

# Experiments in human-computer cooperation for the semantic annotation of Portuguese corpora\*



<http://www.linguateca.pt/ACDC/>

<http://www.linguateca.pt/ACDC/corte-e-costura/>

Diana Santos [Diana.Santos@sintef.no](mailto:Diana.Santos@sintef.no)

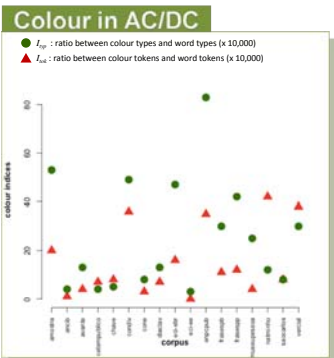
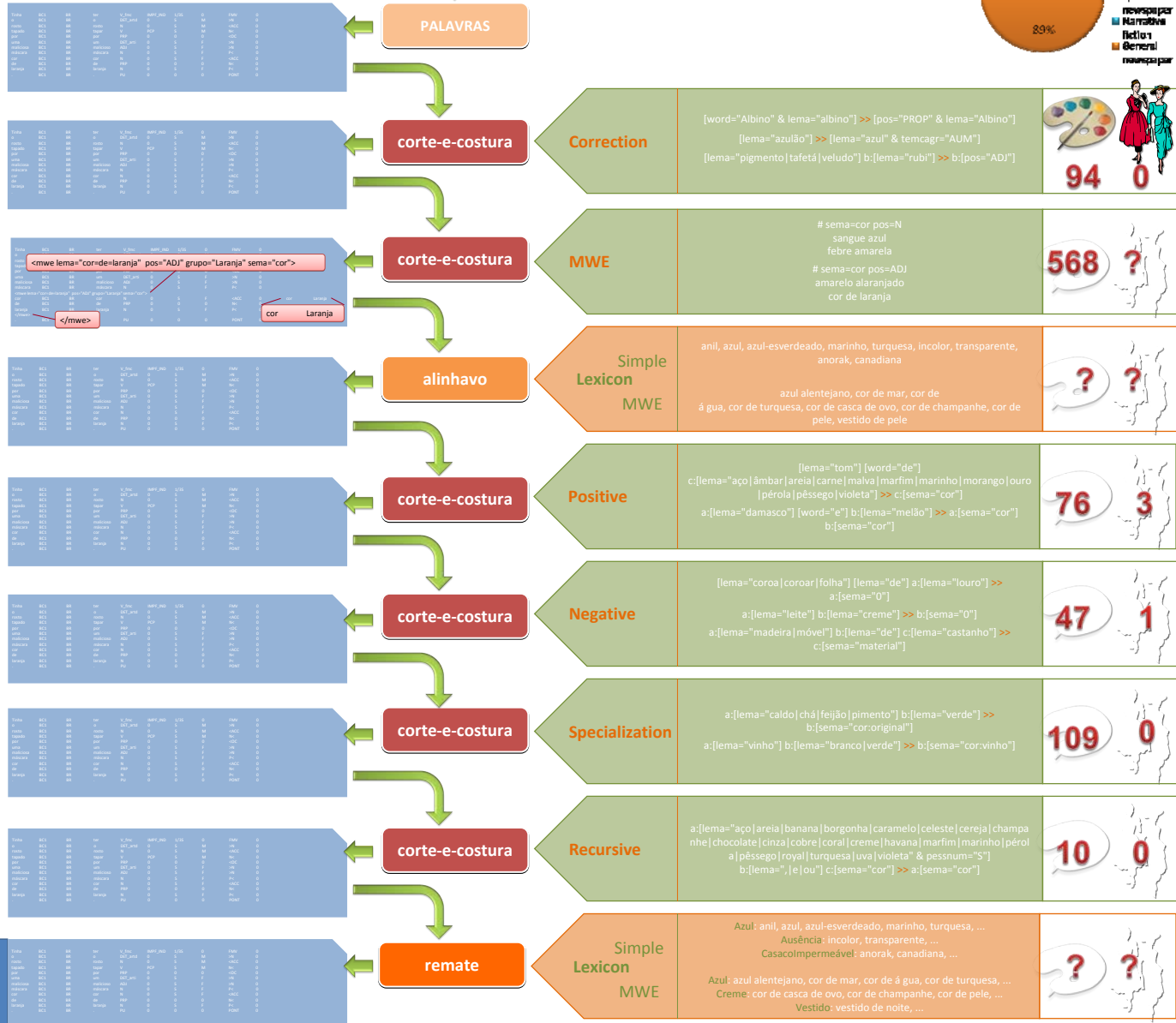
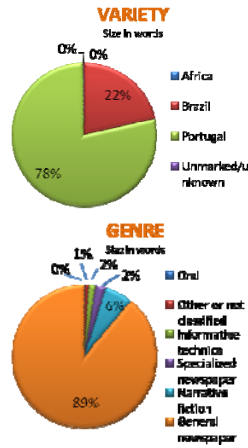
Cristina Mota [cmota@ist.utl.pt](mailto:cmota@ist.utl.pt)

## Motivation

- Every corpus a different formalism and search expression
- Every corpus different undocumented assumptions and terminology
- No replication or addition to already existing corpora possible
- No revision/feedback by users possible

AC/DC answer: cooperative annotation of already richly-annotated corpora

AC/DC Corpora					
Corpus	Units	Words	Sentences	Var.	Short description
Amostra	137.832	98.796	4.955	B	Pos-tagged sample of NILC Corpus
ANCIB	1.690.376	1.258.764	80.992	B	(Moderated) Brazilian discussion list on librarianship
Avante!	7.766.418	6.501.257	204.414	P	Portuguese party-political newspaper Avante!, 1997-2002
CD HAREM	222.407	147.077	8.185	all	Golden collection of the First HAREM
CETEMPúblico	232.543.379	189.575.095	7.665.410	P	Major Portuguese daily newspaper, 1991-1998
CHAVE	123.868.725	99.478.954	4.740.448	P B	News from Público and Folha de São Paulo, 1994-1995
ClassiPPE	1.922.601	1.304.282	74.690	P	Portuguese fiction, drama and poetry, 15th and 19th century
CONDIVort	7.088.775	5.577.632	328.214	P B	Sports, fashion and health magazines (50s, 70s and 2000)
CoNE	925.230	685.225	31.561	any	Spam or general e-mail messages
DiaCLAV	7.758.467	6.651.549	232.152	P	Four Portuguese regional daily newspapers
ECl-EBR	917.127	724.015	44.381	B	Corpus Borba-Ramoses of Brazilian Portuguese
ECl-EE	32.034	27.140	839	P	Call for the EU ESPRIT program
ENPCPub	92.693	72.389	4.371	P B	Translated fiction from English, subset of the ENPC corpus
FrasesPB	23.313	19.162	653	B	Individual sentences in Brazilian Portuguese
FrasesPP	20.049	16.233	594	P	Individual sentences pos-tagged
MuseuPessoa	517.747	375.158	27.388	P B	Transcriptions of oral interviews from Museu da Pessoa
Natura/Minho	2.156.187	1.749.083	68.910	P	Unedited version of regional Portuguese newspaper
Natura/Público	7.369.349	6.274.542	225.752	P	Two first paragraphs of each Público article, 1991-1994
NILC	42.608.038	32.342.456	1.963.795	B	Newspaper, commercial letters, textbooks, etc.
Vercial	18.854.273	14.315.992	996.869	P	Portuguese fiction, poetry and drama, 16th to 20th century
Total	457.707.958	368.107.685	16.342.734		



Rule activation & Corpora changes						
Corpus	Rules fired	Words changed	Word types changed	Attr. changed	New MWE	MWE changed
Amostra	30	312	142	570	5	3
ANCIB	93	388	106	641	21	3
Avante!	707	4237	372	7641	40	34
CHAVE	19965	132244	1998	242486	2796	700
CoNE	27	393	105	642	1	0
DiaCLAV	1973	7513	481	13116	444	25
ECl-EBR	150	1969	412	3699	25	3
ENPCPub	225	377	162	637	6	2
MuseuPessoa	71	751	137	1451	7	1
Natura/Minho	553	1388	224	2093	39	2
Vercial	2921	50246	1553	93831	546	99

## Conclusions

- It is possible with little effort and with a large user community behind
- Remains to be seen if the remaining annotation corpora are used

Note: Reuse of OpenCWB and (Open)PALAVRAS a must!  
Without them we would again have started from scratch...

- OpenCWB <http://cwb.sourceforge.net/>
- (Open)PALAVRAS <http://visl.sdu.dk/visl/pt>

\* All numbers dated May 5th, 2010