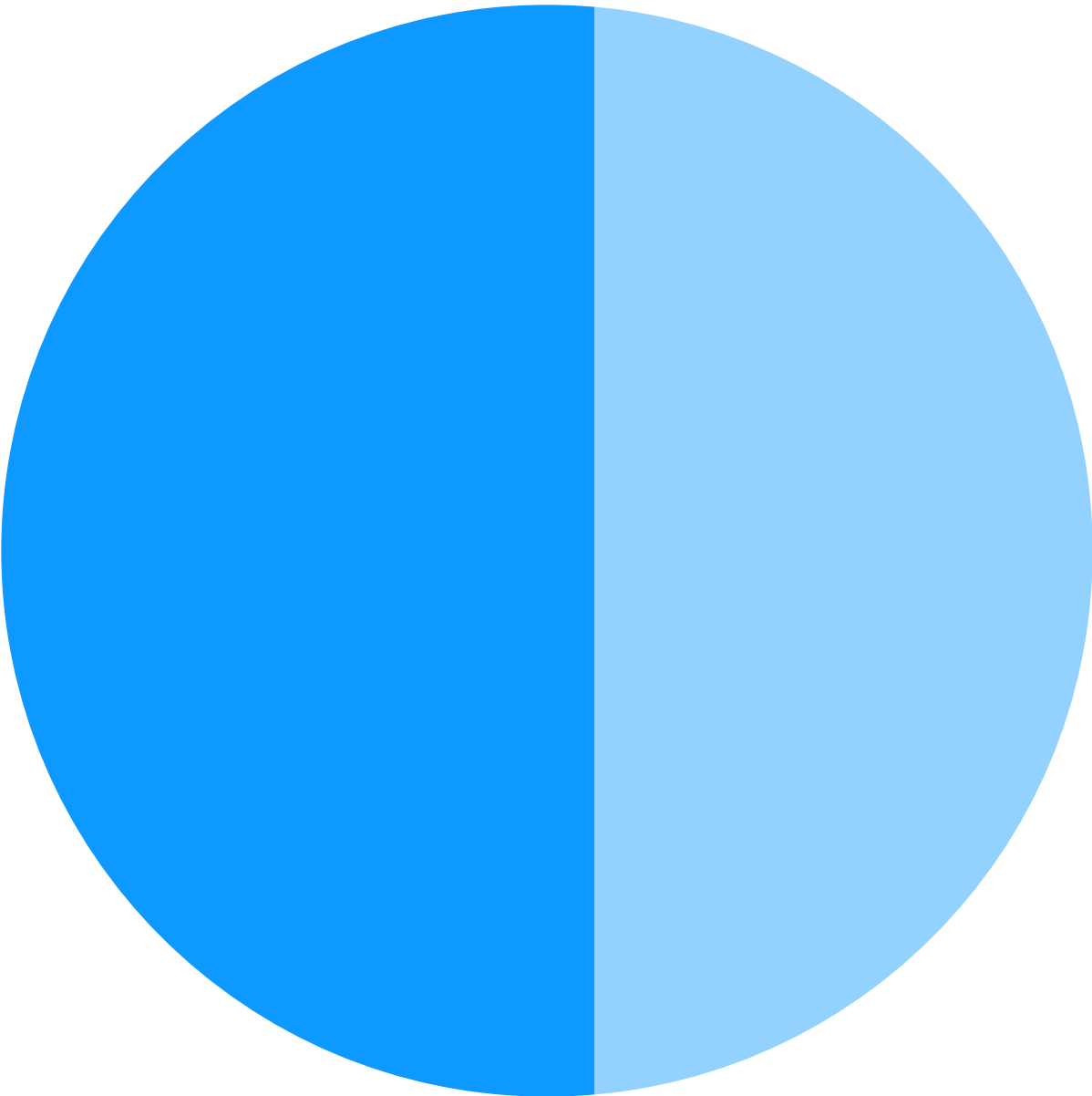


# SPRINT 1

Mockingjay [Brastel Co., Ltd.](#)



# O que fizemos?

Pr-e-processamento e análise exploratória de dados

O que os dados têm a nos dizer?

# Pré-processamento

Desenvolver uma pipeline de pré-processamento de texto para normalizar dados textuais, facilitando análises posteriores.

# Pré-processamento

- Remoção de sequências de escape
- Conversão para minúsculo
- Remoção de números
- Remoção de pontuação
- Tokenização
- Remoção de Stopwords
- Lematização
- Stemming

# Testes

As funções dos testes unitários para serem testados e seus resultados esperados, estão logo abaixo:

- `TestRemoveEscapeSequences` (remove sequências de escape): "Bom dia, quanto da 150.000 yenes no Brasil hj ?" -> "Bom dia, quanto da 150.000 yenes no Brasil hj ?"

- `TestLowercase` (sentença em caixa baixa):

"Bom dia, quanto da 150.000 yenes no Brasil hj ?" -> "bom dia, quanto da 150.000 yenes no brasil hj ?"

- `TestRemoveNumbers` (Remoção de números na sentença):

"Bom dia, quanto da 150.000 yenes no Brasil hj ?" -> "Bom dia, quanto da . yenes no Brasil hj ?"

- `TestRemovePunctuation` (Remoção de sinais de pontuação):

"Bom dia, quanto da 150.000 yenes no Brasil hj ?" -> "Bom dia quanto da 150000 yenes no Brasil hj"

- `TestTokenization` (Segmentação do texto em unidades menores chamadas tokens) :

"Bom dia, quanto da 150.000 yenes no Brasil hj ?" -> ['Bom', 'dia', ',', 'quanto', 'da', '150.000', 'yenes', 'no', 'Brasil', 'hj', '?']

- `TestRemoveStopwords` (Palavras que são removidas do texto para reduzir a dimensionalidade e focar nas palavras de maior relevância):

['Bom', 'dia', ',', 'quanto', 'da', '150.000', 'yenes', 'no', 'Brasil', 'hj', '?'] -> ['Bom', 'dia', ',', 'quanto', '150.000', 'yenes', 'Brasil', 'hj', '?']

- `TestStemming` (Conversão de palavras para seu radical base):

['Bom', 'dia', ',', 'quanto', 'da', '150.000', 'yenes', 'no', 'Brasil', 'hj', '?'] -> ['Bom', 'dia', ',', 'quant', 'da', '150.000', 'yen', 'no', 'Brasil', 'hj', '?']

- `TestLemmatization` (Conversão de palavras para seus lemas):

['Bom', 'dia', ',', 'quanto', 'da', '150.000', 'yenes', 'no', 'Brasil', 'hj', '?'] -> ['Bom', 'dia', ',', 'quanto', 'de o', '150.000', 'yene', 'em o', 'Brasil', 'hj', '?']

# Análise exploratória

Analisar interações de clientes via chat para identificar padrões e preparar os dados para desenvolvimento de uma solução de chatbot.

# Quais perguntas fazer?

- Distribuição dos tipos de perguntas (dúvidas, problemas, solicitações).
- Complexidade das perguntas e como elas variam na base de dados.
- Respostas mais frequentemente fornecidas pelos atendedores.
- Padrões de respostas repetidas ou padronizadas.
- Correlação entre tipo de pergunta e necessidade de escalonamento.

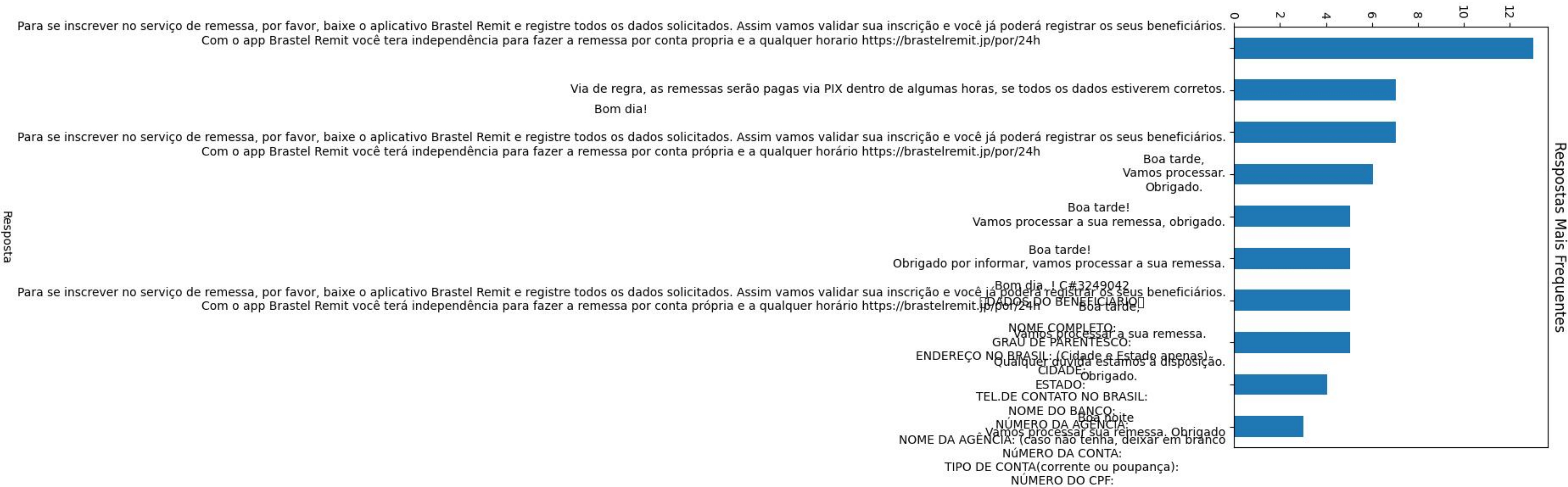
# Perguntas







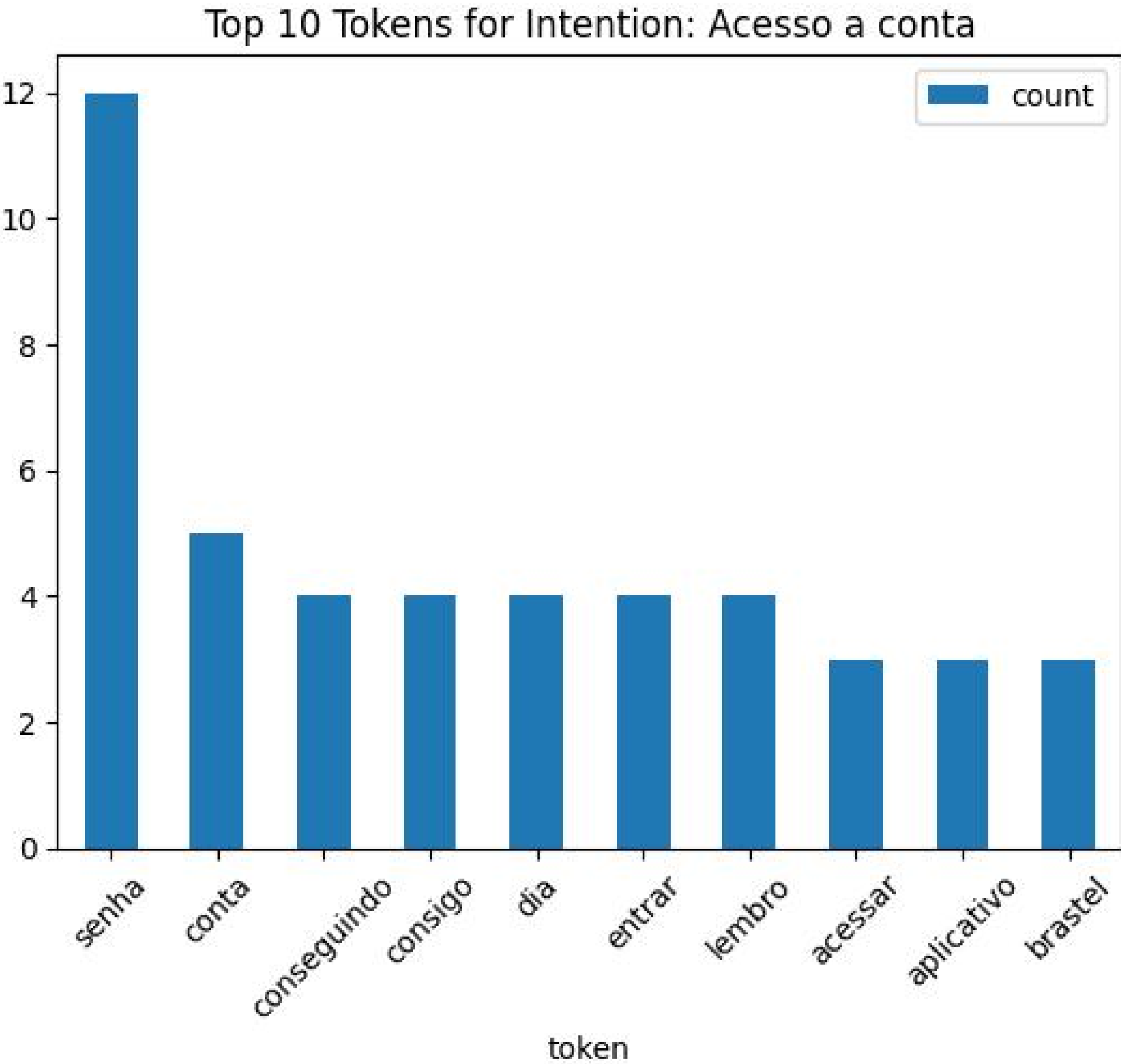
# Respostas frequentes



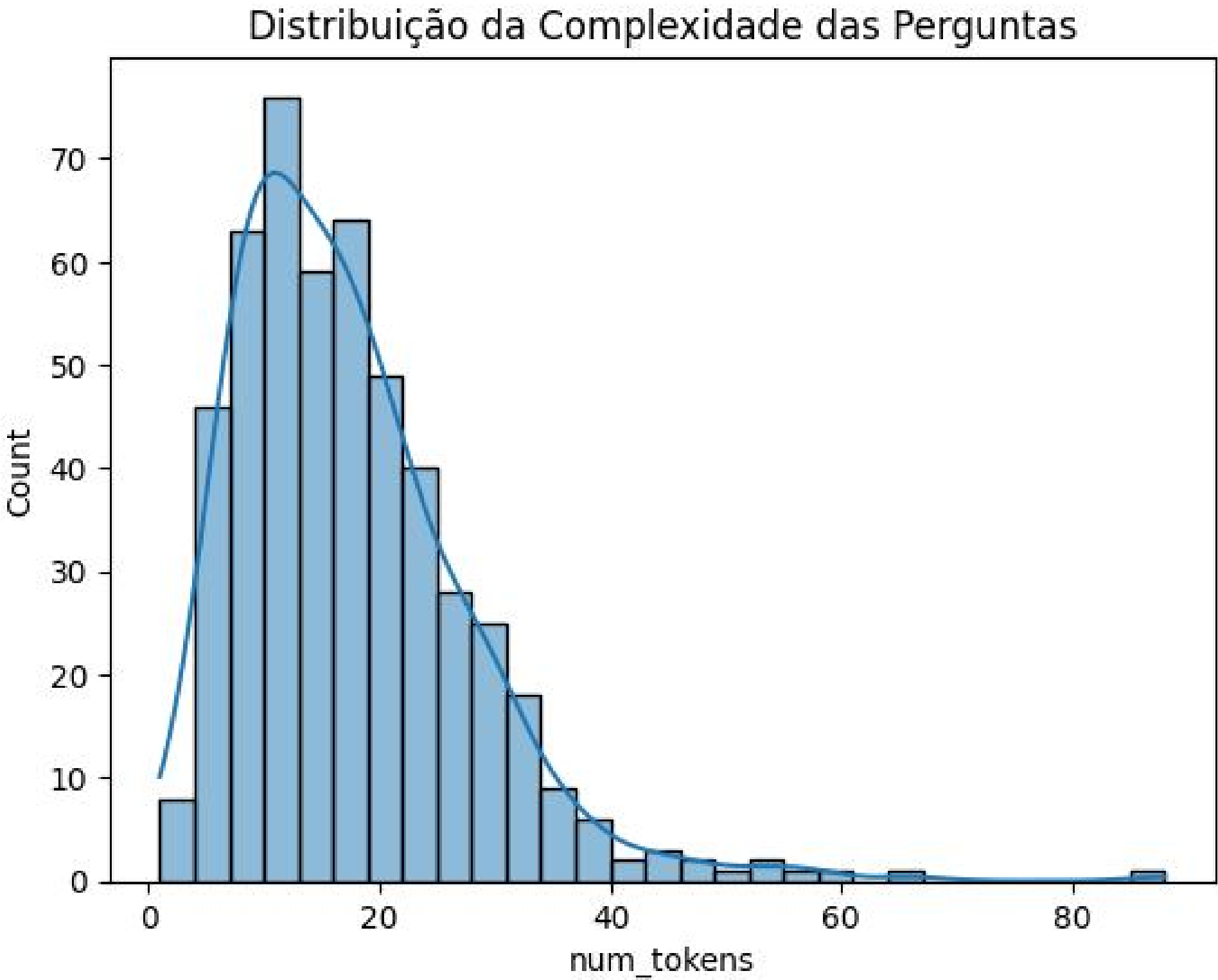
# Respostas



# Respostas



# Tokens por pergunta





# Próximos passos

- Implementação de um modelo baseline
- Implementação de modelo com rede neural
- Continuação do artigo
- Melhoria da pipeline

7 de agosto de 2024

Instituto de Tecnologia e Liderança

# Nosso time

Grupo Mockingjay



Allan



Elias



Giovana



Rafael



Gábrio



Cristiane



Melyssa

7 de agosto de 2024

Instituto de Tecnologia e Liderança

# Obrigado

