

Nos bastidores da Gramateca: uma série de serviços

ALBERTO SIMÕES E DIANA SANTOS

<http://www.linguateca.pt>

MOTIVAÇÃO

Demasiado texto anotado para podermos conhecê-lo sem ferramentas informáticas

- potenciar a procura e apresentar a repartição
- permitir comparações
- permitir análises alternativas
- visualizar sob diferentes ângulos

TEMAS DE INVESTIGAÇÃO

- Usabilidade
- Comparação de anotadores sintáticos
- Gramática baseada em corpos
- Estudo de diferenças entre teorias linguísticas

AC/DC

Acesso a corpos, disponibilização de corpos

- servindo 27 corpos diferentes, 1,5 mil milhões de palavras
- anotada com muita informação associada

Há 15 anos servindo as comunidades do processamento computacional da língua portuguesa e da linguística, graças ao PALAVRAS e ao OpenCWB (<http://cwb.sourceforge.net/>).

SERVIÇOS

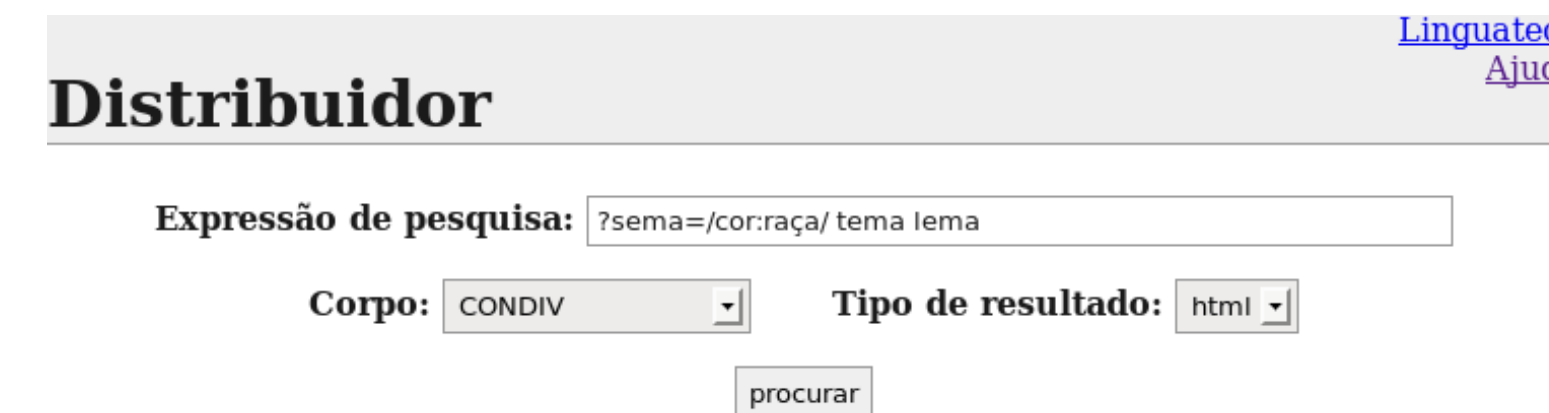
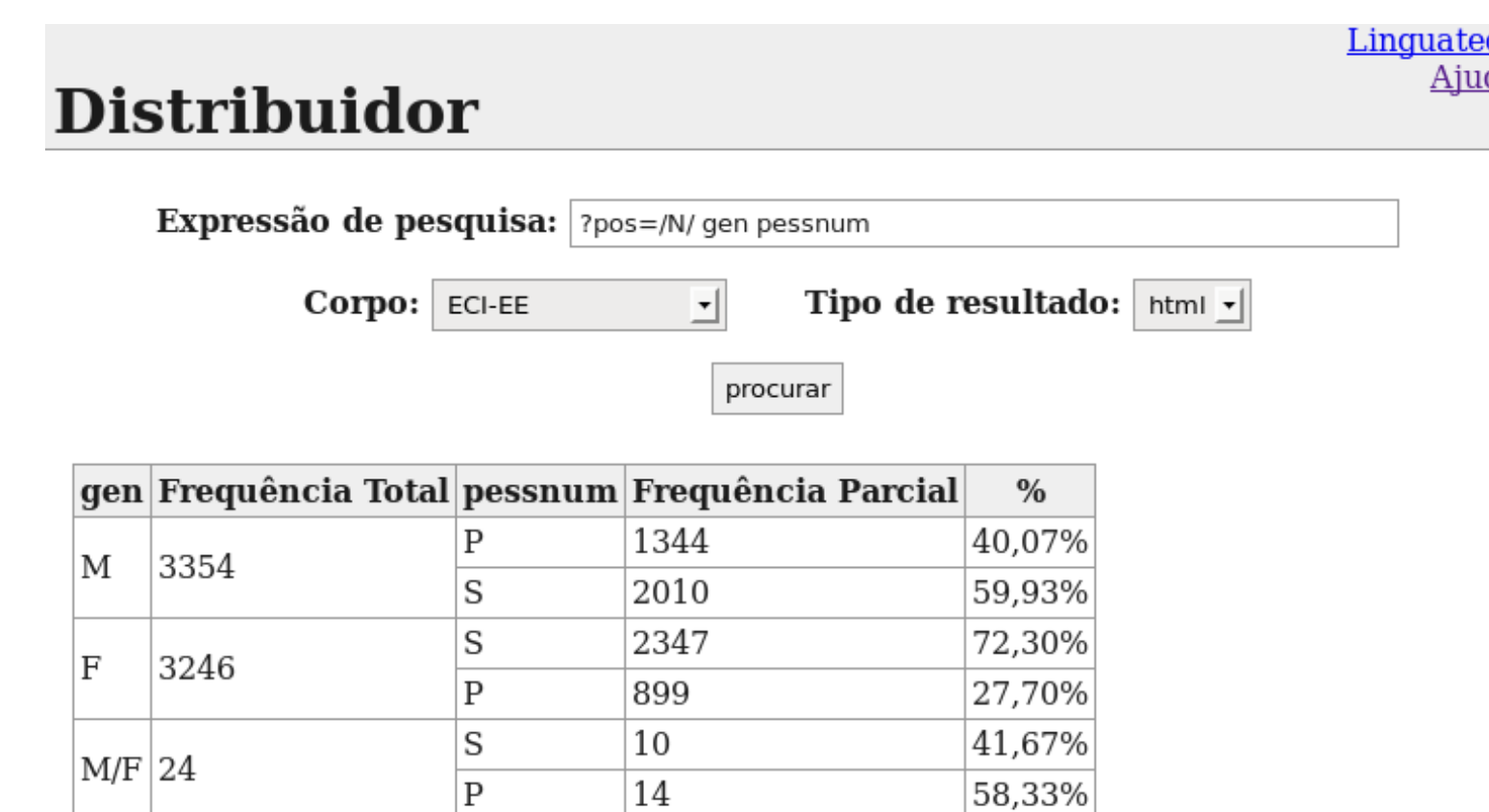
- Procura** <http://www.linguateca.pt/ACDC/>
Distribuidor <http://www.linguateca.pt/Distribuidor/>
Comparador <http://www.linguateca.pt/Comparador/>
Ordenador <http://www.linguateca.pt/Ordenador/>
VARRA <http://www.linguateca.pt/VARRA/>
Ensinador <http://www.linguateca.pt/Ensinador/>
Ensinador paralelo <http://dinis2.linguateca.pt/paraensinador/>
Rêve <http://www.linguateca.pt/Reve/>
Repositório GitHub <https://github.com/linguateca>

REFERÊNCIAS

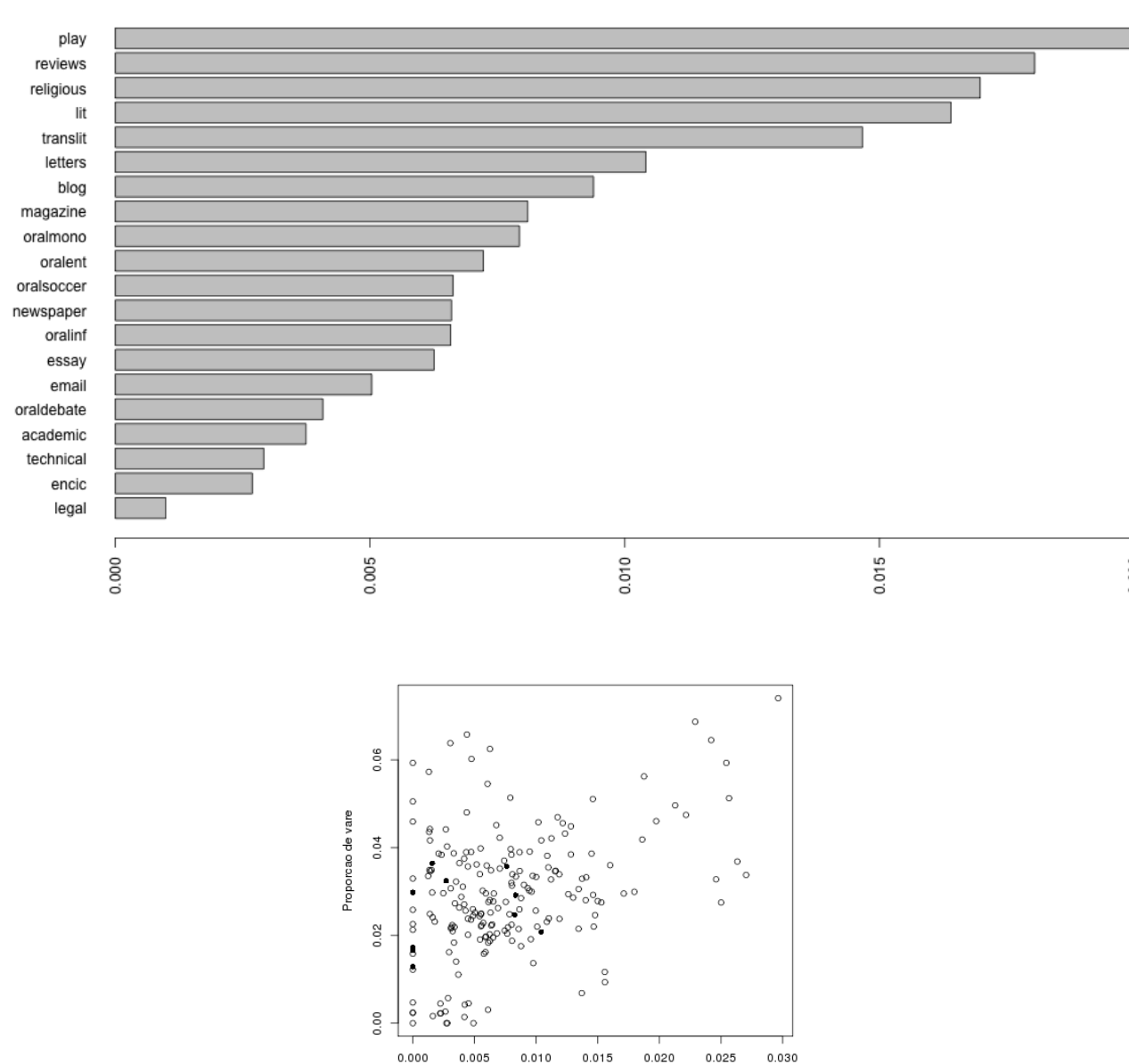
- Luís Costa, Diana Santos e Paulo Alexandre Rocha. 2014. Estudando o português tal como é usado: o serviço AC/DC. In *The 7th Brazilian Symposium in Information and Human Language Technology (STIL 2009)* (São Carlos, Brasil, 8-11 de Setembro 2009).
- Eckhard Bick. 2000. *The Parsing System "Palavras": Automatic Grammatical Analysis of Portuguese in a Constraint Grammar Framework*. Aarhus University Press.
- Cláudia Freitas, Diana Santos, Hugo Gonçalo Oliveira & Violeta Quental. 2014. VARRA: Validação, Avaliação e Revisão de Relações semânticas no AC/DC. In Simone Sarmento, Tony Berber Sardinha, Livia Pretto Mottin & Ana Maria T. Ibaños. *Pesquisas e perspectivas em linguística de corpus*, Mercado de Letras, pp. 199-232.
- Diana Santos. 2014. Corpora at Linguateca: Vision and roads taken. In Tony Berber Sardinha & Telma de Lurdes São Bento Ferreira (eds.), *Working with Portuguese Corpora*, Bloomsbury, 2014, pp. 219-236.
- Alberto Simões & Diana Santos. 2011. Ensinador: corpus-based Portuguese grammar exercises. In *Procesamiento del Lenguaje Natural*, 47, pp. 301-309.

EXEMPLOS DA SUA INVOCACÃO

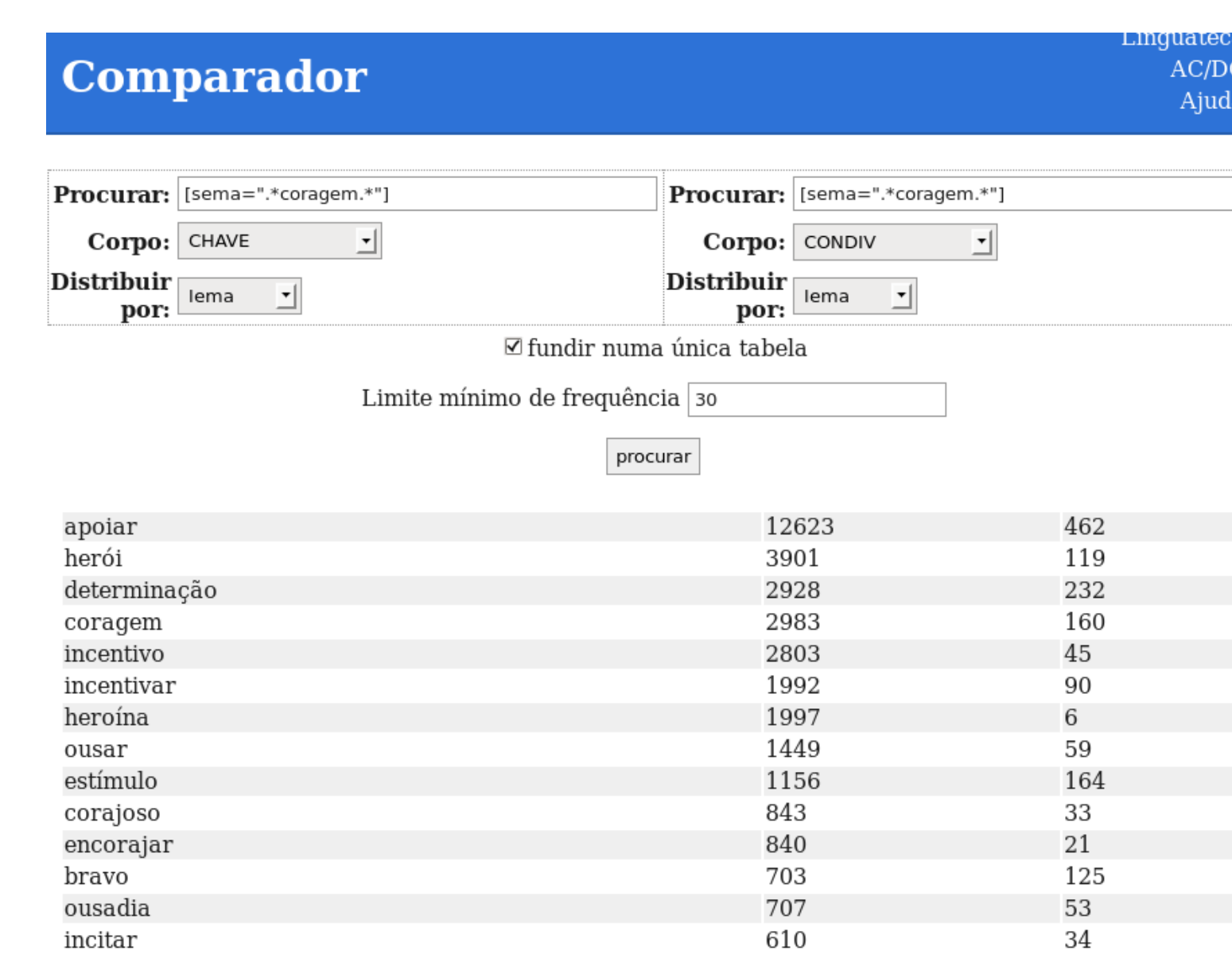
Distribuição mais fina



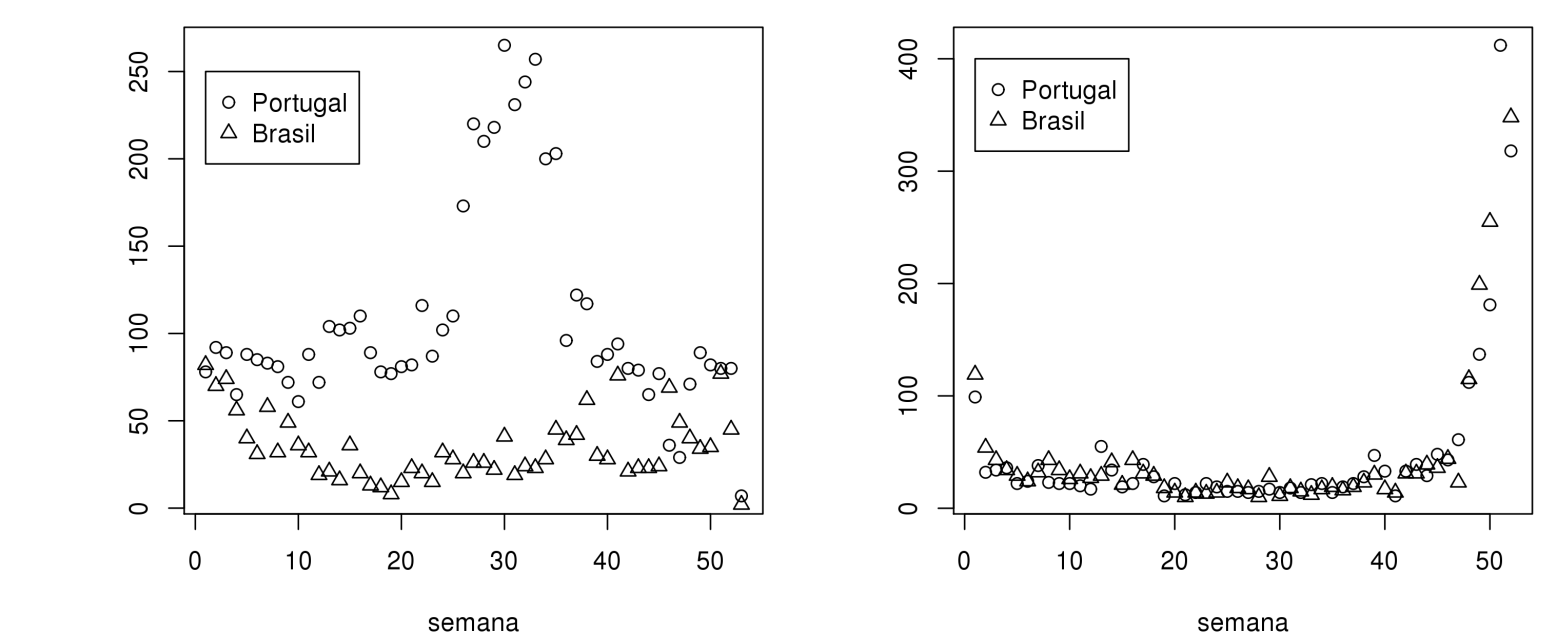
Capacidades de visualização



Comparação de procuras



Inspeção fina de correlações



Criação de material didático

Ensinador

A procurar " [word="quandoise"]- [pos!="V.*"]- [temagr=".*SUBJ"]lema" no corpus C-Oral-Brasil v. 3.2 [121 entradas.]

Seleccione as concordâncias que deseja usar.

- Ah, é, uai // 24 se n morrer antes deles. (morrer)
- Lá em casa, aviso se cê for lá. (ser)
- Se eu fizesse // 280 se eu fizesse. (fazer)
- Eu, se eu fosse cê fazia isso não, Carão // 204 se eu fosse cê matava aqui o \$ / / 205 eu te avisei, Zé // 206 p 'cê + CAR (ser)

Ensinador Paralelo

Eu não sei o que se passa, mas há coisa misteriosa que eu não consigo adivinhar. er. Far har holdt seg inestetligt sammen med fetter Baltasar bebo.

Falou-me em criados mortos; mas eu não entendi... Jeg kunne ikke høre helt hva hun sa til meg.

Prometa partir para Coimbra logo que o _____ fazer sem já ai Teresa. behandling. Han loved å dra til Coimbra så snart han var vis ikke ville bli utsatt for noe ubehag i hans fravær.

Um apenas adiantara coisa que não alimiar a justiça. et til. Like etter. En av dem la til. uten at det hjalp de som e vicia a ser que o mató. nas vizinhanças do local. fora startet. at småskogen i nærheten av funnstedet var hogget chapotado.

Nesta escuridade a justiça não dar passo algum. Utten flere opplysninger kunne ikke myndighetene utrette noe som helst.

Tanto ao velho como ao morgado convinha apagar algum indício que _____ envolvesse no mistério daquelas duas. Det passet både Tadesu og Baltasar godt at enhver opplysning som kunne blande dem inn i drapsmysteriet ble undertrykt.

Reanotação de dados anotados

Rêve

Partes do corpo

- Classifique a palavra segundo o seu sentido:
- referindo-se ao corpo humano (corpo):
 - usando uma palavra do corpo para outros fins:
 - corpo:faculdade
 - corpo:outros
 - corpo:animal
 - corpo:partidoobjeto
 - corpo:sentimento
 - corpo:vegetal
 - corpo:grupo
 - corpo:doença
 - corpo:posição
 - corpo:emoção
 - corpo:lugar
 - corpo:movimento
 - ou não se referindo de todo ao corpo humano: 0

Se a palavra fizer parte de uma expressão maior relacionada com o corpo, indique essa conexão no campo de comentário, com o sufixo EVP, por exemplo *corpoEVP* ou *corpo:lugarEVP*. Quando considerar que mais de uma classificação é possível, indique ambas no comentário ligadas por ...

- E001-PT332: E na aula recebem-me com uma salva de **palmas**. corpo:outros
- E001-PT905: Eu andei com caixotes às **costas**, todos nós. corpo
- E003-PT24: A ferramenta é esta que tenho na **mão**: a colmeadeira, o lareiro e a estaca. corpo
- E103-PT39: Eu enjoo muito e como fui de carreira, ia muito "amorrinhada", uma mulher que lá atrás de mim, num sítio qualquer bateu-me nas **costas** e disse-me: "Saia aqui!". corpo
- E105-PT91: Os guardas passaram por mim e eu levava o alforçes às **costas** e eles não me viram. corpo
- E133-BR-1060: Meu irmão já até levantou a **mão** para minha mãe, ela tem medo que ele faça alguma besteira. corpo:outros corpo:outrosEVP, etc

AGRADECIMENTOS

A Linguateca foi financiada pelo governo português e pela União Europeia através do MCTES, a UMIC, e a FCT de 1998 até 2010, e gerida e apoiada pela infraestrutura informática da FCCN desde 2000. O trabalho mais recente foi financiado pela Universidade de Oslo. Estamos também gratos ao Research Computing Group desta universidade.

