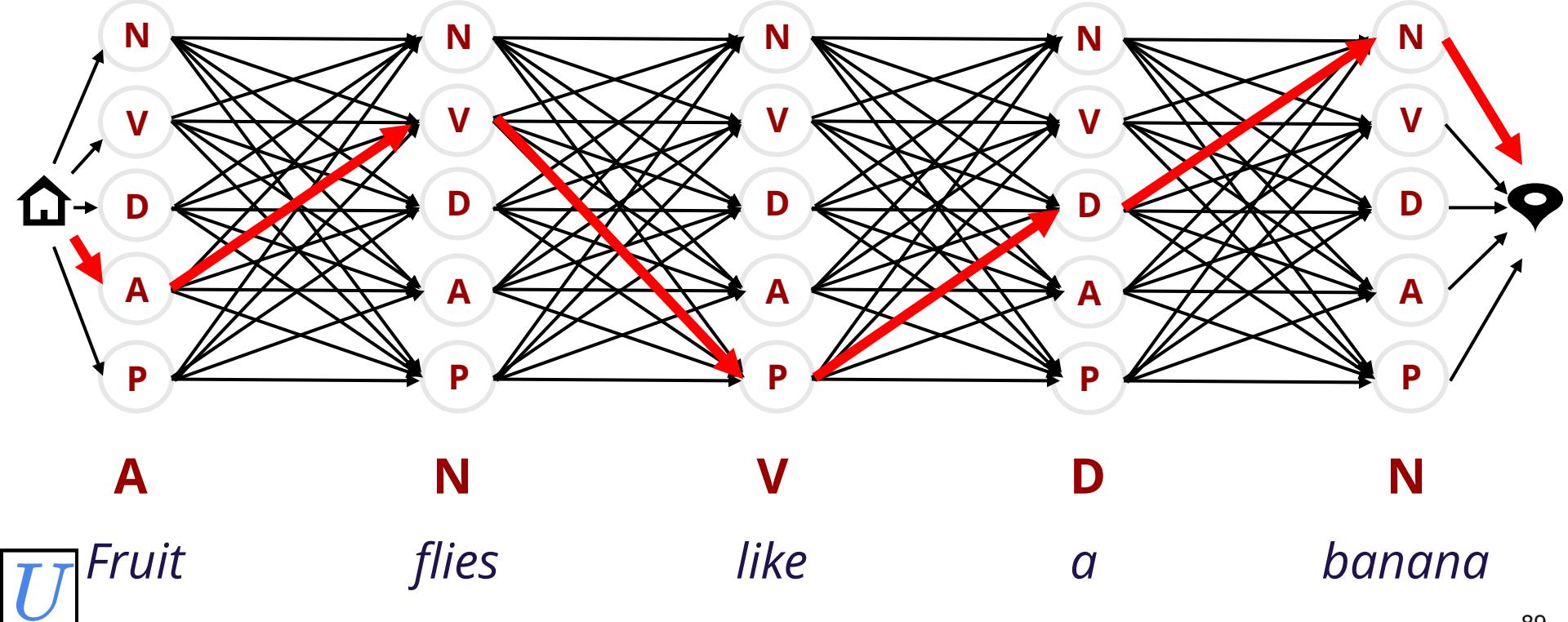
# Structured Perceptron

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \max_y (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$



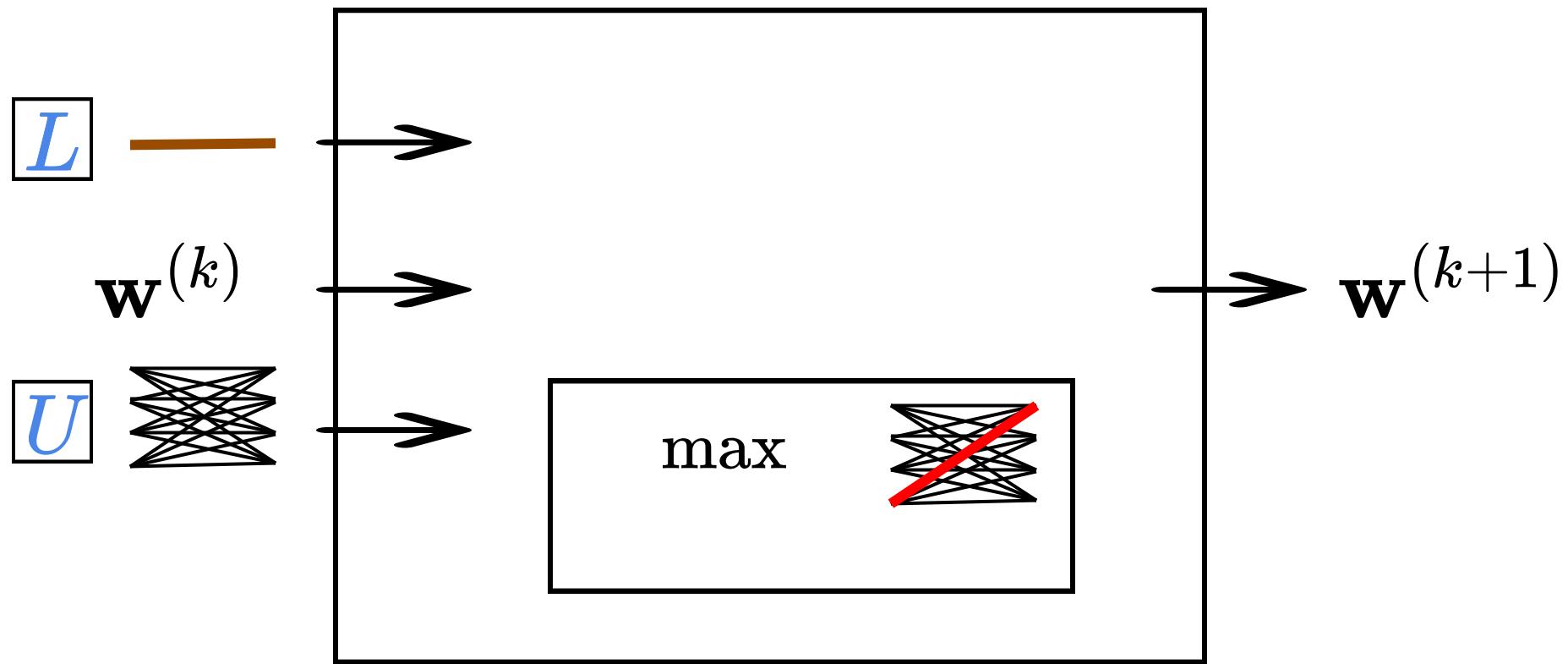
# Structured Perceptron

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \max_y (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$



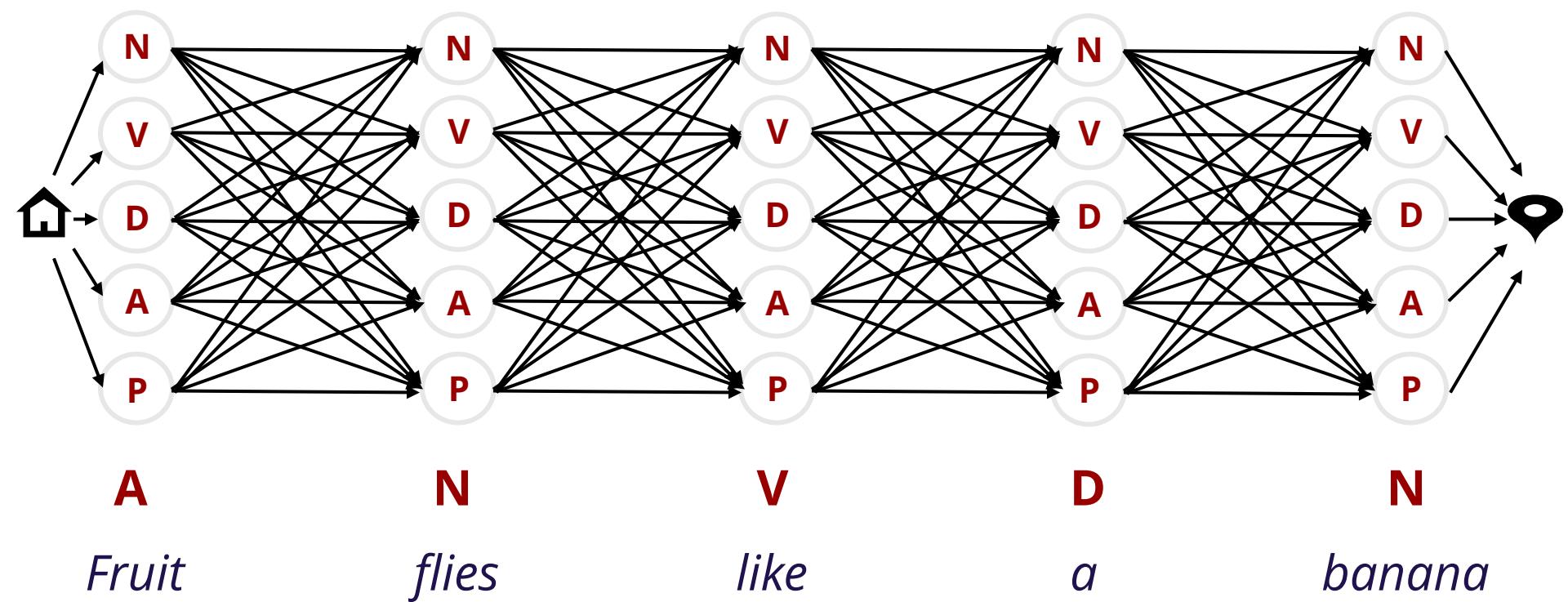
# Learning

## Structured Perceptron



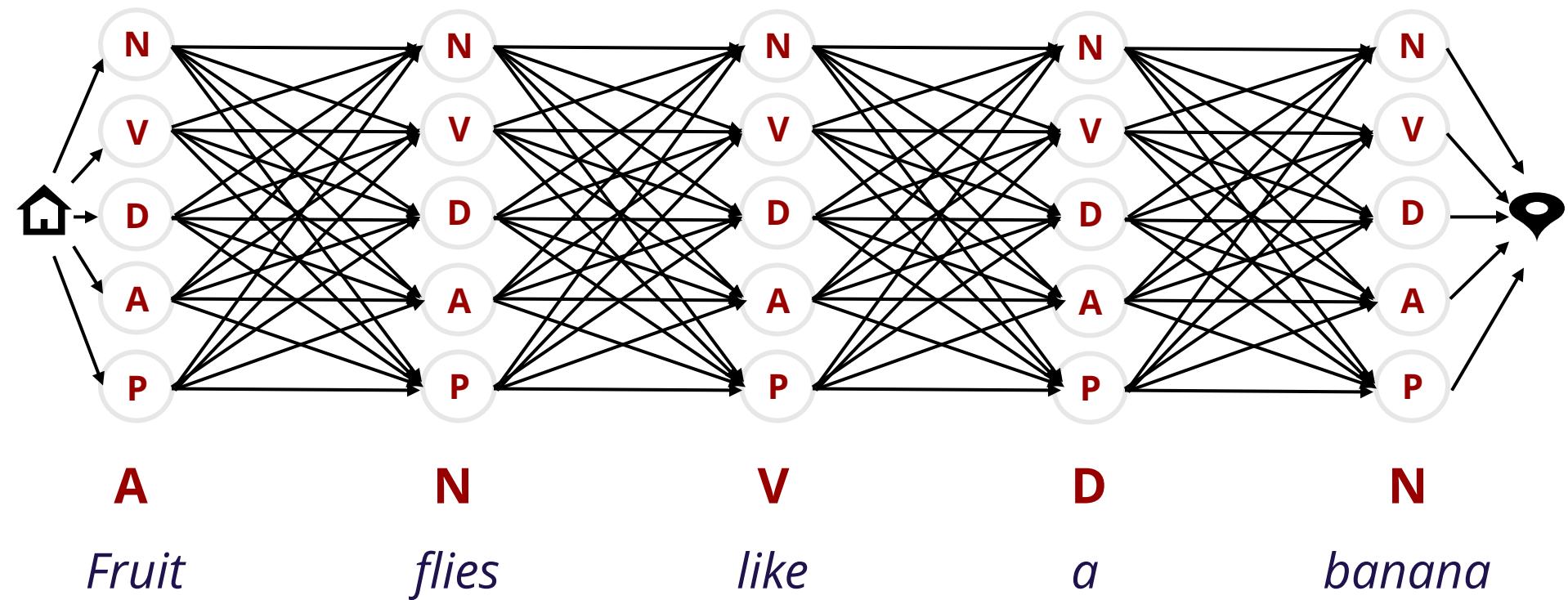
# A Probabilistic View

$$\max_w \sum_i \log p(y_i|x_i)$$



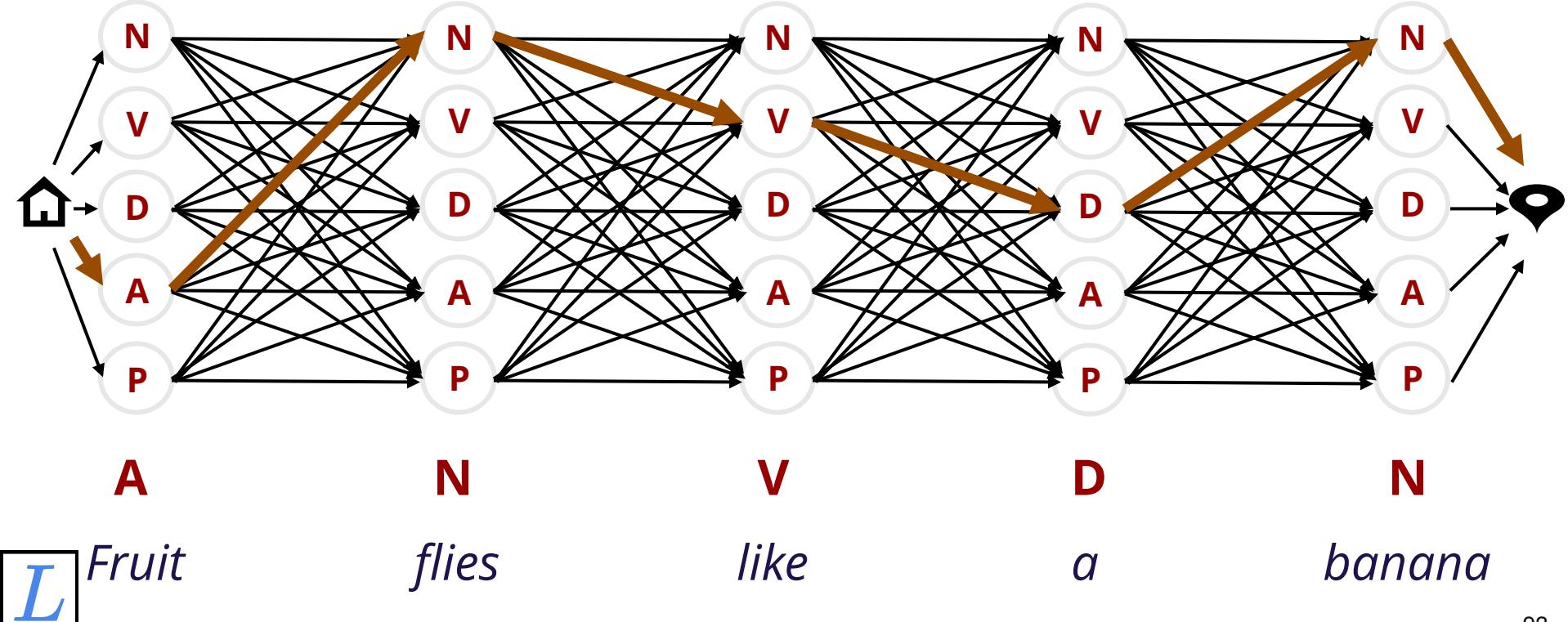
# Conditional Random Fields

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$



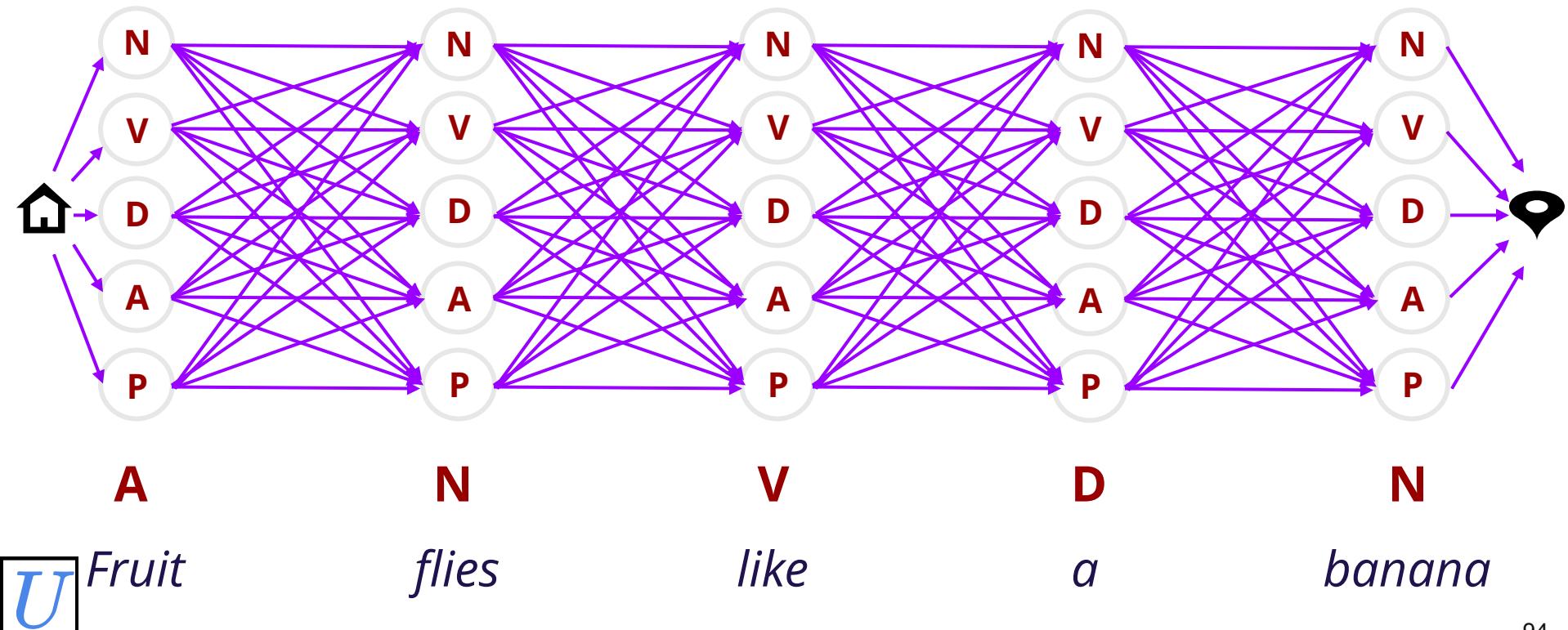
# Conditional Random Fields

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$



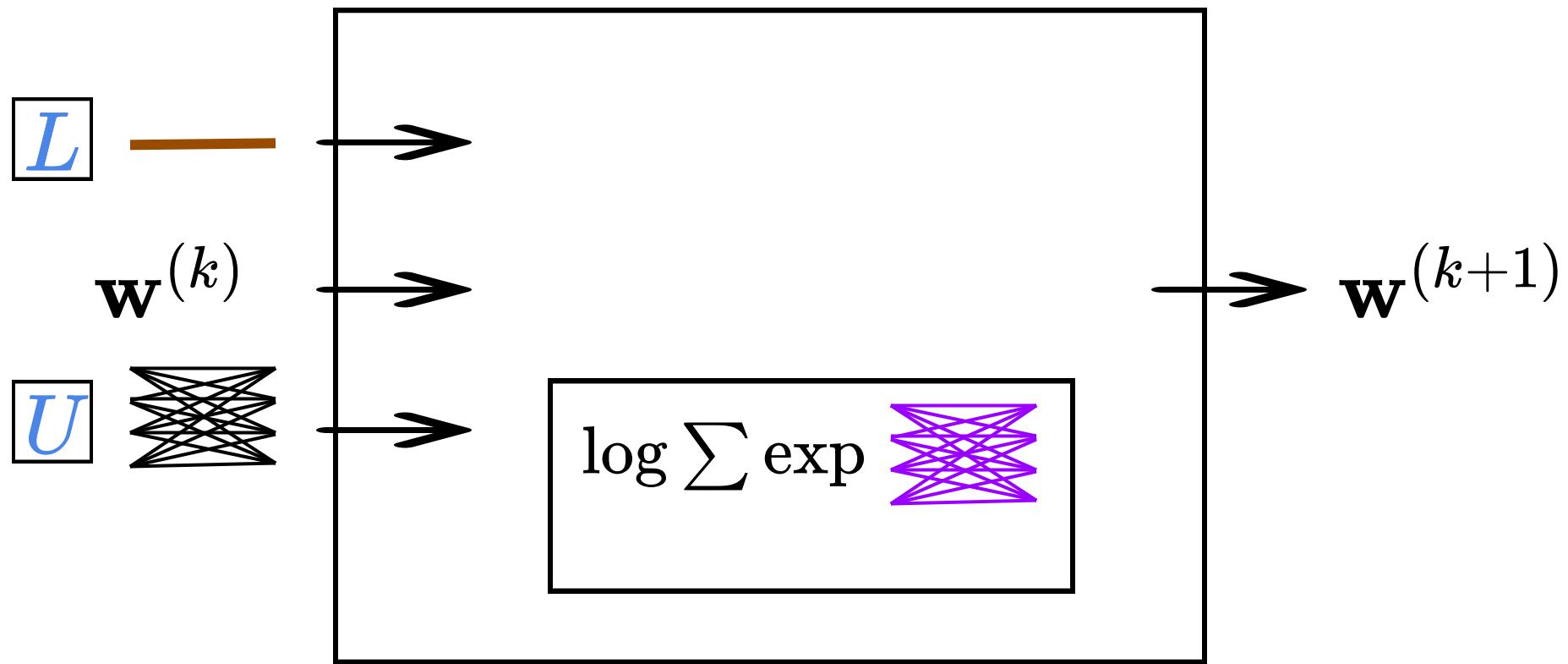
# Conditional Random Fields

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$



# Learning

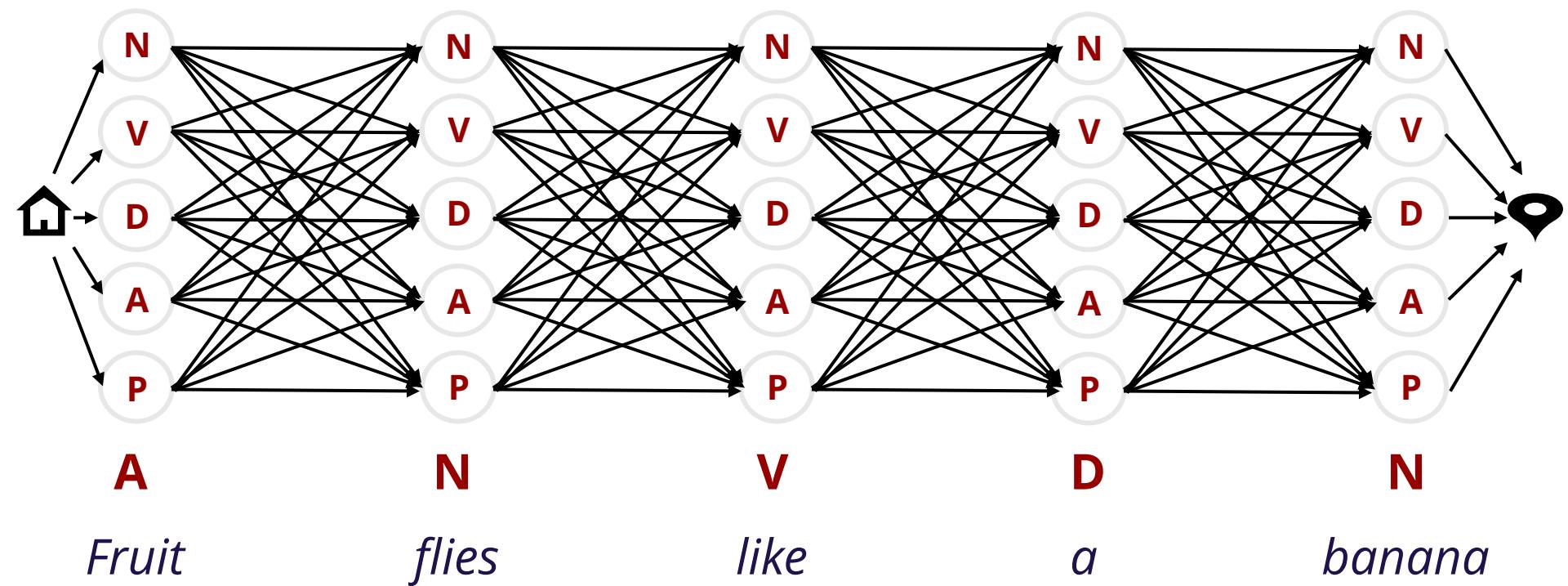
CRF



# Max-Margin

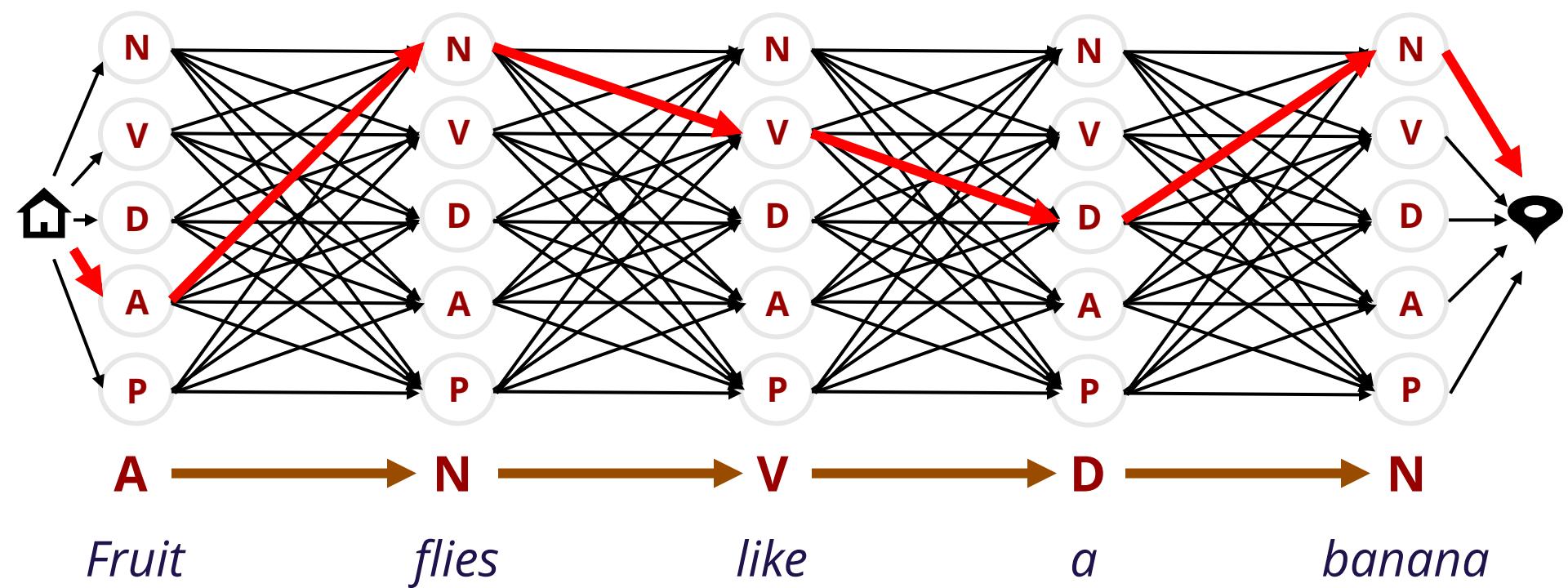
$$\min_w \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \max_y (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$
$$\min_w \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \max_y (\Delta(y_i, y) + \mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$

$$\Delta(\mathbf{D}, \mathbf{P}) = 1, \Delta(\mathbf{N}, \mathbf{N}) = 0, \Delta(\mathbf{N}, \mathbf{V}) = 10$$



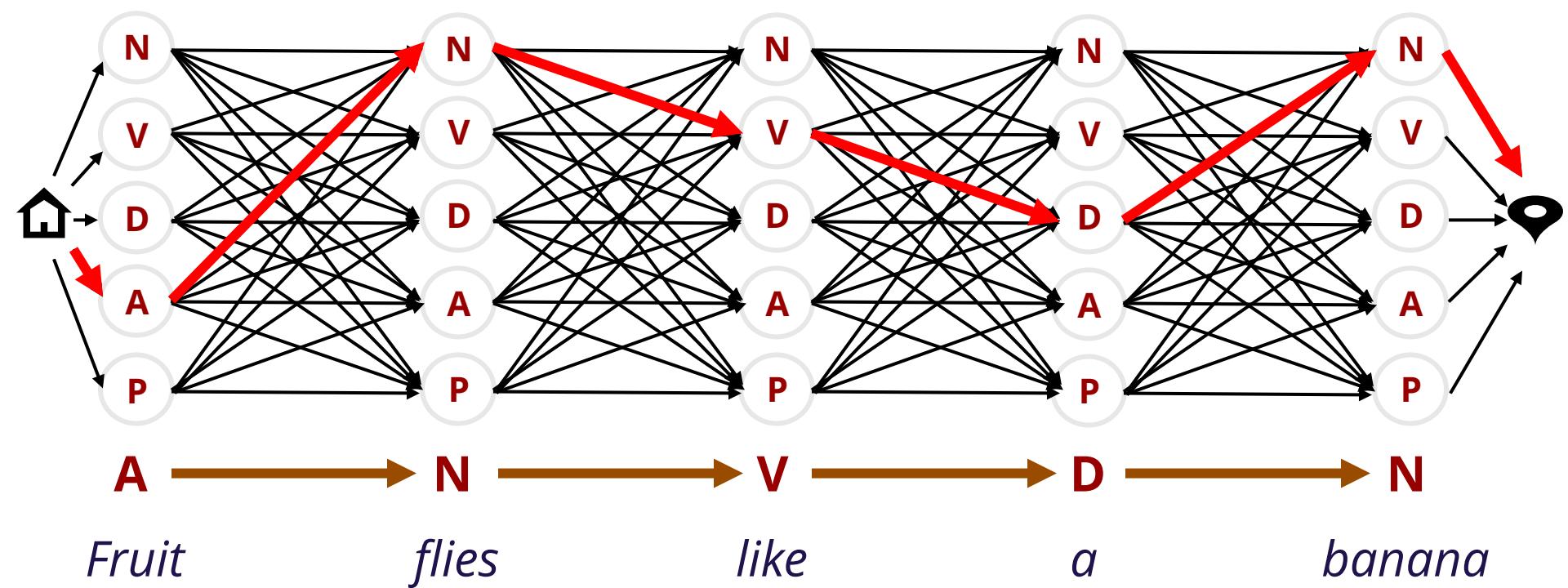
# Decode with Oracle

$$\begin{aligned} \min_{\mathbf{w}} \sum_i & \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \max_y (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right) \\ \min_{\mathbf{w}} \sum_i & \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \max_y (\Delta(y_i, y) + \mathbf{w} \cdot \mathbf{f}(x_i, y)) \right) \\ \Delta(\mathbf{D}, \mathbf{P}) = 1, \Delta(\mathbf{N}, \mathbf{N}) = 0, \Delta(\mathbf{N}, \mathbf{V}) = 10 \end{aligned}$$



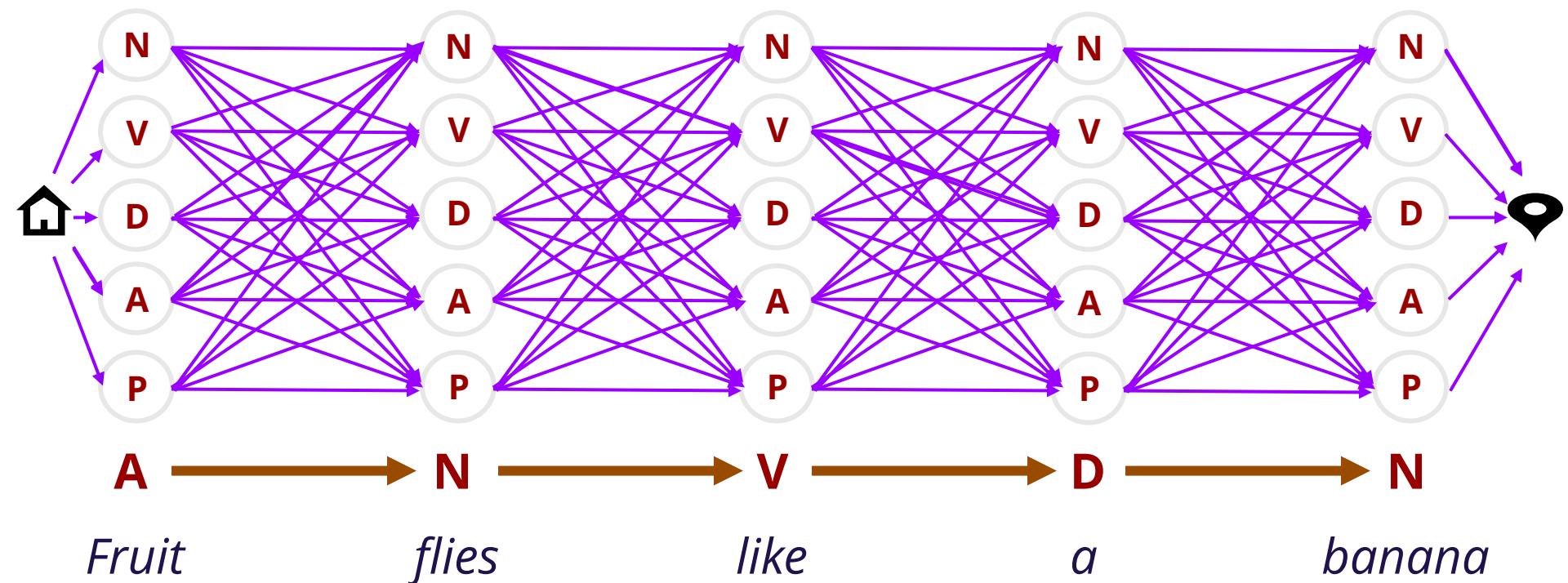
# Structural SVM

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \max_y (\Delta(y_i, y) + \mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$



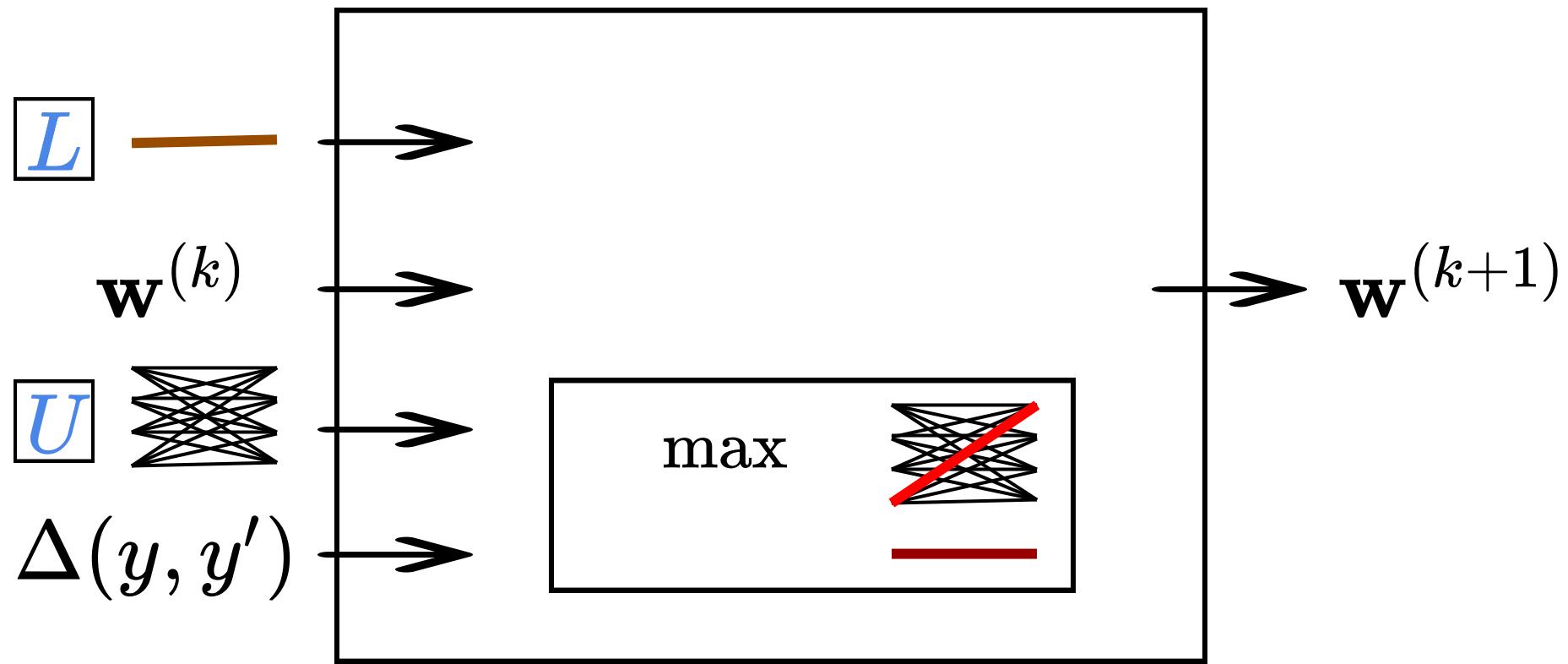
# "Soft"Max-Margin CRF

$$\min_w \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$
$$\min_w \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\Delta(y_i, y) + \mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$



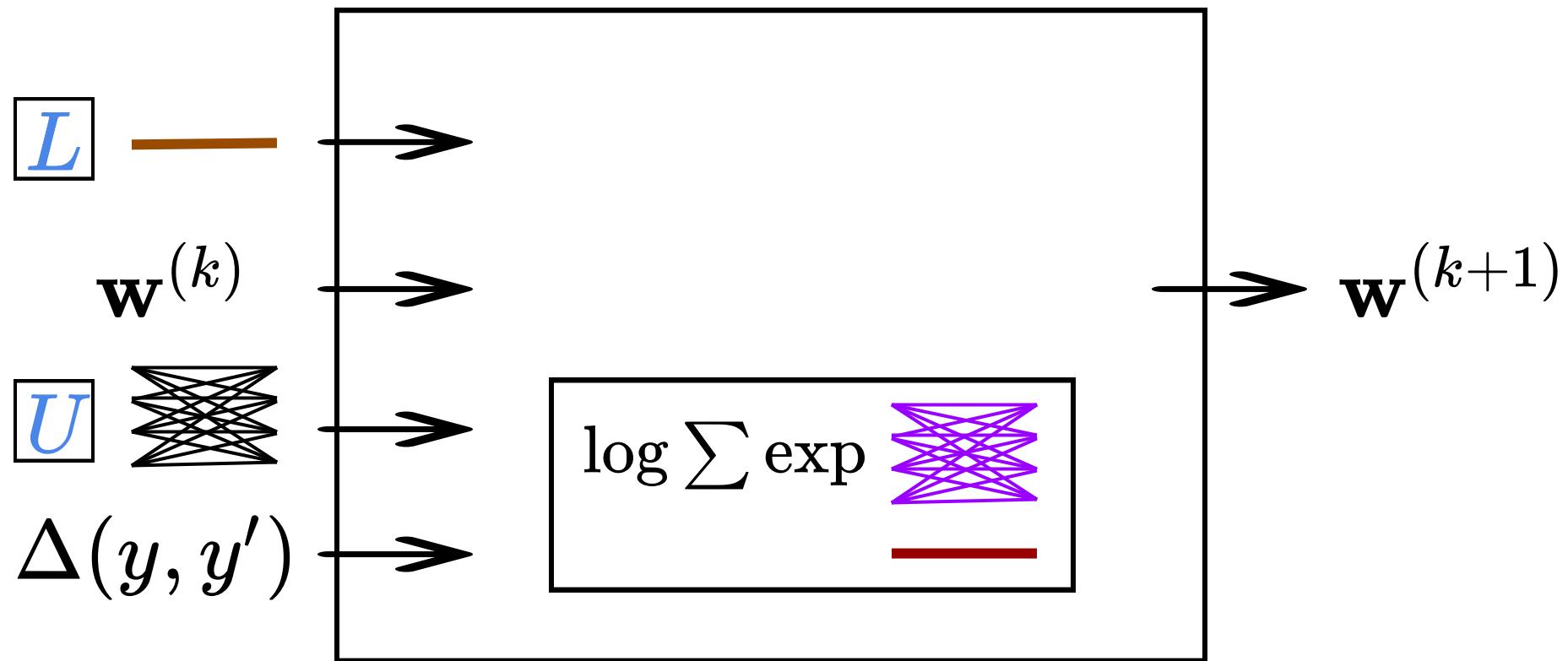
# Learning

## Structured Perceptron, Structural SVM



# Learning

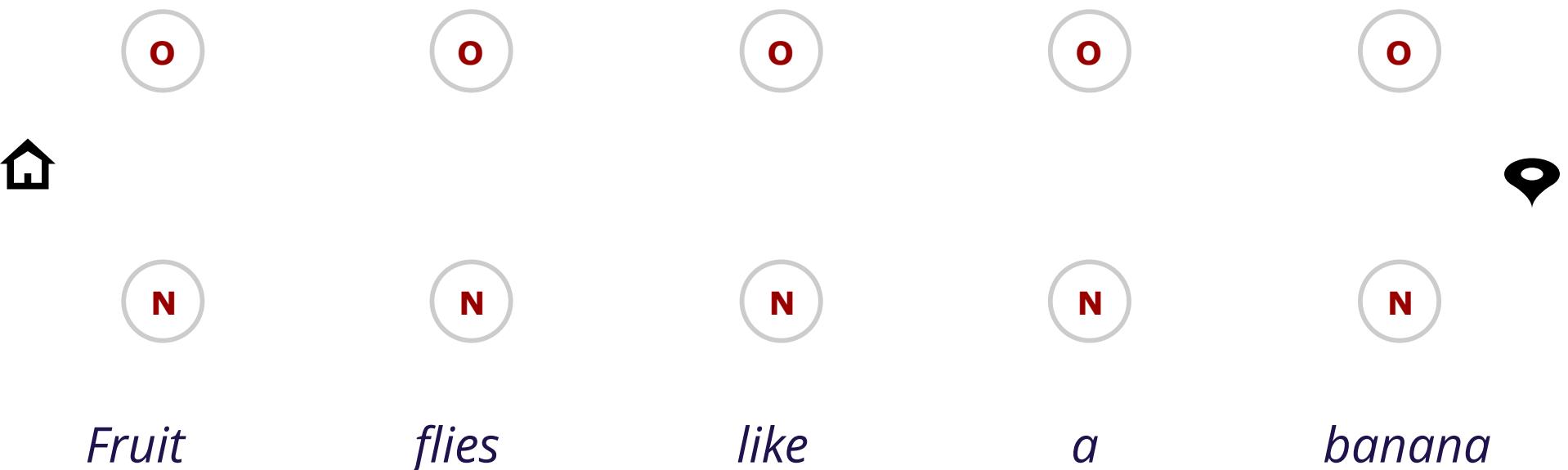
CRF, "Soft"Max-Margin CRF



# Semi-Markov CRF

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$

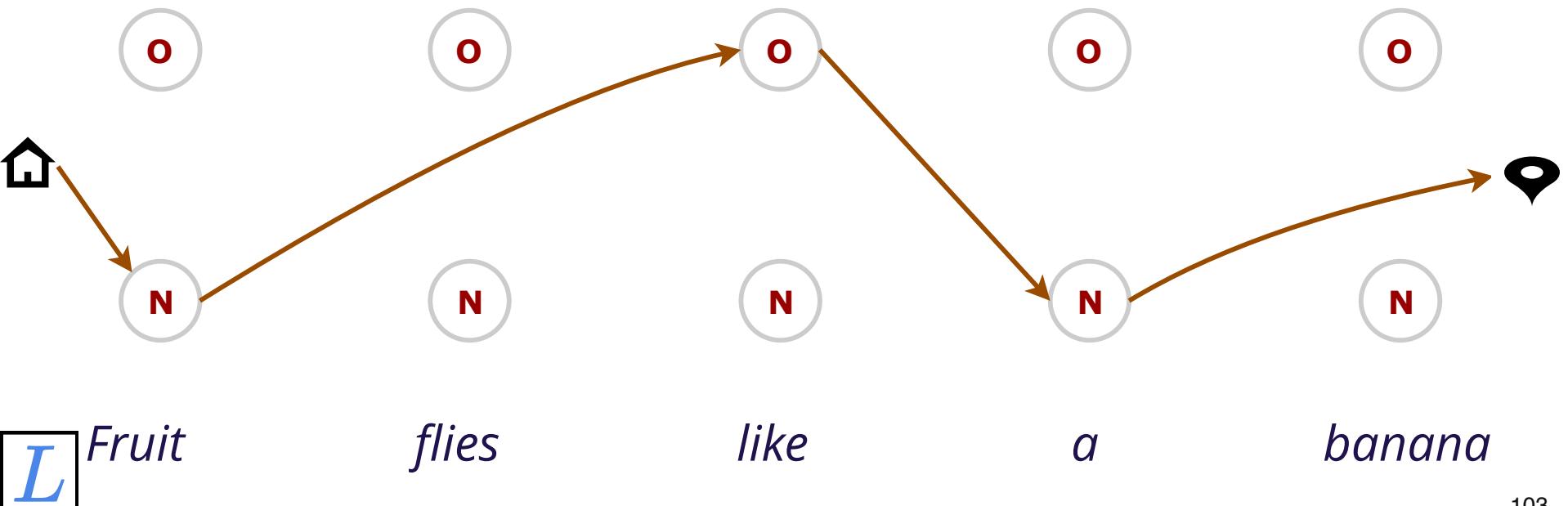
$$\mathbf{f}_1(x, [y^{0-1} = \mathbf{N}, y^2 = \mathbf{O}])$$



# Semi-Markov CRF

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$

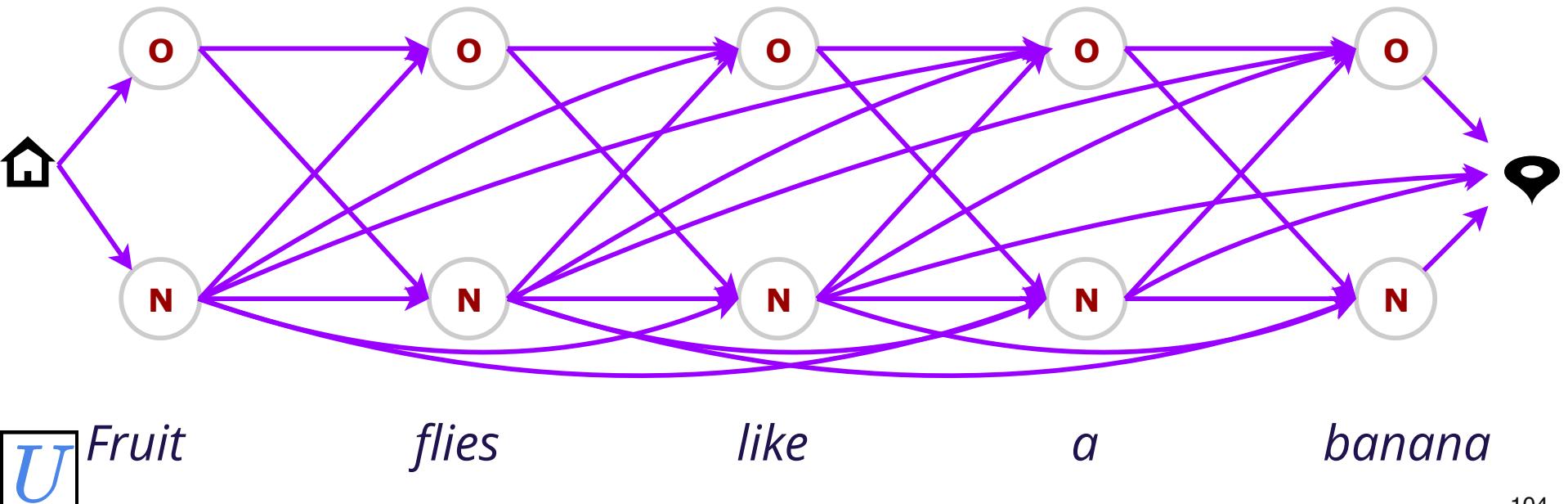
$$\mathbf{f}_1(x, [y^{0-1} = \mathbf{N}, y^2 = \mathbf{O}])$$



# Semi-Markov CRF

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$

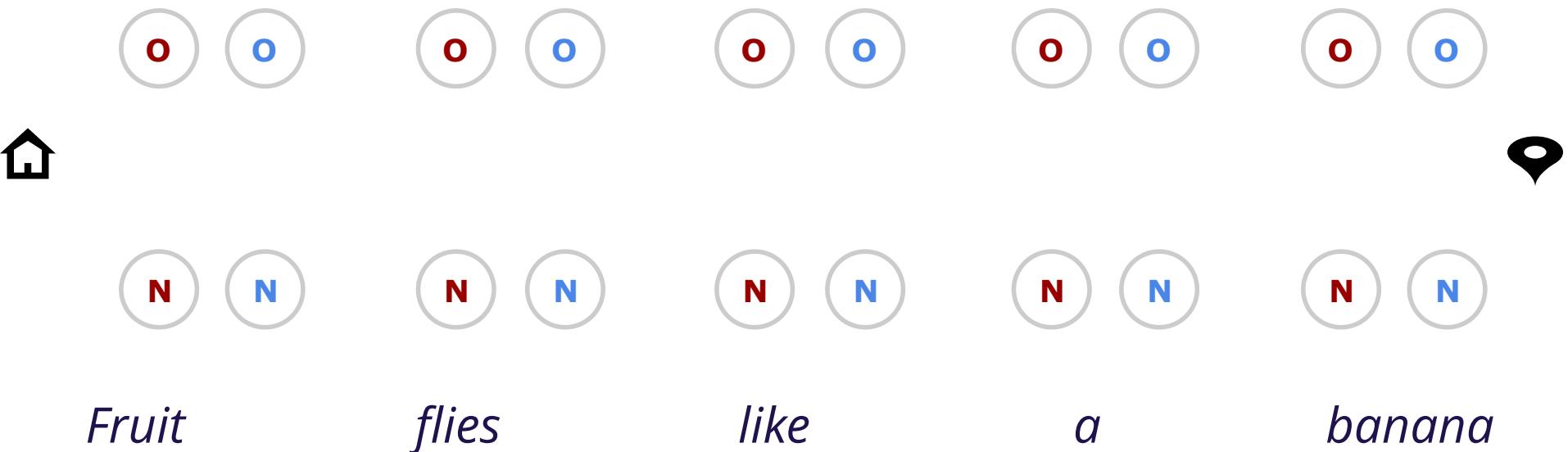
$$\mathbf{f}_1(x, [y^{0-1} = \mathbf{N}, y^2 = \mathbf{O}])$$



# Weak Semi-Markov CRF

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$

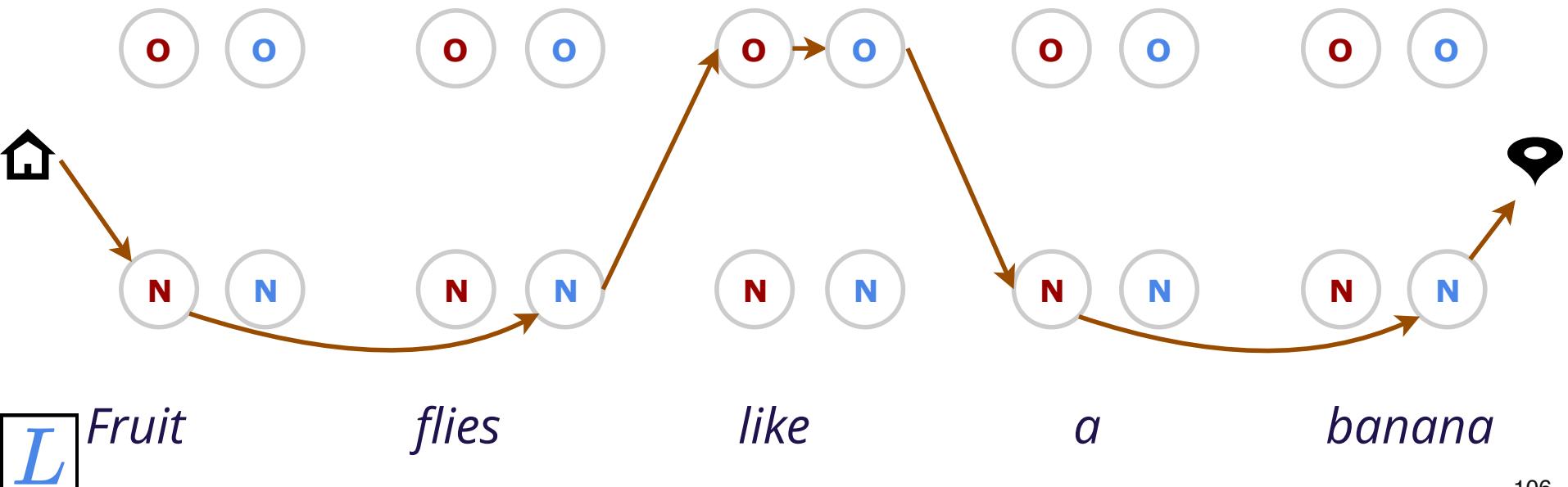
$$\mathbf{f}_1(x, [y^{0-1} = \mathbf{N}]) \quad \mathbf{f}_2(x, [y^1 = \mathbf{N}, y^2 = \mathbf{O}])$$



# Weak Semi-Markov CRF

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$

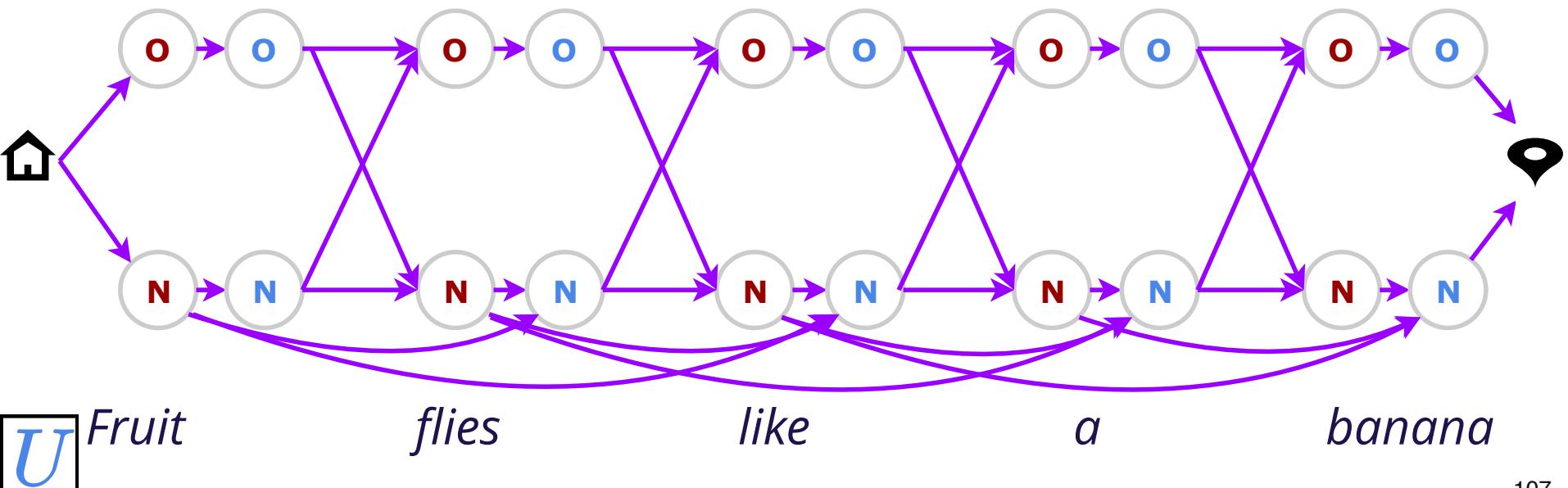
$$\mathbf{f}_1(x, [y^{0-1} = \mathbf{N}]) \quad \mathbf{f}_2(x, [y^1 = \mathbf{N}, y^2 = \mathbf{O}])$$



# Weak Semi-Markov CRF

$$\min_{\mathbf{w}} \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$

$$\mathbf{f}_1(x, [y^{0-1} = \mathbf{N}]) \quad \mathbf{f}_2(x, [y^1 = \mathbf{N}, y^2 = \mathbf{O}])$$



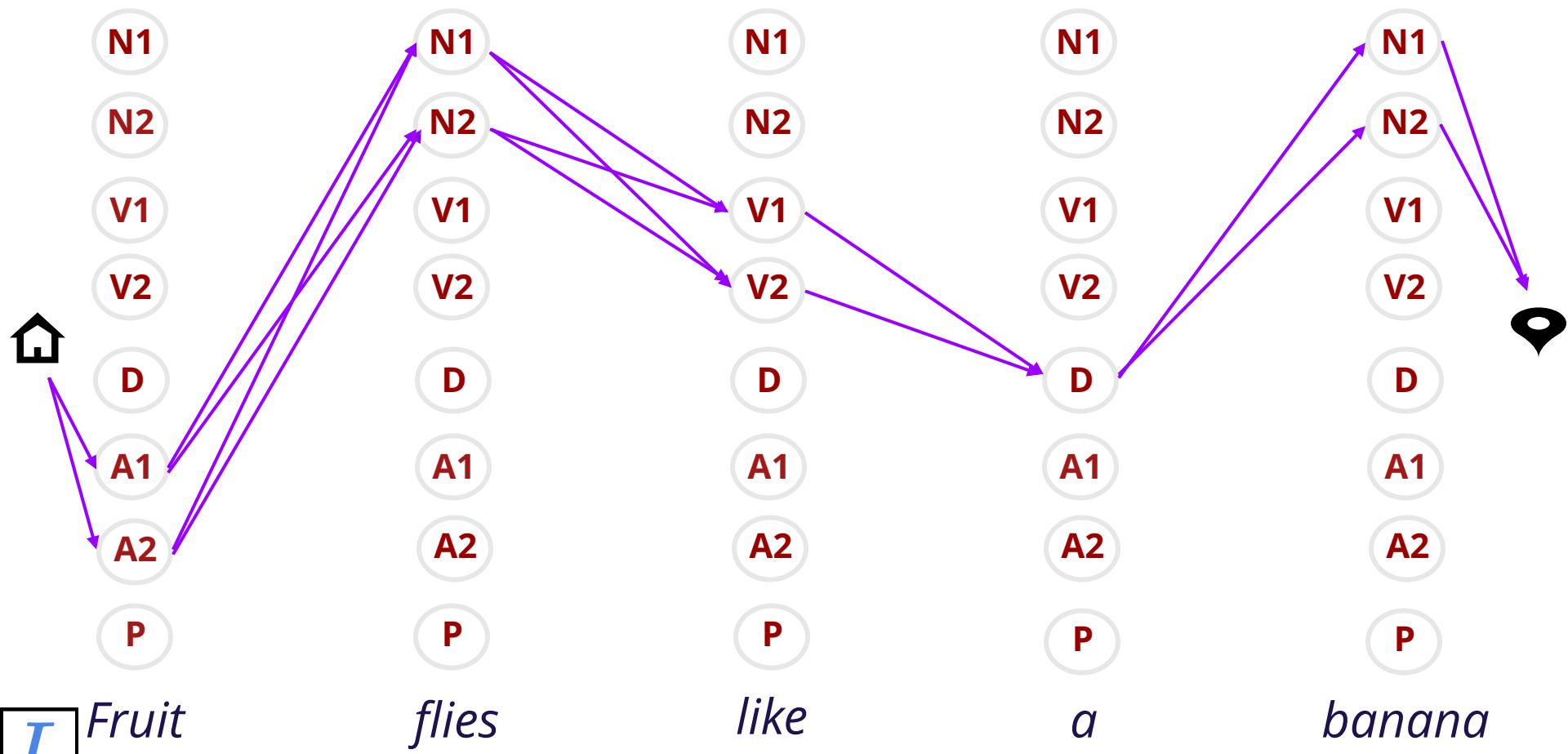
# Latent-Variable CRF

$$\min_{\mathbf{w}} \sum_i \left( -\log \sum_h \exp(\mathbf{w} \cdot \mathbf{f}(x_i, h, y_i)) + \log \sum_{h',y} \exp(\mathbf{w} \cdot \mathbf{f}(x_i, h', y)) \right)$$



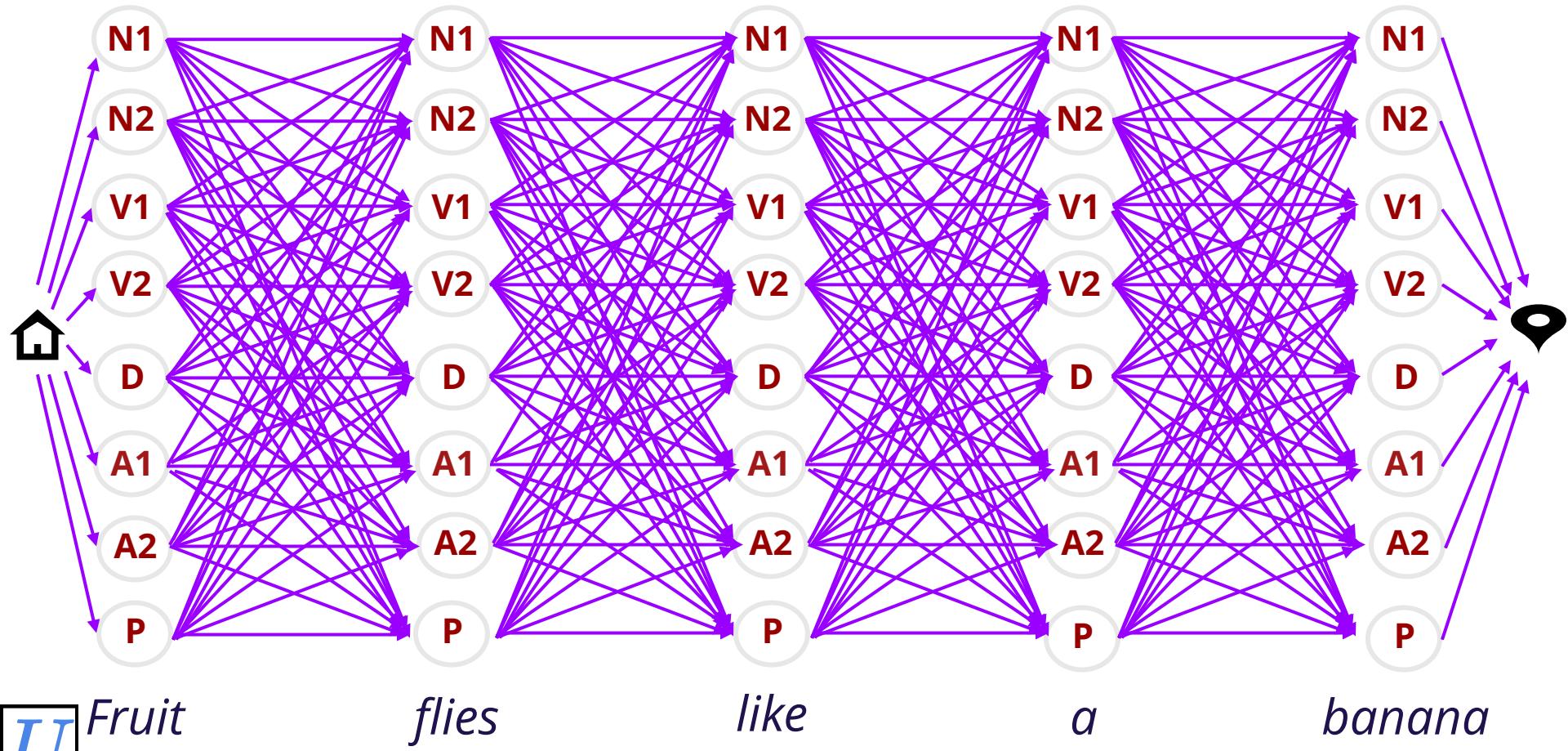
# Latent-Variable CRF

$$\min_{\mathbf{w}} \sum_i \left( -\log \sum_h \exp (\mathbf{w} \cdot \mathbf{f}(x_i, h, y_i)) + \log \sum_{h',y} \exp (\mathbf{w} \cdot \mathbf{f}(x_i, h', y)) \right)$$



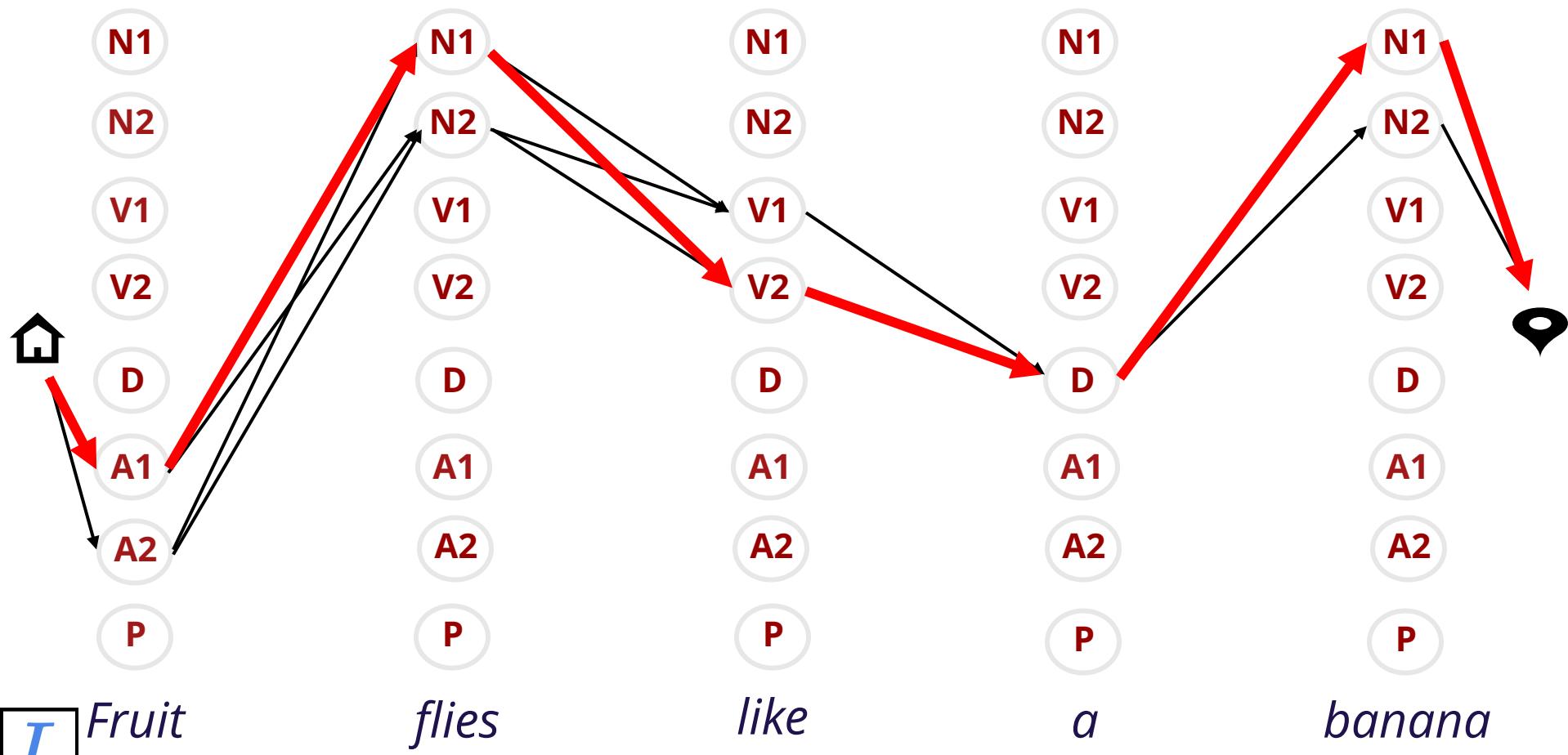
# Latent-Variable CRF

$$\min_{\mathbf{w}} \sum_i \left( -\log \sum_h \exp (\mathbf{w} \cdot \mathbf{f}(x_i, h, y_i)) + \log \sum_{h',y} \exp (\mathbf{w} \cdot \mathbf{f}(x_i, h', y)) \right)$$



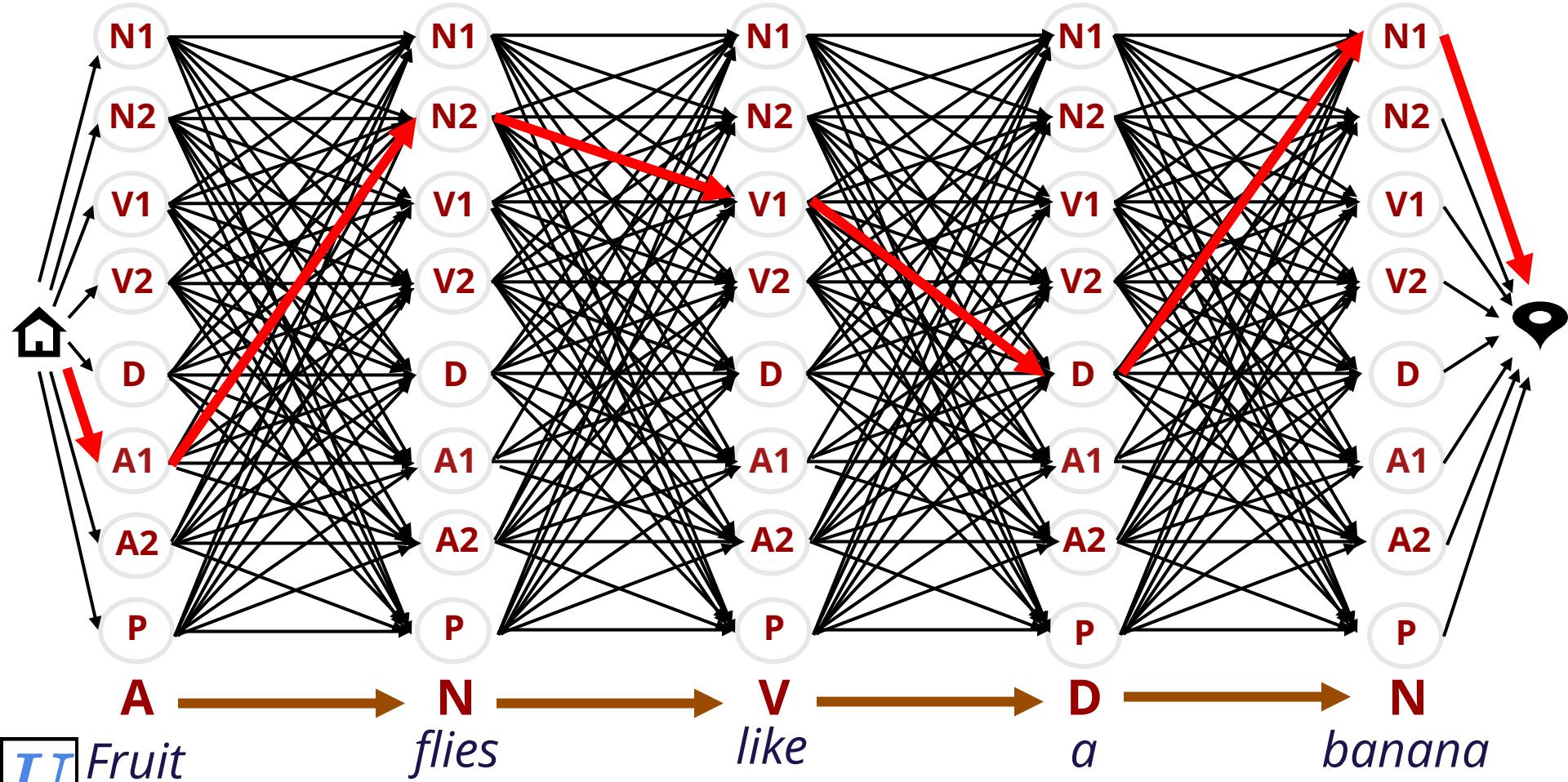
# Latent-Variable SSVM

$$\min_{\mathbf{w}} \sum_i \left( -\max_h (\mathbf{w} \cdot \mathbf{f}(x_i, h, y_i)) + \max_{h', y} (\Delta(y_i, y, h') + \mathbf{w} \cdot \mathbf{f}(x_i, h', y)) \right)$$



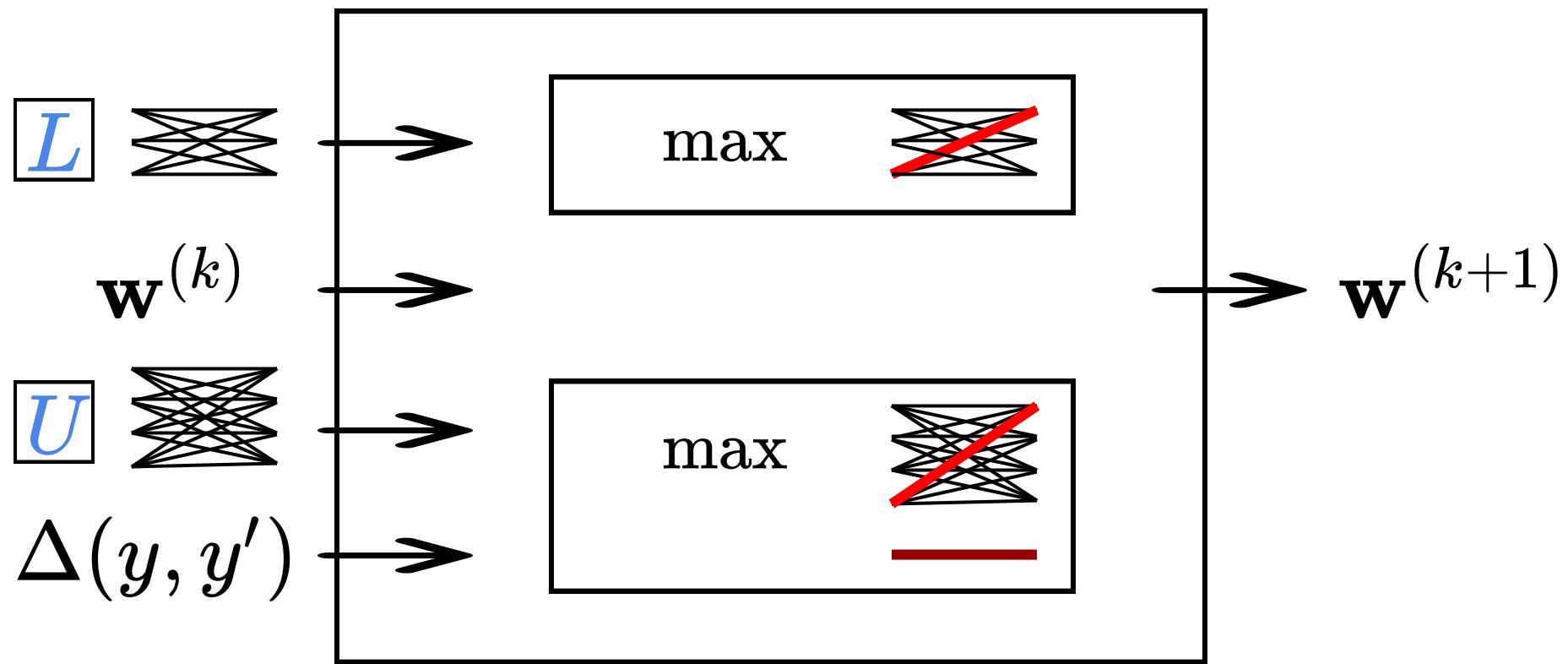
# Latent-Variable SSVM

$$\min_{\mathbf{w}} \sum_i \left( -\max_h (\mathbf{w} \cdot \mathbf{f}(x_i, h, y_i)) + \max_{h', y} (\Delta(y_i, y, h') + \mathbf{w} \cdot \mathbf{f}(x_i, h', y)) \right)$$



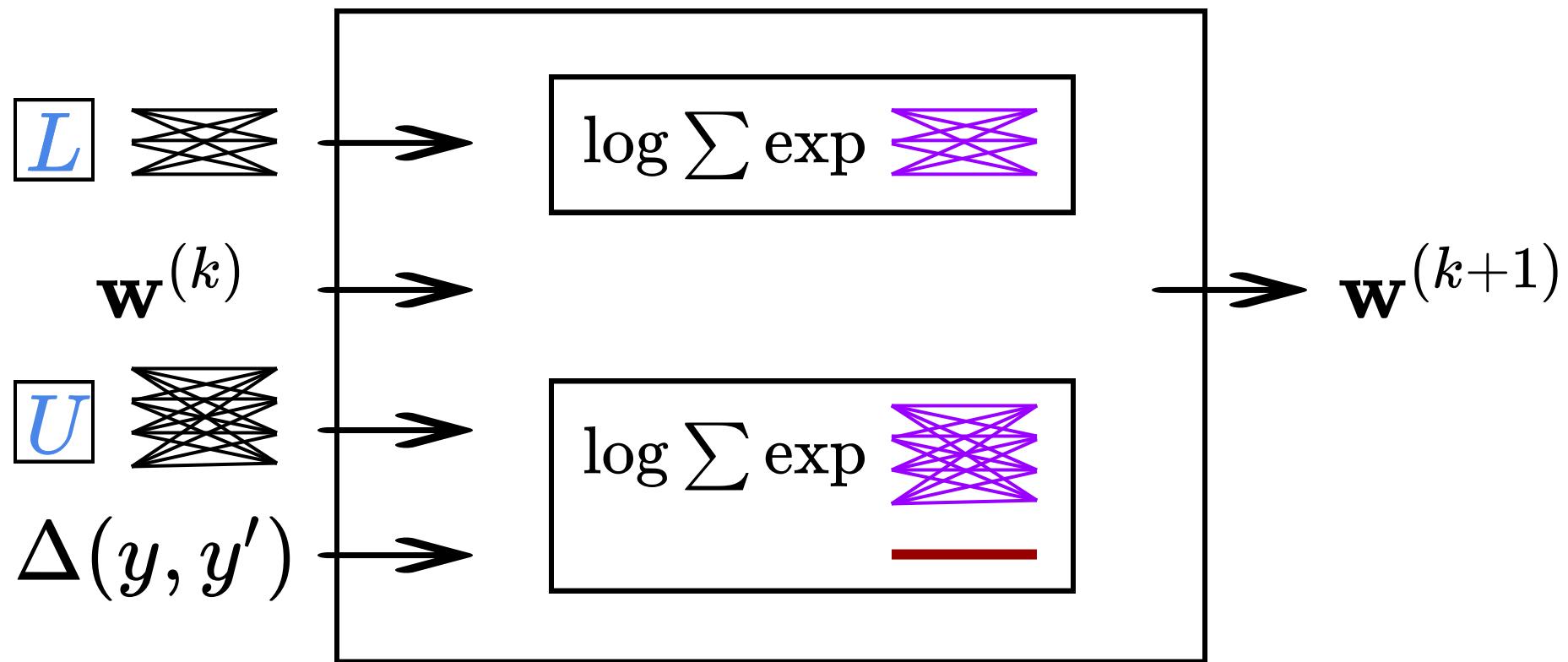
# Learning

Structured Perceptron, SSVM, Latent SSVM, ...

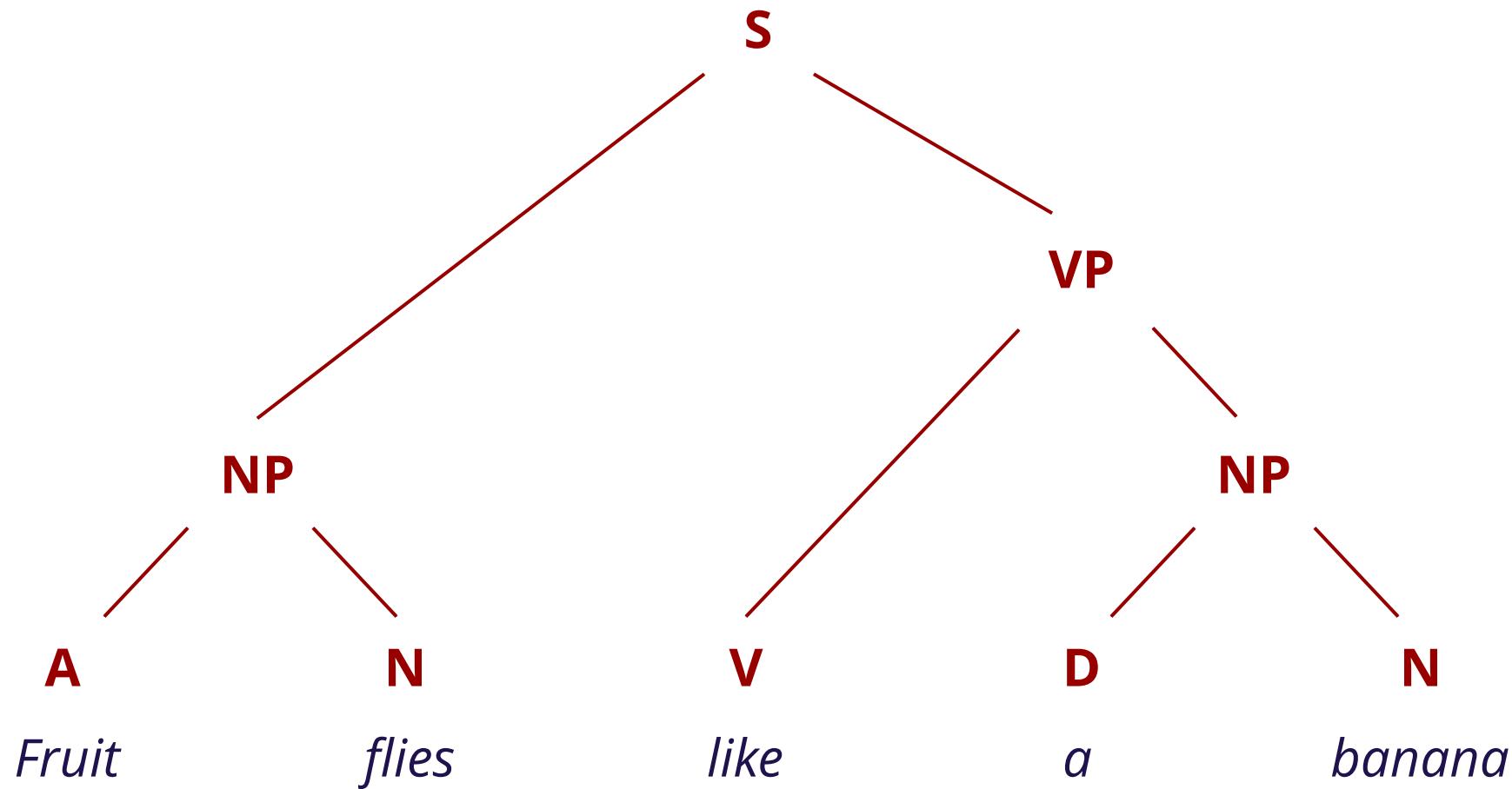


# Learning

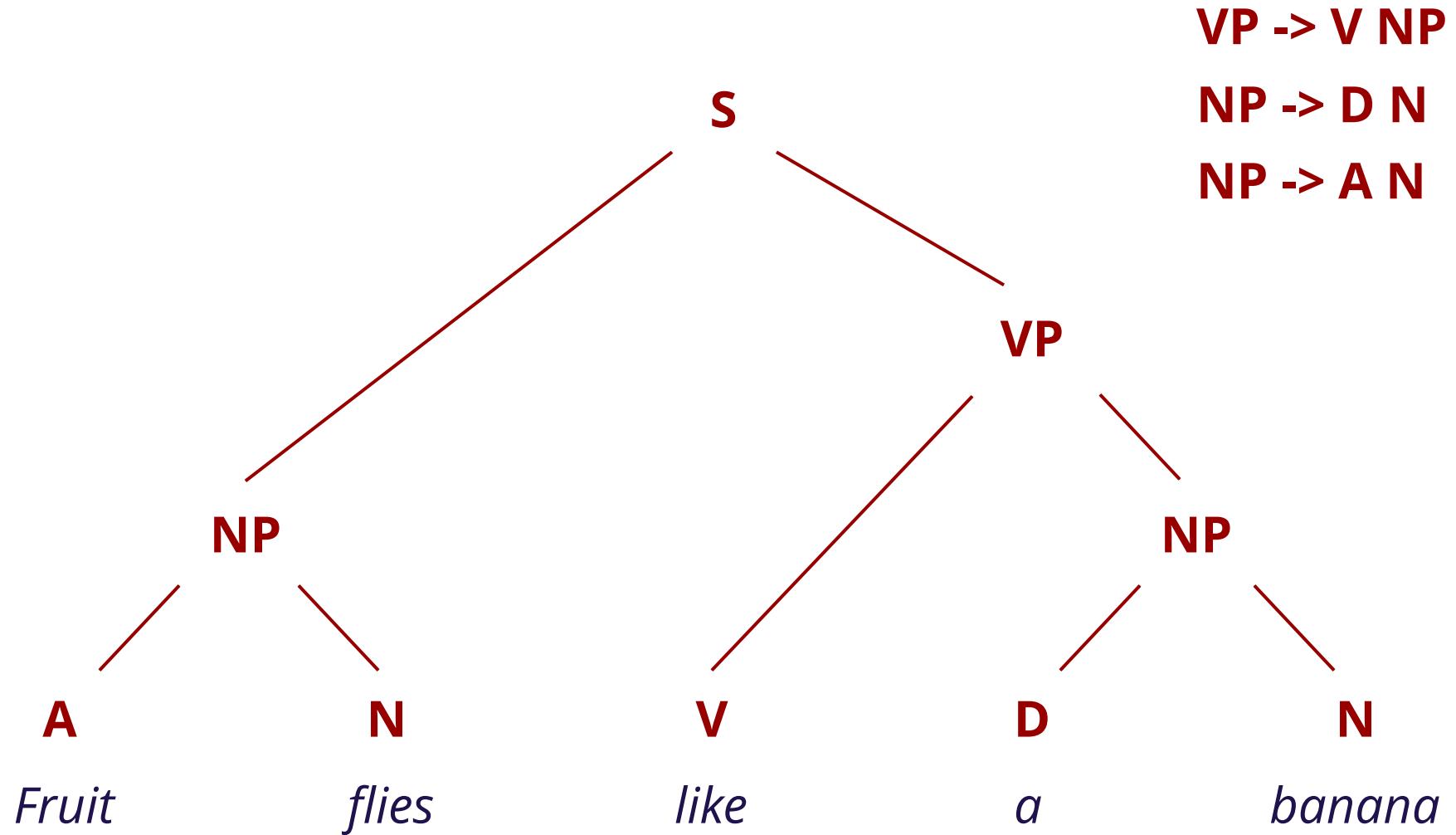
Linear/Semi/Latent/Softmax-margin CRF, ...



# Constituency Parsing

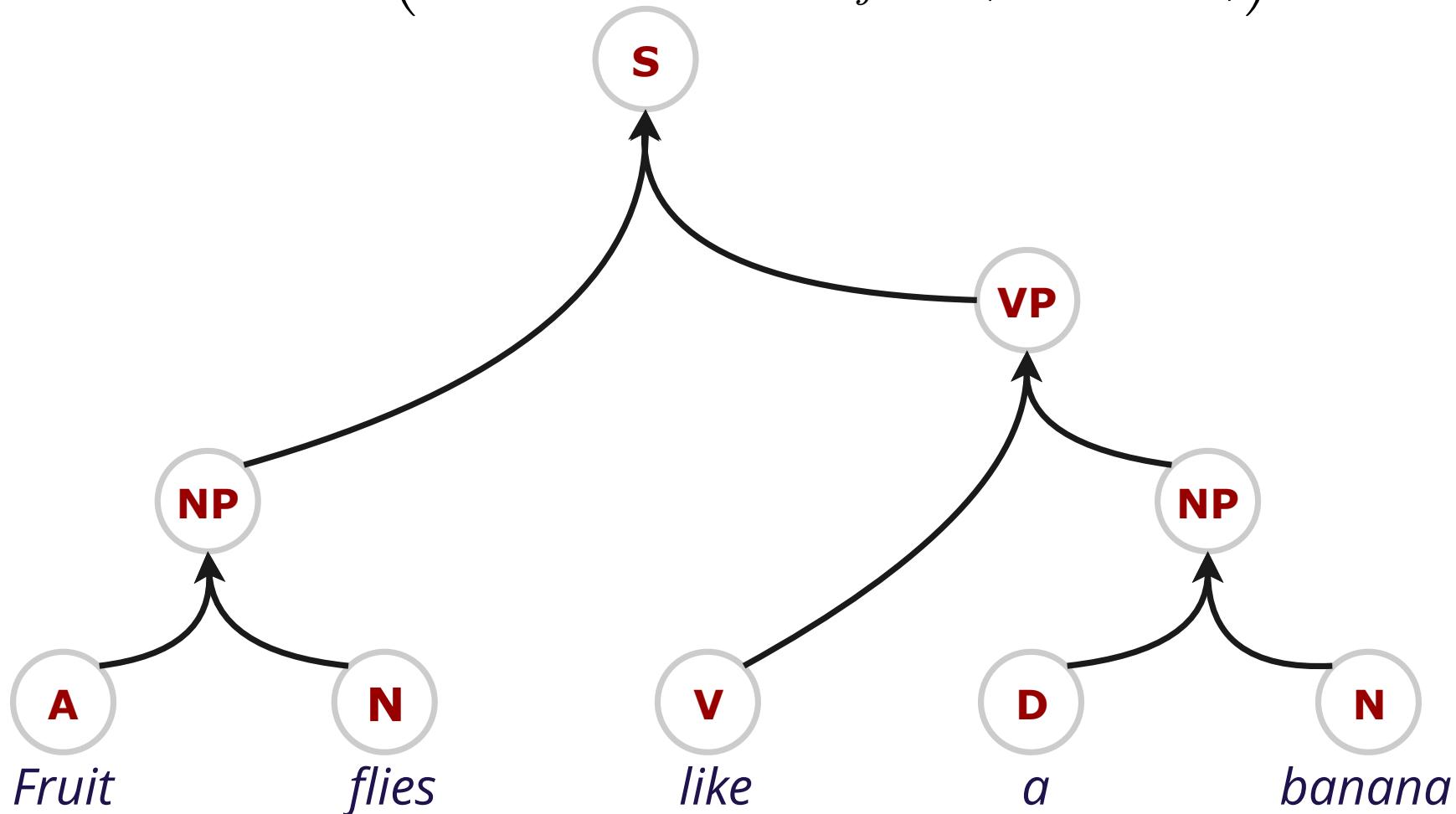


# Constituency Parsing



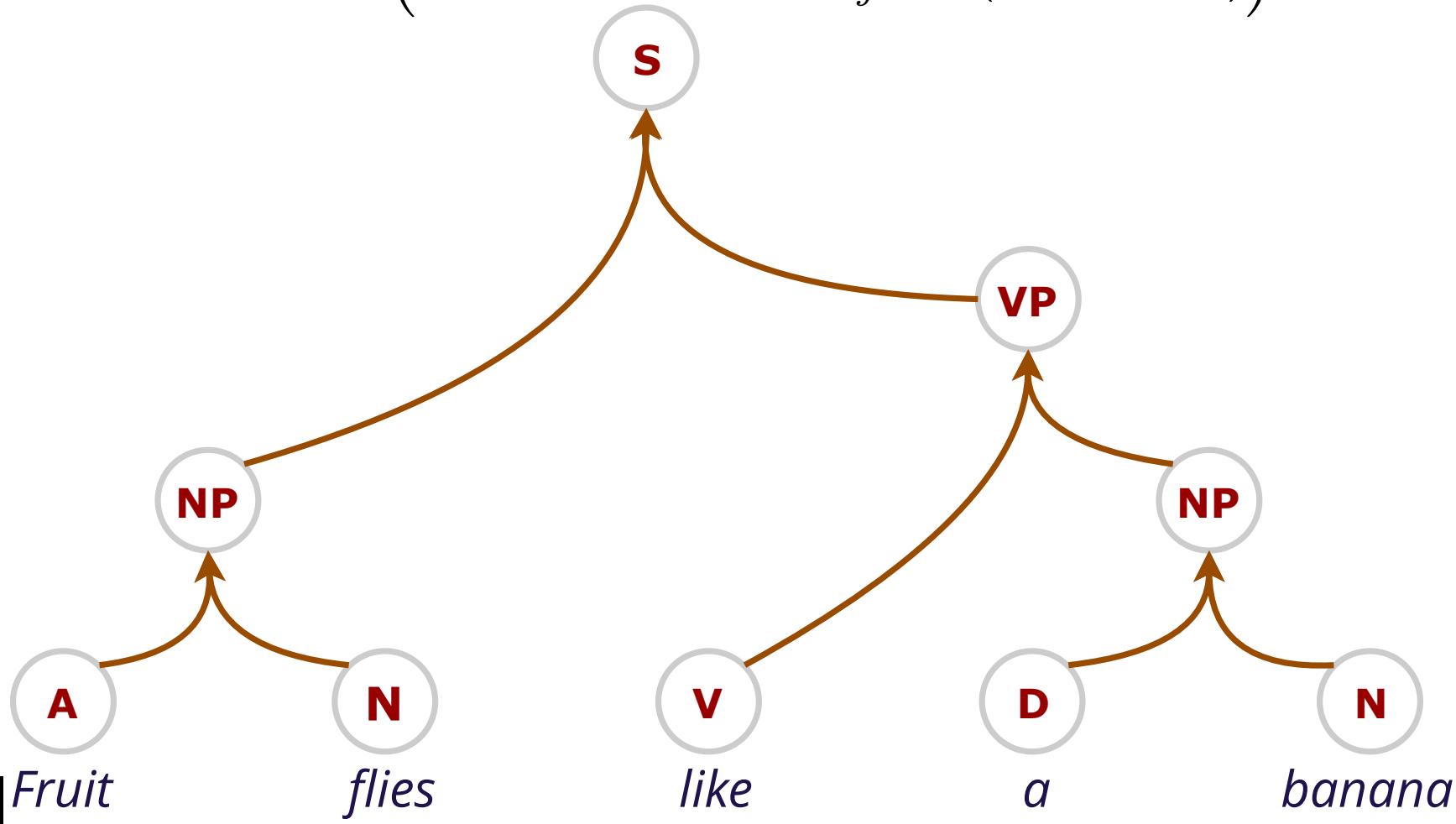
# Parsing with CRF

$$\min_w \sum_i \left( -w \cdot f(x_i, y_i) + \log \sum_y \exp(w \cdot f(x_i, y)) \right)$$



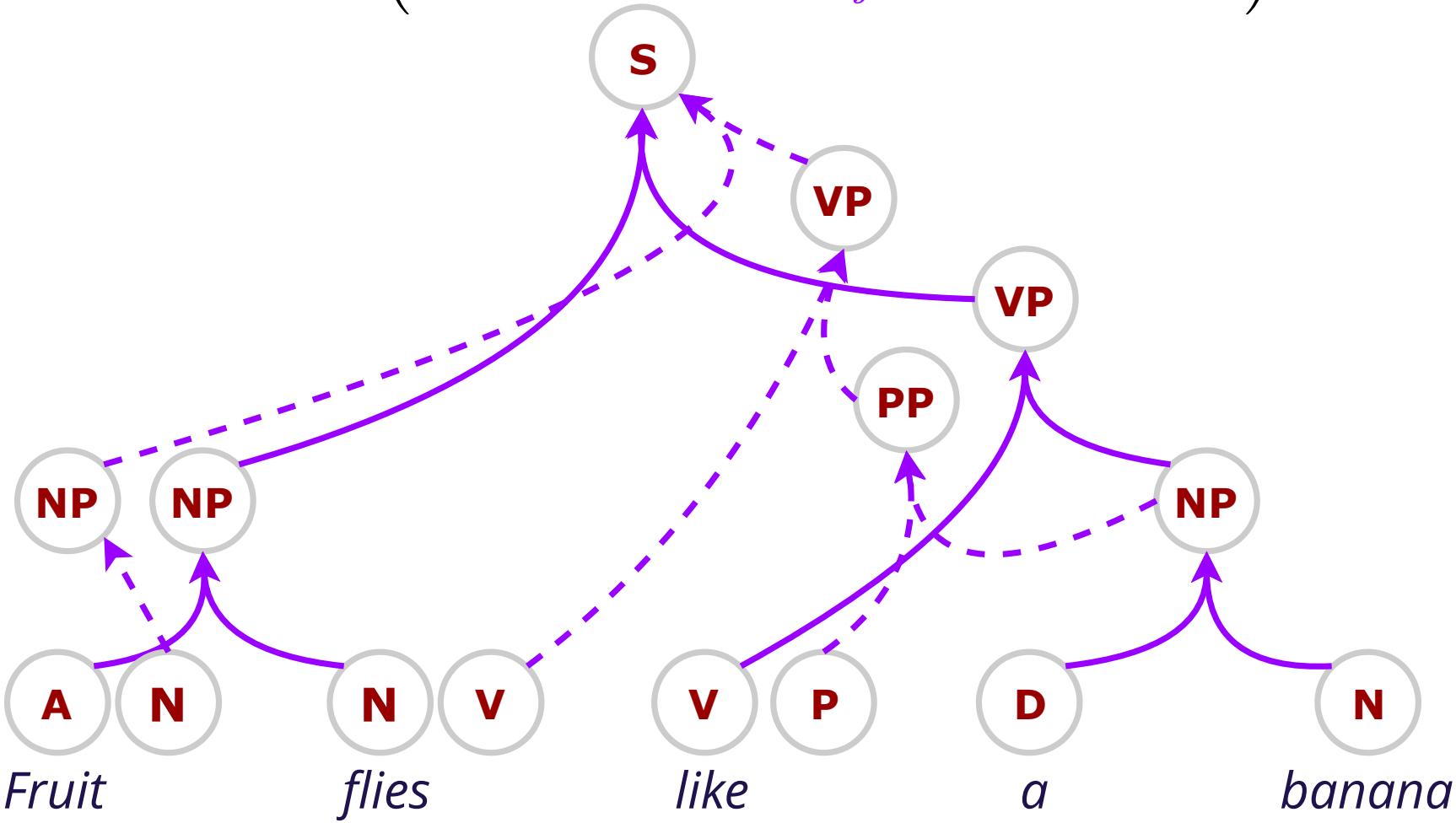
# Parsing with CRF

$$\min_w \sum_i \left( -\mathbf{w} \cdot \mathbf{f}(x_i, y_i) + \log \sum_y \exp (\mathbf{w} \cdot \mathbf{f}(x_i, y)) \right)$$



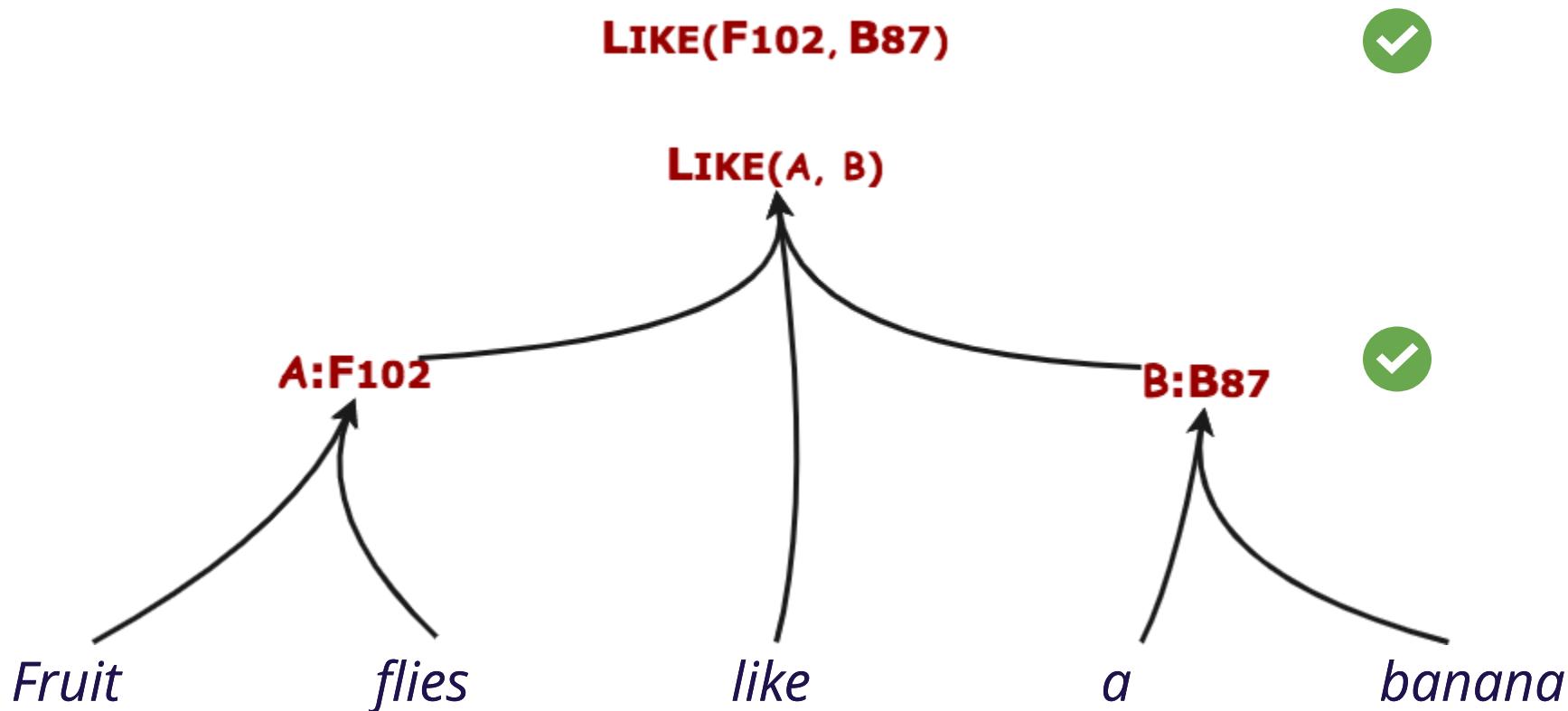
# Parsing with CRF

$$\min_w \sum_i \left( -w \cdot f(x_i, y_i) + \log \sum_y \exp(w \cdot f(x_i, y)) \right)$$



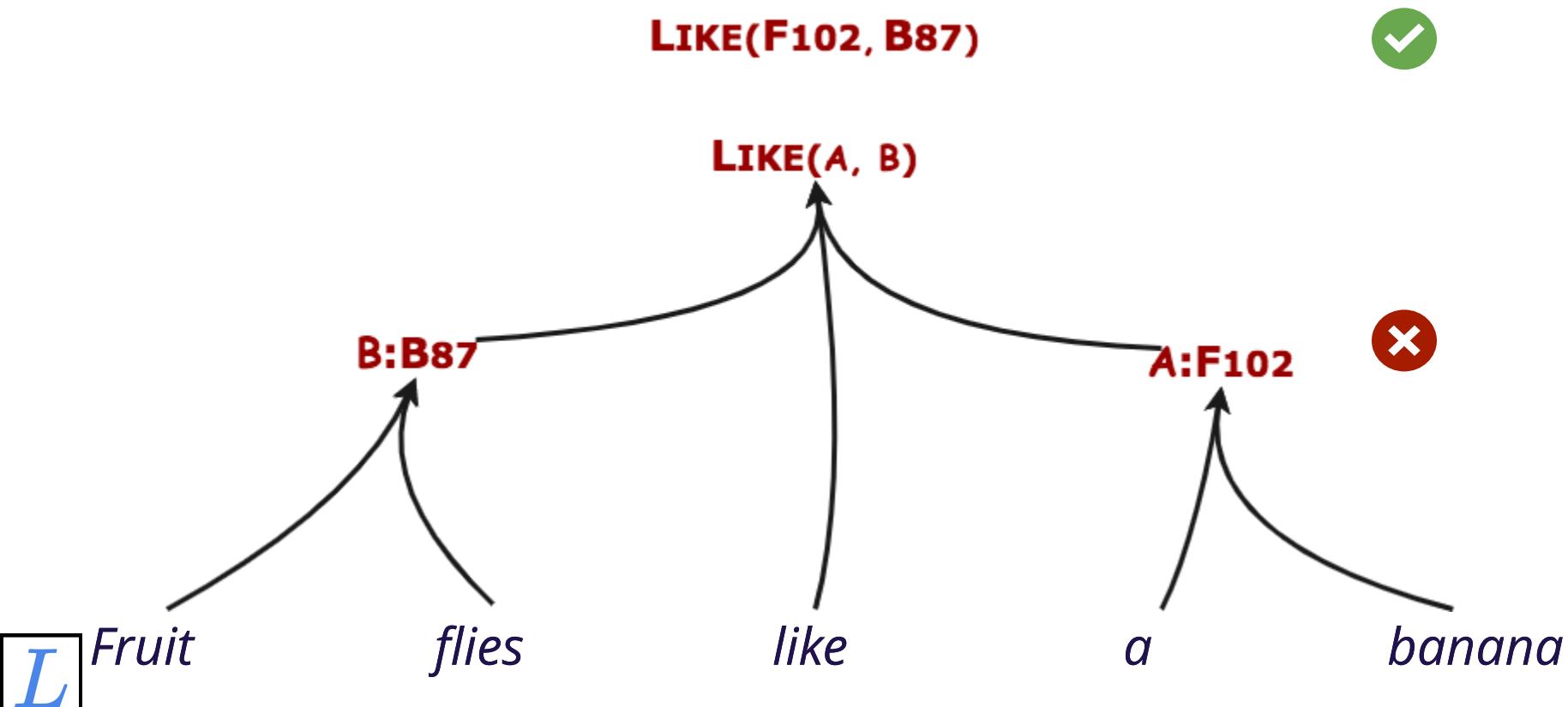
# Hybrid Tree

$$\min_w \sum_i \left( -\log \sum_h \exp(w \cdot f(x_i, h, y_i)) + \log \sum_{h',y} \exp(w \cdot f(x_i, h', y)) \right)$$



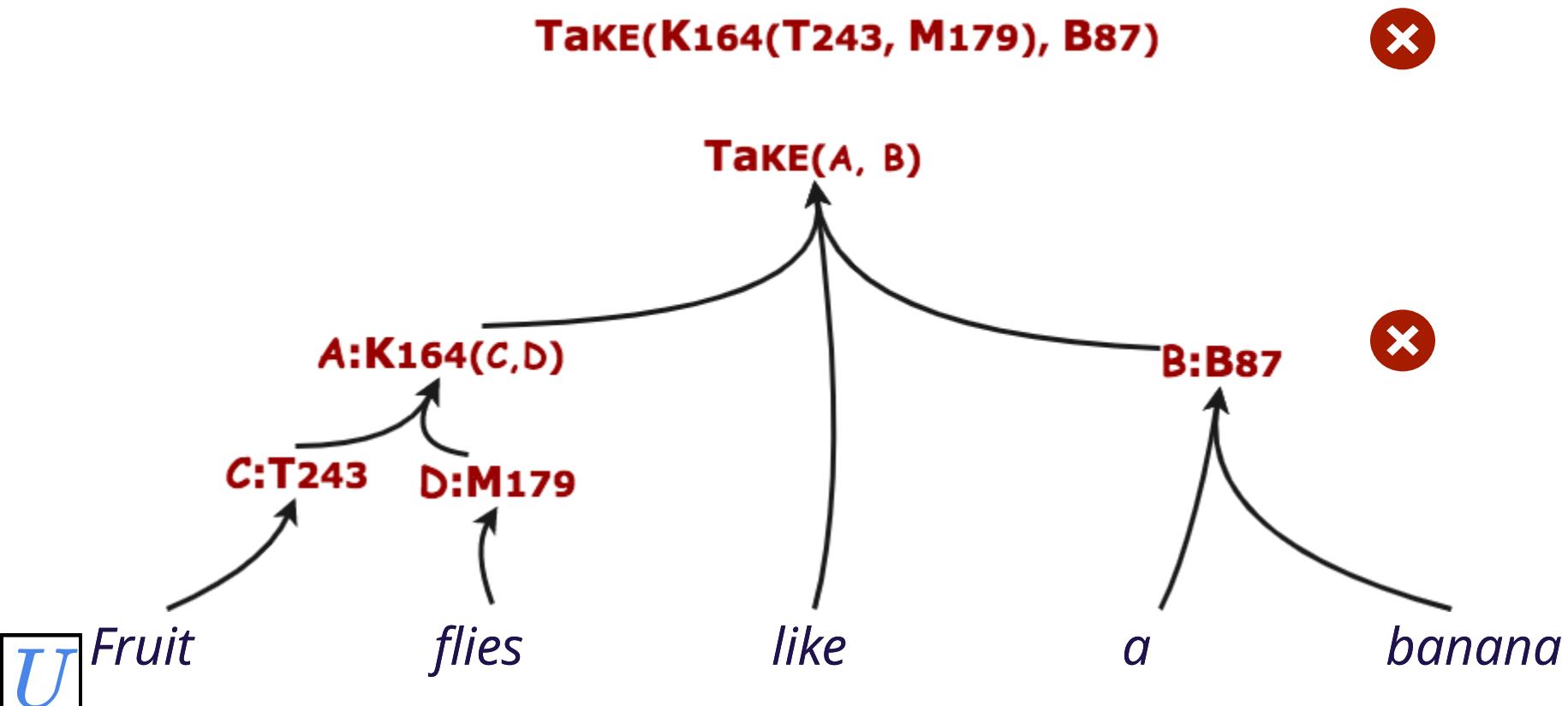
# Hybrid Tree

$$\min_w \sum_i \left( -\log \sum_h \exp(w \cdot f(x_i, h, y_i)) + \log \sum_{h',y} \exp(w \cdot f(x_i, h', y)) \right)$$

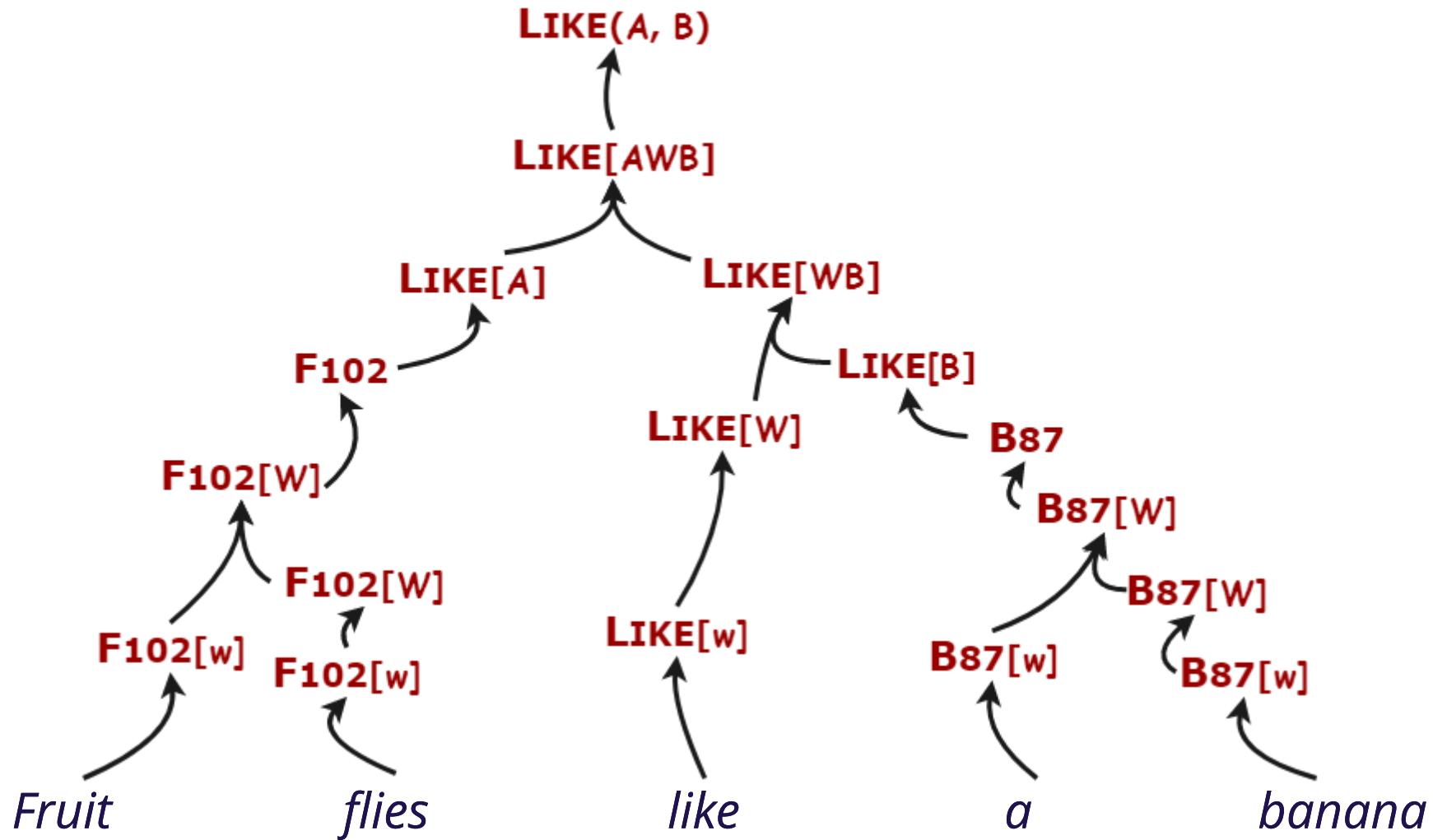


# Hybrid Tree

$$\min_w \sum_i \left( -\log \sum_h \exp(w \cdot f(x_i, h, y_i)) + \log \sum_{h',y} \exp(w \cdot f(x_i, h', y)) \right)$$



# Hybrid Tree



# Overlapping Structures

NX

NX

NX

NX

*Fruit*

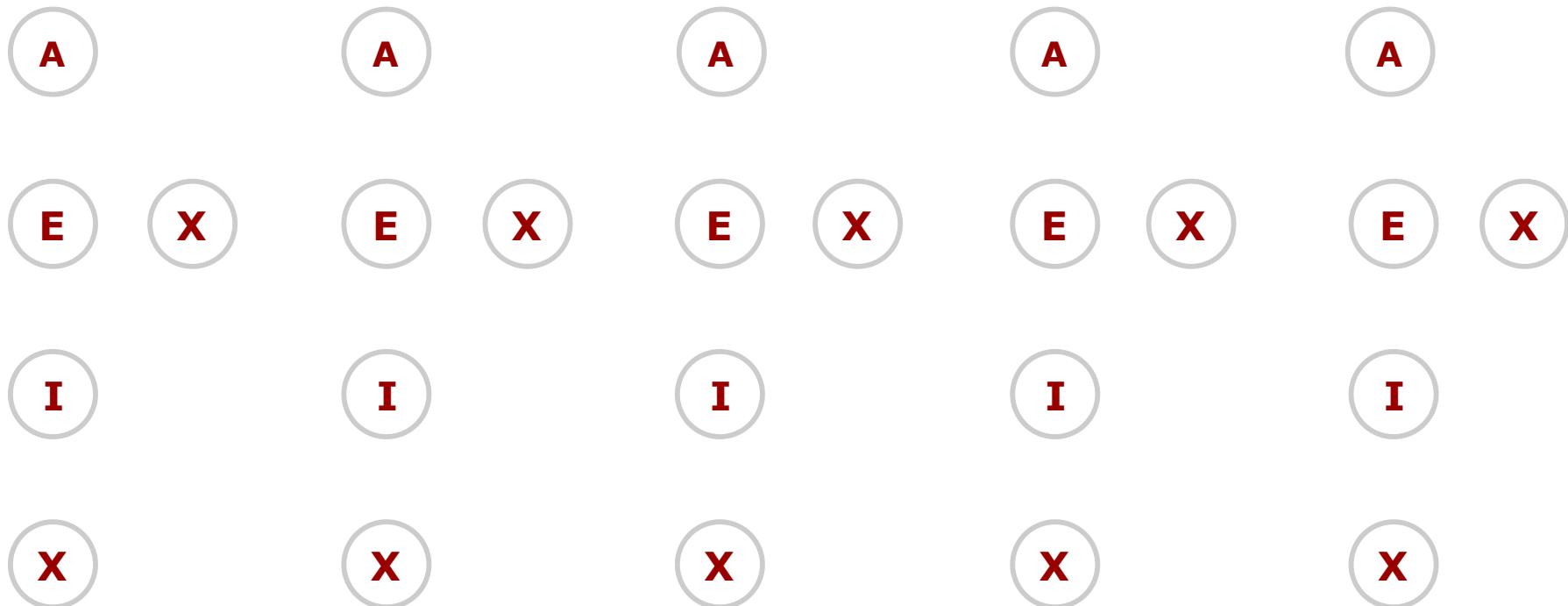
*flies*

*like*

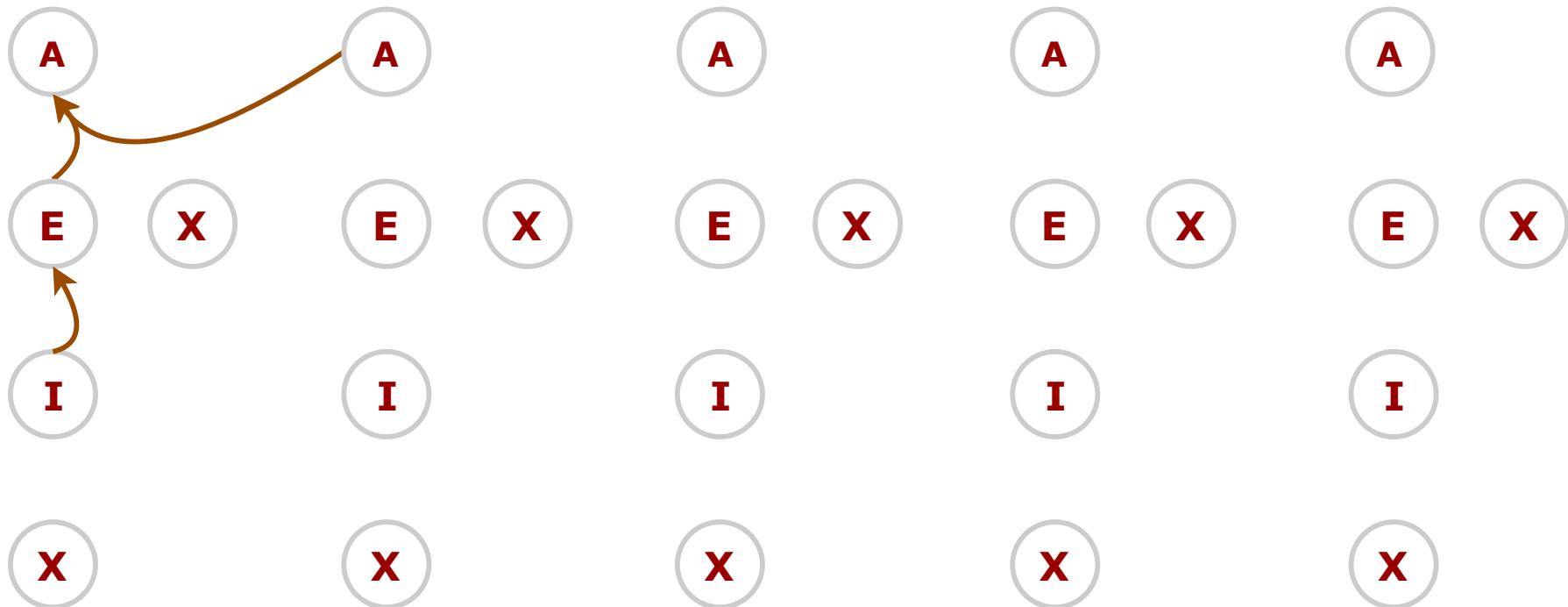
*a*

*banana*

# Mention Hypergraphs



# Mention Hypergraphs



L

Fruit

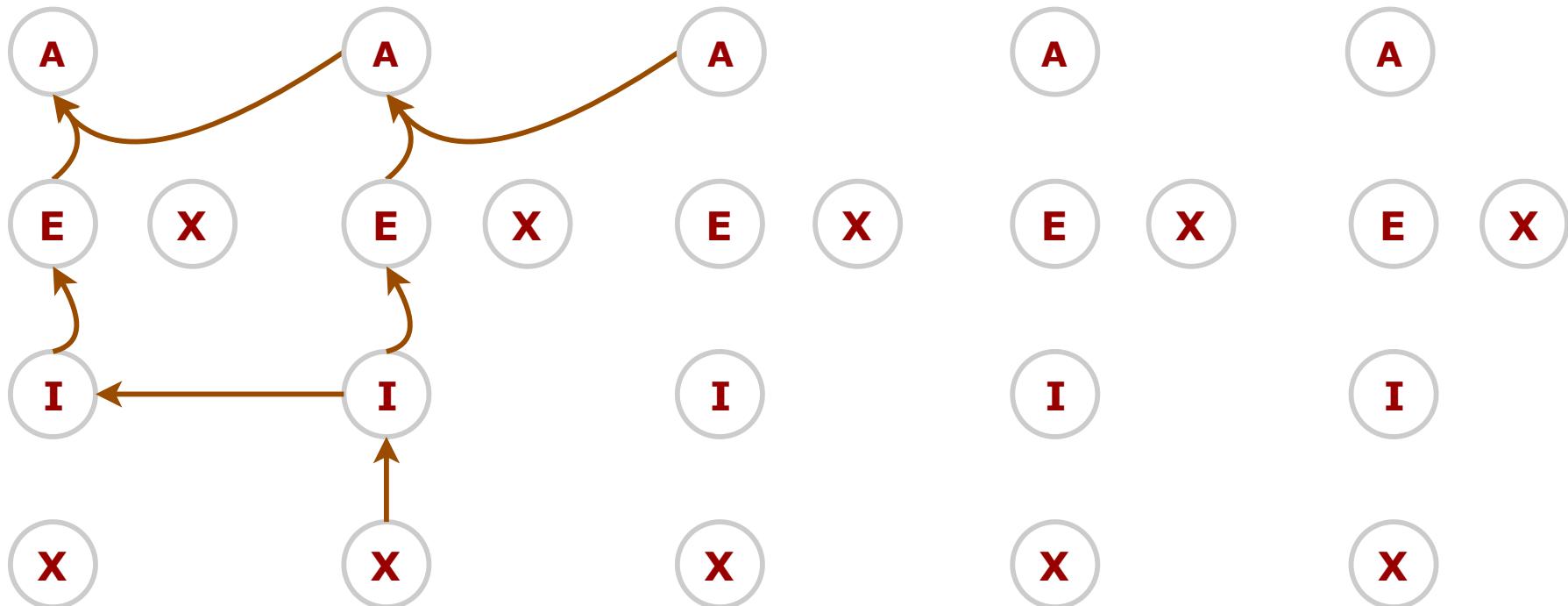
flies

like

a

banana

# Mention Hypergraphs



L

Fruit

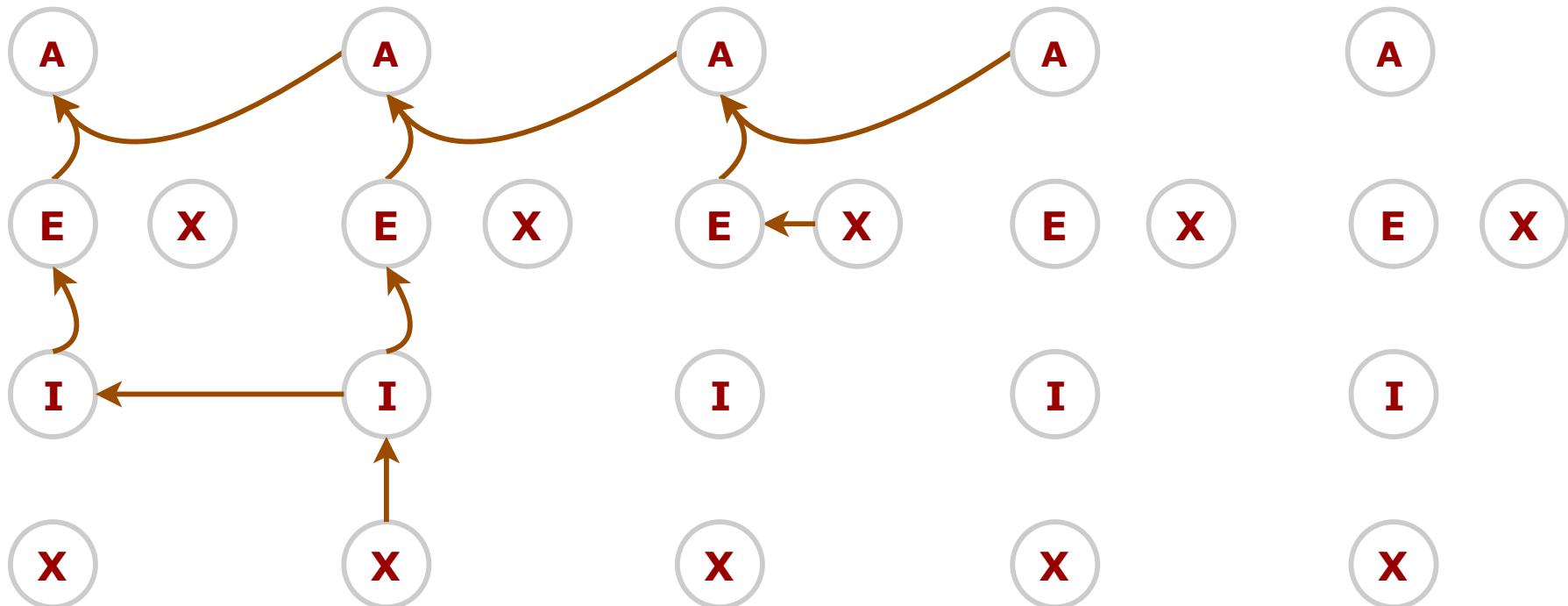
flies

like

a

banana

# Mention Hypergraphs



L

Fruit

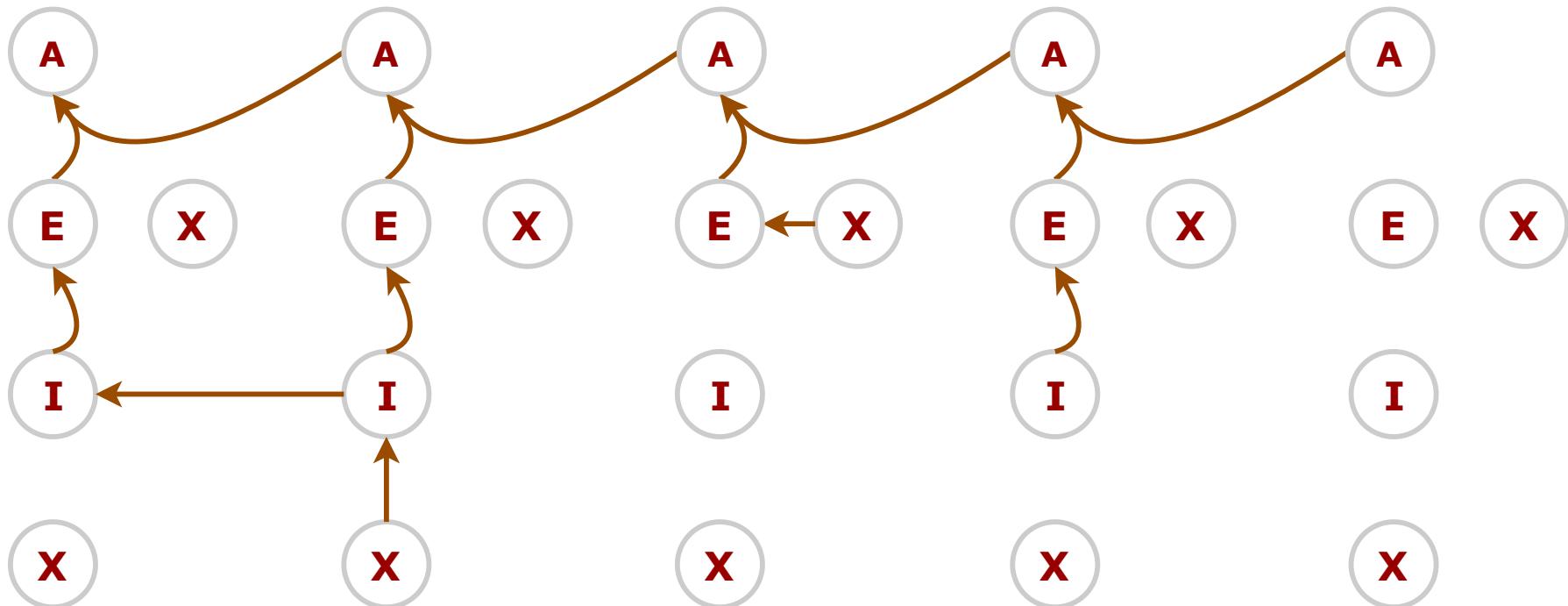
*flies*

*like*

*a*

*banana*

# Mention Hypergraphs



L

Fruit

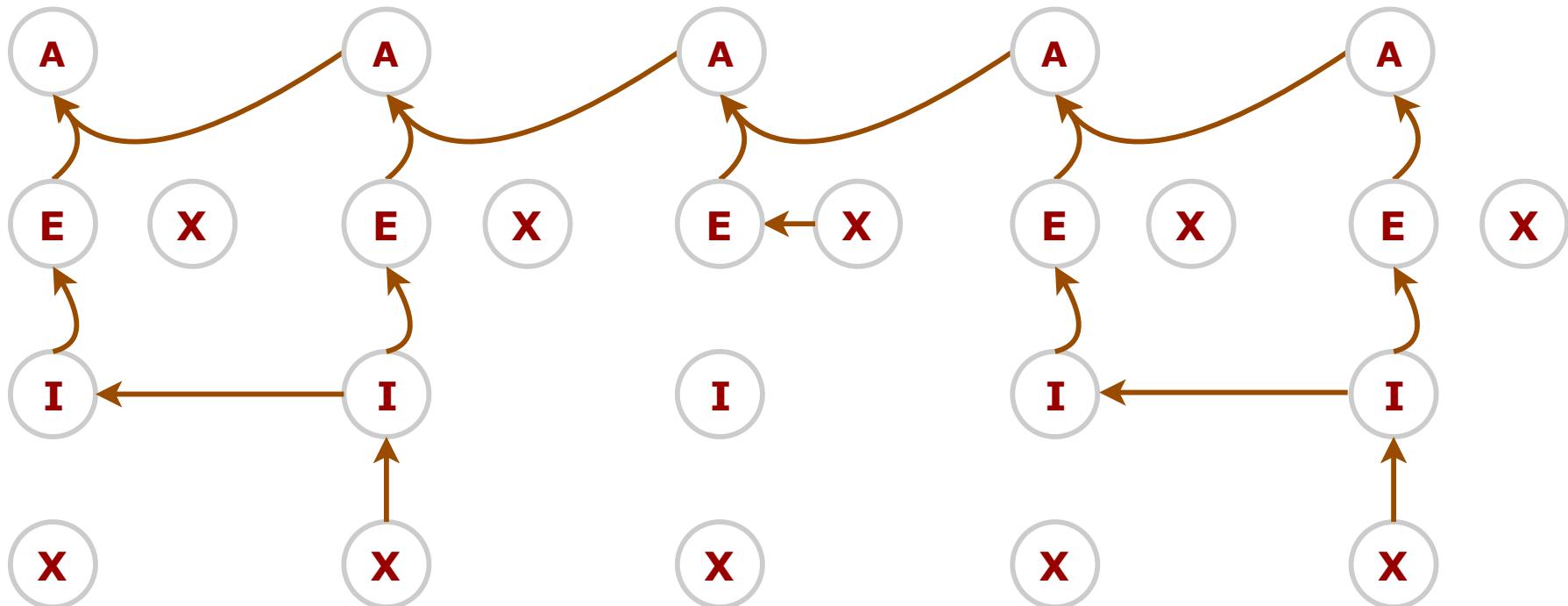
*flies*

*like*

*a*

*banana*

# Mention Hypergraphs



L

Fruit

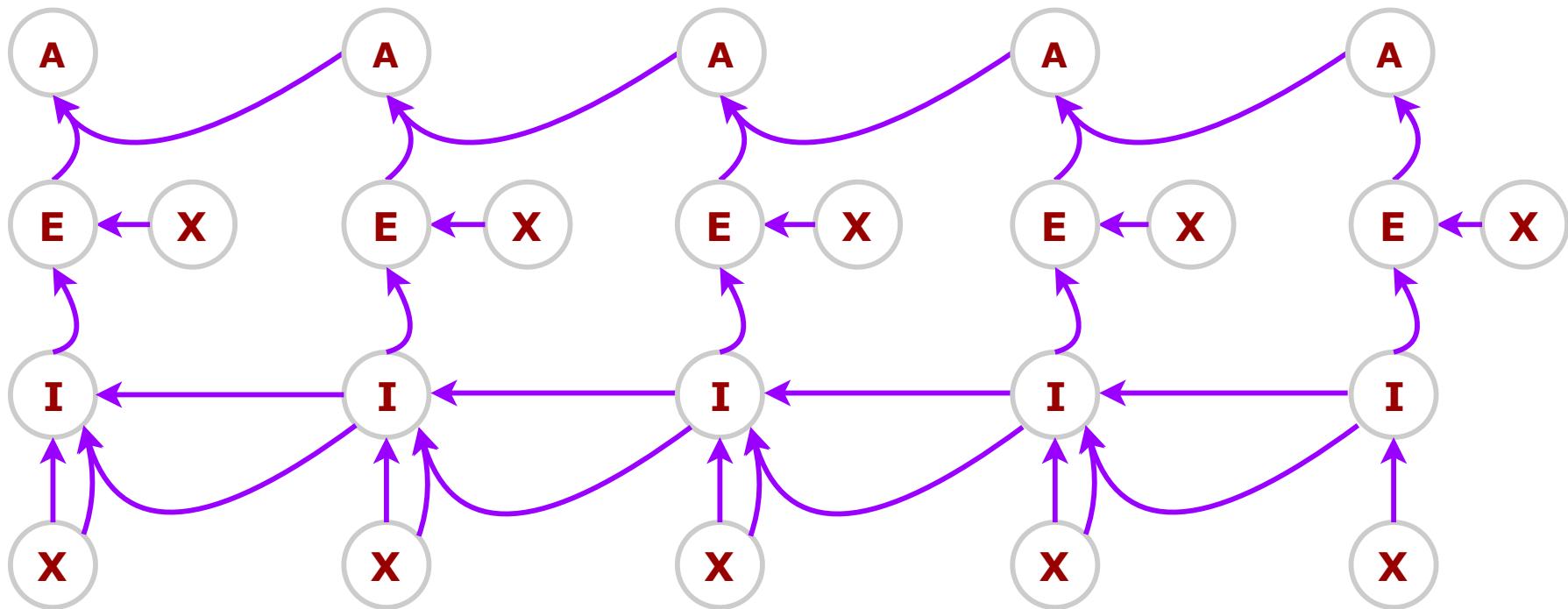
*flies*

*like*

*a*

*banana*

# Mention Hypergraphs



U

Fruit

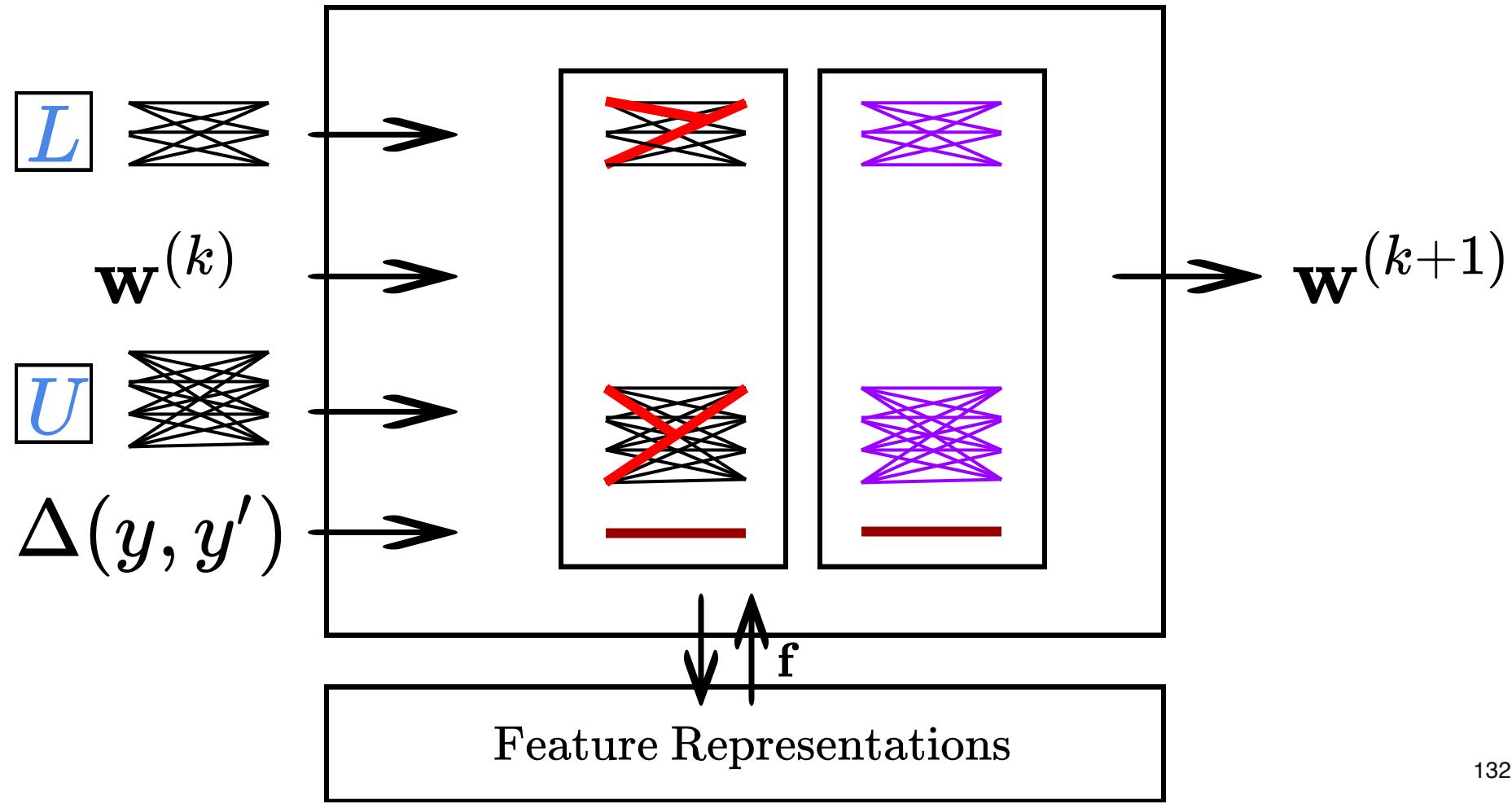
*flies*

*like*

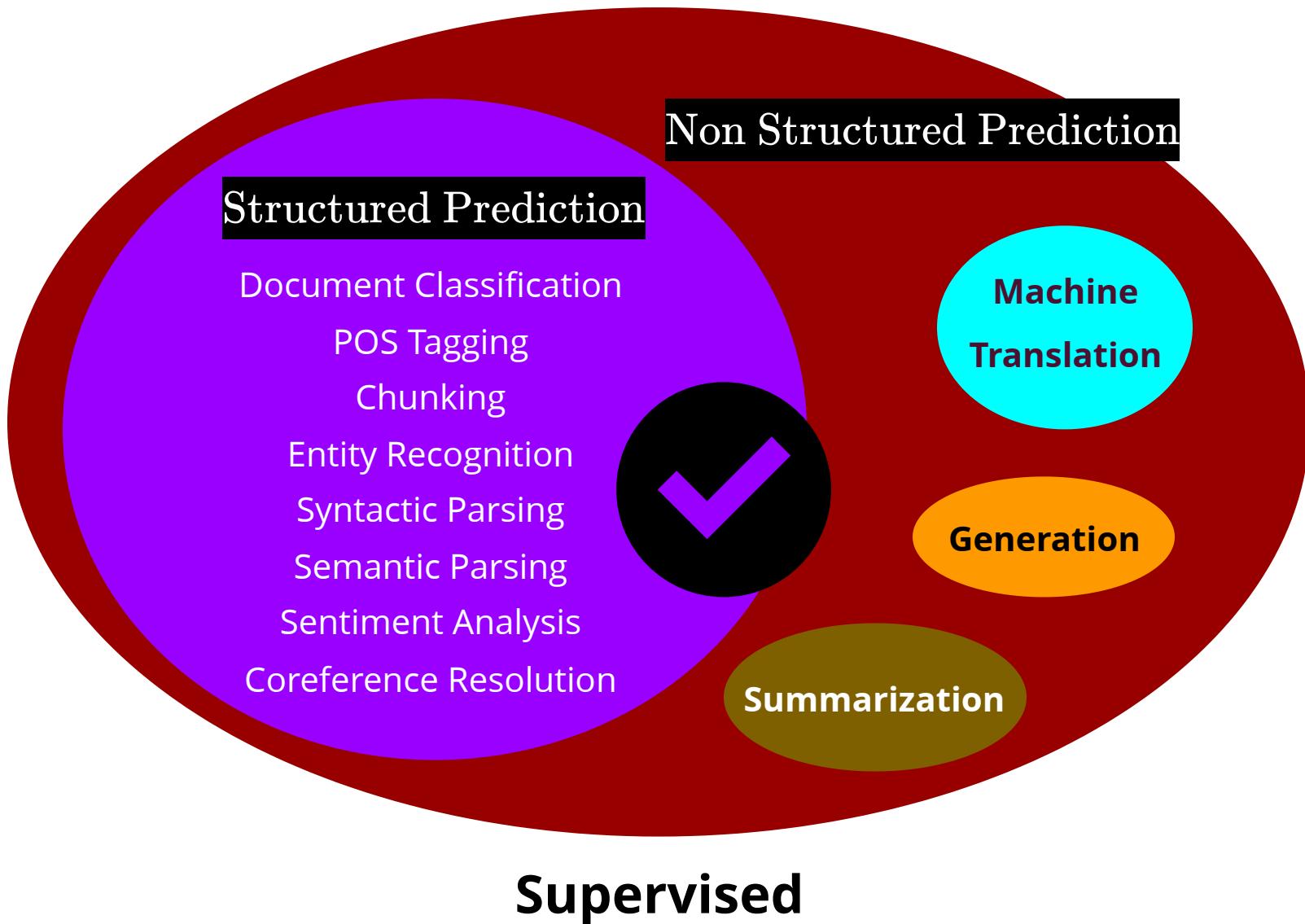
*a*

*banana*

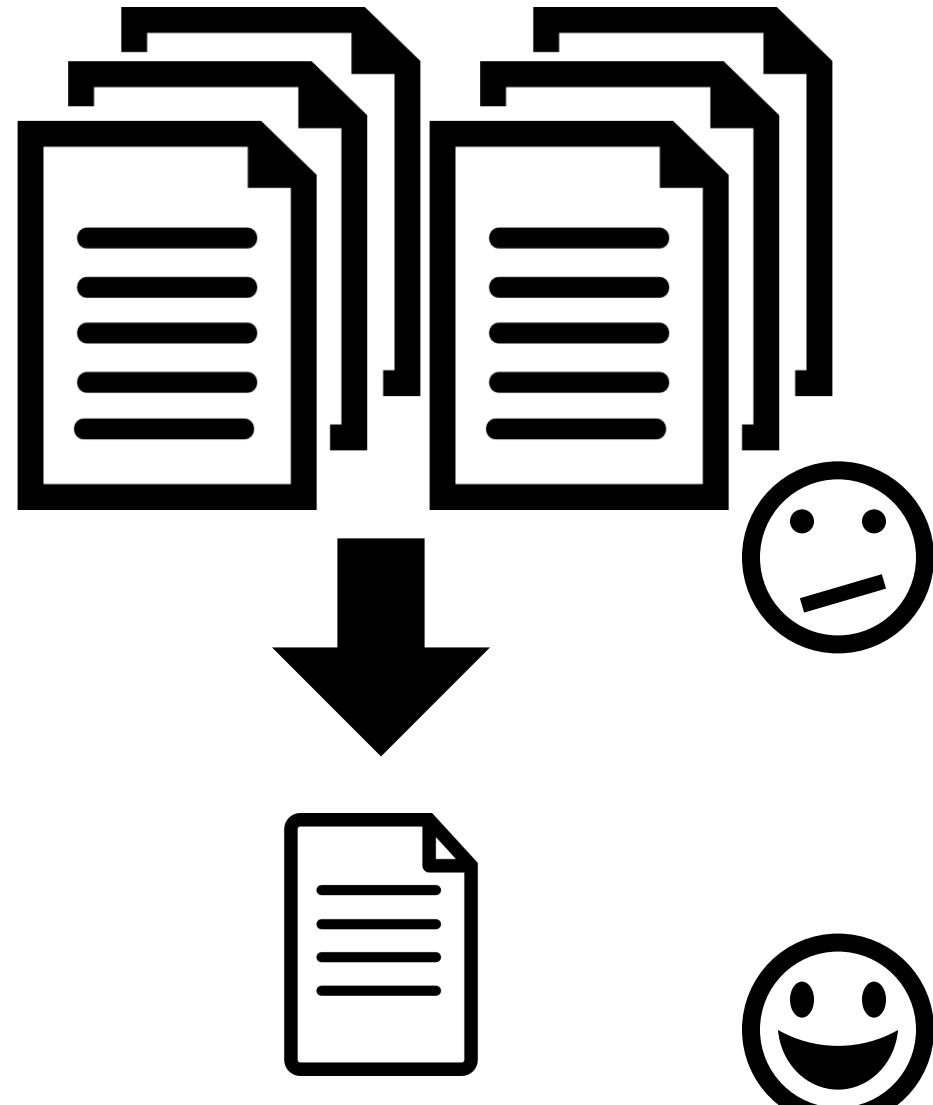
# A Unified Framework for Structured Prediction



# Tasks in NLP



# Text Summarization



# Other Generation Tasks

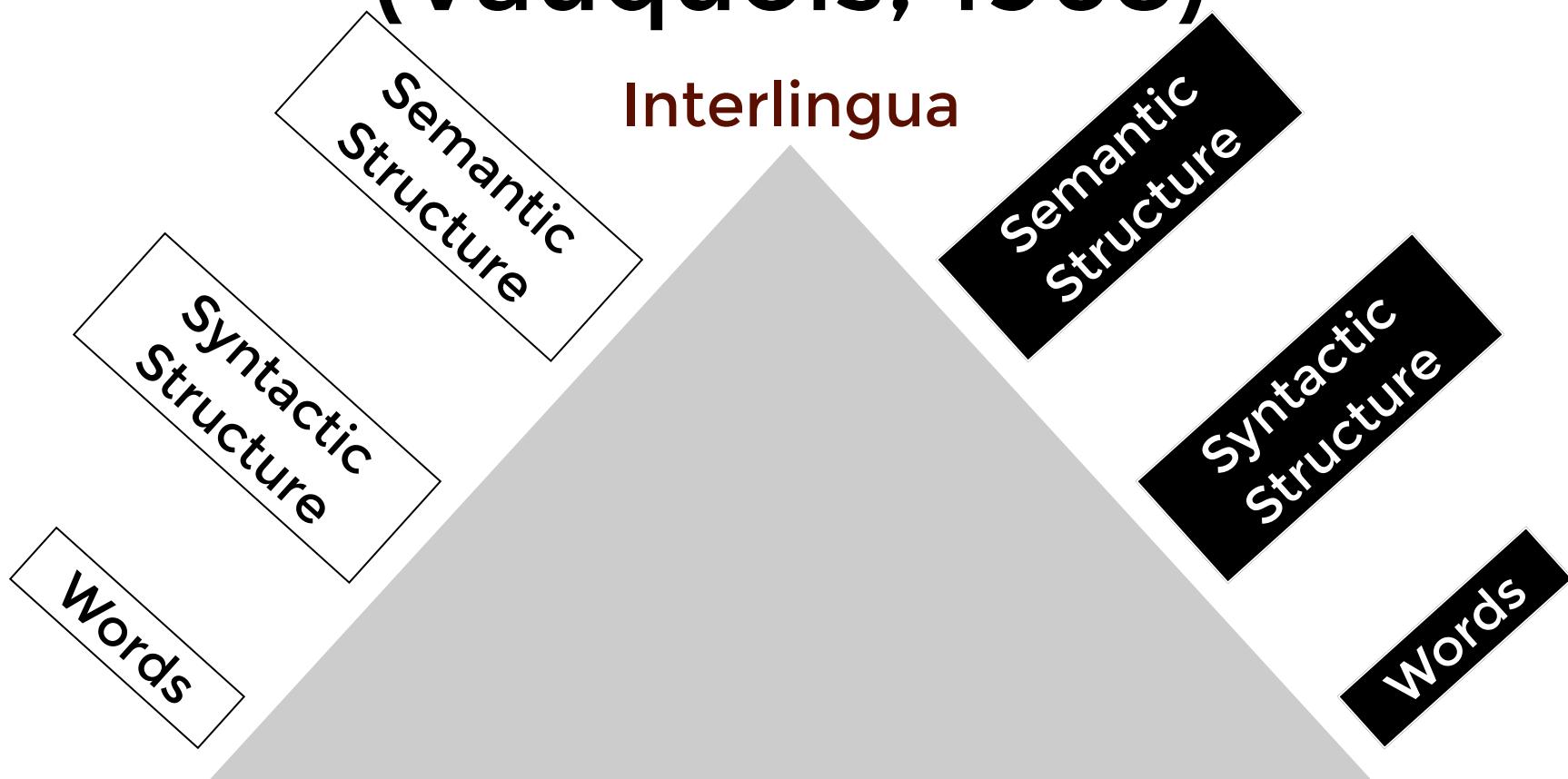
Dialogue  
Systems

Question  
Answering

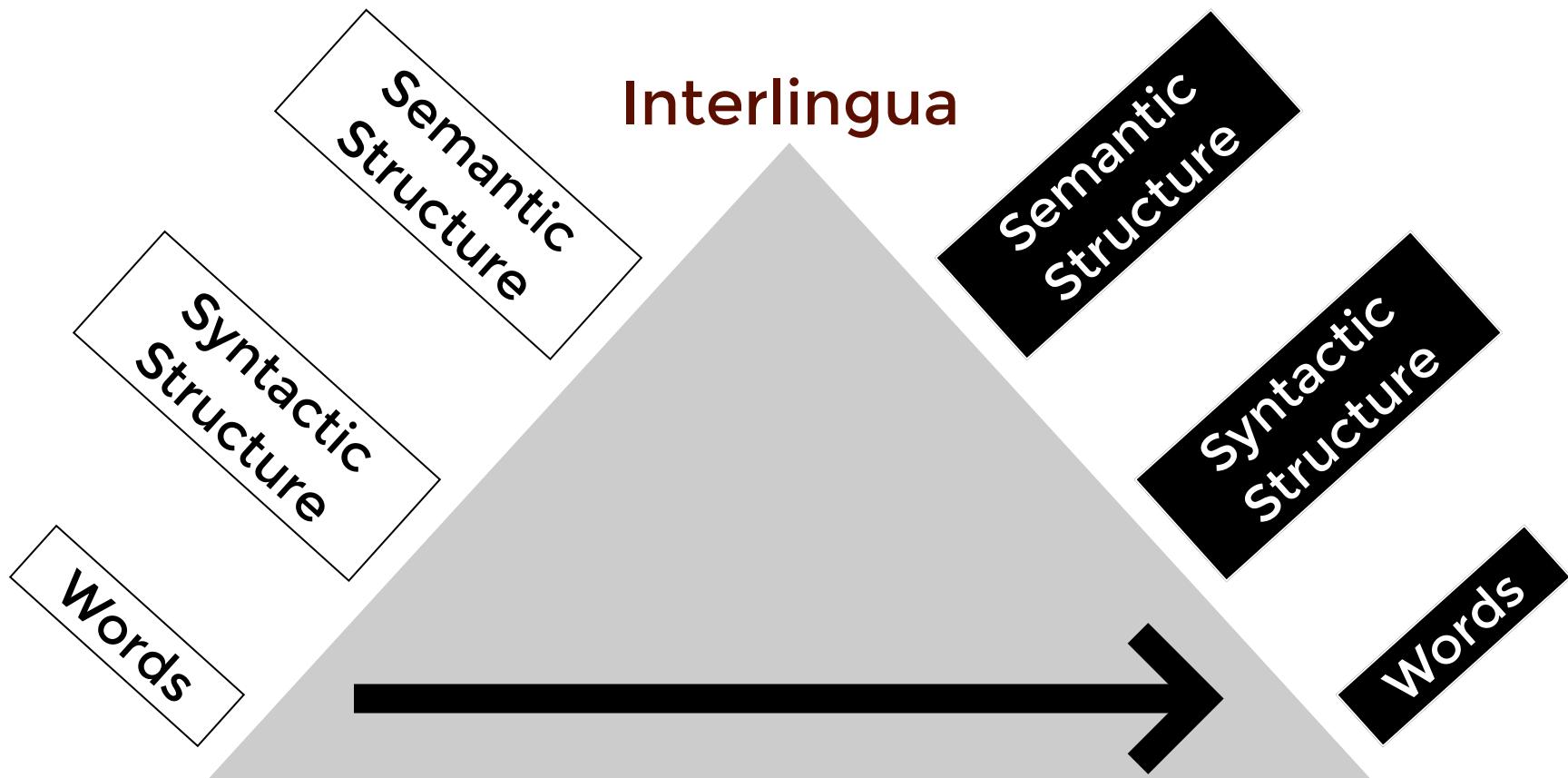
Machine  
Translation

They are all text-to-text problems!

# Machine Translation (Vauquois, 1968)



# Machine Translation



Text-to-text Problem

# IBM Model 1

$$p(f_1, f_2, \dots, f_m | e_0, e_1, e_2, \dots, e_n, m)$$

We would like to maximize:

$$\sum_{i=1}^m \log \left( \sum_{j=0}^n t(f_i | e_j) \right)$$

Ignores the word ordering information.  
Learn with Expectation Maximization.

NULL SUTD is the only university in the East .

新加坡 科技 设计 大学 是 东部 唯一 的 一 所 大学。

# IBM Model 2

$$p(f_1, f_2, \dots, f_m | e_0, e_1, e_2, \dots, e_n, m)$$

$$p(f_1, f_2, \dots, f_m, a_1, a_2, \dots, a_m | e_0, e_1, e_2, \dots, e_n, m)$$

$$\prod_{i=1}^m q(a_i | i, n, m) t(f_i | e_{a_i})$$

It additionally learns a distribution that captures (absolute) word reordering

NULL SUTD is the only university in the East .

新加坡 科技 设计 大学 是 东部 唯一 的 一 所 大学。

# IBM Models

Pioneering work on word-level  
translation

Not particularly useful models for  
translation themselves, but can  
yield useful alignment  
information for the training set

A first step towards building other  
advanced models such as phrase-  
based and syntax-based models

# Phrase-based Translation

**LM : Language Model**

How well the translated sentence reads

**TM : Phrase Translation Model**

How faithful the translation is to the original

**DM : Distortion Model**

How much efforts on "moving the eyes" in translation is required

# Phrase-based Translation

$$p_5 = (6, 6, \text{in the East})$$

SUTD is the only university **in the East**

新加坡 科技 设计 大学 是 **东部** 唯一 的 一所 大学 。

$$\underbrace{\log q(\text{in}|\text{only, university}) + \log q(\text{the}|\text{university, in}) + \log q(\text{East}|\text{in, the})}_{\text{Language model}}$$

$$+ \underbrace{\log t(\text{东部}|\text{in the East})}_{\text{Phrase translation model}}$$

$$+ \underbrace{\eta \times 5}_{\text{Distortion model}}$$

# Phrase-based Translation

We know how to score a translation derivation, but how do we search for the most optimal derivation?

The position of the last French word in the previous French phrase translated

The score of the partial derivation so far

A state

$$s = (e_1, e_2, b, r, \alpha)$$

The last two English words in the previous translated English phrase

A bit string indicating which words in French are (not yet) translated.

# Phrase-based Translation

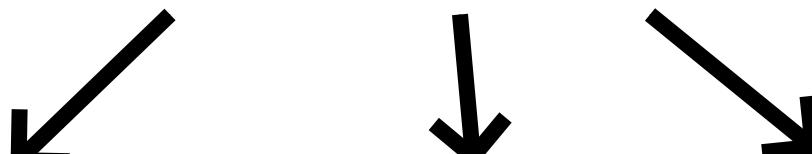
$$p_2 = (5, 5, \text{is})$$

SUTD is

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

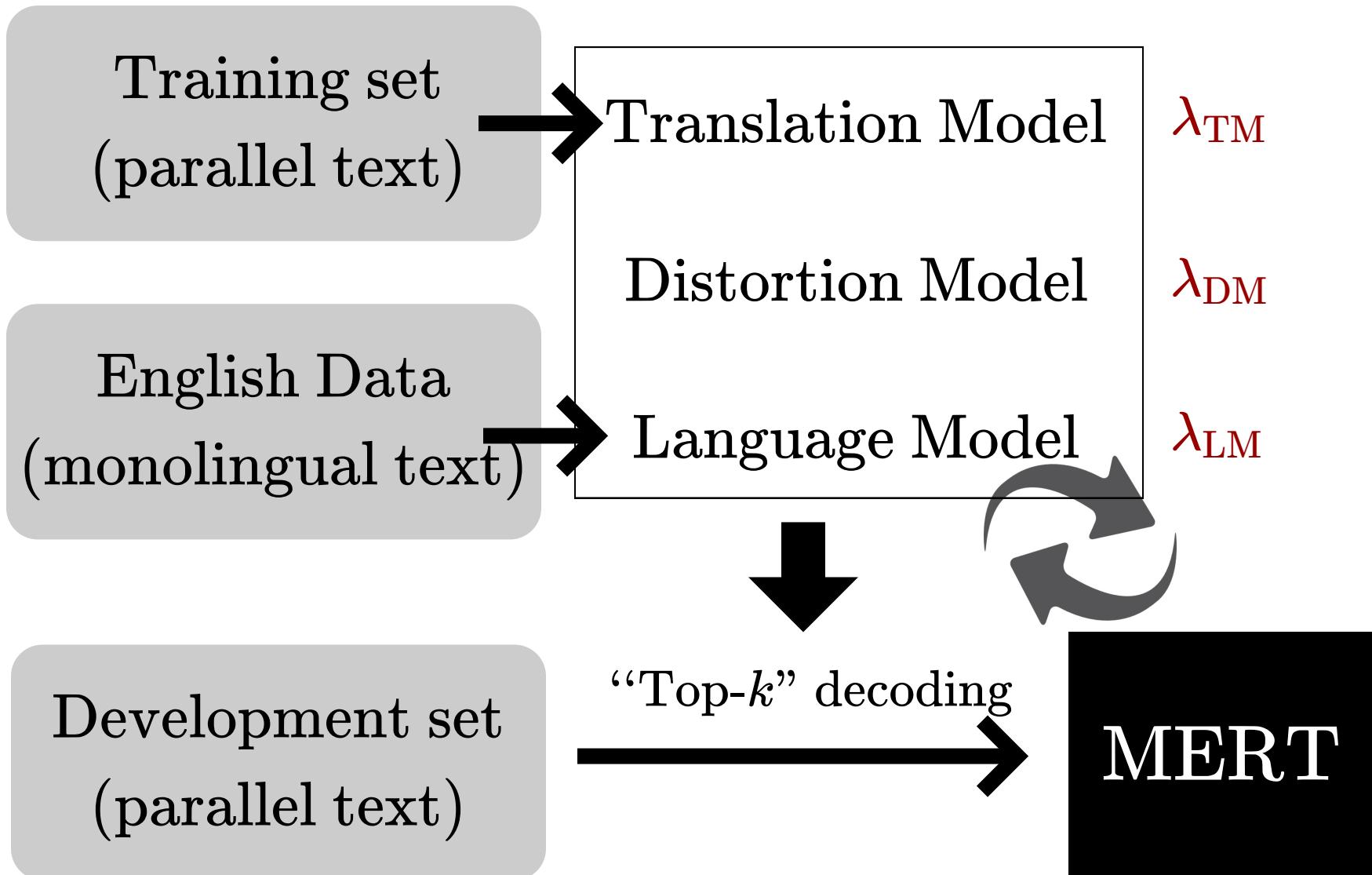
Each  $p_k$  is essentially an action!

(⟨START⟩, SUTD, 11111000000, 5, 8.9)



(..., 11111000010, ...)   ...   (... , 11111000111, ...)

# Phrase-based Translation



# BLEU



**BLEU: a Method for Automatic Evaluation of Machine Translation**

**Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu**  
IBM T. J. Watson Research Center  
Yorktown Heights, NY 10598, USA  
[{papineni,roukos,toddward,weijing}@us.ibm.com](mailto:{papineni,roukos,toddward,weijing}@us.ibm.com)

**Abstract**

Human evaluations of machine translation are extensive but expensive. Human evaluations can take months to finish and involve human labor that can not be reused. We propose a method of automatic machine translation evaluation that is quick, inexpensive, and language-independent, that correlates highly with human evaluation, and that has little marginal cost per run. We present this method as an automated understudy to skilled human judges which substitutes for them when there is need for quick or frequent evaluations.<sup>1</sup>

**1 Introduction**

**1.1 Rationale**

Human evaluations of machine translation (MT) weigh many aspects of translation, including *adequacy, fidelity, and fluency* of the translation (Hovy, 1999; White and O'Connell, 1994). A comprehensive catalog of MT evaluation techniques and their rich literature is given by Reeder (2001). For the most part, these various human evaluation approaches are quite expensive (Hovy, 1999). Moreover, they can take weeks or months to finish. This is a big problem because developers of machine translation systems need to monitor the effect of daily changes to their systems in order to weed out bad ideas from good ideas. We believe that MT progress stems from evaluation and that there is a logjam of fruitful research ideas waiting to be released from BLEU.

<sup>1</sup>So we call our method the bilingual evaluation understudy.

the evaluation bottleneck. D  
fit from an inexpensive autom  
quick, language-independent, a  
with human evaluation. We prop  
tion method in this paper.

**1.2 Viewpoint**

How does one measure translation p  
The closer a machine translation is to a p  
human translation, the better it is. This n  
tral idea behind our proposal. To judge th  
of a machine translation, one measures its cl  
to one or more reference human translations ac  
ing to a numerical metric. Thus, our MT evalua  
system requires two ingredients:

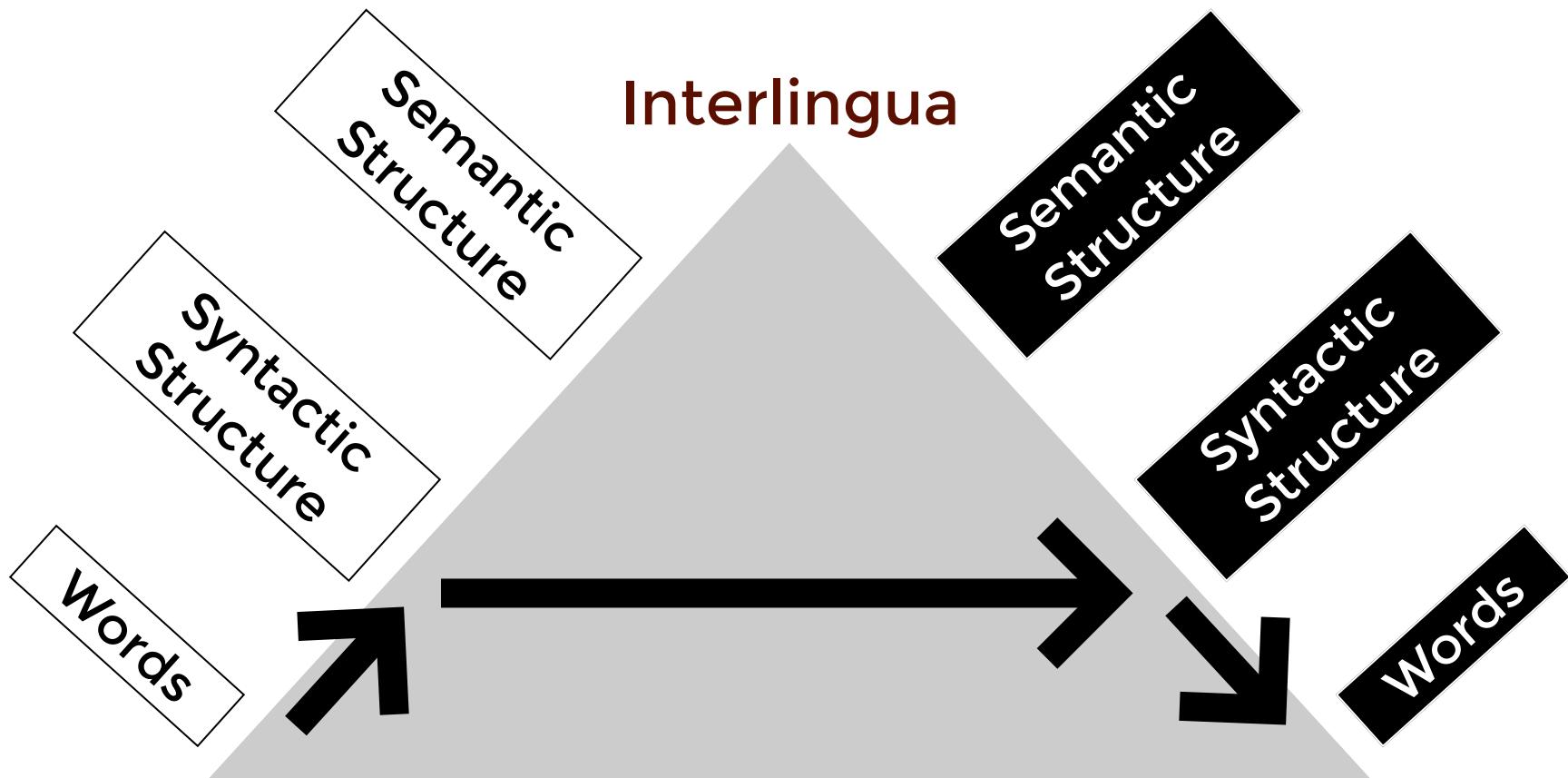
1. a numerical "translation closeness" metric
2. a corpus of good quality human reference trans  
lations

We fashion our closeness metric after the highly suc  
cessful word error rate metric used by the speech  
recognition community, appropriately modified for  
multiple reference translations and allowing for le  
gitimate differences in word choice and word or  
der. The main idea is to use a weighted average of  
variable length phrase matches against the reference  
translations. This view gives rise to a family of met  
rics using various weighting schemes. We have se  
lected a promising baseline metric for

In Section 2, we d  
detail. In Sectio  
BLEU. In Sectio  
experiment. In S  
metric performan

The most widely adopted evaluation metric for measuring MT quality.

# Machine Translation



# Synchronous Grammar



**Stochastic Inversion Transduction Grammars and Bilingual Parsing of Parallel Corpora**

Dekai Wu\*  
Hong Kong University of Science and Technology

We introduce (1) a novel stochastic inversion transduction grammar formalism for bilingual language modeling of sentence-pairs, and (2) the concept of bilingual parsing with a variety of parallel corpus analysis applications. Aside from the bilingual orientation, three major features distinguish the formalism from the finite-state transducers more traditionally found in computational linguistics: it skips directly to a context-free rather than finite-state base, it permits a minimal extra degree of ordering flexibility, and its probabilistic formulation admits an efficient maximum-likelihood bilingual parsing algorithm. A convenient normal form is shown to exist. Analysis of the formalism's expressiveness suggests that it is particularly well suited to modeling ordering shifts between languages, balancing needed flexibility against complexity constraints. We discuss a number of examples of how stochastic inversion transduction grammars bring bilingual constraints to bear upon problematic corpus analysis tasks such as segmentation, bracketing, phrasal alignment, and parsing.

**1. Introduction**

We introduce a general formalism for modeling of bilingual sentence pairs, known as an **inversion transduction grammar**, with potential application in a variety of corpus analysis areas. Transduction grammar models, especially of the finite-state family, have long been known. However, the imposition of identical ordering constraints upon both streams severely restricts their applicability, and thus transduction grammars have received relatively little attention in language-modeling research. The inversion transduction grammar formalism skips directly to a context-free, rather than finite-state, base and permits one extra degree of ordering flexibility, while retaining properties necessary for efficient computation, thereby sidestepping the limitations of traditional transduction grammars.

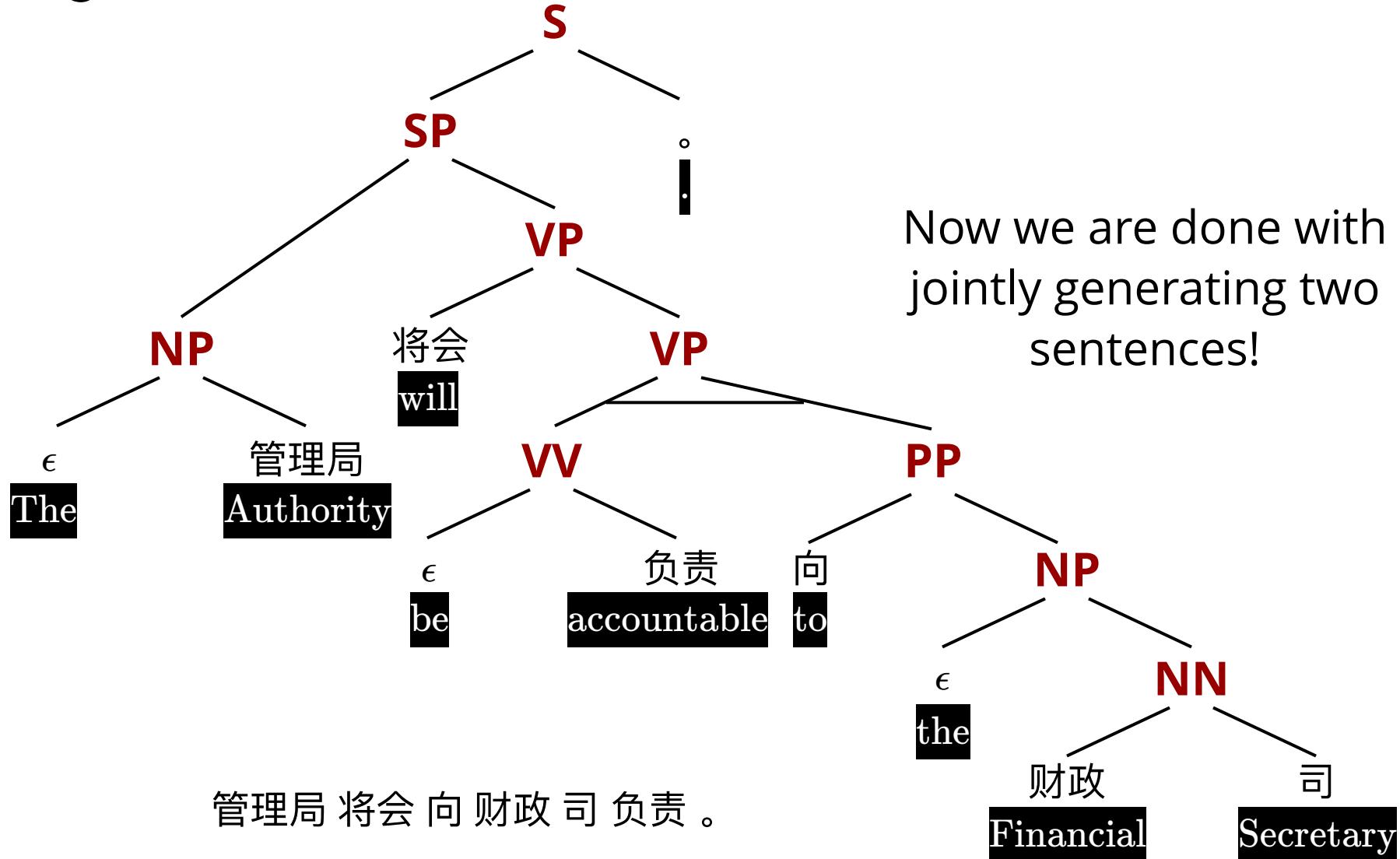
In tandem with the concept of bilingual language-modeling, we propose the concept of bilingual parsing, where the input is a sentence-pair rather than a sentence. Though inversion transduction grammars remain inadequate as full-fledged translation models, bilingual parsing with simple inversion transduction grammars turns out to be very useful for parallel corpus analysis when the true grammar is not fully known. Parallel bilingual corpora have been shown to provide a rich source of constraints for statistical analysis (Brown et al. 1990; Gale and Church 1991; Gale, Church, and Yarowsky 1992; Church 1993; Brown et al. 1993; Dagan, Church, and Gale 1993;

\* Department of Computer Science, University of Science and Technology, Clear Water Bay, Hong Kong  
E-mail: dekai@cs.ust.hk

© 1997 Association for Computational Linguistics

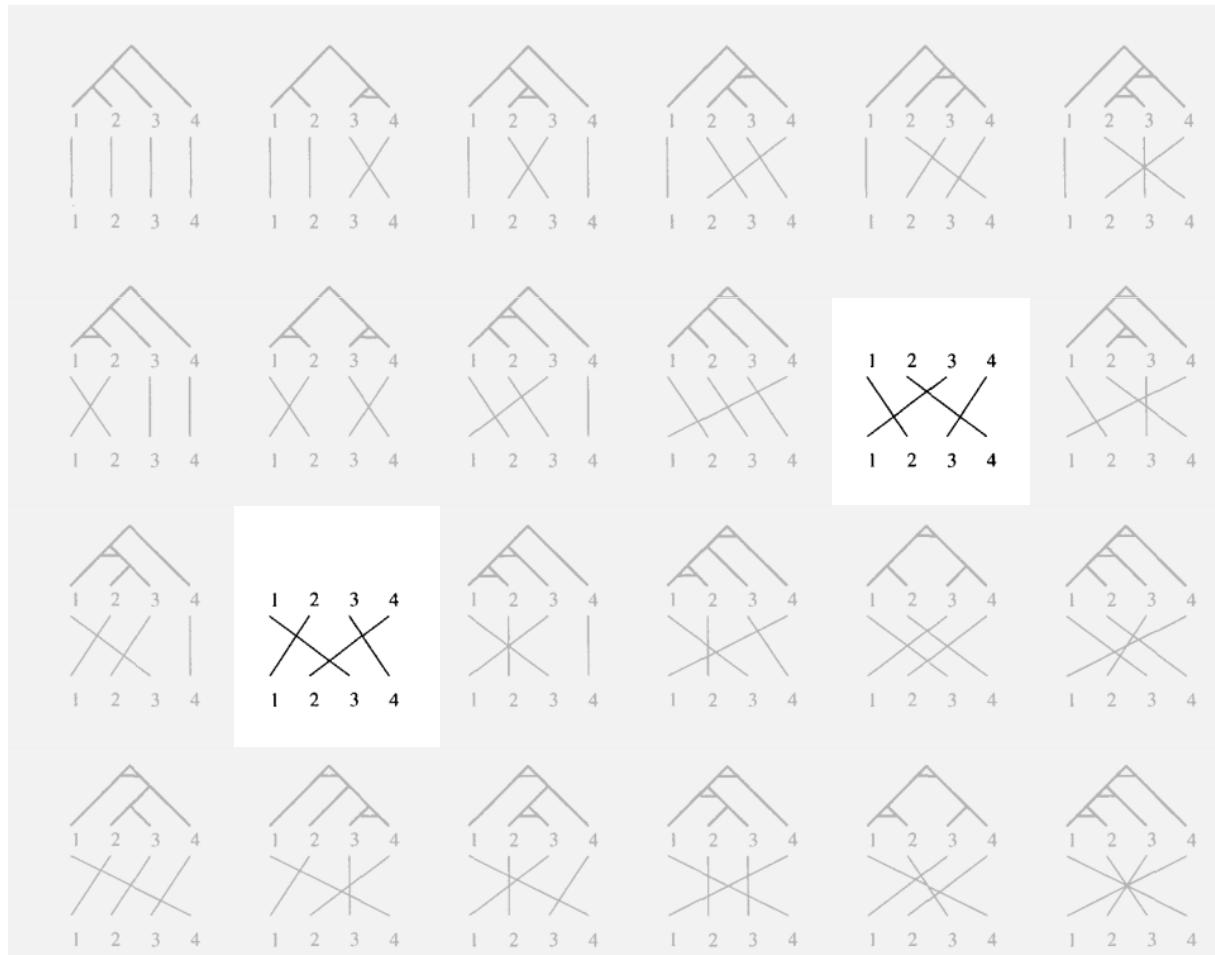
**Inversion Transduction Grammar**

# Synchronous Derivation



The Authority will be accountable to the Financial Secretary .

# Limitations with ITG



2 out of 24 cases where the ITG is unable to handle.

# Formal Syntax

## Hierarchical Phrase-Based Translation

David Chiang\*  
Information Sciences Institute  
University of Southern California

We present a statistical machine translation model that uses hierarchical phrases—phrases that contain subphrases. The model is formally a synchronous context-free grammar but learned from a parallel text without any syntactic annotations. Thus it can be seen as combining fundamental ideas from both syntax-based translation and phrase-based translation. We describe our system's training and decoding methods in detail, and evaluate it for translation speed and translation accuracy. Using BLEU as a metric of translation accuracy, we find that our system performs significantly better than the Alignment Template System, a state-of-the-art phrase-based system.

### 1. Introduction

The alignment template translation model (Och and Ney 2004) and related phrase-based models advanced the state of the art in machine translation by expanding the basic unit of translation from words to **phrases**, that is, substrings of potentially unlimited size (but not necessarily phrases in any syntactic theory). These phrases allow a model to learn local reorderings, translations of multiword expressions, or insertions and deletions that are sensitive to local context. This makes them a simple and powerful mechanism for translation.

The basic phrase-based model is an instance of the noisy-channel approach (Brown et al. 1993). Following convention, we call the source language “French” and the target language “English”; the translation of a French sentence  $f$  into an English sentence  $e$  is modeled as:

$$\arg \max_e P(e | f) = \arg \max_e P(e, f) \quad (1)$$

$$= \arg \max_e (P(e) \times P(f | e)) \quad (2)$$

The phrase-based translation model  $P(f | e)$  “encodes”  $e$  into  $f$  by the following steps:

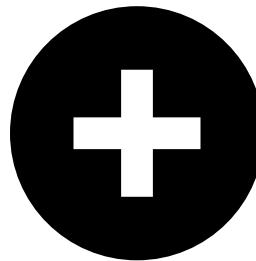
1. segment  $e$  into phrases  $\bar{e}_1 \dots \bar{e}_l$ , typically with a uniform distribution over segmentations;



\* 4676 Admiralty Way, Suite 1001, Marina del Rey, CA 90292, USA. E-mail: chiang@isi.edu. Much of the research presented here was carried out while the author was at the University of Maryland Institute for Advanced Computer Studies.  
Submission received: 1 May 2006; accepted for publication: 3 October 2006.

# Formal Syntax

Phrase-based  
Translation



Synchronous  
Grammar

It combines the idea of  
phrase-based translation  
and synchronous parsing.

# Phrase-based SCFG

澳洲 是 与 北韩 有 邦交 的 少数 国家 之一

Australia is one of the few countries that have diplomatic relations with North Korea

$\mathbf{X} \rightarrow (\text{与 } \mathbf{X}_1 \text{ 有 } \mathbf{X}_2, \text{ have } \mathbf{X}_2 \text{ with } \mathbf{X}_1) + 7.2$

$\mathbf{X} \rightarrow (\text{北韩}, \text{ North Korea}) - 0.7$

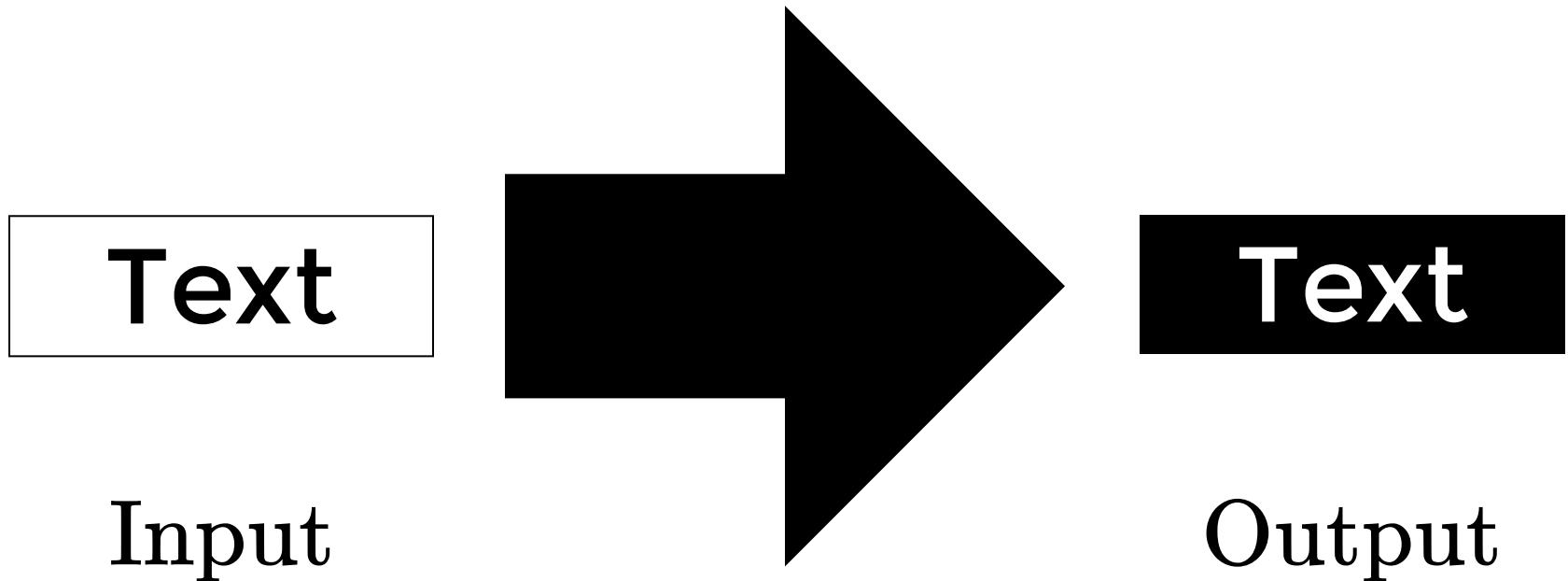
$\mathbf{X} \rightarrow (\text{邦交}, \text{ diplomatic relations}) + 3.9$

We are interested in a weighted  
phrase-level synchronous  
context-free grammar (SCFG)

# Score of Derivation

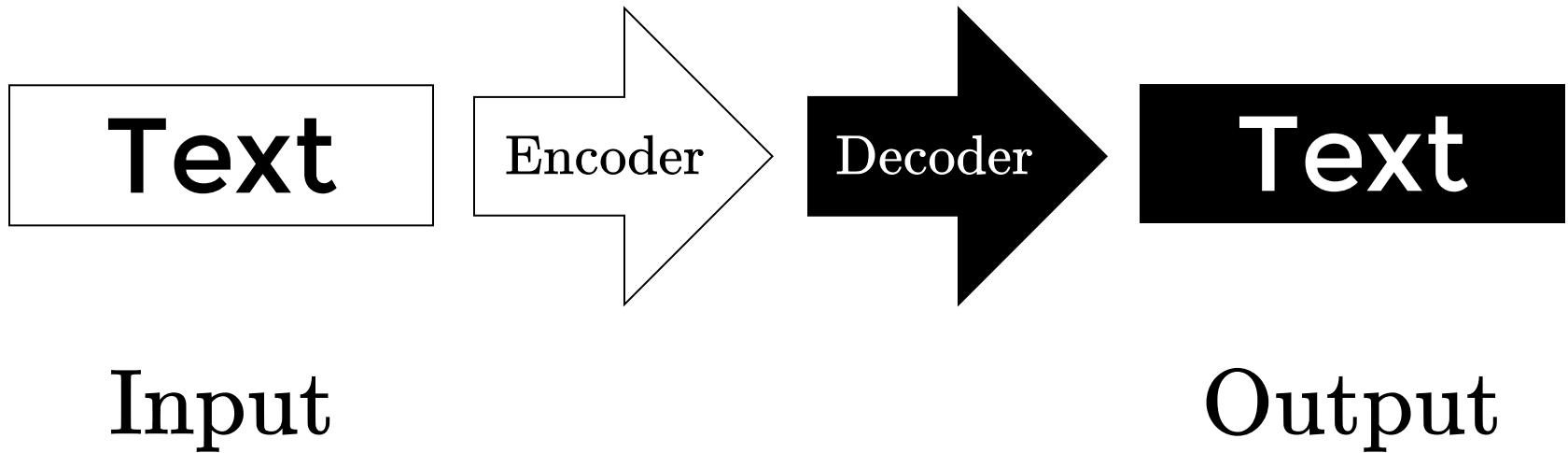
$$score(D) = \underbrace{\sum_{\mathbf{X} \rightarrow \langle \gamma, \alpha \rangle \in D} score(\mathbf{X} \rightarrow \langle \gamma, \alpha \rangle)}_{\text{Grammar rules}} + \underbrace{\log q(e)}_{\text{Language model}}$$

# Text-to-Text Generation



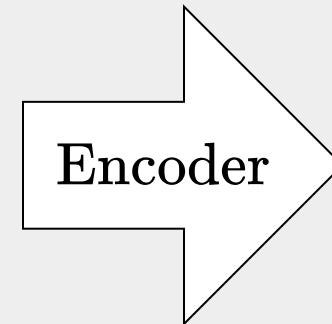
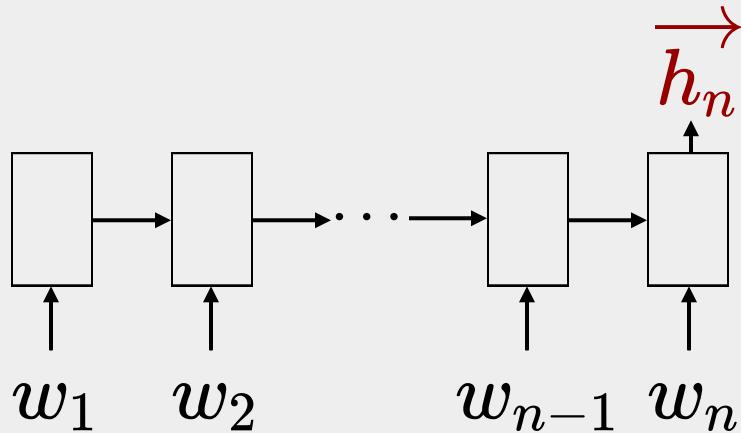
Text-to-text problem with **neural networks**

# Encoder-Decoder

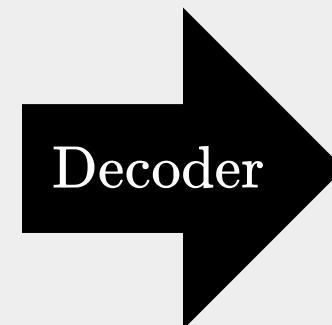
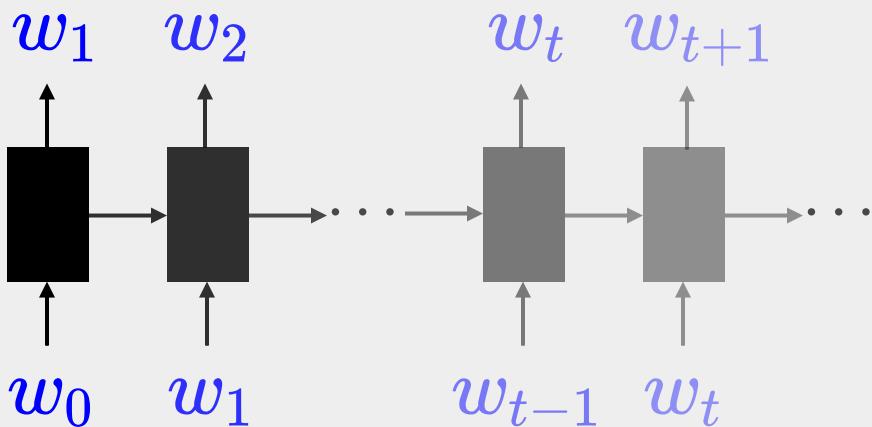


# LSTM

## Two applications of LSTM

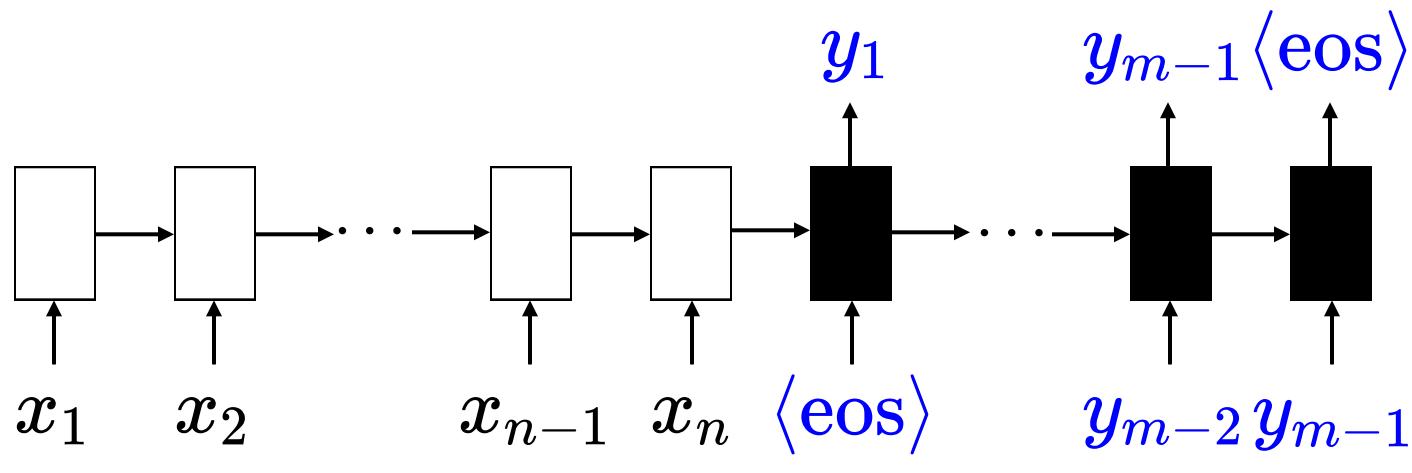


Sentence Embedding



Sentence Generation

# Sequence to Sequence

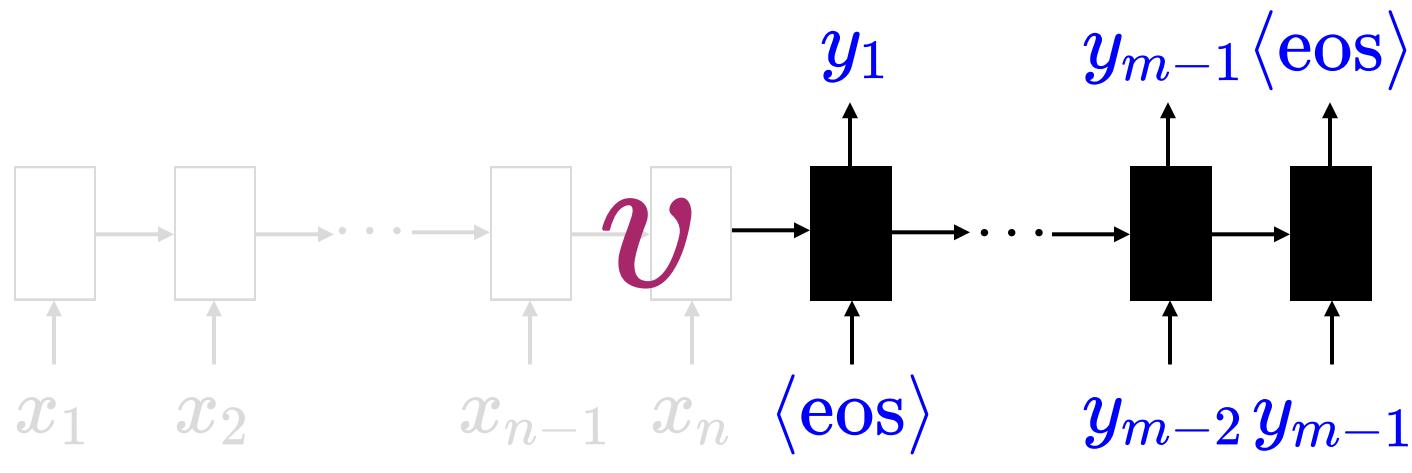


Encoder

Decoder

$$p(y_1, \dots, y_m | x_1, \dots, x_n)$$

# Sequence to Sequence

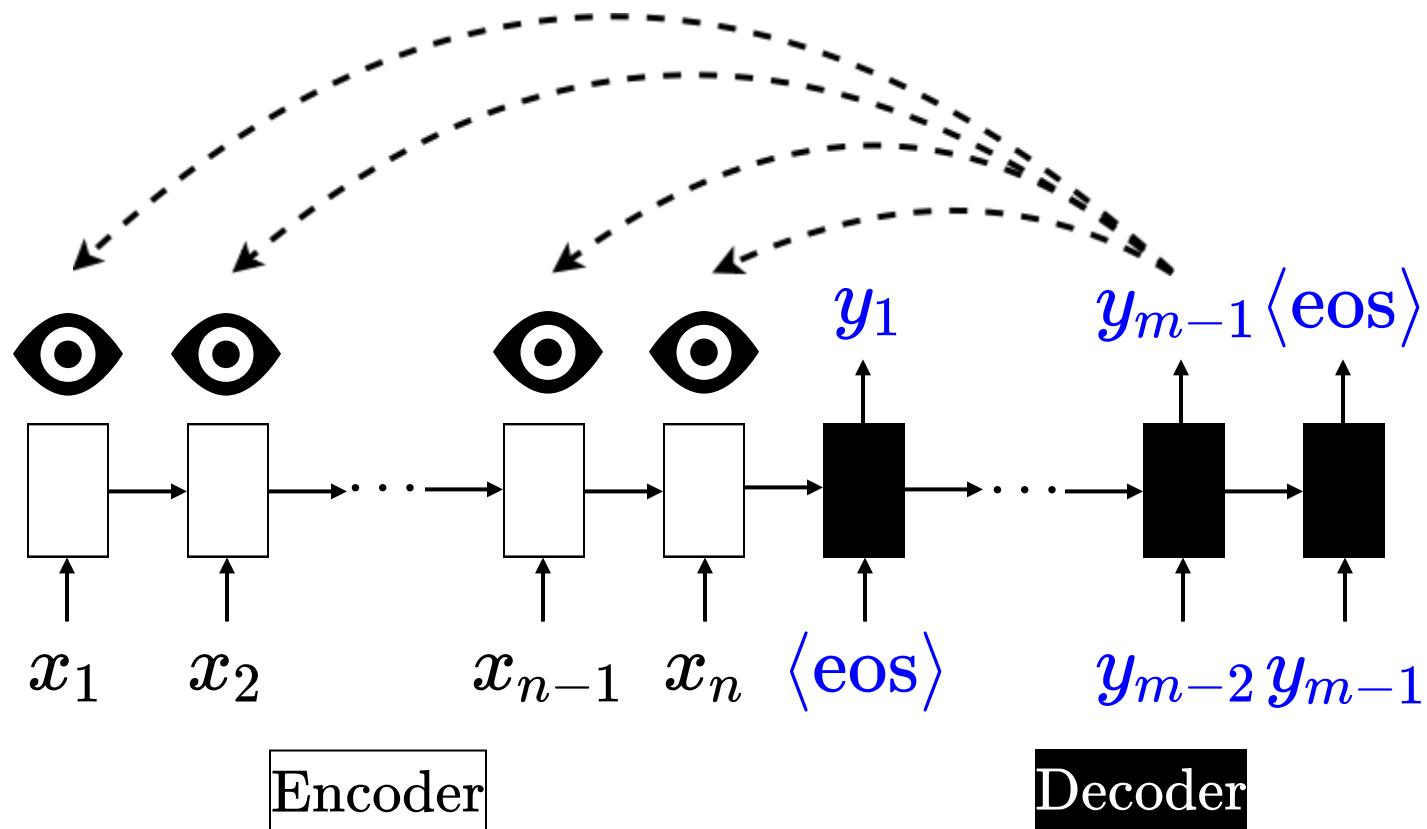


Encoder

Decoder

$$p(y_1, \dots, y_m | x_1, \dots, \mathcal{V}, x_n)$$

# Seq2Seq with Attention (Bahdanau et al. 2015)



We may need to look at the input sequence  
when generating the outputs...

# Attention

## Cosine

$$\frac{\mathbf{d}_i^T \mathbf{e}_j}{\|\mathbf{d}_i\| \cdot \|\mathbf{e}_j\|}$$

(Graves et al. 2014)

## Concatenation

$$\mathbf{v}^T \tanh(\mathbf{W}[\mathbf{e}_j; \mathbf{d}_i])$$

(Bahdanau et al. 2015)

## Scaled Product

$$\frac{\mathbf{d}_i^T \mathbf{e}_j}{\sqrt{\delta}}$$

(Vaswani et al. 2017)

## Dot-Product

$$\mathbf{d}_i^T \mathbf{e}_j$$

(Luong et al. 2015)

## General Product

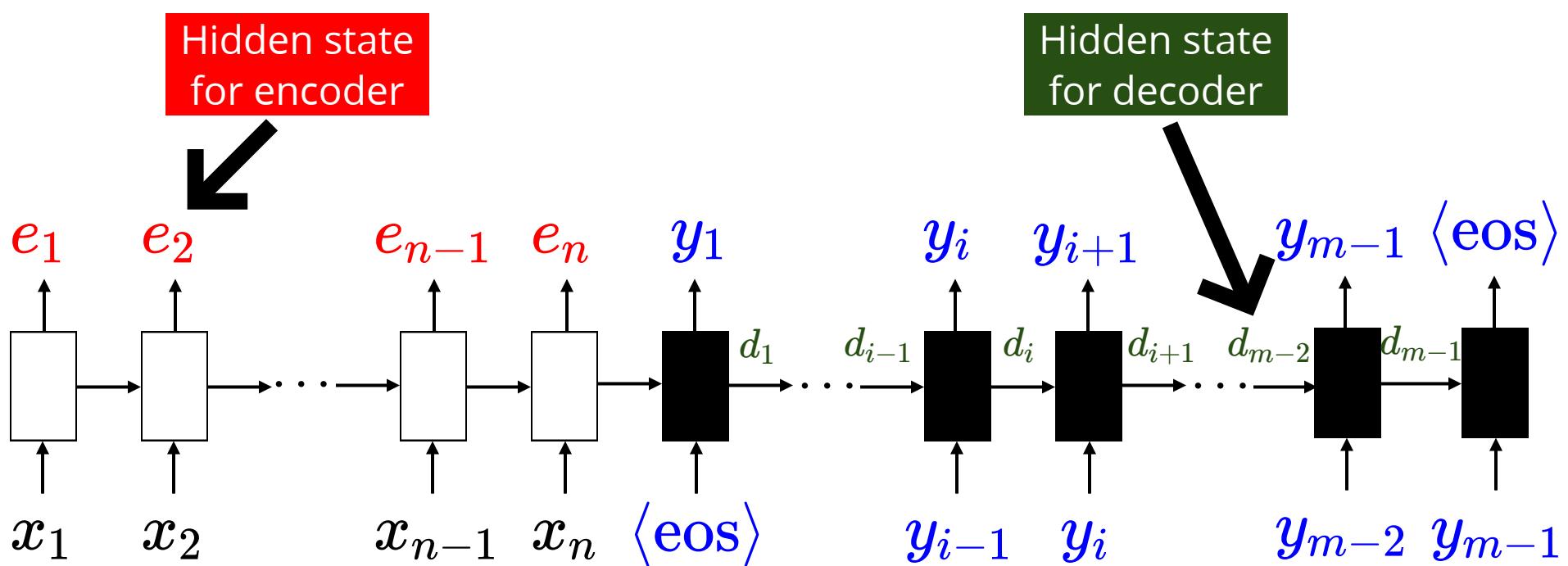
$$\mathbf{d}_i^T \mathbf{W} \mathbf{e}_j$$

(Luong et al. 2015)

# Seq2Seq with Attention

$$u_j^i = v^T \tanh(W_1 e_j + W_2 d_i)$$

Learnable parameters



Encoder

Decoder

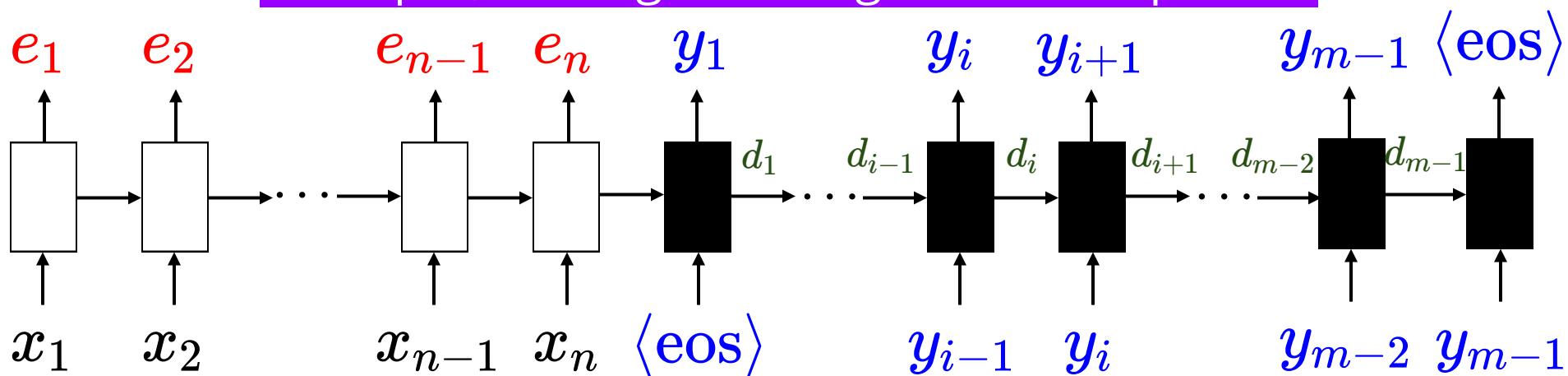
# Seq2Seq with Attention

$$u_j^i = v^T \tanh(W_1 e_j + W_2 d_i)$$

$$a_j^i = \text{softmax}(u_j^i)$$



How much I would like to focus on the j-th input, when generating the i-th output



Encoder

Decoder

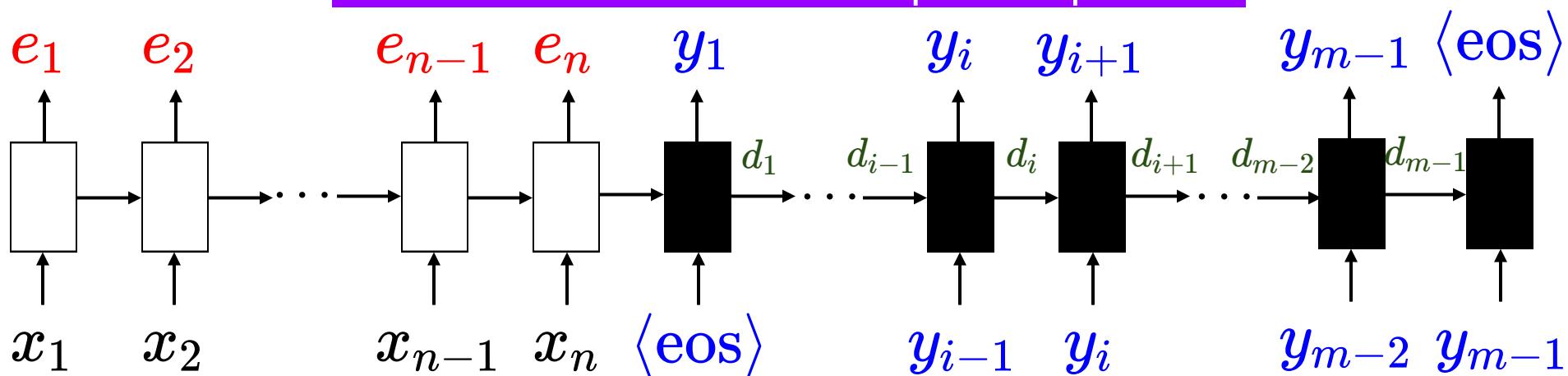
# Seq2Seq with Attention

$$u_j^i = v^T \tanh(W_1 e_j + W_2 d_i)$$

$$a_j^i = \text{softmax}(u_j^i)$$

$$d'_i = \sum_{j=1}^n a_j^i e_j$$

The weighted combination of hidden states from the entire input sequence



Encoder

Decoder

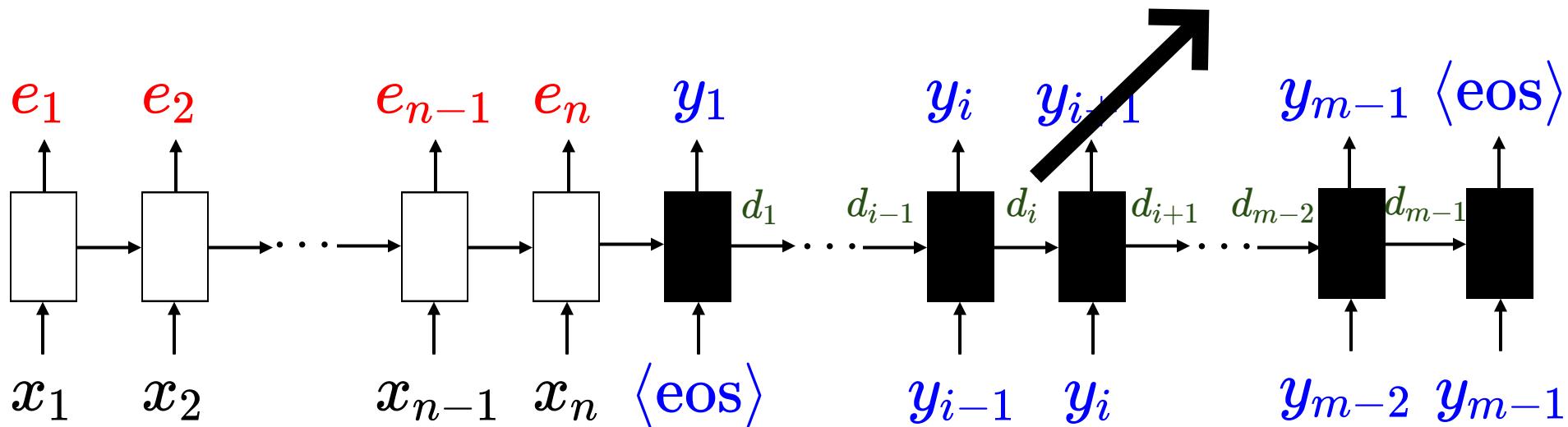
# Seq2Seq with Attention

$$u_j^i = v^T \tanh(W_1 e_j + W_2 d_i)$$

$$a_j^i = \text{softmax}(u_j^i)$$

$$d'_i = \sum_{j=1}^n a_j^i e_j$$

When predicting  $y_i$   
use  $[d_i, d'_i]$  instead



Encoder

Decoder

# Neural Models

**seq2seq**

recurrent network with attention

**fairseq**

convolutional network based encoder

**transformer**

self-attention network based encoder

# Neural Models

seq2seq

recurrent network with attention

fairseq

convolutional network based encoder

transformer

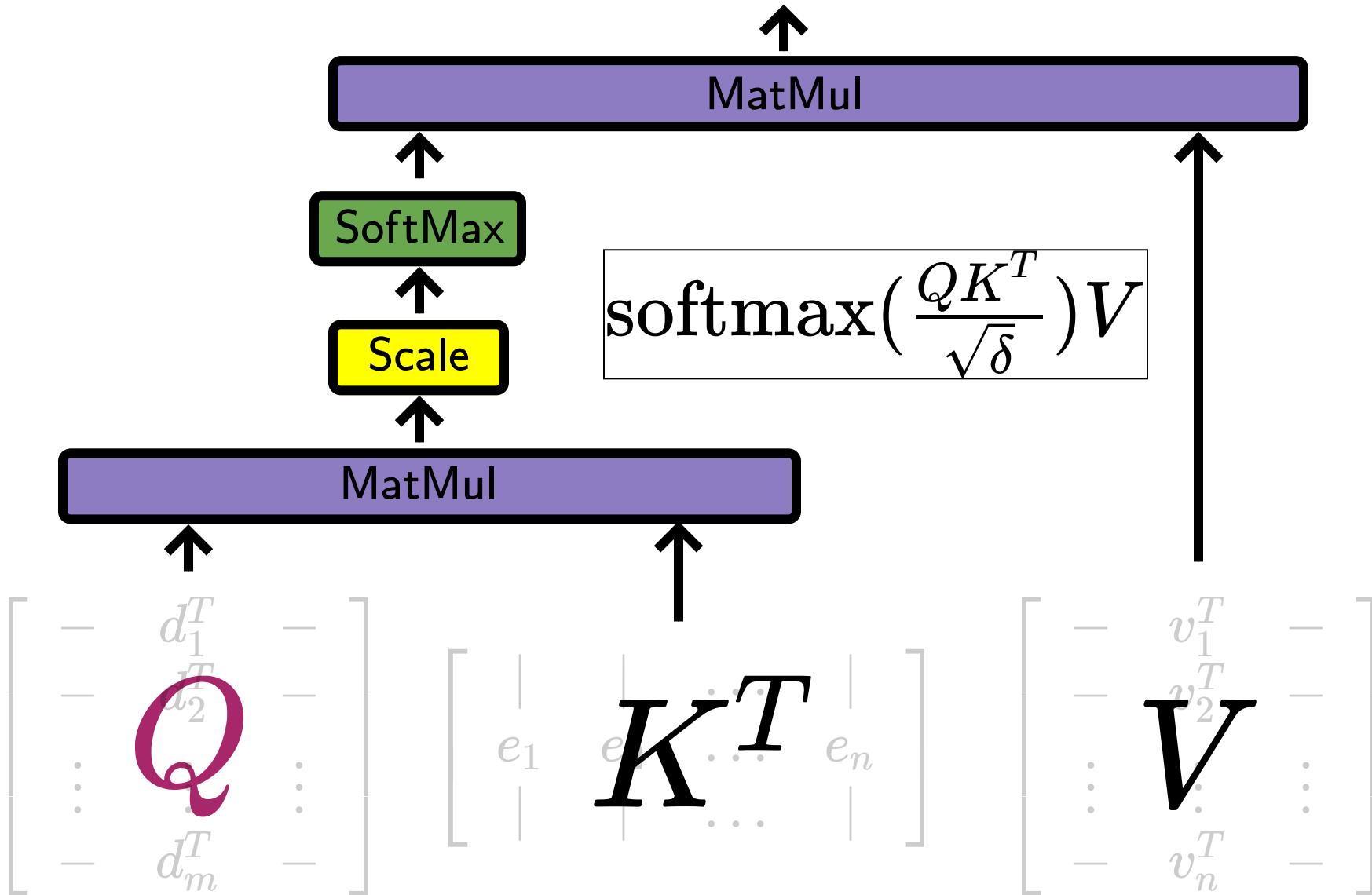
self-attention network based encoder

# Attention

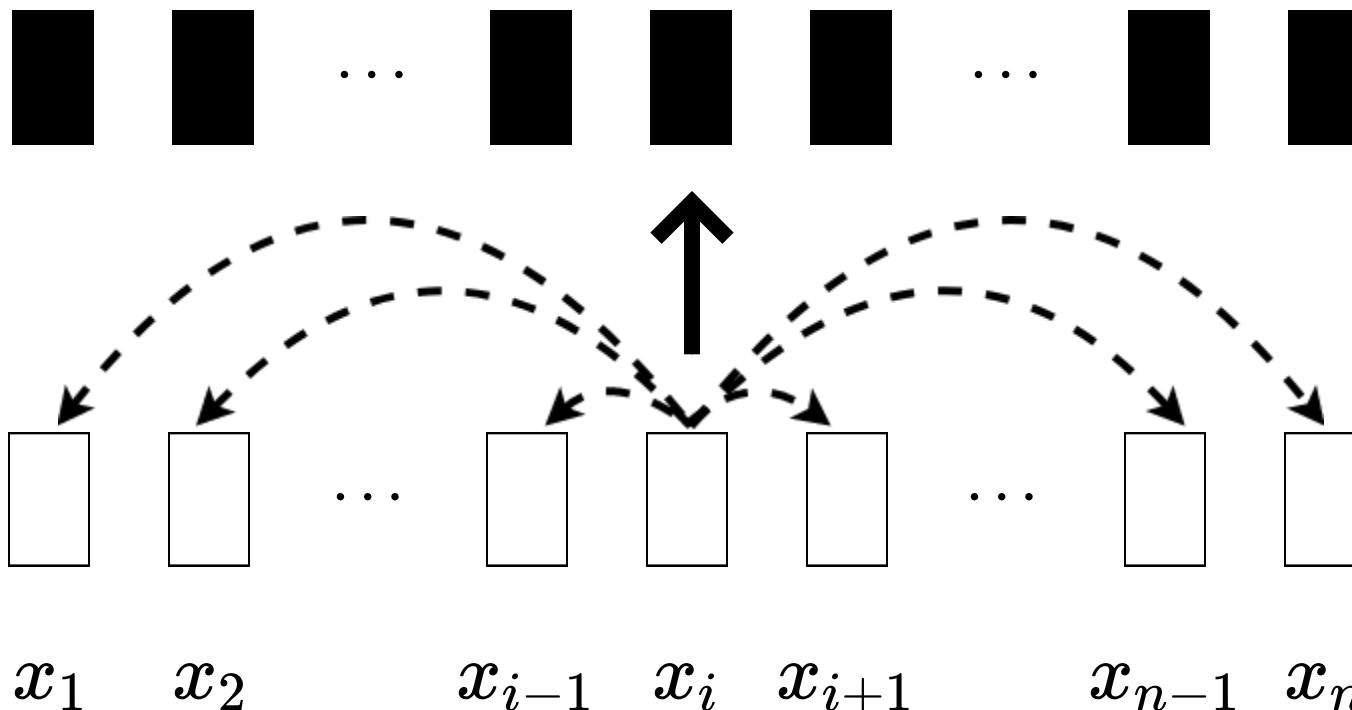
Query	Key	Value
$\begin{bmatrix} - & d_1^T & - \\ - & d_2^T & - \\ \vdots & \vdots & \vdots \\ - & d_m^T & - \end{bmatrix}$	$\begin{bmatrix} - & e_1^T & - \\ - & e_2^T & - \\ \vdots & \vdots & \vdots \\ - & e_n^T & - \end{bmatrix}$	$\begin{bmatrix} - & v_1^T & - \\ - & v_2^T & - \\ \vdots & \vdots & \vdots \\ - & v_n^T & - \end{bmatrix}$
$Q$	$K$	$V$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{\delta}}\right)V$$

# Scaled Dot Product



# Self Attention



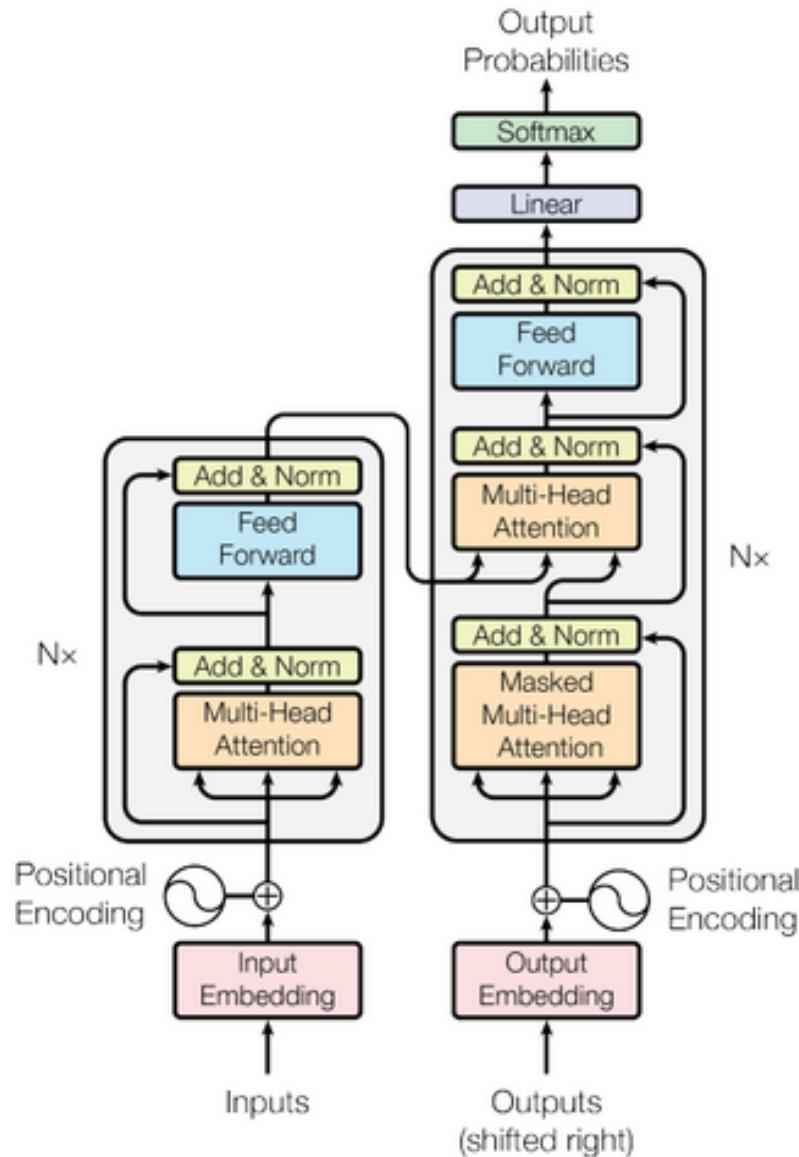
$$Q = K = V$$

Query, key and value matrices are the same!

# Transformer

Encoder

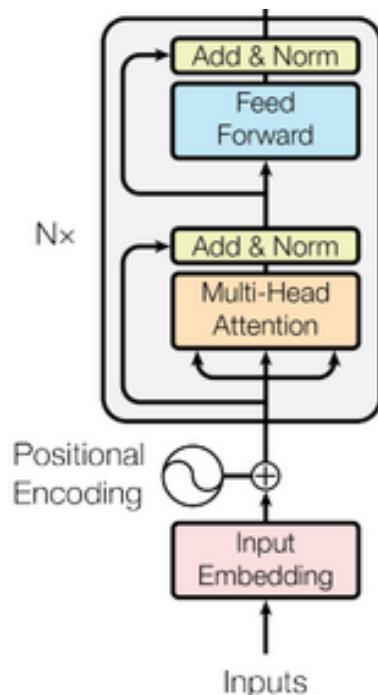
Decoder



# BERT

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	# #ing	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{my}$	$E_{dog}$	$E_{is}$	$E_{cute}$	$E_{[SEP]}$	$E_{he}$	$E_{likes}$	$E_{play}$	$E_{##ing}$	$E_{[SEP]}$
Segment Embeddings	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
Position Embeddings	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$

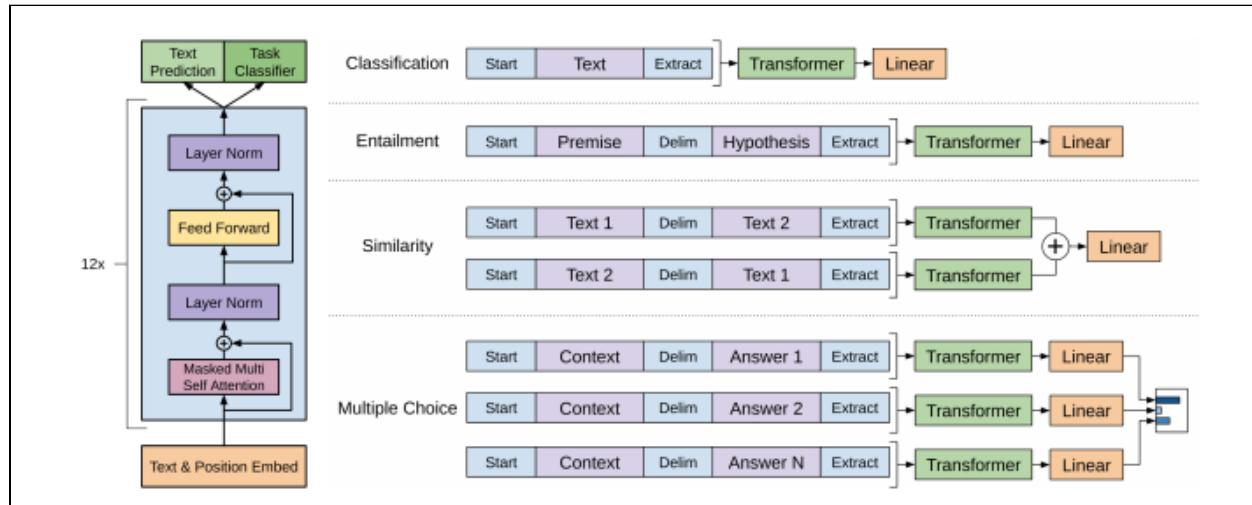
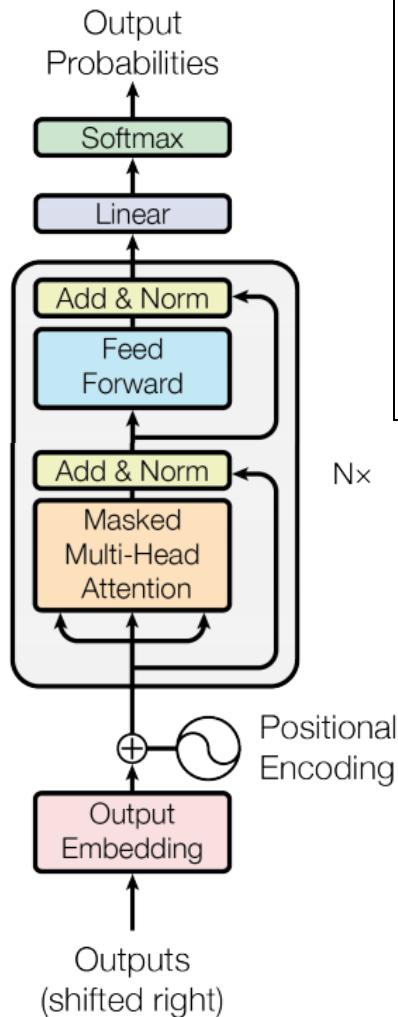
Encoder



Task #1: Masked LM  
Useful for learning context representations within a sequence

Task #2: Next Sentence Prediction  
Useful for tasks that involve identifying relations between multiple sentences

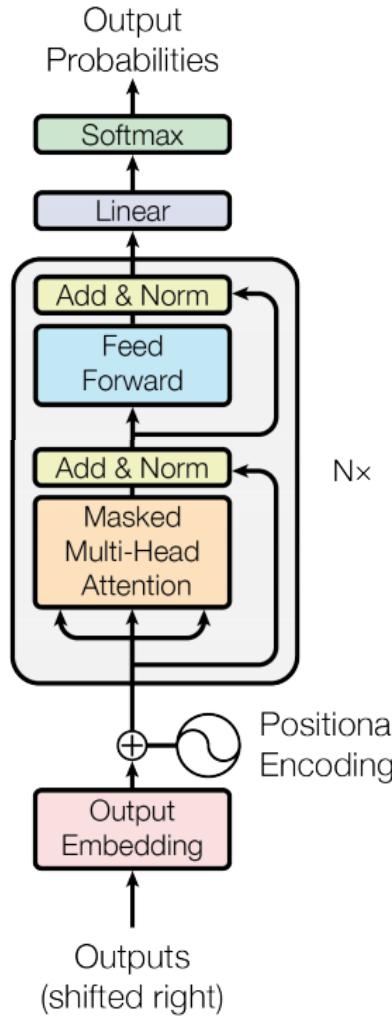
# GPT



Decoder

GPT: relies on task-specific fine-tuning,  
similar to BERT

# GPT-2, GPT-3



Model	Size (# Params)
GPT-2 (Open AI)	1.5 Billion
Megatron (Nvidia)	8 Billion
Turing NLG (Microsoft)	17 Billion
GPT-3 (Open AI)	175 Billion

Decoder

GPT-2: LM as unsupervised multitask learner  
GPT-3: LM as few-short learner

# Tasks in NLP



Document Classification

POS Tagging

Chunking

Entity Recognition

Syntactic Parsing

Semantic Parsing

Sentiment Analysis

Coreference Resolution

Natural Language Generation

Text Summarization

Machine Translation

Word Clusters

GloVe, word2vec

Topic Modeling

Language Modeling

ELMo, BERT, XLNet

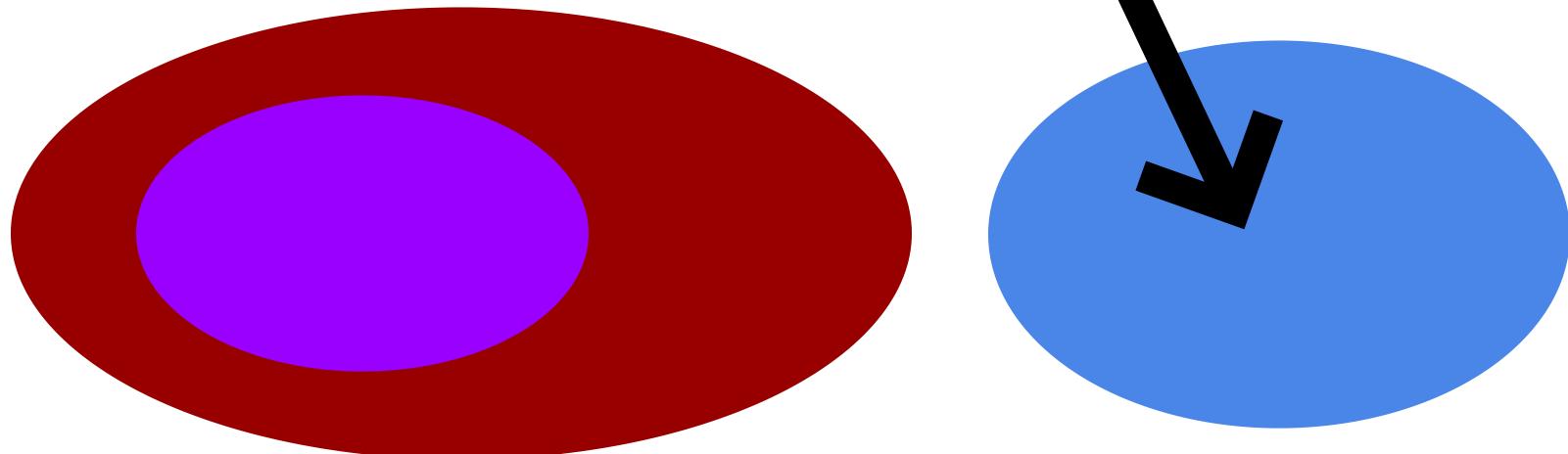


**Supervised**

**Unsupervised**

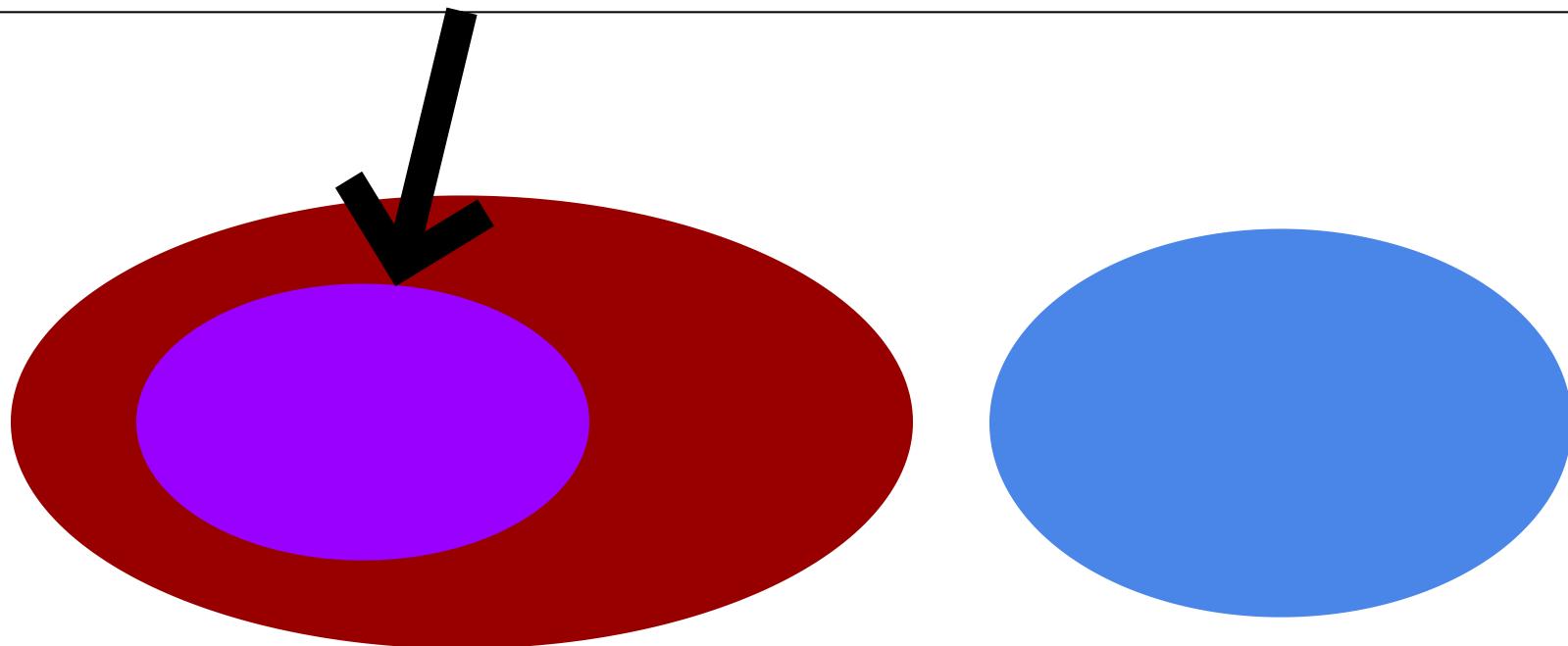
# Some Future Directions

Better methods for  
unsupervised learning of  
word/phrase/sentence  
representations



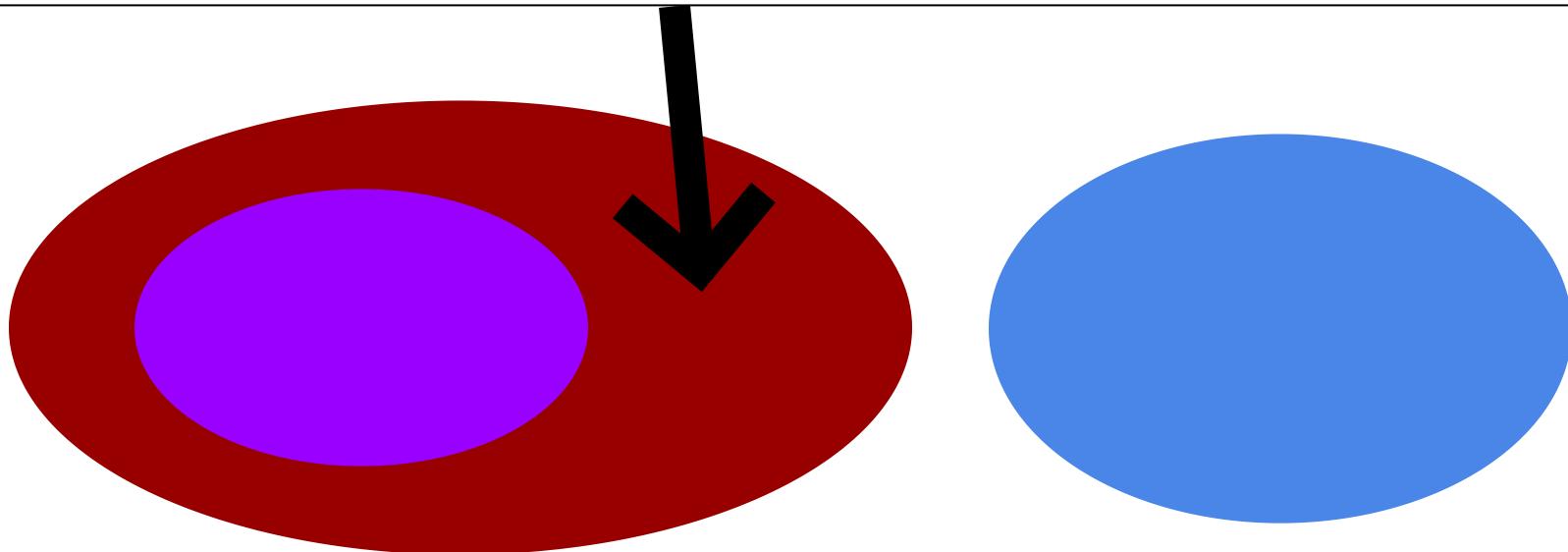
# Some Future Directions

Better understanding of the  
connection between Different  
**Structured Prediction Models**



# Some Future Directions

Better ways of training non  
structured prediction models



# Thank You

[statnlp.org](http://statnlp.org)

# We Are Hiring



清华大学  
Tsinghua University



浙江大学  
ZHEJIANG UNIVERSITY



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE



NYU



香港科技大学  
THE HONG KONG  
UNIVERSITY OF SCIENCE  
AND TECHNOLOGY



SINGAPORE UNIVERSITY OF  
TECHNOLOGY AND DESIGN



UMass  
Amherst



上海交通大学  
SHANGHAI JIAO TONG UNIVERSITY



National University  
of Singapore



A screenshot of the StatNLP website. At the top, there is a navigation bar with links to 'About', 'People', 'Publications', and 'Tools'. Below the navigation bar, there is a large image of a complex network graph with red and blue nodes and connecting lines. In the center of the page, there is a box containing the text 'Our Vision' and 'Conduct fine research &amp; Nurture world class researchers'. A large black arrow points downwards from the top of this section towards the bottom of the page.



COLUMBIA UNIVERSITY  
IN THE CITY OF NEW YORK



THE UNIVERSITY  
of EDINBURGH



Stanford  
University

Carnegie  
Mellon  
University



Cornell University

TEXAS  
The University of Texas at Austin

ByteDance

Tencent 腾讯

StatNLP.org