

# Psych 156A/ Ling 150: Acquisition of Language II

## Lecture 17 Learning Language Structure

### Announcements

Work on structure review questions

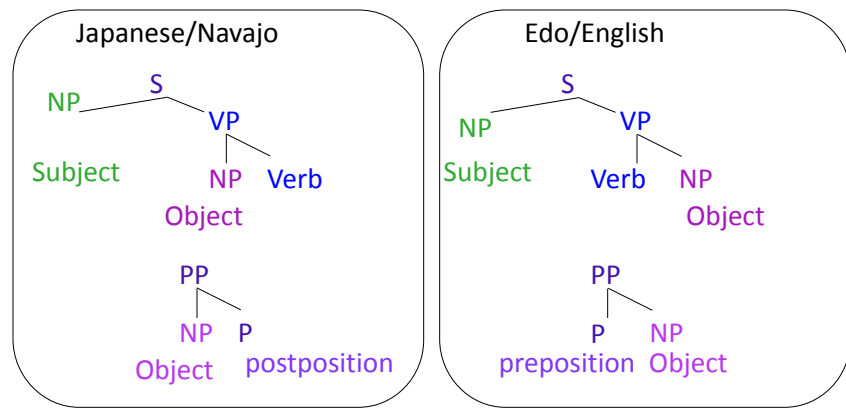
Final review this Thursday 6/5/14

Final exam next Thursday 6/12/14 between 1:30 and 3:30pm (taken online through EEE).

Consider taking more language science classes in the future!

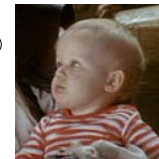
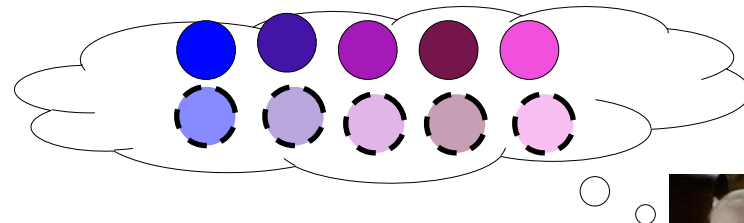
### Language variation: Recap from before

While languages may differ on many levels, they have many similarities at the level of language structure (syntax). Even languages with no shared history seem to share similar structural patterns.



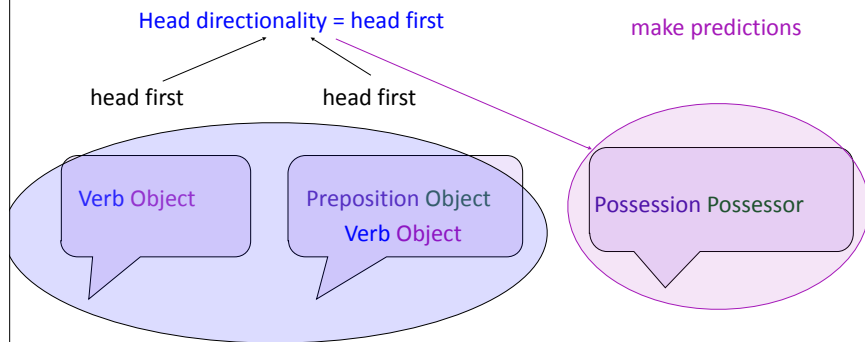
### Language variation: Recap from before

One way for children to learn the complex structures of their language is to have them already be aware of the ways in which human languages can vary. Linguistic nativists believe this is knowledge contained in Universal Grammar. Then, children listen to their native language data to decide which patterns their native language follows.



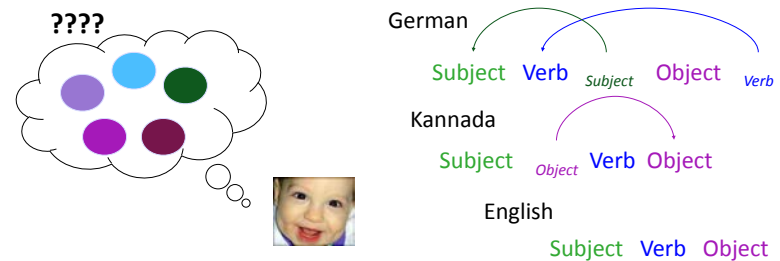
## Language variation: Recap from before

Languages can be thought to vary structurally on a number of linguistic parameters. One purpose of parameters is to explain how children learn some hard-to-notice structural properties.



## Issue from last time: Learning parameter values

The observable data are often the result of a **combination of interacting parameters**. That is, the observable data are the result of some unobservable process, and the child has to reverse engineer the observable data to figure out what parameter values might have produced the observable data - even if the child already knows what the parameters are!



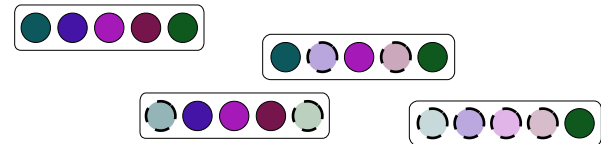
## Learning structure with statistical learning: Linguistic parameters and probability



## Linguistic knowledge for learning structure

Parameters = constraints on language variation. Only certain rules/patterns are possible. This is linguistic knowledge.

A language's grammar  
= combination of language rules  
= combination of parameter values



Idea: use statistical learning to learn which value (for each parameter) that the native language uses for its grammar. This is a combination of using linguistic knowledge & statistical learning.

## Yang 2004: Variational learning

Idea taken from evolutionary biology:

In a population, individuals compete against each other. The fittest individuals survive while the others die out.

How do we translate this to learning language structure?

## Yang 2004: Variational learning

Idea taken from evolutionary biology:

In a population, individuals compete against each other. The fittest individuals survive while the others die out.

How do we translate this to learning language structure?

Individual = grammar (combination of parameter values that represents the structural properties of a language)



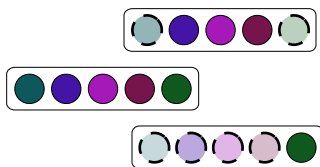
Fitness = how well a grammar can analyze the data the child encounters

## Yang 2004: Variational learning

Idea taken from evolutionary biology:

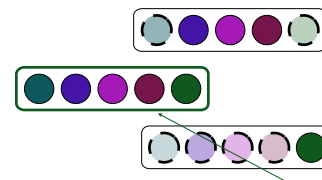
A child's mind consists of a population of grammars that are competing to analyze the data in the child's native language.

Population of grammars



## Yang 2004: Variational learning

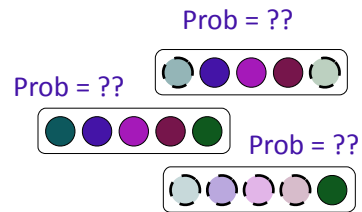
Intuition: The most successful (fittest) grammar will be the native language grammar because it can analyze all the data the child encounters. This grammar will "win", once the child encounters enough native language data because none of the other competing grammars can analyze all the data.



If this is the native language grammar, this grammar can analyze all the input while the other two can't.

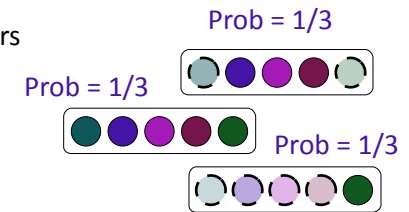
## Variational learning details

At any point in time, a grammar in the population will have a probability associated with it. This represents the child's belief that this grammar is the correct grammar for the native language.



## Variational learning details

Before the child has encountered any native language data, all grammars are equally likely. So, initially all grammars have the same probability, which is 1 divided the number of grammars available.

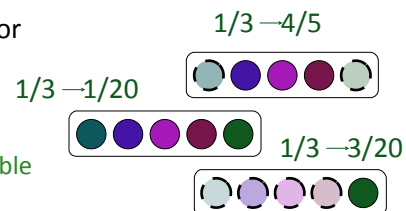


If there are 3 grammars, the initial probability for any given grammar =  $1/3$

## Variational learning details

As the child encounters data from the native language, some of the grammars will be more fit because they are better able to account for the structural properties in the data.

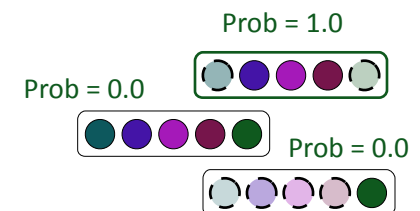
Other grammars will be less fit because they cannot account for some of the data encountered.



Grammars that are more compatible with the native language data will have their probabilities increased while grammars that are less compatible will have their probabilities decreased over time.

## Variational learning details

After the child has encountered enough data from the native language, the native language grammar should have a probability near 1.0 while the other grammars have a probability near 0.0.



## The power of unambiguous data

Unambiguous data from the native language can only be analyzed by grammars that use the native language's parameter value.

This makes **unambiguous data very influential data** for the child to encounter, since these data are **only compatible with the parameter value that is correct for the native language.**

## Unambiguous data

**Problem: Do unambiguous data exist for entire grammars?**

This requires data that are incompatible with every other possible parameter value of every other possible grammar....

This seems unlikely for real language data because parameters connect with different types of patterns, which may have nothing to do with each other.

## Unambiguous issues

Parameter 1: subject-drop



Spanish: +subject-drop



Patterns allowed:

Vamos

*go-1<sup>st</sup>-pl-pres*

"We go"

Subject dropped

Nosotros vamos

*1<sup>st</sup>-pl go-1<sup>st</sup>-pl-pres*

"We go"

Subject spoken

## Unambiguous issues

Parameter 1: subject-drop



English: -subject-drop



Patterns allowed:

~~*go-1<sup>st</sup>-pl-pres*~~

~~"go" ≠ "we go"~~

~~Subject dropped~~

*1<sup>st</sup>-pl*

"We"

*go-pres*

"go"

Subject spoken

## Unambiguous issues

Parameter 2: Head-directionality



Edo/English: Head first



Basic word order:

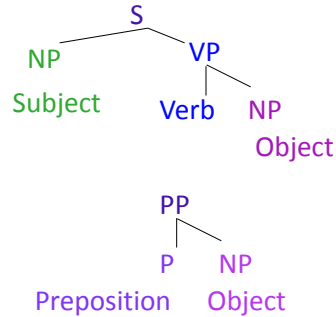
Subject Verb Object [SVO]

Prepositions:

Preposition Noun Phrase

Possessed before Possessor

Possession Possessor



## Unambiguous issues

Parameter 2: Head-directionality



Japanese/Navajo: Head-final



Basic word order:

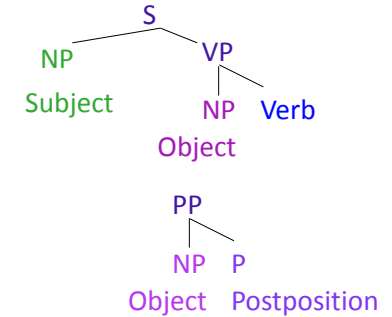
Subject Object Verb [SOV]

Postpositions:

Noun Phrase Postposition

Possessor before Possessed

Possessor Possession

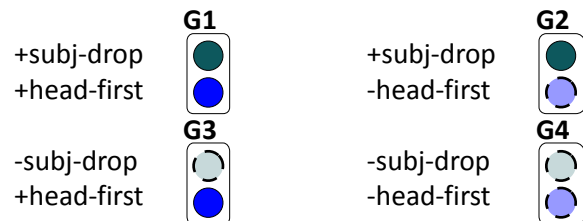


## Unambiguous issues

Data point:

Subject Object Verb

Grammars available:

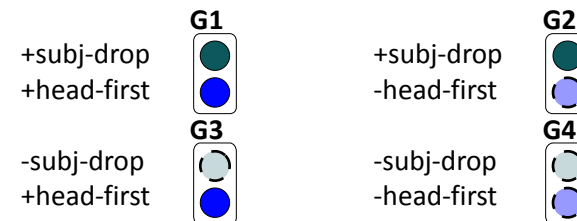


## Unambiguous issues

Data point:

Subject Object Verb

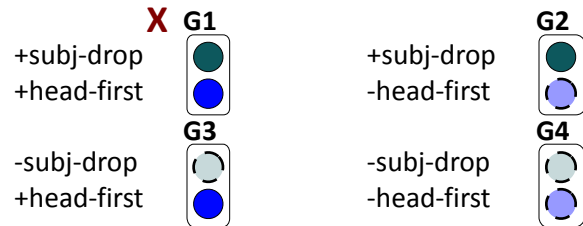
Which grammars can analyze this data point?



## Unambiguous issues

Data point:            Subject   Object   Verb

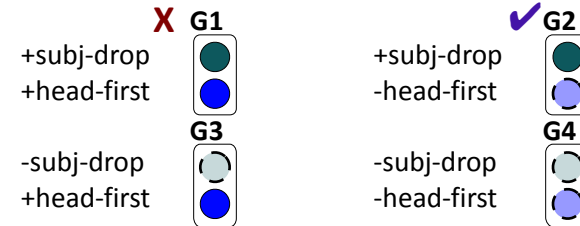
**G1?**    ✓ +subj-drop allows Subject to be spoken  
          X +head-first predicts SVO



## Unambiguous issues

Data point:            Subject   Object   Verb

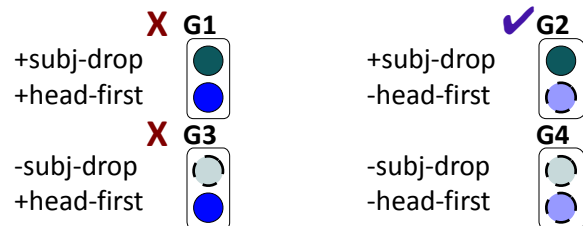
**G2?**    ✓ +subj-drop allows Subject to be spoken  
          ✓ -head-first predicts SOV



## Unambiguous issues

Data point:            Subject   Object   Verb

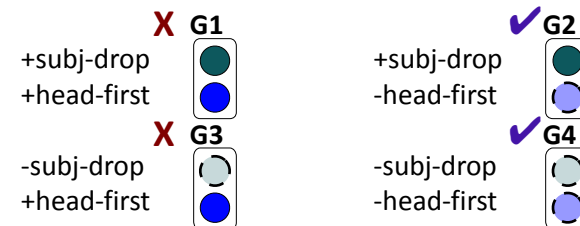
**G3?**    ✓ -subj-drop requires Subject to be spoken  
          X +head-first predicts SVO



## Unambiguous issues

Data point:            Subject   Object   Verb

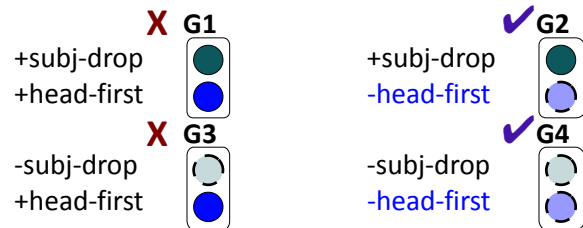
**G4?**    ✓ -subj-drop requires Subject to be spoken  
          ✓ -head-first predicts SOV



## Unambiguous issues

Data point:            Subject   Object   Verb

There's more than one grammar compatible with this data point...even though we feel like it should be informative for head directionality.

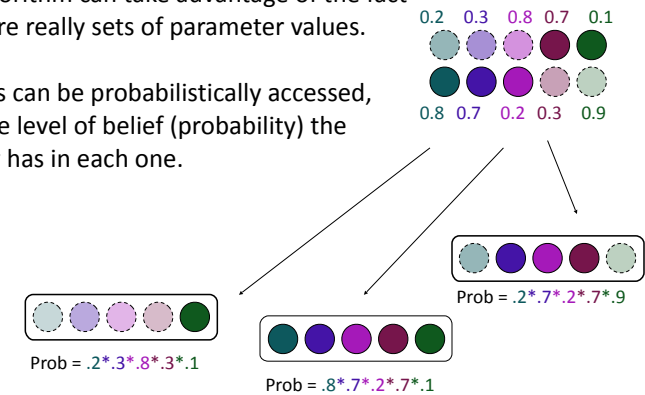


## Using parameters

### Parameterized grammars

Yang (2004)'s algorithm can take advantage of the fact that grammars are really sets of parameter values.

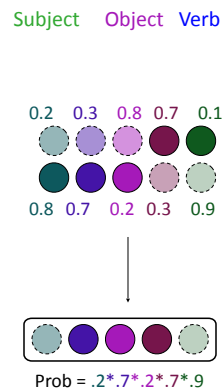
Parameter values can be probabilistically accessed, depending on the level of belief (probability) the learner currently has in each one.



## The learning algorithm

For each data point encountered in the input...

- (1) Choose a grammar to test out on a particular data point. Select a grammar by choosing a set of parameter values, based on the probabilities associated with each parameter value.

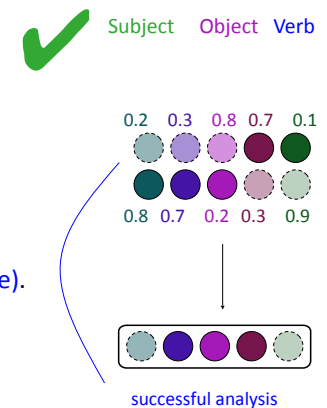


Denison, Bonawitz, Gopnik, & Griffiths 2013:  
Experimental evidence from 4 and 5-year-olds suggests that children are sensitive to the probabilities of complex representations (such as parameters), and so this kind of sampling is not unrealistic.

## The learning algorithm

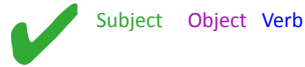
- (2) Try to analyze the data point with this grammar.

If this grammar **can** analyze the data point, increase the probability of all participating parameters values slightly (**reward each value**).





## The learning algorithm

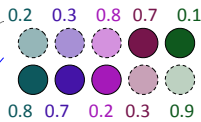


(2) Try to analyze the data point with this grammar.

If this grammar can analyze the data point, increase the probability of all participating parameters values slightly (reward each value).

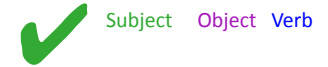
Actual update equation for reward:

$p_v$  = previous value of successful parameter value = .2  
 $p_o$  = previous value of opposing parameter value = .8



successful analysis

## The learning algorithm



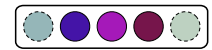
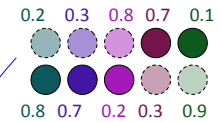
(2) Try to analyze the data point with this grammar.

If this grammar can analyze the data point, increase the probability of all participating parameters values slightly (reward each value).

Actual update equation for reward:

$p_{v\_updated} = p_v + \gamma(1 - p_v)$   
 $p_{o\_updated} = (1 - \gamma)p_o$

$\gamma$  = learning rate (ex:  $\gamma = .125$ )



successful analysis

## The learning algorithm



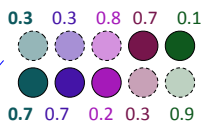
(2) Try to analyze the data point with this grammar.

If this grammar can analyze the data point, increase the probability of all participating parameters values slightly (reward each value).

Actual update equation for reward:

If  $p_v = .2$  and  $p_o = .8$ ...  
 $p_{v\_updated} = .2 + .125(1 - .2) = .3$   
 $p_{o\_updated} = (1 - .125).8 = .7$

$\gamma$  = learning rate (ex:  $\gamma = .125$ )



successful analysis

= .2  $\rightarrow$  .3  
 = .8  $\rightarrow$  .7

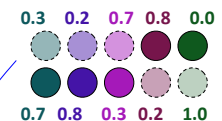
## The learning algorithm



(2) Try to analyze the data point with this grammar.

If this grammar can analyze the data point, increase the probability of all participating parameters values slightly (reward each value).

Do this for each parameter value in the chosen grammar.



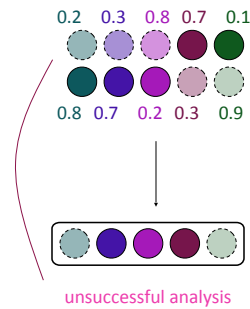
successful analysis

## The learning algorithm

X Subject Object Verb

(2) Try to analyze the data point with this grammar.

If this grammar **cannot** analyze the data point, decrease the probability of all participating parameters values slightly (**punish each value**).



## The learning algorithm

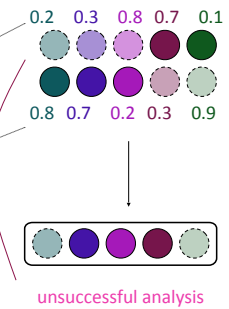
X Subject Object Verb

(2) Try to analyze the data point with this grammar.

If this grammar **cannot** analyze the data point, decrease the probability of all participating parameters values slightly (**punish each value**).

Actual update equation for reward:

$p_v$  = previous value of successful parameter value ● = .2  
 $p_o$  = previous value of opposing parameter value ● = .8



## The learning algorithm

X Subject Object Verb

(2) Try to analyze the data point with this grammar.

If this grammar **cannot** analyze the data point, decrease the probability of all participating parameters values slightly (**punish each value**).

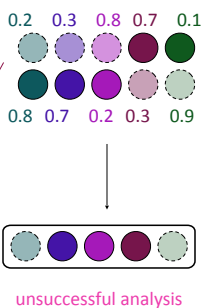
Actual update equation for reward:

$$p_{v\_updated} = (1-\gamma)p_v$$

$$p_{o\_updated} = \gamma + (1-\gamma)p_o$$

$\gamma$  = learning rate (ex:  $\gamma = .125$ )

● = .2  
● = .8



## The learning algorithm

X Subject Object Verb

(2) Try to analyze the data point with this grammar.

If this grammar **cannot** analyze the data point, decrease the probability of all participating parameters values slightly (**punish each value**).

Actual update equation for reward:

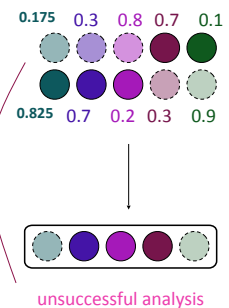
If  $p_v = .2$  and  $p_o = .8$ ...

$$p_{v\_updated} = (1-.125).2 = .175$$

$$p_{o\_updated} = .125 + (1-.125).8 = .825$$

$\gamma$  = learning rate (ex:  $\gamma = .125$ )

● = .2 → .175  
● = .8 → .825



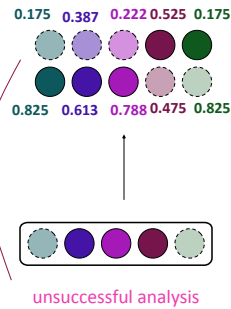
## The learning algorithm

**X** Subject Object Verb

(2) Try to analyze the data point with this grammar.

If this grammar **cannot** analyze the data point, decrease the probability of all participating parameters values slightly (**punish each value**).

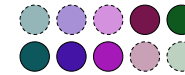
Do this for each parameter value in the chosen grammar.



## Unambiguous data

Problem ameliorated!

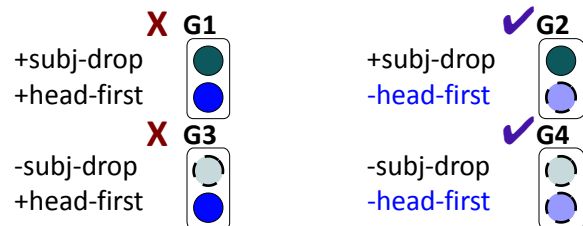
Unambiguous data are much more likely to exist for individual parameter values instead of entire grammars.



## Unambiguous issues – no more!

Data point: Subject Object Verb

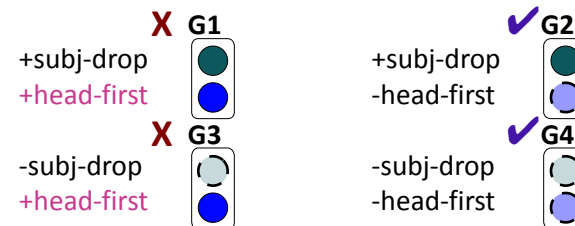
In this case, if either G2 or G4 were selected, -head-first would be rewarded (in addition to whichever subj-drop value was used).



## Unambiguous issues – no more!

Data point: Subject Object Verb

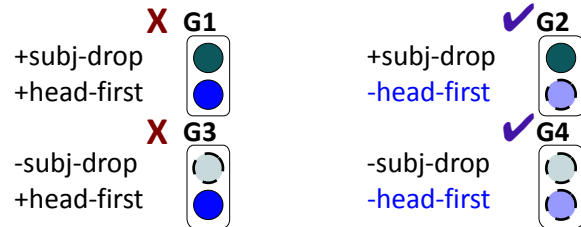
In this case, if either G1 or G3 were selected, +head-first would be punished (in addition to whichever subj-drop value was used).



## Unambiguous issues – no more!

Data point:            Subject    Object    Verb

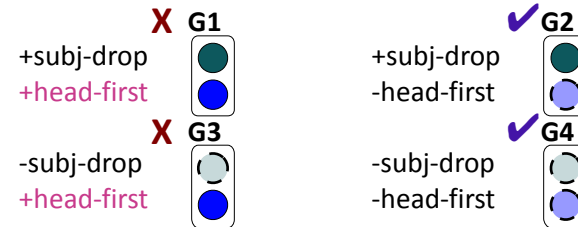
Because this data point is unambiguous for **-head-first**, grammars using that value would be rewarded and its probability as a parameter value would become 1.0 over time.



## Unambiguous issues – no more!

Data point:            Subject    Object    Verb

Meanwhile, grammars using **+head-first** would be punished every time, and its probability as a parameter value would approach 0.0 over time.



## Unambiguous data

Idea from Yang 2004: The more unambiguous data there are, the faster the native language's parameter value will "win" (reach a probability near 1.0). This means that the child will learn the associated structural pattern faster.

Example: the more unambiguous +subject-drop data the child encounters, the faster a child should learn that the native language allows subjects to be dropped.

## Unambiguous data

Idea from Yang 2004: The more unambiguous data there are, the faster the native language's parameter value will "win" (reach a probability near 1.0). This means that the child will learn the associated structural pattern faster.

Question: Is it true that the amount of unambiguous data the child encounters for a particular parameter determines when the child learns that structural property of the language?

Yang 2004, 2011:  
Unambiguous data learning examples

Wh-fronting for questions

Wh-word moves to the front (like English)

Sarah will see **who**?

*Underlying form of the question*

Yang 2004, 2011:  
Unambiguous data learning examples

Wh-fronting for questions

Wh-word moves to the front (like English)

**Who** will Sarah will see who?

*Observable (spoken) form of the question*

Yang 2004, 2011:  
Unambiguous data learning examples

Wh-fronting for questions

Wh-word moves to the front (like English)

**Who** will Sarah will see who?

Wh-word stays "in place" (like Chinese)

Sarah will see **who**?

*Observable (spoken) form of the question*

Yang 2004, 2011:  
Unambiguous data learning examples

Wh-fronting for questions

Parameter: +/- wh-fronting

Native language value (English): +wh-fronting

Unambiguous data: any (normal) wh-question, with **wh-word in front** (ex: "Who will Sarah see?")

Frequency of unambiguous data to children: 25% of input

Age of +wh-fronting acquisition: very early (before 1 yr, 8 months)

## Yang 2004, 2011: Unambiguous data learning examples

### Topic drop

Chinese (+topic-drop): can drop NP (subject or object) if it is the understood topic of the discourse

Understood topic: Jareth

*Speakers had been talking about Jareth*

## Yang 2004, 2011: Unambiguous data learning examples

### Topic drop

Chinese (+topic-drop): can drop NP (subject or object) if it is the understood topic of the discourse

Understood topic: Jareth

Mingtian guiji hui xiayu.

Tomorrow estimate will rain

'It is tomorrow that (Jareth) believes it will rain'

*Speaker doesn't have to say "Jareth"*

## Yang 2004, 2011: Unambiguous data learning examples

### Topic drop

Chinese (+topic-drop): can drop NP (subject or object) if it is the understood topic of the discourse

Understood topic: Jareth

Mingtian guiji hui xiayu.

Tomorrow estimate will rain

'It is tomorrow that (Jareth) believes it will rain'

English (-topic-drop): can't drop topic NP

*Speaker has to say "Jareth"*

\*It is tomorrow that believes it will rain.

It is tomorrow that Jareth believes it will rain.

## Yang 2004, 2011: Unambiguous data learning examples

### Topic drop

Parameter: +/- topic-drop

Native language value (Chinese): +topic-drop

Unambiguous data: any utterance where the object NP is dropped because it is the topic

Frequency of unambiguous data to children: 12% of input

Age of +topic-drop acquisition: very early (before 1 yr, 8 months)

## Yang 2004, 2011: Unambiguous data learning examples

### Subject drop

Italian (+subject-drop): can drop the subject

Verrá?  
*3<sup>rd</sup>-sg-will-come*  
“Will s/he come?”

English (-subject-drop): can't drop subject NP

\*Will come?  
Will he come?

## Yang 2004, 2011: Unambiguous data learning examples

### Subject drop

Parameter: +/- subject-drop

Native language value (Italian): +subject-drop

Unambiguous data: Dropped subjects in questions

Frequency of unambiguous data to children: 10% of input

Age of +subject-drop acquisition: very early (before 1 yr, 8 months)

## Yang 2004, 2011: Unambiguous data learning examples

### Subject drop

Parameter: +/- subject-drop

Native language value (English): -subject-drop

Unambiguous data: Expletive subjects (ex: *It* seems he's going to come after all.)

Frequency of unambiguous data to children: 1.2% of input

Age of -subject-drop acquisition: 3 years old

## Yang 2004, 2011: Unambiguous data learning examples

### Verb raising

Verb moves “above” (before) the *adverb/negative word* (French)

Jean *souvent* voit Marie  
Jean *often* sees Marie

Jean *pas* voit Marie  
Jean *not* sees Marie

*Underlying form of the sentence*

## Yang 2004, 2011: Unambiguous data learning examples

### Verb raising

Verb moves “above” (before) the **adverb/negative word** (French)

Jean **voit** **souvent** voit Marie  
 Jean **often** **sees** Marie

Jean **voit** **pas** voit Marie  
 Jean **not** **sees** Marie

*Observable (spoken) form of the sentence*

## Yang 2004, 2011: Unambiguous data learning examples

### Verb raising

Verb moves “above” (before) the **adverb/negative word** (French)

Jean **voit** **souvent** voit Marie  
 Jean **often** **sees** Marie

Jean **voit** **pas** voit Marie  
 Jean **not** **sees** Marie

Verb stays “below” (after) the **adverb/negative word** (English)

Jean **often** **sees** Marie.  
 Jean **does not see** Marie.

*Observable (spoken) form of the sentence*

## Yang 2004, 2011: Unambiguous data learning examples

### Verb raising

Parameter: **+/- verb-raising**

Native language value (French): **+verb-raising**

Unambiguous data: data points that have both a **verb** and an **adverb/negative word** in them, where the positions of each can be seen (“Jean **voit souvent** Marie”)

Frequency of unambiguous data to children: **7% of input**

Age of +verb-raising acquisition: **1 yr, 8 months**

## Yang 2004, 2011: Unambiguous data learning examples

### Verb Second

Verb moves to second phrasal position, **some other phrase** moves to the first position (German)

Sarah **liest** Sarah **das Buch** liest  
 Sarah **reads** **the book** “Sarah reads the book.”

Das Buch **liest** Sarah **das Buch** liest  
 The book **reads** Sarah “Sarah reads the book.”

Verb does not move (English)

Sarah **reads** **the book**.

*Observable (spoken) form of the sentence*



## Yang 2004, 2011: Unambiguous data learning examples

### Verb Second

Parameter: +/- verb-second

Native language value (German): +verb-second

Unambiguous data: Object Verb Subject data points in German (“Das Buch liest Sarah”), since they show the Object and the Verb in front of the Subject

Frequency of unambiguous data to children: 1.2% of input

Age of +verb-second acquisition: ~3 yrs

## Yang 2004, 2011: Unambiguous data learning examples

### Intermediate wh-words in complex questions

(Hindi, some German) *Observable (spoken) form of the question*

Wer glaubst du wer Recht hat?

Who think-2nd-sg you who right has

“Who do you think has the right?”

## Yang 2004, 2011: Unambiguous data learning examples

### Intermediate wh-words in complex questions

(Hindi, some German)

Wer glaubst du wer Recht hat?

Who think-2nd-sg you who right has

“Who do you think has the right?”

No intermediate wh-words in complex questions (English)

Who do you think has the right?

*Observable (spoken) form of the question*

## Yang 2004, 2011: Unambiguous data learning examples

### Intermediate wh-words in complex questions

Parameter: +/- intermediate-wh

Native language value (English): -intermediate-wh

Unambiguous data: complex questions of a particular kind that show the absence of a wh-word at the beginning of the embedded clause (“Who do you think has the right?”)

Frequency of unambiguous data to children: 0.2% of input

Age of -intermediate-wh acquisition: > 4 yrs

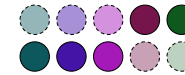
## Yang 2004, 2011: Unambiguous data learning examples

Parameter value	Frequency of unambiguous data	Age of acquisition
+wh-fronting (English)	25%	Before 1 yr, 8 months
+topic-drop (Chinese)	12%	Before 1 yr, 8 months
+subject-drop (Italian)	10%	Before 1 yr, 8 months
+verb-raising (French)	7%	1 yr, 8 months
+verb-second (German)	1.2%	3 yrs
-subject-drop (English)	1.2%	3 yrs
-intermediate-wh (English)	0.2%	> 4 yrs

The quantity of unambiguous data available in the child's input seems to be a good indicator of when they will acquire the knowledge. The more there is, the sooner they learn the right parameter value for their native language.

## Summary: Variational learning for language structure

Big idea: When a parameter is set depends on **how frequent the unambiguous data are** in the data the child encounters. This can be captured easily with the variational learning idea, since unambiguous data are very influential: They always reward the native language grammar and always punish grammars with the non-native parameter value.



## Summary: Variational learning for language structure

Predictions of variational learning:

Parameters set early: more unambiguous data available

Parameters set late: less unambiguous data available

These predictions seem to be born out by available data on when children learn certain structural patterns (parameter values) about their native language.

## Questions?



You should be able to do all the questions on the structure review questions. Remember to bring questions to the final exam review next class!