



The
University
Of
Sheffield.

Natural Language Processing for the Semantic Web

Isabelle Augenstein

Department of Computer Science, University of Sheffield, UK

i.augenstein@sheffield.ac.uk

Linked Data for NLP Tutorial, ESWC Summer School 2014

August 24, 2014



The British Museum

- Visiting
- > What's on
- Explore
- Research
- Learning
- About us
- Membership
- Support us
- Channel
- Blog
- Shop

Search the website

What's on > Exhibitions > Ming >

[Tickets](#) [Visiting](#) [Highlight objects](#) [Supporters](#) [Schools](#) [Events](#) [Restaurant packages](#)

English 中文

The BP exhibition
Ming
50 years that changed China

18 September 2014 –
5 January 2015
[Tickets on sale](#)

[Book online >](#)

Members free / Open late Fridays

This major exhibition will explore a golden age in China's history.

Between AD 1400 and 1450, China was a global superpower run by one family – the Ming dynasty – who established Beijing as the capital and built the Forbidden City. During this period, Ming China was thoroughly connected with the outside world. Chinese artists absorbed many fascinating influences, and created some of the most beautiful objects and paintings ever made.

The exhibition will feature a range of these spectacular objects – including exquisite porcelain, gold, jewellery, furniture, paintings, sculptures and textiles – from museums across China and the rest of the world. Many of them have only very recently been discovered and have never been seen outside China.



The British Museum

Visiting
> What's on
Explore
Research
Learning
About us

Membership
Support us
Channel
Blog
Shop

Search the website

semi-structured information

What's on > Exhibitions > Ming >

Tickets Visiting Highlight objects Supporters Schools Events Restaurant packages

English 中文

unstructured information

The BP exhibition
Ming
50 years that changed China

18 September 2014 –
5 January 2015
Tickets on sale

Book online >

Members free / Open late Friday

This major exhibition will explore a golden age in China's history.

Between AD 1400 and 1450, China was a global superpower run by one family – the Ming dynasty – who established Beijing as the capital and built the Forbidden City. During this period, Ming China was thoroughly connected with the outside world. Chinese artists absorbed many fascinating influences, and created some of the most beautiful objects and paintings ever made.

The exhibition will feature a range of these spectacular objects – including exquisite porcelain, gold, jewellery, furniture, paintings, sculptures and textiles – from museums across China and the rest of the world. Many of them have only very recently been discovered and have never been seen outside China.



The British Museum

- Visiting
- Membership
- Support us
- Channel
- Blog
- Shop

Search the website

What's on > Exhibitions > Ming >

Tickets Visiting Highlight objects Supporters Schools Events Restaurant packages

English 中文

semi-structured information

unstructured information

How to link this information to a knowledge base automatically?

The BP exhibition
Ming
50 years that changed China

18 September 2014 –
5 January 2015
Tickets on sale

Book online >

Members free / Open late Friday

This major exhibition will explore a golden age in China's history. Between AD 1400 and 1450, China was a global superpower run by one family – the Ming dynasty – who established Beijing as the capital and built the Forbidden City. During this period, Ming China was thoroughly connected with the outside world. Chinese artists absorbed many fascinating influences, and created some of the most beautiful objects and paintings ever made.

The exhibition will feature a range of these spectacular objects – including exquisite porcelain, gold, jewellery, furniture, paintings, sculptures and textiles – from museums across China and the rest of the world. Many of them have only very recently been discovered and have never been seen outside China.



The British Museum

- Visiting
- > What's on
- Explore
- Research
- Learning
- About us
- Membership
- Support us
- Channel
- Blog
- Shop

Search the website

semi-structured information

What's on > Exhibitions > Ming >

Tickets Visiting Highlight objects Supporters Schools Events Restaurant packages

English 中文

unstructured information

The BP exhibition
Ming
50 years that
changed China

18 September 2014 –
5 January 2015

Tickets on sale

Book online >

Members free / Open late Friday

This major exhibition will explore a golden age in China's history.

Between AD 1400 and 1450, China was a global superpower run by one family – the Ming dynasty – who established Beijing as the capital and built the Forbidden City. During this period, Ming China was thoroughly connected with the outside world. Chinese artists absorbed many fascinating influences, and created some of the most beautiful objects and paintings ever made.

The exhibition will feature a range of these spectacular objects – including exquisite porcelains, jade, jewellery, furniture, paintings, sculptures and textiles – from museums across China and the rest of the world. Many of them have only very recently been discovered and have never been seen outside China.

How to link this information to a knowledge base automatically?

Information Extraction!



The
University
Of
Sheffield.

Information Extraction

6

Between AD 1400 and 1450, China was a global superpower run by one family – the Ming dynasty – who established Beijing as the capital and built the Forbidden City.

Between AD 1400 and 1450, **China** was a global superpower run by one family – the **Ming dynasty** – who established **Beijing** as the capital and built the **Forbidden City**.

Named Entity Recognition

Between AD 1400 and 1450, **China** was a global superpower run by one family – the **Ming dynasty** – who established Beijing as the capital and built the **Forbidden City**.

Named Entity Recognition (NER)

Named Entity Classification (NEC):

China: /location/country

Ming dynasty: /royalty/royal_line

Beijing: /location/city

Forbidden City: /location/city

Between AD 1400 and 1450, **China** was a global superpower run by one family – the **Ming dynasty** – who established Beijing as the capital and built the **Forbidden City**.

Named Entity Recognition

Named Entity Classification:

China: /location/country

Ming dynasty: /royalty/royal_line

Beijing: /location/city

Forbidden City: /location/city

Named Entity Linking:

China: /m/0d05w3

Ming dynasty: /m/0bw_m

Beijing: /m/01914

Forbidden City: /m/0j0b2

Named Entities: Proper nouns, which refer to real-life entities

Named Entity Recognition: Detecting boundaries of named entities (NEs)

Named Entity Classification: Assigning classes to NEs, such as PERSON, LOCATION, ORGANISATION, or fine-grained classes such as ROYAL LINE

Named Entity Linking / Disambiguation: Linking NEs to concrete entries in knowledge base, example:

China -> LOCATION: Republic of China, country in East Asia

-> LOCATION: China proper, core region of China during Qing dynasty

-> LOCATION: China, Texas

-> PERSON: China, Brazilian footballer born in 1964

-> MUSIC: China, a 1979 album by Vangelis

-> ...



Relations

Between AD 1400 and 1450, **China** was a global
superpower run by one family – the **Ming dynasty** – who
established **Beijing** as the capital and built the **Forbidden City**.

/location/country/capital (arrow from Beijing to China)

/royalty/royal_line/kingdom_s_ruled (arrow from Ming dynasty to China)

/location/location/containedby (arrow from Forbidden City to Beijing)

Named Entity Recognition

Relation Extraction



Relations and Time Expressions ¹²

Between AD 1400 and 1450, China was a global

1400-XX-XX -- 1450-XX-XX

/royalty/royal_line/kingdom_s_ruled

superpower run by one family – the Ming dynasty – who

/location/country/capital

/location/location/containedby

established Beijing as the capital and built the Forbidden City.

Named Entity Recognition

Relation Extraction

Temporal Extraction



Relations, Time Expressions and Events

Between AD 1400 and 1450, China was a global
 1400-XX-XX -- 1450-XX-XX
 superpower run by one family – the Ming dynasty – who
 /location/country/capital established Beijing as the capital and built the Forbidden City.
 /royalty/royal_line/kingdom_s_ruled
 /location/location/containedby

Named Entity Recognition

Event Extraction

Relation Extraction

Event: /royalty/royal_line: Ming dynasty

Temporal Extraction

/royalty/royal_line/ruled_from: 1400-XX-XX
 /royalty/royal_line/ruled_to: 1450-XX-XX
 /royalty/royal_line/kingdom_s_ruled: China

Relations: Two or more entities which relate to one another in real life

Relation Extraction: Detecting relations between entities and assigning relation types to them, such as CAPITAL-OF

Temporal Extraction: Recognising and normalising time expressions: times (e.g. “3 in the afternoon”), dates (“tomorrow”), durations (“since yesterday”), and sets (e.g. “twice a month”)

Events: Real-life events that happened at some point in space and time, e.g. kingdom, assassination, exhibition

Event Extraction: Extracting events consisting of the name and type of event, time and location

- Information extraction (IE) methods such as named entity recognition (NER), named entity classification (NEC), named entity linking, relation extraction (RE), temporal extraction, and event extraction can help to add markup to Web pages
- Information extraction approaches can serve two purposes:
 - Annotating every single mention of an entity, relation or event, e.g. to add markup to Web pages
 - Aggregating those mentions to populate knowledge bases, e.g. based on confidence values and majority voting

China	LOCATION	0.9
China	LOCATION	0.8
China	PERSON	0.4
→ China	LOCATION	

- Focus of the rest of the tutorial and hands-on session: Named entity recognition and classification (NERC)
- Possible methodologies
 - Rule-based approaches: write manual extraction rules
 - Machine learning based approaches
 - Supervised learning: manually annotate text, train machine learning model
 - Unsupervised learning: extract language patterns, cluster similar ones
 - Semi-supervised learning: start with a small number of language patterns, iteratively learn more (bootstrapping)
 - Gazetteer-based method: use existing list of named entities
 - Combination of the above



The
University
Of
Sheffield.

Information Extraction: Methods ¹⁷

Developing a NERC involves programming based around APIs..



Developing a NERC involves programming based around APIs..

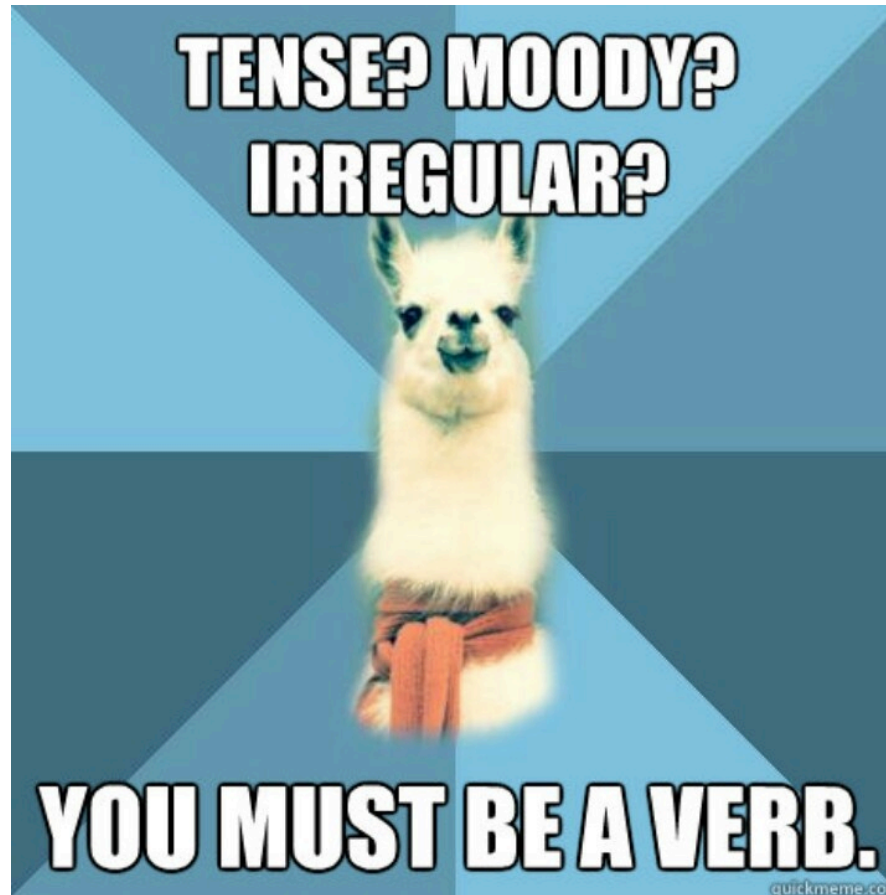
NEVER HAVE I FELT SO
CLOSE TO ANOTHER SOUL
AND YET SO HELPLESSLY ALONE
AS WHEN I GOOGLE AN ERROR
AND THERE'S ONE RESULT
A THREAD BY SOMEONE
WITH THE SAME PROBLEM
AND NO ANSWER
LAST POSTED TO IN 2003

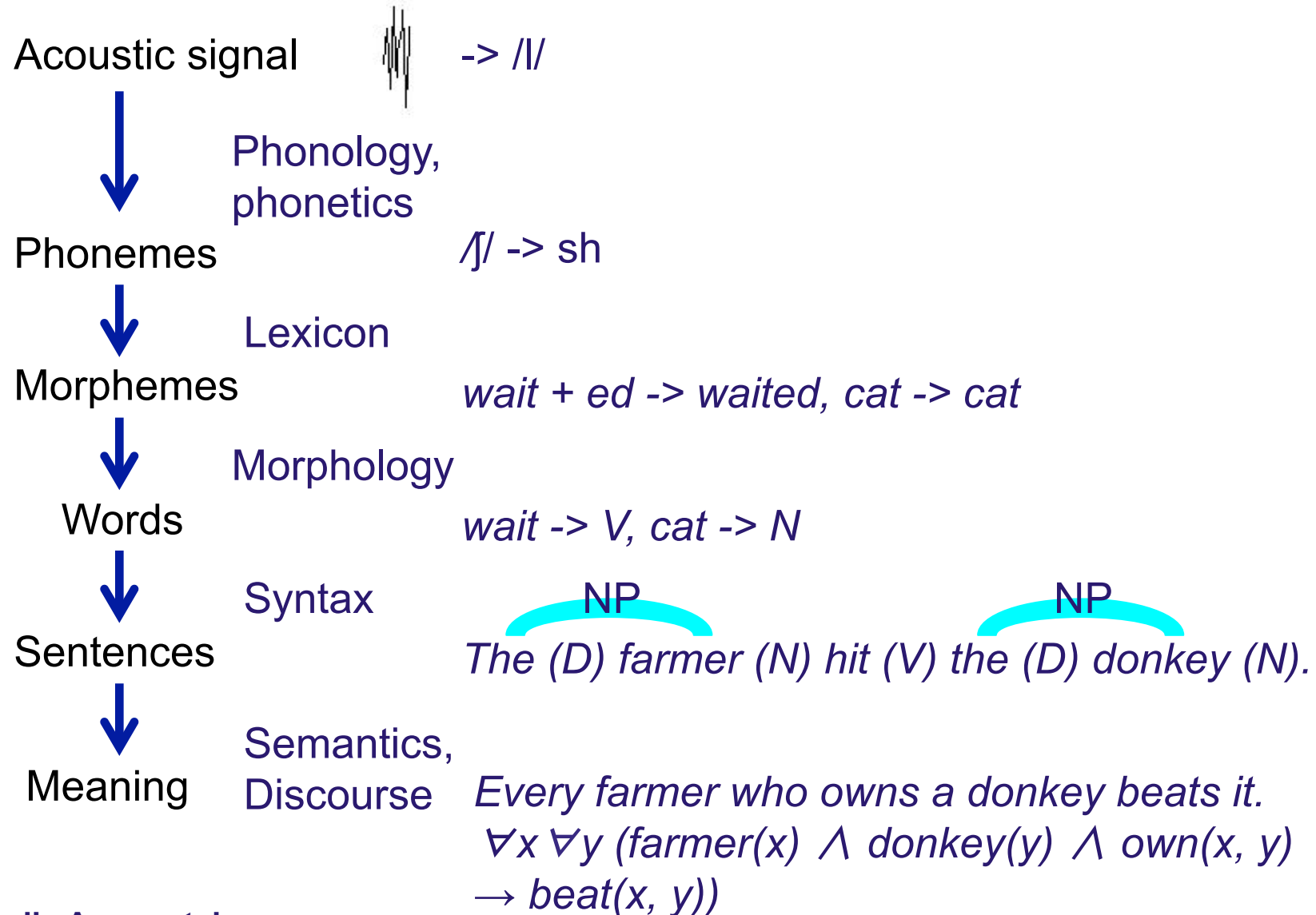


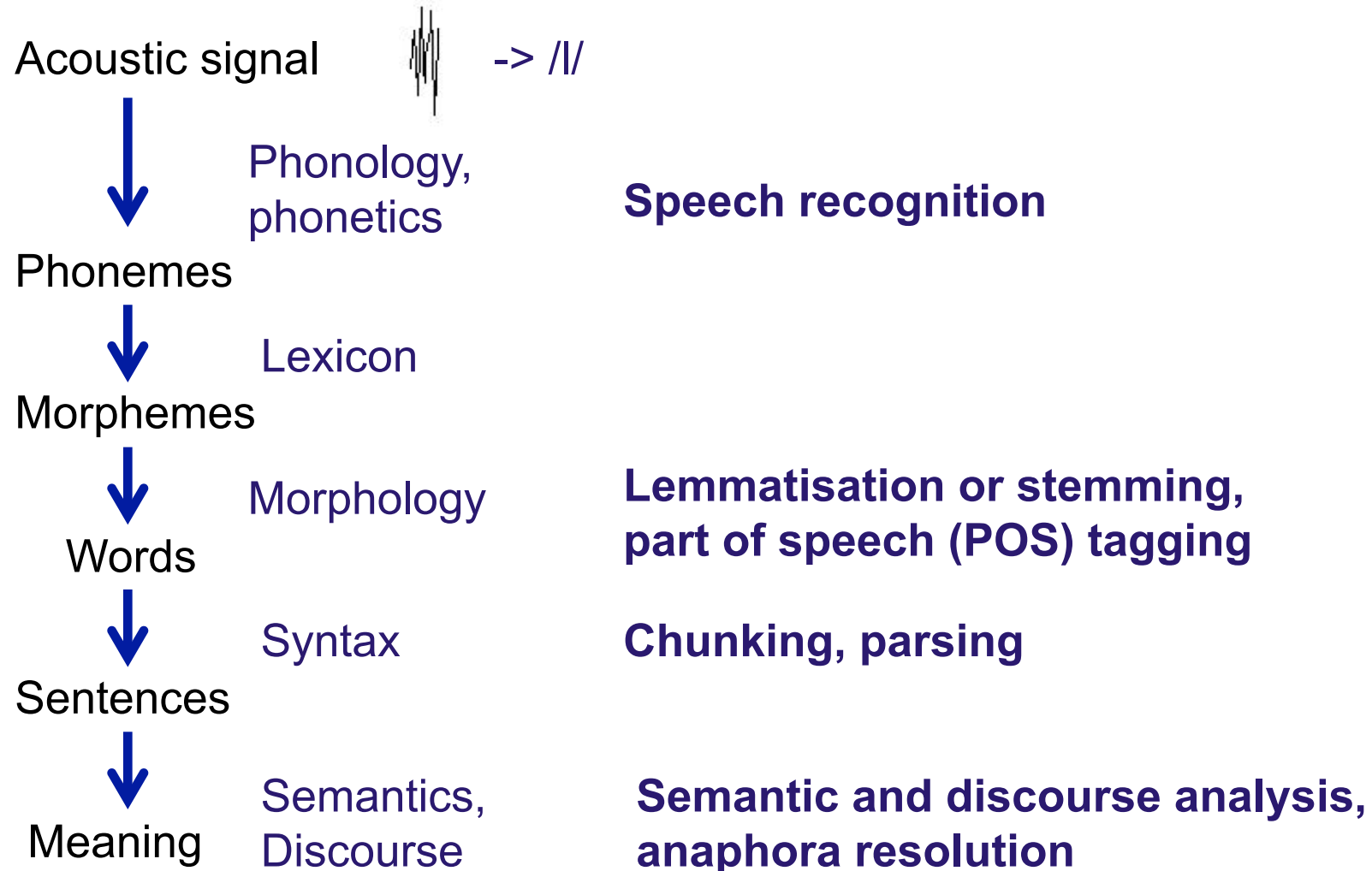
which can be frustrating at times

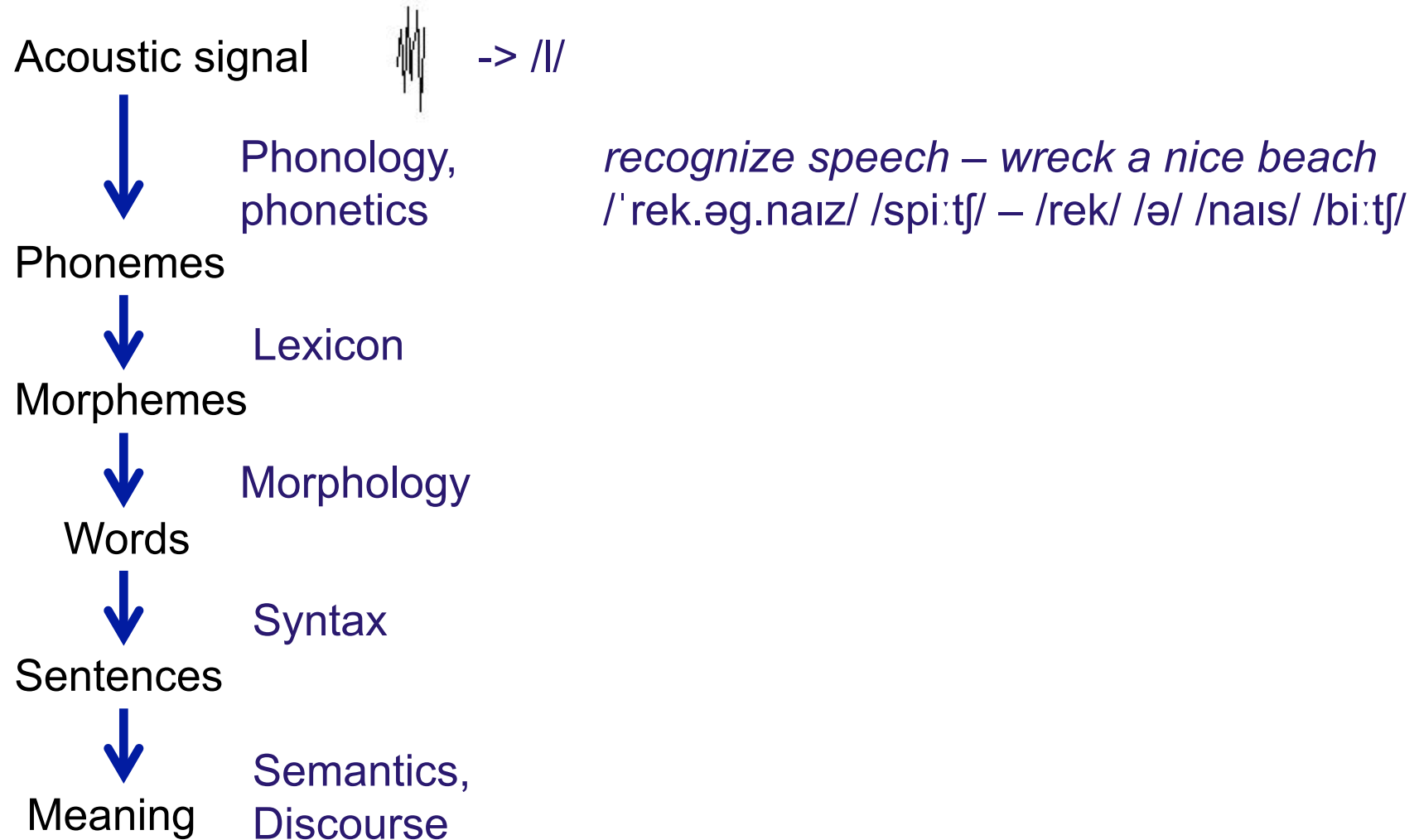


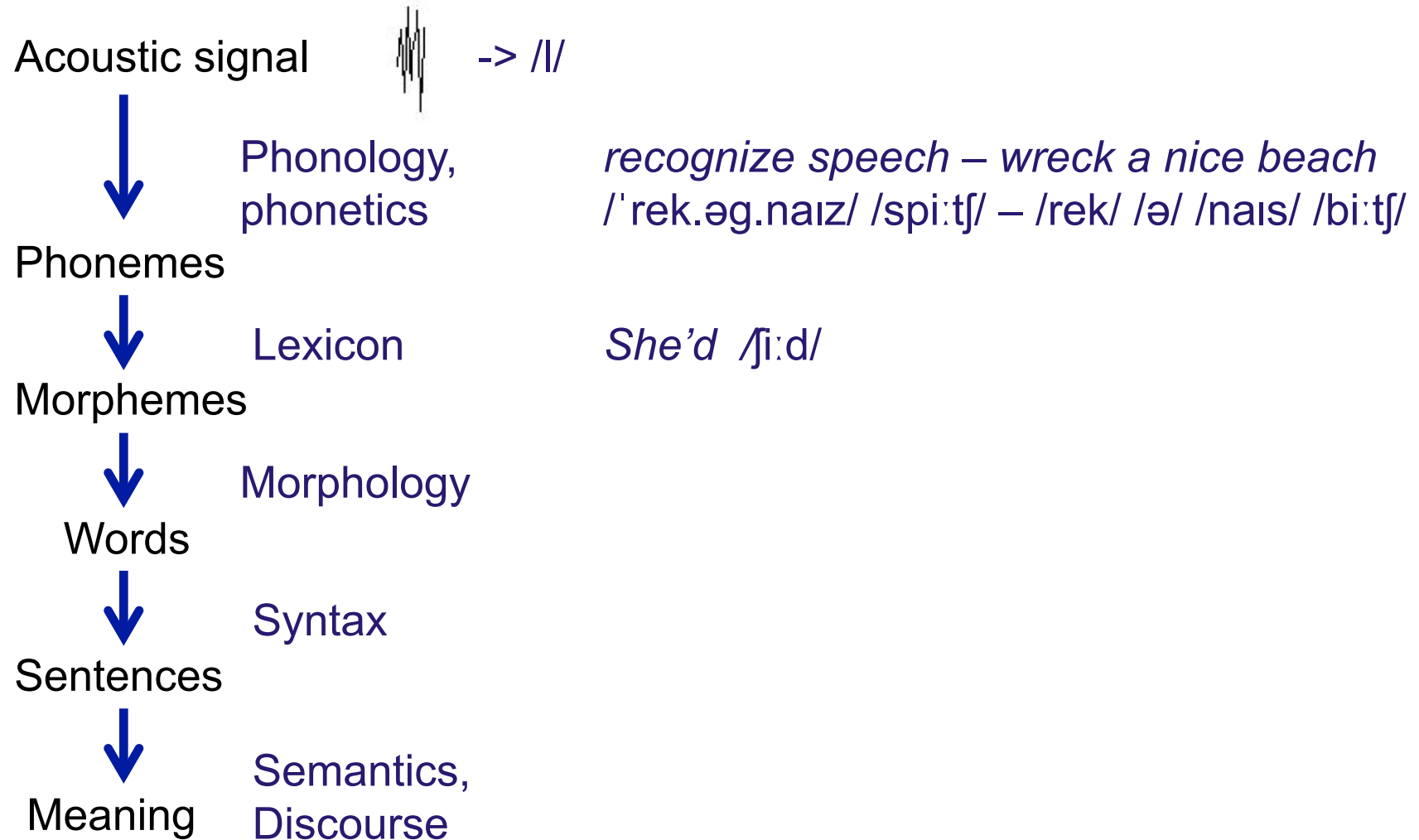
and (at least basic) knowledge about linguistics

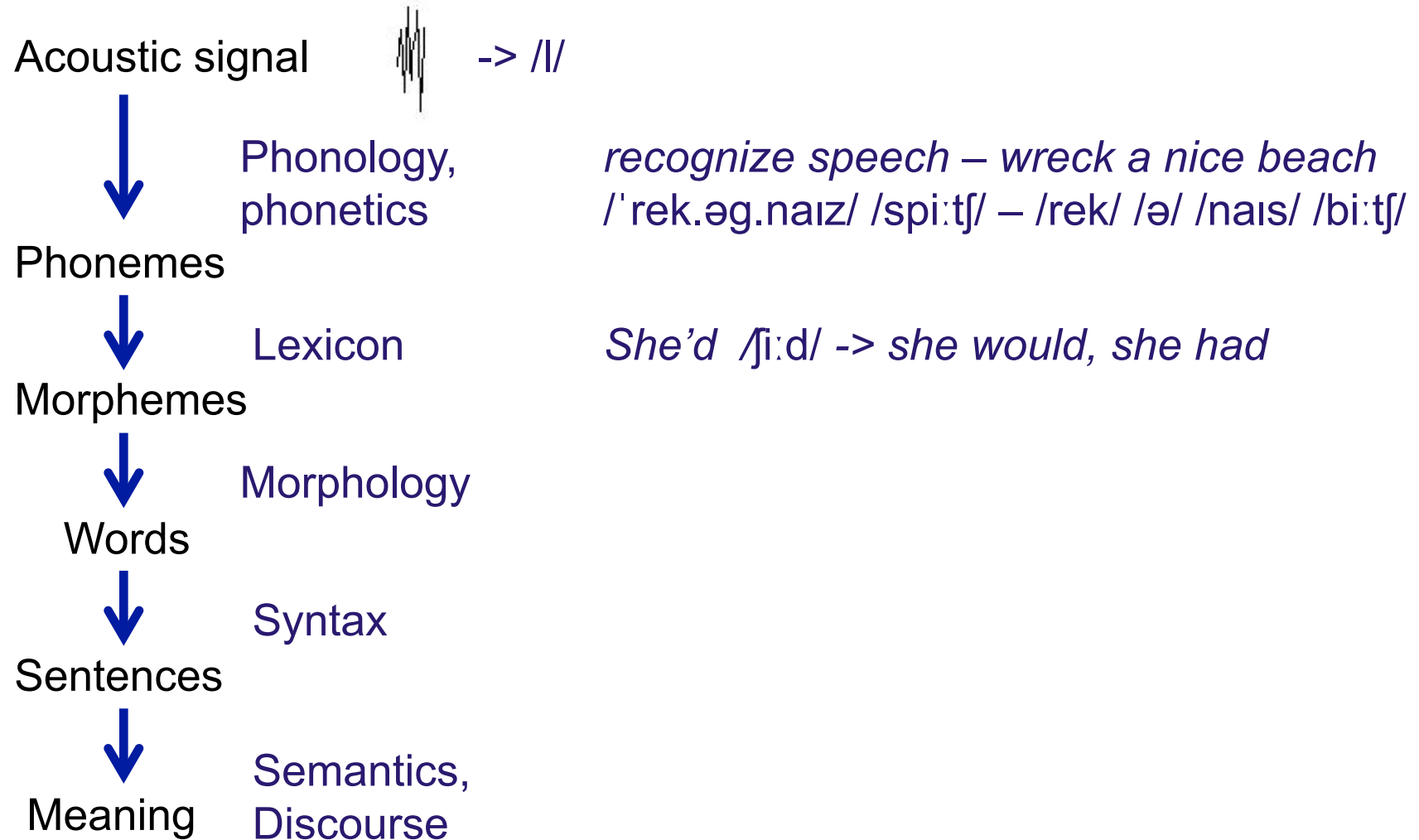


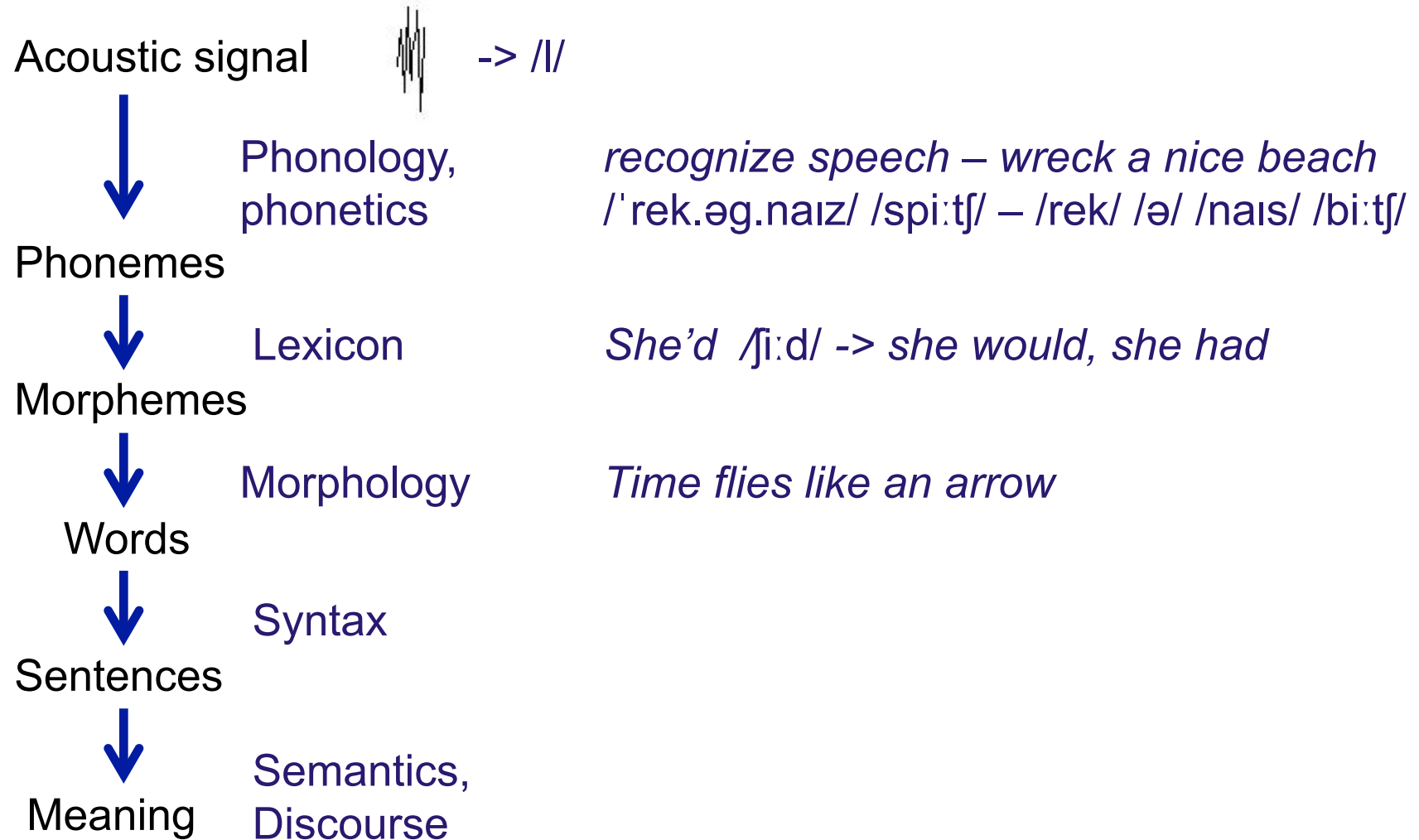


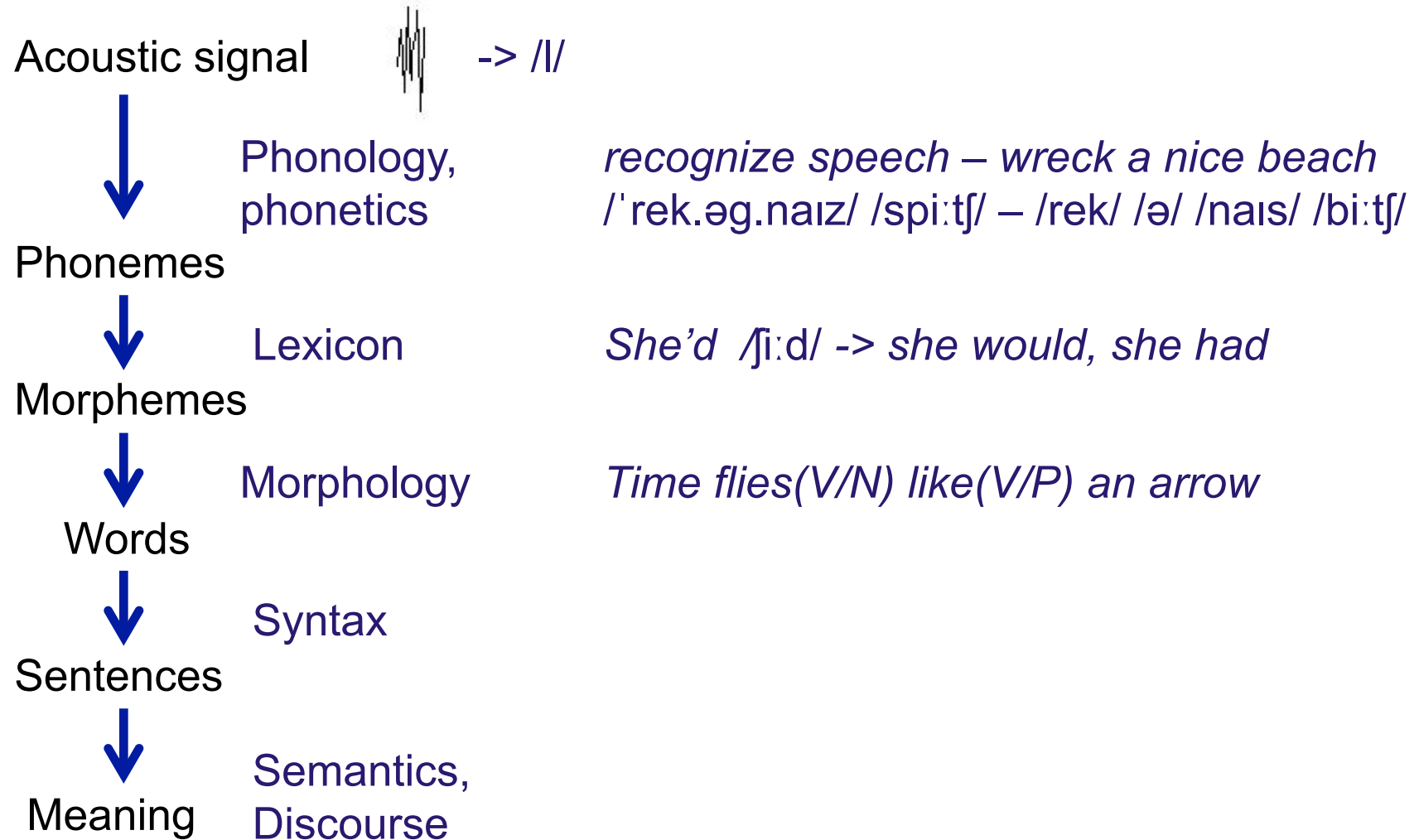


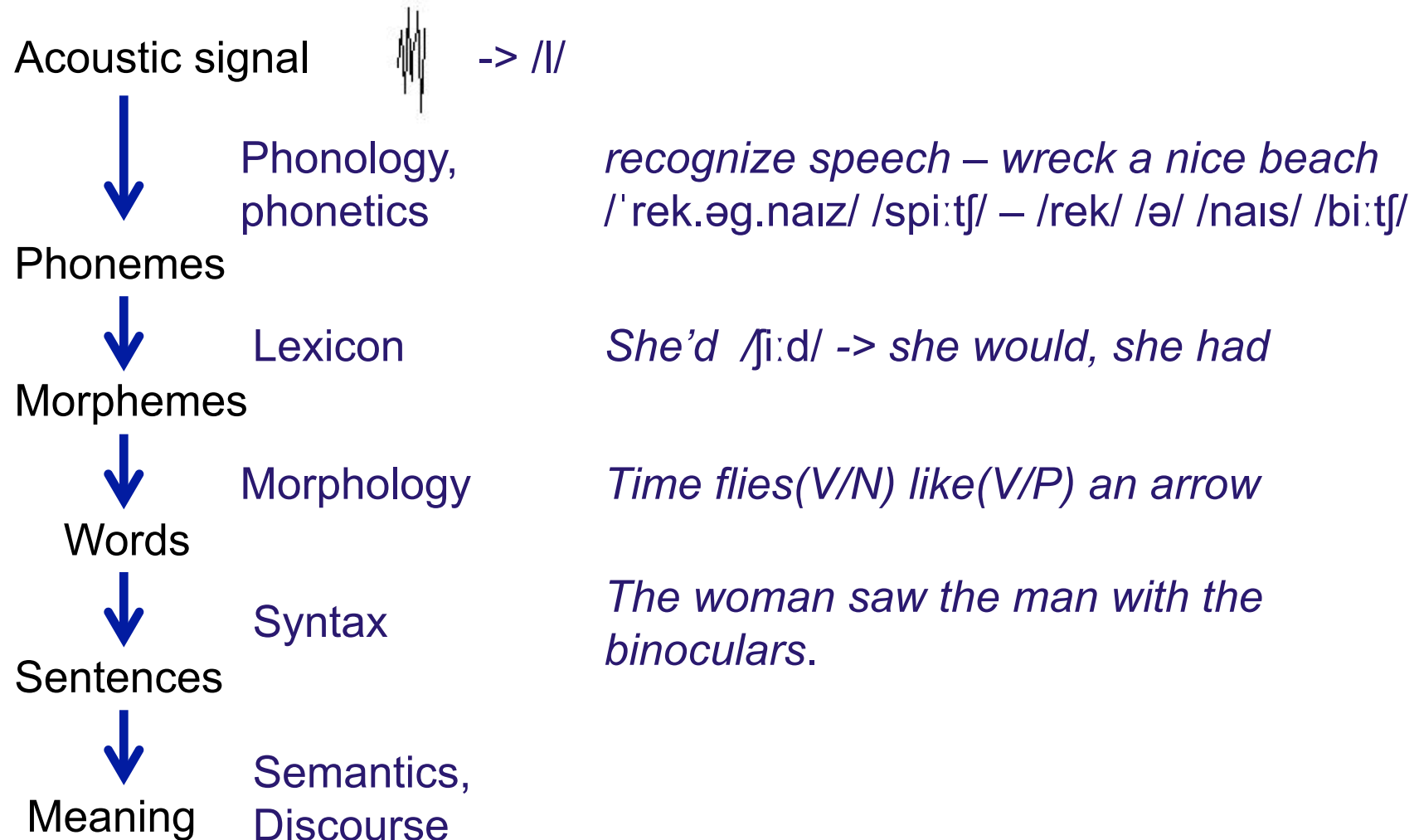


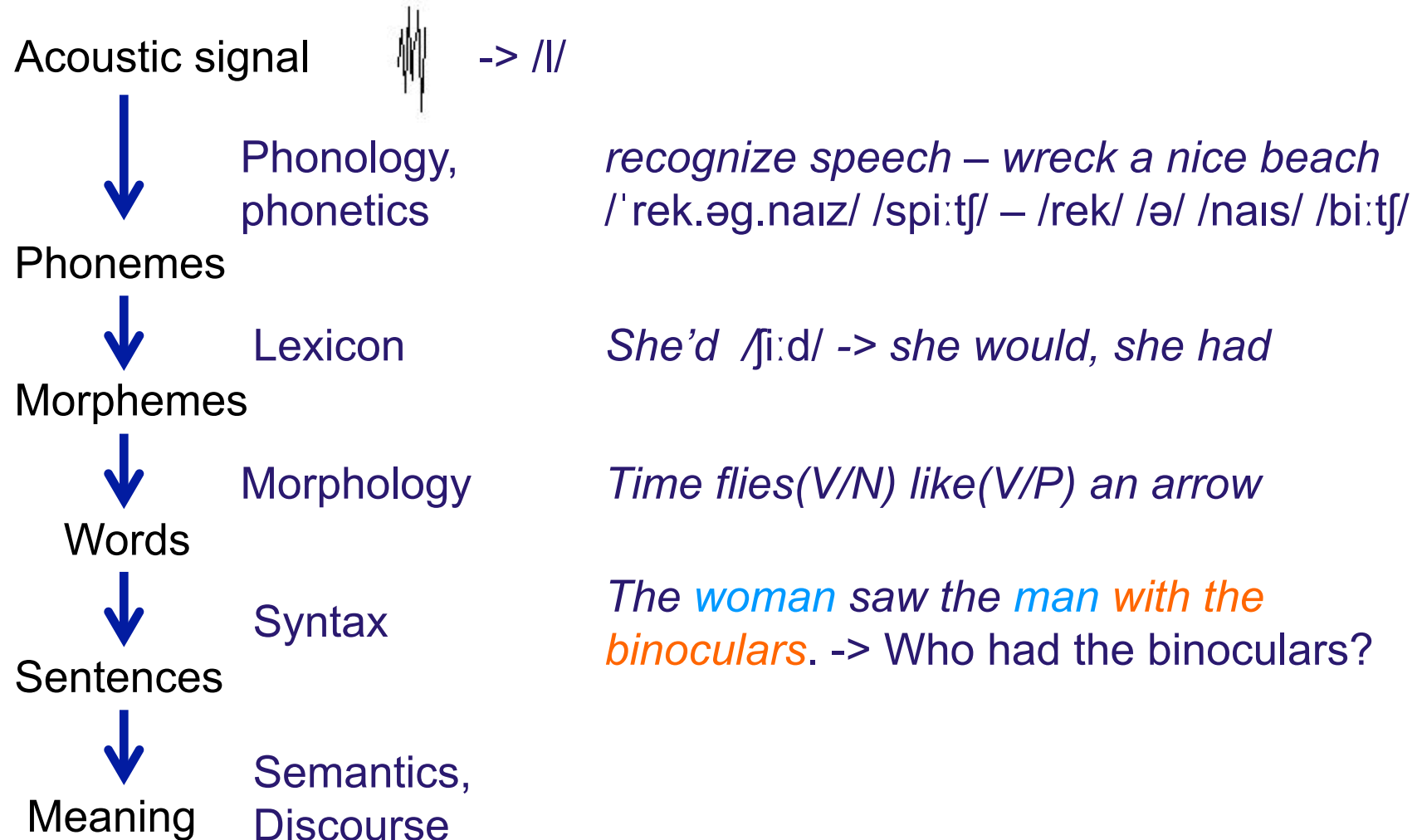


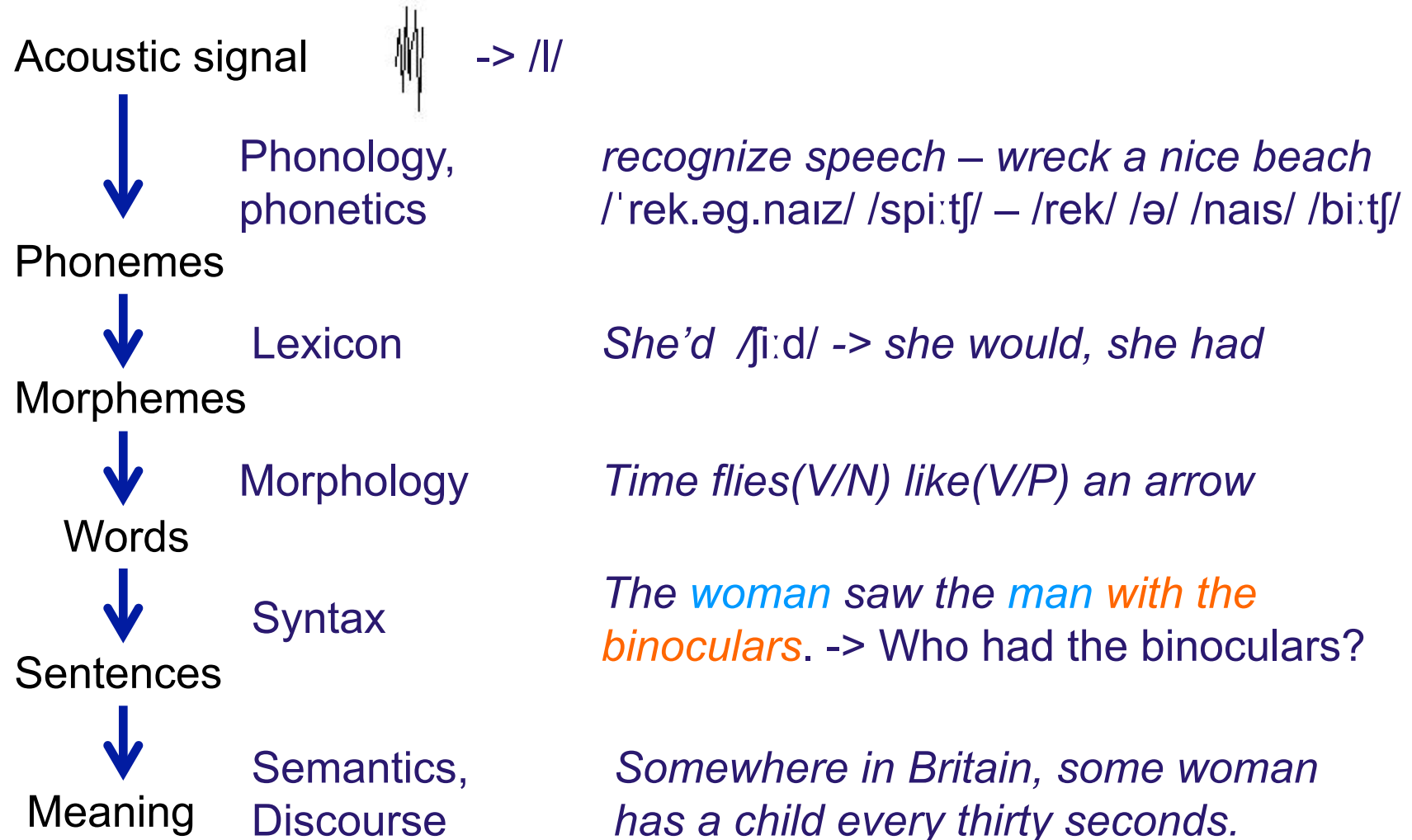


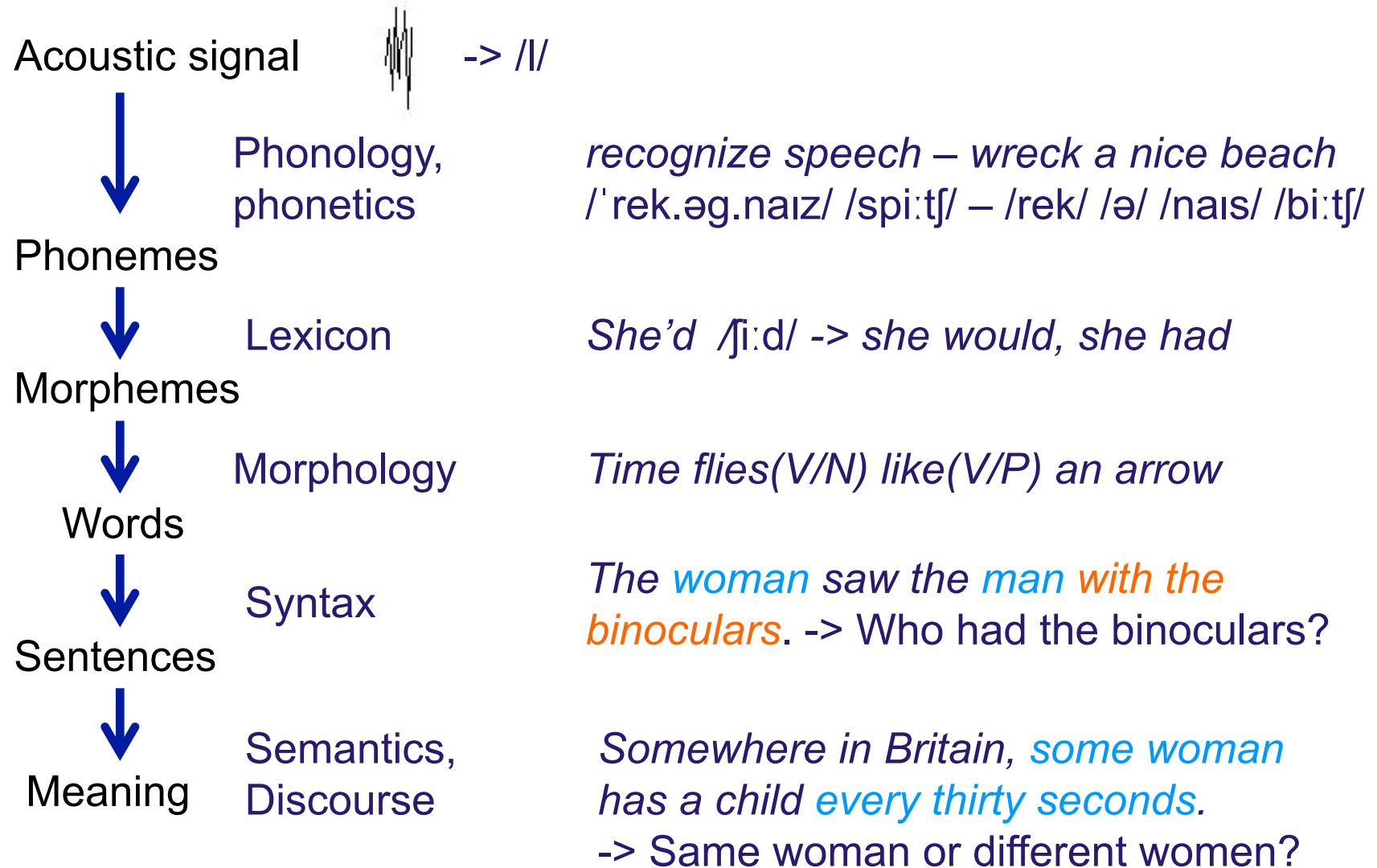


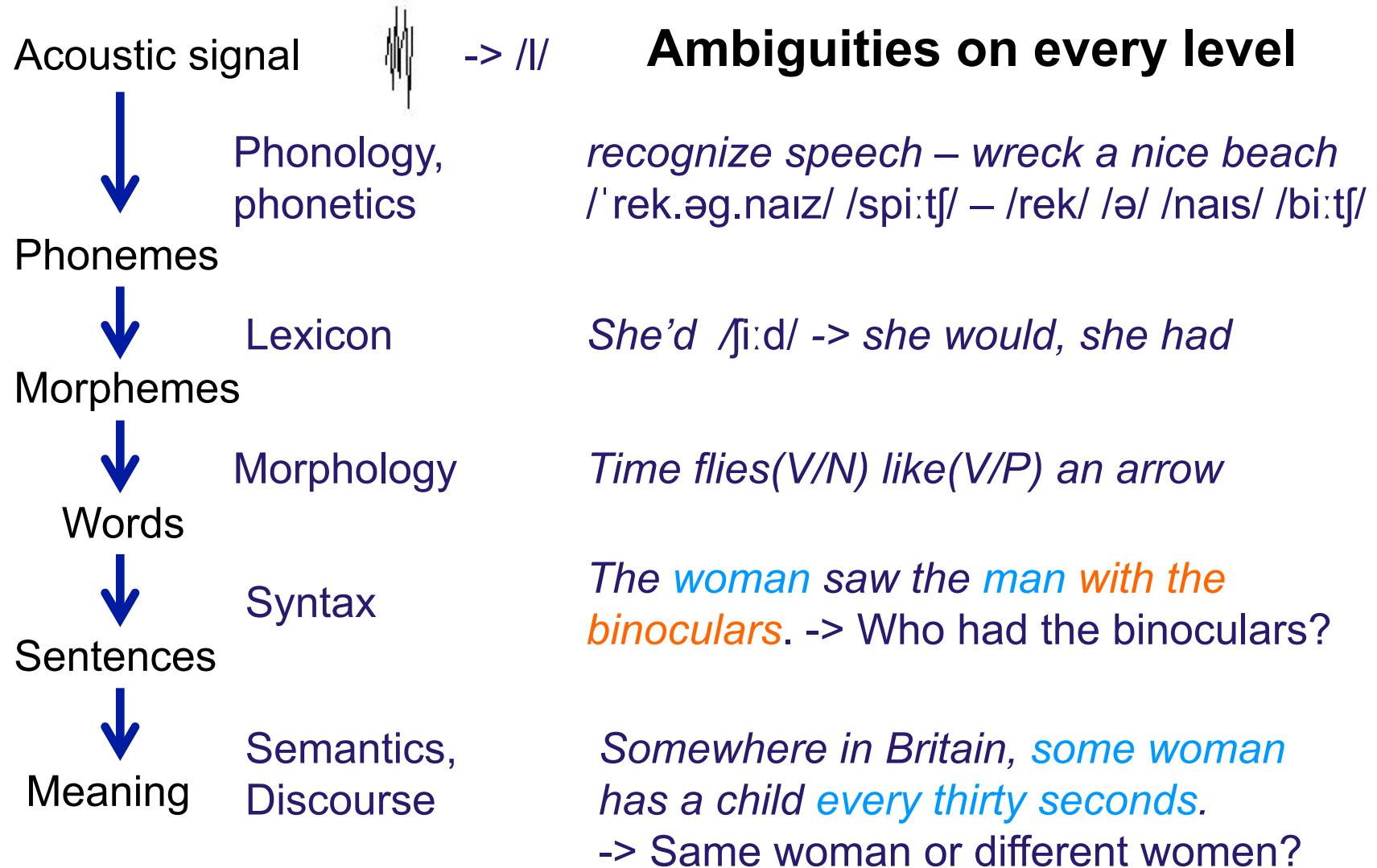














Acoustic signal



-> //

Ambiguities on every level



Phonemes



Morpheme



Words

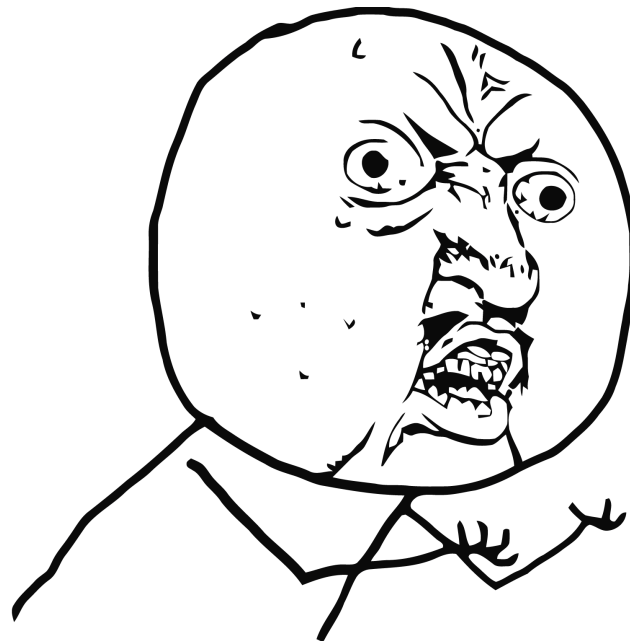


Sentences



Meaning

Y U SO



AMBIGUOUS?

vreck a nice beach
- /rek/ /ə/ /nais/ /bi:tʃ/

ould, she had

/P) an arrow

man with the
ad the binoculars?

n, some woman
ty seconds.
lifferent women?

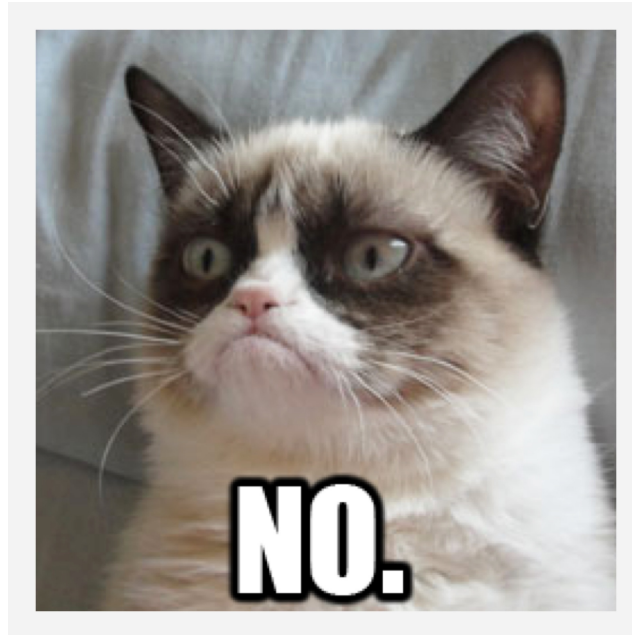
Language is ambiguous..

Can we still build named entity extractors that extract all entities from unseen text correctly?



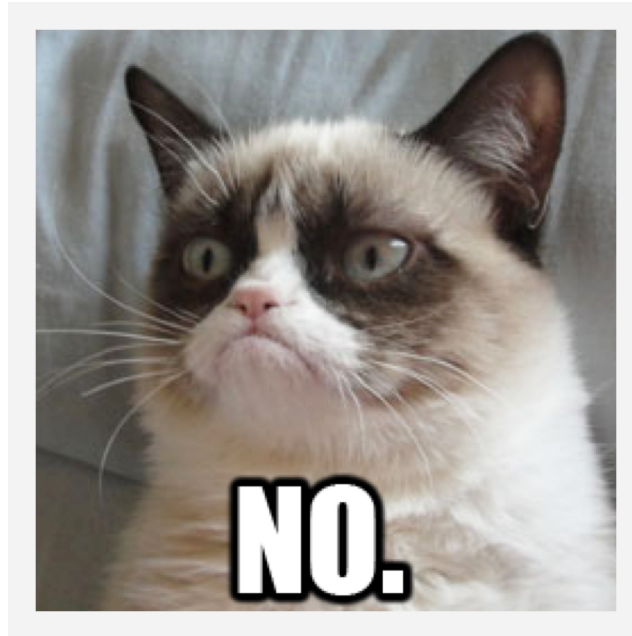
Language is ambiguous..

Can we still build named entity extractors that extract all entities from unseen text correctly?



Language is ambiguous..

Can we still build named entity extractors that extract all entities from unseen text correctly?



However, we can try to extract most of them correctly
using linguistic cues and background knowledge!

What can help to recognise and/or classify named entities?

- Words:
 - Words in window before and after mention
 - Sequences
 - Bags of words

Between AD 1400 and 1450, **China** was a global superpower

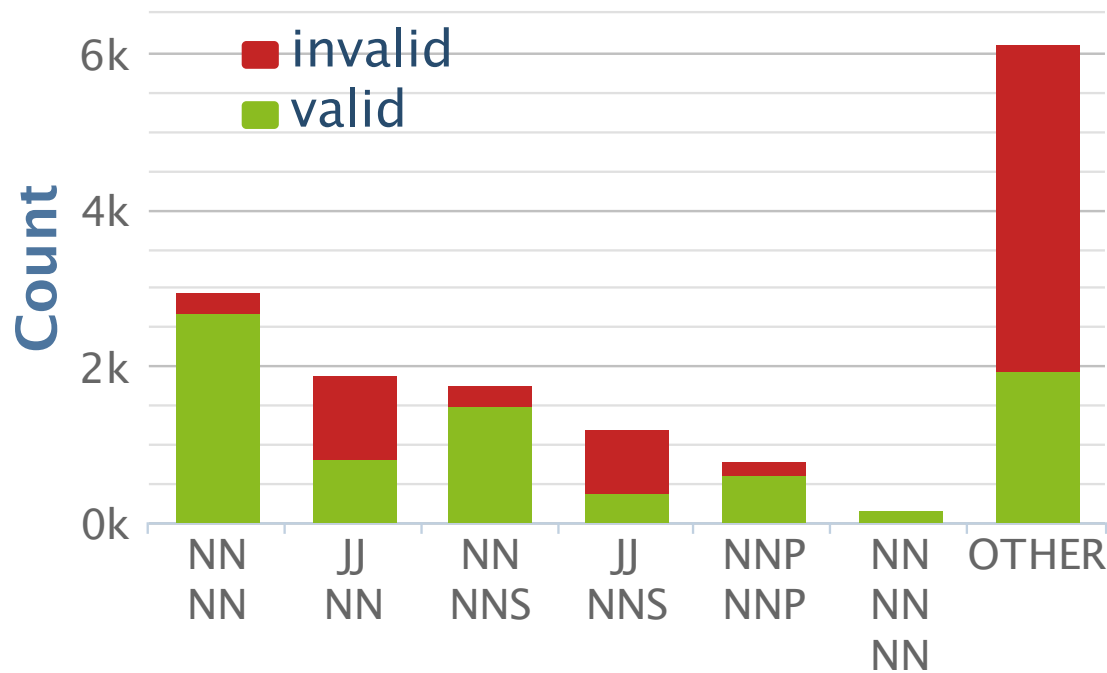
w: China w-1: , w-2: 1450 w+1: was w+2: a
seq[-]: 1450, seq[+]: was a
bow: China bow[-]: , bow[-]: 1450 bow[+]: was bow[+]: a

What can help to recognise and/or classify named entities?

- Morphology:
 - Capitalisation: is upper case (*China*), all upper case (*IBM*), mixed case (*eBay*)
 - Symbols: contains \$, £, €, roman symbols (*IV*), ..
 - Contains period (*google.com*), apostrophe (*Mandy's*), hyphen (*speed-o-meter*), ampersand (*Fisher & Sons*)
 - Stem or Lemma (*cats -> cat*), prefix (*disadvantages -> dis*), suffix (*cats -> s*), interfix (*speed-o-meter -> o*)

What can help to recognise and/or classify named entities?

- POS (part of speech) tags
 - Most named entities are nouns



Prokofyev (2014)

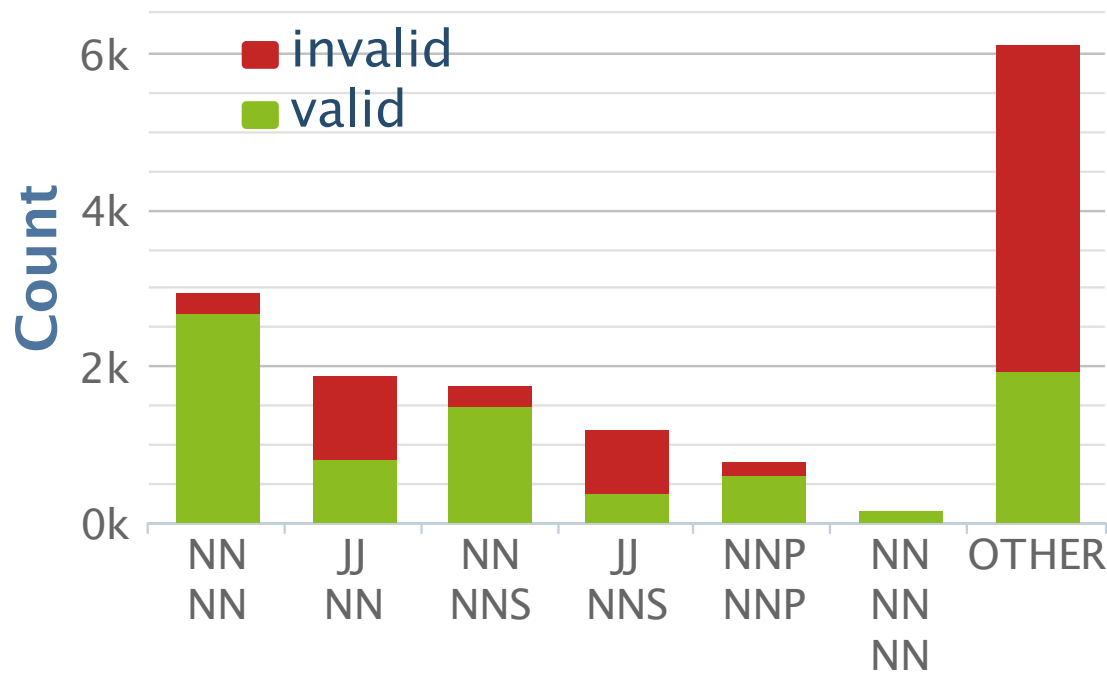
Number	Tag	Description
1.	CC	Coordinating conjunction
2.	CD	Cardinal number
3.	DT	Determiner
4.	EX	Existential <i>there</i>
5.	FW	Foreign word
6.	IN	Preposition or subordinating conjunction
7.	JJ	Adjective
8.	JJR	Adjective, comparative
9.	JJS	Adjective, superlative
10.	LS	List item marker
11.	MD	Modal
12.	NN	Noun, singular or mass
13.	NNS	Noun, plural
14.	NNP	Proper noun, singular
15.	NNPS	Proper noun, plural
16.	PDT	Predeterminer
17.	POS	Possessive ending
18.	PRP	Personal pronoun
19.	PRP\$	Possessive pronoun
20.	RB	Adverb
21.	RBR	Adverb, comparative
22.	RBS	Adverb, superlative
23.	RP	Particle
24.	SYM	Symbol
25.	TO	<i>to</i>
26.	UH	Interjection
27.	VB	Verb, base form
28.	VBD	Verb, past tense
29.	VBG	Verb, gerund or present participle
30.	VBN	Verb, past participle
31.	VBP	Verb, non-3rd person singular present
32.	VBZ	Verb, 3rd person singular present
33.	WDT	Wh-determiner
34.	WP	Wh-pronoun
35.	WP\$	Possessive wh-pronoun
36.	WRB	Wh-adverb

Number	Tag	Description			
1.	CC	Coordinating conjunction		18.	PRP Personal pronoun
2.	CD	Cardinal number		19.	PRP\$ Possessive pronoun
3.	DT	Determiner	Adjectives (all start with J)	20.	RB Adverb
4.	EX	Existential <i>there</i>		21.	RBR Adverb, comparative
5.	FW	Foreign word		22.	RBS Adverb, superlative
6.	IN	Preposition or subordinating conjunction		23.	RP Particle
7.	JJ	Adjective		24.	SYM Symbol
8.	JJR	Adjective, comparative		25.	TO <i>to</i>
9.	JJS	Adjective, superlative		26.	UH Interjection
10.	LS	List item marker		27.	VB Verb, base form
11.	MD	Modal		28.	VBD Verb, past tense
12.	NN	Noun, singular or mass	Nouns (all start with N)	29.	VBG Verb, gerund or present participle
13.	NNS	Noun, plural		30.	VBN Verb, past participle
14.	NNP	Proper noun, singular		31.	VBP Verb, non-3rd person singular present
15.	NNPS	Proper noun, plural		32.	VBZ Verb, 3rd person singular present
16.	PDT	Predeterminer		33.	WDT Wh-determiner
17.	POS	Possessive ending		34.	WP Wh-pronoun
				35.	WP\$ Possessive wh-pronoun
				36.	WRB Wh-adverb

**Verbs (all
start with V)**

What can help to recognise and/or classify named entities?

- POS (part of speech) tags
 - Most named entities are nouns



Prokofyev (2014)

What can help to recognise and/or classify named entities?

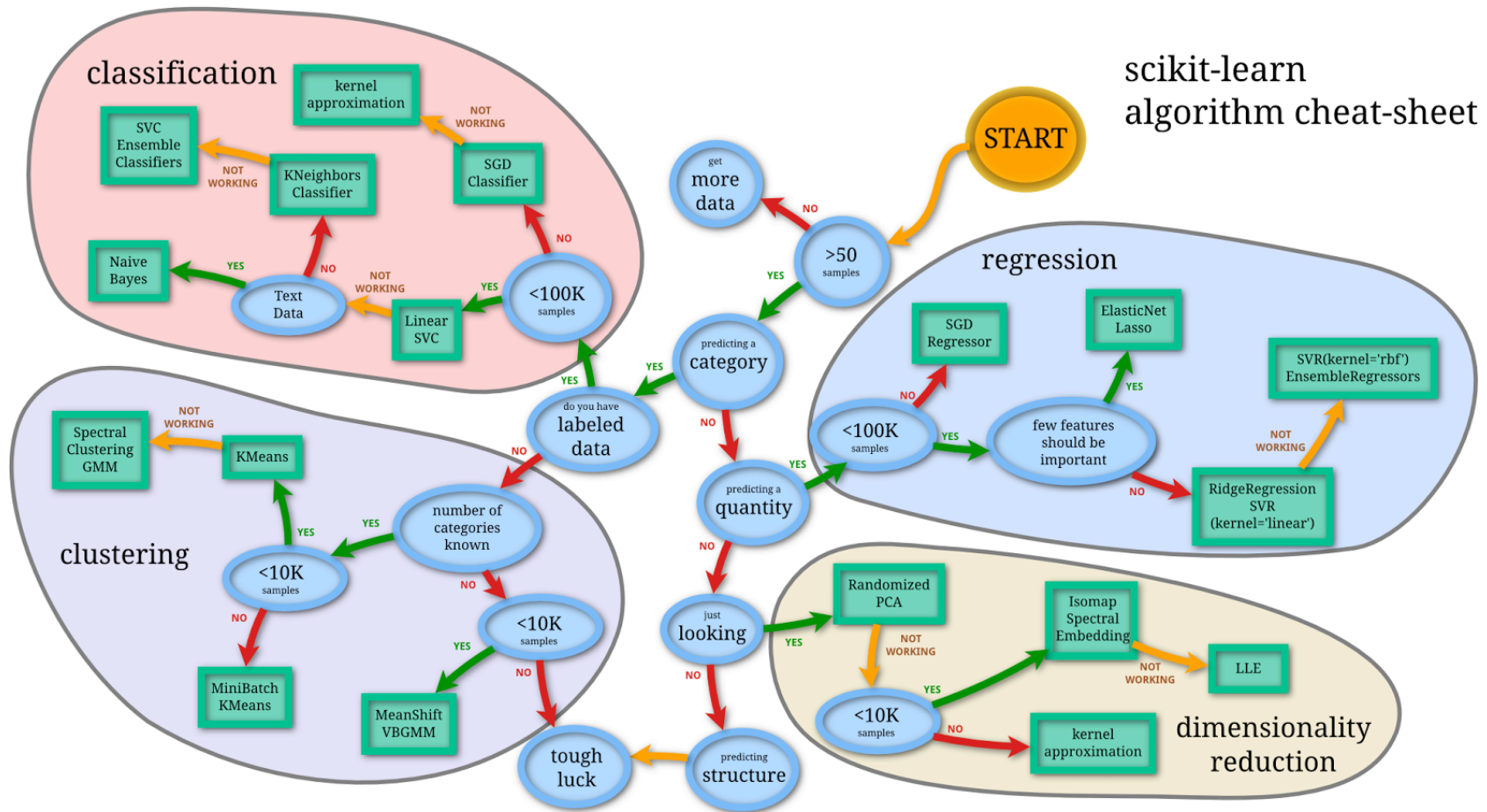
- Gazetteers
 - Retrieved from HTML lists or tables

Religion

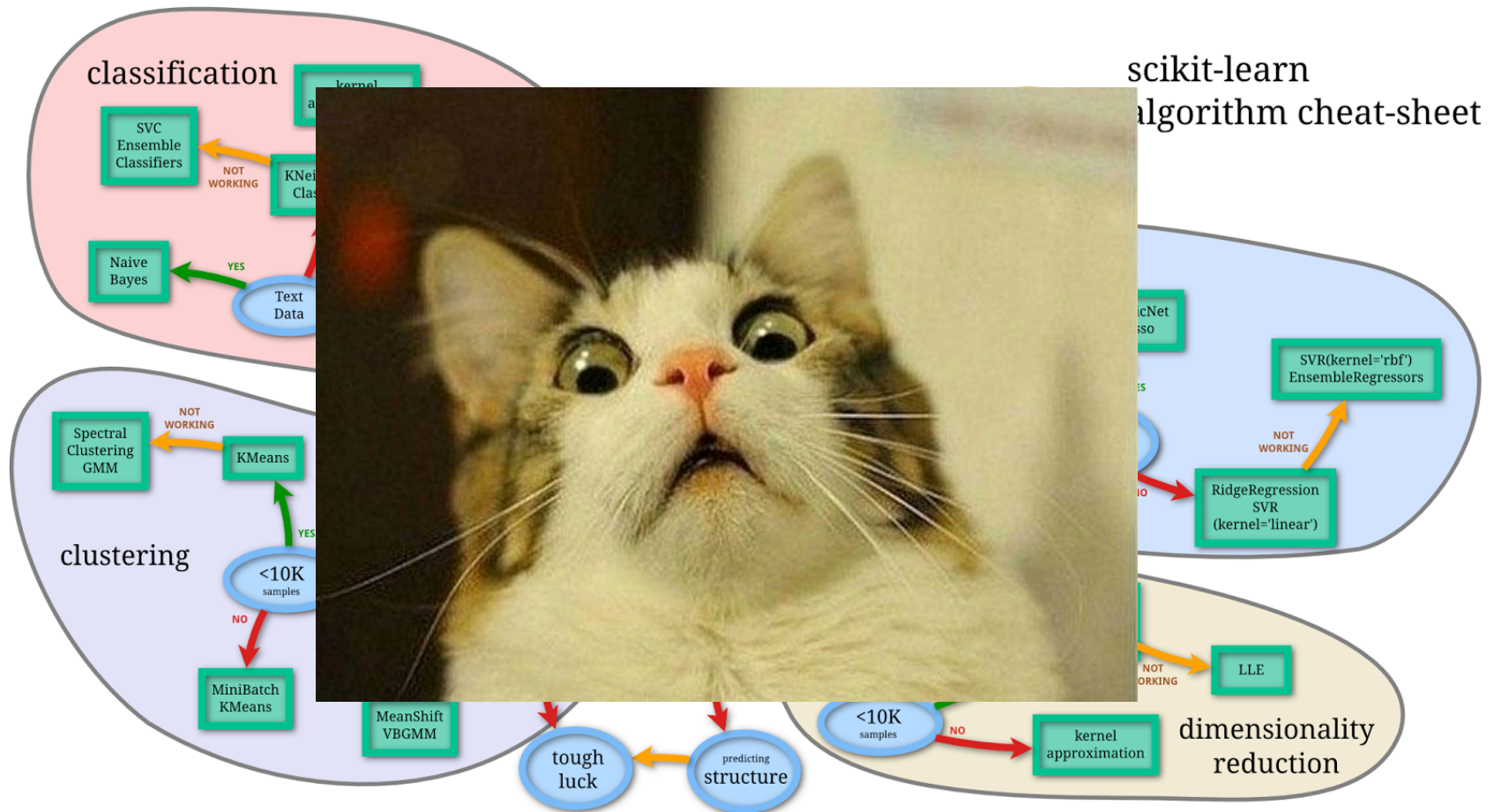
Heaven worship,
Taoism, Confucianism,
Buddhism, Chinese
folk religion, Islam

- Using regular expressions patterns and search engines (e.g. “*Religions such as **”)
- Retrieved from knowledge bases

Extensive choice of machine learning algorithms for training NERCs



Extensive choice of machine learning algorithms for training NERCs



- Unfortunately, there isn't enough time to explain machine learning algorithms in detail
- **CRFs** (conditional random fields) are one of the most widely used algorithms for NERC
 - Graphical models, view NERC as a sequence labelling task
 - Named entities consist of a beginning token (*B*), inside tokens (*I*), and outside tokens (*O*)
China (*B-LOC*) built (*O*) the (*O*) Forbidden (*B-LOC*) City (*I-LOC*) . (*O*)
- For now, we will focus on **rule- and gazetteer-based** NERC
- It is fairly easy to write manual extraction rules for NEs, can achieve a high performance when combined with gazetteers
 - This can be done with the GATE software (general architecture for text engineering) and Jape rules
-> **Hands-on session**

Natural Language Processing:

- GATE (general purpose architecture, includes other NLP and ML software as plugins)
- Stanford NLP (Java)
- OpenNLP (Java)
- NLTK (Python)

Machine Learning:

- scikit-learn (Python, rich documentation, highly recommended!)
- Mallet (Java)
- WEKA (Java)
- Alchemy (graphical models, Java)
- FACTORIE (graphical models, Scala)
- CRFSuite (efficient implementation of CRFs, Python)

Ready to use NERC software:

- ANNIE (rule-based, part of GATE)
- Wikifier (based on Wikipedia)
- FIGER (based on Wikipedia, fine-grained Freebase NE classes)

Almost ready to use NERC software:

- CRFSuite (already includes Python implementation for feature extraction, you just need to feed it with training data, which you can also download)

Ready to use RE software:

- ReVerb (Open IE, extracts patterns for any kind of relation)
- MultiR (Distant supervision, relation extractor trained on Freebase)

Web Content Extraction software:

- Boilerpipe (extract main text content from Web pages)
- Jsoup (traverse elements of Web pages individually, also allows to extract text)



Thank you for
your attention!

Questions?