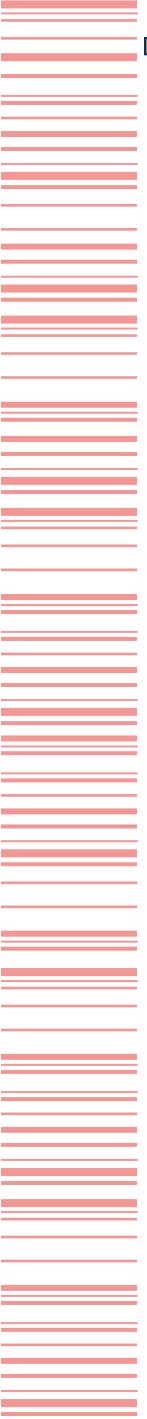


# A Step towards Natural-Sounding Speech Generation

Maarika Traat  
Tartu University  
Institute of Computer Science



# Contents

- ◆ Natural Speech vs. Computer Speech
  - Intonation
- ◆ What is Information Structure?
  - Relation to Intonation
- ◆ Discourse Representation Theory (DRT)
- ◆ Information Structure in DRT
- ◆ Categorial Grammars
  - Combinatory Categorial Grammar
  - Unification Categorial Grammar
- ◆ Unification-based Combinatory Categorial Grammar (UCCG)
- ◆ Information Structure in UCCG

# Natural Speech vs. Computer Speech

- ◆ Text-to-speech vs. content-to-speech
- ◆ Formant synthesis (rule-based synthesis)
  - an acoustic model.
  - parameters: fundamental frequency, voicing, and noise levels varied over time to create a waveform of artificial speech (Bell Labs 63) 
- ◆ Articulatory synthesis
  - computational models of the human vocal tract; the articulation processes occurring there
  - mostly of academic interest
- ◆ Concatenative Synthesis
  - concatenation (or stringing together) of segments of recorded speech
  - generally the most natural sounding synthesized speech
  - but sometimes audible glitches in the output, detracting from the naturalness of the synthesized speech



# Natural Speech vs. Computer Speech

- ◆ Text-to-speech vs. content-to-speech
- ◆ Formant synthesis (rule-based synthesis)
  - an acoustic model.
  - parameters: fundamental frequency, voicing, and noise levels varied over time to create a waveform of artificial speech (Bell Labs 63)
- ◆ Articulatory synthesis
  - computational models of the human vocal tract; the articulation processes occurring there
  - mostly of academic interest
- ◆ Concatenative Synthesis
  - concatenation (or stringing together) of segments of recorded speech
  - generally the most natural sounding synthesized speech
  - but sometimes audible glitches in the output, detracting from the naturalness of the synthesized speech

# Natural Speech vs. Computer Speech

- ◆ Examples of concatenative speech

- AT&T:

US 

UK 

IN 

I like all kinds of ice cream.

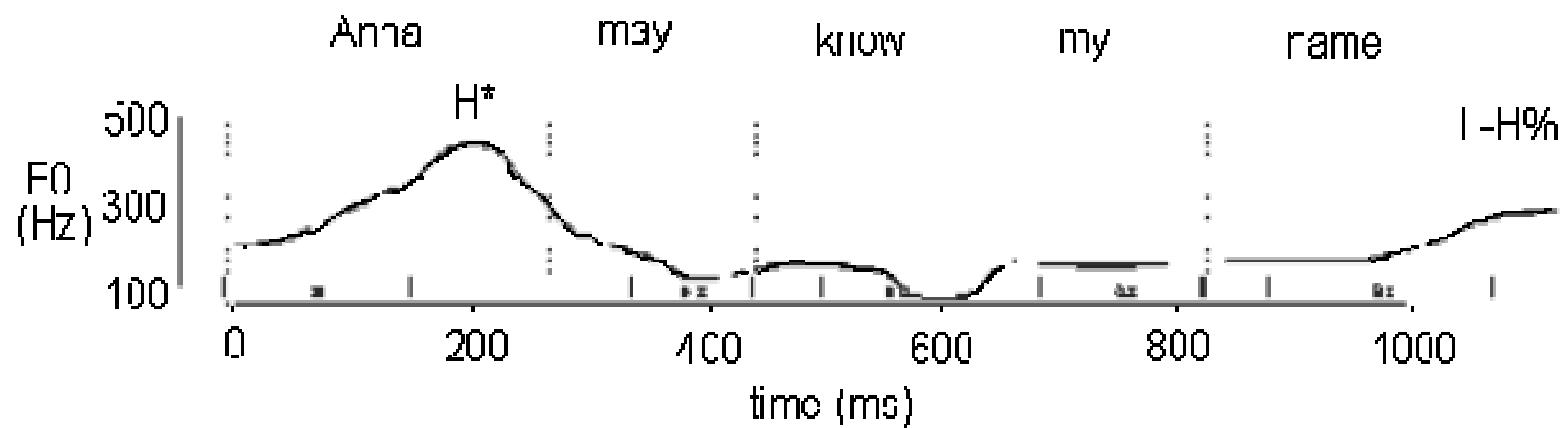
Usually chocolate flavour is my favourite, but sometimes simple vanilla ice cream is even better.

Which one do you like better: chocolate or vanilla ice cream?

- Institute of Cybernetics at TTU: EST 

# Intonation

- ◆ Intonation, prosody, speech melody
- ◆ Intonation contour: F0 contour
  - The fundamental tone: the lowest frequency in a harmonic series. In addition to F0 there are overtones (faster waves).

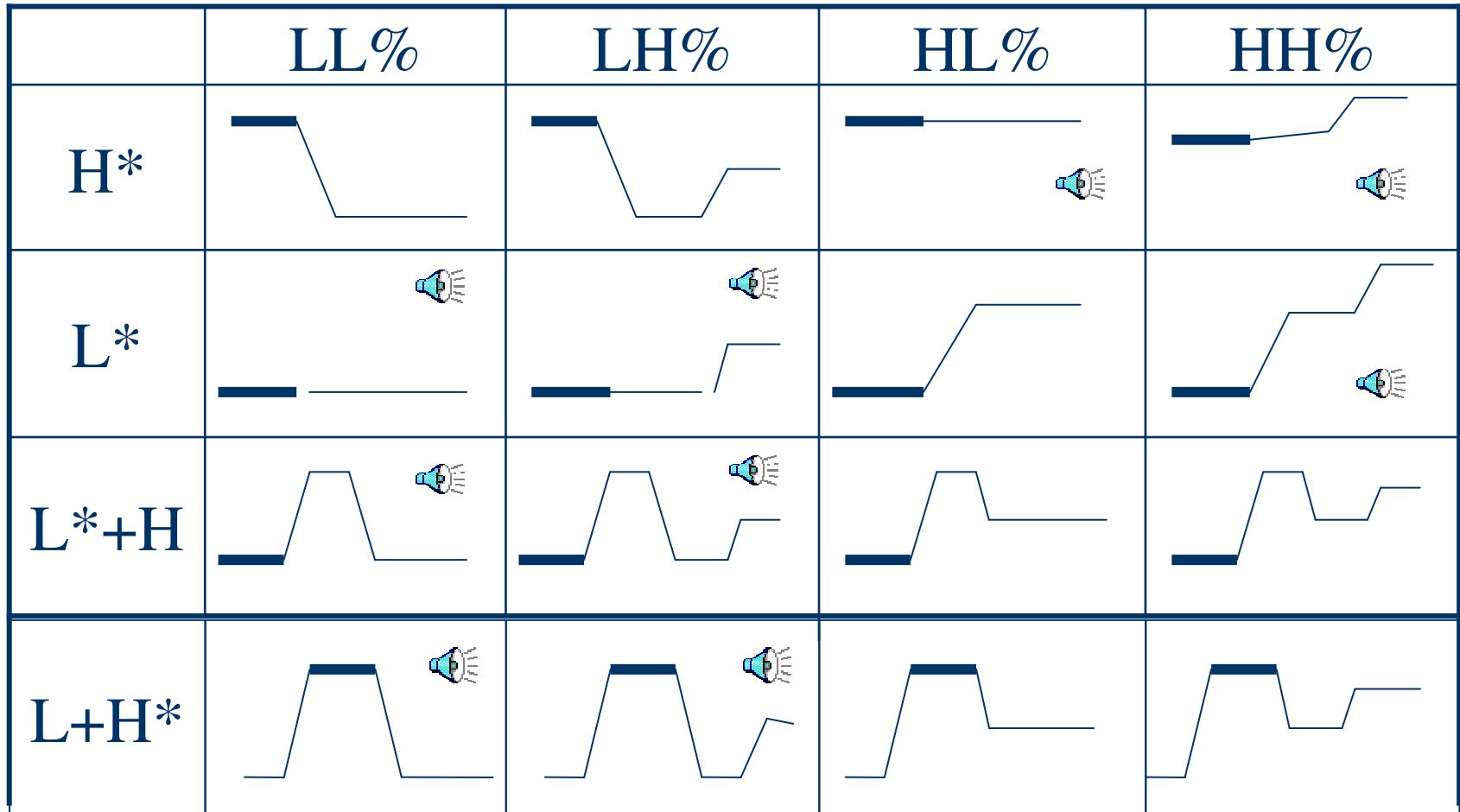


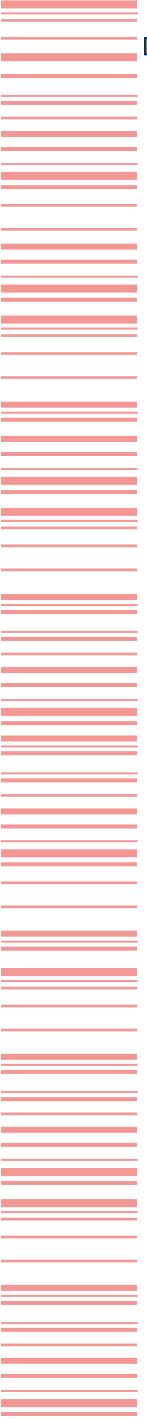


# Pierrehumbert's System of Prosodic Elements (mod.)

- ◆ Pitch accents
  - On words: makes words stand out, emphasises them
    - H\*, L\*, L+H\*, L\*+H, H+L\*, H\*+L, H\*+H
- ◆ Boundary tones
  - At the end of prosodic phrases
    - LL%, LH%, HH%, HL%

# Prosodic Elements Contd.





# Information Structure

- Aristotle taught young ALEXANDER.
- \* ARISTOTLE taught young Alexander.

# Information Structure

- Who did Aristotle teach?
- Aristotle taught young ALEXANDER
- \* ARISTOTLE taught young Alexander.

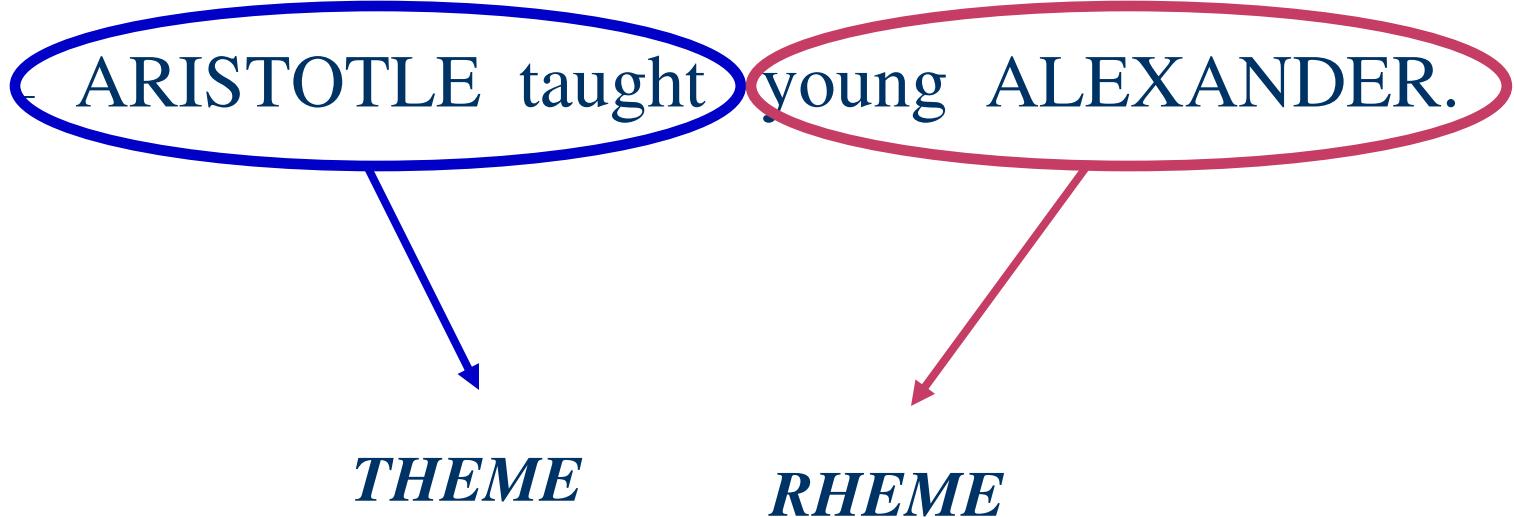
**THEME**

**RHEME**

**FOCUS**

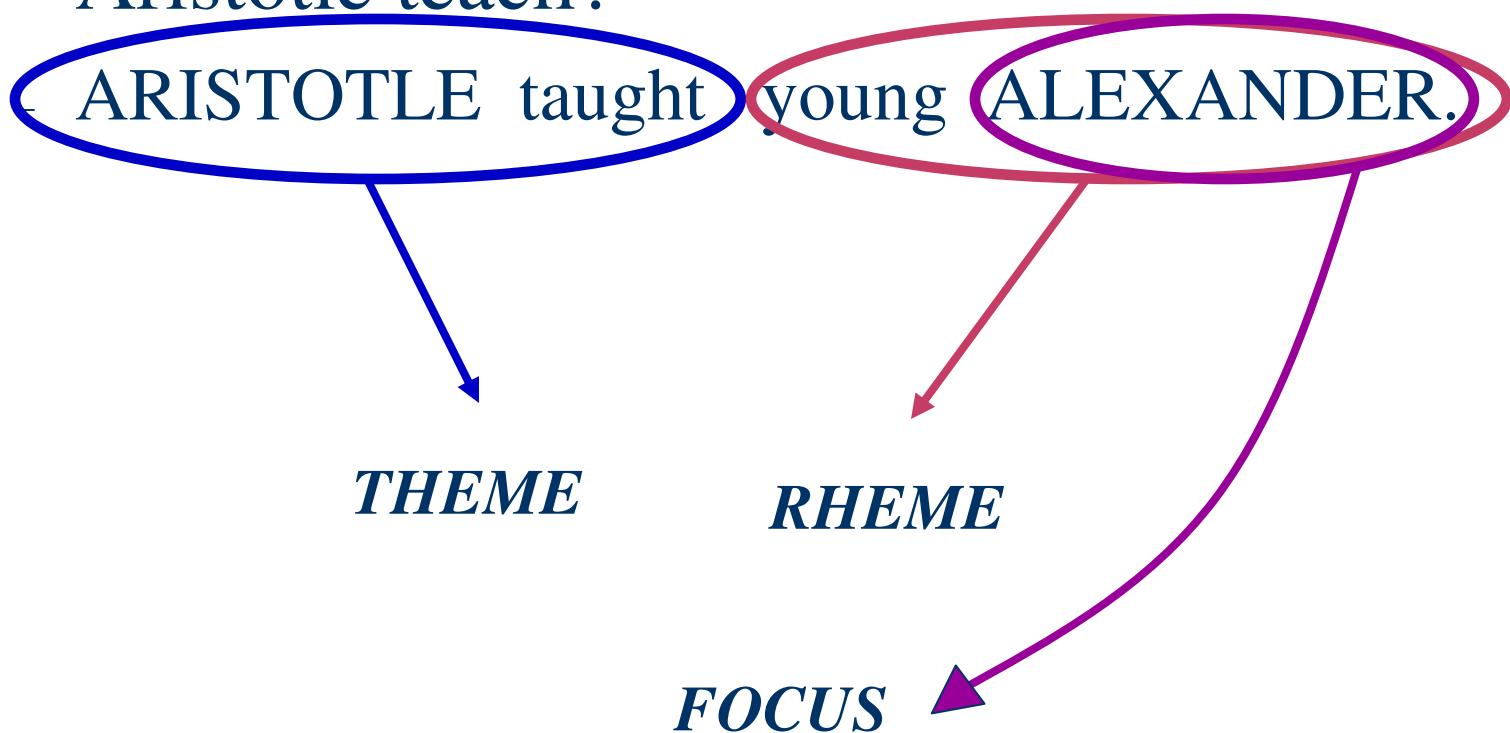
# Information Structure

- I know who you teach, but who did Aristotle teach?



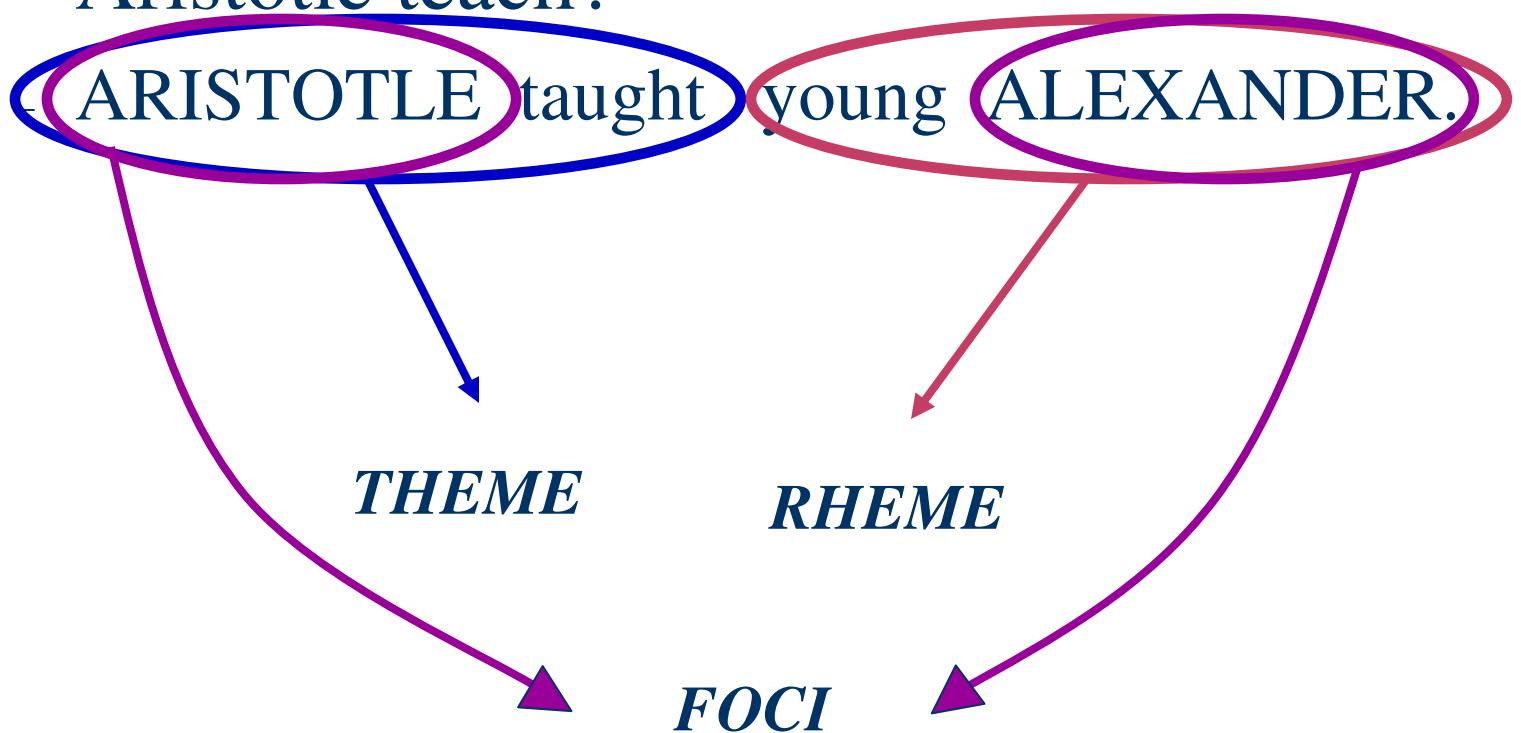
# Information Structure

- I know who you teach, but who did Aristotle teach?



# Information Structure

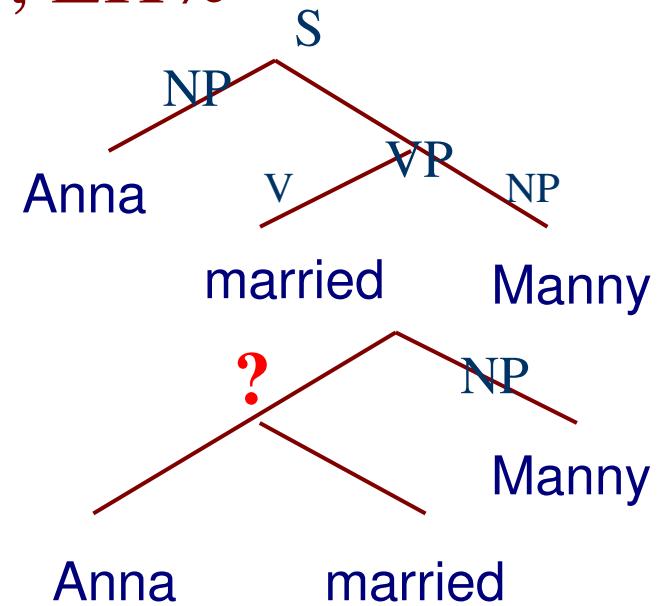
- I know who you teach, but who did Aristotle teach?



# Steedman's Prosodic Approach to Information Structure

- ◆ Theme pitch accents, e.g. L+H\*
- ◆ Rheme pitch accents, e.g. H\*
- ◆ Boundary tones, e.g. LL%, LH%

- *Anna married Manny.*  
H\* LL%                    L+H\* LH%
- *Anna married Manny.*  
L+H\*            LH%      H\* LL%



# Discourse Representation Theory (DRT)

- ◆ Discourse Representation Structure (DRS):
  - Universe: discourse referents
  - Discourse conditions

Anna likes ice cream. She eats it often.

X Y E
anna(X)
female(X)
like(E)
ice_cream(Y)
thing(Y)
agent(E,X)
patient(E,Y)



Z W E <sub>1</sub>
Z=?
female(?)
eat(E <sub>1</sub> )
W=? <sub>1</sub>
thing(? <sub>1</sub> )
agent(E <sub>1</sub> ,Z)
patient(E <sub>1</sub> ,W)
often(E <sub>1</sub> )



X Y Z W E E <sub>1</sub>
anna(X)
female(X)
like(E)
ice_cream(Y)
thing(Y)
agent(E,X)
patient(E,Y)
Z=X
W=Y
eat(E <sub>1</sub> )
agent(E <sub>1</sub> ,Z)
patient(E <sub>1</sub> ,W)
often(E <sub>1</sub> )



# Goals



- ◆ Combine
    - information structure
    - DRT semantics
  - ◆ Provide a mechanism to get from
    - an intonationally annotated text to our semantic representation marked with information structure, and vice versa
- 
- IS-DRS



# Our approach: IS-DRS

- ◆ Information structure flags on discourse conditions
- ◆ Suitable both for natural language generation and automatic inference

Anna H\* LL% married  
Manny L+H\* LH%.

theme  $\theta$   
rheme  $\rho$   
focus  $+/-$

X,Y,E	
anna(X)	$\rho+$
manny(Y)	$\theta+$
marry(E)	$\theta-$
agent(E,X)	
patient(E,Y)	



# Categorial Grammars

- ◆ Grammar: lexicon and rules for constructing more complex structures from simpler ones
- ◆ Categorial Grammars:
  - Lexicalised theories of grammar
  - Category – functional type associated with an entry in the lexicon

man :- n

a :- np/n

- ◆ Combinatory rules

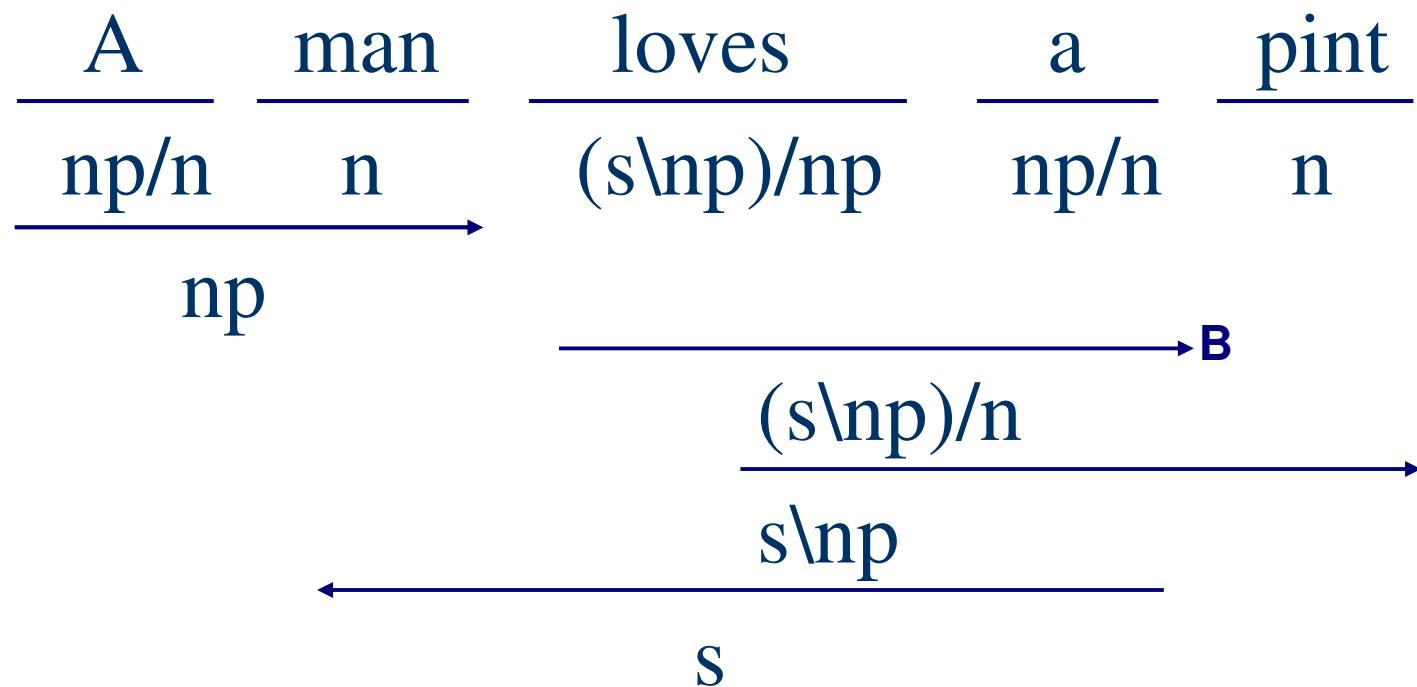
$$\frac{a \quad \text{man}}{\text{np/n} \quad \text{n}} \Rightarrow \frac{}{\text{np}}$$

# Combinatory Categorial Grammar (CCG)

- ◆ Categories
  - Basic categories: s, n, np
  - Complex categories:  $s\backslash np$ ,  $np/n$ ,  $(n\backslash n)/(s\backslash np)$ 
    - Directional slashes: \ and /
- ◆ Combinatory rules:
  - Forward application:  $X/Y \ Y \rightarrow X$
  - Backward application:  $Y \ X\backslash Y \rightarrow X$
  - Forward composition:  $X/Y \ Y/Z \rightarrow X/Z$
  - Backward composition:  $Y\backslash Z \ X\backslash Y \rightarrow X\backslash Z$
  - Etc.



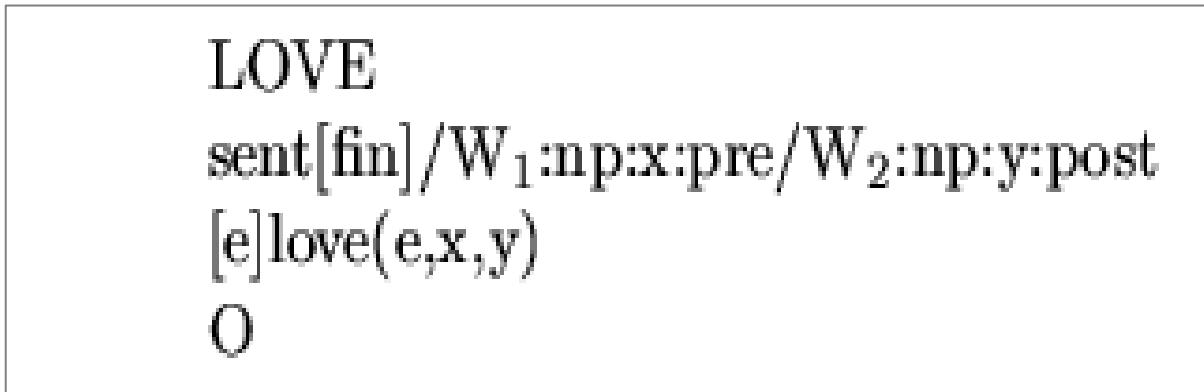
# Combinatory Categorial Grammar





# Unification Categorial Grammar

- ◆ Unification
- ◆ Feature structures called signs
- ◆ Forward and backward application



LOVE  
sent[fin]/W<sub>1</sub>:np:x:pre/W<sub>2</sub>:np:y:post  
[e]love(e,x,y)  
O

# Unification-based Combinatory Categorial Grammar

- ◆ Combinatory rules
- ◆ Directional slashes
- ◆ Unification
- ◆ Signs
- ◆ Compositional semantics
- ◆ Basic categories: s, n, vp
- ◆ DRT semantics (1<sup>st</sup> order)
- ◆ Information structure in semantics

}

CCG

}

UCG

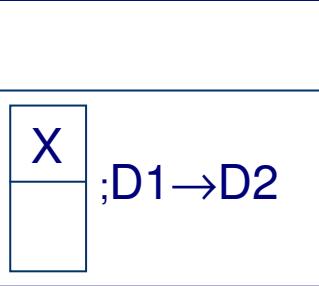


# Signs

- ◆ Basic signs

**PHO:** child  
**CAT:** n  
**VAR:** X  
**DRS:**   
child(X)

- ◆ Complex sign

**PHO:** every+W1+W2  
**CAT:** s  
**DRS:**   
X  
;D1→D2

/

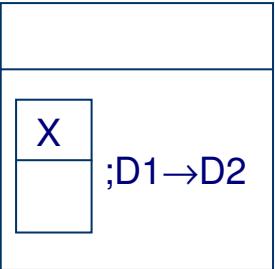
**PHO:** W2  
**CAT:** vp  
**AGR:** fin  
**VAR:** X  
**DRS:** D2

) /

**PHO:** W1  
**CAT:** n  
**VAR:** X  
**DRS:** D1

# Combining of Signs

PHO: every+W1+W2  
CAT: s  
DRS:



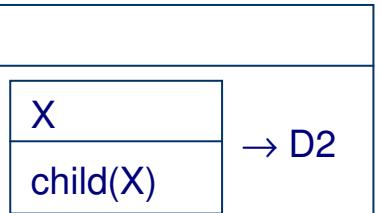
PHO: W2  
CAT: vp  
AGR: fin  
VAR: X  
DRS: D2

) /  
PHO: W1  
CAT: n  
VAR: X  
DRS: D1

PHO: child  
CAT: n  
VAR: Y  
DRS:



PHO: every+child+W2  
CAT: s  
DRS:



PHO: W2  
CAT: vp  
AGR: fin  
VAR: X  
DRS: D2





# Theme and Rheme in UCCG

- ◆ Themeness or rhemeness of a lexical item depends on the type of pitch accent that occurs on it
- ◆ If no pitch accent on the item, its information structure value is a variable → gets its value from a marked item with which the current item is combined
  - theme marked item can only combine with
    - another theme marked item
    - an unmarked item
  - rheme marked item can only combine with
    - another rheme marked item
    - an unmarked item

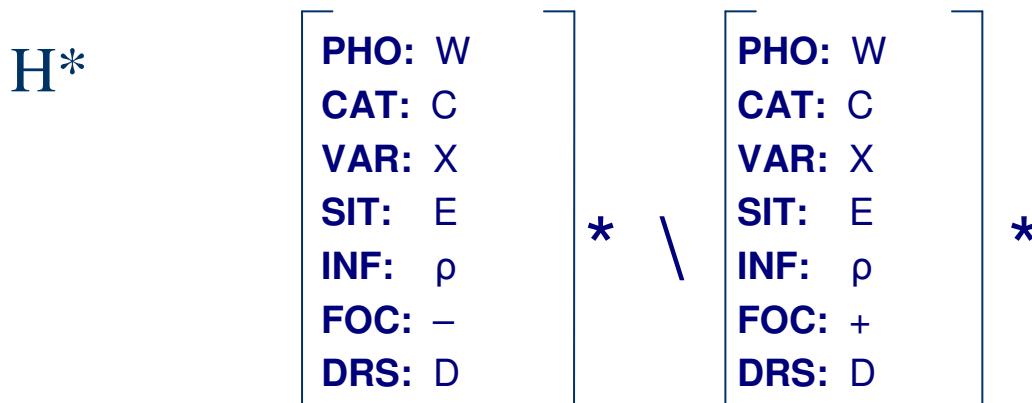
# Information Structure in UCCG Feature Structures

- ◆ INF feature
  - theme, rheme, phrase
- ◆ FOC feature
  - + the word carries a pitch accent
  - the word does not have a pitch accent on it

<b>PHO:</b> manny + W	/	<b>PHO:</b> W	
<b>CAT:</b> s		<b>CAT:</b> vp	
<b>INF:</b> I		<b>VAR:</b> X	
<b>FOC:</b> F		<b>SIT:</b> E	
<b>DRS:</b> <table border="1" style="display: inline-table;"><tr><td>X</td></tr><tr><td>Manny(X) I F</td></tr></table> $\otimes D$		X	Manny(X) I F
X			
Manny(X) I F			
<b>FOC:</b> F		<b>FOC:</b> F	
<b>DRS:</b> D		<b>DRS:</b> D	

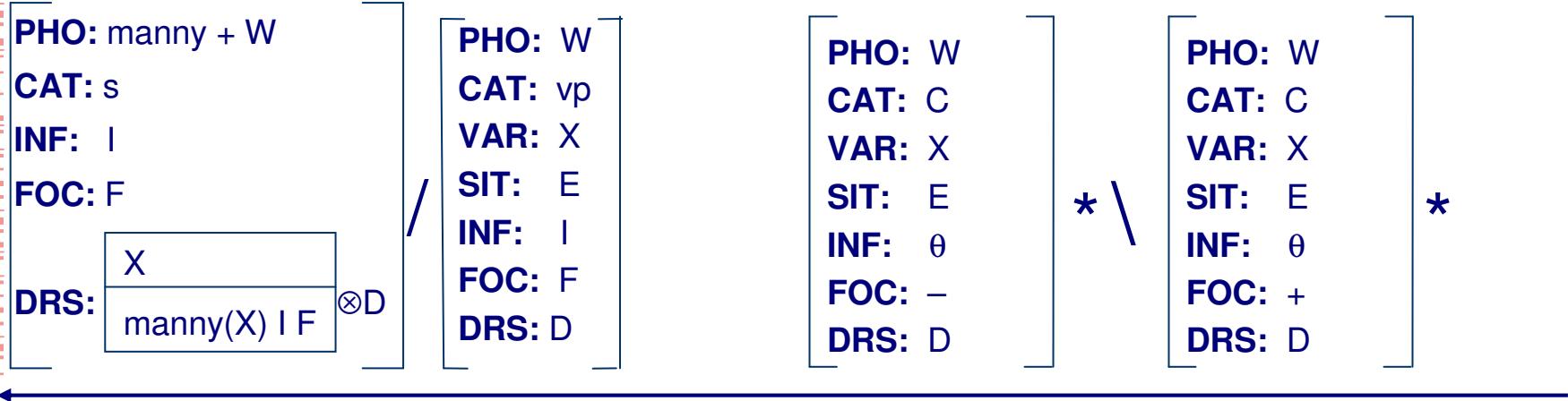
# Pitch Accent

- ◆ Appears as an individual lexical item
- ◆ Copies the sign on its left, and instantiates the information structure feature value to theme ( $\theta$ ) or rheme ( $\rho$ )
- ◆ Via unification marks the semantics of the lexical item as being focused (+) (different values in active & result parts)



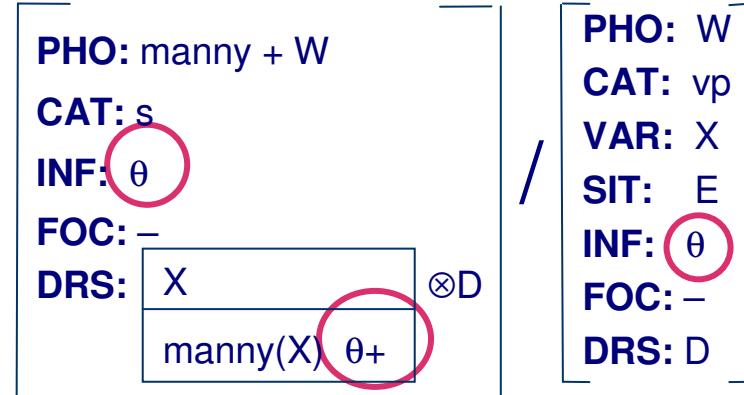
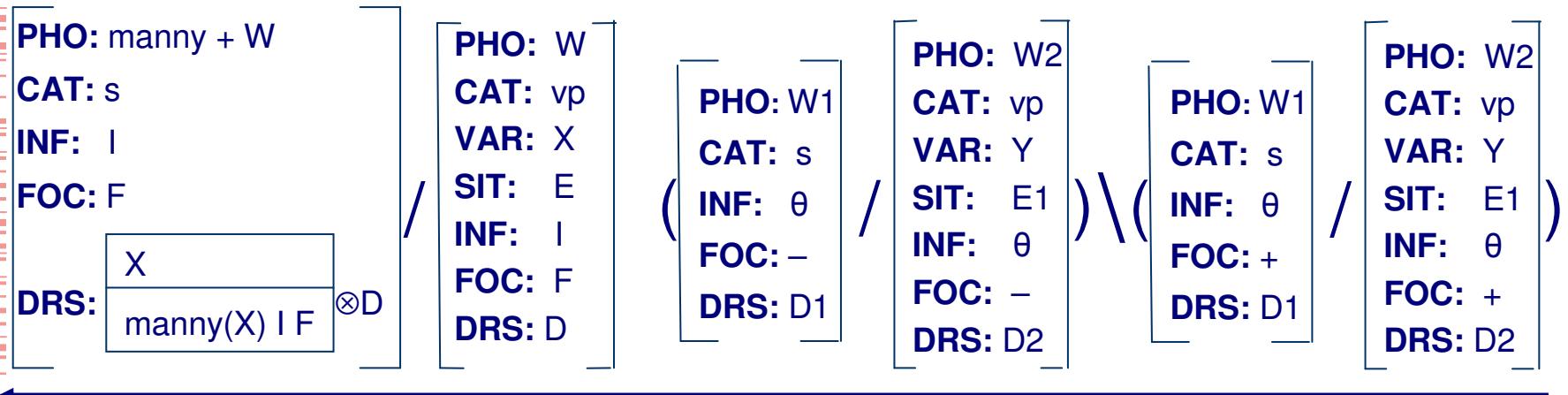
# Combining with a Pitch Accent

## ◆ Manny L+H\*



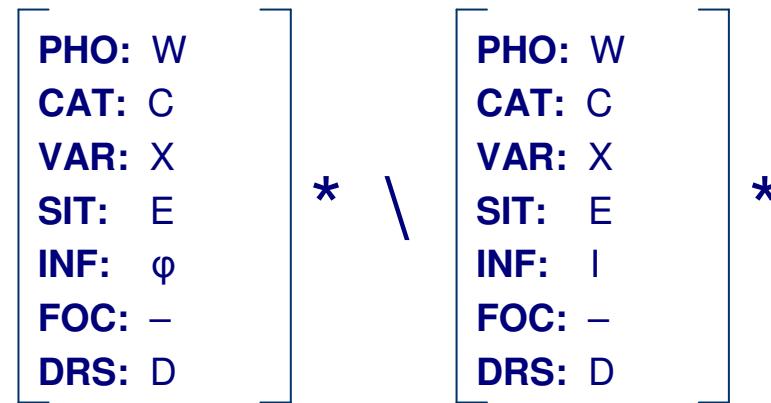
# Combining with a Pitch Accent

## ◆ Manny L+H\*



# Boundary Tones

- ◆ Boundary tones “freeze” intonational phrases:  
cannot combine with anything else but a complete phrase



# Combining with a Boundary Tone

- ◆ married Manny L+H\* LH%

<b>PHO:</b> married + manny
<b>CAT:</b> vp
<b>VAR:</b> Y
<b>SIT:</b> E
<b>INF:</b> θ
<b>FOC:</b> –
<b>DRS:</b> X, E manny(X): θ + married(E): θ – patient(E,X) agent(E,Y)

<b>PHO:</b> W
<b>CAT:</b> C
<b>VAR:</b> Z
<b>SIT:</b> E1
<b>INF:</b> φ
<b>FOC:</b> –
<b>DRS:</b> D

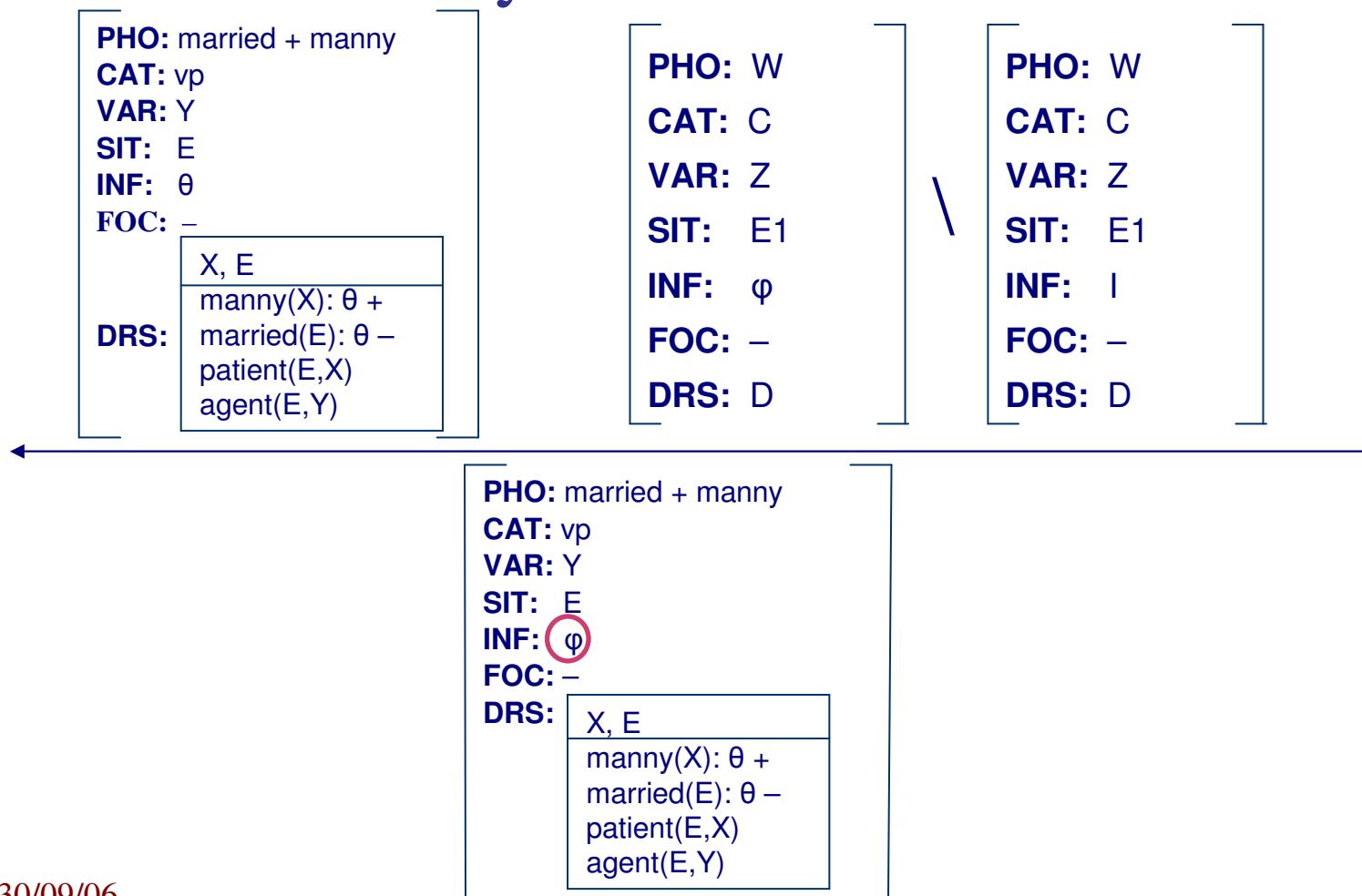
\* \

<b>PHO:</b> W
<b>CAT:</b> C
<b>VAR:</b> Z
<b>SIT:</b> E1
<b>INF:</b> I
<b>FOC:</b> –
<b>DRS:</b> D



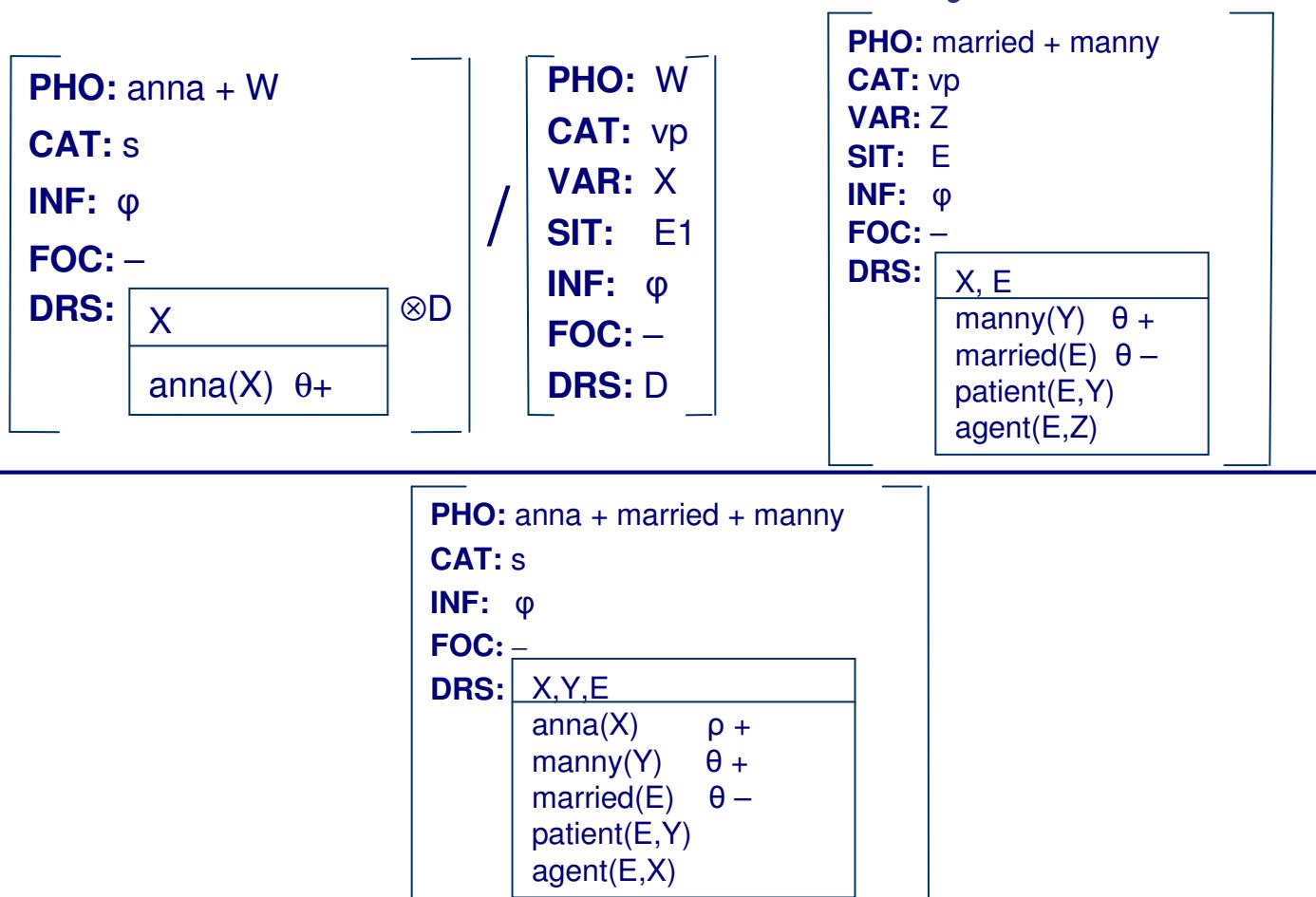
# Combining with a Boundary Tone

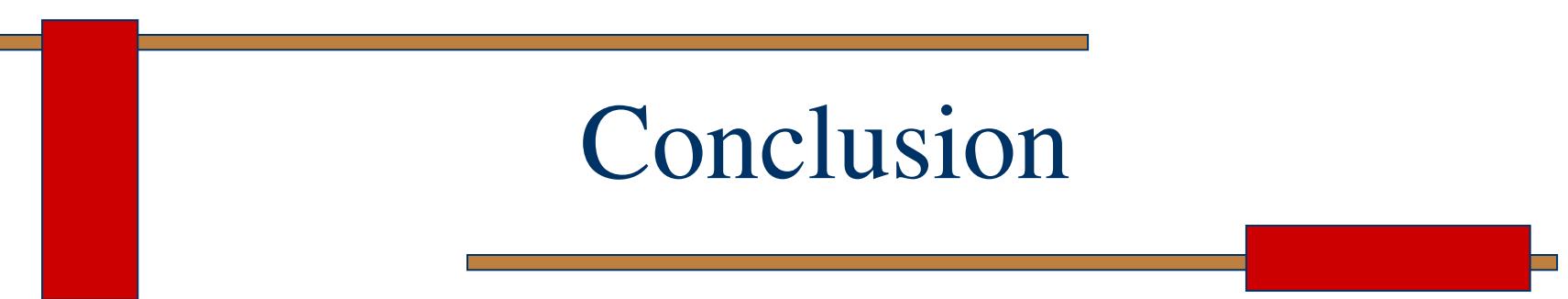
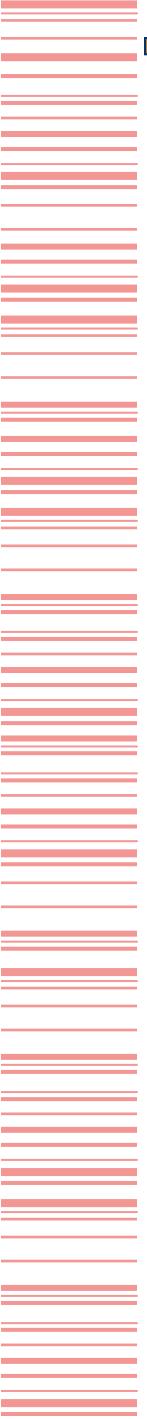
- ◆ married Manny L+H\* LH%



# Combining of Full Intonational Phrases

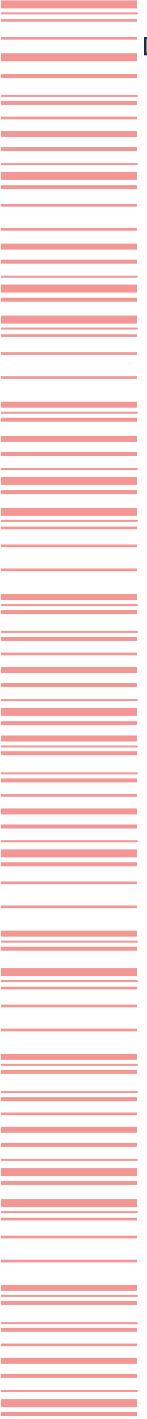
- Anna H\* LL% married Manny L+H\* LH%





# Conclusion

- ◆ Integrated DRT, IS and CCG into a single framework
  - IS-DRS
    - DRT combined with a theory of information structure
  - Unification-based Combinatory Categorial Grammar
    - parsing and generating prosodically annotated text
    - compositional analysis of information structure



# Future Work



- ◆ Enhance the existing UCCG parser
- ◆ Implement a generation component based on UCCG
- ◆ Use the formalism in existing spoken dialogue systems

For a more complete account see:

<http://www.iccs.inf.ed.ac.uk/~s0129610/thesis.ps.gz>