

Trees!

Jacob Andreas / MIT 6.804-6.864 / Spring 2020

Admin

All grades are now Pass / No Record.
II / TQE status is not affected.

Midterm has been turned into a homework assignment,
and no peer grading on future assignments.
Deadlines are flexible.

Let us know if you're having trouble accessing videos.

Stay safe!

Recap: labels and sequences

Predicting labels

$$S = W_2^T f(W_1^T x)$$

$$f(W_1^T x) = h_1$$

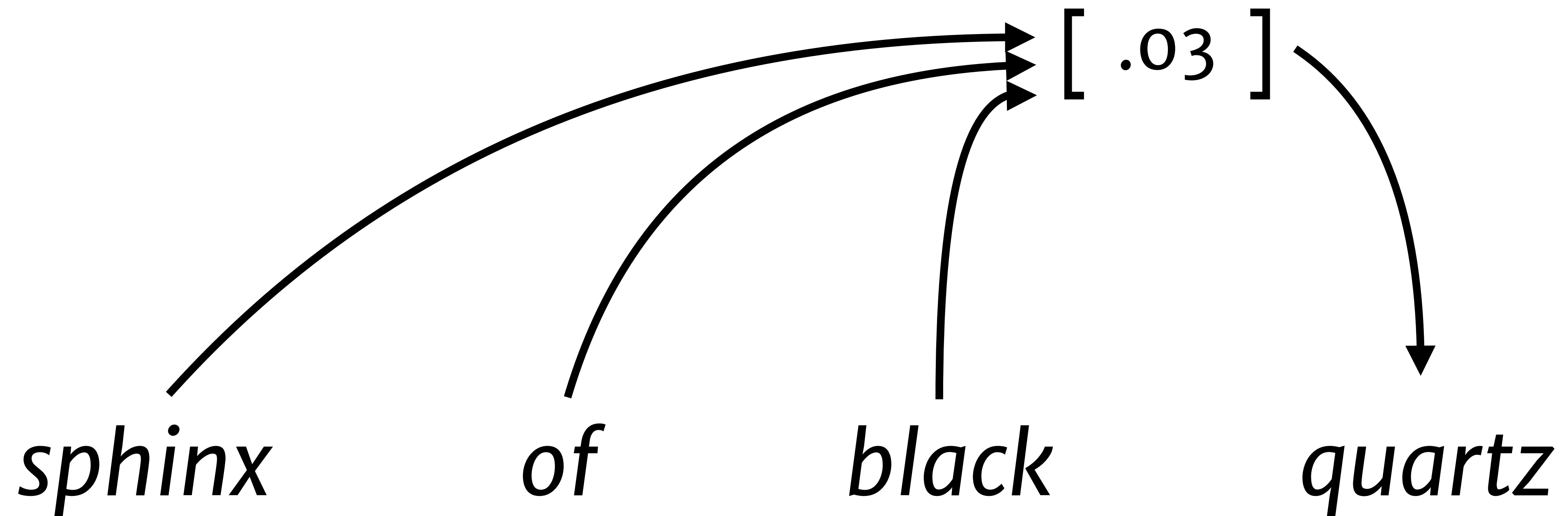
$$W_2^T h_1 = S$$

$$\begin{matrix} x \\ \left[\begin{array}{c} 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{array} \right] \text{input} \end{matrix}$$

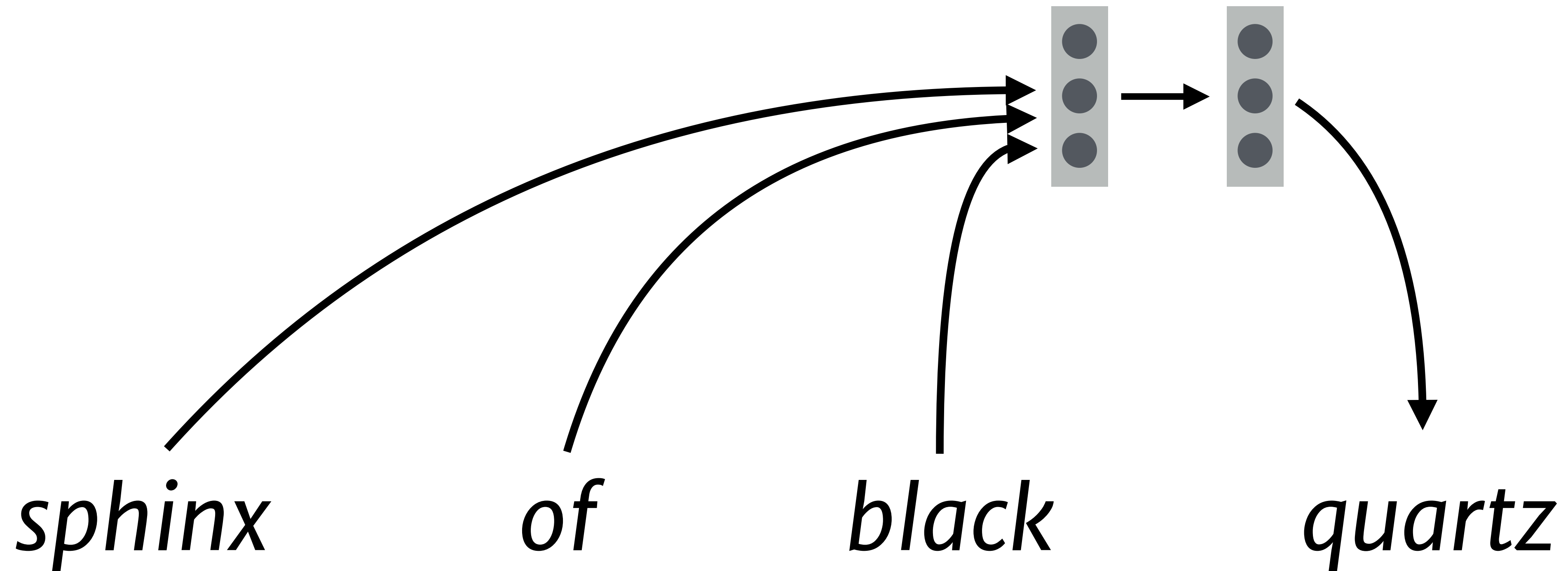
$$\left[\begin{array}{c} 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{array} \right] \text{“hidden layer”}$$

$$\left[\begin{array}{c} 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{array} \right] \text{output}$$

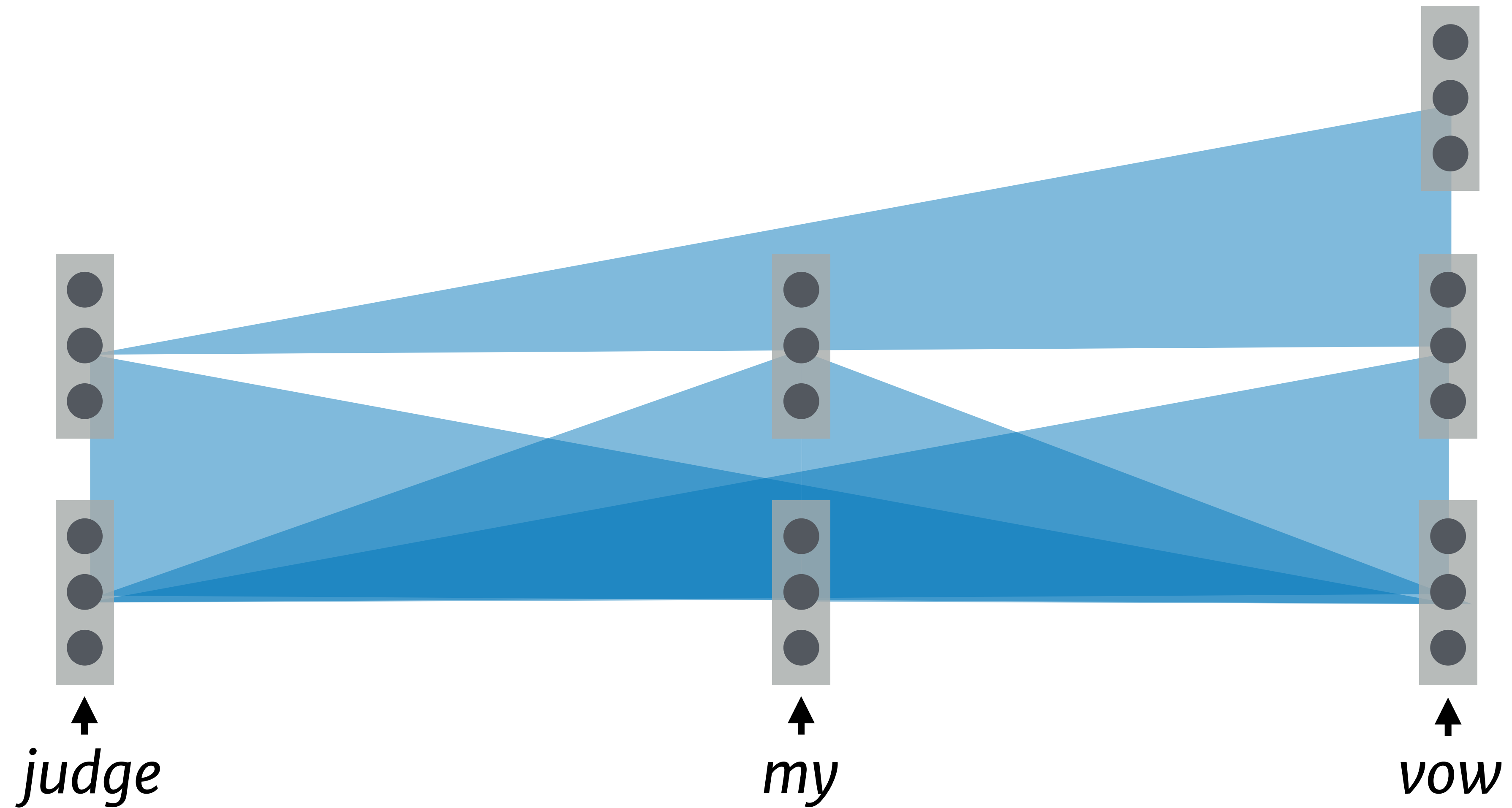
Predicting sequences: n-gram models



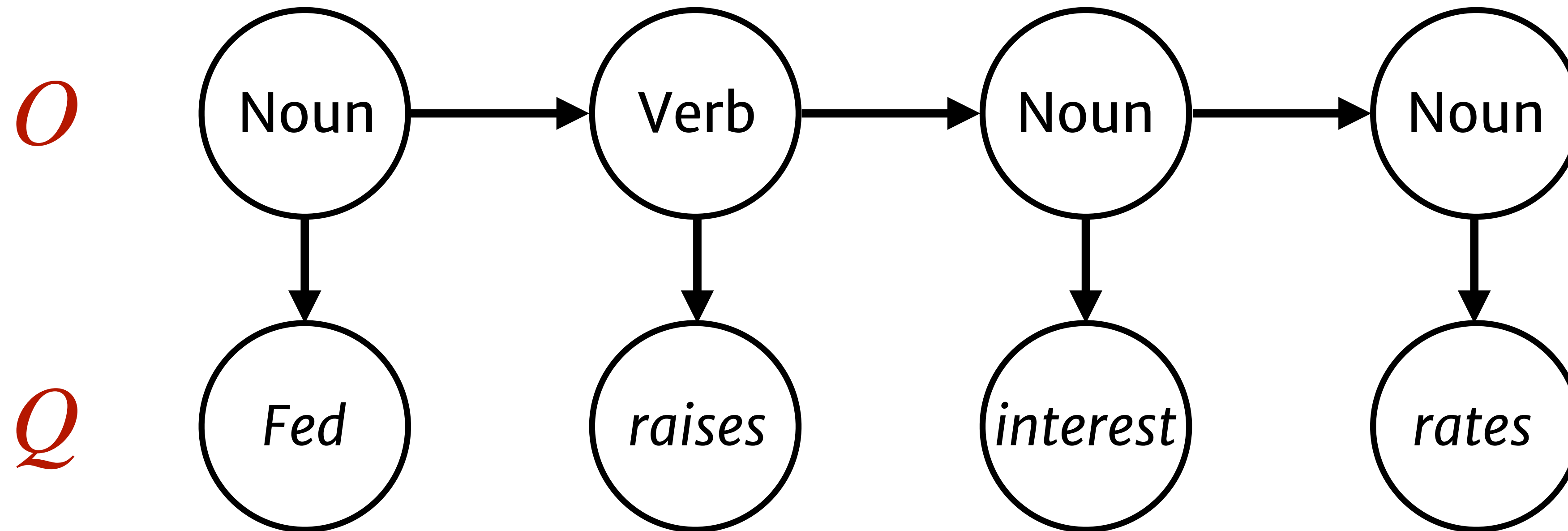
Predicting sequences: n-gram models



Predicting sequences: neural networks

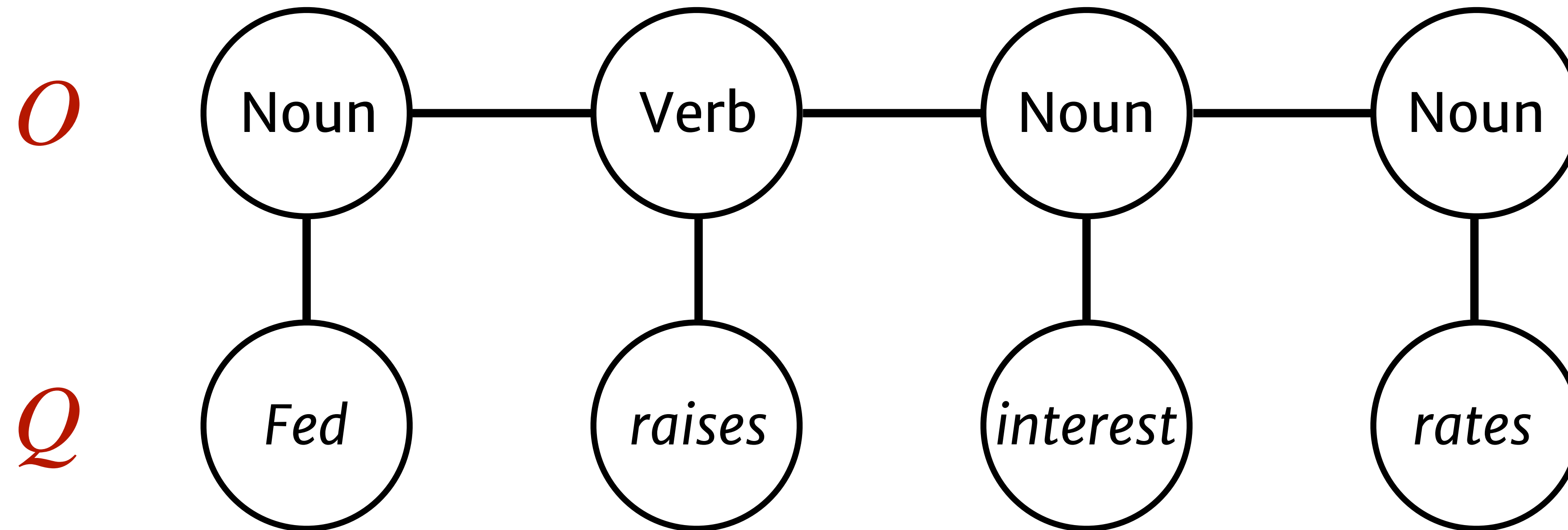


Labeling sequences: HMMs & CRFs



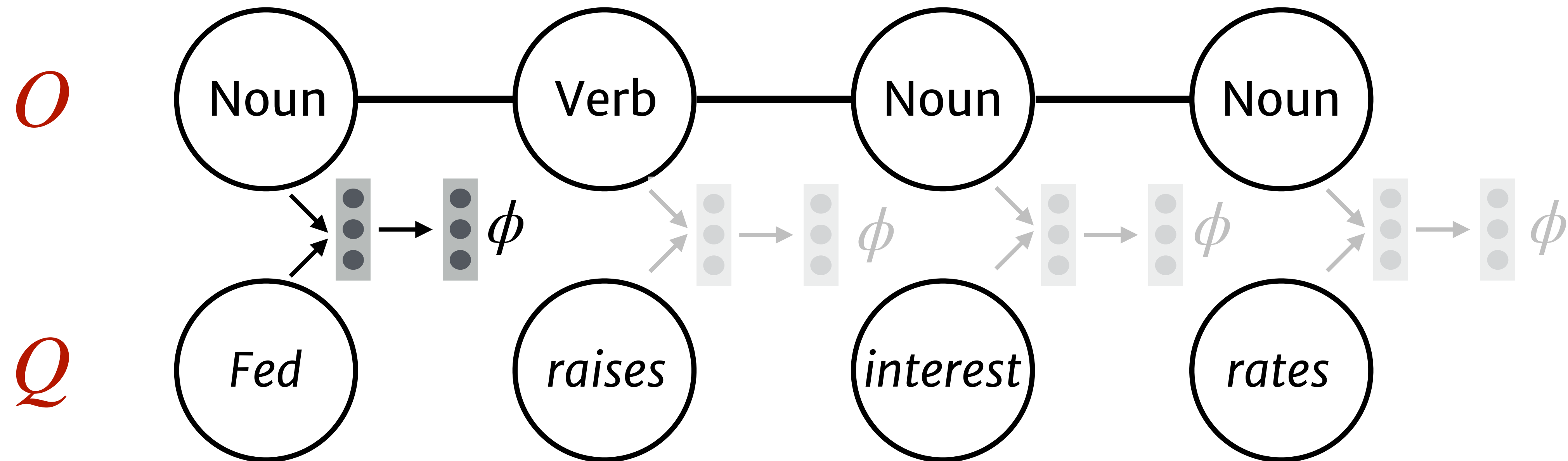
$$\text{HMM: } p(O, Q) = \prod_t p(q_t | q_{:t-1}) p(o_t | q_t)$$

Labeling sequences: HMMs & CRFs



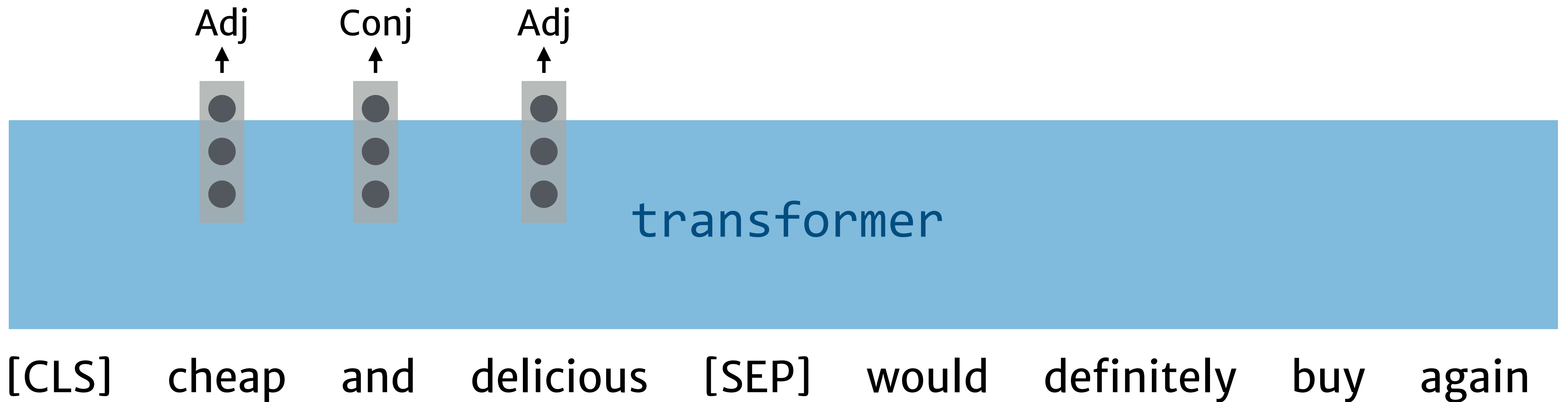
CRF:
$$p(O, Q) = \frac{1}{Z} \exp \left\{ \sum_t a^\top \phi(q_t, q_{:t-1}) + b^\top \phi(o_t | q_t) \right\}$$

Labeling sequences: HMMs & CRFs

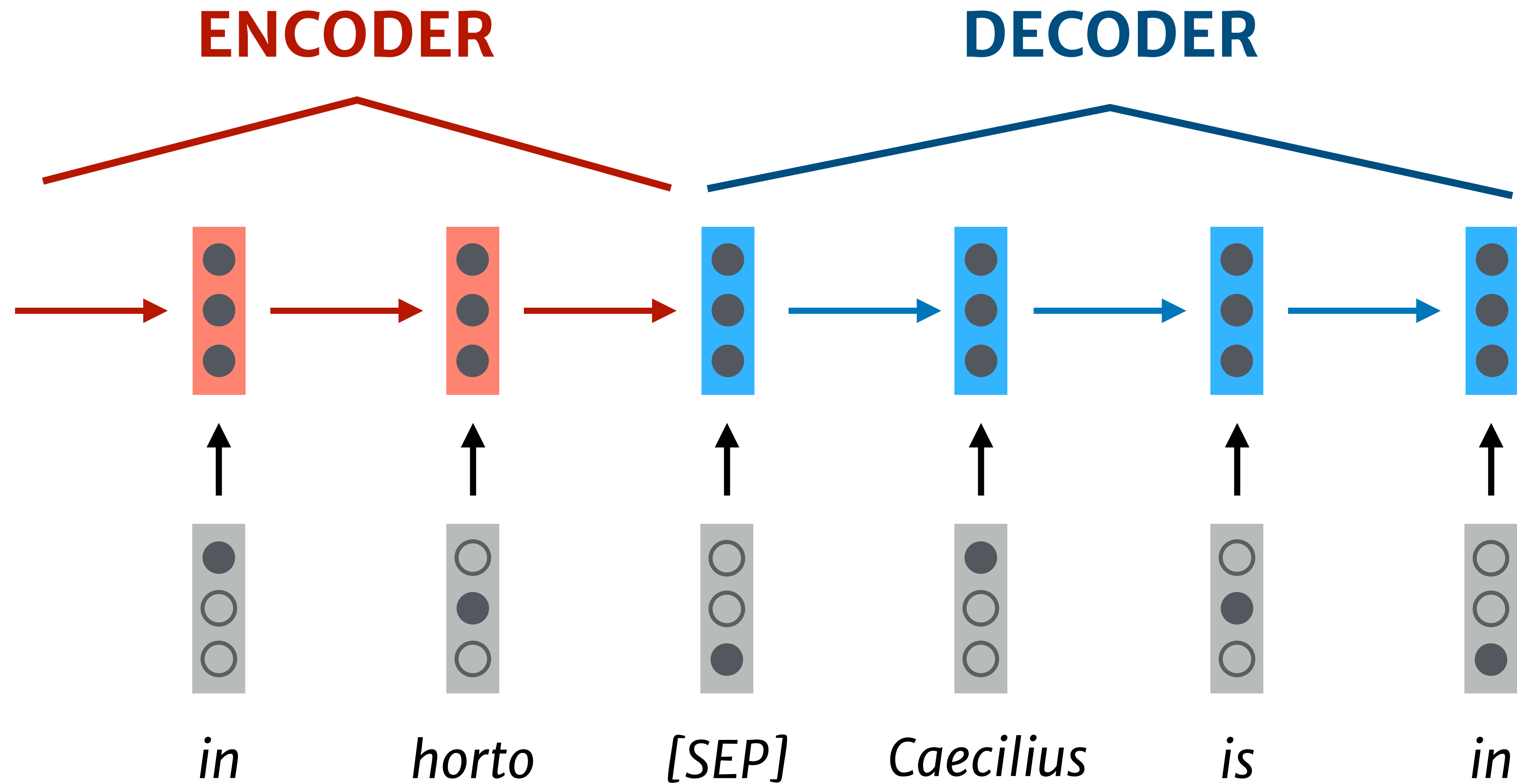


$$\text{CRF: } p(O, Q) = \frac{1}{Z} \exp \left\{ \sum_t a^\top \phi(q_t, q_{:t-1}) + b^\top \phi(o_t | q_t) \right\}$$

Labeling sequences: neural networks

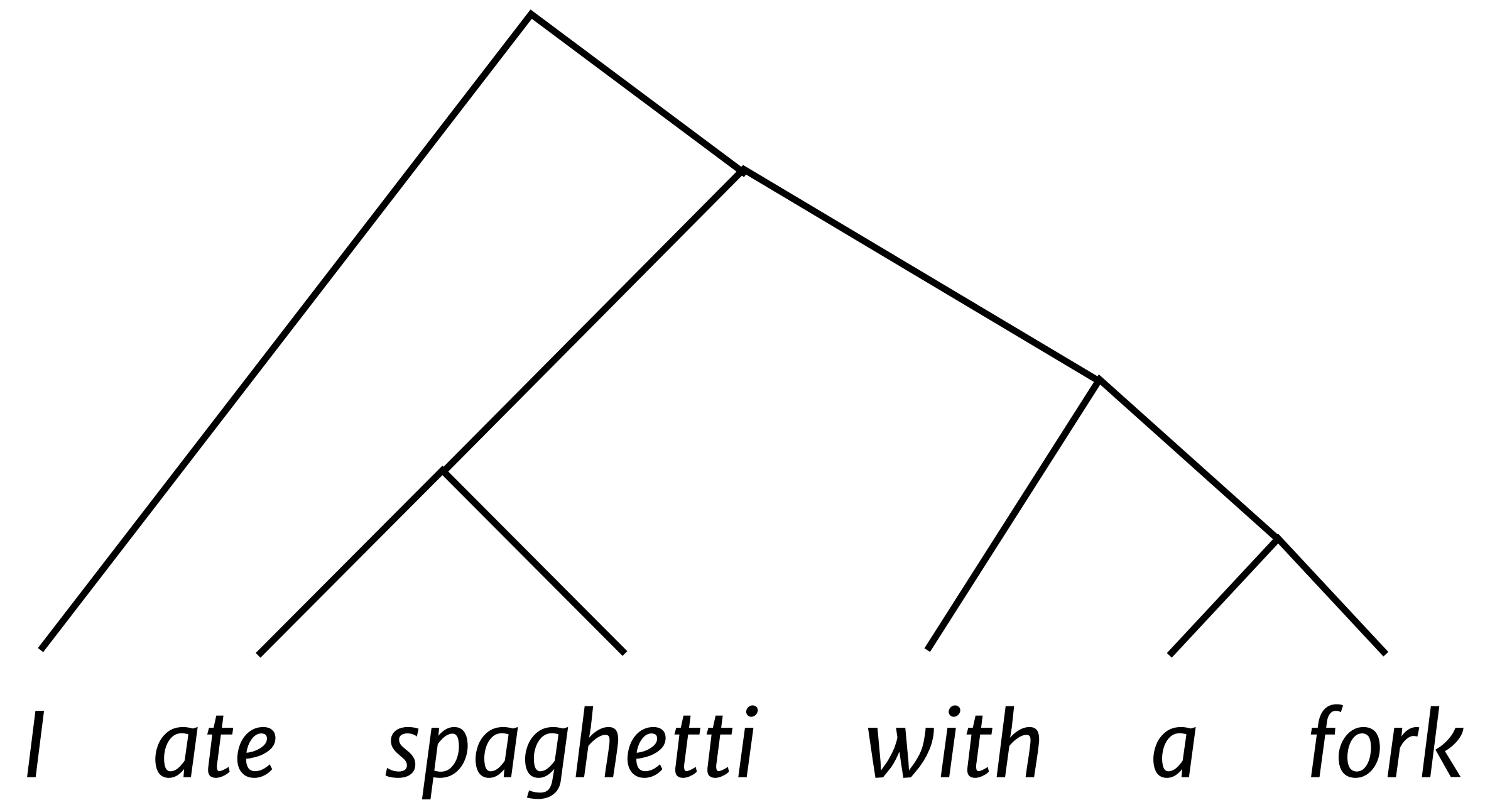
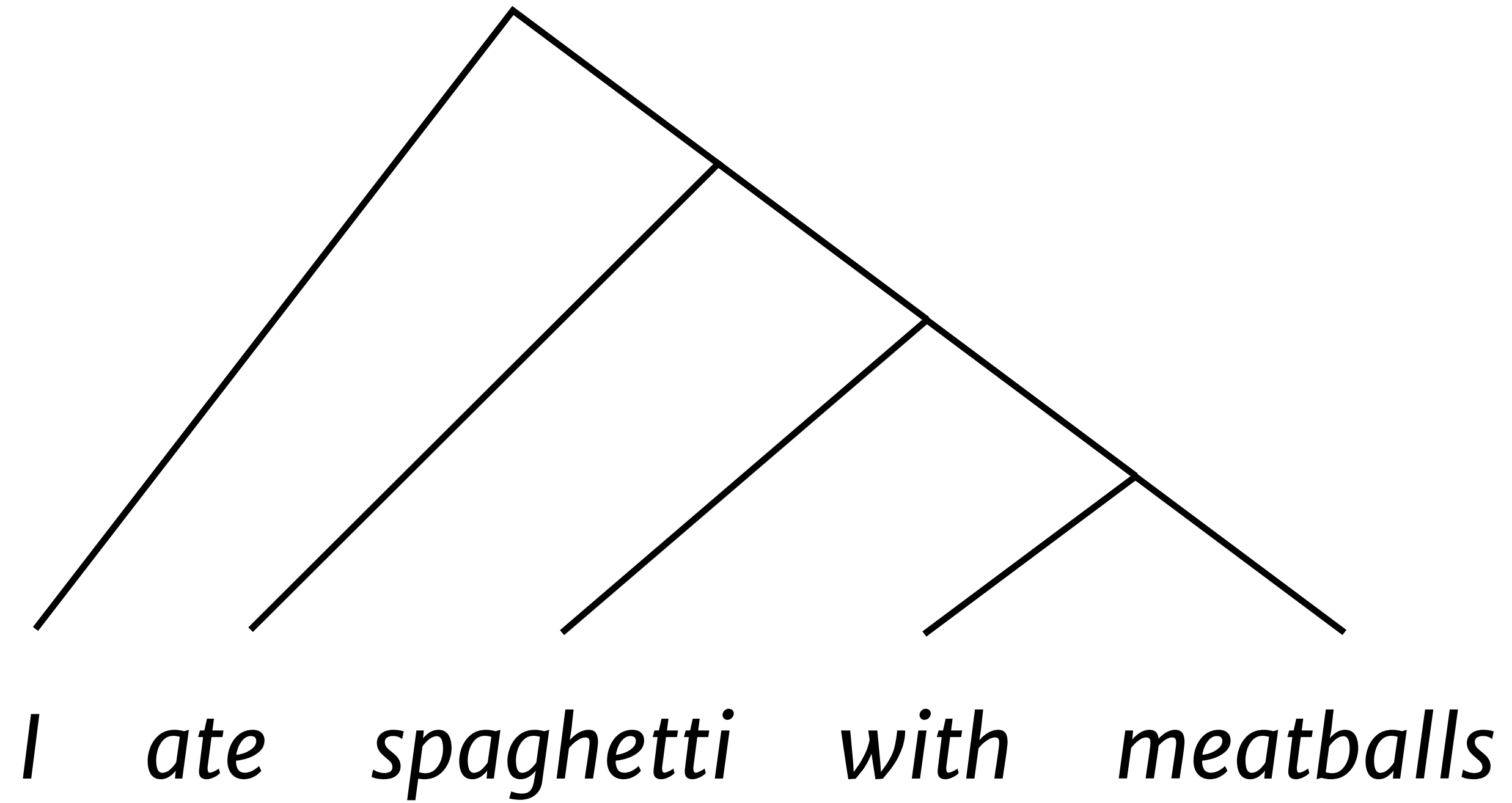


Sequence-to-sequence models

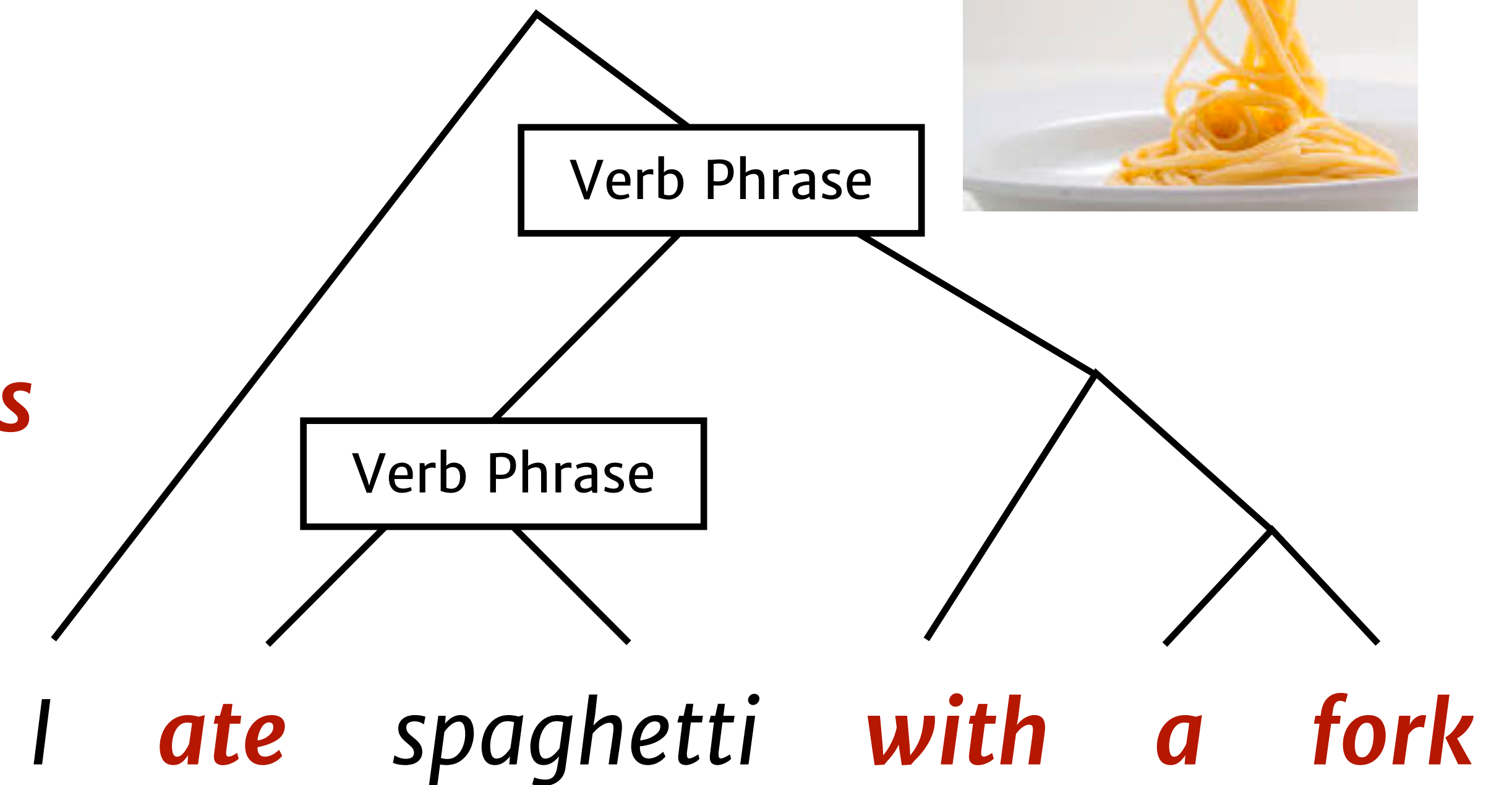
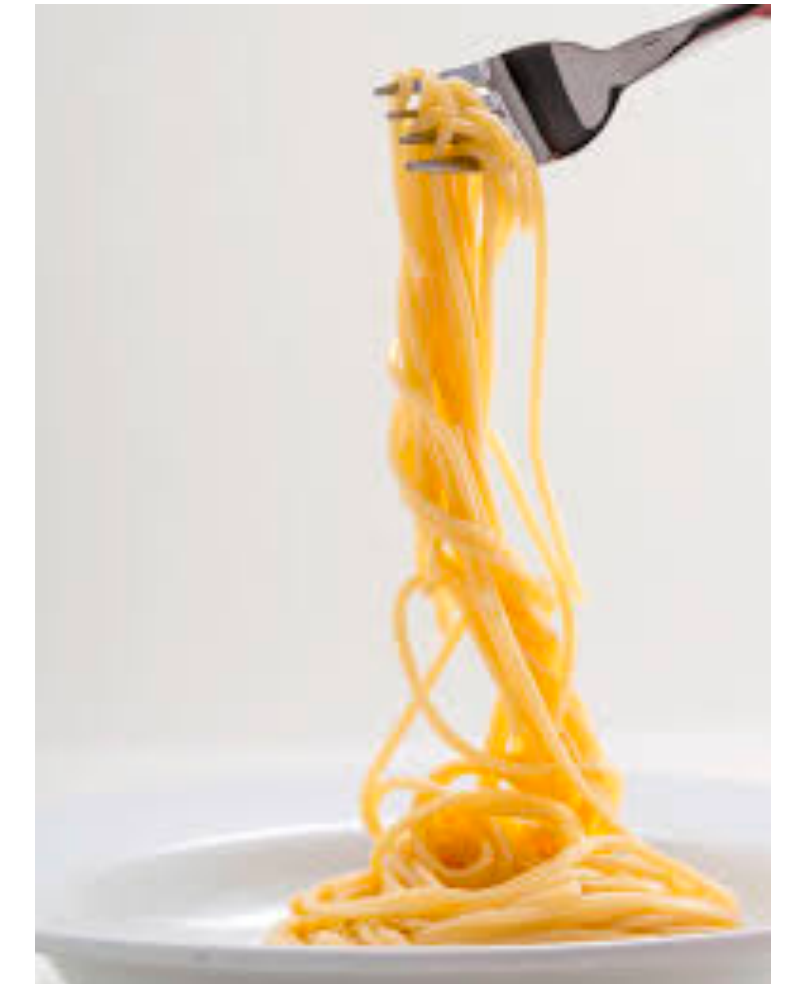
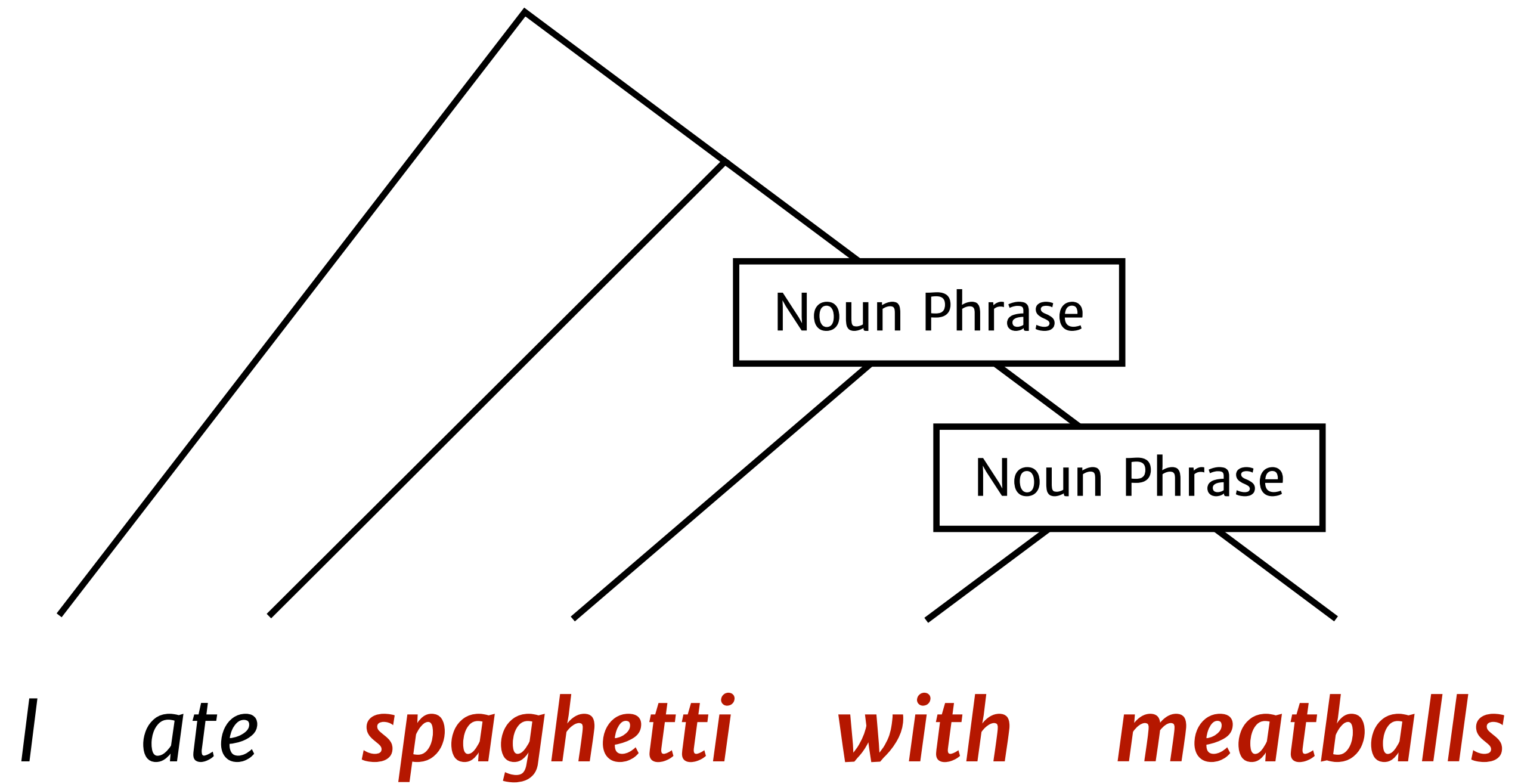


Other structures

Syntax

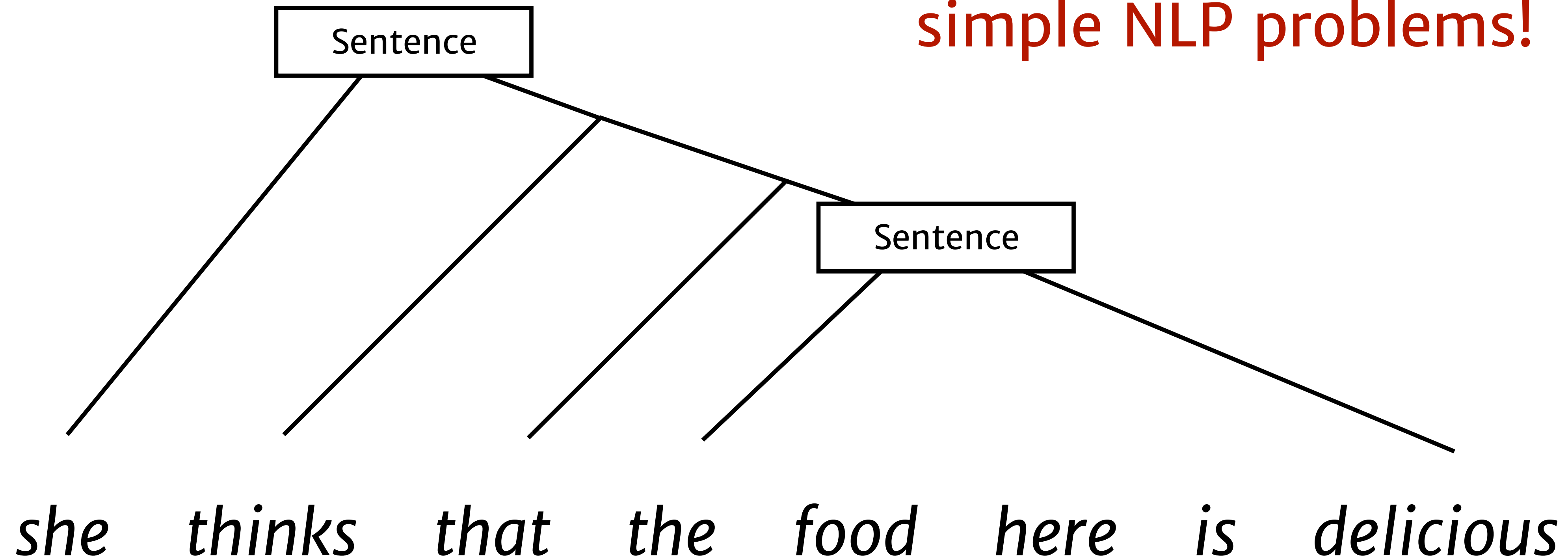


Syntax



Syntax

Useful to distinguish between statements and beliefs, even in simple NLP problems!



Semantics

```
forbid(send(find(type=email), find(keyword=dean)))
```

don't send the email to the dean

Semantics

`forbid(send(find(type=email), find(keyword=dean)))`

`forbid`

don't

`send`

send

`find(type=email)`

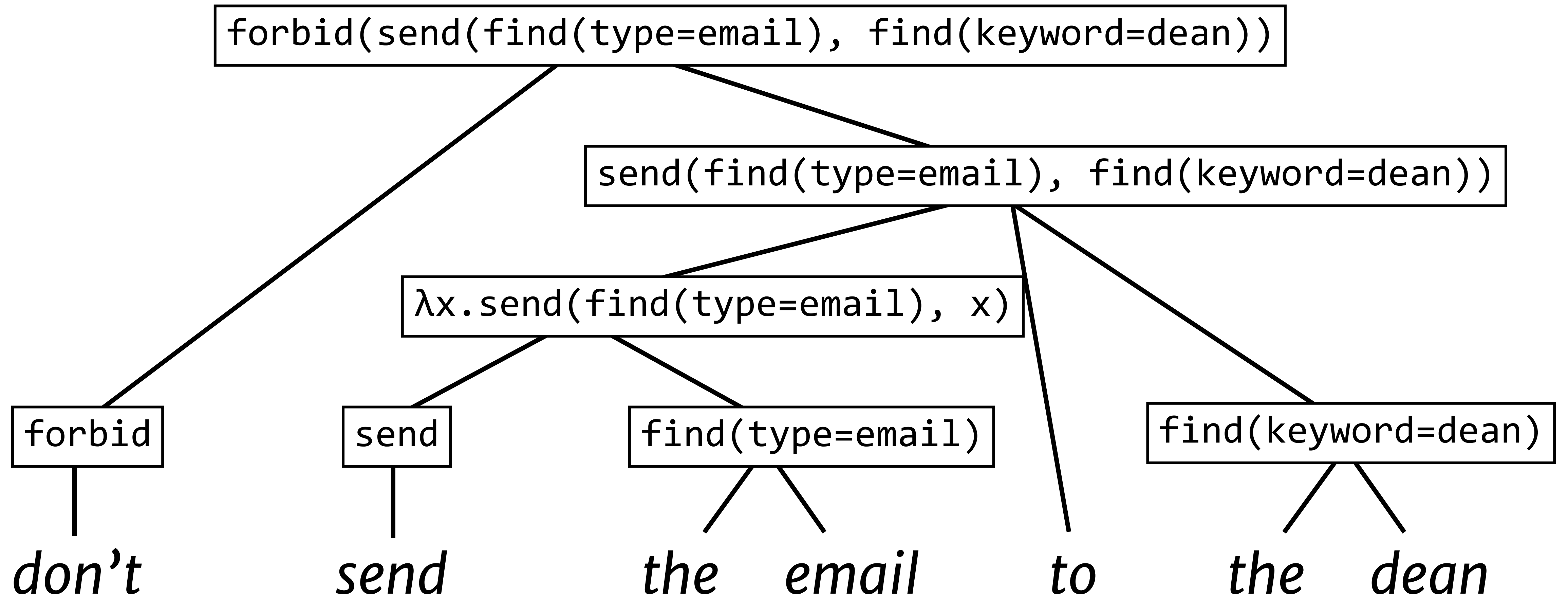
the email

to

`find(keyword=dean)`

the dean

Semantics



Discourse

COMPARISON

The projections are in the neighborhood of 50 cents a share to 75 cents,

CIRCUMSTANCE

compared with a restated \$1.65 a share a year earlier,

when profit was \$107.8 million on sales of \$435.5 million.

Why trees?

“Simplest” formal generative process that provides **hierarchical relationships and long-distance dependencies:**

My aunt gave me a microscope.

My aunt's sister gave me a microscope.

My aunt's sister, who works at the NIH, gave me a microscope.

My aunt's sister, who works at a little-known constituent institute of th

Syntax in ten minutes

Constituents

Key idea from previous examples: some sentence fragments “stick together”—can be moved around, replaced, and modified without affecting meaning / grammaticality:

I ate spaghetti with meatballs

I ate

I ate it

It was spaghetti with meatballs that I ate

Constituents

Some fragments are harder to manipulate:

I ate spaghetti with meatballs

I ate meatballs ✗ (meaning changed)

It was ate spaghetti with that I meatballs ✗ (not grammatical)

Constituents

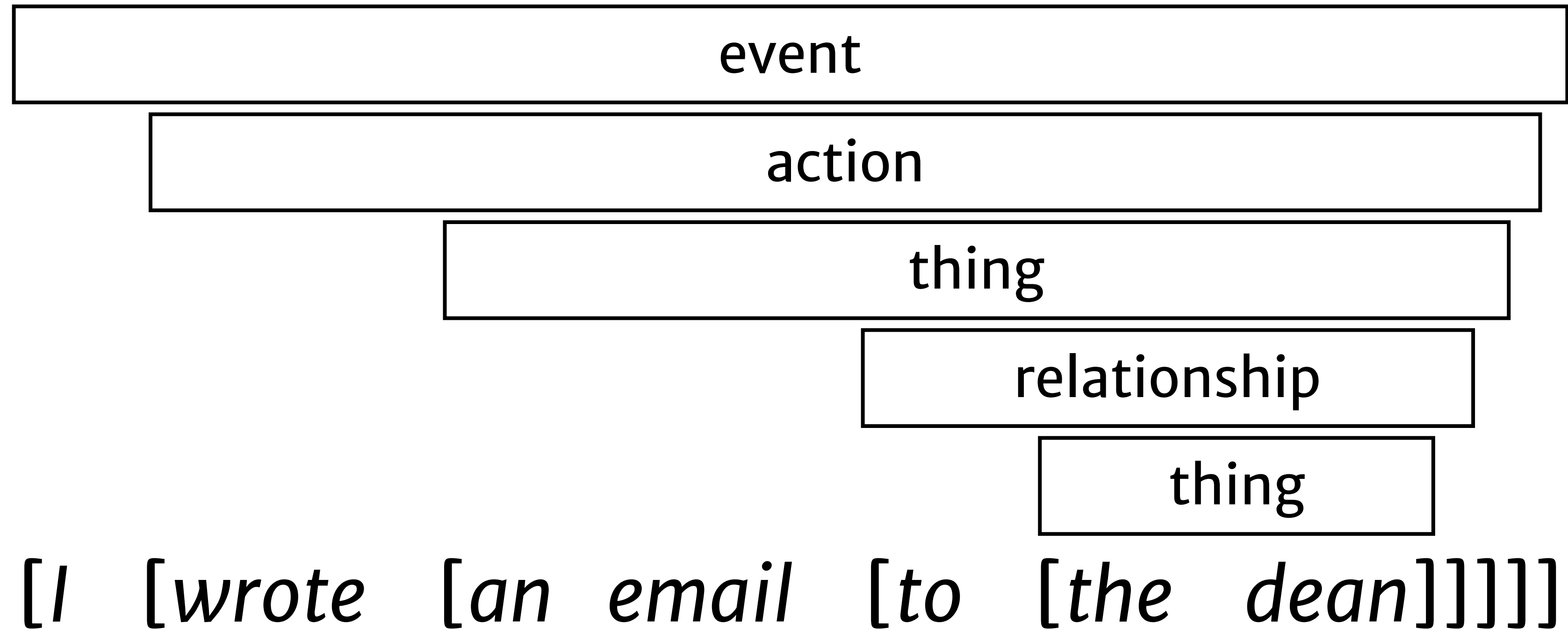
Not just things:

*I ate spaghetti **with a fork***

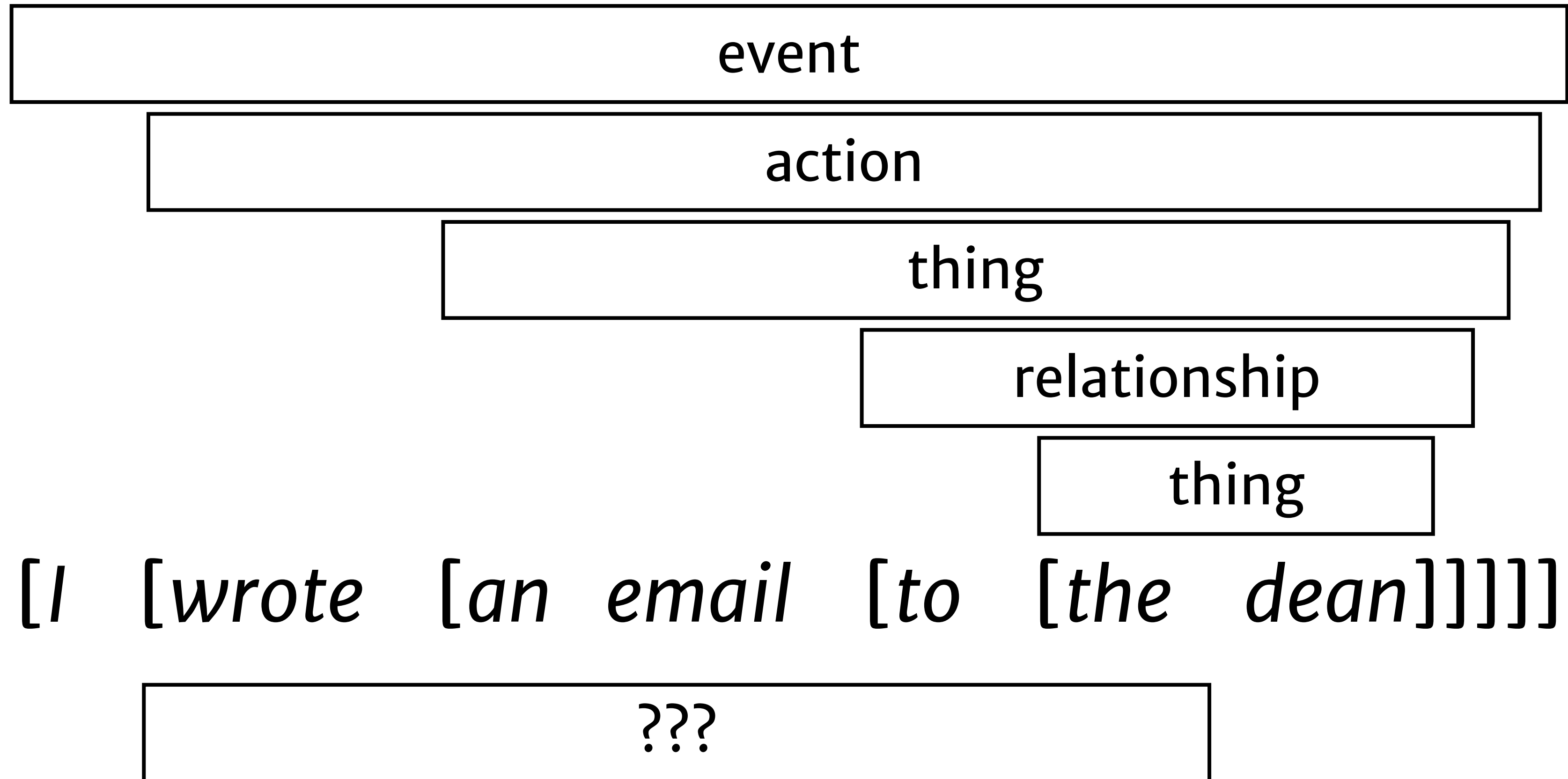
I ate spaghetti

*It was **with a fork** that I ate spaghetti*

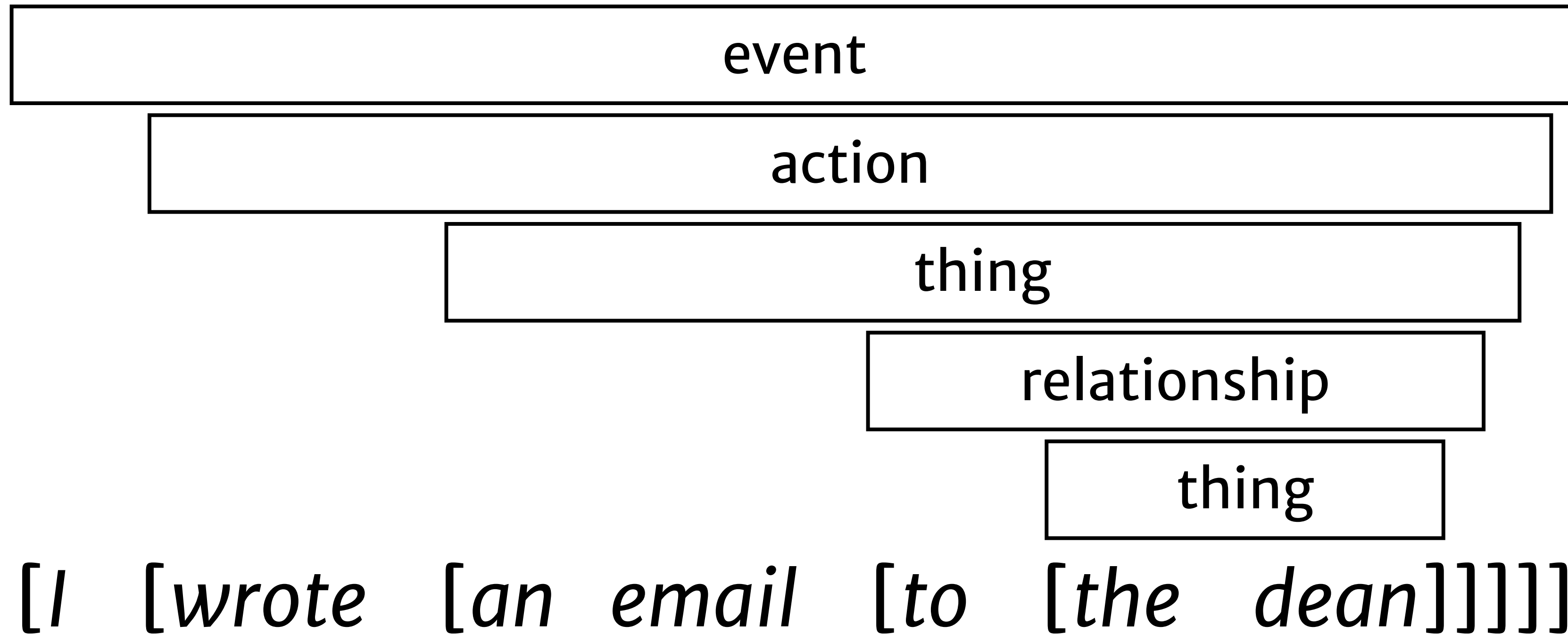
Constituents & Types



Constituents & Types

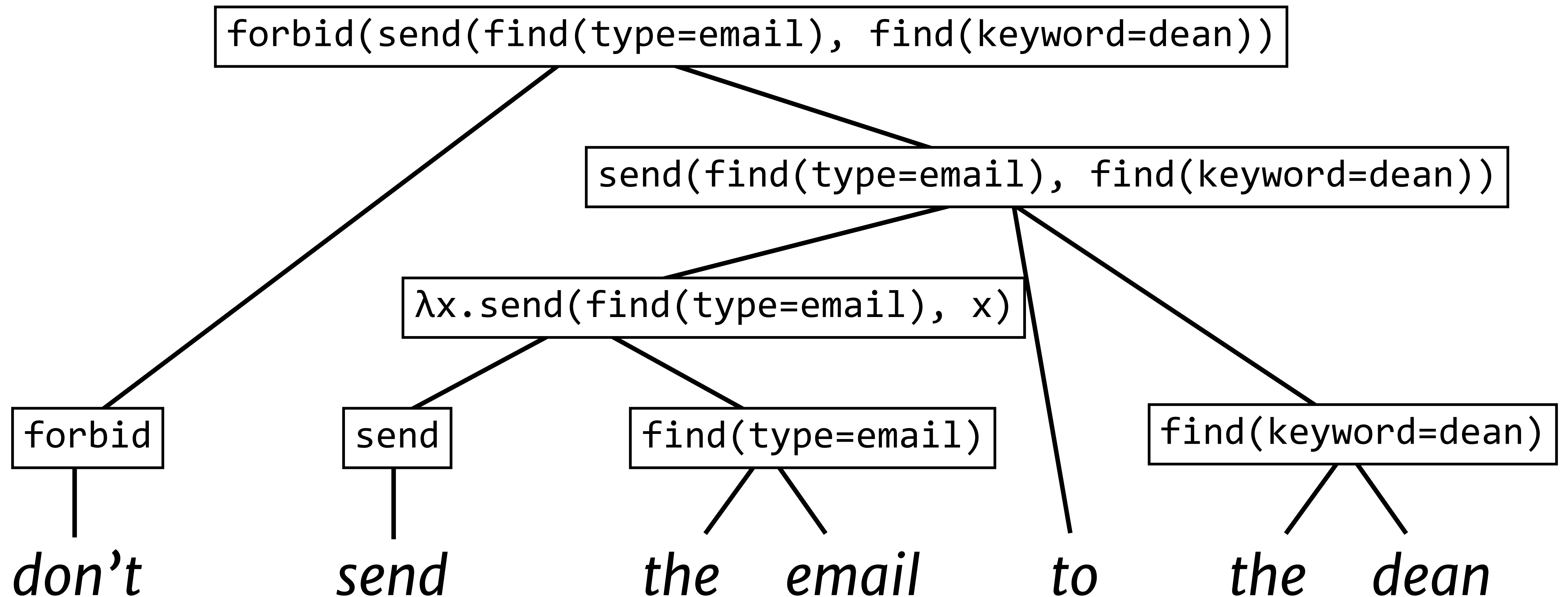


Constituents & Types

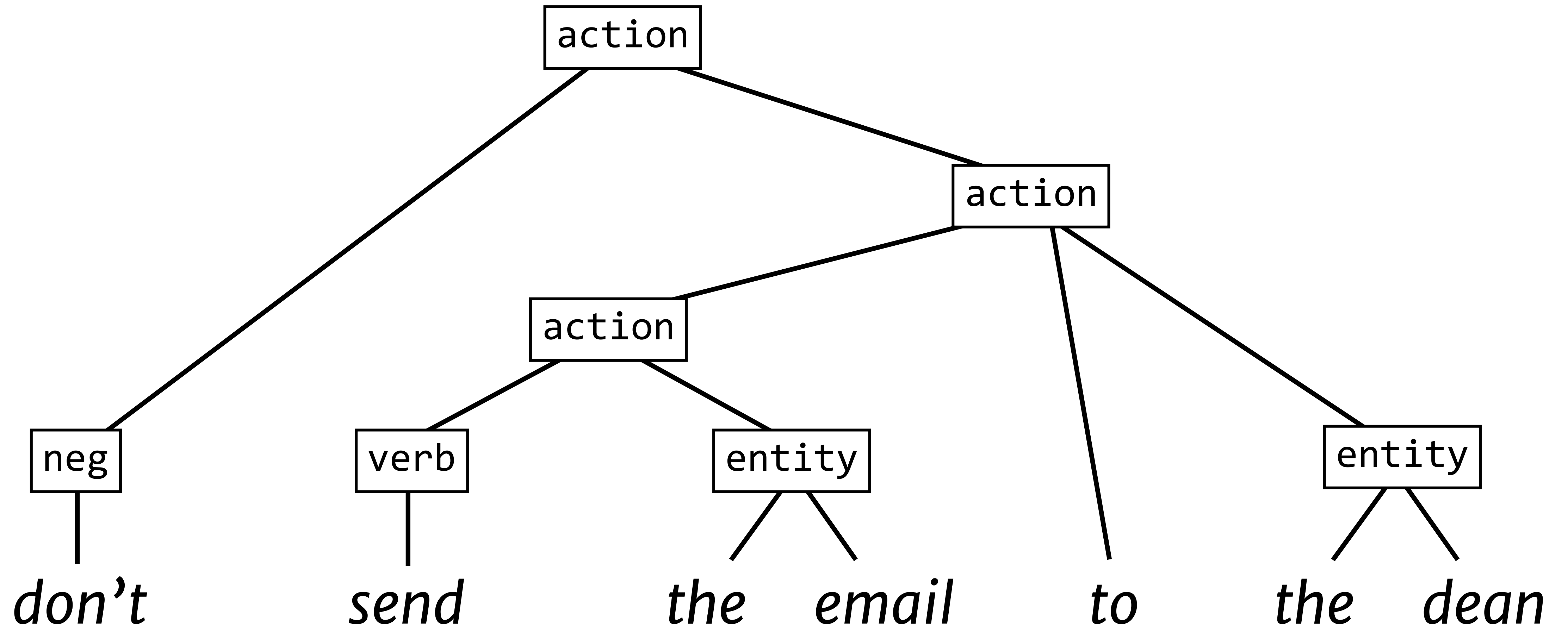


Lots of research on the exact form of this hierarchy.
For most NLP applications: **entities, events, relations.**

Types & semantics



Types & semantics



Context free grammars

Just like in HMMs, we'd like to define some joint distribution over sentences and underlying structures, and reason about marginals and conditionals.

What's the right distribution over trees and sentences?

Context free grammars

A **sentence** might consist of an entity and an action.

[I] [*swallowed the spider*]

Context free grammars

A **sentence** might consist of an entity and an action.

[I] [swallowed the spider]

$S \rightarrow NP \ VP$ a Noun Phrase followed by a Verb Phrase make a Sentence

A **sentence** might just consist of an action.

[eat the spider]

$S \rightarrow VP$ a Verb Phrase makes a Sentence

Context free grammars

$S \rightarrow NP VP \mid VP$

“or”

the followed by a noun makes an entity

$NP \rightarrow \textit{the} N$

$N \rightarrow \textit{cat} \mid \textit{dog} \mid \textit{spider} \mid \textit{cheesecake} \mid \textit{democracy}$

a verb and an optional entity make an action

$VP \rightarrow V \mid V NP$

$V \rightarrow \textit{eat} \mid \textit{eats} \mid \textit{run} \mid \textit{differentiate} \mid \dots$

A sample from our CFG

S → NP VP | VP
NP → *the* N
N → *cat* | *dog* | *spider* | *cheesecake* | *democracy*
VP → V | V NP
V → eat | eats | run | differentiate | ...

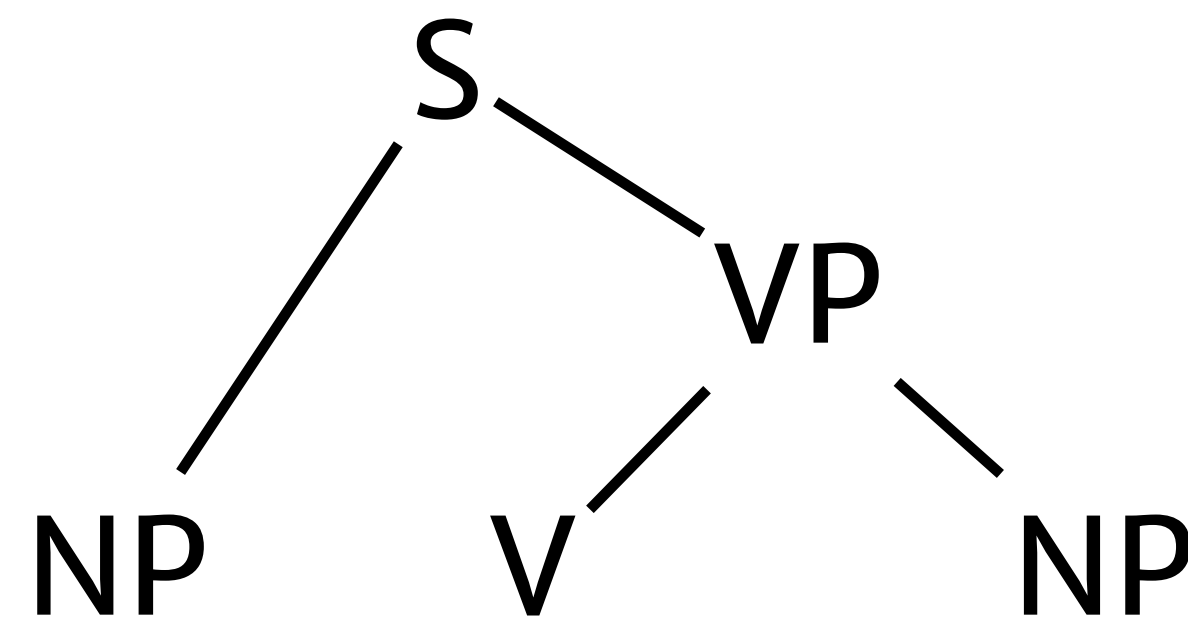
A sample from our CFG

S → NP VP | VP
NP → *the* N
N → *cat* | *dog* | *spider* | *cheesecake* | *democracy*
VP → V | V NP
V → *eat* | *eats* | *run* | *differentiate* | ...

S	<i>the cat</i> VP	<i>the cat eat the N</i>
NP VP	<i>the cat V NP</i>	<i>the cat eat the democracy</i>
<i>the</i> N VP	<i>the cat eat NP</i>	

A sample from our CFG

S → NP VP | VP
NP → *the* N
N → *cat* | *dog* | *spider* | *cheesecake* | *democracy*
VP → V | V NP
V → *eat* | *eats* | *run* | *differentiate* | ...



the cat eat the democracy

What about other languages?

a taky na to většinou nemá peníze

and also for it generally hasn't money

and in most cases he has no money for it either

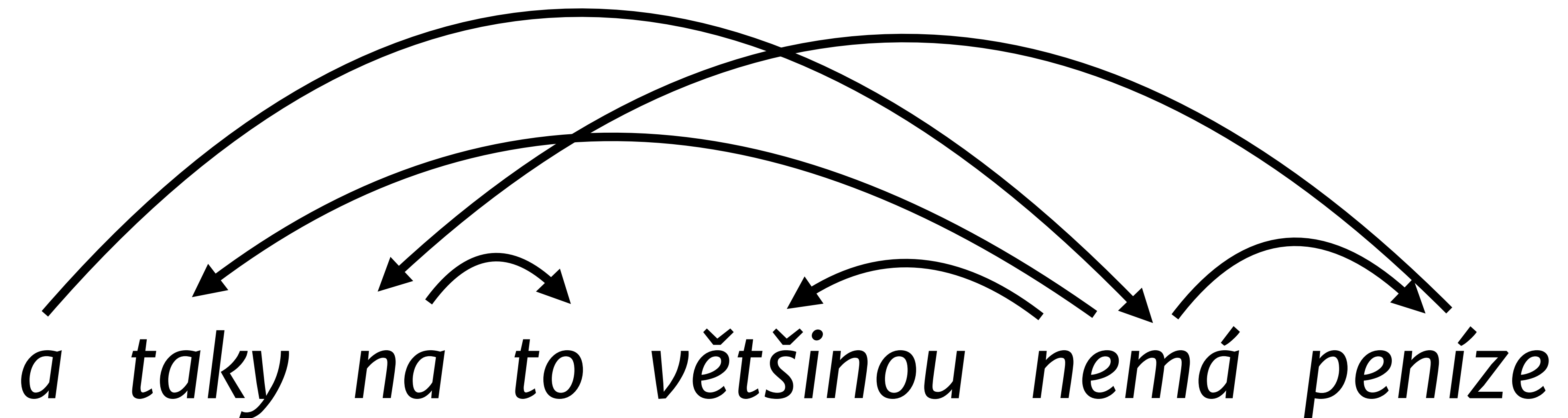
What about other languages?

a taky na to většinou nemá peníze
and also for it generally hasn't money

and in most cases he has no money for it either

can't draw a constituency tree!

Dependency grammar

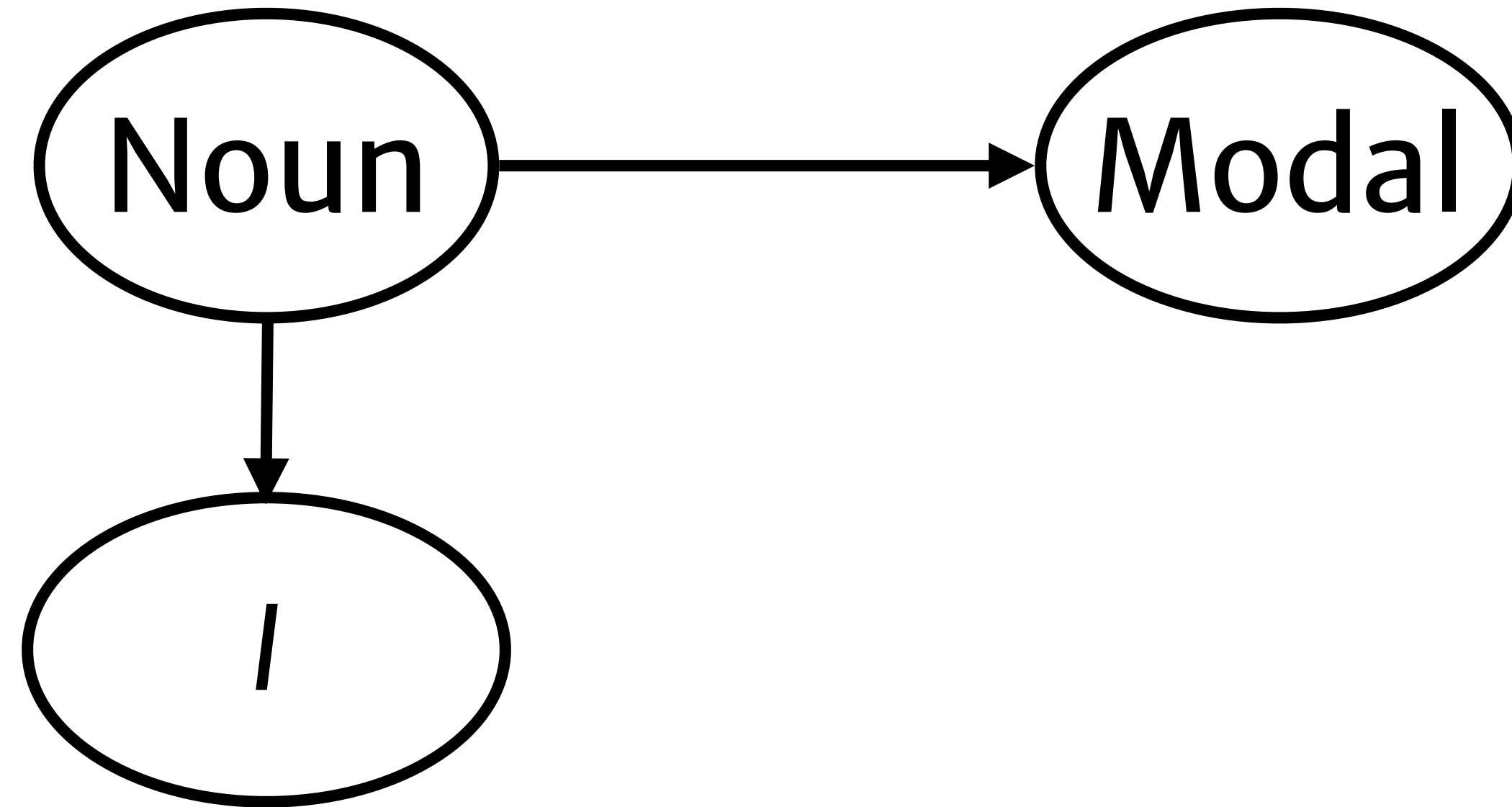


Probabilistic grammars

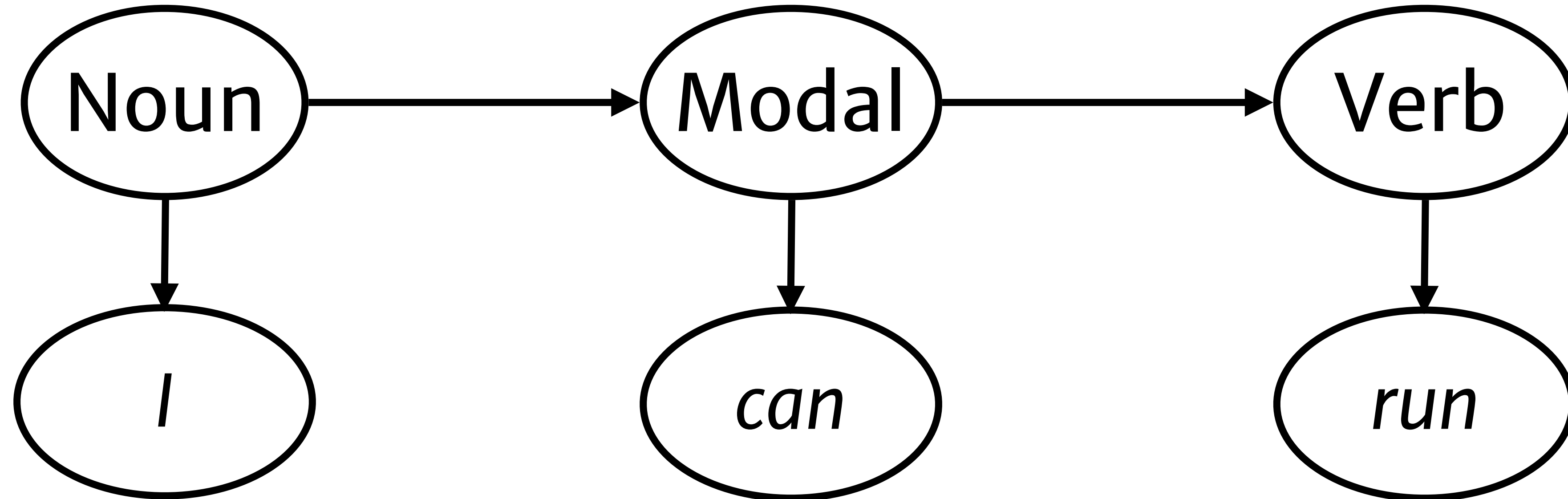
CFGs as generative models

Noun

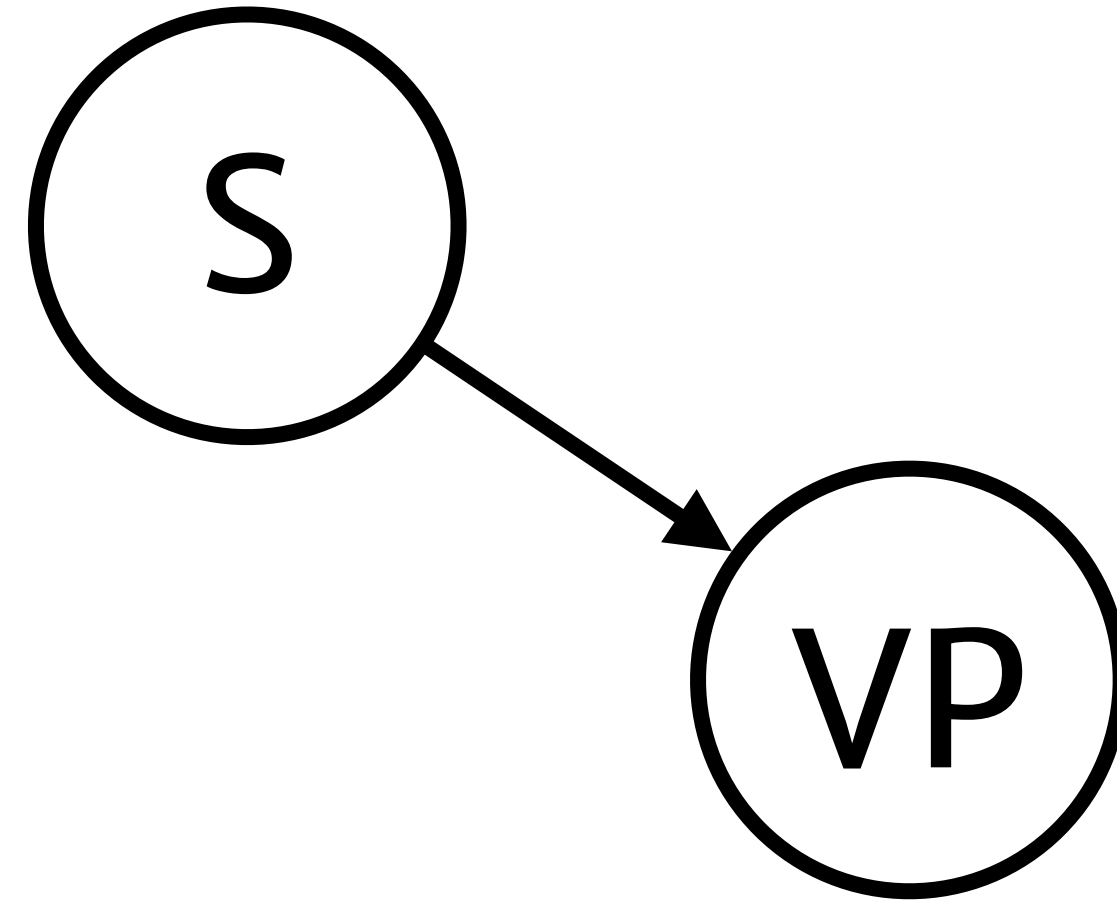
CFGs as generative models



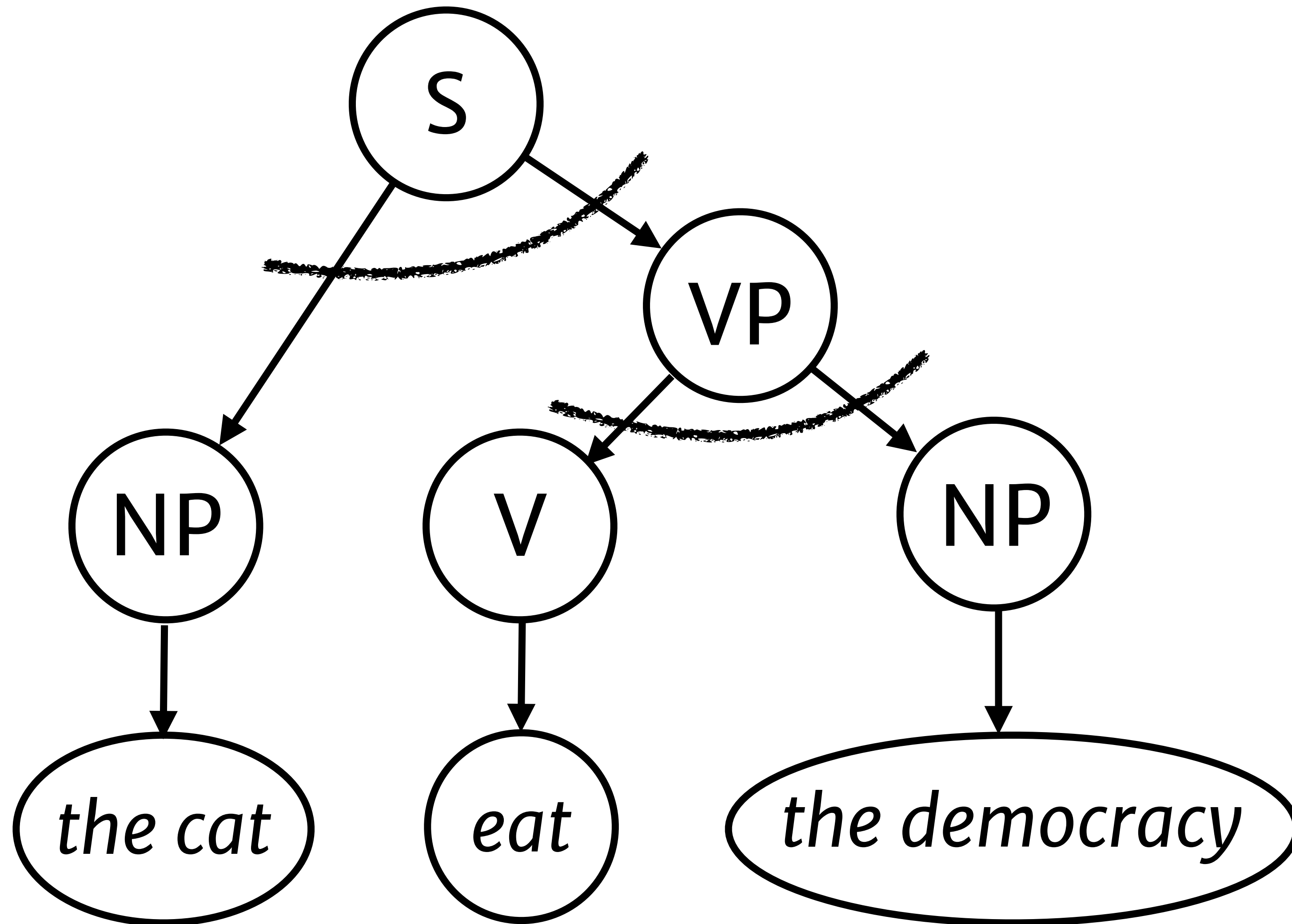
CFGs as generative models



CFGs as generative models



CFGs as generative models



Probabilistic CFGs

A probabilistic context free grammar (PCFG) consists of

1. A set of nonterminal symbols N
2. A set of terminal symbols T
3. A set of rules R
4. A set of rule probabilities $p(r \in R \mid n \in N)$

Probabilistic CFGs

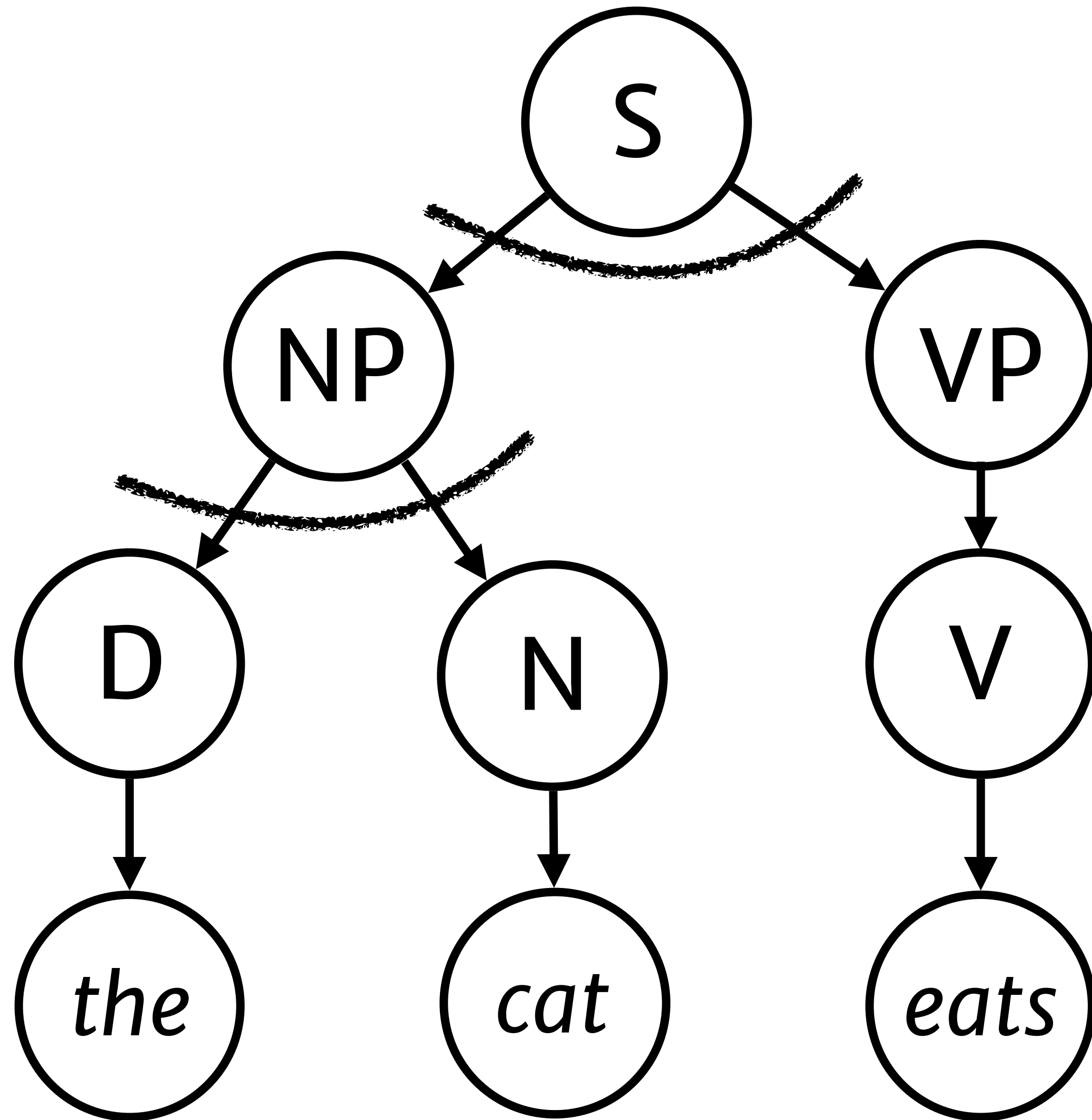
A rule consists of

1. A left hand symbol
2. A sequence of right-hand symbols

LHS		RHS	prob
S	→	NP VP	0.75

such that $\sum_{\text{rules with LHS symbol A}} p(\text{rule} \mid \text{LHS}) = 1$

Queries: joint probability



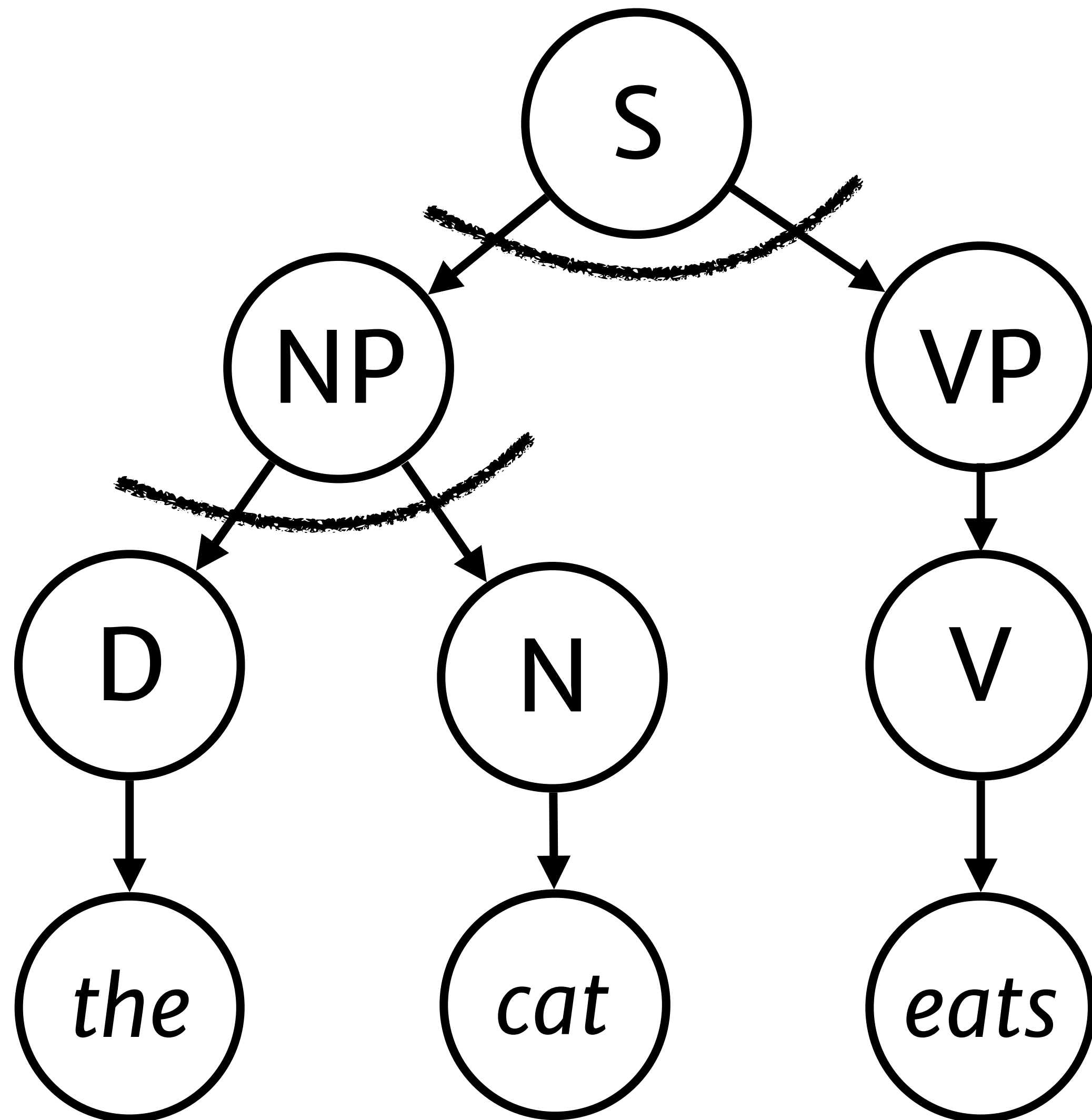
$$p(T, S)$$

joint prob. of tree and sentence

$$p(S \rightarrow NP VP) \cdot p(NP \rightarrow D N)$$

$$p(D \rightarrow the) \cdot p(N \rightarrow cat) \dots$$

Queries: best tree



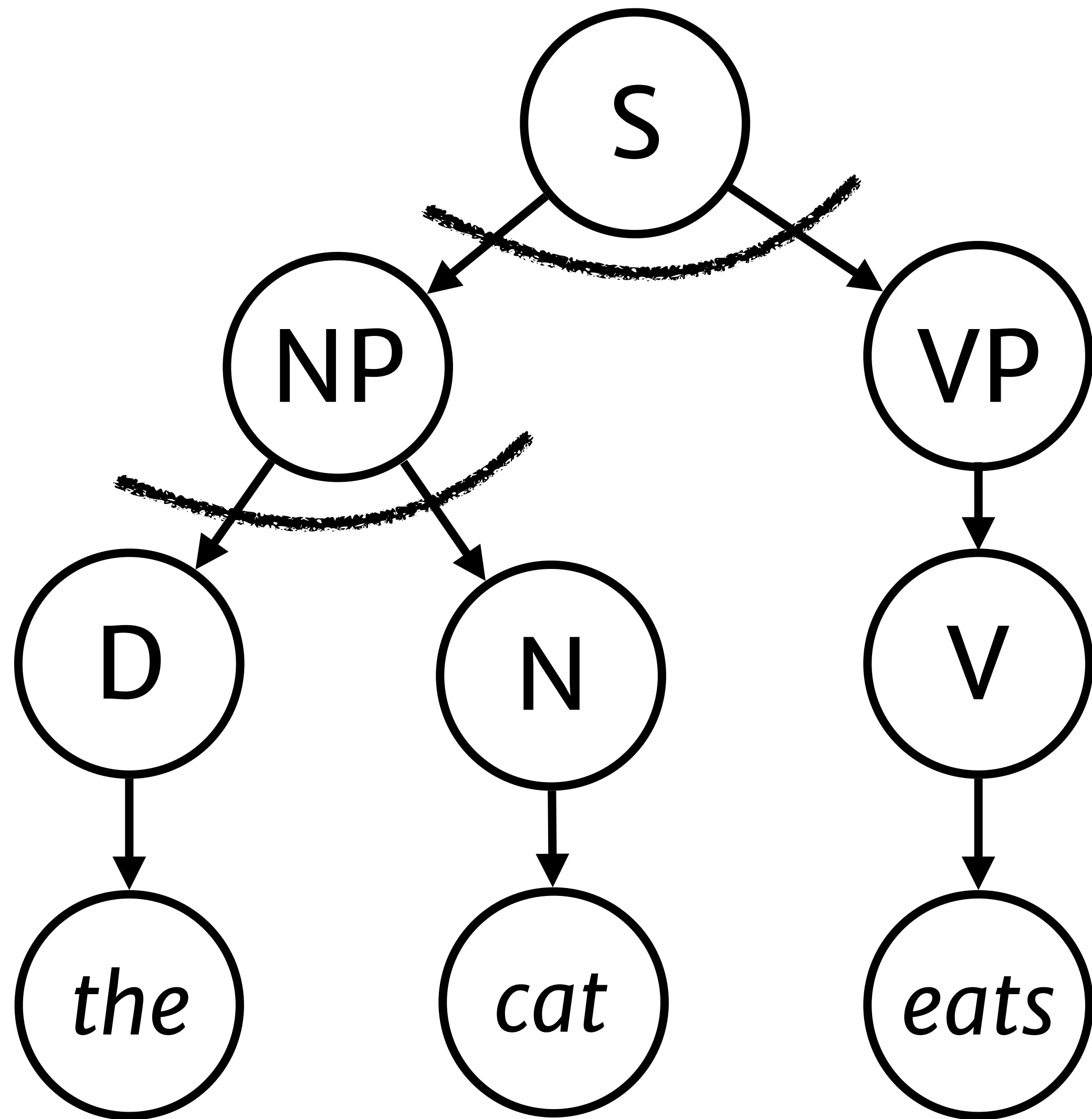
$$\operatorname{argmax}_T p(T | S)$$

most probable tree given sent.

$$p(S \rightarrow NP VP) \cdot p(NP \rightarrow D N)$$

$$p(D \rightarrow \text{the}) \cdot p(N \rightarrow \text{cat}) \dots$$

Queries: sentence marginal



$p(S)$

prob. of sentence under **any** tree

(there are $\frac{(2n)!}{(n+1)!n!}$ *unlabeled* binary trees over n words...)

Parsing

Chomsky normal form

Notational convenience: only binary trees.

Every rule has one of these forms:

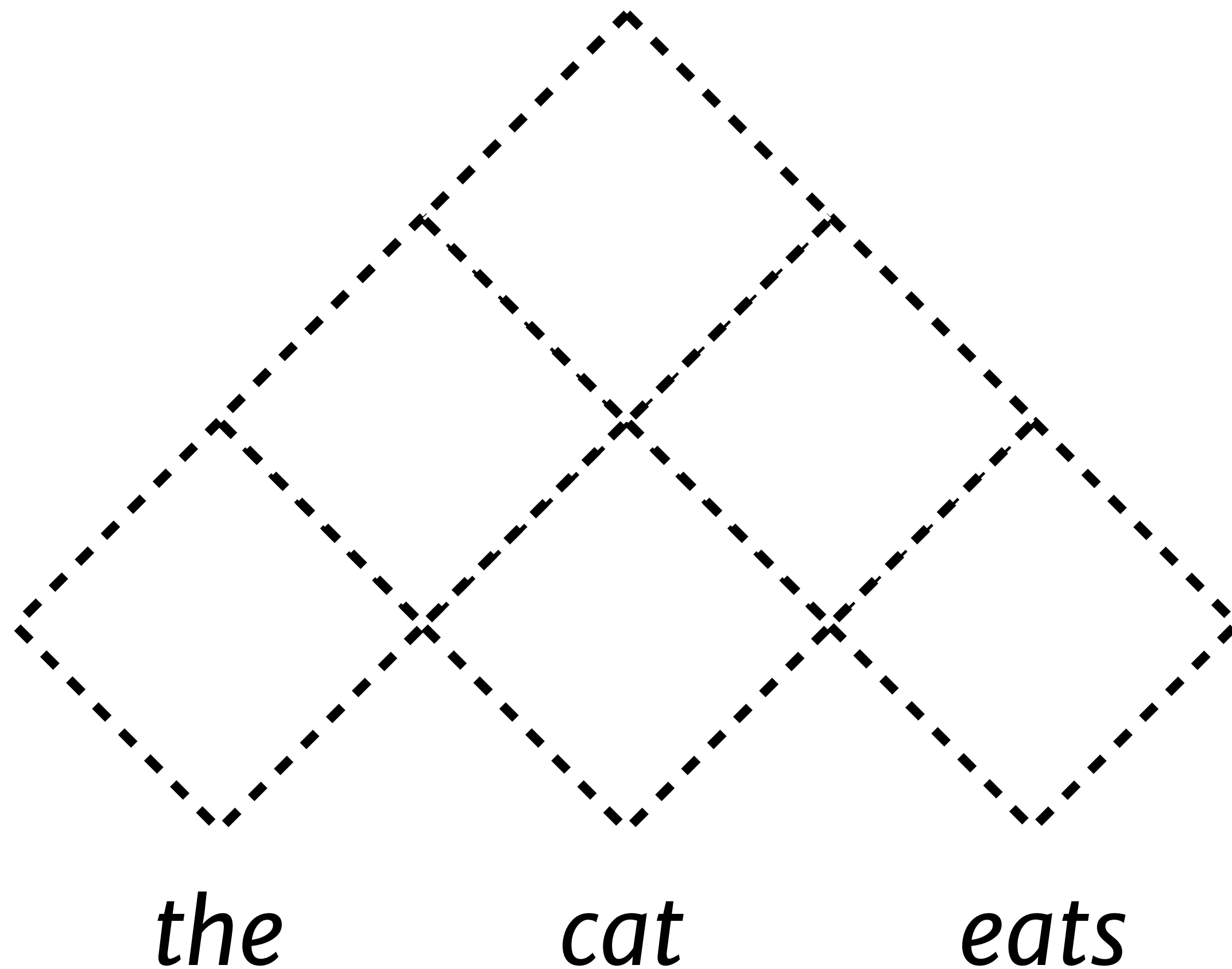
Nonterminal \rightarrow Terminal

Nonterminal \rightarrow Nonterminal Nonterminal

(Can always get rules into this form by introducing new NTs)

Warmup: the CKY algorithm

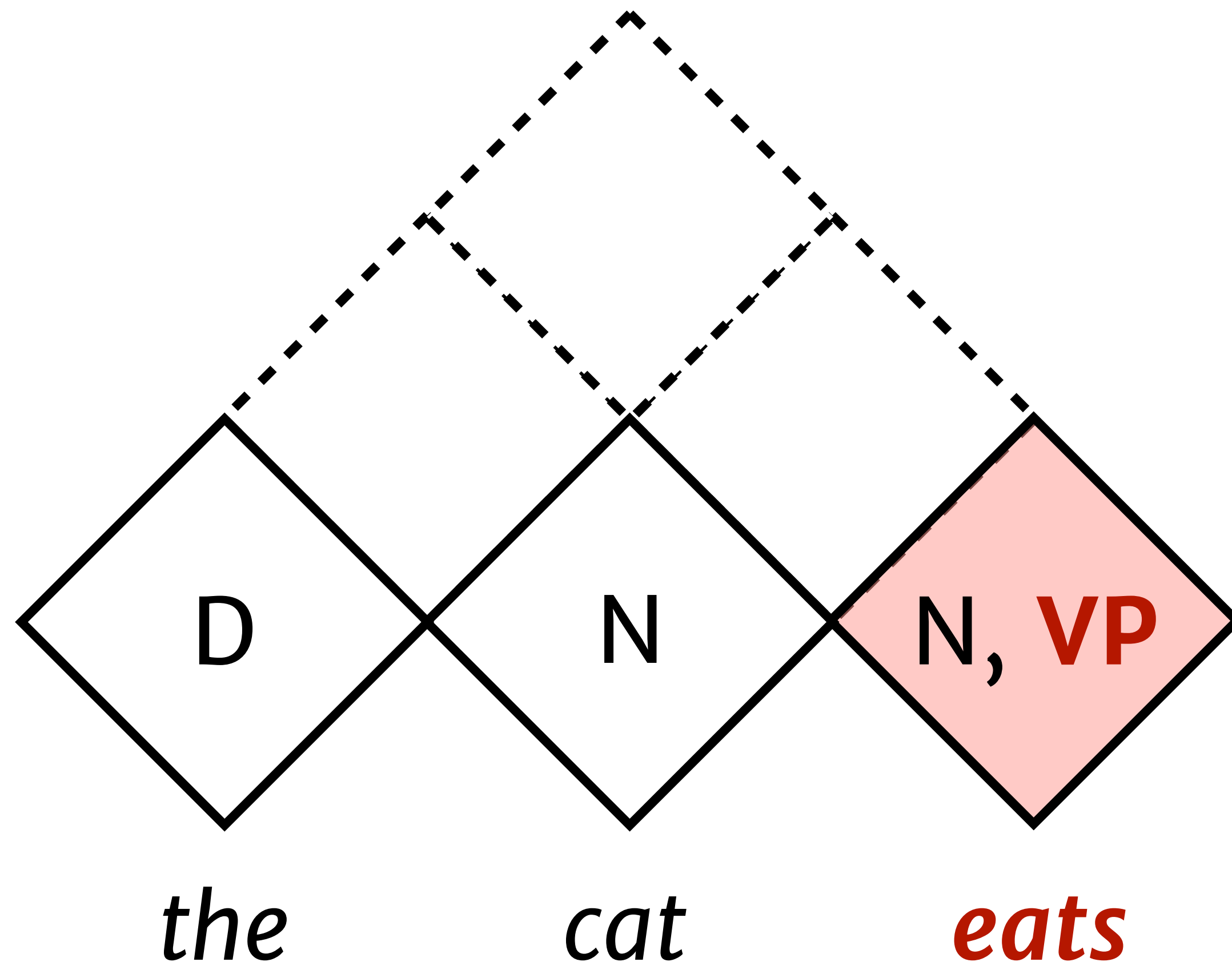
Is the string S generated by CFG G ?



$S \rightarrow NP VP$
 $NP \rightarrow D N$
 $N \rightarrow \text{cat} \mid \text{eats}$
 $VP \rightarrow \text{eats} \mid \text{sings}$
 $D \rightarrow \text{the} \mid \text{a}$

Warmup: the CKY algorithm

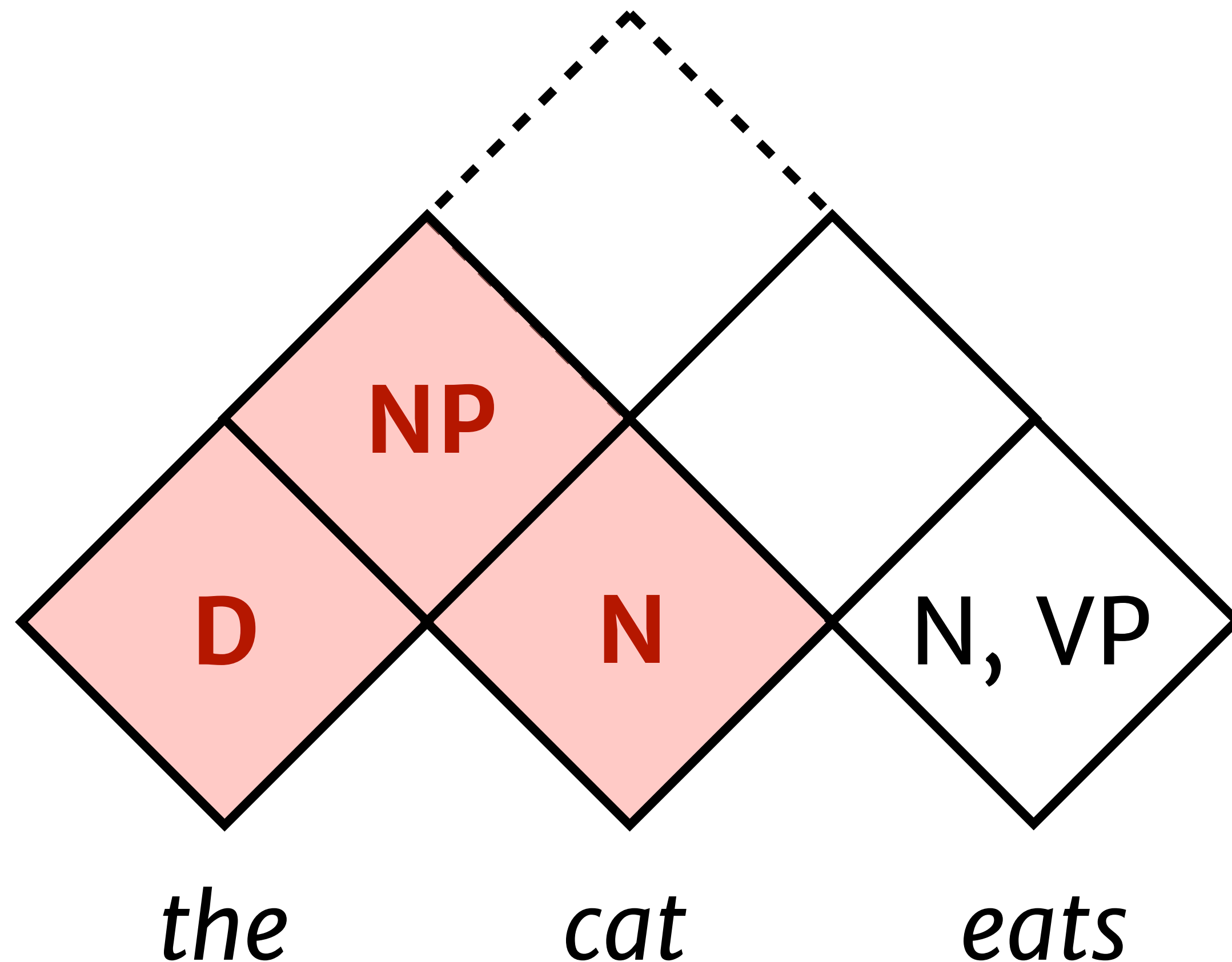
1. Fill in bottom row with NTs that can generate observed words



S → NP VP
NP → D N
N → cat | eats
VP → eats | sings
D → the | a

Warmup: the CKY algorithm

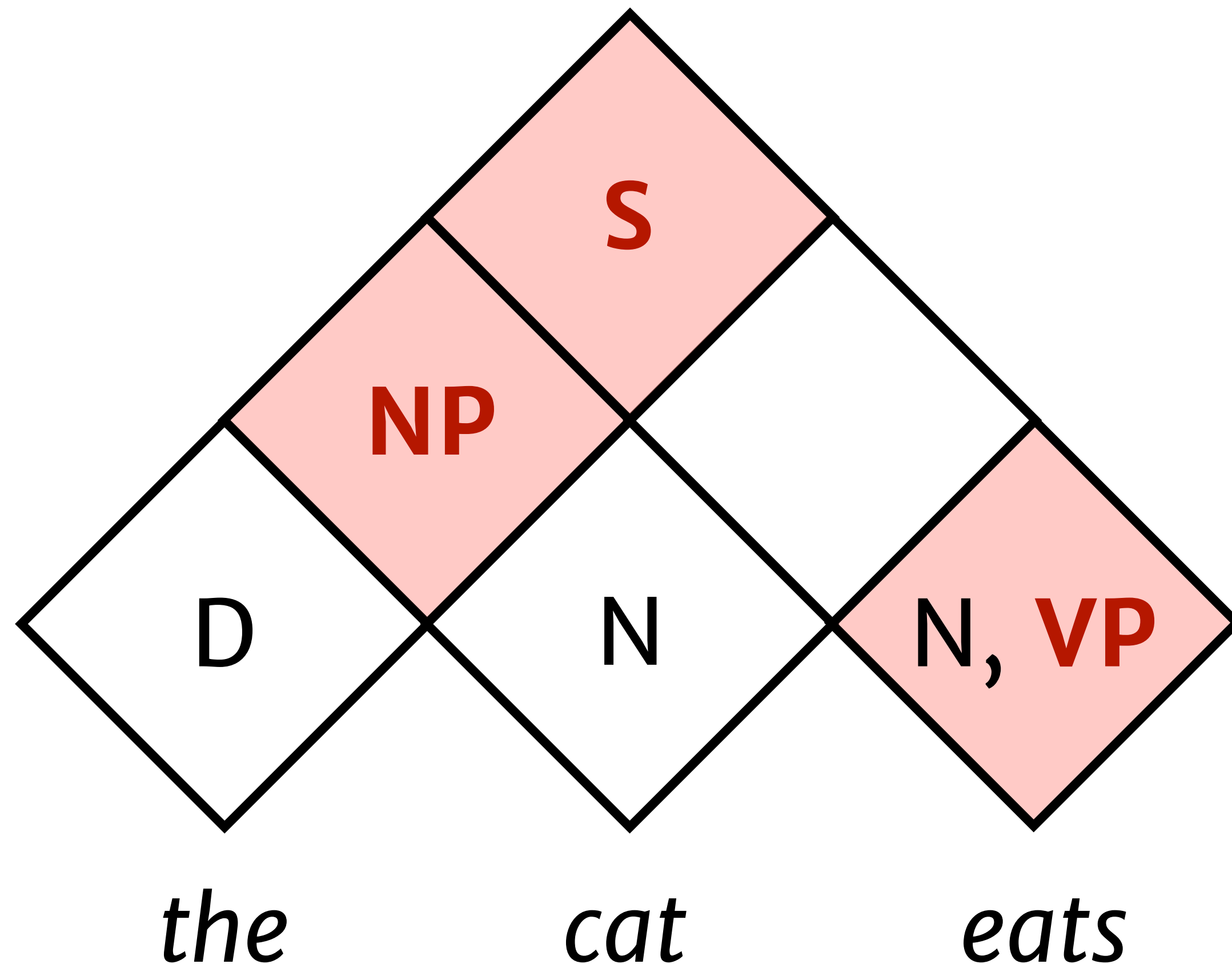
2. Fill in second row with NTs that generate a symbol in **each** child



S → NP VP
NP → **D N**
N → cat | eats
VP → eats | sings
D → the | a

Warmup: the CKY algorithm

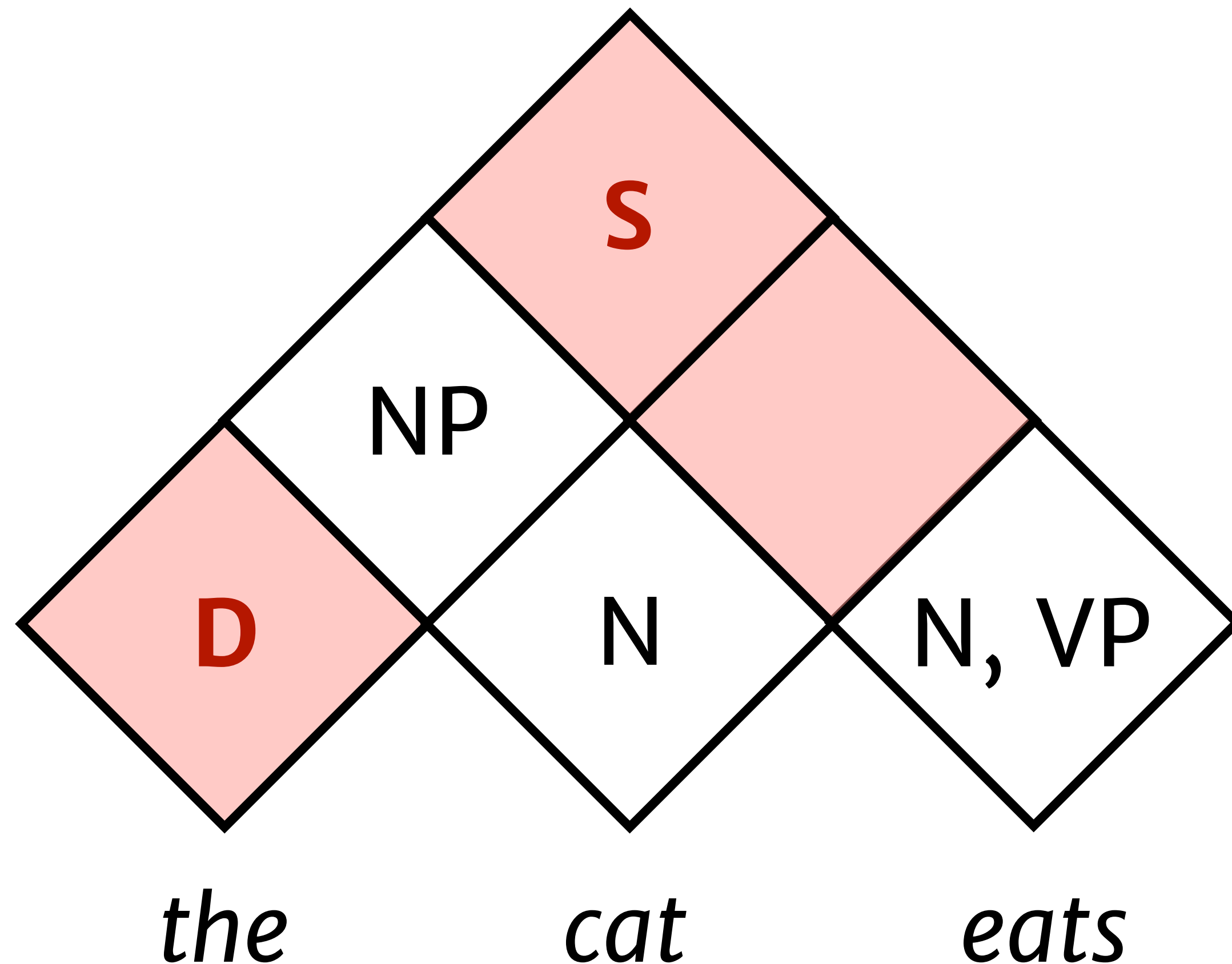
3. Fill in higher rows with NTs that generate a symbol **any pair of non-overlapping children**



S → **NP VP**
NP → D N
N → cat | eats
VP → eats | sings
D → the | a

Warmup: the CKY algorithm

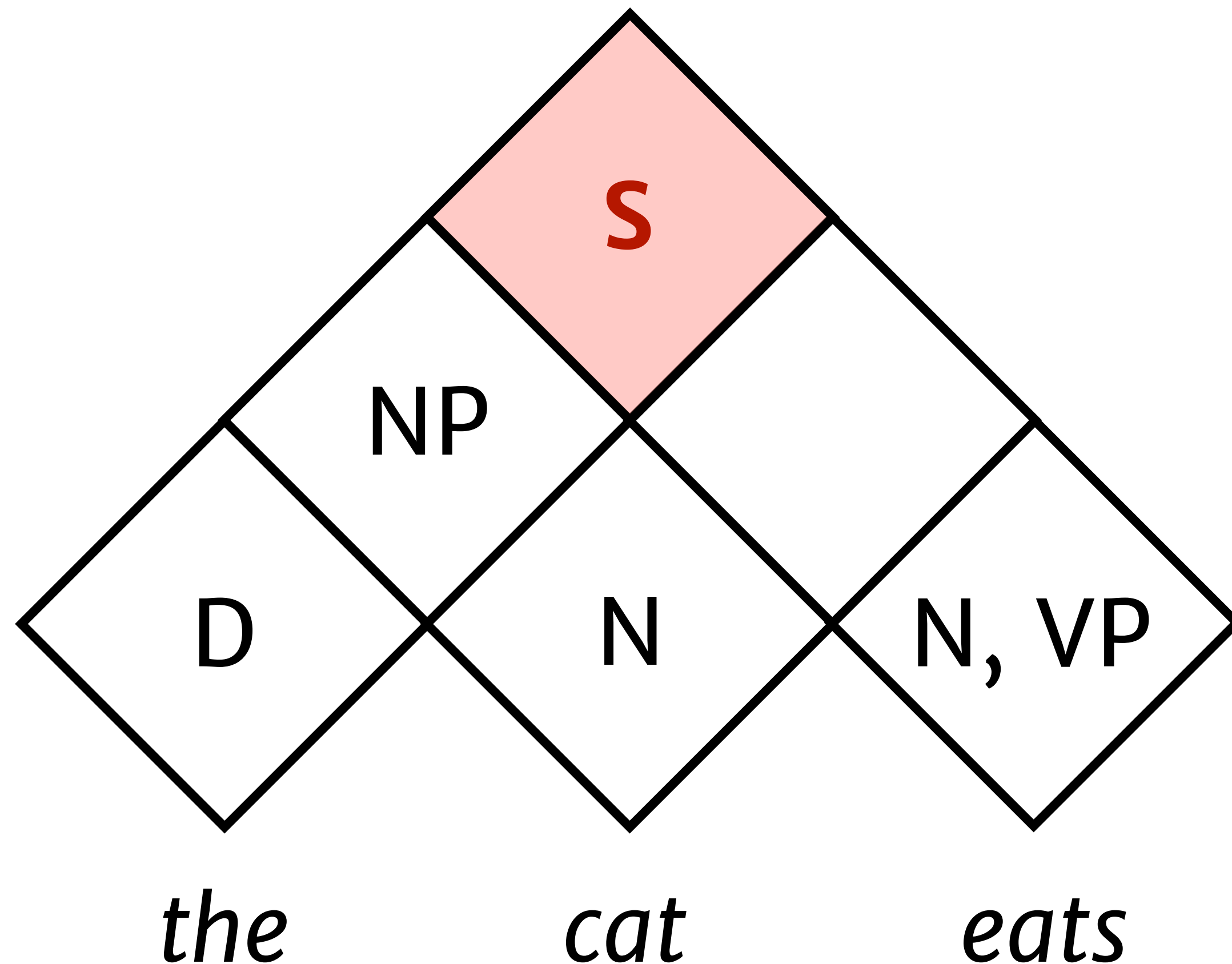
3. Fill in higher cells with NTs that generate symbols in **any pair of non-overlapping children**



S → **NP VP**
NP → D N
N → cat | eats
VP → eats | sings
D → the | a

Warmup: the CKY algorithm

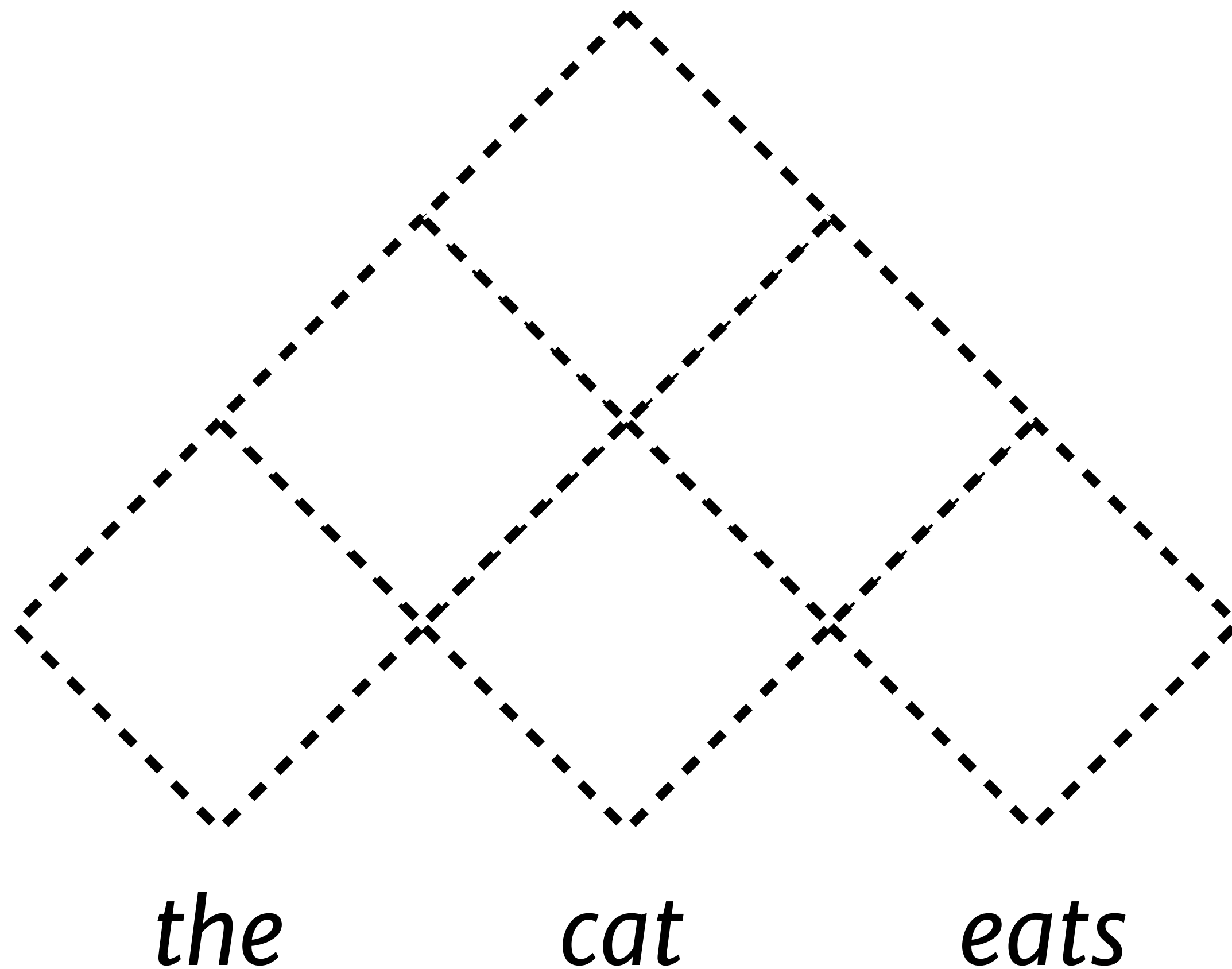
4. If the top cell contains the start symbol, the string is generated.



$S \rightarrow NP VP$
 $NP \rightarrow D N$
 $N \rightarrow \text{cat} \mid \text{eats}$
 $VP \rightarrow \text{eats} \mid \text{sings}$
 $D \rightarrow \text{the} \mid \text{a}$

Highest-scoring parse

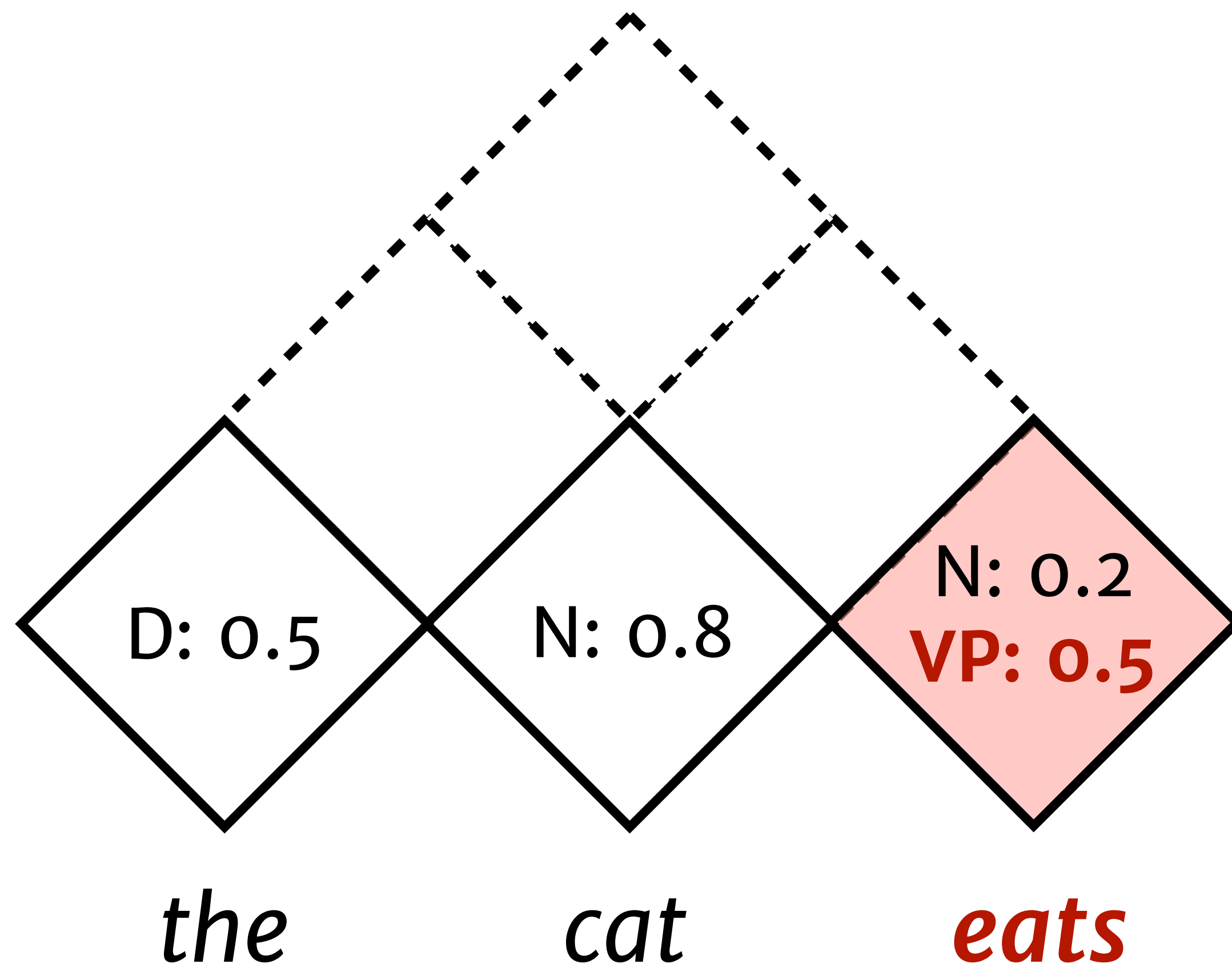
What parse assigns highest prob. to S under the PCFG G ?



S	\rightarrow	0.9		0.1
		NP	VP	NP N
NP	\rightarrow	1		
		D	N	
N	\rightarrow	0.8		0.2
		cat	eats	
VP	\rightarrow	0.5		0.5
		eats	sings	
D	\rightarrow	0.5		0.5
		the	a	

Highest-scoring parse

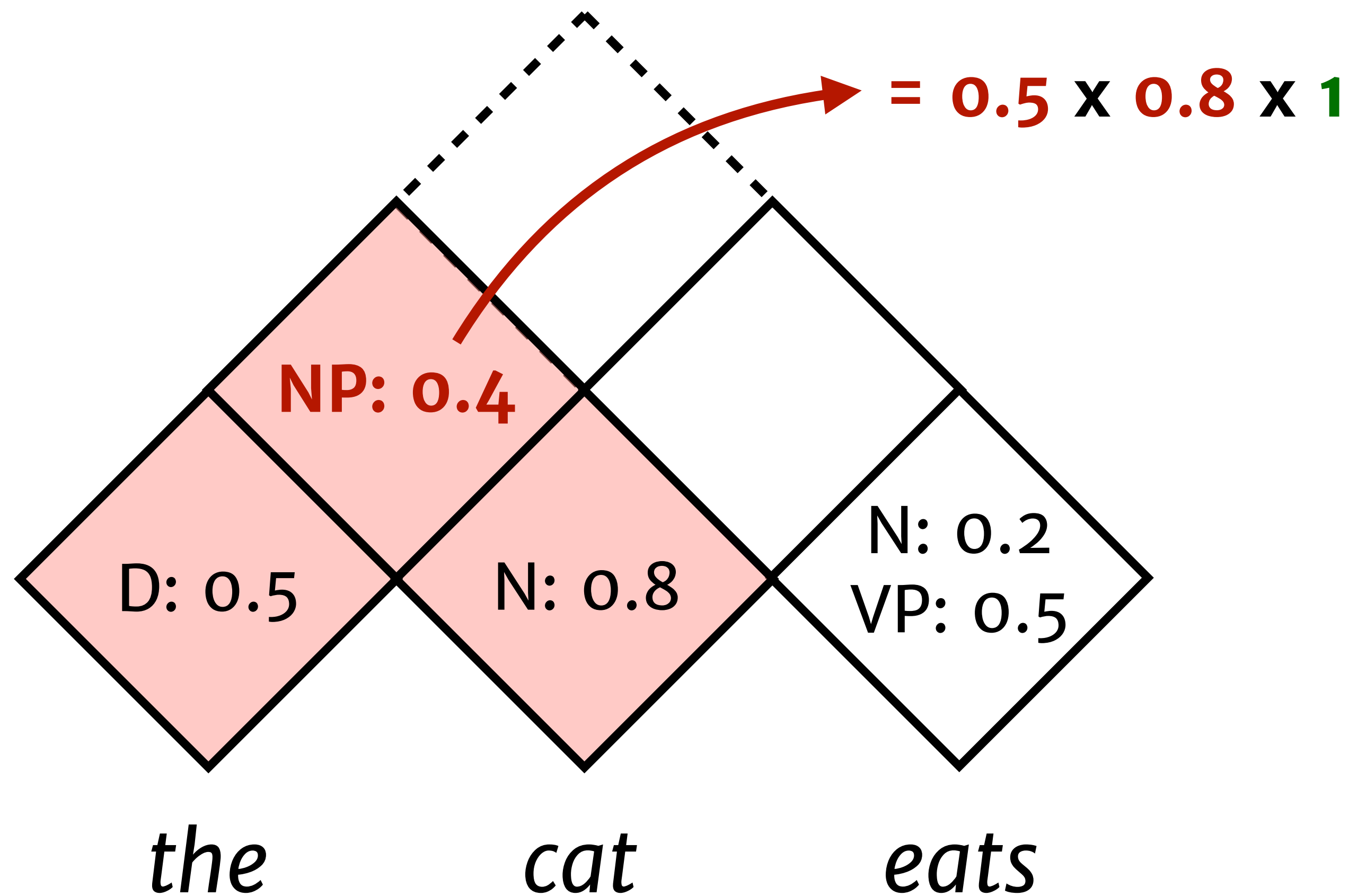
1. Fill in bottom row with prob. that each NT generates word



S	→	0.9		0.1
		NP VP		NP N
NP	→	1		
		D N		
N	→	0.8		0.2
		cat		eats
VP	→	0.5		0.5
		eats		sings
D	→	0.5		0.5
		the		a

Highest-scoring parse

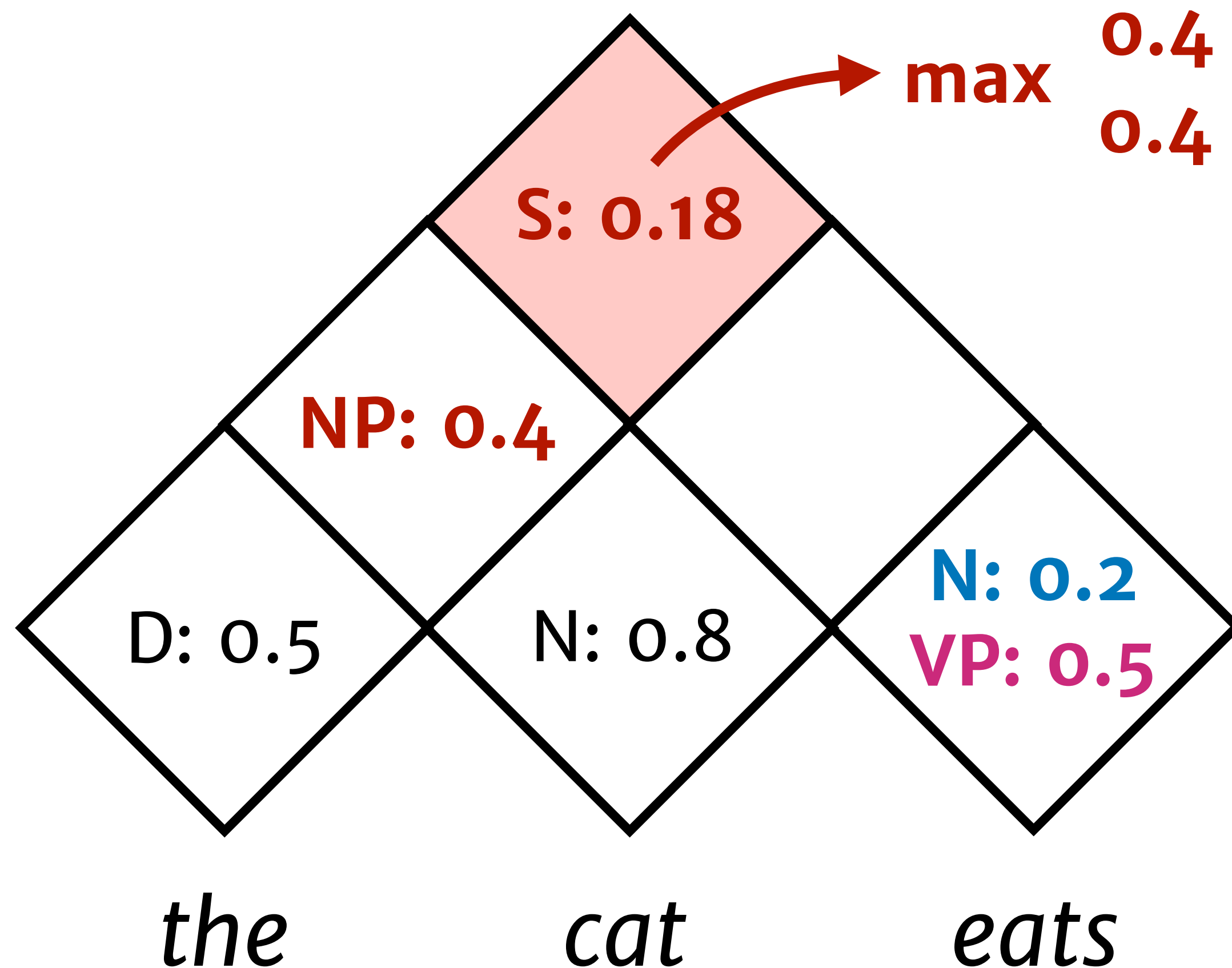
2. Fill in higher rows with highest-scoring product of child probs. times rule prob.



S	→	NP VP		NP N
		0.9		0.1
NP	→	D N		
		1		
N	→	cat		eats
		0.8		0.2
VP	→	eats		sings
		0.5		0.5
D	→	the		a
		0.5		0.5

Highest-scoring parse

2. Fill in higher rows with highest-scoring product of child probs. times rule prob.

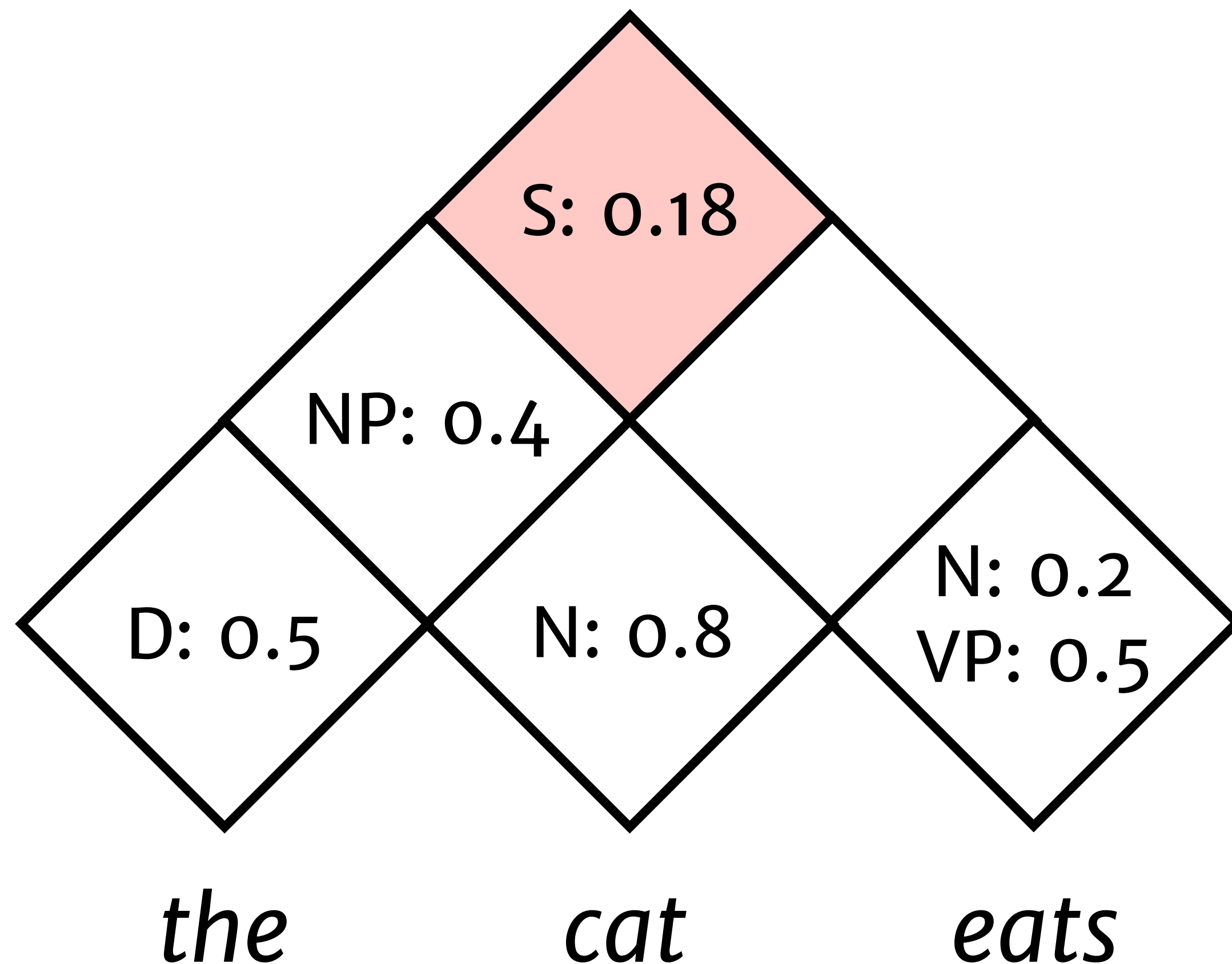


$0.4 \times 0.5 \times 0.9,$
 $0.4 \times 0.2 \times 0.1$

$S \rightarrow \text{NP VP} \mid \text{NP N}$
 $\text{NP} \rightarrow \text{D N}$
 $\text{N} \rightarrow \text{cat} \mid \text{eats}$
 $\text{VP} \rightarrow \text{eats} \mid \text{sings}$
 $\text{D} \rightarrow \text{the} \mid \text{a}$

Highest-scoring parse

3. The score for *S* in the top cell is the score of the best parse.



<i>S</i>	→	0.9		0.1
		<i>NP VP</i>		<i>NP N</i>
<i>NP</i>	→	1		
		<i>D N</i>		
<i>N</i>	→	0.8		0.2
		<i>cat</i>		<i>eats</i>
<i>VP</i>	→	0.5		0.5
		<i>eats</i>		<i>sings</i>
<i>D</i>	→	0.5		0.5
		<i>the</i>		<i>a</i>

The Viterbi algorithm for CFGs

Q: what is the **most probable** tree for a given sentence?

$$\max_T p(T, S)$$

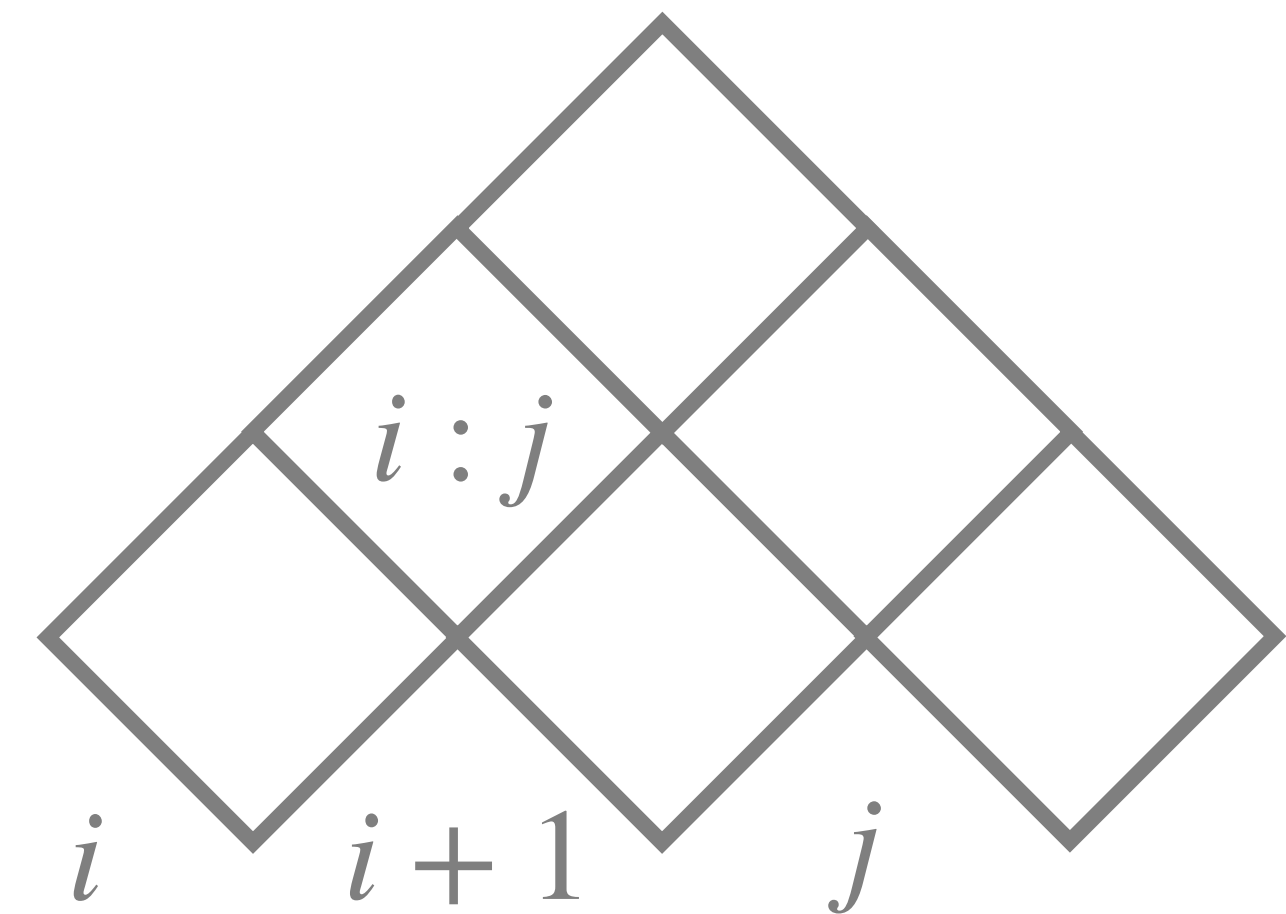
$\delta(s, i, j)$ highest-scoring tree with root s covering words $i:j$

base case:

$$\delta(s, i, i + 1) = p(s \rightarrow w_i)$$

inductive case:

$$\delta(s, i, j) = \max_{k \in [i+1, j-1]} \max_{s', s''} p(s \rightarrow s' s'') \delta(s', i, k) \delta(s'', k, j)$$



The Viterbi algorithm for CRFs

Q: what is the **most probable** assignment of tags to observations?

$$\operatorname{argmax}_Q p(Q | O)$$

$$\delta(t, j) = \max_i \delta(t - 1, i) a_{ij} b_j(o_t) \quad \delta(1, j) = \pi(j) b_j(o_1)$$

The inside algorithm for CFGs

Q: what is the marginal probability of a sentence given a tree?

$$\sum_T p(T, S)$$

$\beta(s, i, j)$ probability of all parses with root s covering words $i:j$

base case:

$$\beta(s, i, i + 1) = p(s \rightarrow w_i)$$

inductive case:

$$\beta(s, i, j) = \sum_{k \in [i+1, j-1]} \sum_{s', s''} p(s \rightarrow s' s'') \beta(s', i, k) \beta(s'', k, j)$$

Tree-structured CRFs

Instead of scores $p(A \rightarrow B \ C)$, use an arbitrary scorer

$$w^\top \phi(A, B, C, i, j)$$

(can be different in each cell & look at full sentence sentence)

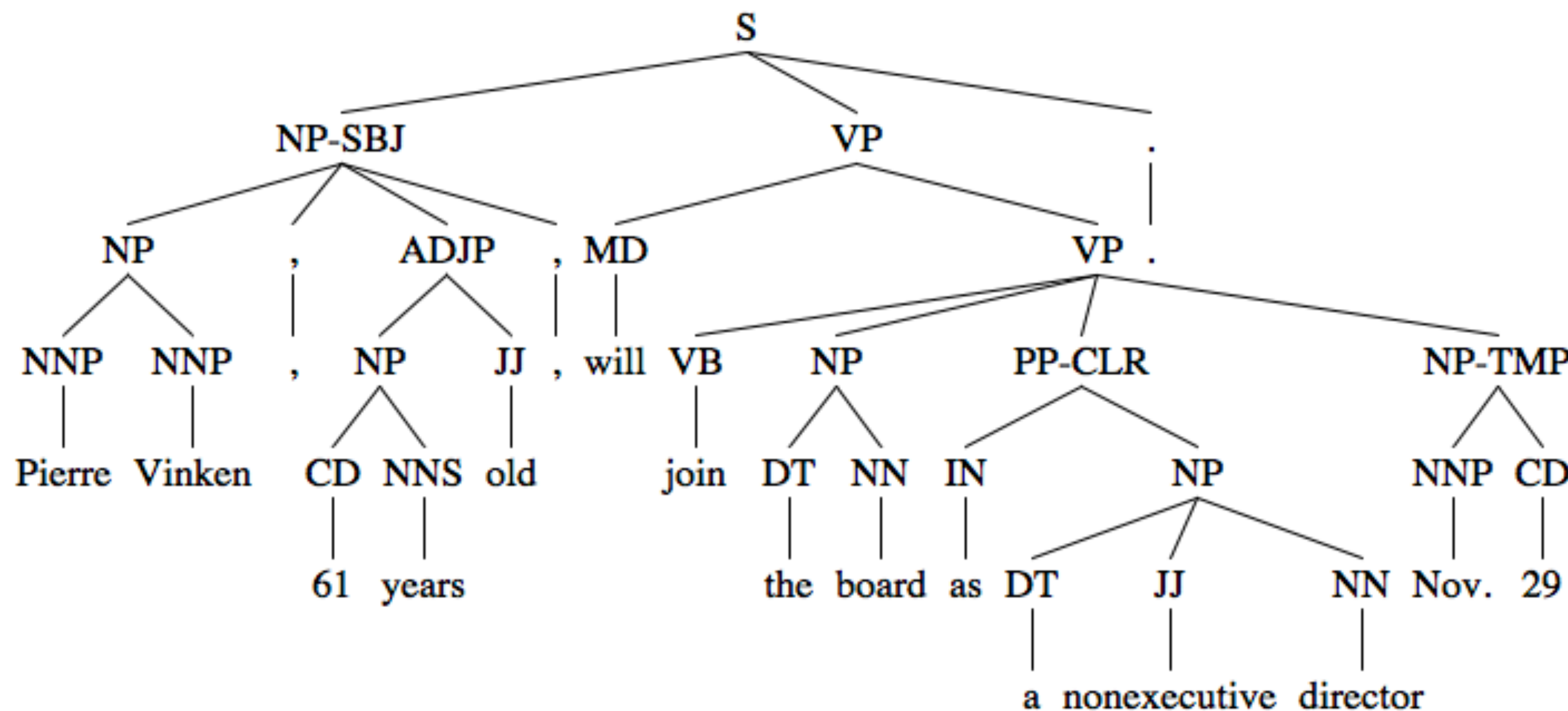
Works just like the HMM version!
 $\beta(S, 0, |S|)$ is the partition function

Learning

Supervised learning

For PCFGs—given a **treebank**,

estimate by counting:



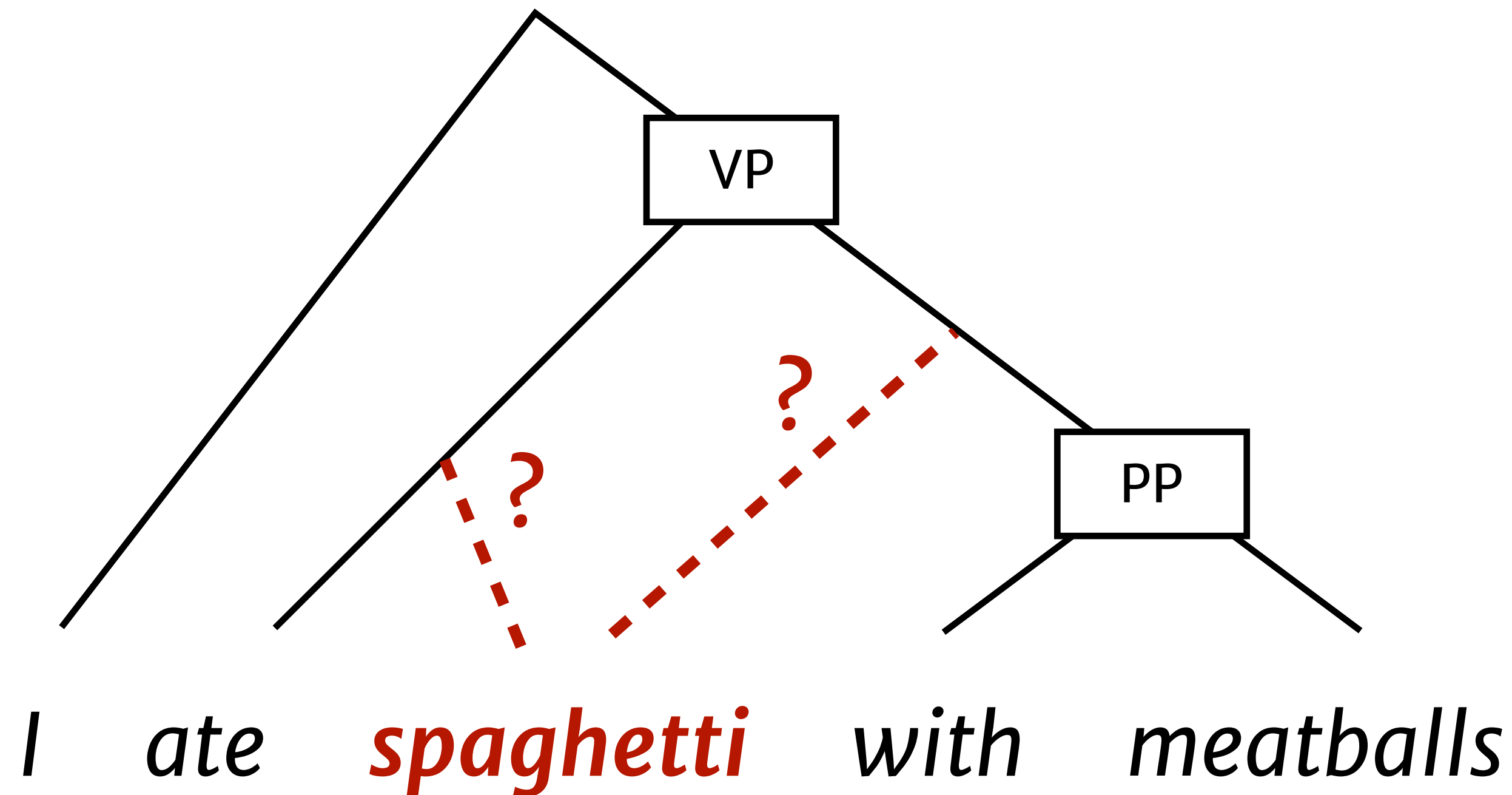
$$p(S \rightarrow NP VP)$$

$$= \frac{\#(S \rightarrow NP VP)}{\#(S)}$$

Supervised learning

For PCFGs—given a **treebank**,

estimate by counting:

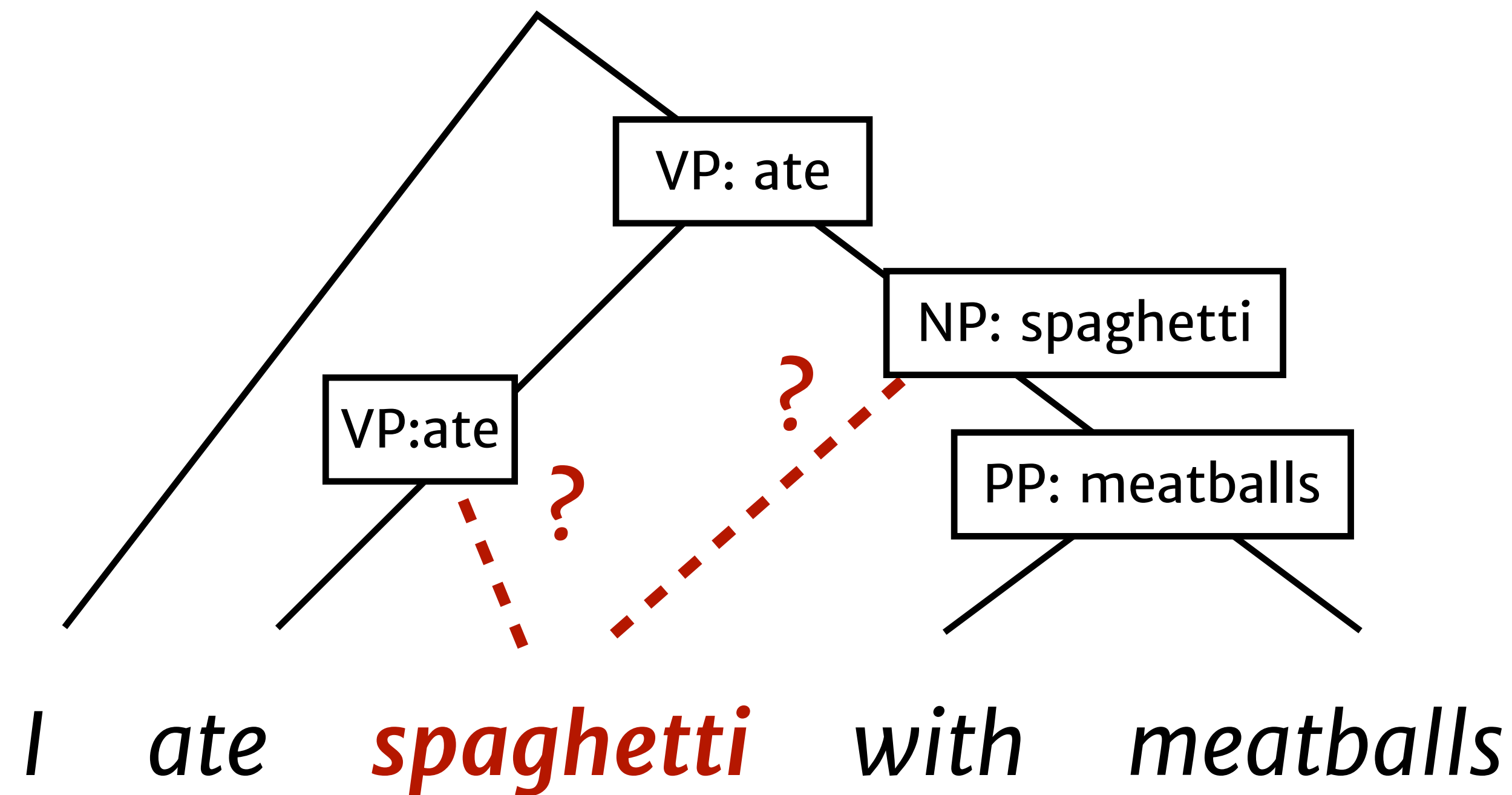


$$p(S \rightarrow NP VP) = \frac{\#(S \rightarrow NP VP)}{\#(S)}$$

This doesn't work very well: basic syntactic categories are too coarse.

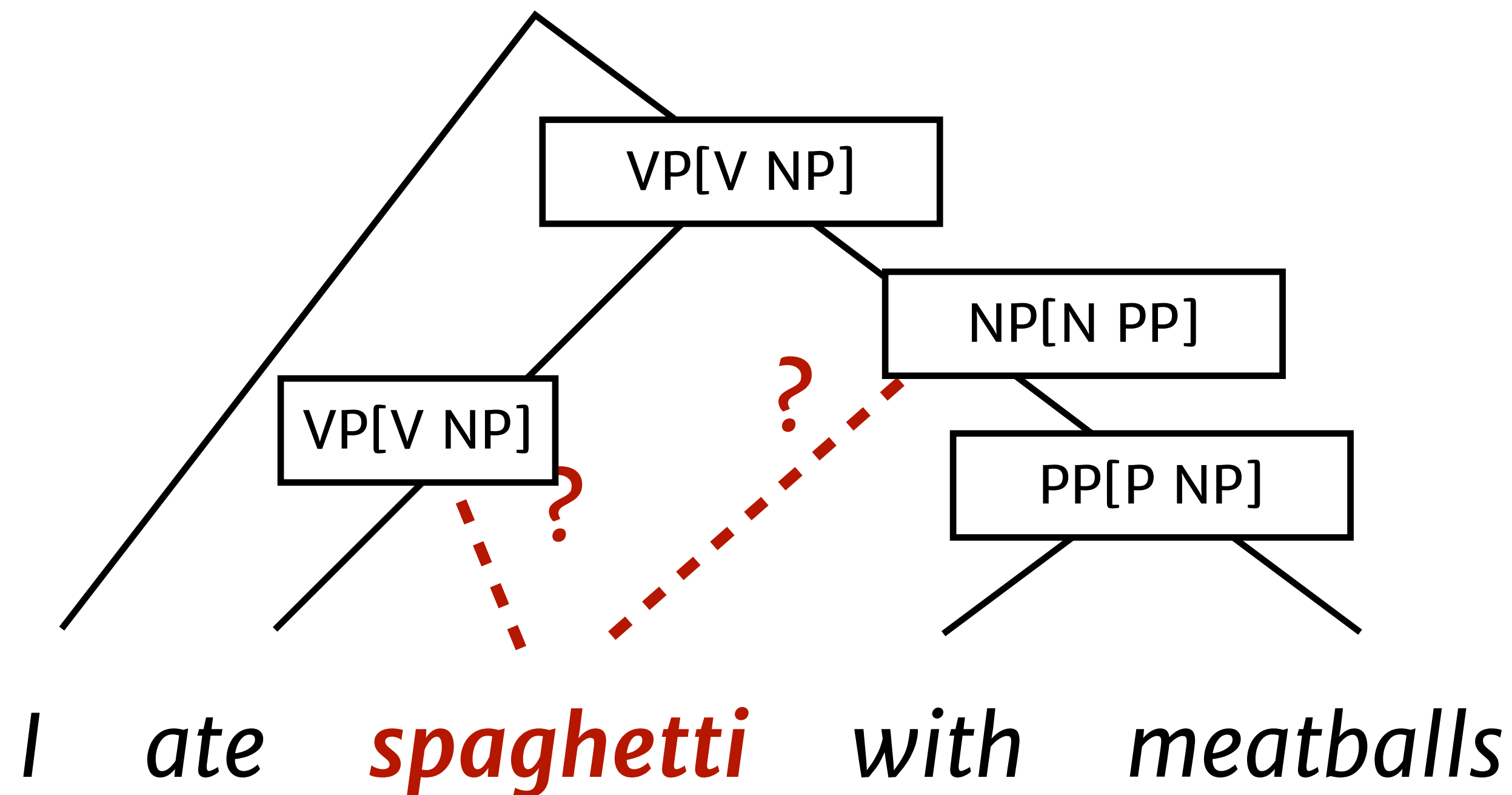
Supervised learning: lexicalization

Idea: enrich nonterminal alphabet with information about the most important **word** underneath:



Supervised learning: Markovization

Idea: enrich nonterminal alphabet with more information about the local tree structure:



Supervised learning: features & NNs

Idea: Use the CRF version

$$P(T) \propto \exp \left\{ \sum_{(A \rightarrow B \ C, i, k, j)} w^\top \phi(A, B, C, i, k, j) \right\}$$

and give ϕ features like “A = NP and $j:k$ contains *fork*”
(or make it a neural network)

Supervised parsing: what's still hard?

Error Type	Occurrences	Nodes	
		Involved	Ratio
PP Attachment	846	1455	1.7
Single Word Phrase	490	490	1.0
Clause Attachment	385	913	2.4
Adverb and Adjective Modifier Attachment	383	599	1.6
Different Label	377	754	2.0
Unary	347	349	1.0
NP Attachment	321	597	1.9
NP Internal Structure	299	352	1.2
Coordination	209	557	2.7
Unary Clause Label	185	200	1.1
VP Attachment	64	159	2.5
Parenthetical Attachment	31	74	2.4
Missing Parenthetical	12	17	1.4
Unclassified	655	734	1.1

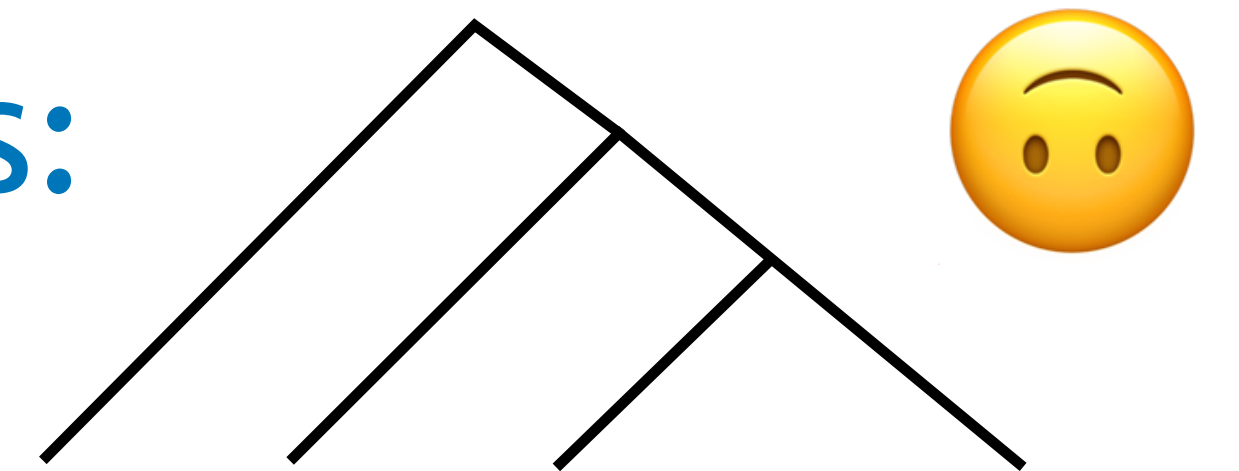
→ *spaghetti with a fork*

→ *[world oil] prices*

Unsupervised learning

Model	F_1	Training/Test PPL
Random Trees	19.5	—
Right Branching	39.5	—
Scalar PCFG (unsupervised)	< 35.0	> 350

worse than assuming every tree looks like this:



Unsupervised learning: embeddings

Model	F_1	Training/Test PPL
Random Trees	19.5	—
Right Branching	39.5	—
Scalar PCFG	< 35.0	> 350
Neural PCFG	52.6	≈ 250

“Grammar embeddings”: $p(A \rightarrow B C) \propto \exp\{v_A^\top f(v_B, v_C)\}$

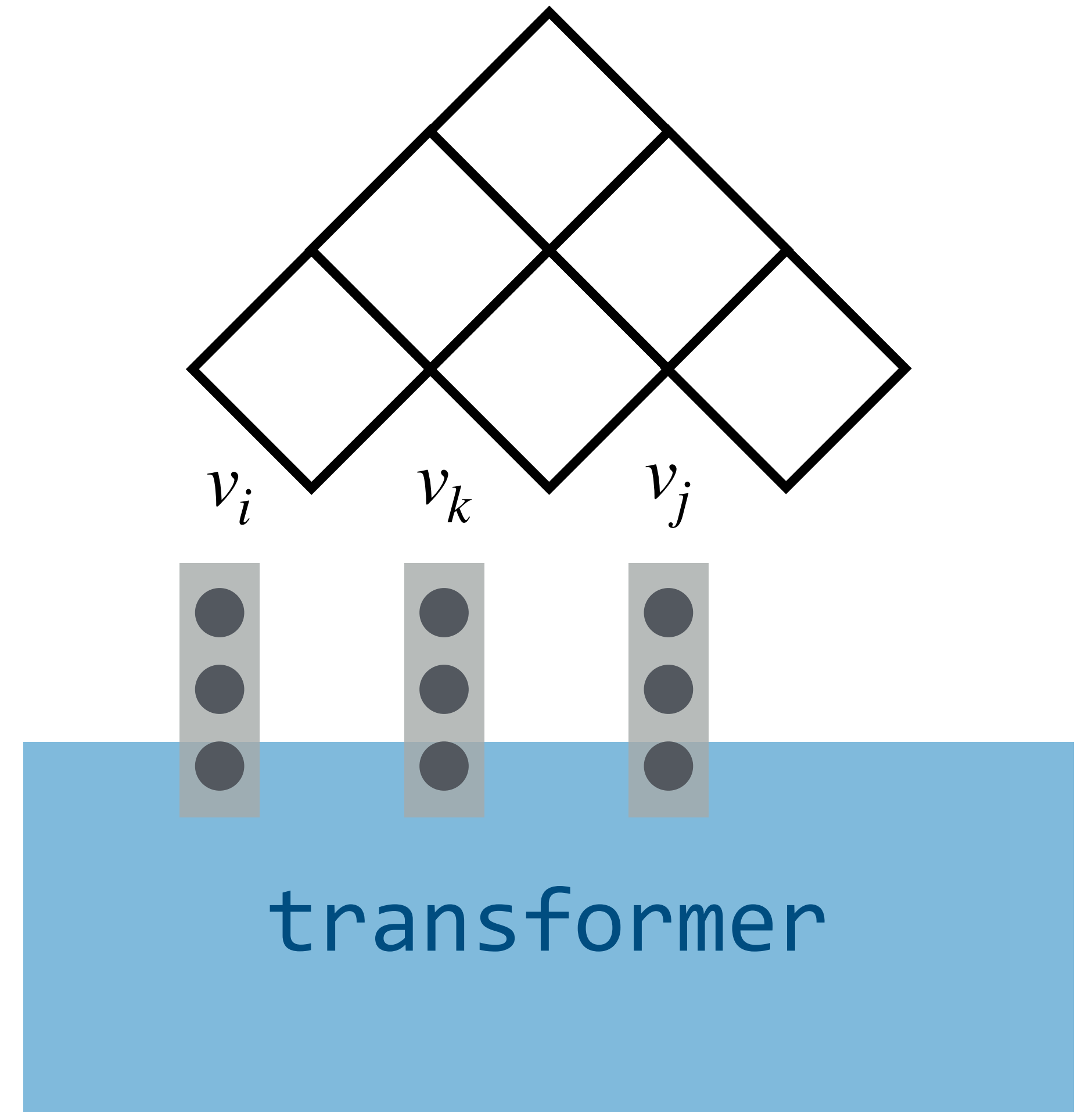
What about neural nets?

best:

$$\exp \left\{ \sum_{(A \rightarrow B \ C, i, k, j)} w^T \phi(A, B, C, v_i, v_k, v_j) \right\} / Z$$

almost as good:

train an independent cell classifier
that takes $v_{i,j,k}$ as input



Next class: advanced language modeling