



Oficina: ChatGPT com Dados Privados usando Elasticsearch

Alex Salgado
Developer Advocate @ Elastic

- | | |
|--|---|
|  @alexsalgadoprof |  @alexsalgadoprof |
|  salgado |  /in/alex-salgado/ |





Alex Salgado
Developer Advocate
LATAM

@alexsalgadoprof

salgado

@alexsalgadoprof

/in/alex-salgado/

- **Mestre** em Ciência da Computação pela UFF (Games)
- **MBA** UFF
- **PhD Candidate UFF: Robótica/Visão Computacional**

- **+ 20 anos** de experiência na área de desenvolvimento de software
- Ocupei diversos cargos, trabalhando em **startups**, pequenas e grandes empresas como Oracle, CSN, BRQ/IBM, **Chemtech/Siemens (9 anos)**.
- **8 anos** como professor universitário





Preocupações em torno da IA Generativa.

80%

Dados mundiais são não-estruturados

[KPMG Generative AI Survey](#)

[The Prompt: Generative AI survey | Google Cloud Blog](#)



Preocupações em torno da IA Generativa.

50%

60%

64%

DESENVOLVEDORES

Vê escassez das habilidades necessárias para a IA Generativa em sua organização.

EXECUTIVOS

Sentem um alto grau de urgência para adotar a IA generativa.

EXECUTIVOS

Acreditam que estão a 1-2 anos de implementar a IA Generativa.



Three solutions powered by one stack

3 solutions



Enterprise Search



Observability



Security

Powered by
the Elastic Stack

Kibana

Elasticsearch

Agent

Beats

Logstash

Deployed
anywhere



Elastic Cloud



Elastic Cloud
Enterprise

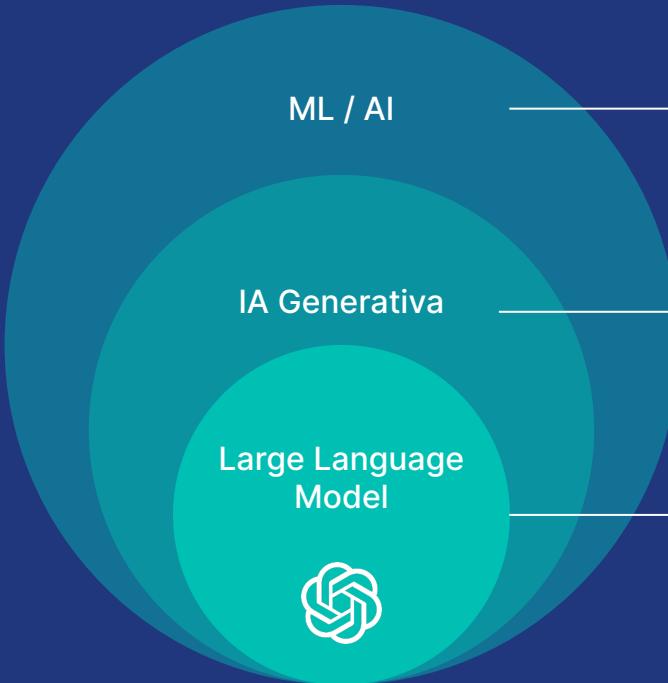


Elastic Cloud
on Kubernetes

SaaS

Orchestration

Conceitos básicos de ML, IA Generativa e LLMs



O que?

Algoritmos programados para aprender o comportamento dos dados e fazer previsões

Casos de uso

Detecção de anomalias, forecasting, reconhecimento de imagem, PLN

Algoritmos programados para criar novos dados

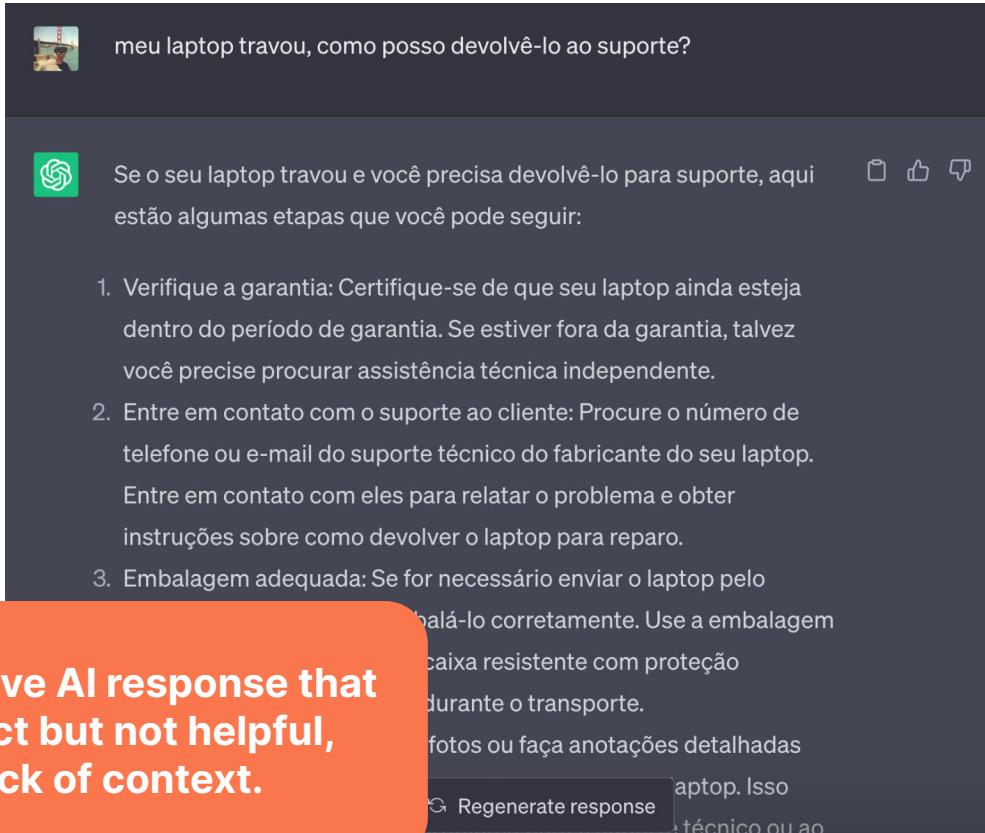
Chatbots, geradores de texto, imagem e música

Algoritmos (Deep Learning) treinados com grandes volumes de dados e programados para criar novos dados

Chatbots, geradores de texto, tradutores, geradores de código, aplicativos de pergunta e resposta

Example

Getting Some Help



meu laptop travou, como posso devolvê-lo ao suporte?

Se o seu laptop travou e você precisa devolvê-lo para suporte, aqui estão algumas etapas que você pode seguir:

1. Verifique a garantia: Certifique-se de que seu laptop ainda esteja dentro do período de garantia. Se estiver fora da garantia, talvez você precise procurar assistência técnica independente.
2. Entre em contato com o suporte ao cliente: Procure o número de telefone ou e-mail do suporte técnico do fabricante do seu laptop. Entre em contato com eles para relatar o problema e obter instruções sobre como devolver o laptop para reparo.
3. Embalagem adequada: Se for necessário enviar o laptop pelo correio, embalá-lo corretamente. Use a embalagem certa e uma caixa resistente com proteção adicional durante o transporte.

A generative AI response that was correct but not helpful, due to a lack of context.

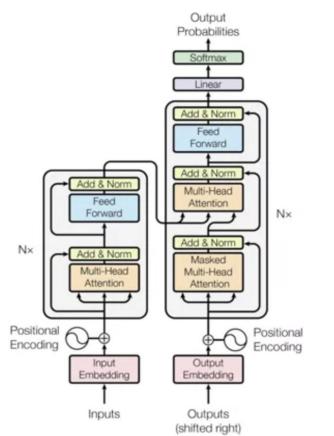
G Regenerate response

aptop. Isso é útil para um técnico ou ao

Generative AI

É um tipo de inteligência artificial capaz de gerar texto, imagens ou outros tipos de mídia em resposta a estímulos, aprendendo os padrões e estruturas de dados de treinamento e gerando novos dados com características semelhantes.

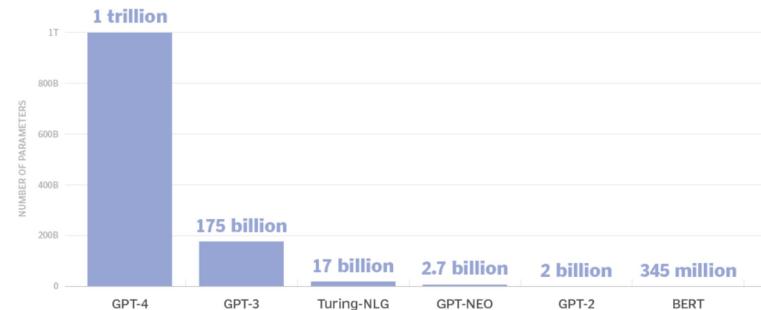
Transformers



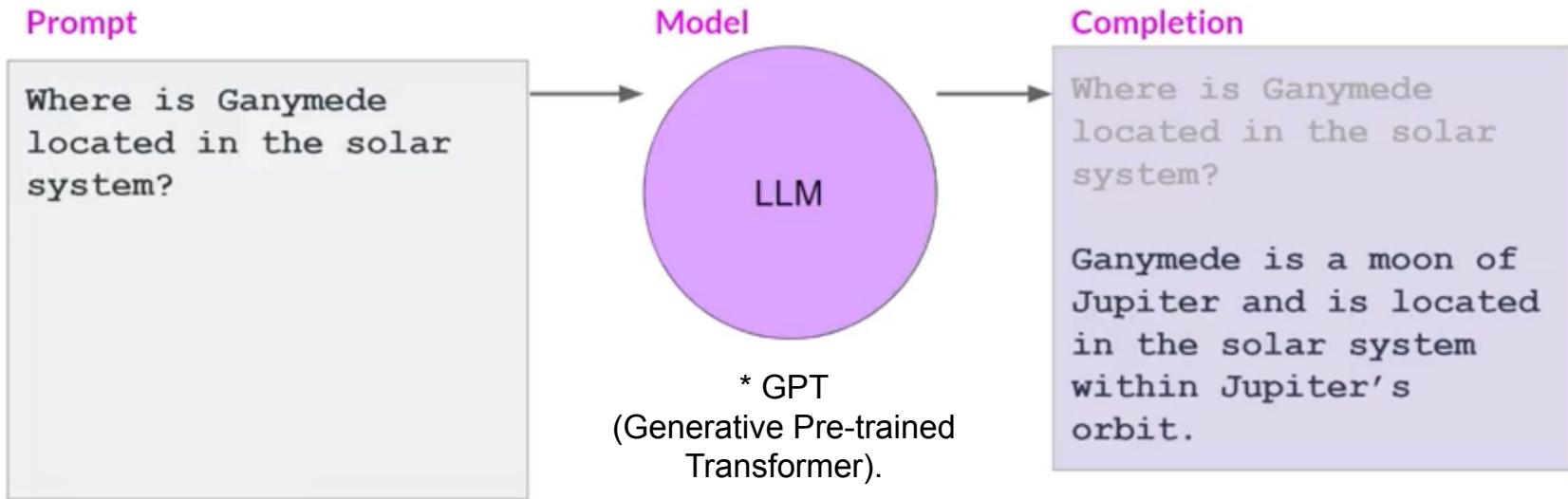
LLM

Large Language Model (LLM) é um tipo de algoritmo de inteligência artificial (IA) que utiliza técnicas de aprendizado profundo e conjuntos de dados enormes para compreender, resumir, gerar e prever novo conteúdo.

Parameters of transformer-based language models



Prompts and completions



Context window

- typically a few 1000 words.

Elasticsearch: You Know, for Search



How fast should my internet be?

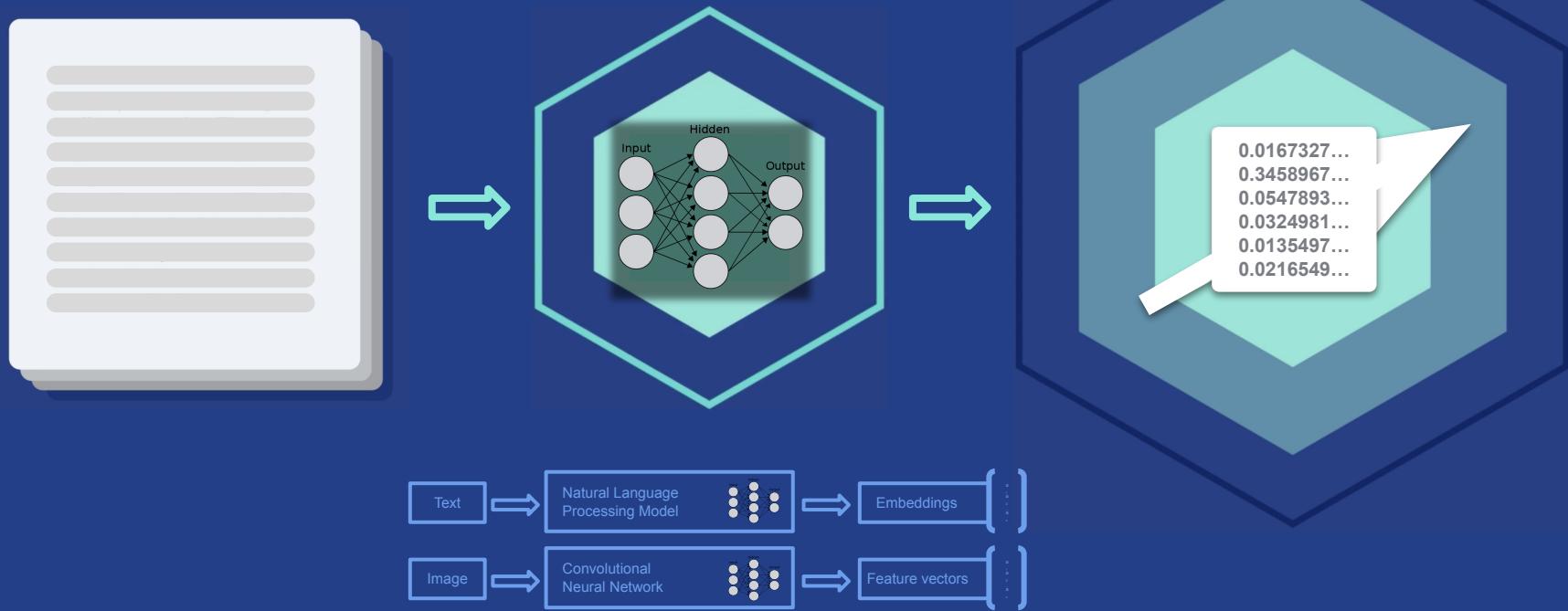
In order to stream from our service you will need a high quality connection. The required connection speed for using the service will vary depending on the quality of

you wish
vice. For

recommend at least...

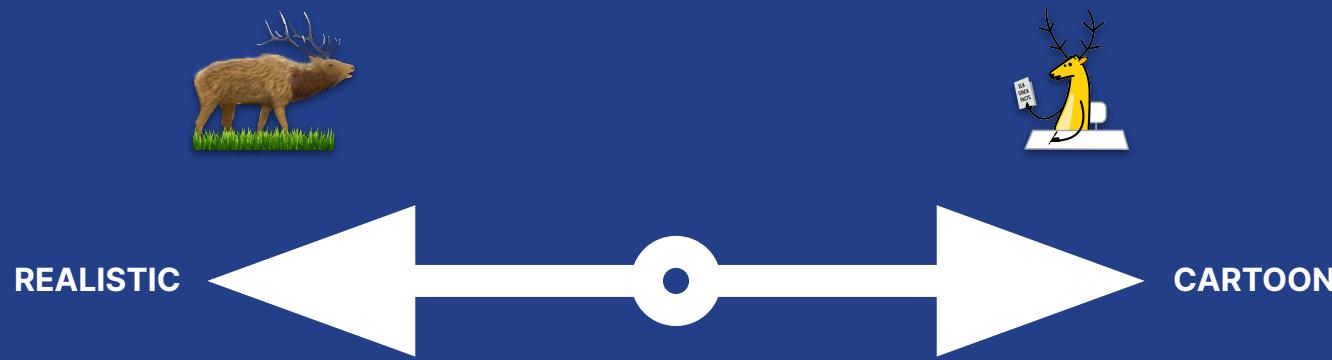
O que é similaridade de vetores?

Converte dados em representações vetoriais onde as distâncias representam similaridade.

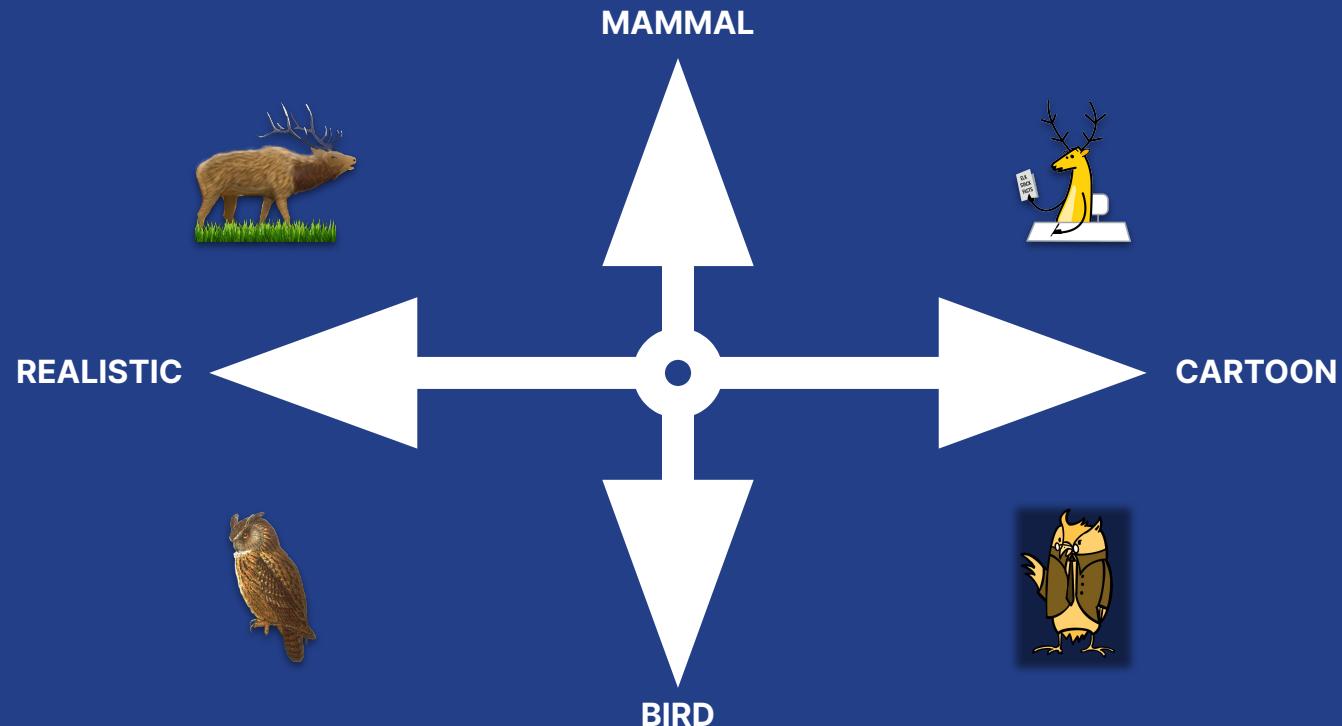


Embedding é a representação do dado no espaço vetorial

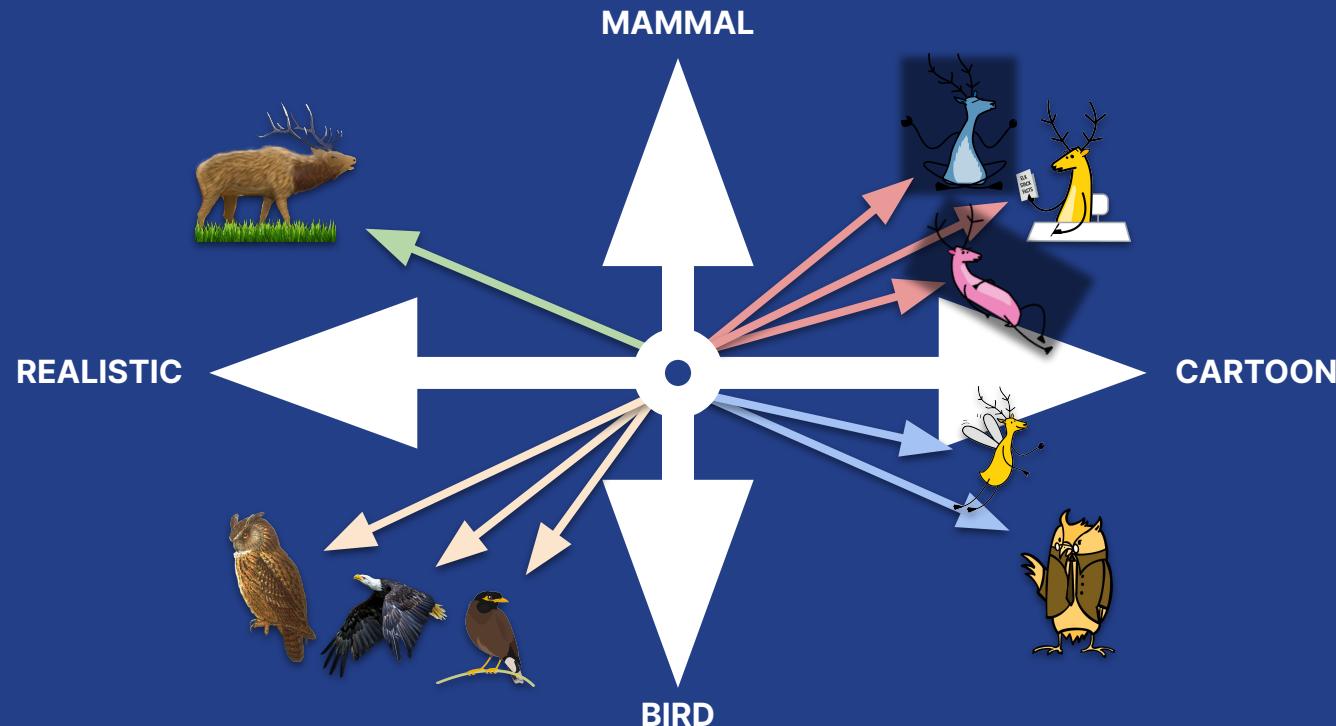
Example: 1-dimensional vector



Múltiplas dimensões representam diferentes aspectos dos dados.

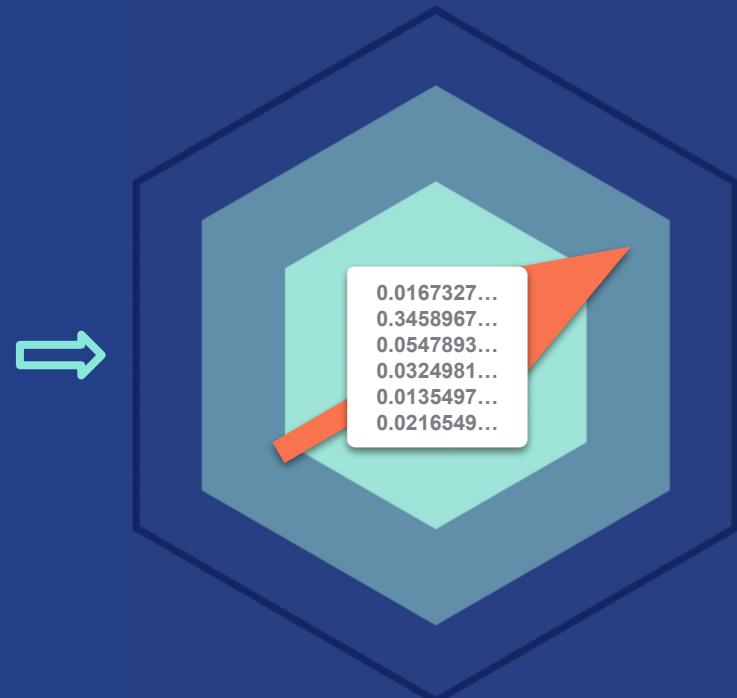
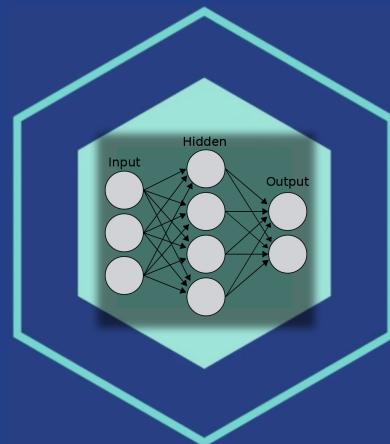


No “Espaço vetorial”, dados similares são agrupados juntos.

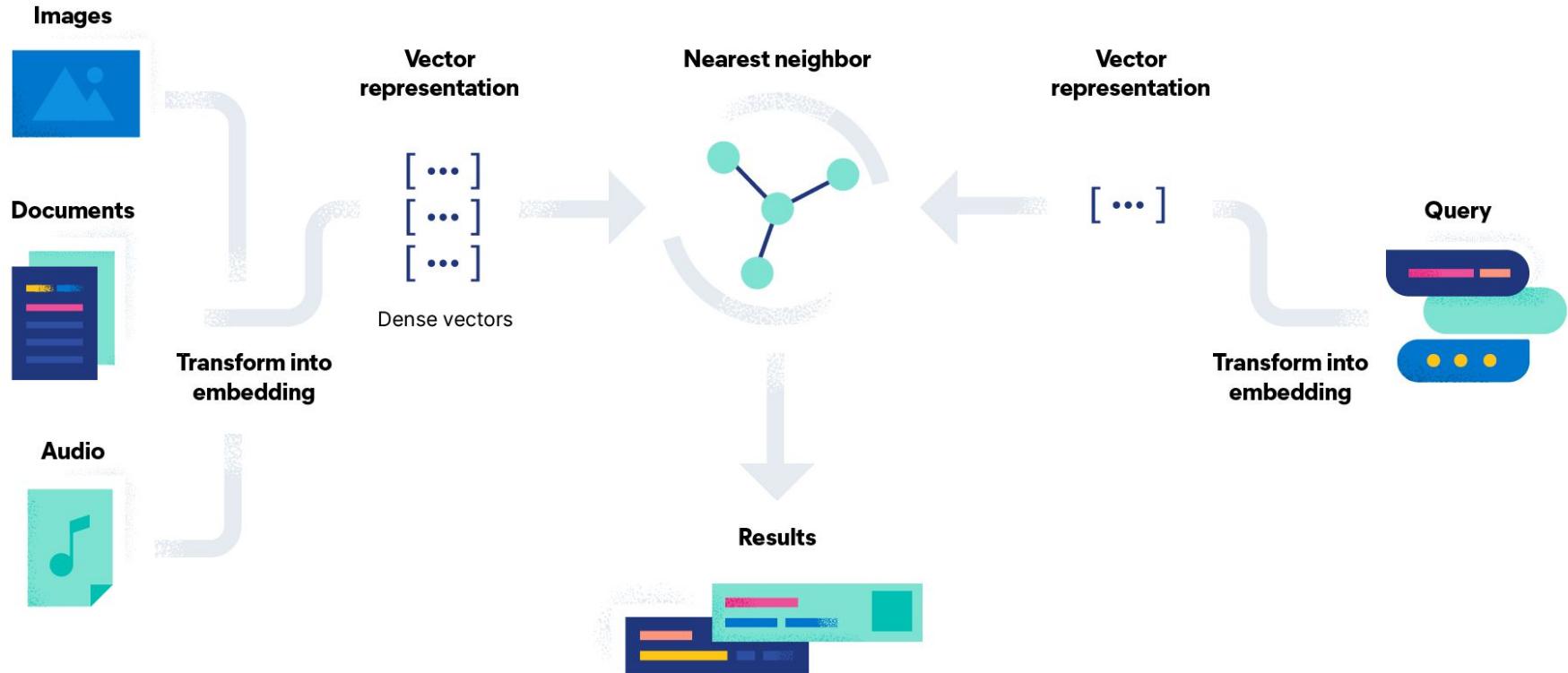


O que é similaridade de vetores?

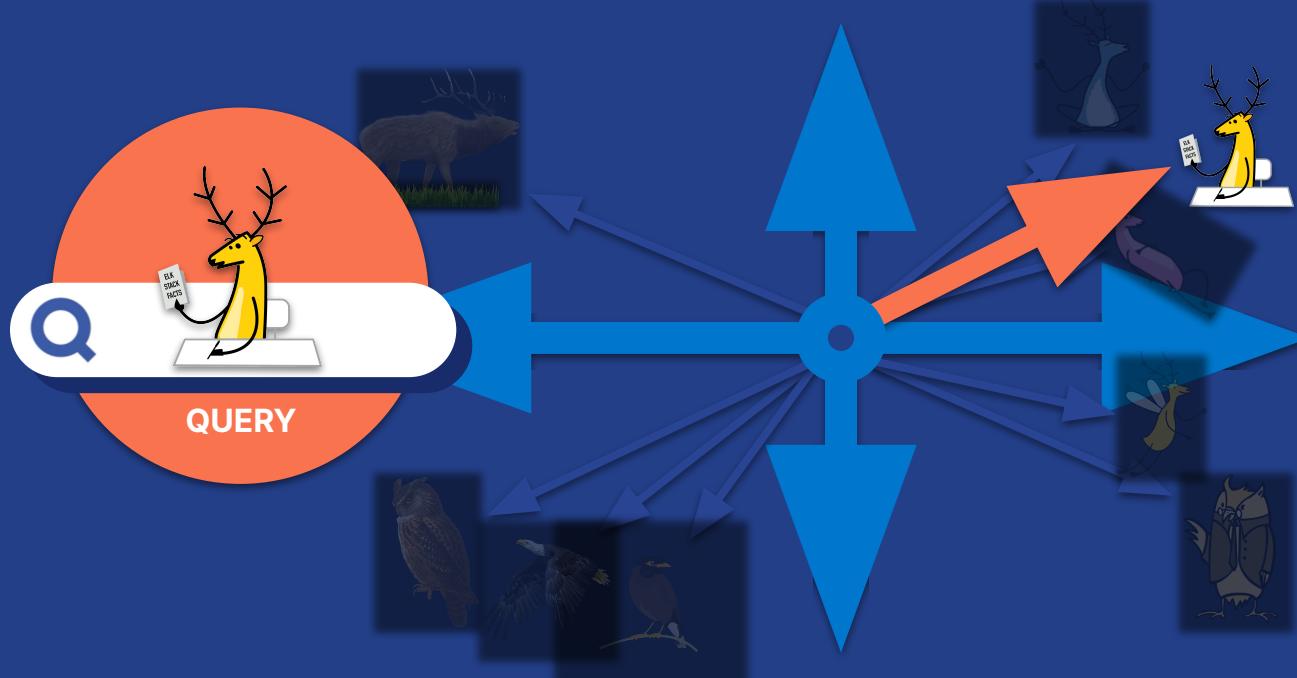
A consulta também é vetorizada (embeddings)



Consultas também são vetorizadas



Busca vetorial ranqueia os objetos baseados em similaridades (relevância) com a consulta



Relevance	Result
Query	
1	
2	
3	
4	
5	

Prática

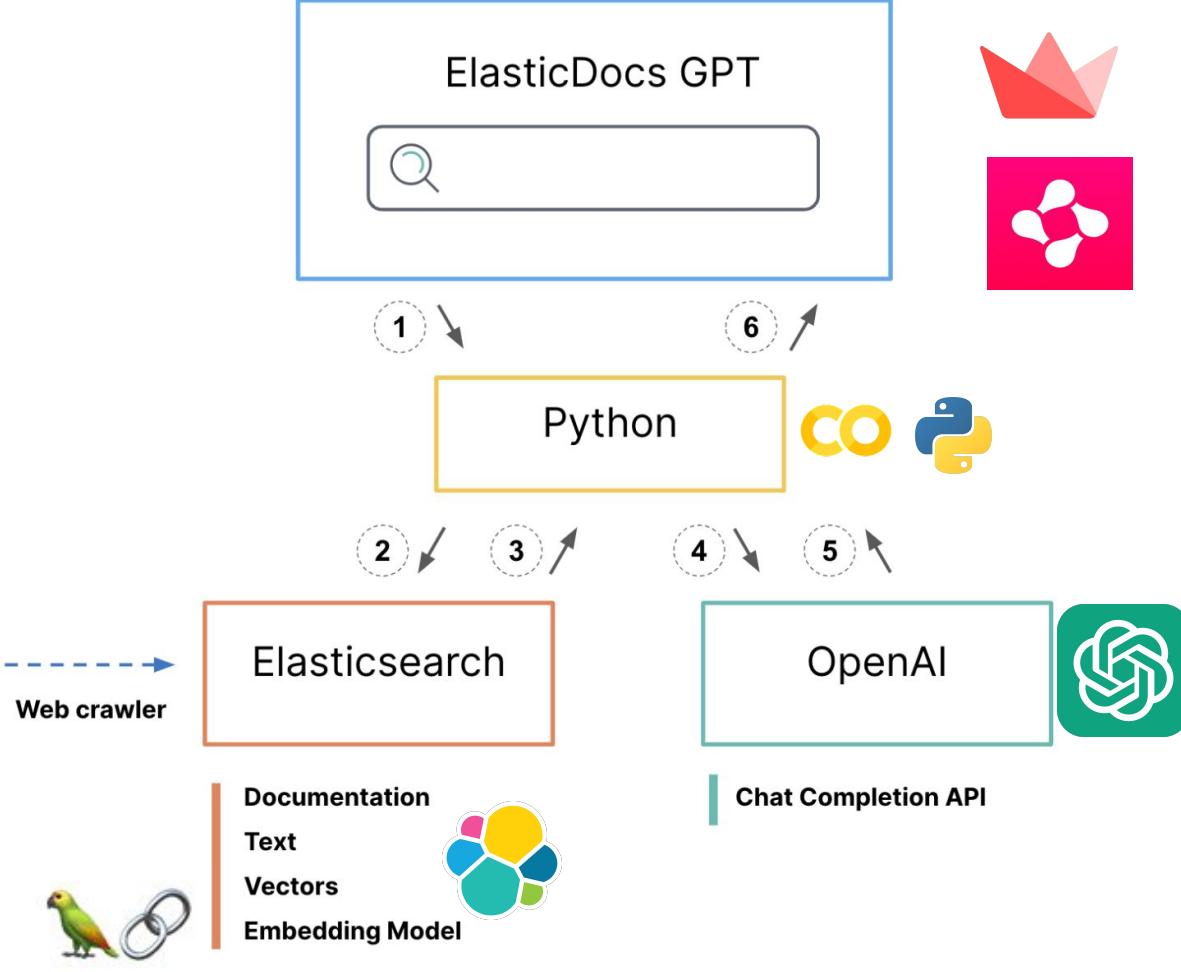
Vamos construir nossa aplicação
LLM para chat?

How to use ChatGPT with Elasticsearch

Retrieval Augmented
Generation (RAG)



this can be replaced
by your own data source!



Documentation
Text
Vectors
Embedding Model



Demo

https://github.com/salgado/ElasticDocs_GPT

<https://www.elastic.co/blog/chatgpt-elasticsearch-openai-meets-private-data>

Conecte-se com a comunidade Elastic no Brasil

Encontre um grupo



Veja os eventos locais



• <https://community.elastic.co/>

<https://ela.st/brvirtual>

Obrigado

Alex Salgado



@alexsalgadoprof



salgado



@alexsalgadoprof



/in/alex-salgado/

Link workshop-mackenzie

<https://github.com/prof-alex/workshop-mackenzie>

Blogs

<https://www.elastic.co/search-labs/chatgpt-elasticsearch-openai-meets-private-data>

Blogs de Alex

<https://www.elastic.co/pt/blog/author/alex-salgado>

<https://dev.to/salgado>

Para casa

Recriar o seu próprio ChatGPT com dados privados

Usando a sua conta trial:

- Replicar o passo a passo desse workshop usando como referência o blog :
<https://www.elastic.co/search-labs/chatgpt-elasticsearch-openai-meets-private-data>
- Definir uma página de crawler que o chatgpt não saiba responder
- Mudar o título da aplicação

Durante o processo, havendo qualquer dúvida, pergunte a comunidade usando estes canais:

Aqui é nosso discuss: <https://discuss.elastic.co/>

Aqui nosso slack:

https://join.slack.com/t/elasticstack/shared_invite/zt-24evvtls4-08ZzDILtGQNDQduk15Mw0w

Para casa

Recriar o seu próprio ChatGPT com dados privados

Para quem conseguir terminar:

- Crie um vídeo de 1 minuto fazendo uma busca na sua aplicação
- Publique no Linkedin com sua percepção
- Coloque a hashtag #WorkShopElastic #Elastic #Mackenzie #ChatGPT
- Marque **eu** e o professor **Rogério**

Aguarde a surpresa !!!!