

Overview: Levels of Analysis in NLP

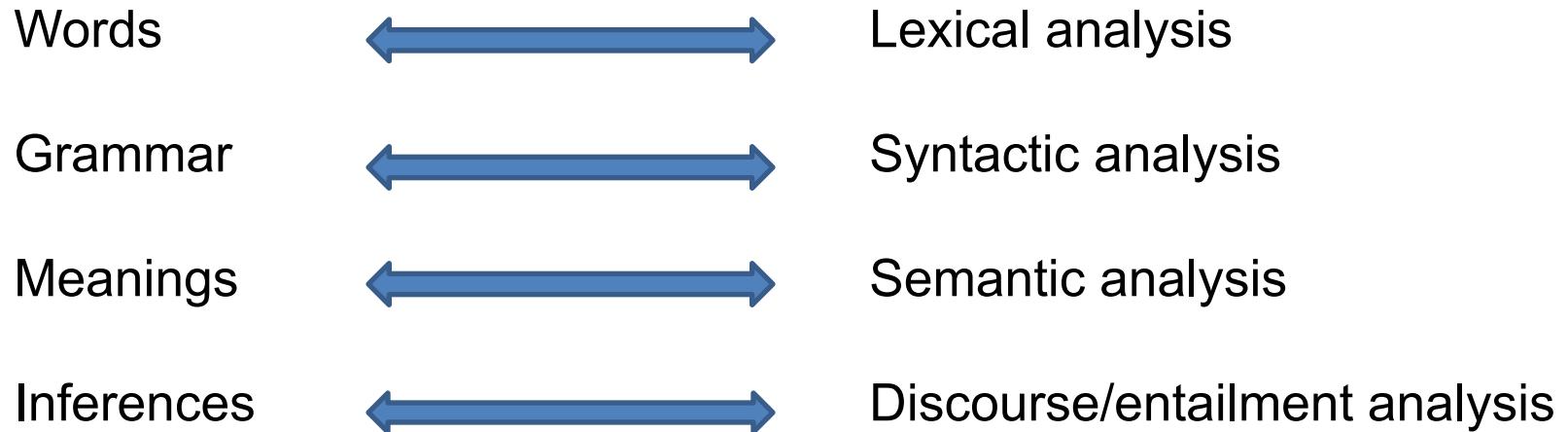
Natural Language Processing

Levels of Language Learning

- When we are growing up, we learn language at multiple levels:
 - Words: what's a real word, how it is spelled
 - Grammar: how words combine to make sentences
 - Meanings: what words and sentences mean in context
 - Inferences: how people draw more out of what is said than what is literally stated

Levels of NLP Analysis

- Similarly, we have different levels in NLP.



Increasing Difficulty

As human beings, we find it gets harder to learn a language as you progress to the next level.

Suppose you are an English speaker learning German for the first time.

Word: “Attention!”



“Achtung!”

Grammar: “It’s raining.”



“Es regnet.”
(lit: “It rains.”)

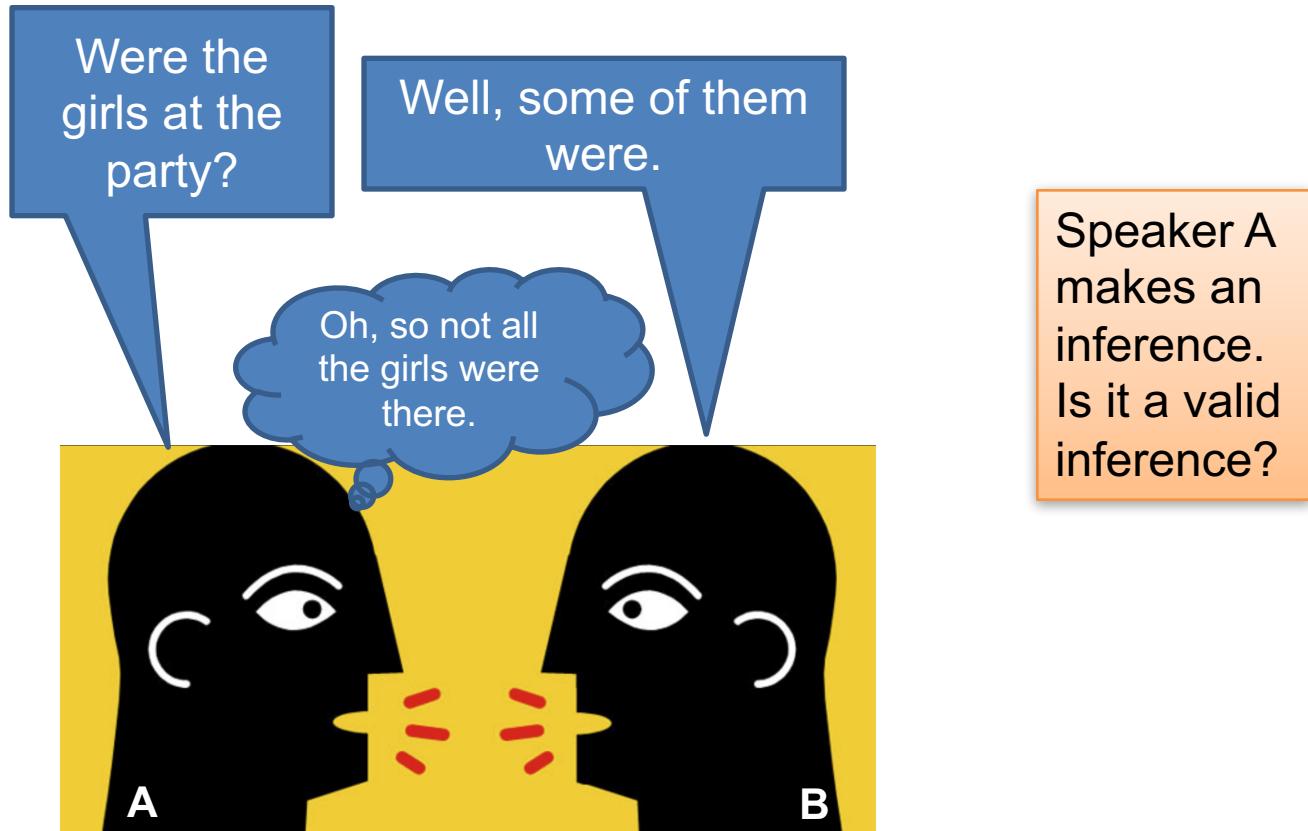
Meanings: “What time is
it? Quarter ’til nine!”



“Wieviel Uhr ist es? Dreiviertel
zehn!”
(lit: How much clock is it?
Three-quarters ten!)

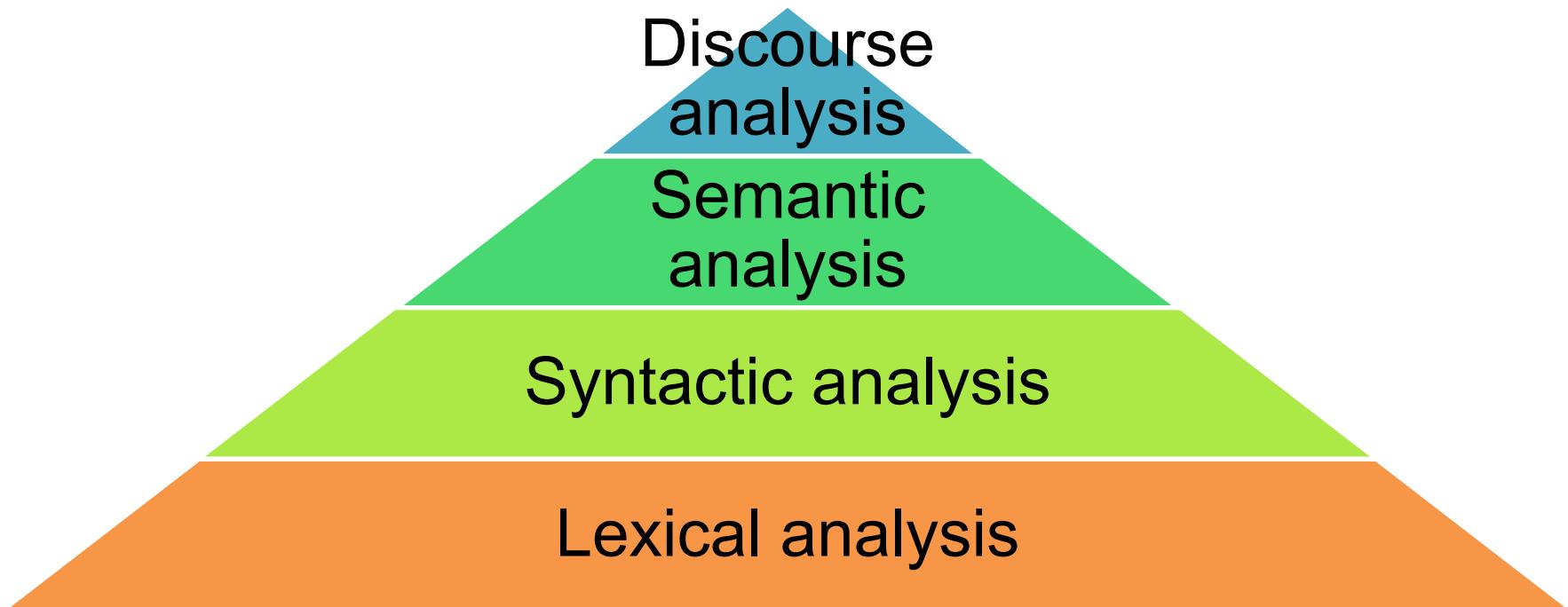
Analysis at the Level of Discourse

The inferences we make in discourse are the hardest of all. Consider the following:



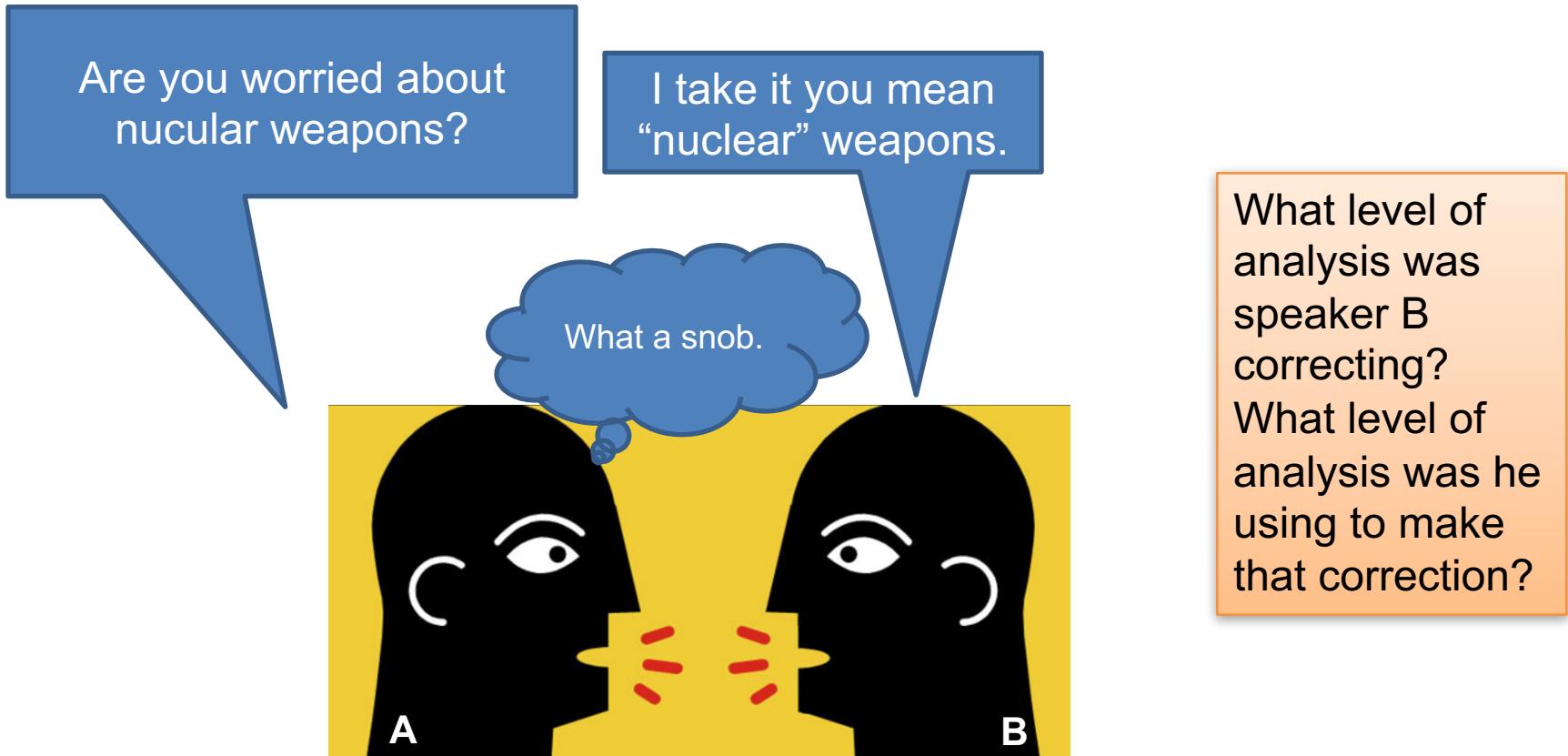
Levels of Analysis

It is tempting to think of these levels of analysis as layers, such that each one builds upon the lower levels beneath it.



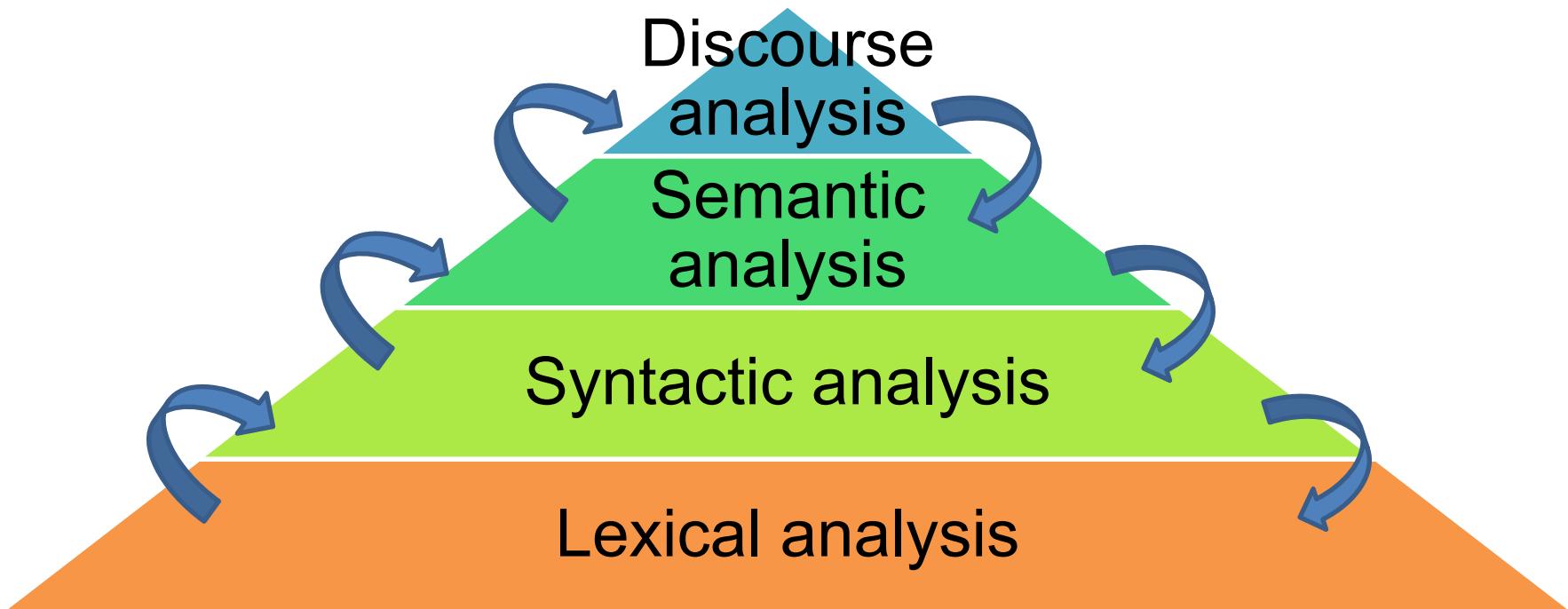
Multilevel Analysis

However, in real life, we use all the levels to help us understand all the levels:



Levels of Analysis

So, for convenience, we can continue to think of these as higher and lower levels as long as we remember that they have every manner of interconnection.



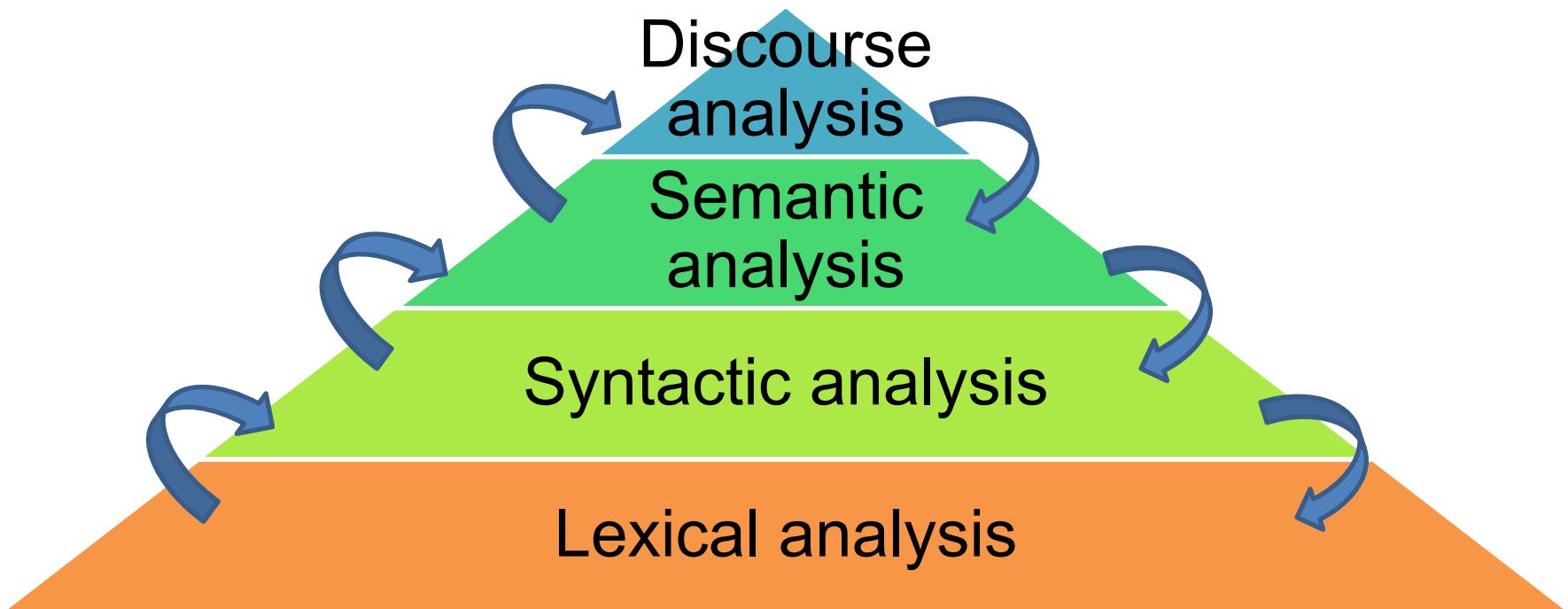
DataScience@SMU

Overview: Lexical Analysis in NLP

Natural Language Processing

Levels of Analysis

For convenience, we can think of there being higher and lower levels of analysis in NLP as long as we remember they are interconnected.



Morphology and Stemming

Morphology means the study of “morphemes,” which are the units that our words are made of.

- Take the word “running.” The root word “run” is a morpheme, and the often-used “-ing” ending is itself a morpheme.

Making an algorithm that when inputted “running” automatically gives you the root morpheme “run” is called a “stemmer.”

What applications could benefit from a stemmer?

Attachment of Metadata

Corpus-derived metadata of various types can make valuable enhancements to lexical resources. For each headword consider adding:

- Word-frequency score
- Common collocations
- Common co-occurring words, context words, etc.



Collocations

Collocations can make a big impact on the meaning of the main word. Consider all these collocations of “take.”

take a shower
take care
take a class
take damage
take a role
take initiative
take a break
take the lead
take a chance

take charge
take responsibility
take a look
take an exam
take a nap
take a claim
take a dump
take a rest
take a class

take a seat
take heart
take an opportunity
take notes
take notice
take a picture
take a taxi
take action
take advantage

It may be better to have a headword entry in the lexicon for each of these—then make sure that lexical look-ups are done right!

Enumeration of Word Senses

- Most words are “polysemous,” meaning they can have more than one meaning.
- Usually, a lexicon tries to list word senses in order, from most common to least common in occurrences.

table (tā'bl), *n.* a flat smooth board, furnished with legs; flat surface; a flat stretch of country; table-land; perspective plane; palm of the hand; upper facet of a gem; table supply; hence, food; party around a table; hence, entertainment; tablet; index or syllabus: *pl.* collection of many particulars brought into one view; collection of numbers or figures methodically arranged; the Ten Commandments: *adj.* pertaining to a table: *v.t.* to catalogue or index; lay or place on a table (as a report) for future consideration.

Enumeration of Word Senses

If WordNet is “the Bible” here, then many real-world word senses are “apocryphal.”

WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Display options for sense: (gloss) "an example sentence"

Noun

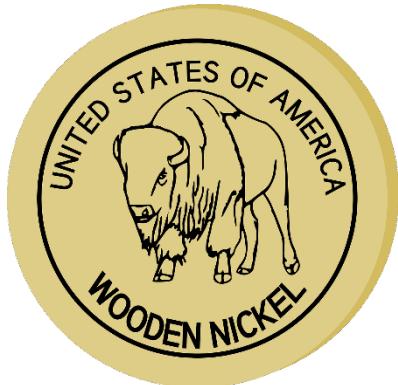
- **S: (n)** [run](#), [tally](#) (a score in baseball made by a runner touching all four bases safely) "*the Yankees scored 3 runs in the bottom of the 9th*"; "*their first tally came in the 3rd inning*"
- **S: (n)** [test](#), [trial](#), [run](#) (the act of testing something) "*in the experimental trials the amount of carbon was measured separately*"; "*he called each flip of the coin a new trial*"
- **S: (n)** [footrace](#), [foot race](#), [run](#) (a race run on foot) "*she broke the record for the half-mile run*"
- **S: (n)** [streak](#), [run](#) (an unbroken series of events) "*had a streak of bad luck*"; "*Nicklaus had a run of birdies*"
- **S: (n)** [run](#), [running](#), [running play](#), [running game](#) ((American football) a play in which a player attempts to carry the ball through or past the opposing team) "*the defensive line braced to stop the run*"; "*the coach put great emphasis on running*"
- **S: (n)** [run](#) (a regular trip) "*the ship made its run in record time*"
- **S: (n)** [run](#), [running](#) (the act of running; traveling on foot at a fast pace) "*he broke into a run*"; "*his daily run keeps him fit*"
- **S: (n)** [run](#) (the continuous period of time during which something (a machine or a factory) operates or continues in operation) "*the assembly line was on a 12-hour run*"
- **S: (n)** [run](#) (unrestricted freedom to use) "*he has the run of the house*"
- **S: (n)** [run](#) (the production achieved during a continuous period of operation (of a machine or factory etc.)) "*a daily run of 100,000 gallons of paint*"

Domain Relations

- Once we have word senses, we realize some senses are linked to, or much stronger in, one domain vs. another.
- Consider “nickel” as a metal, a coin, a sports term.

----- DOMAINS -----

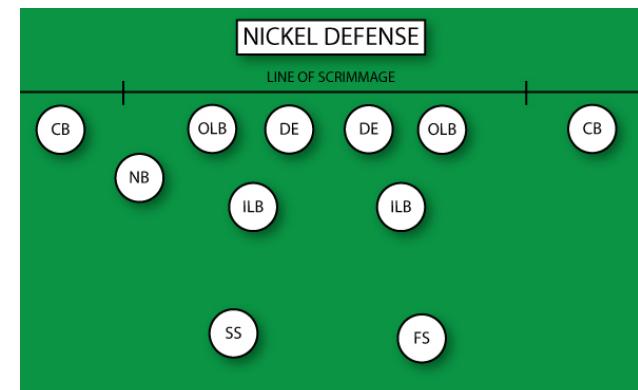
MONEY



METALLURGY



SPORTS



Text and Corpus Analytics

At the lexical level, we can already do a lot of very interesting analysis of large texts, and even of whole corpora.

Examples:

- Spell correction
- Terminology extraction
- Lexical diversity measurement

```
>>> lexical_diversity(text3)
0.06230453042623537
>>> lexical_diversity(text5)
0.13477005109975562
```

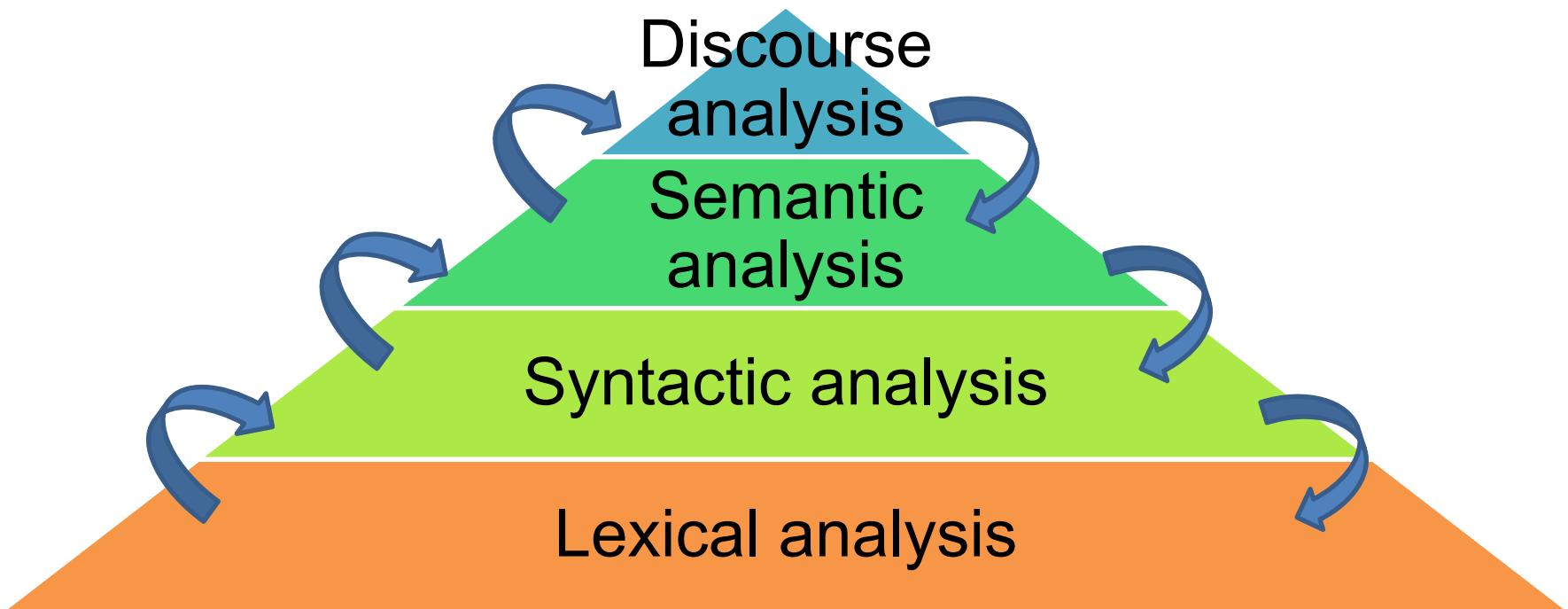
DataScience@SMU

Overview: Lexical Analysis in NLP

Natural Language Processing

Levels of Analysis

For convenience, we can think of there being higher and lower levels of analysis in NLP as long as we remember they are interconnected.



Morphology and Stemming

Morphology means the study of “morphemes,” which are the units that our words are made of.

- Take the word “running.” The root word “run” is a morpheme, and the often-used “-ing” ending is itself a morpheme.

Making an algorithm that when inputted “running” automatically gives you the root morpheme “run” is called a “stemmer.”

What applications could benefit from a stemmer?

Attachment of Metadata

Corpus-derived metadata of various types can make valuable enhancements to lexical resources. For each headword consider adding:

- Word-frequency score
- Common collocations
- Common co-occurring words, context words, etc.

A word-cloud built from co-occurrence frequencies with the word “MORALS” might look like this.



Collocations

Collocations can make a big impact on the meaning of the main word. Consider all these collocations of “take.”

take a shower
take care
take a class
take damage
take a role
take initiative
take a break
take the lead
take a chance

take charge
take responsibility
take a look
take an exam
take a nap
take a claim
take a dump
take a rest
take a class

take a seat
take heart
take an opportunity
take notes
take notice
take a picture
take a taxi
take action
take advantage

It may be better to have a headword entry in the lexicon for each of these—then make sure that lexical look-ups are done right!

Enumeration of Word Senses

- Most words are “polysemous,” meaning they can have more than one meaning.
- Usually, a lexicon tries to list word senses in order, from most common to least common in occurrences.

table (tā'bl), *n.* a flat smooth board, furnished with legs; flat surface; a flat stretch of country; table-land; perspective plane; palm of the hand; upper facet of a gem; table supply; hence, food; party around a table; hence, entertainment; tablet; index or syllabus: *pl.* collection of many particulars brought into one view; collection of numbers or figures methodically arranged; the Ten Commandments: *adj.* pertaining to a table: *v.t.* to catalogue or index; lay or place on a table (as a report) for future consideration.

Enumeration of Word Senses

If WordNet is “the Bible” here, then many real-world word senses are “apocryphal.”

WordNet Search - 3.1
- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations
Display options for sense: (gloss) "an example sentence"

Noun

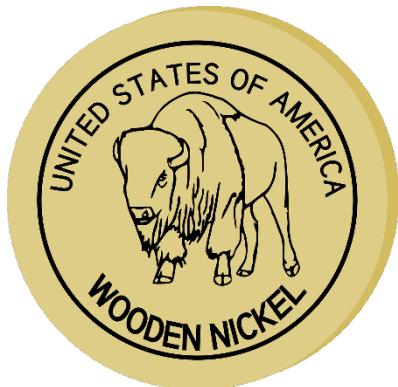
- **S: (n) run, tally** (a score in baseball made by a runner touching all four bases safely)
"the Yankees scored 3 runs in the bottom of the 9th"; "their first tally came in the 3rd inning"
- **S: (n) test, trial, run** (the act of testing something) *"in the experimental trials the amount of carbon was measured separately"; "he called each flip of the coin a new trial"*
- **S: (n) footrace, foot race, run** (a race run on foot) *"she broke the record for the half-mile run"*
- **S: (n) streak, run** (an unbroken series of events) *"had a streak of bad luck"; "Nicklaus had a run of birdies"*
- **S: (n) run, running, running play, running game** ((American football) a play in which a player attempts to carry the ball through or past the opposing team) *"the defensive line braced to stop the run"; "the coach put great emphasis on running"*
- **S: (n) run** (a regular trip) *"the ship made its run in record time"*
- **S: (n) run, running** (the act of running; traveling on foot at a fast pace) *"he broke into a run"; "his daily run keeps him fit"*
- **S: (n) run** (the continuous period of time during which something (a machine or a factory) operates or continues in operation) *"the assembly line was on a 12-hour run"*
- **S: (n) run** (unrestricted freedom to use) *"he has the run of the house"*
- **S: (n) run** (the production achieved during a continuous period of operation (of a machine or factory etc.)) *"a daily run of 100,000 gallons of paint"*

Domain Relations

- Once we have word senses, we realize some senses are linked to, or much stronger in, one domain vs. another.
- Consider “nickel” as a metal, a coin, a sports term.

----- DOMAINS -----

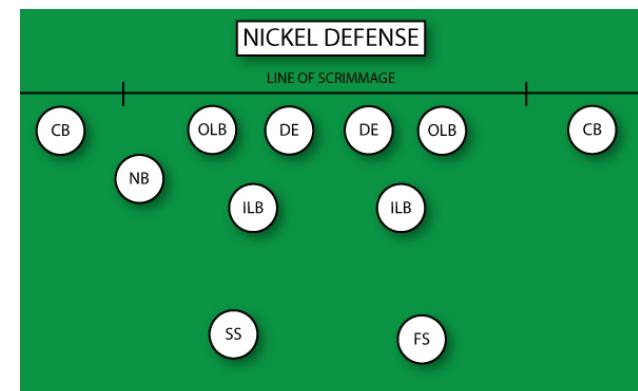
MONEY



METALLURGY



SPORTS



Text and Corpus Analytics

At the lexical level, we can already do a lot of very interesting analysis of large texts, and even of whole corpora.

Examples:

- Spell correction
- Terminology extraction
- Lexical diversity measurement

```
>>> lexical_diversity(text3)
0.06230453042623537
>>> lexical_diversity(text5)
0.13477005109975562
```

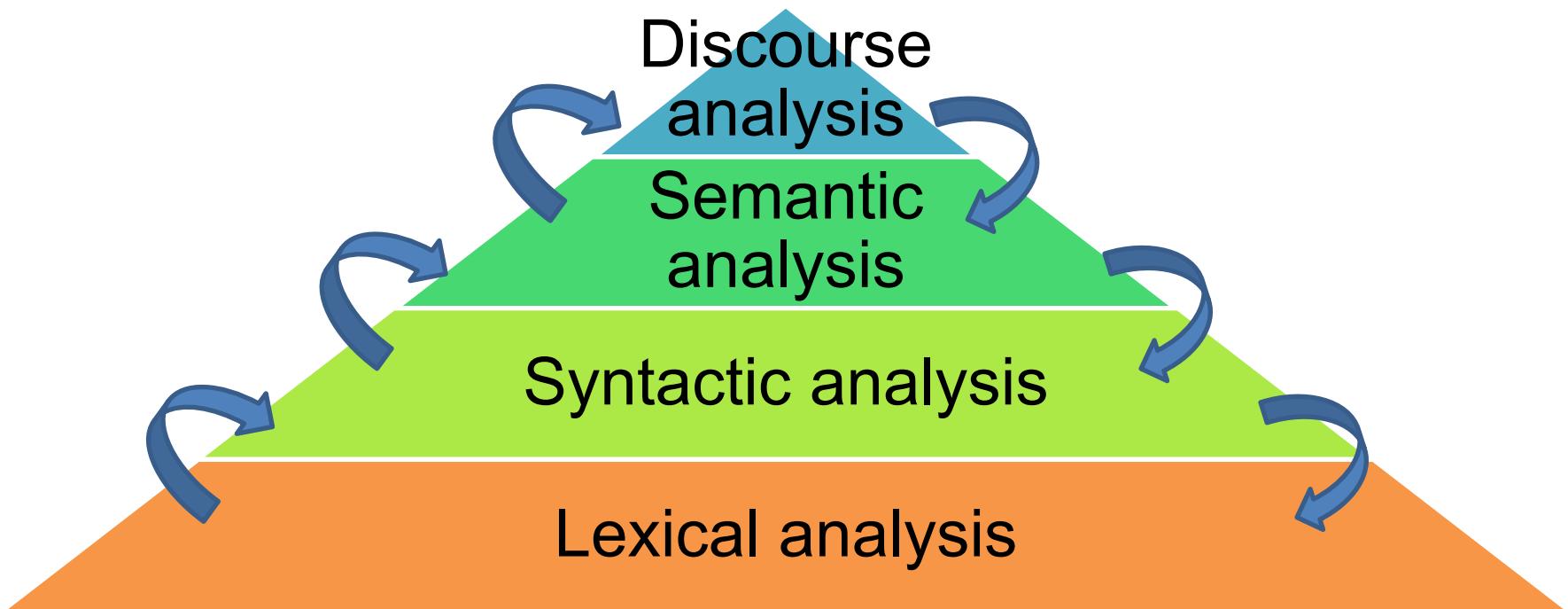
DataScience@SMU

Overview: Syntactic Analysis in NLP

Natural Language Processing

Levels of Analysis

For convenience, we can think of there being higher and lower levels of analysis in NLP as long as we remember they are interconnected.

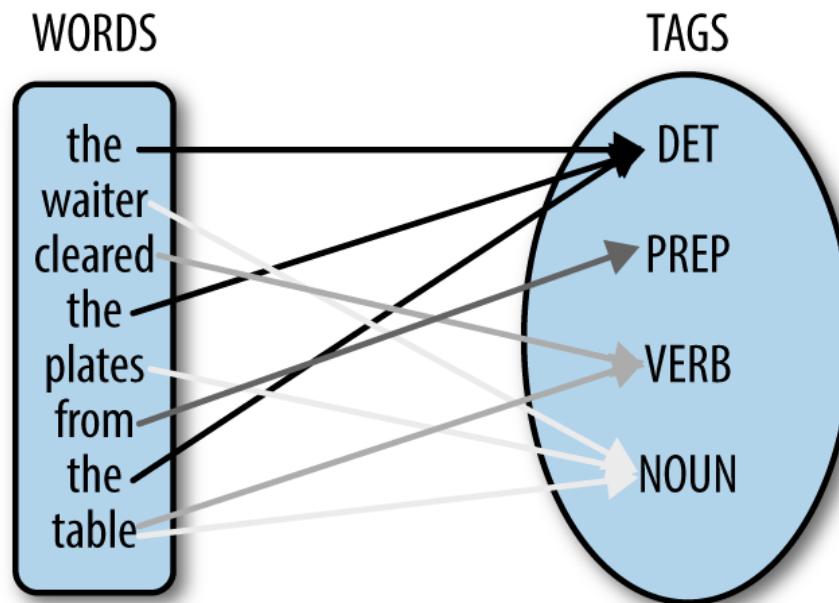


Sentence Boundary Detection

- The most common sentence boundary marker, the period, “.”, is used unfortunately for abbreviations; therefore, sentence boundary detection is nontrivial in English.

POS Tagging

- Most words can have more than one part of speech; therefore, POS tagging is nontrivial.



POS Tagging

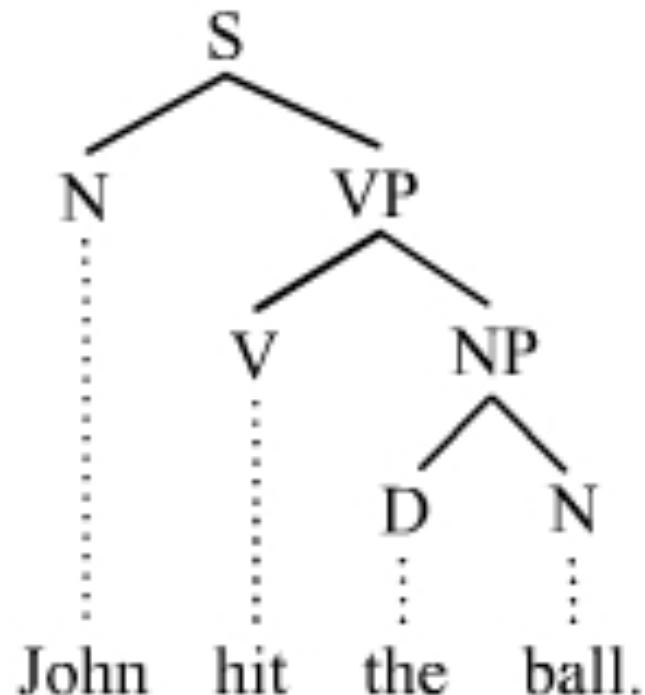
One of the most commonly used tag sets:

Penn Treebank Tagset

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &</i>
CD	Cardinal number	<i>one, two, three</i>	TO	"to"	<i>to</i>
DT	Determiner	<i>a, the</i>	UH	Interjection	<i>ah, oops</i>
EX	Existential 'there'	<i>there</i>	VB	Verb, base form	<i>eat</i>
FW	Foreign word	<i>mea culpa</i>	VBD	Verb, past tense	<i>ate</i>
IN	Preposition/sub-conj	<i>of, in, by</i>	VBG	Verb, gerund	<i>eating</i>
JJ	Adjective	<i>yellow</i>	VBN	Verb, past participle	<i>eaten</i>
JJR	Adj., comparative	<i>bigger</i>	VBP	Verb, non-3sg pres	<i>eat</i>
JJS	Adj., superlative	<i>wildest</i>	VBZ	Verb, 3sg pres	<i>eats</i>
LS	List item marker	<i>1, 2, One</i>	WDT	Wh-determiner	<i>which, that</i>
MD	Modal	<i>can, should</i>	WP	Wh-pronoun	<i>what, who</i>
NN	Noun, sing. or mass	<i>llama</i>	WP\$	Possessive wh-	<i>whose</i>
NNS	Noun, plural	<i>llamas</i>	WRB	Wh-adverb	<i>how, where</i>
NNP	Proper noun, singular	<i>IBM</i>	\$	Dollar sign	<i>\$</i>
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	<i>#</i>
PDT	Predeterminer	<i>all, both</i>	"	Left quote	<i>(‘ or “)</i>
POS	Possessive ending	<i>'s</i>	"	Right quote	<i>(‘ or ”)</i>
PRP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	<i>([, (, {, <)</i>
PRP\$	Possessive pronoun	<i>your, one's</i>)	Right parenthesis	<i>(],), }, >)</i>
RB	Adverb	<i>quickly, never</i>	,	Comma	<i>,</i>
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	<i>(. ! ?)</i>
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	<i>(: ; ... – -)</i>
RP	Particle	<i>up, off</i>			

Parsing

- No parser is perfect, but for most applications we can live with less-than-perfect parsing.



Lemmatization

- A lemma is the canonical (conventional) form that represents a set of related word forms, e.g., for *run, runs, ran, running*, the lemma is “run.”
- Often, a context-free stemming gives us the lemma, but in some cases not. Example:

wetter → wet

but

better → bet?

Usually not!

better → good

- If we parse “better” as an adjective, not a noun, then we can determine its correct lemma.

Discrete Text Field Analysis

- Unitizing
- Normalizing
- “Smart ETL”

The screenshot shows an Amazon product page for an Acer Aspire E 15 laptop. The top navigation bar includes the Amazon logo, a search bar with 'acer laptop', and various menu options like Electronics, Departments, and Account & Lists. A promotional banner at the top right says '\$10 & under with FREE shipping'. The main product image is a black Acer Aspire E 15 laptop with its screen open, showing the Windows 10 desktop. To the left of the main image is a vertical sidebar with six smaller thumbnail images of different laptops. The product title is 'Acer Aspire E 15 E5-575-33BM 15.6-Inch FHD Notebook (Intel Core i3-7100U 7th Generation, 4GB DDR4, 1TB 5400RPM HD, Intel HD Graphics 620, Windows 10 Home), Obsidian Black'. It has a 4.5-star rating from 3,410 reviews. The price is listed as \$349.99 with FREE Shipping. The product is described as 'In Stock' and offers 'One-Day Shipping'. Key features listed include a 7th Generation Intel Core i3-7100U Processor, 15.6" Full HD Widescreen ComfyView LED-backlit Display, 4GB DDR4 Memory, and Windows 10 Home.

Discrete Text Field Analysis

The screenshot shows an Amazon product page for an Acer Aspire E 15 laptop. The product is displayed prominently in the center, showing its dark frame and screen displaying the Windows 10 desktop. To the left is a vertical sidebar with icons for different laptop models. The main title is "Acer Aspire E 15 E5-575-33BM 15.6-Inch FHD Notebook (Intel Core i3-7100U 7th Generation, 4GB DDR4, 1TB 5400RPM HD, Intel HD Graphics 620, Windows 10 Home), Obsidian Black". Below the title, it says "#1 Best Seller in Traditional Laptop Computers". The price is listed as \$349.99 & FREE Shipping. A red box highlights the product description text. To the right, a large red callout box contains the text: "All of this text is crammed into a single data field—the “product name” field. Marketers do this all the time for SEO." An arrow points from this callout box up towards the highlighted product description text.

Electronics > Computers & Accessories > Computers & Tablets > Laptops > Traditional Laptops

Computers Laptops Tablets Desktops Monitors Computer Accessories PC Components

Deliver to sanfrancisco 94102 Departments Your Recommendations Today's Deals Gift Cards Registry Sell Help

EN Hello, Sign in Account & Lists

Electronics > Computers & Accessories > Computers & Tablets > Laptops > Traditional Laptops

Acer Aspire E 15 E5-575-33BM 15.6-Inch FHD Notebook (Intel Core i3-7100U 7th Generation, 4GB DDR4, 1TB 5400RPM HD, Intel HD Graphics 620, Windows 10 Home), Obsidian Black

#1 Best Seller in Traditional Laptop Computers

Price: \$349.99 & FREE Shipping. Details

In Stock. Want it tomorrow, April 5? Order within 17 hrs 53 mins and choose One-Day Delivery. Ships from and sold by Amazon.com. Gift-wrap available.

Style: Laptop Only

Laptop + Microsoft Office
\$434.98

Laptop Only
\$349.99

- 7th Generation Intel Core i3-7100U Processor (2.4GHz, 3MB L3 cache)
- 15.6" Full HD Widescreen ComfyView LED-backlit Display supporting Acer TrueColor
- 4GB DDR4 Memory, 1TB 5400RPM HDD
- Windows 10 Home
- Up to 12-hours Battery Life

› See more product details

All of this text is crammed into a single data field—the “product name” field. Marketers do this all the time for SEO.

These text strings obviously do not follow regular grammar, so a standard parser is no help to us. We might write a custom parser using a catalog of product names and features.

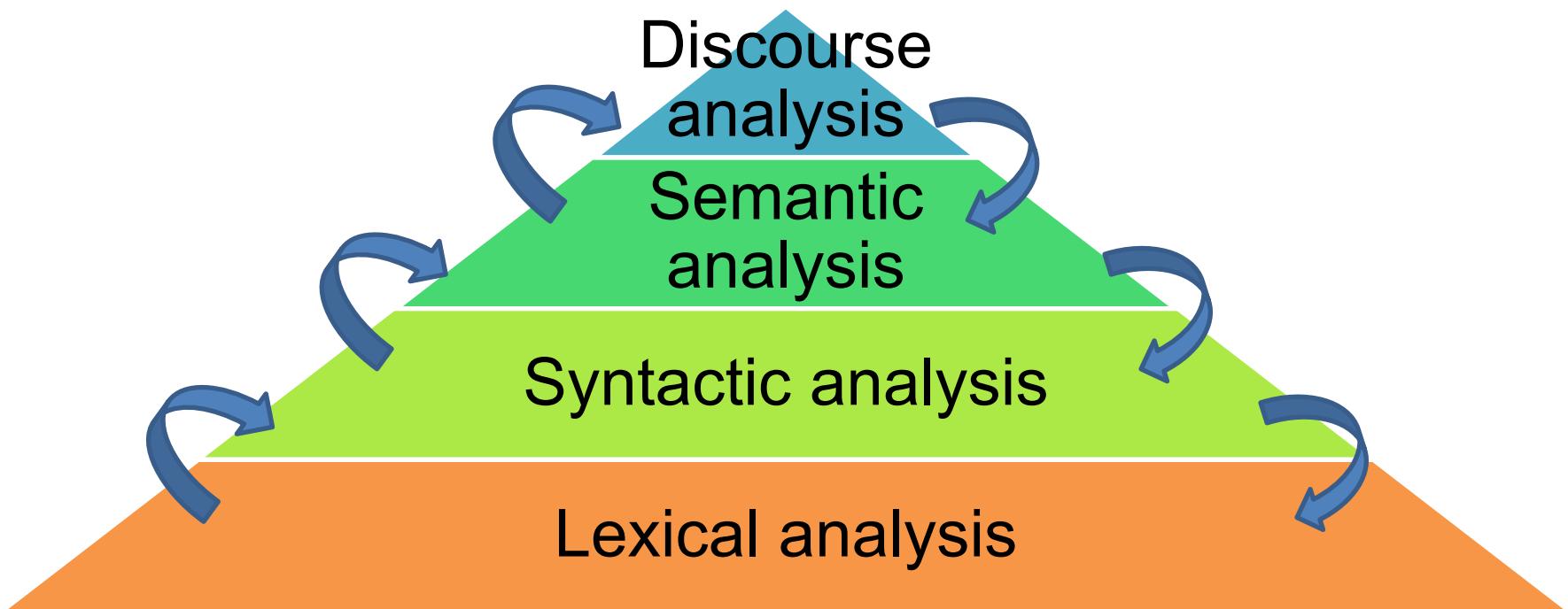
DataScience@SMU

Overview: Semantic Analysis in NLP

Natural Language Processing

Levels of Analysis

For convenience, we can think of there being higher and lower levels of analysis in NLP as long as we remember they are interconnected.



Semantic Analysis

Some major types:

- Named entity extraction (NEE a.k.a. NER)
- Relationship extraction (between NEs)
- Word sense disambiguation (WSD)
- Classification
- Tagging
- Topic segmentation
- Sentiment analysis

Named-Entity Extraction

- Also called “named-entity recognition” (abbreviated NEE or NER).
- Simple NEE just recognizes entities without typing.
- More sophisticated NEE does basic typing:
 - Persons
 - Organizations
 - Places
 - Events

Named-Entity Extraction

Good NEE will cluster together many variants, including epithets.



Hillary Clinton Donald Trump

Hillary Rodham Clinton The Donald

Senator Clinton Donald J. Trump

Secretary Clinton Donald Jackass Trump

Hitlery Clinton Donald Drumpf

Crooked Hillary Herr Trump

Relationship Extraction (between NEs)

- What relation does one NE have to another? Sometimes we can answer this during the NEE process.

This time, the CEO of Apple Tim Cook, said that “the DACA situation is one that I am truthfully, as an American, deeply offended by.”

These are sometimes called “triples” because there are three elements: two NEs and one relation.

Collections of triples are called “triple stores.”

Word Sense Disambiguation (WSD)

WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Display options for sense: (gloss) "an example sentence"

Noun

- S: (n) chair (a seat for one person, with a support for the back) "*he put his coat over the back of the chair and sat down*"
- S: (n) professorship, chair (the position of professor) "*he was awarded an endowed chair in economics*"
- S: (n) president, chairman, chairwoman, chair, chairperson (the officer who presides at the meetings of an organization) "*address your remarks to the chairperson*"
- S: (n) electric chair, chair, death chair, hot seat (an instrument of execution by electrocution; resembles an ordinary seat for one person) "*the murderer was sentenced to die in the chair*"
- S: (n) chair (a particular seat in an orchestra) "*he is second chair violin*"

Verb

- S: (v) chair, chairman (act or preside as chair, as of an academic department in a university) "*She chaired the department for many years*"
- S: (v) moderate, chair, lead (preside over) "*John moderated the discussion*"

Word Sense Disambiguation (WSD)

WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Display options for sense: (gloss) "an example sentence"

Noun

- S: (n) chair (a seat for one person, with a support for the back) "*he put his coat over the back of the chair and sat down*"
- S: (n) professorship, chair (the position of professor) "*he was awarded an endowed chair in economics*"
- S: (n) president, chairman, chairwoman, chair, chairperson (the officer who presides at the meetings of an organization) "*address your remarks to the chairperson*"
- S: (n) electric chair, chair, death chair, hot seat (an instrument of execution by electrocution; resembles an ordinary seat for one person) "*the murderer was sentenced to die in the chair*"
- S: (n) chair (a particular seat in an orchestra) "*he is second chair violin*"

Context words give us a chance to disambiguate intended senses.

Verb

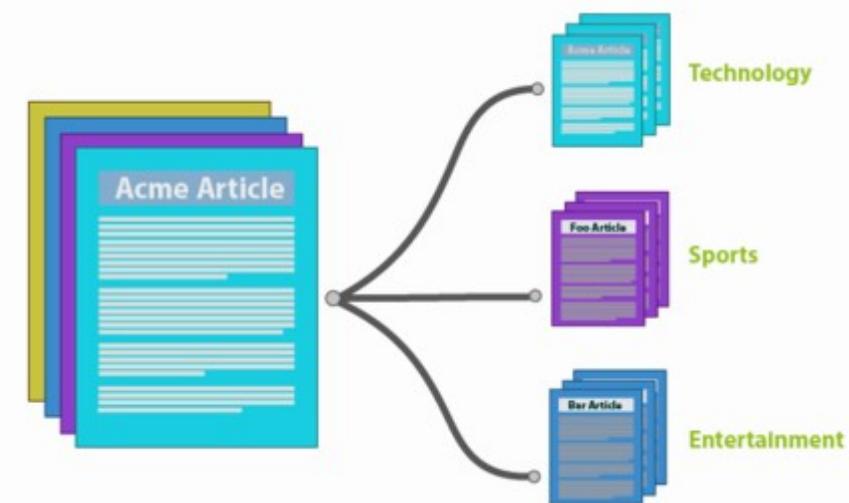
- S: (v) chair, chairman (act or preside as chair, as of an academic department in a university) "*She chaired the department for many years*"
- S: (v) moderate, chair, lead (preside over) "*John moderated the discussion*"

Doing this at scale is hard! It's still an open research area of NLP (i.e., unsolved).

Classification

- In classification, we use a tree-structured graph to place documents into categories.
- Often we employ machine learning methods, such as SVM*.

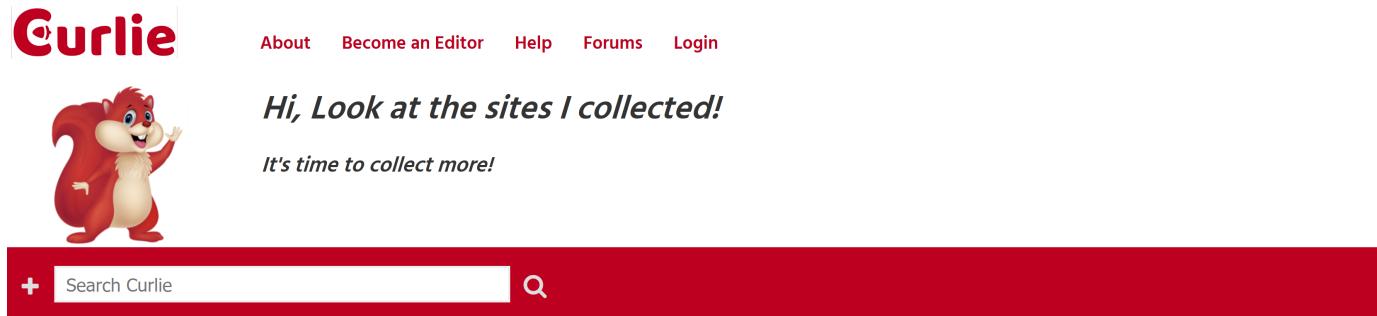
We always start from a preexisting category schema, often called a “taxonomy.”



*support vector machine

Classification

Curlie.org is the successor of the once-famous DMOZ directory of web pages.



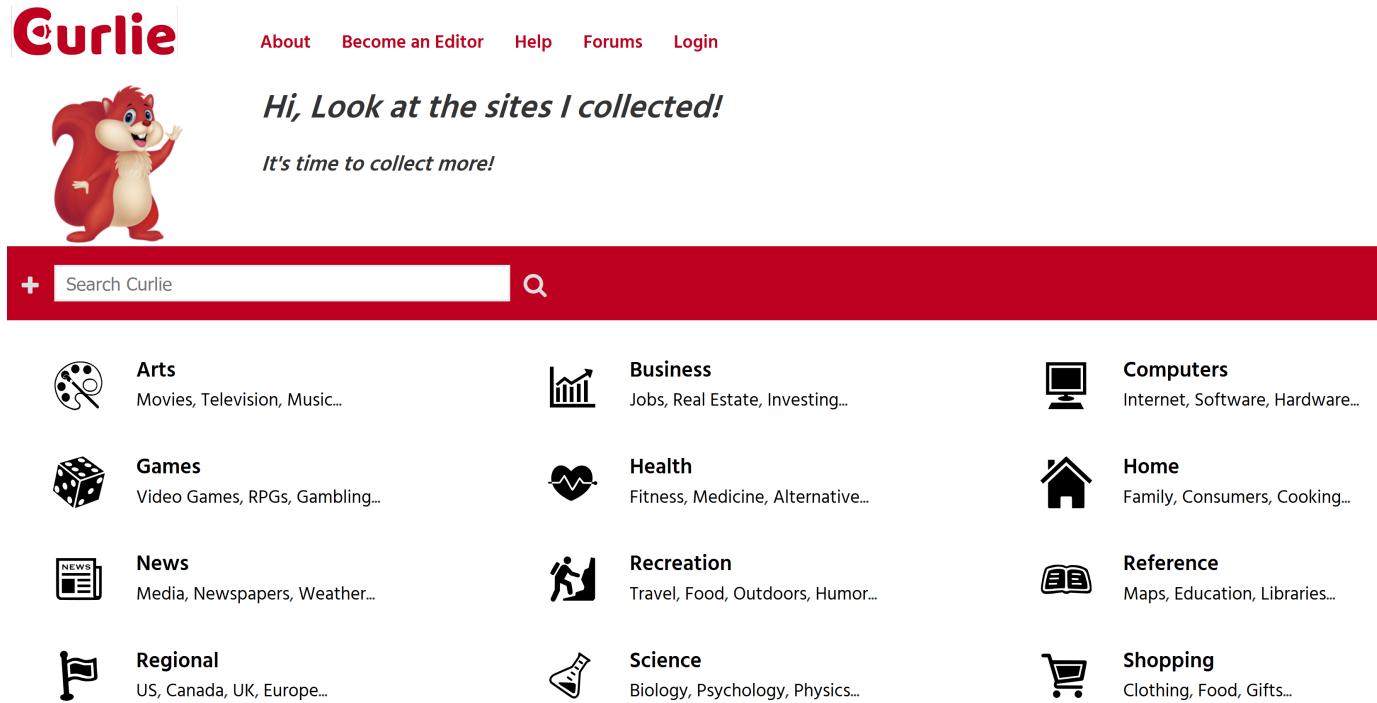
The screenshot shows the Curlie.org homepage. At the top left is the "Curlie" logo with a red squirrel icon. To its right are links for "About", "Become an Editor", "Help", "Forums", and "Login". Below the logo is a large red squirrel icon with the text "Hi, Look at the sites I collected!" and "It's time to collect more!". A red search bar with a magnifying glass icon and a plus sign is centered below the squirrel. The main content area displays ten categories with icons and descriptions:

 Arts Movies, Television, Music...	 Business Jobs, Real Estate, Investing...	 Computers Internet, Software, Hardware...
 Games Video Games, RPGs, Gambling...	 Health Fitness, Medicine, Alternative...	 Home Family, Consumers, Cooking...
 News Media, Newspapers, Weather...	 Recreation Travel, Food, Outdoors, Humor...	 Reference Maps, Education, Libraries...
 Regional US, Canada, UK, Europe...	 Science Biology, Psychology, Physics...	 Shopping Clothing, Food, Gifts...

Think about classifying web pages into the above categories.
Why is it a misnomer to call this a “taxonomy”?

Classification

Curlie.org is the successor of the once-famous DMOZ directory of web pages.



The screenshot shows the Curlie.org homepage. At the top left is the red "Curlie" logo with a white outline. To its right are links for "About", "Become an Editor", "Help", "Forums", and "Login". Below the logo is a cartoon squirrel waving. Next to the squirrel is the text "Hi, Look at the sites I collected!" and "It's time to collect more!". A red search bar with a white input field containing "Search Curlie" and a magnifying glass icon is centered below the squirrel. Below the search bar are nine category cards arranged in a grid. Each card has an icon on the left and a title and description on the right. The categories are: Arts (Movies, Television, Music...), Business (Jobs, Real Estate, Investing...), Computers (Internet, Software, Hardware...), Games (Video Games, RPGs, Gambling...), Health (Fitness, Medicine, Alternative...), Home (Family, Consumers, Cooking...), News (Media, Newspapers, Weather...), Recreation (Travel, Food, Outdoors, Humor...), Reference (Maps, Education, Libraries...), Regional (US, Canada, UK, Europe...), Science (Biology, Psychology, Physics...), and Shopping (Clothing, Food, Gifts...).

 Arts Movies, Television, Music...	 Business Jobs, Real Estate, Investing...	 Computers Internet, Software, Hardware...
 Games Video Games, RPGs, Gambling...	 Health Fitness, Medicine, Alternative...	 Home Family, Consumers, Cooking...
 News Media, Newspapers, Weather...	 Recreation Travel, Food, Outdoors, Humor...	 Reference Maps, Education, Libraries...
 Regional US, Canada, UK, Europe...	 Science Biology, Psychology, Physics...	 Shopping Clothing, Food, Gifts...

Why is it a misnomer to call this a “taxonomy”?
Because there will be plenty of cases of dual classification.
In a real taxonomy, every item is classified into one and only category path.

Classification

Q: What are some other big (famous, influential) document taxonomies on the web?

?

Classification

Q: What are some other big (famous, influential) document taxonomies on the web?

A: IMDB, IAB, About.com, the Google Product Taxonomy, and many website taxonomies such as Amazon, CNET, etc.

Other Examples

Additional types
of semantic
analysis include:

- Tagging
- Topic segmentation
- Sentiment analysis

Topic:
Knee
injuries

Topic:
Rehab

Tags: sports, injuries, coaching, soccer

Lorem ipsum dolor sit amet, consectetur adipiscing elit. posuere tortor vitae elit. Sed vitae metus a elit bibendum malesuada cras pulvinar. Quisque pellentesque nibh in sem. Curabitur ligula. Suspendisse potenti. Duis sit amet augue eu arcu ultrices auctor. Suspendisse elementum, nunc ut molestie elementum, neque augue vulputate elit, eu blandit enim velit vitae nulla. Duis sed.

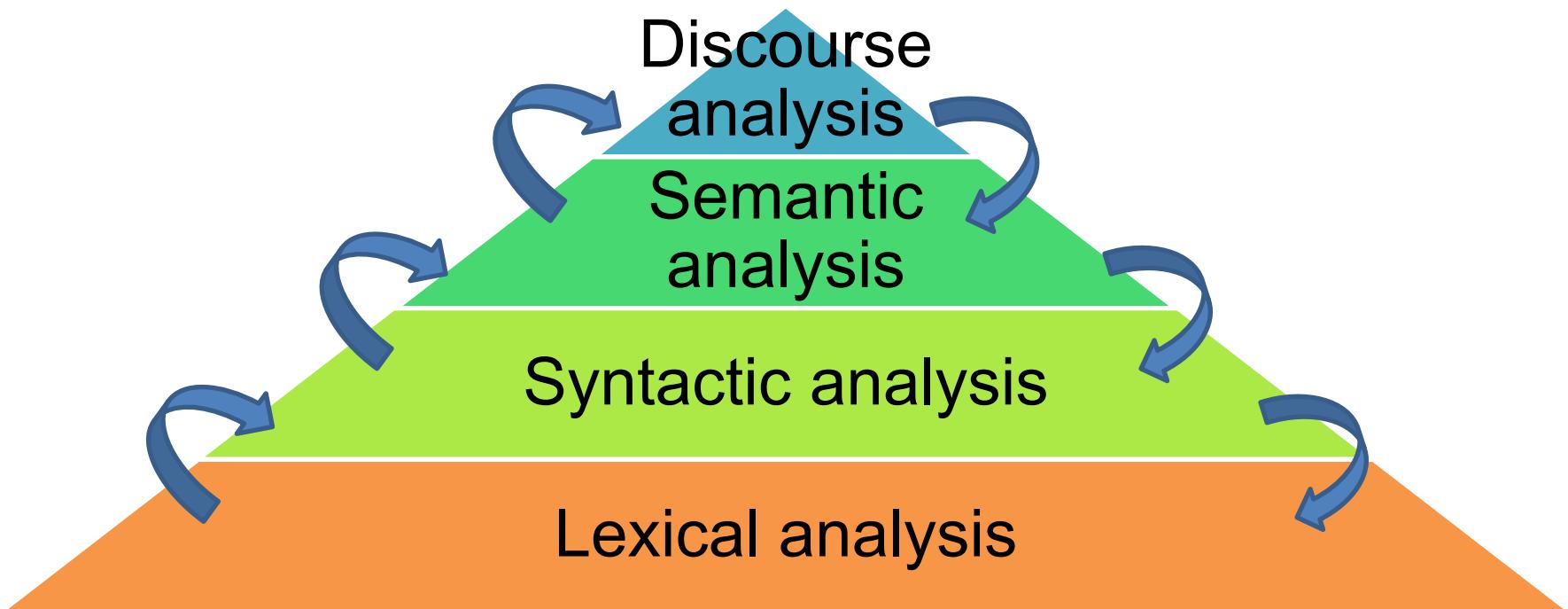
DataScience@SMU

Overview: Semantic Analysis in NLP

Natural Language Processing

Levels of Analysis

For convenience, we can think of there being higher and lower levels of analysis in NLP as long as we remember they are interconnected.



Semantic Analysis

Some major types:

- Named entity extraction (NEE a.k.a. NER)
- Relationship extraction (between NEs)
- Word sense disambiguation (WSD)
- Classification
- Tagging
- Topic segmentation
- Sentiment analysis

Named-Entity Extraction

- Also called “named-entity recognition” (abbreviated NEE or NER).
- Simple NEE just recognizes entities without typing.
- More sophisticated NEE does basic typing:
 - Persons
 - Organizations
 - Places
 - Events

Named-Entity Extraction

Good NEE will cluster together many variants, including epithets.



Hillary Clinton Donald Trump

Hillary Rodham Clinton The Donald

Senator Clinton Donald J. Trump

Secretary Clinton Donald Jackass Trump

Hitlery Clinton Donald Drumpf

Crooked Hillary Herr Trump

Relationship Extraction (between NEs)

- What relation does one NE have to another? Sometimes we can answer this during the NEE process.

This time, the CEO of Apple Tim Cook, said that “the DACA situation is one that I am truthfully, as an American, deeply offended by.”

These are sometimes called “triples” because there are three elements: two NEs and one relation.

Collections of triples are called “triple stores.”

Word Sense Disambiguation (WSD)

WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Display options for sense: (gloss) "an example sentence"

Noun

- S: (n) chair (a seat for one person, with a support for the back) "*he put his coat over the back of the chair and sat down*"
- S: (n) professorship, chair (the position of professor) "*he was awarded an endowed chair in economics*"
- S: (n) president, chairman, chairwoman, chair, chairperson (the officer who presides at the meetings of an organization) "*address your remarks to the chairperson*"
- S: (n) electric chair, chair, death chair, hot seat (an instrument of execution by electrocution; resembles an ordinary seat for one person) "*the murderer was sentenced to die in the chair*"
- S: (n) chair (a particular seat in an orchestra) "*he is second chair violin*"

Verb

- S: (v) chair, chairman (act or preside as chair, as of an academic department in a university) "*She chaired the department for many years*"
- S: (v) moderate, chair, lead (preside over) "*John moderated the discussion*"

Word Sense Disambiguation (WSD)

WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Display options for sense: (gloss) "an example sentence"

Noun

- S: (n) chair (a seat for one person, with a support for the back) "*he put his coat over the back of the chair and sat down*"
- S: (n) professorship, chair (the position of professor) "*he was awarded an endowed chair in economics*"
- S: (n) president, chairman, chairwoman, chair, chairperson (the officer who presides at the meetings of an organization) "*address your remarks to the chairperson*"
- S: (n) electric chair, chair, death chair, hot seat (an instrument of execution by electrocution; resembles an ordinary seat for one person) "*the murderer was sentenced to die in the chair*"
- S: (n) chair (a particular seat in an orchestra) "*he is second chair violin*"

Verb

- S: (v) chair, chairman (act or preside as chair, as of an academic department in a university) "*She chaired the department for many years*"
- S: (v) moderate, chair, lead (preside over) "*John moderated the discussion*"

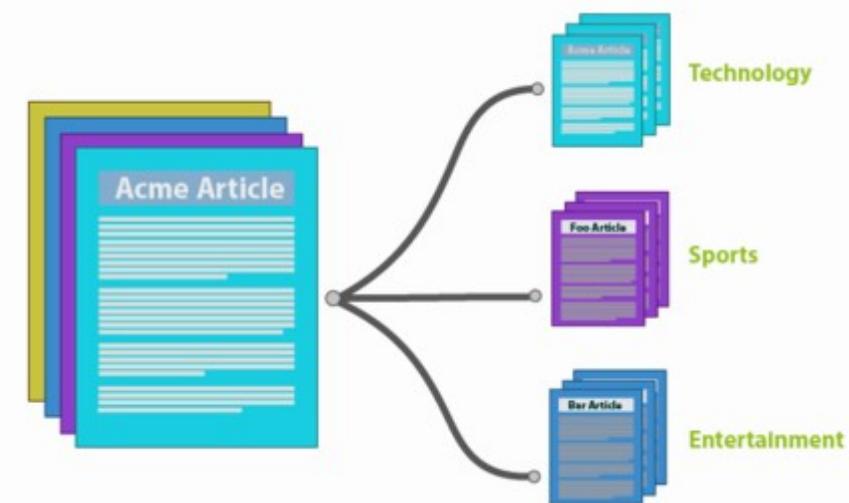
Context words give us a chance to disambiguate intended senses.

Doing this at scale is hard!
It's still an open research area of NLP (i.e., unsolved).

Classification

- In classification, we use a tree-structured graph to place documents into categories.
- Often we employ machine learning methods, such as SVM*.

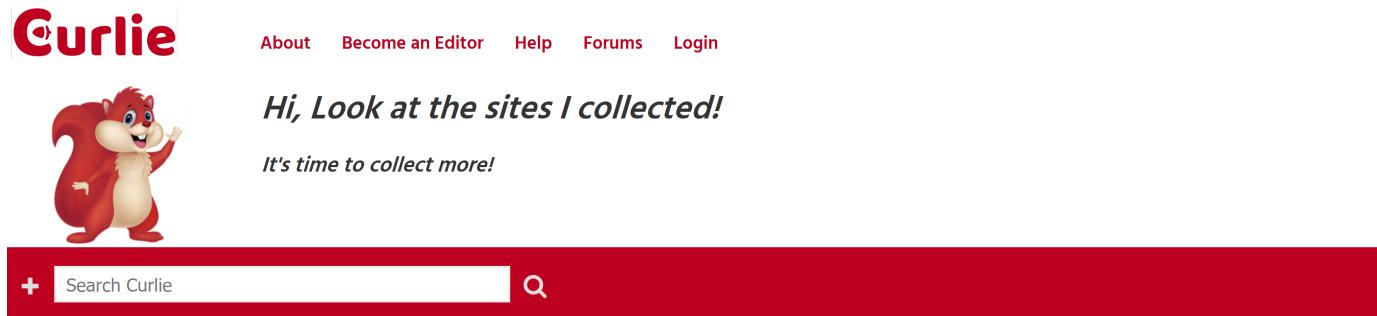
We always start from a preexisting category schema, often called a “taxonomy.”



*support vector machine

Classification

Curlie.org is the successor of the once-famous DMOZ directory of web pages.



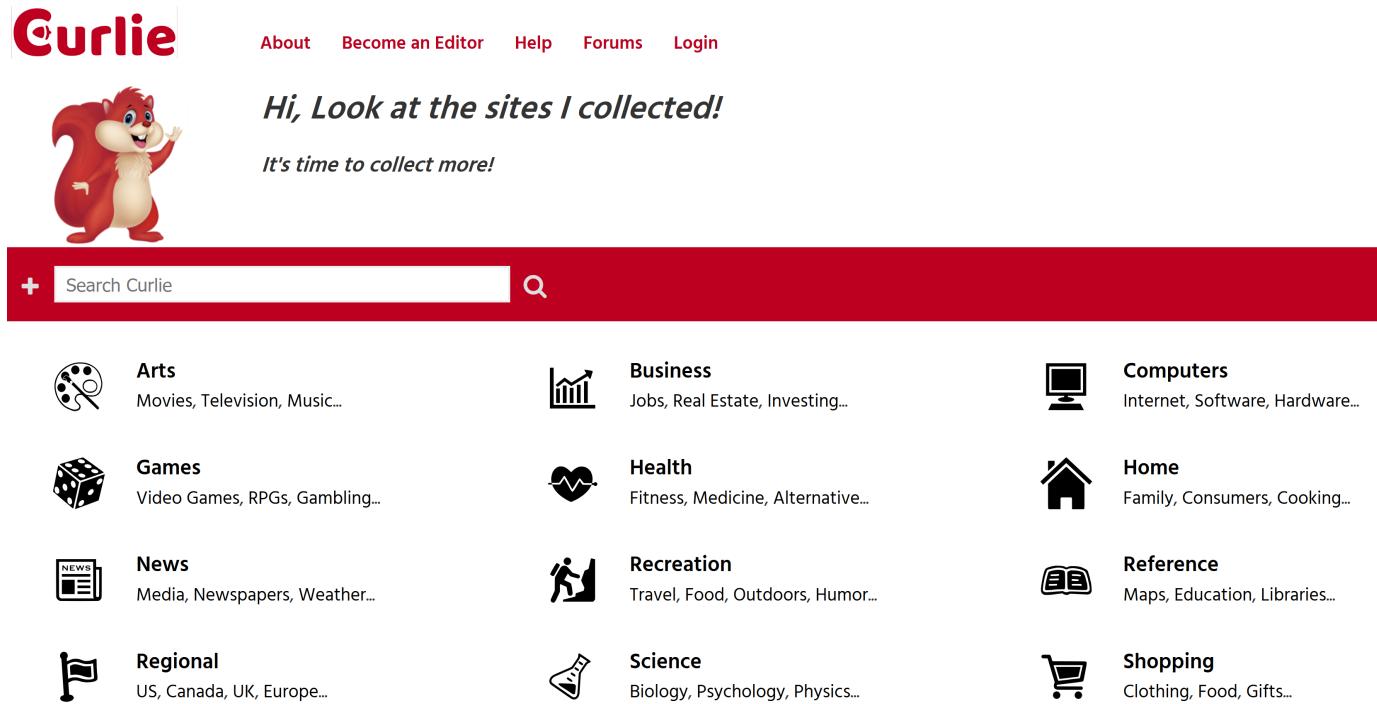
The screenshot shows the Curlie.org homepage. At the top left is the "Curlie" logo with a red squirrel icon. To its right are links for "About", "Become an Editor", "Help", "Forums", and "Login". Below the logo is a large red squirrel icon with the text "Hi, Look at the sites I collected!" and "It's time to collect more!". A red search bar with a magnifying glass icon and a plus sign is centered below the squirrel. The main content area displays ten categories with icons and descriptions:

 Arts Movies, Television, Music...	 Business Jobs, Real Estate, Investing...	 Computers Internet, Software, Hardware...
 Games Video Games, RPGs, Gambling...	 Health Fitness, Medicine, Alternative...	 Home Family, Consumers, Cooking...
 News Media, Newspapers, Weather...	 Recreation Travel, Food, Outdoors, Humor...	 Reference Maps, Education, Libraries...
 Regional US, Canada, UK, Europe...	 Science Biology, Psychology, Physics...	 Shopping Clothing, Food, Gifts...

Think about classifying web pages into the above categories.
Why is it a misnomer to call this a “taxonomy”?

Classification

Curlie.org is the successor of the once-famous DMOZ directory of web pages.



The screenshot shows the Curlie.org homepage. At the top left is the red "Curlie" logo with a white outline. To its right are links for "About", "Become an Editor", "Help", "Forums", and "Login". Below the logo is a cartoon squirrel waving. Next to the squirrel is the text "Hi, Look at the sites I collected!" and "It's time to collect more!". A red search bar with a white input field containing "Search Curlie" and a magnifying glass icon is centered below the squirrel. Below the search bar are nine category cards arranged in a grid. Each card has an icon on the left and a title and description on the right. The categories are: Arts (Movies, Television, Music...), Business (Jobs, Real Estate, Investing...), Computers (Internet, Software, Hardware...), Games (Video Games, RPGs, Gambling...), Health (Fitness, Medicine, Alternative...), Home (Family, Consumers, Cooking...), News (Media, Newspapers, Weather...), Recreation (Travel, Food, Outdoors, Humor...), Reference (Maps, Education, Libraries...), Regional (US, Canada, UK, Europe...), Science (Biology, Psychology, Physics...), and Shopping (Clothing, Food, Gifts...).

 Arts Movies, Television, Music...	 Business Jobs, Real Estate, Investing...	 Computers Internet, Software, Hardware...
 Games Video Games, RPGs, Gambling...	 Health Fitness, Medicine, Alternative...	 Home Family, Consumers, Cooking...
 News Media, Newspapers, Weather...	 Recreation Travel, Food, Outdoors, Humor...	 Reference Maps, Education, Libraries...
 Regional US, Canada, UK, Europe...	 Science Biology, Psychology, Physics...	 Shopping Clothing, Food, Gifts...

Why is it a misnomer to call this a “taxonomy”?
Because there will be plenty of cases of dual classification.
In a real taxonomy, every item is classified into one and only category path.

Classification

Q: What are some other big (famous, influential) document taxonomies on the web?

?

Classification

Q: What are some other big (famous, influential) document taxonomies on the web?

A: IMDB, IAB, About.com, the Google Product Taxonomy, and many website taxonomies such as Amazon, CNET, etc.

Other Examples

Additional types
of semantic
analysis include:

- Tagging
- Topic segmentation
- Sentiment analysis

Topic:
Knee
injuries

Topic:
Rehab

Tags: sports, injuries, coaching, soccer

Lorem ipsum dolor sit amet, consectetur adipiscing elit. posuere tortor vitae elit. Sed vitae metus a elit bibendum malesuada cras pulvinar. Quisque pellentesque nibh in sem. Curabitur ligula. Suspendisse potenti. Duis sit amet augue eu arcu ultrices auctor. Suspendisse elementum, nunc ut molestie elementum, neque augue vulputate elit, eu blandit enim velit vitae nulla. Duis sed.

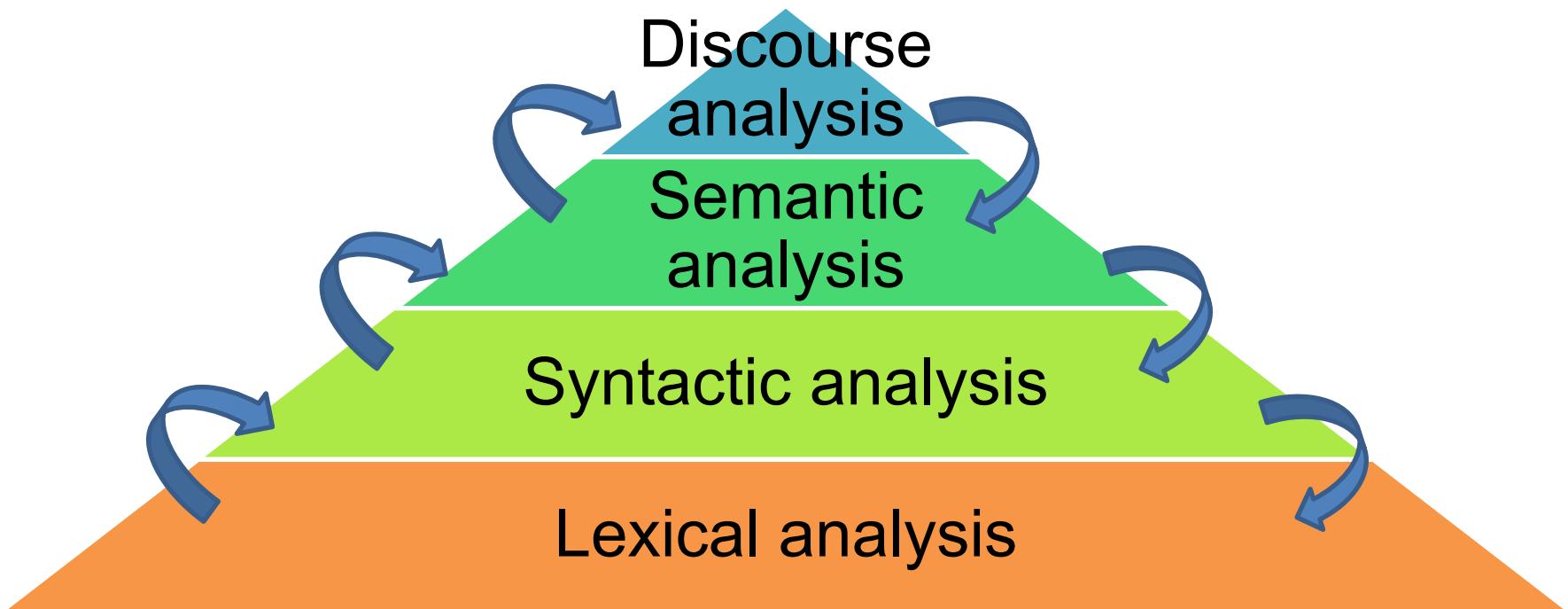
DataScience@SMU

Overview: Discourse Analysis in NLP

Natural Language Processing

Levels of Analysis

For convenience, we can think of there being higher and lower levels of analysis in NLP as long as we remember they are interconnected.



Discourse Analysis

Some major types:

- Anaphora resolution
- Discourse modeling
- Question answering
- Textual entailment
- Pragmatic analysis

Anaphora Resolution

- A pronoun always refers to its most recent antecedent, right?
- No. Pronouns don't always behave that way

“Jim told Bob he would give him the quarterly report next Monday.”

Who is “he” and
who is “him”?

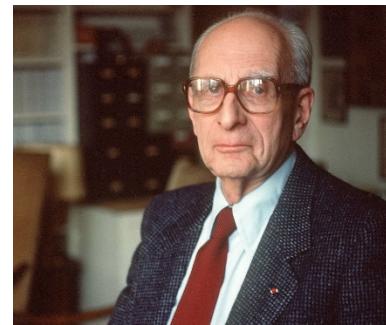
Discourse Modeling

Waiter: Introduce.Self.Name
Waiter: Introduce.Self.Role
Waiter: Offer.Beverage
Customer: Order.Beverage
Waiter: Offer.Appetizer
Customer: Order.Appetizer
Etc.

Roger Schank's work on "scripts" in the '70s...

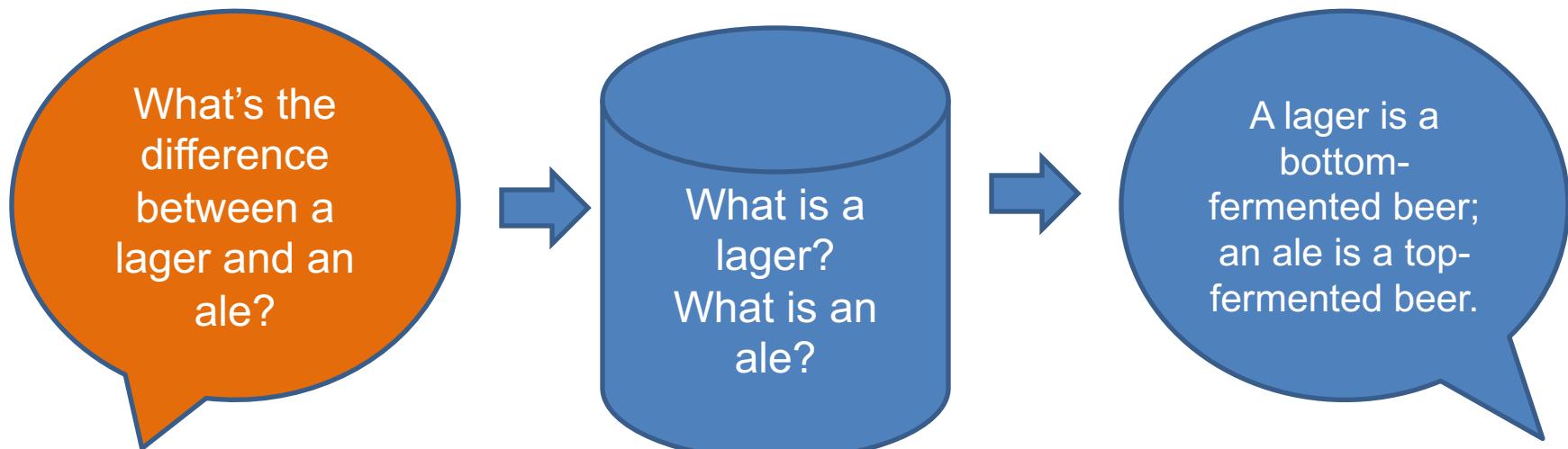


...was anticipated by Claude Lévi-Strauss's "structuralism" in the '40s.



Question Answering

- Most straightforward approach is matching preexisting Q-A materials
- The matching is often not straightforward. We can build models that adapt question formats.



We may not have the exact question in our db but can combine other questions (with answers) to make a response.

Applying Inference to Text

- Textual entailment

“All men are mortal.”

“Socrates is a man.”

Inference: Socrates is mortal.

This inference follows from strict logic.

- Pragmatic analysis

“Did you like the play?”

“It was, um, interesting.”

Inference: He didn’t like the play.

This inference **doesn’t** follow from strict logic.

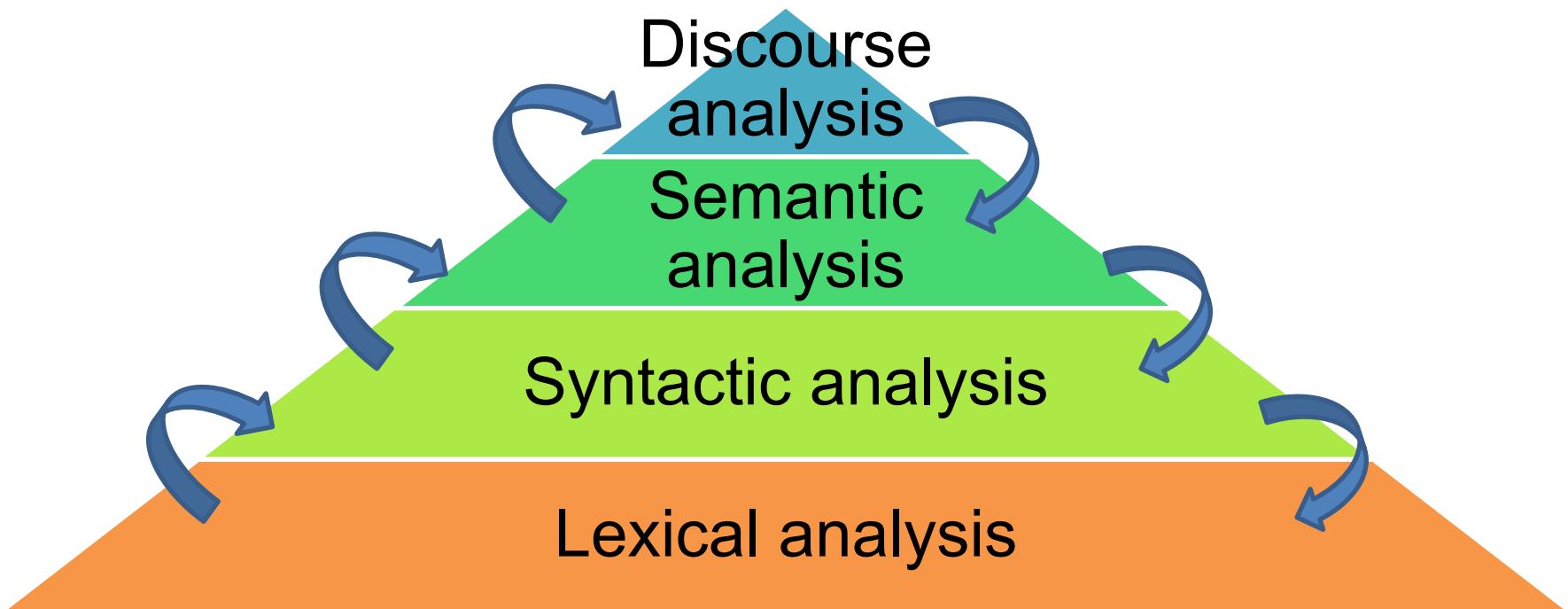
DataScience@SMU

Review: Levels of Analysis and Applications in NLP

Natural Language Processing

Levels of Analysis

For convenience, we can think of there being higher and lower levels of analysis in NLP as long as we remember they are interconnected.



Levels of Analysis and Implementations

Generally, we use lower levels also when we move to “higher” levels of analysis.

This table just shows *typical* levels of analysis employed in NLP.

One can often enhance results by adding a higher level than is strictly necessary, e.g., adding syntactic analysis will improve spell correction.

Minimum level of analysis generally used for implementations	Lexical Analysis	Syntactic Analysis	Semantic Analysis	Discourse Analysis
Stemming	X			
Keyword Tagging	X			
Spell Correction	X			
Terminology Extraction	X			
Reading Level Estimation	X			
Sentence-Boundary Detection	X	X		
POS-Tagging	X	X		
Parsing	X	X		
Lemmatization	X	X		
Unitizing/Normalizing	X	X		
NEE	X	X	X	
WSD	X	X	X	
Conceptual Tagging	X	X	X	
Classification	X	X	X	
Topic Segmentation	X	X	X	
Sentiment Analysis	X	X	X	
Anaphora Resolution	X	X	X	X
Discourse Modeling	X	X	X	X
Question Answering	X	X	X	X
Textual Entailment	X	X	X	X
Pragmatic Analysis	X	X	X	X

DataScience@SMU