

# DeepDict

**A Graphical Corpus-based Dictionary  
of Word Relations**

Eckhard Bick

University of Southern Denmark & GrammarSoft ApS

# Motivation – the lexicographer's view

- a) lexicography: better data
  - corpus data have better coverage and better legitimacy than introspection or chance quotes from literature
  - Quantity: the more the better
  - Authenticity: real running data from mixed sources (Internet), or at least a source with mixed topics (Wikipedia, news text)
  - given a corpus, a grammatically annotated one is best for the extraction of lexical patterns and statistics
  - lexical patterns are best based on linguistic relations (subject, object etc.) rather than mere adjacency in text
  - even given a corpus that satisfies all of the above, it is cumbersome to search for and quantify lexical patterns
  - even a statistics-integrating interface like our CorpusEye will only provide data for one pattern at a time

# Motivation – the user's view

- b) lexical information: better accessibility
  - electronic vs. paper: no size limitations, easy searching, “depth-on-demand” (QuickDict vs. DeepDict)
  - passive (“definitional”) vs. active (“productive-contextual”)
  - Advanced Learner's Dictionary: information on how to use a word in context – syntactic and semantic restrictions and combinatorial aspects, e.g. **A gives x to B** (A,B = +HUM, x,y = -HUM)
  - “live” examples
  - but how to show, for a given entry word, all constructions and examples, and how not to forget any?

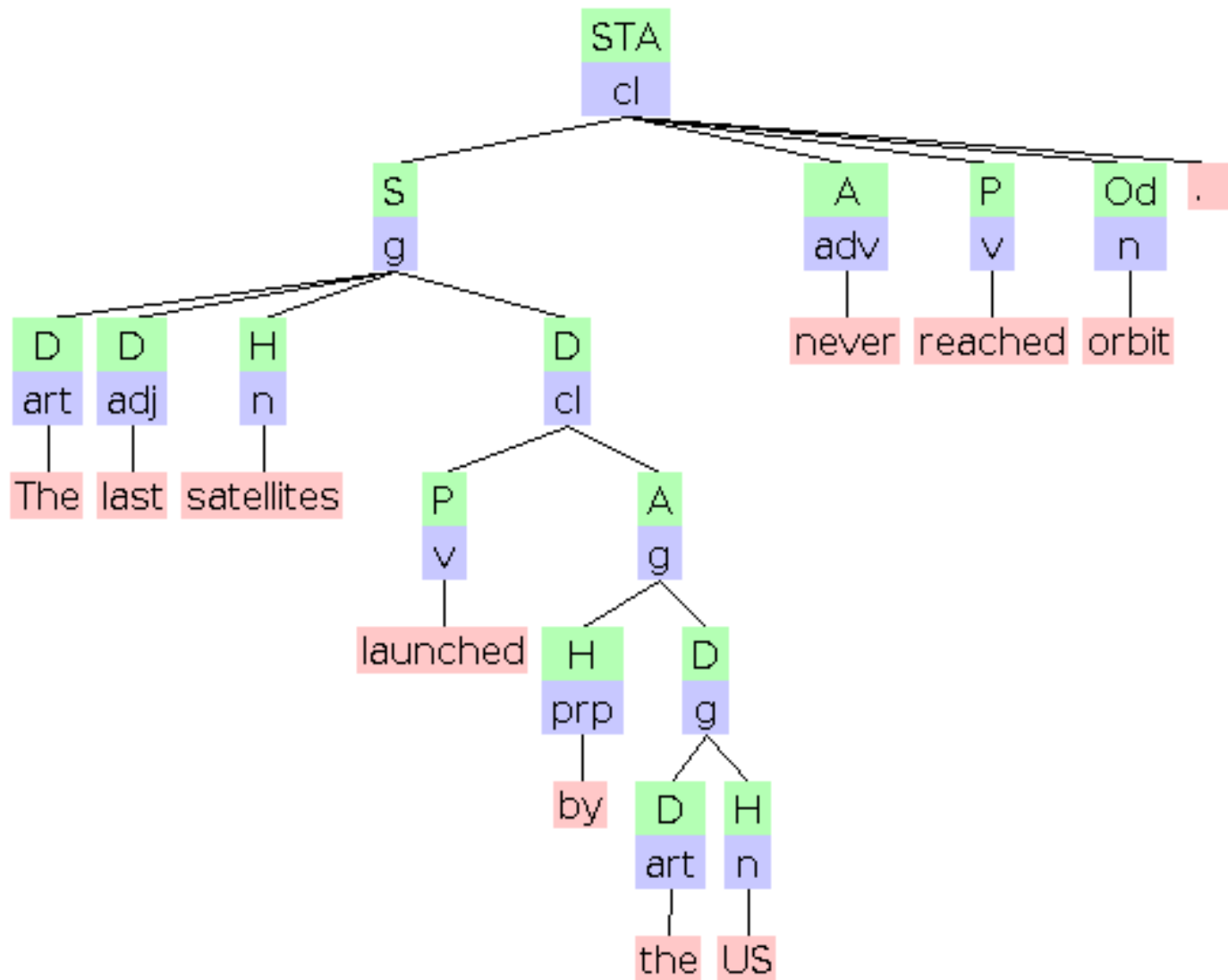
# Idea: graphical presentation of **lexical complements** based on **dependency statistics**

- annotate a corpus
  - Peter [Peter] @SUBJ ate a couple of apples ["apple"] @ACC
  - Cats ["cat"] @SUBJ eat mice ["mouse"] @ACC
- generalize by collecting and counting dependency pairs of lemmas (simplified):
  - PROP\_SUBJ -> eat, cat\_SUBJ -> eat
  - apple\_ACC -> eat, mouse\_ACC -> eat
- present the result in list form
  - {PROP,cat} SUBJ -> eat <- {apple,mouse} ACC

# Dependency tree annotation

The	<def>	ART	@>N	#1->3
last	<num-ord>	ADJ	@>N	#2->3
satellites	<Vair>	N P NOM	@SUBJ>	#3->9
launched		V PCP2 PAS	@ICL-N<	#4->3
by		PRP	@<PASS	#5->4
the	<def>	ART	@>N	#6->7
US	<civ>	PROP F S	@P<	#7->5
never	<atemp>	ADV	@ADVL>	#8->9
reached		V PAST	@FMV	#9->0
orbit	<L>	N S NOM	@<ACC	#10->9
\$.				#11->0

# equivalent constituent tree



# how to distinguish between typical and non-informative complements?

- use frequency counts for dependency pairs
- normalize for lexical frequency
- $C * \log(p(a \rightarrow b) ^2 / (p(a) * p(b)))$
- use thresholds for minimum co-occurrence strength and minimum absolute number of occurrences
  - ask for strong positive correlation (mutual information)
  - $\log_2$  frequency classes: 1 (1), 2 (2-4), 3 (5-8), 4 (9-16) ....
- for a few special word classes, use generalizations:
  - PROP/hum (names)
  - NUM (numbers)
- separate treatment of pronouns (only relative to each other)

# Data production

**raw text:**  
- Wikipedia  
- newspaper  
- Internet  
- Europarl

**corpus:**  
encoding cleaning  
sentence separation  
id-marking

DanGram  
EngGram  
...

Comp. lexica  
CG grammars  
Dep grammar

example  
concordance  
\* .....  
\* .....  
\* .....

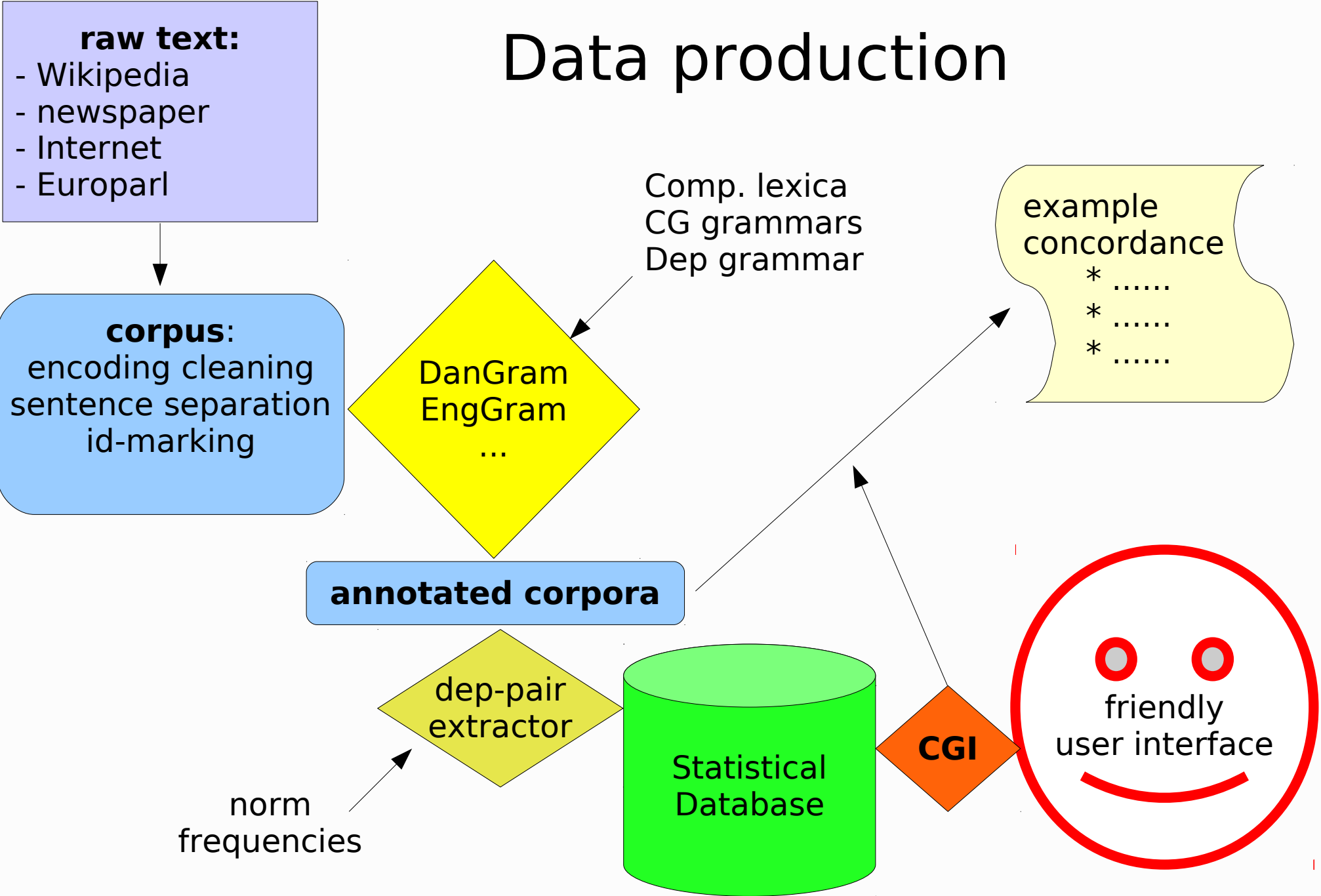
**annotated corpora**

dep-pair  
extractor

norm  
frequencies

Statistical  
Database

**CGI**















# the applicational environment

- DeepDict is hosted at **www.gramtrans.com** where it is part of an integrated suite of translation tools
- it was developed as a spin-off from many years of
  - a) parser development
  - b) corpus annotation projects (Corpus Eye at corp.hum.sdu.dk)
  - c) lexicography and MT research
- The primary languages are the Germanic and Romance languages, but the method is largely language-independent given a lemmatized and dependency annotated corpus in that language
- Few similar approaches of turning corpus data into lexicography:
  - Sketch Engine (Kilgariff et al. 2004)
  - Leipzig Wortschatz project (Biemann et al. 2004)

# Corpus sources and parsers

	<b>Parser</b>	<b>Lexicon</b>	<b>Grammar</b>	<b>Corpora</b>
	<a href="#">DanGram</a>	100.000 lexemes, 40.000 names	8.400 rules	ca. 159 M words (mixed)
	<a href="#">PALAVRAS</a>	70.000 lexemes, 15.000 names	7.500 rules	ca. 210 M words (news) [+170 mill. wiki a.o.]
	<a href="#">HISPAL</a>	73.000 lexemes	4.900 rules	ca. 90 M words (Wiki, Europarl, Internet)
	<a href="#">EngGram</a>	81.000 val/sem	4.500 rules	ca. 210 M words (mixed) [+106 M email & chat]
	<a href="#">SweGram</a>	65.000 val/sem	8.400 rules	ca. 60 M words (news, Europarl) [+ Wiki]
	<a href="#">NorGram</a>	OBT / via DanGram	OBT / via DanGram	ca. 50 M words (Wikipedia) [+ internet]
	<a href="#">FrAG</a>	57.000 lexemes	1.400 rules	-  [+67 mill Wiki, Europarl]
	<a href="#">GerGram</a>	25.000 val/sem	LS+1.300 rules	ca. 44 M words (Wiki, Europarl) [+ internet]
	<a href="#">EspGram</a>	30.000 lexemes	2.600 rules	ca. 58 M words (mixed)
	<a href="#">ItaGram</a>	30.600 lexemes	1.600 rules	46 M (Wiki, Europarl)

cut-and-paste text

Languages:  
en, da, no,  
se, pt, eo

URL

Tools:  
- WAP  
- sms  
- Firefox plugin  
- Docs

Skip Navigation

### Navigation

- Home
- News
- Products
- Free
- Personal
- Commercial Lite
- Commercial Standard
- Commercial Direct
- Schools
- Languages
- System Features
- Donate Words
- GramTrans
- Contact
- Feedback
- Research

### Tools

- Document Translation
- Mozilla Firefox Extension
- Texting (SMS)
- Statistics

### Controls

- Site Admin
- My Profile
- Logout

### News

- New lexicon feature
- New registration policy
- Substantial improvements for English - Danish

## Text Translation

Danish to English    Length 0

## Web Page Translation

http://

## News

### 2008-04-11: New lexicon feature

Many people use GramTrans not just for machine translation, but also - or primarily - as a dictionary service, and for those users we have long provided the **QuickDict** feature - a simple lexicon lookup for individual words (including inflected words), where alternative translations can be preserved rather than contextually discarded. We have now gone one better and offer the **DeepDict** lexifier, a full "Learner's Dictionary" lookup providing usage information for the target language words, as well as example sentences and statistical information. Just click on the QuickDict translation (in red).

DeepDict is also available as such, of course, for direct lookup, and has its [own page](#) on the GramTrans site. For now, the QuickDict->DeepDict link works for the English side of Danish->English translations, while direct lookups also cover Danish and Spanish DeepDict, with more languages in the pipeline. As a target group we envision not just translators, but also schools, researchers and lexicographers. DeepDict is based on the grammatical and statistical analysis of huge text corpora, and can thus offer a unique and authentic overview of lexical usage, far beyond the scope of ordinary dictionaries.

2008-01-24: [New registration policy](#)

2008-01-03: [Substantial improvements for English - Danish](#)

2007-10-22: [Testing Document Translation](#)

Ads by Google

**Learn English in 10 days**  
Learn English in 10 days with top language learning software  
[www.natively.com/Englis](http://www.natively.com/Englis)

**Learn Danish For Free**  
Download your free Danish BYKI™ software and learn Danish fast!  
[www.BYKI.com](http://www.BYKI.com)

**One Click Translation**  
Translate 50 Different Languages  
Translate translation  
[www.babylon.com](http://www.babylon.com)

**Freelance Translation**  
Get quotes from freelance translators! All languages.  
[www.TranslatorsTown.co](http://www.TranslatorsTown.co)

# The first dictionary layer: QuickDict

The screenshot shows the GramTrans website interface. At the top, there is a navigation bar with links for Home, News, Products, Languages, and Deep. Below this is a large banner with the text "GRAMTRANS" in a stylized white font on a dark background with diagonal stripes. Underneath the banner, there are flags for various languages and a "Skip Navigation" link. On the left side, there is a "Navigation" menu with links for Home, News, Products, Free, Personal, Commercial Lite, Commercial Standard, Commercial Direct, Schools, Languages, and Donate Words. The main content area is divided into two sections: "Translated Result" and "Dictionary Lookup Result".

**Translated Result**

voice

textual choice

[Printable View](#)

*Information: Translation took 0.728 seconds which is 8.245 CpS.*

**Dictionary Lookup Result**

The QuickDict lookup provides a few alternative translations. Touch a translation for frequency information, or click it for contextual information ([DeepDict](#)).

DeepDict link

mouse-over frequency

**stemme**

stemme (n) – voice; vote

stemme (v) – vote; tune; be correct

TL/SL-inflected forms

# The second dictionary layer: DeepDict

## Implicit semantics from lexical relations

### Quick Reference

An entry consists of the following elements, taking "8.37:6 comic" as example:

- **8.37**: The co-occurrence strength between the lookup word and a given relation (relative frequency).
- **6**: Dual logarithmic value of absolute frequency. Scale is from 1 to 9.
- **comic**: The co-occurring word. In bold face if the absolute frequency value is 4 or higher.

Red numbers can be clicked to show examples in concordance form, if available.

For more information, please see the [DeepDict Reference page](#).

### Change Lookup Parameters

#### Word to look up:

#### Word class:

- Noun
- Verb
- Adverb
- Adjective

#### Lookup language:

- Danish
- English
- French
- German
- Portuguese
- Spanish

#### Lexical frequency threshold:

- High
- Medium
- Low
- None

Minimum occurrence:

Minimum relative frequency:

Show top:

# DeepDict: nouns 1

## voice (noun)

countable

<b>Premodifiers:</b> 6.73:7 <b>loud</b> · 6.57:7 <b>NUM</b> · 6.41:6 <b>distinctive</b> · 5.05:7 <b>deep</b> · 6.64:5 <b>soprano</b> · 7.46:4 <b>gravelly</b> · 7.44:4 <b>husky</b> · 4.34:7 <b>single</b> · 5.21:6 <b>inner</b> · 4.2:7 <b>own</b> · 6.59:4 <b>baritone</b> · 5.58:5 <b>passive</b> · 6.52:4 <b>hoarse</b> · 4.46:6 <b>soft</b> · 5.46:5 <b>authoritative</b> · 4.32:6 <b>quiet</b> · 4.28:6 <b>human</b> · 6.28:4 <b>squeaky</b> · 6.16:4 <b>narrative</b> · 5.98:4 <b>gruff</b>	<b>PP postmodifiers:</b> 8.72:8 <b>rel-INDP</b> 1.35:3 <b>interr-INDP</b> 2.58:5 <b>of character</b> 2.43:5 <b>of reason</b> 2.25:5 <b>of god</b> 3.08:4 <b>from behind</b> 1.75:5 <b>of america</b> 2.7:4 <b>of conscience</b> 2.43:3 <b>of dissent</b>	<b>Modifier of:</b> 6.04:7 <b>actor</b> · 5.94:4 <b>telephony</b> · 3.58:5 <b>actress</b> · 4.91:3 <b>coil</b> · 2.04:5 <b>communication</b> · 2.88:4 <b>talent</b> · 2.88:4 <b>recorder</b> · 3.52:3 <b>choir</b> · 2.19:4 <b>transmission</b> · 1.02:5 <b>vote</b> · 1.76:4 <b>channel</b> · 2.7:3 <b>characterization</b> · 3.59:2 <b>synthesizer</b> · 3.47:2 <b>inflection</b> · 3.41:2 <b>synthesis</b> · 0.31:5 <b>system</b> · 1.27:4 <b>message</b> · 1.6:3 <b>directive</b> · 0.51:4 <b>call</b> · 1.43:3 <b>lesson</b>
---	---	--

<b>one can ...</b>	14.93:2 <b>modulate</b> · 12.4:2 <b>murmur</b> · 8.62:5 <b>recognise</b> · 11.36:1 <b>hush</b> · 10.96:1 <b>shriek</b> · 3.44:8 <b>hear</b> · 9.28:2 <b>amplify</b> · 8.34:2 <b>imitate</b> · 3.72:6 <b>lower</b> · 6.7:3 <b>obey</b> · 7.2:2 <b>mimic</b> · 4.9:4 <b>lend</b> · 5.02:3 <b>possess</b> · 4.68:3 <b>dub</b> · 0.53:7 <b>raise</b> · 5.52:2 <b>heed</b> · 4.36:3 <b>drown</b> · 6.14:1 <b>clip</b> · 5.08:2 <b>equal</b> · 6.06:1 <b>sharpen</b>  8.54:4 <b>creep into</b> · 10.09:2 <b>exclaim in</b> · 9.09:2 <b>mutter in</b> · 2.63:8 <b>speak with</b> · 4.18:5 <b>sing in</b> · 8.07:1 <b>retort in</b> · 6.42:2 <b>whisper in</b> · 3.65:4 <b>reply in</b> · 4.24:3 <b>cry in</b> · 6.06:1 <b>recite in</b> · 5.78:1 <b>startle by</b> · 4.77:2 <b>inject into</b> · 2.77:4 <b>listen to</b> · 0.54:6 <b>speak in</b> · 3.39:3 <b>detect in</b> · 3.34:3 <b>consist of</b> · 2.91:3 <b>shout in</b> · 2.04:3 <b>sing with</b> · 0.78:3 <b>sound like</b>	<b>a voice</b>
<b>a voice can ...</b>	14.14:3 <b>muffle</b> · 12.44:4 <b>tremble</b> · 9.54:4 <b>whisper</b> · 11.44:2 <b>crackle</b> · 11.32:2 <b>growl</b> · 6.13:7 <b>sound</b> · 10.61:1 <b>wobble</b> · 9.49:2 <b>drip</b> · 9.49:2 <b>thicken</b> · 10.29:1 <b>squeak</b> · 7.85:3 <b>falter</b> · 5.82:5 <b>echo</b> · 8.69:2 <b>waver</b> · 7.43:3 <b>harden</b> · 8.24:2 <b>reverberate</b> · 8.7:1 <b>exclaim</b> · 5.38:4 <b>fade</b> · 5.25:4 <b>shout</b> · 6.17:3 <b>deepen</b> · 7.17:2 <b>startle</b>	
<b>a voice can be</b>	13.33:3 <b>muffle</b> · 12.39:2 <b>hush</b> · 8.1:3 <b>clip</b> · 9.75:1 <b>tinge</b> · 8.7:1 <b>amplify</b> · 1.58:7 <b>hear</b> · 4.68:3 <b>dub</b> · 5.12:2 <b>choke</b> · 4.09:2 <b>drown</b> · 3.62:2 <b>strain</b>	<b>... 'ed</b>

# DeepDict: nouns 2 - the “word field” side benefit

## língua (noun)

countable

<b>Premodifiers:</b> 1.91:7 próprio · 0.63:6 segundo · 0.33:4 só	<b>PP postmodifiers:</b> 1.8:5 de areia 0.01:6 de trabalho 0.4:3 de difusão 0.12:1 de gringo	<b>Adjectival postmodifiers:</b> 7.16:9 português · 6.6:9 oficial · 6.58:9 inglês · 3.85:8 francês · 5.79:6 castelhano · 3.91:7 alemã · 4.39:5 gestual · 5.17:4 veicular · 4.13:5 nativo · 2.71:6 chinês · 3.56:5 latino · 1.37:7 nacional · 2.17:6 comum · 2.12:6 espanhol · 2.01:6 diferente · 2.88:5 albanês · 3.76:4 falado · 3.57:4 berbere · 2.44:5 curdo · 1.29:6 regional · 1.99:5 galego · 1.74:5 original · 0.69:6 local · 2.47:4 natal · 2.44:4 eslavo
<b>se pode ...</b>	10.44:1 vivificar · 9.9:1 escovar · 2.79:7 <b>aprender</b> · 6.45:2 manejar · 6.44:2 afiar · -0.88:9 <b>falar</b> · 1.11:7 <b>dominar</b> · 4.88:3 morder · 4.64:3 desatar · 2.53:5 <b>ensinar</b> · 3.48:4 <b>soltar</b>  6.02:3 verter para · 1.53:7 <b>traduzir em</b> · 4.01:4 <b>redigir em</b> · 1.35:5 <b>expressar em</b> · 0.17:6 <b>cantar em</b> · 0.09:5 <b>editar em</b> · 0.94:4 <b>imprimir em</b> · 0.78:4 <b>expressar em</b>	<b>uma língua</b>
<b>uma língua pode ...</b>	2.16:3 <b>soltar</b> · 0.46:4 <b>ensinar</b>	

# Linked example concordances and “word sketches”

Forms	Abs Freq	Rel Freq
Total	6749	100.00%
vægt -> lægger	3608	53.46%
vægt -> lægge	1270	18.82%
vægt -> lagde	820	12.15%
vægt -> lagt	815	12.08%
vægten -> lægger	100	1.48%
vægten -> lægge	68	1.01%

[Show all forms...](#)

word sketch:

lægger	stor mest afgørende meget	vægt	på {, / . / , fordi } {bag / på denne / på ved ansættelser / på det }
--------	------------------------------------	------	--

## Concordances for: vægt\_N -> lægge\_V

ID	Text
inf30-33978	« Markedsføring på værdier og ansvar skal være en vigtig kilde til at øget konkurrenceevne på verdensmarkedet , fordi den globale forbruger i stadig større udstrækning <b>lægger vægt</b> på den slags , når han eller hun køber ind .
inf100-35624	I sidste ende giver den diplomatiske balancegang dog mulighed for , at den tyrkiske regering for en tid kan <b>lægge</b> mest <b>vægt</b> på sin egen tillægserklæring af indenrigspolitiske årsager , mens EU fokuserer på toldaftalen .
inf40-148330	« I dag <b>lægger</b> vi megen <b>vægt</b> på , at tillidsmanden skal arbejde med uddannelse og kunne svare på spørgsmål om pension .
c2000- dmfrinab114	Netop den pind har vi fået ind , og det <b>lægger</b> vi stor <b>vægt</b> på , fordi vi netop vil gøre



# collocation sketches 2

## Form Statistics for: adoption\_N -> give\_V

Forms	Abs Freq	Rel Freq
Total	23	100.00%
adoption -> gave for	12	52.17%
adoption -> give for	7	30.43%
adoption -> gives for	3	13.04%
adoption -> given for	1	4.35%

*word sketch:*

for            him up  
              up            for adoption ->    { . / . / , / and }  
              up her baby

## Concordances for: adoption\_N -> give\_V

ID	Text
w7-1803406	The agencies will cover the costs of delivery and the medical care for any woman who <b>gives</b> up her baby <b>for adoption</b> .
w7-400518	Soon after release , she was impregnated by a man she barely knew and gave birth to a baby girl , which she had to <b>give up for adoption</b> .

# DeepDict: Verb + Complements

## caress (verb)

total of 527 relations

### Subjects:

**PERS:** we, he, they, she

6.21:2 PROP · 4.79:2 finger ·  
4.62:1 breeze · 4.44:1 thumb ·  
2.89:2 hand · 1.47:1 eye

### Accusative objects:

**PERS:** her, one another

6.62:2 cheek · 5.12:2 skin · 5.83:1 fingertip · 4.74:2 hair · 4.24:2 breast ·  
4.7:1 spine · 3.45:2 face · 4.42:1 jaw · 3.86:1 neck · 2.71:2 body ·  
3.71:1 PROP · 2.59:2 back · 3:1 length · 0.25:1 head

caress ...	5.54:2 gently · 3.71:1 sensuously
caress <b>to</b> ...	4.48:1 waist
caress <b>with</b> ...	4.01:1 tongue · 1.5:1 hand
caress <b>in</b> ...	0.23:1 way

# Special treatment of word classes

- Each major PoS has its own lexicogram setup
  - verb + arguments/adjuncts, noun/adjective + modifiers
- PROP and NUM are generalized to avoid noise
- pronouns are very frequent and can't be directly compared to other lexical material
- but pronouns are carriers of abstracted semantic information (cp. Odense pronominal valency approach)
  - ± human: *who, what*
  - male / female: *he, she, him, her*
  - place, direction: *der, derhen, her, herhen (Danish)*
  - countable / mass: *much, many*

# Pronouns as semantic classifiers

- “drikke” (drink): +anim vs. quantity
  - jeg (I), vi (we), han (he), ..... den (UTR-it)  
<=> den (UTR-it), meget (much)
- “marry”: +hum/male vs. +hum/female, 2ps > 1ps
  - they, he, she, you, who, we <=> you, her, him, them, who, me
- “learn”: +hum <=> -hum/abstract
  - we, they, you, he, I, she <=> something, what, a lot, them, much
  - subclause complements: that-KS, interrogatives

## learn (verb)

total of 63496 relations

Hide Frequencies

### Subjects:

**PERS:** we, they, you, he, I, she, i, who, one, one, everyone  
**9.92:9 PROP** · **3.9:8 child** · **4.72:5 pupil** · **3.35:6 student** ·

### Subclauses:

**4.95:9 that**  
**5.42:6 interr**

### Accusative objects:

**PERS:** many, something, what, a lot, them, much, anything, that,  
nothing, all, which

# Verb - adverb collocations

- (a) free adverb(ial)s: time, place, manner ...
- (b) valency bound adverb(ial)s
  - feel *how* (manner argument)
  - live *where* (place argument)
  - go *where* (direction argument)
- (c) verb-integrated particles
  - give *up*, fall *apart*. ? cut *out* (object predicative?)
- Since DeepDict is a lexicographical rather than a syntactic tool, we only keep verb-integrated particles separate (to allow sub-lemmatization), and lump everything else in an umbrella category (brown field)

# “run” + adverbs

## Verbal particles:

4.21:9 out · 4.62:8 off · 3.56:8 down · 2.68:7 over ·  
0.97:4 through

- known verbal particles

6.11:9 away · 5.64:7 unsuccessfully · 5.99:6 counter · 6.52:5 aground · 5.23:6 midway · 6.19:5 concurrently · 3.75:7 fast ·  
4.42:6 homely · 4.32:6 north-south · 4.68:5 amok · 4.61:5 upstairs · 2.51:7 back · 5.4:4 northwards · 3.32:6 through ·  
3.12:6 south · 3.06:6 east · 4.02:5 smoothly · 2.98:6 west · 1.94:7 now · 3.77:5 east-west · 4.76:4 firstly · 1.76:7 well ·  
2.57:6 north · 3.53:5 quickly · 1.51:7 lately · 1.03:7 up · 0.7:6 down · 0.29:5 in

- new verb-integrated particles: *run amok, run counter (to)*
  - direction valency: *away, north-south, back, northwards*
  - free manner adverbs: *unsuccessfully, quickly, smoothly*
- mirrored by semantically distinct object complementation classes:

- *run (the) length / course (of) --- move adv.tr.*
- *run (n) miles --- move itr.*
- *run (the) risk --- fixed expression*
- *run (a) finger (along/over) --- move np.tr.*
- *run (a) program / system --- tr. “operate”*
- *run (a) school / centre / business -- tr. “organize”*

# Verb - preposition collocates

- maybe the most ignored piece of usage information in dictionaries
- very difficult for learners, since the choice of preposition is more syntactic than semantic (cp. also aphasia research, Broca vs. Wernicke centres)
- like adverbs, prepositions (or rather pp's) can either be valency bound or free complements. It's near-impossible to make the distinction automatically, but know valencies are \*-marked
- the binary dependency link has to be extended from the syntactic to the semantic head of the pp, storing 3-part links in the database

# known valency-bound pp complements

<b>run in ...</b>	3.83:8 <b>election</b> · 3.93:5 <b>mode</b> · 2.73:6 <b>direction</b> · 0.19:3 <b>median</b> · 0.34:2 <b>groove</b> · 0.32:2 <b>vein</b> · 0.88:1
<b>run at ...</b>	6.47:5 <b>racecourse</b> · 0.11:4 <b>speed</b>
<b>run *for ...</b>	5.21:6 <b>re-election</b> · 4.81:6 <b>auditor</b> · 3.55:7 <b>gover</b> 3.27:6 <b>episode</b> · 1.14:8 <b>year</b> · 3.1:6 <b>commissione</b> 1.38:6 <b>leadership</b> · 1.33:6 <b>senate</b> · 1.32:6 <b>office</b> 1.49:3 <b>governorship</b> · 1.02:3 <b>touchdown</b> · 0.68:2 <b>kilo</b>
<b>run *into ...</b>	3.39:7 <b>trouble</b> · 3.58:6 <b>difficulty</b> · 1.33:6 <b>problem</b>
<b>run *as ...</b>	3.2:7 <b>candidate</b> · 1.09:3 <b>independent</b>
<b>run at ...</b>	3.53:6 <b>speed</b> · 1.43:4 <b>theatre</b> · 2.98:2 <b>racecourse</b> ·
<b>run *on ...</b>	3.25:6 <b>platform</b> · 1.83:5 <b>ticket</b> · 2.77:4 <b>petrol</b> · 0. 0.95:5 <b>track</b> · 1.94:4 <b>processor</b> · 0.7 0.63:3 <b>mainframe</b> · 0.6:3 <b>architecture</b>

<b>run through ...</b>	2.93:6 <b>hair</b> · 1.33:6 <b>term</b> · 1.03:6 <b>town</b> · 1.67: 0.88:1 <b>shire</b> · 0.09:1 <b>curl</b>
<b>run up ...</b>	2.64:5 <b>stairs</b> · 0.56:3 <b>steps</b> · 0.19:1 <b>overdraft</b>
<b>run over ...</b>	2.1:5 <b>distance</b>
<b>run along ...</b>	1.69:5 <b>edge</b> · 0.28:5 <b>side</b> · 0.02:5 <b>line</b> · 0.58:
<b>run under ...</b>	1.35:5 <b>window</b> · 1.42:4 <b>banner</b> · 0.38:3 <b>moti</b>
<b>run from ...</b>	3.25:3 <b>terminus</b> · 1.67:4 <b>may</b> · 1.56:3 <b>junction</b>
<b>run down ...</b>	2.02:4 <b>stairs</b>
<b>run on ...</b>	0.96:5 <b>basis</b> · 0.48:1 <b>microcomputer</b> · 0.19:1 <b>p</b>

free adverbial pp  
complements



# “drikke” (drink) + pp: implicit action frame

<b>drikke med ...</b>	3.13:5 <b>PROP-hum</b> · 2.07:1 svend · 0.
<b>drikke i ...</b>	3.02:4 <b>slurk</b> · 1.86:4 <b>PROP-top</b> · 0.1
<b>drikke fra ...</b>	1.04:5 <b>samling</b> · 0.35:4 <b>sans</b>
<b>drikke til ...</b>	0.08:4 <b>mad</b> · 0.46:1 aftensmÅl tid
<b>drikke af ...</b>	2.07:1 <b>plasticrus</b> · 0.43:2 vandhane ·
<b>drikke som ...</b>	1.27:1 <b>aperitif</b>

social act of drinking together

manner: in sips,  
where: place names

fixed expression:  
“drink s.o. unconscious”

drinking context: dinner

ritualized drink types:  
starter drink

**fra sig fra sans -> og samling** {aftenen før brylluppet / og ryger til udpumpning til Vejle havn}  
**og** {samling}

# DeepDict: Verb + Prep.

vote ...	6.18:8 <b>therefore</b> · 4.87:8 <b>today</b> · 5.85:7 <b>tomorrow</b> · 4.24:6 <b>unanimously</b> · 3.62:5 <b>against</b> · 3.61:5 <b>overwhelmingly</b> · 1.49:7 <b>justly</b> · 2.31:5 <b>yesterday</b> · 4.22:3 firstly · 3.21:4 <b>accordingly</b> · 0.76:6 <b>now</b> · 3.47:3 tactically · 2.33:4 <b>differently</b> · 1.6:4 <b>separately</b> · 0.34:5 <b>however</b> · 0.16:5 <b>so</b> · 1.57:3 namely · 2.51:2 wholeheartedly · 1.37:3 actually · 0.26:4 <b>simple</b> · 0.93:3 freely · 0.89:3 naturally · 0.74:3 except · 1.53:2 intelligently · 0.43:3 narrowly
vote <b>in</b> ...	12.04:9 <b>favour</b> · 3.36:8 <b>election</b> · 3.66:6 <b>referendum</b> · 4.52:4 <b>entirety</b> · 1:6 <b>committee</b> · 0.61:5 <b>vote</b> · 0.07:4 <b>chamber</b> · 0.25:3 may · 0.06:1 plebiscite
vote <b>against</b> ...	5.36:9 <b>report</b> · 3.7:7 <b>resolution</b> · 3.66:7 <b>amendment</b> · 2.59:7 <b>proposal</b> · 3.58:6 <b>motion</b> · 3.03:4 <b>paragraph</b> · 2.33:4 <b>directive</b> · 1.6:4 <b>text</b> · 0.84:4 <b>recommendation</b> · 0.49:3 PROP-hum · 0.17:3 appointment · 1.16:2 accession · 0.53:2 discharge · 0.03:2 ratification
vote <b>*for</b> ...	5.06:9 <b>report</b> · 4.35:7 <b>resolution</b> · 3.67:7 <b>amendment</b> · 4.08:6 <b>motion</b> · 2.56:6 <b>reason</b> · 3.51:5 <b>directive</b> · 1.66:6 <b>candidate</b> · 3.22:4 <b>censure</b> · 1.12:5 <b>proposal</b> · 1.6:4 <b>text</b> · 2.24:3 accession · 0.54:4 <b>bush</b> · 0.03:4 <b>regulation</b> · 1.84:2 incorporation · 0.83:3 paragraph · 0.49:3 PROP-hum · 0.85:2 abolition · 0.51:2 postponement · 0.22:2 discharge · 0.01:2 tomorrow · 0.99:1 E · 0.82:1 deletion · 0.7:1 continuance · 0.7:1 guillotine · 0.5:1 assent
vote <b>in</b> ...	6.54:7 <b>favour</b> · 0.35:5 <b>poll</b>
vote <b>*on</b> ...	3.21:6 <b>amendment</b> · 1.97:7 <b>report</b> · 1.92:6 <b>resolution</b> · 2.38:5 <b>text</b> · 2.23:4 <b>directive</b> · 1.02:5 <b>proposal</b> · 1.01:5 <b>matter</b> · 1.34:4 <b>motion</b> · 1.72:3 paragraph · 0.53:4 <b>basis</b> · 1.42:3 enlargement · 0.29:3 statute · 0.25:3 may · 0.85:2 accession · 0.48:2 recital · 0.99:1 website · 0.01:1 scoreboard

# structural symptoms of semantical differences

- meaning change between pre-modifying and post-modifying position (Romance languages)
- meaning change depending on head
  - *ill child/horse/relative* (state)
  - *ill omen/fate/fortune* (quality)
  - *ill wind/temper/humour* (intention?)
- adverbial premodifiers for adjective classification
  - intensity: *very*
  - measure: *<unit> noun modifiers*
  - state: *temporal adverbs: often*
  - $\pm$ control: *intentionally*
  - result: *from, by, with*

# DeepDict: Adjectives

## pesado (adjective)

### Pre-modifiers:

10.18:9 mais · 2.25:7 muito · 2.05:6 tão ·  
1.54:6 menos · 2.32:5 demasiado ·  
2.61:4 cada vez mais · 1.02:5 bastante ·  
0.98:4 algo · 1.4:3 excessivamente ·  
0.33:4 extremamente · 0.83:2 um pouco ·  
0.16:2 de tal forma

### Premodifier of:

5.13:7 herança · 5.07:7 derrota · 4.49:6 multa · 5.21:5 fardo ·  
3.06:6 pena · 2.92:5 encargo · 1.6:6 responsabilidade ·  
2.38:5 sanção · 3.23:4 tributo · 1.95:5 carga · 0.54:5 estrutura ·  
1.22:4 indemnização · 1.13:4 perda · 1.04:4 condenação ·  
0.98:4 factura · 1.91:3 sérvia · 1.51:3 ónus · 0.48:4 silêncio ·  
0.4:4 tarefa · 0.18:4 custo · 1.13:3 coima · 1.1:3 hum · 0.05:4 dívida  
· 0.64:3 bombardeamento · 0.57:3 burocracia

### Postmodifier of:

7.34:8 artilharia · 7.1:8 metal · 6.3:8 veículo · 4.87:8 arma ·  
4.84:6 armamento · 3.38:6 peso · 4.17:5 metralhadora ·  
2.91:6 pena · 3.81:5 comercial · 2.77:6 viatura · 2.16:6 estrutura ·  
2.88:5 herança · 3.63:4 maquinaria · 2.45:5 camião · 2.24:5 carga ·  
2.2:5 condutor · 1.81:5 derrota · 1.76:5 droga · 1.22:5 mão ·  
1.9:4 multa · 0.88:5 equipamento · 1.59:4 colisão · 1.39:4 motorista  
· 0.39:5 terreno · 1.33:4 castigo

# ill (adjective)

total of 4061 relations

Hide Frequencies

## Pre-modifiers:

9.03:7 **terminally** · 7.84:6 **seriously** · 6.11:6 **seriously** ·  
6.99:5 **gravely** · 4.37:7 **very** · 5.86:5 **critically** · 3.75:6 **too** ·  
4.67:5 **severely** · 4.8:4 **really** · 4.44:4 **physically** ·  
5.16:3 **chronically** · 4.14:4 **desperately** · 4.32:3 **violently** ·  
1.93:5 **so** · 2.76:4 **quite** · 3.47:3 **dangerously** · 3.02:3 **critically** ·  
2.92:3 **dangerously** · 2.92:3 **extremely** · 2.69:3 **chronically** ·  
3.46:2 **acutely** · 1.92:3 **increasingly** · 2.11:2 **mortally** ·  
1.07:3 **extremely** · 1.75:2 **terribly**

## Post-modifiers:

4.23:4 **with fever**  
2.66:4 **with cancer**  
3.32:3 **with tuberculosis**  
2.87:2 **with pneumonia**  
3.42:1 **with indigestion**  
3.42:1 **with gout**  
1.99:2 **with malaria**  
2.91:1 **with cyst**  
1.75:2 **with pleurisy**

## Post-modifiers:

4.23:4 **with fever**  
2.66:4 **with cancer**  
3.32:3 **with tuberculosis**  
2.87:2 **with pneumonia**  
3.42:1 **with indigestion**  
3.42:1 **with gout**  
1.99:2 **with malaria**  
2.91:1 **with cyst**  
1.75:2 **with pleurisy**

## Premodifier of:

5.28:7 **health** · 7.86:3 **repute** · 4.7:6 **effect** ·  
4.67:6 **patient** · 5.4:4 **effects** · 5.91:3 **omen** ·  
3.37:5 **treatment** · 3.48:4 **fortune** · 5.33:2 **humour** ·  
2.55:3 **luck** · 2.49:3 **feeling** · 1.39:4 **person** · 0.59:4 **child**  
· 2.52:2 **temper** · 1.51:3 **wind** · 2.02:2 **deed** · 0.97:2 **fame**  
· 0.58:2 **fate** · 0.26:2 **intention** · 0.06:2 **relative** ·  
0.73:1 **behoves** · 0.53:1 **judgement** · 0.2:1 **fit** ·  
0.11:1 **reception** · 0.08:1 **grace**

# Semantically motivated collocation restrictions: usage of *big* / *high* / *large*

high ...	big ...	large ...
<ul style="list-style-type: none"> <li>• level</li> <li>• [school]</li> <li>• concentration</li> <li>• speed</li> <li>• proportion</li> <li>• altitude</li> <li>• elevation</li> <li>• temperature</li> </ul>	<ul style="list-style-type: none"> <li>• [bang, band]</li> <li>• hit</li> <li>• problematic</li> <li>• break</li> <li>• difference</li> <li>• brother</li> <li>• star, bird</li> <li>• man, city</li> </ul>	<ul style="list-style-type: none"> <li>• number</li> <li>• quantity</li> <li>• amount</li> <li>• proportion</li> <li>• sum</li> <li>• portion, part</li> <li>• city, island</li> <li>• population</li> </ul>
<p>&lt;degree&gt; &lt;measure&gt;</p>	<p>&lt;size&gt; &lt;importance&gt;</p>	<p>&lt;extension&gt; &lt;quantity-mass&gt;</p>

# direct collocates vs. dependency collocates

- easy to catch: multi-word expressions (MWE), or what in English amounts to “compounds”, often with phonetic stress on the first word
  - big bang, big band
  - high tide, high society
- more difficult, often non-adjacent, profits from dependency relations
  - high temperature
    - **high** room **temperature**
    - **ambient** temperature was rather **high** when ...

# Bilingual polysemy / equivalence check: “caress” objects in Danish / Swedish

- DANISH: kærtegne

- **body parts:** bryst, krop, kind, hud, balder, mave, inderlår, brystvorte, hår, ansigt, klitoris, lår, sexbombe, nosse, røvhul, nakke, hals, kropsdel, bagdel
- **surfaces:** silkestof, græsbane
- **PROP-hum**

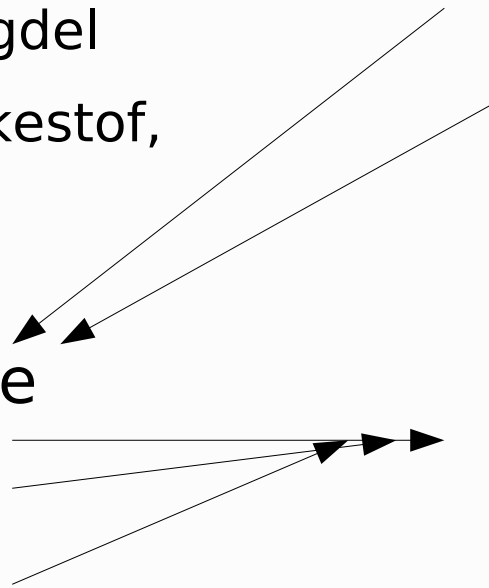
DAN: stryge

- \* swipe
- \* iron
- \* remove
- \* move ...

- SWEDISH: smeka

- **body parts:** kind, kønsorgan, bryst, stjärt, klitoris, kropp
- **PROP-hum**
- **ball:** boll, passning, tennisboll
- **instrument:** elgitarr
- **things:** lack, rännil, instrmentpanel, julle, murbrok, vidunder

SWE: stryka





# DeepDict: Verb + semantic classes (only compiled for Danish, Norwegian and Esperanto)

- uses ~ 200 semantic prototype classes for nouns
- correlations are computed the same way as for words
- offers a level of abstraction, and can compensate for sparse data

## drikke (verb)

total of 41235 relations

Hide Frequencies

### Subjects:

**PERS:** jeg, vi, han, de, man, du, der, hun, den

7.48:8 <H> · 7.38:8 **PROP-hum** · 5.9:7 <Hprof> ·

5.6:6 <Hnat> · 5.4:6 <Hfam> · 5.38:6 <HH> ·

5.48:3 PROP · 1.25:5 **dansker** · 1.74:4 <Hideo> ·

1.62:4 <Azo> · 1.62:4 <A> · 0.12:5 **ung** ·

1.21:3 PROP-org · 2.07:2 alkoholafhængig ·

1.05:3 <Adom> · 0.68:3 <Lh> · 0.24:3 <food> ·

1.37:1 medhustru · 1.37:1 nescafé · 0.39:1 Forældrene

### Accusative objects:

**PERS:** den, meget

11.48:9 <drink-h> · 9.83:9 <drink-m-h> · 9.2:9 <drink-c-h> · 8.88:9 <cm-liq> ·

8.78:9 <occ> · 7.24:9 **kaffe** · 7.47:8 <cm-liq-h> · 7.2:8 <B> · 6.93:8 øl-2 · 8.8:5 **PROP** ·

5.73:8 **vin** · 6.51:7 øl-1 · 6.31:7 <amount> · 6.24:7 <drink> · 6.23:7 **te** · 5.52:7 **alkohol** ·

6.13:6 **bajer** · 4.05:8 **vand** · 5.44:6 <con> · 5.29:6 **PROP-hum** · 5.28:6 **cola** ·

5.05:6 <cc-h> · 4.91:6 <ac> · 4.87:6 **rødvin** · 4.68:6 **whisky** · 4.46:6 <drink-m> ·

4.16:5 <mat-h> · 4.13:5 <cc> · 3.9:5 <drink-c> · 3.73:5 **PROP-org** · 3.45:5 <H> ·

3.4:5 <HH> · 3.4:5 <food-h> · 3.19:5 <am> · 3.01:5 **PROP-top** · 3.01:5 <Hprof> ·

2.67:4 <temp> · 2.59:4 <cm-h> · 2.26:4 <food-m-h> · 2.17:4 <cm> · 2.17:4 <unit> ·

1.86:4 <anorg> · 1.62:4 <per> · 1.05:3 <Hfam> · 1.05:3 <f-psych>











### Verbal particles:

4.92:6 **ud** · 3.03:6 **ihjel** · 1.28:6 **sammen** · 0.45:5 **op** · 1.47:3 **ned** · 3.41:1 **bort**

# “drikke” (drink) - object classes

- Words:
  - kaffe
  - vin
  - øl
  - te
  - alkohol
  - bajer
  - vand
  - cola
  - rødvin
  - whisky
- Semantic prototypes
  - <drink-h>, <drink-m-h>, <drink-c-h>, <drink>, <drink-m>, <drink-c>
  - <cm-liq> (*vand, urin*)
  - <amount> (*masse, mundfuld*)
  - <con> (*kop, flaske*), <unit> (*liter*)
- ambiguity artifacts:
  - kaffe <occ> <Lh>, glas <mat-h>
- metaphor:
  - <anorg> (“drink one's brain out”)
- cross-class:
  - <food-h> (*gift*)

# The background: 10 dependency parsers

Language	Parser	Lexicon	Analyzer	Grammar	Levels	Applications
	DanGram	100.000 lexemes, 40.000 names	Full	8.400 rules	morph., syntax, dep., psg, sem. roles	Teaching, corpus annotation, MT, Spell/Grammar checker, QA-systems, NER
	PALAVRAS	70.000 lexemes, 15.000 names	Full	7.500 rules	morph., syntax, dep., psg, sem. roles	Teaching, corpus annotation, MT, QA-systems, NER
	HISPAL	73.000 lexemes	Full	4.900 rules	morph., syntax, dep., psg, sem. roles	Teaching, corpus annotation
	EngGram	81.000 val/sem	Full	4.500 rules	morph./syntax, dep., psg	Teaching, corpus annotation, MT
	SweGram	65.000 val/sem	Full	8.400 rules	morph./syntax, dep., psg	Teaching, corpus annotation, MT
	NorGram	OBT / via DanGram	Full	OBT / via DanGram	morph./syntax, dep., psg	Teaching, corpus annotation, MT
	FrAG	57.000 lexemes	DTT + analysis	1.400 rules	morph.-correction, syntax, dep., psg	Teaching, corpus annotation
	GerGram	25.000 val/sem	Full (Lingsoft)	LS+1.300 rules	morph. (Lingsoft), syntax, dep., psg	Teaching, corpus annotation
	EspGram	30.000 lexemes	Full	2.600 rules	morph., syntax, dep.	Teaching, corpus annotation, MT
	ItaGram	30.600 lexemes	DTT + analysis	1.600 rules	morph., syntax, dep.	Teaching, corpus annotation

# How to use DeepDict 1

- as a lexicographer
  - find inspiration as to **complementation patterns** (selection restrictions) for dictionary entries
  - find the **most typical** (not just the most common!) **example** of a certain construction
  - find candidates for **multi-word expressions**
  - find candidates for **metaphorical usage** (often high correlation index because one of the parts is infrequent on its own)
  - find **semantic distinctions** and **subsenses** not otherwise obvious, and triggered by head or dependent words (e.g. mistænksom - mistænkelig, high - big - large)

# How to use DeepDict 2

- for teaching
  - create **lexical fields**
    - e.g. edibles, drinkables etc., via 'eat', 'drink')
    - a list of languages? countries? professions?
    - all about language
    - horse/share/oil-related words for an essay: what does it do (SUBJ), what do you do with it (ACC), how is it characterized (prenominals)
  - find **phrasal verbs / prepositional complements**
  - describe **usage differences** between near synonyms
  - distinguish between **literal** and **abstract** uses (e.g. heavy - tung - schwer) ... are some of these cross-language phenomena?
  - find **metaphors** (caress\_V)
  - find **gender differences** through pronouns (caress\_V)

# Sociolinguistic exercise: *my new neighbour is a ...* refugee - fugitive - immigrant - foreigner

- **immigrant:** illegal, German, Irish, Italian, Jewish, NUM, Chinese, recent, European, Polish, legal, undocumented, Mexican
  - subculture, (il)legality
- **foreigner:** NUM, unauthorized, undesirable, untidy, barbarian, meddlesom, friendless, unreliable, stateless, naturalized
  - (negative) focus on (lack of) assimilation
- **refugee:** Afghan, Palestinian, Jewish, genuine, Vietnamese, Kurdish, Politcal, protestant Albanian, NUM, huguenot, Palestinian, Rwandan
  - recent, reason focus
- **fugitive:** wanted, al-Qaeda, ephemeral, NUM, hunted, highest-ranking, royalist, exhausted, Russian, harried, displaced
  - process/context focus (war, factions, persecution)

# Perspectives 1: Lexicography

- DeepDict shows how syntactically related word pairs can be “harvested” from dependency- and function-annotated corpora
- It allows the lexicographer
  - not only to **find examples and frequencies** for certain (known) collocations and lexical constructions, but also
  - to **compile new lists** of such collocations and constructions
- DeepDict is a language-independent method, and could be built not only for further languages, but also for specialized or customer-built corpora (genre-variation, diachronic variation, spoken language, specific author)

# Perspectives 2: Grammars

- Better parsers, specifically, better parser lexica
- DeepDict databases can be used to harvest
  - **valency** patterns
  - semantical **selection restrictions** (generalized from words or prototypes)
  - likelihood thresholds for semantical fillers of syntactic slot
- Cyclical interplay between DeepDict-style corpus information and the parser(s) that provided it

e.g. Portuguese lexicon entry:

- pensar ('think'): <fSUBJ/H:74>, <fSUBJ/org:25>

used in grammars, e.g. +hum marking in anaphor grammar:

- ADD (£hum) TARGET PERS + @P<  
(p @PIV LINK 0 PRP-COM LINK p (<fPRP-com/H>70>))



# Perspectives 3: FrameNet

- current DeepDict: shows one relation at a time, i.e. computes e.g. subject and object fields independently of each other, which is fine for many applications, but could be taken a step further using:
- lexico-semantic frames (Berkeley FrameNet)
  - @SUBJ / <hum> ==> “**read**” <vt> ==> @ACC / <sem-r>
- same corpus annotation needs as for DeepDict
  - verbal subsenses need to be specified and semantically classified
  - senses can be structurally corroborated or supplemented interactively from corpus data
- project site: [www.framenet.dk](http://www.framenet.dk)

DeepDict

@

gramtrans.com

CorpusEye: [corp.hum.sdu.dk](http://corp.hum.sdu.dk)

Parsers: [visl.sdu.dk](http://visl.sdu.dk)

eckhard.bick@mail.dk

# Spin-offs: WebPainter

- live in-line markup of web pages
- mouse-over translations while reading

File Edit View History Bookmarks Tools Help

← Tr... My Acc... Welcom... http... Norweg... Bicycle ... Google ... GramTr...

GramTrans WebPainter  English to Danish [Return to front page](#)


WebPaint

**mouse-over translation:**  
trekanter (trekant)

optional grammar (here: **SUBJ** and **prep**)

All of the above bicycle races involve diamond frame bicycles of two triangles . An alternative is the recumbent , a bicycle on which the rider sits back with the legs horizontal . This puts the body in a position where there is less wind drag . Proponents claim it provides more portable riding , with no weight on the wrists . The recumbent is a more aerodynamic design of bicycle , and world speed records were set with them .

Bicycle races are popular all over the world , especially in Europe . The most devoted countries are Italy , Spain , Belgium , Germany , France , the Netherlands and Switzerland , although the United States has international standing , as does Australia . The USA boasts three-time Tour de France and first American winner , Greg LeMond as well as seven-time winner Lance Armstrong . Australia has seen success through Michael Rogers ( World Road Time Trial Champion , 2002 , 2004 , and



In Europe bicycle racing expresses nation prestige : German Democratic Republic postage stamp depicting Gustav Adolf Schur ( Täve ) , 1960