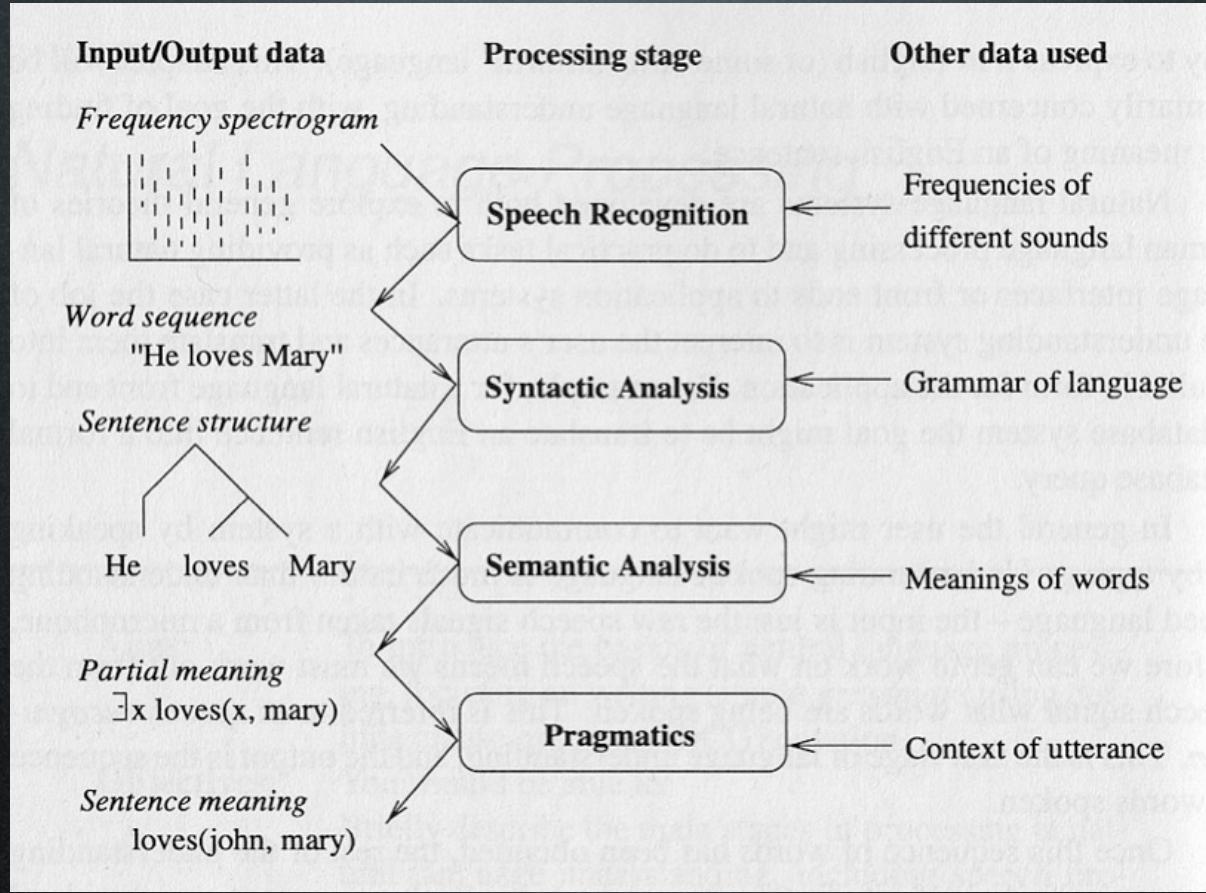


More of the Same Boring AI Stuff

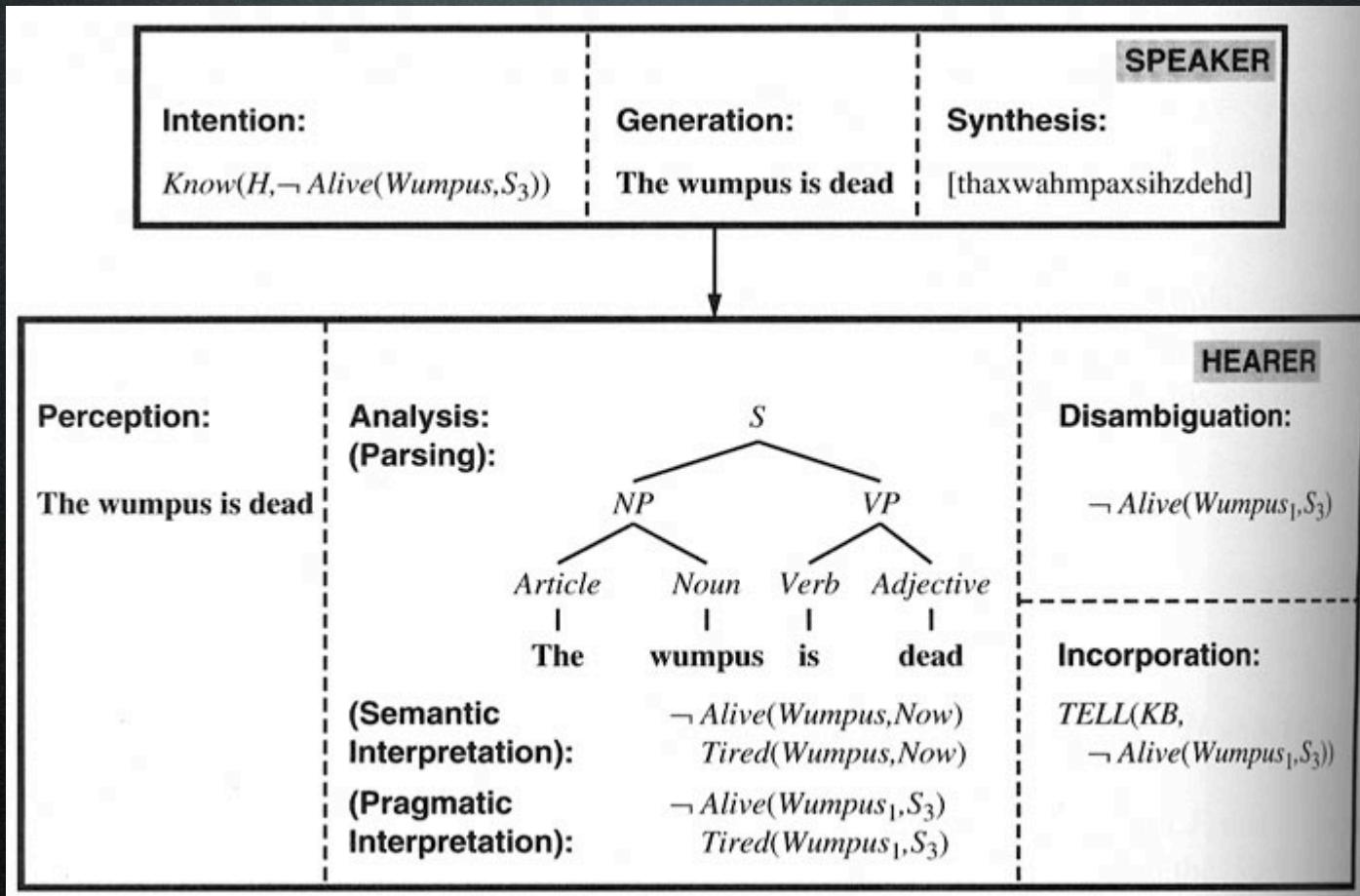
M. Anthony Kapolka III
CS 340 AI Fall 2019
Wilkes University

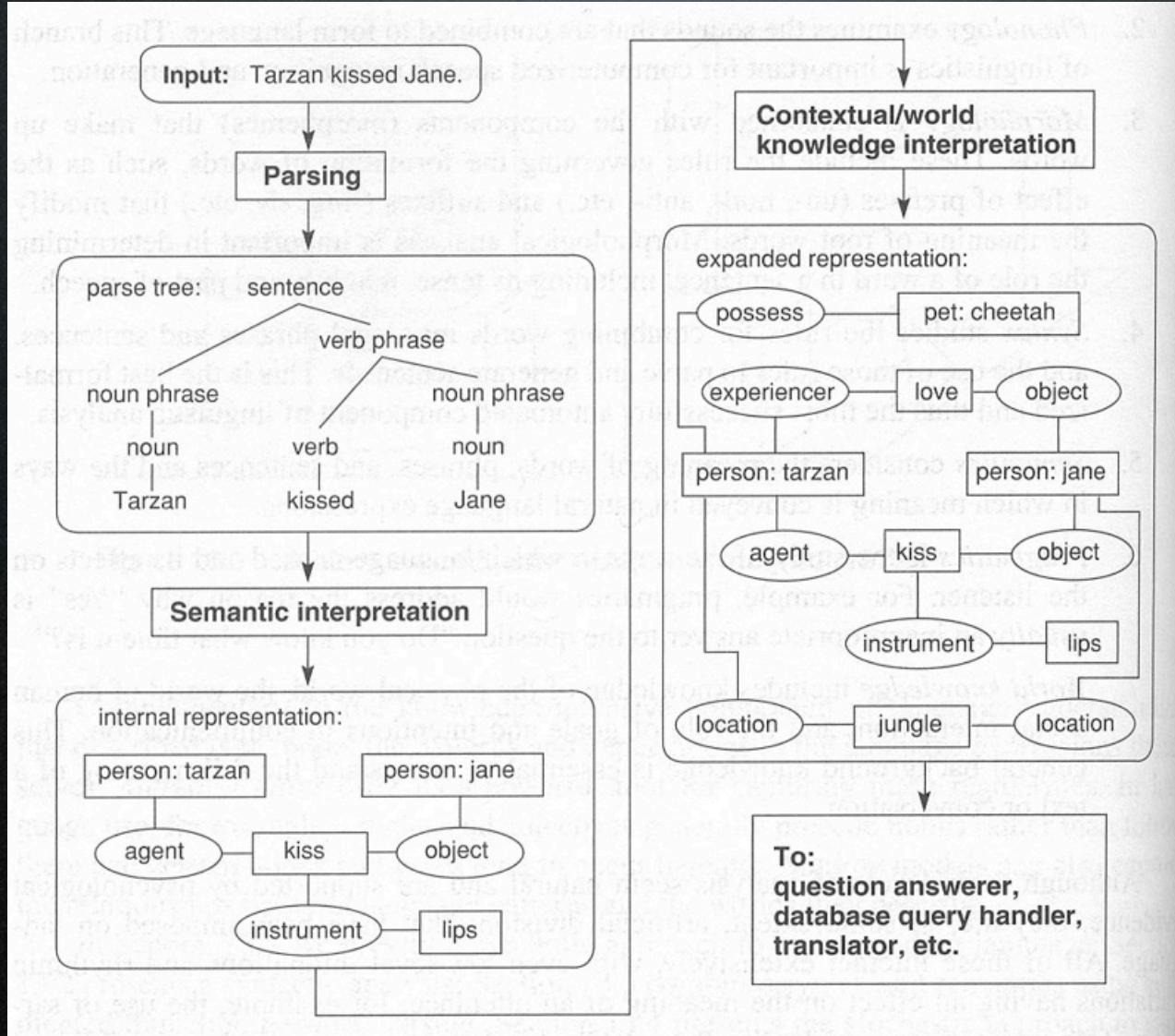


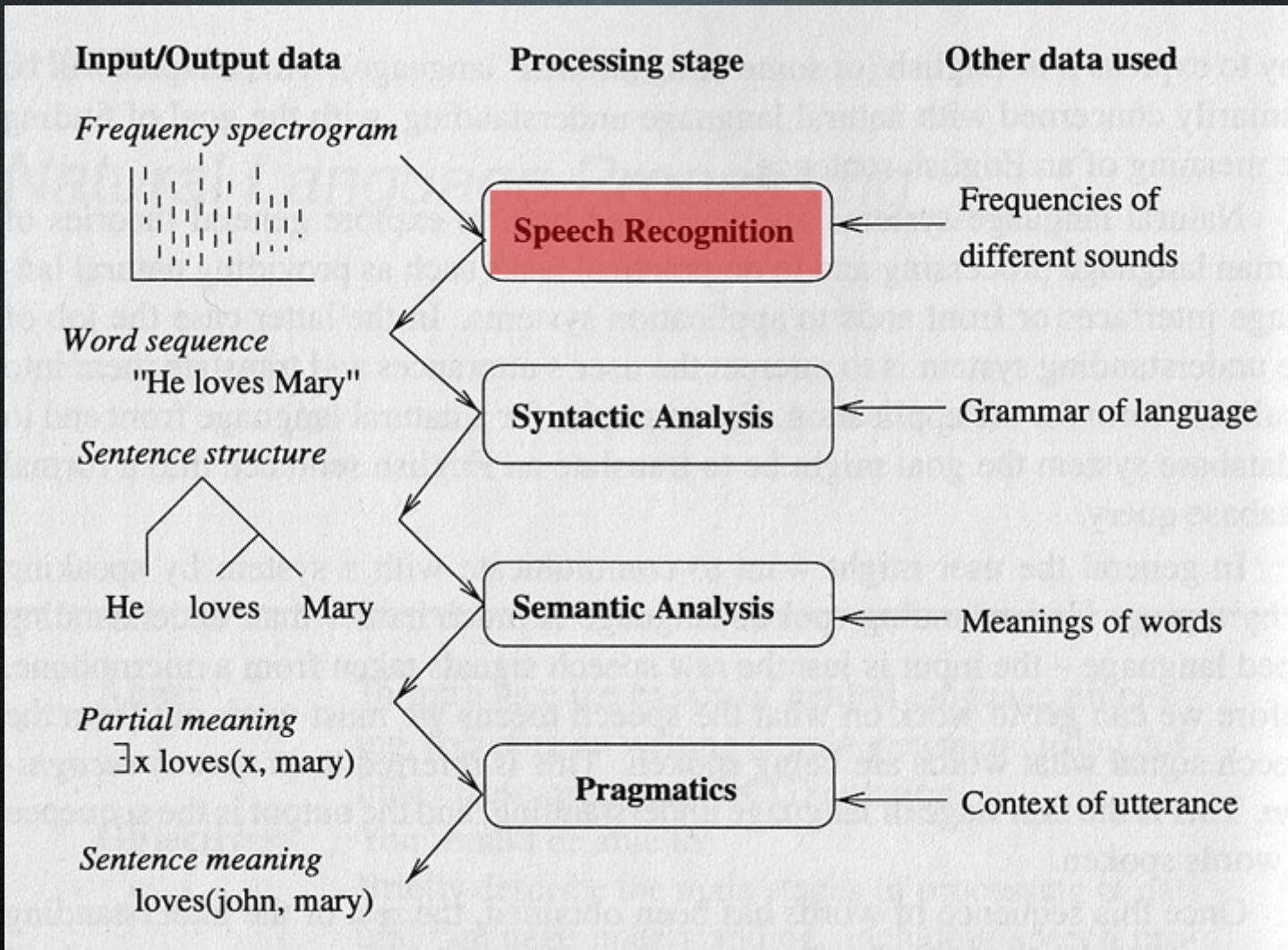
How to wreck a nice beach.



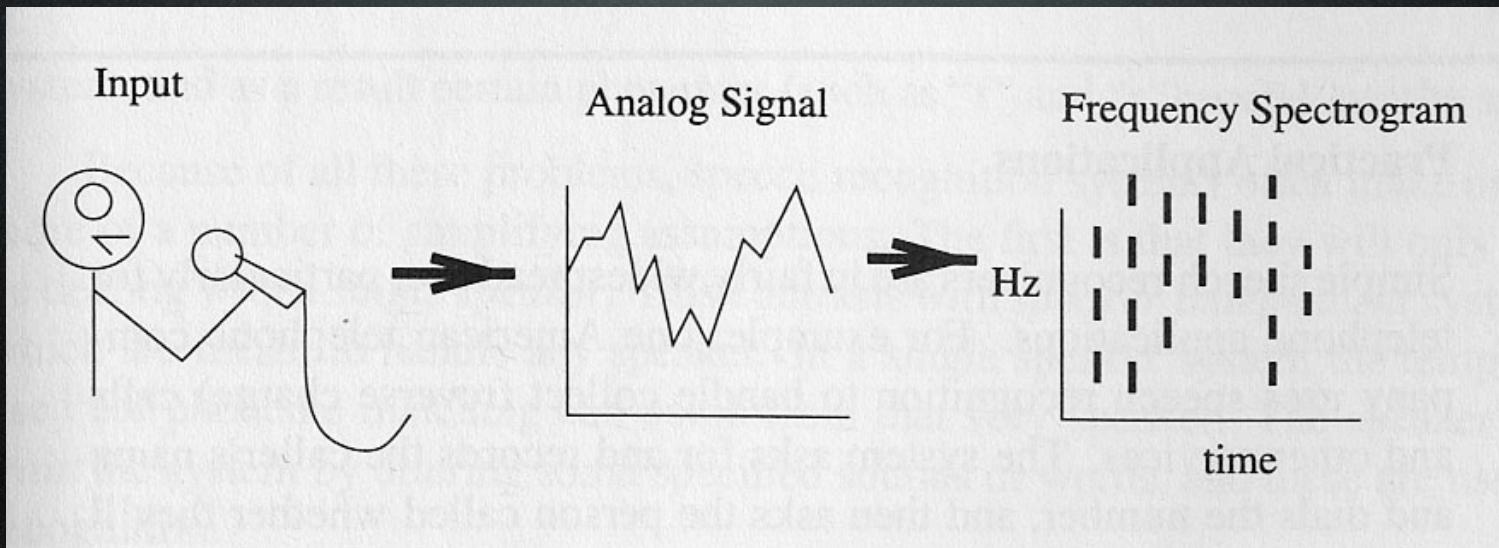
How to recognize speech.







Signal Processing



Signal Processing

- Sampling rate
 - how often a data point is recorded
 - typically 8-16 KHz
- Quantization (Precision)
 - Number of bits for each sample point
 - typically 8-12 bits

Lots of Data

8 KHz, 8 bit sampling =
about 1/2 MB per minute

Speech Recognition

Mapping of digitally encoded acoustic signal into a string of characters

(e.g. SPEECH -> TEXT.)

Speech Comprehension

- Goal: ability to understand meaning
- Harder than recognition - domain-specific
- Easier than recognition - close enough is ok
- Can use context to identify words

Speech Recognition

- Constrained
 - many commercial vendors, good success
 - lightweight apps running on PDAs & cellphones

constrained voice recognition

WED 30 NOV 2005

Get Your Own R2

READ MORE: R2-D2, ROBOTS



I know you've been waiting for this—don't deny it. Finally, R2-D2 can be yours, in your own home, and for only \$119.95. Fully functional, R2 can roam around and **obeys over 40 voice commands** (please refrain from asking him to project Obi-Wan please). He also plays children's games like tag with an IR sensor that can be set to detect motion (which can also be set to sound an alarm for a watchdog effect). Available at Hammacher Schlemmer.

[Star Wars R2-D2 Robot Replica On Sale \[i4u\]](#)

[R2-D2 Robot Product Page \[Hammacher Schlemmer\]](#)

Speech Recognition

- Unconstrained (e.g. Dictation)
 - few vendors, commercial systems
 - constrained by dictionary
 - adaptive, requires training
 - error rates 1 order of magnitude > human
 - most users disappointed
 - human-like performance projected 2050.

Speech Recognition

- Human speech consists of 40-50 phonemes
- A phoneme is a basic speech sound,
- Some letters/letter combinations map to multiple phonemes.

English Phonemes

[iy] , <u>beat</u>	[ih] , <u>bit</u>	[ey] , <u>bet</u>	[ae] , <u>bat</u>
[ah] , <u>but</u>	[ao] , <u>bought</u>	[ow] , <u>boat</u>	[uh] , <u>book</u>
[ux] , <u>beauty</u>	[er] , <u>Bert</u>	[ay] , <u>buy</u>	[oy] , <u>boy</u>
[axr] , <u>diner</u>	[aw] , <u>down</u>	[ax] , <u>about</u>	[ix] , <u>roses</u>
[aa] , <u>cot</u>	[b] , <u>bet</u>	[ch] , <u>Chet</u>	[d] , <u>debt</u>
[f] , <u>fat</u>	[g] , <u>get</u>	[hh] , <u>hat</u>	[hv] , <u>high</u>
[jh] , <u>jet</u>	[k] , <u>kick</u>	[l] , <u>let</u>	[el] , <u>bottle</u>
[m] , <u>met</u>	[em] , <u>bottom</u>	[n] , <u>net</u>	[en] , <u>button</u>
[ng] , <u>sing</u>	[eng] , <u>Washington</u>	[p] , <u>pet</u>	[r] , <u>rat</u>
[s] , <u>set</u>	[sh] , <u>shoe</u>	[t] , <u>ten</u>	[th] , <u>thick</u>
[dh] , <u>that</u>	[dx] , <u>butter</u>	[v] , <u>vet</u>	[w] , <u>wet</u>
[wh] , <u>which</u>	[y] , <u>yet</u>	[z] , <u>zoo</u>	[zh] , <u>measure</u>

Detecting Phonemes

- Often the signal is processed with a sliding window (called a frame)
- typical width 10 milliseconds

Process signal...

- Identify features (frequency, amplitude) to identify phonemes.
- Using dictionary, assemble phonemes into words

Process signal...

- Identify features (frequency, amplitude) to identify phonemes.
- Using dictionary, assemble phonemes into words

[k], [ae], [t]  cat

Problems

- homophones (homonyms)
- segmentation (identifying word breaks
 - not always present in fluent speech)

Solutions

- Extract most likely string of words
- Extract multiple possible words
- Invoke Parser / Semantic Analysis

Can use a probabilistic model

This model is broken down into components
using Bayes rule.

$P(\text{words} \mid \text{signal}) =$

$$(P(\text{words}) \times P(\text{signal} \mid \text{words})) / P(\text{signal})$$

Given a signal, the task is to find the sequence
of words that maximizes $P(\text{words} \mid \text{signal})$.

$$(P(\text{words}) \times P(\text{signal} \mid \text{words})) / P(\text{signal})$$

$P(\text{signal} \mid \text{words})$ is the acoustic model.

[k], [ae], [t]  cat

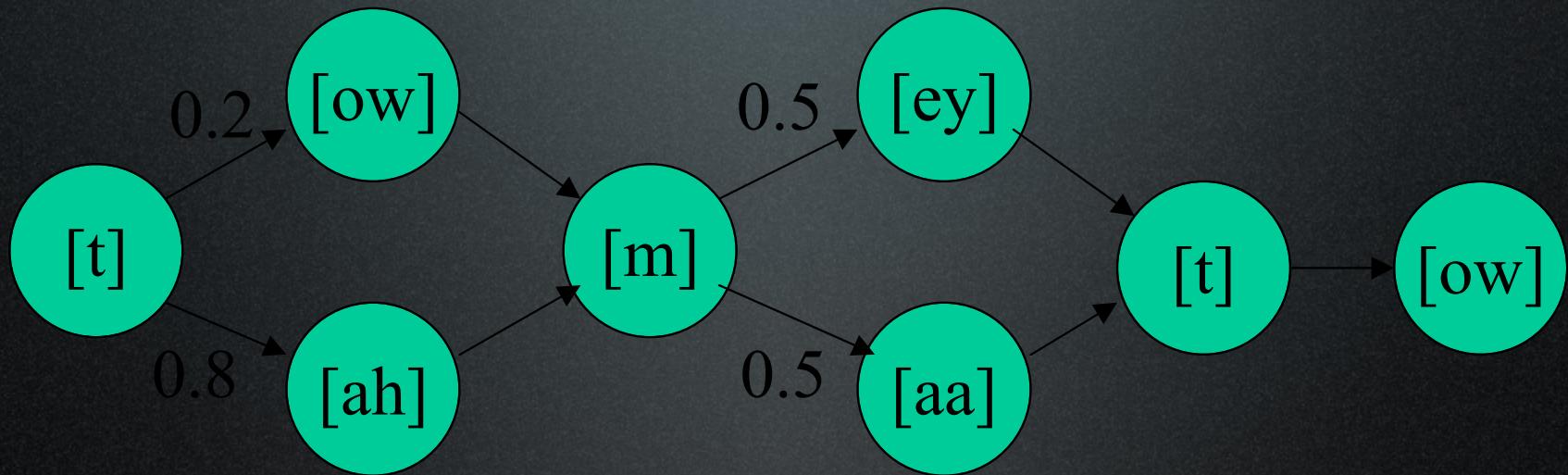
Problems:

- 1) dialects (different pronunciations)
- 2) coarticulation (quick, slurred speech)



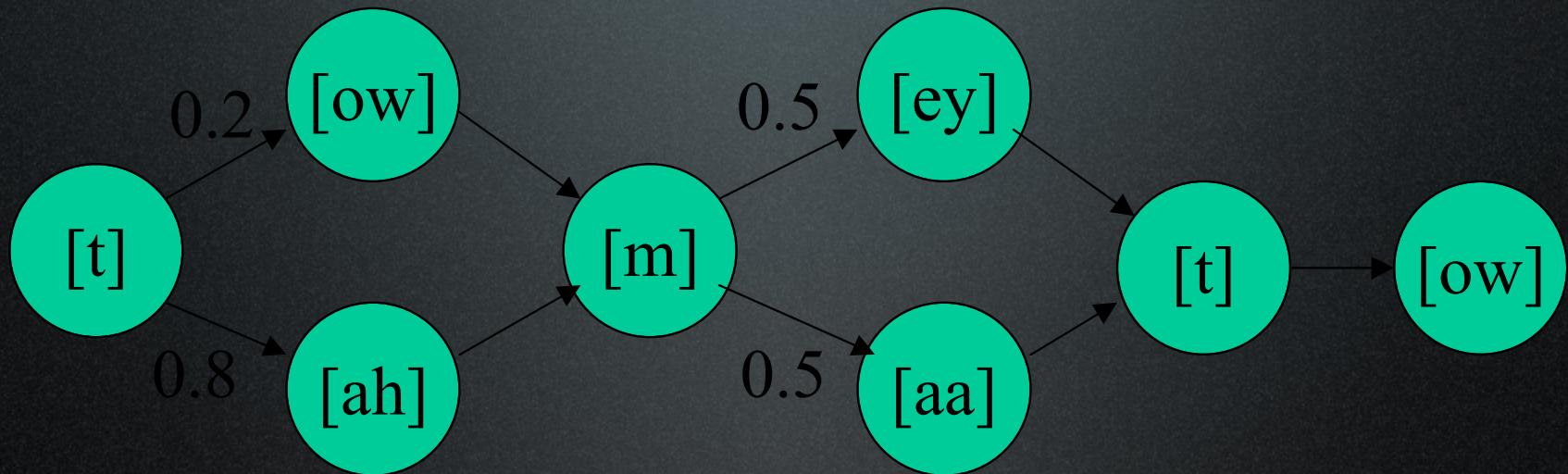
Josh Diak, Citizen of Australia

you say tomato...



example of a markov model showing all phoneme combinations for tomato (dialect & coarticulation)

you say tomato...



$$P([\text{towmeytow}] \mid \text{"tomato"}) = 0.1$$

$$P([\text{tahmaatow}] \mid \text{"tomato"}) = 0.4$$



I say...

t ah mey t ae hh

$$(P(\text{words}) \times P(\text{signal} \mid \text{words})) / P(\text{signal})$$

$P(\text{words})$ is the language model.

That is, probability based
on language context.

$$(P(\text{words}) \times P(\text{signal} \mid \text{words})) / P(\text{signal})$$



the man is saying “I have a ...”

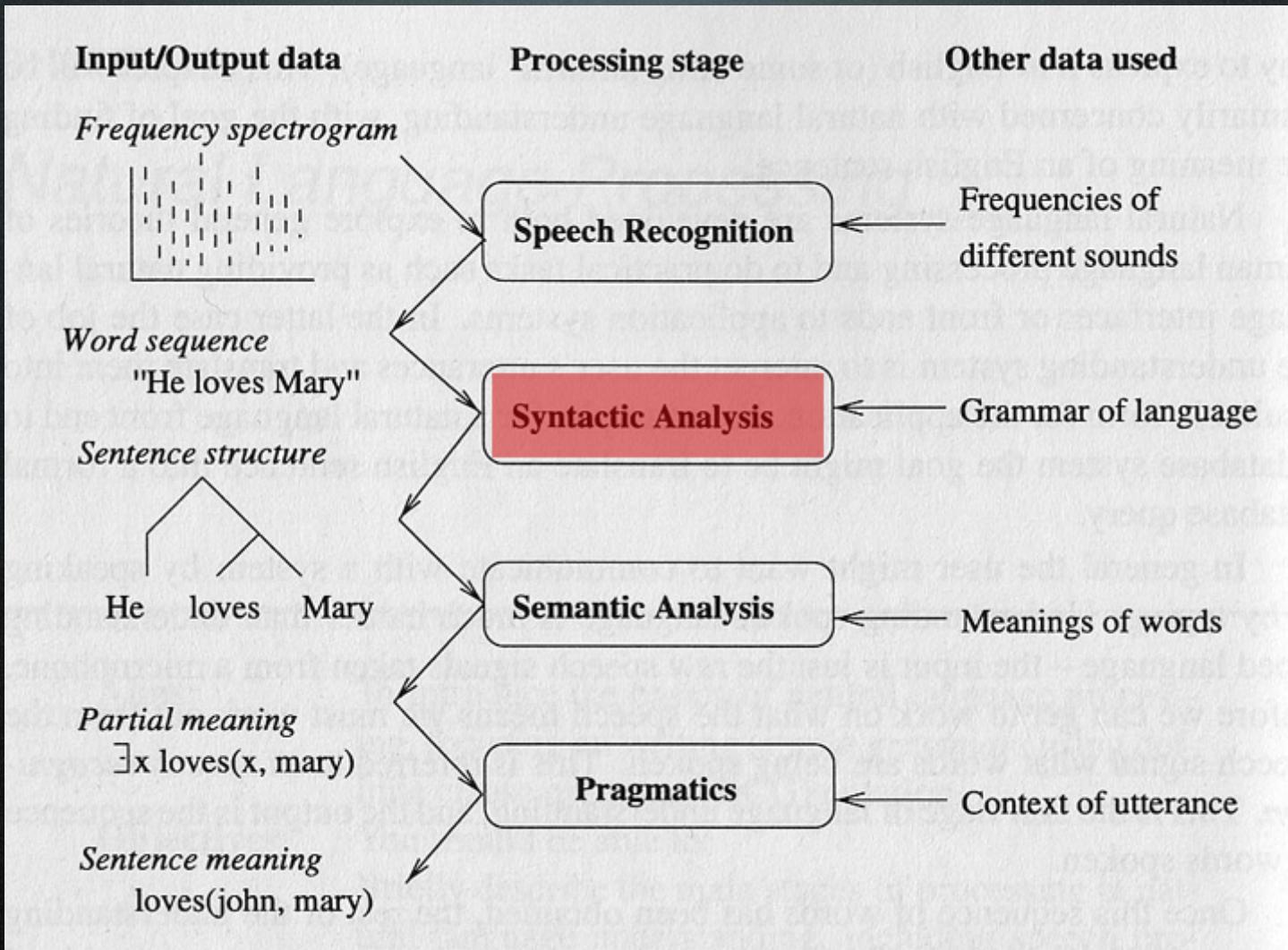
$$(P(\text{words}) \times P(\text{signal} \mid \text{words})) / P(\text{signal})$$



$$P(\text{signal} \mid \text{gum}) = 0.7 \quad P(\text{signal} \mid \text{gun}) = 0.4$$

$$(P(\text{words}) \times P(\text{signal} \mid \text{words})) / P(\text{signal})$$

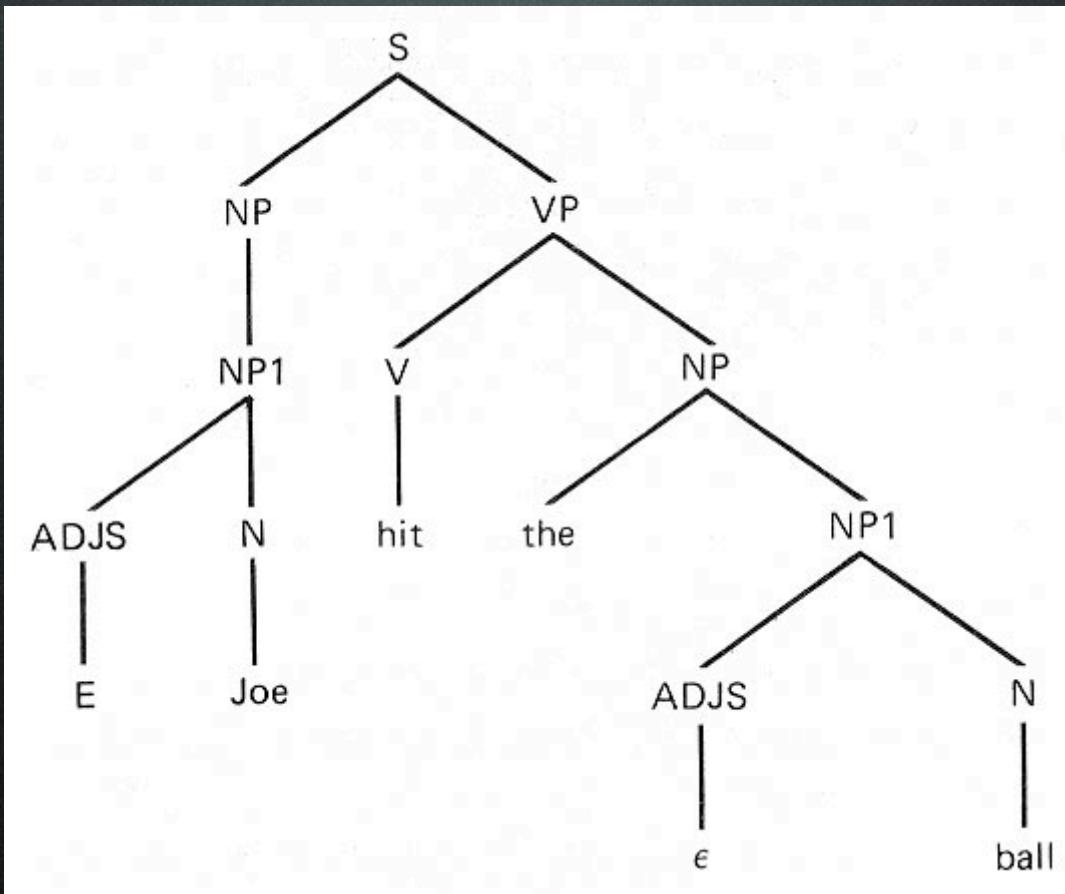
- $P(\text{words})$ is a factor to add weight to more likely utterances.
- How to compute $P(\text{words})$??
- Coincident words
- Feedback from speech comprehension.



Grammar

```
S → NP VP
NP → the NPl
NP → NPl
NPl → ADJS N
ADJS → ε | ADJ ADJS
VP → V
VP → V NP
N → Joe | boy | ball
ADJ → little | big
V → hit | ran
```

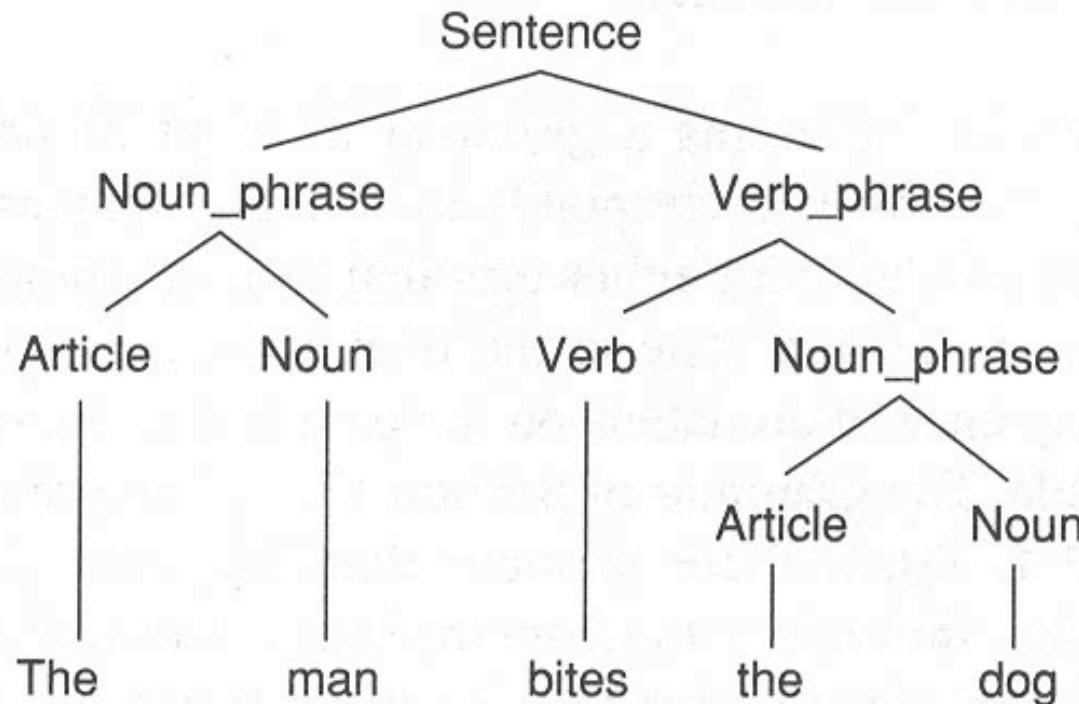
Parse Tree



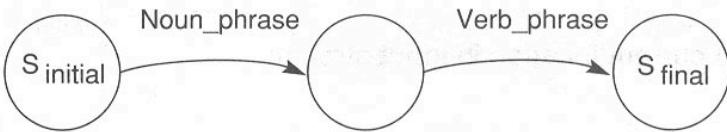
A Dog Grammar

1. sentence \leftrightarrow noun_phrase verb_phrase
2. noun_phrase \leftrightarrow noun
3. noun_phrase \leftrightarrow article noun
4. verb_phrase \leftrightarrow verb
5. verb_phrase \leftrightarrow verb noun_phrase
6. article \leftrightarrow a
7. article \leftrightarrow the
8. noun \leftrightarrow man
9. noun \leftrightarrow dog
10. verb \leftrightarrow likes
11. verb \leftrightarrow bites

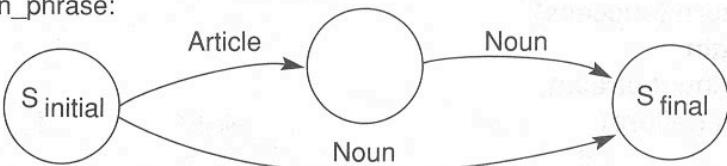
And Parse Tree



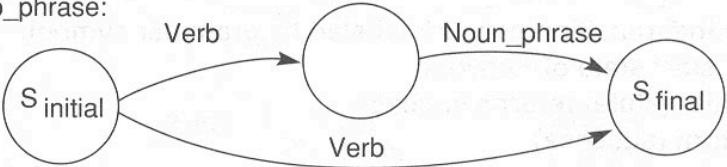
Sentence:



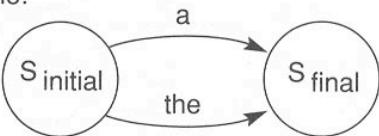
Noun_phrase:



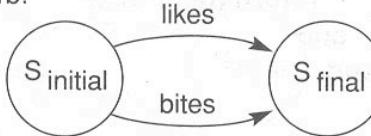
Verb_phrase:



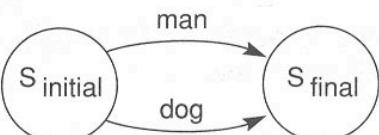
Article:



Verb:



Noun:



sentence:

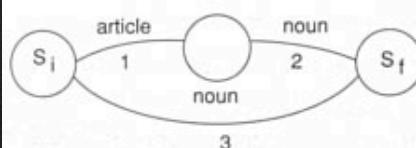


```
function sentence-1;
begin
    NOUN_PHRASE := structure returned by
        noun_phrase network;
    SENTENCE.SUBJECT := NOUN_PHRASE;
end.
```

```
function sentence-2;
begin
    VERB_PHRASE := structure returned by
        verb_phrase network;

    if NOUN_PHRASE.NUMBER =
        VERB_PHRASE.NUMBER
        then begin
            SENTENCE.VERB_PHRASE := VERB_PHRASE;
            return SENTENCE
        end
        else fail
    end.
```

noun_phrase:

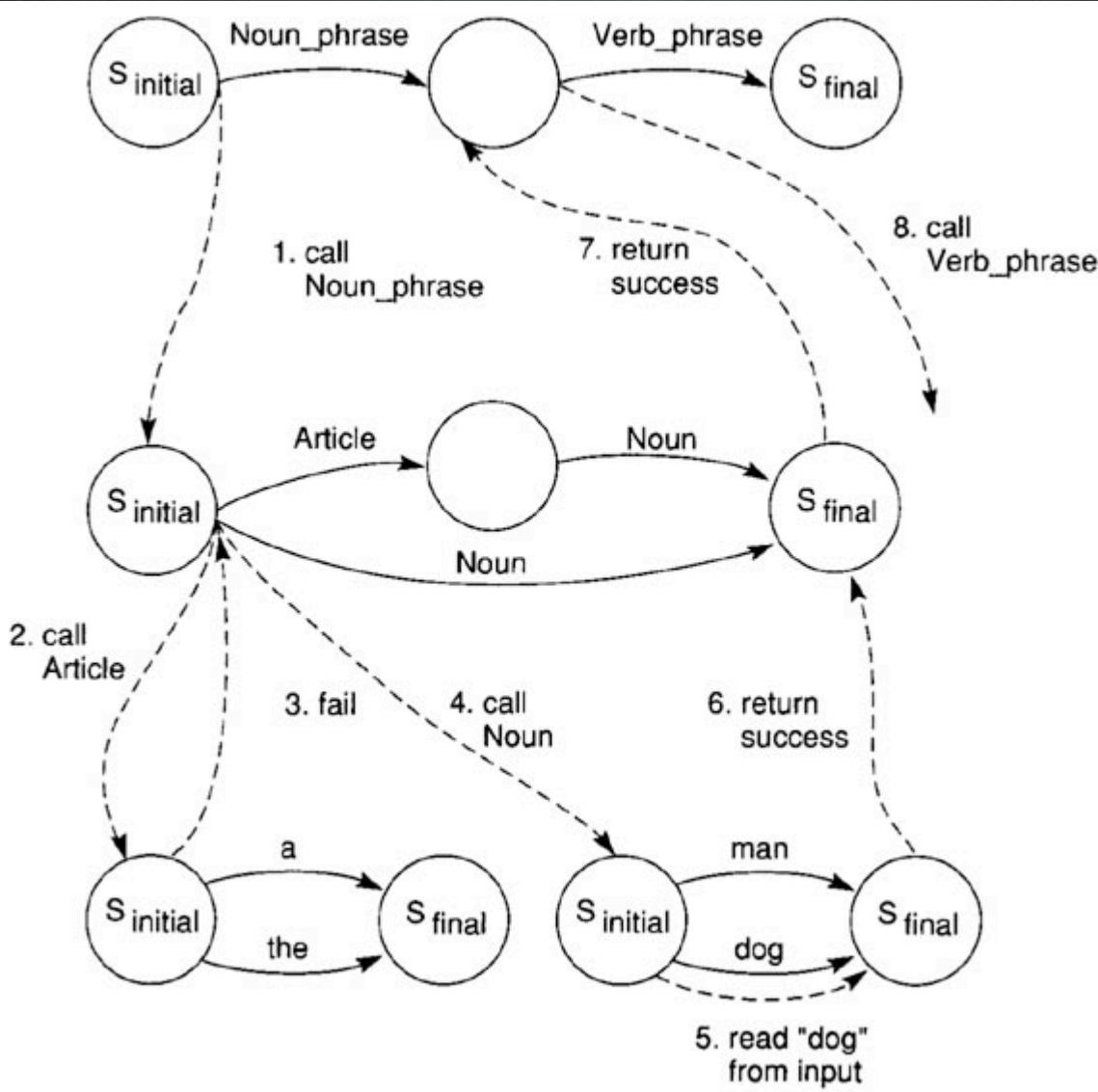


```
function noun_phrase-1;
begin
    ARTICLE := definition frame for next word of input;

    if ARTICLE.PART_OF_SPEECH=article
        then NOUN_PHRASE.DETERMINER := ARTICLE
        else fail
    end.
```

```
function noun_phrase-2;
begin
    NOUN := definition frame for next word of input;

    if NOUN.PART_OF_SPEECH=noun and
        NOUN.NUMBER agrees with
            NOUN_PHRASE.DETERMINER.NUMBER
        then begin
            NOUN_PHRASE.NOUN := NOUN;
            NOUN_PHRASE.NUMBER := NOUN.NUMBER
            return NOUN_PHRASE
        end
        else fail
    end.
```



Classic Jokes

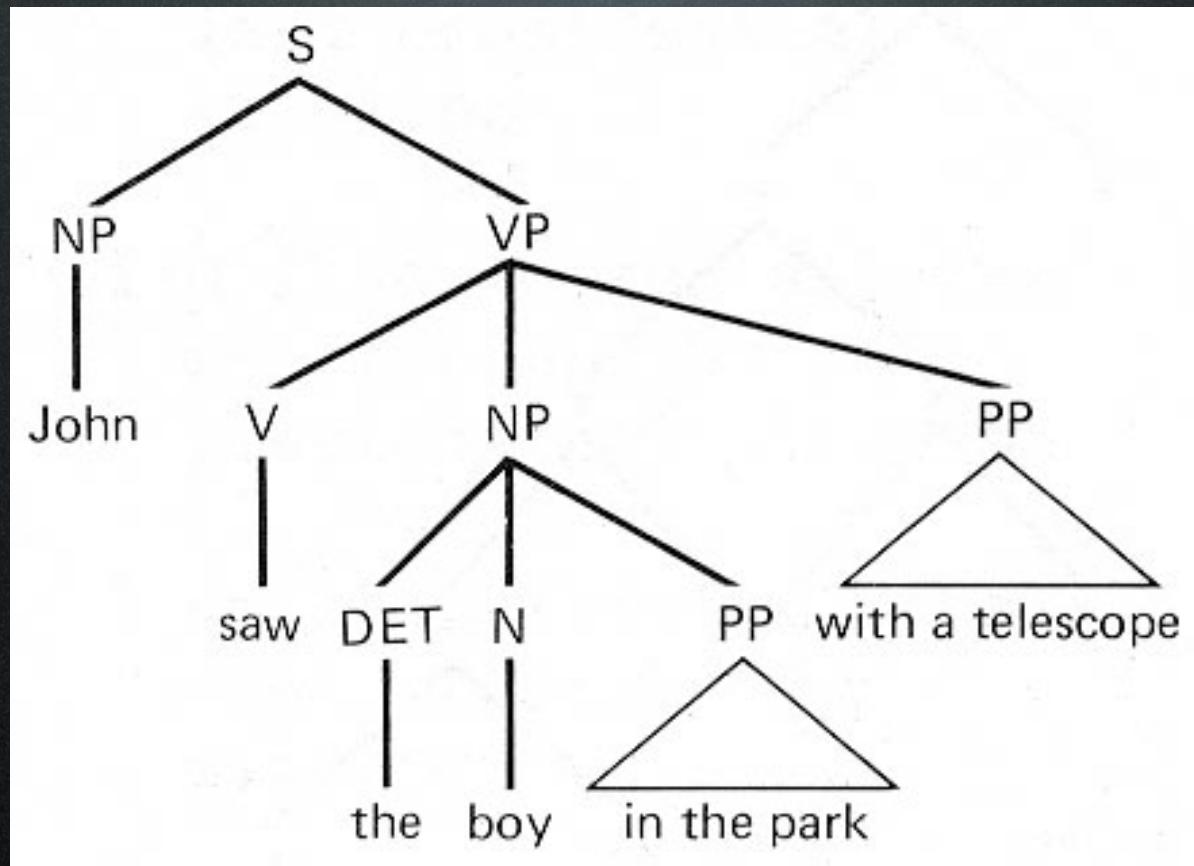
Time flies like an arrow.

Classic Jokes

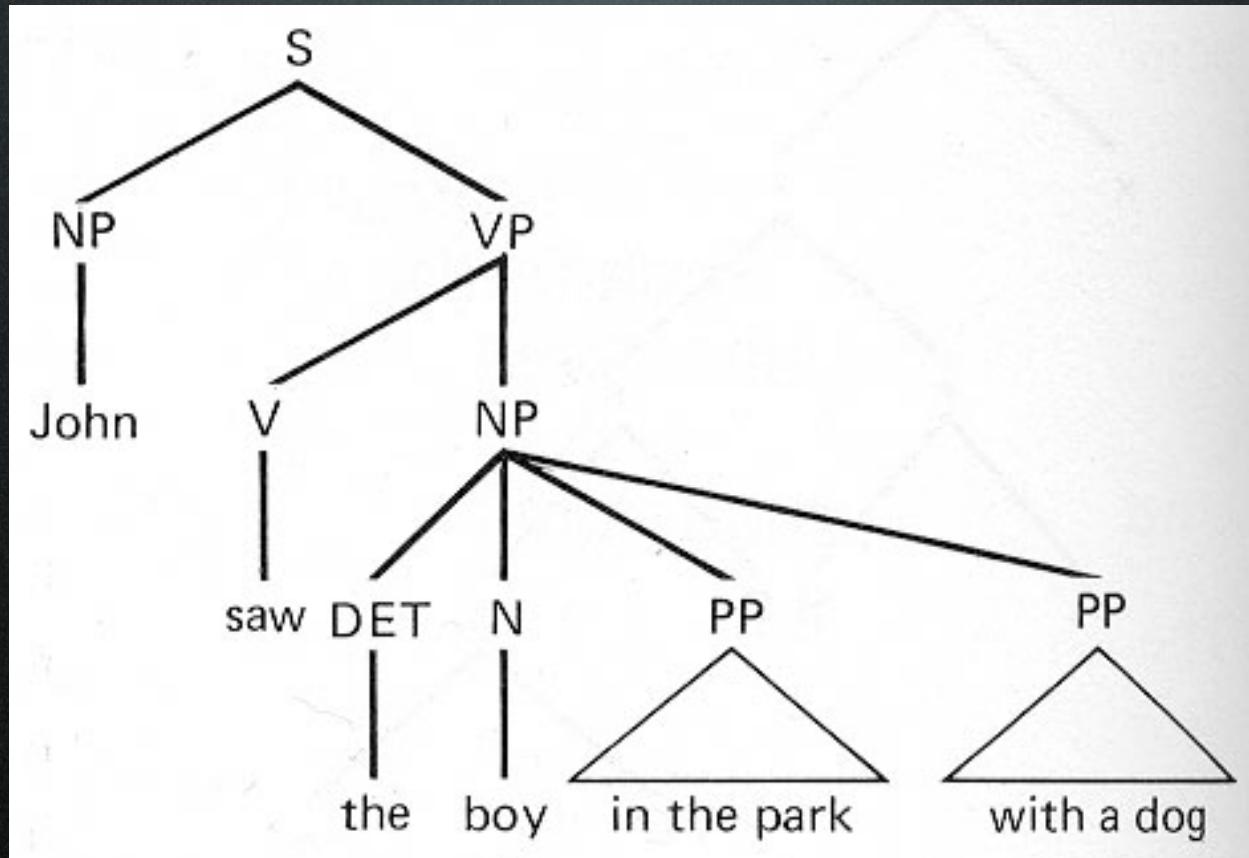
Time flies like an arrow.

Fruit flies like a banana.

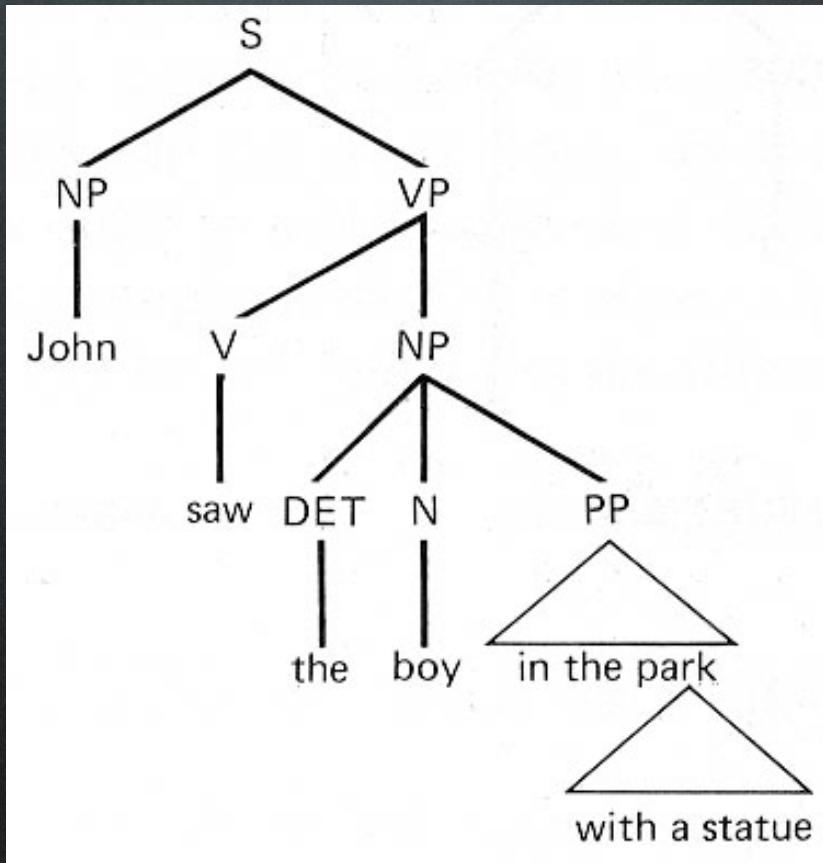
Syntactic Ambiguity



Syntactic Ambiguity

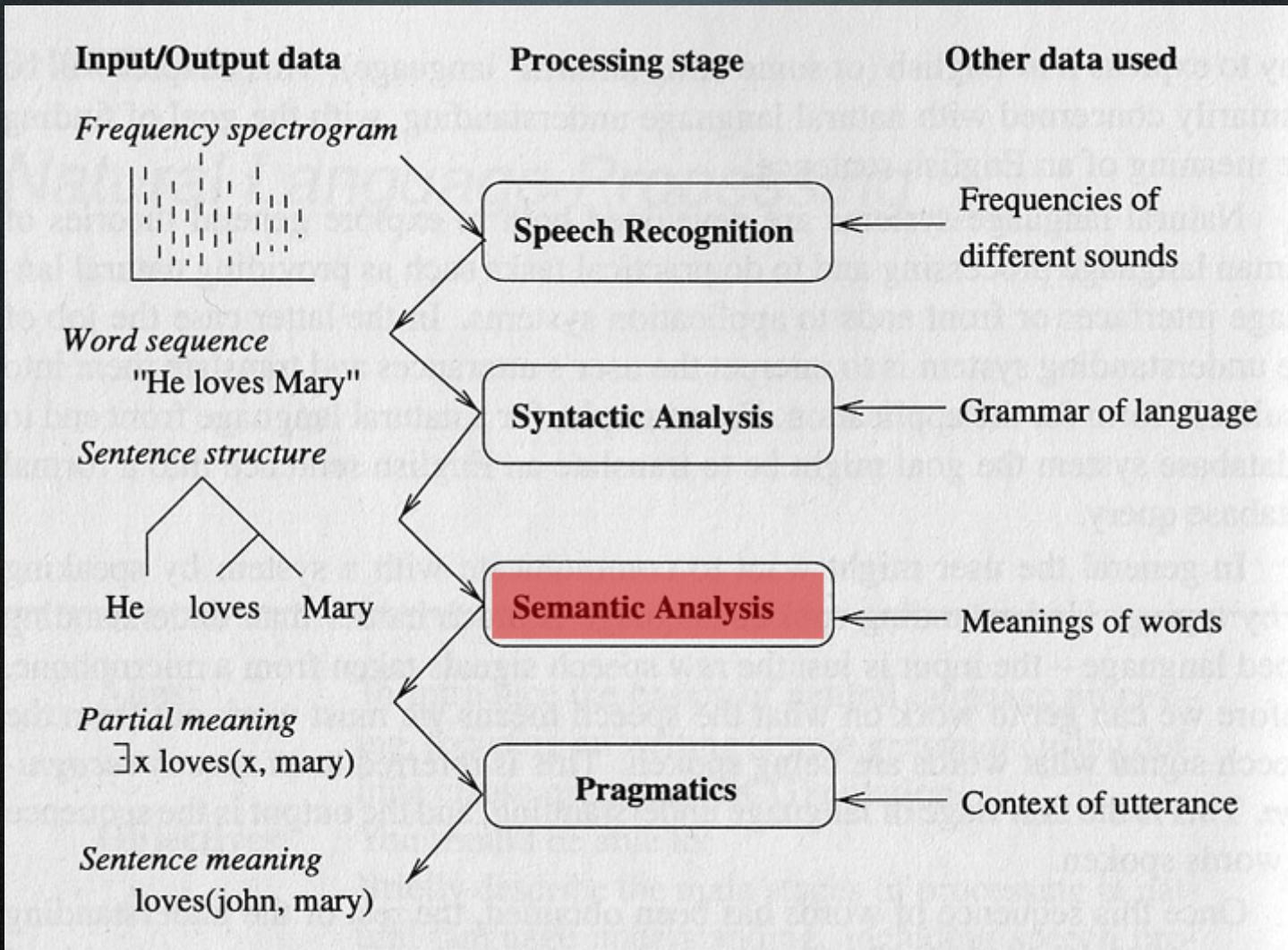


Syntactic Ambiguity



Syntactic Ambiguity

- How do we disambiguate?
- Context
- Feedback from speech comprehension



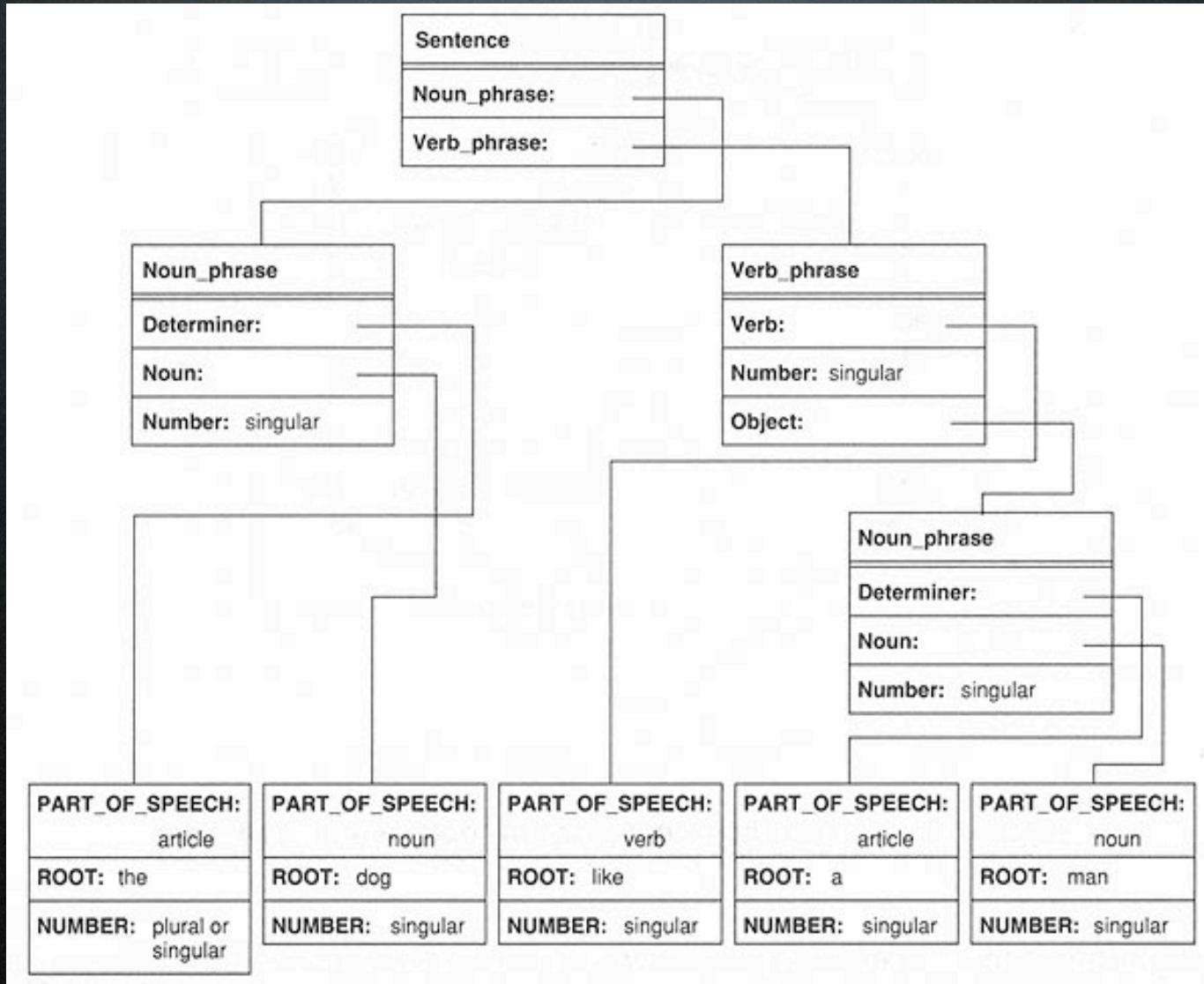
Can use a slot-filler ds

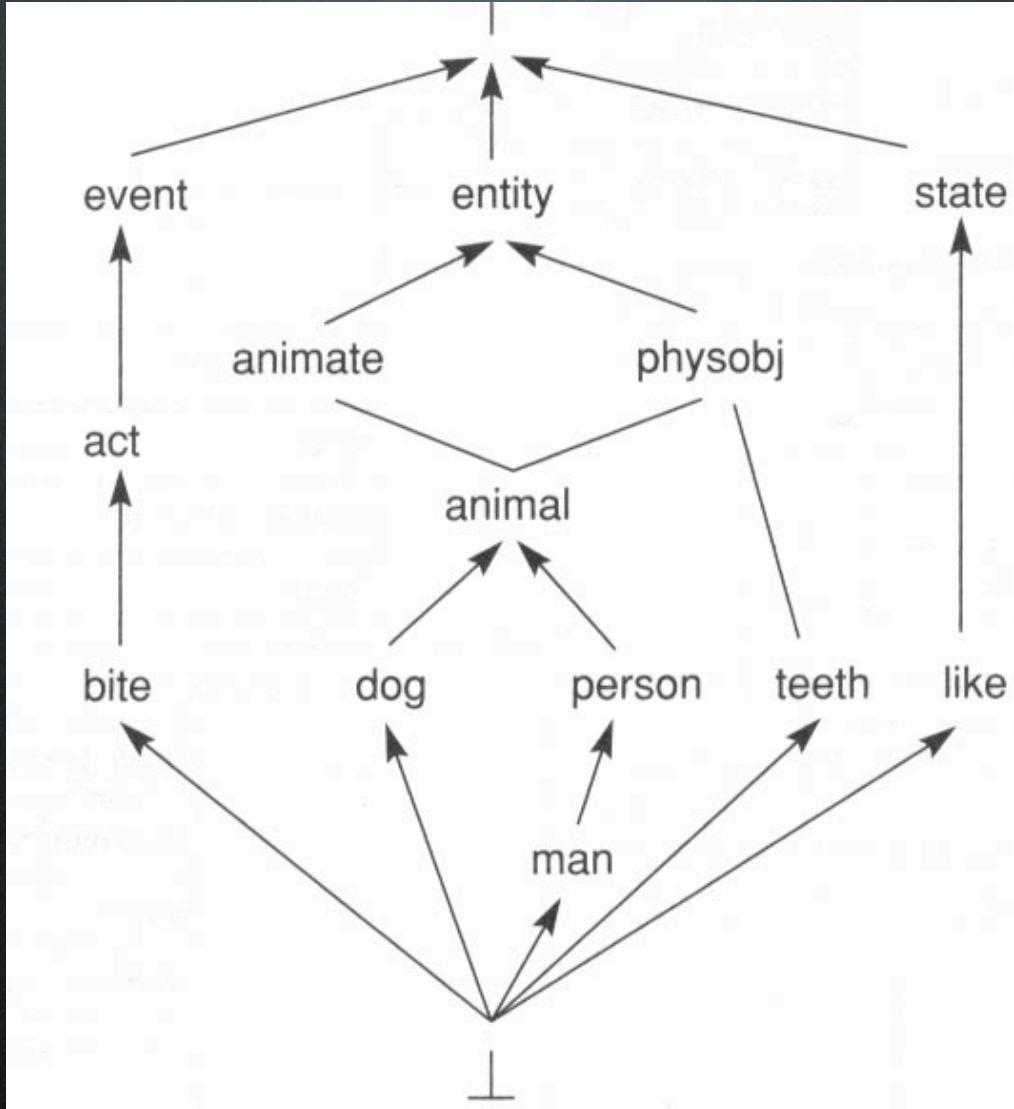
Sentence
Noun phrase:
Verb phrase:

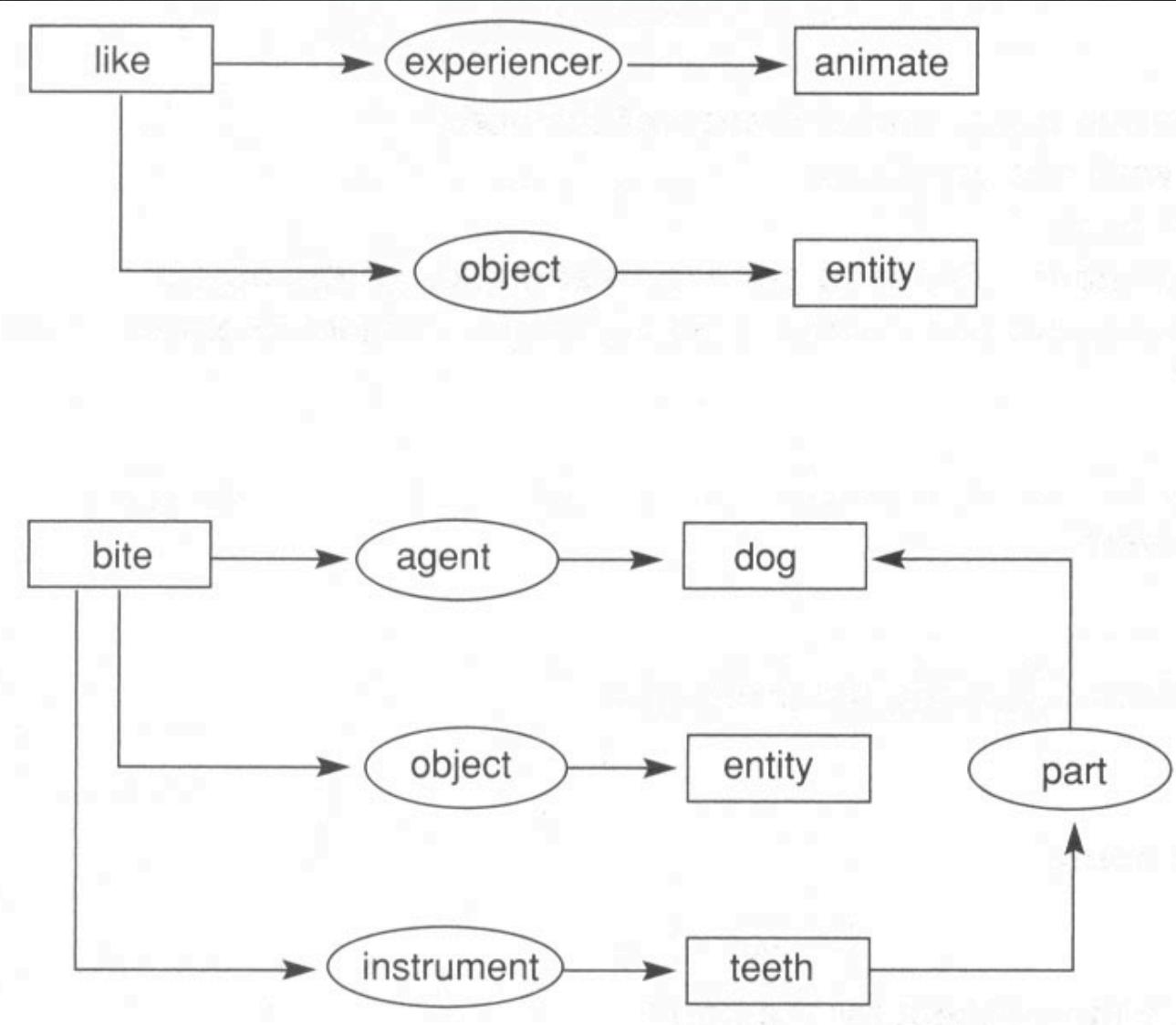
Noun phrase
Determiner:
Noun:
Number:

Verb phrase
Verb:
Number:
Object:

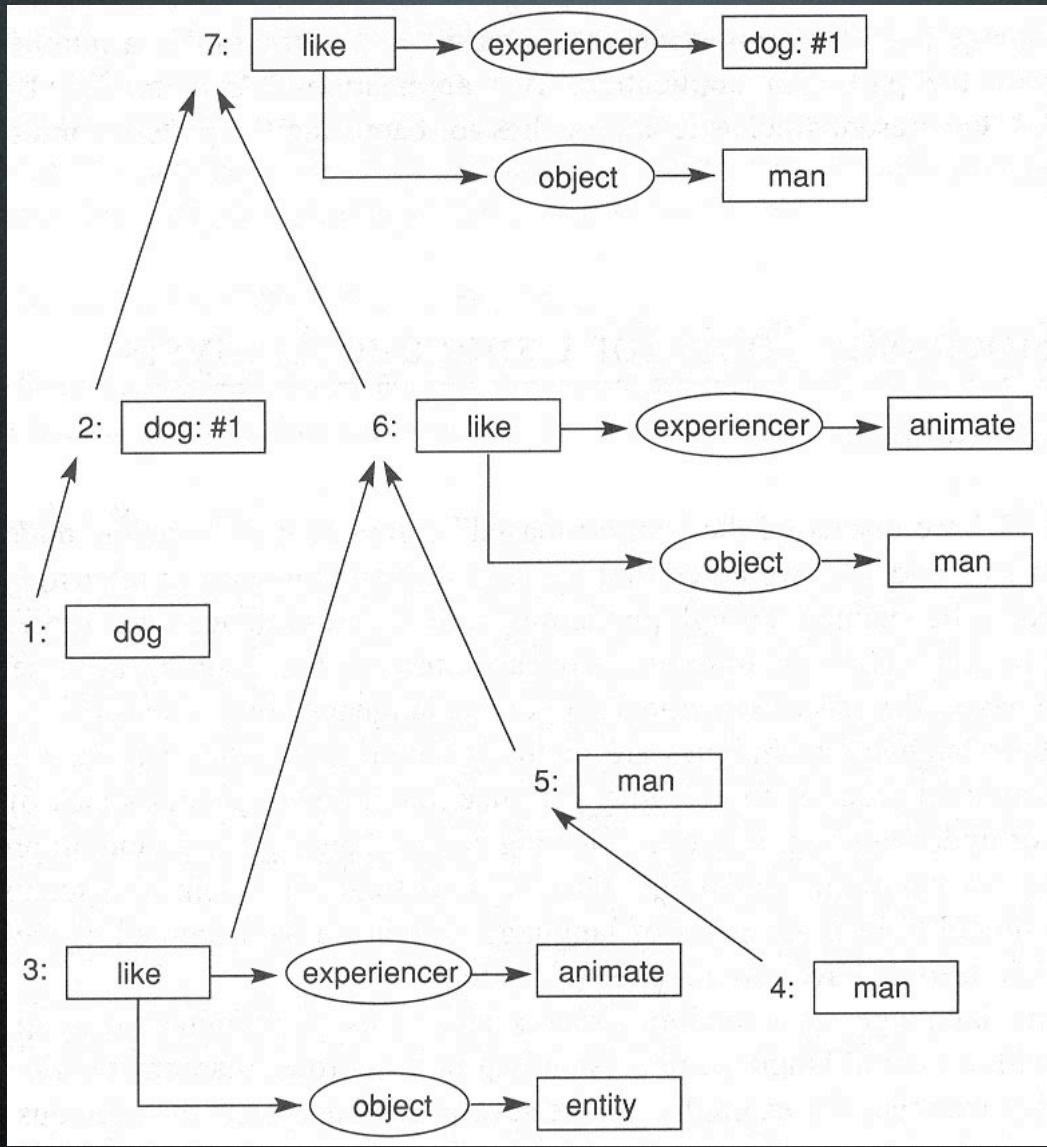
Word	Definition	Word	Definition
a	PART_OF_SPEECH: article ROOT: a NUMBER: singular	like	PART_OF_SPEECH: verb ROOT: like NUMBER: plural
bite	PART_OF_SPEECH: verb ROOT: bite NUMBER: plural	likes	PART_OF_SPEECH: verb ROOT: like NUMBER: singular
bites	PART_OF_SPEECH: verb ROOT: bite NUMBER: singular	man	PART_OF_SPEECH: noun ROOT: man NUMBER: singular
dog	PART_OF_SPEECH: noun ROOT: dog NUMBER: singular	men	PART_OF_SPEECH: noun ROOT: man NUMBER: plural
dogs	PART_OF_SPEECH: noun ROOT: dog NUMBER: plural	the	PART_OF_SPEECH: article ROOT: the NUMBER: plural or singular

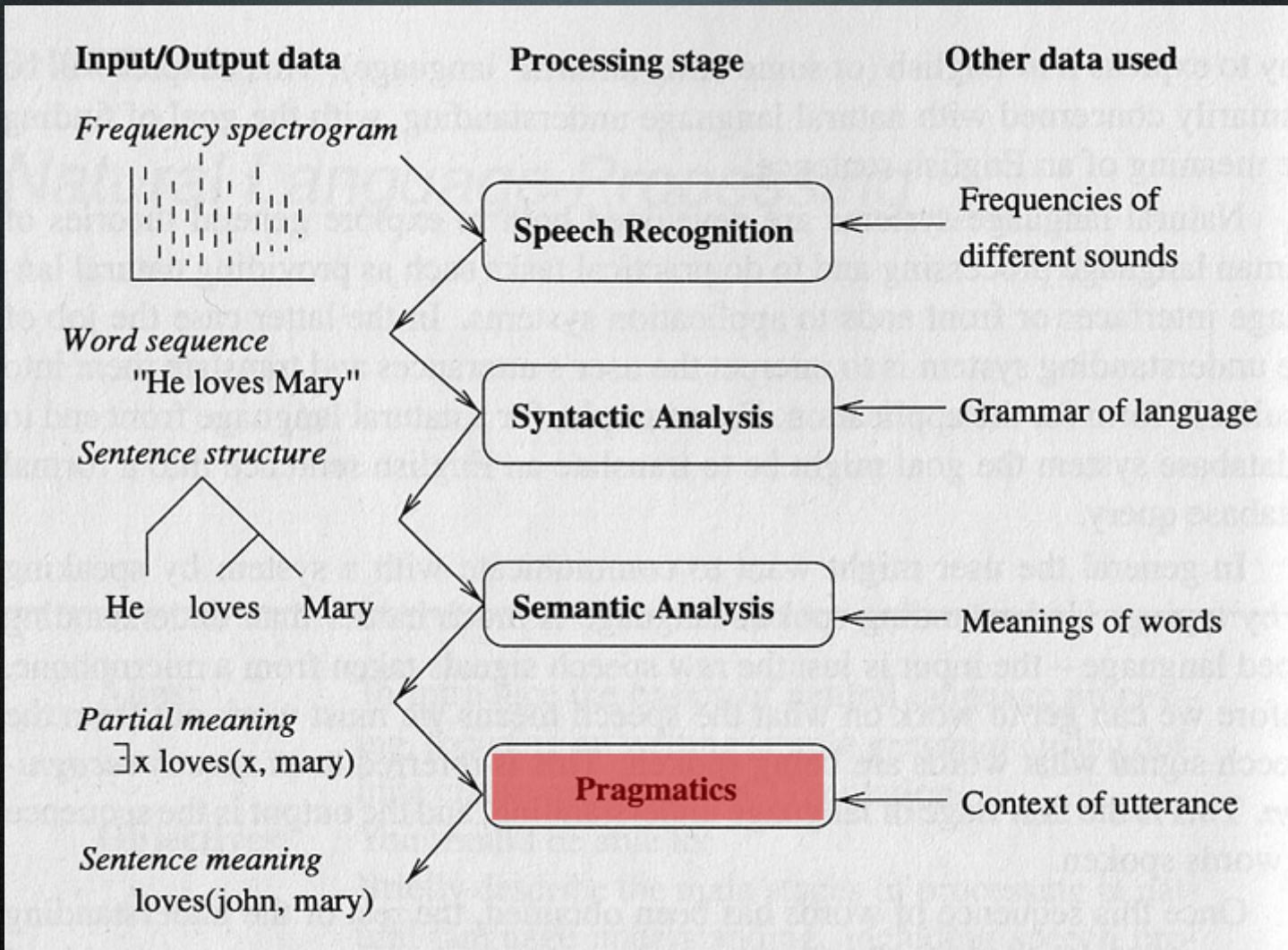






Augment our parse tree with
semantic meaning





Pragmatics

Language expresses more than just facts.

Language can convey a deeper purpose.

When I ask...

Where is the coffee??

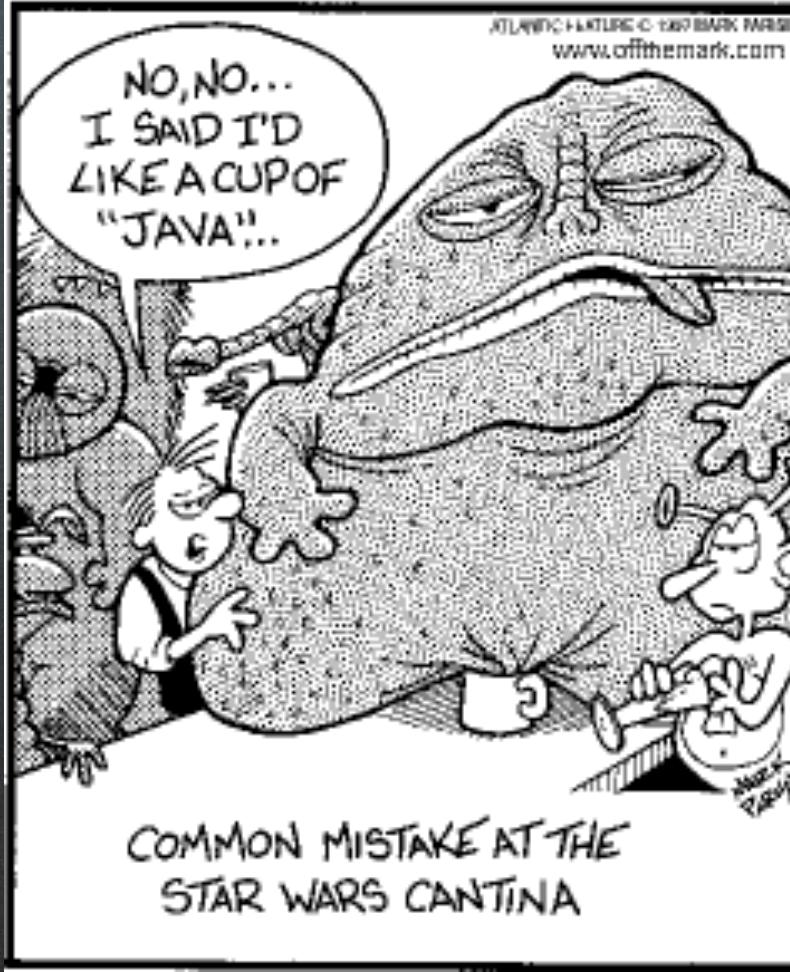
off the mark

www.offthemark.com

by Mark Parisi

ATLANTIC PICTURES © 1989 MARK PARISI

www.offthemark.com



When I ask...

Where is the coffee??



I might actually mean...

I want to make some coffee.

When I ask...

Where is the coffee??



I might actually mean...

Would you make some coffee?

When I ask...

Where is the coffee??



I might actually mean...

Where did you put my coffee?

Look at the bigger picture

c.f restaurant scripts

Thematic Roles

Think of context as a scene.

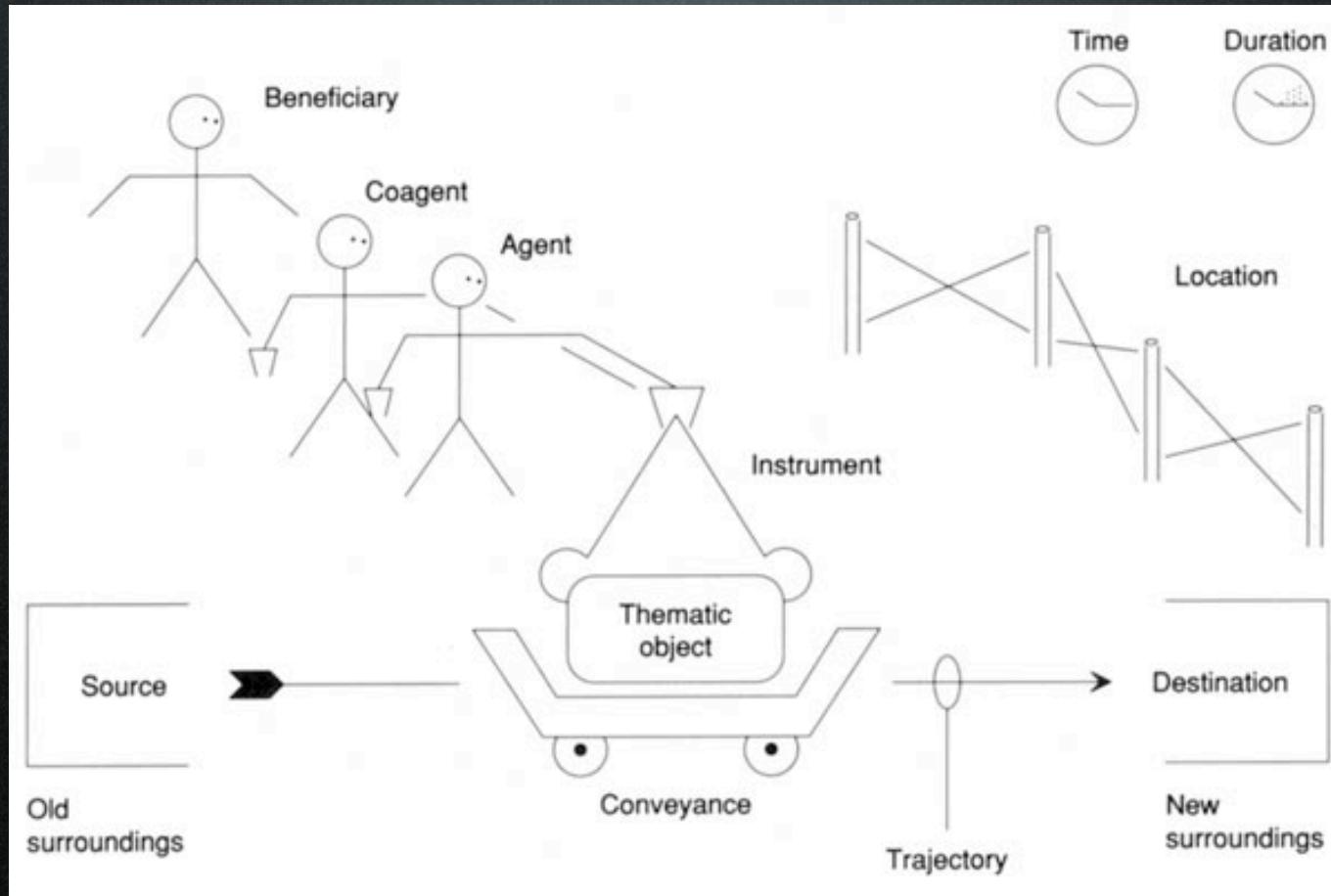
Parts of the scene represented by thematic roles.

Exact number of thematic roles open to debate.

Example Roles

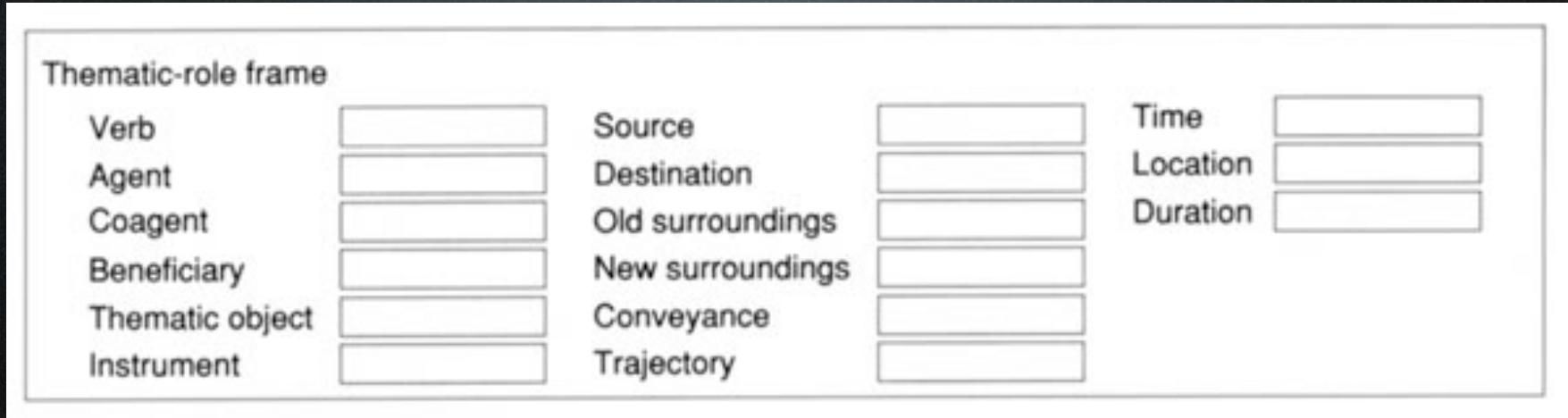
- Agent - causes action to occur
- Coagent - assists causing action
- Beneficiary - for whom action is performed
- Thematic object - thing the sentence is about
- Instrument - tool used by agent

Thematic Roles



[Winston 1992]

Thematic Roles



could be implemented as a frame system
roles filled by sentence analysis
carried as context to next sentence
used to build a script, answer questions