

# 1.Processamento de Linguagem Natural e LLM

Davi Bezerra Barros

## Introdução

NLP (Processamento de Linguagem Natural) é um campo da inteligência artificial que permite a máquinas entender, interpretar e gerar linguagem humana, facilitando a interação entre computadores e humanos por meio de texto ou fala.

### Dificuldades da NLP

- As linguagens são ambíguas, a nlp tem dificuldade em identificar ironia, por exemplo.
- As palavras tem significados diferentes, dependendo do contexto e posição da palavra na frase.
- modelos de NLP processam letras, palavras, sentenças, textos, documentos e corpos compostos de muitos documentos

### Gramática:

não é possível descrever formalmente as regras de uma linguagem, devido à sua complexidade e natureza dinâmica. Como as linguagens estão em constante mudança, elas requerem atualizações frequentes nos modelos.

## NLP com SpaCy

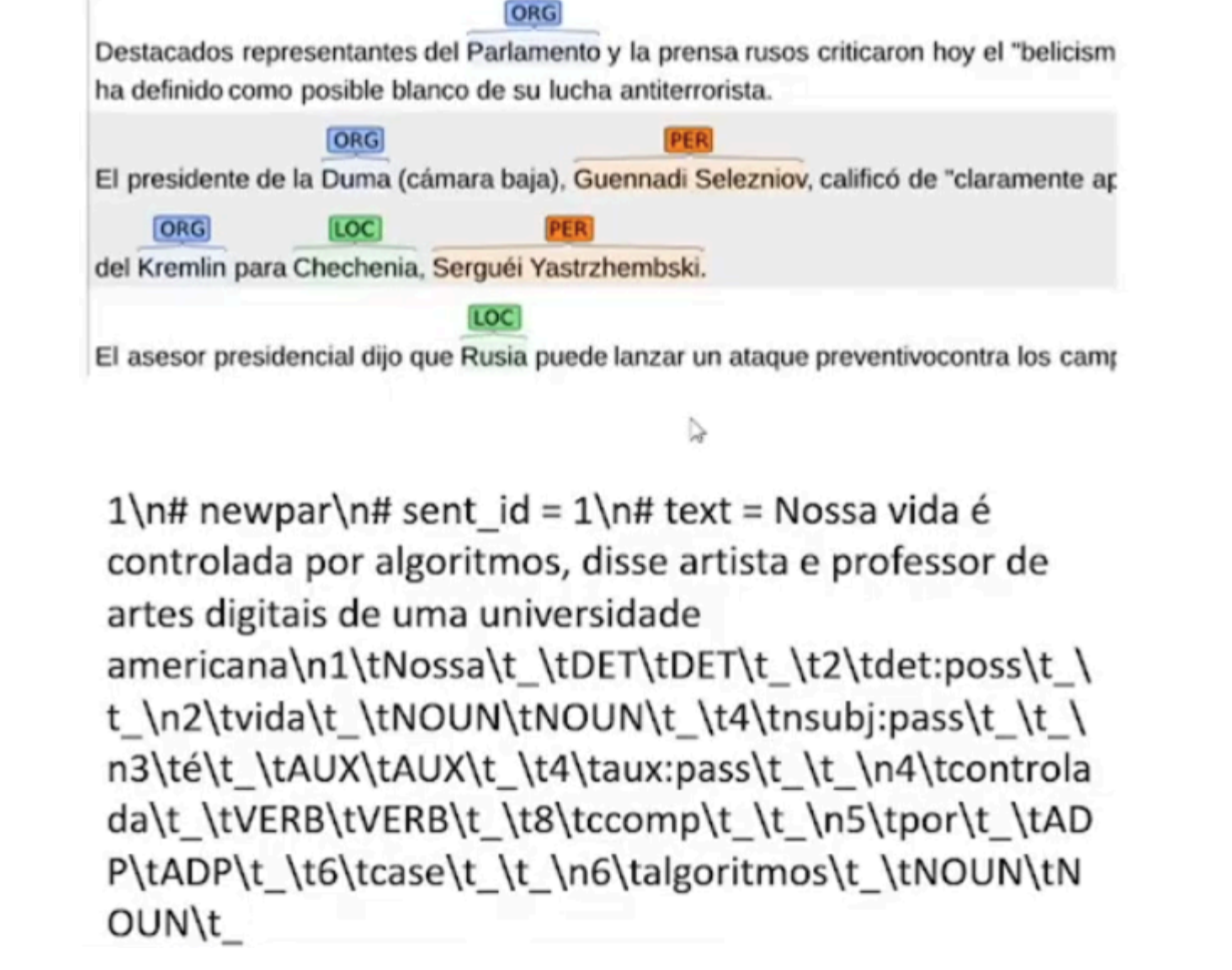
Os modelos pré treinados no spacy vem em três tamanhos:

- Pequeno: ~21MB
- Médio: ~48MB
- Grande: ~440 MB

## Conceitos básicos

**Corpus:** Um conjunto de documentos, textos não estruturados.

**Anotações:** Metadados adicionados ao texto, como tags gramaticais, relações sintáticas ou identificação de entidades.



**Tokenização:** Processo de separar a sentença em suas partes constituintes: palavras, pontos, símbolos etc.

**Parts-of-Speech Tagging(POS):** Adiciona tags a cada token, o caracterizando. Ferramentas são capazes de anotar tokens dentro do contexto, atribuindo significados diferentes.

**Lemma:** Traz a palavra na sua flexão, de modo que possam ser analisadas juntas. Cria uma representação única daquela palavra.

**Stemming:** Redução de palavras a suas raízes (ex: "correndo" → "corr"), menos preciso que a lematização.

**Dependency Parsing:** Encontra a relação de dependência hierárquica entre palavras

**Ngram:** Trata de palavras consecutivas. Bigrama trata duas palavras e trigramas, três. Mais palavras não serão processadas devido a esparsidade.

**Modelo:** É um banco de dados linguístico, específico de cada idioma. As bibliotecas de NLP utilizam seus próprios modelos, ou de terceiros.

## Word Embedding

Por conta da incapacidade do computador de processar informação não estruturada, é necessário criar formas de representar o texto para que ele entenda. Embeddinga técnica utilizada para mapear cada token a um vetor espacial, permitindo o processamento semântico das palavras ao relacionar seu significado com posições espaciais. É um processo fundamental para o processamento de linguagem natural, por que transforma dados categóricos, como palavras, em uma forma que pode ser processada por uma rede neural.

**One-hot encoding** é utilizado para criar vetores de características para identificar as palavras dentro de um conjunto de dados de um vocabulário / linguagem. O problema desta solução é que os vetores precisariam ser do tamanho do corpus, e não há nenhuma representação de contexto.

**TF-IDF:** É uma forma de representar as palavras de acordo com sua frequência no texto, trazendo o peso e a importância da palavra no documento.

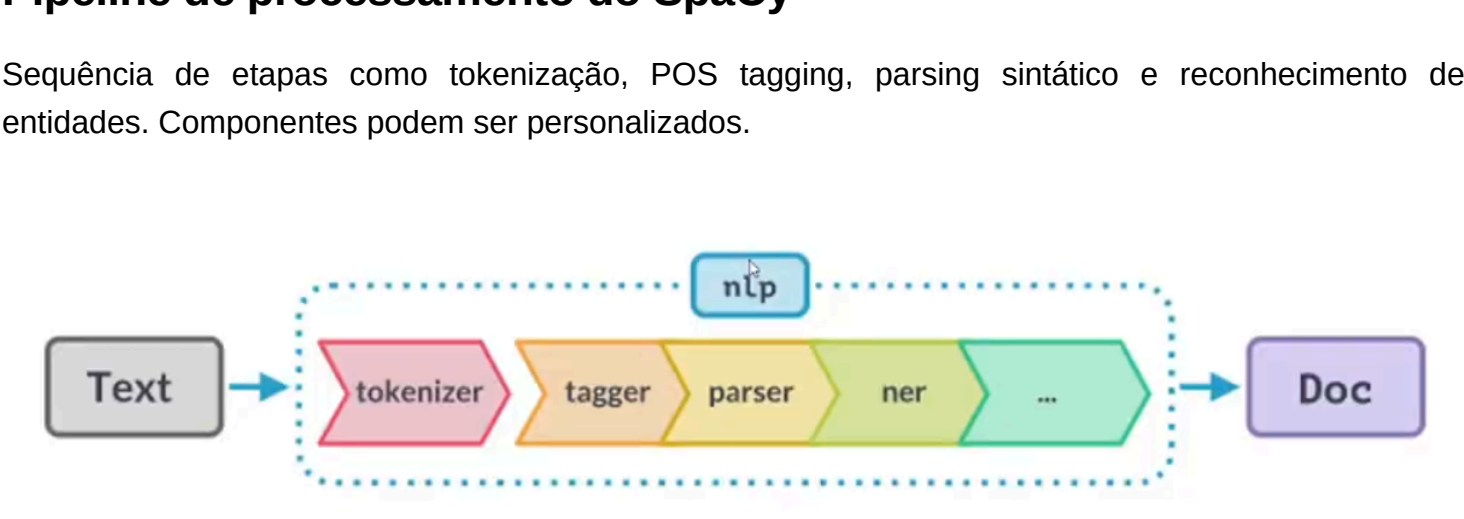
**Word2Vec:** Mostra além da frequência, a relação entre as palavras. através de um treinamento, a word2vec produz um vetor que demonstra matematicamente a relação entre as palavras. Tem duas formas principais:

- CBOW: busca prever uma palavra central no contexto
- Skip-gram: busca prever o contexto a partir de uma palavra central

O word2vec cria um vetor para cada palavra, que indic o nível de semelhaça semântica entre as palavras representadas por aquele vetor.

## Pipeline de processamento do SpaCy

Sequência de etapas como tokenização, POS tagging, parsing sintático e reconhecimento de entidades. Componentes podem ser personalizados.

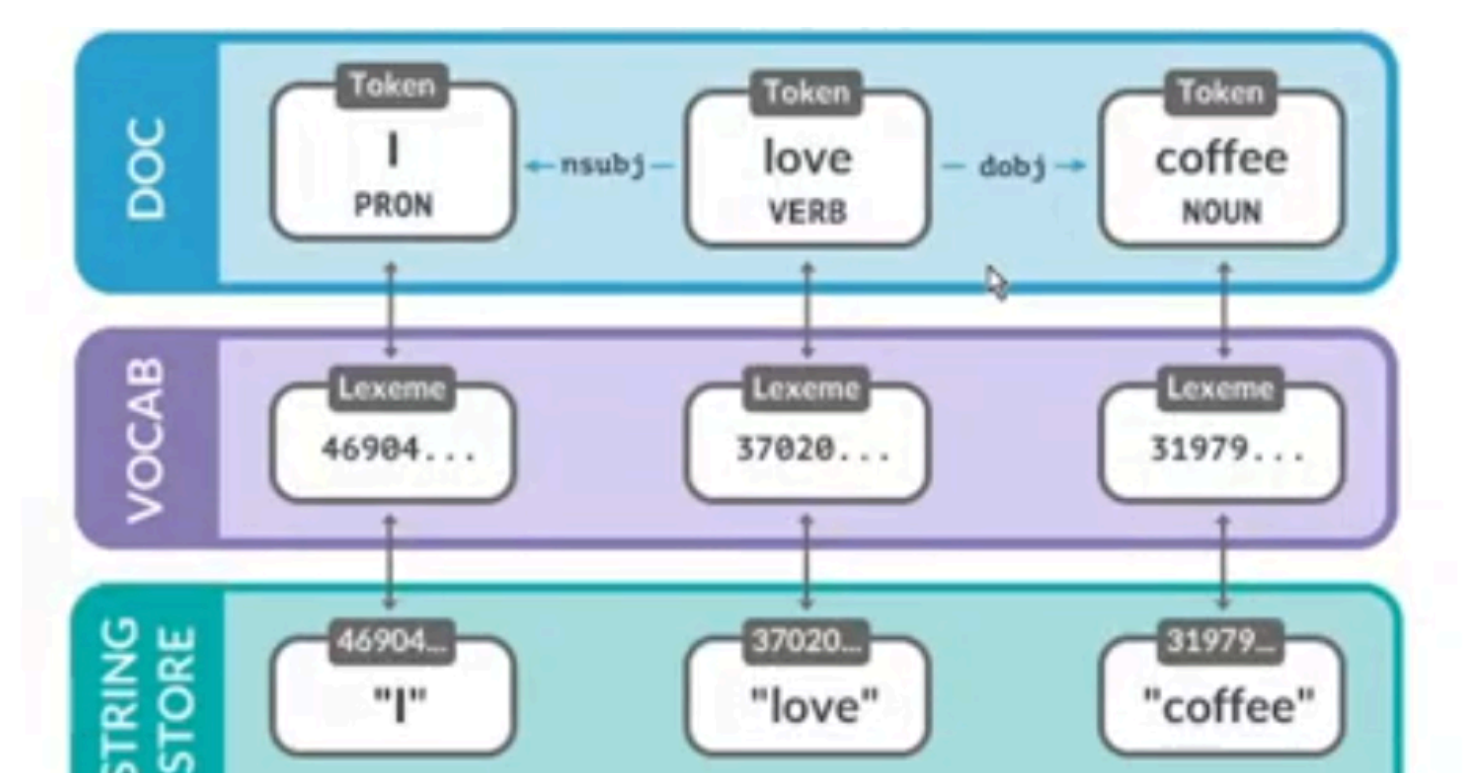


Os componentes do pipeline podem ser adicionados, removidos, e personalizados.

```
1 print(nlp.pipe_names)

['tok2vec', 'morphologizer', 'parser', 'lemmatizer', 'attribute_ruler', 'ner']
```

## Produção de Tokens



### Propriedades do Token:

- **text:** O texto original do token.
- **lemma:** Forma canônica (lema) da palavra, baseada em dicionário.
- **pos:** Classe gramatical no contexto.
- **tag:** Tag detalhada (ex: tempo verbal, número).
- **dep:** Relação sintática na árvore de dependências.
- **enttype:** Tipo de entidade nomeada.
- **isalpha:** Indica se o token contém apenas caracteres alfabéticos.
- **isstop:** Indica se o token é uma "stop word".
- **shape:** Forma do token .
- **vector:** Vetor de embedding .
- **similarity():** Método para calcular similaridade semântica com outro token/doc.
- **etc..**

pratica:



```
1 for token in doc:
2     print(token.text)
```

As  
ações  
do  
Magazine  
Luiza  
S.A.  
,  
Franca  
,  
Brasil  
,  
acumularam  
baixa  
de  
70  
%  
ao  
ano  
.  
Assim  
já  
devolveram  
todos  
os  
ganhos  
do  
período  
da  
pandemia  
.

Exemplos de atributos:

```
1 #atributos
2 print("Tokens:",[token.text for token in doc])
3 print("Stop word:",[token.is_stop for token in doc])
4 print("Alfanumérico",[token.is_alpha for token in doc])
5 print("Maiúculo:",[token.is_upper for token in doc])
6 print("Pontuação:",[token.is_punct for token in doc])
7 print("Número:",[token.like_num for token in doc])
8 print("Sentença inicial:",[token.is_sent_start for token in doc])
9
10
```

Tokens: ['As', 'ações', 'do', 'Magazine', 'Luiza', 'S.A.', ',', 'Franca', ',', 'Brasi  
Stop word [True, False, True, False, False, False, False, False, False, False, False, T  
Alfanumérico [True, True, True, True, True, False, False, True, False, False, T  
Maiúculo: [False, False, False, False, False, False, True, False, False, False, False  
Pontuação: [False, False, False, False, False, True, False, True, False, False, True  
Número: [False, False, False, False, False, False, False, False, False, False, False  
Sentença inicial: [True, False, False, False, False, False, False, False, False, Fals

## Pos-Taggin e Dependências

- POS:** Classifica palavras em categorias .
- Dependências:** Define relações hierárquicas entre palavras.

```
1 for token in doc:
2     print(token.text, "-", token.pos_, "-", token.dep_, "-", token.lemma_, '-', token.shape_)
```

As - DET - det - o - Xx  
ações - NOUN - nsubj - ação - xxxx  
do - ADP - case - de o - xx  
Magazine - PROPN - nmod - Magazine - Xxxxx  
Luiza - PROPN - flat:name - Luiza - Xxxxx  
S.A. - PROPN - flat:name - S.A. - X.X.  
, - PUNCT - punct - , -  
Franca - PROPN - conj - Franca - Xxxxx  
, - PUNCT - punct - , -  
Brasil - PROPN - conj - Brasil - Xxxxx  
, - PUNCT - punct - , -  
acumularam - VERB - ROOT - acumular - xxxx  
baixa - NOUN - obj - baixa - xxxx  
de - ADP - case - de - xx  
70 - NUM - nummod - 70 - dd  
% - SYM - nmod - % - %  
ao - ADP - case - a o - xx  
ano - NOUN - nmod - ano - xxx  
. - PUNCT - punct - . - .  
Assim - ADV - mark - assim - Xxxxx  
já - ADV - advmod - já - xx  
devolveram - VERB - ROOT - devolver - xxxx  
todos - DET - det - todo - xxxx  
os - DET - fixed - o - xx  
ganhos - NOUN - obj - ganho - xxxx  
do - ADP - case - de o - xx  
período - NOUN - nmod - período - xxxx  
da - ADP - case - de o - xx  
pandemia - NOUN - nmod - pandemia - xxxx  
. - PUNCT - punct - . - .

## Listando Entidades Nomeadas:

Identifica e classifica entidades como pessoas, locais ou datas no texto.

```
1 #Tag
2 for ent in doc.ents:
3     print(ent.text, " - ", ent.label_)
```

Magazine Luiza S.A. - ORG  
Franca - LOC  
Brasil - LOC

## Gerenciamento de stop Words

Permite adicionar ou remover palavras irrelevantes

```
1 for token in doc:
2     if token.is_stop:
3         print("Stop word:", token.text)
```

Stop word: As  
Stop word: do  
Stop word: de  
Stop word: Assim  
Stop word: já  
Stop word: todos  
Stop word: os  
Stop word: do  
Stop word: da

## Vocabulário

Conjunto de palavras conhecidas pelo modelo, com atributos como vetores e frequência. As palavras armazenadas tem seu proprio hash, e sao armazenadas como lexemas

```
1 print("hash:", nlp.vocab.strings["dados"])
2 print("hash:", doc.vocab.strings["dados"])
3 print("string:", nlp.vocab.strings["6013848609874238634"])
```

hash: 6013848609874238634  
hash: 6013848609874238634  
string: 15189847787397396957

```
1 #lexemas, palavras que não eestão representadas com uma dependencia de contexto
2 lex = nlp.vocab["dados"]
3 print(lex.text, "-", lex.orth, "-", lex.is_alpha, "-", lex.is_lower)
```

dados - 6013848609874238634 - True - True

## Busca de Similaridade

Compara vetores de palavras/documentos para medir proximidade semântica, com base na semelhança cosenoidal entre os vetores.

```
1 doc1 = nlp("ele viaja regularmente de carro")
2 doc2 = nlp("ele viaja, nunca de avião")
3
4 print(doc1.similarity(doc2))
```

0.9020987749099731

```
1 doc3 = nlp("é pave ou pacomé")
2
3 tokenA = doc3[1]
4 tokenB = doc3[3]
5
6 print(tokenA.similarity(tokenB))
7
```

0.0

## Busca de expressões com matching

Encontra padrões específicos usando regras baseadas em tokens para identificar sequências específicas. Parece um pouco com expressões regulares.

```
1 from spacy.matcher import Matcher
2
3 doc5 = nlp("Você pode ligar para (63) - 984021557 ou (41) 81322096")
4
5 #inicializando matcher
6 matcher = Matcher(nlp.vocab)
7 padrao = [{"ORTH": "("},
8           {"SHAPE": "dd"},
9           {"ORTH": ")"},
10          {"ORTH": "-", "OP": "?"},
11          {"IS_DIGIT": True}]
12
13 matcher.add("telefone", [padrao])
14 matches = matcher(doc5)
15
16 for match_id, start, end in matches:
17     print(doc5[start:end])
```

(63) - 984021557  
(41) 81322096

## Visualização com Displacy

Gera gráficos interativos para dependências sintáticas e entidades.

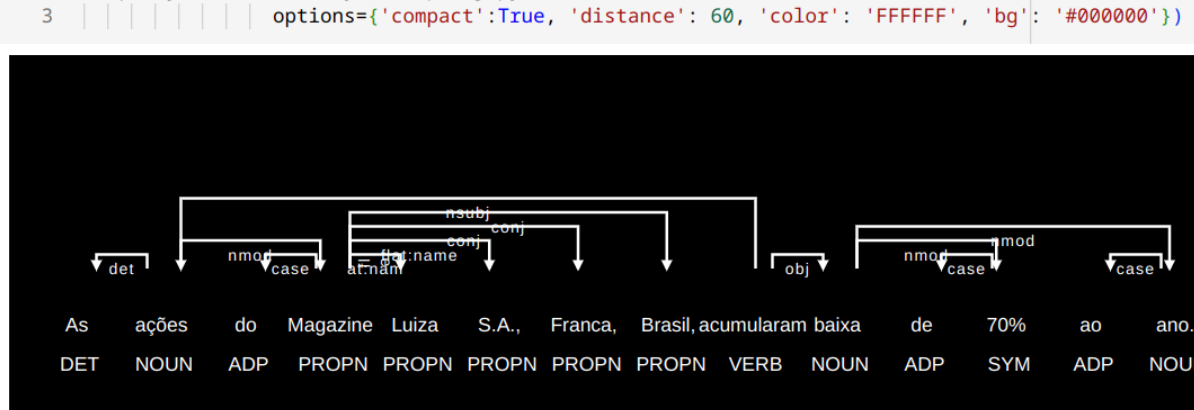
visualização de entidades:

```
1 from spacy import displacy
2
3 displacy.render(doc,style="ent", jupyter=True)
```

As ações do Magazine Luiza S.A. Franca LOC Brasil LOC , acumularam baixa de 70% ao ano. Assim já devolveram todos os ganhos do período da pandemia.

visualizando dependências:

```
1 doc.user_data["title"] = "Exemplo"
2 displacy.render(doc,style="dep", jupyter=True,
3                 options={'compact':True, 'distance': 60, 'color': 'FFFFFF', 'bg': '#000000'})
```



## Gerenciando Pipelines

Personalização de fluxos de processamento

```
[ ] 1 print('Pipeline Normal:', nlp.pipe_names)

Pipeline Normal: ['morphologizer', 'parser', 'lemmatizer', 'attribute_ruler', 'ner']

[ ] 1 nlp.remove_pipe("tok2vec")

('tok2vec', <spacy.pipeline.tok2vec.Tok2Vec at 0x7de483a2fc50>)

[ ] 1 print('Pipeline sem tokenizer:', nlp.pipe_names)

Pipeline sem tokenizer: ['morphologizer', 'parser', 'lemmatizer', 'attribute_ruler', 'ner']

[ ] 1 nlp.add_pipe("tok2vec", after='morphologizer')
2

<spacy.pipeline.tok2vec.Tok2Vec at 0x7de483970470>
```

## NLP com NLTK

Biblioteca **Natural Language Toolkit**. No NLTK, o pipeline de processamento não vem pronto como no spacy, as etapas tem de ser feitas manualmete ou definidas em um pipeline personalizado.

## Teste com a música "O papa perdoa tom zé, de Tom Zé:"

Produção de tokens

Divide textos em unidades menores usando métodos como word\_tokenize.

