

Sharpening the empirical claims of generative syntax through formalization

Tim Hunter

University of Minnesota, Twin Cities

NASSLLI, June 2014

Part 1: Grammars and cognitive hypotheses

What is a grammar?

What can grammars do?

Concrete illustration of a target: Surprisal

Parts 2–4: Assembling the pieces

Minimalist Grammars (MGs)

MGs and MCFGs

Probabilities on MGs

Part 5: Learning and wrap-up

Something slightly different: Learning model

Recap and open questions

Sharpening the empirical claims of generative syntax
through formalization

Tim Hunter — NASSLLI, June 2014

Part 1

Grammars and Cognitive Hypotheses

Outline

- 1 What we want to do with grammars
- 2 How to get grammars to do it
- 3 Derivations and representations
- 4 Information-theoretic complexity metrics

Outline

- 1 What we want to do with grammars
- 2 How to get grammars to do it
- 3 Derivations and representations
- 4 Information-theoretic complexity metrics

Claims made by grammars

What are grammars used for?

- “Mostly” for accounting for acceptability judgements
- But there are other ways a grammar can figure in claims about cognition

Claims made by grammars

What are grammars used for?

- “Mostly” for accounting for acceptability judgements
- But there are other ways a grammar can figure in claims about cognition

Often tempting to draw a distinction between “linguistic evidence” (where grammar lives) and “experimental evidence” (where cognition lives)

- One need not make this distinction
- We will proceed without it, i.e. it's all linguistic (and/or all experimental)

Claims made by grammars

There's a "boring" sense in which every syntax paper makes a cognitive claim, i.e. a claim testable via acceptability facts.

Claims made by grammars

There's a "boring" sense in which every syntax paper makes a cognitive claim, i.e. a claim testable via acceptability facts.

- For one thing, this is not a cop-out, even if it might seem like it is!
 - Why does it seem like a cop-out?
 - Lingering externalism/Platonism?
 - Perhaps partly because it's just relatively rare to see anything being tested by other measures

Claims made by grammars

There's a "boring" sense in which every syntax paper makes a cognitive claim, i.e. a claim testable via acceptability facts.

- For one thing, this is not a cop-out, even if it might seem like it is!
 - Why does it seem like a cop-out?
 - Lingering externalism/Platonism?
 - Perhaps partly because it's just relatively rare to see anything being tested by other measures
- For another, we can incorporate grammars into claims that are testable by other measures.
 - **This is the main point of the course!**
 - The claims/predictions will depend on internal properties of grammars, not just what they say is good and what they say is bad
 - And we'll do it without seeing grammatical derivations as real-time operations

Claims made by grammars

If we accept — as I do — ... that the rules of grammar enter into the processing mechanisms, then evidence concerning production, recognition, recall, and language use in general can be expected (in principle) to have bearing on the investigation of rules of grammar, on what is sometimes called “grammatical competence” or “knowledge of language”.

(Chomsky 1980: pp.200-201)

[S]ince a competence theory must be incorporated in a performance model, evidence about the actual organization of behavior may prove crucial to advancing the theory of underlying competence.

(Chomsky 1980: p.226)

Claims made by grammars

If we accept — as I do — ... that the rules of grammar enter into the processing mechanisms, then evidence concerning production, recognition, recall, and language use in general can be expected (in principle) to have bearing on the investigation of rules of grammar, on what is sometimes called “grammatical competence” or “knowledge of language”.

(Chomsky 1980: pp.200-201)

[S]ince a competence theory must be incorporated in a performance model, evidence about the actual organization of behavior may prove crucial to advancing the theory of underlying competence.

(Chomsky 1980: p.226)

Evidence about X can only advance Y if Y makes claims about X!

Preview

What we will do:

- Put together a chain of linking hypotheses that bring “experimental evidence” to bear on “grammar questions”
 - e.g. reading times, acquisition patterns
 - e.g. move as distinct operation from merge vs. unified with merge
- Illustrate with some toy examples

What we will not do:

- Engage with state-of-the-art findings in the sentence processing literature
- End up with claims that one particular set of derivational operations is empirically better than another

Teasers

We'll take pairs of equivalent grammars that differ only in the move/re-merge dimension.

- They will make different predictions about sentence comprehension difficulty.
- They will make different predictions about what a learner will conclude from a common input corpus.

Teasers

We'll take pairs of equivalent grammars that differ only in the move/re-merge dimension.

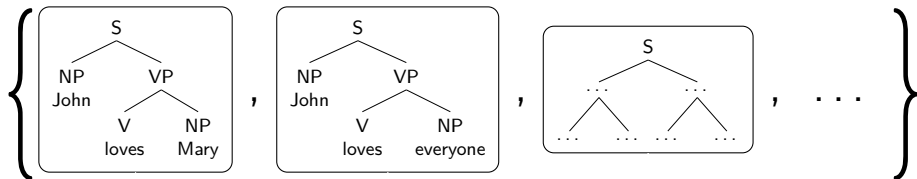
- They will make different predictions about sentence comprehension difficulty.
- They will make different predictions about what a learner will conclude from a common input corpus.

The issues become “distant but empirical questions”. That's all we're aiming for, for now.

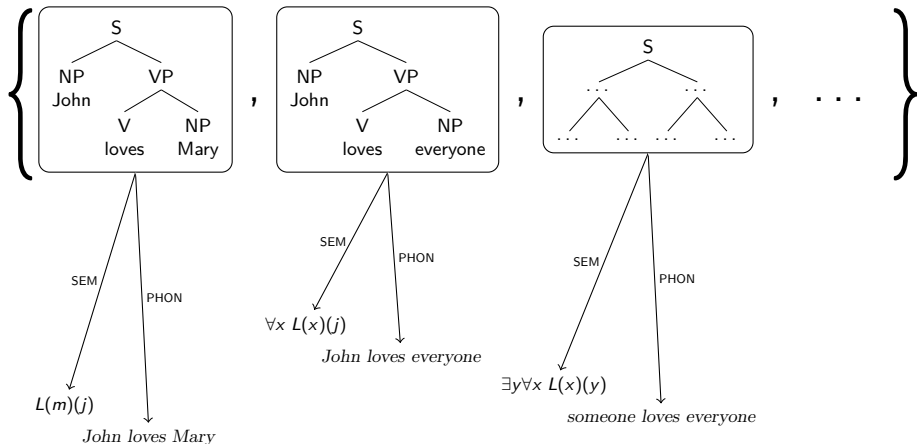
Outline

- 1 What we want to do with grammars
- 2 How to get grammars to do it**
- 3 Derivations and representations
- 4 Information-theoretic complexity metrics

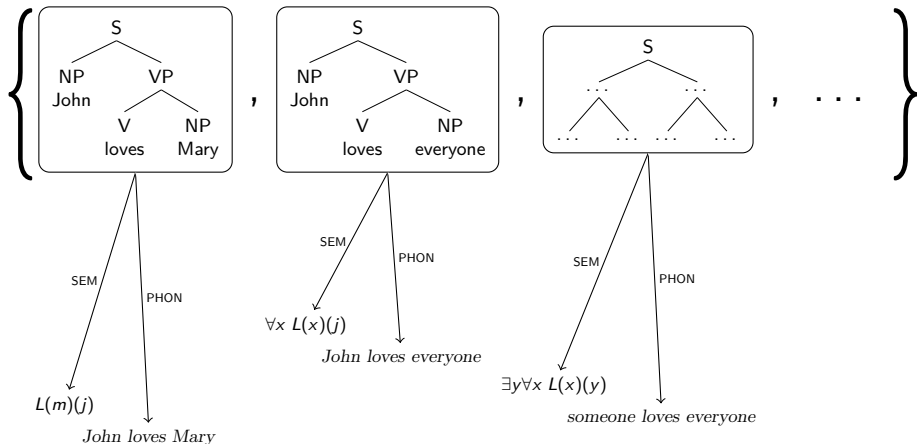
Interpretation functions



Interpretation functions



Interpretation functions



Caveats:

- Maybe we're interested in the finite specification of the set
- Maybe there's no clear line between observable and not
- Maybe some evidence is based on relativities among interpretations

Telling grammars apart

So, what if we have two different grammars — systems that define different sets of objects — that we can't tell apart via the sound and meaning interpretations?

(Perhaps because they're provably equivalent, or perhaps because the evidence just happens to be unavailable.)

Telling grammars apart

So, what if we have two different grammars — systems that define different sets of objects — that we can't tell apart via the sound and meaning interpretations?

(Perhaps because they're provably equivalent, or perhaps because the evidence just happens to be unavailable.)

- Option 1: Conclude that the differences are irrelevant to us (or “they're not actually different”).
- Option 2: Make the differences matter ... somehow ...

What are syntactic representations for?

Morrill (1994) in favour of Option 1:

*The construal of a language as a collection of signs [sound-meaning pairs] presents as an investigative task the characterisation of this collection. This is usually taken to mean the specification of a set of "structural descriptions" (or: "syntactic structures"). Observe however that on our understanding a sign is an association of prosodic [phonological] and semantic properties. It is these properties that can be observed and that are to be modelled. There appears to be no observation which bears directly on syntactic as opposed to prosodic and/or semantic properties, and this implies an asymmetry in the status of these levels. **A structural description is only significant insofar as it is understood as predicting prosodic and semantic properties (e.g. in interpreting the yield of a tree as word order). Attribution of syntactic (or prosodic or semantic) structure does not of itself predict anything.***

What are syntactic representations for?

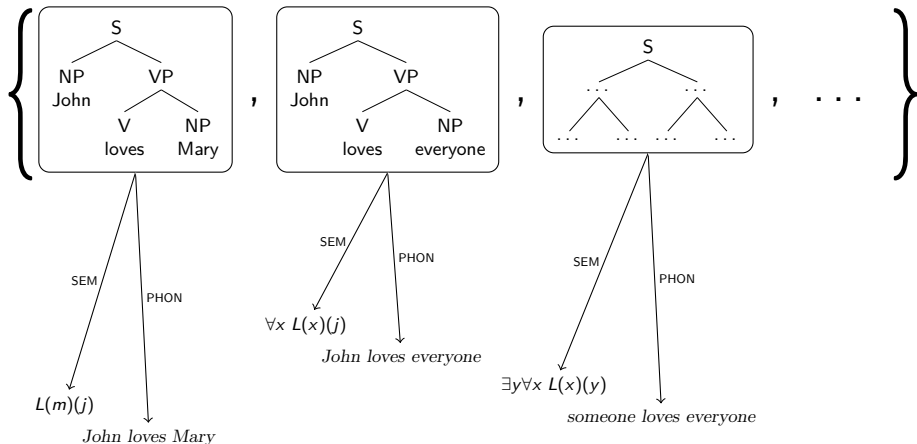
Morrill (1994) in favour of Option 1:

*The construal of a language as a collection of signs [sound-meaning pairs] presents as an investigative task the characterisation of this collection. This is usually taken to mean the specification of a set of “structural descriptions” (or: “syntactic structures”). Observe however that on our understanding a sign is an association of prosodic [phonological] and semantic properties. It is these properties that can be observed and that are to be modelled. There appears to be no observation which bears directly on syntactic as opposed to prosodic and/or semantic properties, and this implies an asymmetry in the status of these levels. **A structural description is only significant insofar as it is understood as predicting prosodic and semantic properties (e.g. in interpreting the yield of a tree as word order). Attribution of syntactic (or prosodic or semantic) structure does not of itself predict anything.***

Where might we depart from this (to pursue Option 2)?

- Object that syntactic structure **does** matter “of itself”
- Object that prosodic and semantic properties are **not** the only ones we can observe

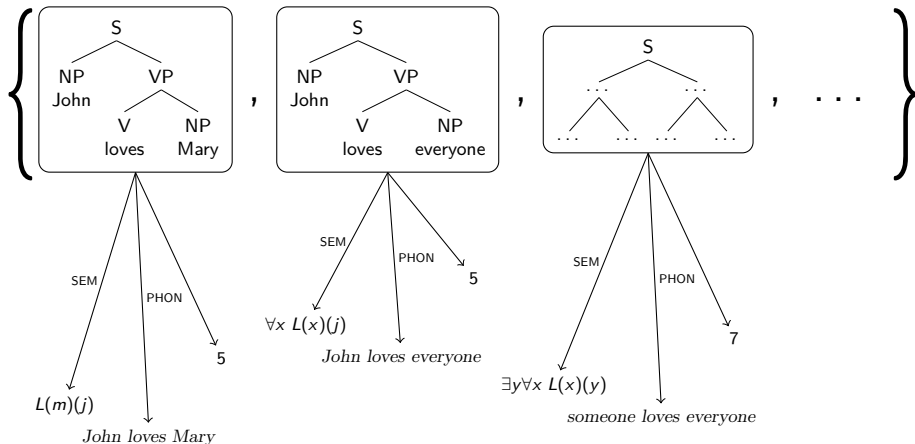
Interpretation functions



Caveats:

- Maybe we're interested in the finite specification of the set
- Maybe there's no clear line between observable and not
- Maybe some evidence is based on relativities among interpretations

Interpretation functions



Caveats:

- Maybe we're interested in the finite specification of the set
- Maybe there's no clear line between observable and not
- Maybe some evidence is based on relativities among interpretations

Interpretation functions for “complexity”

What are some other interpretation functions?

- number of nodes

Interpretation functions for “complexity”

What are some other interpretation functions?

- number of nodes
- ratio of total nodes to terminal nodes (Miller and Chomsky 1963)

Ratio of total nodes to terminal nodes

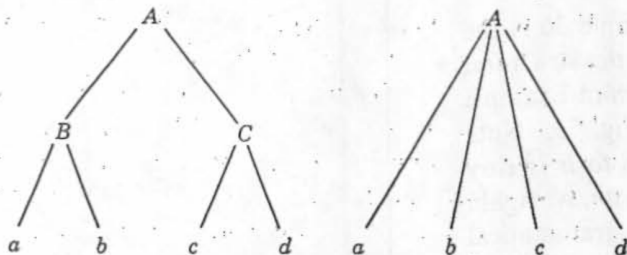


Fig. 8. Illustrating a measure of structural complexity. $N(Q)$ for the P -marker (a) is $7/4$; for (b), $N(Q) = 5/4$.

Ratio of total nodes to terminal nodes

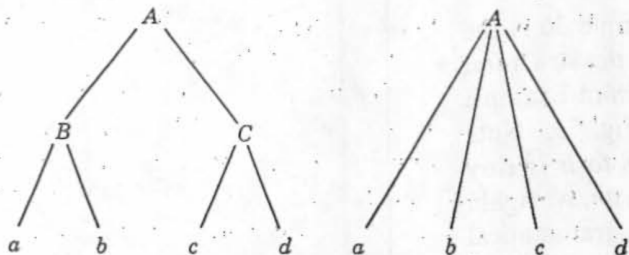


Fig. 8. Illustrating a measure of structural complexity. $N(Q)$ for the P -marker (a) is $7/4$; for (b), $N(Q) = 5/4$.

Won't distinguish center-embedding from left- and right-embedding

- | | | |
|-----|--|----------|
| (1) | The mouse [the cat [the dog bit] chased] died. | (center) |
| (2) | The dog bit the cat [which chased the mouse [which died]]. | (right) |
| (3) | [[the dog] 's owner] 's friend | (left) |

Interpretation functions for “complexity”

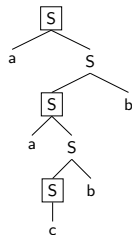
What are some other interpretation functions?

- number of nodes
- ratio of total nodes to terminal nodes (Miller and Chomsky 1963)
- degree of self-embedding (Miller and Chomsky 1963)

Degree of (centre-)self-embedding

A tree's degree of self-embedding is m iff:

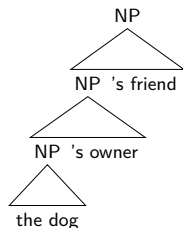
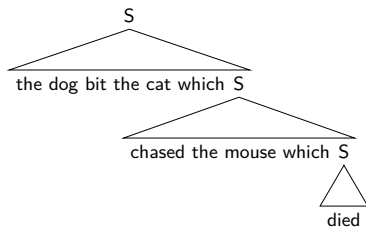
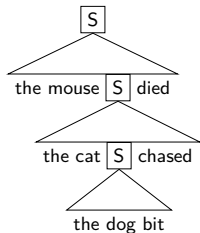
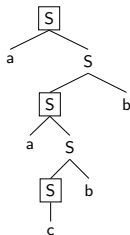
“there is ... a continuous path passing through $m + 1$ nodes N_0, \dots, N_m , each with the same label, where each N_i ($i \geq 1$) is fully self-embedded (with something to the left and something to the right) in the subtree dominated by N_{i-1} ”



Degree of (centre-)self-embedding

A tree's degree of self-embedding is m iff:

"there is ... a continuous path passing through $m + 1$ nodes N_0, \dots, N_m , each with the same label, where each N_i ($i \geq 1$) is fully self-embedded (with something to the left and something to the right) in the subtree dominated by N_{i-1} "



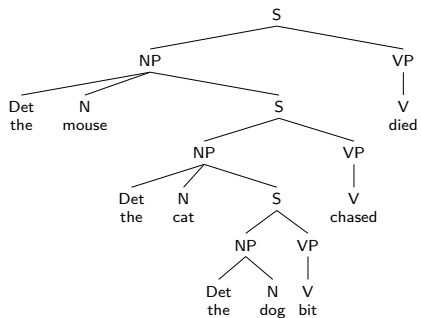
Interpretation functions for “complexity”

What are some other interpretation functions?

- number of nodes
- ratio of total nodes to terminal nodes (Miller and Chomsky 1963)
- degree of self-embedding (Miller and Chomsky 1963)
- “depth” of memory required by a top-down parser (Yngve 1960)

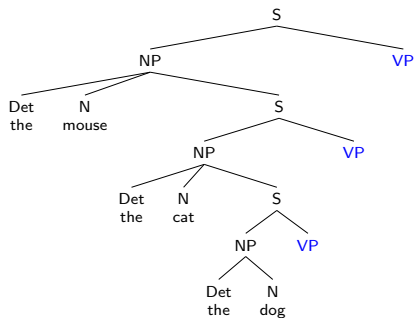
Yngve's depth

Number of constituents expected but not yet started:



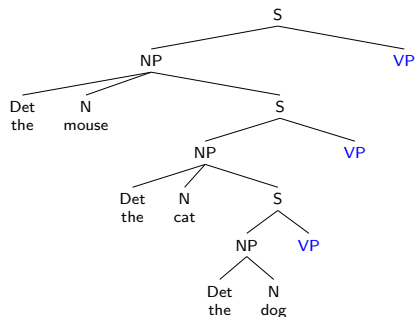
Yngve's depth

Number of constituents expected but not yet started:



Yngve's depth

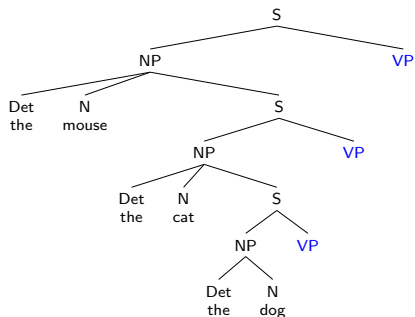
Number of constituents expected but not yet started:



- Unlike (center-)self-embedding, right-embedding doesn't create such large lists of expected constituents (because the expected stuff is all part of **one** constituent).
- But left-embedding does.

Yngve's depth

Number of constituents expected but not yet started:



- Unlike (center-)self-embedding, right-embedding doesn't create such large lists of expected constituents (because the expected stuff is all part of **one** constituent).
- But left-embedding does.
- Yngve's theory was set within — perhaps justified by — a procedural story, but we can arguably detach it from that and treat depth as just another property of trees.

Interpretation functions for “complexity”

What are some other interpretation functions?

- number of nodes
- ratio of total nodes to terminal nodes (Miller and Chomsky 1963)
- degree of self-embedding (Miller and Chomsky 1963)
- “depth” of memory required by a top-down parser (Yngve 1960)
- minimal attachment, late closure, etc.?

Reaching conclusions about grammars

complexity metric + grammar \longrightarrow prediction

Typically, arguments hold the grammar fixed and present evidence in favour of a metric.

Reaching conclusions about grammars

complexity metric + grammar \longrightarrow prediction

Typically, arguments hold the grammar fixed and present evidence in favour of a metric.

We can flip this around: hold the metric fixed and present evidence in favour of a grammar.

If we accept — as I do — ... that the rules of grammar enter into the processing mechanisms, then evidence concerning production, recognition, recall, and language use in general can be expected (in principle) to have bearing on the investigation of rules of grammar, on what is sometimes called “grammatical competence” or “knowledge of language”.

(Chomsky 1980: pp.200-201)

Reaching conclusions about grammars

complexity metric + grammar \longrightarrow prediction

Example: hold [self-embedding](#) fixed as the complexity metric.

- (4) That [the food that [John ordered] tasted good] pleased him.
- (5) That [that [the food was good] pleased John] surprised Mary.

Grammar question: Does a relative clause have a node labeled S?

Reaching conclusions about grammars

complexity metric + grammar \rightarrow prediction

Example: hold **self-embedding** fixed as the complexity metric.

- (4) That [the food that [John ordered] tasted good] pleased him.
- (5) That [that [the food was good] pleased John] surprised Mary.

Grammar question: Does a relative clause have a node labeled S?

Proposed answer	(4) structure	(5) structure	Prediction
Yes	$\dots [S \dots [S \dots]]$	$\dots [S \dots [S \dots]]$	(4) & (5) same
No	$\dots [S \dots [RC \dots]]$	$\dots [S \dots [S \dots]]$	(5) harder

Reaching conclusions about grammars

complexity metric + grammar \rightarrow prediction

Example: hold **self-embedding** fixed as the complexity metric.

(4) That [the food that [John ordered] tasted good] pleased him.

(5) That [that [the food was good] pleased John] surprised Mary.

Grammar question: Does a relative clause have a node labeled S?

Proposed answer	(4) structure	(5) structure	Prediction
Yes	$\dots [S \dots [S \dots]]$	$\dots [S \dots [S \dots]]$	(4) & (5) same
No	$\dots [S \dots [RC \dots]]$	$\dots [S \dots [S \dots]]$	(5) harder

Conclusion: The fact that (5) is harder supports the “No” answer.

Outline

- 1 What we want to do with grammars
- 2 How to get grammars to do it
- 3 Derivations and representations**
- 4 Information-theoretic complexity metrics

Derivations and representations

Question

But these metrics are all properties of a final, **fully-constructed tree**.

How can anything like this be sensitive to differences in the derivational operations that build these trees? (e.g. TAG vs. MG, whether move is re-merge)

Interpretation functions for “complexity”

What are some other interpretation functions?

- number of nodes
- ratio of total nodes to terminal nodes (Miller and Chomsky 1963)
- degree of self-embedding (Miller and Chomsky 1963)
- “depth” of memory required by a top-down parser (Yngve 1960)
- minimal attachment, late closure, etc.?
- “nature, number and complexity of” transformations (Miller and Chomsky 1963)

“nature, number and complexity of the grammatical transformations involved”

*The psychological plausibility of a transformational model of the language user would be strengthened, of course, if it could be shown that our performance on tasks requiring an appreciation of the structure of transformed sentences is **some function of the nature, number and complexity of the grammatical transformations involved.***

(Miller and Chomsky 1963: p.481)

Derivations and representations

Question

But these metrics are all properties of a final, **fully-constructed tree**.

How can anything like this be sensitive to differences in the derivational operations that build these trees? (e.g. TAG vs. MG, whether move is re-merge)

Derivations and representations

Question

But these metrics are all properties of a final, **fully-constructed tree**.

How can anything like this be sensitive to differences in the derivational operations that build these trees? (e.g. TAG vs. MG, whether move is re-merge)

Answer

The relevant objects on which the interpretation functions are defined encode a complete **derivational history**.

Derivations and representations

Question

But these metrics are all properties of a final, **fully-constructed tree**.

How can anything like this be sensitive to differences in the derivational operations that build these trees? (e.g. TAG vs. MG, whether move is re-merge)

Answer

The relevant objects on which the interpretation functions are defined encode a complete **derivational history**.

e.g. The function which, given a complete “recipe” for carrying out a derivation, returns the number of movement steps called for by the recipe.

(Maybe only useful when we’re holding a grammar fixed)

Full derivation recipes?

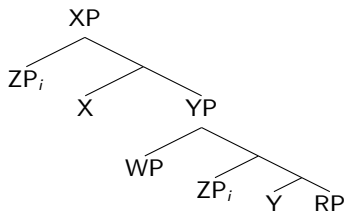
Are the inputs to these functions **really** full derivation recipes?

For minimalist syntax it's hard to tell, because the final **derived object** very often uniquely identifies a **derivational history/recipe**.

Full derivation recipes?

Are the inputs to these functions **really** full derivation recipes?

For minimalist syntax it's hard to tell, because the final **derived object** very often uniquely identifies a **derivational history/recipe**.



- merge Y with RP
- merge the result with ZP
- merge the result with WP
- merge X with the result
- move ZP

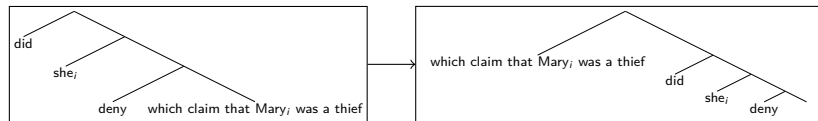
Full derivation recipes?

A few cases reveal that (we must all be already assuming that) it's full derivations/recipes that count.

Full derivation recipes?

A few cases reveal that (we must all be already assuming that) it's full derivations/recipes that count.

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

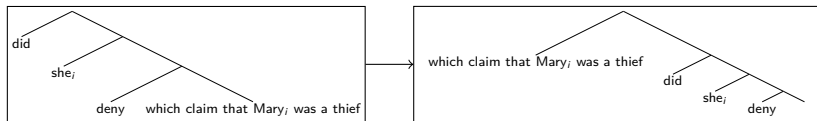


(7) Which claim [that $Mary_i$ made] did she_i deny?

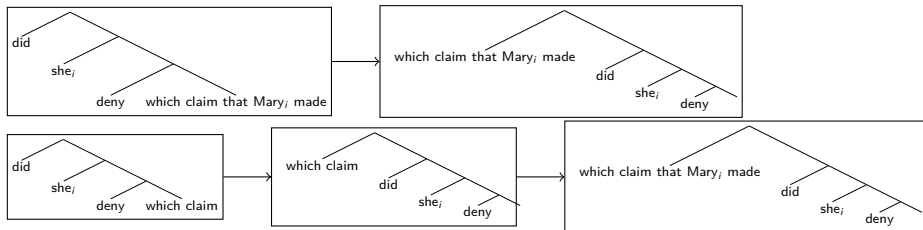
Full derivation recipes?

A few cases reveal that (we must all be already assuming that) it's full derivations/recipes that count.

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



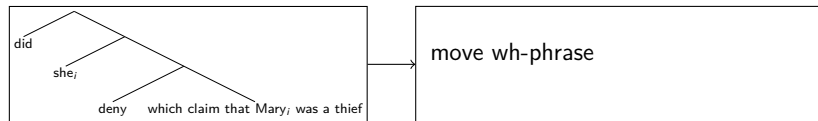
(7) Which claim [that Mary_i made] did she_i deny?



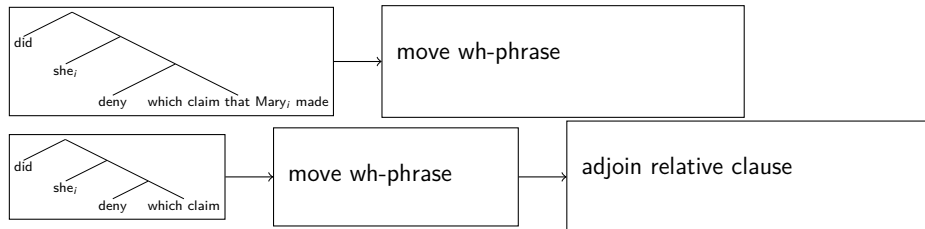
Full derivation recipes?

A few cases reveal that (we must all be already assuming that) it's full derivations/recipes that count.

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



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



Full derivation recipes?

Also:

- subjacency effects without traces
- compare categorial grammar

Full derivation recipes?

Also:

- subjacency effects without traces
- compare categorial grammar

And this is not a new idea!

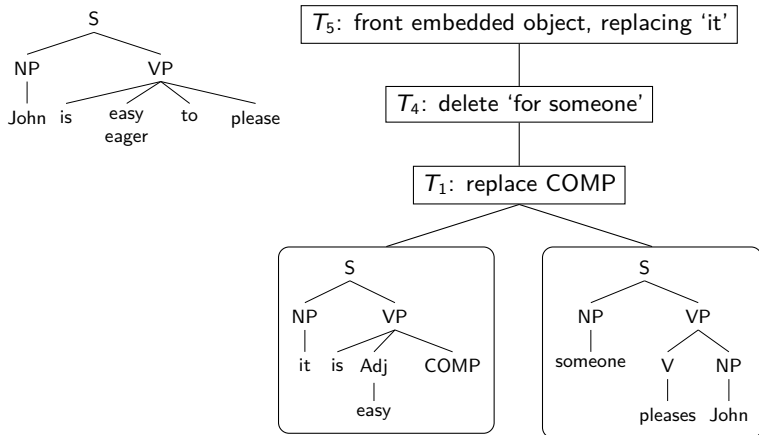
[The perceptual model] will utilize the full resources of the transformational grammar to provide a structural description, consisting of a set of P-markers and a transformational history

Miller and Chomsky (1963: p.480)

Full derivation recipes?

[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

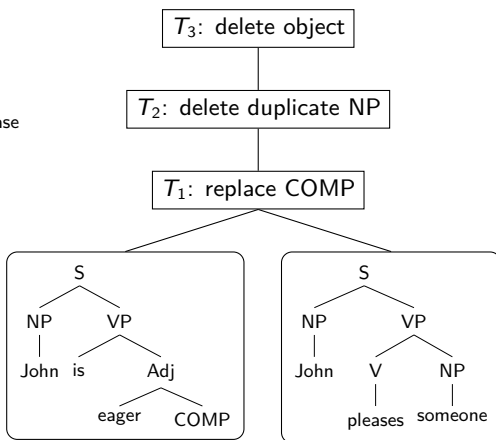
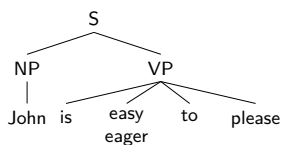
Miller and Chomsky (1963: p.480)



Full derivation recipes?

[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

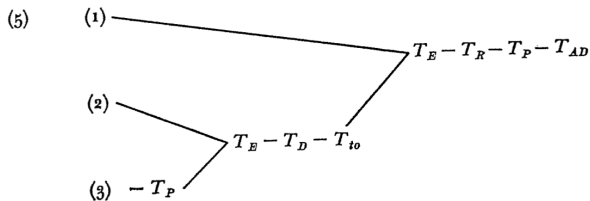
Miller and Chomsky (1963: p.480)



Full derivation recipes?

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

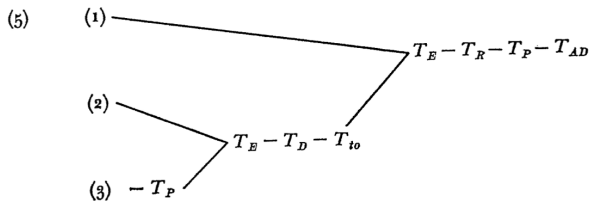
The “transformational history” of (4) by which it is derived from its basis might be represented, informally, by the diagram (5).



Full derivation recipes?

- (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).



Differences these days:

- We'll have things like **merge** and **move** at the internal nodes instead of T_P , T_E , etc.
- We'll have lexical items at the leaves rather than base-derived trees.

Outline

- 1 What we want to do with grammars
- 2 How to get grammars to do it
- 3 Derivations and representations
- 4 Information-theoretic complexity metrics**

Surprisal and entropy reduction

Why these complexity metrics?

- Partly just for concreteness, to give us a goal.
- They are **formalism neutral** to a degree that others aren't.
- They are **mechanism neutral** (Marr level one).
- The pieces of the puzzle that we need to get there (e.g. probabilities) seem likely to be usable in other ways.

Surprisal and entropy reduction

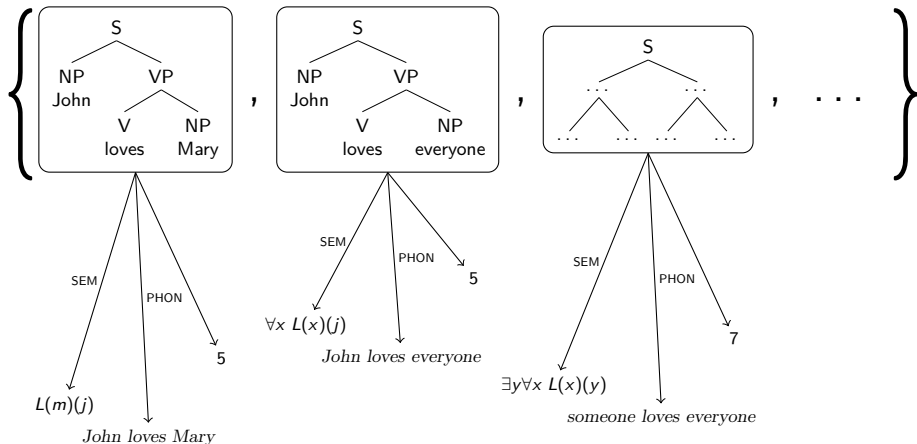
Why these complexity metrics?

- Partly just for concreteness, to give us a goal.
- They are **formalism neutral** to a degree that others aren't.
- They are **mechanism neutral** (Marr level one).
- The pieces of the puzzle that we need to get there (e.g. probabilities) seem likely to be usable in other ways.



John Hale, Cornell Univ.

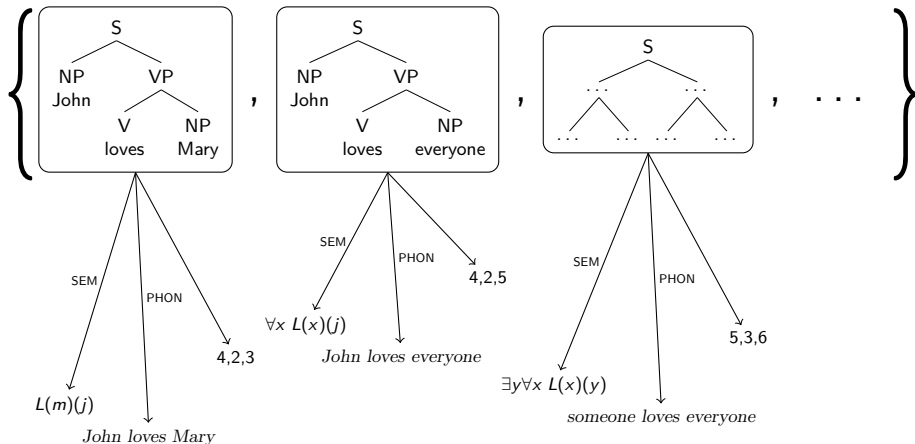
Interpretation functions



Caveats:

- Maybe we're interested in the finite specification of the set
- Maybe there's no clear line between observable and not
- Maybe some evidence is based on relativities among interpretations

Interpretation functions



Caveats:

- Maybe we're interested in the finite specification of the set
- Maybe there's no clear line between observable and not
- Maybe some evidence is based on relativities among interpretations

Surprisal

Given a sentence $w_1 w_2 \dots w_n$:

$$\text{surprisal at } w_i = -\log P(W_i = w_i \mid W_1 = w_1, W_2 = w_2, \dots, W_{i-1} = w_{i-1})$$

Surprisal

0.4	John ran
0.15	John saw it
0.05	John saw them
0.25	Mary ran
0.1	Mary saw it
0.05	Mary saw them

What predictions can we make about the difficulty of comprehending
'John saw it'?

Surprisal

0.4	John ran
0.15	John saw it
0.05	John saw them
0.25	Mary ran
0.1	Mary saw it
0.05	Mary saw them

What predictions can we make about the difficulty of comprehending
'John saw it'?

$$\begin{aligned}\text{surprisal at 'John'} &= -\log P(W_1 = \text{John}) \\ &= -\log(0.4 + 0.15 + 0.05) \\ &= -\log 0.6 \\ &= 0.74\end{aligned}$$

Surprisal

0.4	John ran
0.15	John saw it
0.05	John saw them
0.25	Mary ran
0.1	Mary saw it
0.05	Mary saw them

What predictions can we make about the difficulty of comprehending
'John saw it'?

$$\begin{aligned}\text{surprisal at 'John'} &= -\log P(W_1 = \text{John}) \\ &= -\log(0.4 + 0.15 + 0.05) \\ &= -\log 0.6 \\ &= 0.74\end{aligned}$$

$$\begin{aligned}\text{surprisal at 'saw'} &= -\log P(W_2 = \text{saw} \mid W_1 = \text{John}) \\ &= -\log \frac{0.15 + 0.05}{0.4 + 0.15 + 0.05} \\ &= -\log 0.33 \\ &= 1.58\end{aligned}$$

Surprisal

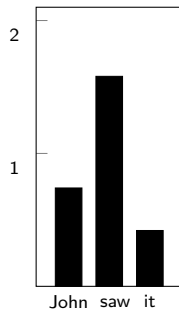
0.4	John ran
0.15	John saw it
0.05	John saw them
0.25	Mary ran
0.1	Mary saw it
0.05	Mary saw them

What predictions can we make about the difficulty of comprehending 'John saw it'?

$$\begin{aligned}
 \text{surprisal at 'John'} &= -\log P(W_1 = \text{John}) \\
 &= -\log(0.4 + 0.15 + 0.05) \\
 &= -\log 0.6 \\
 &= 0.74
 \end{aligned}$$

$$\begin{aligned}
 \text{surprisal at 'saw'} &= -\log P(W_2 = \text{saw} \mid W_1 = \text{John}) \\
 &= -\log \frac{0.15 + 0.05}{0.4 + 0.15 + 0.05} \\
 &= -\log 0.33 \\
 &= 1.58
 \end{aligned}$$

$$\begin{aligned}
 \text{surprisal at 'it'} &= -\log P(W_3 = \text{it} \mid W_1 = \text{John}, W_2 = \text{saw}) \\
 &= -\log \frac{0.15}{0.15 + 0.05} \\
 &= -\log 0.75 \\
 &= 0.42
 \end{aligned}$$



Accurate predictions made by surprisal

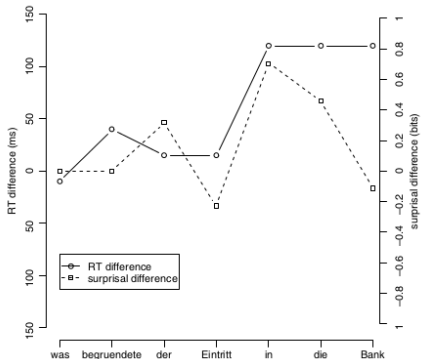


(Hale 2001)

Accurate predictions made by surprisal

- (8) The reporter [who ____ attacked the senator] left the room. (easier)
- (9) The reporter [who the senator attacked ____] left the room. (harder)

Difference between object-initial and subject-initial reading times and surprisals of (11)



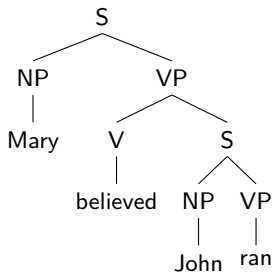
An important distinction

Using surprisal as a complexity metric says nothing about the form of the knowledge that the language comprehender is using!

- We're asking "what's the probability of w_i , given that we've seen $w_1 \dots w_{i-1}$ in the past".
- This does not mean that the comprehender's knowledge takes the form of answers to this kind of question.
- The linear nature of the metric reflects the **task**, not the **knowledge being probed**.

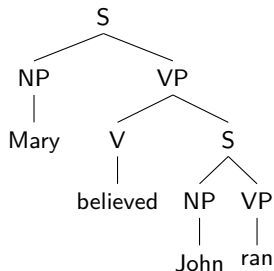
Probabilistic CFGs

1.0	$S \rightarrow NP VP$
0.3	$NP \rightarrow \text{John}$
0.7	$NP \rightarrow \text{Mary}$
0.2	$VP \rightarrow \text{ran}$
0.5	$VP \rightarrow V NP$
0.3	$VP \rightarrow V S$
0.4	$V \rightarrow \text{believed}$
0.6	$V \rightarrow \text{knew}$



Probabilistic CFGs

1.0	$S \rightarrow NP VP$
0.3	$NP \rightarrow \text{John}$
0.7	$NP \rightarrow \text{Mary}$
0.2	$VP \rightarrow \text{ran}$
0.5	$VP \rightarrow V NP$
0.3	$VP \rightarrow V S$
0.4	$V \rightarrow \text{believed}$
0.6	$V \rightarrow \text{knew}$



$$\begin{aligned}
 P(\text{Mary believed John ran}) &= 1.0 \times 0.7 \times 0.3 \times 0.4 \times 1.0 \times 0.3 \times 0.2 \\
 &= 0.00504
 \end{aligned}$$

Surprisal with probabilistic CFGs

Goal: Calculate step-by-step surprisal values for 'Mary believed John ran'

surprisal at 'John' = $-\log P(W_3 = \text{John} \mid W_1 = \text{Mary}, W_2 = \text{believed})$

Surprisal with probabilistic CFGs

Goal: Calculate step-by-step surprisal values for 'Mary believed John ran'

surprisal at 'John' = $-\log P(W_3 = \text{John} \mid W_1 = \text{Mary}, W_2 = \text{believed})$

0.098	Mary believed Mary
0.042	Mary believed John
0.012348	Mary believed Mary knew Mary
0.01176	Mary believed Mary ran
0.008232	Mary believed Mary believed Mary
0.005292	Mary believed Mary knew John
0.005292	Mary believed John knew Mary
0.00504	Mary believed John ran
...	...

Surprisal with probabilistic CFGs

Goal: Calculate step-by-step surprisal values for 'Mary believed John ran'

surprisal at 'John' = $-\log P(W_3 = \text{John} \mid W_1 = \text{Mary}, W_2 = \text{believed})$

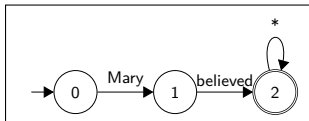
0.098	Mary believed Mary
0.042	Mary believed John
0.012348	Mary believed Mary knew Mary
0.01176	Mary believed Mary ran
0.008232	Mary believed Mary believed Mary
0.005292	Mary believed Mary knew John
0.005292	Mary believed John knew Mary
0.00504	Mary believed John ran
...	...

There are an **infinite number of derivations** consistent with input at each point!

$$\begin{aligned} \text{surprisal at 'John'} &= -\log P(W_3 = \text{John} \mid W_1 = \text{Mary}, W_2 = \text{believed}) \\ &= -\log \frac{0.042 + 0.005292 + 0.00504 + \dots}{0.098 + 0.042 + 0.12348 + 0.01176 + 0.008232 + \dots} \end{aligned}$$

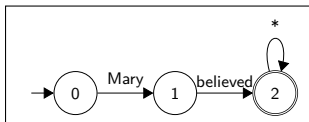
Intersection grammars

1.0	S → NP VP
0.3	NP → John
0.7	NP → Mary
0.2	VP → ran
0.5	VP → V NP
0.3	VP → V S
0.4	V → believed
0.6	V → knew

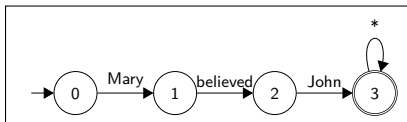
 \cap

 $=$
 G_2

Intersection grammars

1.0 S \rightarrow NP VP
 0.3 NP \rightarrow John
 0.7 NP \rightarrow Mary
 0.2 VP \rightarrow ran
 0.5 VP \rightarrow V NP
 0.3 VP \rightarrow V S
 0.4 V \rightarrow believed
 0.6 V \rightarrow knew

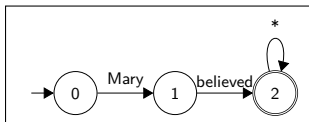
 \cap  $= G_2$

1.0 S \rightarrow NP VP
 0.3 NP \rightarrow John
 0.7 NP \rightarrow Mary
 0.2 VP \rightarrow ran
 0.5 VP \rightarrow V NP
 0.3 VP \rightarrow V S
 0.4 V \rightarrow believed
 0.6 V \rightarrow knew

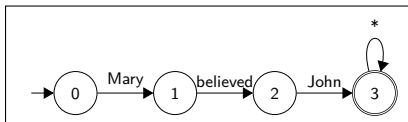
 \cap  $= G_3$

Intersection grammars

1.0 S \rightarrow NP VP
 0.3 NP \rightarrow John
 0.7 NP \rightarrow Mary
 0.2 VP \rightarrow ran
 0.5 VP \rightarrow V NP
 0.3 VP \rightarrow V S
 0.4 V \rightarrow believed
 0.6 V \rightarrow knew

 \cap  $= G_2$

1.0 S \rightarrow NP VP
 0.3 NP \rightarrow John
 0.7 NP \rightarrow Mary
 0.2 VP \rightarrow ran
 0.5 VP \rightarrow V NP
 0.3 VP \rightarrow V S
 0.4 V \rightarrow believed
 0.6 V \rightarrow knew

 \cap  $= G_3$

surprisal at 'John' = $-\log P(W_3 = \text{John} \mid W_1 = \text{Mary}, W_2 = \text{believed})$

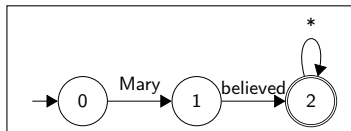
$$= -\log \frac{\text{total weight in } G_3}{\text{total weight in } G_2}$$

$$= -\log \frac{0.0672}{0.224}$$

$$= 1.74$$

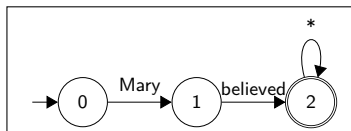
Grammar intersection example (simple)

1.0	$S \rightarrow NP VP$
0.3	$NP \rightarrow \text{John}$
0.7	$NP \rightarrow \text{Mary}$
0.2	$VP \rightarrow \text{ran}$
0.5	$VP \rightarrow V NP$
0.3	$VP \rightarrow V S$
0.4	$V \rightarrow \text{believed}$
0.6	$V \rightarrow \text{knew}$



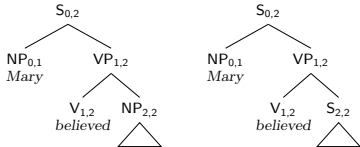
Grammar intersection example (simple)

1.0 $S \rightarrow NP VP$
 0.3 $NP \rightarrow John$
 0.7 $NP \rightarrow Mary$
 0.2 $VP \rightarrow ran$
 0.5 $VP \rightarrow V NP$
 0.3 $VP \rightarrow V S$
 0.4 $V \rightarrow believed$
 0.6 $V \rightarrow knew$



1.0 $S_{0,2} \rightarrow NP_{0,1} VP_{1,2}$
 0.7 $NP_{0,1} \rightarrow Mary$
 0.5 $VP_{1,2} \rightarrow V_{1,2} NP_{2,2}$
 0.3 $VP_{1,2} \rightarrow V_{1,2} S_{2,2}$
 0.4 $V_{1,2} \rightarrow believed$

1.0 $S_{2,2} \rightarrow NP_{2,2} VP_{2,2}$
 0.3 $NP_{2,2} \rightarrow John$
 0.7 $NP_{2,2} \rightarrow Mary$
 0.2 $VP_{2,2} \rightarrow ran$
 0.5 $VP_{2,2} \rightarrow V_{2,2} NP_{2,2}$
 0.3 $VP_{2,2} \rightarrow V_{2,2} S_{2,2}$
 0.4 $V_{2,2} \rightarrow believed$
 0.6 $V_{2,2} \rightarrow knew$

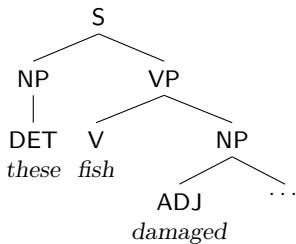
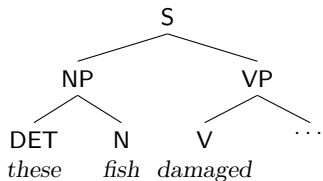


NB: Total weight in this grammar is not one! (What is it? Start symbol is $S_{0,2}$.)
 Each derivation has the weight "it" had in the original grammar.

Grammar intersection example (more complicated)

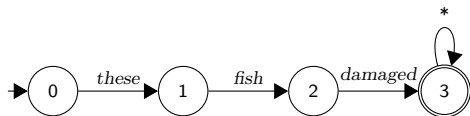
$S \rightarrow NP VP$ $V \rightarrow \textit{fish}$
 $VP \rightarrow V NP$ $V \rightarrow \textit{damaged}$
 $NP \rightarrow DET$ $DET \rightarrow \textit{these}$
 $NP \rightarrow DET N$ $N \rightarrow \textit{fish}$
 $NP \rightarrow ADJ N$ $ADJ \rightarrow \textit{damaged}$

These fish damaged ...



Grammar intersection example (more complicated)

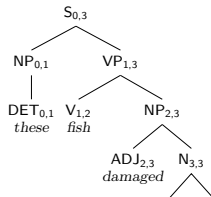
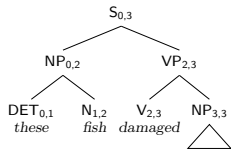
S	→	NP VP	V	→	<i>fish</i>
VP	→	V NP	V	→	<i>damaged</i>
NP	→	DET	DET	→	<i>these</i>
NP	→	DET N	N	→	<i>fish</i>
NP	→	ADJ N	ADJ	→	<i>damaged</i>



$S_{0,3}$	→	$NP_{0,2}$	$VP_{2,3}$
$NP_{0,2}$	→	$DET_{0,1}$	$N_{1,2}$
$VP_{2,3}$	→	$V_{2,3}$	$NP_{3,3}$
$DET_{0,1}$	→	<i>these</i>	
$N_{1,2}$	→	<i>fish</i>	
$V_{2,3}$	→	<i>damaged</i>	

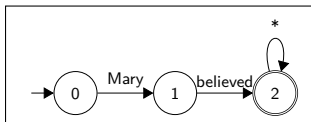
$S_{0,3}$	→	$NP_{0,1}$	$VP_{1,3}$
$NP_{0,1}$	→	$DET_{0,1}$	
$VP_{1,3}$	→	$V_{1,2}$	$NP_{2,3}$
$NP_{2,3}$	→	$ADJ_{2,3}$	$N_{3,3}$
$V_{1,2}$	→	<i>fish</i>	
$ADJ_{2,3}$	→	<i>damaged</i>	

$NP_{3,3}$	→	$ADJ_{3,3}$	$N_{3,3}$
$NP_{3,3}$	→	$DET_{3,3}$	$N_{3,3}$
$NP_{3,3}$	→	$DET_{3,3}$	
$N_{3,3}$	→	<i>fish</i>	
$DET_{3,3}$	→	<i>these</i>	
$ADJ_{3,3}$	→	<i>damaged</i>	

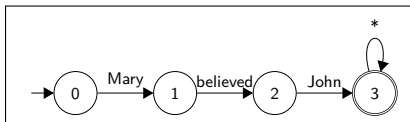


Intersection grammars

1.0 S \rightarrow NP VP
 0.3 NP \rightarrow John
 0.7 NP \rightarrow Mary
 0.2 VP \rightarrow ran
 0.5 VP \rightarrow V NP
 0.3 VP \rightarrow V S
 0.4 V \rightarrow believed
 0.6 V \rightarrow knew

 \cap

 $= G_2$

1.0 S \rightarrow NP VP
 0.3 NP \rightarrow John
 0.7 NP \rightarrow Mary
 0.2 VP \rightarrow ran
 0.5 VP \rightarrow V NP
 0.3 VP \rightarrow V S
 0.4 V \rightarrow believed
 0.6 V \rightarrow knew

 \cap

 $= G_3$

surprisal at 'John' = $-\log P(W_3 = \text{John} \mid W_1 = \text{Mary}, W_2 = \text{believed})$

$$= -\log \frac{\text{total weight in } G_3}{\text{total weight in } G_2}$$

$$= -\log \frac{0.0672}{0.224}$$

$$= 1.74$$

Computing sum of weights in a grammar (“partition function”)

$$Z(A) = \sum_{A \rightarrow \alpha} (p(A \rightarrow \alpha) \cdot Z(\alpha))$$

$$Z(\epsilon) = 1$$

$$Z(a\beta) = Z(\beta)$$

$$Z(B\beta) = Z(B) \cdot Z(\beta) \quad \text{where } \beta \neq \epsilon$$

(Nederhof and Satta 2008)

$$1.0 \quad S \rightarrow NP VP$$

$$0.3 \quad NP \rightarrow \text{John}$$

$$0.7 \quad NP \rightarrow \text{Mary}$$

$$0.2 \quad VP \rightarrow \text{ran}$$

$$0.5 \quad VP \rightarrow V NP$$

$$0.4 \quad V \rightarrow \text{believed}$$

$$0.6 \quad V \rightarrow \text{knew}$$

$$Z(V) = 0.4 + 0.6 = 1.0$$

$$Z(NP) = 0.3 + 0.7 = 1.0$$

$$\begin{aligned} Z(VP) &= 0.2 + (0.5 \cdot Z(V) \cdot Z(NP)) \\ &= 0.2 + (0.5 \cdot 1.0 \cdot 1.0) = 0.7 \end{aligned}$$

$$\begin{aligned} Z(S) &= 1.0 \cdot Z(NP) \cdot Z(VP) \\ &= 0.7 \end{aligned}$$

Computing sum of weights in a grammar (“partition function”)

$$Z(A) = \sum_{A \rightarrow \alpha} (p(A \rightarrow \alpha) \cdot Z(\alpha))$$

$$Z(\epsilon) = 1$$

$$Z(a\beta) = Z(\beta)$$

$$Z(B\beta) = Z(B) \cdot Z(\beta) \quad \text{where } \beta \neq \epsilon$$

(Nederhof and Satta 2008)

1.0	$S \rightarrow NP VP$	$Z(V) = 0.4 + 0.6 = 1.0$
0.3	$NP \rightarrow \text{John}$	$Z(NP) = 0.3 + 0.7 = 1.0$
0.7	$NP \rightarrow \text{Mary}$	
0.2	$VP \rightarrow \text{ran}$	$Z(VP) = 0.2 + (0.5 \cdot Z(V) \cdot Z(NP))$
0.5	$VP \rightarrow V NP$	$= 0.2 + (0.5 \cdot 1.0 \cdot 1.0) = 0.7$
0.4	$V \rightarrow \text{believed}$	$Z(S) = 1.0 \cdot Z(NP) \cdot Z(VP)$
0.6	$V \rightarrow \text{knew}$	$= 0.7$

1.0	$S \rightarrow NP VP$	
0.3	$NP \rightarrow \text{John}$	$Z(V) = 0.4 + 0.6 = 1.0$
0.7	$NP \rightarrow \text{Mary}$	$Z(NP) = 0.3 + 0.7 = 1.0$
0.2	$VP \rightarrow \text{ran}$	
0.5	$VP \rightarrow V NP$	$Z(VP) = 0.2 + (0.5 \cdot Z(V) \cdot Z(NP)) + (0.3 \cdot Z(V) \cdot Z(S))$
0.3	$VP \rightarrow V S$	$Z(S) = 1.0 \cdot Z(NP) \cdot Z(VP)$
0.4	$V \rightarrow \text{believed}$	
0.6	$V \rightarrow \text{knew}$	

Things to know

Technical facts about CFGs:

- Can intersect with a “prefix FSA”
- Can compute the total weight (and the entropy)

Things to know

Technical facts about CFGs:

- Can intersect with a “prefix FSA”
- Can compute the total weight (and the entropy)

More generally:

- Intersecting a grammar with a prefix produces a new grammar which is a representation of the comprehender’s sentence-medial state
- So we can construct a [sequence of grammars](#) which represents the comprehender’s sequence of knowledge-states
- Ask “what changes” (or “how much changes”, etc.) at each step

The general approach is compatible with many very different grammar formalisms (any grammar formalism?) — provided the technical tricks can be pulled off.

Looking ahead

Wouldn't it be nice if we could do all that for minimalist syntax?

The average syntax paper shows **illustrative derivations**, not a **fragment**.

What would we need?

- An explicit characterization of the set of possible derivations
- A way to “intersect” that with a prefix
- A way to define probability distributions over the possibilities

This will require certain idealizations. (But what's new?)

- Billot, S. and Lang, B. (1989). The structure of shared forests in ambiguous parsing. In *Proceedings of the 1989 Meeting of the Association of Computational Linguistics*.
- Chomsky, N. (1965). *Aspects of the Theory of Syntax*. MIT Press, Cambridge, MA.
- Chomsky, N. (1980). *Rules and Representations*. Columbia University Press, New York.
- Ferreira, F. (2005). Psycholinguistics, formal grammars, and cognitive science. *The Linguistic Review*, 22:365–380.
- Gärtner, H.-M. and Michaelis, J. (2010). On the Treatment of Multiple-Wh Interrogatives in Minimalist Grammars. In Hanneforth, T. and Fanselow, G., editors, *Language and Logos*, pages 339–366. Akademie Verlag, Berlin.
- Gibson, E. and Wexler, K. (1994). Triggers. *Linguistic Inquiry*, 25:407–454.
- Hale, J. (2006). Uncertainty about the rest of the sentence. *Cognitive Science*, 30:643–672.
- Hale, J. T. (2001). A probabilistic early parser as a psycholinguistic model. In *Proceedings of the Second Meeting of the North American Chapter of the 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. and Dyer, C. (2013). Distributions on minimalist grammar derivations. In *Proceedings of the 13th Meeting on the Mathematics of Language*.
- Koopman, H. and Szabolcsi, A. (2000). *Verbal Complexes*. MIT Press, Cambridge, MA.

- Lang, B. (1988). Parsing incomplete sentences. In *Proceedings of the 12th International Conference on Computational Linguistics*, pages 365–371.
- Levy, R. (2008). Expectation-based syntactic comprehension. *Cognition*, 106(3):1126–1177.
- Michaelis, J. (2001). Derivational minimalism is mildly context-sensitive. In Moortgat, M., editor, *Logical Aspects of Computational Linguistics*, volume 2014 of *LNCS*, pages 179–198. Springer, Berlin Heidelberg.
- 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.
- Morrill, G. (1994). *Type Logical Grammar: Categorical Logic of Signs*. Kluwer, Dordrecht.
- Nederhof, M. J. and Satta, G. (2008). Computing partition functions of pcfgs. *Research on Language and Computation*, 6(2):139–162.
- Seki, H., Matsumara, T., Fujii, M., and Kasami, T. (1991). On multiple context-free grammars. *Theoretical Computer Science*, 88:191–229.
- 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.
- Stabler, E. P. (2011). Computational perspectives on minimalism. In Boeckx, C., editor, *The Oxford Handbook of Linguistic Minimalism*. Oxford University Press, Oxford.
- Stabler, E. P. and Keenan, E. L. (2003). Structural similarity within and among languages. *Theoretical Computer Science*, 293:345–363.

- Vijay-Shanker, K., Weir, D. J., and Joshi, A. K. (1987). Characterizing structural descriptions produced by various grammatical formalisms. In *Proceedings of the 25th Meeting of the Association for Computational Linguistics*, pages 104–111.
- Weir, D. (1988). *Characterizing mildly context-sensitive grammar formalisms*. PhD thesis, University of Pennsylvania.
- Yngve, V. H. (1960). A model and an hypothesis for language structure. In *Proceedings of the American Philosophical Society*, volume 104, pages 444–466.