



Criação de um Ambiente para o Processamento de Córpus de Português Histórico

Arnaldo Candido Junior

**Orientadora:
Sandra Maria Aluísio**

**Apoio:
Conselho Nacional de Desenvolvimento Científico e Tecnológico**

São Carlos - 02 de abril - NILC - ICMC - USP

Roteiro

- Introdução sobre córpus
- O projeto DHPB
- Construção do córpus
- Criação de glossários
- Acesso ao córpus
- Redação de verbetes
- Ambiente para processamento de córpus de Português Histórico

Córpus

“Um conjunto de dados lingüísticos (pertencentes ao uso oral ou escrito da língua ou a ambos), sistematizados segundo determinados critérios, suficientemente extensos em amplitude e profundidade, (...) que possam ser processados por computador, com a finalidade de propiciar resultados vários e úteis para a descrição e análise.”

(Sardinha, 2002, p. 18)

Córpus ⁽²⁾

- Diversas tipologias (Atkins et al., 1992; Sinclair, 1996; Giouli e Piperidis., 2002)
- Internacionais:
 - The Bank Of English (524 milhões),
 - FRANTEXT (150 milhões)
- Nacionais:
 - Córpus do NILC (41 milhões)
 - Tycho-Brahe (2 milhões)

Construção de corpus

- **Projeto:** escolha dos textos do corpus
 - Representatividade, balanceamento, reusabilidade, extensibilidade
- **Compilação:** coleta dos textos
 - Obtenção de direitos, coleta, limpeza, conversão de formato

Construção de corpus (2)

- **Anotação:** tratamento de metadados
 - Administrativos, editoriais, analíticos e descritivos (Wynne, 2005)
 - Padrões internacionais (TEI, XCES, ...)
- **Uso:** diferentes enfoques (língua, conteúdo e mídia)
 - Profissionais: lexicógrafos, lingüistas computacionais, historiadores, desenvolvedores de software

Ferramentas de apoio

- **Compilação:** reconhecedores óticos
- **Anotação**
 - **Manual:** editores XML
 - **Automática:** etiquetadores morfossintáticos
 - **Semi-automática:** revisão manual da anotação automática
- **Acesso a córpus:** concordanceadores
- **Extração de conhecimento:** sumarizadores

O Projeto DHPB

- Dicionário Histórico do Português do Brasil (Vale et al., 2008)
- Primeiro projeto do gênero
- No século XVI: PB divergia PP (cultura, fauna, flora)
- 1500-1808 (período pré-imprensa)

O Projeto DHPB (2)

- Textos
 - Cartas dos Jesuítas
 - Documentos dos bandeirantes
 - Relatos dos sertanistas
 - Etc.
- Autores: Brasileiros ou Portugueses (que viveram no Brasil por um longo período)

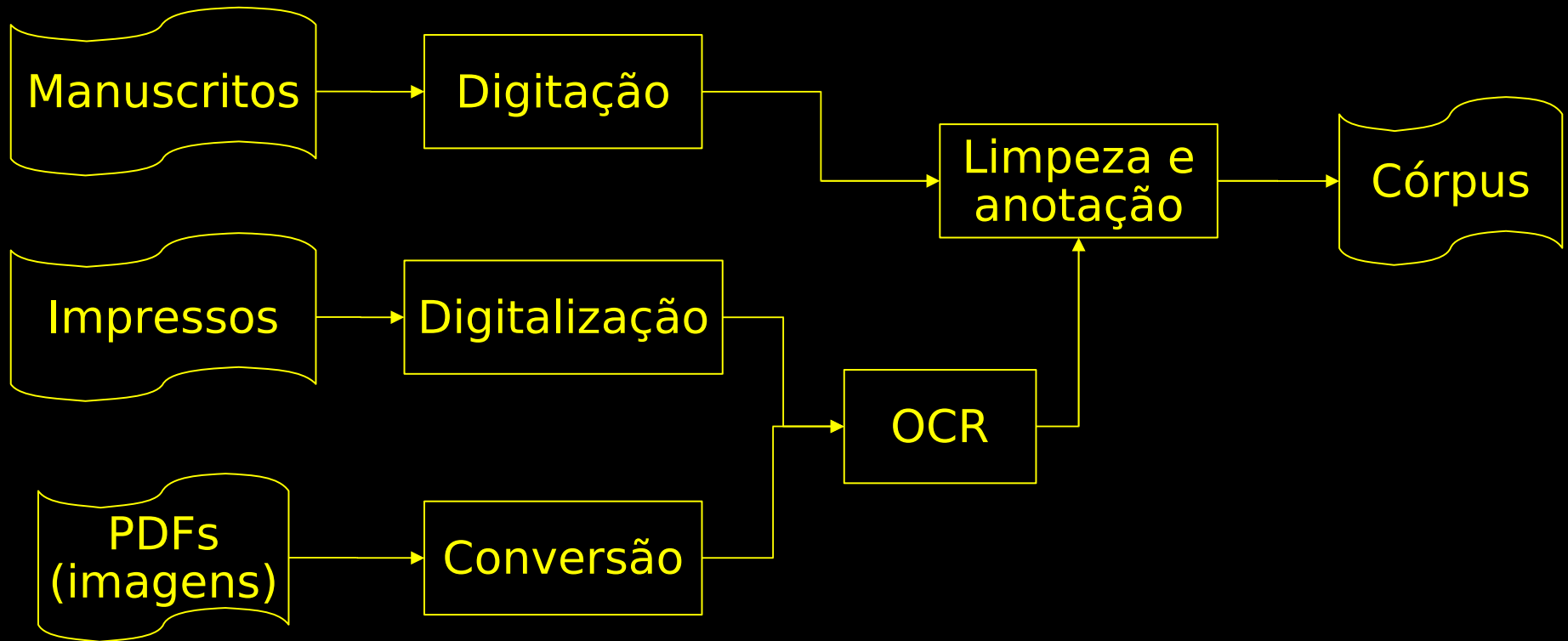
Motivação

- Diversos recursos / metodologias / ferramentas para construção de **córpus contemporâneos**
 - Poucos atendem às necessidades de **córpus históricos**
- Diversas ferramentas para dicionários contemporâneos
 - Poucas atendem às necessidades de **dicionários históricos**

Objetivos

- Levantar as necessidades do projeto DHPB
- Desenvolver metodologias, recursos, e ferramentas para atender estas necessidades
- Criar um ambiente para o processamento de *córpus* de Português Histórico

Compilação do corpus



Digitalização

7. - BAÍA 9 DE AGOSTO DE 1549

127

convertidos, onde estaremos Vicente Rodriguez e eu, e hum soldado¹⁹ que se meteo comnosco para nos servir, e está agora em Exercicios, de que eu estou muy contente. Faremos nossa igreja, onde insinemos os nossos novos christãos, e aos domingos e festas visitarey a Cidade e pregarey. ²⁰⁵

O Padre Antonio Pirez e o P.^e Navarro estaram em outras Aldeas longe, onde já lhes fazem casas. E portanto hé necessario V. R. mandar officiaes, e am-de vir já com a paga, porque cá diz ho Governador que, ainda que venha Alvará de S. A. para nos dar o necessario, que nom o averá ²¹⁰ hi para isto. Os officiaes que cá estão tem muito que fazer, e que o nom tenham, estão com grande saudade do Reyno, porque deixão lá suas molheres e filhos, e nom aceitaram a nossa obra depois que cumprirem com S. A., e tambem ho trabalho que tem com as viandas e o mais os tira disso. ²¹⁵ Portanto me parece que avião de vir de lá, e, se possivel

OCR

7. - BAÍA 9 DE AGOSTO DE 1549 127

convertidos, onde estaremos Vicente Rodriguez e eu, e hum soldado¹⁹ que se meteo comnosco para nos servir, e está agora em Exercicios, de que eu estou muy contente, Faremos nossa igreja, onde insinemos os nossos novos christãos, e aos domingos e festas visitarey a Cidade e pregarey. 205

O Padre Antonio Pirez e o P.^e Navarro estaram em outras Aldeas longe, onde já lhes fazem casas. E portanto hé necessario V. R. mandar officiaes, e am-de vir já com a paga, porque cá diz ho Governador que, ainda que venha Alvará de S. A. para nos dar o necessario, que nom o averá 210 hi para isto. Os officiaes que cá estão tem muito que fazer, e que o nom tenham, estão com grande saudade do Reyno, porque deixão lá suas mulheres e filhos, e nom aceitaram a nossa obra depois que cumprirem com S. A., e tambem ho trabalho que tem com as viandas e o mais os tira disso. 215 Portanto me parece que avião de vir de lá, e, se possivel fosse, com suas mulheres e filhos, e alguns que fação taipas e carpinteiros. Cá está hum Mestre para as obras, que hé hum sobrinho²⁰ de Luis Diaz, mestre das obras d'El-Rey, ho qual veo con trinta mil reis de partido.

(...)

¹⁹ Simão Gonçalves. LEITE I 573.

²⁰ Este «bom oficial», sobrinho de Luís Dias, era Diogo Peres. LEITE I 22.

Limpeza e anotação

<p> {7. - BAÍA 9 DE AGOSTO DE 1549 127 - A00_0002.txt,.N} </p>

<p> convertidos, onde estaremos Vicente Rodriguez e eu, e hum soldado <note place="foot"n="19"> Simão Gonçalves. LEITE I 573. </note> que se meteo comnosco para nos servir, e está agora em Exercicios, de que eu estou muy contente, Faremos nossa igreja, onde insinemos os nossos novos christãos, e aos domingos e festas visitarey a Cidade e pregarey. </p>

<p> O Padre Antonio Pirez e o P.^e Navarro estaram em outras Aldeas longe, onde já lhes fazem casas. E portanto hé necessario V. R. mandar officiaes, e am-de vir já com a paga, porque cá diz ho Governador que, ainda que venha Alvará de S. A. para nos dar o necessario, que nom o averá hi para isto. Os officiaes que cá estão tem muito que fazer, e que o nom tenham, estão com grande saudade do Reyno, porque deixão lá suas mulheres e filhos, e nom aceitaram a nossa obra depois que cumprirem com S. A., e tambem ho trabalho que tem com as viandas e o mais os tira disso. Portanto me parece que avião de vir de lá, e, se possivel fosse, com suas mulheres e filhos, e alguns que fação taipas e carpinteiros. Cá está hum Mestre para as obras, que hé hum sobrinho <note place="foot"n="20"> Este «bom oficial», sobrinho de Luís Dias, era Diogo Peres. LEITE I 22. </note> de Luis Diaz, mestre das obras d'El-Rey, ho qual veo con trinta mil reis de partido. Este nom hé necessario porque abasta ho tio para as obras de S. A.; a este avião de dar o cuidado do nosso collegio; hé bom official. Serão cá muito necessarias pessoas que teção algodão, que há muito, e outros officiaes. </p>

(...)

Ficha catalográfica

Tipologias

1. Tipo da Fonte: EDIÇÃO IMPRESSA

2.1 Domínio Discursivo/Subdomínio Discursivo:

2.2 Gênero/Subgênero:

3a. Tipologia de Assuntos:

3b. Características Sociolingüísticas do Autor:

4. Descrição: CARTAS JESUÍTICAS DISPOSTAS EM ORDEM CRONOLÓGICA, ORGANIZADAS E, QUANDO PRECISO, TRADUZIDAS E ANOTADAS PELO P.^e SERAFIM LEITE (1538-1553) - 3 VOLUMES

5: Localização da Obra: UNESP - CAMPUS DE ARARAQUARA

Fonte

6 Nome do Autor do Texto: P. MANUEL DA NÓBREGA

7: Título do Texto: CARTA DO P. MANUEL DA NÓBREGA AO P. SIMÃO RODRIGUES, BAÍA 9 DE AGOSTO 1549

8. Data em que o Texto foi produzido pelo Autor: BAÍA 9 DE AGOSTO 1549

9. Amostra: INTEGRAL

10. Título da Obra: CARTAS DOS PRIMEIROS JESUÍTAS DO BRASIL

11. Editor: SERAFIM LEITE S. J

12. Organizador/Coordenador (coletânea/livro): SERAFIM LEITE S. J

13. Editora: COMISSÃO DO IV CENTENÁRIO DA CIDADE DE SÃO PAULO

14. Local da Edição: SÃO PAULO

(...)

Anotação do cabeçalho

```
<?xml version="1.0" encoding="UTF-8" standalone="no"?>
<!DOCTYPE TEI.2 SYSTEM "http://docsouth.unc.edu/dtds/teixlite.dtd">
<TEI.2>
<teiHeader>
  <fileDesc>
    <titleStmt>
      <title> CARTA DO P. MANUEL DA NÓBREGA AO P. SIMÃO RODRIGUES, BAÍA 9
DE AGOSTO 1549 </title>
      <author>
        <name> P. MANUEL DA NÓBREGA </name>
        <date> BAÍA 9 DE AGOSTO 1549 </date>
      </author>
    </titleStmt>
    <sourceDesc>
      <biblStruct>
        <monogr>
          <imprint>
            <pubDate> 1956</pubDate>
          </imprint>
        </monogr>
      </biblStruct>
    </sourceDesc>
  </fileDesc>
</teiHeader>
```

Compilação do corpus

- Fases
 - Codificação do corpus
 - Limpeza
 - Anotação
 - Formatos e versões
- Ferramentas: Protej e Protew
- Estatísticas do corpus compilado

Codificação do corpus

Acentos combinados

Símbolo	Unicode	Freqüência	Exemplo
^	0302	2	quar [^] y
~	0303	24.568	com [~] ercio
-	0304	596	caca ⁻ o
¨	0308	48	mu [¨] y
˘	0309	1.804	s [˘] o [˘] mente
'	0313	371	tinha ['] o
´	0301	2	qua [´] e [´] s
ˇ	0306	2	apanh ^ˇ e

Símbolos em geral

Símbolo	Unicode	Freqüência	Exemplo
Æ	00C6	41	Æthyopia
æ	00E6	1.378	græti
œ	0153	116	Cœteris
§	00A7	1.131	§ (parágrafo)
‡	2132	4	‡ixit
f	017F	928	Defcobrio (*)
f	0192	149.909	feito (*)
e	0250	4	passade
&	0026	20.649	&c. (etc.)
@	0040	192	@nrique

* Freqüência sujeita a revisão

Limpeza

Formatação, hifenização, numeração de linhas e de parágrafos.

(...) , E mortes que sem se IHe dar Cauza tem **ex=**
ecutado no Rio da Madeira o Gentio da Nasçaô **Mu=**
ra , impedindo o Comersio dos Moradores naquelle Rio,
E pondo em temor, E consternaçãô as Missoenz **esta=**
blecidas nelle. Ordeno ao Dout.or Ouvidor geral desta (...)

(...) , E mortes que sem se IHe dar Cauza tem **executado** no
Rio da Madeira o Gentio da Nasçaô **Mura**, impedindo o
Comersio dos Moradores naquelle Rio,E pondo em temor, E
consternaçãô as Missoenz **establecidas** nelle. Ordeno ao
Dout.or Ouvidor geral desta (...)

Anotação

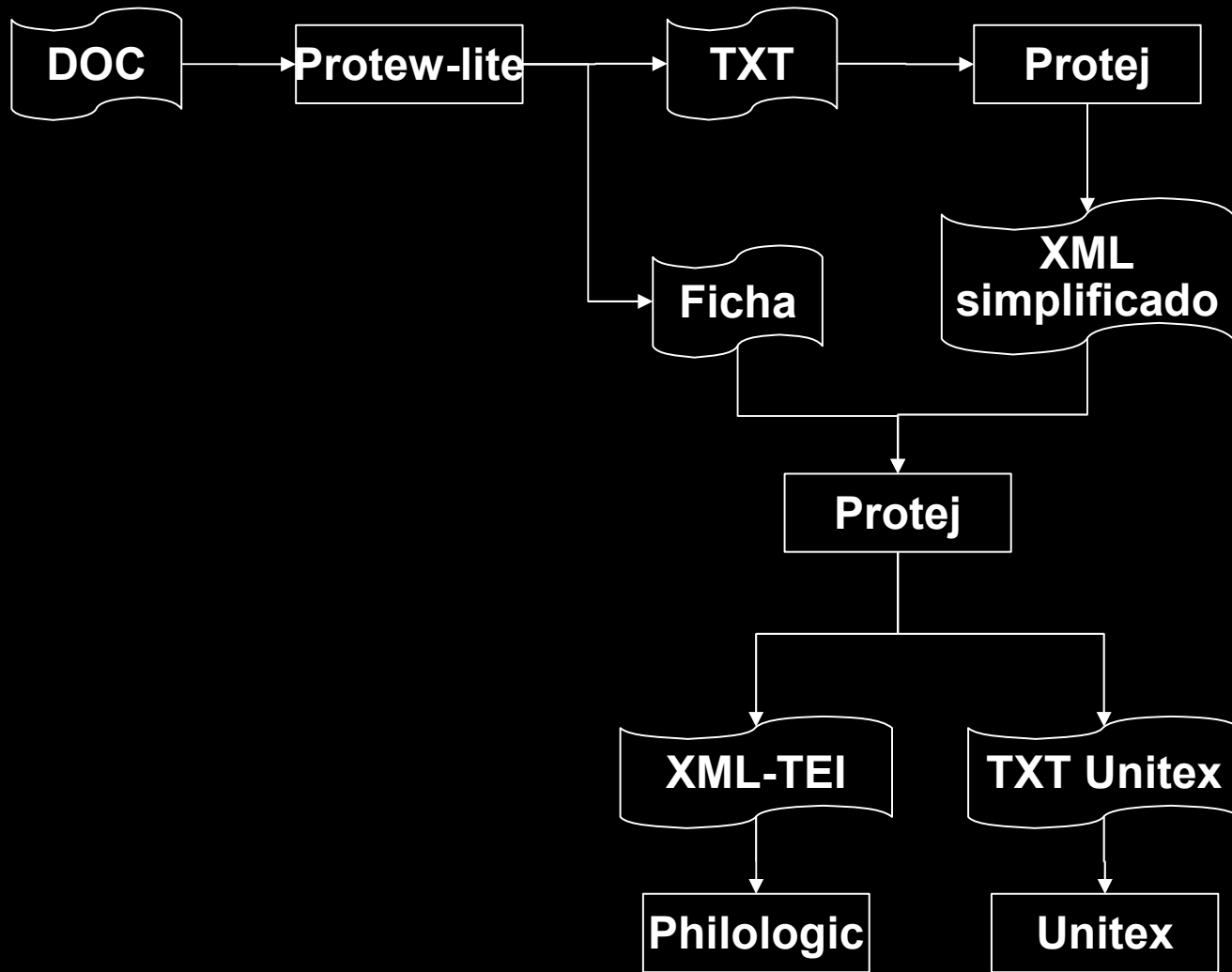
Edição de texto e inserção de etiquetas

(...) da pose vyrem: Como, no anno do nacimiento
5 de Nosso Senhor Jesu Christo de mil e quynhentos
e sessenta 1 anos, haos dose dias do mes d'Aguosto do dyto
(...)
10 campo e borda do matto, Fernão Jorge juiz hordinario
dyta vila e campo, ante my apareceu ho Irmão (...)

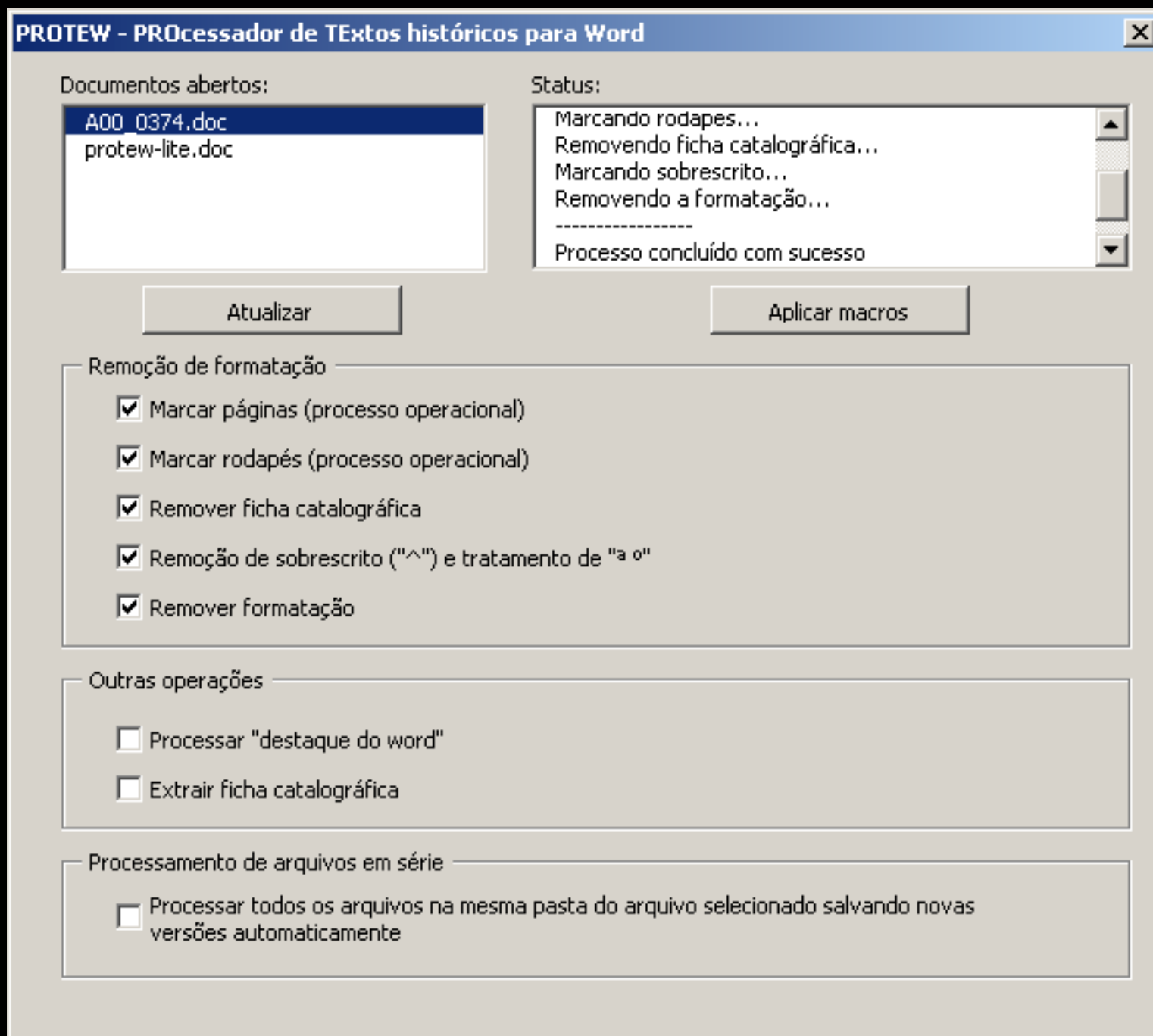
6 Impresso setenta em vez de sessenta. (...)

(...) da pose vyrem: Como, no anno do nacimiento
5 de Nosso Senhor Jesu Christo de mil e quynhentos
e sessenta 1 anos, haos dose dias do mes d'Aguosto do dyto
<note place="foot"n="6" type="line">Impresso setenta (...) </note>
(...)
10 campo e borda do matto, Fernão Jorge juiz hordinario
dyta vila e campo, ante my apareceu ho Irmão (...)

Geração de versões



Protew-lite



Protej

The image shows a screenshot of the Protej software interface. The main window displays a document titled "PROTEJ - /home/arnaldo/projetos/milenio/txts-utf16le/corpus3/A00_0515.txt". The menu bar includes "Arquivo", "Editar", and "Tarefas". The toolbar contains icons for navigation (back, forward, search), settings (gears), and editing (undo, redo, delete). A green heading reads "Tarefa: converter notas que referenciam palavras para XML". The document text includes "CAPITULO 10 O SUCESSO DA MISSÃO DOS PADRES MISSIONARIOS PARA JAGOAGUARA E GURUPATYBA, ONDE ULTIMAMENTE FIZERAM SUA RESIDENCIA...".

Overlaid on the document is a dialog box titled "Protej - Selecionar tarefas". It has "ok" and "cancelar" buttons. The dialog is divided into three tabs: "Pré-processamento inicial", "Tratamento de notas" (which is active), and "Outras operações". Under "Tratamento de notas", there are two sections:

- Etiquetagem de notas**
 - Verificar numeração marginal
 - Converter notas com asteriscos
 - Converter notas alfabéticas
 - Converter notas entre parenteses
 - Converter notas que ocupam múltiplas páginas
 - Converter notas que referenciam palavras para XML
 - Notas no fim da página
 - Notas no fim do documento
 - Converter notas que referenciam linhas para XML
 - Remover numeração marginal
- Processo operacional**
 - Remover marcadores de página e rodapé

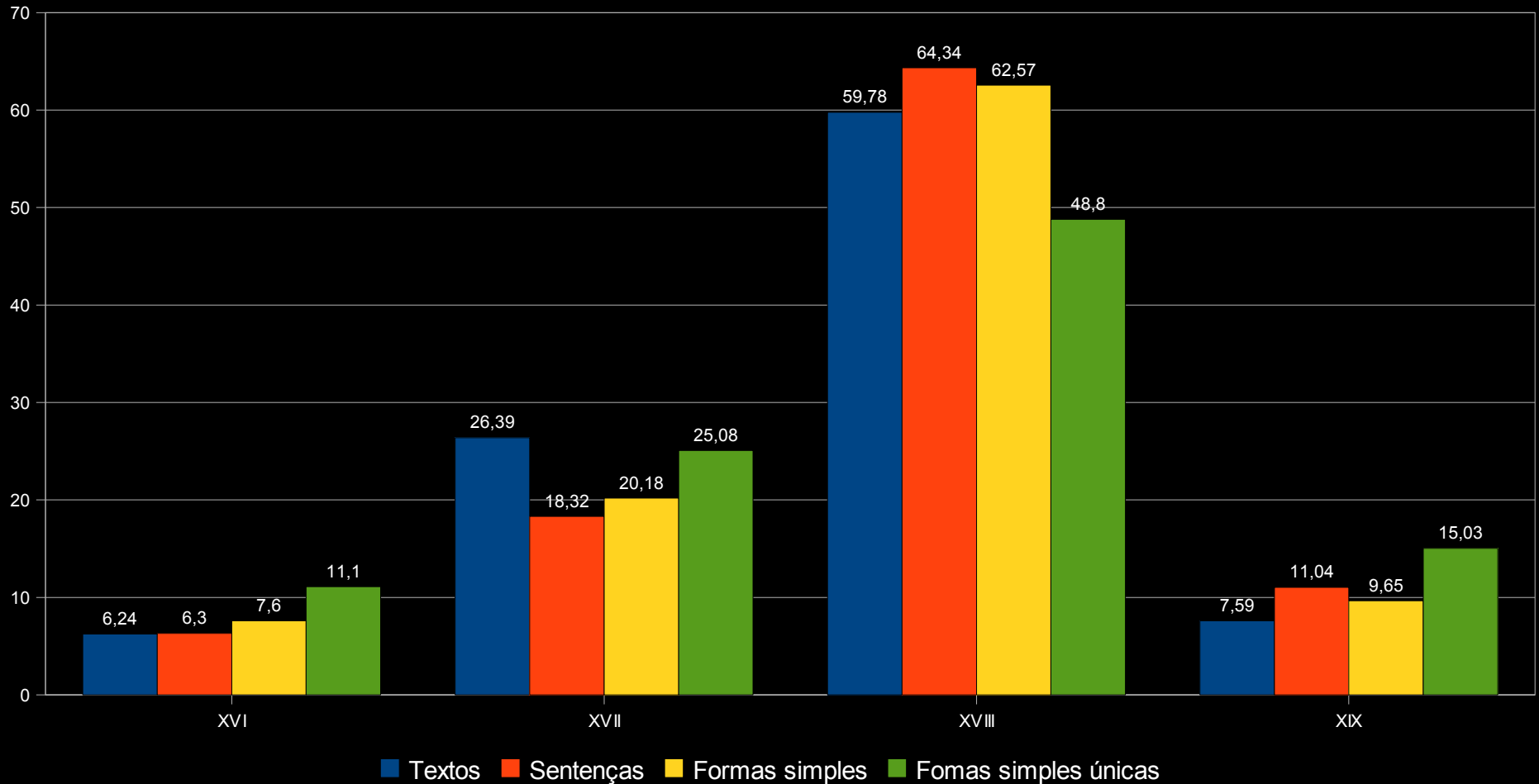
Protej - C rpus Philologic

Inserindo marca o de par grafo...
3 inser es realizadas
Inserindo ficha catalogr fica...
Ficha inserida
Processando A00_0093.txt
Extraindo ficha catalogr fica...
Ficha extraida
Removendo notas em xml
0 notas removidas
Inserindo marca o de par grafo...
4 inser es realizadas
Inserindo ficha catalogr fica...
Ficha inserida
Processando A00_0335.txt
Extraindo ficha catalogr fica...
Ficha extraida
Removendo notas em xml
0 notas removidas
Inserindo marca o de par grafo...
44 inser es realizadas
Inserindo ficha catalogr fica...
Ficha inserida
Processo concluido com sucesso

Dados do corpus

Dados	Valores
<i>Tokens</i>	16.505.808
<i>Types</i>	368.850
Formas simples	7.492.473
Formas simples únicas	368.529
Sentenças	287.570
Textos	2.458
Tamanho em MegaBytes (UTF-16)	82

Textos por século (%)



Glossários

- Fenômenos observados: abreviaturas, junções e variantes de grafia
- Dificultam a contagem de frequências e criação de concordâncias
 - Dificultam a compreensão do texto
- Soluções: substituição, anotação (manual ou automática), criação de glossários

Junções

- Levantadas por uma bolsista do Projeto para separação automática de Junções
- Casos comuns: artigos e preposições
- Exemplos: asmesmas, doestillo, emhumapontadetera, serraniasque, sobpena, ...
- 10.369 junções

Abreviaturas

- Expansões de B^o: bairro, Bartolomeu, bastardo, beco, bento, Bernardo, etc.
- Abreviaturas de janeiro: jan., jan.^{ro}, jan^{ro}, janr.^o, jan.^o, etc.
- Glossário F: extraído de Flexor (1991)
- Glossário C: extraído do corpus
- Protej usado para pré-processamento

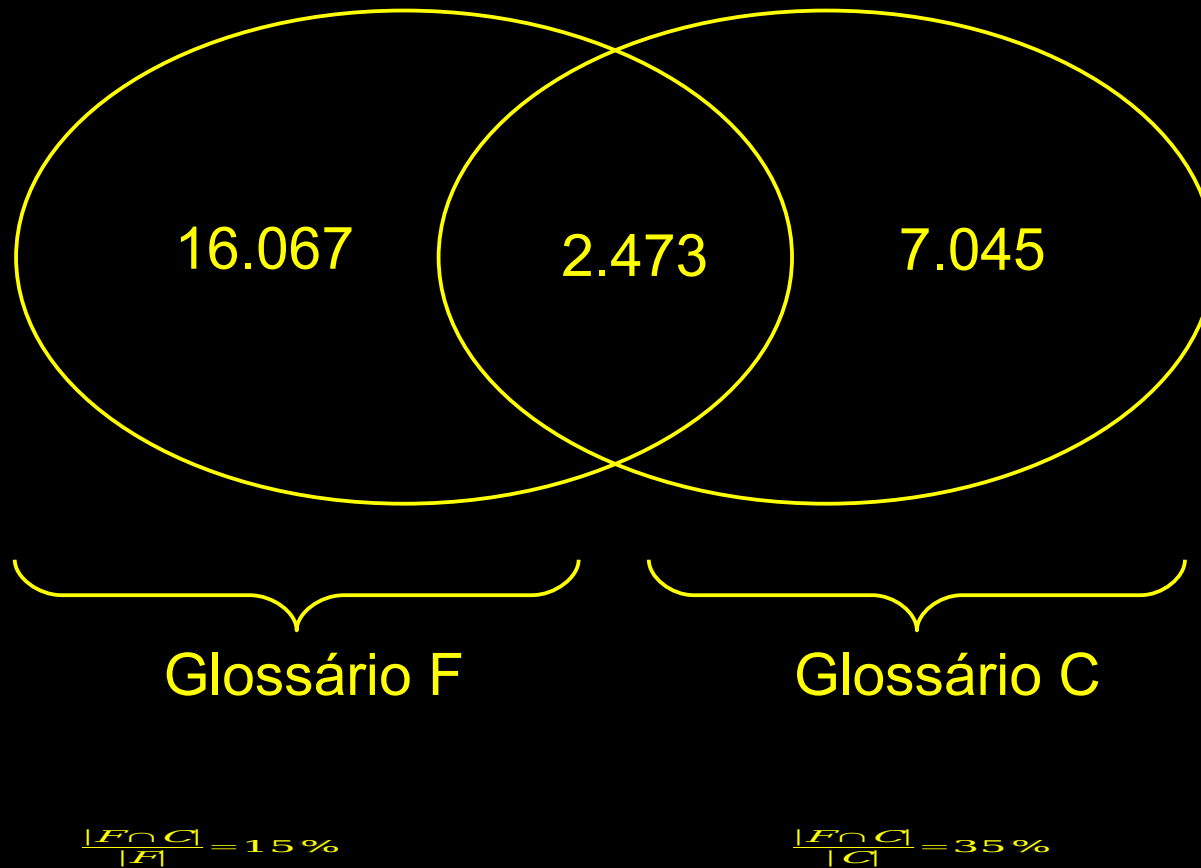
Glossário F

- Convertido para formato digital por uma bolsista do projeto
- Anotação morfossintática (A, B e C)
- Anotação semântica para reconhecimento de entidades nomeadas (A, B, e C)
- 21.869 abreviaturas com 8.721 expansões
- 16.067 abreviaturas simples e 5.635 expansões
- Abreviaturas simples no corpus: 18,92%

Glossário C

- Apenas abreviaturas (sem expansões)
- 7.045 abreviaturas
 - Presença de sobrescrito: ant.^o, cid.^e, p.^a (61%)
 - Ponto interno sucedido por até 4 símbolos: cid.e, embg.e, ex.mo (24%)
 - Palavras terminadas por algumas consoantes: cap, reg, liv, v (15%)

Glossário C versus glossário F



Variantes de grafia

- Inexistência de um sistema ortográfico em textos históricos
- Exemplo: villa, vyla, vjlla, vylla, vjla
- Diversas abordagens: Hirohashi (2004), Rayson et al. (2005), Archer et al. (2005)
- Giusti et al. (2007): regras de transformação (Siaconf) e distância de edição (Philologic + Agrep)

Regras de transformação

- Triplas (E_1 E_2 S), exemplo: (e[ao] e ei)
 - "e[ao]" cobre "aldeia", "meo", "cheas", etc
 - "e" subcadeia a ser substituída (aldeia, meo, cheas, etc)
 - "ei" cadeia de substituição (resultado em aldeia, meio, cheias, etc)
- 43 regras de transformação, 18.082 agrupamentos e 41.170 variantes

Philologic x Siaconf

Técnica	Verdadeiros Positivos	Falsos Positivos	Precisão	Cobertura Comparativa
Regras de transformação (siacnf)	36	0	100%	72%
Distância de edição (Philologic)	41	196	20.92%	84%

Fenômenos combinados

- Não foram considerados

Exemplo	Abreviatura	Junção	Variante
Sarg.^José (Sargento José)	X	X	
abafe (a base)		X	X
supp^te (supostamente)	X		X
héalagadacomm^tos (é alagada com muitos)	X	X	X

Acesso a corpus

- Processadores analisados
 - Philologic: TEI + Interface Web + variantes
 - Unitex: glossários DELA
 - GATE
 - Tenka Text
 - Xaira
- Critério: ISO 9126 (UNIVERSITÉ DE GENÈVE, 2006)

Philologic

The screenshot shows a Mozilla Firefox browser window titled "Philologic Results - Mozilla Firefox". The address bar contains "http://moodle.ici" and the search engine is set to Google. The browser's menu bar includes "Arquivo", "Editar", "Exibir", "Histórico", "Favoritos", "Ferramentas", and "Ajuda". The page content features a header with the word "Philologic" in a stylized font, overlaid on a background of Latin text. Below the header, there is a navigation bar with links: "home", "the ARTFL project", "download", "documentation", and "sample databases". The main content area displays the message "Found 5 matches, shown with frequencies in entire database." and a search interface with "SEARCH" and "CLEAR" buttons. The search results are as follows:

9	<input type="checkbox"/>	giboia
3	<input type="checkbox"/>	giboias
1	<input type="checkbox"/>	giboja
6	<input type="checkbox"/>	giboya
4	<input type="checkbox"/>	gyboia

At the bottom of the browser window, the status bar shows "Concluído".

Unitex

Unitex 2.0beta (October 4, 2007) - current language is Histc

Text DELA FSGraph Lexicon-Grammar Edit File Edition Windows Info

Concordance: /home/arnaldo/unitex/Historical Portuguese (Brazil)/C...

Ihe resulte, offerecendo seu sangue, & [cabedal](#), forças, & industria por instrumentos da conser
ortante Frôta, que em numero de náos, & [cabedal](#) de fazendas, enriqueceo este Reyno. A cargo do
cerca, & fe tornou a pedir mais gente, & [cabedal](#) pera paffar ao Maranhão, enuiando entre tato a d
am um flamengo, homem do maior crédito, [cabedal](#) e sciência de esquipagem, de quantos têm estas
tem tão fervoroso esperito, [ilegível] [cabedal](#) de virtudes que se possa dar por seguro, certo
a desp.ªza. E oq.ªal sendo hum só tenha [cabedal](#) p.ªa ella. E ajustarão a quantid.ªe de Indios, q
cia o captiveiro : como pobre não tinha [cabedal](#) para o resgate: e como a Justiça Divina tinha c
rios a practicar os salvages; Não terá [cabedal](#) para os gastos, e finalmente apenas poderá vive
o nosso, que quase o iguala na grandeza [cabedal](#) de ágoas, ainda agora corre murmurando da sente
usca por neççar de mais industria, e [cabedal](#), mas aseguraõ aver delle, e de Prata muitas Min
ustral, que por serem de menos conta, e [cabedal](#), são pouco nomeados, segue-se o Rio Anajá, que
. 6 v.). Depois de consumido o tempo, e [cabedal](#), que consta dos mandados do caderno do ról de p
ente Rio Tocantins, pelo grande peso, e [cabedal](#) de ágoas, que traz, como já dissemos; e segundo
treter para não lhe ir tomar o passo, e [cabedal](#) de tanta gente rica sem nenhũa defensa, pelo t
pedir mais gente, {A00_2079 - 186,.N} e [cabedal](#) pera a conquista, que o Governador dilatou athé
podem fazer isso, pela grande fábrica e [cabedal](#) que é necessário para se poderem granjeiar canav
os que faziam todos, com a substância e [cabedal](#) dos índios, porque dêles era tudo o que se gast
sca por nessessitar de mais industria e [cabedal](#) mas aseguraõ aver delle e de Prata m.ªtas minas
e ueuia Com toda a sua famillia Caza e [cabedal](#) e fora fugido e retirado pera o Sertam roubado
morou e que ueuia Contada a ssua Caza e [cabedal](#) e foi fogisido e retirado pera o seu Curral e r
bem no mar, sepultando infinita gente e [cabedal](#). Por diligencias do Presidente do Senado, o Dr.
gem de sul. O rio Baures, de extensão e [cabedal](#) d'aguas igual ao Guaporé, de que é o maior conf
gem do sul. O rio Baures, de extensão e [cabedal](#) de aguas igual ao Guaporé, de que é o maior con
costa de Hespanha, e tê seu comercio e [cabedal](#) em conquistas ultramarinas. A 8.ªa Que pode sah

Suporte a edição de verbetes

- Procorph: edição e formatação de verbetes
 - Variantes de grafia
 - Datação
 - Referências ao corpus
 - Acepções (definições)
 - Abonações (excertos do corpus)

Procorph

- Dados gerais
- Variantes

Alterar verbete

Verbetes:
Classe e atributo:
Situação:
Redator:
Data de criação: 2008-01-24

Variantes de grafia

Variante: <input type="text" value="prejuizo"/>	Acima / Abaixo / Remover
Variante: <input type="text" value="preiuzo"/>	Acima / Abaixo / Remover
Variante: <input type="text" value="preioizo"/>	Acima / Abaixo / Remover
Variante: <input type="text" value="preijuzo"/>	Acima / Abaixo / Remover
Variante: <input type="text" value="preyuzo"/>	Acima / Abaixo / Remover
Variante: <input type="text" value="preyoizo"/>	Acima / Abaixo / Remover
Variante: <input type="text" value="prejoizo"/>	Acima / Abaixo / Remover
Variante: <input type="text" value="prejuiso"/>	Acima / Abaixo / Remover
Variante: <input type="text" value="perjuizo"/>	Acima / Abaixo / Remover
Variante: <input type="text" value="prejuifo"/>	Acima / Abaixo / Remover

Procorph

- Abonações
- Acepções
- Referências

Acepções

Acepção:	Dano ou perda de qualquer natureza, como a honra, a saúde. Detrimento.
Atributos:	
Abonação:	Chamei o Padre Antão Gonçalves e examinando bem o caso, achei que tudo fôra uma equivocação, que causara uma má intelligencia das palavras no animo do capitão-mór. Comtudo, como já se não haviam de dar bem ambos, e podia de lá seguir algum *prejuizo* á missão, tirei o padre de Tapeacorú, pondo em seu lugar o Padre Manoel Rabello e o Padre João [p.531] de Avellar, o qual por. mais velho, na religião, ficou por missionario principal daquella residencia, sendo por suas muitas virtudes ambos
Texto:	A00_0548 padre. joão felippe betendorf [1699]. <i>capitulo 15 - acudo á aldea dos guajajaras no marcú e remedeio a um desgosto do capitão-mor do tapeacorú, no maranhão</i>
Página:	530 Acima / Abaixo / Remover
Acepção:	Dano financeiro, perda de dinheiro.
Atributos:	
	Em 21 do passado me recolhi a esta sua caza de volta da minha jornada, e em aliqua couza aproveittou respeito a cobrar aliqua couza do que se deve a esta caza, e mais o havia de ser se o não tivesse empedido a falta de sulimão naquella caza de moeda, que ja passa de dous mezes, q. não

Procorph

- Relacionados
- Observações
- Primeira datação

Verbetes relacionados

[Adicionar verbete relacionado](#)

Observações

Observação:

A forma 'prejuizo' é a mais frequente.

[Acima](#) / [Abaixo](#) / [Remover](#)

[Adicionar observação](#)

Primeira datação

Primeira datação:

E todos tem bem o que hão mister, e a necesidade lhes não fará *prejuizo* algum. Estão espantados de ver a magestade com que entramos e estamos, e temem-nos muito, o que tambem ajuda.]

Texto:

A00_0694 padre manuel da nobrega . [1549]. *carta que o padre manonel da nobrega, preposito provincial da companhia de jesus, em o brasil, escreveu ao padre mestre simão o anno de 1549. (ms. copiado da livraria publica)*

Página:

460

Procorph

- **Visualização (açúcar)**

açúcar: substantivo feminino.

Variantes: assúcar, asúcar

1. Substância doce fabricada industrialmente, extraída, em geral, da cana de açúcar.

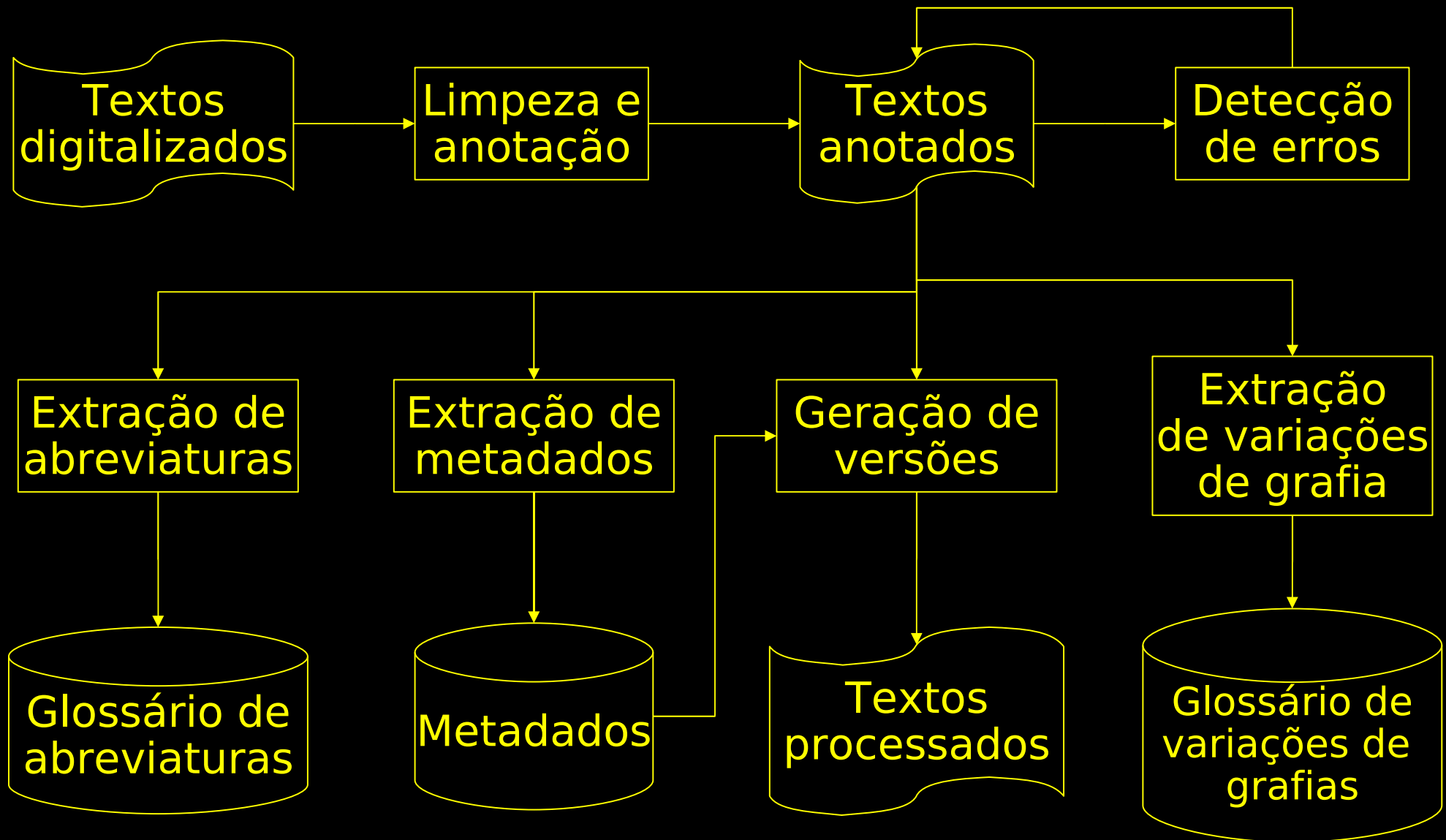
Porque quem vos mostrara, há hoje trezentos anos, uma cana de que se faz o **açúcar** e vos dissera que daquela cana se havia de formar com a indústria humana um pão de **açúcar** tão formoso como hoje o vemos, tê-lo-íeis por coisa ridícula; e, pelo conseguinte, se vos fôsse mostrado um pedaço de pano velho de linho e vos afirmassem que daquele pano se havia de fazer o papel em que escrevemos, quem duvida que o teríeis por zombaria? **ambrósio fernandes brandão [1618]. diálogo primeiro, p. 1 .**

Primeira datação: A renda, que El-Rei cá tem nesta Baía, é esta, scilicet: as miunças que rendem cento e vinte mil réis em que andam arrendadas; o peixe e mandioca e algodão andam em cento e trinta mil réis; pagos em ordenado, que é um terço menos, pode valer em dinheiro oitenta mil réis; o **açúcar** do Engenho anda em cento e cinquenta cruzados. **p. manuel da nóbrega [1558]. carta do p. manuel da nóbrega ao p. miguel de torres, baía 8 de maio 1558, p. 26.**

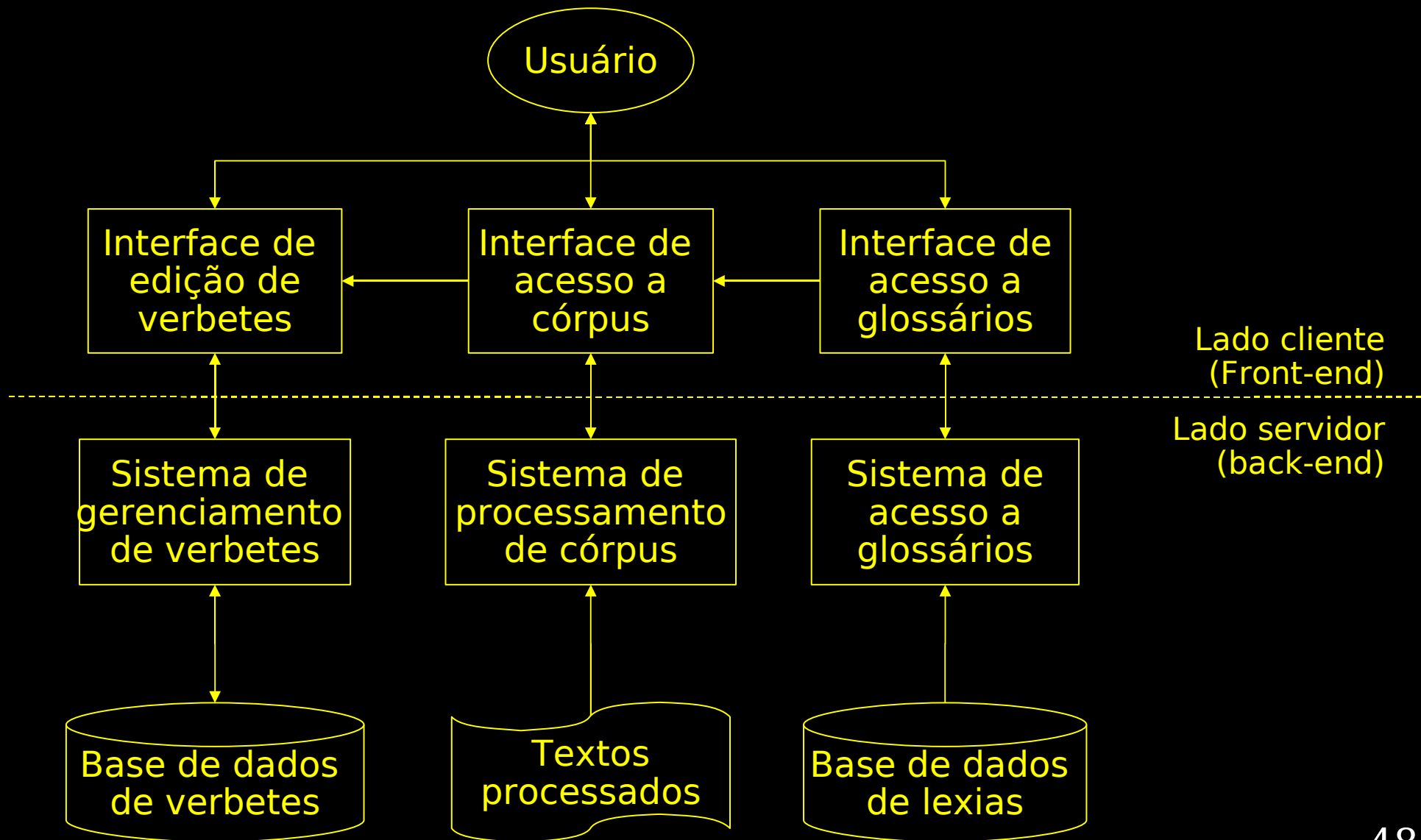
Ambiente para processamento de **córpus** de Português Histórico

- **Arquitetura 1:**
 - **Compilação do **córpus****
 - **Construção de glossários**
 - **Geração de versões**
- **Arquitetura 2:**
 - **Acesso a **córpus****
 - **Acesso a glossários**
 - **Edição de verbetes**

Arquitetura 1



Arquitetura 2



Trabalhos futuros

- Separação automática de junções
- Extração dos metadados domínio e gênero
- Detecção de erros de OCR
- Novas etiquetas TEI
- Novas regras de transformação
- Acesso a glossários no Procorph
- Acesso a corpúscos no Procorph

Conclusões

- Este trabalho faz parte do projeto DHPB
- Contribuições:
 - Metodologia
 - Ferramentas (Protew, Protej, Procorph)
 - Recursos (córpus e glossários)
- As contribuições podem ser usadas em outros projetos

Referências citadas na apresentação

ATKINS, S.; CLEAR, J.; OSTLER, N. Corpus design criteria. *Journal of Literary and Linguistic Computing*, v. 7, n. 1, 1992.

ARCHER, D., ERNST-GERLACH A., KEMPKEN S., PILZ T., RAYSON P. The identification of spelling variants in English and German historical texts: manual or automatic. In: *Digital Humanities*, 2006, Paris: Sorbonne, 2006. p. 3-5.

FLEXOR, M. H. O. *Abreviaturas: Manuscritos dos séculos xvi ao xix*. 2. ed. [S.l.]: UNESP, 1991. 468 p.

GIOULI, V.; PIPERIDIS, S.. *Corpora and HLT: Current trends in corpus processing and annotation*. Disponível em: <http://www.larflast.bas.bg/balric/eng_files/corpus_deliverable_final.htm>. Acesso em: 25 fev. 2008.

GIUSTI, R.; CANDIDO JR, A.; MUNIZ, M. C. M.; CUCATTO, L. A.; ALUÍSIO, S. M. Automatic detection of spelling variation in historical corpus: An application to build a Brazilian Portuguese spelling variants dictionary. In: *Corpus Linguistics*, 2007, Londres. Corpus Linguistics, 2007.

HIROHASHI, A. S. *Aprendizado de regras de substituição para normatização de textos históricos*. Dissertação (Mestrado) — Instituto de Matemática e Estatística, USP, São Paulo, 2004.

Referências citadas na apresentação

RAYSON, P., D. ARCHER AND N. SMITH. VARD versus Word: A comparison of the UCREL variant detector and modern spell checkers on English historic corpora, In Proceedings of Corpus Linguistics 2005, vol. 1, no. 1. Birmingham: Birmingham University.

SARDINHA, T. B. Lingüística de Corpus. Barueri, SP: Manole, 2004.

SINCLAIR, J. Preliminary recommendations on Corpus Typology. EAGLES, 1996. Disponível em: <<http://www.ilc.cnr.it/EAGLES96/pub/eagles/corpora/corpus typ.ps.gz>>. Acesso em: 16 fev. 2007.

UNIVERSITÉ DE GENÈVE. The ISO 9126 Standard. 2006. Disponível <<http://www.issco.unige.ch/ewg95/node1.html>>. Acesso em: 14 nov. 2006.

VALE, O. A. ; CANDIDO JUNIOR, A. ; Muniz ; BENGTON, C. G. ; Cucatto ; ALMEIDA, G. M. B. ; BIDERMAN, M. T. ; Aluísio . Building a large dictionary of abbreviations for named entity recognition in Portuguese historical corpora. In: American Association for Corpus Linguistics, 2008, ACL 2008.

WYNNE, M. (Ed.). Developing Linguistic Corpora: a guide to good practice. Oxford: Oxbow Books, 2005. Disponível em: <<http://ahds.ac.uk/linguistic-corpora/>>. Acesso em: 23 fev. 2007.