

The empirical significance of derivational operations

Tim Hunter

University of California, Los Angeles

February 15, 2017

Overview

Take-home message

Distinct hypotheses about derivational operations can have distinct consequences for our theories' predictions about speakers' linguistic behavior.

Overview

Take-home message

Distinct hypotheses about derivational operations can have distinct consequences for our theories' predictions about speakers' linguistic behavior.

Illustrative case study:

- distinct implementations of movement: move vs. remerge
- distinct predictions when plugged into models of
 - sentence comprehension difficulty, via surprisal
 - grammar selection by a learner, via simple maximum likelihood learning

Overview

Take-home message

Distinct hypotheses about derivational operations can have distinct consequences for our theories' predictions about speakers' linguistic behavior.

Illustrative case study:

- distinct implementations of movement: move vs. remerge
- distinct predictions when plugged into models of
 - sentence comprehension difficulty, via surprisal
 - grammar selection by a learner, via simple maximum likelihood learning

Key idea: Derivations encode the relationship between

- primitive, memorized chunks of knowledge; and
- consequences computed from those

Outline

- 1 What are grammars?
- 2 Derivationally distinct implementations of merge
- 3 Telling them apart
- 4 Historical perspective

Outline

- 1 What are grammars?
- 2 Derivationally distinct implementations of merge
- 3 Telling them apart
- 4 Historical perspective

Grammars

Q: What is a grammar?

A: A (finite) collection of memorized statements that allows a speaker to recognize an (infinite) collection of expressions

Grammars

Q: What is a grammar?

A: A (finite) collection of memorized statements that allows a speaker to recognize an (infinite) collection of expressions

Why do we want a grammar rather than just a list of expressions?

- The infinitely many expressions can't be encoded directly in a finite mind
- We recognize as well-formed (and understand, and ...) expressions we've never encountered before

Grammars

Q: What is a grammar?

A: A (finite) collection of memorized statements that allows a speaker to recognize an (infinite) collection of expressions

Why do we want a grammar rather than just a list of expressions?

- The infinitely many expressions can't be encoded directly in a finite mind
- We recognize as well-formed (and understand, and ...) expressions we've never encountered before

Q: What is a derivation?

A: A particular **chaining-together** of interacting, individually-memorized chunks.

Derivations

Q: What is a derivation?

A: A particular [chaining-together](#) of interacting, individually-memorized chunks.

Derivations

Q: What is a derivation?

A: A particular **chaining-together** of interacting, individually-memorized chunks.

How does a speaker recognize 'painters' as well-formed?

Derivations

Q: What is a derivation?

A: A particular **chaining-together** of interacting, individually-memorized chunks.

How does a speaker recognize 'painters' as well-formed?

A speaker recognizes 'painters' as a well-formed expression by

- knowing that $N + \text{'-s'}$ is well-formed if N is well-formed
- recognizing that 'painter' is well-formed

Derivations

Q: What is a derivation?

A: A particular **chaining-together** of interacting, individually-memorized chunks.

How does a speaker recognize 'painters' as well-formed?

A speaker recognizes 'painters' as a well-formed expression by

- knowing that $N + \text{'-s'}$ is well-formed if N is well-formed
- recognizing that 'painter' is well-formed

A speaker recognizes 'painter' as a well-formed expression by

- knowing that $V + \text{'-er'}$ is well-formed if V is well-formed
- knowing that 'paint' is well-formed

Derivations

Q: What is a derivation?

A: A particular **chaining-together** of interacting, individually-memorized chunks.

How does a speaker recognize 'painters' as well-formed?

A speaker recognizes 'painters' as a well-formed expression by

- knowing that $N + \text{'-s'}$ is well-formed if N is well-formed (a rule)
- recognizing that 'painter' is well-formed

A speaker recognizes 'painter' as a well-formed expression by

- knowing that $V + \text{'-er'}$ is well-formed if V is well-formed (a rule)
- knowing that 'paint' is well-formed (a lexical item)

Phillips and Lewis (2013)

Three views of derivations: literalist, formalist, extensionalist

These correspond to different views about which linking hypotheses expose the the derivational claims of a theory to empirical testing.

Phillips and Lewis (2013)

Three views of derivations: literalist, formalist, extensionalist

These correspond to different views about which linking hypotheses expose the the derivational claims of a theory to empirical testing.

- Literalist (one extreme):
Derivational operations are real-time operations, so observing real-time operations bears directly on derivations.

Phillips and Lewis (2013)

Three views of derivations: literalist, formalist, extensionalist

These correspond to different views about which linking hypotheses expose the the derivational claims of a theory to empirical testing.

- Literalist (one extreme):
Derivational operations are real-time operations, so observing real-time operations bears directly on derivations.
- Extensionalist (other extreme):
Derivational operations are abstract to a degree that makes no linking hypotheses available.

Phillips and Lewis (2013)

Three views of derivations: literalist, formalist, extensionalist

These correspond to different views about which linking hypotheses expose the the derivational claims of a theory to empirical testing.

- Literalist (one extreme):
Derivational operations are real-time operations, so observing real-time operations bears directly on derivations.
- Extensionalist (other extreme):
Derivational operations are abstract to a degree that makes no linking hypotheses available.
- Formalist (middle ground):
Derivational operations are cognitive hypotheses testable by certain (less direct) linking hypotheses.

Phillips and Lewis (2013)

*When we are told, for example, that the wh-word 'what' is initially merged with a verb and subsequently moved to a left peripheral position in the clause, **what claim is this making about the human language system?***

Phillips and Lewis (2013)

Phillips and Lewis (2013)

*When we are told, for example, that the wh-word 'what' is initially merged with a verb and subsequently moved to a left peripheral position in the clause, **what claim is this making about the human language system?***

Phillips and Lewis (2013)

My answer:

- One component of the human language system is the knowledge of a systematic relationship ("merge") that holds between
 - the well-formedness of the expression 'ate what', and
 - the well-formedness of the expressions 'ate' and 'what'.

Phillips and Lewis (2013)

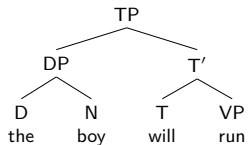
*When we are told, for example, that the wh-word 'what' is initially merged with a verb and subsequently moved to a left peripheral position in the clause, **what claim is this making about the human language system?***

Phillips and Lewis (2013)

My answer:

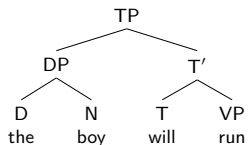
- One component of the human language system is the knowledge of a systematic relationship ("merge") that holds between
 - the well-formedness of the expression 'ate what', and
 - the well-formedness of the expressions 'ate' and 'what'.
- One component of the human language system is the knowledge of a systematic relationship ("move") that holds between
 - the well-formedness of the expression 'what John ate t_i ', and
 - the well-formedness of the expression 'John ate what'.

Contemporary syntactic derivations



How does a speaker recognize this as a well-formed expression?

Contemporary syntactic derivations

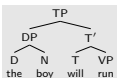


How does a speaker recognize this as a well-formed expression?

Wrong answer: By knowing that this is a well-formed expression.

So some [chaining-together](#) of other chunks of knowledge is required, i.e. a derivation.

A speaker recognizes that



is a well-formed expression by

- knowing that



is well-formed if $\underline{\text{DP}}$ and $\underline{\text{T}'}$ are well-formed

- recognizing that



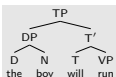
is well-formed

- recognizing that



is well-formed



A speaker recognizes that



is a well-formed expression by

- knowing that



is well-formed if  and  are well-formed

- recognizing that



is well-formed

- recognizing that



is well-formed



A speaker recognizes that



is a well-formed expression by

- knowing that



is well-formed if  and  are well-formed

- knowing that



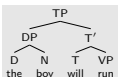
is well-formed

- knowing that



is well-formed

A speaker recognizes that



is a well-formed expression by

- knowing that



is well-formed if $\underline{\text{DP}}$ and $\underline{\text{T}'}$ are well-formed

- recognizing that



is well-formed

- recognizing that



is well-formed

A speaker recognizes that



is a well-formed expression by

- knowing that



is well-formed if $\underline{\text{D}}$ and $\underline{\text{N}}$ are well-formed

- knowing that



is well-formed

- knowing that



is well-formed

A speaker recognizes that



is a well-formed expression by

- knowing that



is well-formed if $\underline{\text{T}}$ and $\underline{\text{VP}}$ are well-formed

- knowing that



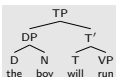
is well-formed

- knowing that





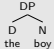
is well-formed

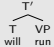
A speaker recognizes that



is a well-formed expression by

- knowing that  is well-formed if  and  are well-formed (MERGE)

- recognizing that  is well-formed


- recognizing that  is well-formed


A speaker recognizes that



is a well-formed expression by

- knowing that  is well-formed if  and  are well-formed (MERGE)

- knowing that  is well-formed

- knowing that  is well-formed


A speaker recognizes that



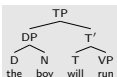
is a well-formed expression by

- knowing that  is well-formed if  and  are well-formed (MERGE)

- knowing that  is well-formed

- knowing that  is well-formed

A speaker recognizes that

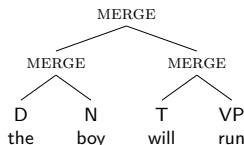


is a well-formed expression by

- knowing that  is well-formed if  and  are well-formed (MERGE)

- recognizing that  is well-formed

- recognizing that  is well-formed




A speaker recognizes that



is a well-formed expression by

- knowing that  is well-formed if  and  are well-formed (MERGE)

- knowing that  is well-formed

- knowing that  is well-formed

A speaker recognizes that



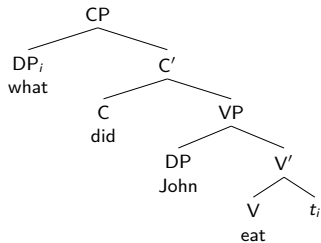
is a well-formed expression by

- knowing that  is well-formed if  and  are well-formed (MERGE)

- knowing that  is well-formed

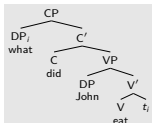
- knowing that  is well-formed

How does a speaker recognize this as a well-formed expression?



(Same wrong answer as before ...)

A speaker recognizes that



is a well-formed expression by

- knowing that

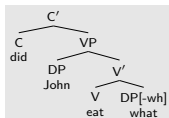


is well-formed if



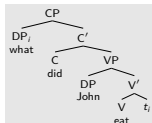
is well-formed

- recognizing that



is well-formed

A speaker recognizes that



is a well-formed expression by

- knowing that

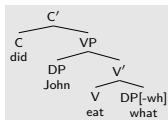


is well-formed if



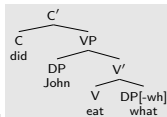
is well-formed

- recognizing that



is well-formed

A speaker recognizes that



is a well-formed expression by

- knowing that



is well-formed if



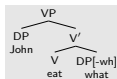
are well-formed

- knowing that



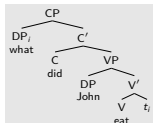
is well-formed

- recognizing that



is well-formed

A speaker recognizes that



is a well-formed expression by

- knowing that

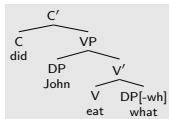


is well-formed if



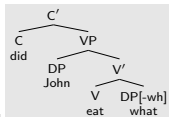
is well-formed (MOVE)

- recognizing that



is well-formed

A speaker recognizes that



is a well-formed expression by

- knowing that



is well-formed if



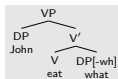
are well-formed (MERGE)

- knowing that



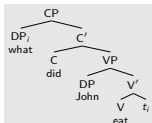
is well-formed

- recognizing that



is well-formed

A speaker recognizes that



is a well-formed expression by

- knowing that

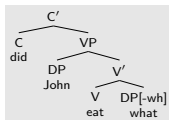


is well-formed if



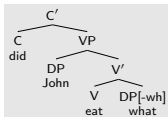
is well-formed (MOVE)

- recognizing that



is well-formed

A speaker recognizes that



is a well-formed expression by

- knowing that



is well-formed if



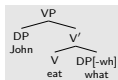
are well-formed (MERGE)

- knowing that

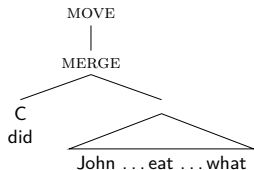


is well-formed

- recognizing that

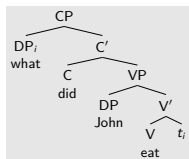


is well-formed



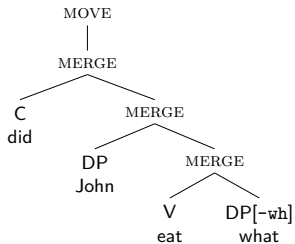
Chaining together memorized chunks

A speaker recognizes that



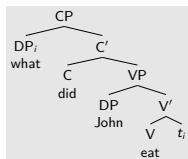
is a well-formed expression by ...

... identifying these relationships between memorized chunks of knowledge:



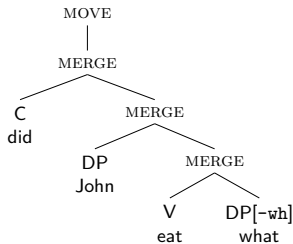
Chaining together memorized chunks

A speaker recognizes that



is a well-formed expression by ...

... identifying these relationships between memorized chunks of knowledge:



But how can we investigate the chained-together chunks rather than just their results?

Distinct divisions of labour

Suppose we have a black box that recognizes a triple of numbers iff each number is drawn from $\{1, 2, 3, 4, 5, 6\}$.

We have two hypotheses about the “grammar” inside the black box.

Distinct divisions of labour

Suppose we have a black box that recognizes a triple of numbers iff each number is drawn from $\{1, 2, 3, 4, 5, 6\}$.

We have two hypotheses about the “grammar” inside the black box.

Hypothesis #1 involves a blue die and a red die :

- a triple (i, j, k) is well-formed if i is well-formed, j is well-formed, and k is well-formed

Distinct divisions of labour

Suppose we have a black box that recognizes a triple of numbers iff each number is drawn from $\{1, 2, 3, 4, 5, 6\}$.

We have two hypotheses about the “grammar” inside the black box.

Hypothesis #1 involves a blue die and a red die :

- a triple (i, j, k) is well-formed if i is well-formed, j is well-formed, and k is well-formed

Hypothesis #2 involves a green die and a yellow die :

- a triple (i, j, k) is well-formed if i is well-formed, j is well-formed, and k is well-formed

Distinct divisions of labour

Suppose we have a black box that recognizes a triple of numbers iff each number is drawn from $\{1, 2, 3, 4, 5, 6\}$.

We have two hypotheses about the “grammar” inside the black box.

Hypothesis #1 involves a blue die and a red die :

- a triple (i, j, k) is well-formed if i is well-formed, j is well-formed, and k is well-formed

Hypothesis #2 involves a green die and a yellow die :

- a triple (i, j, k) is well-formed if i is well-formed, j is well-formed, and k is well-formed

These are different “divisions of labour”, ways of breaking down the work into (finitely many) chainable chunks

Division of labour

What do these hypotheses say about the triple $(4, 5, 6)$?

Hypothesis #1:

$(4, 5, 6)$ is well-formed if

- 4 is well-formed
- 5 is well-formed
- 6 is well-formed

Hypothesis #2:

$(4, 5, 6)$ is well-formed if

- 4 is well-formed
- 5 is well-formed
- 6 is well-formed

Division of labour

What do these hypotheses say about the triple (4, 5, 6)?

Hypothesis #1:

(4, 5, 6) is possible to the extent that

- 4 is possible
- 5 is possible
- 6 is possible

Hypothesis #2:

(4, 5, 6) is possible to the extent that

- 4 is possible
- 5 is possible
- 6 is possible

Division of labour

What do these hypotheses say about the triple $(4, 5, 6)$?

Hypothesis #1:

$(4, 5, 6)$ is probable to the extent that

- 4 is probable
- 5 is probable
- 6 is probable

Hypothesis #2:

$(4, 5, 6)$ is probable to the extent that

- 4 is probable
- 5 is probable
- 6 is probable

Division of labour

What do these hypotheses say about the triple (4, 5, 6)?

Hypothesis #1:

(4, 5, 6) is probable to the extent that

- 4 is probable
- 5 is probable
- 6 is probable

So it should pattern with (5, 4, 6).

Both have probability: $P(4) \cdot P(5) \cdot P(6)$

Hypothesis #2:

(4, 5, 6) is probable to the extent that

- 4 is probable
- 5 is probable
- 6 is probable

So it should pattern with (4, 6, 5).

Both have probability: $P(4) \cdot P(5) \cdot P(6)$

Division of labour

What do these hypotheses say about the triple (4, 5, 6)?

Hypothesis #1:

(4, 5, 6) is probable to the extent that

- 4 is probable
- 5 is probable
- 6 is probable

So it should pattern with (5, 4, 6).

Both have probability: $P(4) \cdot P(5) \cdot P(6)$

- Same **possibilities**: both recognize the set $\{1, 2, 3, 4, 5, 6\}^3$
- Different **probabilities**: distinct ranges of probability distributions over this set
- More generally: distinct similarity relations over this set

Hypothesis #2:

(4, 5, 6) is probable to the extent that

- 4 is probable
- 5 is probable
- 6 is probable

So it should pattern with (4, 6, 5).

Both have probability: $P(4) \cdot P(5) \cdot P(6)$

Outline

- 1 What are grammars?
- 2 Derivationally distinct implementations of merge
- 3 Telling them apart
- 4 Historical perspective

Roadmap

Plan for this section and the next:

- Introduce two grammatical systems that differ only in their derivational operations (i.e. their “chainable” chunks)
 - merge and move as distinct derivational primitives (Stabler 1997, Keenan and Stabler 2003)
 - merge and move implemented by a single derivational primitive (Stabler 2006, Hunter 2011)

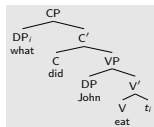
Roadmap

Plan for this section and the next:

- Introduce two grammatical systems that differ only in their derivational operations (i.e. their “chainable” chunks)
 - merge and move as distinct derivational primitives (Stabler 1997, Keenan and Stabler 2003)
 - merge and move implemented by a single derivational primitive (Stabler 2006, Hunter 2011)
- Show that they produce different empirical predictions when plugged in to common probabilistic modeling settings
 - sentence comprehension difficulty via surprisal
 - selection among candidate grammars by a learner

IMG derivations

A speaker recognizes that

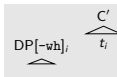


is a well-formed expression by

- knowing that

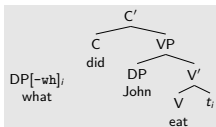


is well-formed if



is well-formed (MRG)

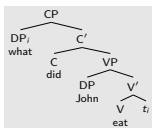
- recognizing that



is well-formed

IMG derivations

A speaker recognizes that



is a well-formed expression by

- knowing that is well-formed if is well-formed (MRG)

- recognizing that is well-formed

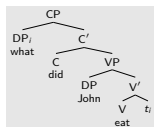
A speaker recognizes that this is a well-formed expression by

- knowing that is well-formed if is well-formed (MRG)

- recognizing that is well-formed

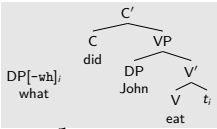
IMG derivations

A speaker recognizes that



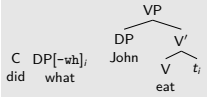
is a well-formed expression by

- knowing that  is well-formed if  is well-formed (MRG)

- recognizing that  is well-formed

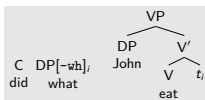
A speaker recognizes that this is a well-formed expression by

- knowing that  is well-formed if  is well-formed (MRG)

- recognizing that  is well-formed

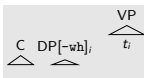
IMG derivations


A speaker recognizes that

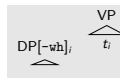


is well-formed by

- knowing that



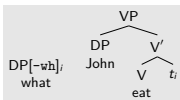
is well-formed if  and



are well-formed (INSERT)

- knowing that  is well-formed

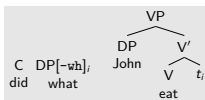
- recognizing that



is well-formed

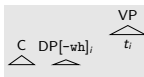
IMG derivations

A speaker recognizes that

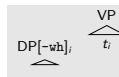


is well-formed by

- knowing that



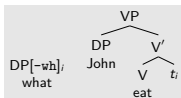
is well-formed if  and



are well-formed (INSERT)

- knowing that  is well-formed

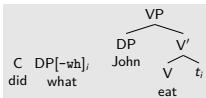
- recognizing that



is well-formed

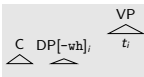
IMG derivations

A speaker recognizes that



is well-formed by

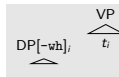
- knowing that



is well-formed if



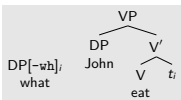
and



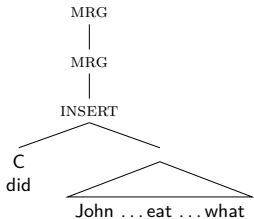
are well-formed (INSERT)

- knowing that  is well-formed

- recognizing that

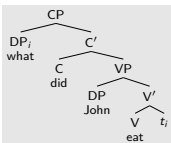


is well-formed



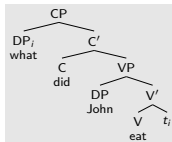
Distinct derivations

Which primitive chunks of knowledge are chained together to produce this expression?

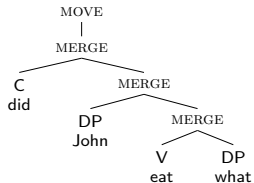


Distinct derivations

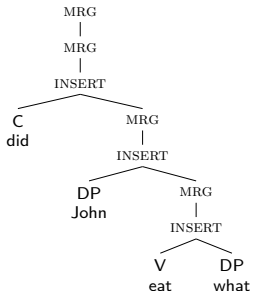
Which primitive chunks of knowledge are chained together to produce this expression?



MG hypothesis: **MERGE** and **MOVE**



IMG hypothesis: **MRG** and **INSERT**



(NB: Don't be distracted by the difference in the **number** of steps.)

Phillips and Lewis (2013)

*When we are told, for example, that the wh-word 'what' is initially merged with a verb and subsequently moved to a left peripheral position in the clause, **what claim is this making about the human language system?***

Phillips and Lewis (2013)

Phillips and Lewis (2013)

When we are told, for example, that the wh-word 'what' is initially merged with a verb and subsequently moved to a left peripheral position in the clause, what claim is this making about the human language system?

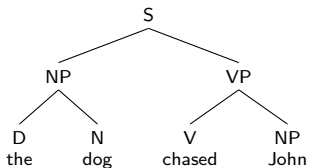
Phillips and Lewis (2013)

When we are told, for example, that the wh-word 'what' is initially inserted and subsequently merged with a verb and then merged into a left peripheral position in the clause, what claim is this making about the human language system?

Outline

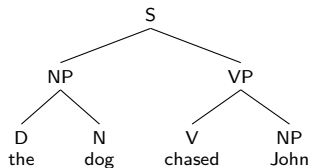
- 1 What are grammars?
- 2 Derivationally distinct implementations of merge
- 3 Telling them apart**
- 4 Historical perspective

Probabilities on a CFG



1.0	S	→	NP VP
0.3	NP	→	John
0.2	NP	→	he
0.5	NP	→	D N
0.3	D	→	the
0.7	D	→	a
0.6	N	→	dog
0.4	N	→	cat
0.8	VP	→	V NP
0.2	VP	→	V
0.9	V	→	chased
0.1	V	→	ate

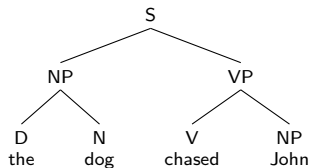
Probabilities on a CFG



$$P(T) = 1.0 \times (0.5 \times 0.3 \times 0.6) \times (0.8 \times 0.9 \times 0.3)$$

1.0	S	→	NP VP
0.3	NP	→	John
0.2	NP	→	he
0.5	NP	→	D N
0.3	D	→	the
0.7	D	→	a
0.6	N	→	dog
0.4	N	→	cat
0.8	VP	→	V NP
0.2	VP	→	V
0.9	V	→	chased
0.1	V	→	ate

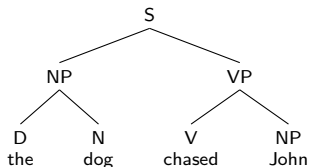
Probabilities on a CFG



$$P(T|\lambda) = \lambda_1 \times (\lambda_4 \times \lambda_5 \times \lambda_7) \times (\lambda_9 \times \lambda_{11} \times \lambda_2)$$

λ_1	S	→	NP VP
λ_2	NP	→	John
λ_3	NP	→	he
λ_4	NP	→	D N
λ_5	D	→	the
λ_6	D	→	a
λ_7	N	→	dog
λ_8	N	→	cat
λ_9	VP	→	V NP
λ_{10}	VP	→	V
λ_{11}	V	→	chased
λ_{12}	V	→	ate

Probabilities on a CFG



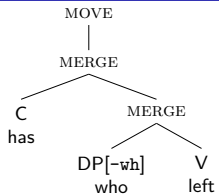
$$P(T|\lambda) = \lambda_1 \times (\lambda_4 \times \lambda_5 \times \lambda_7) \times (\lambda_9 \times \lambda_{11} \times \lambda_2)$$

λ_1	S	→	NP VP
λ_2	NP	→	John
λ_3	NP	→	he
λ_4	NP	→	D N
λ_5	D	→	the
λ_6	D	→	a
λ_7	N	→	dog
λ_8	N	→	cat
λ_9	VP	→	V NP
λ_{10}	VP	→	V
λ_{11}	V	→	chased
λ_{12}	V	→	ate

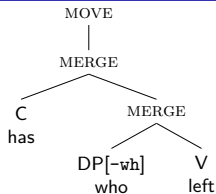
Training question: What values of $\lambda_1, \lambda_2, \dots, \lambda_{12}$ maximize the likelihood of the training data $P(D|\lambda)$?

Note that the choice of [grammatical rules](#) (division of labour) told us what the [parameters](#) were, i.e. defined a space of probability distributions to explore.

Probabilities on minimalist grammars



Probabilities on minimalist grammars

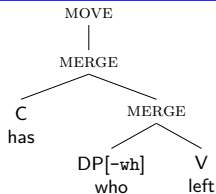


(Things are more complicated because the applicability of a particular rule can't be determined by looking at the directly neighbouring rules in the derivation; cf. n-grams vs. HMMs)

(Hunter and Dyer 2013)

$$\begin{aligned}
 P = & \frac{\exp(\lambda_{\text{MOVE}} + \lambda_{\text{wh}})}{\exp(\lambda_{\text{MOVE}} + \lambda_{\text{wh}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})} \times \frac{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})}{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})} \times \frac{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}})}{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{c}})} \\
 & \times \frac{\exp(\lambda_{\text{left}})}{\exp(\lambda_{\text{left}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}})} \times \frac{\exp(\lambda_{\text{has}})}{\exp(\lambda_{\text{has}}) + \exp(\lambda_{\text{will}})} \times \frac{\exp(\lambda_{\text{who}})}{\exp(\lambda_{\text{who}}) + \exp(\lambda_{\text{what}})}
 \end{aligned}$$

Probabilities on minimalist grammars



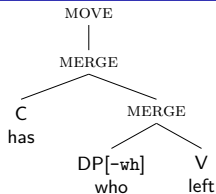
(Things are more complicated because the applicability of a particular rule can't be determined by looking at the directly neighbouring rules in the derivation; cf. n-grams vs. HMMs)

(Hunter and Dyer 2013)

$$\begin{aligned}
 P = & \frac{\exp(\lambda_{\text{MOVE}} + \lambda_{\text{wh}})}{\exp(\lambda_{\text{MOVE}} + \lambda_{\text{wh}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})} \times \frac{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})}{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})} \times \frac{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}})}{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{c}})} \\
 & \times \frac{\exp(\lambda_{\text{left}})}{\exp(\lambda_{\text{left}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}})} \times \frac{\exp(\lambda_{\text{has}})}{\exp(\lambda_{\text{has}}) + \exp(\lambda_{\text{will}})} \times \frac{\exp(\lambda_{\text{who}})}{\exp(\lambda_{\text{who}}) + \exp(\lambda_{\text{what}})}
 \end{aligned}$$

Training question: What values of λ_{MERGE} , λ_{MOVE} , λ_{wh} , λ_{d} , ... maximize the likelihood of the training data $P(D|\lambda)$? (think **blue** and **red**)

Probabilities on minimalist grammars



(Things are more complicated because the applicability of a particular rule can't be determined by looking at the directly neighbouring rules in the derivation; cf. n-grams vs. HMMs)

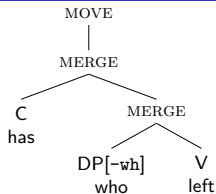
(Hunter and Dyer 2013)

$$\begin{aligned}
 P = & \frac{\exp(\lambda_{\text{MOVE}} + \lambda_{\text{wh}})}{\exp(\lambda_{\text{MOVE}} + \lambda_{\text{wh}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})} \times \frac{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})}{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})} \times \frac{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}})}{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{c}})} \\
 & \times \frac{\exp(\lambda_{\text{left}})}{\exp(\lambda_{\text{left}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}})} \times \frac{\exp(\lambda_{\text{has}})}{\exp(\lambda_{\text{has}}) + \exp(\lambda_{\text{will}})} \times \frac{\exp(\lambda_{\text{who}})}{\exp(\lambda_{\text{who}}) + \exp(\lambda_{\text{what}})}
 \end{aligned}$$

Training question: What values of λ_{MERGE} , λ_{MOVE} , λ_{wh} , λ_{d} , ... maximize the likelihood of the training data $P(D|\lambda)$? (think **blue** and **red**)

And with IMGs, things will be different: λ_{MRG} , λ_{INSERT} , λ_{wh} , λ_{d} , ... (think **green** and **yellow**)

Probabilities on minimalist grammars



(Things are more complicated because the applicability of a particular rule can't be determined by looking at the directly neighbouring rules in the derivation; cf. n-grams vs. HMMs)

(Hunter and Dyer 2013)

$$\begin{aligned}
 P = & \frac{\exp(\lambda_{\text{MOVE}} + \lambda_{\text{wh}})}{\exp(\lambda_{\text{MOVE}} + \lambda_{\text{wh}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})} \times \frac{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})}{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{v}})} \times \frac{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}})}{\exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{c}})} \\
 & \times \frac{\exp(\lambda_{\text{left}})}{\exp(\lambda_{\text{left}}) + \exp(\lambda_{\text{MERGE}} + \lambda_{\text{d}})} \times \frac{\exp(\lambda_{\text{has}})}{\exp(\lambda_{\text{has}}) + \exp(\lambda_{\text{will}})} \times \frac{\exp(\lambda_{\text{who}})}{\exp(\lambda_{\text{who}}) + \exp(\lambda_{\text{what}})}
 \end{aligned}$$

Training question: What values of λ_{MERGE} , λ_{MOVE} , λ_{wh} , λ_{d} , ... maximize the likelihood of the training data $P(D|\lambda)$? (think **blue** and **red**)

And with IMGs, things will be different: λ_{MRG} , λ_{INSERT} , λ_{wh} , λ_{d} , ... (think **green** and **yellow**)

So each system has a different **range of probability distributions** to explore.

A toy minimalist lexicon

$\epsilon :: =t \ c$
 $\epsilon :: =t \ +wh \ c$
 $will :: =v \ =subj \ t$
 $shave :: v$
 $shave :: =obj \ v$
 $boys :: subj$
 $who :: subj \ -wh$

$boys :: =x \ =det \ subj$
 $\epsilon :: x$
 $some :: det$

 $themselves :: =ant \ obj$
 $\epsilon :: =subj \ ant \ -subj$
 $will :: =v \ +subj \ t$

boys will shave
 boys will shave themselves
 who will shave
 who will shave themselves
 some boys will shave
 some boys will shave themselves

Some details:

- Subject is base-generated in SpecTP; no movement for Case
- Transitive and intransitive versions of 'shave'
- 'some' is a determiner that optionally combines with 'boys' to make a subject
 - Dummy feature x to fill complement of 'boys' so that 'some' goes on the left
- 'themselves' can appear in object position via a movement theory of reflexives
 - A subj can be turned into an ant -subj
 - 'themselves' combines with an ant to make an obj
 - 'will' can attract its subject by move as well as merge

Distinct ranges of probability distributions

Take a single training corpus . . .

Possible sentences	Training corpus frequency
boys will shave	10
boys will shave themselves	2
who will shave	3
who will shave themselves	1
some boys will shave	5
some boys will shave themselves	0

Distinct ranges of probability distributions

Take a single training corpus ...

Possible sentences	Training corpus frequency
boys will shave	10
boys will shave themselves	2
who will shave	3
who will shave themselves	1
some boys will shave	5
some boys will shave themselves	0

Separately ask:

- what values for λ_{MERGE} , λ_{MOVE} , λ_{d} , λ_{wh} , ... best fit this training data?
- what values for λ_{MRG} , λ_{INSERT} , λ_{d} , λ_{wh} , ... best fit this training data?

And what do the two results say about the common set of sentences?

Distinct ranges of probability distributions

Take a single training corpus ...

Possible sentences	Training corpus frequency
boys will shave	10
boys will shave themselves	2
who will shave	3
who will shave themselves	1
some boys will shave	5
some boys will shave themselves	0

Separately ask:

- what values for $\lambda_{\text{MERGE}}, \lambda_{\text{MOVE}}, \lambda_{\text{d}}, \lambda_{\text{wh}}, \dots$ best fit this training data?
- what values for $\lambda_{\text{MRG}}, \lambda_{\text{INSERT}}, \lambda_{\text{d}}, \lambda_{\text{wh}}, \dots$ best fit this training data?

And what do the two results say about the common set of sentences?

MG, i.e. MERGE and MOVE	
0.35478	boys will shave
0.35478	some boys will shave
0.14801	who will shave
0.05022	boys will shave themselves
0.05022	some boys will shave themselves
0.04199	who will shave themselves

IMG, i.e. MRG and INSERT	
0.35721	boys will shave
0.35721	some boys will shave
0.095	who will shave
0.095	who will shave themselves
0.04779	boys will shave themselves
0.04779	some boys will shave themselves

Distinct ranges of probability distributions

Take a single training corpus ...

Possible sentences	Training corpus frequency
boys will shave	10
boys will shave themselves	2
who will shave	3
who will shave themselves	1
some boys will shave	5
some boys will shave themselves	0

Separately ask:

- what values for λ_{MERGE} , λ_{MOVE} , λ_{d} , λ_{wh} , ... best fit this training data?
- what values for λ_{MRG} , λ_{INSERT} , λ_{d} , λ_{wh} , ... best fit this training data?

And what do the two results say about the common set of sentences?

MG, i.e. MERGE and MOVE	
0.35478	boys will shave
0.35478	some boys will shave
0.14801	who will shave
0.05022	boys will shave themselves
0.05022	some boys will shave themselves
0.04199	who will shave themselves

IMG, i.e. MRG and INSERT	
0.35721	boys will shave
0.35721	some boys will shave
0.095	who will shave
0.095	who will shave themselves
0.04779	boys will shave themselves
0.04779	some boys will shave themselves

Choice points in the MG-derived MCFG

Question or not?

$\langle c \rangle_0 \rightarrow \langle =t c \rangle_0 \quad \langle t \rangle_0 \quad \exp(\lambda_{\text{MERGE}} + \lambda_t)$

$\langle c \rangle_0 \rightarrow \langle +wh c, -wh \rangle_0 \quad \exp(\lambda_{\text{MOVE}} + \lambda_{wh})$

Non-wh antecedent lexical or complex?

$\langle \text{ant -subj} \rangle_0 \rightarrow \langle =\text{subj ant -subj} \rangle_1 \quad \langle \text{subj} \rangle_0 \quad \exp(\lambda_{\text{MERGE}} + \lambda_{\text{subj}})$

$\langle \text{ant -subj} \rangle_0 \rightarrow \langle =\text{subj ant -subj} \rangle_1 \quad \langle \text{subj} \rangle_1 \quad \exp(\lambda_{\text{MERGE}} + \lambda_{\text{subj}})$

Non-wh subject merged and complex, merged and lexical, or moved?

$\langle t \rangle_0 \rightarrow \langle =\text{subj } t \rangle_0 \quad \langle \text{subj} \rangle_0 \quad \exp(\lambda_{\text{MERGE}} + \lambda_{\text{subj}})$

$\langle t \rangle_0 \rightarrow \langle =\text{subj } t \rangle_0 \quad \langle \text{subj} \rangle_1 \quad \exp(\lambda_{\text{MERGE}} + \lambda_{\text{subj}})$

$\langle t \rangle_0 \rightarrow \langle +\text{subj } t, -\text{subj} \rangle_0 \quad \exp(\lambda_{\text{MOVE}} + \lambda_{\text{subj}})$

Wh-phrase same as subject or separated because of doubling?

$\langle t, -wh \rangle_0 \rightarrow \langle =\text{subj } t \rangle_0 \quad \langle \text{subj -wh} \rangle_1 \quad \exp(\lambda_{\text{MERGE}} + \lambda_{\text{subj}})$

$\langle t, -wh \rangle_0 \rightarrow \langle +\text{subj } t, -\text{subj}, -wh \rangle_0 \quad \exp(\lambda_{\text{MOVE}} + \lambda_{\text{subj}})$

Choice points in the IMG-derived MCFG

Question or not?

$\langle -c \rangle_0$	\rightarrow	$\langle +t -c, -t \rangle_1$	$\exp(\lambda_{\text{MRG}} + \lambda_t)$
$\langle -c \rangle_0$	\rightarrow	$\langle +wh -c, -wh \rangle_0$	$\exp(\lambda_{\text{MRG}} + \lambda_{wh})$

Non-wh antecedent lexical or complex?

$\langle +subj -ant -subj, -subj \rangle_0$	\rightarrow	$\langle +subj -ant -subj \rangle_0$	$\langle -subj \rangle_0$	$\exp(\lambda_{\text{INSERT}})$
$\langle +subj -ant -subj, -subj \rangle_0$	\rightarrow	$\langle +subj -ant -subj \rangle_0$	$\langle -subj \rangle_1$	$\exp(\lambda_{\text{INSERT}})$

Non-wh subject merged and complex, merged and lexical, or moved?

$\langle +subj -t, -subj \rangle_0$	\rightarrow	$\langle +subj -t \rangle_0$	$\langle -subj \rangle_0$	$\exp(\lambda_{\text{INSERT}})$
$\langle +subj -t, -subj \rangle_0$	\rightarrow	$\langle +subj -t \rangle_0$	$\langle -subj \rangle_1$	$\exp(\lambda_{\text{INSERT}})$
$\langle +subj -t, -subj \rangle_0$	\rightarrow	$\langle +v +subj -t, -v, -subj \rangle_1$		$\exp(\lambda_{\text{MRG}} + \lambda_v)$

Wh-phrase same as subject or separated because of doubling?

$\langle -t, -wh \rangle_0$	\rightarrow	$\langle +subj -t, -subj -wh \rangle_0$	$\exp(\lambda_{\text{MRG}} + \lambda_{\text{subj}})$
$\langle -t, -wh \rangle_0$	\rightarrow	$\langle +subj -t, -subj, -wh \rangle_0$	$\exp(\lambda_{\text{MRG}} + \lambda_{\text{subj}})$

Learned weights on the MG

$$\lambda_t = 0.094350 \quad \exp(\lambda_t) = 1.0989$$

$$\lambda_{\text{subj}} = -5.734063 \quad \exp(\lambda_{\text{subj}}) = 0.0032$$

$$\lambda_{\text{wh}} = -0.094350 \quad \exp(\lambda_{\text{wh}}) = 0.9100$$

$$\lambda_{\text{MERGE}} = 0.629109 \quad \exp(\lambda_{\text{MERGE}}) = 1.8759$$

$$\lambda_{\text{MOVE}} = -0.629109 \quad \exp(\lambda_{\text{MOVE}}) = 0.5331$$

$$P(\text{antecedent is lexical}) = 0.5$$

$$P(\text{antecedent is non-lexical}) = 0.5$$

$$P(\text{wh is reflexivized}) = \frac{\exp(\lambda_{\text{MOVE}})}{\exp(\lambda_{\text{MERGE}}) + \exp(\lambda_{\text{MOVE}})} = 0.2213$$

$$P(\text{wh not reflexivized}) = \frac{\exp(\lambda_{\text{MERGE}})}{\exp(\lambda_{\text{MERGE}}) + \exp(\lambda_{\text{MOVE}})} = 0.7787$$

$$P(\text{question}) = \frac{\exp(\lambda_{\text{MOVE}} + \lambda_{\text{wh}})}{\exp(\lambda_{\text{MERGE}} + \lambda_t) + \exp(\lambda_{\text{MOVE}} + \lambda_{\text{wh}})} = 0.1905$$

$$P(\text{non-question}) = \frac{\exp(\lambda_{\text{MERGE}} + \lambda_t)}{\exp(\lambda_{\text{MERGE}} + \lambda_t) + \exp(\lambda_{\text{MOVE}} + \lambda_{\text{wh}})} = 0.8095$$

$$P(\text{non-wh subject merged and complex}) = \frac{\exp(\lambda_{\text{MERGE}})}{\exp(\lambda_{\text{MERGE}}) + \exp(\lambda_{\text{MERGE}}) + \exp(\lambda_{\text{MOVE}})} = 0.4378$$

$$P(\text{non-wh subject merged and lexical}) = \frac{\exp(\lambda_{\text{MERGE}})}{\exp(\lambda_{\text{MERGE}}) + \exp(\lambda_{\text{MERGE}}) + \exp(\lambda_{\text{MOVE}})} = 0.4378$$

$$P(\text{non-wh subject moved}) = \frac{\exp(\lambda_{\text{MOVE}})}{\exp(\lambda_{\text{MERGE}}) + \exp(\lambda_{\text{MERGE}}) + \exp(\lambda_{\text{MOVE}})} = 0.1244$$

$$P(\text{who will shave}) = 0.1905 \times 0.7787 = 0.148$$

$$P(\text{boys will shave themselves}) = 0.5 \times 0.8095 \times 0.1244 = 0.050$$

Learned weights on the IMG

$$\lambda_t = 0.723549$$

$$\exp(\lambda_t) = 2.0617$$

$$P(\text{antecedent is lexical}) = 0.5$$

$$\lambda_v = 0.440585$$

$$\exp(\lambda_v) = 1.5536$$

$$P(\text{antecedent is non-lexical}) = 0.5$$

$$\lambda_{wh} = -0.723459$$

$$\exp(\lambda_{wh}) = 0.4850$$

$$P(\text{wh-phrase reflexivized}) = 0.5$$

$$\lambda_{\text{INSERT}} = 0.440585$$

$$\exp(\lambda_{\text{INSERT}}) = 1.5536$$

$$P(\text{wh-phrase non-reflexivized}) = 0.5$$

$$\lambda_{\text{MRG}} = -0.440585$$

$$\exp(\lambda_{\text{MRG}}) = 0.6437$$

$$P(\text{question}) = \frac{\exp(\lambda_{\text{MRG}} + \lambda_{wh})}{\exp(\lambda_{\text{MRG}} + \lambda_t) + \exp(\lambda_{\text{MRG}} + \lambda_{wh})} = \frac{\exp(\lambda_{wh})}{\exp(\lambda_t) + \exp(\lambda_{wh})} = 0.1905$$

$$P(\text{non-question}) = \frac{\exp(\lambda_{\text{MRG}} + \lambda_t)}{\exp(\lambda_{\text{MRG}} + \lambda_t) + \exp(\lambda_{\text{MRG}} + \lambda_{wh})} = \frac{\exp(\lambda_t)}{\exp(\lambda_t) + \exp(\lambda_{wh})} = 0.8095$$

$$P(\text{non-wh subject merged and lexical}) = \frac{\exp(\lambda_{\text{INSERT}})}{\exp(\lambda_{\text{INSERT}}) + \exp(\lambda_{\text{INSERT}}) + \exp(\lambda_{\text{MRG}} + \lambda_v)} = 0.4412$$

$$P(\text{non-wh subject merged and complex}) = \frac{\exp(\lambda_{\text{INSERT}})}{\exp(\lambda_{\text{INSERT}}) + \exp(\lambda_{\text{INSERT}}) + \exp(\lambda_{\text{MRG}} + \lambda_v)} = 0.4412$$

$$P(\text{non-wh subject moved}) = \frac{\exp(\lambda_{\text{MRG}} + \lambda_v)}{\exp(\lambda_{\text{INSERT}}) + \exp(\lambda_{\text{INSERT}}) + \exp(\lambda_{\text{MRG}} + \lambda_v)} = 0.1176$$

$$P(\text{who will shave}) = 0.5 \times 0.1905 = 0.095$$

$$P(\text{boys will shave themselves}) = 0.5 \times 0.8095 \times 0.1176 = 0.048$$

Surprisal predictions

Grammar: MG, i.e. MERGE and MOVE

Sentence: 'who will shave themselves'

MG, i.e. MERGE and MOVE	
0.35478	boys will shave
0.35478	some boys will shave
0.14801	who will shave
0.05022	boys will shave themselves
0.05022	some boys will shave themselves
0.04199	who will shave themselves

Surprisal predictions

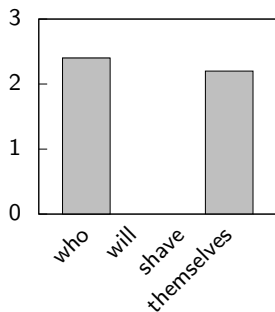
Grammar: MG, i.e. MERGE and MOVE

Sentence: 'who will shave themselves'

$$\begin{aligned} \text{surprisal at 'who'} &= -\log P(W_1 = \text{who}) \\ &= -\log(0.15 + 0.04) \\ &= -\log 0.19 \\ &= 2.4 \end{aligned}$$

$$\begin{aligned} \text{surprisal at 'themselves'} &= -\log P(W_4 = \text{themselves} \mid W_1 = \text{who}, \dots) \\ &= -\log \frac{0.04}{0.15 + 0.04} \\ &= -\log 0.21 \\ &= 2.2 \end{aligned}$$

MG, i.e. MERGE and MOVE	
0.35478	boys will shave
0.35478	some boys will shave
0.14801	who will shave
0.05022	boys will shave themselves
0.05022	some boys will shave themselves
0.04199	who will shave themselves



Surprisal predictions

Grammar: IMG, i.e. MRG and INSERT

Sentence: 'who will shave themselves'

IMG, i.e. MRG and INSERT	
0.35721	boys will shave
0.35721	some boys will shave
0.095	who will shave
0.095	who will shave themselves
0.04779	boys will shave themselves
0.04779	some boys will shave themselves

Surprisal predictions

Grammar: IMG, i.e. MRG and INSERT

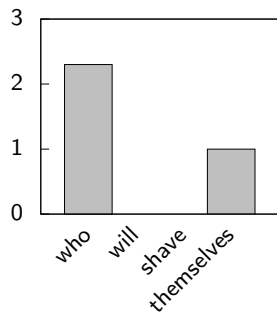
Sentence: 'who will shave themselves'

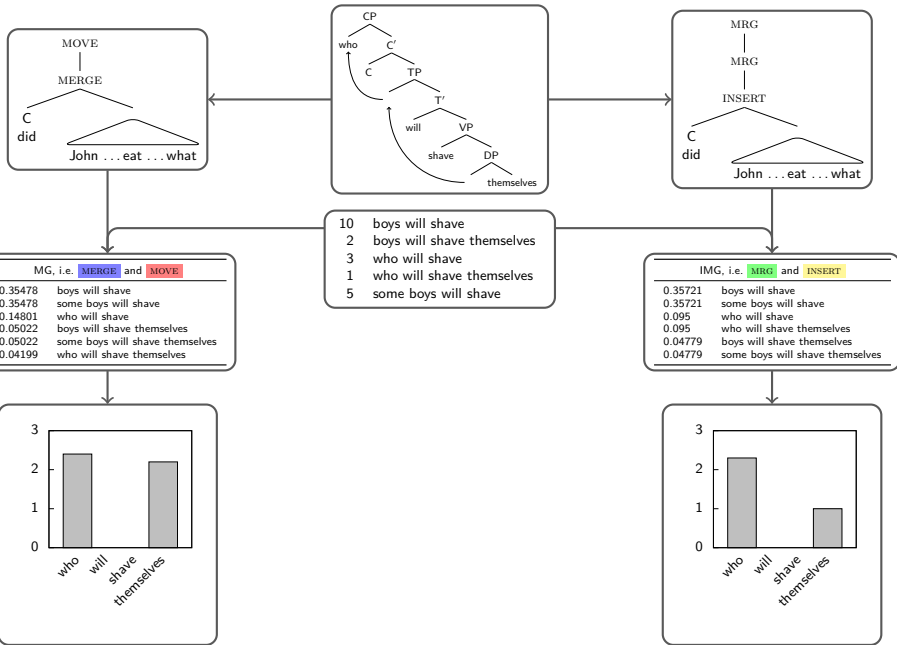
IMG, i.e. MRG and INSERT

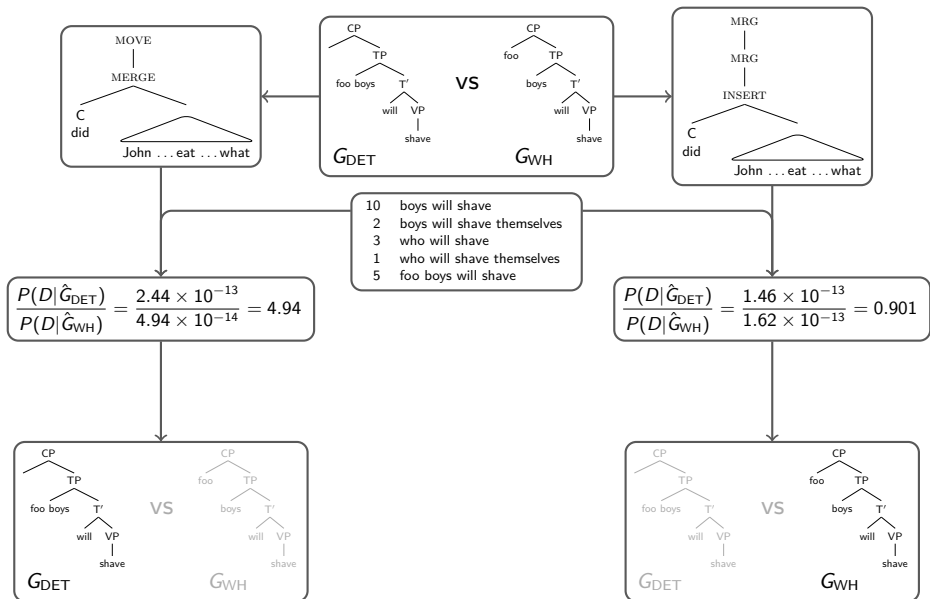
0.35721	boys will shave
0.35721	some boys will shave
0.095	who will shave
0.095	who will shave themselves
0.04779	boys will shave themselves
0.04779	some boys will shave themselves

$$\begin{aligned}
 \text{surprisal at 'who'} &= -\log P(W_1 = \text{who}) \\
 &= -\log(0.10 + 0.10) \\
 &= -\log 0.2 \\
 &= 2.3
 \end{aligned}$$

$$\begin{aligned}
 \text{surprisal at 'themselves'} &= -\log P(W_4 = \text{themselves} \mid W_1 = \text{who}, \dots) \\
 &= -\log \frac{0.10}{0.10 + 0.10} \\
 &= -\log 0.5 \\
 &= 1
 \end{aligned}$$







Outline

- 1 What are grammars?
- 2 Derivationally distinct implementations of merge
- 3 Telling them apart
- 4 Historical perspective**

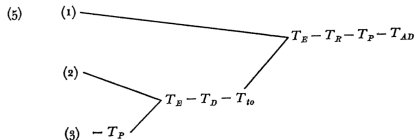
Back in the good old (heavily derivational) days

[The perceptual model] will utilize the full resources of the transformational grammar to provide a structural description, consisting of a set of P-markers and a transformational history

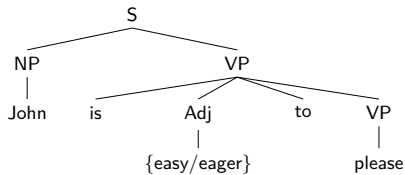
(Miller and Chomsky 1963: p.480)

(4) the man who persuaded John to be examined by a specialist
was fired

The “transformational history” of (4) by which it is derived from its basis might be represented, informally, by the diagram (5).



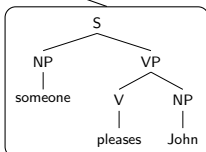
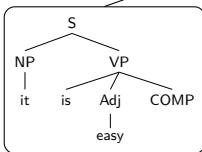
Full derivational histories



T_5 : front embedded object, replacing 'it'

T_4 : delete 'for someone'

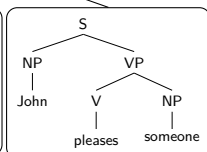
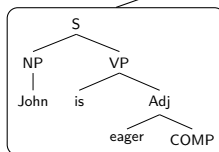
T_1 : replace COMP

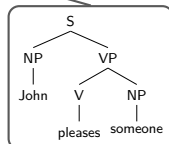
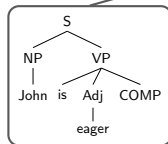
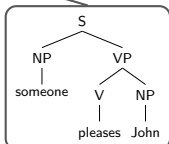
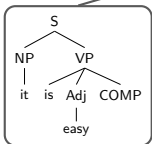
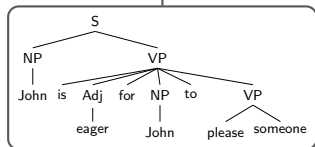
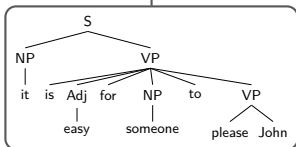
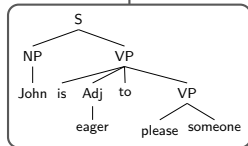
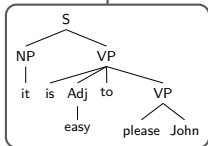
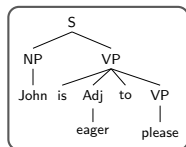
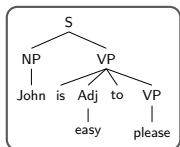


T_3 : delete object

T_2 : delete duplicate NP

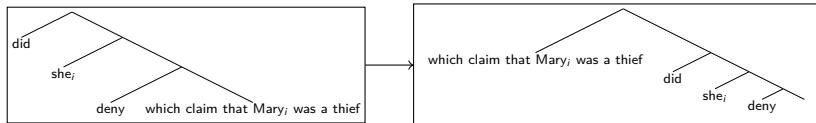
T_1 : replace COMP





Full derivational histories

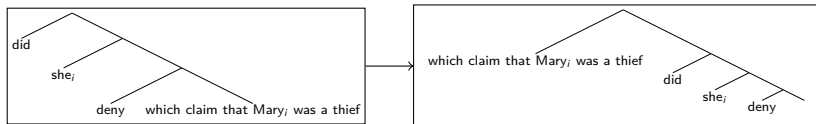
(1) * Which claim [that Mary_i was a thief] did she_i deny?



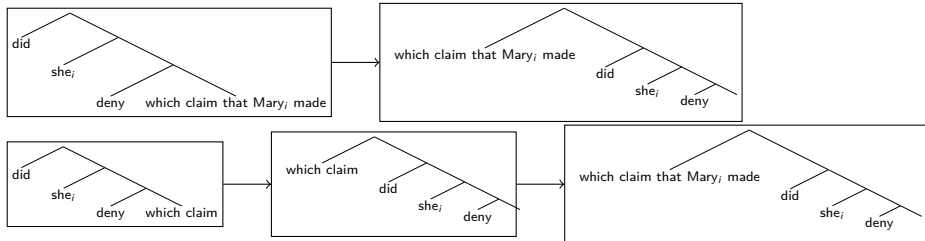
(2) Which claim [that Mary_i made] did she_i deny?

Full derivational histories

(1) * Which claim [that Mary_i was a thief] did she_i deny?

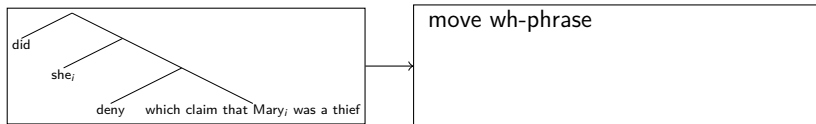


(2) Which claim [that Mary_i made] did she_i deny?

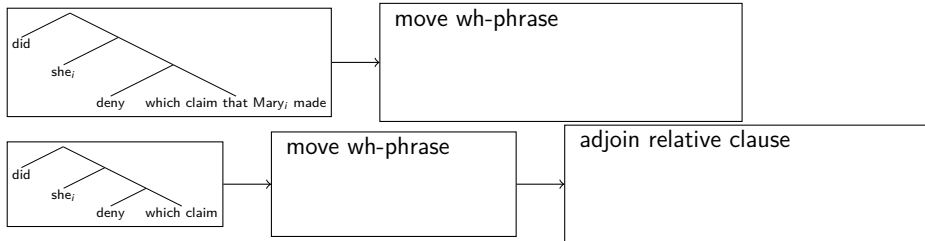


Full derivational histories

(1) * Which claim [that Mary_i was a thief] did she_i deny?



(2) Which claim [that Mary_i made] did she_i deny?



Conclusion and open issues

We can formulate a theory where

- the choice of derivational operations has empirically-testable consequences (so we are **not extensionalists**), and
- this does not happen by taking derivational operations to be real-time operations (because we are **not literalists**).

Conclusion and open issues

We can formulate a theory where

- the choice of derivational operations has empirically-testable consequences (so we are **not extensionalists**), and
- this does not happen by taking derivational operations to be real-time operations (because we are **not literalists**).

Unanswered question: “So **what are** the real-time operations?”

- For complementary proposals on this see Stabler (2013), Kobele et al. (2013), Graf et al. (2015), Hunter (forthcoming)
- . . . but all of these take the form of procedures for **identifying a derivation tree/T-marker**.
- Distinct from the question of what the chunks to be (somehow) chained together are.

References I

- Graf, T., Fodor, B., Monette, J., Rachiele, G., Warren, A., and Zhang, C. (2015). A refined notion of memory usage for minimalist parsing. In *Proceedings of the 14th Meeting on the Mathematics of Language (MoL 2015)*, pages 1–14. Association for Computational Linguistics.
- Hunter, T. (2011). Insertion Minimalist Grammars: Eliminating redundancies between merge and move. In Kanazawa, M., Kornai, A., Kracht, M., and Seki, H., editors, *The Mathematics of Language (MOL 12 Proceedings)*, volume 6878 of *LNCS*, pages 90–107, Berlin Heidelberg. Springer.
- Hunter, T. (forthcoming). Left-corner parsing of minimalist grammars. In Berwick, B. and Stabler, E., editors, *Minimalist Parsing*. Oxford University Press.
- Hunter, T. and Dyer, C. (2013). Distributions on minimalist grammar derivations. In *Proceedings of the 13th Meeting on the Mathematics of Language*.
- Keenan, E. L. and Stabler, E. P. (2003). *Bare Grammar*. CSLI Publications, Stanford, CA.
- Kobele, G. M., Gerth, S., and Hale, J. (2013). Memory resource allocation in top-down minimalist parsing. In Morrill, G. and Nederhof, M.-J., editors, *Formal Grammar 2012/2013*, volume 8036 of *Lecture Notes in Computer Science*, pages 32–51. Springer.
- Miller, G. A. and Chomsky, N. (1963). Finitary models of language users. In Luce, R. D., Bush, R. R., and Galanter, E., editors, *Handbook of Mathematical Psychology*, volume 2. Wiley and Sons, New York.

References II

- Phillips, C. and Lewis, S. (2013). Derivational order in syntax: evidence and architectural consequences. *Studies in Linguistics*, 6:11–47.
- Stabler, E. (2013). Two models of minimalist, incremental syntactic analysis. *Topics in Cognitive Science*, 5(3):611–633.
- Stabler, E. P. (1997). Derivational minimalism. In Retoré, C., editor, *Logical Aspects of Computational Linguistics*, volume 1328 of *LNCS*, pages 68–95, Berlin Heidelberg. Springer.
- Stabler, E. P. (2006). Sideways without copying. In Wintner, S., editor, *Proceedings of The 11th Conference on Formal Grammar*, pages 157–170, Stanford, CA. CSLI Publications.