

ANTONIS ANASTASOPOULOS
CS499 INTRODUCTION TO NLP

CONDITIONAL RANDOM

FIELDS



<https://cs.gmu.edu/~antonis/course/cs499-spring21/>

With adapted slides by Graham Neubig

STRUCTURE OF THIS LECTURE

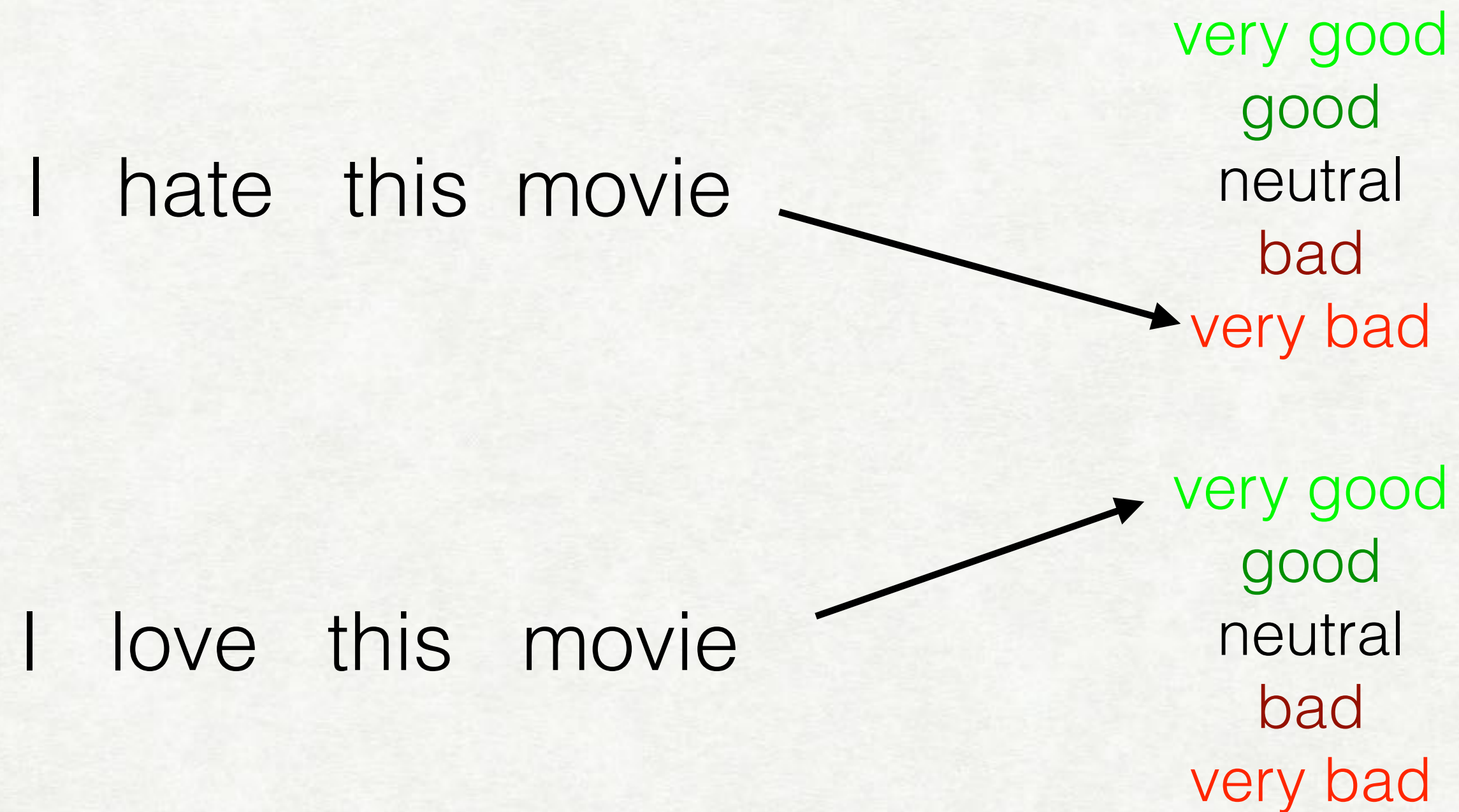
1 Structured Prediction

2 Structured Perceptron

3 Conditional Random Fields


4 Viterbi

A PREDICTION PROBLEM

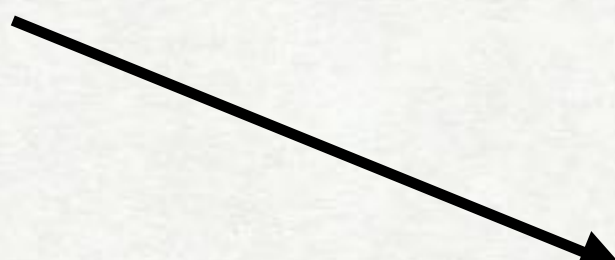


TYPES OF PREDICTION

Two classes (binary classification)

I hate this movie  positive
negative

- Multiple classes (**multi-class classification**)

I hate this movie  very good
good
neutral
bad
very bad

- Exponential/infinite labels (**structured prediction**)

I hate this movie  PRP VBP DT NN

I hate this movie  *kono eiga ga kirai*

WHY CALL IT "STRUCTURED" PREDICTION?

Classes are too numerous to enumerate

Need some sort of method to exploit the *problem structure* to learn efficiently

Example of "structure", the following two outputs are similar:

PRP VBP DT NN

PRP VBP VBP NN

MANY VARIETIES OF STRUCTURED PREDICTION!

Models:

RNN-based decoders

Convolution/self attentional decoders

CRFs w/ local factors

Training algorithms:

Structured perceptron, structured large margin

Sampling corruptions of data

Exact enumeration with dynamic programs

Reinforcement learning/minimum risk training

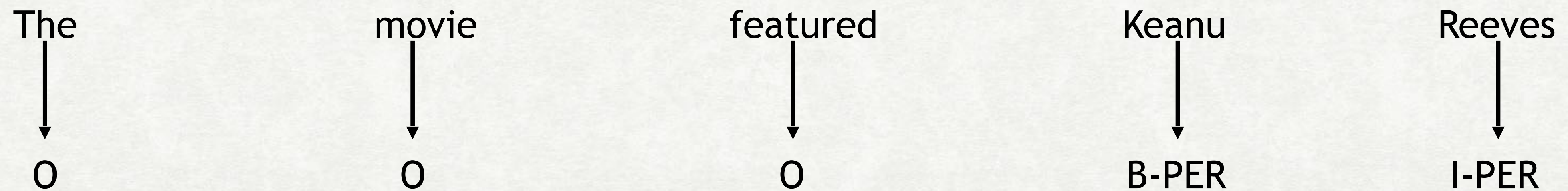
SEQUENCE LABELING

One tag for one word

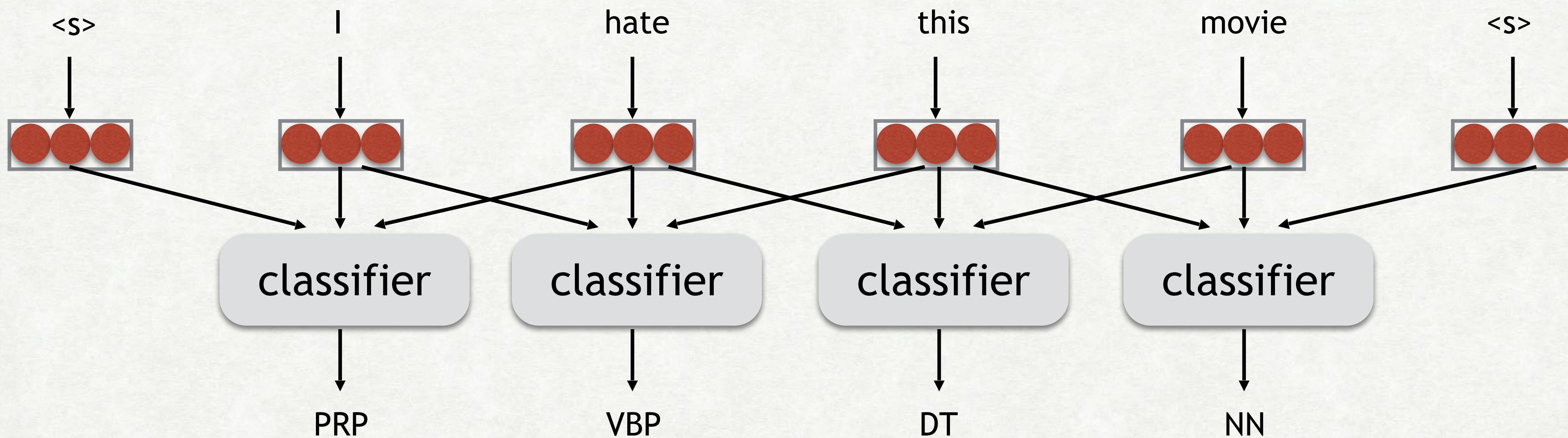
e.g. Part of speech tagging



• e.g. Named entity recognition

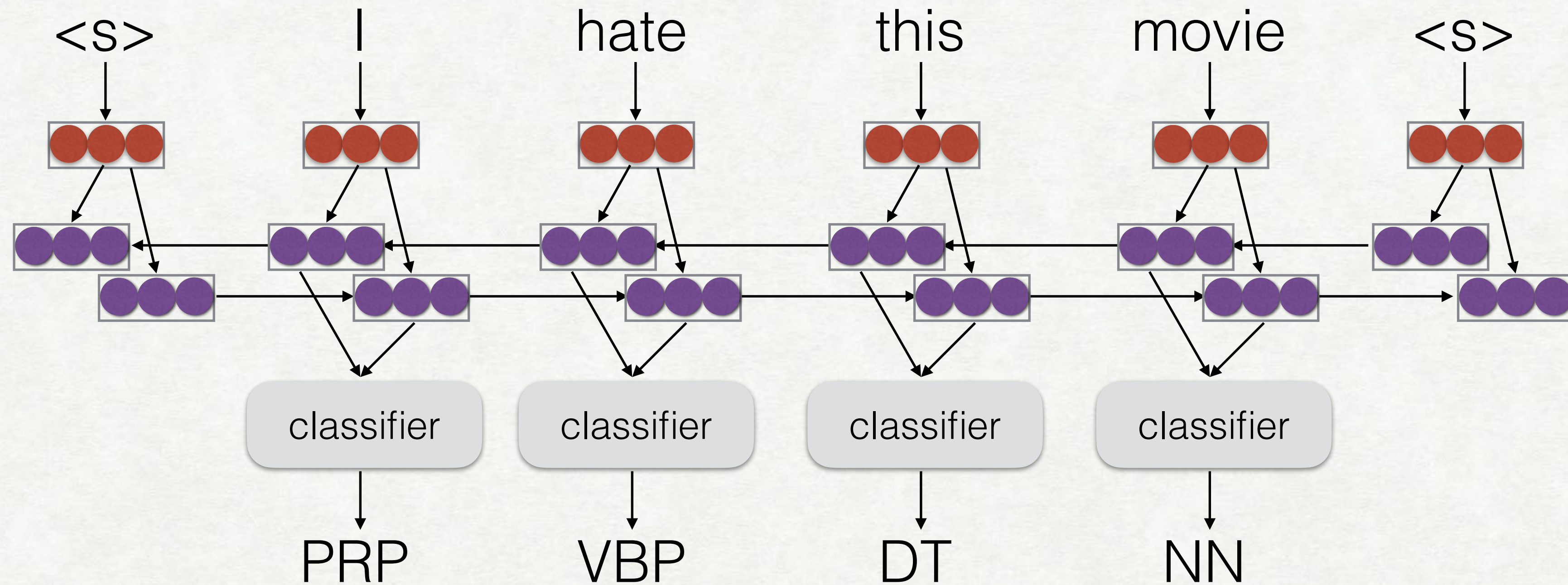


SEQUENCE LABELING AS INDEPENDENT CLASSIFICATION



SEQUENCE LABELING W/ BILSTM

Still not modeling output structure! Outputs are independent



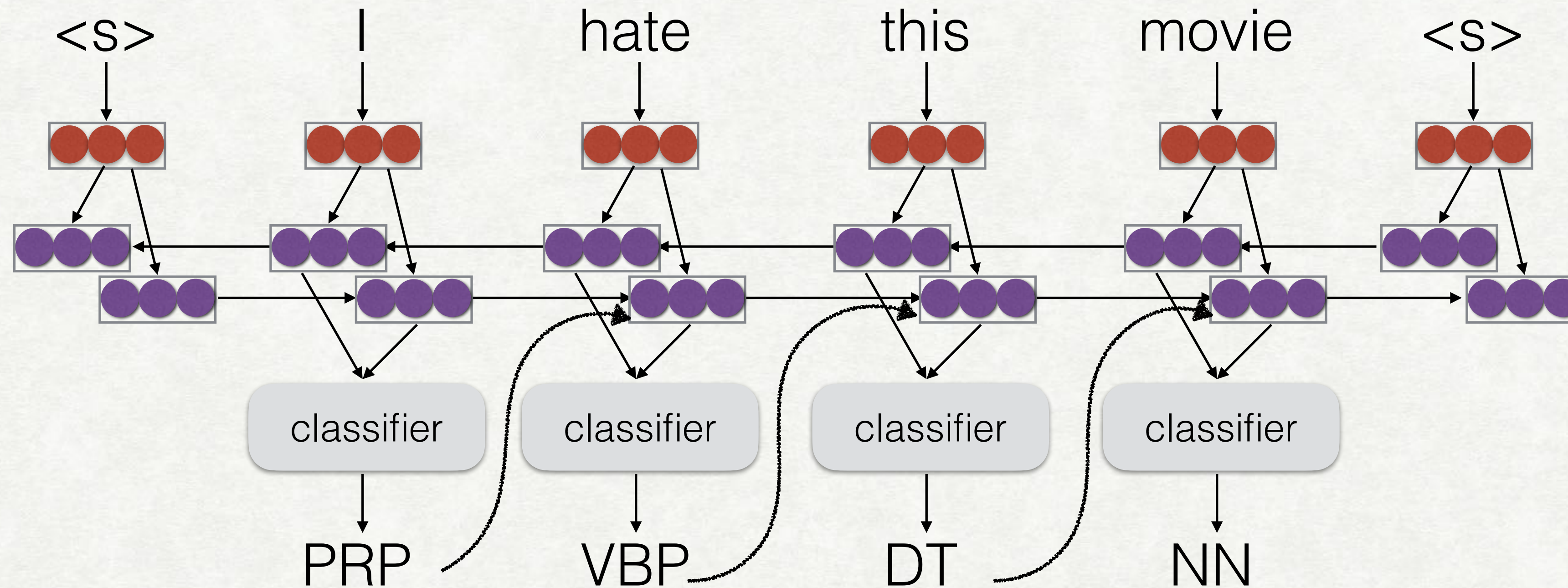
WHY MODEL INTERACTIONS IN OUTPUT?

Consistency is important!

time	flies	like	an	arrow	
N	V	Prep	DT	N	(time moves similarly to an arrow)
N	NS	V	DT	N	("time flies" are fond of arrows)
V	N	Prep	DT	N	(please measure the time of flies similarly to how an arrow would)
		↓	max frequency		
N	NS	Prep	DT	N	("time flies" that are similar to an arrow)

A TAGGER CONSIDERING OUTPUT STRUCTURE

Tags are inter-dependent



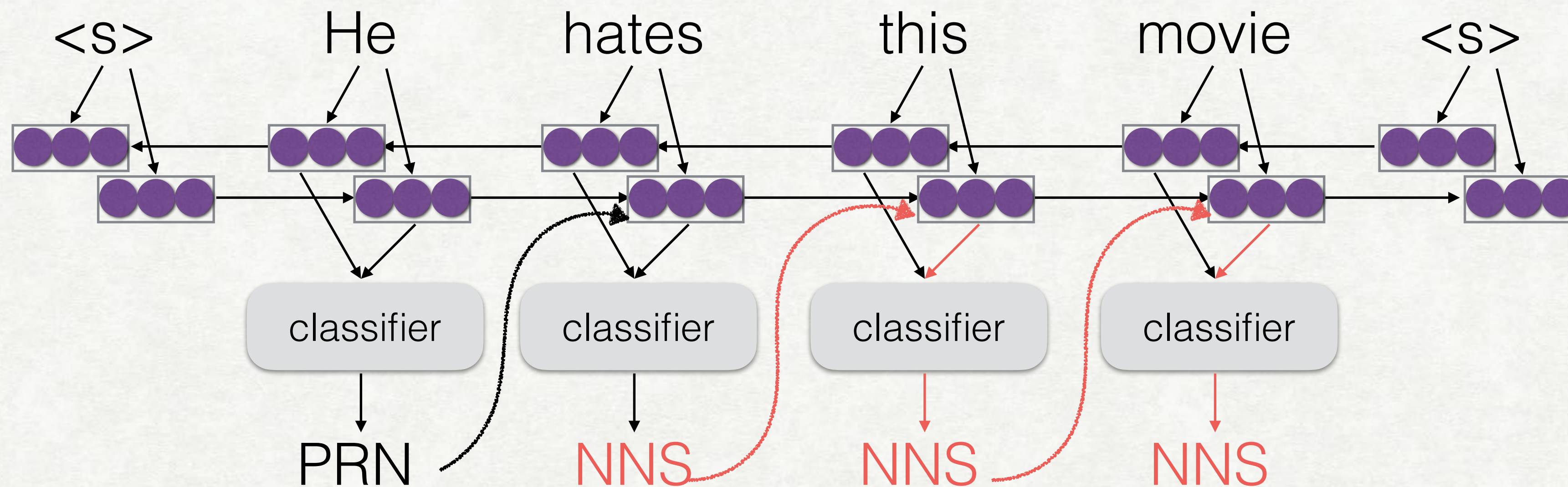
TRAINING STRUCTURED MODELS

Simplest training method "teacher forcing"

Just feed in the *correct* previous tag

TEACHER FORCING AND EXPOSURE BIAS

Teacher forcing assumes feeding correct previous input, but at test time we may make mistakes that propagate



- **Exposure bias:** The model is not exposed to mistakes during training, and cannot deal with them at test

LOCAL NORMALIZATION VS. GLOBAL NORMALIZATION

Locally normalized models: each decision made by the model has a probability that adds to one

$$P(Y | X) = \prod_{j=1}^{|Y|} \frac{e^{S(y_j | X, y_1, \dots, y_{j-1})}}{\sum_{\tilde{y}_j \in V} e^{S(\tilde{y}_j | X, y_1, \dots, y_{j-1})}}$$

Globally normalized models (a.k.a. energy-based models): each sentence has a score, which is not normalized over a particular decision

$$P(Y | X) = \frac{e^{\sum_{j=1}^{|Y|} S(y_j | X, y_1, \dots, y_{j-1})}}{\sum_{\tilde{Y} \in V^*} e^{\sum_{j=1}^{|\tilde{Y}|} S(\tilde{y}_j | X, \tilde{y}_1, \dots, \tilde{y}_{j-1})}}$$

PROBLEMS TRAINING GLOBALLY NORMALIZED MODELS

Problem: the denominator is too big to expand naively

We must do something tricky:

$$P(Y | X) = \frac{e^{\sum_{j=1}^{|Y|} S(y_j | X, y_1, \dots, y_{j-1})}}{\sum_{\tilde{Y} \in V^*} e^{\sum_{j=1}^{|\tilde{Y}|} S(\tilde{y}_j | X, \tilde{y}_1, \dots, \tilde{y}_{j-1})}}$$

Consider only a subset of hypotheses

Design the model so we can efficiently enumerate all hypotheses

STRUCTURED PERCEPTRON

THE STRUCTURED PERCEPTRON ALGORITHM

An extremely simple way of training (non-probabilistic) global models

Find the one-best, and if it's score is better than the correct answer, adjust parameters to fix this

$$\hat{Y} = \operatorname{argmax}_{\tilde{Y} \neq Y} S(\tilde{Y} | X; \theta) \quad \leftarrow \text{Find one best}$$

if $S(\hat{Y} | X; \theta) \geq S(Y | X; \theta)$ **then** \leftarrow If score better than reference

$$\theta \leftarrow \theta + \alpha \left(\frac{\partial S(Y|X;\theta)}{\partial \theta} - \frac{\partial S(\hat{Y}|X;\theta)}{\partial \theta} \right) \quad \leftarrow \text{Increase score of ref, decrease score of one-best (here, SGD update)}$$

end if

STRUCTURED PERCEPTRON LOSS

Structured perceptron can also be expressed as a loss function!

$$\ell_{\text{percept}}(X, Y) = \max(0, S(\hat{Y} | X; \theta) - S(Y | X; \theta))$$

- Resulting **gradient looks like perceptron algorithm**

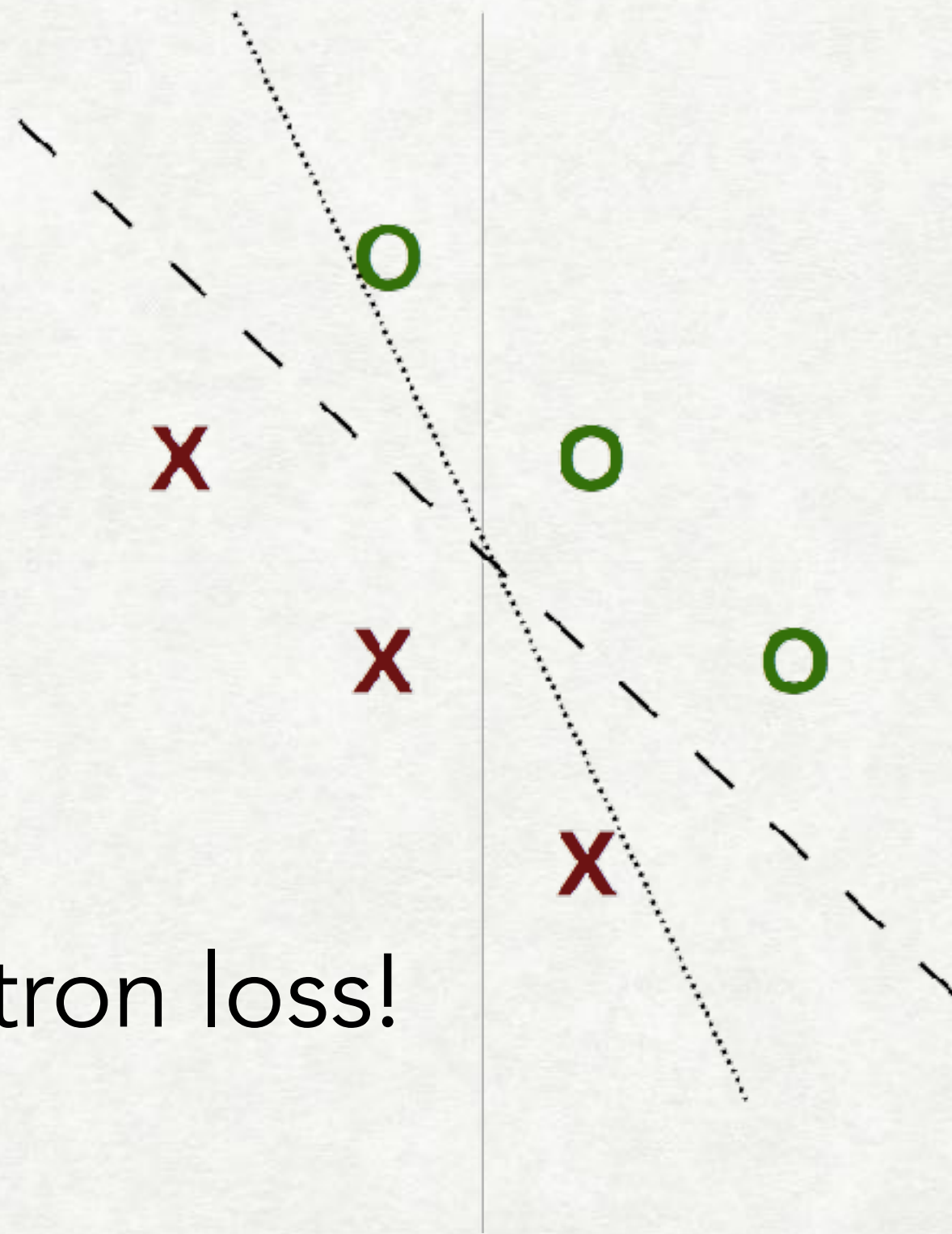
$$\frac{\partial \ell_{\text{percept}}(X, Y; \theta)}{\partial \theta} = \begin{cases} \frac{\partial S(Y|X; \theta)}{\partial \theta} - \frac{\partial S(\hat{Y}|X; \theta)}{\partial \theta} & \text{if } S(\hat{Y} | X; \theta) \geq S(Y | X; \theta) \\ 0 & \text{otherwise} \end{cases}$$

- This is a normal loss function, **can be used in NNs**
- But! Requires finding the argmax in addition to the true candidate: must **do prediction during training**

HINGE LOSS AND COST-SENSITIVE TRAINING

PERCEPTRON AND UNCERTAINTY

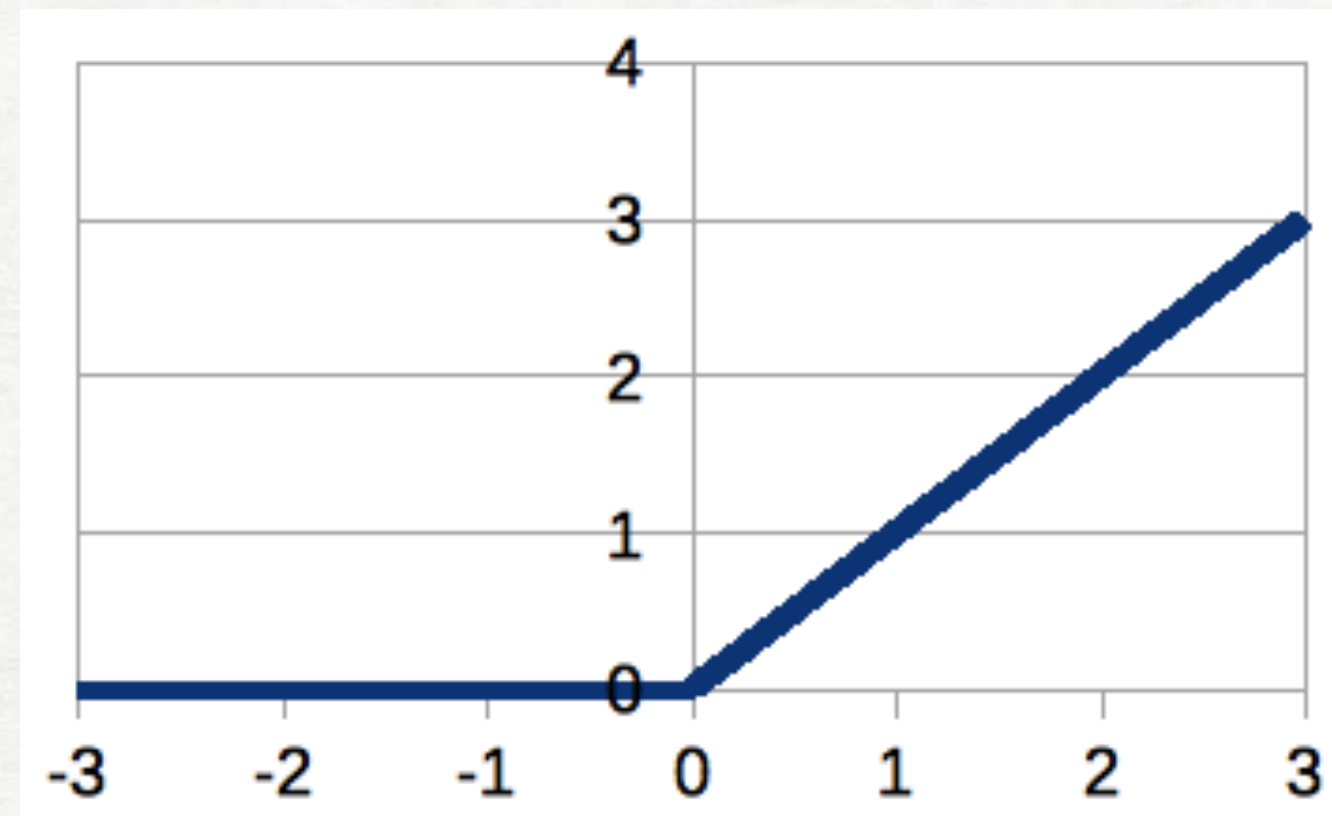
Which is better, dotted or dashed?



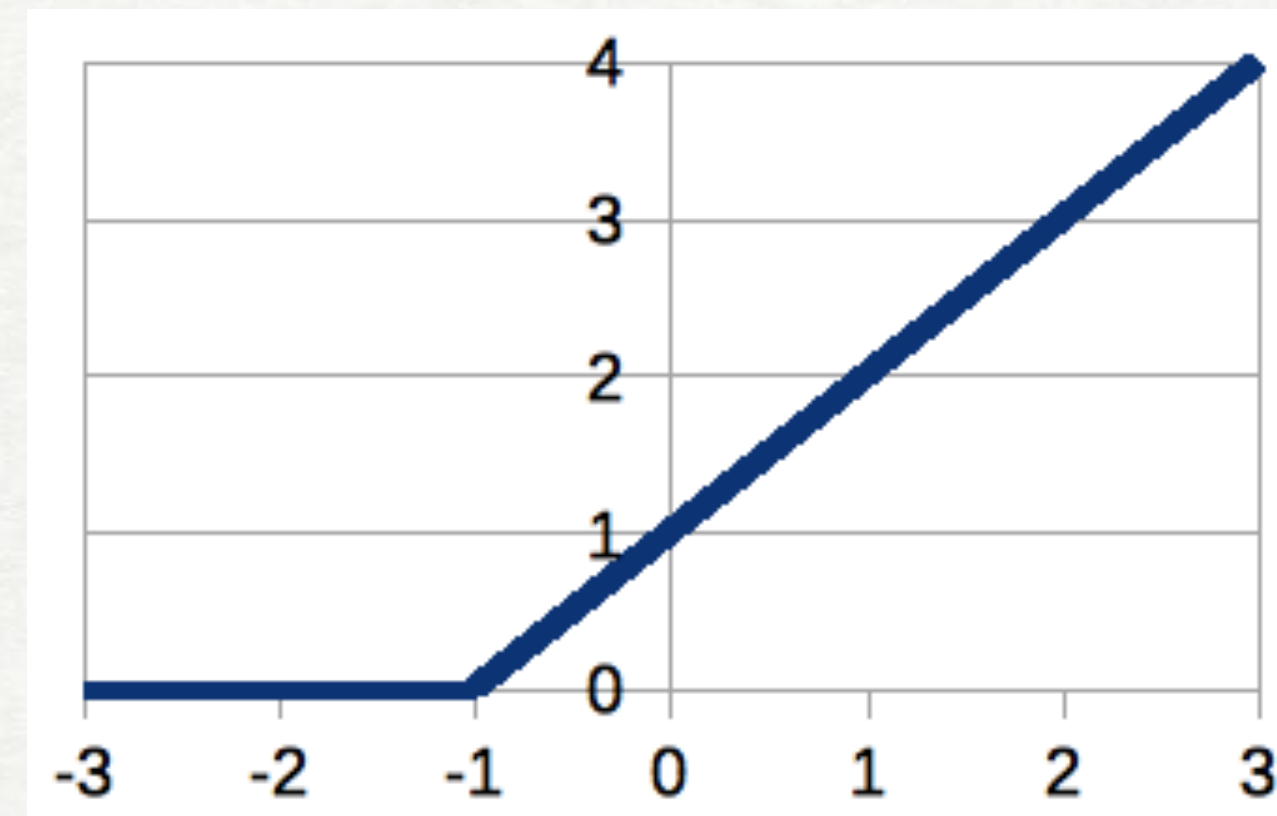
Both have zero perceptron loss!

ADDING A "MARGIN" WITH HINGE LOSS

Penalize when incorrect answer is within margin m



Perceptron

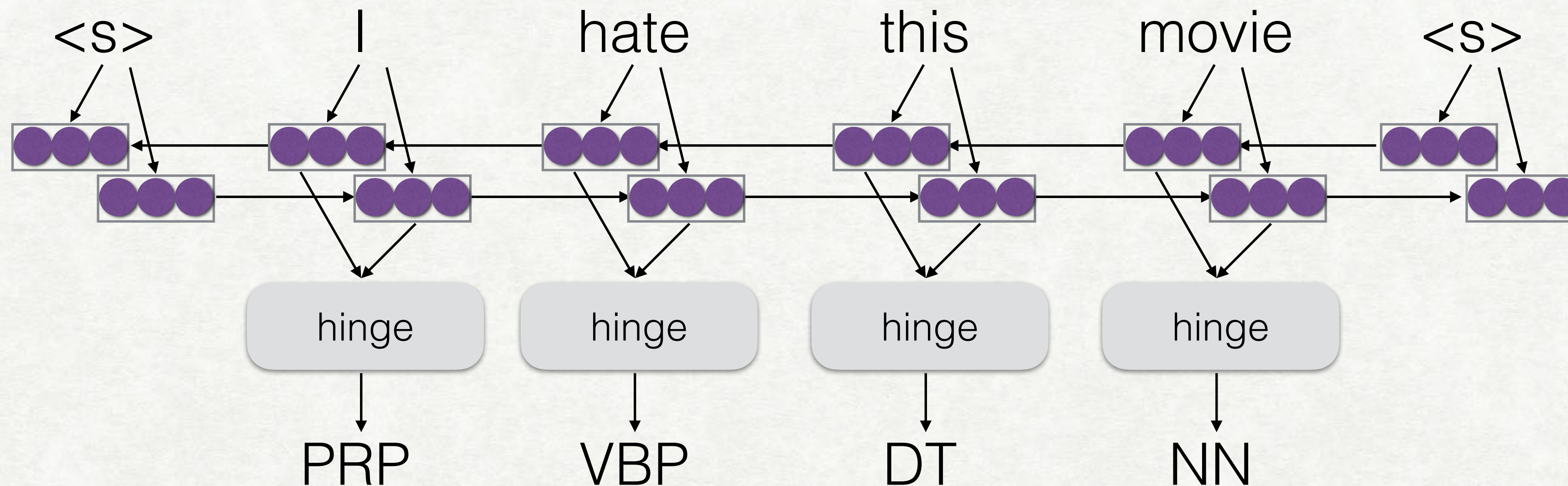


Hinge

$$\ell_{\text{hinge}}(x, y; \theta) = \max(0, m + S(\hat{y} | x; \theta) - S(y | x; \theta))$$

HINGE LOSS FOR ANY CLASSIFIER!

We can swap cross-entropy for hinge loss anytime



```
loss = dy.pickneglogsoftmax(score, answer)
```

```
↓  
loss = dy.hinge(score, answer, m=1)
```


LOCALLY-DEPENDENT MODELS

PROBLEMS

- Independent classification models

- Strong independent assumption

$$P(Y|X) = \prod_{i=1}^L P(y_i|X)$$

- No guarantee of valid (consistent) structured outputs
 - BIO tagging scheme in NER

- Locally normalized models (e.g. history-based RNN, seq2seq)

- Prior order

$$P(Y|X) = \prod_{i=1}^L P(y_i|X, y_{<i})$$

- Approximating decoding
 - Greedy search
 - Beam search
- Label bias

MODELS W/ LOCAL DEPENDENCIES

Some independence assumptions, but not entirely independent (local dependencies)

Exact and optimal decoding/training via dynamic programs

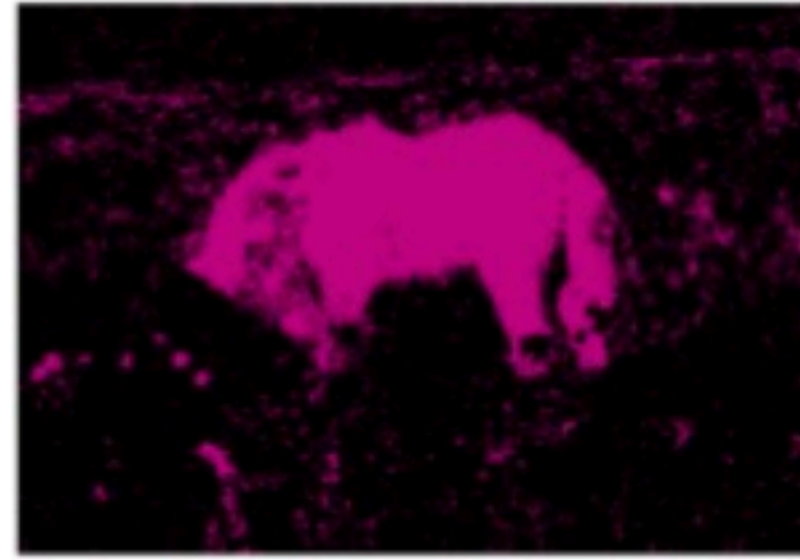
Conditional Random Fields! (CRFs)



LOCAL VS LOCALLY NORMALIZED



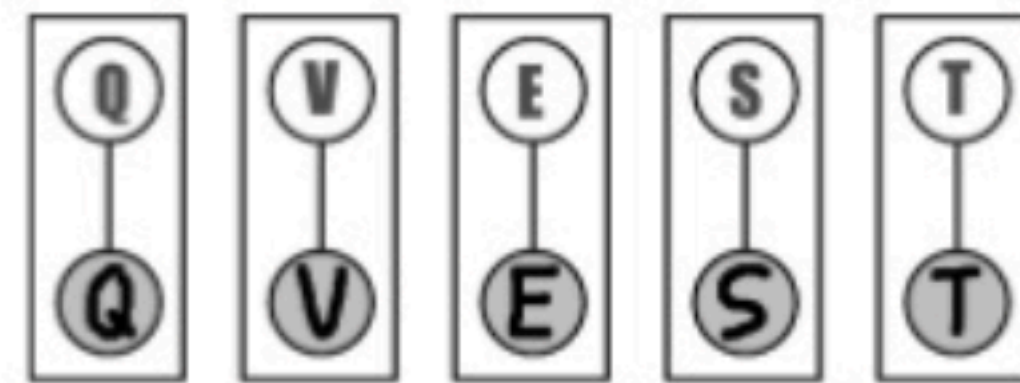
original



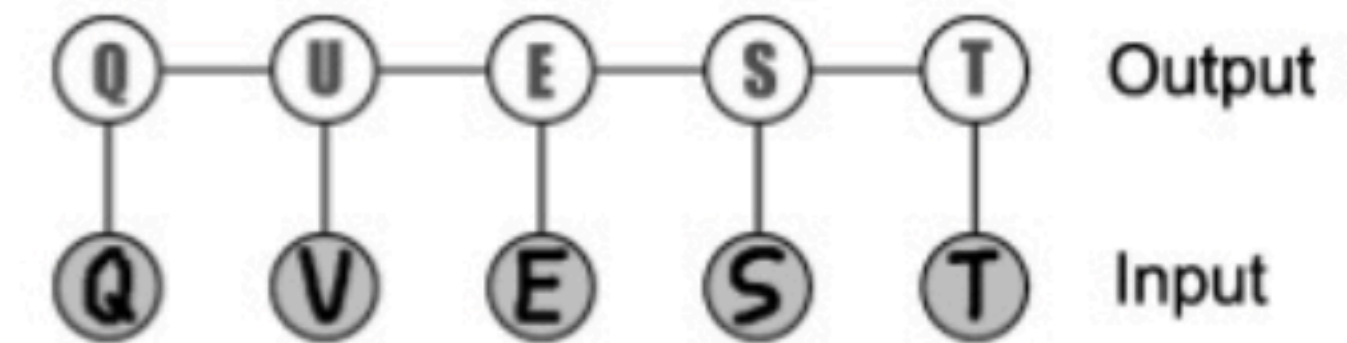
local classification



local + smoothness



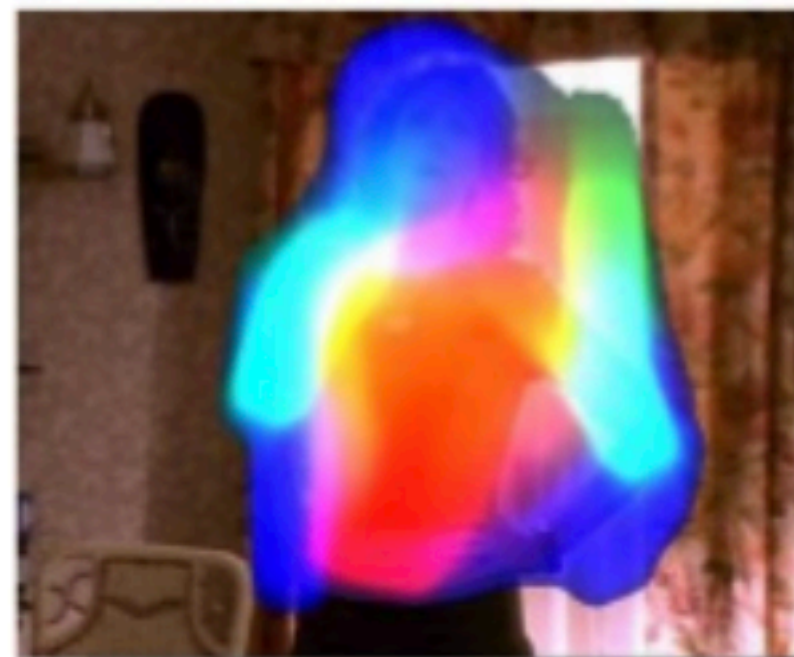
local classification



local + "correction"



original



local classification



local + geometry

REMINDER: GLOBALLY NORMALIZED MODELS

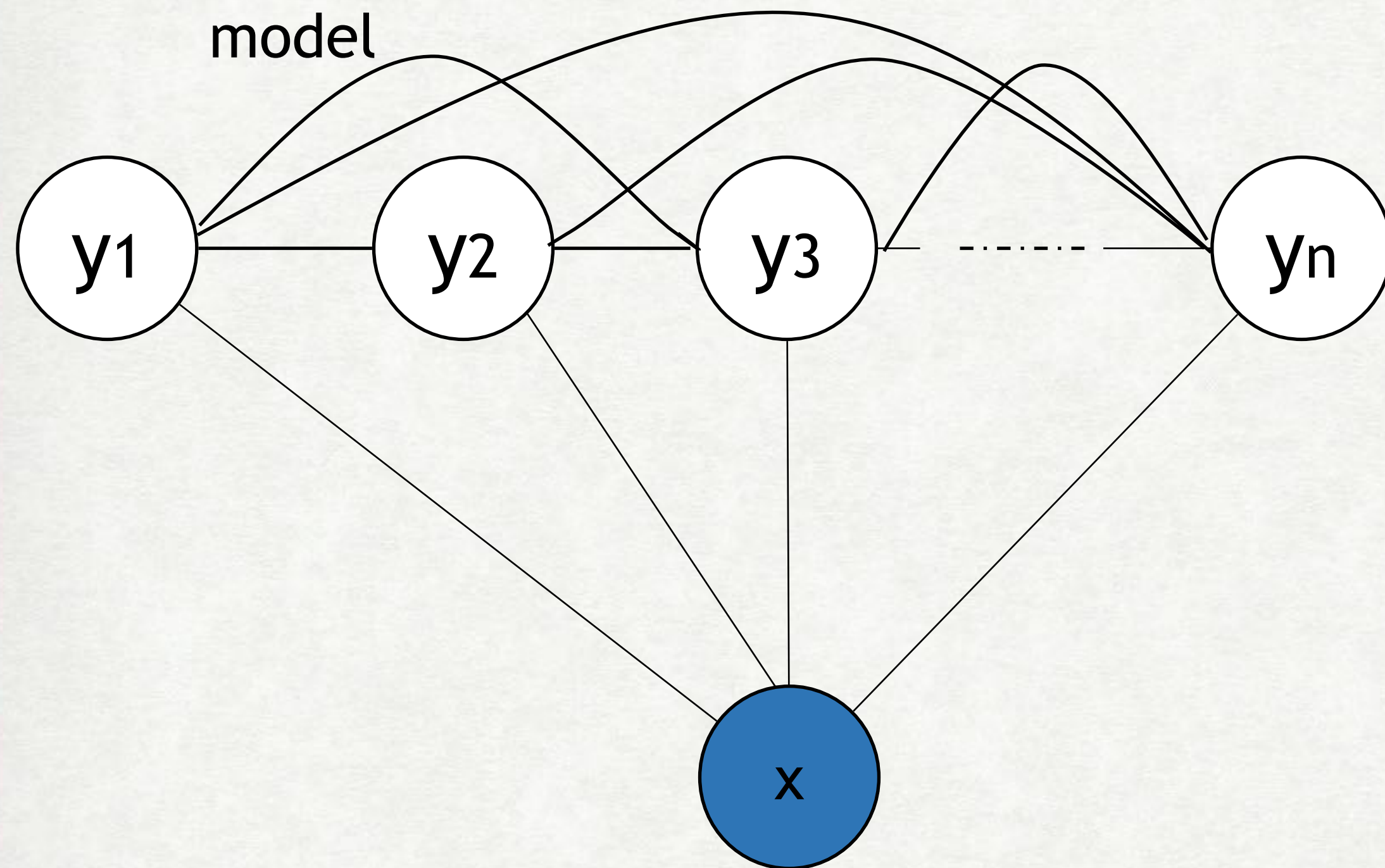
- Each output sequence has a score, which is not normalized over a particular decision

$$P(Y|X) = \frac{\exp(S(Y, X))}{\sum_{Y'} \exp(S(Y', X))} = \frac{\psi(Y, X)}{\sum_{Y'} \psi(Y', X)}$$

where $\psi(Y, X)$ are potential functions.

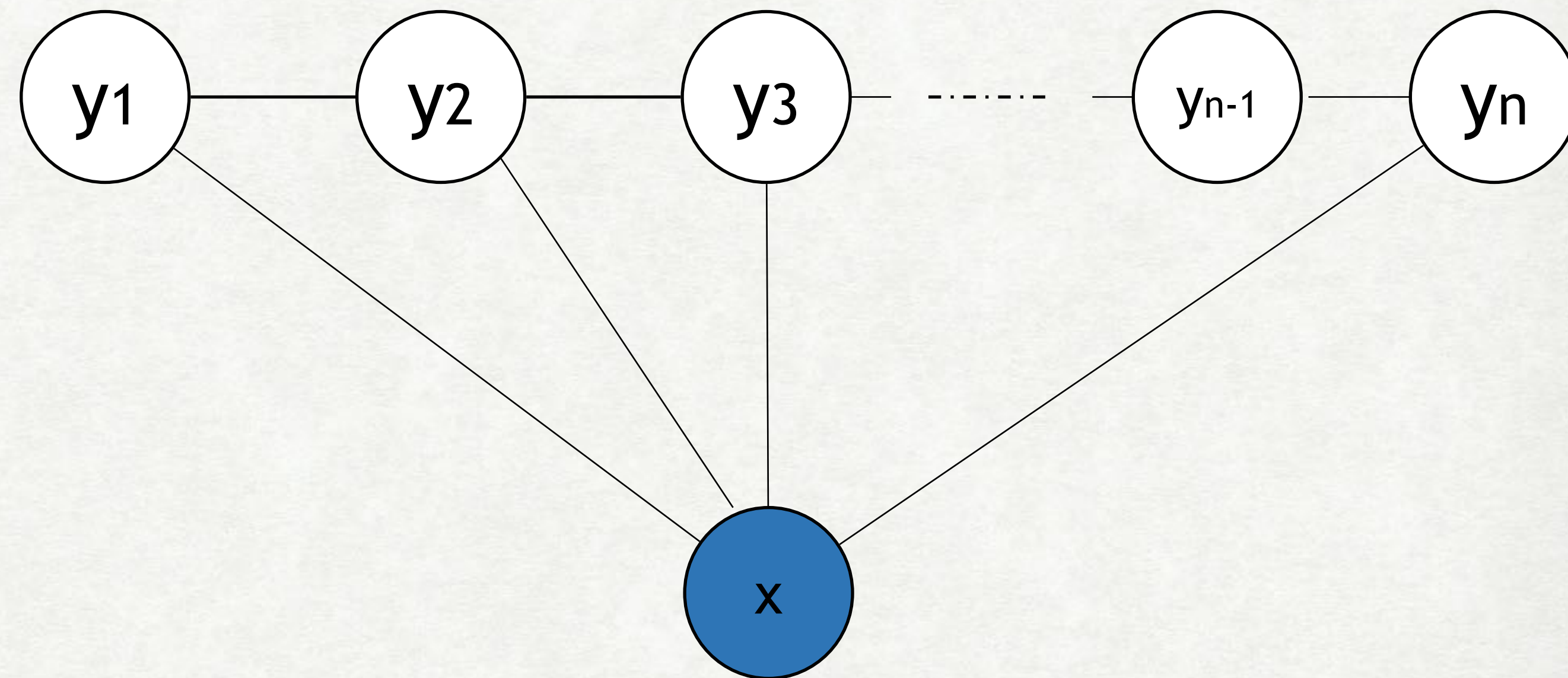
CONDITIONAL RANDOM FIELDS

General form of globally normalized model



$$P(Y|X) = \frac{\psi(Y, X)}{\sum_{Y'} \psi(Y', X)}$$

First-order linear CRF



$$P(Y|X) = \frac{\prod_{i=1}^L \psi_i(y_{i-1}, y_i, X)}{\sum_{Y'} \prod_{i=1}^L \psi_i(y'_{i-1}, y'_i, X)}$$

POTENTIAL FUNCTIONS

$$\psi_i(y_{i-1}, y_i, X) = t(y_{i-1}, y_i, X) \times e(y_i, X)$$

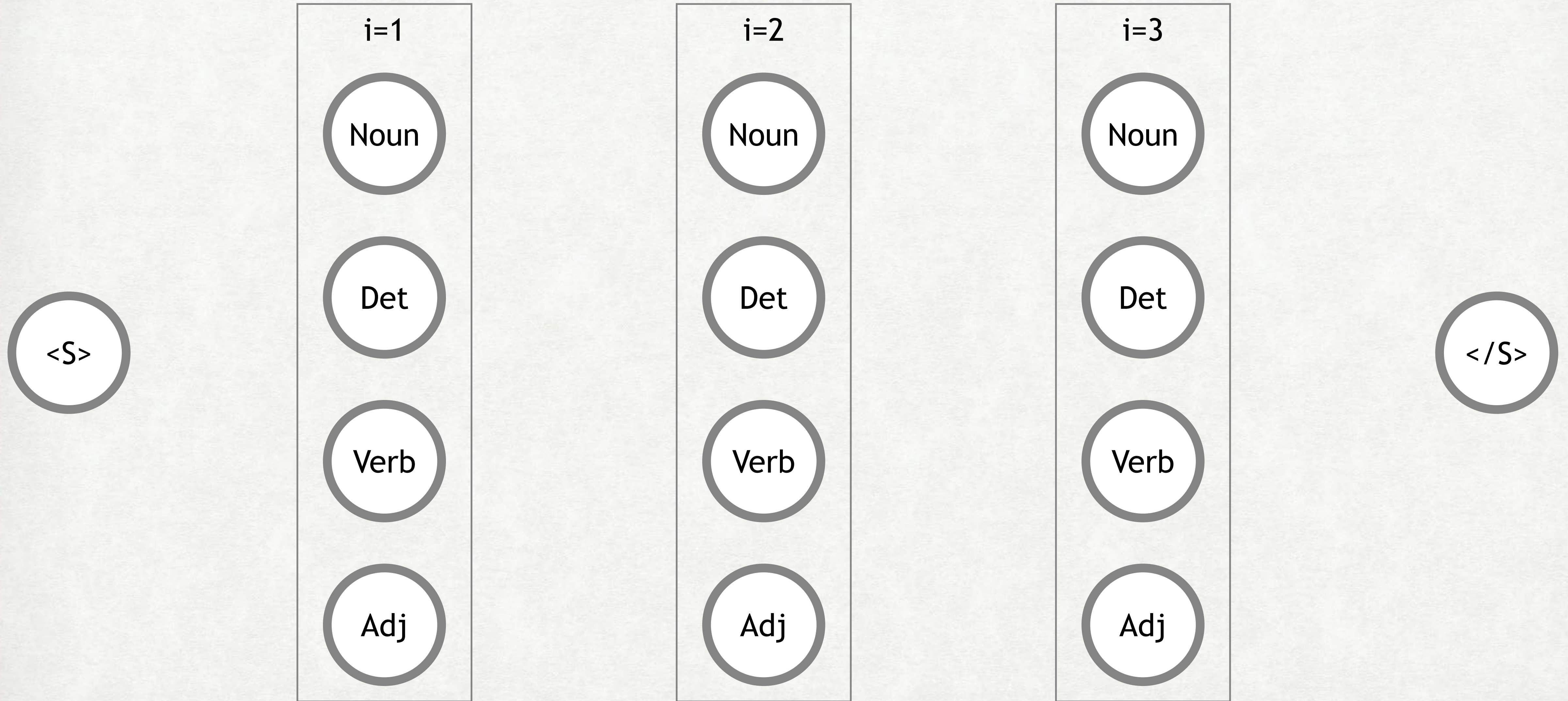
“Transition” “Emission”

The diagram shows the equation $\psi_i(y_{i-1}, y_i, X) = t(y_{i-1}, y_i, X) \times e(y_i, X)$. The term $t(y_{i-1}, y_i, X)$ is enclosed in an orange rounded rectangle, and the term $e(y_i, X)$ is enclosed in a green rounded rectangle. An arrow points from the text "Transition" below to the orange box, and another arrow points from the text "Emission" below to the green box.

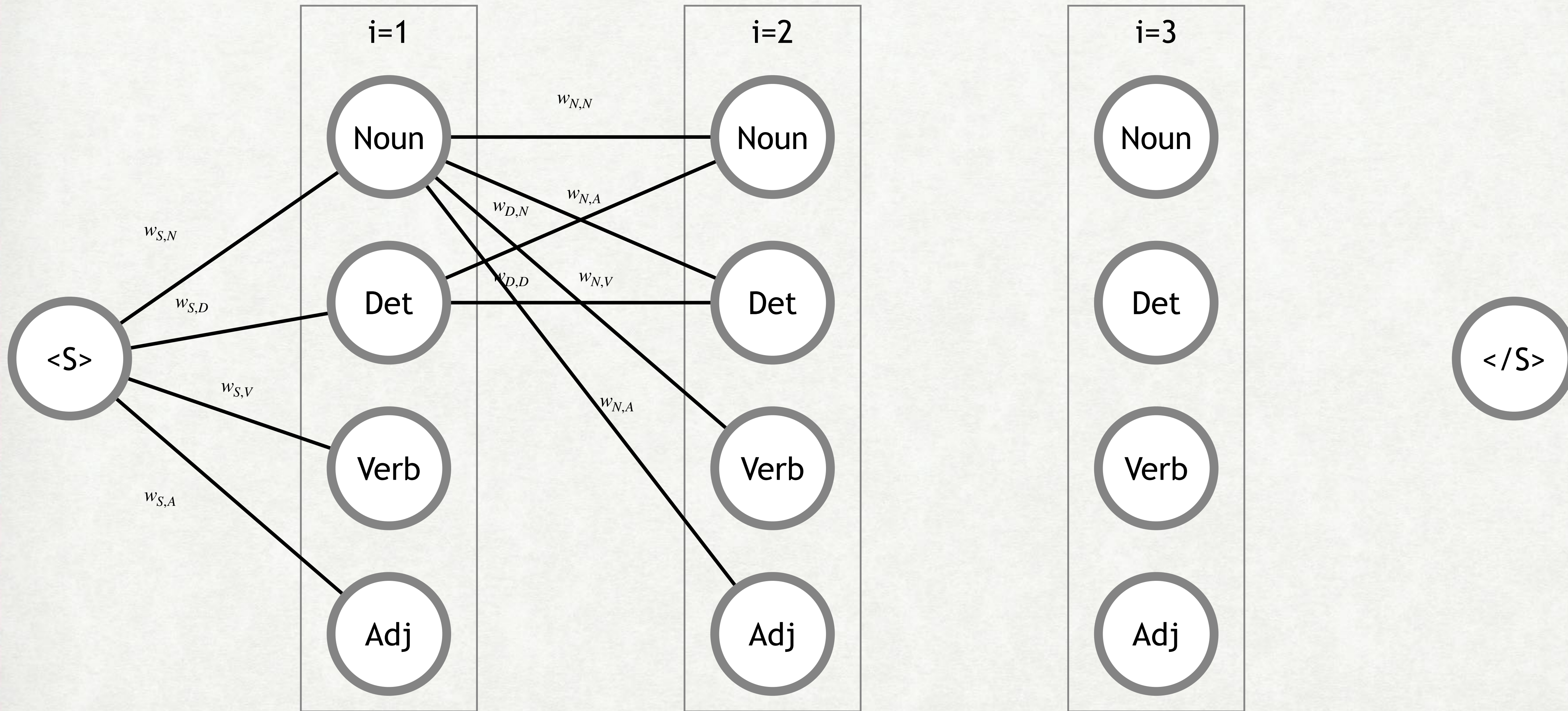
Simpler Version

$$\psi_i(y_{i-1}, y_i, X) = t(y_{i-1}, y_i) \times e(y_i, X)$$

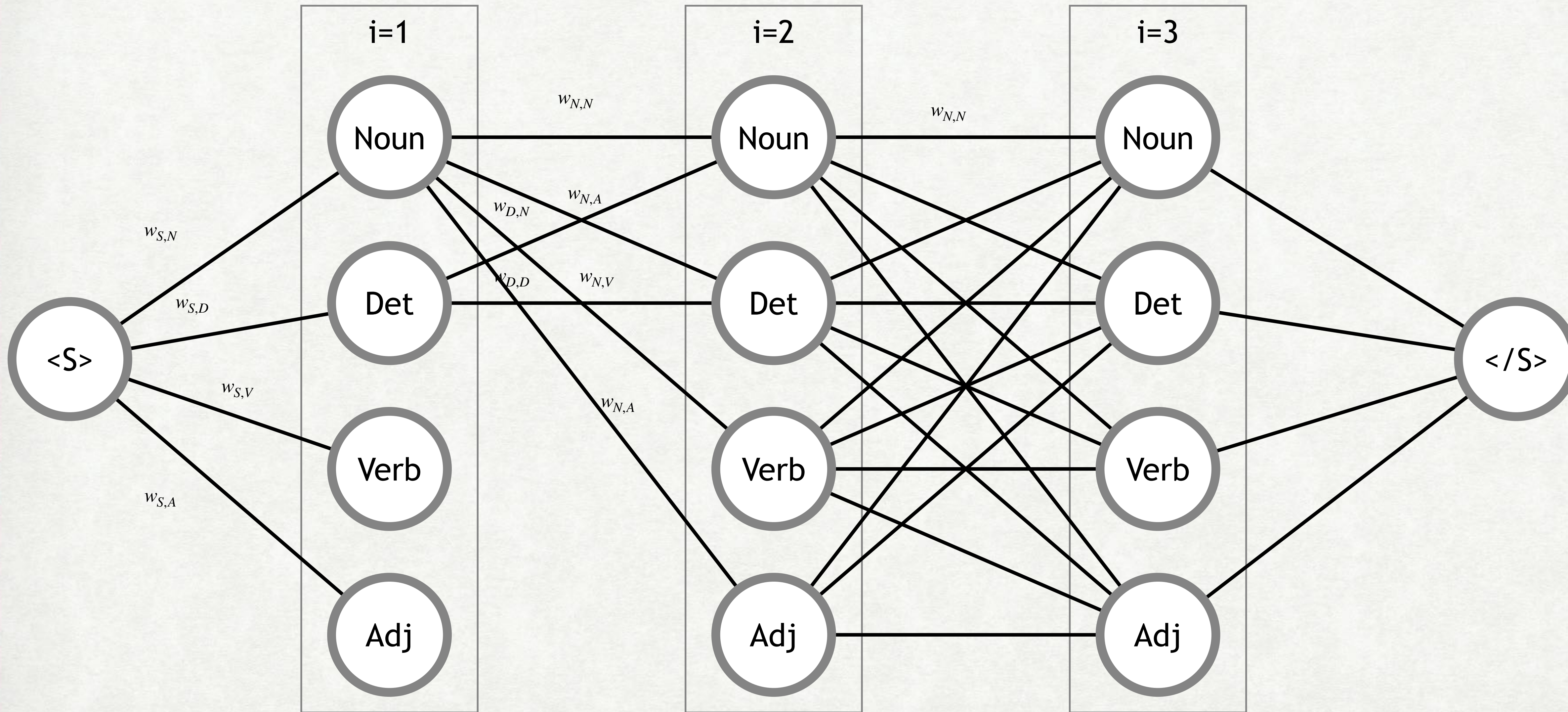
EXAMPLE



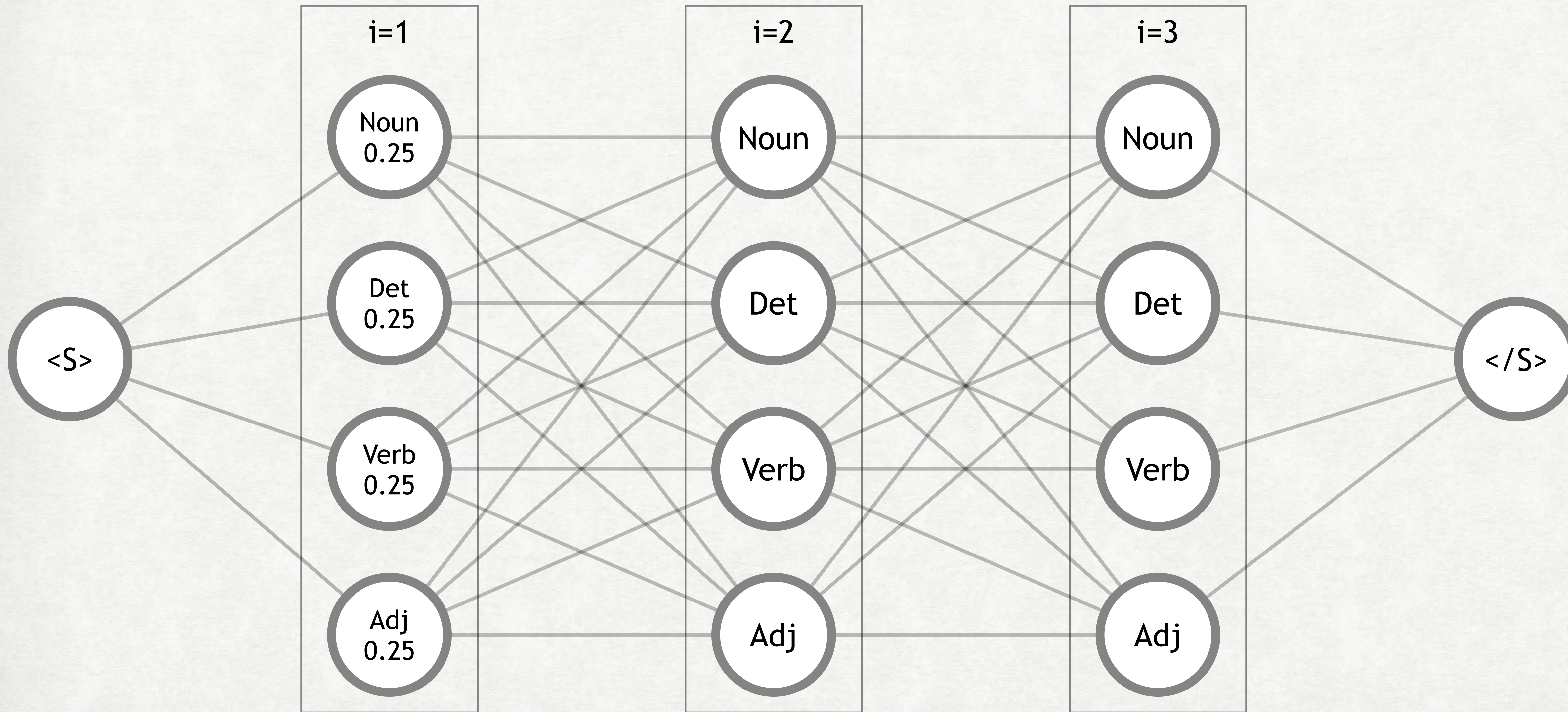
TRANSITION PARAMETERS



TRANSITION PARAMETERS

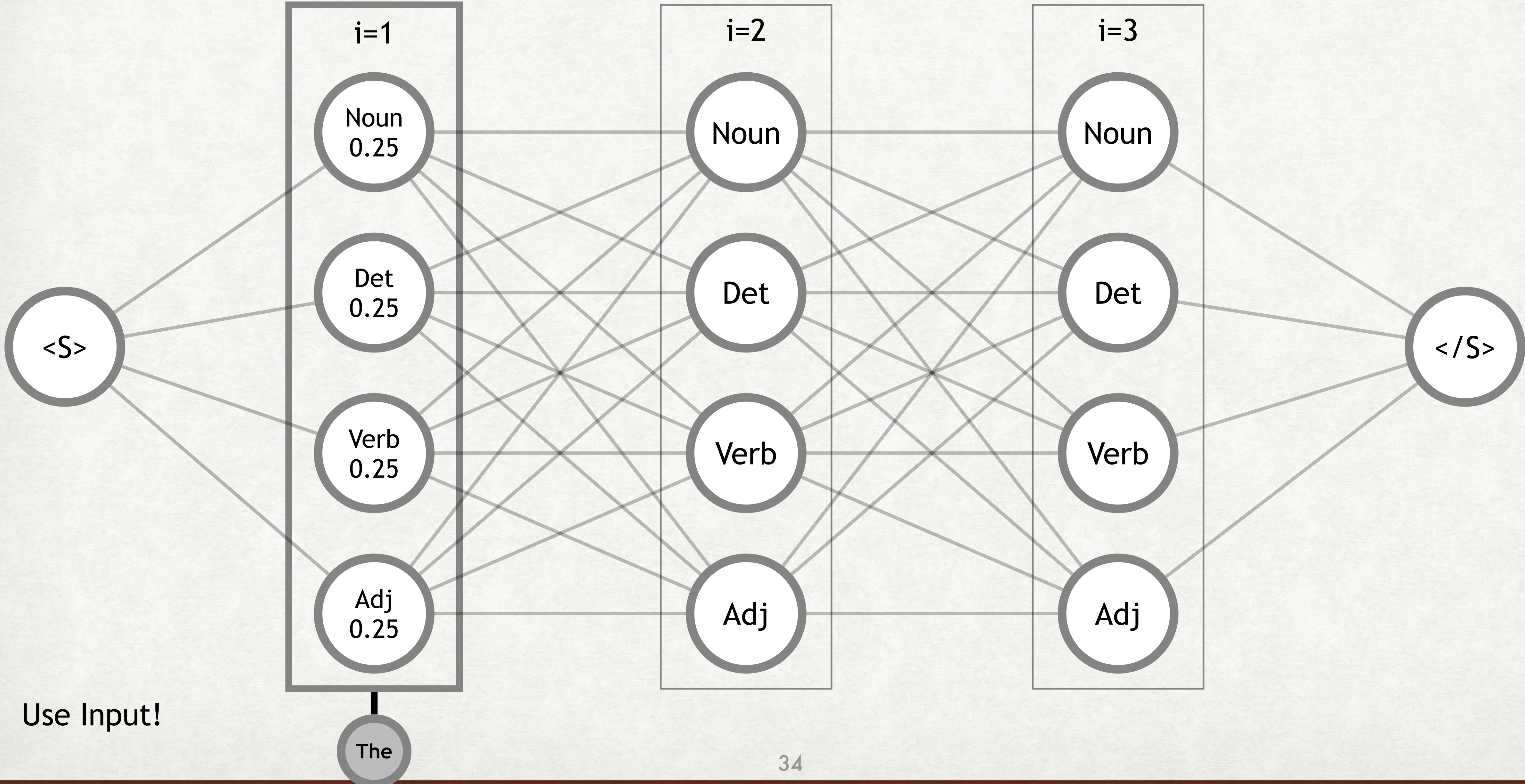


EMISSION PROBABILITIES



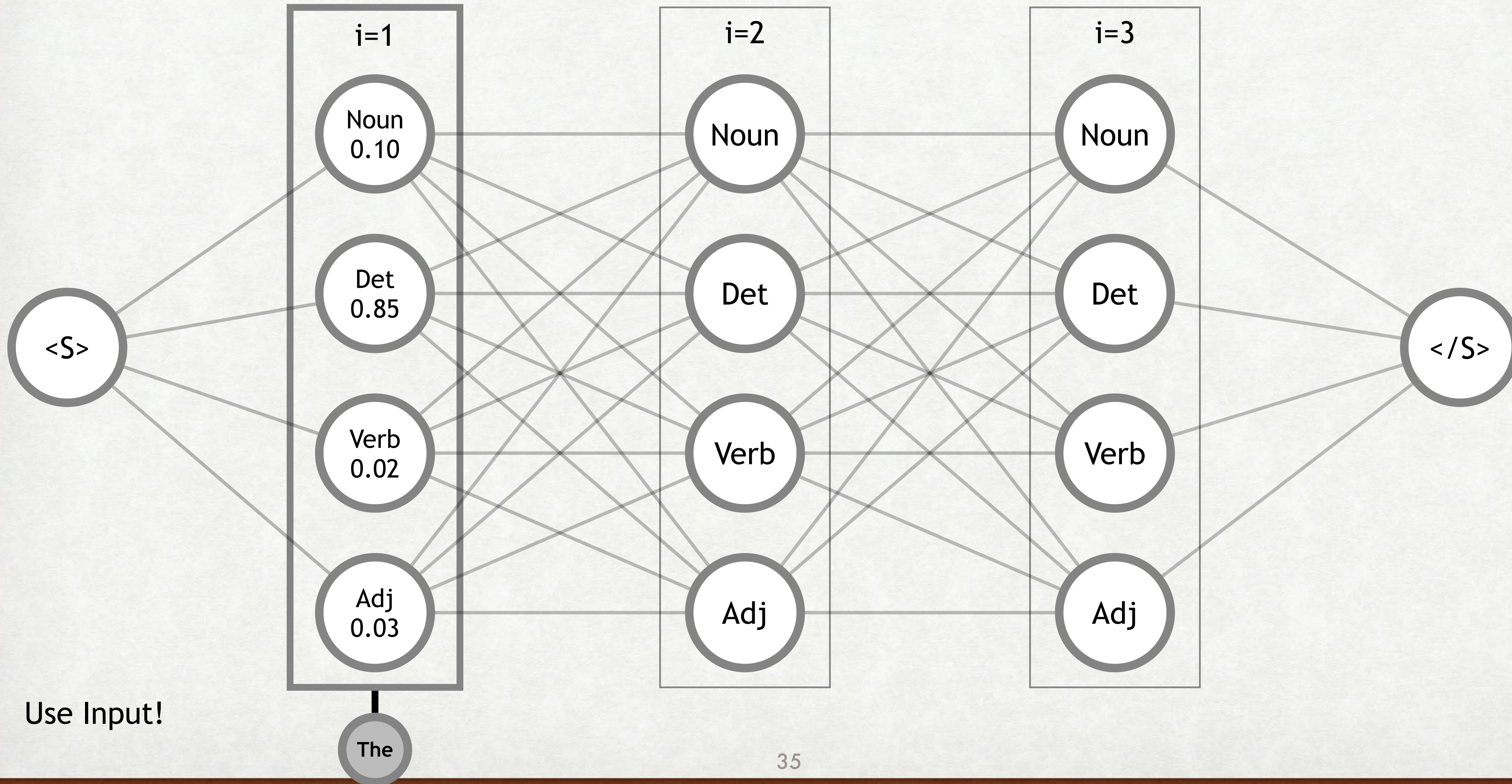
Equally likely?

EMISSION PROBABILITIES

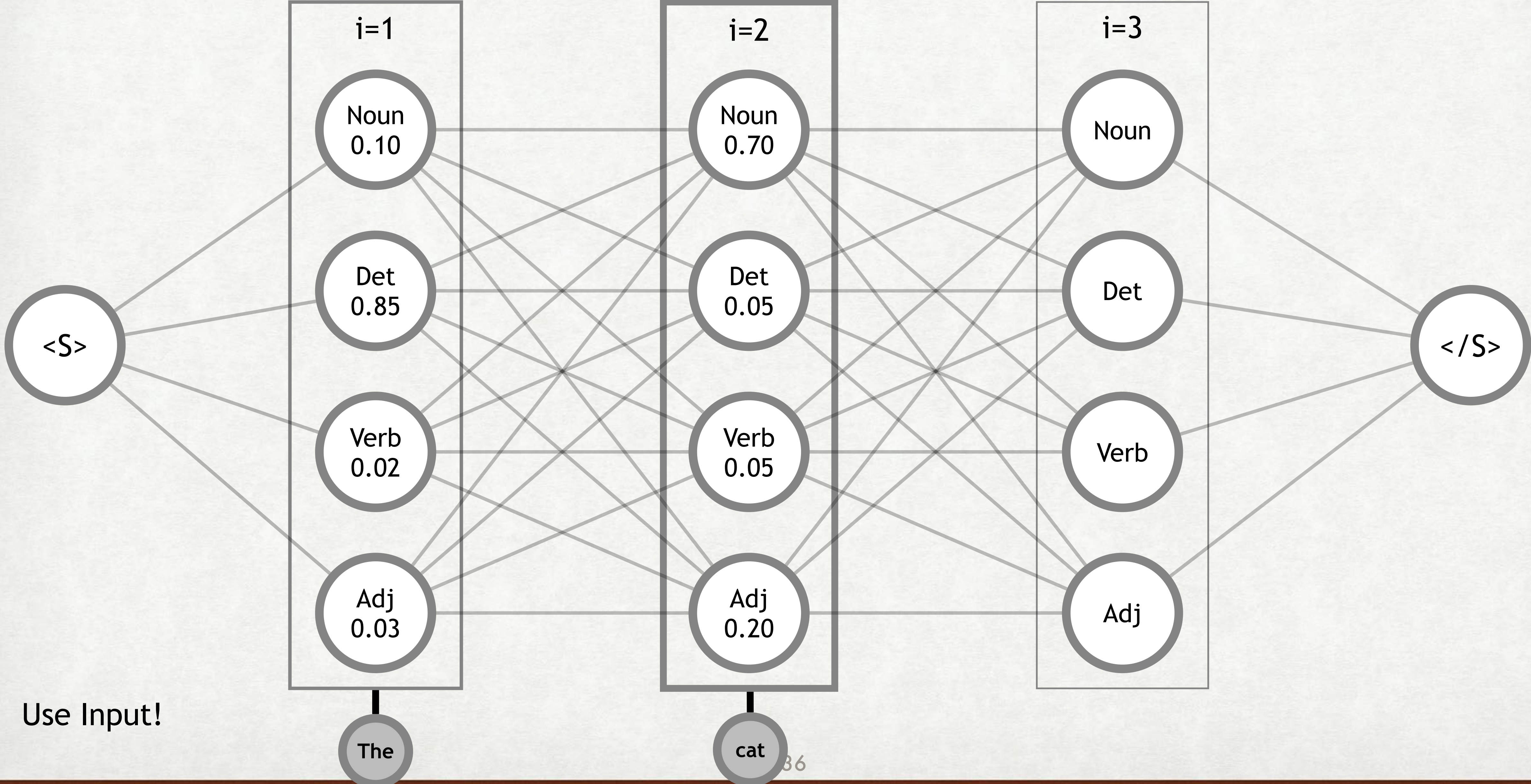


Use Input!

EMISSION PROBABILITIES

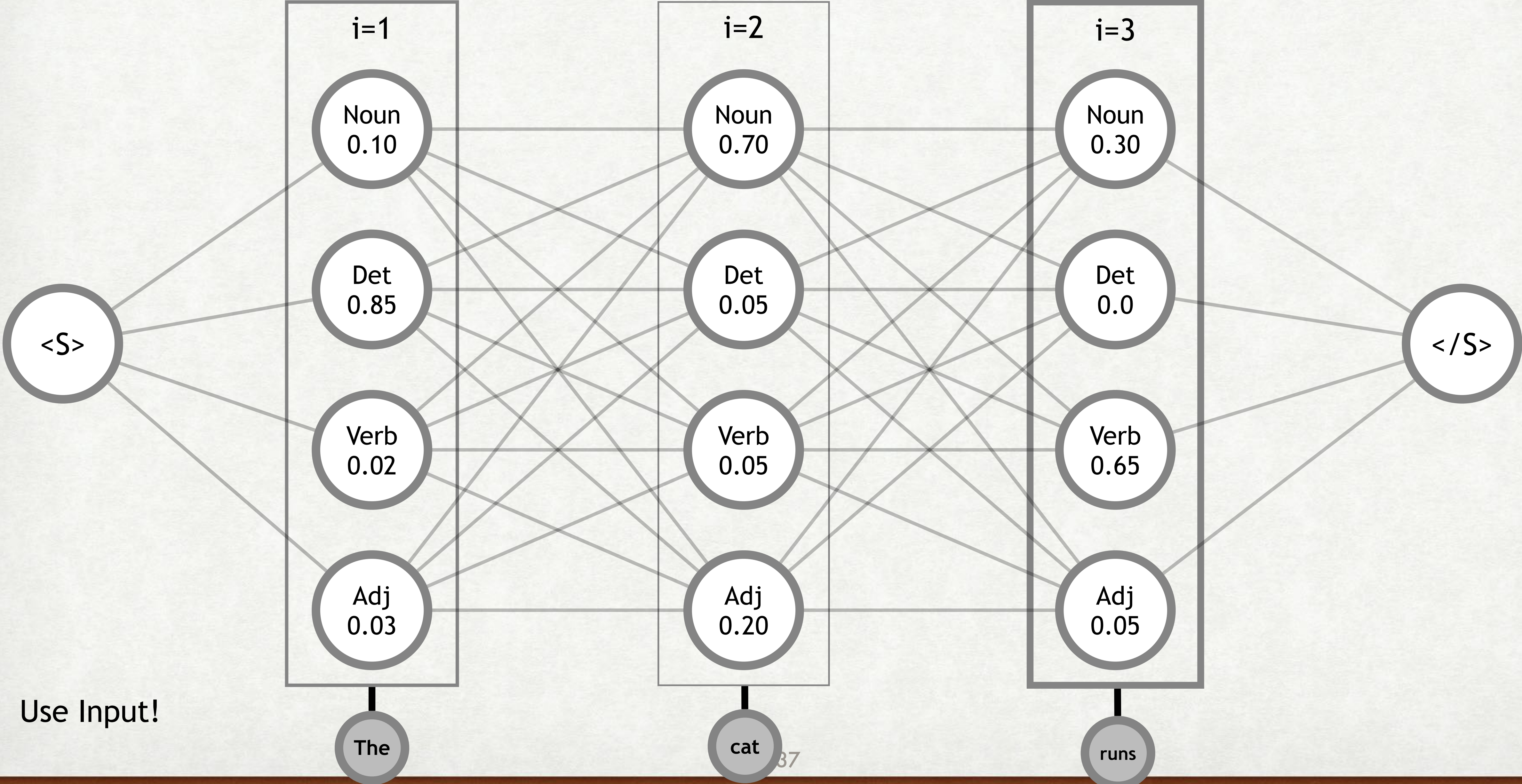


EMISSION PROBABILITIES



Use Input!

EMISSION PROBABILITIES



TRAINING

We want to maximize the probability of the correct output sequence

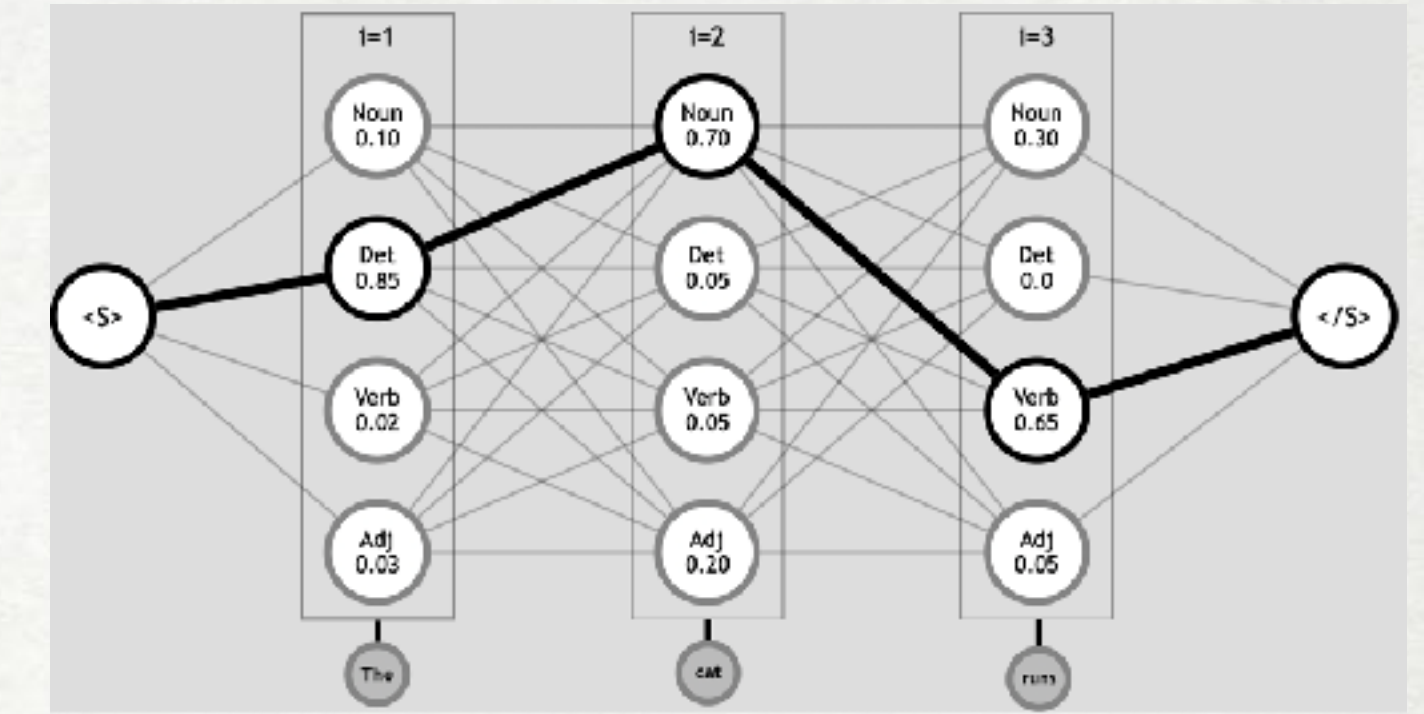
$$\bullet P(Y|X) = \frac{\prod_{i=1}^L \psi_i(y_{i-1}, y_i, X)}{\sum_{Y'} \prod_{i=1}^L \psi_i(y'_{i-1}, y'_i, X)} = \frac{\prod_{i=1}^L \psi_i(y_{i-1}, y_i, X)}{Z(X)}$$

- Training: computing the partition function $Z(X)$

$$Z(X) = \sum_Y \prod_{i=1}^L \psi_i(y_{i-1}, y_i, X)$$

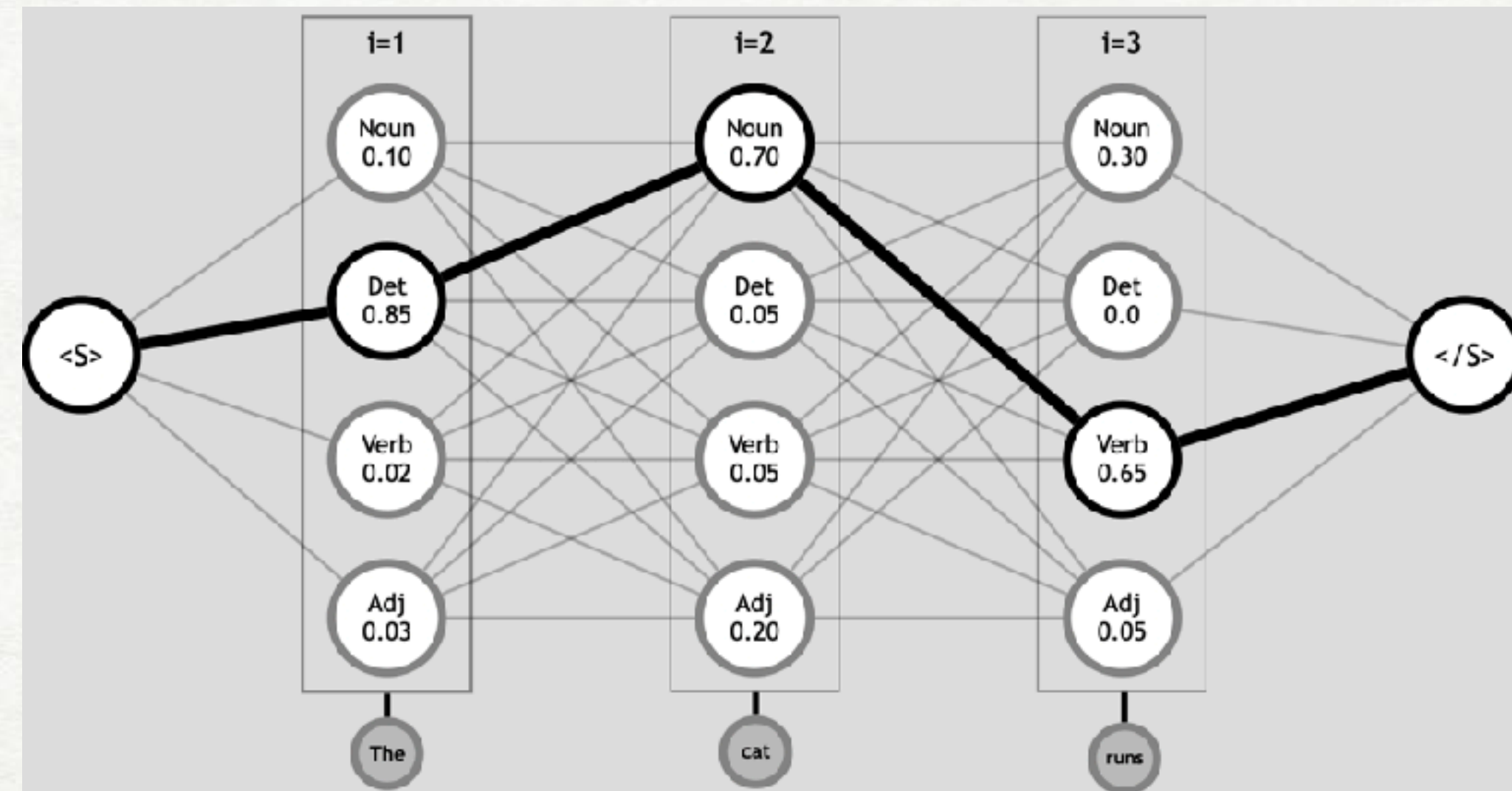
- Decoding

$$y^* = \operatorname{argmax}_Y P(Y|X)$$



TRAINING

We want to maximize the probability of the correct output sequence



Traditionally:

1. Extract features from the input words/context/sentence
2. Train with features as input, target sequences as desired output
3. Use some optimization technique to weight the feature importance for each position

A bunch of algorithms are nicely implemented in scikit-learn:

<https://sklearn-crfsuite.readthedocs.io/en/latest/>

TRY IT!

DOWNLOAD THE NOTEBOOK FROM THE CLASS
WEBSITE

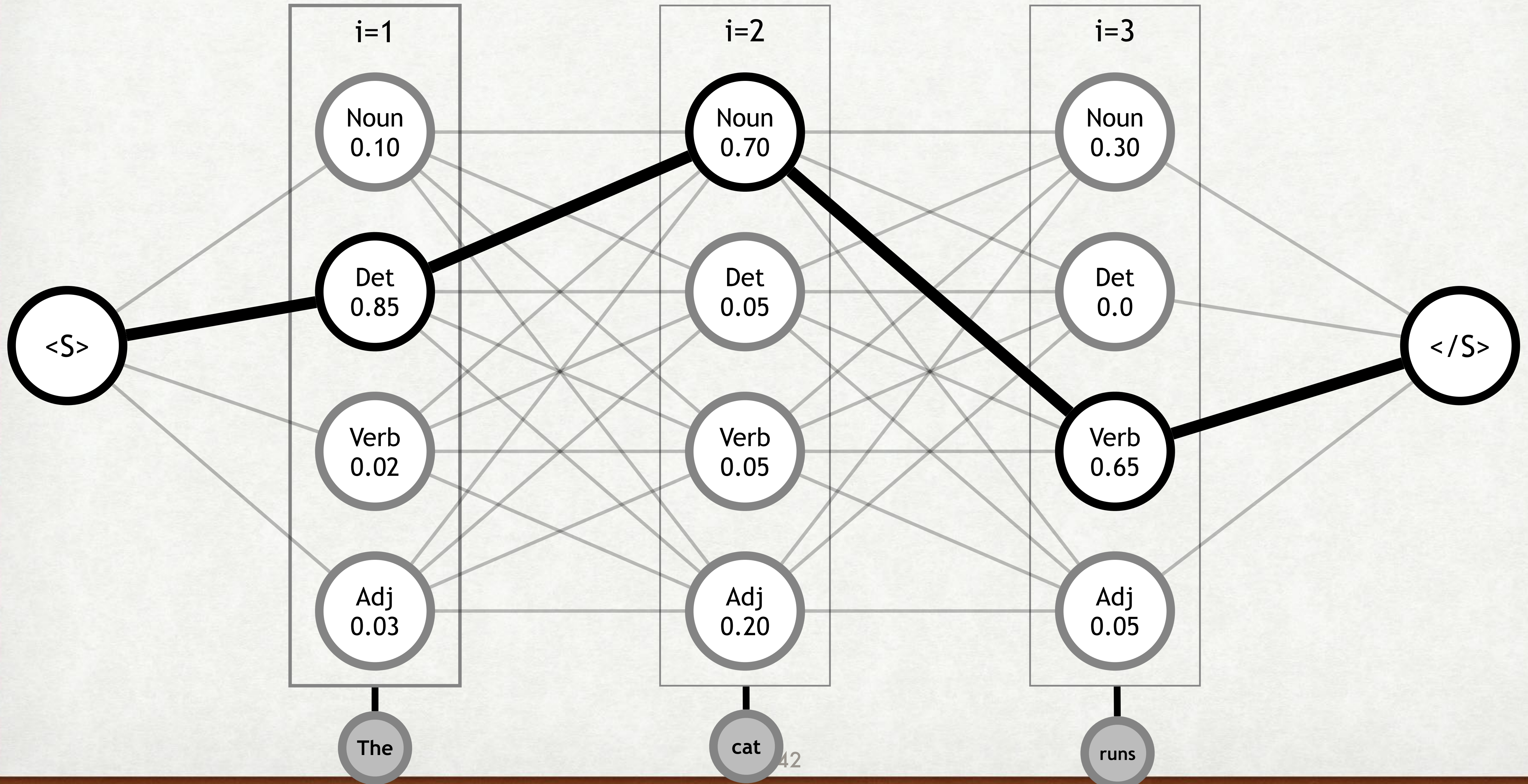
Break-Out Room Exercise [around 20 minutes]:

https://docs.google.com/document/d/1ifTqeqmK6cG2Zk-f5kDMvU_baNh3x-nbxQXz1d549HY/edit?usp=sharing

TRAINING & DECODING OF CRF: VITERBI/FORWARD BACKWARD ALGORITHM

TRAINING

$$P(Y|X) = \frac{\prod_{i=1}^L \psi_i(y_{i-1}, y_i, X)}{\sum_{Y'} \prod_{i=1}^L \psi_i(y'_{i-1}, y'_i, X)} = \frac{\prod_{i=1}^L \psi_i(y_{i-1}, y_i, X)}{Z(X)}$$



INTERACTIONS

- each label depends on the input and nearby labels
 - but, given *adjacent* labels, the others do not matter!
 - If we knew the score of every sequence y_1, y_2, \dots, y_{n-1}
we could easily compute the score of every sequence $y_1, y_2, \dots, y_{n-1}, y_n$
 - So, we only really need to know the score of all the sequences ending in each y_{n-1}
- (Think of that as a “*pre-calculation*” that happens before we think about y_n)

STEP 1: INITIAL

First, calculate transition from $\langle S \rangle$ and emission of the first word for every POS

natural



STEPS: MIDDLE PARTS

For middle words, calculate the scores for all possible previous POS tags

natural language

1:NN → 2:NN

1:JJ → 2:JJ

1:VB → 2:VB

1:LRB → 2:LRB

1:RRB → 2:RRB

...

...

$\text{score}["2 \text{ NN}"] = \log_sum_exp(\text{score}["1 \text{ NN}"] + T(\text{NN}|\text{NN}) + S(\text{language} | \text{NN}),$
 $\text{score}["1 \text{ JJ}"] + T(\text{NN}|\text{JJ}) + S(\text{language} | \text{NN}),$
 $\text{score}["1 \text{ VB}"] + T(\text{NN}|\text{VB}) + S(\text{language} | \text{NN}),$
 $\text{score}["1 \text{ LRB}"] + T(\text{NN}|\text{LRB}) + S(\text{language} | \text{NN}),$
 $\text{score}["1 \text{ RRB}"] + T(\text{NN}|\text{RRB}) + S(\text{language} | \text{NN}),$
...)

$\text{score}["2 \text{ JJ}"] = \log_sum_exp(\text{score}["1 \text{ NN}"] + T(\text{JJ}|\text{NN}) + S(\text{language} | \text{JJ}),$
 $\text{score}["1 \text{ JJ}"] + T(\text{JJ}|\text{JJ}) + S(\text{language} | \text{JJ}),$
 $\text{score}["1 \text{ VB}"] + T(\text{JJ}|\text{VB}) + S(\text{language} | \text{JJ}),$
...)

$$\log \text{sum exp}(x, y) = \log(\exp(x) + \exp(y))$$

STEPS: FINAL PART

Finish up the sentence with the sentence final symbol

science

/:NN → /+1:</S>

/:JJ

/:VB

/:LRB

/:RRB

...

score["/ +1 </S>"] = log_sum_exp(
score["/ NN"] + T(</S>|NN),
score["/ JJ"] + T(</S>|JJ),
score["/ VB"] + T(</S>|VB),
score["/ LRB"] + T(</S>|LRB),
score["/ NN"] + T(</S>|RRB),
...
)