

# ***Dicionário Histórico do Português do Brasil (séculos XVI, XVII e XVIII)***

**Ferramentas computacionais para possibilitar o uso universal da Base Informatizada sobre o Brasil dos séculos XVI, XVII e XVIII**

Sandra Maria Aluísio



# Roteiro

- Como fizemos o Portal do Projeto Lácio-Web (LW)
  - Definição dos cabeçalho dos textos
    - Tipologia Textual quadripartida: gênero, tipos de texto, domínio e meio de distribuição
  - Editor de cabeçalhos e a marcação em XML
  - Base de dados que dá suporte às buscas para montagem de subcórpus
- O que aprendemos: o que faríamos diferente
- Proposta para o Projeto do Dic\_Hist
  - Proposta para o Portal e mão-de-obra necessária
  - Proposta para o gerenciamento do projeto: uso de uma Web colaborativa – coteia do ICMC



# Projeto Lácio-Web

- 30 meses de projeto
- 4 tipos diferentes de corpus: Lácio-Ref, Par-C, Comp-C, Mac-Morpho (embora 6 fossem previstos)
- Obtida a Autorização de Uso – via assinatura de um termo – para todos os textos dos corpus
- 16 bolsistas entre DTI e ITI
- 3 pesquisadores – 2 computação, 1 lingüista de corpus

# Lácio-Web



Compilação de Corpus do Português do Brasil e Implementação de Ferramentas para Análises Lingüísticas

:: Conteúdo ::

:: [Descrição](#)

:: [Corpus](#)

:: [Ferramentas](#)

:: [Lançamentos](#)

:: [Manuais](#)

:: [Downloads](#)

:: [Publicações](#)

:: [Por que se cadastrar?](#)

:: [Colaboradores](#)

:: [Como contribuir?](#)

:: [FAQ](#)

:: [Equipe](#)

:: [Apoio](#)

:: [Contato](#)

:: Página Principal ::

*Última Flor do Lácio*

Olavo Bilac, 1914

*Última flor do Lácio, inculta e bela,  
É, a um tempo, esplendor e sepultura:  
Ouro nativo, que na ganga impura  
A bruta mina entre os cascalhos vela...*

*Amo-te assim, desconhecida e obscura,  
Tubo de aço clanger, luto sinfonia,  
Que és o meu - o meu - o meu - o meu - o meu - o meu -  
E o arrolho da saudade e da ternura!*

*Amo o teu viço agreste e o teu aroma  
De virgens selvas e de oceano largo!  
Amo-te, ó rude e doloroso idioma,*

*em que da voz materna ouvi: "meu filho!",  
E em que Camões chorou, no exílio amargo,  
O gênio sem ventura e o amor sem brilho!*



IME - Instituto de  
Matemática e Estatística



*Bem-vindo ao Lácio-Web!*

**[CADASTRE-SE E ACESSE OS CORPUS!](#)**

Usuário:

Senha:

Entrar

[Esqueci minha senha!](#)

Melhor visualizado na resolução: 1024 x 768 Fonte: Média | Última Atualização : 28/06/2004



Iniciar



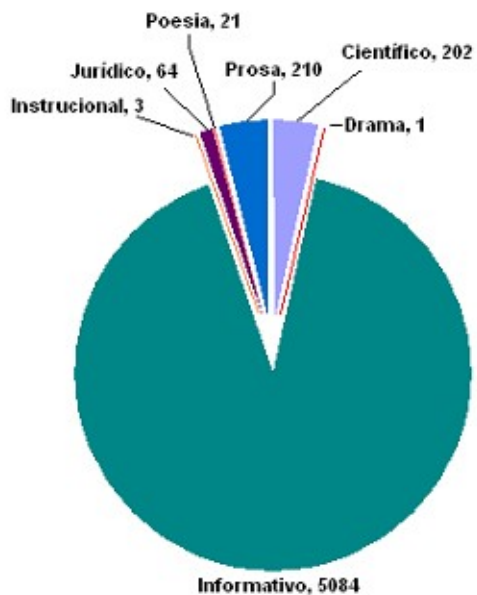
Lácio-Web - Microsoft...

Internet

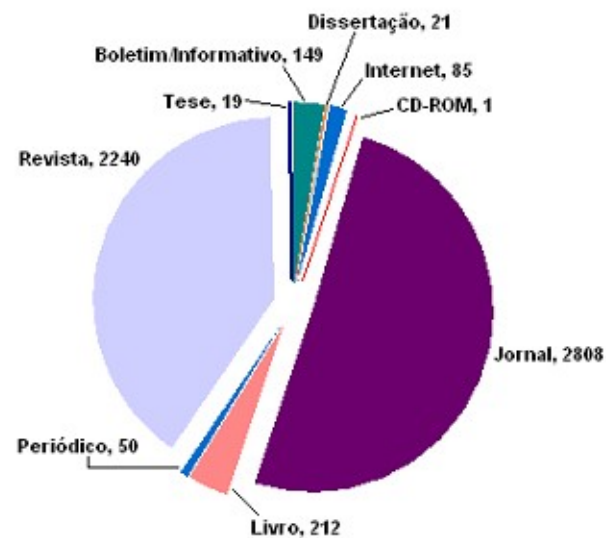


18:22

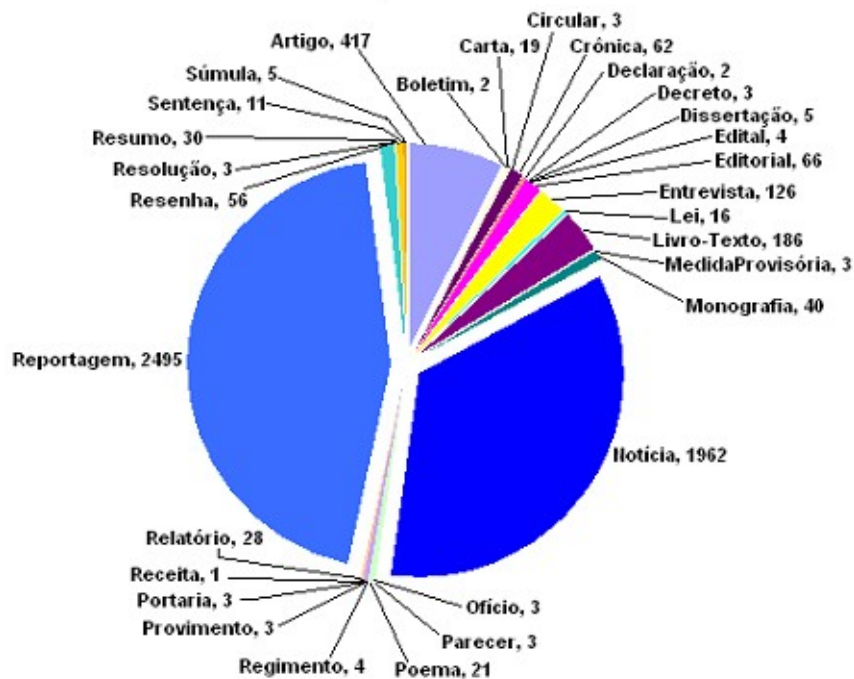
### CLASSIFICAÇÃO DE GÊNEROS



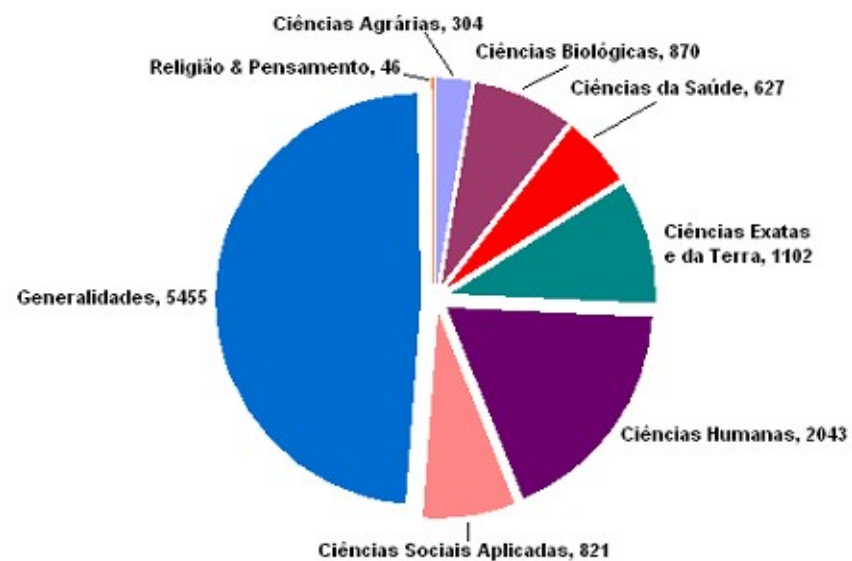
### CLASSIFICAÇÃO DE MEIOS DE DISTRIBUIÇÃO



### CLASSIFICAÇÃO DE TIPOS DE TEXTOS



### CLASSIFICAÇÃO DE DOMÍNIOS

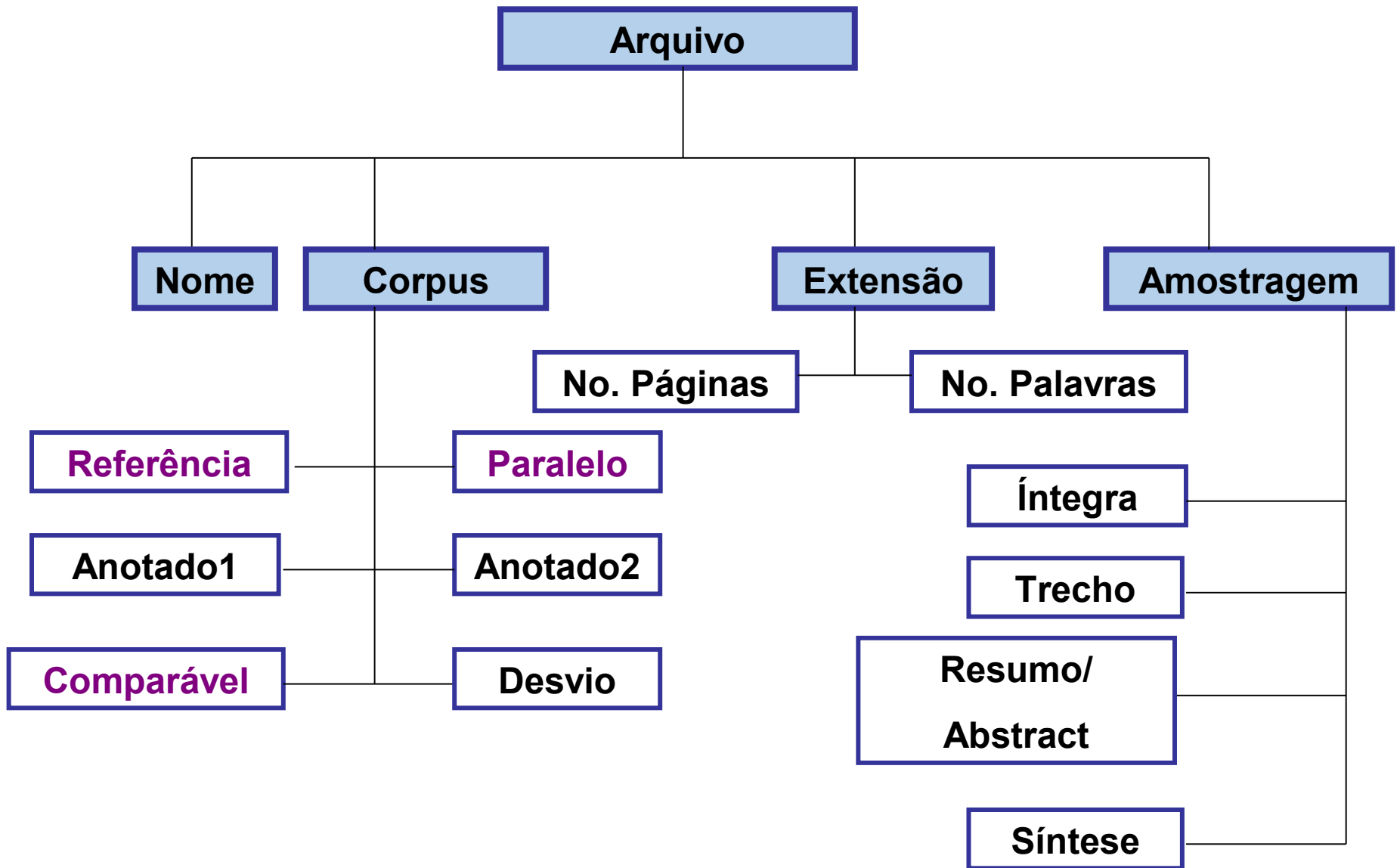


# Como fizemos o Portal do LW

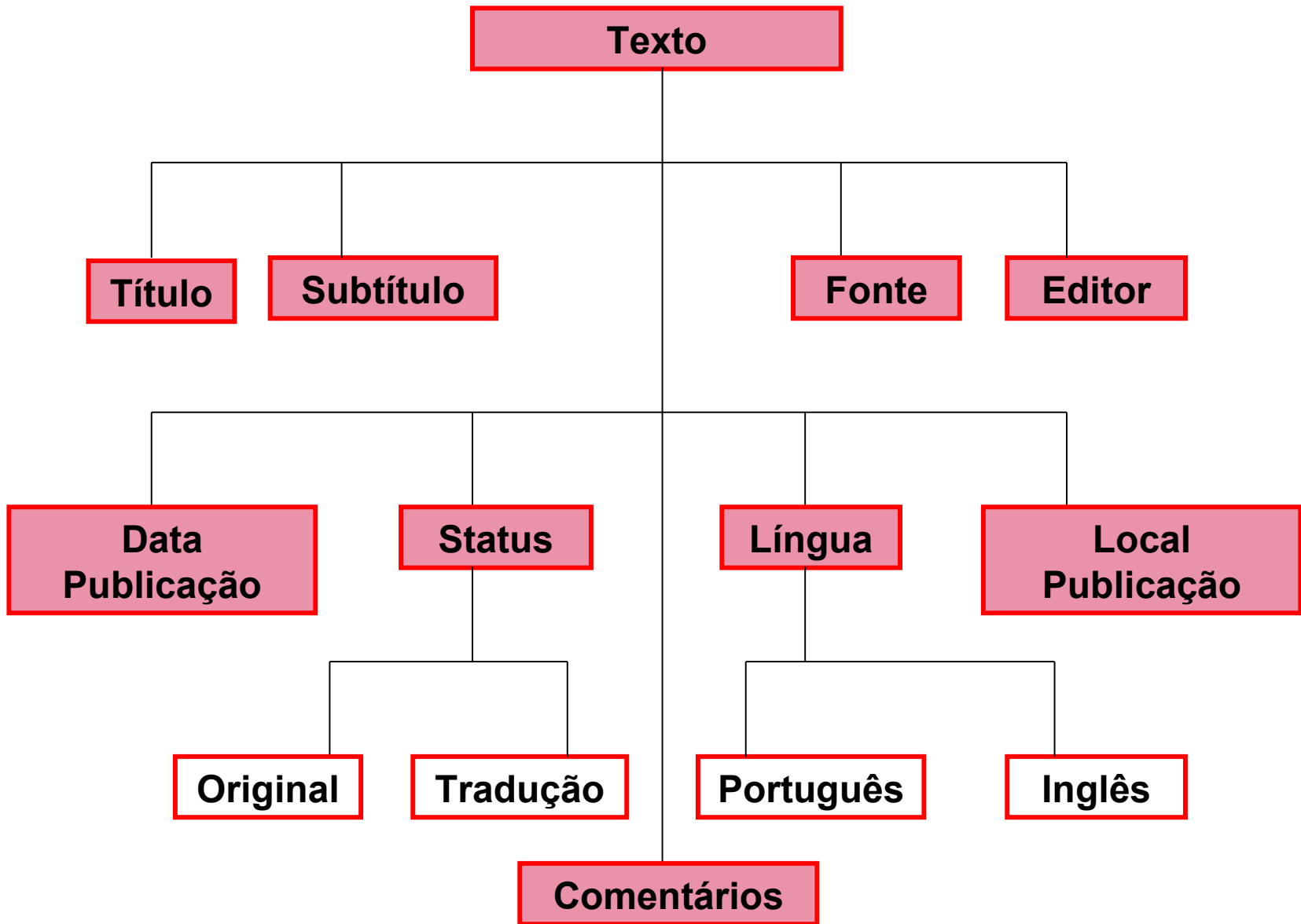
- Decisão sobre as informações do cabeçalho
  - dados bibliográficos comuns,
  - dados de catalogação como tamanho do arquivo, tipo da autoria, resumo do texto (se houver), e a **tipologia textual**
- Nomeação e formatação dos arquivos
- Edição dos cabeçalhos com ajuda de um EDITOR
  - Obtém dados dos nomes (preenche automaticamente alguns campos), ajuda na edição e gera um texto em XML
- Subida dos **dados** acima na base de dados que dá suporte às buscas

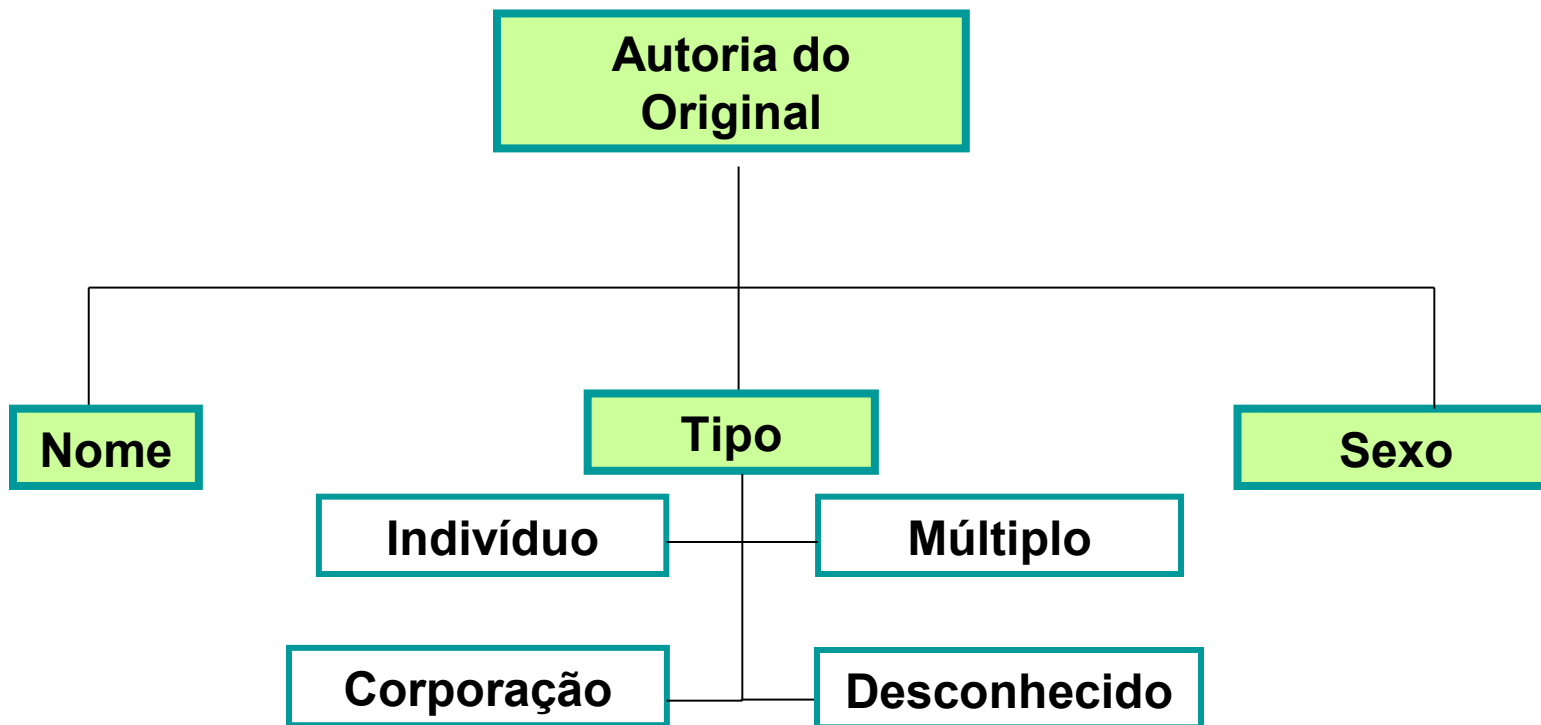


# Informações do cabeçalho

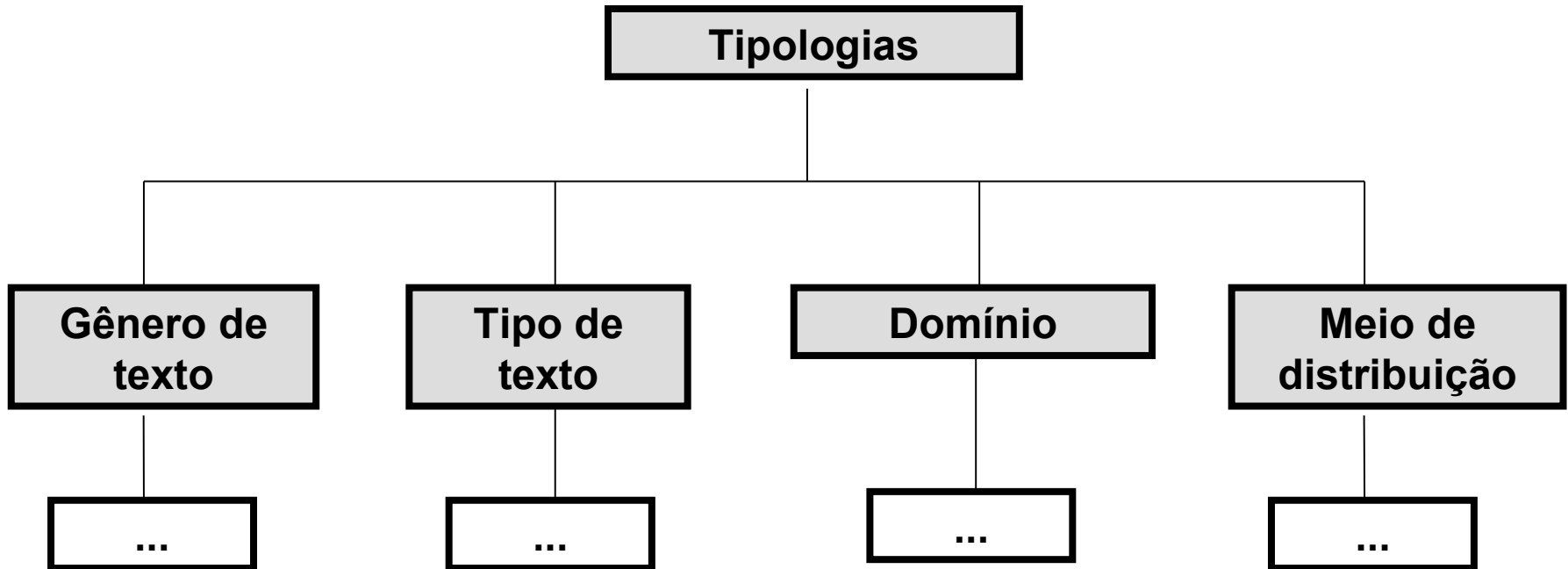




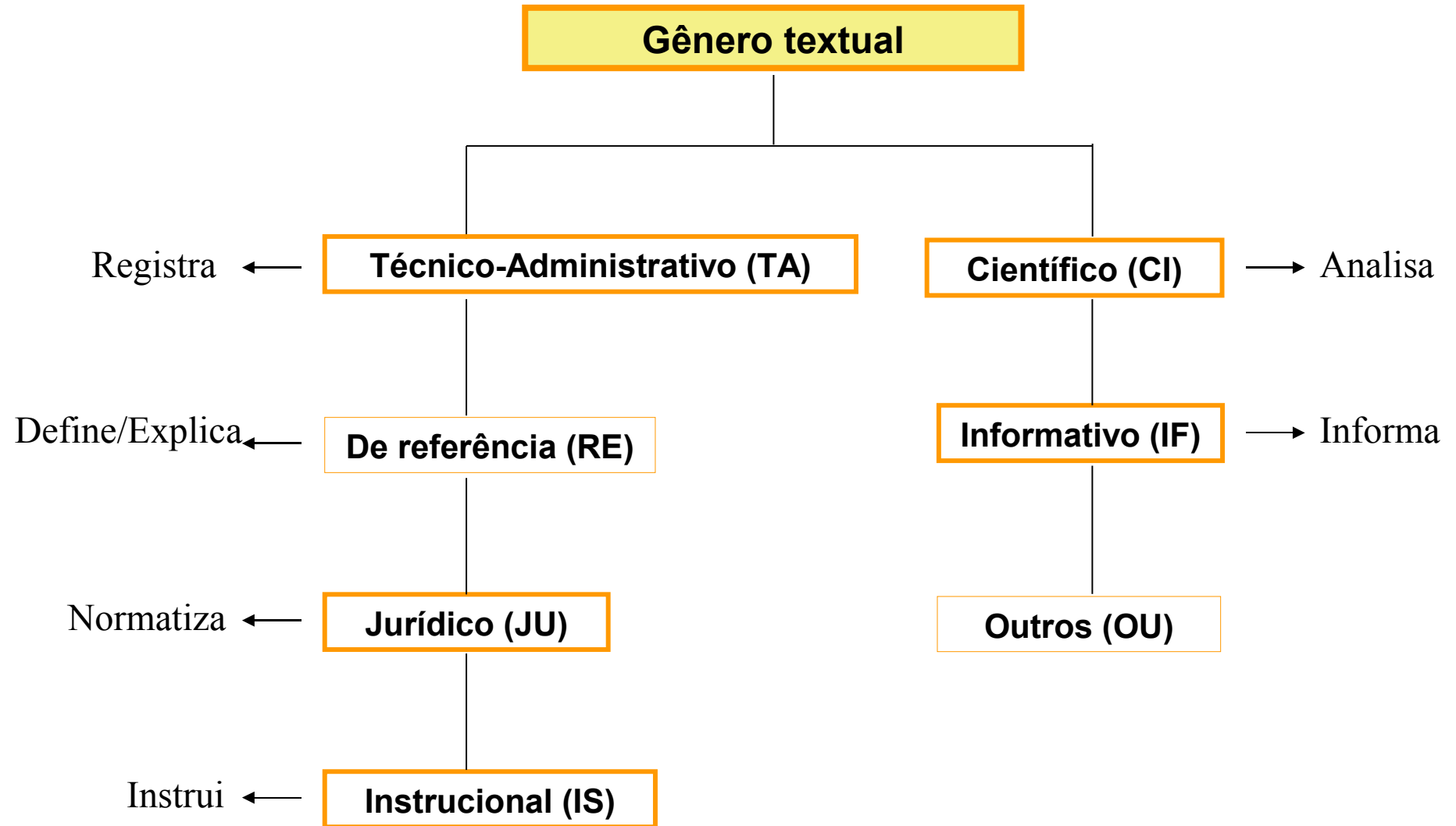




**OBS:** Este esquema deve ser produzido para duas outras categorias de autoria: 1) autor da síntese 2) autor da resenha e 3) orientador. Neste último caso, a categoria “corporação” não deve constar dentre as possibilidades.

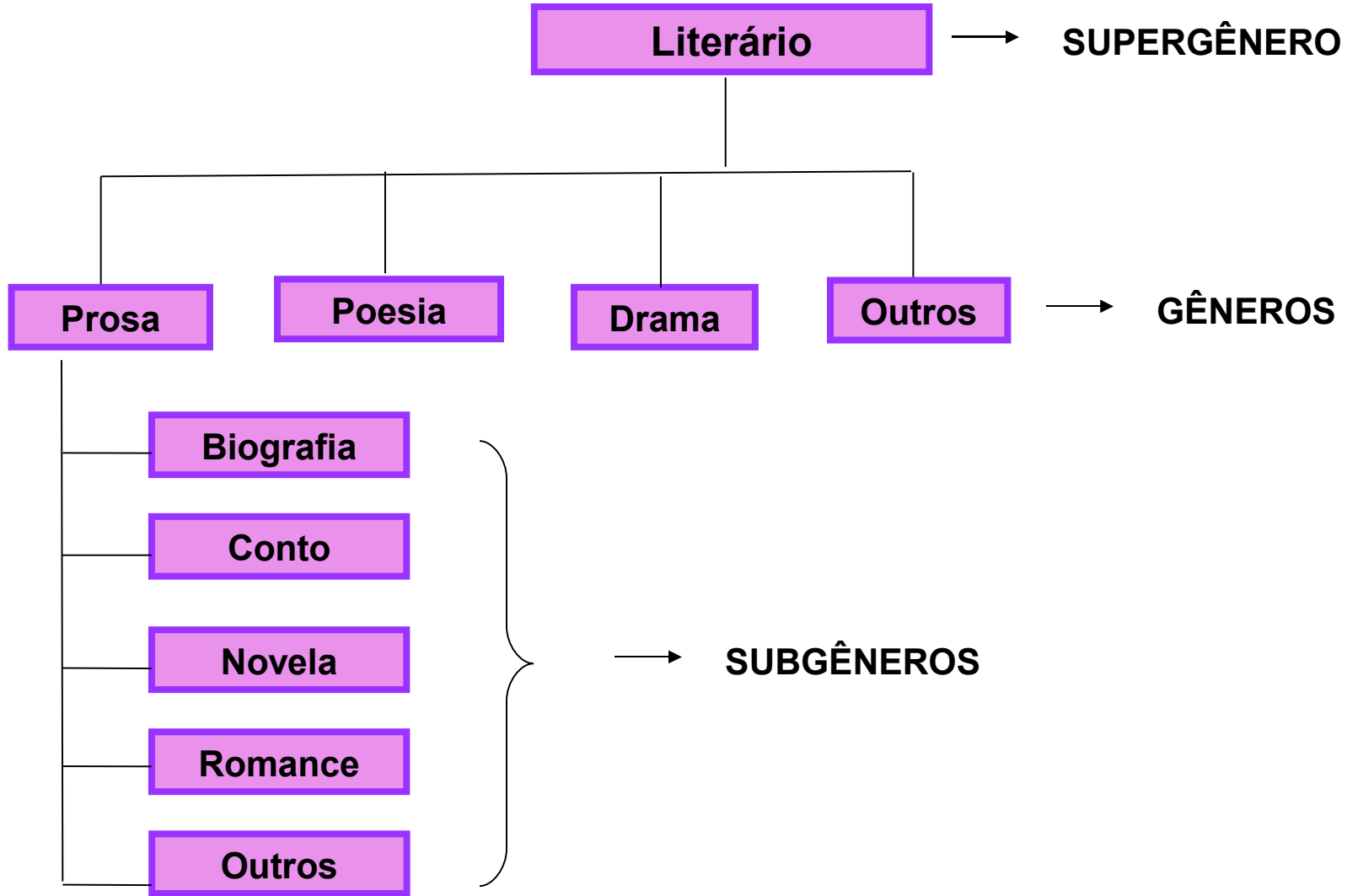


## *Gêneros de texto*



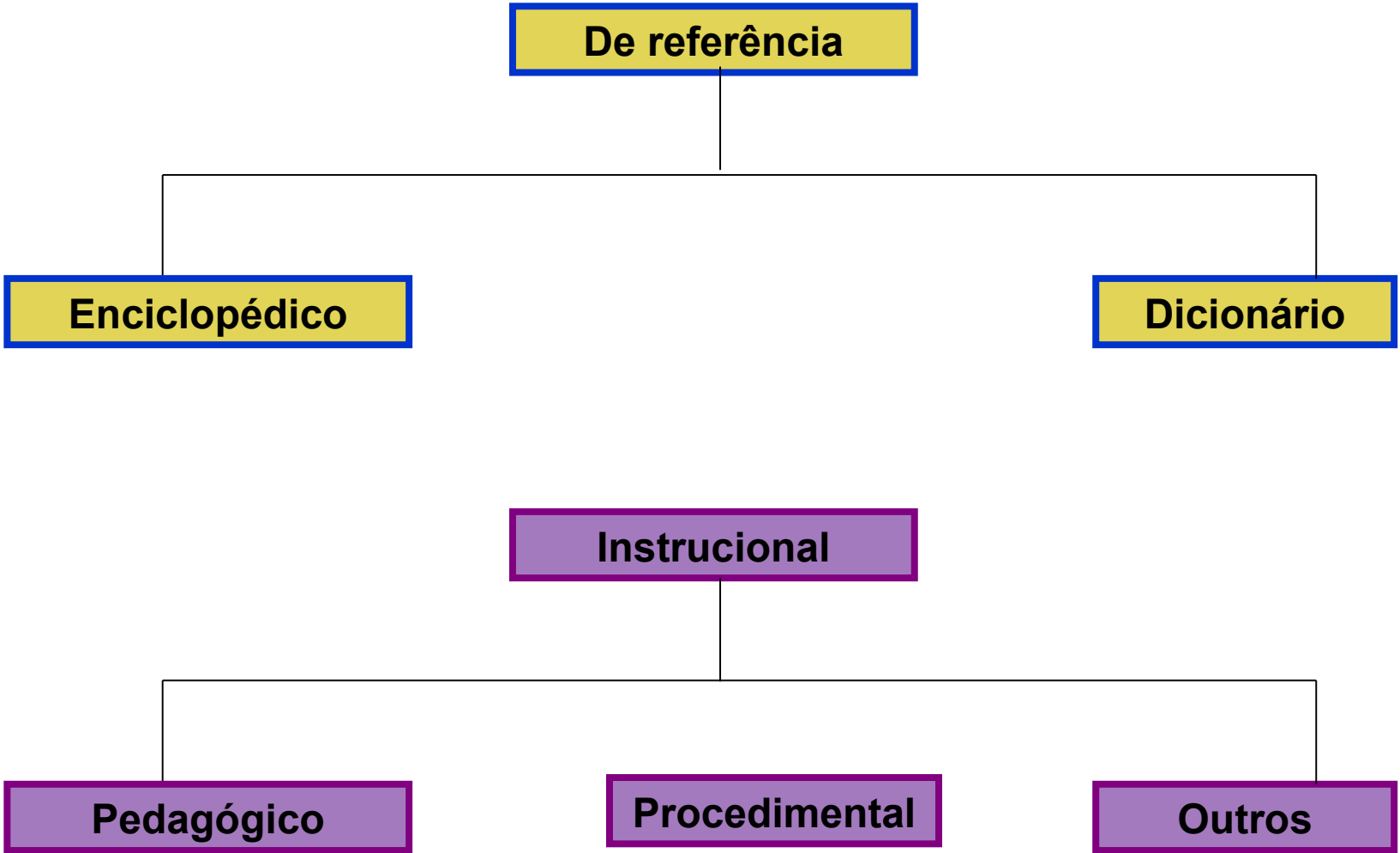
**Critérios distintivos:** formação do gênero (Bakhtin) / comunidade discursiva (Swales) / Intenção comunicativa

# *Supergênero, gêneros e subgêneros de texto*

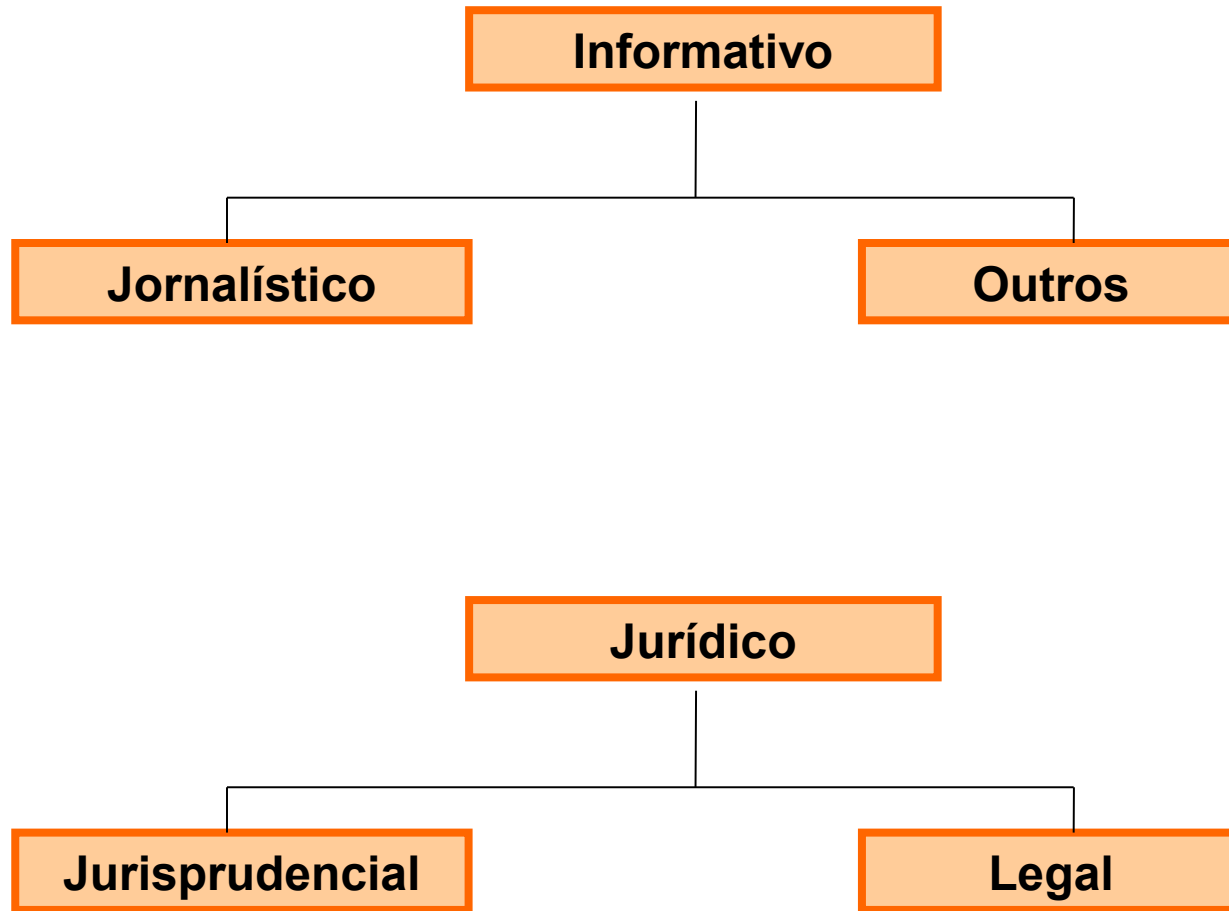


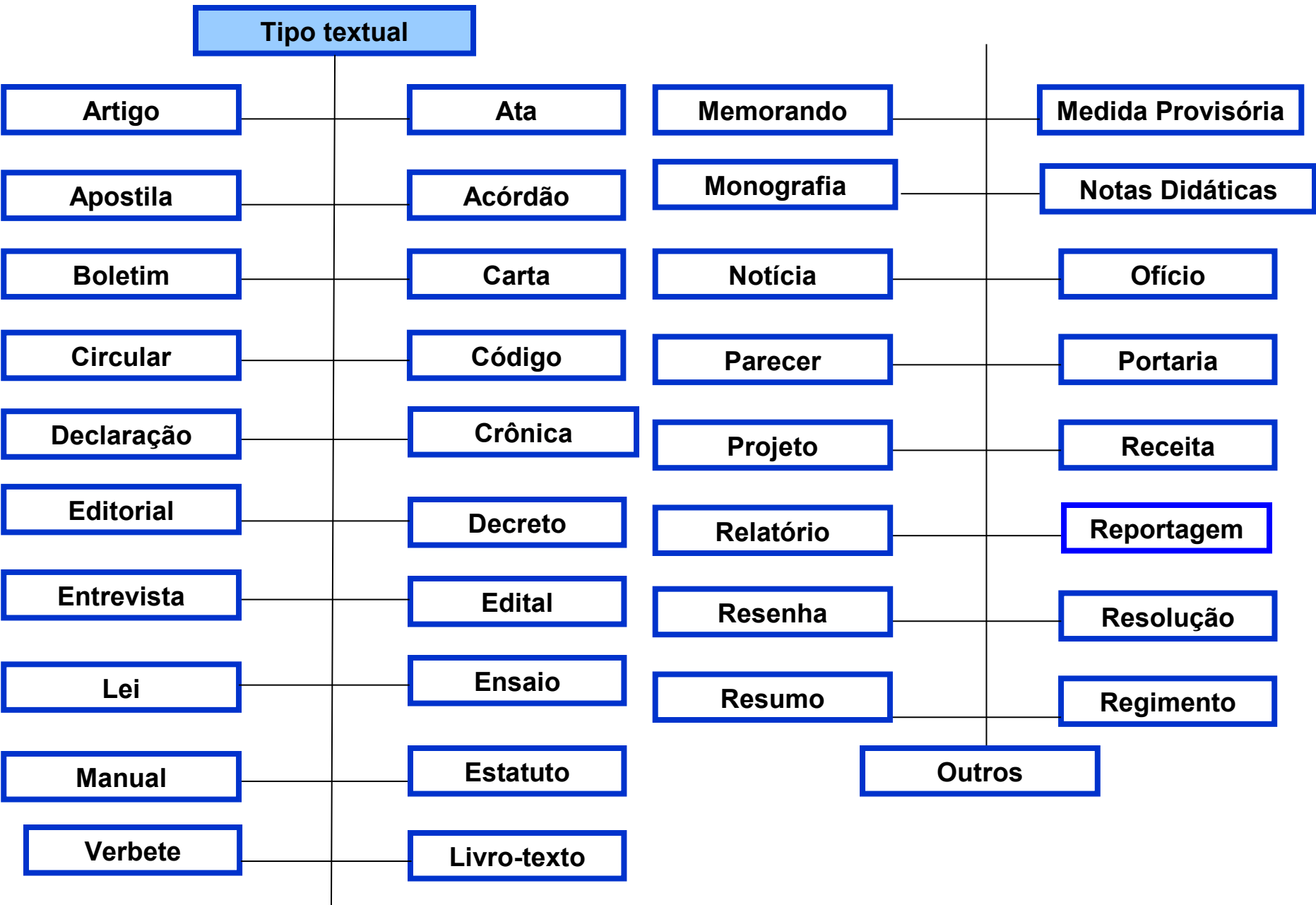
**Critérios distintivos:** Divisão histórica clássica, já consagrada por teóricos literários

*Gêneros e subgêneros de texto*



# *Gêneros e subgêneros de texto*







**Domínios**

**Científico**

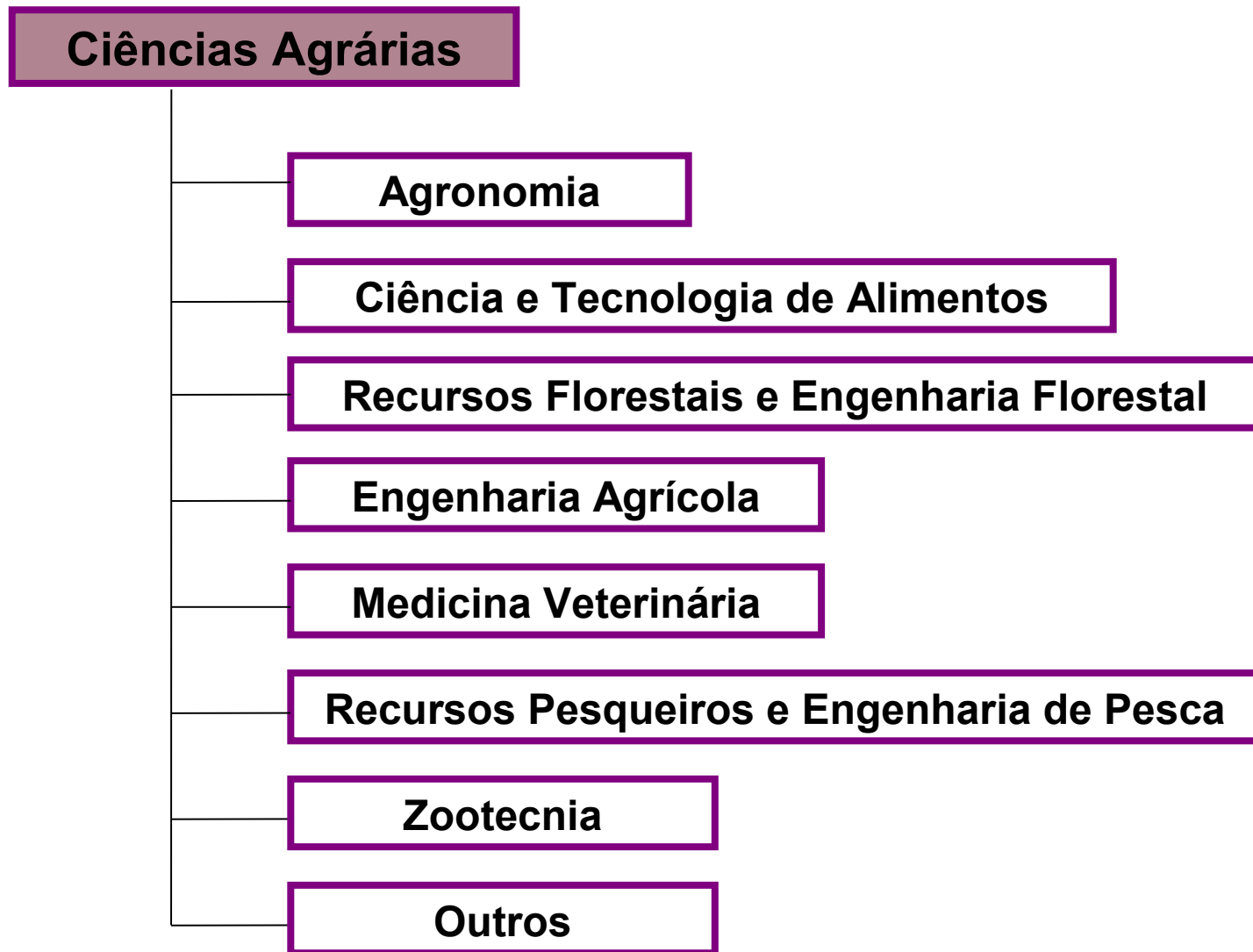
**Religião &  
Pensamento**

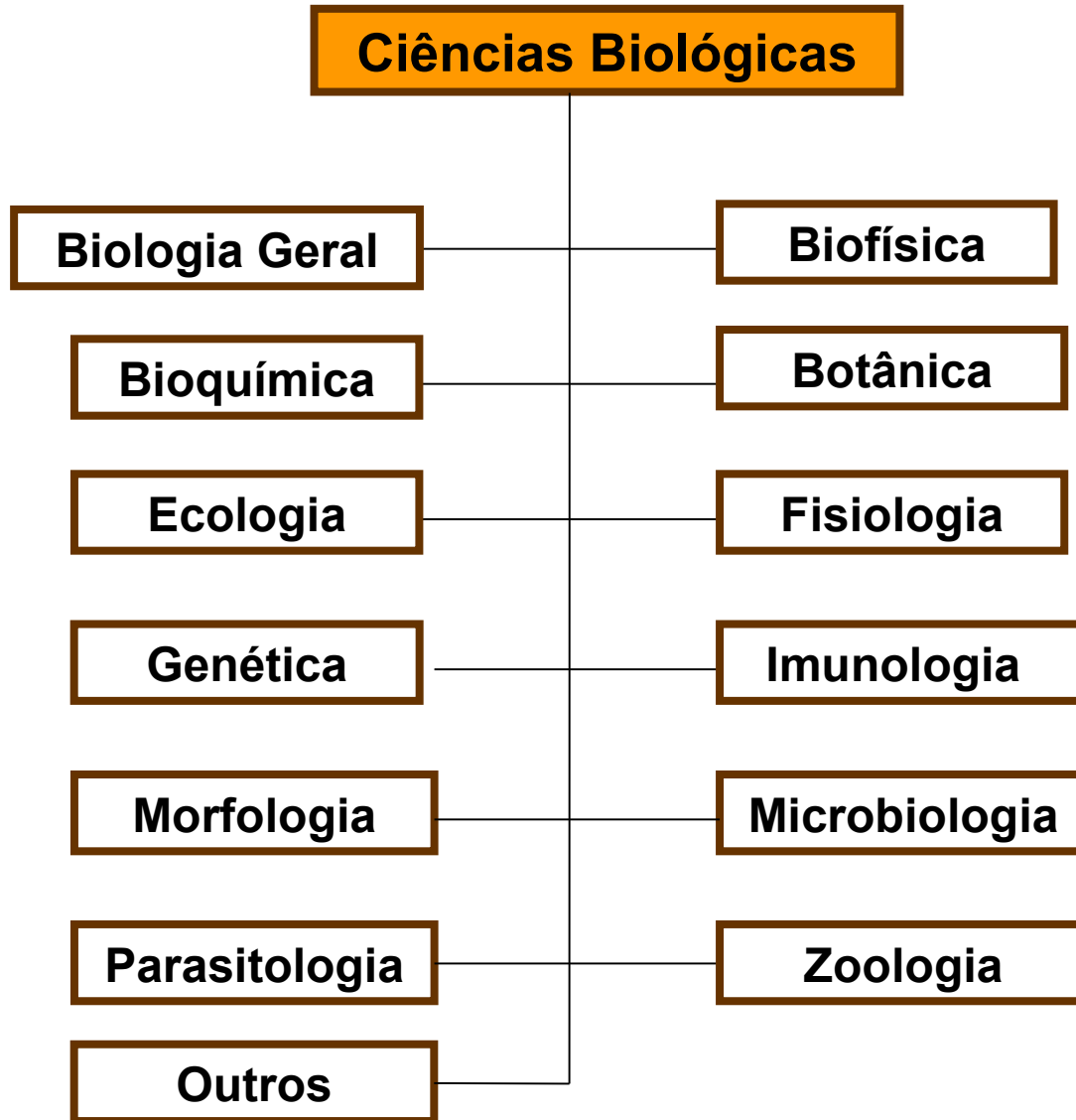
**Generalidades**

...

...

...





**Ciências da Saúde**

```
graph TD; A[Ciências da Saúde] --- B[Educação Física]; A --- C[Enfermagem]; A --- D[Farmácia]; A --- E[Fisioterapia]; A --- F[Fonoaudiologia]; A --- G[Medicina]; A --- H[Nutrição]; A --- I[Odontologia]; A --- J[Saúde Coletiva]; A --- K[Terapia Ocupacional]; A --- L[Outros];
```

**Educação Física**

**Enfermagem**

**Farmácia**

**Fisioterapia**

**Fonoaudiologia**

**Medicina**

**Nutrição**

**Odontologia**

**Saúde Coletiva**

**Terapia Ocupacional**

**Outros**

# Ciências Exatas e da Terra



# Ciências Humanas

Antropologia

Arqueologia

Ciências Contábeis

Ciência Política

Filosofia

Geografia

História

Jornalismo

Lingüística, Letras e Artes

Pedagogia

Psicologia

Publicidade e Propaganda

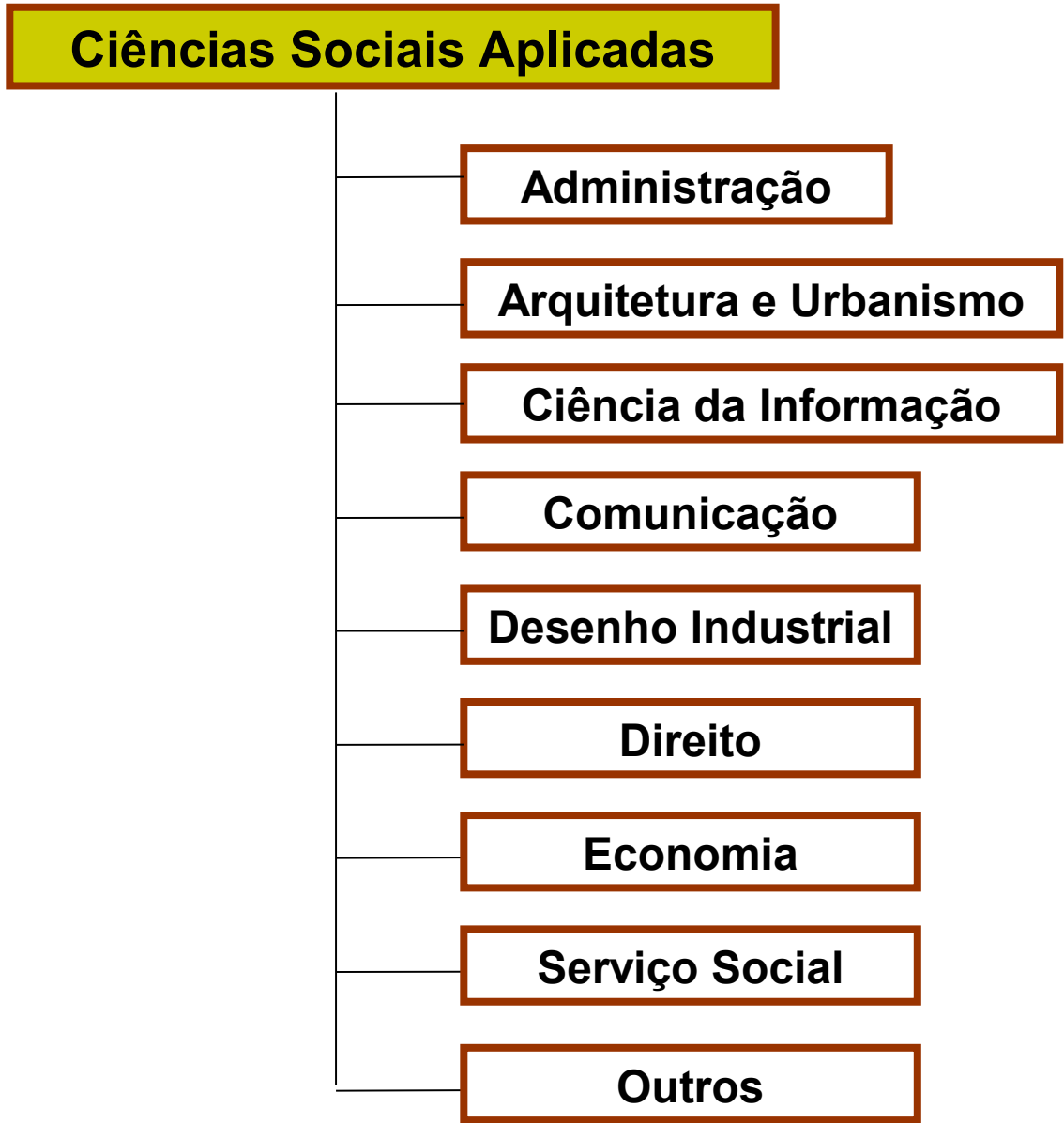
Relações Públicas

Sociologia

Turismo

Teologia

Outros



**Religião & Pensamento**

```
graph TD; A[Religião & Pensamento] --- B[Auto-ajuda]; A --- C[Magia e Bruxaria]; A --- D[Religião]; A --- E[Outros];
```

**Auto-ajuda**

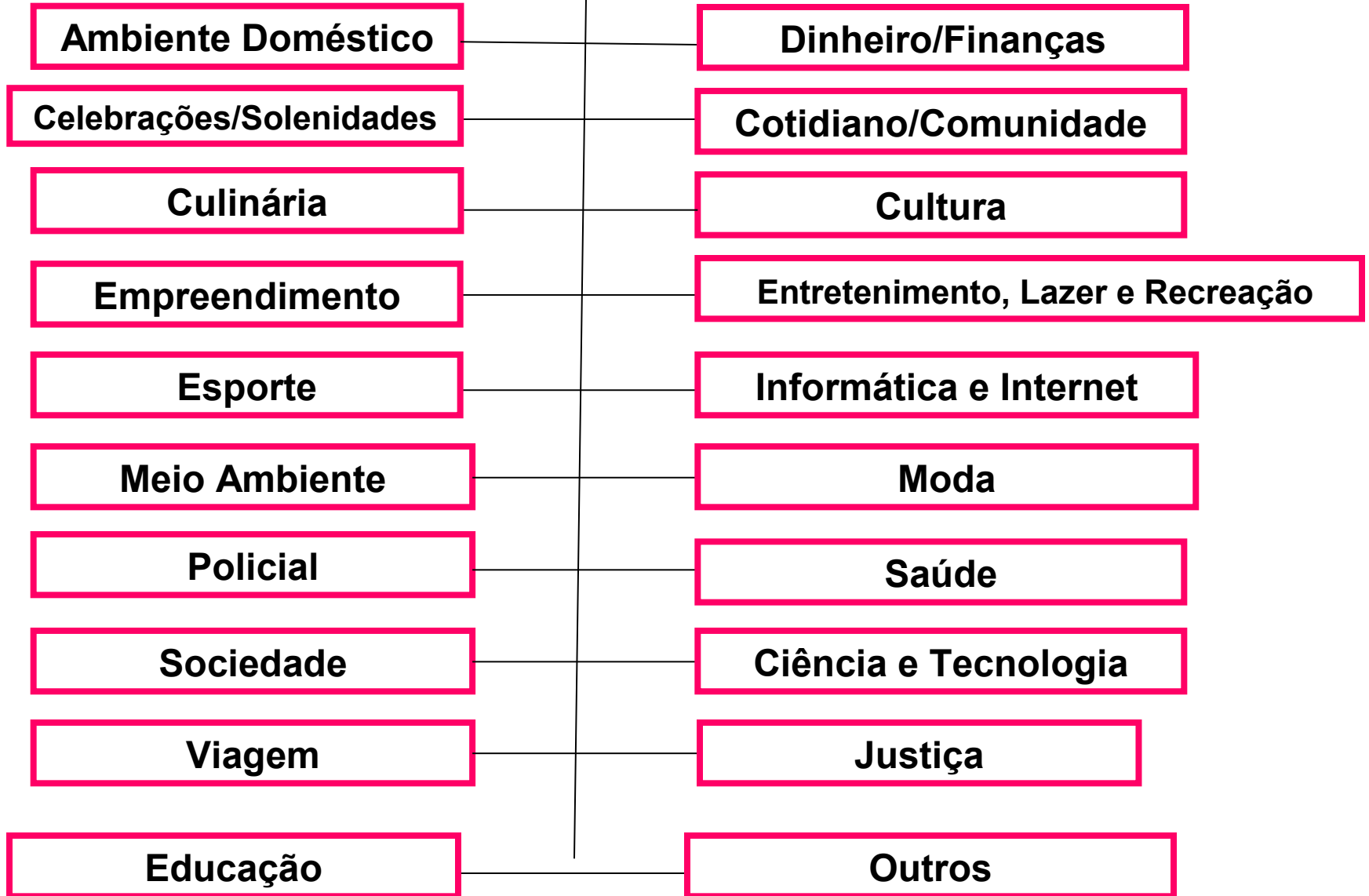
**Magia e Bruxaria**

**Religião**

**Outros**



## Generalidades



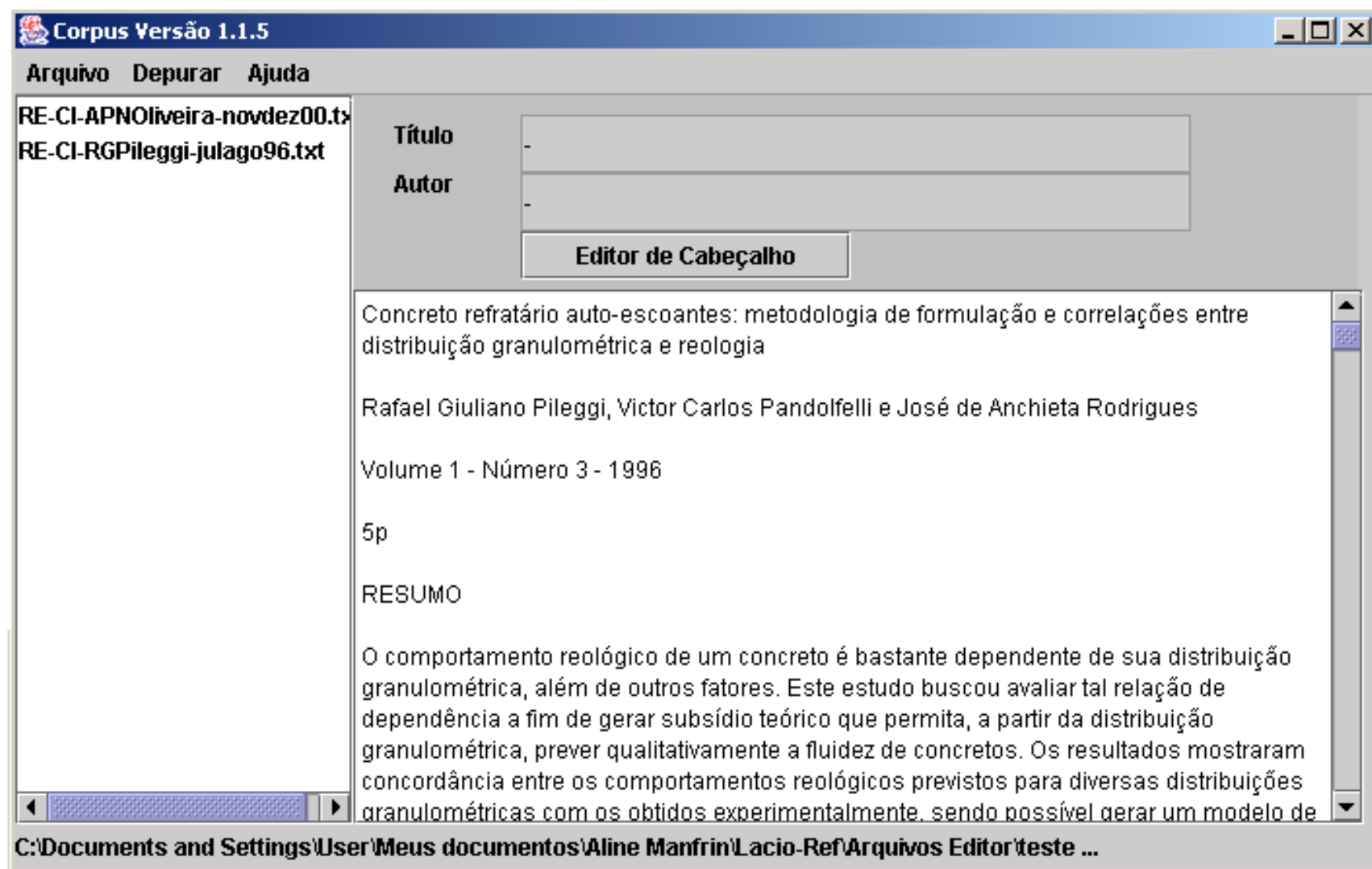
## Meio de Distribuição



# Nomeação e formatação dos arquivos

- Nomeação:
  - Meio, super-gênero e gênero e informações particulares do tipo de texto:
    - IN-JU-SE-03mar99.txt
- Formatação
  - Transformação para txt, marcação da presença de elementos gráficos

# Editor de Cabeçalho – após nomeação e formatação





Cabeçalho

**Nome do arquivo**

nfrin\Lacio-Ref\Arquivos Editor\teste 115\RE-CI-RGPileggi-julago96.txt

**Corpus**

Referência



**Número de Páginas**

6

**No. Palavras**

2254

**Amostra**

Íntegra



Arquivo

Texto

Autoria

Tipologias

### Cabeçalho

**Título**

formulação e correlações entre distribuição granulométrica e reologia

**Subtítulo**

**Língua**

Português do Brasil (PB) ▼

**Fonte**

**Editor**

**Local de Publicação**

**Data**

07-08.1996

**Status**

- ▼

**Comentários**

Arquivo

Texto

Autoria

Tipologias



Cabeçalho

**Autoria de**

Texto ▼

**Tipo de Autoria/ texto**

Múltiplo ▼

**Nome do autor do texto**

Giuliano Pileggi, Victor Carlos Pandolfelli, José de Anchieta Rodrigues

**Sexo do autor do texto**

Masculino, Masculino, Masculino

**Sexo de Rafael Giuliano Pileggi**

Masculino ▼

**Sexo de Victor Carlos Pandolfelli**

Masculino ▼

**Sexo de José de Anchieta Rodrig...**

Masculino ▼

Arquivo

Texto

Autoria

Tipologias



### Cabeçalho

**Gênero**

Científico ▼

**Tipo Textual**

Artigo ▼

**Domínio geral**

Ciências Exatas e da Terra ▼  
Ciências Exatas e da Terra ▼  
- ▼

**Domínio específico**

Engenharia Civil ▼  
Engenharia de Materiais ▼

**Definição**

Anotador ▼  
Auto-def ▼

**Distribuição**

Revista (RE) ▼

Arquivo

Texto

Autoria

Tipologias



```
<?xml version="1.0" encoding="ISO-8859-1"?>
```

```
<!-- v:1.3.2-->
```

```
<text>
```

```
<header>
```

```
  <title>
```

```
    <fileName>IN-JU-SE-03mar99.txt</fileName>
```

```
    <corpus>Referência</corpus>
```

```
    <nPages>4</nPages>
```

```
    <nWords>1171</nWords>
```

```
    <sample>Íntegra</sample>
```

```
  </title>
```

```
<sourceText>
```

```
  <titleOfText>ACÓRDÃO nº 23.048</titleOfText>
```

```
  <language>Português do Brasil (PB)</language>
```

```
  <source>Supremo Tribunal Regional</source>
```

```
  <pubPlace>Natal</pubPlace>
```

```
  <date>03.03.1999</date>
```

```
  <status>Original</status>
```

```
  <comments>Agravo Regimental nº 03-00722/99-0</comments>
```

```
</sourceText>
```

Editor salva no  
formato XML

```
<author>
  <textAuthor>
    <name>Maria do Perpétuo Socorro Wanderley de Castro, Ezequiel
    Escolástico Bezerra, Nicodemos Fabrício Maia</name>
    <gender>Feminino, Masculino, Masculino</gender>
    <position>Juíza Presidente, Juiz, Procurador do Trabalho</position>
    <type>Múltiplo</type>
  </textAuthor>
</author>
<textClassification>
  <textGenre>
    <genre>Jurídico</genre>
  </textGenre>
  <textType>Sentença</textType>
  <domain>
    <generalDomain defined="annotador-
    def">Científico/Ciências Sociais Aplicadas</generalDomain>
    <specificDomain>Direito</specificDomain>
  </domain>
  <distribution>Internet</distribution>
</textClassification>
</header>
```

<body>  
<omit desc="head">  
ACÓRDÃO nº 23.048

Agravo Regimental nº 03-00722/99-0

Juiz Redator: Ezequiel Escolástico Bezerra

Agravante: Ministério Público do Trabalho

Agravado: Petróleo Brasileiro S/A – Petrobrás

Procedência: TRT 21ª Região/RN

Maria do Perpétuo Socorro Wanderley de Castro  
Juíza Presidente

Ezequiel Escolástico Bezerra  
Juiz designado para redigir o Acórdão

Nicodemos Fabrício Maia  
Procurador do Trabalho  
</omit>

Agravo Regimental em Mandado de Segurança. Concessão de liminar em Ação Civil Pública.  
Pertinência da medida concedida pelo órgão de Primeiro Grau. ....

# Base de dados que dá suporte às buscas para montagem de subcorpúpus

- Todas as informações do cabeçalho são guardadas numa base de dados, permitindo buscas:
  - **Pesquisa Simples:** meio, super-gênero e gênero
  - **Pesquisa Avançada:** além dos campos da Pesquisa Simples, permite buscar por dados de catalogação bibliográficas, como *Nome de Autor*, *Nome do Periódico* e *Caderno*.
  - **Pesquisa Personalizada:** o usuário deve definir detalhadamente o recorte de sua investigação. Campos de seleção como: o *Tipo de Amostragem*, o *Tamanho da Amostra*, o *Tipo de Autoria*, o *Tipo Textual* e o Domínio são apresentados ao usuário, sendo que a grande maioria deles possui conteúdos dinâmicos.

# O que aprendemos: o que faríamos diferente

- Nomeação dos arquivos mais simples
- Faríamos um editor de cabeçalho para a Web que após edição do cabeçalho automaticamente colocasse as informações no banco de dados
  - No LW o processo é dividido em 2 etapas pois o editor funciona no Windows e a subida à base é feita com ajuda de programas acionados manualmente
- Avaliaríamos:
  - o Plonetaryum da FAPESP (<http://plonetaryum.incubadora.fapesp.br/portal>) para fazer o mesmo Portal do LW com a infra-estrutura e servidores da incubadora de conteúdos digitais (<http://incubadora.fapesp.br/>)

# Proposta para o Projeto do Dic\_Hist

- Proposta para o Portal e mão-de-obra necessária
  
  
  
  
  
  
  
  
  
  
- Proposta para o gerenciamento do projeto:  
uso de uma Web colaborativa – coteia do  
ICMC

# Portal

- 1) Implementação de um site do projeto com informações sobre o projeto, equipe, publicações, manuais de anotação, etc. em português e também em inglês para ampliar a divulgação da pesquisa.
  - Esse site será integrado o site interno; os dois terão uniformidade de *design* gráfico.
  - 4 meses de pagamento equivalente a uma bolsa ITI A do CNPq
- 2) Modelagem de um banco de dados para guardar os cabeçalhos dos textos do *córpus* e permitir pesquisas rápidas no *córpus*.
  - Posterior inclusão dos dados a partir da planilha criada pelo Projeto Resgate e criação da estrutura de diretórios que abrigará fisicamente o *córpus* no servidor Web.
  - 1 ano de pagamento equivalente a uma bolsa ITI A do CNPq

- 3) Adaptação das interfaces Web de pesquisa e montagem de subcorpús de estudo do projeto Lácio-Web que comporão a parte interna do site do item 1), a princípio restrita aos membros do projeto durante a sua duração (3 anos).
  - Após esse prazo a parte interna ficará disponível a outros pesquisadores. Inclusão das ferramentas básicas como concordanceadores e contadores de freqüência utilizadas no projeto Lácio-Web na parte interna do site do corpús do projeto atual.
  - 4 meses de pagamento equivalente a uma bolsa ITI A do CNPq
- 4) Inclusão das ferramentas desenvolvidas em outros projetos do NILC como lematizadores, geradores de n-gramas, extratores de termos, separadores de sentenças, interfaces com o léxico do NILC que pode ser utilizado como dicionário de exclusão na verificação de arcaísmos na parte interna do site do corpús do projeto atual.
  - Treinamento para os lingüistas do projeto quanto ao uso das ferramentas e interpretação dos dados.
  - 1 ano de pagamento equivalente a uma bolsa DTI H do CNPq

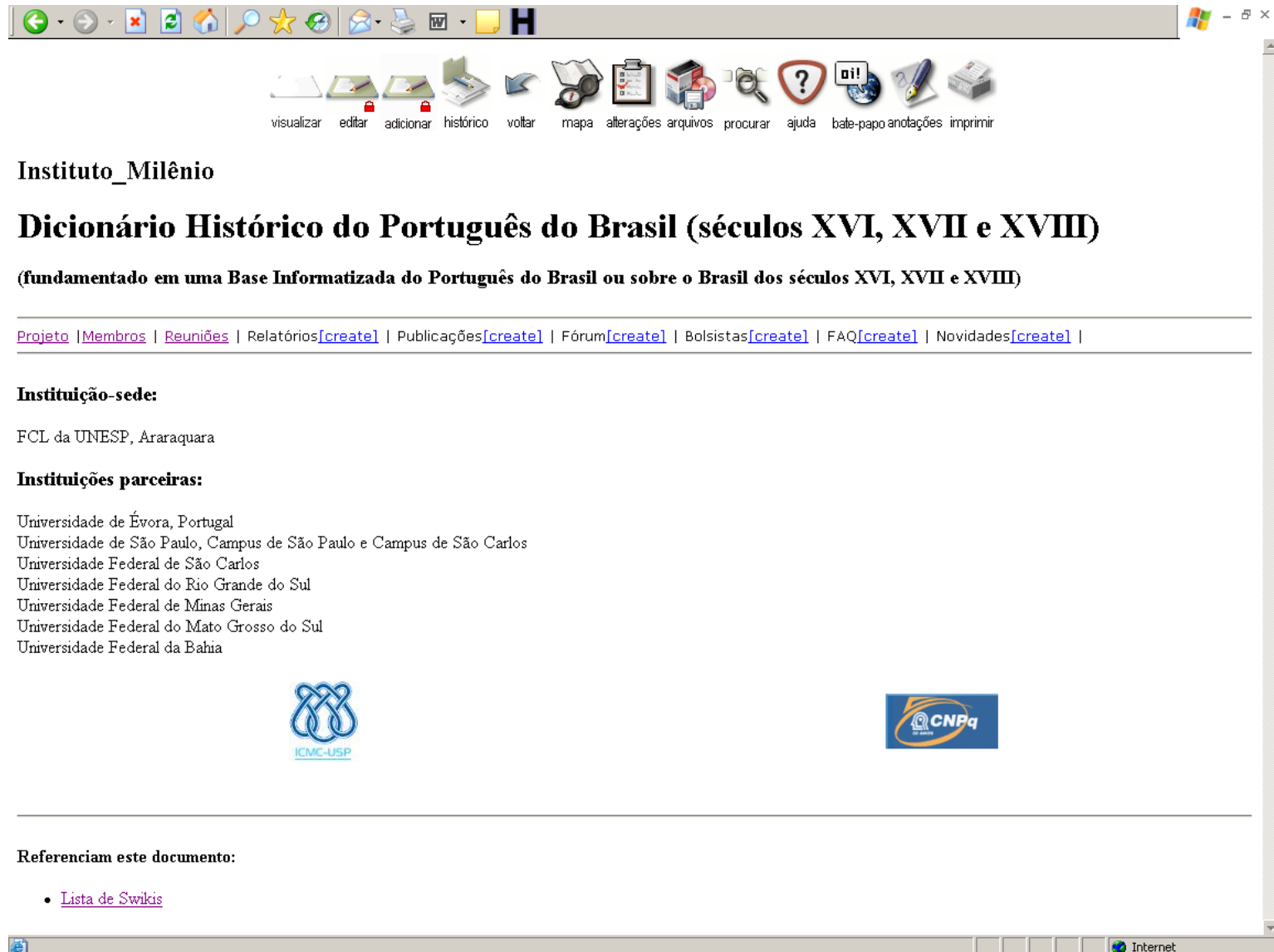


- 5) Anotação semi-automática (feita por ferramentas e corrigida por lingüistas) de uma parte do córpus de textos em português dos séculos XVI, XVII e XVIII (a princípio prevemos que esse córpus anotado tenha 500 mil palavras).
  - Essa tarefa prevê a avaliação e possível adaptação do conjunto de etiquetas utilizadas no projeto Lácio-Web para o córpus de português histórico.
  - Trabalho durante 1 ano de uma equipe contratada a partir do início do projeto que será formada por:
    - um graduado em letras com experiência de 2 anos (pagamento equivalente a uma bolsa DTI G do CNPq)
    - 3 alunos da graduação em letras (pagamento equivalente a uma bolsa ITI A do CNPq)
- 6) Treinamento de etiquetadores morfossintáticos com o córpus de 500 mil palavras anotado.
  - Inclusão dos etiquetadores criados na parte interna do site, além de concordanceadores que façam a pesquisa por palavra e etiqueta.
  - 4 meses de pagamento equivalente a uma bolsa ITI A do CNPq

- 7) Codificação do corpus na linguagem de marcação XML.
  - Essa codificação prevê a marcação do cabeçalho dos textos (cujas informações ficarão no banco de dados) e
  - a marcação da estrutura geral - capítulos, parágrafos, títulos e subtítulos, notas de rodapé e elementos gráficos e também a marcação da estrutura de subparágrafos - elementos que são de interesse lingüístico, tais como sentenças, citações, palavras, abreviações, nomes, referências, datas e ênfase.
  - Essa anotação do corpus seguindo padrões internacionais permitirá a sua expansão, facilidade de uso e reuso em outros projetos.
  - De interesse principal aqui estão a marcação dos parágrafos e sentenças para geração de estatísticas sobre os textos do corpus e também abreviações e transcrições ortográficas. Essa anotação explícita dos elementos dos textos permitirá a geração automática de listas de abreviações

1 ano de pagamento equivalente a uma bolsa DTI D do CNPq

# Proposta para o gerenciamento do projeto: Web colaborativa do ICMC



The screenshot shows a web browser window with a toolbar at the top containing icons for back, forward, home, search, and other functions. Below the toolbar is a row of icons for various actions: visualizar, editar, adicionar, histórico, voltar, mapa, alterações, arquivos, procurar, ajuda, bate-papo, anotações, and imprimir. The main content area of the browser displays the following text:

**Instituto\_Milênio**

**Dicionário Histórico do Português do Brasil (séculos XVI, XVII e XVIII)**

**(fundamentado em uma Base Informatizada do Português do Brasil ou sobre o Brasil dos séculos XVI, XVII e XVIII)**

---

[Projeto](#) | [Membros](#) | [Reuniões](#) | [Relatórios\[create\]](#) | [Publicações\[create\]](#) | [Fórum\[create\]](#) | [Bolsistas\[create\]](#) | [FAQ\[create\]](#) | [Novidades\[create\]](#) |

---

**Instituição-sede:**

FCL da UNESP, Araraquara

**Instituições parceiras:**

Universidade de Évora, Portugal  
Universidade de São Paulo, Campus de São Paulo e Campus de São Carlos  
Universidade Federal de São Carlos  
Universidade Federal do Rio Grande do Sul  
Universidade Federal de Minas Gerais  
Universidade Federal do Mato Grosso do Sul  
Universidade Federal da Bahia

Below the text, there are two logos: the ICMC-USP logo on the left and the CNPq logo on the right.

---

**Referenciam este documento:**

- [Lista de Swikis](#)

The browser's status bar at the bottom shows the address bar and the text "Internet".

# Endereço da Coteia

- <http://coteia.icmc.usp.br/coteia/>  
– Instituto\_Milênio
  
- <http://coteia.icmc.usp.br/coteia/mostra.php>