



Ciclo de  
Encontros:  
4 x PLN

# Humanidades Digitais e PLN

Renata Vieira, Fernanda Olival,  
Helena Cameron

28 de outubro 2024

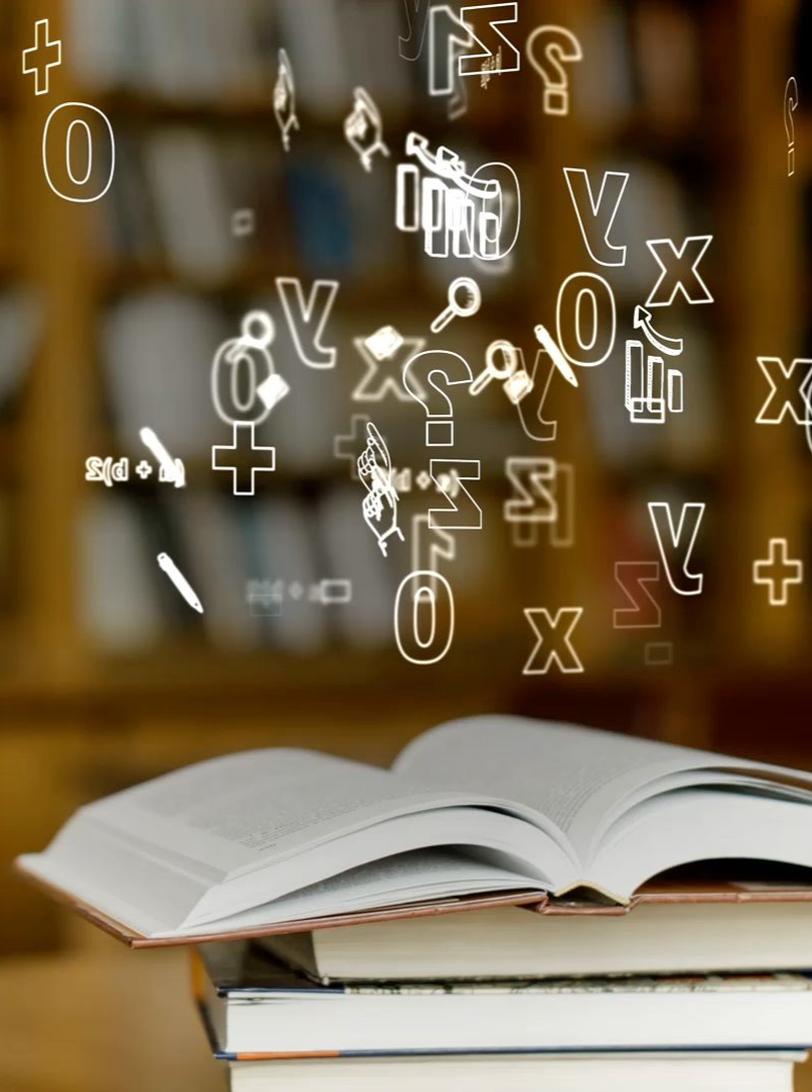
Fátima Farrica  
Maria José Bocorny Finatto  
Ana Paula Banza  
Ana Sofia Ribeiro  
Cassia Trojahn

# Sumário

---

- ❑ Introdução
- ❑ Preparação das fontes
- ❑ Transformação de textos em dados
- ❑ PLN, HD e Língua Portuguesa: Projetos
- ❑ Considerações finais

# Introdução



# Introdução

---

Na área de HD, em fontes textuais, encontramos uma grande variação

- nos **períodos históricos** das fontes,
- no seu suporte (**manuscritos em papel, impressos, fotografados**, etc),
- no seu estágio de digitalização, que pode variar entre **imagens digitais, textos em PDF e textos digitalizados** em outros formatos.

Todas essas variações adicionam esforços extras de processamento.



# Introdução

---

O capítulo discute os requisitos de **preparação e organização das fontes**, objetos de análise que depois de transcritas e **digitalizadas**, podem ser submetidas a **processamentos mais avançados**.

O objetivo é mostrar não apenas como o **PLN é útil e relevante nesse domínio**, mas também como a área de **HD é rica em despertar novas questões** para o desenvolvimento do PLN.



# Preparação das Fontes



Digitalização



Normalização



Metadados



Anotação

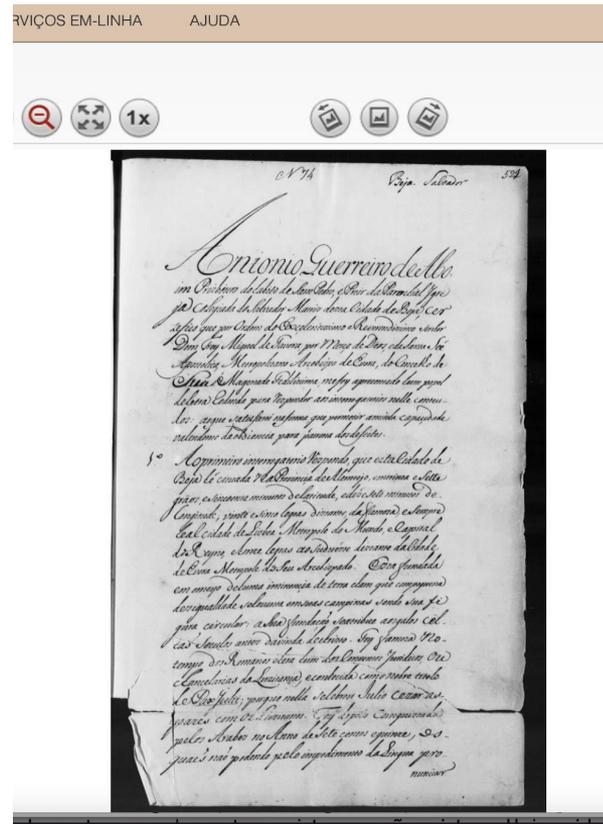
# Digitalização

Muitos projetos na área de HD necessitam lidar com a digitalização de **manuscritos** originais.

A área de paleografia, em particular, lida com a leitura e transcrição de versões manuscritas para um suporte atual.

Paleografia manual x automática (ex. Transkribus)

Digitalizar impressos (ex. eScriptorium, vFlat scan)



# Metadados

Uma vez digitalizada a fonte ou corpus de estudo, é necessário pensar na **organização dos seus metadados**.

O material digital deve conter a informação sobre a qual acervo pertence, identificar unicamente cada arquivo e, quando pertinente, associar autoria, data e outros elementos pertinentes.

Os metadados podem descrever a estrutura do documento, identificar volumes, capítulos, páginas, cabeçalhos, notas de rodapé, numeração de páginas ou comentários adicionados aos originais.



# Normalização

---

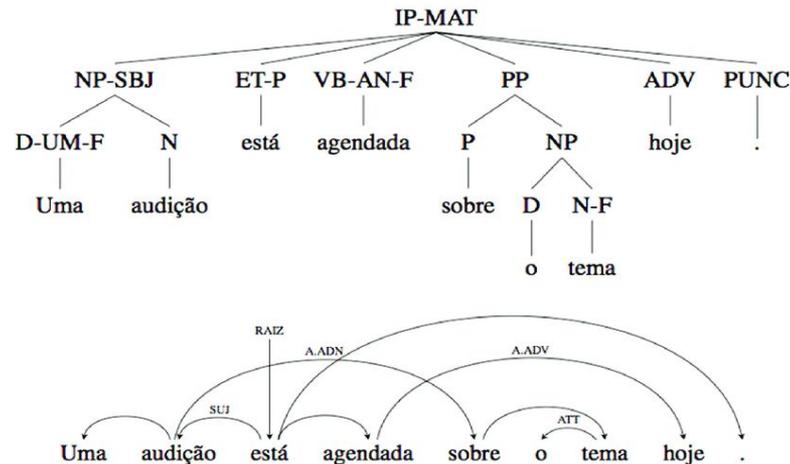
Em HD, as fontes textuais, quer manuscritas, quer impressas, sendo de uma época anterior ao estágio atual da língua, apresentam **variações ortográficas ou morfossintáticas**, não só em relação ao padrão atual como também dentro da mesma época.



# Anotação

A anotação textual ocorre em diferentes níveis, **sintático, semântico lexical, e discursivo** para identificar elementos de interesse e embasar diferentes estudos.

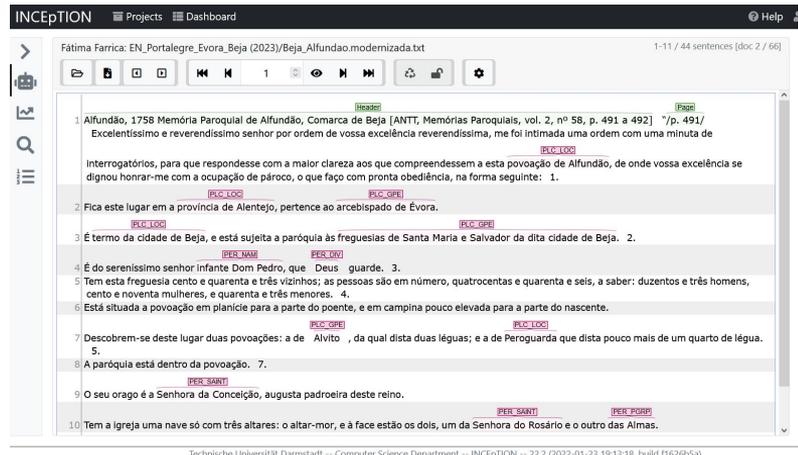
Geralmente a anotação sintática permite um tipo de análise diferenciada da anotação semântica.



# Anotação

A anotação de **entidades** é frequentemente adotada, com utilidade para estudos de literatura, história, geografia...

Conforme a complexidade do fenômeno e a respectiva disponibilidade de recursos, a anotação pode ser manual ou automática.



The screenshot displays the INCEPTION software interface. The top bar shows 'INCEPTION' and navigation options like 'Projects' and 'Dashboard'. The main window displays a document titled 'Fátima Farrica: EN\_Portalegre\_Evora\_Beja (2023)/Beja\_Alfundao.modernizada.txt' with 1-11 / 44 sentences. The document text is annotated with various entity types in red and green boxes, such as 'PERSON', 'LOCATION', and 'DATE'. The interface includes a search bar on the left and a toolbar with navigation and editing tools at the top.

Technical University Darmstadt -- Computer Science Department -- INCEPTION -- 22.7.2022-01-23 10:13:18 build f16206763

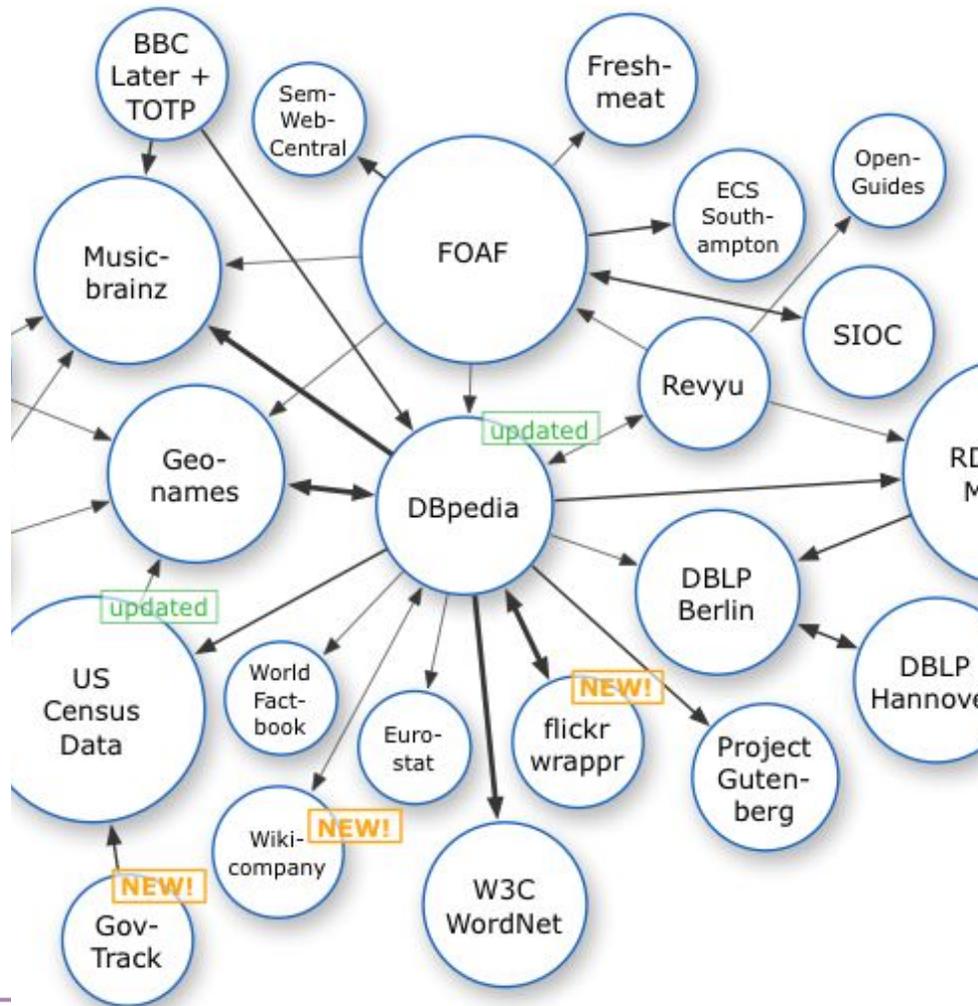


# Transformação de Texto em Dados

Extração de Informação

Representação de conhecimento e Ontologias

Dados Ligados e Dados FAIR



# Extração de Informação

---

O PLN atua como área responsável pela transformação de **textos em dados**.

Em textos das áreas das humanidades já digitalizados e normalizados, podem ser aplicadas ferramentas de PLN já desenvolvidas para diferentes tarefas:

- reconhecimento de **entidades nomeadas**,
- extração de **eventos**,
- resolução de **correferências**, e
- sistemas de **respostas a perguntas**

Amplia a **capacidade de análise**, pois pode-se organizar essa informação de forma ágil e rápida em estruturas bem definidas.



# Extração de Informação

Código	Nome	Ligação ao SO	Ofício	Morada_prc
	Diogo Pereira de Sampaio	Processo incompleto		Beira
14819	Domingos Fernandes	Familiar SO	Mercador de panos	Beira
	Domingos João	Familiar SO	Mercador	Beira
	Estêvão Fernandes	Familiar SO	Barbeiro	Estremadura
20437	Estêvão Lopes	Familiar SO	Tanoeiro	Estremadura
	Estêvão Pina	Familiar SO	Boticário	Estremadura
	Filipe Lourenço, Padre	Notário SO	Escrivão do Comissário SO	Beira
	Francisco de Araújo	Oficial SO	Dispenseiro SO	
	Francisco Gomes de Melo, Padre	Notário SO	Notário SO junto do comiss	Beira
	Francisco Gonçalves	Familiar SO	Alfaiate	Estremadura

# Representação do Conhecimento e Ontologias

---

Ontologias têm sido usadas em diversos domínios do conhecimento em diferentes tarefas:

- representar conhecimento de **domínios** complexos;
- organizar e anotar grandes quantidades de dados (**anotação semântica**);
- integrar dados de diversas fontes (**integração** semântica);
- dar ancoragem para o **raciocínio automático**;
- suportar maior **eficiência semântica** em buscas.



# Dados Ligados e Dados FAIR

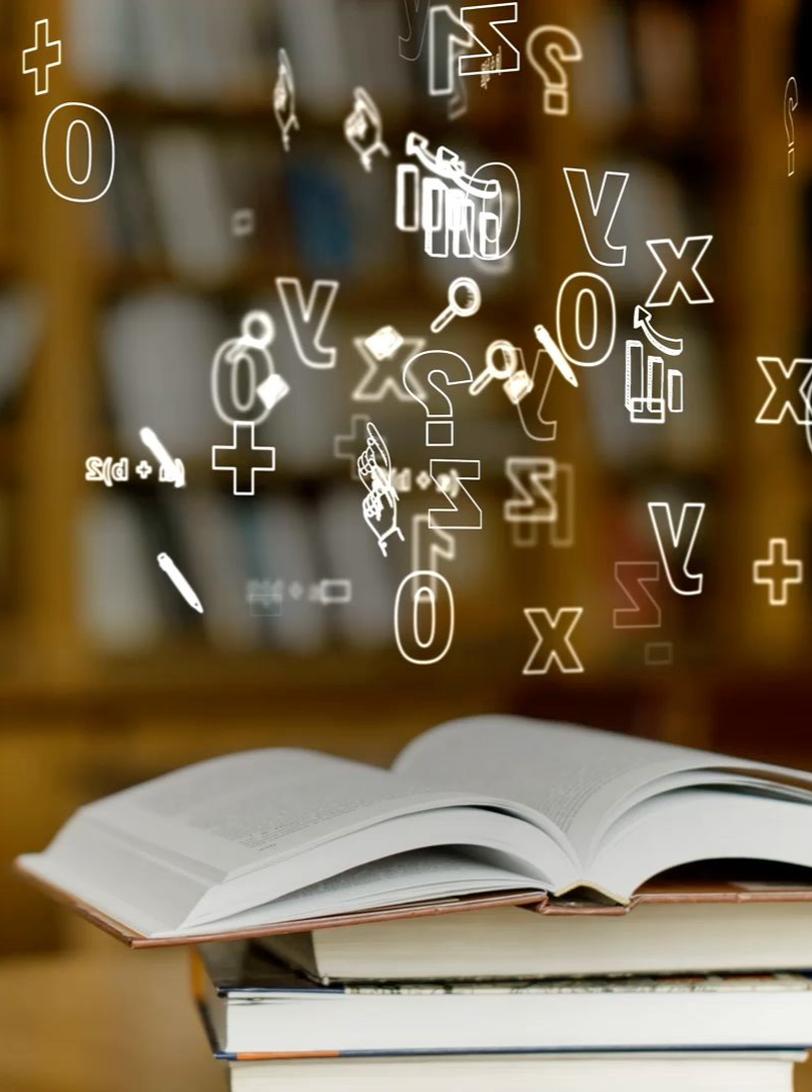
---

Findable, Accessible, Interoperable, Reusable

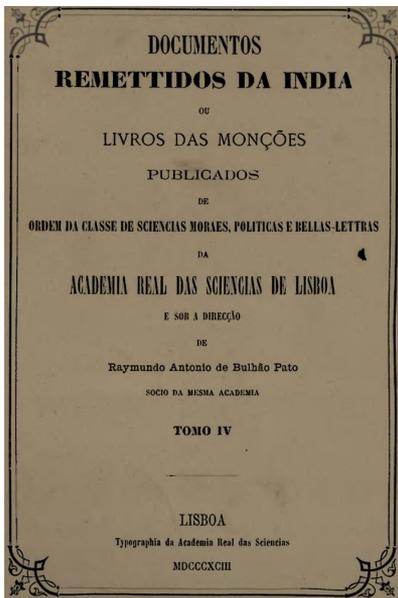
- (F) atribuição de **identificadores exclusivos** e persistentes a conjuntos de dados, descrevendo-os com metadados ricos que permitem sua indexação e descoberta;
- (A) uso de **protocolos abertos e de padrões** para acesso a conjuntos de dados;
- (I) uso de linguagens formais e de **vocabulários padronizados** FAIR para representar (meta)dados;
- (R) uso de metadados ricos sobre **licença de uso, proveniência e qualidade** de dados.



PLN, HD e Língua  
Portuguesa:  
Projetos

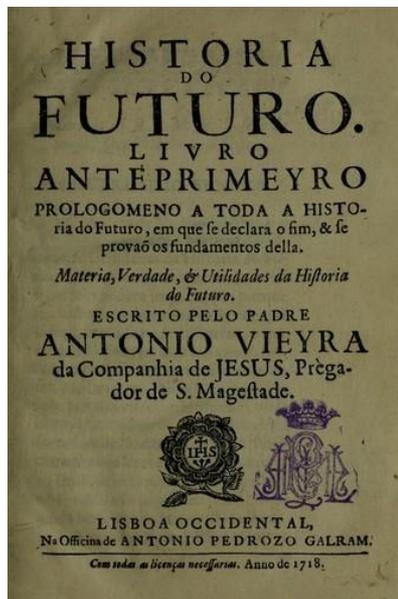






Extração de Eventos  
Análise de eventos de conflitos  
Anotação de Entidades Nomeadas  
Pesquisa por Entidades  
Relacionamento Eventos x Entidades

# Livro das Monções - Ana Sofia Ribeiro



Normalização

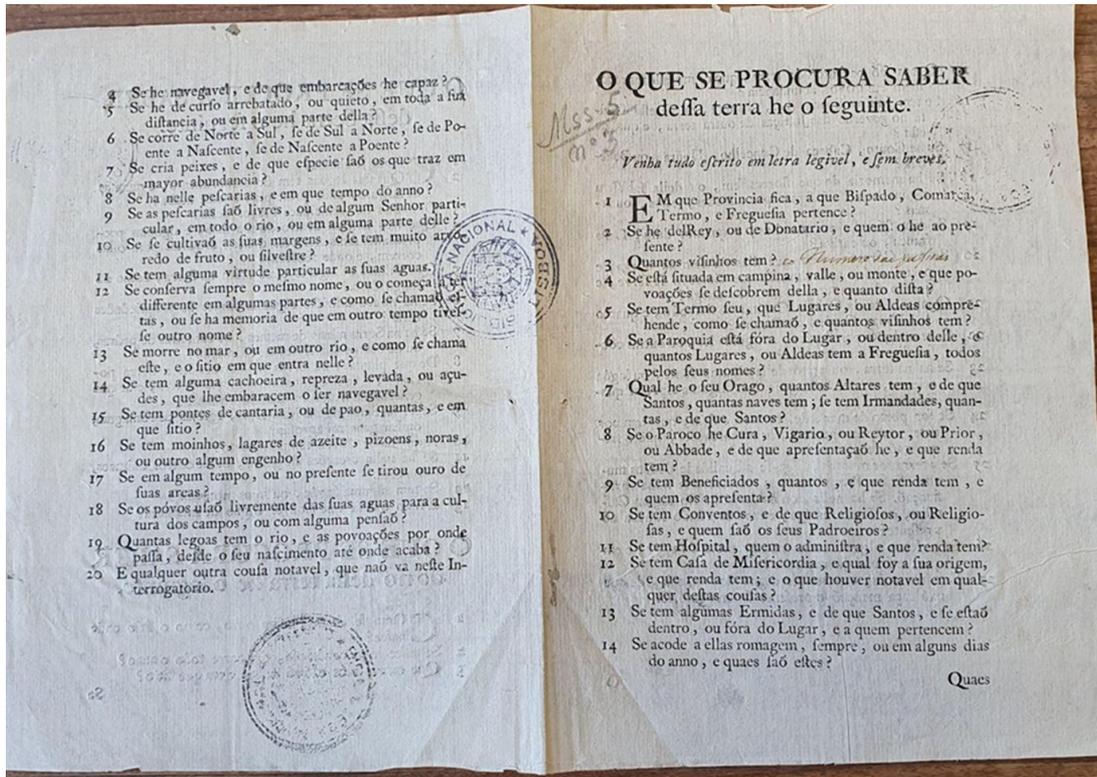
Análise morfológica e sintática

Elaboração do lexicon

Para breve: anotação semântica

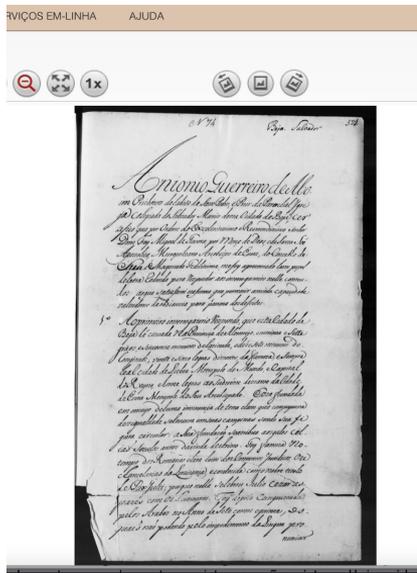
---

# História do Futuro - Ana Paula Banza



# Memórias Paroquiais - Fernanda Olival

BNP, Ms. 5, nº 3



Digitalização (manual)

Normalização (manual)

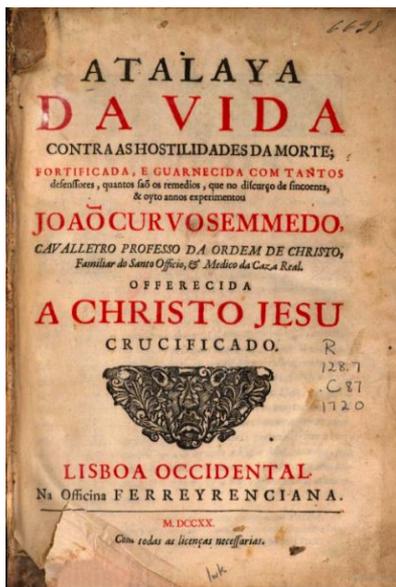
Anotação de Entidades Nomeadas

Reconhecimento de Entidades Nomeadas (BerTimbau, Albertina)

para expandir anotação para toda a coleção

Análise da fonte a partir de Entidades

# Memórias Paroquiais - Fernanda Olival

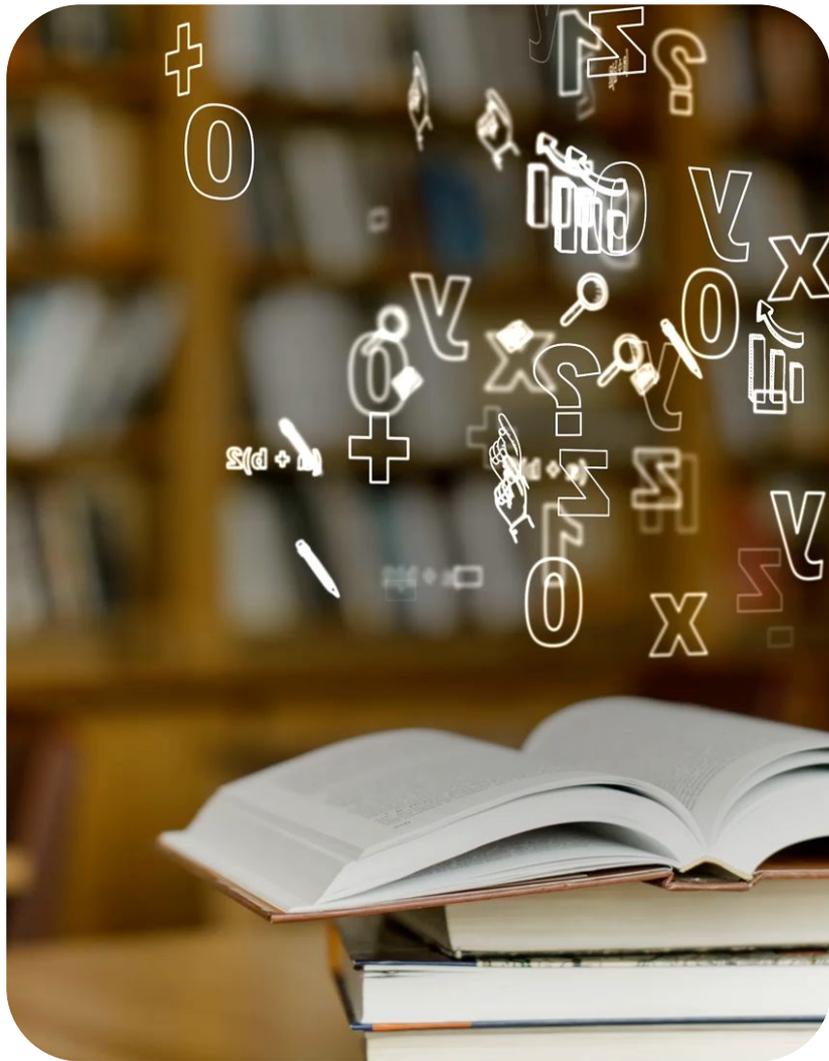


Mapear as antigas terminologias de Saúde em seus equivalentes da atualidade.

Estabelecer quais as áreas da Medicina, da Anatomia e que tipo de doenças e medicamentos eram conhecidos daquela época.

Terminologias usadas nos manuais de Curvo Semedo contrastadas com o vocabulário de Enfermagem

## Textos Médicos Séc XVIII - Maria José Finatto



## Considerações Finais

### **PLN e HD**

uma correlação de crescimento mútuo



# Considerações Finais

---

É preciso integrar:

- os objetivos de um pesquisador da área de humanidades, como linguística e literatura, ciências sociais, história,
- os objetivos de um pesquisador em PLN de encontrar, aplicar, desenvolver soluções apropriadas para cada tipo de problema, maximizando a correção e a performance das soluções encontradas de acordo com a disponibilidade de recursos.

Celebrando 3 edições do **Workshop de PLN e Humanidades Digitais** junto ao Propor



# Considerações Finais

---

A tecnologia faz parte da definição do **ser humano**

Quem produz tecnologia deve ter uma formação humanística

Atenta às origens e consequências do que produz

Por exemplo a ética, mas não só

Computação não pode ser um mistério acessível a poucos (como a leitura no início da imprensa)

Humanidades, como formação deve incluir conhecimentos de computação





<https://brasileiraspln.com/livro-pln/2a-edicao/parte-dominios/cap-humanidades-digitais>

Obrigada!



# Elementos úteis

