

# Lecture 8

# Corpus analysis and computational linguistics

- © Kushneruk Svetlana Leonidovna
- Doctor of Philology, Professor of
- Chelyabinsk State University

# 1. Corpus analyses of vocabulary

a set of lexical forms  
having the **same stem** and  
belonging to the **same**  
major **word class**, differing  
only in inflection and/or  
spelling



ACE  
MON  
STELLAR  
SOPHAGUS  
SPEED  
ALGEBRA  
ARTISTIC  
CONTOUR  
ARTS  
BUILDINGS  
SKV  
DT  
A  
SMA  
BREAM  
NAMED  
GREAT  
BERMUDA  
AQUA  
CLARITY  
WORDS  
CONCERNING  
ENGLAND  
DEFINITION  
AWARDS  
FARM  
STINGRAY  
NORFOLK  
WIND  
LIPIDS  
NORFOLK  
ARTS  
NORFOLK  
RATHON  
ARENA  
LATERAL  
ICAL  
SCALES  
LOWER  
NIGHT  
THEOREMS  
AZURE  
BENT  
BRIGHT  
BILAYER  
VEGAS  
CAGE  
FORMULAS  
MATHEMATICS  
SAND  
BODY  
BLOSSOMS

a base form of a word + its inflected forms + derivatives

- **LEMMA** : inflected forms of the verb
  - *develops, developed, developing*
- **WORD FAMILY** additionally includes representatives of other word classes
  - adjectives: *undeveloped, underdeveloped*
  - nouns: *development, developments, developer, developers*



# 1.1. Polysemy



**Polysemy** – the existence of several meanings in an individual word



**The more frequent a word is, the more meanings it is likely to have**

the polysemy of the word *bear*

## Corpus of Contemporary American English

SEARCH

FREQUENCY

[List](#) [Chart](#) [Word](#) [Browse](#) +

bear\_v\*

[POS]?

Find matching strings

Reset

Sections [Texts/Virtual](#) [Sort/Limit](#) [Options](#)

Corpus of Contemporary American English

SEARCH FREQUENCY CONTEXT

ON CLICK: [CONTEXT](#) [TRANSLATE \(??\)](#) [GOOGLE](#) [IMAGE](#) [PRON/VIDEO](#) [BOOK](#) (HELP)

HELP ALL FORMS (SAMPLE): 100 200 500 FREQ

1	<input type="checkbox"/>	BEAR	27665
---	--------------------------	------	-------

Corpus of Contemporary American English

SEARCH FREQUENCY CONTEXT ACCOUNT

CLICK FOR MORE CONTEXT NEW SAVE TRANSLATE ANALYZE

ID	Year	Type	Source	Actions	Text
1	2012	BLOG	addinginfo.org		, my friends, if you are (wo) man enough to grin and bear it, can be very enlightening. # *rimshot* Try the veal, I'll
2	2012	BLOG	dailypaul.com		Here's a little help researching the Bronfmans: http: **35;1818;TOOLONG... # And bear in mind, many of the really big names go to great effort to make
3	2012	BLOG	...tionalgeographic.com		, we generally take precautions to avoid interactions or attacks (hiking in groups, bear spray, a gun, etc). However, in other places in the
4	2012	BLOG	...tionalgeographic.com		comparison, during the 20th century there have been 71 fatal grizzly (brown) bear attacks in North America. Each year in the United States, 16-18 people d
5	2012	BLOG	...uality.wordpress.com		pointless discrimination hurts all of us, from those impacted directly to those helping to bear the burden of the societal costs. # As much as Mitt Romney a
6	2012	BLOG	freethoughtblogs.com		God's. Men and women are made in the image of God, they bear his image... so therefore they owe a something to God. That's the
8	2012	BLOG	core77.com (1)		probably for legal reasons, so I will not quote anyone directly. Also bear in mind I'm going off of an anecdote I heard only once, about
9	2012	BLOG	...mbreen.wordpress.com		out concerns, even if we do not see the risks ourselves. We must bear in mind that perpetrators are often charming, sociable characters and just because
10	2012	BLOG	...ithchip.blogspot.com		are " traders, like everyone else in a free society, and they should bear that title proudly, considering the crucial importance of the services they offer. "
11	2012	BLOG	...ocetah.wordpress.com		something like " Yanaweti " or " Yonaweti, " from " yonah ", bear, and " uweti, " old. The actual pronunciation would have been something
12	2012	BLOG	...esternjournalism.com		# We have an Ace-in-the-Hole. Just one. He has no credentials that can bear scrutiny. His " sealed " past has to be Un-Sealed by a military force
13	2012	BLOG	...nstructionlitmag.com		Knicks. If his early returns -- his signature three-pointers, a positive attitude -- bear out over the course of the season, he may end up usurping the title
14	2012	BLOG	blackwindmetal.com		feelings I got from each album; The Time Of The Oath just seems to bear the mark of a confident band. Again, it's not all too different
15	2012	BLOG	blog.adw.org		adjustment. # Some of us who got married were simply unable to conceive and bear children. I don't think that absolved us from the sacramental respons
16	2012	BLOG	blog.adw.org		their love can be a sacrament of Christ's love for the Church which can bear fruit in other ways. I know you said May's argument was directed at
17	2012	BLOG	blog.adw.org		the form it does does not mean that couples who for some reason can not bear a child are excluded. It only means that they can not purposefully, or
18	2012	BLOG	...le-east.blogspot.com		US State Department should ask the Congress to allocate funds so that these efforts could bear fruition. While the term has been overused, the need to w
19	2012	BLOG	joyfullyjobless.com		the home office gallery gathered by Judy Heminsley. You'll see wonderful environments that bear no resemblance to a cubicle. 3 Responses to A Sense of I
20	2012	BLOG	...yabsoff.blogspot.com		interest in kissing the un-cute boy in front of me. # i couldn't bear to let the opportunity pass me by altogether, so i turned to my best

## 1.2. Synonymy as a lexical relation



**Synonymy** – a bilateral or symmetrical sense relation in which two or more linguistic forms have similar meaning: *issue* and *problem*



no two words can be considered perfect synonyms as **corpus data** reveal *important differences* in the phraseological patterns



SEARCH

WORD

CONTEXT

ANALYZE TEXT

-SAMPLES- MY TEXTS

SEARCH CLEAR

HELP	<a href="#">Compare to previous WordAndPhrase (PDF)</a>
	<a href="#">Overview</a>
	<a href="#">Word-oriented functions</a>
	<a href="#">Phrase-oriented functions</a>



You can always get to this page by clicking this icon above.

You can enter any text that you would like in the form at the left -- for example, a paper that you've written, or a newspaper article that you've copied from another website. After inputting the text, you can then see useful information about words and phrases in that text, based on data from COCA.

First, it will highlight all of the medium and lower-frequency words in your text and create lists of these words that you can use offline. This frequency data can help language learners focus on new words, and it can allow you to see "what the text is about" (i.e. text-specific words).

Second, you can click through the words in the text to see a detailed "word sketch" of any of the words -- showing their definition and their translation (in more than 100 languages); links to pronunciation, images, and videos; related topics, collocates, "clusters" (2, 3, and 4-word phrases); and concordance lines.

Finally, you can do powerful searches on selected phrases in your text, to show related phrases in COCA. In this way, this resource is like a "collocational thesaurus" to see what related phrases are most likely in different styles of English.

Just enter some text (or analyze a random text via SAMPLES), and there will be more help files on the next page.



# TRUMP inaugural address

Corpus of Contemporary American English

BROWSE/RANDOM      WORD      CONTEXT      ANALYZE TEXT

EDIT TEXT	SAVE TEXT	● WORD	● PHRASE
FREQ RANGE	1-500	501-3000	> 3000
1466 WORDS	65 %	11 %	10 %

CLICK ON ANY WORD BELOW FOR A FULL WORD SKETCH

Chief Justice Roberts, President Carter, President Clinton, President Bush, President Obama, fellow Americans, and people of the world, thank you. We the citizens of America are now joined in a great national effort to rebuild our country and restore its promise for all of our people. Together we will determine the course of America, and the world, for many, many years to come. We will face challenges. We will confront hardships, but we will get the job done. Every four years, we gather on these steps to carry out the orderly and peaceful transfer of power, and we are grateful to President Obama and First Lady Michelle Obama for their gracious aid throughout this transition. They have been magnificent. Thank you. Today's ceremony, however, has very special meaning, because today we are not merely transferring power from one administration to another, or from one party to another, but we are transferring power from Washington, D.C., and giving it back to you, the people. For too long, a small group in our nation's capital has reaped the rewards of government, while the people have borne the cost. Washington flourished, but the people did not share in its wealth. Politicians prospered, but the jobs left and the factories closed. The establishment protected itself, but not the citizens of our country. Their victories have not been your victories. Their triumphs have not been your triumphs, and while they celebrated in our nation's capital, there was little to celebrate for struggling families all across our land. That all changes, starting right here and right now, because this moment is your moment -- it belongs to you. It belongs to everyone gathered here today, and everyone watching, all across America. This is your day. This is your celebration, and this is the United States of

(CLICK ANY WORD FOR FULL WORD SKETCH)

LOW FREQ	MID FREQ	HIGH FREQ
4: wealth 3: bless, factories 2: allegiance, destiny, glorious, loyalty, oath, righteous, shine, trillions, triumphs 1: alliances, almighty, assembled, bedrock, bleed, cannot, carnage, celebrate, celebrated, celebration, ceremony, civilized, complaining, confidence, confront, constantly, conviction, courage, creator, crucial, decay, decree, depletion, deprived, disagreements, disrepair, dissipated, enforcement, enriched, eradicate, establishment, expense, flag, flourished, flush, forever, friendship, gangs, goodness, goodwill, govern, gracious, grateful, hardships, harness, heal, highways, historic,	5: dreams, protected 4: citizens, everyone, heart 3: borders, capital, foreign 2: action, anyone, belongs, breath, exists, fight, forgotten, forward, industry, mountain, movement, ocean, pain, politicians, safe, seek, share, strength, success, transferring, victories, workers 1: accept, administration, affairs, aid, airports, armies, arrives, beautiful, benefit, birth, blood, born, borne, bridges, brown, carry, cash, challenge, challenges, chief, completely, crime, debate, decades, decision, defend, defended, define, demands, destroying, determine, disease,	75: and 73: the 50: we 49: our 48: of 43: will 37: to 20: is 15: a, for, you 14: all, in 13: be, but 12: are 11: from, it, their, your 10: not, people, that, this 9: again, country, nation 8: one, with 7: every 6: across, at, back, by, great, has, make, never, new, no, now, on, same, while, world 5: have, many, other, right, they, today 4: as, been, bring, day,

EDIT TEXT SAVE TEXT  WORD  PHRASE

FREQ RANGE	1-500	501-3000	> 3000
1 WORDS	100 %	0 %	0 %

(CLICK ANY WORD FOR FULL WORD SKETCH)

LOW FREQ	MID FREQ	HIGH FREQ
		1: problem

CLICK ON ANY WORD BELOW FOR A FULL WORD SKETCH

[problem](#)

[See in iWeb](#) [🏠](#) [Collocates](#) [Clusters](#) [Topics](#) [Texts](#) [KWIC](#) [⬇️](#) [HELP](#)

**problem** (NOUN) ★ 🔄 #187 +



1. a state of difficulty that needs to be resolved 2. a source of difficulty 3. a question raised for consideration or solution

D M O C G **E**

[YouGlish](#) [PlayPhrase](#) [Yarn](#)  
[Translate: choose language](#)

SYNONYMS NEW: DEFIN +SPEC +GENL

**difficulty** difficulty, drawback, glitch, hindrance, hitch, obstacle, obstruction, problem, setback, snag **puzzle** conundrum, enigma, poser, problem, puzzle, question, riddle

TOPICS (more)

[solve](#), [solution](#), [eg](#), [student](#), [solving](#), [fix](#), [ie](#), [symptom](#), [treatment](#), [cause](#), [variable](#), [difficulty](#), [mathematical](#), [computer](#), [correlate](#), [behavior](#), [condition](#), [finding](#), [classroom](#), [score](#)

COLLOCATES (more)

**NOUN** [solution](#), [solving](#), [america](#), [us](#), [alcohol](#), [drinking](#), [obama](#), [pollution](#)

**VERB** [solve](#), [cause](#), [face](#), [address](#), [fix](#), [resolve](#), [arise](#), [pose](#)

**ADJ** [serious](#), [environmental](#), [mental](#), [severe](#), [fundamental](#), [behavioral](#), [psychological](#), [underlying](#)

**ADV** [eg](#), [ie](#), [therein](#), [head-on](#), [twofold](#)

RELATED WORDS

[problematic](#), [problem-solving](#), [problem-solver](#), [problem-based](#), [problem-solve](#), [problematical](#), [problem-free](#), [unproblematic](#), [problematically](#), [problematic](#), [sub-problem](#)

# Information about clusters, text types and concordance lines

Corpus of Contemporary American English

SEARCH WORD CONTEXT ANALYZE TEXT

## CLUSTERS (more)

problem •	problem with • problems with • problems in • problem in • problem for • problem here • problems for • problem solving
• problem	health problems • real problem • big problem • only problem • have problems • one problem • biggest problem • serious problem
problem ••	problem is not • problem with it • problems associated with • problem is n't • problems such as • problem that we • problems that we • problem at all
•• problem	have a problem • to the problem • not a problem • have no problem • got a problem • solve the problem • had a problem • with the problem
problem •••	problem is that we • problem is that it • problem is that you • problem is that there • problem is that i • problem is that they • problem i have with • problem is that he
••• problem	part of the problem • you have a problem • n't have a problem • one of the problems • we have a problem • i have no problem • there is a problem • i have a problem

























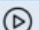



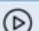









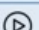


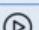



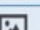




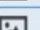


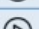
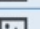


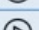
## TEXTS / VIRTUAL CORPORA (more)

ACAD:LearningDisability • ACAD:Canadian J Exper Psychology • WEB:wisegEEK.com • ACAD:EducTreatmenChildren • BLOG:blogs.suntimes.com • WEB:ilto.wordpress.com • BLOG:blog.mrmeyer.com • ACAD:DrugIssues • BLOG:wattsupwiththat.com • WEB:chronicle.com • ACAD:EducTreatmenChildren • WEB:supermemo.com • ACAD:AI Magazine • WEB:hanselman.com • BLOG:...square.blogs.cnn.com • ACAD:ExceptionalChildren • ACAD:EmotBehavDis • WEB:slideshare.net • ACAD:LearningDisability • ACAD:SchoolPsych • BLOG:feministe.us • WEB:freethoughtblogs.com • WEB:history.navy.mil • BLOG:control.com • BLOG:freethoughtblogs.com •

## CONCORDANCE LINES (more)

1	MAG: 1995: TotalHealth	11582 Most adults suffer from acute low back problems at some time in their lives . Many patients traditionally have
2	ACAD: 2001: SchoolPsych	two intervals divided by the total number of intervals of problem behavior and then multiplied by 100% . # Hypothesis
3	ACAD: 2001: SchoolPsych	accuracy percentage was reported to be approximately 70% . Problem behavior included inappropriate teacher engagement and
4	ACAD: 2001: SchoolPsych	environmental variables that may control or maintain problem behavior . Although ABA approaches are extremely useful in
5	ACAD: 2001: SchoolPsych	analysis outcomes also confirmed the hypothesis for Sal 's problem behavior . During low peer attention conditions , Sal was more

# SYNONYMS

	RANK	FREQ	Word	PoS	Audio	Video	Image	LANG?
1	187	504175	problem	NOUN				
2	206	457301	question	NOUN				
3	972	99455	challenge	NOUN				
4	2039	43423	difficulty	NOUN				
5	4099	16342	obstacle	NOUN				
6	4501	14181	catch	NOUN				
7	4652	13476	puzzle	NOUN				
8	8052	5407	setback	NOUN				
9	8292	5126	obstruction	NOUN				
10	10043	3689	drawback	NOUN				
11	11643	2850	glitch	NOUN				
12	12058	2688	riddle	NOUN				
13	14198	1969	conundrum	NOUN				
14	14667	1841	hitch	NOUN				
15	15413	1669	enigma	NOUN				
16	18376	1171	snag	NOUN				
17	18646	1133	hindrance	NOUN				
18	29329	392	poser	NOUN				

# issue

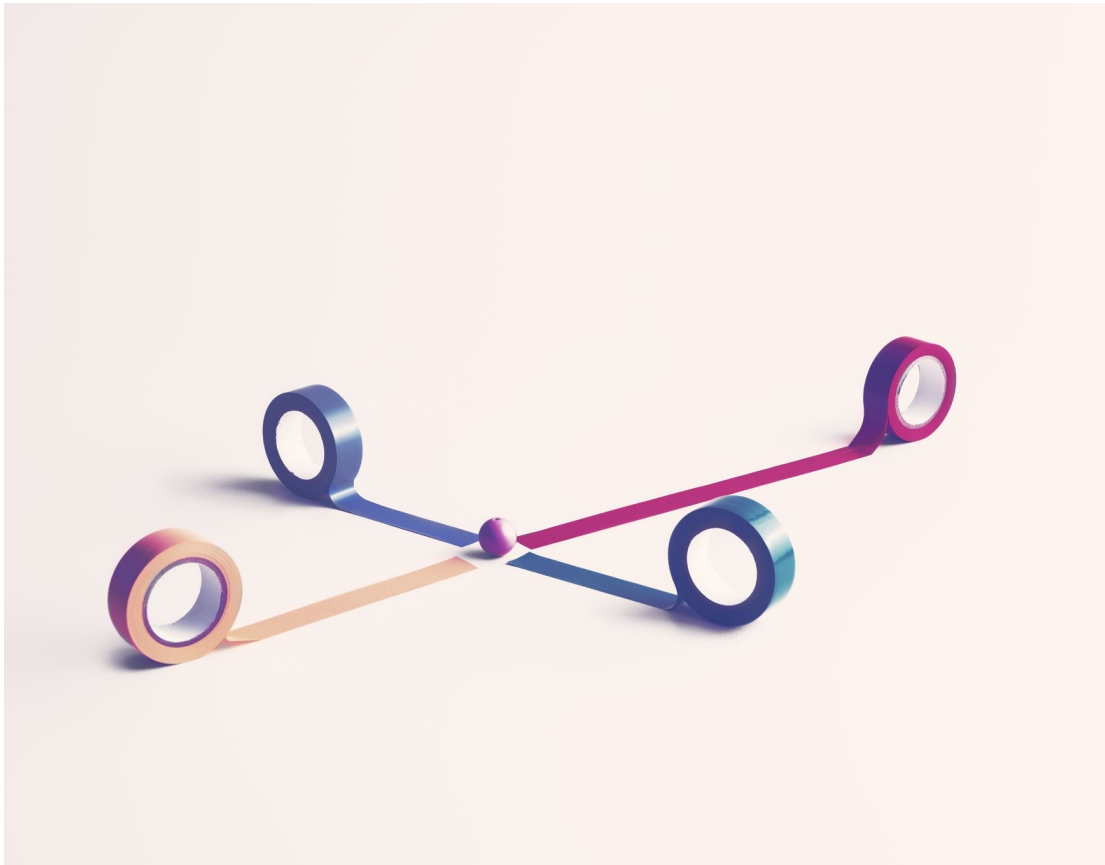
1	187	504175	problem	NOUN				
2	206	457301	question	NOUN				
3	211	444412	number	NOUN				
4	261	362137	issue	NOUN				
5	575	167351	matter	NOUN				
6	803	120773	subject	NOUN				
7	876	110976	concern	NOUN				
8	1037	93507	production	NOUN				
9	1610	58410	topic	NOUN				
10	1865	48365	copy	NOUN				
11	1878	47967	release	NOUN				
12	2322	36338	distribution	NOUN				
13	2394	34656	publication	NOUN				
14	2494	32842	edition	NOUN				
15	2750	28903	delivery	NOUN				
16	3369	22027	dispute	NOUN				
17	5891	9267	circulation	NOUN				
18	6980	6884	offspring	NOUN				
19	9082	4422	installment	NOUN				
20	16233	1512	issuance	NOUN				

# *problem*

	RANK	FREQ	Word
1	187	504175	problem
2	206	457301	question
3	972	99455	challenge
4	2039	43423	difficulty
5	4099	16342	obstacle
6	4501	14181	catch
7	4652	13476	puzzle
8	8052	5407	setback
9	8292	5126	obstruction
10	10043	3689	drawback
11	11643	2850	glitch
12	12058	2688	riddle
13	14198	1969	conundrum
14	14667	1841	hitch
15	15413	1669	enigma
16	18376	1171	snag
17	18646	1133	hindrance
18	29329	392	poser

# *issue*

1	187	504175	problem
2	206	457301	question
3	211	444412	number
4	261	362137	issue
5	575	167351	matter
6	803	120773	subject
7	876	110976	concern
8	1037	93507	production
9	1610	58410	topic
10	1865	48365	copy
11	1878	47967	release
12	2322	36338	distribution
13	2394	34656	publication
14	2494	32842	edition
15	2750	28903	delivery
16	3369	22027	dispute
17	5891	9267	circulation
18	6980	6884	offspring
19	9082	4422	installment
20	16233	1512	issuance



## 1.3. Metaphoricity and idiomaticity of words

**Metaphoric meaning** is a type of extended meaning when words are used in an extended non-literal way to explain something in a comparative relationship


**hook, line and sinker**

*I fell for his story hook, line and sinker.*

SEARCH		FREQUENCY		CONTEXT		OVERVIEW	
<a href="#">CLICK FOR MORE CONTEXT</a>				<span>NEW</span>		<span>SAVE</span> <span>TRANSLATE</span> <span>ANALYZE</span>	
1	KCX	S_conv	🔍	🔍	Q	Timeshare. (SP:PS1FC) Oh (pause) no! She falls for it I-- hook line and <b>sinker</b> . Then I forget that I've been joking with her (pause) then she tells	
2	KP1	S_conv	🔍	🔍	Q	lot have pulled on us and everybody's fell for it hook, line and <b>sinker</b> (SP:PS50U) (unclear) (SP:PS50T) ha (SP:PS50U) the first year was horrible but last year was (	
3	ARK	W_fict_prose	🔍	🔍	Q	chair.' Special Branch bought it, swallowed it -- hook, line and <b>sinker</b> ." And Tweed?' his deputy enquired. Gareth Morgan was forty-two	
4	FR3	W_fict_prose	🔍	🔍	Q	, to become just another off-the-peg person dangling on the idiomatic hook (line and <b>sinker</b> ), my voting habits purely a function of minute alterations in fiscal pc	
5	G02	W_fict_prose	🔍	🔍	Q	lavatory seat is up, maroon bowl stained white, unflushed paper, a dark <b>sinker</b> lurking like a sub in the U-bend.' Baby barf NOW, I	
6	HTG	W_fict_prose	🔍	🔍	Q	but it could only be as chairman, they swallowed it hook, line and <b>sinker</b> . Glastonbury was out. I don't think he even noticed. I was	
7	JY5	W_fict_prose	🔍	🔍	Q	had no one but herself to blame if she'd fallen hook, line and <b>sinker</b> in love. She slid him a sideways glance, only to be all but	
8	K1Y	W_news_script	🔍	🔍	Q	The winning line.Ben's a scrambling champ at six. And, hook line and <b>sinker</b> . The bait that Neville Chamberlain was happy to fall for. Change the law	
9	AJA	W_newsp_brdshst_nat_misc	🔍	🔍	Q	nodding off during the sermon with a smart crack on the head with a lead <b>sinker</b> . But few came odder than Joseph Neeld. Up to the age of 39	
10	AJY	W_newsp_brdshst_nat_misc	🔍	🔍	Q	language: addicts, paramours and dupes of all kinds fall hook, line and <b>sinker</b> to the cunning and the dangerous. The hook is the epitome of angling,	
11	A1Y	W_newsp_brdshst_nat_report	🔍	🔍	Q	who doubts.' They just wanted me to swallow it hook, line and <b>sinker</b> ,' a disillusioned student said. Society leaders argue that their single issue approach	
12	CF9	W_newsp_other_sports	🔍	🔍	Q	a wind-up.National and local papers, television and radio crews fell hook, line and <b>sinker</b> , for the 29-year-old Sizewell B steel erector's story, but the EADT decide	
13	A6R	W_pop_lore	🔍	🔍	Q	my broly as a pout started to rattle my rod tip when a 5 oz <b>sinker</b> ripped a hole through my broly. The trace from a chap 20 yards up	
14	CK3	W_pop_lore	🔍	🔍	Q	the ribber. The colour changers will fit standard and find gauge machines. The <b>sinker</b> plate supplied with the YC6 for single bed knitting has wheels and brushes	
15	CK3	W_pop_lore	🔍	🔍	Q	with the YC6 for single bed knitting has wheels and brushes beneath it like any <b>sinker</b> plate. To use it on a fin gauge machine, the cog wheels have	
16	CK3	W_pop_lore	🔍	🔍	Q	replaced with ones to fit the fine gauge. (The cogs must fit the <b>sinker</b> posts which are closer together on a fine gauge than on a standard of course	
17	CK3	W_pop_lore	🔍	🔍	Q	and slip stripes, the yarn goes in feeder 1 of the YC6 colour changer <b>sinker</b> plate. If you have the YC5 model, as we saw last month,	
18	CK3	W_pop_lore	🔍	🔍	Q	knitting single bed Fair Isle. The second yarn goes in feeder two of the <b>sinker</b> plate, while it is the yarn in feeder one that changes. As with	
19	CK3	W_pop_lore	🔍	🔍	Q	while it is the yarn in feeder one that changes. As with the standard <b>sinker</b> plate, it is the yarn in feeder two that knits the pattern and the	

- British National Corpus (BNC)
- concordances for the word '*sinker*'





quantitative analysis of frequencies of words +  
interpretative discourse-oriented analysis

Computers and concordancers are not able to identify idiomatic meanings of words and therefore ***automatic searches need to be combined with qualitative forms of analysis*** based on the reading of concordance lines

---

## 1.4. Register variation: jargon, slang and appropriateness

### Register

it focuses on the varied use of language as dependent on factors such as communicative situation, formality level, the age and gender of language users and relationships between them



# Corpus of Contemporary American English

SEARCH

ERROR

List Chart Word Browse +



Sections Texts/Virtual Sort/Limit Options

The most frequent nouns across different registers in **COCA** corpus

# Corpus of Contemporary American English



SEARCH

FREQUENCY

CONTEXT

OVERVIEW

[CONTEXT](#)
[TRANSLATE \(RU\)](#)
[GOOGLE](#)
[IMAGE](#)
[PRON/VIDEO](#)
[BOOK](#)
[\(HELP\)](#)

	ALL	BLOG	WEB-GENL	TV/MOVIES	SPOKEN	FICTION	MAGAZINE	NEWSPAPER	ACADEMIC	1990-1994	1995-1999	2000-2004	2005
PEOPLE	1707619	318038	267974	138749	415659	101645	167611	198660	99283	192320	188237	181568	183
TIME	1601568	236428	228079	232871	218237	200699	185038	162623	137593	179867	195515	191206	189
WAY	1066264	163441	148363	165626	161449	152288	116071	91908	67118	122389	130210	126490	125
YEARS	1017103	136505	131781	76935	153612	88610	155761	178846	95053	129462	129993	128160	123
MAN	665748	55141	73797	185470	85148	150487	47310	47961	20434	98151	99316	89968	87
LIFE	679447	95869	103849	100139	83905	82242	87317	62671	63455	78489	86179	81956	81
DAY	665337	90617	89513	97742	96838	97284	84429	79557	29357	74005	81782	83870	84
WORLD	678541	105582	114646	53575	77455	53124	101784	86059	86316	81855	79182	84808	72
YEAR	648697	100126	91629	7425	87036	34127	108310	172353	47691	79375	78144	75755	74
THINGS	604458	109431	94568	79775	127560	69906	54095	49673	19450	62327	69004	65640	66
THING	551285	84190	65234	135318	117139	70333	41778	37293		59964	68659	66050	67
SCHOOL	474568	53031	53873	26507	41375	40753	38894	105973	114162	49646	54585	61015	69
HOUSE	432152	42459	49826	55297	92884	78693	45936	67057		54023	59649	53596	54
PRESIDENT	438560	51142	60278		207258		32663	85140	2079	58045	49284	48594	38
FAMILY	411396	47718	50309	44738	62904	40637	58567	63595	42928	43740	53199	52038	57
NIGHT	392940	34053	35727	92729	59913	94881	36006	39631		56775	57705	54594	51
CHILDREN	395055	45461	51627		54401	36547	54207	59527	93285	55908	55403	48927	47
MONEY	408460	72084	57751	65332	70701	34657	44085	61082	2768	51111	51538	45486	47
STATE	414892	67617	78162		48523		35259	116024	69307	44949	39806	40867	38

ON CLICK: [CONTEXT](#) [TRANSLATE \(RU\)](#) [GOOGLE](#) [IMAGE](#) [PRON/VIDEO](#) [BOOK](#) (HELP)

HELP		ALL	BLOG	WEB-GENL	TV/MOVIES	SPOKEN	FICTION	MAGAZINE	NEWSPAPER	ACADEMIC	1990-1994	1995-1999	2000-2004	2005-2009	2010-2014	2015-2019	
1	<input type="checkbox"/>	FLUKE	2383	582	502	216	172	223	242	357	89	179	284	236	198	247	155

## Informal words in a corpus


*fluke* – an unlikely chance occurrence, especially a surprising piece of luck





















EN (172)

100

/ 2 > >>






















FOR MORE CONTEXT

NEW 

OK	ABC_20/20			Q	sees this as just almost divine. White tigers were rare, kind of a fluke in nature. NARRATOR-MALE) : All the v
OK	CNN_Newsroom			Q	added only 20,000 jobs. Now we will find out whether that was just a fluke or the beginning of a hiring slow
OK	NPR_Saturday			Q	, I have to say that's - in a way, it's a fluke. But of course, I've been thinking about the issues that it discusse
OK	NPR_Saturday			Q	. And they're going to make life difficult. But this is not a fluke because in any sort a tournament, you can w
OK	PBS_Newshour			Q	limits to being anti-Trump, second, that the Trump phenomenon was not just a fluke, that it's based on so
OK	Fox: The Five			Q	, President Obama just releases one. And you thought Bill Ayers was just a fluke. The lesson is when the lef
OK	CBS: This Morning			Q	it was the second nomination. And I thought, maybe it's not a fluke. GAYLE-KING# Mm-Hm. JULIA-ROBERTS
OK	ABC: 20/20			Q	. Teacher said my hands were too small. So really, it was a fluke that I ended up with the violin. AMY-ROBA
OK	PBS_Newshour			Q	detected a gravitational wave, and it helps prove the first time wasn't a fluke. JUDY-WOODRUFF# Still to co
OK	PBS_Newshour			Q	Bernie Sanders campaign really hopes tomorrow that they can show that Michigan was not a fluke. They h

# The spoken section – 172 examples

---

ONTEXT			NEW 
Legl Tropic Diseases			, Min DY, et al. Risk factors for <i>Opisthorchis viverrini</i> and minute intestinal fluke infections in Lao PDR, 2009-20
JNE			dataset, stroking was identified when oscillation on the high-frequency component of surge accelerations ind
bia Law Review			in more enduring factors affecting the long-run success of the firm, rather than to fluke factors unlikely to rep
Gamma			my first article rejection came, I dismissed it. It had to be a fluke. However, when more rejections came, I was
Educ			a well-developed tonal ability, correct notes become the fruit of audiation rather than the fluke of technique.
alPlantSci			hepatica and / or <i>F. gigantica</i> (Soulsby, 1987). The common liver fluke, <i>F. hepatica</i> is a trematode and widely d
alPlantSci			. Cyst wall is dissolved in the gastrointestinal tract of the host and the young fluke emerges. It penetrates and
alPlantSci			of the above description is also applicable to <i>F. gigantica</i> , another species of liver fluke, which is restricted to c
alPlantSci			linked to the presence of the appropriate mater bodies and on adequate climate characteristics enabling fluk
alPlantSci			, 2004). The variations pertaining to the prevalence and seasonal fluctuations of various fluke species in a par

The academic section – 89 examples



SEARCH

FREQUENCY

CONTEXT

OVERVIEW

SECTION: MAGAZINE (242)

FIND SAMPLE: [100](#) [200](#)

PAGE: &lt;&lt; &lt; 1/3 &gt; &gt;&gt;

CLICK FOR MORE CONTEXT

NEW

SAVE

TRANSLATE

ANAL

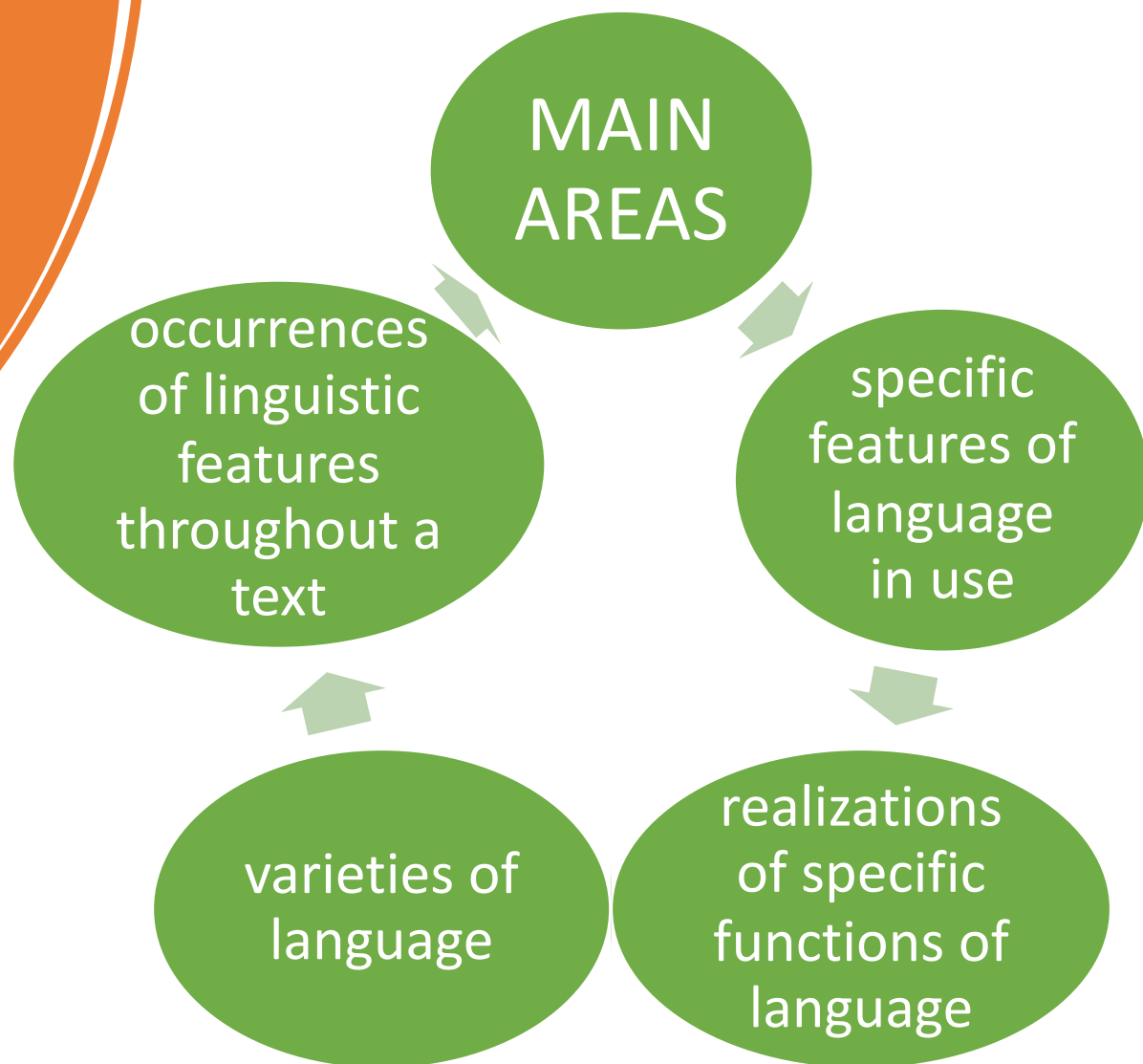
1	2019	MAG	Salon				I did everything I was supposed to do, and still, there was a <b>fluke</b> . I could have died. My baby could have died or been seriously delayed
2	2019	MAG	Bleacher Report				on Team USA in Tokyo to prove the FIBA World Cup slip-up was merely a <b>fluke</b> . 81164284 Terrible News: Stephen Miller Reportedly
3	2019	MAG	Salon				and that had been represented entirely by Republican senators since 1992. Was that a <b>fluke</b> event? # It was not a fluke event. We're going to prove
4	2019	MAG	Salon				senators since 1992. Was that a fluke event? # It was not a <b>fluke</b> event. We're going to prove that in 2020. # Advertisement: #
5	2019	MAG	Salon				the silent dark felt secure enough that this influx of small rodents wasn't a <b>fluke</b> . They also made for tasty eating. # The man could not remember how
6	2019	MAG	ESPN				the radar, as if his 2017 London Marathon victory was nothing more than a <b>fluke</b> . He'd entered the race as the new kid on the block, the
7	2018	MAG	ESPN				. Health is a meaningful skill for NFL quarterbacks. Every passer is susceptible to <b>fluke</b> injuries, and wear and tear eventually slows down every quarterback
8	2018	MAG	ESPN				either side. Tiger has 21, the second most. That's not a <b>fluke</b> , not a small sample size (they've played 83 matches between them)
9	2018	MAG	Sports Illustrated				. The tournament would also want to describe the incident as a tragic but nonetheless <b>fluke</b> accident. This is important because reasonable measures do not
10	2018	MAG	Sports Illustrated				but nonetheless fluke accident. This is important because reasonable measures do not always prevent <b>fluke</b> accidents. Courts, in turn, typically do not hold

# The magazine section – 242 examples

---



## 2. Corpus analysis of discourse



## 2.1. Lexical features of discourse

The notion of **lexical cohesion** is of central importance to the study of discourse. It can be defined as the ways *'in which the components of the surface texts are mutually connected within a sequence'*.



## a. REITERATION



**Reiteration** involves either the direct repetition of a word or the use of a related item that carries the same meaning



synonym, near-synonym or superordinate



'furniture' can be replaced by superordinates  
'item', 'object', 'thing'

[In a shoe-shop: <\$1> is the customer, <\$2> is the assistant]

<\$2> Probably needs adjusting but I'll check that.

<\$1> Oh right. That's lovely.

<\$2> | Okay.

<\$1> | Yeah | that's nice.

<\$2> | They're nice aren't they.

<\$1> Yeah they are nice.

<\$2> Very very nice.

<\$1> Thank you.

<\$2> They feel right?

<\$1> Yeah.

<\$2> Does it?

<\$1> That feels pretty good actually.

<\$2> Yes smashing fit.

<\$1> Yeah. (CANCODE © Cambridge University Press)

**b. RELEXICALIZATION** – a specific form of paraphrasing in which speakers take up one another's vocabulary as they participate in conversations

### c. INTERTEXTUALITY

explores how a discourse makes reference to prior and future discourses

### d. DISCOURSE MARKERS

*'fine', 'good', 'like', 'now', 'okay',  
'right', 'so', 'actually', 'well', 'I mean',  
'you know', 'so to speak', 'in other  
words'*

A **discourse marker** is a word or phrase – *a conjunction, adverbial, comment clause, interjection* – that is uttered with the primary function of *bringing the listener's attention to a particular kind of linkage* of the upcoming utterance with the immediate discourse context

discourse  
marker  
*'so to speak'*



SEARCH

FREQUENCY

List Chart Word Browse +

so to speak [POS]?

Find matching strings Reset

Sections Texts/Virtual Sort/Limit Options

1	IGNORE ----- TV/MOVIES BLOG WEB-GENL SPOKEN FICTION MAGAZINE NEWSPAPER ACADEMIC	2	IGNORE ----- TV/MOVIES BLOG WEB-GENL SPOKEN FICTION MAGAZINE NEWSPAPER ACADEMIC
---	------------------------------------------------------------------------------------------------------------	---	------------------------------------------------------------------------------------------------------------

SORTING RELEVANCE SEC1 : SEC2

MINIMUM FREQUENCY  10 0



---

Exploring how different **forms of vocabulary are configured** over longer stretches of discourse contributes to a better understanding of their role in the creation and maintenance of ***discourse flow***

## 2.2. Exemplary study

---

- Baker, P. (2013). Corpora and discourse analysis, in K. Hyland (ed.) *Discourse Studies Reader: Essential excerpts*. London: Bloomsbury Academic, 11–34.
- **12 holiday leaflets** (17,865 words)

# DISCOURSE STUDIES READER

Essential Excerpts





Using [WordSmith tools](#), Baker started his analysis by producing a **list of the most frequent words** from his holiday corpus and compared them with data from the BNC

*Table 9.1* Examples of the most frequent words in the holiday corpus and BNC expressed in percentages (adapted from Baker 2013)

	<i>Word</i>	<i>% frequency in the holiday corpus</i>	<i>% frequency in the BNC</i>
1	The	5.55	6.20
2	And	3.62	2.68
3	To	2.64	2.66
4	A	2.44	2.21
5	Of	1.96	3.12
6	You	1.95	0.68
7	For	1.38	0.90
8	In	1.37	1.97
9	On	1.15	0.74
10	All	1.04	0.28

# Observations

**1. the personal pronoun 'you' is typically regarded as a feature of spoken language**

- the high frequency of 'you' suggests a 'personal style' of writing, where the writer is directly addressing the reader



**2. the frequency of function words**

- tend to have a high frequency, irrespective of the type of texts analyzed

the most frequent lexical items from

## the holiday corpus

- ‘beach’, ‘pool’, ‘studios’
  - represent the type of language that is used when people describe their holidays
  - analysis of clusters (n-grams) formed around words such as ‘bar’, ‘club’
  - explores the occurrence of verbs as another source of information about the *discourse of tourism*





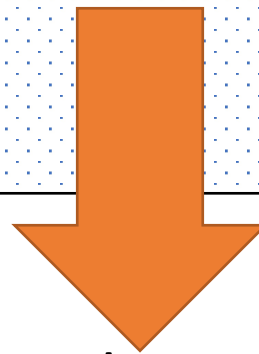
### 3. Computational linguistics and corpus linguistics

- **Computational linguistics** is the field of computer science that looks at how computer systems can be created that work with language in some way

# DATA MINING

identification of patterns and extraction of information across very large datasets

## TEXT MINING



- programs for orthography and grammar correction
- information retrieval from document databases
- translation from one natural language into another

**Computational linguistics** might be considered as a synonym of **automatic processing of natural language**, since the main task of computational linguistics is just the **construction of computer programs to process words and texts in natural language**



# Overlaps

- **Corpus linguistics** is ultimately about **finding out about the nature and usage of language**.
- **Computational linguistics** may also be concerned with **modelling the nature of language computationally**, it is *in addition* focused on **solving technical problems involving language** [McEnery and Hardie 2012: 228].



The screenshot shows the Linguistic Data Consortium (LDC) website. The header includes the LDC logo and the text 'Linguistic Data Consortium'. A navigation menu on the left lists categories such as ABOUT, MEMBERS, COMMUNICATIONS, LANGUAGE RESOURCES, Data, Obtaining Data, Catalog, By Year, Top Ten Corpora, Projects, Search, Memberships, Data Scholarships, Tools, Papers, LR Wiki, DATA MANAGEMENT, and COLLABORATIONS. The main content area displays the title 'English Gigaword' and a list of metadata fields and values.

Home > Language Resources > Data	
<h2>English Gigaword</h2>	
<i>Item Name:</i>	English Gigaword
<i>Author(s):</i>	David Graff, Christopher Cieri
<i>LDC Catalog No.:</i>	LDC2003T05
<i>ISBN:</i>	1-58563-260-0
<i>ISLRN:</i>	953-543-425-922-6
<i>DOI:</i>	<a href="https://doi.org/10.35111/0z6y-q265">https://doi.org/10.35111/0z6y-q265</a>
<i>Release Date:</i>	January 28, 2003
<i>Member Year(s):</i>	2003
<i>DCMI Type(s):</i>	Text
<i>Data Source(s):</i>	newswire
<i>Project(s):</i>	TIDES, GALE, EARS
<i>Application(s):</i>	natural language processing, language mod
<i>Language(s):</i>	English
<i>Language ID(s):</i>	eng
<i>License(s):</i>	<a href="#">LDC User Agreement for Non-Members</a>
<i>Online Documentation:</i>	<a href="#">LDC2003T05 Documents</a>
<i>Licensing Instructions:</i>	<a href="#">Subscription &amp; Standard Members, and Non-</a>
<i>Citation:</i>	Graff, David, and Christopher Cieri. English ( Philadelphia: Linguistic Data Consortium, 20

## • Computational linguistics

- ❖ makes extensive use of corpora and other sorts of digital ‘language resources’
- ❖ many large corpora are constructed mainly for use in computational linguistics

## English Gigaword corpus





---

## Computational linguistics

seeks to develop the computational machinery needed for an agent to exhibit various forms of linguistic behavior [Mani 2013: 466].

---


*Agent*: both human beings and artificial agents such as computer programs

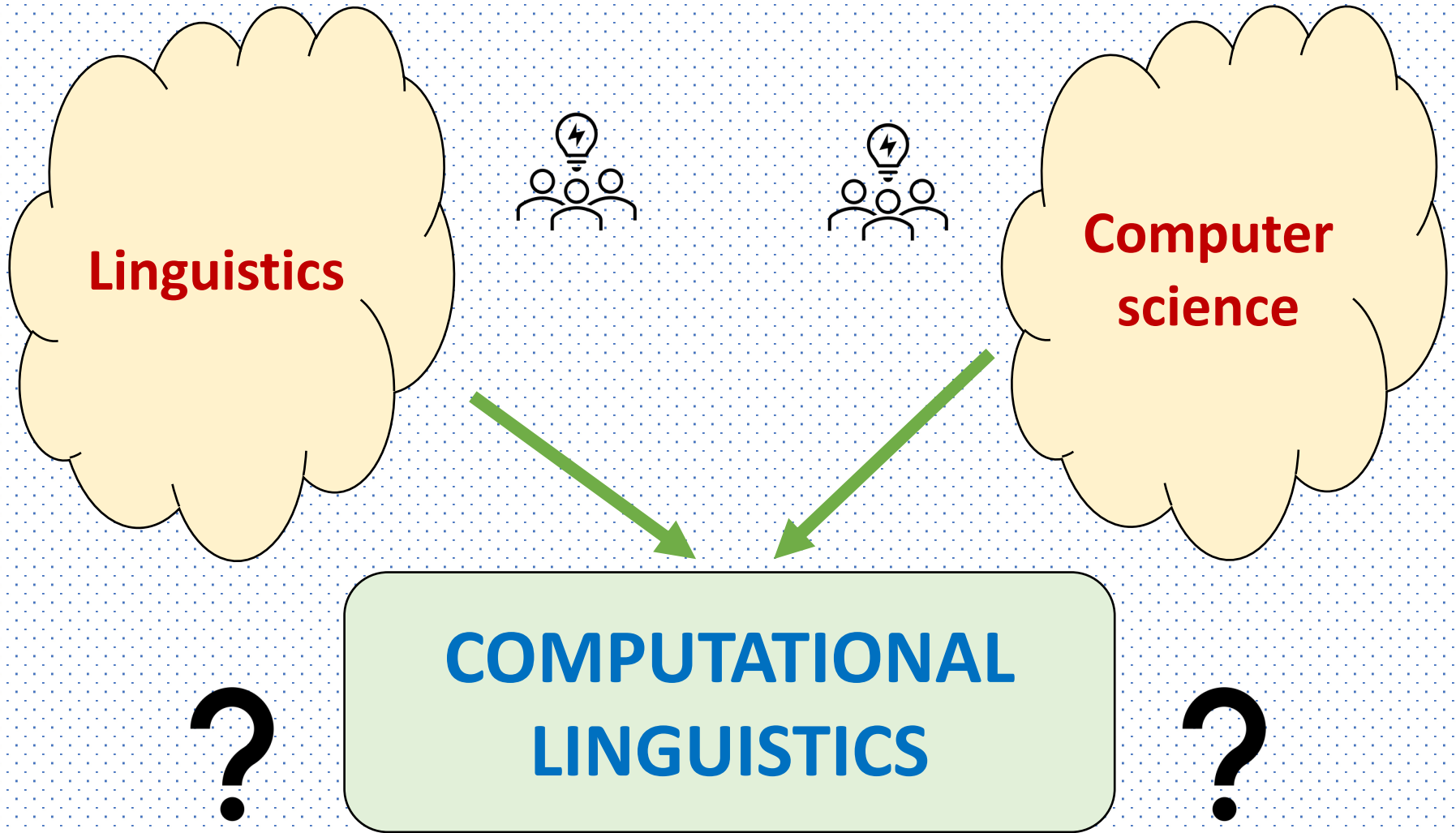
*Machinery*: computer programs as well as the linguistic knowledge that they contain

---

## Computational linguistics

involves designing and developing programs to carry out linguistic tasks





- how language is represented and processed by the human mind
- what is the role of statistics in language understanding

# 3.2. Classification of applied linguistic systems



## Text preparation



– *Automatic hyphenation* of words in natural language texts



## – *Spell checking*

detection and correction of typographic and spelling errors



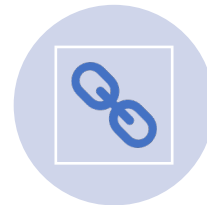
## – *Grammar checking*

detection and correction of grammatical errors


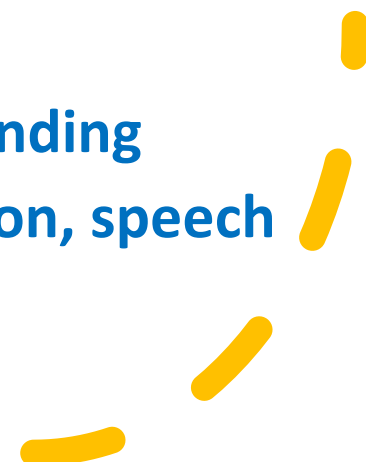


## – *Style checking*

detection and correction of stylistic errors



– *Referencing* specific words, word combinations, and semantic links between them

- 
- **Information retrieval** in scientific, technical, and business document databases
  - **Automatic translation** from one natural language to another
  - **Natural language interfaces** to databases and other systems
  - **Extraction of factual data** from business or scientific texts
  - **Text generation** from pictures and formal specifications
  - **Natural language understanding**
  - **Optical character recognition, speech recognition**
- 

## 4. Morphological processing

**Morphology** is the study of the structure of words

The task of an automatic morphological analyzer is to take a word in a language and break it down into its **stem form** along with any **affixes** that it may have attached to that stem

## *Lily reads well*

---

- *Lily* as a proper name
- *reads* as the 3d p.s. present form of the verb *read* (*read + s*)
- *well* as either an adverb or a singular noun



# 4.1. Tokenization

the first step in morphological analysis  
it is to identify separate words



## 4.2. Morphological analysis and synthesis

---

- a. An *automatic morphological analyzer* takes a word and breaks it down into its component morphemes (stems and affixes)
- Sometimes, instead of a full morphological analysis, a simple **stemming** algorithm is used which strips off suffixes to arrive at a stem form



---

b. Another strategy is to use a fully inflected lexicon, which includes all the possible affixed forms of every word in the language



The analyzer simply looks up the word in the list. Such lists are almost inevitably **incomplete**, and they can become too large and unwieldy for computers to handle

# 5. Syntactic processing

- Given a set of linguistic rules that describe how elements of a sentence can be put together, a computer program called a *syntactic parser* will try to find the best grammatical analysis of a sentence



*I can fish.*

## 5.1. Context- free grammars

- *phrase structure rules* break up a sentence into its constituent parts, consisting of syntactic phrases or words

S → NP VP  
VP → Aux V  
VP → V NP  
VP → V  
VP → Aux V NP  
NP → D N  
NP → N  
NP → Pronoun

V → can  
V → fish  
V → dance  
Aux → can  
D → the  
N → fish  
N → dance  
Pronoun → I

# “Toy grammar” abbreviations

---

<b>Aux</b>	<b>Auxiliary (can)</b>
<b>D</b>	determiner (the)
<b>N</b>	noun
<b>NP</b>	noun phrase
<b>S</b>	sentence
<b>V</b>	verb
<b>VP</b>	verb phrase

---

*I can fish.*

## 5.2. Parsing



How can we get a computer to analyze the syntactic structure of a sentence?



A **parser** takes an input sentence and produces one or more syntactic representations of it

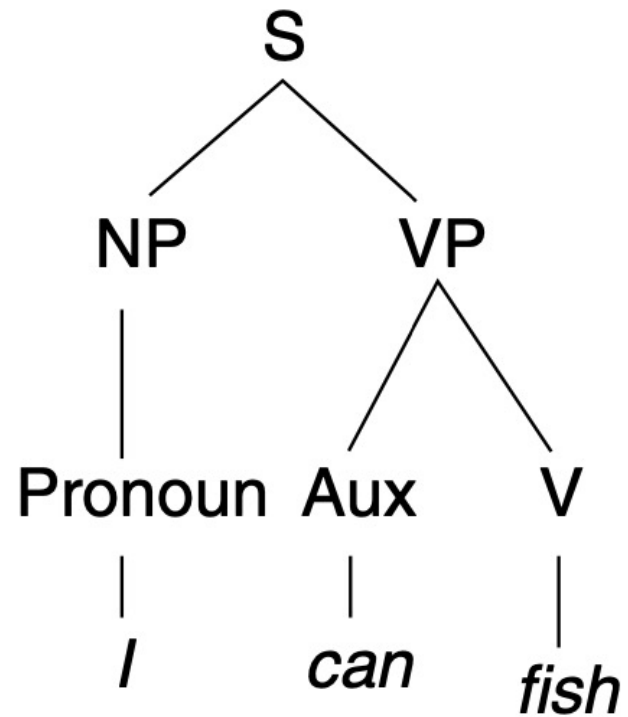
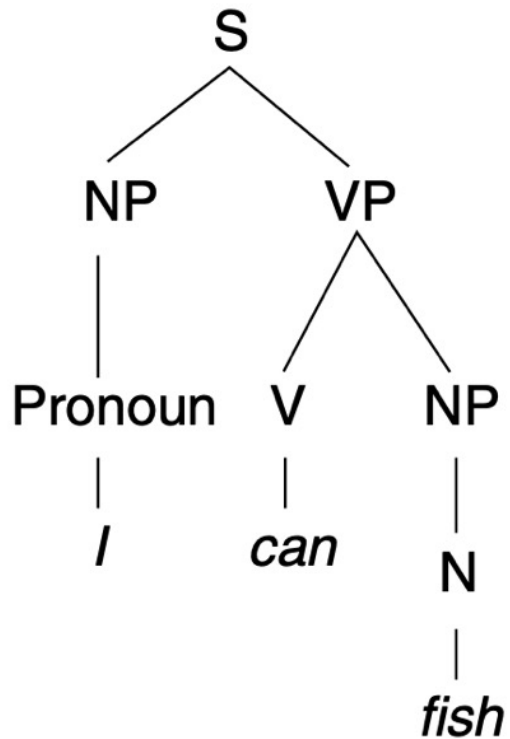


It produces a single representation if the sentence is syntactically **unambiguous**, but more than one representation if there is syntactic **ambiguity**

*I can fish.*

## PARSE TREE

top-down parser builds the parse trees from the top-down



## 5.3. Part-of-speech tagging



To each word that has more than one part of speech, the tagger assigns the **most likely part-of-speech**

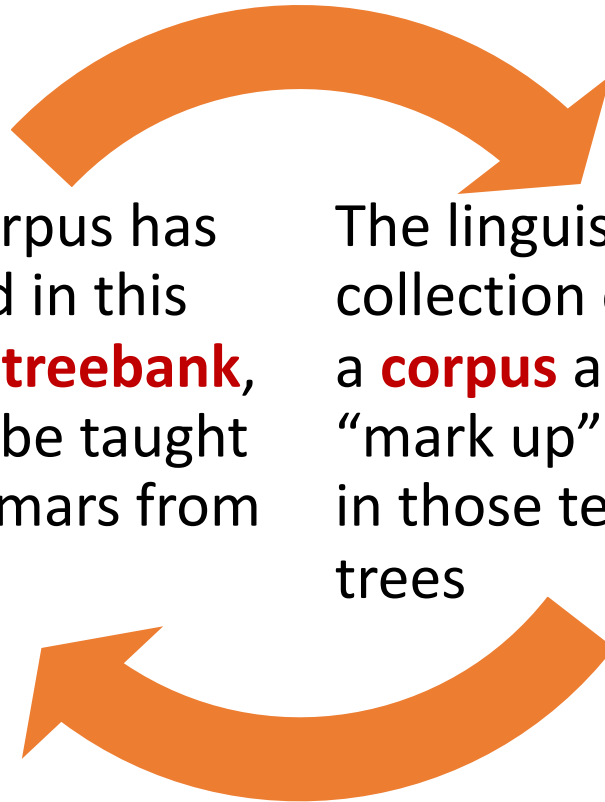


This is based on context-based rules derived by **human intuition** (after *the, fish* is likely to be a **noun**) **machines** that learn from a collection of example sentences that have been tagged already with parts of speech

## 5.4. Statistical parsing

Once a large corpus has been annotated in this way, creating a **treebank**, computers can be taught to induce grammars from it

The linguists take a collection of texts called a **corpus** and analyze and “mark up” the sentences in those texts with parse trees





## 6. Semantic processing

*How can a computer decide which meaning is intended?*

- **word-sense disambiguator**
- It can use the **context of neighboring words** in the sentence as well as other words in the document to figure out which meaning of a given word is most likely.



the verb  
*cook*

- A verb's meaning includes
- syntactic *subcategorization*  
⇒  
the syntactic elements, or  
“arguments,” it combines with
- *thematic roles*  
⇒  
• the semantic relations  
between a verb and its  
arguments

*She cooked meatloaf*

*She cooked Mary a great dinner*

a. COOK (Theme<sub>NP</sub>)

b. COOK (Recipient<sub>NP</sub>, Theme<sub>NP</sub>)

c. COOK (Theme<sub>NP</sub>, Recipient<sub>for-PP</sub>)

*She cooked a great dinner for Mary*

# Sentence meaning

*How can a computer combine meanings?*

- One way to approach this is to have **semantic rules** that accompany syntactic rules in the grammar.
- The syntactic rule (a) below is augmented with a semantic rule (b) that says the VP's meaning is constructed by applying the V's meaning to the NP's meaning

a.  $VP \rightarrow V NP$

b.  $VP.meaning = \text{Apply}(V.meaning, NP.meaning)$