

# Setting up for Corpus Lexicography

Adam Kilgarriff, Jan Pomikalek, Pete  
Whitelock  
LCL & OUP

# Premise

- Corpus technology can support lexicography making it
  - *more accurate*
  - *more consistent*
  - *faster*

Rundell and Kilgarriff 2011

*Automating the creation of dictionaries*  
in Sylviane Granger's Festschrift

- This paper
  - A case study

# A new Portuguese dictionary

- OUP
- Pt-En and En-Pt
- 40,000 headwords on each side
- Pt-En starts from
  - Dictionary
    - Medium-sized Pt-Dutch
  - Corpus
    - blank sheet

# Agenda

1. Collect corpus
2. Process with best tools
3. From parser output to corpus system input
4. Finding good examples
5. Regional variants

# Status

1. Collect corpus
2. Process with best tools
3. From parser output to corpus system input
4. Finding good examples
5. Regional variants

# Corpus collection

- Big and diverse
- 100m not big enough
  - 40,000 headwords
  - 40,000th word in BNC: 27 hits

# dogfish *(noun)* British National Corpus freq = 27 (0.2 per million)

<u>modifier</u>	<u>7</u>	<u>1.0</u>
spotted	<u>1</u>	6.19
odd	<u>3</u>	4.55
spur	<u>1</u>	4.41
lesser	<u>1</u>	4.09
breeding	<u>1</u>	3.75

<u>modifies</u>	<u>2</u>	<u>0.3</u>
skin	<u>1</u>	1.38

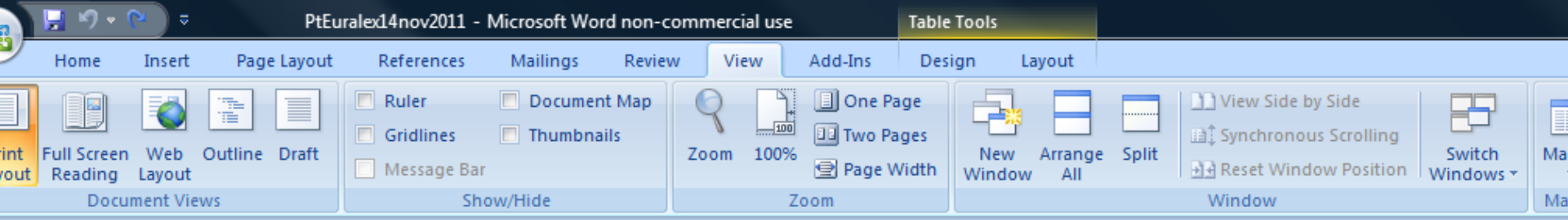
<u>and/or</u>	<u>18</u>	<u>5.5</u>
gurnard	<u>1</u>	10.25
thornback	<u>1</u>	10.19
dogfish	<u>2</u>	10.14
pollock	<u>1</u>	10.09
Scuba	<u>1</u>	9.42
pout	<u>3</u>	9.33
codling	<u>1</u>	7.83
whiting	<u>1</u>	7.48
dab	<u>1</u>	6.86
cod	<u>2</u>	6.8
ling	<u>1</u>	6.64
crab	<u>1</u>	5.21

<u>pp obj for-p</u>	<u>3</u>	<u>10.1</u>
Mwnt	<u>1</u>	12.68
mark	<u>1</u>	1.72

<u>pp obj with-p</u>	<u>2</u>	<u>9.2</u>
daylight	<u>1</u>	4.73
poor	<u>1</u>	1.17

<u>pp to-p</u>	<u>1</u>	<u>2.9</u>
lb	<u>1</u>	5.16

<u>pp obj o</u>
Catch
Nerve



dogfish (*noun*) ukWaC freq = 462 (0.3 per million) |

<u>object of</u>	<u>77</u>	<u>1.3</u>
white	<u>2</u>	8.59
dab	<u>6</u>	8.12
pout	<u>3</u>	7.63
catch	<u>8</u>	1.24

<u>modifier</u>	<u>202</u>	<u>1.5</u>
lesser-spotted	<u>9</u>	10.22
deep-sea	<u>14</u>	8.53
spiny	<u>7</u>	7.88
<u>huss</u>	<u>2</u>	7.81
wrasse	<u>11</u>	7.74

<u>modifies</u>	<u>88</u>	<u>0.6</u>
coalfish	<u>2</u>	7.61
<u>pollack</u>	<u>3</u>	7.15
conger	<u>3</u>	6.92
wrasse	<u>3</u>	5.97
ling	<u>2</u>	5.9

<u>and/or</u>
<u>Huss</u>
<u>bullhuss</u>
<u>nursehound</u>
wrasse
gulper



# Where from?

- Web
- Quantity
  - Yeah
- Quality
  - As good or better
    - Keller and Lapata 2003, Sharoff 2006, Baroni et al 2009

# How?

- New linguistics-specialist crawler
  - Was Heritrix, next time: Spiderling
    - Other talk
- Cleaning *including language-identification*
  - jusText
    - Pomikalek thesis
- Deduplication
  - Onion
    - Pomikalek thesis

# Processing tools

- Reviewed options
  - Best: ***Palavras***
    - Bick 2000
    - + ongoing development since
  - Contacted author, negotiated licence
  - Installed
  - Applied to 2b words

# Vast process

- Parsing is usually slow - would it take years?
- Parallelised in 12 processes
- Many bugs encountered, resolved with developers
- Crashed on many input files – leave them out
- Final run: 15 days

# Corpus creation stats

	European	Brazilian
HTML data downloaded	1.10 TB	1.37 TB
Unique URLs	31.5 million	39.1 million
Crawling time	8 days	10 days
Final corpus size (words)	0.7 billion	1.0 billion

# From dependency parse to word sketch

- Palavras: *dependency parser*
- Output for each word
  - Lemma, pos tag
  - “my governor is word N”
  - “relation is ...”
- Like CONLL output
- SKEW-2, Siva Reddy
  - Word sketches from CONLL format data

# To get better word sketches

- Parser output and lexicographic word sketches
  - *Not quite the same*
- Anomalies in parser output
- Large project

# Preposition-article contractions

→ satisfação [satisfação] <cjt> <act> <percep-f> N F S @<ACC #19->17  
→ de [de] <sam-> <np-close> PRP @N< #20->19  
→ os [o] <-sam> <artd> DET M P @>N #21->23  
nossos [nosso] <poss 1P> DET M P @>N #22->23  
→ clientes [cliente] <Hattr> N M P @P< #23->20

19 satisfação satisfação N F:S 14,V obj %w\_N/%w\_V obj  
→ 20 dos de PRP 19,NIL/%w\_N dep PRP  
21 nossos nosso DET M:P 22,NIL/DET spec\_of %w\_N  
22 clientes cliente N M:P 19,N\_de\_%w\_N/%w\_N\_de\_N



# Verb form reconstruction

Deveria- [dever] <\*> <hyfen> <fmc> <aux> V COND 3S VFIN @FS-STA #1->0  
se- [se] <hyfen> PERS M/F 3S/P ACC @<SUBJ #2->1  
começar [começar] <vH> <mv> V INF @ICL-AUX< #3->1

- 1 Dever-se-ia dever V COND:3S:VFIN 1,REFL-SUBJ
- 2 começar começar V INF 1,V comp %w\_V/%w\_V comp V

# Multi-word unpacking

A=Comunidade=de=Direitos=Humanos

[A=Comunidade=de=Direitos=Humanos] <tit> <\*> PROP F P @NPHR #2->0



<mwe parsed="yes" pos="PROP">


2 A o DET F:S 3,NIL/DET spec\_of %w\_N



3 Comunidade comunidade N F:S



4 de de PRP 3,NIL/%w\_N dep PRP



5 Direitos direito N M:P 3,N\_de\_%w\_N/%w\_N\_de\_N



6 Humanos humano ADJ M:P 5,N mod %w\_ADJ/%w\_N mod ADJ

</mwe>

# Trinary Relations/Coordination

um [um] <arti> DET M S @>N #17->18

simulador [simulador] <H> N M S @<SC #18->0

de [de] <np-close> PRP @N< #19->18

inclinação [inclinação] <cjt-head> <percep-f> <am> N F S @P< #20->19

e [e] KC @CO #21->20

direção [direção] <cjt> <HH> <dir> <Ltop> N F S @P< #22->20

17 um um DET M:S 18,NIL/DET spec\_of %w\_N

18 simulador simulador N M:S

19 de de PRP 18,NIL/%w\_N dep PRP

20 inclinação inclinação N F:S 18,N\_de\_%w\_N/%w\_N\_de\_N

21 e e KC

22 direção direção N F:S 18,N\_de\_%w\_N/%w\_N\_de\_N;  
20,N e|ou %w\_N/%w\_N e|ou N

# Control relations

não [não] ADV @ADVL> #3->4

é [ser] <vK> <fmc> <mv> V PR 3S IND VFIN @FS-STA #4->0

viável [viável] <nh> ADJ F S @<SC #5->4

sua [seu] <poss 3S> DET F S @>N #6->7

aplicação [aplicação] <act> <sem-r> N F S @<SUBJ #7->4

3 não não ADV 4,%w\_ADV mod\_of V/ADV mod\_of %w\_V


4 é ser V PR:3S:IND:VFIN 7,N subj\_of %w\_V/%w\_N subj\_of V

5 viável viável ADJ F:S 4,V dep %w\_ADJ/%w\_V dep ADJ;  
7,N subj\_of %w\_ADJ/%w\_N subj\_of ADJ


6 sua seu DET F:S 7,NIL/DET spec\_of %w\_N

7 aplicação aplicação N F:S

# Reanalysis



variados [variar] V PCP M P @>N #21->22  
aspectos [aspecto] <cjt-head> <ac-cat> N M P @<SUBJ #22->17  
de [de] <sam-> <np-close> PRP @N< #23->9  
a [o] <-sam> <artd> DET F S @>N #24->25  
tecnologia [tecnologia] <domain> N F S @P< #25->23



20 variados variar V PCP:M:P  
21 aspectos aspecto N M:P 35,N subj\_of %w\_N/%w\_N subj\_of N  
22 da de PRP 21,NIL/%w\_N dep PRP  
23 tecnologia tecnologia N F:S 21,N\_de\_%w\_N/%w\_N\_de\_N

# Lemmatization

Old spelling	New spelling
acto	ato
carbónico	carbônico
cabeça-de-burro	cabeça de burro
concetual	conceptual
auto-sugestão	autossugestão

Female form	Male form
amiga	amigo

# GDEX

- Good dictionary example finder
- Customise for Portuguese
  - Follow Slovene lead

# Regional variation

- European vs Brazilian
- Method
  - Keyword list of each vs other
  - If in top 1%: add note to word sketch