

Instituto Federal do Espírito Santo
Bacharelado em Sistemas de Informação

Sistema *web* de arquitetura modular para
processamento de corpora

Ádler Oliveira Silva Neves

Orientador: Me. Ernani Leite Ribeiro Filho

Sumário

- 1 Introdução ao tópico
- 2 Objetivos
 - Geral
 - Específicos
- 3 Desenvolvimento
 - Um *overview* sobre os objetivos
 - Aprofundando nos objetivos
- 4 Resultados obtidos
 - Treino do etiquetador
 - O sistema desenvolvido
- 5 Conclusão
 - Trabalhos futuros

O que falta nos atuais?

Software	Interface	Gratuito	Corpus fornecido pelo usuário	Tokenizador	Etiquetador	Concordanciador
WordSmith 5	Desktop	Não	Sim	Sim	Não	Sim
Unitex/GramLab	Desktop	Sim	Sim	Sim	Francês	Sim
CorpusEye	Web	Sim	Não	Sim	Sim	Sim
COCA Online Corpus	Web	Sim	Não	Sim	Sim	Sim
NoSketch Engine	Web	Sim	Sim	Sim	Não	Sim
Sketch Engine	Web	Não	Sim	Sim	Sim	Sim
Corpus Slayer	Web	Sim	Sim	Sim	Sim	Sim

Tabela 1: Tabela comparativa resumida de softwares de Processamento de Linguagem Natural

Fonte: O autor

Objetivo geral

Desenvolver uma aplicação *web* de código aberto para marcação e busca de partes do discurso em corpora, visando ampliar as funcionalidades em relação a *softwares* similares existentes e com interface amigável ao usuário.

Objetivos específicos

- Analisar comparativamente os recursos das ferramentas WordSmith, CorpusEye, COCA Online Corpus, Unitex/GramLab e Sketch Engine;
- Desenvolver ou adaptar um módulo extrator de sentenças;
- Desenvolver ou adaptar um módulo extrator de lista de palavras;
- Desenvolver ou adaptar um módulo etiquetador de partes do discurso que atue sobre sentenças;
- Desenvolver ou adaptar um módulo concordanciador que suporte busca por etiquetas;
- Integrar os módulos desenvolvidos ou adaptados numa aplicação web;
- Disponibilizar uma ferramenta livre para uso educacional.

Análise comparativa

- Tabela 4, seção 3.1; (p. 29)
- Várias funcionalidades desejáveis:
 - Não há tempo hábil para implementar todas;
- A tabela 1 era uma versão resumida desta;
- Deu origem aos objetivos subsequentes.

Um *overview* sobre os objetivos

Separador de sentenças

Separador de frases

O que faz: Separa um texto em frases;

Desafios: Siglas, abreviações e abreviaturas; (Sr., Sra., V.Exa.)

Implementação: Adaptada do Unitex/Gramlab;

```
Senhora.snt (/home/adler/workspace/Unitex-GramLa...
5655 sentence delimiters, 176549 (12019 diff) tokens, 78229 ...
77592 occurrences (21198 DLF entries) simple words, 1562 oc...

Senhora
{S}José de Alencar
{S}Projeto Lácio-WEB
{S}Disponível em:
http://www.nilc.icmc.usp.br/lacioweb/
{S}Ao Leitor
{S}Este livro, como os dois que o precederam, não
são da própria lavra do escritor, a quem
geralmente os atribuem.
{S}A história é verdadeira:{S} e a narração vem de
pessoa que recebeu diretamente, e em
circunstâncias que ignoro, a confidência dos
principais atores deste drama curioso.
{S}O suposto autor não passa rigorosamente de
editor.{S} É certo que tomando a si o encargo de
corrigir a forma e dar-lhe um lavor literário, de
```

Figura 3: Lista de sentenças do Unitex/GramLab sobre o livro Senhora de José de Alencar

Fonte: O autor

Um *overview* sobre os objetivos

Extrator de lista de palavras

O que faz: Identifica as palavras do texto e as classifica como simples ou composta, seu lema, e suas possíveis flexões;

Implementação: Adaptada do Unitex/Gramlab;

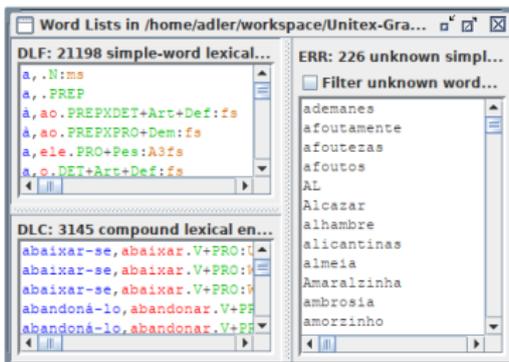


Figura 4: Lista de palavras do Unitex/GramLab sobre o livro Senhora de José de Alencar

Concordanciador

O que faz: “extrai todas as ocorrências de uma palavra de busca num corpus juntamente com seu cotexto [...]”⁴;

Implementação: Própria.

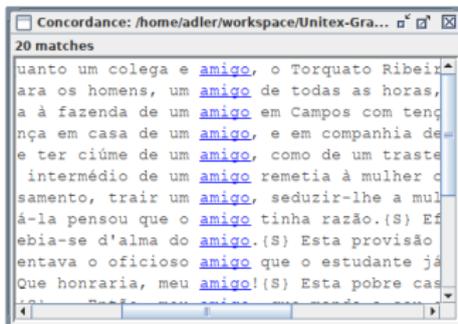


Figura 5: Concordâncias para a palavra “amigo” gerada pelo UniteX/GramLab sobre o livro Senhora de José de Alencar

Fonte: O autor

⁴TAGNIN, S. E. **Glossário de linguística de corpus**. São Paulo: HUB Editorial, 2010. p. 358.

Etiquetador de partes do discurso

O que faz: Atribui etiquetas de partes do discurso a cada palavra da sentença;

Desafios: Ambiguidades; (“casa” é verbo ou substantivo?)

Implementação: Se o treinamento do Unitex/Gramlab obtiver precisão⁵ maior que 75%, será utilizado o etiquetador deste; caso contrário, serão utilizados os etiquetadores treinados por Aires⁶, priorizados por precisão.

⁵“denota a proporção de casos preditos como positivos que são considerados positivos reais” POWERS, D. Evaluation: From Precision, Recall and F-Measure To Roc, Informedness, Markedness Correlation. **Journal of Machine Learning Technologies**, v. 2, n. 1, p. 37–63, 2011. ISSN 2229-3981. p. 38.

⁶AIRES, R. V. X. **Implementação, adaptação, combinação e avaliação de etiquetadores para o português do Brasil**. 154 p. Tese (Doutorado) — Universidade de São Paulo, São Carlos, 2000.

Um *overview* sobre os objetivos

Etiquetador de partes do discurso

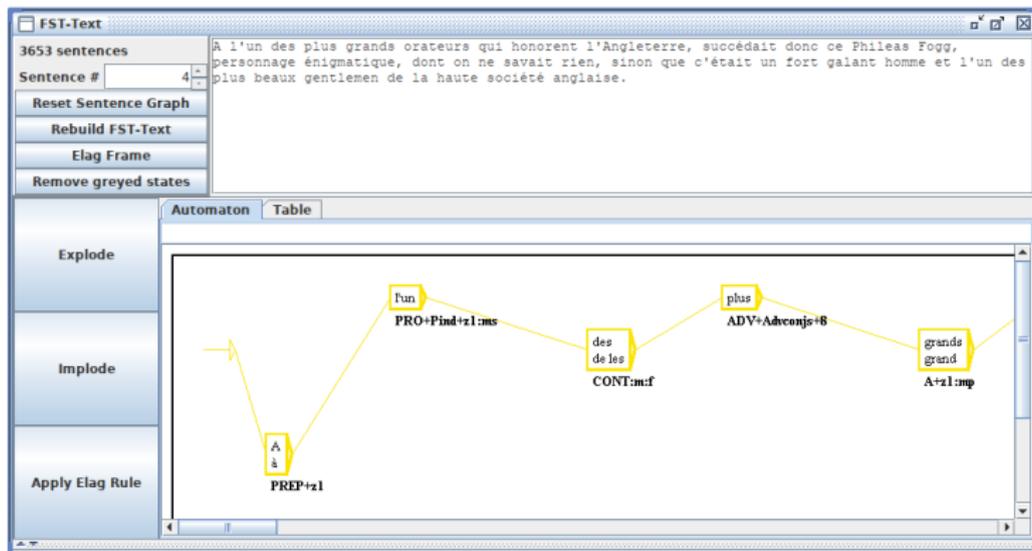


Figura 6: Etiquetador do Unitex/GramLab sobre o livro A Volta ao Mundo em 80 Dias de Júlio Verne

Fonte: O autor

Concordanciador

- Problemas a serem resolvidos:

The screenshot shows the search interface of the COCA Online Corpus concordancer. At the top, there are navigation tabs: "List", "Chart", "Collocates", "Compare", and "KWIC", with "KWIC" being the active tab. Below the tabs is a search input field, currently empty, with a "[POS]" label to its right. Underneath the input field is a row of buttons: "L", followed by six hyphens, "R", and an asterisk. Below this row is a "Keyword in Context (KWIC)" button and a "Reset" button.

Figura 7: Tela de busca do concordanciador do COCA Online Corpus

Fonte: O autor

- Qual a sintaxe disso?
- Como busco por etiquetas?
- Como o servidor vai entender o que digitei?
- O que eu queria buscar é o que o servidor me retornou?

Concordanciador

- Ações possíveis:
 - Busca por etiqueta;
 - Busca por palavra exata ou partes desta;
 - Intervalo de fixo ou variável de palavras a ignoradas;
 - Combinação das anteriores.

Treino do etiquetador

1 de 3

Corpus: Floresta Sintática⁷ (ordem de milhão de amostras)

Etiquetador: *Unitex/GramLab*

Problema: Qual o significado das etiquetas de saída?

- Documentação incompleta;
- Dos 5 artigos citados, apenas um era de acesso público e não trazia dados sobre o significado das etiquetas.

Problema: O resultado obtido se compara a quê? O que seria um resultado ruim?

⁷LINGUATECA. Projecto Floresta Sinta(c)tica. 2010.

Treino do etiquetador

2 de 3

Desenvolvido outro etiquetador, para ser o parâmetro de ruim:

Corpus: Floresta Sintática⁸ (ordem de milhão de amostras)

Etiquetador: *YAS-Tagger*

Funcionamento: 5 tabelas associativas de trigramas, bigramas e unigramas para etiqueta;

Resultados inesperados levaram à dúvida: “qual seria o impacto se o corpus fosse uma ordem de grandeza menor?”

⁸LINGUATECA. Projecto Floresta Sinta(c)tica. 2010.

Treino do etiquetador

3 de 3

Corpus: Aires⁹ (ordem de centena de milhar de amostras)

Etiquetador: *YAS-Tagger*

⁹AIRES, R. V. X. **Implementação, adaptação, combinação e avaliação de etiquetadores para o português do Brasil**. 154 p. Tese (Doutorado) — Universidade de São Paulo, São Carlos, 2000.

Treinos dos etiquetadores

- 1 Floresta Sintática + Unitex/GramLab
- 2 Floresta Sintática + YAS-Tagger
- 3 Aires + YAS-Tagger

Floresta Sintática + Unitex/GramLab

- Precisão $\approx 60,76\%$;
- Precisão concentrada em 3 das 4 etiquetas mais frequentes:
 - PREP;
 - PRON;
 - V;
- A frequência das etiquetas da saída do etiquetador não apresenta uma clara correlação com a frequência no treino e teste.

Floresta Sintática + Unitex/GramLab

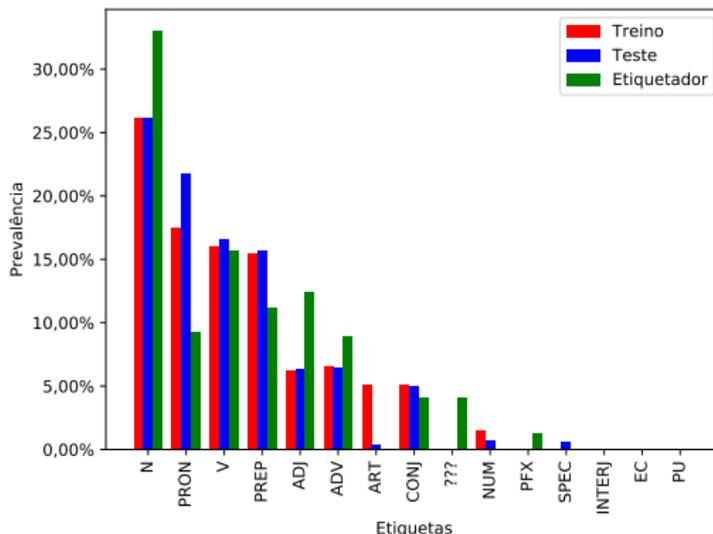


Figura 8: Prevalência das 15 etiquetas mais frequentes do conjunto de treino, teste e etiquetado pelo *Unitex/GramLab*

Fonte: O autor

Floresta Sintática + YAS-Tagger

- Precisão $\approx 76,93\%$;
- Precisão distribuída mais uniformemente por etiqueta;
- A frequência no conjunto de saída do etiquetador é sempre menor que a frequência desta no treino e teste;

Floresta Sintática + YAS-Tagger

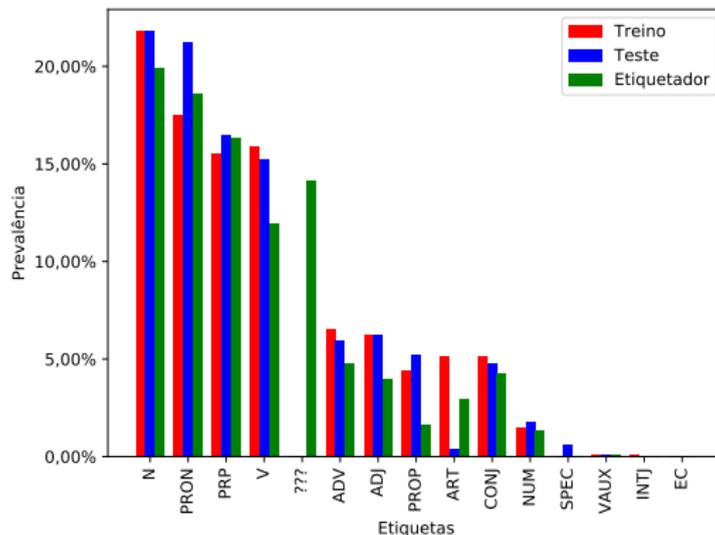


Figura 9: Prevalência das 15 etiquetas mais frequentes do conjunto de treino, teste e etiquetado pelo *YAS-Tagger*

Fonte: O autor

Aires + YAS-Tagger

- Precisão $\approx 53,40\%$ (queda de $23,56\%$);
- A etiqueta “???” agora representa $33,74\%$ da saída do etiquetador:
 - aumento de 138%

Aires + YAS-Tagger

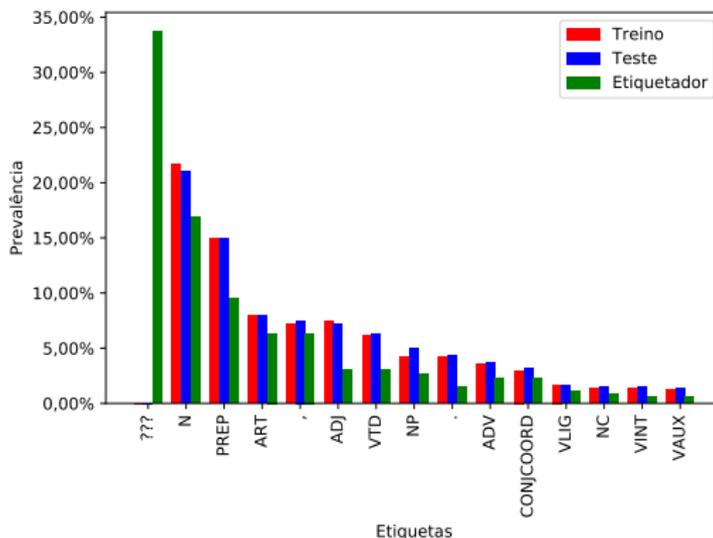


Figura 10: Prevalência das 15 etiquetas mais frequentes do conjunto de treino, teste e etiquetado pelo *YAS-Tagger* sobre o copus de Aires

Fonte: O autor

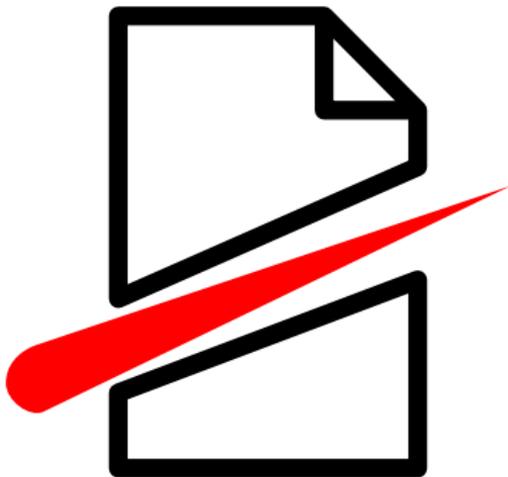
Etiquetador	Precisão
MXPOST	89,66%
Brill Tagger	88,76%
Tree Tagger	88,47%
YAS-Tagger	76,93%
Unitex/GramLab	60,76%

Tabela 2: Comparação da precisão entre os etiquetadores *MXPOST*, *Brill Tagger*, *Tree Tagger*, *YAS-Tagger* e *Unitex/GramLab*

Fonte: Aires¹⁰ e o autor.

¹⁰AIRES, R. V. X. **Implementação, adaptação, combinação e avaliação de etiquetadores para o português do Brasil**. 154 p. Tese (Doutorado) — Universidade de São Paulo, São Carlos, 2000. p. 82.

O sistema desenvolvido



<https://corpusslayer.com>

*Os corpora aparecerão aqui
depois que você entrar na plataforma*

Entrar

Usuário

Senha

[Entrar](#)

[Esqueci minha senha](#)

[Registrar](#)

Figura 11: Página inicial do sistema desenvolvido

Fonte: O autor

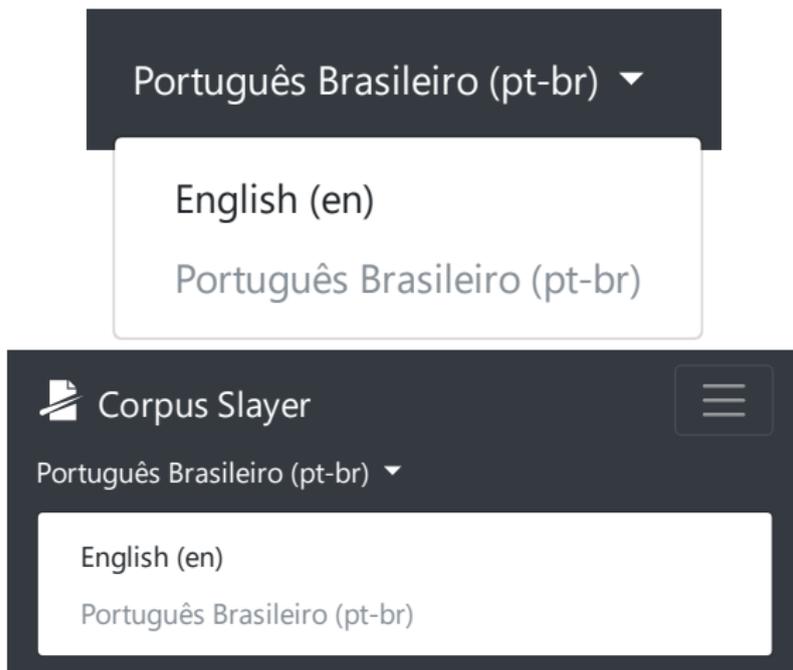


Figura 12: Detalhe do seletor de idiomas, *desktop* em cima e em dispositivos móveis em baixo

The screenshot shows the main interface of the 'Corpus Slayer' application. At the top, there is a navigation bar with the application name 'Corpus Slayer', a language dropdown set to 'Português Brasileiro (pt-br)', and links for 'Configurações' and 'Sair'. The main content area displays three corpora cards: 'Oriente Médio Grande' (48 documents), 'Fireworks' (1 document), and 'Oriente Médio Pequeno' (1 document). Each card includes a 'Detalhes' button, creation and modification dates, and a '+ Adicionar corpus' button at the bottom. To the right, a large grey box contains a 'Bem Vindo!' message and a 'Conectado como adler' status. The footer contains the text 'Site esculpido à mão com amor usando Django' and links for 'Ajuda', 'Privacidade', and 'Termos'.

Figura 13: Página inicial do sistema desenvolvido, mostrando a listagem dos corpora dum usuário

Fonte: O autor

Oriente Médio Grande
48 documentos

Detalhes →

Criado 2 meses, 1 semana atrás

Modificado 2 meses, 1 semana atrás

Fireworks
1 documentos

Detalhes →

Criado 2 meses, 1 semana atrás

Modificado 2 meses, 1 semana atrás

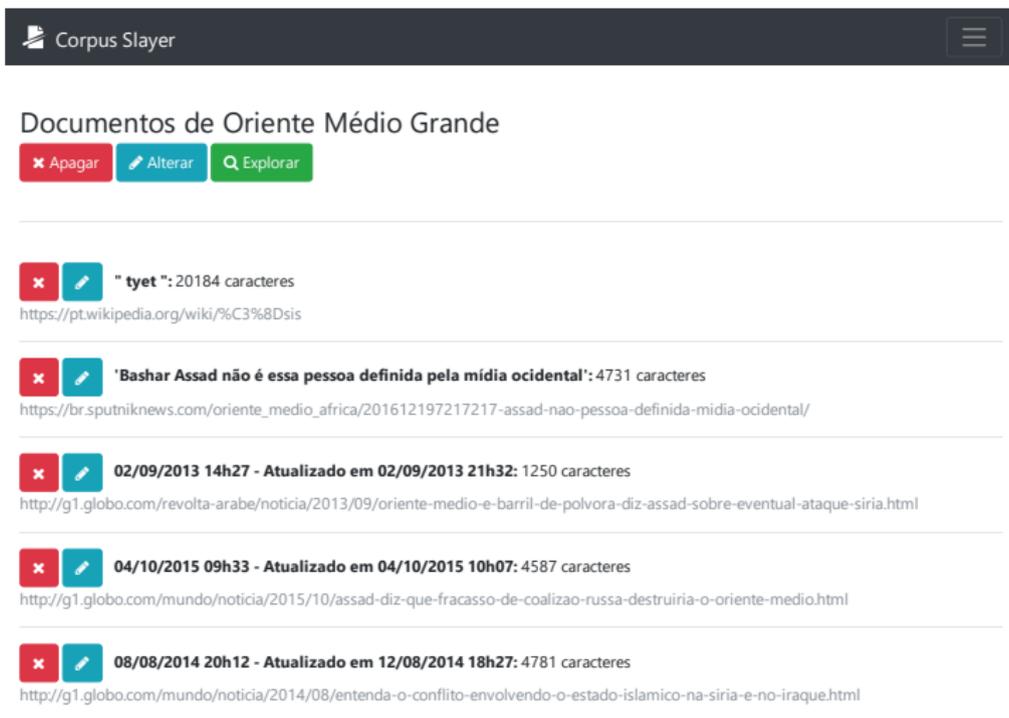
Oriente Médio Pequeno
1 documentos

Detalhes →

Criado 2 meses, 1 semana atrás

Modificado 2 meses, 1 semana atrás

+ Adicionar corpus



The screenshot shows the 'Corpus Slayer' application interface. At the top, there is a dark header with the application name and a menu icon. Below the header, the main content area is titled 'Documentos de Oriente Médio Grande'. Underneath the title, there are three action buttons: 'Apagar' (red), 'Alterar' (blue), and 'Explorar' (green). The main area displays a list of documents, each with a red 'x' icon, a blue edit icon, and a text snippet. The first document is titled '"tyet ": 20184 caracteres' with a URL from pt.wikipedia.org. The second is "'Bashar Assad não é essa pessoa definida pela mídia ocidental': 4731 caracteres' with a URL from sputniknews.com. The third is '02/09/2013 14h27 - Atualizado em 02/09/2013 21h32: 1250 caracteres' with a URL from globo.com. The fourth is '04/10/2015 09h33 - Atualizado em 04/10/2015 10h07: 4587 caracteres' with a URL from globo.com. The fifth is '08/08/2014 20h12 - Atualizado em 12/08/2014 18h27: 4781 caracteres' with a URL from globo.com.

Corpus Slayer

Documentos de Oriente Médio Grande

✖ Alterar Explorar

✖ "tyet ": 20184 caracteres
<https://pt.wikipedia.org/wiki/%C3%8Dsis>

✖ 'Bashar Assad não é essa pessoa definida pela mídia ocidental': 4731 caracteres
https://br.sputniknews.com/oriente_medio_africa/201612197217217-assad-nao-pessoa-definida-midia-ocidental/

✖ 02/09/2013 14h27 - Atualizado em 02/09/2013 21h32: 1250 caracteres
<http://g1.globo.com/revolta-arabe/noticia/2013/09/oriente-medio-e-barril-de-polvora-diz-assad-sobre-eventual-ataque-siria.html>

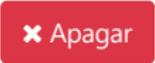
✖ 04/10/2015 09h33 - Atualizado em 04/10/2015 10h07: 4587 caracteres
<http://g1.globo.com/mundo/noticia/2015/10/assad-diz-que-fracasso-de-coalizacao-russa-destruiria-o-oriente-medio.html>

✖ 08/08/2014 20h12 - Atualizado em 12/08/2014 18h27: 4781 caracteres
<http://g1.globo.com/mundo/noticia/2014/08/entenda-o-conflito-envolvendo-o-estado-islamico-na-siria-e-no-iraque.html>

Figura 14: Lista de documentos dum corpus

Fonte: O autor

Documentos de Oriente Médio Grande

 x Apagar Alterar Explorar

" **tyet** ": 20184 caracteres

<https://pt.wikipedia.org/wiki/%C3%8Dsis>



'**Bashar Assad não é essa pessoa definida pela mídia ocidental**': 4731 caracteres

https://br.sputniknews.com/oriente_medio_africa/201612197217217-assad-nao-pessoa-definida-n



Figura 15: Opções de análises para um corpus

Fonte: O autor

Análise - Oriente Médio Grande

+ Enviar corpus construído com o BootCaT [bootcat_upload]

✘ Descartar dados processados [Unitex/GramLab]

✘ Descartar dados processados [MXTERMINATOR]

✘ Descartar dados processados [MXPOST]

✘ Descartar dados processados [TreeTagger]

→ Processar Corpus [Unitex/GramLab]

→ Processar Corpus [MXTERMINATOR]

→ Processar Corpus [MXPOST]

→ Processar Corpus [TreeTagger]

☰ Lista de Sentenças [Unitex/GramLab]

☰ Lista de Sentenças [MXTERMINATOR]

☰ Lista de Sentenças [Unitex/GramLab]

☰ Lista de Sentenças [MXTERMINATOR]

☰ Frequência de Palavras [Unitex/GramLab]

☰ Lista de Palavras [Unitex/GramLab]

☰ Autômato de Sentença [Unitex/GramLab]

☰ Autômato do Texto [Unitex/GramLab]

👉 Texto Etiquetado [MXPOST]

👉 Texto Etiquetado [TreeTagger]

☰ Concordanciador [concord]



Lista de Sentenças - Oriente Médio Pequeno

1.  Frente dos rebeldes sírios em Aleppo desmorona e leva pânico a refugiados
2.  Ofensiva lançada pelas tropas de Bashar al Assad com apoio russo encurrala refugiados
3.  Recomandar no Facebook
4.  Um guia para entender quem é quem no complexo conflito da Síria
5.  As linhas rebeldes no norte da Síria desmoronam rapidamente nos últimos dias, por causa da intensificação da ofensiva contra Aleppo realizada por forças leais ao presidente Bashar al Assad , apoiadas por combatentes xiitas e pela aviação russa.
6.  Na última segunda-feira, as fileiras do regime se situavam nos arredores de Tal Rifat, a apenas 20 quilômetros da passagem fronteiriça de Öncüpınar/Bab al Salam, entre a Turquia e a Síria.

Figura 16: Lista de sentenças dum corpus

Fonte: O autor

 Corpus Slayer

 Lista de Palavras - Oriente Médio Pequeno

Palavras não reconhecidas

Afrin Ahrar Alepo Aleppo Ansari Assad Bab Bashar Corps EI EL Facebook Hassan
Hawa Hazer Idlib jihad jihadistas Kazim Kenyo Kerem Kilis Kinik Marea Mercy Nusra
Nyirjesy Ôncüpinar Qaeda Recomanar Reyhanli Rifat salafista Salam Salih Sham

Palavra Composta	Lema	Gramática & Semântica	Inflexão
 porta-voz	porta-voz	N VN	
 queixaram-se	queixar	V PRO	

Figura 17: Lista de palavras dum corpus

Fonte: O autor



Lista de Sentenças Etiquetadas -

Oriente Médio Pequeno

- 
Frente N dos NPROP rebeldes N sírios N em PREP
Aleppo N desmorona N e KC leva V pânico ADJ a ART
refugiados N
- 
Ofensiva ADJ lançada PCP peLas N tropas N de PREP
Bashar NPROP aL NPROP Assad NPROP com PREP apoio N russo ADJ

Figura 18: Corpus processado pelo *Tree Tagger*

Fonte: O autor

Concordanciador - Oriente Médio Grande

Corpus etiquetado

TreeTagger

Este campo é obrigatório.

Visibilidade da vizinhança

4 palavras de cada lado

Busca

.ado__VERBO {0,1} terrorista {0,*} .mic.

Este campo é obrigatório.

Buscar

Decompondo consulta

termina com (até) VERBO | pula de (até) Palavras | e (terrorista)

pula de (até) Palavras | contém (mic)

Referência da notação de busca

Observe a tabela abaixo:

Busca	Significado
abc	A palavra é "abc"
\\	A palavra é "\"
.\.\.	A palavra é "..."

Figura 19: Tela de busca do concordanciador

Fonte: O autor

Concordanciador - **Oriente Médio Grande**

Corpus etiquetado

TreeTaqqr

Este campo é obrigatório.

Visibilidade da vizinhança

4 palavras de cada lado

Busca

Este campo é obrigatório.

Decompondo consulta

termina com ado

VERBO

pula de 0 até 1 Palavras

é terrorista

pula de 0 até any Palavras

contém mic

Referência da notação de busca

Observe a tabela abaixo:

Busca	Significado
abc	A palavra é "abc"
\\	A palavra é "\"
.\.\\.	A palavra é "..."



Resultados do Concordanciador - Oriente Médio Grande

- apenas PDEN | atua V | conforme PREP | o ART | planejado N | . P | na PREP
 propaganda N | terrorista ADJ | gravada PCP | com PREP | duas NUM | câmeras N
- Estado N | Islâmico NPROP | e NPROP | já NPROP | considerado PCP | o ART
 grupo N | terrorista ADJ | mais KC | perigoso ADJ | de PREP | todos PRDADJ
- do ADV | executivo ADJ | turco N | e KC | considerado PCP | um ART
 grupo N | terrorista ADJ | por PREP | Ancara NPROP | . P | pelos N

Figura 20: Tela de resultados do concordanciador

Fonte: O autor

Resultados do Concordanciador - Oriente Médio Grande

- ↩ apenas PDEN atua V conforme PREP o ART planejado N , , na PREP
 propaganda N terrorista ADJ gravada PCP com PREP duas NUM câmeras N
- ↩ Estado N Islâmico NPROP é NPROP já NPROP considerado PCP o ART
 grupo N terrorista ADJ mais KC perigoso ADJ de PREP todos PROADJ
- ↩ do ADV executivo ADJ turco N e KC considerado PCP um ART
 grupo N terrorista ADJ por PREP Ancara NPROP , , pelos N

2.  Estado N Islâmico NPROP é NPROP já NPROP considerado PCP o ART
- grupo N terrorista ADJ mais KC perigoso ADJ de PREP todos PROADJ



Corpus Slayer



Documento #88 de Oriente Médio Grande

Território reivindicado pelo Estado Islâmico no mundo.

https://pt.m.wikipedia.org/wiki/Estado_Isl%C3%A2mico_do_Iraque_e_do_Levante

Território reivindicado pelo Estado Islâmico no mundo.

Desde 2004, a principal meta do grupo é a fundação de um Estado islâmico . [117] [118] O
EIII procurou estabelecer-se como um califado, um tipo de Estado islâmico liderado por um

Figura 21: Visualização de documento a partir dum clique na seta para trás

Fonte: O autor



Corpus Slayer



Documentos encontrados em **Oriente Médio Grande**

08/08/2014 20h12 - Atualizado em 12/08/2014 18h27

<http://g1.globo.com/mundo/noticia/2014/08/entenda-o-conflito-envolvendo-o-estado-islamico-na-siria-e-no-iraque.html>

14 pontos-chave sobre o Oriente Médio e o papel do Estado Islâmico

http://brasil.elpais.com/brasil/2015/10/11/internacional/1444563614_586697.html

Figura 22: Listagem de documentos a partir dum clique na seta para trás

Fonte: O autor

- É possível conseguir uma ferramenta comparável às pagas apenas integrando softwares gratuitos existentes;
- Este trabalho deixa uma fonte de inspiração para concordanciadores existentes e futuros uma para aumentar a usabilidade destes por usuários inexperientes;
- Unitex/GramLab tem uma documentação incompleta que cobre apenas o uso da interface.

Trabalhos futuros

- Garantir que o sistema desenvolvido seja acessível por cegos;
- Adicionar interoperabilidade do sistema desenvolvido com outros sistemas que usam este como execução remota de procedimento;
- Adicionar elementos de rede social, de forma a ser possível compartilhar resultados entre pesquisadores;
- Implementar todos os requisitos levantados, mas não concretizados neste trabalho;
- Implementar pontuador automático para auxiliar pessoas com dislexia.

Perguntas?