

EXPLORANDO DISFLUÊNCIAS NO CÓRPUS COBRA-7, UM CÓRPUS EXTENSO DE PROCESS-WRITING DE APRENDIZES DE INGLÊS COMO LÍNGUA ESTRANGEIRA

Pesquisa de mestrado

Autor: Wendel Mendes Dantas

Orientador: Tony Berber Sardinha

Instituição: PUC-SP

Linha de pesquisa: Linguagem, tecnologia e educação

Início da pesquisa: 02/2010

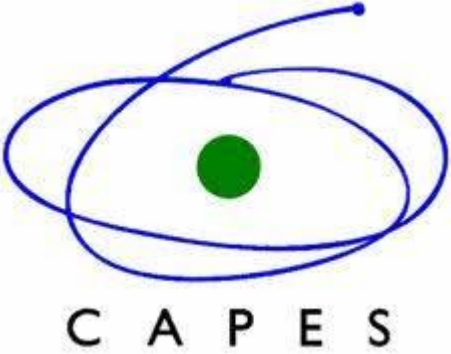
Previsão de defesa: 02/2012

Andamento: análise piloto do corpus disponível

E-mail: wendel@bilinguismo.org

Grupo de estudos: <http://corpuslg.org/gelc/gelc.php>

Agradecimentos





Perguntas de pesquisa

- (1) Quais são os padrões léxicogramaticais mais típicos no cópús COBRA-7?
- (2) Quais são as frequências dos padrões léxicogramaticais típicos do COBRA-7 nos córpora de referência?



Metodologia

1. Escolha dos corpóra de estudo e referência;
2. Compilação do corpús de estudo COBRA-7 (Corpús Brasileiro de Aprendizes de Inglês como Língua Estrangeira – Seven Idiomas);
3. Feitura de uma lista *Index* de palavras do corpús de estudo para cada nível utilizando o comando “*Make/Add to index*” da ferramenta *Wordlist*, do programa computacional Wordsmith Tools versão 5.0 (Scott, 2008) (doravante WST5);
4. Feitura de uma lista *Index* de palavras do BNC (*British National Corpus*) (corpús de referência) utilizando o comando “*Make/Add to index*” da ferramenta *Wordlist*, do programa computacional WST5);



Metodologia

5. Feitura de listas de *clusters* de três palavras para cada nível analisado e para o córpus de referência (BNC), usando a ferramenta *Wordlist*, do programa computacional WST5, comandos “*Compute/Clusters*” .
6. Feitura de listas de *clusters-chave* para cada um dos níveis analisados utilizando como córpus de referência o BNC, por meio da ferramenta *Keywords*, do programa computacional WST5.
7. Padrões dos *clusters-chave*: Observação dos padrões léxicogramaticais das 10 primeiras ocorrências dos *clusters* de três palavras na lista de *clusters-chave* do nível pré-intermediário, comparando-os com os dois outros níveis.



Metodologia

8. Como os 10 primeiros *clusters*-chave da lista correspondente ao nível pré-intermediário não foram todos encontrados nos outros dois níveis pelo aplicativo *Keywords* do programa computacional WST5, realizou-se, então, uma busca desses *clusters*-chave para esses níveis utilizando o aplicativo *Concord*, do WST5.

9. Normalização dos resultados por 1000.

10. Busca dos *clusters*-chave no COCA, conseguidas digitando-os no campo de busca do *site* e selecionando-se a opção CHART, que fornece um gráfico estatístico dos subcórpora nos quais os *clusters* são encontrados.



Metodologia

11. Normalização dos resultados dos corpóra de estudo e referência por 1000 para criar um quadro comparável adequado.

12. Observação dos resultados listados.

Piloto

Padrões Léxicogramaticais presentes nas listas de *clusters-chave* da versão 1 das redações nos três níveis analisados no *córpus de estudo* (valores normalizados por 1000):

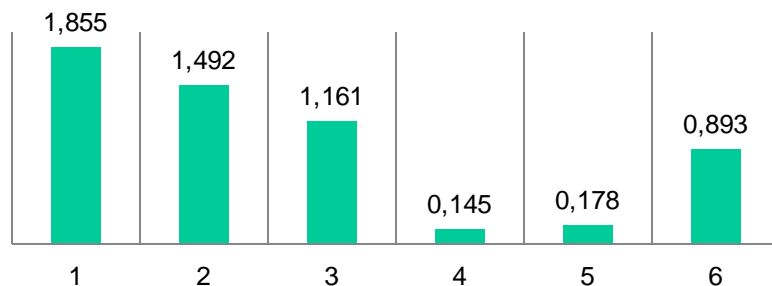
Padrões léxicogramaticais	Freq. COBRA-7			Freq. BNC	Freq. COCA*	
	Pre-inter	Inter.	Inter-sup		Escrito	Oral
Verbo + A lot of + subst. sing/pl	1,855	1,492	1,161	0,145	0,178	0,893
I used to + verbo	0,699	0,056	0,189	0,024	0,013	0,024
It's a + adj./subst.	0,602	0,338	0,162	0,141	0,153	0,528
I want to + verbo/fim de oração	0,554	0,507	0,459	0,048	0,060	0,243
Suj./to + Go to the + subst.	0,530	0,450	0,300	0,032	0,030	0,074
Verbo + To work with + obj..	0,506	0,056	0,081	0,010	0,015	0,026
I don't + verbo	0,506	0,676	0,486	0,357	0,312	1,010
I was very + adj./verbo participio	0,506	0,113	0,189	0,007	0,004	0,018
It was a +subst./adj.	0,506	0,422	0,378	0,134	0,111	0,202
subst./pron./verbo + To go to + obj.	0,482	0,370	0,320	0,040	0,040	0,105

* Este valor foi dividido por 1000, pois o COCA é normalizado por 1 milhão.

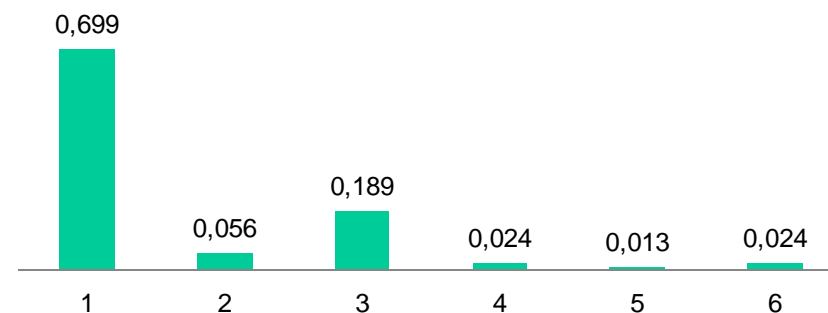
Piloto



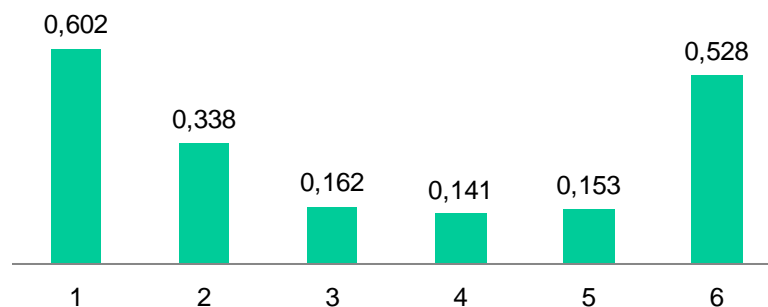
A lot of



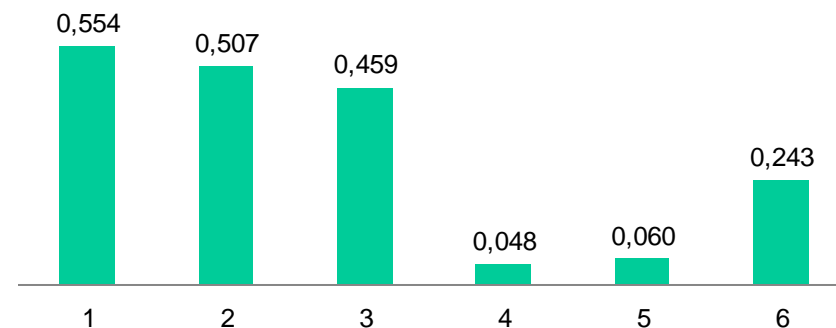
I used to



It's a

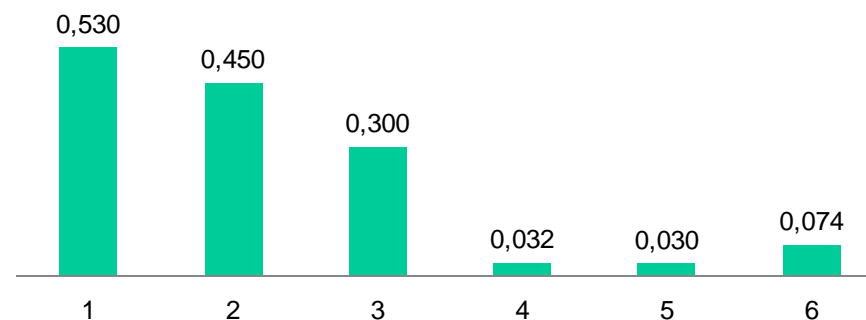


I want to



- 1 = Pre-intermediário
- 2 = Intermediário
- 3 = Intermediário superior
- 4 = BNC
- 5 = COCA (Escrito)
- 6 = COCA (Oral)

Go to the





Piloto

Resultados auferidos até o momento:

Os dados apontam um uso dos *clusters*-chave pesquisados nas composições dos alunos superior àquele observado nas composições dos falantes nativos do inglês, coletados nos *córpore* de referência.

Mostram ainda que o uso dos *clusters*-chave pesquisados no *córpore* de estudo se aproxima mais da oralidade dos falantes nativos do inglês do que da produção escrita.

Porém, observa-se entre os níveis pré-intermediário e intermediário superior uma maior aproximação do padrão de ocorrência revelado pelos *córpore* de referência.

