

Probabilistic Sentence Processing 1

Core ideas and algorithms

Roger Levy

UC San Diego
Department of Linguistics

IPAM summer school: Probabilistic Models of Cognition
18 July 2007

What is “sentence processing”?

- ▶ *Sentence processing* is the study of how humans comprehend and produce sentences (and words within sentences, and sequences of sentences, etc.) in real time.
- ▶ We'll start with *comprehension*, and do a bit of *production* near the end.

Theoretical Desiderata

Realistic models of human sentence processing must account for:

- ▶ Robustness to arbitrary input
- ▶ Accurate disambiguation
- ▶ Inference on basis of incomplete input (Tanenhaus et al 1995, Altmann and Kamide 1999, Kaiser and Trueswell 2004)
- ▶ Processing difficulty is *differential* and *localized*

Robustness

Real linguistic input is not always totally well-formed. . .

I think when she finally came to the realization that, you know, no, I can not, I can not take care of myself.

...

I mean, for somebody who is, you know, for most of their life has, has, uh, not just merely had a farm but had ten children had a farm, ran everything because her husband was away in the coal mines.

And, you know, facing that situation, it's, it's quite a dilemma.

... but usually we come to understand it pretty well anyway.

Robustness

Real linguistic input is not always totally well-formed. . .

I think when she finally came to the realization that, you know, no, I can not, I can not take care of myself.

...

I mean, for somebody who is, you know, for most of their life has, has, uh, not just merely had a farm but had ten children had a farm, ran everything because her husband was away in the coal mines.

And, you know, facing that situation, it's, it's quite a dilemma.

(The woman is facing being put in a resting home.)

... but usually we come to understand it pretty well anyway.

Robustness

Real linguistic input is not always totally well-formed. . .

I think when she finally came to the realization that, you know, no, I can not, I can not take care of myself.

...

I mean, for somebody who is, you know, for most of their life has, has, uh, not just merely had a farm but had ten children had a farm, ran everything because her husband was away in the coal mines.

And, you know, facing that situation, it's, it's quite a dilemma.

(The woman is facing being put in a resting home.)

... but usually we come to understand it pretty well anyway.

Accurate disambiguation

Most sentences are ambiguous in ways we do not even notice:

Mary forgot the pitcher of water sitting near the stove.

Accurate disambiguation

Most sentences are ambiguous in ways we do not even notice:

Mary forgot the pitcher of water sitting near the stove.



Accurate disambiguation

Most sentences are ambiguous in ways we do not even notice:

Mary forgot the pitcher of water sitting near the stove.



Accurate disambiguation

Most sentences are ambiguous in ways we do not even notice:

Mary forgot the pitcher of water sitting near the stove.



That's probably not what you were thinking of...

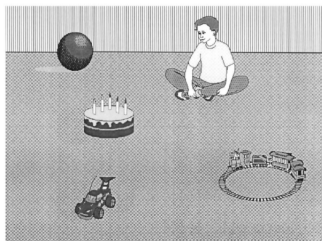
Inference on the basis of incomplete input

Comprehenders do not wait until the whole sentence has been heard to make inferences about what it means or will wind up meaning:

(Altmann and Kamide, 1999)

Inference on the basis of incomplete input

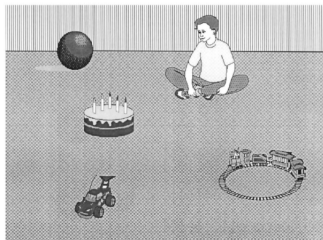
Comprehenders do not wait until the whole sentence has been heard to make inferences about what it means or will wind up meaning:



(Altmann and Kamide, 1999)

Inference on the basis of incomplete input

Comprehenders do not wait until the whole sentence has been heard to make inferences about what it means or will wind up meaning:

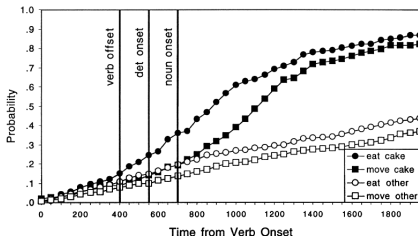
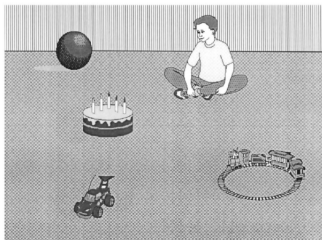


“The boy will **eat**/**move** the cake. . .”

(Altmann and Kamide, 1999)

Inference on the basis of incomplete input

Comprehenders do not wait until the whole sentence has been heard to make inferences about what it means or will wind up meaning:

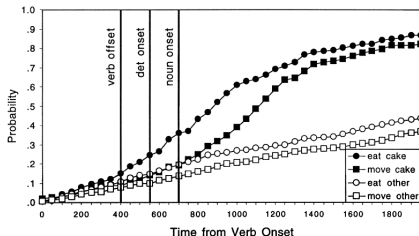
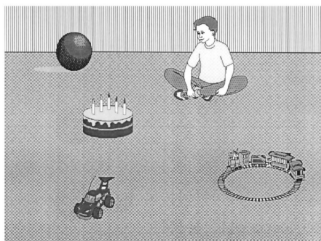


“The boy will **eat**/**move** the cake. . . .”

(Altmann and Kamide, 1999)

Inference on the basis of incomplete input

Comprehenders do not wait until the whole sentence has been heard to make inferences about what it means or will wind up meaning:



“The boy will **eat**/**move** the cake. . . .”

That is, comprehension is *incremental*

(Altmann and Kamide, 1999)

Processing difficulty is differential

Using multiple relative clauses in a sentence can make processing difficult:

This is the malt that the rat that the cat that the dog worried killed ate.

It's not the meaning of the sentence, or the use of relative clauses, that makes it hard:

This is the malt that was eaten by the rat that was killed by the cat that was worried by the dog.

Processing difficulty is differential

Using multiple relative clauses in a sentence can make processing difficult:

This is the malt that the rat that the cat that the dog worried killed ate.

It's not the meaning of the sentence, or the use of relative clauses, that makes it hard:

This is the malt that was eaten by the rat that was killed by the cat that was worried by the dog.

Processing difficulty is differential

Using multiple relative clauses in a sentence can make processing difficult:

This is the malt that the rat that the cat that the dog worried killed ate.

It's not the meaning of the sentence, or the use of relative clauses, that makes it hard:

This is the malt that was eaten by the rat that was killed by the cat that was worried by the dog.

Processing difficulty is localized

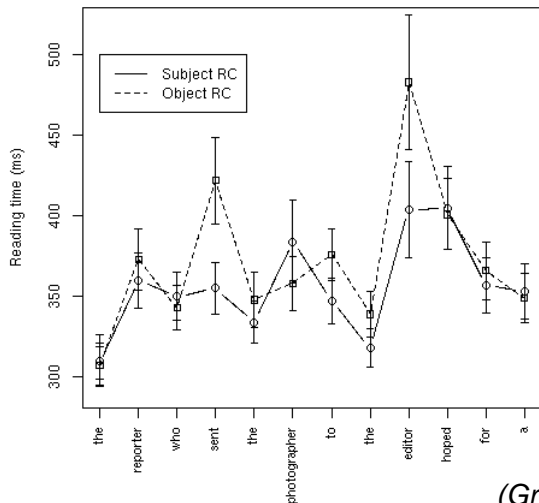
[self-paced reading demo, Example1]

(Grodner and Gibson, 2005)

Processing difficulty is localized

[self-paced reading demo, Example1]

Word-by-word reading times for sentences with different types of relative clauses (RCs)



(Grodner and Gibson, 2005)

Goal of modeling sentence processing

- ▶ We need to account for these properties (robustness, disambiguation, inference from incomplete input, and differential difficulty) together
- ▶ Probabilistic models are an interesting way to do this
- ▶ You are probably already convinced about robustness, disambiguation, and inference!
- ▶ I'll talk largely about some interesting tie-ins with differential difficulty

Goal of modeling sentence processing

- ▶ We need to account for these properties (robustness, disambiguation, inference from incomplete input, and differential difficulty) together
- ▶ Probabilistic models are an interesting way to do this
- ▶ You are probably already convinced about robustness, disambiguation, and inference!
- ▶ I'll talk largely about some interesting tie-ins with differential difficulty

Goal of modeling sentence processing

- ▶ We need to account for these properties (robustness, disambiguation, inference from incomplete input, and differential difficulty) together
- ▶ Probabilistic models are an interesting way to do this
- ▶ You are probably already convinced about robustness, disambiguation, and inference!
- ▶ I'll talk largely about some interesting tie-ins with differential difficulty

Goal of modeling sentence processing

- ▶ We need to account for these properties (robustness, disambiguation, inference from incomplete input, and differential difficulty) together
- ▶ Probabilistic models are an interesting way to do this
- ▶ You are probably already convinced about robustness, disambiguation, and inference!
- ▶ I'll talk largely about some interesting tie-ins with differential difficulty

Try to guess the next word in the sentence

- ▶ Empirically, it's been shown that more highly predictable words are read more quickly (Ehrlich and Rayner, 1981)
- ▶ Why would this be the case?

Try to guess the next word in the sentence

My brother came inside to...

- ▶ Empirically, it's been shown that more highly predictable words are read more quickly (Ehrlich and Rayner, 1981)
- ▶ Why would this be the case?

Try to guess the next word in the sentence

My brother came inside to... chat? get warm? talk? eat? rest?

- ▶ Empirically, it's been shown that more highly predictable words are read more quickly (Ehrlich and Rayner, 1981)
- ▶ Why would this be the case?

Try to guess the next word in the sentence

My brother came inside to... chat? get warm? talk? eat? rest?
The children went outside to...

- ▶ Empirically, it's been shown that more highly predictable words are read more quickly (Ehrlich and Rayner, 1981)
- ▶ Why would this be the case?

Try to guess the next word in the sentence

My brother came inside to... chat? get warm? talk? eat? rest?

The children went outside to... play

- ▶ Empirically, it's been shown that more highly predictable words are read more quickly (Ehrlich and Rayner, 1981)
- ▶ Why would this be the case?

Try to guess the next word in the sentence

My brother came inside to... chat? get warm? talk? eat? rest?

The children went outside to... play

- ▶ Empirically, it's been shown that more highly predictable words are read more quickly (Ehrlich and Rayner, 1981)
- ▶ Why would this be the case?

Try to guess the next word in the sentence

My brother came inside to... chat? get warm? talk? eat? rest?

The children went outside to... play

- ▶ Empirically, it's been shown that more highly predictable words are read more quickly (Ehrlich and Rayner, 1981)
- ▶ Why would this be the case?

Surprisal as a possible metric for processing difficulty

An event's surprisal is simply its negative log conditional probability

$$\log \frac{1}{P(x|\text{Context})}$$

Intuitively, this is a measure of the amount of information contained in the event

Three proposals for surprisal as a measure of processing difficulty/time:

- ▶ Surprisal of a word as primitive measure of processing (Attneave, 1959; Hale, 2001)
- ▶ KL divergence as size of update that the word induces for distribution over interpretations of input (Levy, 2005, 2007)
 - ▶ independently proposed as a measure of surprise in visual scene perception (Itti and Baldi, 2005)
- ▶ Surprisal as an optimal solution to the speed/resource tradeoff in language comprehension (Smith, 2006)

Probabilistic grammars for estimating surprisal

- ▶ Comprehenders' expectations about upcoming words should reflect structural distributional regularities of the language
- ▶ Hence, probabilistic grammars (e.g., PCFGs) are a good candidate

PCFG review

a man arrived yesterday

0.3 S → S CC S

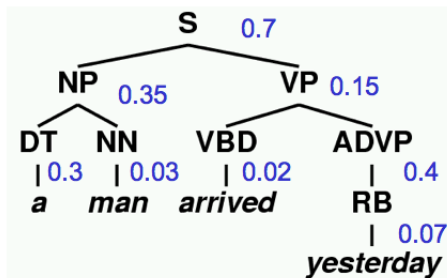
0.7 S → NP VP

0.35 NP → DT NN

0.15 VP → VBD ADVP

0.4 ADVP → RB

...



Total probability: $0.7 \cdot 0.35 \cdot 0.15 \cdot 0.3 \cdot 0.03 \cdot 0.02 \cdot 0.4 \cdot 0.07 = 1.85 \times 10^{-7}$

Algorithms by Lafferty and Jelinek (1992),
Stolcke (1995) give us $p_i(w)$ from a PCFG

PCFG review (2)

- ▶ The *probability of a string* $w_1\dots n$ is the sum of the probabilities of all trees whose yield **is** $w_1\dots n$
- ▶ The *probability of a string prefix* $w_1\dots j$ is the sum of the probabilities of all trees whose yield **begins with** $w_1\dots j$
- ▶ If we had the probabilities of two string prefixes $w_1\dots j-1$ and $w_1\dots j$, we could calculate the conditional probability $P(w_j|w_1\dots j-1)$ as their ratio.

Inference over infinite tree sets

Consider the following noun-phrase grammar:

$\frac{2}{3}$ NP \rightarrow Det N
 $\frac{1}{3}$ NP \rightarrow NP PP
1 PP \rightarrow P NP

1 Det \rightarrow the
 $\frac{2}{3}$ N \rightarrow dog
 $\frac{1}{3}$ N \rightarrow cat
1 P \rightarrow near

Question: given a sentence starting with

the...

what is the probability that the next word is *dog*?

Intuitively, the answers to this question should be

$$P(\text{dog}|\text{the}) = \frac{2}{3}$$

because the second word HAS to be either *dog* or *cat*.

Inference over infinite tree sets

Consider the following noun-phrase grammar:

$\frac{2}{3}$ NP \rightarrow Det N
 $\frac{1}{3}$ NP \rightarrow NP PP
1 PP \rightarrow P NP

1 Det \rightarrow the
 $\frac{2}{3}$ N \rightarrow dog
 $\frac{1}{3}$ N \rightarrow cat
1 P \rightarrow near

Question: given a sentence starting with

the...

what is the probability that the next word is *dog*?

Intuitively, the answers to this question should be

$$P(\text{dog}|\text{the}) = \frac{2}{3}$$

because the second word HAS to be either *dog* or *cat*.

Inference over infinite tree sets

Consider the following noun-phrase grammar:

$\frac{2}{3}$ NP \rightarrow Det N
 $\frac{1}{3}$ NP \rightarrow NP PP
1 PP \rightarrow P NP

1 Det \rightarrow the
 $\frac{2}{3}$ N \rightarrow dog
 $\frac{1}{3}$ N \rightarrow cat
1 P \rightarrow near

Question: given a sentence starting with

the...

what is the probability that the next word is *dog*?

Intuitively, the answers to this question should be

$$P(\text{dog}|\text{the}) = \frac{2}{3}$$

because the second word HAS to be either *dog* or *cat*.

Inference over infinite tree sets

Consider the following noun-phrase grammar:

$\frac{2}{3}$ NP \rightarrow Det N
 $\frac{1}{3}$ NP \rightarrow NP PP
1 PP \rightarrow P NP

1 Det \rightarrow the
 $\frac{2}{3}$ N \rightarrow dog
 $\frac{1}{3}$ N \rightarrow cat
1 P \rightarrow near

Question: given a sentence starting with

the...

what is the probability that the next word is *dog*?

Intuitively, the answers to this question should be

$$P(\text{dog}|\text{the}) = \frac{2}{3}$$

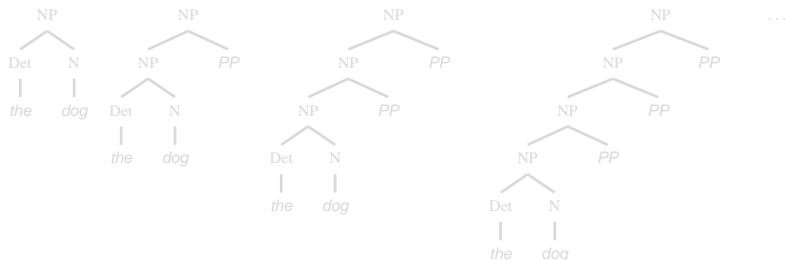
because the second word HAS to be either *dog* or *cat*.

Inference over infinite tree sets (2)

$\frac{2}{3} \rightarrow 1$
NP \rightarrow Det N
NP \rightarrow NP PP
1 PP \rightarrow P NP

1 Det \rightarrow the
 $\frac{2}{3} \rightarrow 1$
N \rightarrow dog
 $\frac{1}{3} \rightarrow 1$
N \rightarrow cat
1 P \rightarrow near

- ▶ We “should” just enumerate the trees that cover *the dog* ..., and divide their total probability by that of *the ...*
- ▶ ... but there are infinitely many trees.

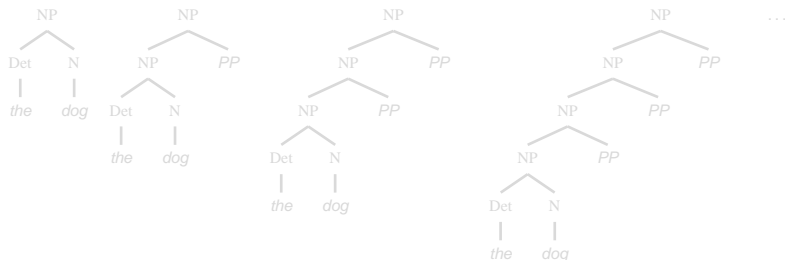


Inference over infinite tree sets (2)

$\frac{2}{3}$ NP \rightarrow Det N
 $\frac{1}{3}$ NP \rightarrow NP PP
1 PP \rightarrow P NP

1 Det \rightarrow the
 $\frac{2}{3}$ N \rightarrow dog
 $\frac{1}{3}$ N \rightarrow cat
1 P \rightarrow near

- ▶ We “should” just enumerate the trees that cover *the dog* ..., and divide their total probability by that of *the ...*
- ▶ ... but there are infinitely many trees.

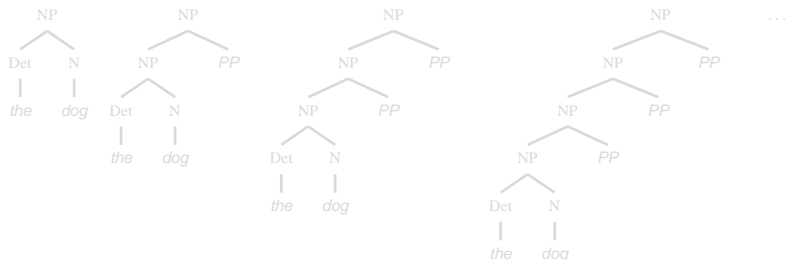


Inference over infinite tree sets (2)

$\frac{2}{3}$ NP \rightarrow Det N
 $\frac{1}{3}$ NP \rightarrow NP PP
1 PP \rightarrow P NP

1 Det \rightarrow the
 $\frac{2}{3}$ N \rightarrow dog
 $\frac{1}{3}$ N \rightarrow cat
1 P \rightarrow near

- ▶ We “should” just enumerate the trees that cover *the dog* ..., and divide their total probability by that of *the* ...
- ▶ ...but there are infinitely many trees.

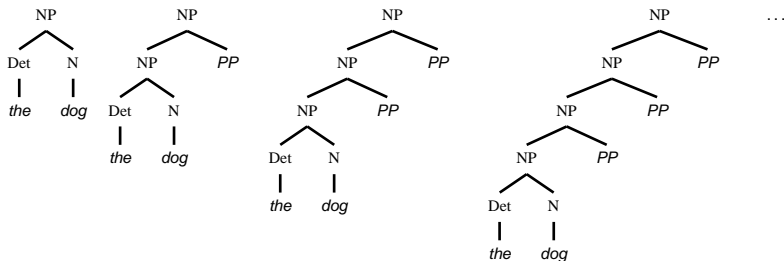


Inference over infinite tree sets (2)

$\frac{2}{3} \rightarrow 1$
2 NP \rightarrow Det N
3 NP \rightarrow NP PP
1 PP \rightarrow P NP

1 Det \rightarrow the
 $\frac{2}{3}$ N \rightarrow dog
 $\frac{1}{3}$ N \rightarrow cat
1 P \rightarrow near

- ▶ We “should” just enumerate the trees that cover *the dog* ..., and divide their total probability by that of *the* ...
- ▶ ...but there are infinitely many trees.

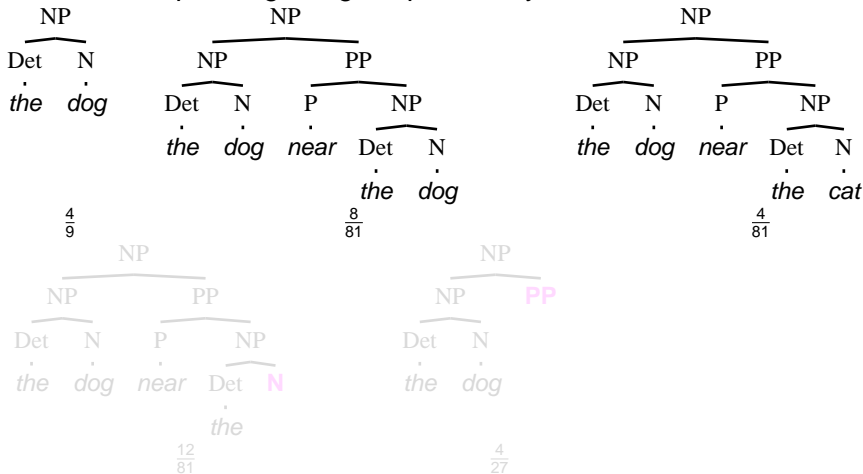


$\frac{2}{3} \rightarrow \frac{1}{3}$
 NP \rightarrow Det N
 NP \rightarrow NP PP
 1 PP \rightarrow P NP

1 Det \rightarrow the
 $\frac{2}{3}$ N \rightarrow dog
 $\frac{1}{3}$ N \rightarrow cat
 1 P \rightarrow near

Shortcut 1: you can think of a *partial* tree as marginalizing over all completions of the partial tree.

It has a corresponding marginal probability in the PCFG.

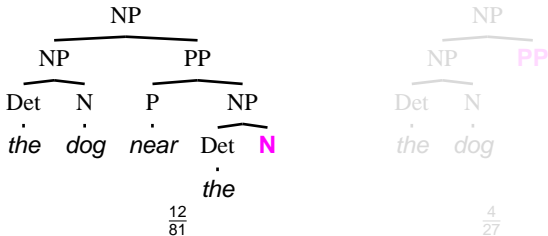
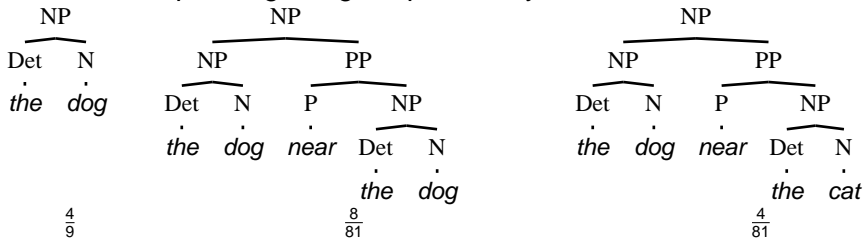


$\frac{2}{3}$ NP \rightarrow Det N
 $\frac{1}{3}$ NP \rightarrow NP PP
 1 PP \rightarrow P NP

1 Det \rightarrow the
 $\frac{2}{3}$ N \rightarrow dog
 $\frac{1}{3}$ N \rightarrow cat
 1 P \rightarrow near

Shortcut 1: you can think of a *partial* tree as marginalizing over all completions of the partial tree.

It has a corresponding marginal probability in the PCFG.

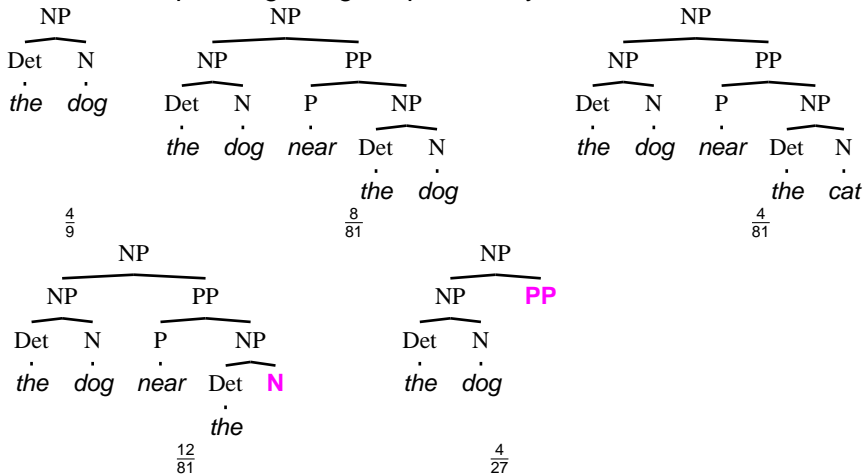


$\frac{2}{3}$ NP \rightarrow Det N
 $\frac{1}{3}$ NP \rightarrow NP PP
 1 PP \rightarrow P NP

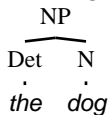
1 Det \rightarrow the
 $\frac{2}{3}$ N \rightarrow dog
 $\frac{1}{3}$ N \rightarrow cat
 1 P \rightarrow near

Shortcut 1: you can think of a *partial* tree as marginalizing over all completions of the partial tree.

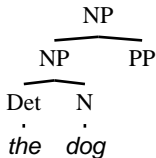
It has a corresponding marginal probability in the PCFG.



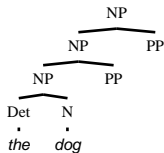
Problem 2: there are still an infinite number of incomplete trees covering a partial input.



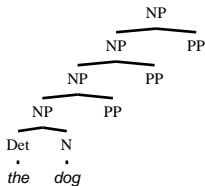
$$\frac{4}{9}$$



$$\frac{4}{27}$$



$$\frac{4}{81}$$



$$\frac{4}{243}$$

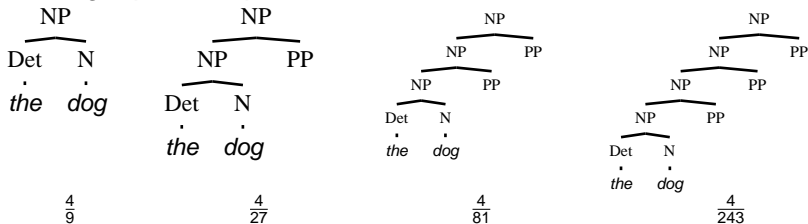
BUT! These tree probabilities form a geometric series:

$$\begin{aligned}
 P(\text{the dog} \dots) &= \frac{4}{9} + \frac{4}{27} + \frac{4}{81} + \frac{4}{243} + \dots \\
 &= \frac{4}{9} \prod_{i=0}^{\infty} \left(\frac{1}{3}\right)^i \\
 &= \frac{2}{3}
 \end{aligned}$$

... which matches the original rule probability

$$\frac{2}{3} N \rightarrow \text{dog}$$

Problem 2: there are still an infinite number of incomplete trees covering a partial input.



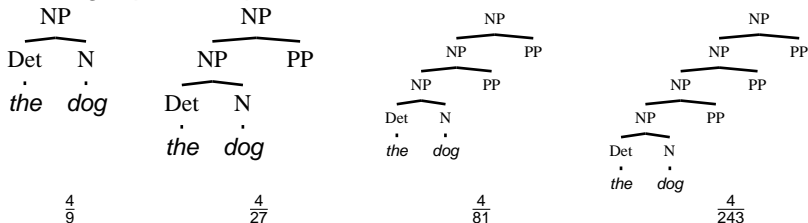
BUT! These tree probabilities form a geometric series:

$$\begin{aligned}
 P(\text{the dog} \dots) &= \frac{4}{9} + \frac{4}{27} + \frac{4}{81} + \frac{4}{243} + \dots \\
 &= \frac{4}{9} \prod_{i=0}^{\infty} \left(\frac{1}{3}\right)^i \\
 &= \frac{2}{3}
 \end{aligned}$$

... which matches the original rule probability

$$\frac{2}{3} N \rightarrow \text{dog}$$

Problem 2: there are still an infinite number of incomplete trees covering a partial input.



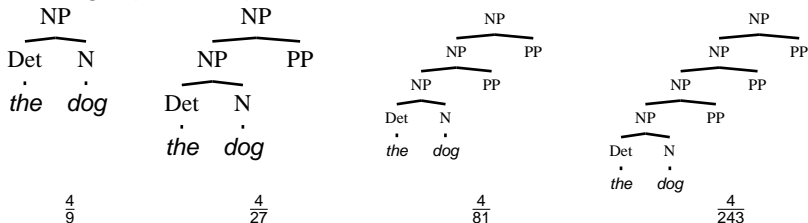
BUT! These tree probabilities form a geometric series:

$$\begin{aligned}
 P(\text{the dog} \dots) &= \frac{4}{9} + \frac{4}{27} + \frac{4}{81} + \frac{4}{243} + \dots \\
 &= \frac{4}{9} \prod_{i=0}^{\infty} \left(\frac{1}{3}\right)^i \\
 &= \frac{2}{3}
 \end{aligned}$$

... which matches the original rule probability

$$\frac{2}{3} N \rightarrow \text{dog}$$

Problem 2: there are still an infinite number of incomplete trees covering a partial input.



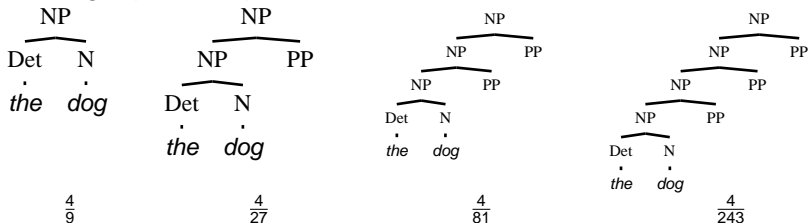
BUT! These tree probabilities form a geometric series:

$$\begin{aligned}
 P(\text{the dog} \dots) &= \frac{4}{9} + \frac{4}{27} + \frac{4}{81} + \frac{4}{243} + \dots \\
 &= \frac{4}{9} \prod_{i=0}^{\infty} \left(\frac{1}{3}\right)^i \\
 &= \frac{2}{3}
 \end{aligned}$$

... which matches the original rule probability

$$\frac{2}{3} N \rightarrow \text{dog}$$

Problem 2: there are still an infinite number of incomplete trees covering a partial input.



BUT! These tree probabilities form a geometric series:

$$\begin{aligned}
 P(\text{the dog} \dots) &= \frac{4}{9} + \frac{4}{27} + \frac{4}{81} + \frac{4}{243} + \dots \\
 &= \frac{4}{9} \prod_{i=0}^{\infty} \left(\frac{1}{3}\right)^i \\
 &= \frac{2}{3}
 \end{aligned}$$

... which matches the original rule probability

$$\frac{2}{3} N \rightarrow \text{dog}$$

Generalizing the geometric series induced by rule recursion

In general, these infinite tree sets arise due to *left recursion* in a probabilistic grammar

$$A \rightarrow B \alpha$$

$$B \rightarrow A \beta$$

We can formulate a stochastic *left-corner matrix* of transitions between categories:

$$P_L = \begin{array}{c|cccc} & A & B & \dots & K \\ \hline A & 0.3 & 0.7 & \dots & 0 \\ B & 0.1 & 0.1 & \dots & 0.2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ K & 0.2 & 0.1 & \dots & 0.2 \end{array}$$

and solve for its closure $R_L = (I - P_L)^{-1}$.
(Stolcke, 1995)

Generalizing the geometric series induced by rule recursion

In general, these infinite tree sets arise due to *left recursion* in a probabilistic grammar

$$A \rightarrow B \alpha$$

$$B \rightarrow A \beta$$

We can formulate a stochastic *left-corner matrix* of transitions between categories:

$$P_L = \begin{array}{c|cccc} & A & B & \dots & K \\ \hline A & 0.3 & 0.7 & \dots & 0 \\ B & 0.1 & 0.1 & \dots & 0.2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ K & 0.2 & 0.1 & \dots & 0.2 \end{array}$$

and solve for its closure $R_L = (I - P_L)^{-1}$.
(Stolcke, 1995)

Generalizing the geometric series induced by rule recursion

In general, these infinite tree sets arise due to *left recursion* in a probabilistic grammar

$$A \rightarrow B \alpha$$

$$B \rightarrow A \beta$$

We can formulate a stochastic *left-corner matrix* of transitions between categories:

$$P_L = \begin{array}{c|cccc} & A & B & \dots & K \\ \hline A & 0.3 & 0.7 & \dots & 0 \\ B & 0.1 & 0.1 & \dots & 0.2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ K & 0.2 & 0.1 & \dots & 0.2 \end{array}$$

and solve for its closure $R_L = (I - P_L)^{-1}$.
(Stolcke, 1995)

Generalizing the geometric series

$\frac{2}{3}$	NP \rightarrow Det N
$\frac{1}{3}$	NP \rightarrow NP PP
1	PP \rightarrow P NP

1	Det \rightarrow the
$\frac{2}{3}$	N \rightarrow dog
$\frac{1}{3}$	N \rightarrow cat
1	P \rightarrow near

The closure of our left-corner matrix is

$$\begin{array}{l} \text{NP} \\ \text{PP} \\ \text{Det} \\ \text{N} \\ \text{P} \end{array} \left(\begin{array}{ccccc} 1.5 & 0 & 1.0 & 0 & 0 \\ 0 & 1.0 & 0 & 0 & 1.0 \\ 0 & 0 & 1.0 & 0 & 0 \\ 0 & 0 & 0 & 1.0 & 0 \\ 0 & 0 & 0 & 0 & 1.0 \end{array} \right)$$

Note that the $\frac{3}{2}$ “bonus” accrued for left-recursion of NPs appears in the (NP,NP) cell of the matrix

Generalizing the geometric series

$\frac{2}{3}$	NP \rightarrow Det N
$\frac{1}{3}$	NP \rightarrow NP PP
1	PP \rightarrow P NP

1	Det \rightarrow the
$\frac{2}{3}$	N \rightarrow dog
$\frac{1}{3}$	N \rightarrow cat
1	P \rightarrow near

The closure of our left-corner matrix is

$$\begin{array}{l} \text{NP} \\ \text{PP} \\ \text{Det} \\ \text{N} \\ \text{P} \end{array} \left(\begin{array}{ccccc} 1.5 & 0 & 1.0 & 0 & 0 \\ 0 & 1.0 & 0 & 0 & 1.0 \\ 0 & 0 & 1.0 & 0 & 0 \\ 0 & 0 & 0 & 1.0 & 0 \\ 0 & 0 & 0 & 0 & 1.0 \end{array} \right)$$

Note that the $\frac{3}{2}$ “bonus” accrued for left-recursion of NPs appears in the (NP, NP) cell of the matrix

Efficient incremental parsing: the probabilistic Earley algorithm

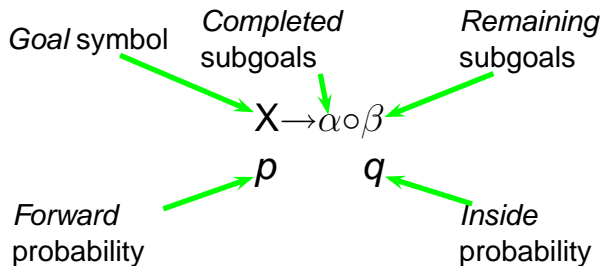
We can use the Earley algorithm (Earley, 1970) in a probabilistic incarnation (Stolcke, 1995) to deal with these infinite tree sets.

The (slightly oversimplified) probabilistic Earley algorithm has two fundamental types of operations:

- ▶ **Prediction:** if X is a possible goal, choose a rule $X \rightarrow Y_\alpha$ and set up Y_α as a new sequence of possible sub-goals of X .
- ▶ **Completion:** if X is a possible goal and we encounter a completed X , absorb it and move on to the next sub-goal in the sequence.

Efficient incremental parsing: the probabilistic Earley algorithm

- ▶ Parsing consists of constructing a *chart* of *states* (items)
- ▶ A state has the following structure:



- ▶ The *forward* probability is the total probability of getting **from** the root at the start of the sentence **through to** this state
- ▶ The *inside* probability is the “bottom-up” probability of the state

Efficient incremental parsing: the probabilistic Earley algorithm

Det \rightarrow o the

1 1

NP \rightarrow o Det N

$\frac{2}{3} + \frac{1}{3}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{2}$ $\frac{1}{3}$

ROOT \rightarrow o NP

1 1

NP \rightarrow NP o PP

$\frac{1}{3}$ $\frac{2}{9}$

NP \rightarrow Det No

$\frac{2}{3}$ $\frac{4}{9}$

N \rightarrow o cat

$\frac{1}{3}$ $\frac{1}{3}$

N \rightarrow o dog

$\frac{2}{3}$ $\frac{2}{3}$

P \rightarrow o near

$\frac{1}{3}$ 1

PP \rightarrow o P NP

$\frac{1}{3}$ $\frac{2}{9}$

Det \rightarrow o the

$\frac{1}{3}$ 1

NP \rightarrow o Det N

$\frac{2}{9} + \frac{1}{9}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{6}$ $\frac{1}{3}$

NP \rightarrow Det o N

1 $\frac{2}{3}$

Det \rightarrow the o

1 1

N \rightarrow dog o

$\frac{2}{3}$ $\frac{2}{3}$

PP \rightarrow P o NP

$\frac{1}{3}$ $\frac{2}{9}$

P \rightarrow near o

$\frac{1}{3}$ 1

NP \rightarrow Det o N

$\frac{1}{3}$ $\frac{2}{3}$

Det \rightarrow the o

$\frac{1}{3}$ 1

the

dog

near

the

Efficient incremental parsing: the probabilistic Earley algorithm

Det \rightarrow o the

1 1

NP \rightarrow o Det N

$\frac{2}{3} + \frac{1}{3}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{2}$ $\frac{1}{3}$

ROOT \rightarrow o NP

1 1

NP \rightarrow NP o PP

$\frac{1}{3}$ $\frac{2}{9}$

NP \rightarrow Det No

$\frac{2}{3}$ $\frac{4}{9}$

N \rightarrow o cat

$\frac{1}{3}$ $\frac{1}{3}$

N \rightarrow o dog

$\frac{2}{3}$ $\frac{2}{3}$

P \rightarrow o near

$\frac{1}{3}$ 1

PP \rightarrow o P NP

$\frac{1}{3}$ $\frac{2}{9}$

Det \rightarrow o the

$\frac{1}{3}$ 1

NP \rightarrow o Det N

$\frac{2}{9} + \frac{1}{9}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{6}$ $\frac{1}{3}$

NP \rightarrow Det o N

1 $\frac{2}{3}$

Det \rightarrow the o

1 1

N \rightarrow dog o

$\frac{2}{3}$ $\frac{2}{3}$

PP \rightarrow P o NP

$\frac{1}{3}$ $\frac{2}{9}$

P \rightarrow near o

$\frac{1}{3}$ 1

NP \rightarrow Det o N

$\frac{1}{3}$ $\frac{2}{3}$

Det \rightarrow the o

$\frac{1}{3}$ 1

the

dog

near

the

Efficient incremental parsing: the probabilistic Earley algorithm

Det \rightarrow o the

1 1

NP \rightarrow o Det N

$\frac{2}{3} + \frac{1}{3}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{2}$ $\frac{1}{3}$

ROOT \rightarrow o NP

1 1

NP \rightarrow NP o PP

$\frac{1}{3}$ $\frac{2}{9}$

NP \rightarrow Det No

$\frac{2}{3}$ $\frac{4}{9}$

N \rightarrow o cat

$\frac{1}{3}$ $\frac{1}{3}$

N \rightarrow o dog

$\frac{2}{3}$ $\frac{2}{3}$

P \rightarrow o near

$\frac{1}{3}$ 1

PP \rightarrow o P NP

$\frac{1}{3}$ $\frac{2}{9}$

Det \rightarrow o the

$\frac{1}{3}$ 1

NP \rightarrow o Det N

$\frac{2}{9} + \frac{1}{9}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{6}$ $\frac{1}{3}$

NP \rightarrow Det o N

1 $\frac{2}{3}$

Det \rightarrow the o

1 1

N \rightarrow dog o

$\frac{2}{3}$ $\frac{2}{3}$

PP \rightarrow P o NP

$\frac{1}{3}$ $\frac{2}{9}$

P \rightarrow near o

$\frac{1}{3}$ 1

NP \rightarrow Det o N

$\frac{1}{3}$ $\frac{2}{3}$

Det \rightarrow the o

$\frac{1}{3}$ 1

the

dog

near

the

Efficient incremental parsing: the probabilistic Earley algorithm

Det \rightarrow o the

1 1

NP \rightarrow o Det N

$\frac{2}{3} + \frac{1}{3}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{2}$ $\frac{1}{3}$

ROOT \rightarrow o NP

1 1

NP \rightarrow NP o PP

$\frac{1}{3}$ $\frac{2}{9}$

NP \rightarrow Det No

$\frac{2}{3}$ $\frac{4}{9}$

N \rightarrow o cat

$\frac{1}{3}$ $\frac{1}{3}$

N \rightarrow o dog

$\frac{2}{3}$ $\frac{2}{3}$

P \rightarrow o near

$\frac{1}{3}$ 1

PP \rightarrow o P NP

$\frac{1}{3}$ $\frac{2}{9}$

Det \rightarrow o the

$\frac{1}{3}$ 1

NP \rightarrow o Det N

$\frac{2}{9} + \frac{1}{9}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{6}$ $\frac{1}{3}$

NP \rightarrow Det o N

1 $\frac{2}{3}$

Det \rightarrow the o

1 1

N \rightarrow dog o

$\frac{2}{3}$ $\frac{2}{3}$

PP \rightarrow P o NP

$\frac{1}{3}$ $\frac{2}{9}$

P \rightarrow near o

$\frac{1}{3}$ 1

NP \rightarrow Det o N

$\frac{1}{3}$ $\frac{2}{3}$

Det \rightarrow the o

$\frac{1}{3}$ 1

the

dog

near

the

Efficient incremental parsing: the probabilistic Earley algorithm

Det → o the

1 1

NP → o Det N

$\frac{2}{3} + \frac{1}{3}$ $\frac{2}{3}$

NP → o NP PP

$\frac{1}{2}$ $\frac{1}{3}$

ROOT → o NP

1 1

NP → NP o PP

$\frac{1}{3}$ $\frac{2}{9}$

NP → Det No

$\frac{2}{3}$ $\frac{4}{9}$

N → o cat

$\frac{1}{3}$ $\frac{1}{3}$

N → o dog

$\frac{2}{3}$ $\frac{2}{3}$

P → o near

$\frac{1}{3}$ 1

PP → o P NP

$\frac{1}{3}$ $\frac{2}{9}$

Det → o the

$\frac{1}{3}$ 1

NP → o Det N

$\frac{2}{9} + \frac{1}{9}$ $\frac{2}{3}$

NP → o NP PP

$\frac{1}{6}$ $\frac{1}{3}$

NP → Det o N

1 $\frac{2}{3}$

Det → the o

1 1

N → dog o

$\frac{2}{3}$ $\frac{2}{3}$

PP → P o NP

$\frac{1}{3}$ $\frac{2}{9}$

P → near o

$\frac{1}{3}$ 1

NP → Det o N

$\frac{1}{3}$ $\frac{2}{3}$

Det → the o

$\frac{1}{3}$ 1

the

dog

near

the

Efficient incremental parsing: the probabilistic Earley algorithm

Det \rightarrow o the

1 1

NP \rightarrow o Det N

$\frac{2}{3} + \frac{1}{3}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{2}$ $\frac{1}{3}$

ROOT \rightarrow o NP

1 1

NP \rightarrow NP o PP

$\frac{1}{3}$ $\frac{2}{9}$

NP \rightarrow Det No

$\frac{2}{3}$ $\frac{4}{9}$

N \rightarrow o cat

$\frac{1}{3}$ $\frac{1}{3}$

N \rightarrow o dog

$\frac{2}{3}$ $\frac{2}{3}$

P \rightarrow o near

$\frac{1}{3}$ 1

PP \rightarrow o P NP

$\frac{1}{3}$ $\frac{2}{9}$

Det \rightarrow o the

$\frac{1}{3}$ 1

NP \rightarrow o Det N

$\frac{2}{9} + \frac{1}{9}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{6}$ $\frac{1}{3}$

NP \rightarrow Det o N

1 $\frac{2}{3}$

Det \rightarrow the o

1 1

N \rightarrow dog o

$\frac{2}{3}$ $\frac{2}{3}$

PP \rightarrow P o NP

$\frac{1}{3}$ $\frac{2}{9}$

P \rightarrow near o

$\frac{1}{3}$ 1

NP \rightarrow Det o N

$\frac{1}{3}$ $\frac{2}{3}$

Det \rightarrow the o

$\frac{1}{3}$ 1

the

dog

near

the

Efficient incremental parsing: the probabilistic Earley algorithm

Det → o the

1 1

NP → o Det N

$\frac{2}{3} + \frac{1}{3}$ $\frac{2}{3}$

NP → o NP PP

$\frac{1}{2}$ $\frac{1}{3}$

ROOT → o NP

1 1

NP → NP o PP

$\frac{1}{3}$ $\frac{2}{9}$

NP → Det No

$\frac{2}{3}$ $\frac{4}{9}$

N → o cat

$\frac{1}{3}$ $\frac{1}{3}$

N → o dog

$\frac{2}{3}$ $\frac{2}{3}$

P → o near

$\frac{1}{3}$ 1

PP → o P NP

$\frac{1}{3}$ $\frac{2}{9}$

Det → o the

$\frac{1}{3}$ 1

NP → o Det N

$\frac{2}{9} + \frac{1}{9}$ $\frac{2}{3}$

NP → o NP PP

$\frac{1}{6}$ $\frac{1}{3}$

NP → Det o N

1 $\frac{2}{3}$

Det → the o

1 1

N → dog o

$\frac{2}{3}$ $\frac{2}{3}$

PP → P o NP

$\frac{1}{3}$ $\frac{2}{9}$

P → near o

$\frac{1}{3}$ 1

NP → Det o N

$\frac{1}{3}$ $\frac{2}{3}$

Det → the o

$\frac{1}{3}$ 1

the

dog

near

the

Efficient incremental parsing: the probabilistic Earley algorithm

Det \rightarrow o the

1 1

NP \rightarrow o Det N

$\frac{2}{3} + \frac{1}{3}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{2}$ $\frac{1}{3}$

ROOT \rightarrow o NP

1 1

NP \rightarrow NP o PP

$\frac{1}{3}$ $\frac{2}{9}$

NP \rightarrow Det No

$\frac{2}{3}$ $\frac{4}{9}$

N \rightarrow o cat

$\frac{1}{3}$ $\frac{1}{3}$

N \rightarrow o dog

$\frac{2}{3}$ $\frac{2}{3}$

P \rightarrow o near

$\frac{1}{3}$ 1

PP \rightarrow o P NP

$\frac{1}{3}$ $\frac{2}{9}$

Det \rightarrow o the

$\frac{1}{3}$ 1

NP \rightarrow o Det N

$\frac{2}{9} + \frac{1}{9}$ $\frac{2}{3}$

NP \rightarrow o NP PP

$\frac{1}{6}$ $\frac{1}{3}$

NP \rightarrow Det o N

1 $\frac{2}{3}$

Det \rightarrow the o

1 1

N \rightarrow dog o

$\frac{2}{3}$ $\frac{2}{3}$

PP \rightarrow P o NP

$\frac{1}{3}$ $\frac{2}{9}$

P \rightarrow near o

$\frac{1}{3}$ 1

NP \rightarrow Det o N

$\frac{1}{3}$ $\frac{2}{3}$

Det \rightarrow the o

$\frac{1}{3}$ 1

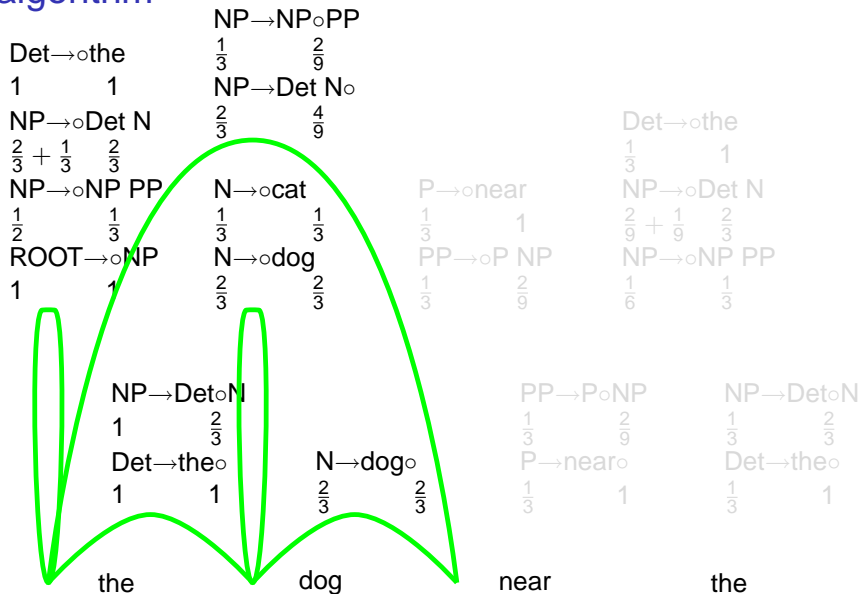
the

dog

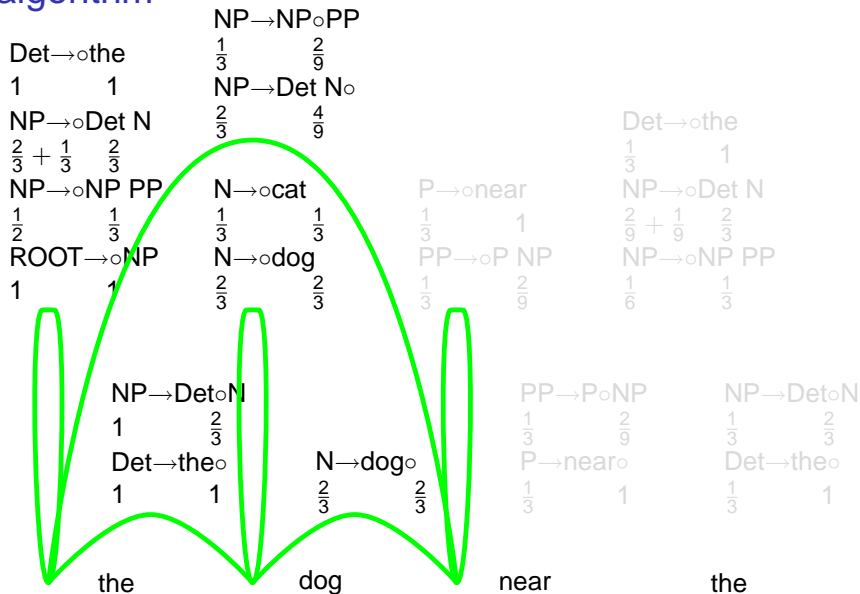
near

the

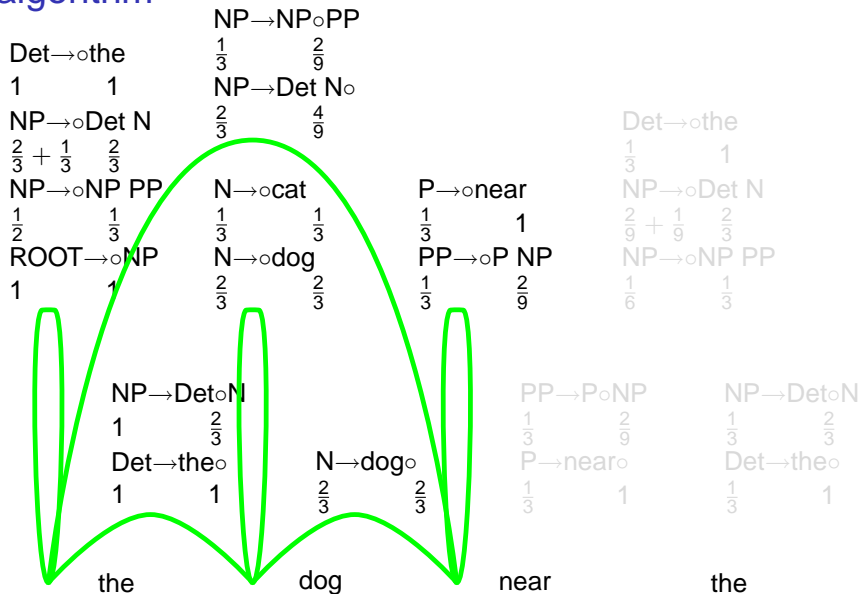
Efficient incremental parsing: the probabilistic Earley algorithm



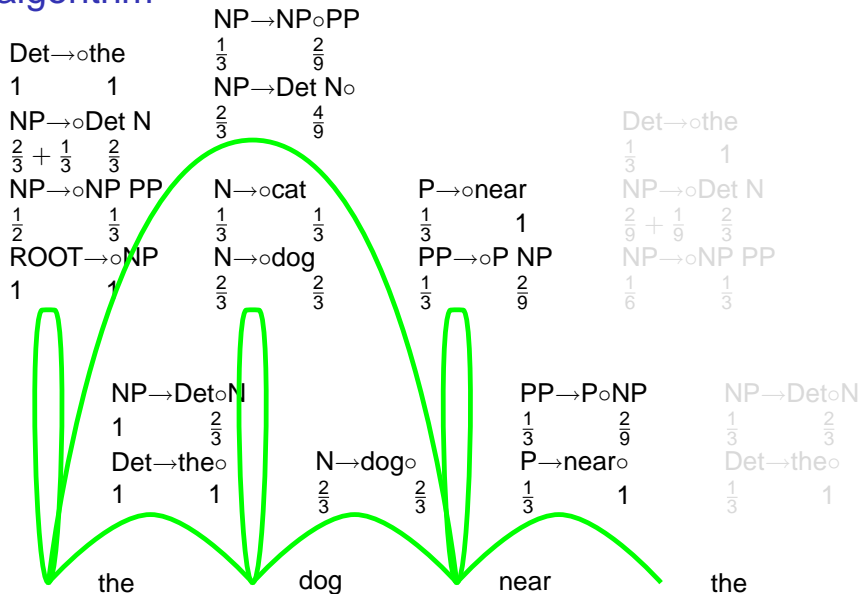
Efficient incremental parsing: the probabilistic Earley algorithm



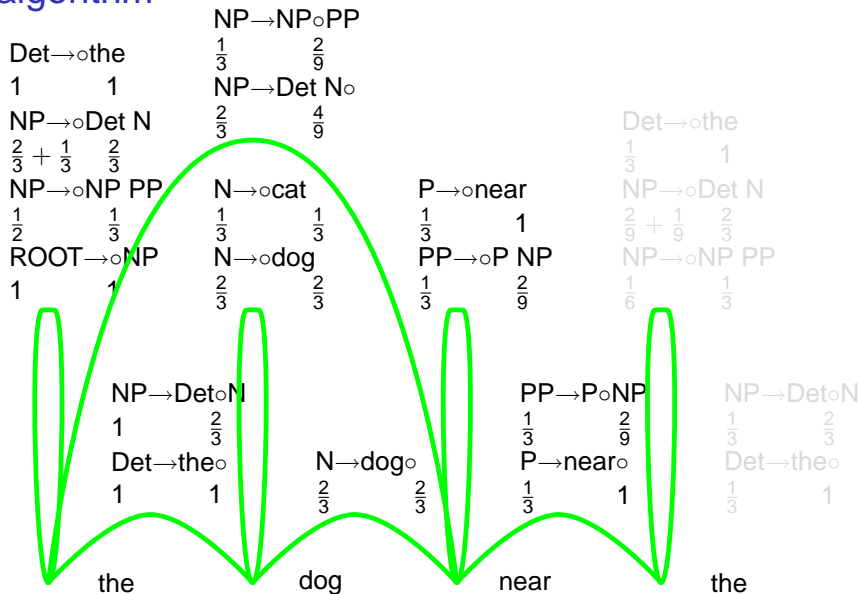
Efficient incremental parsing: the probabilistic Earley algorithm



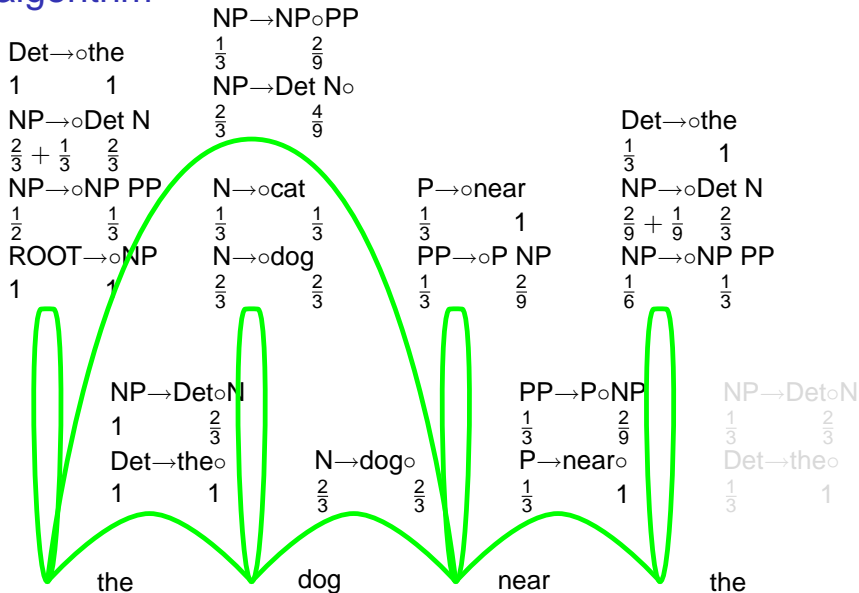
Efficient incremental parsing: the probabilistic Earley algorithm



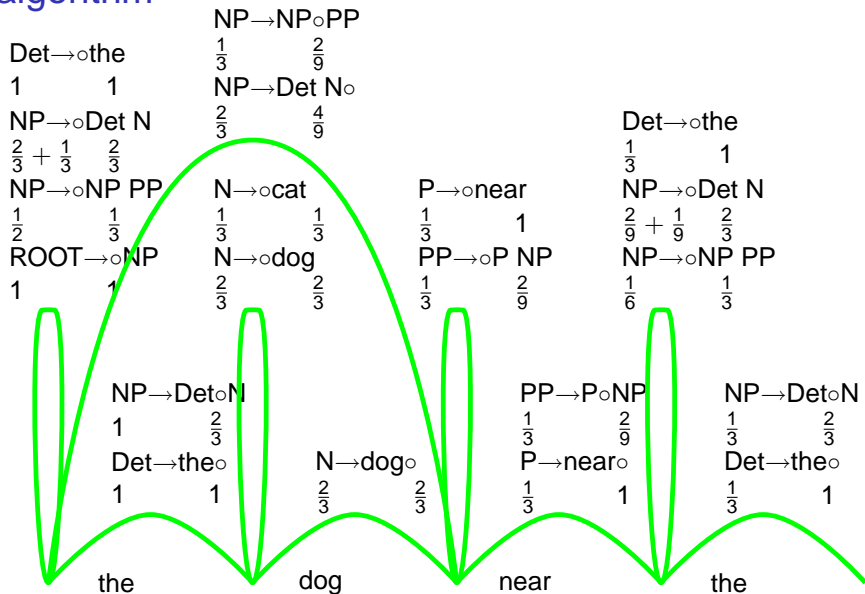
Efficient incremental parsing: the probabilistic Earley algorithm



Efficient incremental parsing: the probabilistic Earley algorithm



Efficient incremental parsing: the probabilistic Earley algorithm



Probabilistic Earley as an “eager” algorithm

- ▶ From the *inside probabilities* of the states on the chart, the posterior distribution on (incremental) trees can be directly calculated
- ▶ This posterior distribution is *precisely* the correct result of the application of Bayes’ rule
- ▶ Hence, probabilistic Earley is also performing rational disambiguation
- ▶ Hale (2001) called this the “eager” property of an incremental parsing algorithm.

Probabilistic Earley algorithm: key ideas

- ▶ We want to use probabilistic grammars for both disambiguation and calculating probability distributions over upcoming events
- ▶ Infinitely many trees can be constructed in polynomial time ($O(n^3)$) and space ($O(n^2)$)
- ▶ The *prefix probability* of the string is calculated in the process
- ▶ By taking the log-ratio of two prefix probabilities, the surprisal of a word in its context can be calculated

Probabilistic Earley algorithm: key ideas

- ▶ We want to use probabilistic grammars for both disambiguation and calculating probability distributions over upcoming events
- ▶ Infinitely many trees can be constructed in polynomial time ($O(n^3)$) and space ($O(n^2)$)
- ▶ The *prefix probability* of the string is calculated in the process
- ▶ By taking the log-ratio of two prefix probabilities, the surprisal of a word in its context can be calculated

Probabilistic Earley algorithm: key ideas

- ▶ We want to use probabilistic grammars for both disambiguation and calculating probability distributions over upcoming events
- ▶ Infinitely many trees can be constructed in polynomial time ($O(n^3)$) and space ($O(n^2)$)
- ▶ The *prefix probability* of the string is calculated in the process
- ▶ By taking the log-ratio of two prefix probabilities, the surprisal of a word in its context can be calculated

Next...

Applications of the idea of surprisal to comprehension and production

References

- Altmann, G. T. and Kamide, Y. (1999). Incremental interpretation at verbs: restricting the domain of subsequent reference. *Cognition*, 73(3):247–264.
- Attneave, F. (1959). *Applications of Information Theory to Psychology: A summary of basic concepts, methods and results*. Holt, Rinehart and Winston.
- Earley, J. (1970). An efficient context-free parsing algorithm. *Communications of the ACM*, 13(2):94–102.
- Ehrlich, S. F. and Rayner, K. (1981). Contextual effects on word perception and eye movements during reading. *Journal of Verbal Learning and Verbal Behavior*, 20:641–655.
- Grodner, D. and Gibson, E. (2005). Some consequences of the serial nature of linguistic input. *Cognitive Science*, 29(2):261–290.
- Hale, J. (2001). A probabilistic Earley parser as a psycholinguistic model. In *Proceedings of NAACL*, volume 2, pages 159–166.
- Itti, L. and Baldi, P. (2005). Bayesian surprise attracts human attention. In *Advances in Neural Information Processing Systems*.
- Levy, R. (2005). *Probabilistic Models of Word Order and Syntactic Discontinuity*. PhD thesis, Stanford University.
- Levy, R. (2007). Expectation-based syntactic comprehension. *Cognition*. In press.
- Smith, N. (2006). Surprisal-based sentence processing as optimal behavior. M.S., UC San Diego.
- Stolcke, A. (1995). An efficient probabilistic context-free parsing algorithm that computes prefix probabilities. *Computational Linguistics*, 21(2):165–201.