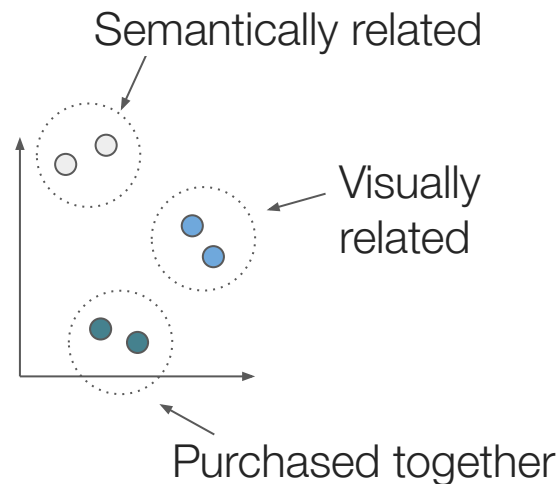


Embeddings as data structures 2.0?

"Data structures 2.0 is written in human unfriendly language, such as the floating point values of an embedding. No human is involved in writing this code ... and coding directly in the floating point values is kind of tedious but possible (I tried)." [2]

$$\operatorname{argmin}_{N_{\theta}, \theta} L(\mathbf{X}, \mathbf{Y}; N_{\theta})$$

Network architecture \nearrow N_{θ}, θ \nwarrow Parameters \nearrow \nwarrow Data describing the similarity to encode



Building Software 2.0 with embeddings

Jesse Clark
Marqo

About

@Marqo - The Embeddings Cloud

Co-founder/CTO 

@jn2clark  

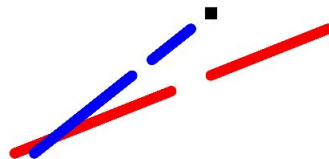
Physics @ UCL & Stanford



ML and OR @ StitchFix



Deep learning and AI @ Amazon



Building Software 2.0 with embeddings

Outline

- What is Software 2.0?
 - Examples
- What are embeddings?
 - Producing embeddings
 - Training models for embeddings
- Building with embeddings
 - Classification, search and RAG

What is Software 2.0?

*"**Software 2.0** is written in human unfriendly language, such as the weights of a neural network. No human is involved in writing this code ... and coding directly in weights is kind of hard (I tried)." [1]*

[1] <https://karpathy.medium.com/software-2-0-a64152b37c35>

What is Software 2.0? - Examples

Visual Recognition



What is Software 2.0? - Examples

Structure prediction



What is Software 2.0? - Examples

Text summarisation

The bottleneck is no longer access to information; now it's our ability to keep up.

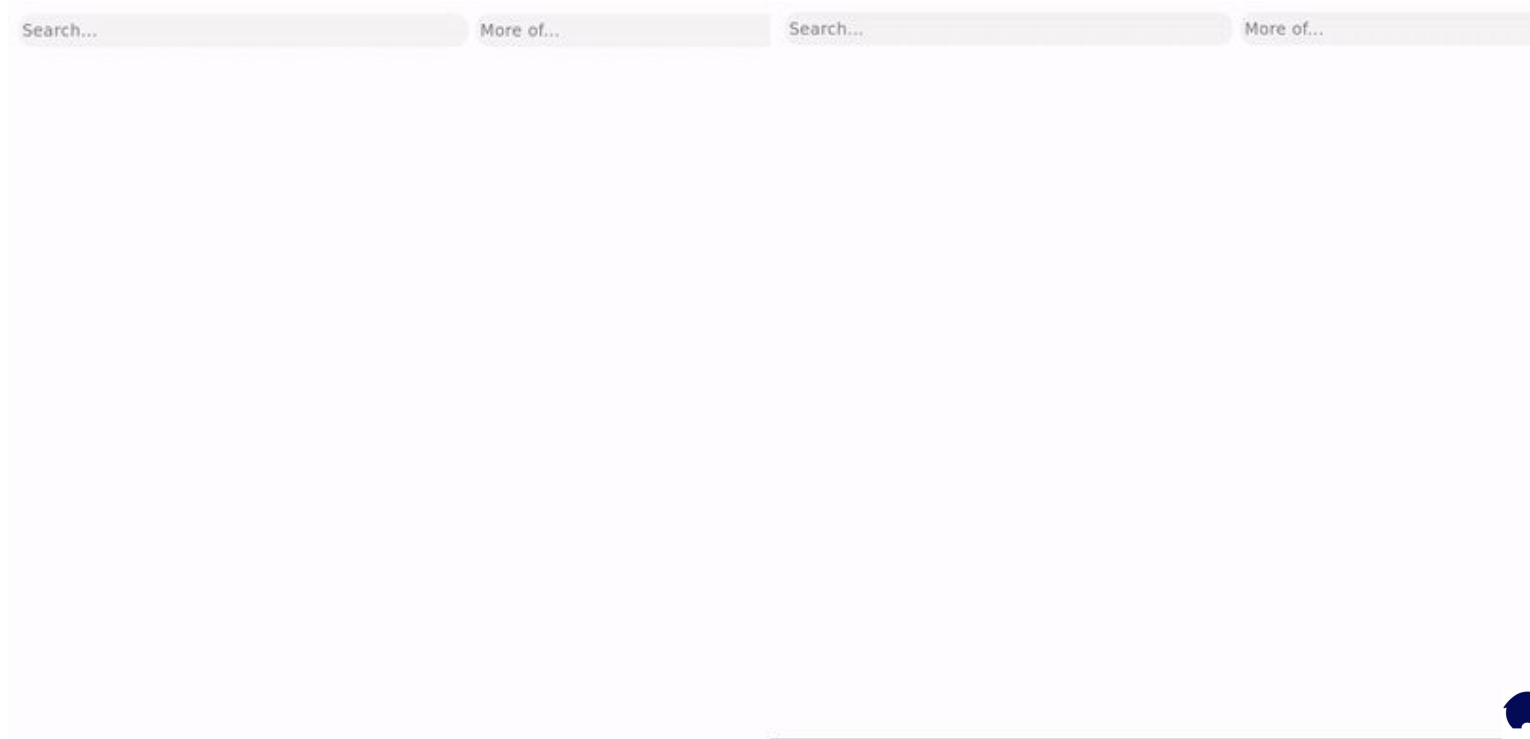
AI can be trained on a variety of different types of texts and summary lengths.

A model that can generate long, coherent, and meaningful summaries remains an open research problem.

The last few decades have witnessed a fundamental change in the challenge of taking in new information. The bottleneck is no longer access to information; now it's our ability to keep up. We all have to read more and more to keep up-to-date with our jobs, the news, and social media. We've looked at how AI can improve people's work by helping with this information deluge and one potential answer is to have algorithms automatically summarize longer texts. Training a model that can generate long, coherent, and meaningful summaries remains an open research problem. In fact, generating any kind of longer text is hard for even the most advanced deep learning algorithms. In order to make summarization successful, we introduce two separate improvements: a more contextual word generation model and a new way of training summarization models via reinforcement learning (RL). The combination of the two training methods enables the system to create relevant and highly readable multi-sentence summaries of long text, such as news articles, significantly improving on previous results. Our algorithm can be trained on a variety of different types of texts and summary lengths. In this blog post, we present the main contributions of our model and an overview of the natural language challenges specific to text summarization.

What is Software 2.0? - Examples

Search and recommendations



GITHUB COPILOT: CHAT



Hi @monalisa, how can I help you?

I'm powered by AI, so surprises and mistakes are possible. Make sure to verify any generated code or suggestions, and share feedback so that we can learn and improve.

Ask a question or type '/' for commands



parse_expenses.py × addresses.rb × sentiments.ts ×

```
1 import datetime
2
3 |
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
```



What is Software 2.0?

"Software 2.0 is written in human unfriendly language, such as the weights of a neural network. No human is involved in writing this code ... and coding directly in weights is kind of hard (I tried)." [1]

Software 1.0

Rules + data → Answers

```
def binary_search(arr, target):
    low = 0
    high = len(arr) - 1

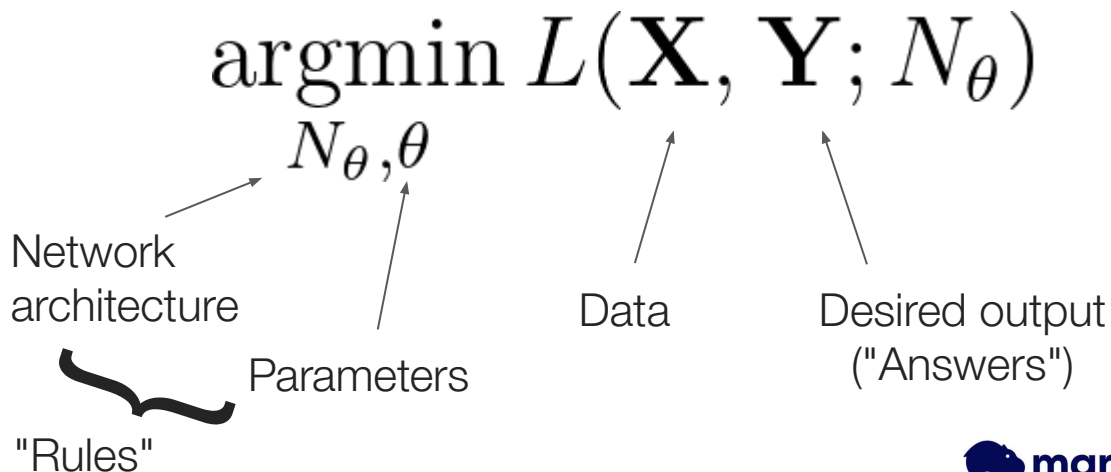
    while low <= high:
        mid = (low + high) // 2

        if arr[mid] == target:
            return mid
        elif arr[mid] < target:
            low = mid + 1
        else:
            high = mid - 1

    return -1
```

Software 2.0

Answers + data → Rules



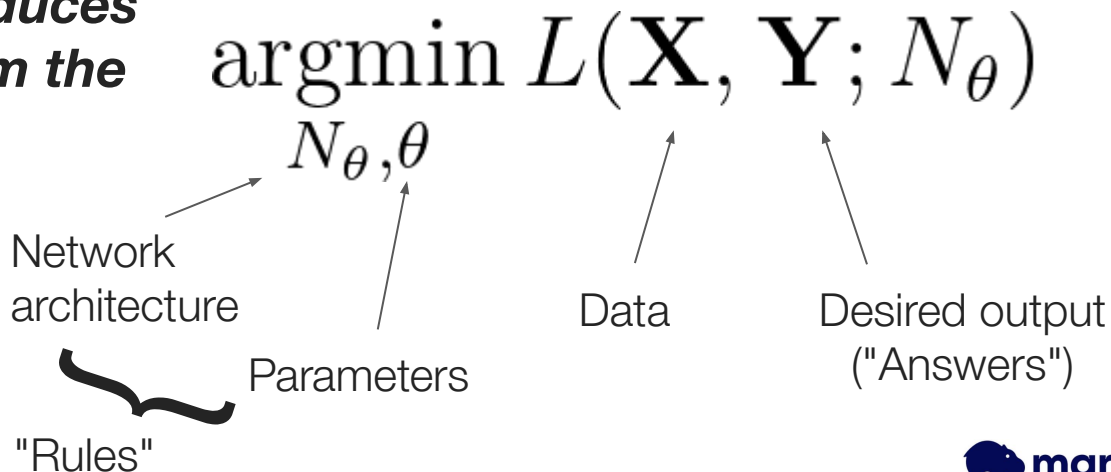
What is Software 2.0?

"Software 2.0 is written in human unfriendly language, such as the weights of a neural network. No human is involved in writing this code ... and coding directly in weights is kind of hard (I tried)." [1]

Software 2.0

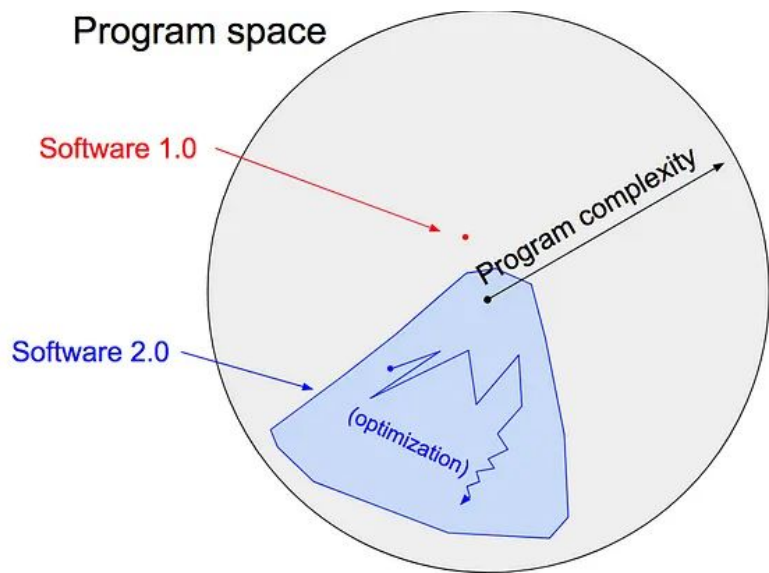
Answers + data \rightarrow Rules

"Find a function that produces the outputs we want from the data we have"



What is Software 2.0?

"Software 2.0 is written in human unfriendly language, such as the weights of a neural network. No human is involved in writing this code ... and coding directly in weights is kind of hard (I tried)." [1]



Software 2.0

Answers + data \rightarrow Rules

$$\operatorname{argmin}_{N_{\theta}, \theta} L(\mathbf{X}, \mathbf{Y}; N_{\theta})$$

Why Software 2.0? - Benefits and considerations

Benefits

- Computationally homogeneous
- Constant running time and memory
- Simple to bake into silicon
- It is highly portable
- Direct optimization
- It is much better

Drawbacks

- Harder to understand
- Soft and unexpected failure modes
- Can be difficult to test
- Requires strict testing and monitoring of data
 - "unit tests for data"

What is Software 2.0? - Summary

Software 1.0:

- Human-engineered source code (e.g., .cpp files)
- Compiled into a binary that does useful work

```
def binary_search(arr, target):
    low = 0
    high = len(arr) - 1

    while low <= high:
        mid = (low + high) // 2

        if arr[mid] == target:
            return mid
        elif arr[mid] < target:
            low = mid + 1
        else:
            high = mid - 1

    return -1
```

Software 2.0:

- Source code comprises:
 - Dataset defining behavior
 - Neural net architecture providing code skeleton
- Many details (weights) to be filled in
- Training process compiles dataset into final neural network binary

$$\operatorname{argmin}_{N_{\theta}, \theta} L(\mathbf{X}, \mathbf{Y}; N_{\theta})$$

Building Software 2.0 with embeddings

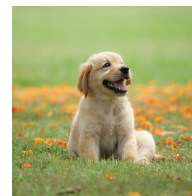
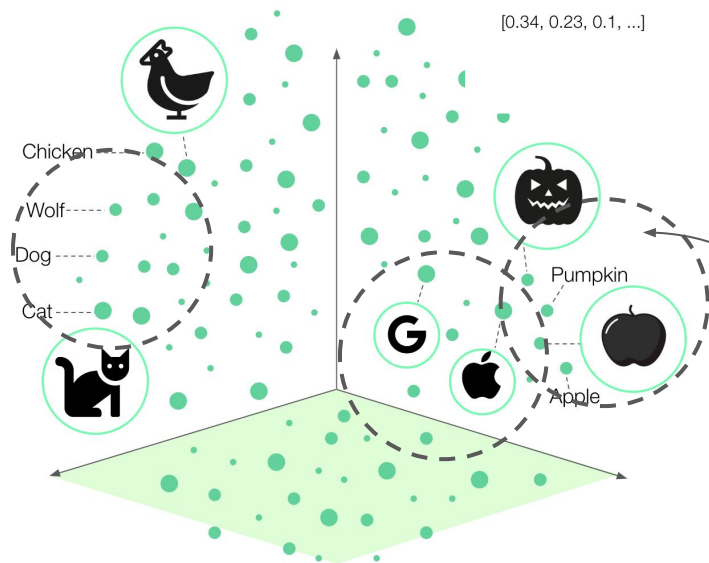
Outline

- What is Software 2.0?
 - Examples
- What are embeddings?
 - Producing embeddings
 - Training models for embeddings
- Building with embeddings
 - Classification, search and RAG

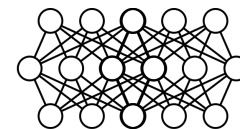
What are embeddings?

Vectors (Embeddings)

- Encode the meaning of objects comprised of text, images, audio, video as vectors
- Embeddings are (learned) vectors
- Proximity in vector space equals "similarity"



A dog sitting on some grass



Model

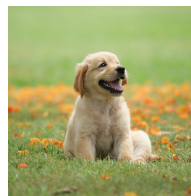
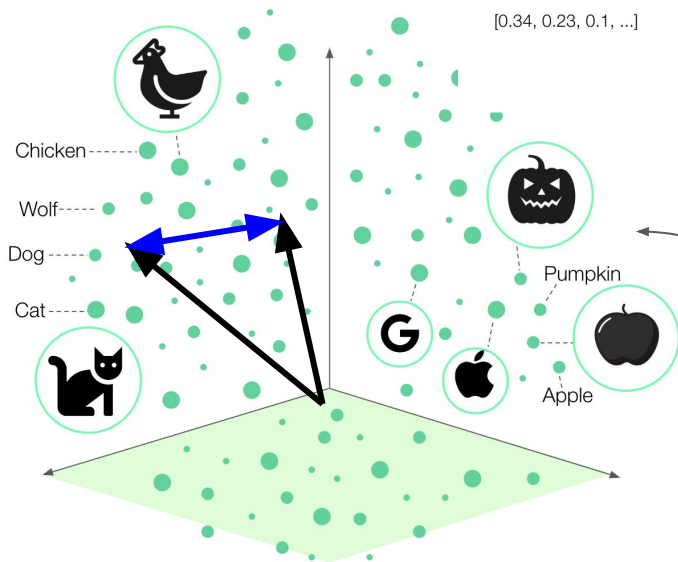
$[0.11, -0.24, \dots, \dots, \dots, \dots, \dots, \dots, \dots, 0.07, 0.16, -0.07]$



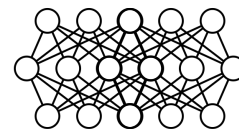
What are embeddings?

Vectors (Embeddings)

- Encode the meaning of objects comprised of text, images, audio, video as vectors
- Embeddings are (learned) vectors
- Proximity in vector space equals "similarity"



A dog sitting on some grass



Model

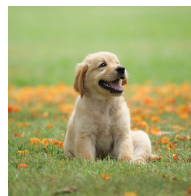
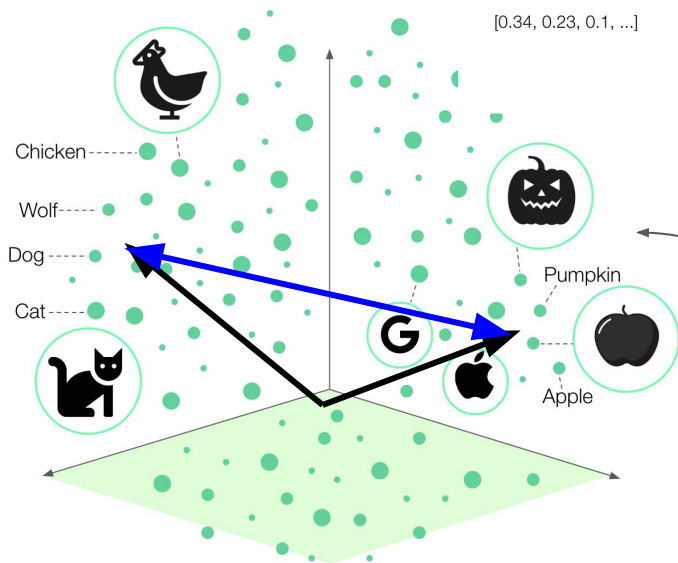
[0.11, -0.24, ..., ..., ..., ..., ..., ..., ..., 0.07, 0.16, -0.07]



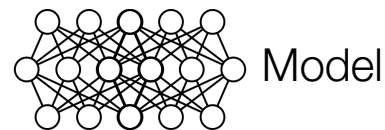
What are embeddings?

Vectors (Embeddings)

- Encode the meaning of objects comprised of text, images, audio, video as vectors
- Embeddings are (learned) vectors
- Proximity in vector space equals "similarity"



A dog sitting on some grass



Model

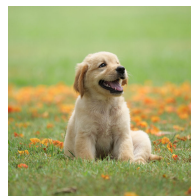
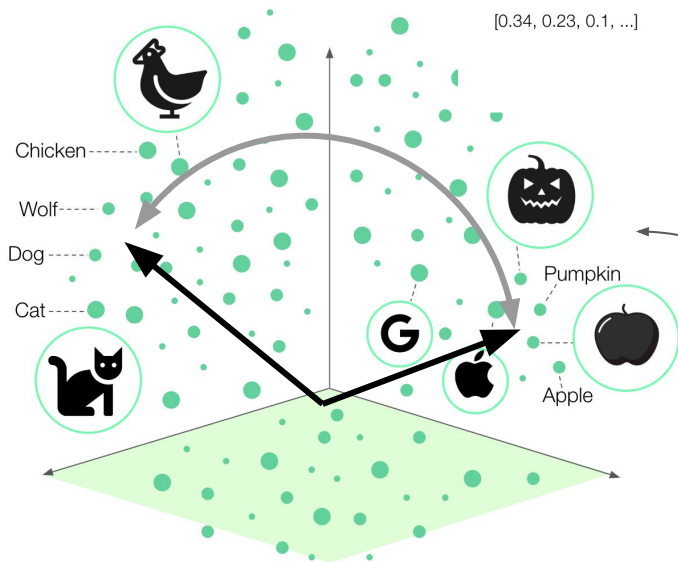
[0.11, -0.24, ..., ..., ..., ..., ..., ..., ..., 0.07, 0.16, -0.07]



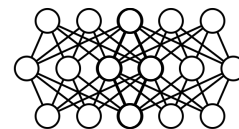
What are embeddings?

Vectors (Embeddings)

- Encode the meaning of objects comprised of text, images, audio, video as vectors
- Embeddings are (learned) vectors
- Proximity in vector space equals "similarity"



A dog sitting on some grass



Model

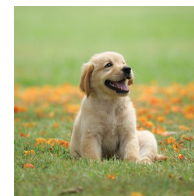
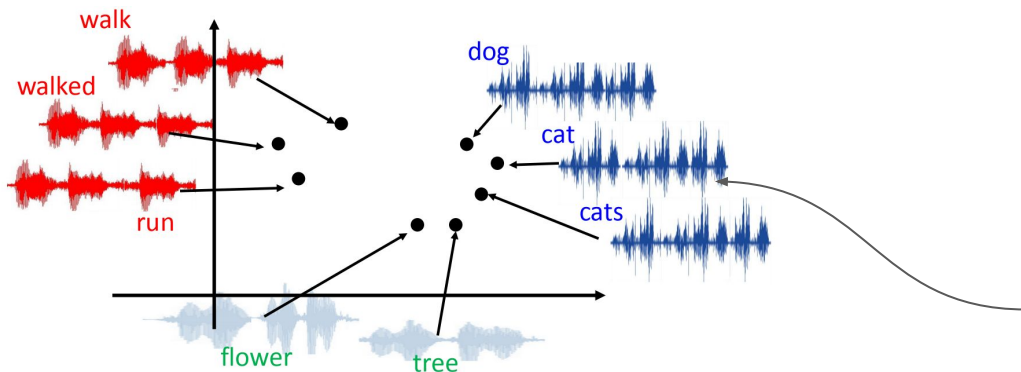
[0.11, -0.24, ..., ..., ..., ..., ..., ..., ..., 0.07, 0.16, -0.07]



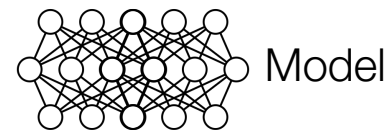
What are embeddings?

Vectors (Embeddings)

- Encode the meaning of objects comprised of text, images, audio, video as vectors
- Embeddings are (learned) vectors
- Proximity in vector space equals "similarity"



A dog sitting on some grass

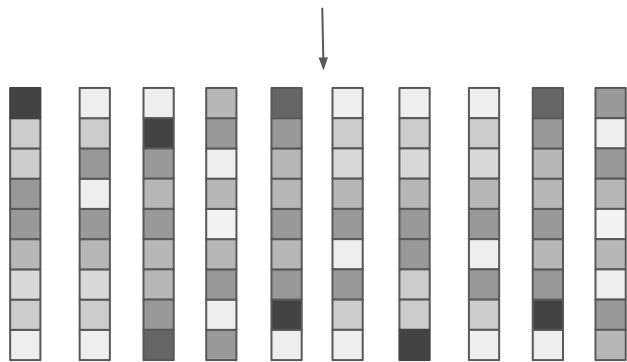
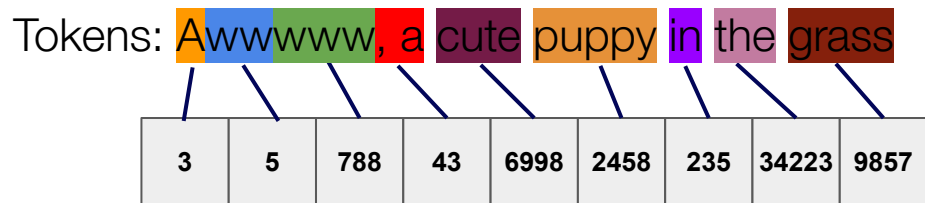


[0.11, -0.24, ..., ..., ..., ..., ..., ..., ..., ..., ..., 0.07, 0.16, -0.07]

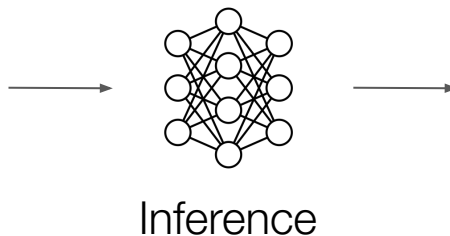


What are embeddings? How models understand text

Query: Awwwww, a cute puppy in the grass



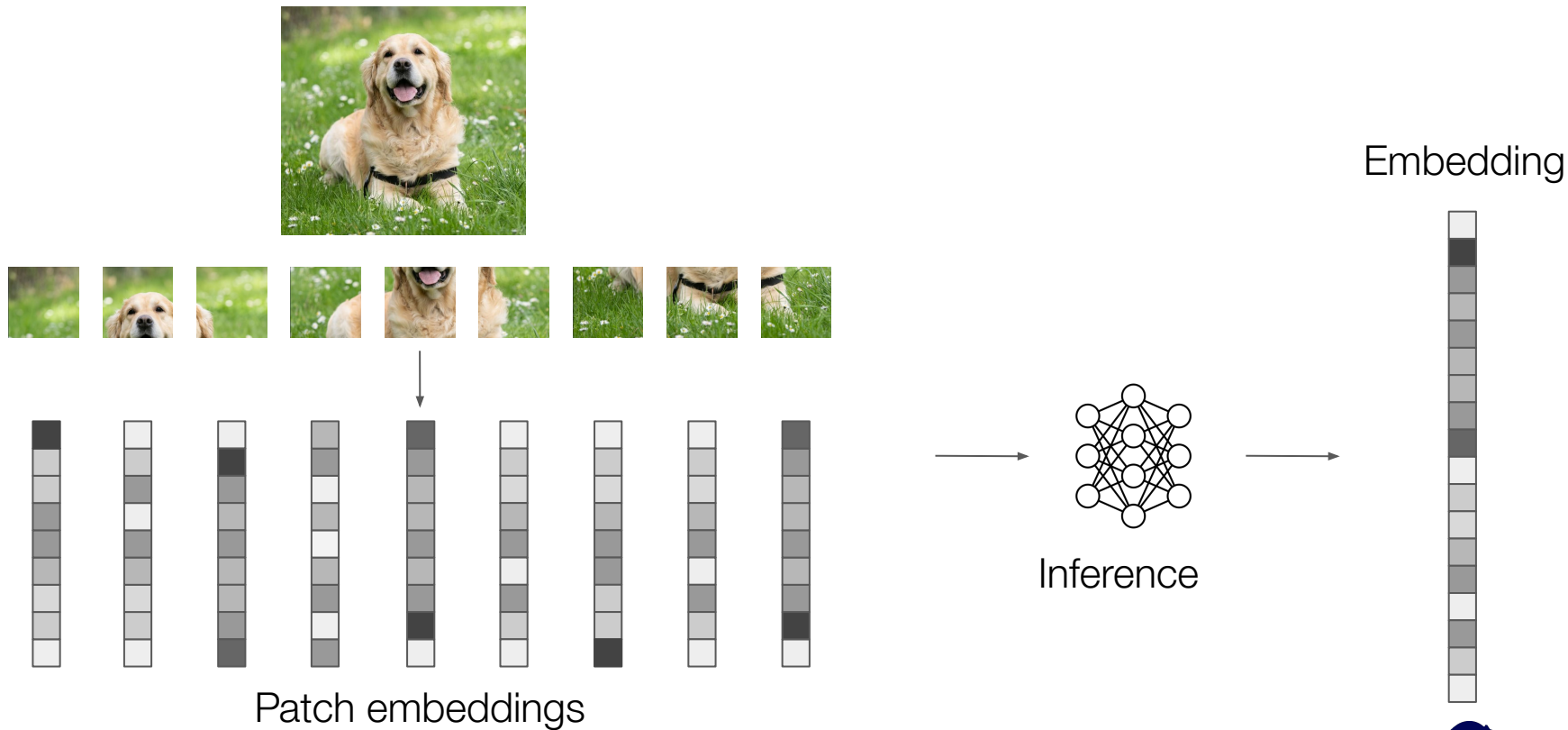
Token embeddings



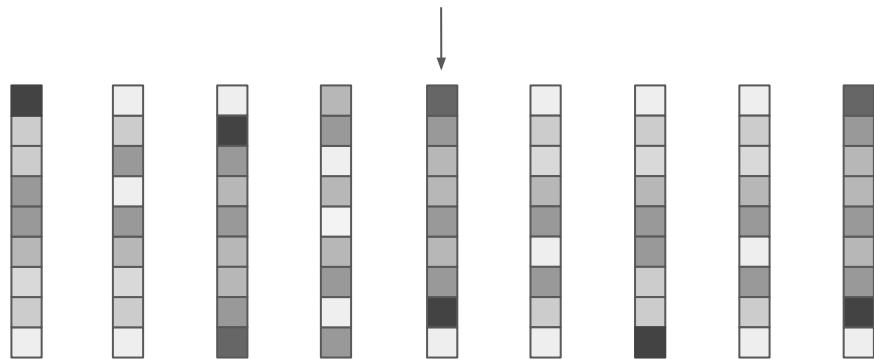
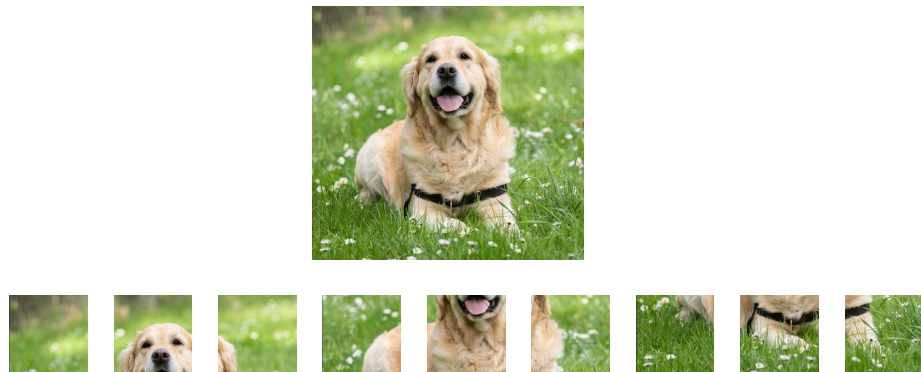
Embedding



What are embeddings? How models understand images



What are embeddings? How models understand images



Patch embeddings

$$f\left(\begin{matrix} m \\ \text{Data} \end{matrix} \cdot \begin{matrix} n \\ \text{"Weights" - to be learnt } (\theta) \end{matrix} = \begin{matrix} n \\ \text{Embedding} \end{matrix} \right)$$

The diagram shows a function f that takes a matrix A (Data) of size $m \times n$ and a matrix B (Weights) of size $n \times n$ as input, resulting in a matrix C (Embedding) of size $m \times n$. The matrices are represented as grids of colored squares: green for A , red for B , and blue for C .

Embedding

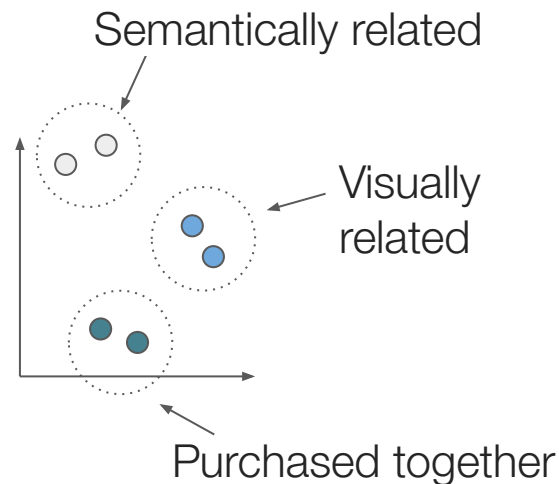


Embeddings as data structures 2.0?

"Data structures 2.0 is written in human unfriendly language, such as the floating point values of an embedding. No human is involved in writing this code ... and coding directly in the floating point values is kind of tedious but possible (I tried)." [2]

$$\operatorname{argmin}_{N_{\theta}, \theta} L(\mathbf{X}, \mathbf{Y}; N_{\theta})$$

Network architecture \nearrow N_{θ}, θ \nwarrow Parameters
Data describing the similarity to encode \nearrow \nwarrow

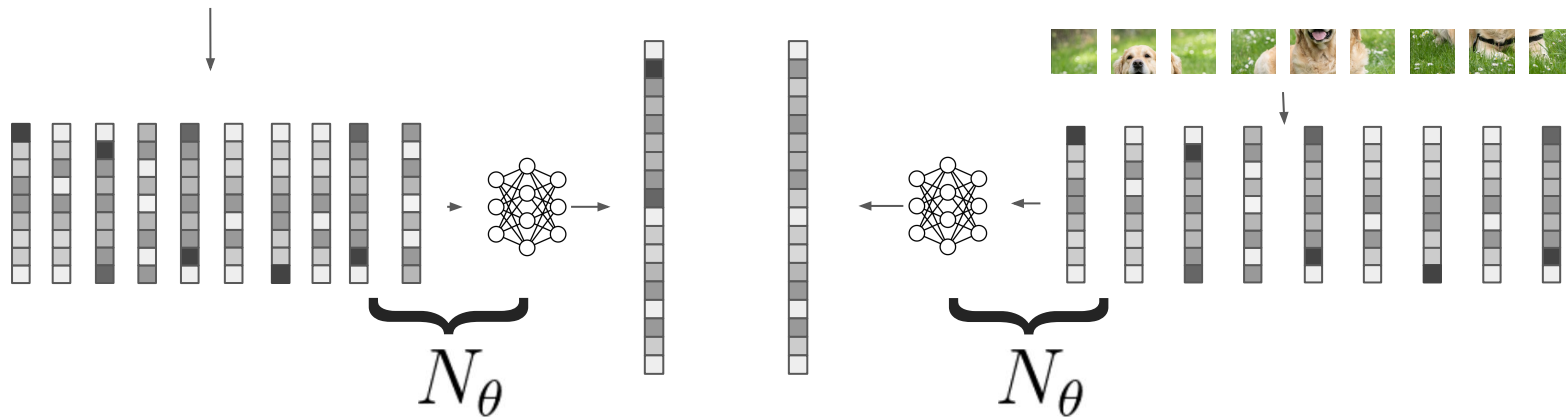
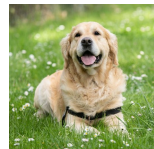


Training embedding models $\operatorname{argmin}_{N_\theta, \theta} L(\mathbf{X}, \mathbf{Y}; N_\theta)$

Query: Awwwww, a cute puppy in the grass

Tokens: Awwwww, a cute puppy in the grass

\mathbf{X} \mathbf{Y}



**Make these embeddings
as close as possible**

$\operatorname{argmin} L$

Marqo FashionCLIP

State-of-the-art embeddings model for fashion search and recommendations



<https://www.marqo.ai/blog/search-model-for-fashion>

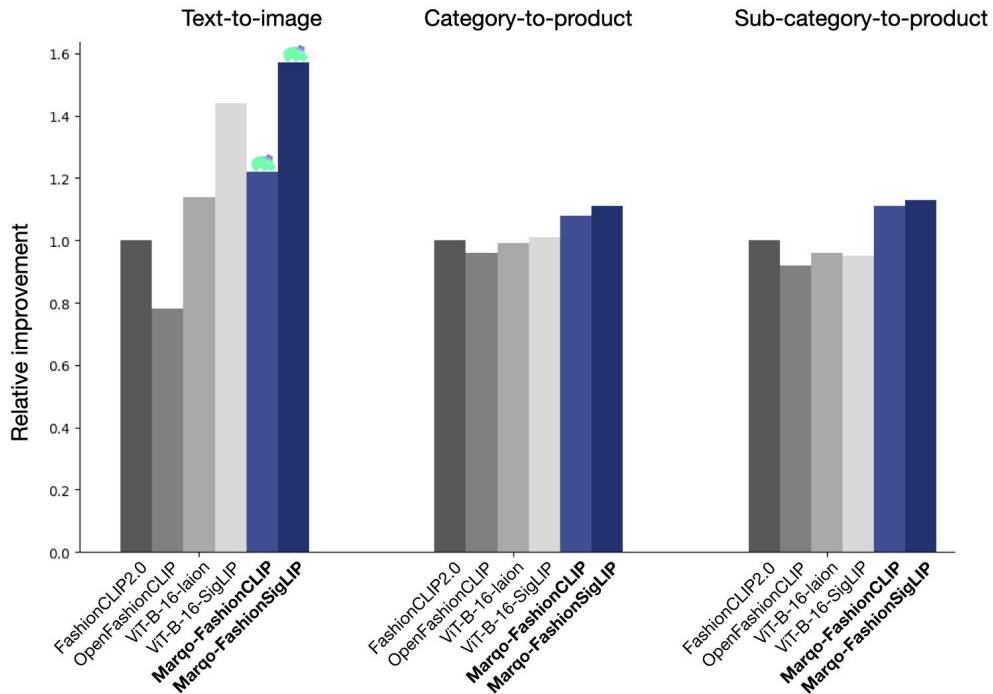


$L(\text{image, text}) +$
 $L(\text{image, colors}) +$
 $L(\text{image, materials}) +$
 $L(\text{image, style}) +$
 $L(\text{image, keywords})$

Midriff Waist Shaper Dress
Colors: blue, purple, pink
Materials: cotton, linen
Style: casual, spring/summer
Keywords: Floral midi dress, blue floral dress

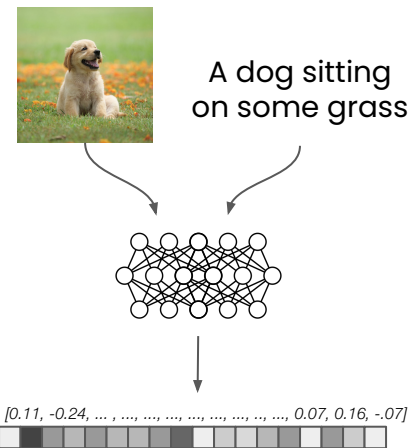
Marqo FashionCLIP

State-of-the-art embeddings model for fashion search and recommendations

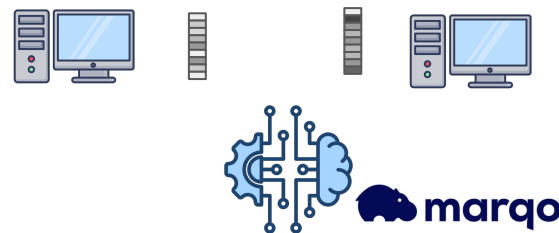


Summary of software 2.0 and embeddings

- Software 2.0
 - learns the functions from data instead of being handcrafted
- Embeddings
 - learn the representation of data instead of being handcrafted
- Embeddings are learned vectors
- Embeddings provide a consistent format for diverse unstructured data



$$\operatorname{argmin}_{N_{\theta}, \theta} L(\mathbf{X}, \mathbf{Y}; N_{\theta})$$

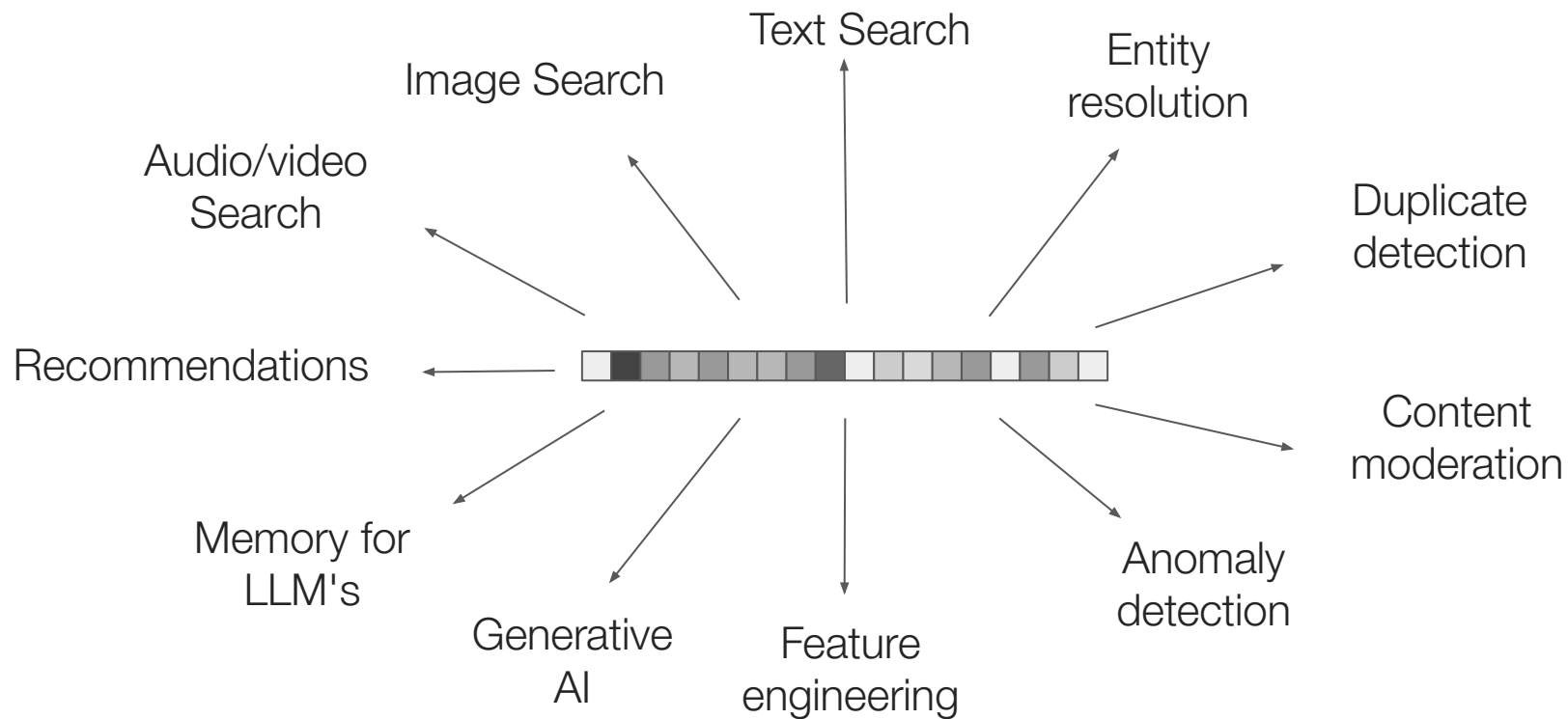


Building Software 2.0 with embeddings

Outline

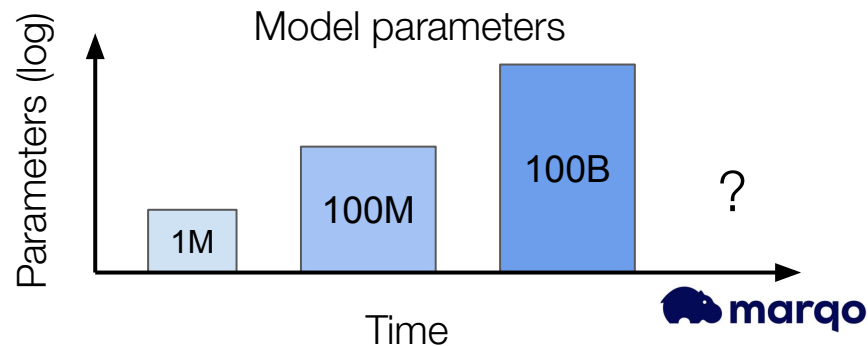
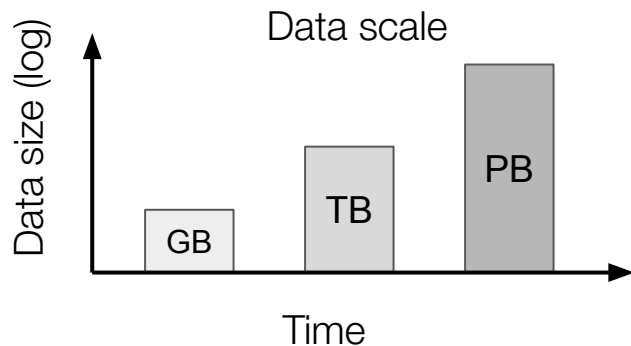
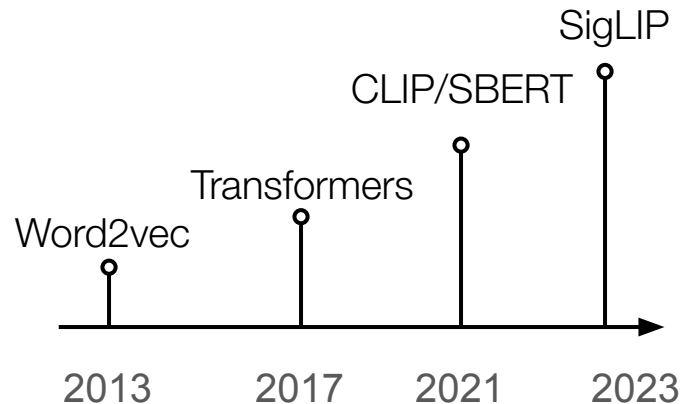
- What is Software 2.0?
 - Examples
- What are embeddings?
 - Producing embeddings
 - Training models for embeddings
- Building with embeddings
 - Classification, search and RAG

Building with Embeddings



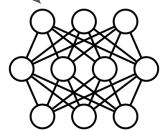
Building with Embeddings - why now?

- Embedding quality has dramatically improved
- Driven by:
 - data quality and quantity
 - training methods
 - hardware improvements
 - architecture improvements



Classification

"Dune is the best movie ever."



Embedding model



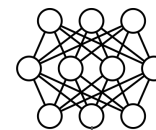
Embedding of the text

Labels

"Cinema"

"Restaurant"

"Clothing"



Embeddings of the labels

Calculate the similarity of
the content and labels $d(c, \mathbf{L})$

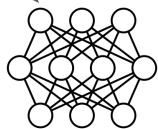
"Cinema": 0.98

"Restaurant": 0.01

"Clothing": 0.01

Classification

"Dune is the best movie ever."



Embedding model

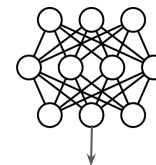


Embedding of the text

Labels

"Positive"

"Negative"



Embeddings of the labels

Calculate the similarity of
the content and labels $d(c, \mathbf{L})$

"Positive": 0.99

"Negative": 0.01

```
# Define the query and classification labels
query = ["Dune is the best movie ever."]
labels = ["Cinema", "Food", "Pet food"]

# Step 1: Get the embedding for the query
query_embedding = get_embedding(query) # type: ignore

# Step 2: Get the embeddings for each classification label
label_embeddings = get_embedding(labels) # type: ignore

# Step 3: Calculate cosine similarity between the query embedding and each label embedding
# This returns a similarity score for each label
similarities = cosine_similarity(query_embedding, label_embeddings)

# Step 4: Find the label with the highest similarity score
max_similarity_index = np.argmax(similarities)
predicted_label = labels[max_similarity_index]

# Step 5: Print the result
print(f"The predicted label for the query is: {predicted_label}")
```



Upload Fashion Item Image

Drop Image Here
- or -
Click to Upload

Or provide an image URL

Classify Clear

Or click on one of the images below to classify it:

Examples

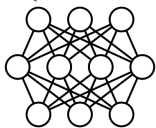
Label

📄

<https://huggingface.co/spaces/Marqo/Marqo-FashionSigLIP-Classification>

Search

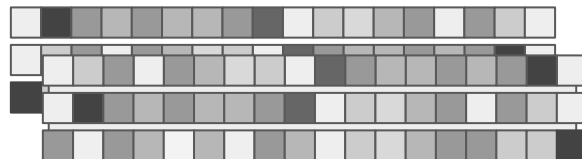
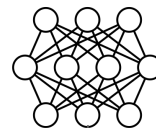
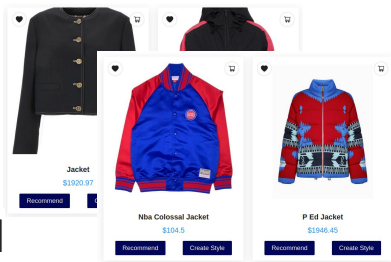
"retro blue jacket"



Embedding model



Embedding of the query

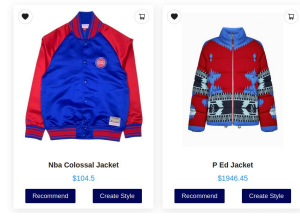


Embeddings of the products

Calculate the similarity of the query and products

$d(q, \mathbf{P})$

"Product4": 0.76
"Product3": 0.74
"Product1": 0.18
"Product2": 0.11



Search for this... More of this... Less of this...  Filters  Score Modifiers 

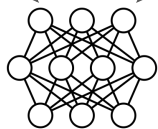
Search as you type Show live recommendations

Search

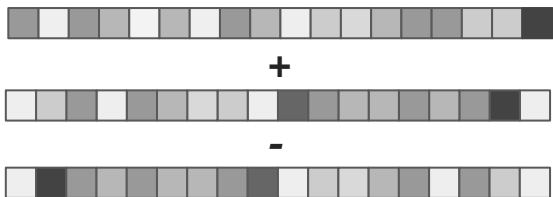
"retro blue jacket"

"sports"

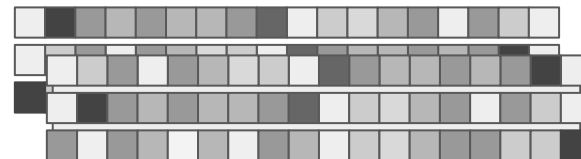
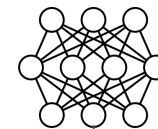
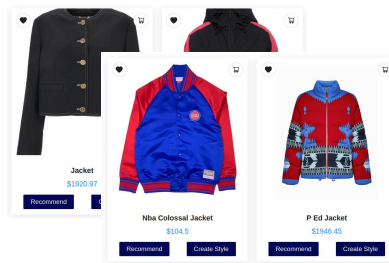
"shiny"



Embedding model



Embedding of the queries

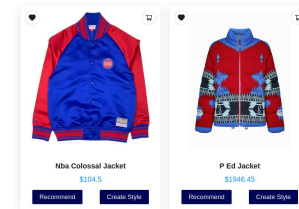


Embeddings of the products

Calculate the similarity of the query and products

$$d(q, \mathbf{P})$$

"Product4": 0.76
"Product3": 0.74
"Product1": 0.18
"Product2": 0.11



```
import marqo
from marqo import Client

# initialize the client
mq = Client("https://api.marqo.ai", api_key=api_key)

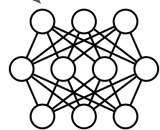
# create an index
mq.create_index(index_name,
                model="open_clip/ViT-B-32/laion2b_s34b_b79k")

# embed content
mq.index(index_name).add_documents(content)

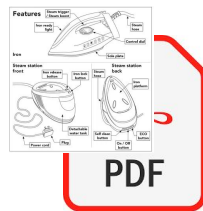
# perform a search query on the index
results = mq.index(index_name).search(q="retro blue jacket")
```

Retrieval augmented generation (RAG)

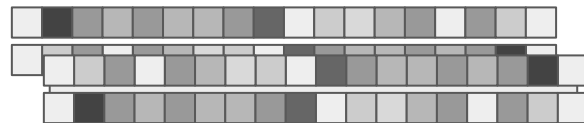
"what is the clicking sound"



Embedding model



Embedding of the query



Embeddings of the document chunks

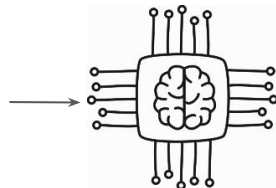
$d(q, \mathbf{C})$

"chunk1": 0.76

"chunk4": 0.74

"chunk2": 0.18

"chunk3": 0.11



Send top chunks to LLM

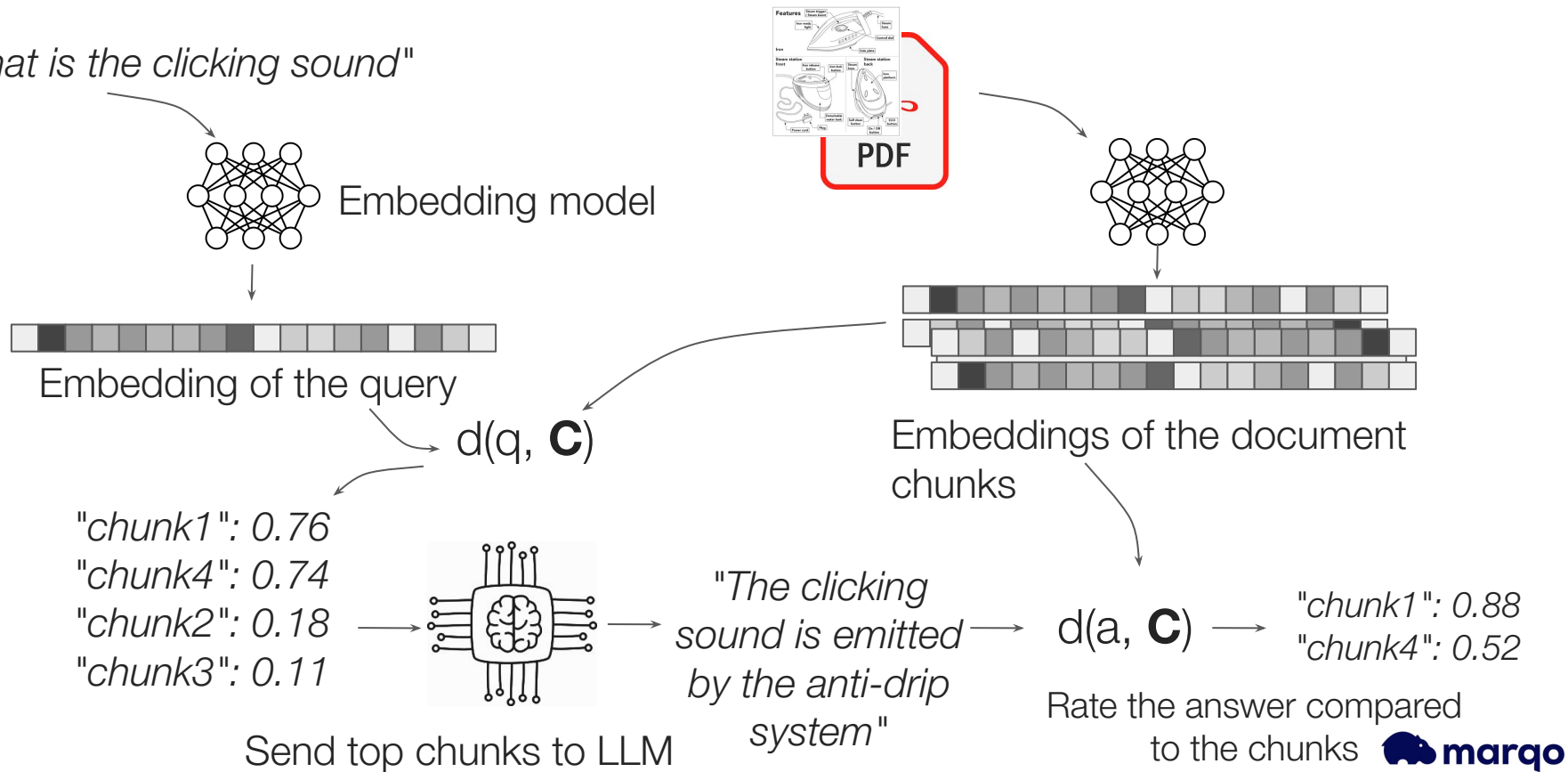
"The clicking sound is emitted by the anti-drip system"



Please enter some text to start searching...

Retrieval augmented generation (RAG)

"what is the clicking sound"



Building with Embeddings

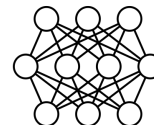
- Embedding based systems share a lot of similarities
- Embedding generation (inference) and comparing (scoring) make up the core
- Embedding generation handled by deep learning models
- The comparison (retrieval) of embeddings is handled by a vector database

Train



$$\underset{N_{\theta}, \theta}{\operatorname{argmin}} L(\mathbf{X}, \mathbf{Y}; N_{\theta})$$

Generate



Retrieve



$$d(q, \mathbf{C})$$

Marqo Embeddings Cloud



Train

Multimodal fine-tuning, optimised for user interaction data.



Generate

Marqo's open source inference engine provides fast image ingestion, inference and search.



Retrieve

Marqo provides a fully-managed Vespa vectorDB engine, deployed with Marqo inference nodes for low latency search and ingestion.

```
model_details =
marqtune_client.train(
    method="weighted_clip",

model="open_clip/ViT-L-14/laion2b_
s32b_b82k",
    data={
        "id": "your_data_uuid",
        "bucket_path":
"path_in_bucket"
    },
    params={
        "param1": "value1",
        "param2": "value2"
    }
)
```

```
mq.create_index("my-index",
model="my-custom-CLIP")

mq.index("my-index").add_documents
(
    [{
        "title": "The Travels of
Marco Polo",
        "image": "my-image-url"
    }],
    tensor_fields=["title",
"image"]
)
```

```
results =
mq.index("my-index").search(
    "Marqo Polo"
)
```

Building Software 2.0 with embeddings

Outline

- What is Software 2.0?
 - Examples
- What are embeddings?
 - Producing embeddings
 - Training models for embeddings
- Building with embeddings
 - Classification, search and RAG

Conclusion

- Software 2.0:
 - Uses optimization to learn functions from data
 - Allows for increasingly complex programs
- Embeddings:
 - The data structures 2.0
 - Uses optimization to learn representations from data
 - Interface with unstructured data
 - Allow for composability
- Embedding models share a lot of similarities with each other
- Embedding based systems share a lot of similarities with each other

$$\operatorname{argmin}_{N_{\theta}, \theta} L(\mathbf{X}, \mathbf{Y}; N_{\theta})$$



Thank you!

jesse@marqo.ai

Marqo

<https://www.marqo.ai/>

Open source

<https://github.com/marqo-ai/marqo>

Classification demo

<https://huggingface.co/spaces/Marqo/Marqo-FashionSigLIP-Classification>

FashionCLIP

<https://huggingface.co/Marqo/marqo-fashionSigLIP>

Blog

<https://www.marqo.ai/blog>

