



Berlin · Beijing · Shenzhen

Neural Search in Action

Representing, transiting & searching **multimodal data**

Han Xiao, Founder of Jina AI



@hxiao



@JinaAI_

About me & Jina AI

Han Xiao, Founder & CEO of Jina AI. Based in Berlin, Germany.

- ML PhD in 2014 TU Munich; Zalando Research; Tencent AI Lab; Creator of Fashion-MNIST.

Jina AI

- Founded in 2020, focus on multimodal AI search & create
- Opensource contributor: Jina, **DocArray (Linux Foundation)**, CLIP-as-service, ...
- 60 people, HQ in Berlin. Offices in Beijing, Shenzhen.



Jina AI Tech Spectrum



Prompt tuning

the process of crafting and refining the input prompts in order to guide its output towards specific, desired responses.

the deployment of fine-tuned models in a production environment, usually requiring substantial resources such as GPU hosting. MLOps, emphasizing the serving of mid-size to large models in a scalable, efficient, and reliable manner.

Model serving

Prompt serving

wrapping and serving prompts through an API, without hosting heavy models. The API calls a public large language model service and handles the orchestration of inputs and outputs in a chain of operations.

Also known as fine-tuning, involves adjusting the parameters of a pre-trained model on a new, often task-specific dataset to improve its performance and adapt it to a specific application.

Model tuning



Prompt tuning

Faster time2market



PromptPerfect



- LLMSearch ◆ DevGPT
- Ensemble ◆ ThinkGPT
- ArtiBanner ◆ AgentChain
- SceneX

Model serving



Inference API

▲ Jina+DocArray

▲ LC-serve

▲ Jcloud

◆ DiscoArt

▲ OpenGPT

◆ DalleFlow

◆ CLIP-as-service

◆ DocsQA

● Rationale

Prompt serving

Long-term invest



Model tuning



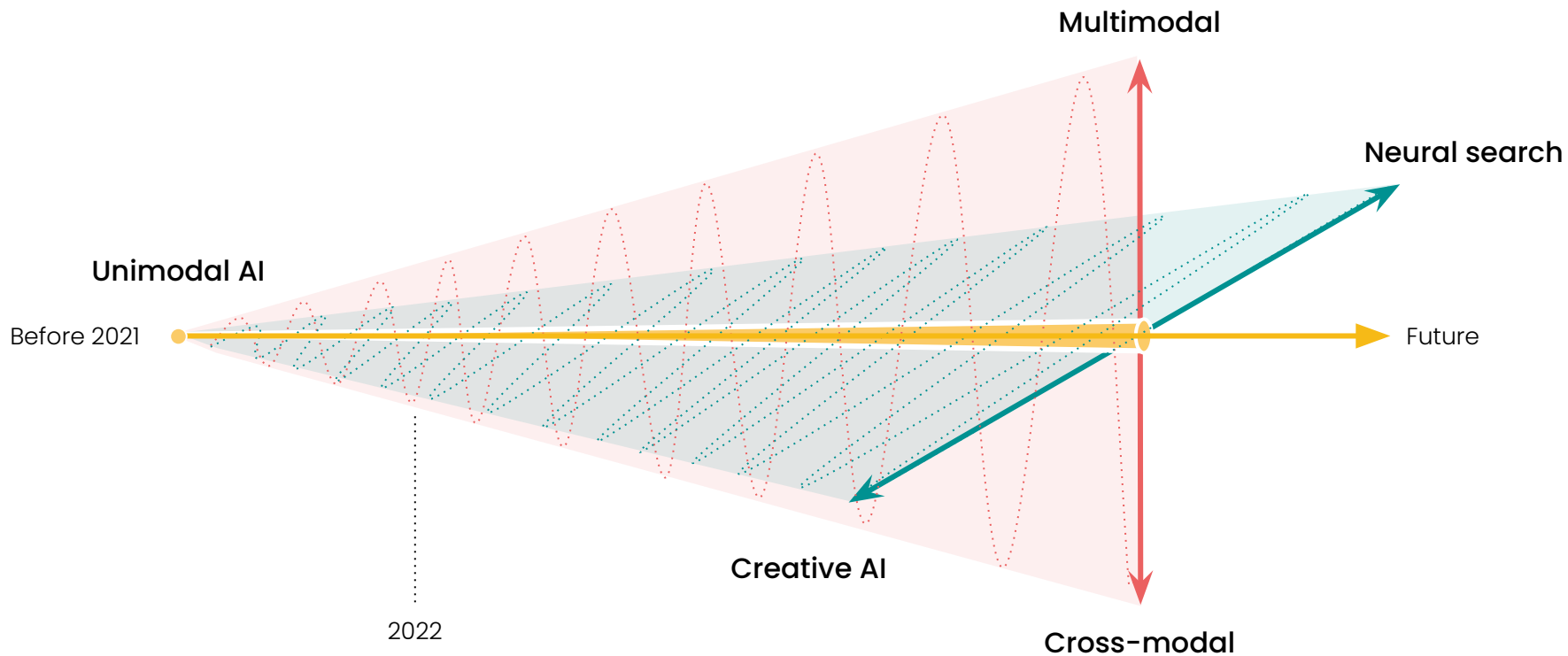
Finetuner

Agenda

- **Preliminary: multimodal AI**
- Opensource package: DocArray
 - Motivation
 - Representing data
 - Transiting data
 - Storing data
 - Retrieving data
- Multimodal at scale in production

This tutorial may require technical knowledge. Familiarity with Python 3.7+ concepts like data classes could be helpful.

Preliminary: from unimodal to multimodal



From unimodal to multimodal

"**modality**" roughly means "**data type**".

- Unimodal AI refers to applying AI to one specific type of data.
- Most early machine learning works fall into this category.
- Even today, when you open any machine learning literature, unimodal AI is still the majority of the content.

Unimodal – NLP

LDA was the 2010's transformer

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

JMLR: Workshop and Conference Proceedings 13: 63-78
2nd Asian Conference on Machine Learning (ACML2010), Tokyo, Japan, Nov. 8–10, 2010.

Efficient Collapsed Gibbs Sampling For Latent Dirichlet Allocation

Han Xiao
Thomas Stibor

Department of Informatics
Technical University of Munich, GERMANY

XIAOH@IN.TUM.DE
STIBOR@IN.TUM.DE

Editor: Masashi Sugiyama and Qiang Yang

Abstract

Collapsed Gibbs sampling is a frequently applied method to approximate intractable integrals in probabilistic generative models such as latent Dirichlet allocation. This sampling method has however the crucial drawback of high computational complexity, which makes it limited applicable on large data sets. We propose a novel *dynamic sampling* strategy to significantly improve the efficiency of collapsed Gibbs sampling. The strategy is explored in terms of efficiency, convergence and perplexity. Besides, we present a straight-forward parallelization to further improve the efficiency. Finally, we underpin our proposed improvements with a comparative study on different scale data sets.

Keywords: Gibbs sampling, Optimization, Latent Dirichlet Allocation

1. Introduction

Latent Dirichlet allocation (LDA) is a generative probabilistic model that was first proposed by Blei et al. (2003) to discover topics in text documents. LDA is based on the

Unimodal tasks in NLP

Adhoc methods for NLP problems

Sentiment analysis

Text classification

Topic modeling

Text summarization

Natural language
generation

Named entity
recognition

Word sense
disambiguation

Parts-of-speech
tagging

Grammatical
parsing

Machine translation

Question answering

Spam filtering

Language modeling

Dialog systems

Information
extraction

Semantic role
labeling

Part-of-speech
induction

Co-reference
resolution

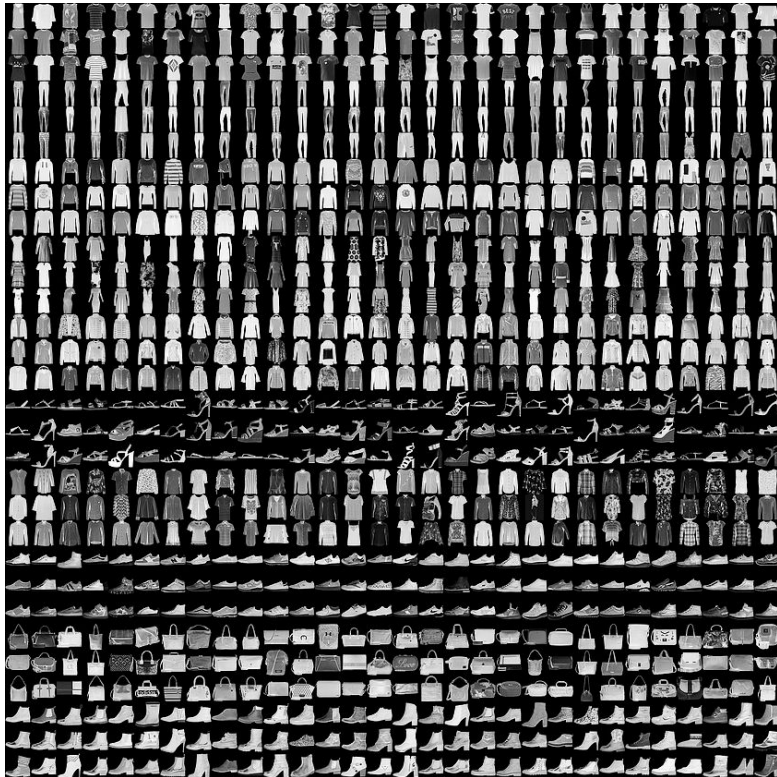
Pronoun resolution

Sentence
segmentation

Textual Modality

Unimodal – CV

Fashion-MNIST, 2017



07747v2 [cs.LG] 15 Sep 2017

Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms

Han Xiao
Zalando Research
Mühlenstraße 25, 10243 Berlin
han.xiao@zalando.de

Kashif Rasul
Zalando Research
Mühlenstraße 25, 10243 Berlin
kashif.rasul@zalando.de

Roland Vollgraf
Zalando Research
Mühlenstraße 25, 10243 Berlin
roland.vollgraf@zalando.de

Abstract

We present Fashion-MNIST, a new dataset comprising of 28×28 grayscale images of 70,000 fashion products from 10 categories, with 7,000 images per category. The training set has 60,000 images and the test set has 10,000 images. Fashion-MNIST is intended to serve as a direct drop-in replacement for the original MNIST dataset for benchmarking machine learning algorithms, as it shares the same image size, data format and the structure of training and testing splits. The dataset is freely available at <https://github.com/zalando-research/fashion-mnist>.

Unimodal tasks in CV

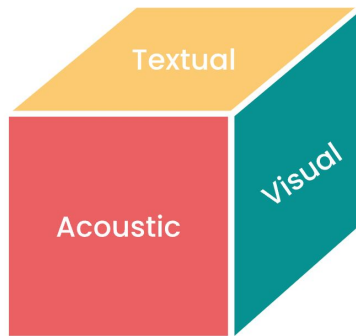
Object classification and detection	Image segmentation	Object tracking	Action recognition	Scene understanding
3D reconstruction	Pose estimation	Depth estimation	Stereo vision	Texture recognition and classification
Material recognition	Object recognition in video	Facial recognition and identification	Human activity recognition	Image super-resolution
Neural style transfer	Image inpainting	Video frame interpolation	Multiple object tracking in 3D	SLAM

Visual Modality

Unimodal tasks in speech & audio

Automatic Speech Recognition	Text-to-Speech	Speaker Recognition	Speaker Diarization	Speech Enhancement
Music Recommendation	Music Genre Recognition	Music Artist Recognition	Music Structure Segmentation	Music Tempo Estimation
Audio Source Separation	Sound Event Detection	Sound Event Classification	Sound Event Localization	Audio Scene Recognition
Audio Captioning	Emotion Recognition	Speech Translation	Voice Activity Detection	Silence Detection
Acoustic Modality				

Unimodal know-how are hardly transferable



- Tasks are specific to just one modality (e.g. textual, visual, acoustic, etc).
- Knowledge is learned from and applied to only one modality (i.e. a visual algorithm can only learn from and be applied to images).

Sentiment analysis	Text classification	Topic modeling	Text summarization	Natural language generation
Named entity recognition	Word sense disambiguation	Parts-of-speech tagging	Grammatical parsing	Machine translation
Question answering	Spam filtering	Language modeling	Dialog systems	Information extraction
Semantic role labeling	Part-of-speech induction	Co-reference resolution	Pronoun resolution	Sentence segmentation
Textual Modality				
Object classification and detection	Image segmentation	Object tracking	Action recognition	Scene understanding
3D reconstruction	Pose estimation	Depth estimation	Stereo vision	Texture recognition and classification
Material recognition	Object recognition in video	Facial recognition and identification	Human activity recognition	Image super-resolution
Neural style transfer	Image inpainting	Video frame interpolation	Multiple object tracking in 3D	SLAM
Visual Modality				
Automatic Speech Recognition	Text-to-Speech	Speaker Recognition	Speaker Diarization	Speech Enhancement
Music Recommendation	Music Genre Recognition	Music Artist Recognition	Music Structure Segmentation	Music Tempo Estimation
Audio Source Separation	Sound Event Detection	Sound Event Classification	Sound Event Localization	Audio Scene Recognition
Audio Captioning	Emotion Recognition	Speech Translation	Voice Activity Detection	Silence Detection
Acoustic Modality				

A detour: cross-modal model

NIPS 2010, Cross-LDA

Toward Artificial Synesthesia: Linking Images and Sounds via Words

Han Xiao, Thomas Stibor
Department of Informatics
Technical University of Munich
Garching, D-85748
{xiaoh, stibor}@in.tum.de

Abstract

We tackle a new challenge of modeling a perceptual experience in which a stimulus in one modality gives rise to an experience in a different sensory modality, termed *synesthesia*. To meet the challenge, we propose a probabilistic framework based on graphical models that enables to link visual modalities and auditory modalities via natural language text. An online prototype system is developed for allowing human judgement to evaluate the model's performance. Experimental results indicate usefulness and applicability of the framework.

1 Introduction

A picture of a golden beach might stimulate human's hearing, probably, by imagining the sound of waves crashing against the shore. On the other hand, the sound of a baaing sheep might illustrate a green hillside in front of your eyes. In neurology, this kind of experience is termed *synesthesia*. That is, a perceptual experience in which a stimulus in one modality gives rise to an experience in a different sensory modality. Without a doubt, the creative process of humans (e.g. painting and composing) is to a large extent attributed to their synesthesia experiences. While cross-sensory links such as sound and vision are quite common to humans, machines do not possess the same

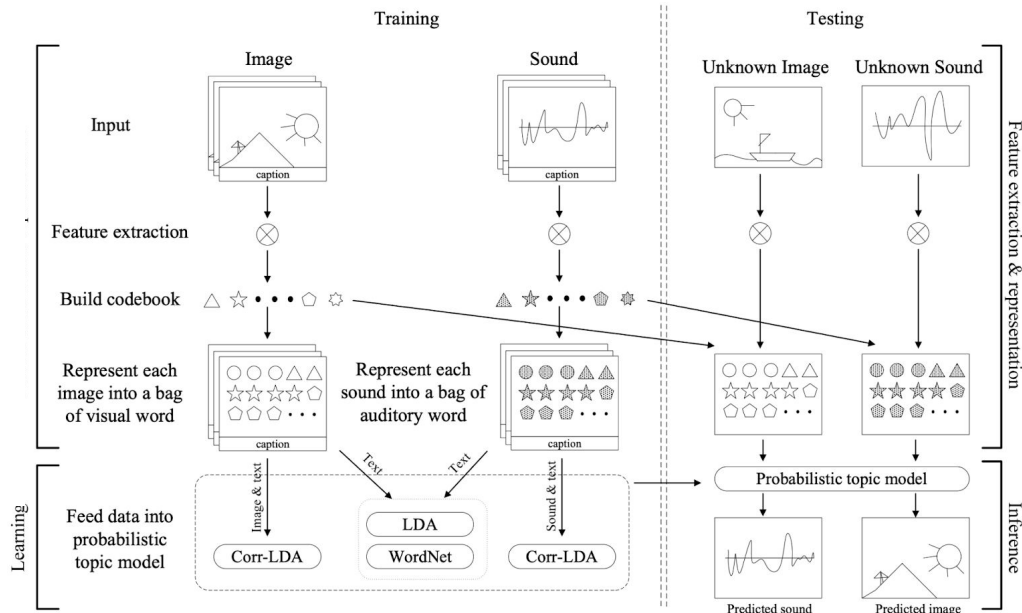
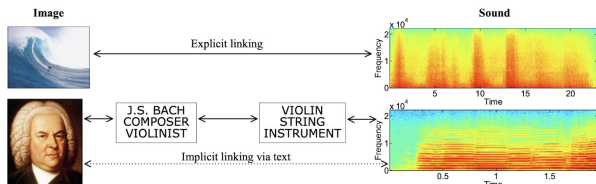


Figure 2: Probabilistic framework for performing the image composition and sound illustration task. The framework is an extension based on the work flow proposed in [8]. Images and sounds are represented in bags-of-words, so that the difference between the two modalities can be omitted. Once we have the algorithm for inferring sounds from an image, we can apply it to infer images from a sound by mirroring the algorithm.

Erase the boundary between modalities



- Tasks are shared and transferred between multiple modalities (so one algorithm can work with images and text and audio).
- Knowledge is learned from and applied to multiple modalities (so an algorithm can learn from textual data and apply that to visual data).

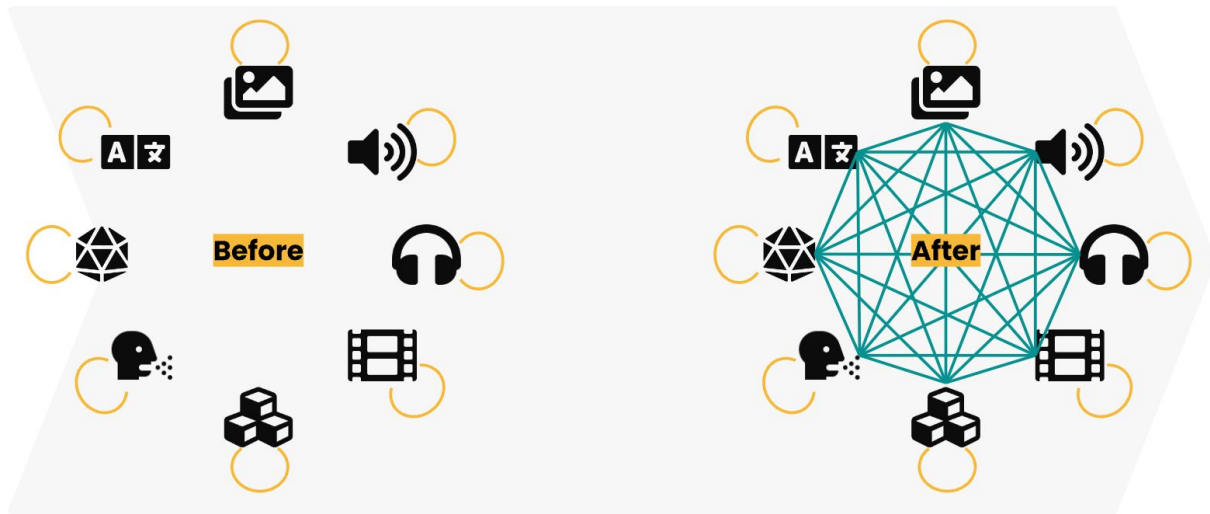
Paradigm shift from unimodal to multimodal

The rise of multimodal AI can be attributed to advances in two machine learning techniques: **Representation learning** and **transfer learning**.

- Representation learning lets models create common representations for all modalities.
- Transfer learning lets models first learn fundamental knowledge, and then fine-tune on specific domains.

CLIP, DALLE, BLIP, Bark, GPT4

We will see more and more AI applications move beyond one data modality and leverage relationships between different modalities



The paradigm shift from single-modal AI to multimodal AI

“An artificial intelligence system trained on words and sentences alone will never approximate human understanding.”

Y. Lecun in 2022 in AI And The Limits Of Language

**Multimodal AI is the future,
but the ML ecosystem is not yet
suited for it.**

Agenda

- Preliminary: multimodal AI
- **Opensource package: DocArray**
 - Motivation
 - Representing data
 - Transiting data
 - Storing data
 - Retrieving data
- Multimodal at scale in production

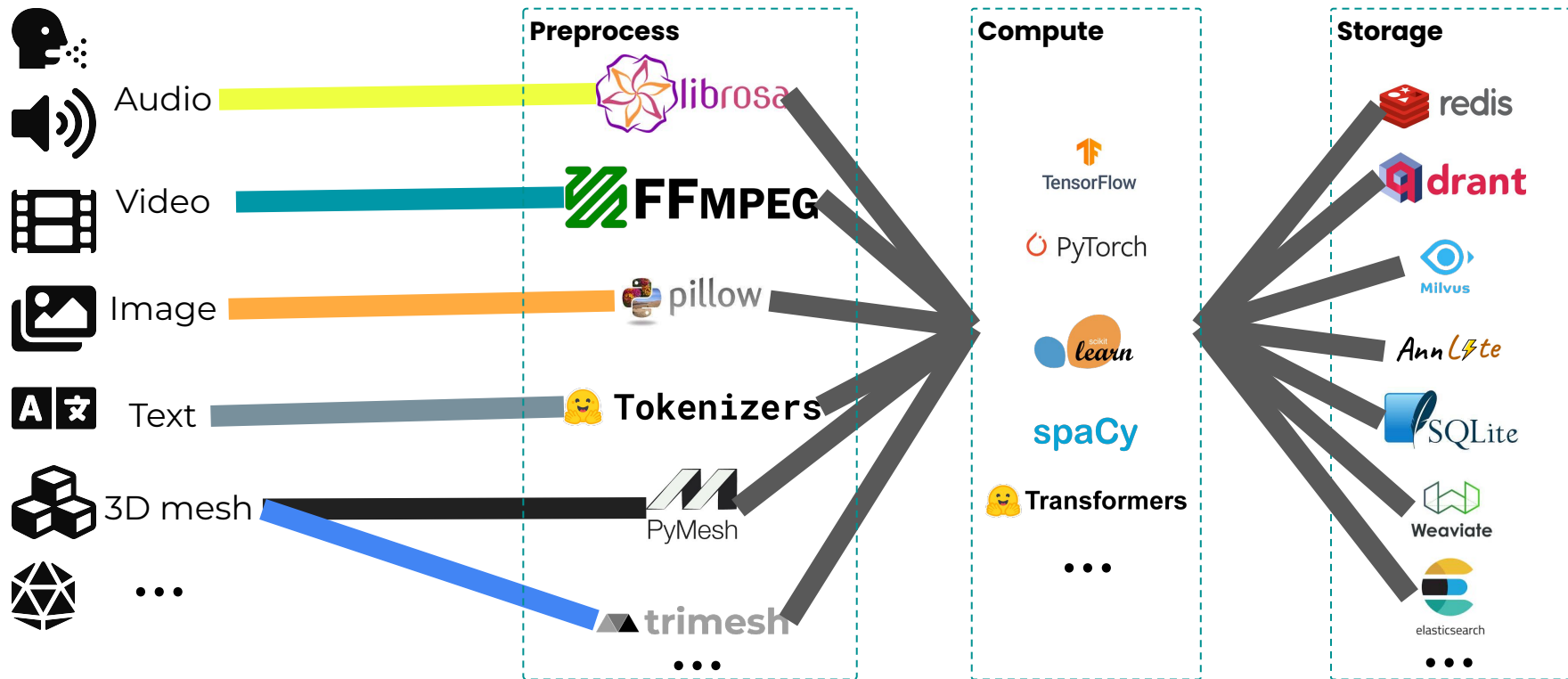
This tutorial may require technical knowledge. Familiarity with Python 3.7+ concepts like data classes could be helpful.

**DocArray for
representing, transiting,
storing, searching
multimodal data**

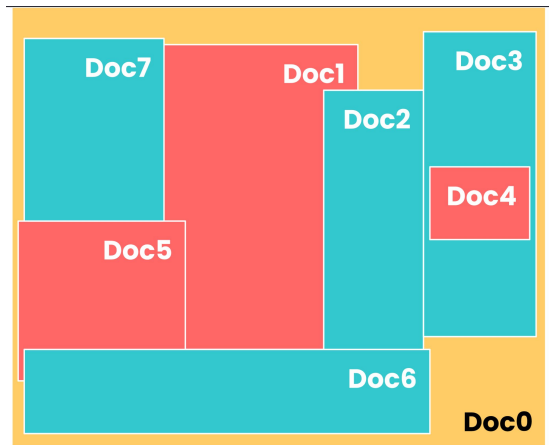
Representing multimodal data is a pain

- Lack of common interface for different modalities makes it difficult to work with multiple modalities at the same time.
- No easy way to represent unstructured and nested multimodal data.

Lack of common interface

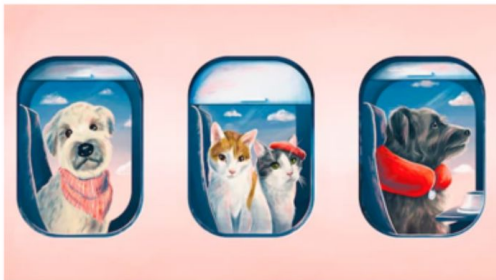


No easy way to represent unstructured nested multimodal data



- Unstructured document
- Nested content
- Different modalities (text, image, ...)

DocArray way of representing multimodal data



By the Way A Post Travel Destination

Everything to know about flying
with pets, from picking your seat to
keeping your animal calm

By Nathan Diller

```
from docarray import dataclass, Document
from docarray.typing import Image, Text, ...
```

```
@dataclass
class WPArticle:
    banner: Image
    headline: Text
    meta: JSON
```

```
a = WPArticle(
    banner='dog-cat-flight.png',
    headline='Everything to know about fl
    meta={
        'author': 'Nathan Diller',
        'column': 'By the Way - A Post Tr
    },
)

doc = Document(a)
```

Frequent data transfer over network is expensive

Multimodal data is processed by multiple models and models are usually deployed in a distributed way.

Data at rest	Data in use	Data in transit
Inactive data under very occasional changes, stored physically in database, warehouse, spreadsheet, archives, etc.	Active data under constant change, stored physically in database, warehouse, spreadsheet, etc.	Traversing a network or temporarily residing in computer memory to be read or updated

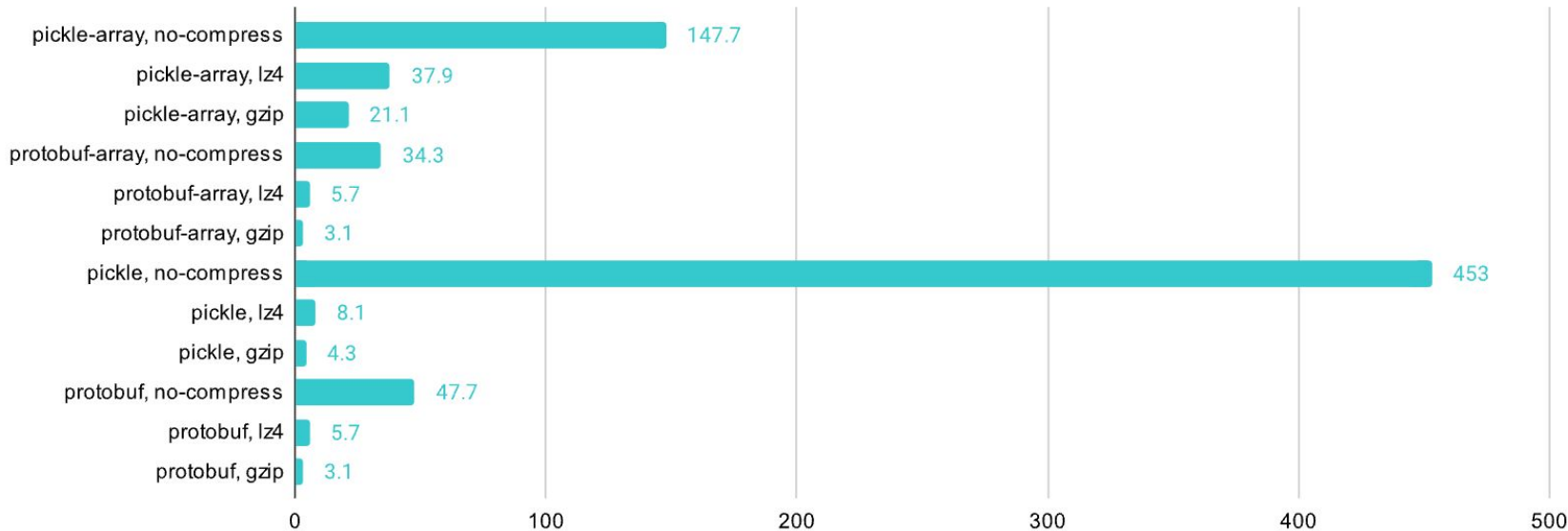
Performant serialization is important

DocArray is designed to be “ready-to-wire” at anytime.

- JSON string: `.from_json() / .to_json()`
 - Pydantic model: `.from_pydantic_model() / .to_pydantic_model()`
- Bytes (compressed): `.from_bytes() / .to_bytes()`
 - Disk serialization: `.save_binary() / .load_binary()`
- Base64 (compressed): `.from_base64() / .to_base64()`
- Protobuf Message: `.from_protobuf() / .to_protobuf()`
- Python List: `.from_list() / .to_list()`
- Pandas Dataframe: `.from_dataframe() / .to_dataframe()`
- Cloud: `.push() / .pull()`

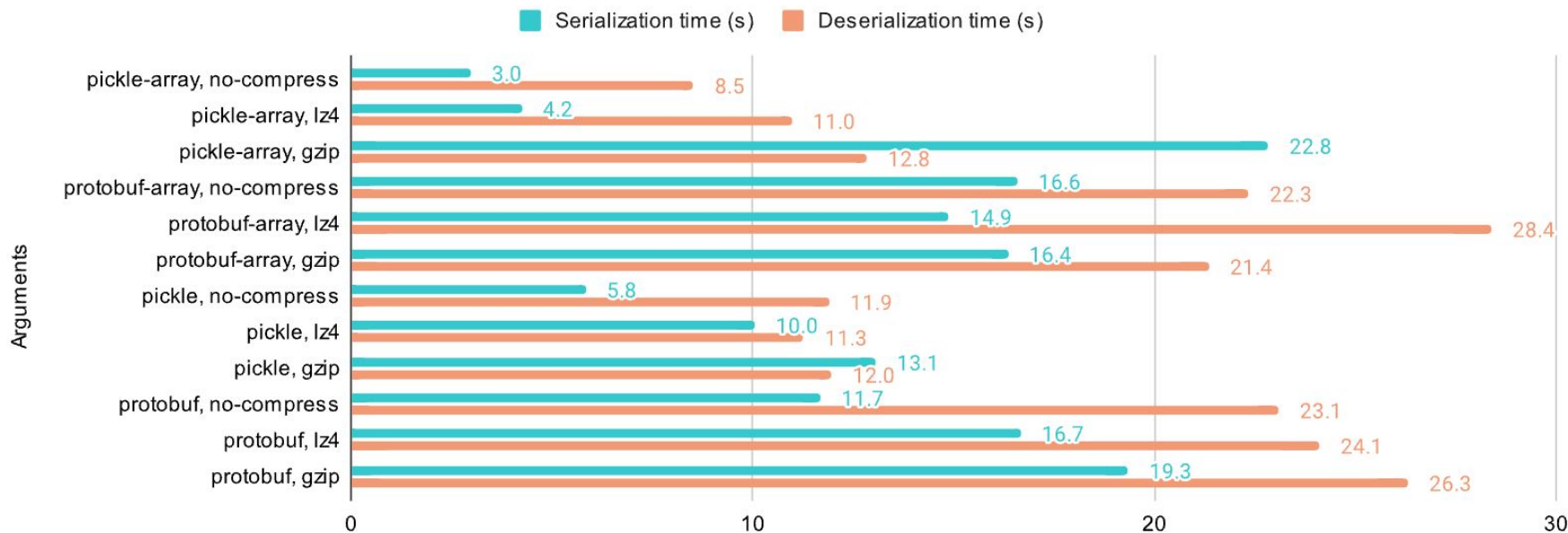
Binary serialization optimized for in-transit & at-rest

Size in MB on 1M Docs



Binary serialization optimized for in-transit & at-rest

Time cost in seconds on 1M Docs



Storing nested data with databases is complicated

- Complex and nested schema are not directly supported in databases
- Explosion in numbers of vector databases with different APIs but no universal client

DocArray way of storing data

DocArray Storage

```
1 from docarray import DocumentArray, Document
2
3 da = DocumentArray(storage='milvus',
4                     config={'connection': 'example.db'})
5
6 with da:
7     da.append(Document())
8 da.summary()
```

DocArray way of storing data

DocArray Storage

```
1 from docarray import DocumentArray, Document
2
3 da = DocumentArray(storage='milvus',
4                     config={'connection':
5
6 with da:
7     da.append(Document())
8 da.summary()
```

```
'milvus'
'qdrant'
'weaviate'
'elasticsearch'
'redis'
'opensearch'
'annlite'
'sqlite'
```

Vector Search via a consistent API

```
1 from docarray import Document, DocumentArray
2 import numpy as np
3
4 n_dim = 3
5 da = DocumentArray(
6     storage='annlite',
7     config={'n_dim': n_dim, 'metric': 'Euclidean'},
8 )
9
10 with da:
11     da.extend([Document(embedding=i * np.ones(n_dim)) for i in range(10)])
12
13 result = da.find(np.array([2, 2, 2]), limit=6)
14 result[:, 'embedding']
```

Vector Search via a consistent API

```

1 from docarray import Document, DocumentArray
2 import numpy as np
3
4 n_dim = 3
5 da = DocumentArray(
6     storage='annlite',
7     config={'n_dim': n_dim, 'metric': 'Euclid
8 )
9
10 with da:
11     da.extend([Document(embedding=i * np.ones
12
13 result = da.find(np.array([2, 2, 2]), limit=6
14 result[:, 'embedding']

```

Name	Construction	Vector search	Vector search + Filter	Filter
In memory	<code>DocumentArray()</code>	✓	✓	✓
SQLite	<code>DocumentArray(storage='sqlite')</code>	✗	✗	✓
Weaviate	<code>DocumentArray(storage='weaviate')</code>	✓	✓	✓
Qdrant	<code>DocumentArray(storage='qdrant')</code>	✓	✓	✓
AnnLite	<code>DocumentArray(storage='annlite')</code>	✓	✓	✓
ElasticSearch	<code>DocumentArray(storage='elasticsearch')</code>	✓	✓	✓
Redis	<code>DocumentArray(storage='redis')</code>	✓	✓	✓
Milvus	<code>DocumentArray(storage='milvus')</code>	✓	✓	✓

Quick Recap

- It's like JSON, but for intensive computation.
- It's like `numpy.ndarray`, but for unstructured data.
- It's like `pandas.DataFrame`, but for nested and mixed media **data with embeddings**.
- It's like Protobuf, but for data scientists and deep learning engineers.



Quick R

- It's like JSON
- It's like `numpy`
- It's like `pandas`
- It's like Protobuf

	DocArray	<code>numpy.ndarray</code>	JSON	<code>pandas.DataFrame</code>	Protobuf
Tensor/matrix data	✓	✓	✗	✓	✓
Text data	✓	✗	✓	✓	✓
Media data	✓	✗	✗	✗	✗
Nested data	✓	✗	✓	✗	✓
Mixed data of the above four	✓	✗	✗	✗	✗
Easy to (de)serialize	✓	✗	✓	✓	✓
Data validation (of the output)	✓	✗	✗	✗	✓
Pythonic experience	✓	✓	✗	✓	✗
IO support for filetypes	✓	✗	✗	✗	✗
Deep learning framework support	✓	✓	✗	✗	✗
multi-core/GPU support	✓	✓	✗	✗	✗
Rich functions for data types	✓	✗	✗	✓	✗

ings.

Hands-on DocArray

Install DocArray

To install `DocArray (0.33)`, you can use the following command:

```
pip install "docarray[full]"
```

<https://docs.docarray.org/>

For old DocArray, more compatibility and features

```
pip install "docarray[full]"==0.21
```


Representing data – Document

At the heart of `DocArray` lies the concept of `BaseDoc`.

The following Python code defines a `BannerDoc` class that can be used to represent the data of a website banner:

```
from docarray import BaseDoc
from docarray.typing import ImageUrl

class BannerDoc(BaseDoc):
    image_url: ImageUrl
    title: str
    description: str
```



Representing data – Document

You can then instantiate a `BannerDoc` object and access its attributes:

```
banner = BannerDoc(  
    image_url='https://example.com/image.png',  
    title='Hello World',  
    description='This is a banner',  
)  
  
assert banner.image_url == 'https://example.com/image.png'  
assert banner.title == 'Hello World'  
assert banner.description == 'This is a banner'
```



Representing multimodal data with nested structure

Let's say you want to represent a YouTube video in your application, perhaps to build a search system for YouTube videos.

A YouTube video is not only composed of a video, but also has a title, description, thumbnail (and more, but let's keep it simple).

All of these elements are from different modalities:

the title and description are text,

the thumbnail is an image,

and the video itself is, well, a video.

DocArray lets you represent all of this multimodal data in a single object.



Year in Review: 2021 in Graphic Design

Linus Boman ✓

119K views • 1 year ago

Representing multimodal data with nested structure

First for the thumbnail image:

```
from docarray import BaseDoc
from docarray.typing import ImageUrl, ImageBytes

class ImageDoc(BaseDoc):
    url: ImageUrl
    bytes: ImageBytes = (
        None # bytes are not always loaded in memory, so we make it optional
    )
```

Representing multimodal data with nested structure

Then for the video itself:

```
from docarray import BaseDoc
from docarray.typing import VideoUrl, VideoBytes

class VideoDoc(BaseDoc):
    url: VideoUrl
    bytes: VideoBytes = (
        None # bytes are not always loaded in memory, so we make it optional
    )
```

Representing multimodal data with nested structure

All the elements that compose a YouTube video are ready:

```
from docarray import BaseDoc

class YouTubeVideoDoc(BaseDoc):
    title: str
    description: str
    thumbnail: ImageDoc
    video: VideoDoc
```

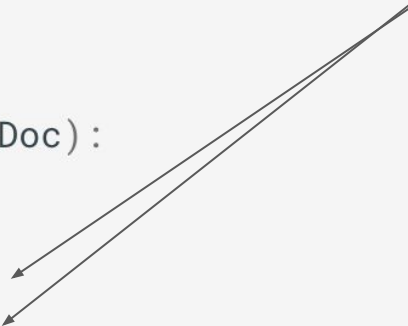


Representing multimodal data with nested structure

All the elements that compose a YouTube video are ready:

```
from docarray import BaseDoc

class YouTubeVideoDoc(BaseDoc):
    title: str
    description: str
    thumbnail: ImageDoc
    video: VideoDoc
```

Two arrows originate from the right side of the code block. One arrow points from the `ImageDoc` type annotation in the `thumbnail` field to the `ImageDoc` class name in the explanatory text box. The other arrow points from the `VideoDoc` type annotation in the `video` field to the `VideoDoc` class name in the same text box.

You see here that `ImageDoc` and `VideoDoc` are also `BaseDoc`, and they are later used inside another `BaseDoc`. This is what we call nested data representation.

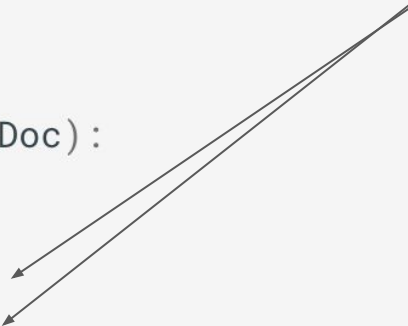
`BaseDoc` can be nested to represent any kind of data hierarchy.

Representing multimodal data with nested structure

All the elements that compose a YouTube video are ready:

```
from docarray import BaseDoc

class YouTubeVideoDoc(BaseDoc):
    title: str
    description: str
    thumbnail: ImageDoc
    video: VideoDoc
```



You see here that `ImageDoc` and `VideoDoc` are also `BaseDoc`, and they are later used inside another `BaseDoc`. This is what we call nested data representation.

`BaseDoc` can be nested to represent any kind of data hierarchy.

This representation can be used to `send` or `store` data. You can even use it directly to `train a machine learning Pytorch` model on this representation.

Recap: representing multimodal data

- "Dataclass" look and feel, for defining the structure
- Strong typing, for defining modality

- Python built-in types
- Numpy types
- URI types

- **Text**
- **Image**
- **Audio**
- **Video**
- **Mesh3D**
- **PointCloud3D**

- Tensor types

- **ImageTensor**
- **AudioTensor**
- **VideoTensor**
- **Embedding**

- **Optional[]**

```
from docarray import BaseDoc
from docarray.typing import ImageUrl, ImageBytes
```

```
class ImageDoc(BaseDoc):
    url: ImageUrl
    bytes: ImageBytes = (
        None # bytes are not always loaded in m
```

Representing an array of multimodal data

The fundamental building block of DocArray is the `BaseDoc` class which represents a *single* document, a *single* datapoint.

However, in machine learning we often need to work with an *array* of documents, and an *array* of data points.

We introduce

- `DocList` which is a **Python list** of `BaseDocs`
- `DocVec` which is a **column-based representation** of `BaseDocs`

Example of DocList

First you need to create a `Doc` class, our data schema. Let's say you want to represent a banner with an image, a title and a description:

```
from docarray import BaseDoc, DocList
from docarray.typing import ImageUrl
```

```
class BannerDoc(BaseDoc):
    image: ImageUrl
    title: str
    description: str
```

Example of DocList

First you need to create a `Doc` class, our data schema. Let's say you want to represent a banner with an image, a title and a description:

```
from docarray import BaseDoc, DocList
from docarray.typing import ImageUrl
```

Let's instantiate several `BannerDoc` s:

```
banner1 = BannerDoc(
    image='https://example.com/image1.png',
    title='Hello World',
    description='This is a banner',
)

banner2 = BannerDoc(
    image='https://example.com/image2.png',
    title='Bye Bye World',
    description='This is (distopic) banner',
)
```

Example of DocList

`DocList` and `DocVec` are both `AnyDocArrays`. The following section will use `DocList` as an example, but the same applies to `DocVec`.

You can now collect them into a `DocList` of `BannerDoc` s:

```
docs = DocList[BannerDoc]([banner1, banner2])  
  
docs.summary()
```



```
DocList Summary  
  
Type      DocList[BannerDoc]  
Length    2
```



```
Document Schema  
  
BannerDoc  
├ image: ImageUrl  
├ title: str  
└ description: str
```

Example of DocList

You can access documents inside it with the usual Python array API:

```
print(docs[0])
```

```
BannerDoc(image='https://example.com/image1.png', title='Hello World', description:
```

or iterate over it:

```
for doc in docs:  
    print(doc)
```

```
BannerDoc(image='https://example.com/image1.png', title='Hello World', description:  
BannerDoc(image='https://example.com/image2.png', title='Bye Bye World', descriptio
```

Accessing member attribute at array level

At the document level:

```
print(banner1.image)
```

```
'https://example.com/image1.png'
```

At the Array level:

```
print(docs.image)
```

```
['https://example.com/image1.png', 'https://example.com/image2.png']
```

Accessing member attribute at array level

At the document level:

```
print(banner1.image)
```

```
https://example.com/image1.png'
```

At the Array level:

```
print(docs.image)
```

```
['https://example.com/image1.png', 'https://example.com/image2.png']
```

You can even access the attributes of the nested `BaseDoc` at the Array level:

```
print(docs.banner.image)
```

```
['https://example.com/image1.png', 'https://example.com/image2.png']
```

This is just the same way that you would do it with `BaseDoc`:

```
print(page1.banner.image)
```

```
'https://example.com/image1.png'
```


DocList[DocType] syntax

`DocList[DocType]` creates a custom `DocList` that can only contain `DocType` Documents.

Non-typing DocList for heterogeneous data

```
from docarray import BaseDoc, DocList
from docarray.typing import ImageUrl, AudioUrl

class ImageDoc(BaseDoc):
    url: ImageUrl

class AudioDoc(BaseDoc):
    url: AudioUrl

docs = DocList(
    [
        ImageDoc(url='https://example.com/image1.png'),
        AudioDoc(url='https://example.com/audio1.mp3'),
    ]
)
```

Strong-typing DocList for homogeneous data

```
try:
    docs = DocList[ImageDoc](
        [
            ImageDoc(url='https://example.com/image1.png'),
            AudioDoc(url='https://example.com/audio1.mp3'),
        ]
    )
except ValueError as e:
    print(e)
```

```
ValueError: AudioDoc(
  id='e286b10f58533f48a0928460f0206441',
  url=AudioUrl('https://example.com/audio1.mp3', host_type='domain')
) is not a <class '__main__.ImageDoc'>
```

DocList vs DocVec

DocList is based on Python Lists. You can append, extend, insert, pop, and so on. In DocList, data is individually owned by each **BaseDoc** collect just different Document references.

Use **DocList** when you want to be able to rearrange or re-rank your data. One flaw of **DocList** is that none of the data is contiguous in memory, so you cannot leverage functions that require contiguous data without first copying the data in a continuous array.

DocVec is a columnar data structure. **DocVec** is always an array of homogeneous Documents. The idea is that every attribute of the **BaseDoc** will be stored in a contiguous array: a column.

DocList vs DocVec

Let's say you want to embed a batch of Images:

```
def embed(image: NdArray['batch_size', 3, 224, 224]):
```

```
    ...
```

DocList vs DocVec

```
from docarray import BaseDoc
from docarray.typing import NdArray

class ImageDoc(BaseDoc):
    image: NdArray[
        3, 224, 224
    ] = None # [3, 224, 224] this just mean we know in advance the shape of the tensor
```

DocList vs DocVec

```
from docarray import BaseDoc
from docarray.typing import NdArray
```

```
class ImageDoc(BaseDoc):
    image: NdArray

docs = DocList[ImageDoc](
    [ImageDoc(image=np.random.rand(3, 224, 224)) for _ in range(10)]
)

embed (np.stack(docs.image))
...
embed (np.stack(docs.image))
```

DocList vs DocVec

```
from docarray import BaseDoc
from docarray.typing import NdArray
```

```
class ImageDoc(BaseDoc):
```

```
    image:
```

```
        3,
```

```
    ] = NdArray
```

```
1 from docarray import DocVec
```

```
2 import numpy as np
```

```
3
```

```
4 docs = DocVec[ImageDoc](
```

```
5     [ImageDoc(image=np.random.rand(3, 224, 224)) for _ in range(10)]
```

```
6 )
```

```
7
```

```
8 embed(docs.image)
```

Access the view of Document in DocVec

If you access a document inside a `DocVec` you will get a document view. A document view is a view of the columnar data structure which looks and behaves like a `BaseDoc` instance. It is a `BaseDoc` instance but with a different way to access the data.

```
from docarray import DocVec

docs = DocVec[ImageDoc](
    [ImageDoc(image=np.random.rand(3, 224, 224)) for _ in range(10)]
)

my_doc = docs[0]

assert my_doc.is_view() # True
```

whereas with `DocList`:

```
docs = DocList[ImageDoc](
    [ImageDoc(image=np.random.rand(3, 224, 224)) for _ in range(10)]
)

my_doc = docs[0]

assert not my_doc.is_view() # False
```

Access the view of Document in DocVec

If you access a document inside a `DocVec` you will get a document view. A document view is a view of the columnar data structure which looks and behaves like a `BaseDoc` instance. It is a `BaseDoc` instance but with a different way to access the data.

you should use `DocVec` when you need to work with contiguous data, and you should use `DocList` when you need to rearrange or extend your data.

```
docs = DocList[ImageDoc](
    [ImageDoc(image=np.random.rand(3, 224, 224)) for _ in range(10)]
)

my_doc = docs[0]

assert not my_doc.is_view() # False
```


Storing & retrieving via Vector Database

```
1 from docarray import DocList, BaseDoc
2 from docarray.index import HnswDocumentIndex
3 import numpy as np
4
5 from docarray.typing import ImageUrl, ImageTensor, NdArray
6
7
8 class ImageDoc(BaseDoc):
9     url: ImageUrl
10    tensor: ImageTensor
11    embedding: NdArray[128]
12
13
14 # create some data
15 dl = DocList[ImageDoc](
16     [
17         ImageDoc(
18             url="https://upload.wikimedia.org/wikipedia/commons/2/2f/Alpamayo.jpg",
19             tensor=np.zeros((3, 224, 224)),
20             embedding=np.random.random((128,)),
21         )
22         for _ in range(100)
23     ]
24 )
25
26 # create a Document Index
27 index = HnswDocumentIndex[ImageDoc](work_dir='/tmp/test_index2')
28
29
30 # index your data
31 index.index(dl)
32
33 # find similar Documents
34 query = dl[0]
35 results, scores = index.find(query, limit=10, search_field='embedding')
```

Storing & retrieving via Vector Database

```
1 from docarray import DocList, BaseDoc
2 from docarray.index import HnswDocumentIndex
3 import numpy as np
4
5 from docarray.typing import ImageUrl, ImageTensor, NdArray
6
7
8 class ImageDoc(BaseDoc):
9     url: ImageUrl
10     tensor: ImageTensor
11     embedding: NdArray[128]
12
13
14 # create some data
15 dl = DocList[ImageDoc](
16     [
17         ImageDoc(
18             url="https://upload.wikimedia.org/wikipedia/commons/2/2f/Alpamayo.jpg",
19             tensor=np.zeros((3, 224, 224)),
20             embedding=np.random.random((128,)),
21         )
22         for _ in range(100)
23     ]
24 )
25
26 # create a Document Index
27 index = HnswDocumentIndex[ImageDoc](work_dir='/tmp/test_index2')
28
29
30 # index your data
31 index.index(dl)
32
33 # find similar Documents
34 query = dl[0]
35 results, scores = index.find(query, limit=10, search_field='embedding')
```



Document Index: ORM for vector DBs

Document Index provides a unified interface to a number of [vector databases](#).

You can think of Document Index as an **ORM for vector databases**.

Currently, DocArray supports the following vector databases:

- [Weaviate](#) | [Docs](#)
- [Qdrant](#) | [Docs](#)
- [Elasticsearch](#) v7 and v8 | [Docs](#)
- [HNSWlib](#) | [Docs](#)

*Old DocArray v0.21 supports Milvus, Redis, Opensearch

Construct a HNSWDocumentIndex

To use [HnswDocumentIndex](#), you need to install extra dependencies with the following command: `pip install "docarray[hnswlib]"`

To create a Document Index, you first need a document that defines the schema of your index:

```
from docarray import BaseDoc
from docarray.index import HnswDocumentIndex
from docarray.typing import NdArray

class MyDoc(BaseDoc):
    embedding: NdArray[128]
    text: str

db = HnswDocumentIndex[MyDoc](work_dir='./my_test_db')
```

Construct a HNSWDocumentIndex

To use `HnswDocumentIndex`, you need to install extra dependencies with the following command: `pip install "docarray[hnswlib]"`

To create a Document Index, you first need a document th

```
from docarray import BaseDoc
from docarray.index import HnswDocumentIndex
from docarray.typing import NdArray

class MyDoc(BaseDoc):
    embedding: NdArray[128]
    text: str

db = HnswDocumentIndex[MyDoc](work_dir='./my_test')
```

In this code snippet, `HnswDocumentIndex` takes a schema of the form of `MyDoc`. The Document Index then creates a column for each field in `MyDoc`.

Construct a HNSWDocumentIndex


To use `HnswDocumentIndex`, you need to install extra dependencies with the following command: `pip install "docarray[hnswlib]"`

To create a Document Index, you first need a document th

```
from docarray import BaseDoc
from docarray.index import HnswDocumentIndex
from docarray.typing import NdArray

class MyDoc(BaseDoc):
    embedding: NdArray[128]
    text: str

db = HnswDocumentIndex[MyDoc](work_dir='./my_test')
```



In this code snippet, `HnswDocumentIndex` takes a schema of the form of `MyDoc`. The Document Index then *creates a column for each field in `MyDoc`*.

The column types in the backend database are determined by the type hints of the document's fields. Optionally, you can *customize the database types for every field*.

Construct a HNSWDocumentIndex

To use `HnswDocumentIndex`, you need to install extra dependencies with the following command: `pip install "docarray[hnswlib]"`

To create a Document Index, you first need a document th

```
from docarray import BaseDoc
from docarray.index import HnswDocumentIndex
from docarray.typing import NdArray

class MyDoc(BaseDoc):
    embedding: NdArray[128]
    text: str

db = HnswDocumentIndex[MyDoc](work_dir='./my_test')
```

In this code snippet, `HnswDocumentIndex` takes a schema of the form of `MyDoc`. The Document Index then *creates a column for each field in `MyDoc`*.

The column types in the backend database are determined by the type hints of the document's fields. Optionally, you can [customize the database types for every field](#).

Most vector databases need to know the dimensionality of the vectors that will be stored. Here, that is automatically inferred from the type hint of the `embedding` field: `NdArray[128]` means that the database will store vectors with 128 dimensions.

Index data

Now that you have a Document Index, you can add data to it, using the `index()` method:

```
import numpy as np
from docarray import DocList

# create some random data
docs = DocList[MyDoc](
    [MyDoc(embedding=np.random.rand(128), text=f'text {i}') for i in range(100)]
)

# index the data
db.index(docs)
```


Index data

Now that you have a Document Index, you can add data to it, using the `index()` method:

```
import numpy as np
from docarray import DocList

# create some random data
docs = DocList[MyDoc](
    [MyDoc(embedding=np.random.rand(128), text=f'text {i}')
     ]
)

# index the data
db.index(docs)

from docarray import BaseDoc
from docarray.index import HnswDocumentIndex
from docarray.typing import NdArray

class MyDoc(BaseDoc):
    embedding: NdArray[128]
    text: str

db = HnswDocumentIndex[MyDoc](work_dir='./my_test_db')
```

As you can see, `DocList[MyDoc]` and `HnswDocumentIndex[MyDoc]` are both parameterized with `MyDoc`. This means that they share the same schema, and in general, the schema of a Document Index and the data that you want to store need to have compatible schemas

Vector search

[Search by Document](#)

[Search by raw vector](#)

```
# create a query Document
query = MyDoc(embedding=np.random.rand(128), text='query')

# find similar Documents
matches, scores = db.find(query, search_field='embedding', limit=5)

print(f'{matches=}')
print(f'{matches.text=}')
print(f'{scores=}')

```

Vector search

Search by Document

Search by raw vector

```
# create a query Document
query = MyDoc(embedding=np.random.rand(128), text='query')

# find similar Documents
matches, scores = db.find(query, search_field='embedding', limit=5)

print(f'{matches=}')
print(f'{matches.text=}')
print(f'{scores=}')

```

Search by Document

Search by raw vector

```
# create a query vector
query = np.random.rand(128)

# find similar Documents
matches, scores = db.find(query, search_field='embedding', limit=5)

print(f'{matches=}')
print(f'{matches.text=}')
print(f'{scores=}')

```

Vector search

Search by Document

Search by raw vector

```
# create a query Document
query = MyDoc(embedding=np.random.rand(128), text='query')

# find similar Documents
matches, scores = db.find(query, search_field='embedding', limit=5)

print(f'{matches=}')
print(f'{matches.text=}')
print(f'{scores=}')

```

Search by Document

Search by raw vector

```
# create a query vector
query = np.random.rand(128)

# find similar Documents
matches, scores = db.find(query, search_field='embedding', limit=5)

print(f'{matches=}')
print(f'{matches.text=}')
print(f'{scores=}')

```

Search by Documents

Search by raw vectors

```
# create some query Documents
queries = DocList[MyDoc](
    MyDoc(embedding=np.random.rand(128), text=f'query {i}') for i in range(3)
)

# find similar Documents
matches, scores = db.find_batched(queries, search_field='embedding', limit=5)

print(f'{matches=}')
print(f'{matches[0].text=}')
print(f'{scores=}')

```

Vector search

Search by Document

Search by raw vector

```
# create a query Document
query = MyDoc(embedding=np.random.rand(128), text='query')

# find similar Documents
matches, scores = db.find(query, search_field='embedding', limit=5)

print(f'{matches=}')
print(f'{matches.text=}')
print(f'{scores=}')
```

Search by Document

Search by raw vector

```
# create a query vector
query = np.random.rand(128)

# find similar Documents
matches, scores = db.find(query, search_field='embedding', limit=5)

print(f'{matches=}')
print(f'{matches.text=}')
print(f'{scores=}')
```

Search by Documents

Search by raw vectors

```
# create some query Documents
queries = DocList[MyDoc](
    MyDoc(embedding=np.random.rand(128), text=f'query {i}') for i in range(3)
)

# find similar Documents
matches, scores = db.find_batched(queries, search_field='embedding', limit=5)

print(f'{matches=}')
print(f'{matches[0].text=}')
print(f'{scores=}')
```

Search by Documents

Search by raw vectors

```
# create some query vectors
query = np.random.rand(3, 128)

# find similar Documents
matches, scores = db.find_batched(query, search_field='embedding', limit=5)

print(f'{matches=}')
print(f'{matches[0].text=}')
print(f'{scores=}')
```

Hybrid search through the query builder

Document Index supports atomic operations for vector similarity search, text search and filter search.

To combine these operations into a single, hybrid search query, you can use the query builder that is accessible through `build_query()`:

```
# prepare a query
q_doc = MyDoc(embedding=np.random.rand(128), text='query')

query = (
    db.build_query() # get empty query object
    .find(query=q_doc, search_field='embedding') # add vector similarity search
    .filter(filter_query={'text': {'$exists': True}}) # add filter search
    .build() # build the query
)

# execute the combined query and return the results
results = db.execute_query(query)
print(f'{results=}')

```

Customize vector DB configuration

```
db = HnswDocumentIndex[MyDoc](work_dir='/tmp/my_db')

db.configure(
    default_column_config={
        np.ndarray: {
            'dim': -1,
            'index': True,
            'space': 'ip',
            'max_elements': 2048,
            'ef_construction': 100,
            'ef': 15,
            'M': 8,
            'allow_replace_deleted': True,
            'num_threads': 5,
        },
        None: {},
    }
)
```

Indexing and searching multimodal data

In the following example you can see a complex schema that contains nested Documents. The `YouTubeVideoDoc` contains a `VideoDoc` and an `ImageDoc`, alongside some "basic" fields:



Year in Review: 2021 in Graphic Design

Linus Boman

119K views · 1 year ago

```

1 from docarray.typing import ImageUrl, VideoUrl, AnyTensor
2
3
4 # define a nested schema
5 class ImageDoc(BaseDoc):
6     url: ImageUrl
7     tensor: AnyTensor = Field(space='cosine', dim=64)
8
9
10 class VideoDoc(BaseDoc):
11     url: VideoUrl
12     tensor: AnyTensor = Field(space='cosine', dim=128)
13
14
15 class YouTubeVideoDoc(BaseDoc):
16     title: str
17     description: str
18     thumbnail: ImageDoc
19     video: VideoDoc
20     tensor: AnyTensor = Field(space='cosine', dim=256)
21
22
23 # create a Document Index
24 doc_index = HnswDocumentIndex[YouTubeVideoDoc](work_dir='/tmp2')
25
26 # create some data
27 index_docs = [
28     YouTubeVideoDoc(
29         title=f'video {i+1}',
30         description=f'this is video from author {10*i}',
31         thumbnail=ImageDoc(url=f'http://example.ai/images/{i}', tensor=np.ones(64)),
32         video=VideoDoc(url=f'http://example.ai/videos/{i}', tensor=np.ones(128)),
33         tensor=np.ones(256),
34     )
35     for i in range(8)
36 ]
37
38 # index the Documents
39 doc_index.index(index_docs)

```


Indexing and searching multimodal data

You can perform search on any nesting level by using the dunder operator to specify the field defined in the nested data.

```
1 # create a query Document
2 query_doc = YouTubeVideoDoc(
3     title=f'video query',
4     description=f'this is a query video',
5     thumbnail=ImageDoc(url=f'http://example.ai/images/1024', tensor=np.ones(64)),
6     video=VideoDoc(url=f'http://example.ai/videos/1024', tensor=np.ones(128)),
7     tensor=np.ones(256),
8 )
9
10 # find by the `youtubevideo` tensor; root level
11 docs, scores = doc_index.find(query_doc, search_field='tensor', limit=3)
12
13 # find by the `thumbnail` tensor; nested level
14 docs, scores = doc_index.find(query_doc, search_field='thumbnail__tensor', limit=3)
15
16 # find by the `video` tensor; neseted level
17 docs, scores = doc_index.find(query_doc, search_field='video__tensor', limit=3)
18
```

Nested DocList with subindex

Documents can be nested by containing a `DocList` of other documents, which is a slightly more complicated scenario than the previous one.

In this case, the nested `DocList` will be represented as a **new sub-index** (or table, collection, etc., depending on the database backend), that is linked with **the parent index** (table, collection, ...).

```

1 class ImageDoc(BaseDoc):
2     url: ImageUrl
3     tensor_image: AnyTensor = Field(space='cosine', dim=64)
4
5
6 class VideoDoc(BaseDoc):
7     url: VideoUrl
8     images: DocList[ImageDoc]
9     tensor_video: AnyTensor = Field(space='cosine', dim=128)
10
11
12 class MyDoc(BaseDoc):
13     docs: DocList[VideoDoc]
14     tensor: AnyTensor = Field(space='cosine', dim=256)
15
16
17 # create a Document Index
18 doc_index = HnswDocumentIndex[MyDoc](work_dir='/tmp3')
19
20 # create some data
21 index_docs = [
22     MyDoc(
23         docs=DocList[VideoDoc](
24             [
25                 VideoDoc(
26                     url=f'http://example.ai/videos/{i}-{j}',
27                     images=DocList[ImageDoc](
28                         [
29                             ImageDoc(
30                                 url=f'http://example.ai/images/{i}-{j}-{k}',
31                                 tensor_image=np.ones(64),
32                             )
33                             for k in range(10)
34                         ]
35                     ),
36                     tensor_video=np.ones(128),
37                 )
38                 for j in range(10)
39             ]
40         ),
41         tensor=np.ones(256),
42     )
43     for i in range(10)
44 ]
45
46 # index the Documents
47 doc_index.index(index_docs)
48

```

Search by subindex

```
1 # find by the `VideoDoc` tensor
2 root_docs, sub_docs, scores = doc_index.find_subindex(
3     np.ones(128), subindex='docs', search_field='tensor_video', limit=3
4 )
5
6 # find by the `ImageDoc` tensor
7 root_docs, sub_docs, scores = doc_index.find_subindex(
8     np.ones(64), subindex='docs__images', search_field='tensor_image', limit=3
9 )
10
```

Transiting data over network

Sending via REST API/JSON -> Backend: FastAPI

Sending via gRPC/ws -> Backend: Jina microservice

```
1 import numpy as np
2 from fastapi import FastAPI
3 from docarray.base_doc import DocArrayResponse
4 from docarray import BaseDoc
5 from docarray.documents import ImageDoc
6 from docarray.typing import NdArray
7
8 class InputDoc(BaseDoc):
9     img: ImageDoc
10     text: str
11
12
13 class OutputDoc(BaseDoc):
14     embedding_clip: NdArray
15     embedding_bert: NdArray
16
17
18 app = FastAPI()
19
20
21 @app.post("/embed/", response_model=OutputDoc, response_class=DocArrayResponse)
22 async def create_item(doc: InputDoc) -> OutputDoc:
23     ## call my fancy model to generate the embeddings
24     doc = OutputDoc(
25         embedding_clip=embed(doc.image), embedding_bert=embed(doc.text))
26     )
27     return doc
28
```

Transiting data over network

Sending via REST API/JSON -> Backend: FastAPI

Sending via gRPC/ws -> Backend: Jina microservice

```
1 import numpy as np
2 from fastapi import FastAPI
3 from docarray.base_doc import DocArrayResponse
4 from docarray import BaseDoc
5 from docarray.documents import ImageDoc
6 from docarray.typing import NdArray
7
8 class InputDoc(BaseDoc):
9     img: ImageDoc
10     text: str
11
12
13 class OutputDoc(BaseDoc):
14     embedding_clip: NdArray
15     embedding_bert: NdArray
16
17
18 app = FastAPI()
19
20
21 @app.post("/e
22 async def cre
23     ## call m
24     doc = Out
25     embed
26 )
27 return do
28
```

```
1 async with AsyncClient(app=app, base_url="http://test") as ac:
2     response = await ac.post("/doc/", data=docs.to_json()) # sending docs as json
3
4 assert response.status_code == 200
5 # You can read FastAPI's response in the following way
6 docs = DocList[TextDoc].from_json(response.content.decode())
```

Transiting data over network

Sending via gRPC/ws -> Backend: Jina microservice

```
1 class WhisperExecutor(Executor):
2     def __init__(self, device: str, *args, **kwargs):
3         super().__init__(*args, **kwargs)
4         self.model = whisper.load_model("medium.en", device=device)
5
6     @requests
7     def transcribe(self, docs: DocList[AudioURL], **kwargs) -> DocList[Response]:
8         response_docs = DocList[Response]()
9         for doc in docs:
10             transcribed_text = self.model.transcribe(str(doc.audio))['text']
11             response_docs.append(Response(text=transcribed_text))
12
13         return response_docs
14
```

Transiting data over network

Sending via gRPC/ws -> Backend: Jina microservice

```

1 class WhisperExecutor(Executor):
2     def __init__(self, device: str, *args, **kwargs):
3         super().__init__(*args, **kwargs)
4         self.model = whisper.load_model("medium.en", device=device)
5
6     @requests
7     def transcribe(self, docs: DocList[AudioURL], **kwargs) -> DocList[Response]:
8         response_docs = DocList[Response]()
9         for doc in docs:
10             transcribed_text = self.model.transcribe(doc.audio)
11             response_docs.append(Response(text=transcribed_text))
12
13     return response_docs
14

```

```

1 dep = Deployment(
2     uses=WhisperExecutor, uses_with={'device': "cpu"}, port=12349,
3     timeout_ready=-1
4 )
5 with dep:
6     docs = d.post(
7         on='/transcribe',
8         inputs=[AudioURL(audio='resources/audio.mp3')],
9         return_type=DocList[Response],
10     )
11
12 print(docs[0].text)
13

```

Agenda

- Preliminary: multimodal AI
- Opensource package: DocArray
 - Motivation
 - Representing data
 - Transiting data
 - Storing data
 - Retrieving data
- **Multimodal at scale in production**

This tutorial may require technical knowledge. Familiarity with Python 3.7+ concepts like data classes could be helpful.

An end to end example

https://docs.docarray.org/how_to/multimodal_training_and_serving/



Berlin · Beijing · Shenzhen

Thanks for your attention



jina.ai



@JinaAI_



han.xiao@jina.ai