

DATASTAX

Concorra a kits da DataStax

A melhor solução para Vector Search e RAG



Workshop: Criando Apps de IA Generativa com NodeJS e NextJS

Samuel Matioli
Solution Engineer at DataStax



» Agenda

- Aplicações com IA Generativa
- Vector Search & DataStax Astra
- O que é RAG?
- Ex1: Implementando RAG *na mão*
- Ex2: Simplificando com LangchainJS
- Ex3: Usando Functions

Ganhe um Badge



Envie o seu DB ID e um screenshot da aplicação:

samuel.matioli@datastax.com

para receber seu Badge!

Sorteio de brindes

Cadastre-se aqui para participar do sorteio.

bit.ly/datastax-tdcsc



Housekeeping - tdclover / #tdclover

Versões: NodeJS 18.

Clonar: **<https://github.com/smatiolids/workshop-genai-js>**

Criar conta em:

OpenAI (<https://platform.openai.com/>)

- Vamos usar os embeddings e o LLM para gerar conteúdo

DataStax Astra (astra.datastax.com) – **Conta GRATUITA!**

- Armazenar dados proprietários e memória



Criando o Vector DB

Create Database

ESC



Select a deployment type

Serverless (Vector)

Recommended

An all-in-one database solution, optimized for Vector and Generative AI workloads

Serverless (Non-Vector)

A more traditional database solution without any of our new vector capabilities

Configuration

Database name*

db_name

Give it a memorable name – this can't be changed later.

Provider*



Google Cloud



Region*



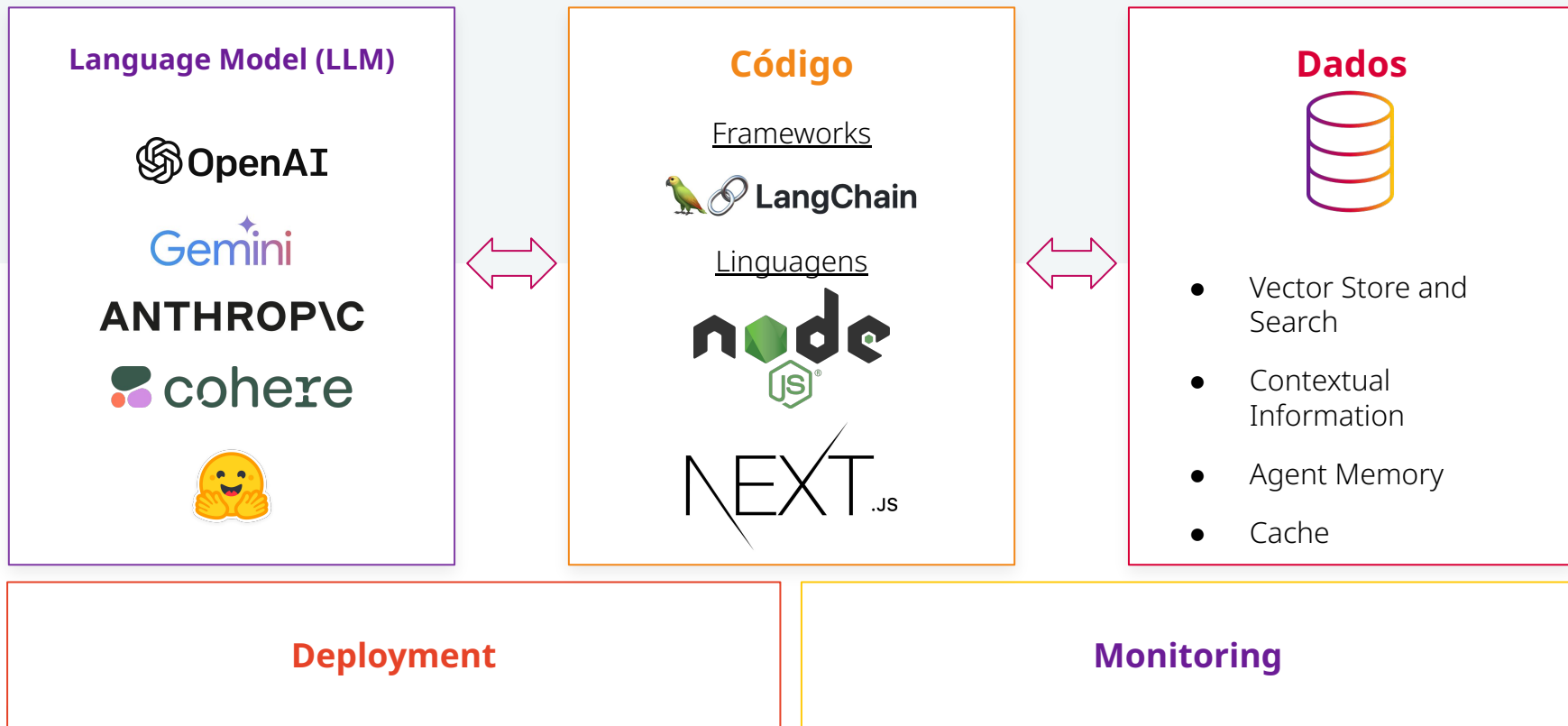
us-east1



Cancel

Create Database

› Como criar aplicações GenAI



› Dificuldades com o RAG em produção

O apps com GenAI "quebram" muito

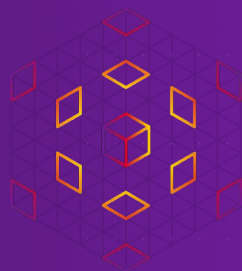
- Muitas dependências
- APIs externas **sem garantias**
- Novas versões de APIs externas **sem aviso**

Escala e resiliência:

- Dados em tempo real
- Milhares de usuários
- Milhões de interações
- Muitos dados...

Necessidade de evolução:

- Agilidade para incrementar e refinar os casos de uso
- Aproveitar novos modelos
- Novas origens de dados
- Ferramentas para a área de negócio



Exercício 1



› Iniciando

Passos:

- Conta no Astra
- Clonar: **<https://github.com/smatiolids/workshop-genai-js>**
- Copiar o .env.example para .env.local
- Atualizar as variáveis
 - OPENAI_API_KEY
 - ASTRA_DB_API_ENDPOINT
 - ASTRA_DB_APPLICATION_TOKEN

› Ex 1: RAG + Vectorize

Passos:

- Acessar localhost:3000/v1/chat e fazer perguntas sobre financiamento imobiliário
- Configurar o \$vectorize
- Criar collection "**real_estate_financing**" no Astra
- Carregar documento em localhost:3000/v1/upload



API Key

Add Integration

ESC



OpenAI

Configure OpenAI to generate embeddings for Astra

API key name*

Name cannot be changed later

OpenAI API key*

Learn more about [adding embedding providers](#) to use in Astra

Add databases to API key scope*



Select up to 10 vector databases to add to key scope.

By adding this integration, you agree to the [DataStax Preview Terms](#).

Cancel

Add integration



Criando a collection

Create collection

ESC



Collection name*

real_estate_financing

☒ Vector-enabled collection ⓘ

Embedding generation method*

OpenAI

Select a configured [embedding provider integration](#) or generate your own embeddings

Api key*

dev_key

Select an api key

Embedding model*

text-embedding-3-small

Select a model to use to generate embeddings

Dimensions*

1536

Vector length in your dataset

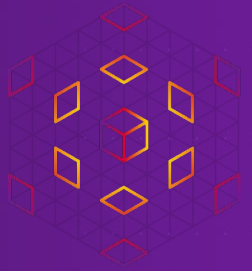
Similarity metric*

Cosine

Method used to calculate vector similarities

› Carregar dados para o Astra

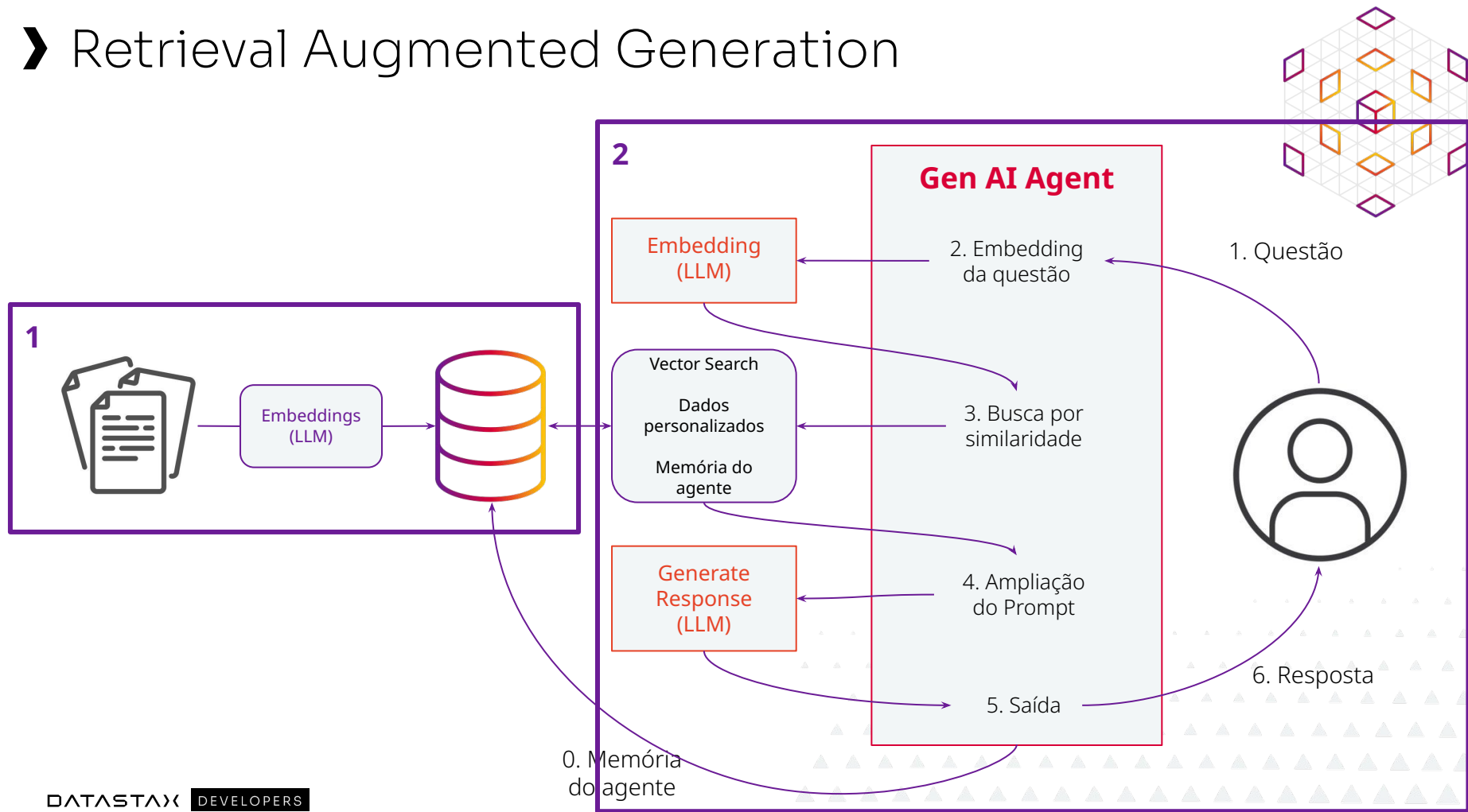
- Criar API Key
 - Vincular DB ao escopo da chave
- Criar collection "**real_estate_financing**"
- Vincular modelo OpenAI "**text-embedding-3-small**"
 - Dimensões: **1536**
 - Similarity Metric: **Cosine**



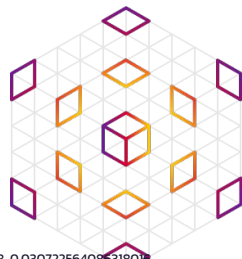
Retrieval Augmented Generation (ou RAG)



➤ Retrieval Augmented Generation



» O que são embeddings



To create a security token that can be used to log into a database, select Token Management from the User Management menu. Then, choose an appropriate role for the user, and click the Generate Token button. Copy the token details to a safe place, as the secret that is shown can never be reproduced in the Astra console for security reasons.

text-embedding-
ada-002

```
[ -0.029334254562854767, 0.06338247656822205, 0.03711941838264465, 0.06770425289869308, 0.030722564086318016,
-0.03855780512094498, 0.05715630576014519, 0.011225797235965729, 0.03209076076745987, 0.018565965816378593,
0.005725057329982519, 0.003278332995250821, 0.019661232829093933, 0.00848303231916428, 0.01011530589312315,
-0.06865981221199036, -0.02742725796997547, 0.004272086545825005, 0.006464742589741945, 0.033381473273038864,
-0.06456394493579865, -0.00168662762735039, -0.02538292482495308, -0.013529211282730103, 0.00674532540125362,
-0.09090574085712433, -0.004533097147941589, -0.011557350866496563, -0.017933078110218048, -0.013565625995397568,
0.037340205162763596, -0.013345368206501007, -0.046318963170051575, 0.018221553415060043, -0.0305146872997283594,
0.06087908893823624, -0.015947293490171432, -0.004384475760161877, 0.011510275304317474, 0.03181201219558716,
0.0004961538943462074, -0.00960769411764431, -0.0026698557194322348, -0.02102011534202099, -0.0244514361227417,
0.018808122724294662, -0.04526928439736366, -0.03507508337497711, 0.011936023831367493, -0.04246317967772484,
-0.04538712650537491, 0.0030571885872632265, -0.02645616978460906, 0.004339783918112516, -0.004989228677004576,
-0.0019529943820089102, -0.015389597043395042, -0.008066490292549133, -0.034285105764865875, -0.06475269049406052,
-0.008249715901911259, 0.031450431793928146, 0.008753931149840355, -0.06284607946872711, -0.00269001311803341,
0.061753395944833755, 0.0314679853618145, 0.005365008022636175, -0.034285105764865875, -0.06475269049406052,
0.06229100224342346, -0.016091199591755867, -0.038237784057855606, -0.01697474904358387, 0.0023320959880948067,
-0.02873411774635315, -0.07216104120016098, 0.04663623124361038, 0.023897146806120872, -0.0282142040193081,
-0.03714695945382118, -0.055613692849874496, -0.0028377221897244453, -0.06574894487857819, -0.061038118572235,
0.06294400244951248, 0.0034130788408219814, 0.07920042425394058, 0.007338271436676025, -0.050656369000997162,
-0.02522268146276474, 0.027450041845440865, -0.01720043271780014, 0.046272337436676025, -0.05018896237015724,
0.015779439359903336, -0.026586400344967842, -0.01974601112306118, -0.00036689057014882565, -0.016816521063447,
-0.025464840233295958, 0.0007100807852111757, 0.04524853080511093, 0.0010508883278816938, -0.00547241186723131,
0.011604292318224907, -0.0427076481282711, -0.02004644088447094, -0.06824997067451477, -0.08084388822317123,
-0.08167271316051483, 0.038480401039123535, -0.0414984643399414, 0.0621405728161335, 0.01636849343776703,
-0.0277055791497898, 0.02410232089459896, 0.021344885230064392, 0.056428126990795135, 0.02979239635169506,
-0.05207456275820732, 0.004299748223274946, 0.03417612612247467, 0.034210722893476486, 0.0001064265313693222,
0.011242502368986607, 0.037193565656074524, -0.004098605364561081, 0.01320237666364845, 0.021659305319190052,
0.03850701451301575, -0.03979567810893059, 0.024909289553761482, 0.0036120524164289236, 0.030269717797636986,
0.03532775491476059, 0.04048445075750351, -0.021236592903733253, 0.05895552039146423, 0.04913758486509323,
-0.047305576503276825, 0.05272332951426506, 0.012154217809438705, -0.02513653226196766, -0.0105582932010293,
-0.049685653299093246, 0.032950107008218765, -0.007436738815158606, -0.0749432702434189, -0.0447106082201004,
0.03816404938697815, -0.029877835884690285, -0.020543526858091354, 0.0253277961647129, 0.011234065571783165,
0.07374250143766403, 0.04288359731435776, 0.03435317426919937, -0.02951200306415558, -0.09385887533426285,
-0.005317367613315582, 0.01705515943467617, -0.00934696663171053, 0.01293235830962658, 0.02108096517622471,
0.030062183737754822, 0.004270109347999096, -0.00579592073362436, 0.006119553931057453, -0.009726069867610931,
-0.0016054088482633233, -0.12823154032320377, 0.005963715258985758, -0.01607099547982216, ]
```

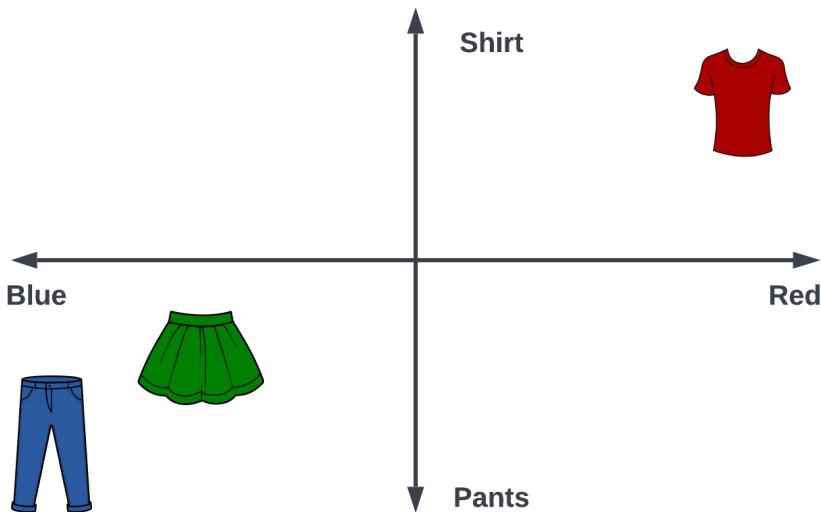
Texto

Embedding

1536 dimensões

» Como funciona o Vector Search?

- Documentos, imagens e outros conteúdos são transformados em **embeddings** (vetores com N dimensões) que representam seu significado.
- Busca semântica é realizada a partir da **similaridade** de suas representações numéricas (os embeddings)
- As métricas de similaridade mais comuns são: Dot Product, Coseno e Euclidiana.



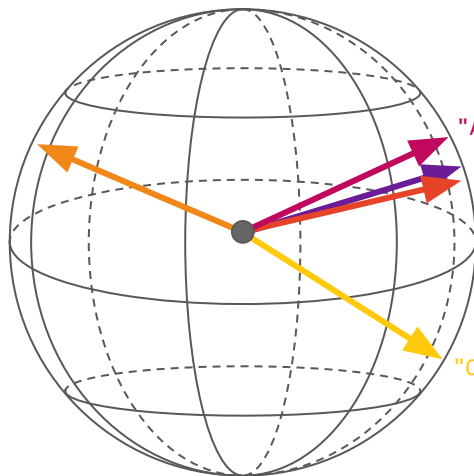
<https://www.datastax.com/guides/what-is-a-vector-embedding>

Vector embeddings

Representação fiel de "coisas" como vetores em um espaço de alta dimensionalidade

- Quais coisas? Principalmente **texto** (também: imagens, vídeos, sons ...)
- Transformar um problema semântico em um problema geométrico

"Aquele que busca a igualdade entre os desiguais ..."

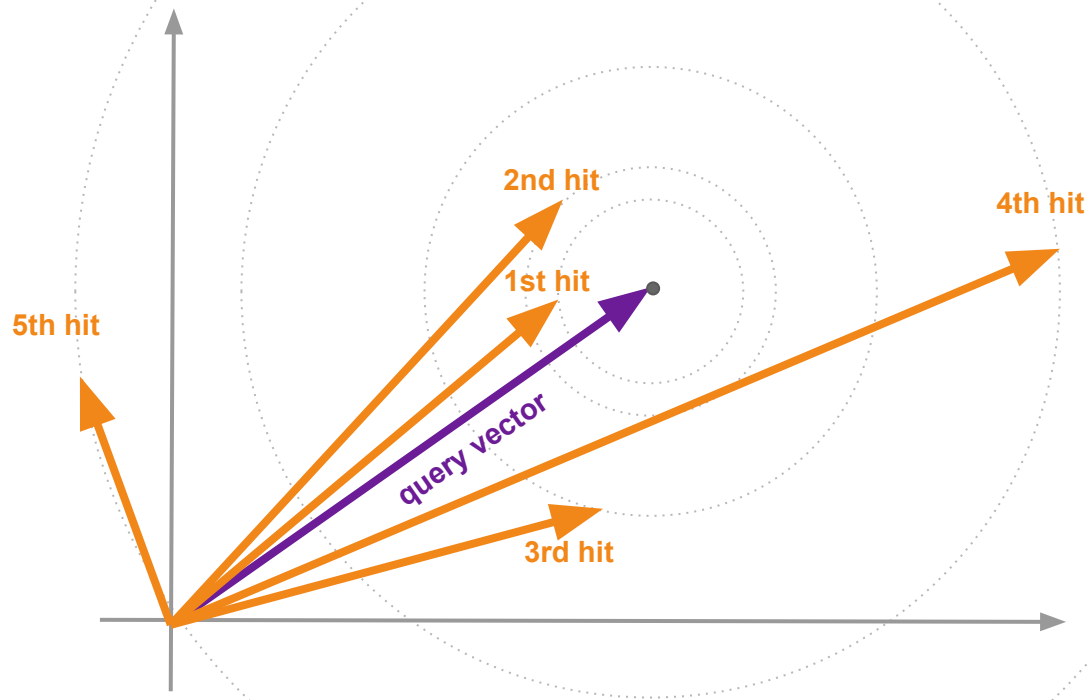


"A escassa satisfação que o homem pode ..."
"Lutamos toda a nossa vida por nada"
"Viver é sofrer, sobreviver é ..."

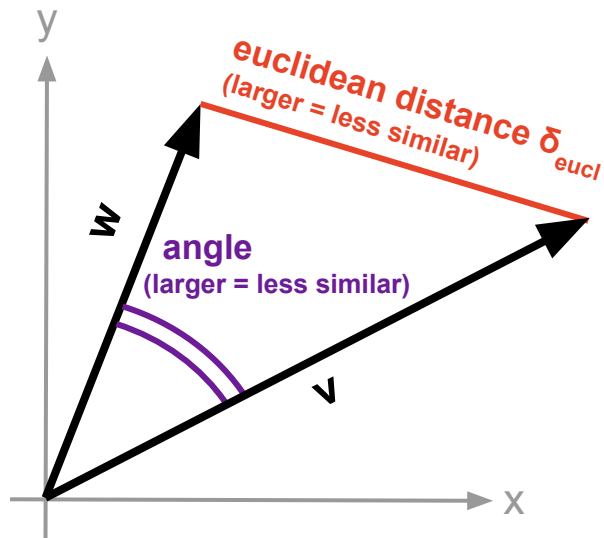
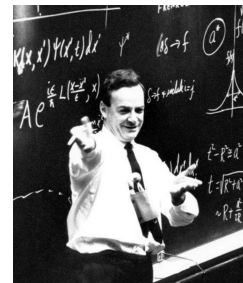
"O nada assombra o Ser."

Embedding vector space

Vector Search em uma imagem



Similaridades



Euclidean similarity

$$S_{\text{eucl}}(\mathbf{v}, \mathbf{w}) = \frac{1}{1 + \sum_i (v_i - w_i)^2} = \frac{1}{1 + \delta_{\text{eucl}}^2(\mathbf{v}, \mathbf{w})}$$

$$\delta_{\text{eucl}}(\mathbf{v}, \mathbf{w}) = \sqrt{\sum_i (v_i - w_i)^2}$$



Cosine similarity (a.k.a. "only the angle counts")

$$S_{\text{cos}}(\mathbf{v}, \mathbf{w}) = \frac{1 + \frac{\sum_i v_i w_i}{|\mathbf{v}| |\mathbf{w}|}}{2} = \frac{1 + \frac{\mathbf{v} \cdot \mathbf{w}}{|\mathbf{v}| |\mathbf{w}|}}{2}$$



DATASTAX

ASTRA

Data Platform all in one:

NoSQL, Vector Search, APIs, Streaming, ZeroOps, Multi Cloud, etc...

**Applications****Data sources**

Web



Applications



IoT

**DATASTAX****Starlight**

JMS, Kafka, RabbitMQ

gRPC, REST, GraphQL,
Document (JSON) and CQL
APIs **PULSAR**Streaming service based
on Apache Pulsar®Database as a service based on
Apache Cassandra®**CDC****Invisible operations**Infrastructure, Security, Scalability,
High Availability, Resiliency**Multi-Cloud & Cloud-Native**Multi-region, Inter-cloud, Serverless
and K8S microservices based**Data ecosystem**

Analytics



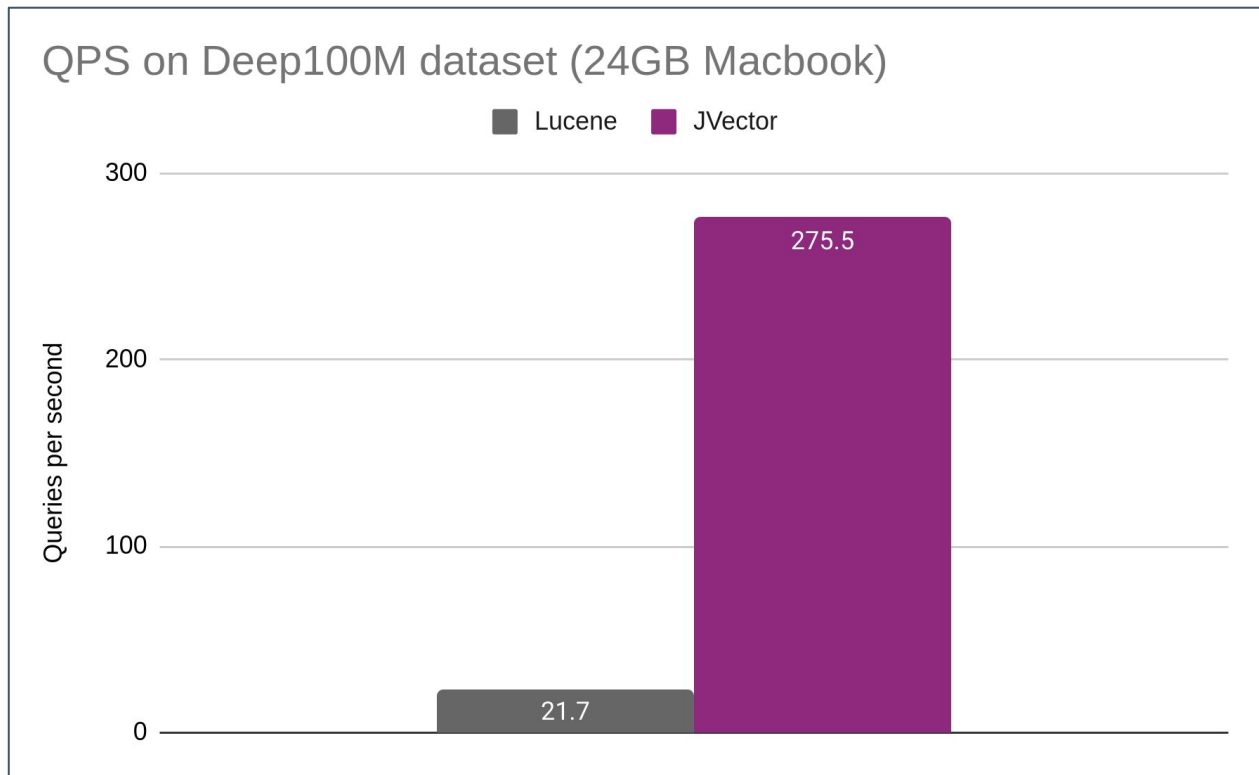
Search

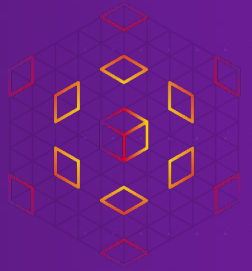


AI

**kubernetes****DATASTAX** DEVELOPERS

JVector performance: read throughput





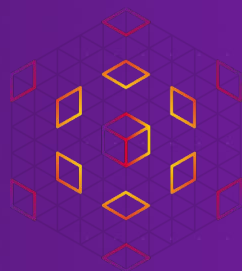
Ex2 - RAG + LangchainJS



› Ex 2: Simplificando com LangchainJS

Passos:

- Acessar localhost:3000/**v2**/chat e fazer perguntas sobre financiamento imobiliário
- Criar collection "**real_estate_financing_langchain**" no Astra
- Carregar documento em localhost:3000/**v2**/upload
- Repetir perguntas em v2/chat



Ex3 - Stocks



› Ex 3: GenAI + DataAPI + NextJS

Passos:

- Acessar localhost:3000/**stocks**
- Criar collection "**stocks**" no Astra
- Carregar dados para a collection
- Perguntar sobre preço de ação: AAPL, GOOG, SBUX

Ganhe um Badge



Envie o seu DB ID e um screenshot da aplicação:

samuel.matioli@datastax.com

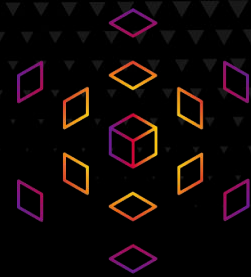
para receber seu Badge!

Sorteio de brindes

Cadastre-se aqui para participar do sorteio.

bit.ly/datastax-tdcsc





Samuel Matioli

Solution Engineer @ DataStax

samuel.matioli@datastax.com

[linkedin.com/in/samuelmatioli](https://www.linkedin.com/in/samuelmatioli)

samuelmatioli.medium.com

[youtube.com/@samuelmatioli](https://www.youtube.com/@samuelmatioli)



Appendix

› Referências adicionais

AstraDB no Langchain

<https://python.langchain.com/docs/integrations/vectorstores/astradb>

RAG na prática

<https://www.youtube.com/watch?v=wc22tG5IGUU>

AstraDB Vector & Vertex AI

<https://www.youtube.com/watch?v=E0Wv4ALVp5w>