

DATASTAX

DEVELOPERS

# AI as an API

Machine Learning with some REST

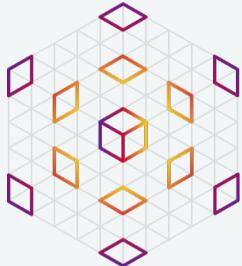


Based on original content by [CodingEnterpreneurs](#)





# › Agenda



---

**01****Housekeeping**

Live and Hands-on

---

**02****Our Goal**

Use case and Technologies

---

**03****Database Setup**

Data model &amp; Astra DB

---

**04****AI**

Train a text classifier

---

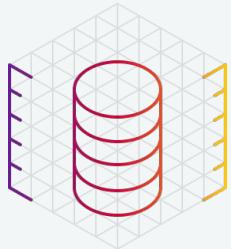
**05****API**

Bring to production

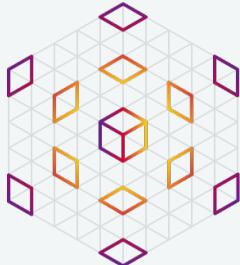
---

**06****What's next?**

Quiz, Homework, Agenda



# › Agenda



**01**

**Housekeeping**  
Live and Hands-on

**02**

**Our Goal**  
Use case and Technologies

**03**

**Database Setup**  
Data model & Astra DB

**04**

**AI**  
Train a text classifier

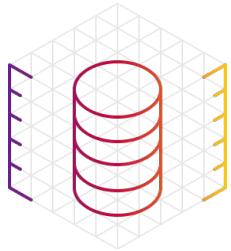
**05**

**API**  
Bring to production

**06**

**What's next?**  
Quiz, Homework, Agenda

# » Live & Interactive



[youtube.com/DataStaxDevs](https://youtube.com/DataStaxDevs)

Live!

Agenda

<b>01</b> PetClinic Architecture & Use Case	<b>02</b> DataStax Astra Cassandra Database The Art of Data Modeling	<b>03</b> Reactive Drivers Reactive vs Async
<b>04</b> Spring Reactive Boot and WebFlux	<b>05</b> User Interface Angular	<b>06</b> Game & Resources

DataStax Developers



YouTube

!menti

How much experience do you have with the **Spring Framework**?

Quiz!

Experience Level	Percentage
Never heard about it.	41%
I know the concepts.	24%
I have already used it.	10%
I use it regularly.	25%

DataStax Developers

!discord

Help!

Discord Developers

#workshop-chat

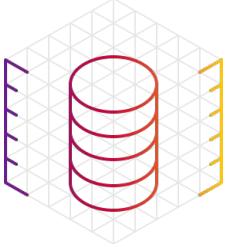


Mentimeter



Discord (#workshop-chat)

# › Hands-on housekeeping

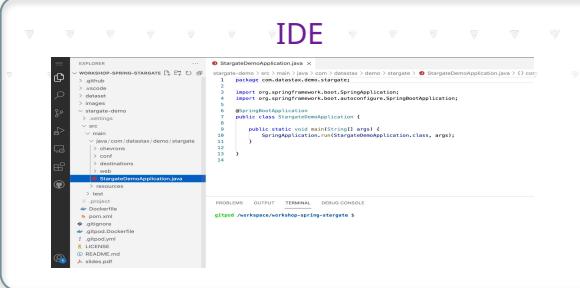


Mini-course + slides + hands-on

**DATASTAX**

ACADEMY

start from  GitHub  !github



IDE

```
gitpod.com/datas... /demo /stargate /StargateDemoApplication.java
```

```
1 import org.springframework.boot.SpringApplication;
2 import org.springframework.boot.autoconfigure.SpringBootApplication;
3
4 @SpringBootApplication
5 public class StargateDemoApplication {
6
7     public static void main(String[] args) {
8         SpringApplication.run(StargateDemoApplication.class, args);
9     }
10 }
```

PROBLEMS OUTPUT TERMINAL DEBUG CONSOLE  
gitpod /workspace/workshop-spring-stargate 3

 Gitpod  !gitpod

Database + Api + Streaming



**DATASTAX**

ASTRA DB

 !astra

CodingEnterpreneurs



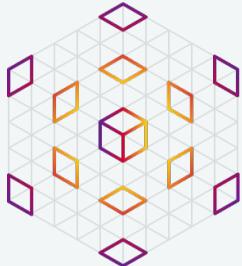
Youtube & blog for more



Nothing to install !



# › Agenda



**01**

## Housekeeping

Live and Hands-on

**02**

## Our Goal

Use case and Technologies

**03**

## Database Setup

Data model & Astra DB

**04**

## AI

Train a text classifier

**05**

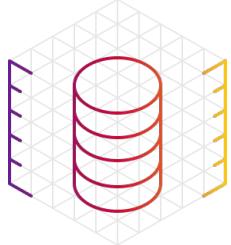
## API

Bring to production

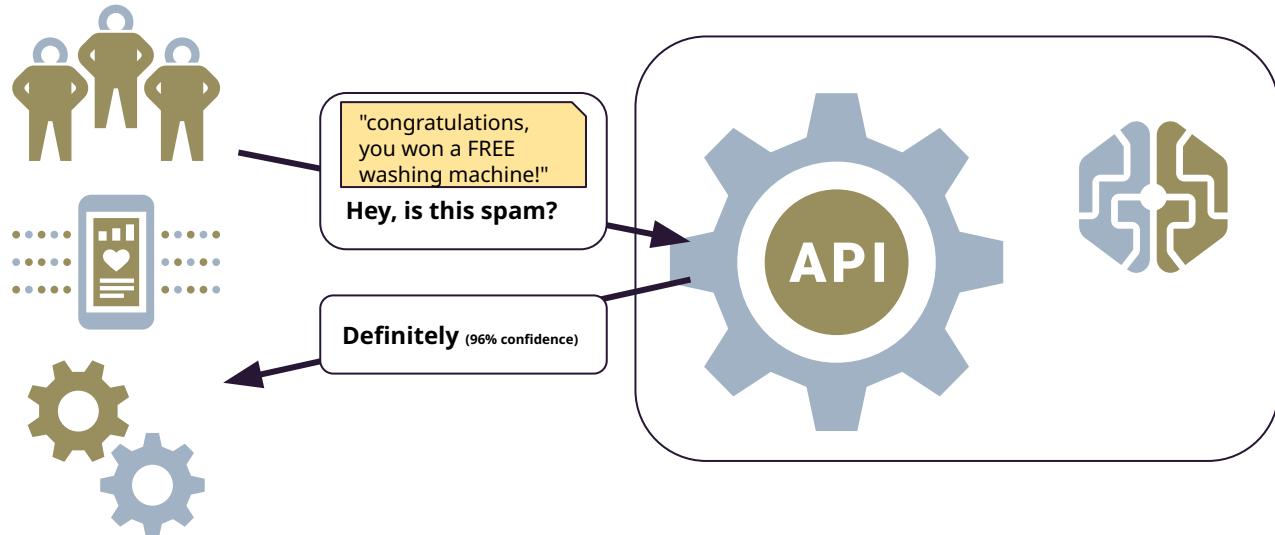
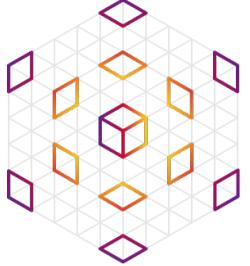
**06**

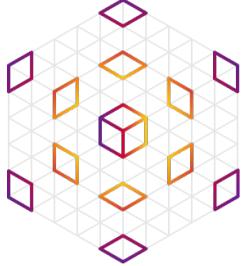
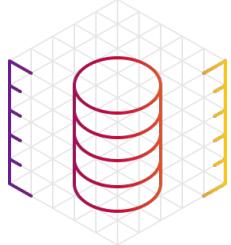
## What's next?

Quiz, Homework, Agenda



# » What we want



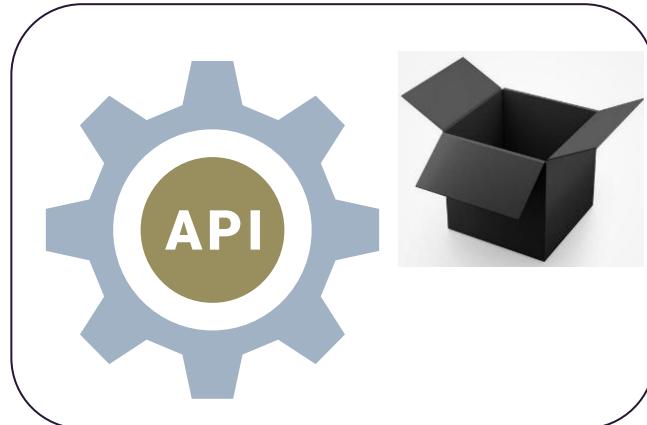


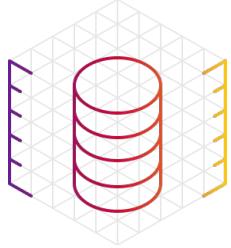
# ➤ Traditional viewpoint

The API sees the AI part as a black box

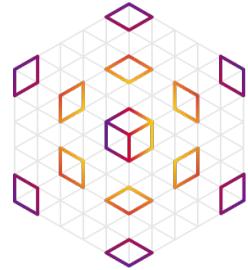
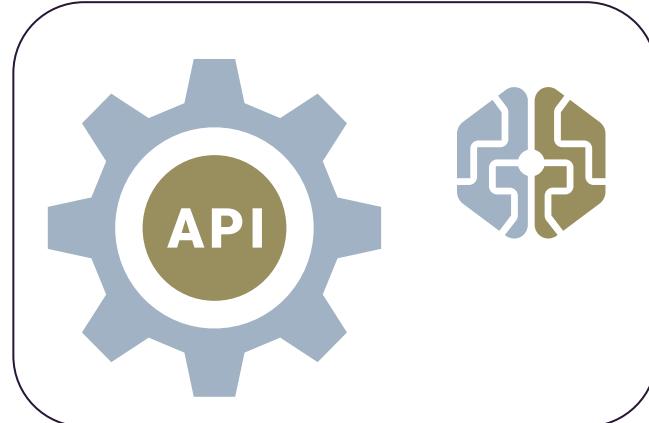
API = concerned with stability & performance

Engineer Stuff™





## ➤ Traditional viewpoint

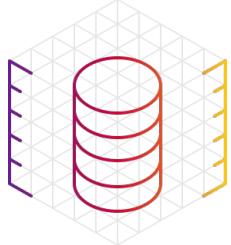


ML components, often messy  
(no offence meant)

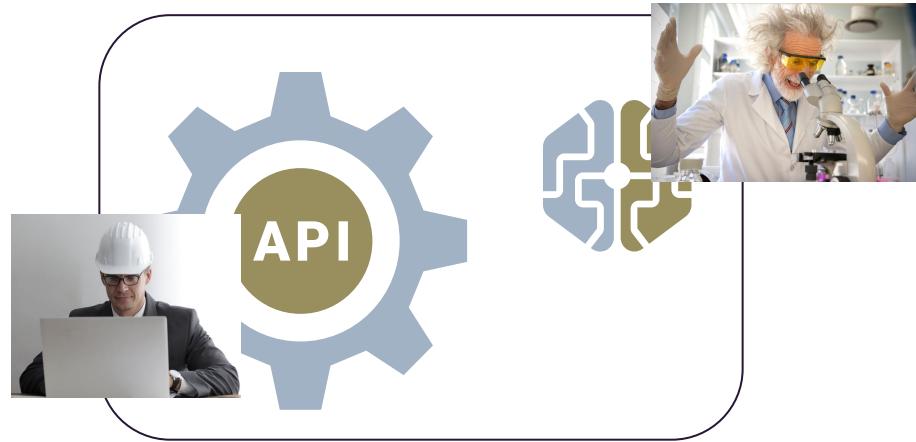
Treasures hidden away (hardly  
reproducible)

A "science" proper  
(experimentation, prototyping...)





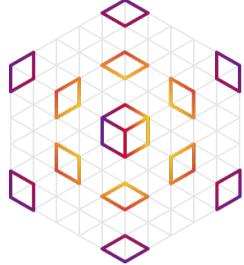
## ➤ Two in one



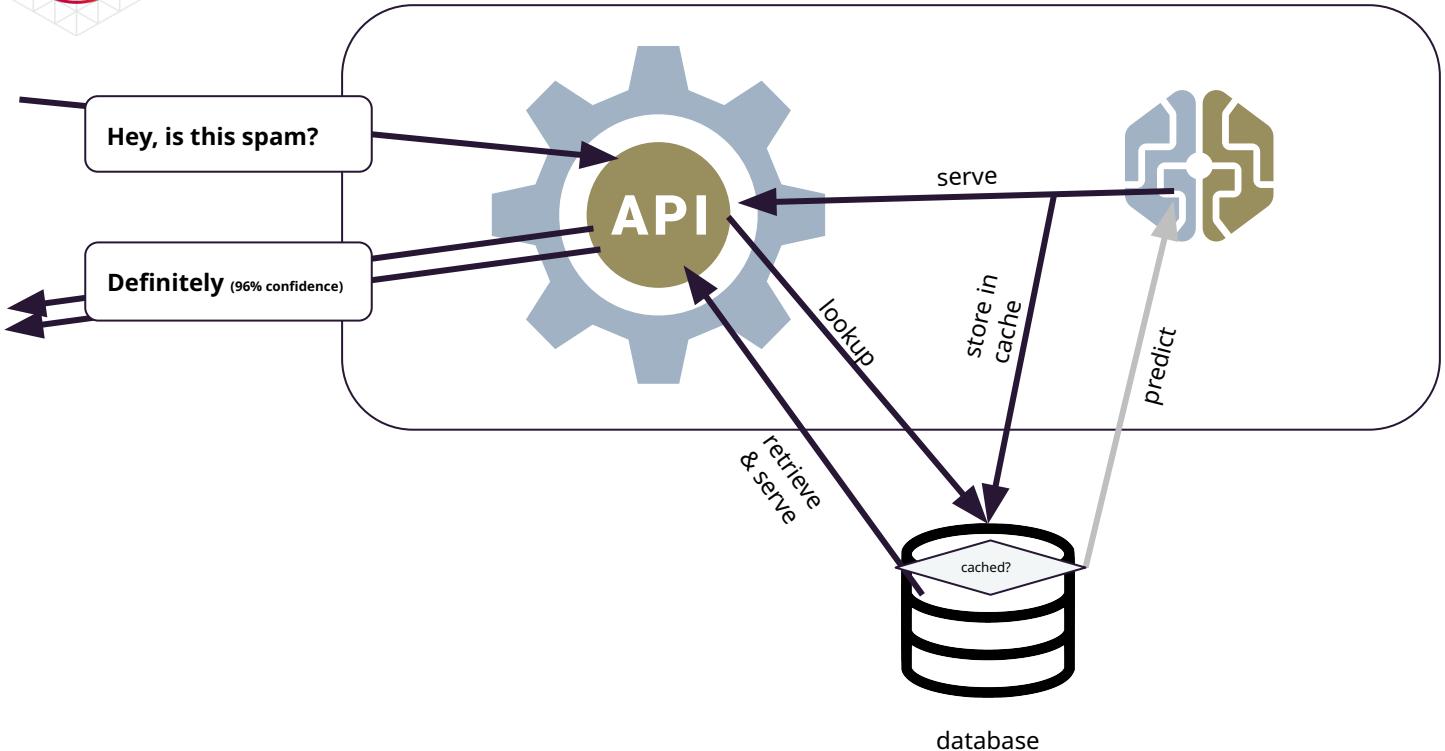
Properly exposing ML elements  
in production is a key skill

Still, two clearly separate phases:

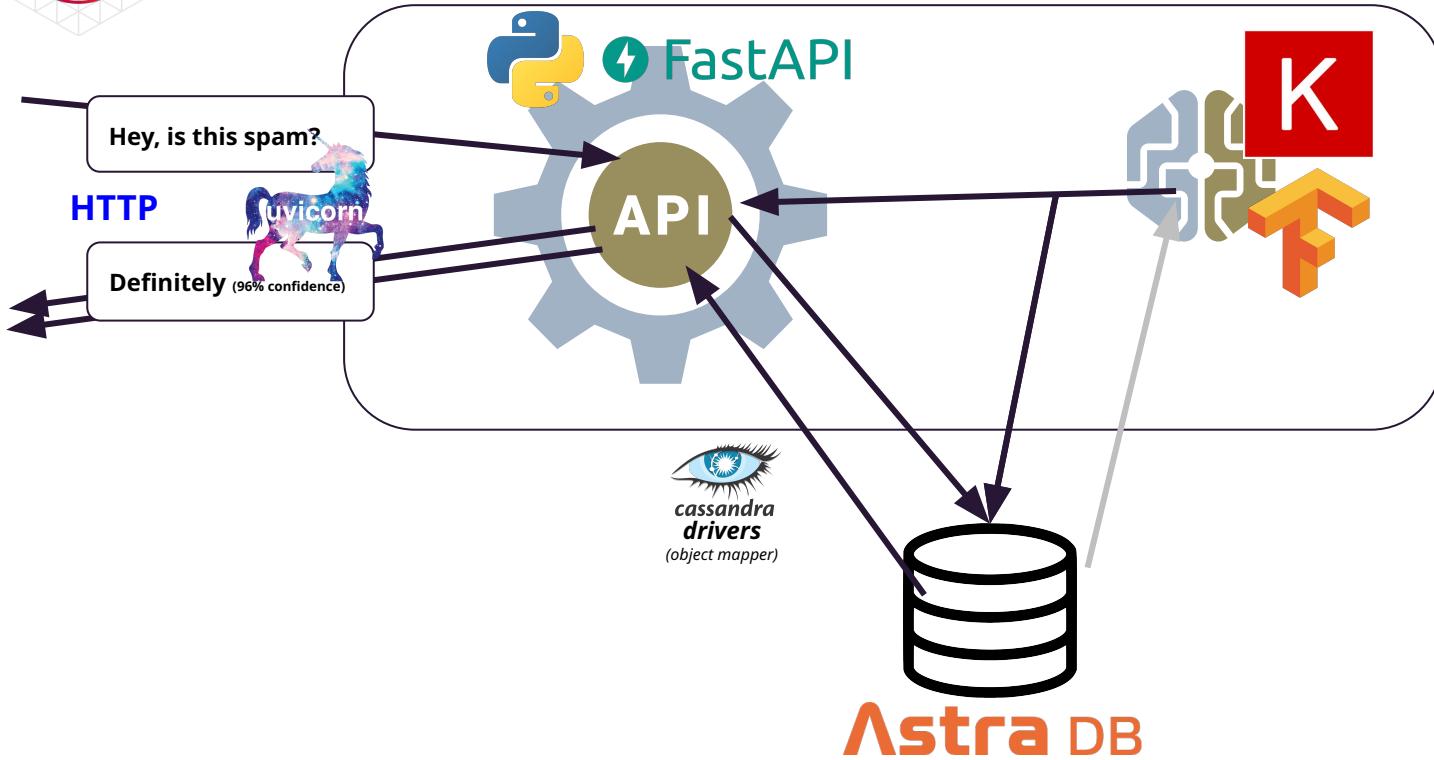
- create the classifier (design, train)
- bring it to production

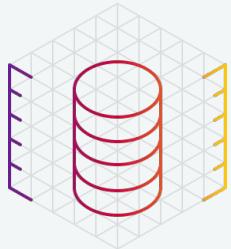


# ➤ Architecture sketch

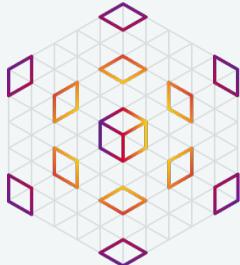


# › Tech stack





# › Agenda



**01**

## Housekeeping

Live and Hands-on

**02**

## Our Goal

Use case and Technologies

**03**

## Database Setup

Data model & Astra DB

**04**

## AI

Train a text classifier

**05**

## API

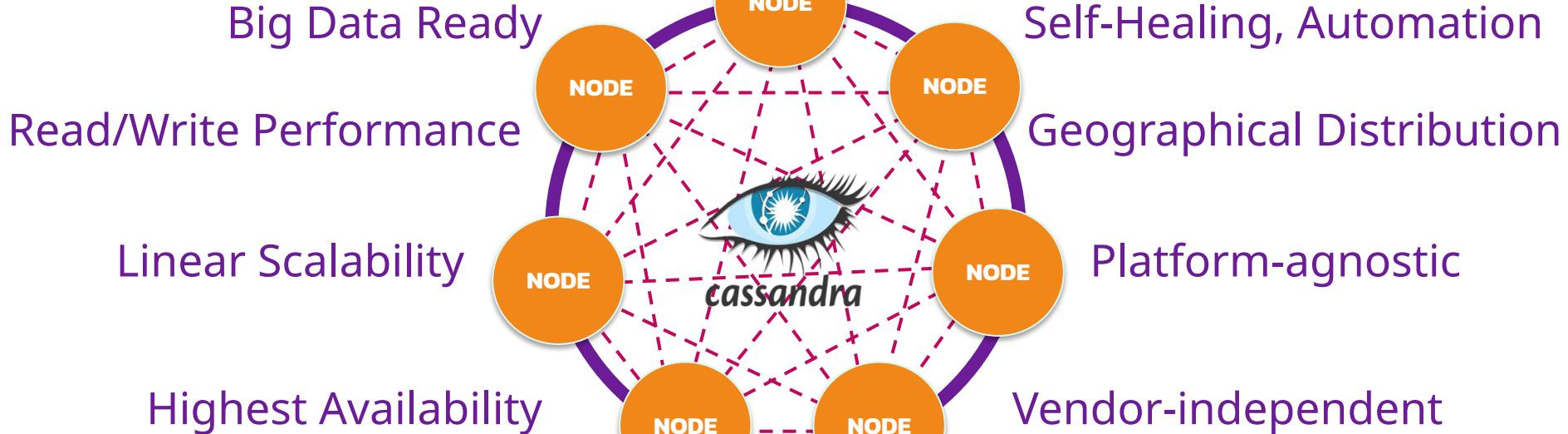
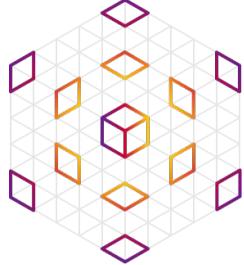
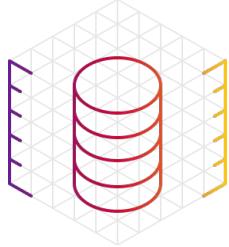
Bring to production

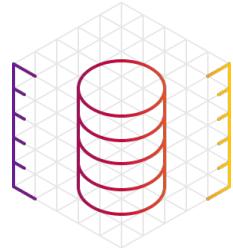
**06**

## What's next?

Quiz, Homework, Agenda

# › Intro to Apache Cassandra

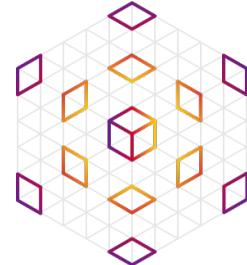




**GRAPHQL API**  
*Discoverability, Flexibility*  
Read only what you need

**CQL API**  
*Compatibility*  
It is Apache Cassandra !

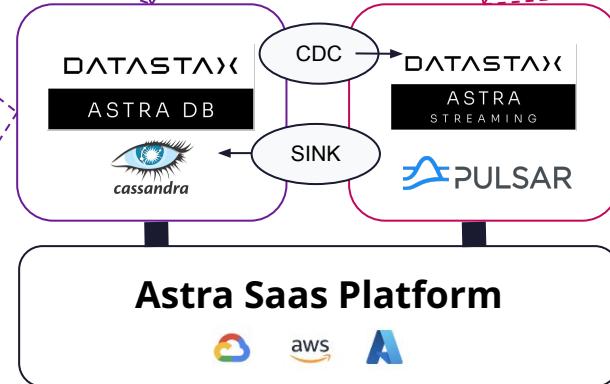
**PULSAR CLIENT**  
*Compatibility*  
It is Apache Pulsar !



**DOCUMENT API**  
*Flexible Schema*  
Document-Oriented usages

**REST API**  
*Stateless, Interoperability*  
CRUD for Cassandra

**GRPC API**  
*Performance, Interoperability*  
Execute cql in the fastest way



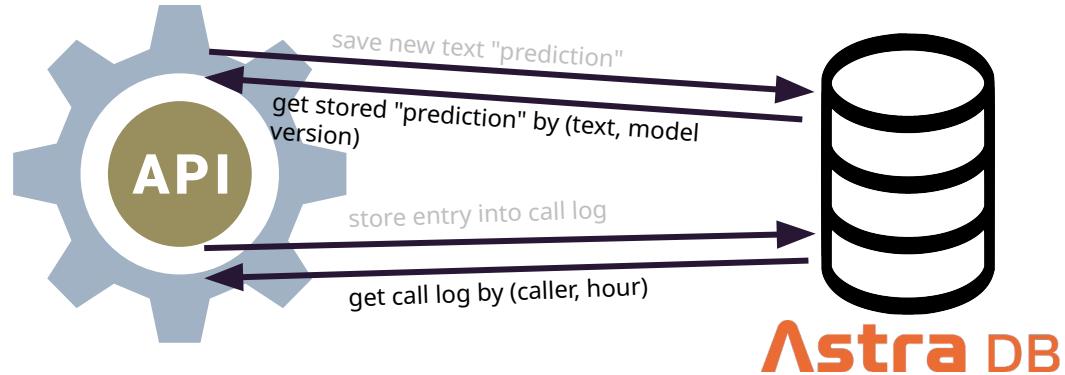
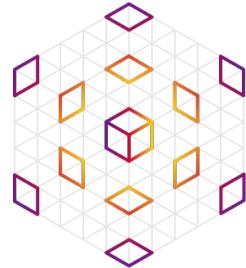
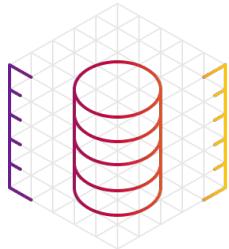
**JAVA MESSAGING SYSTEM**  
*Interoperability*  
Plug any existing Java Application

**AMQP**  
*Interoperability*  
Plug any existing Java Application

**KAFKA Interfaces**  
*Interoperability*  
Migrate Easily

API  
DEVOPS

# ➤ DB usage patterns

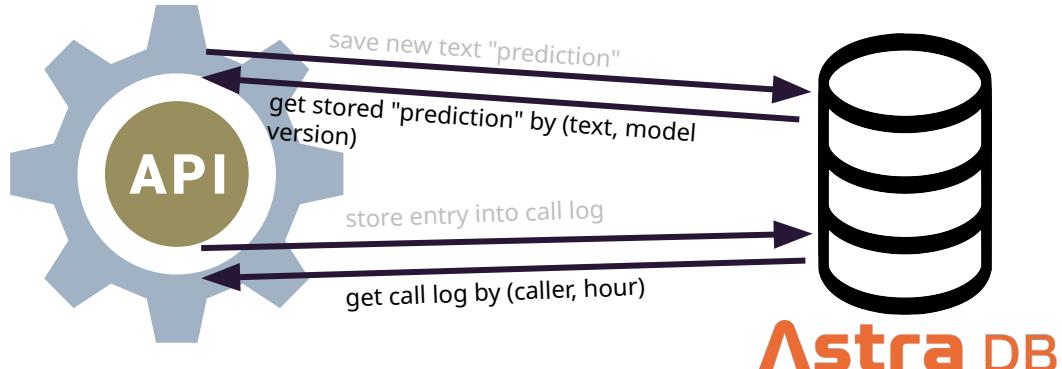
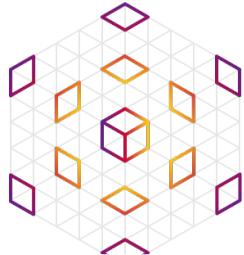
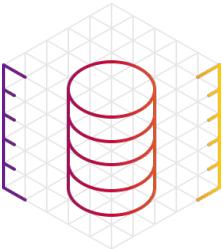


Data modeling, the Cassandra way:

***"design the table after the query"***

(also: tables are partitioned!)

# ➤ Data modeling



Astra DB

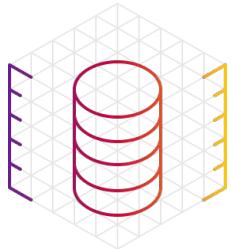
spam_cache_items		
model_version	text	K
input	text	K
confidence	float	
prediction_map	map<text,float>	
result	text	
stored_at	timeuuid	

spam_calls_per_caller		
caller_id	text	K
called_hour	timestamp	K
called_at	timeuuid	C
input_text	text	

"partition key" (rows stored together)

"primary key" (row uniqueness)

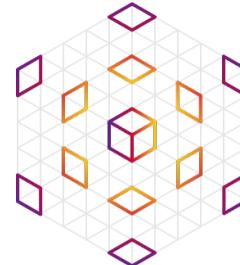
"Chebotko diagrams"



# ➤ DB I/O in practice



**cassandra  
drivers**  
*(object mapper)*



```
class SpamCallItem(Model):
    table_name__ = 'spam_calls_per_caller'
    __keyspace__ = ASTRA_DB_KEYSPACE
    __connection__ = 'my-astra-session'
    caller_id = columns.Text(primary_key=True, par-
    called_hour = columns.DateTime(primary_key=True,
    called_at = columns.TimeUUID(primary_key=True,
    input = columns.Text()
```

**Tables are created  
automatically...**

```
# Database
logging.info('      DB initialization')
DBSession = initSession()
sync_table(SpamCacheItem)
sync_table(SpamCallItem)
logging.info('      API Startup completed.')
```

```
query = SpamCallItem.objects().filter(
    caller_id=caller_id,
    called_hour=called_hour,
)
for item in query:
    yield item
```

```
called_hour = getThisHour()
for input in inputs:
    SpamCallItem.create(
        caller_id=caller_id,
        called_hour=called_hour,
        input=input,
    )
```

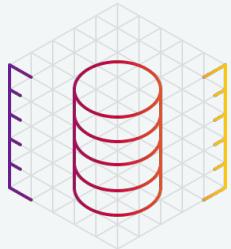


# Lab 1

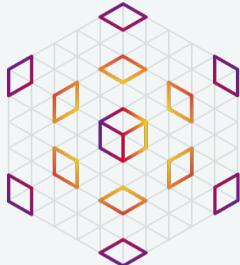
## DB setup

**Create your Astra DB instance**

**Create a "DB Admin" token**



# › Agenda



---

**01****Housekeeping**

Live and Hands-on

**02****Our Goal**

Use case and Technologies

**03****Database Setup**

Data model &amp; Astra DB

---

**04****AI**

Train a text classifier

**05****API**

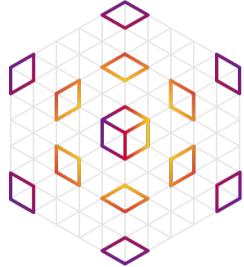
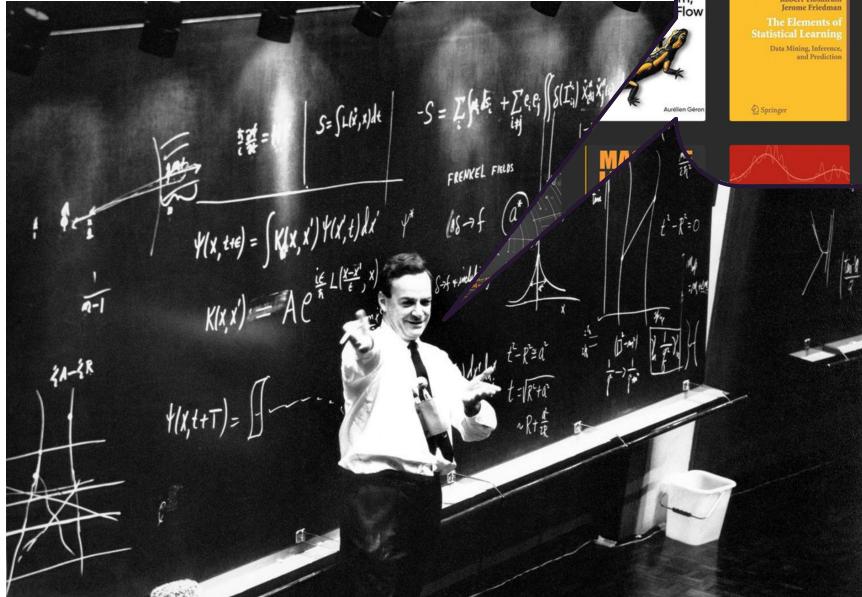
Bring to production

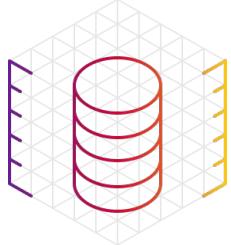
**06****What's next?**

Quiz, Homework, Agenda

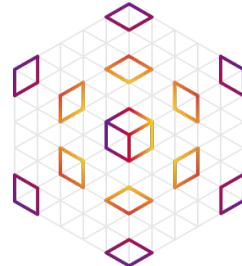
# » Disclaimer

What this section is **not**:





# › AI and ML



## "ML: LSTM RNN for NLP"

This is what we want:

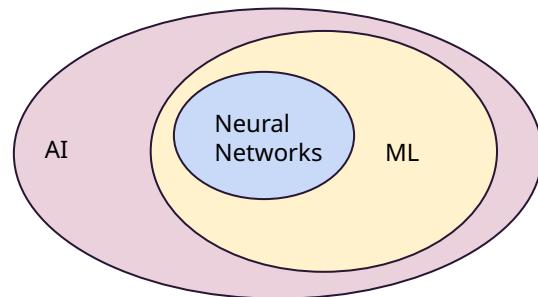


**Machine Learning** = *algorithms that improve by being fed data, without being explicitly instructed what to do.*

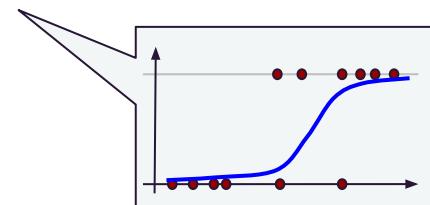
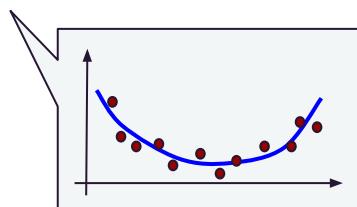
It's essentially statistical inference (with superpowers).

Lots of **math** involved (linear algebra, calculus, prob./statistics).

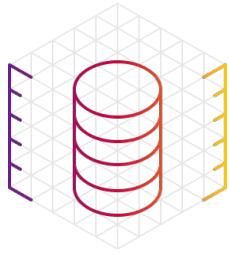
Nowadays accessible as neatly-packaged tools (good for us!)



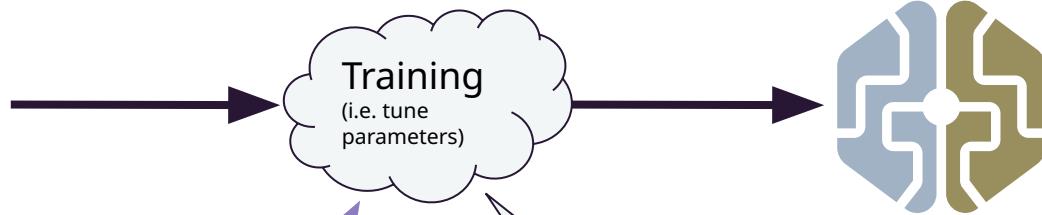
*Simple* examples of ML: least-squares fits, logistic regressions.



# › Supervised learning



"Blank slate" model



**Choice of the model**  
(also "hyperparameters")

**Labeled training**

Go until jurong point Ok lar... Joking wif Free entry in 2 a wk U dun say so early ha Nah I don't think he FreeMsg Hey there da Even my brother is n As per your request WINNER!! As a valued Had your mobile 11 m	X	ham ham spam ham ham spam ham ham spam spam
	Y	

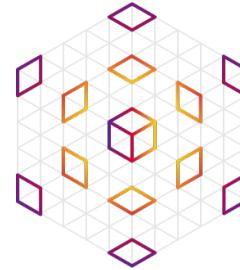
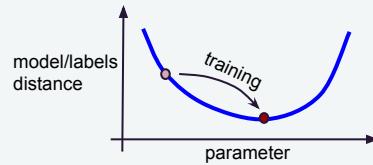


Ready to predict !

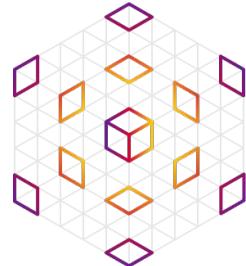
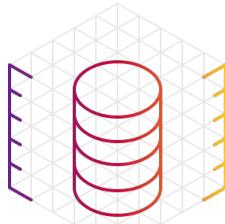
(classify, recognize objects  
in images, do  
speech-to-text, ...)

**"Training":**

tuning model parameters to make  
predictions as close to labeled  
data as possible *(more later)*



# » Neural Networks



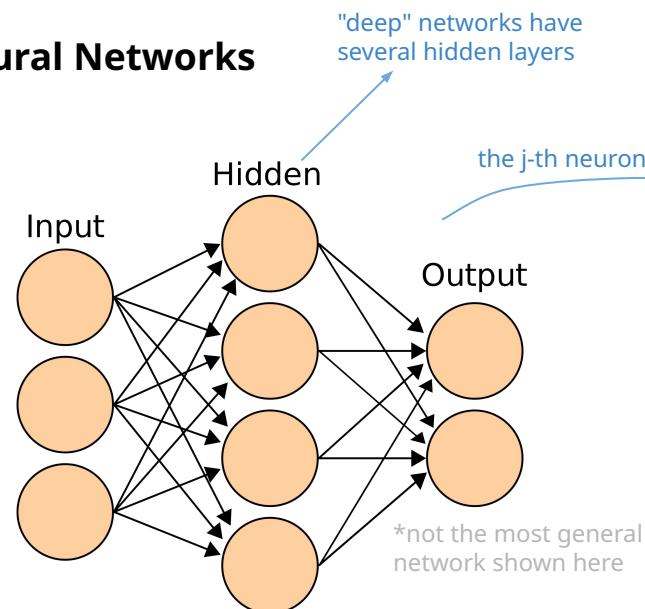
Many algorithms and structures for ML:

linear/logistic regression, support-vector-machines, decision trees/random forests, gradient boosting, ...

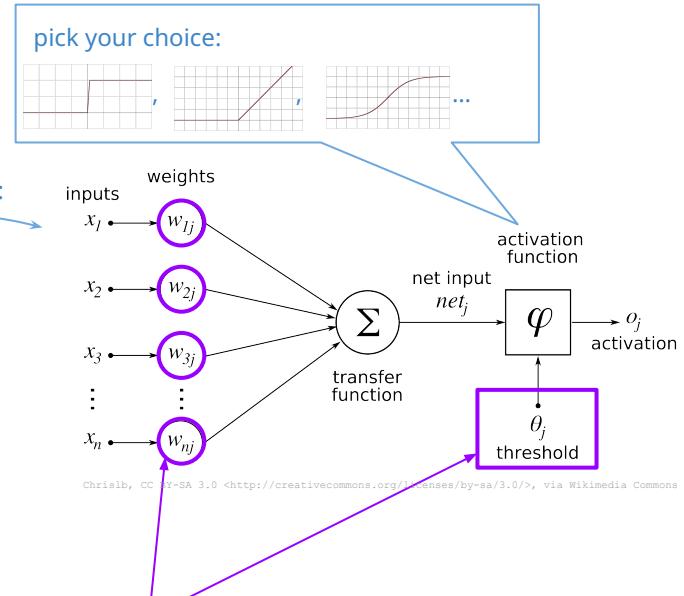


## Focus on Neural Networks

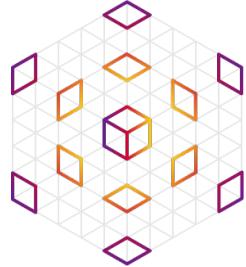
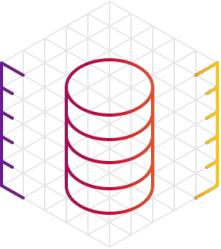
Input (text, image, ...) is encoded and fed here



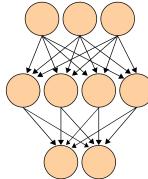
many parameters per neuron:  
training gets CPU- and time-intensive



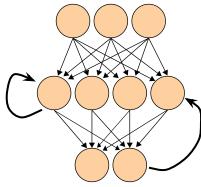
# ➤ RNN and LSTM



Architectures



"Feed-forward"



"Recurrent neural network"

has *loops*: ~ "memory"

better with sequences (time series, speech, **texts**, ...)

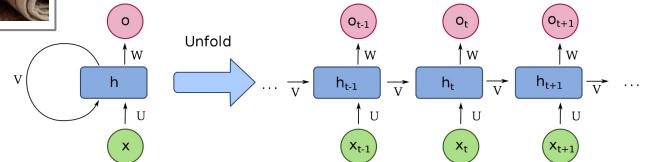
loops "unfold" as time-steps: training gets harder

➤ "Long-short-term-memory"

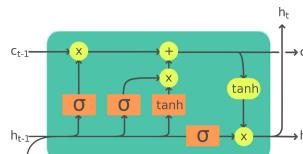
introduces explicit memory elements

overcomes the "vanishing gradient" training problem

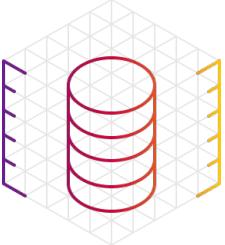
a standard RNN to process "**input sequences**"



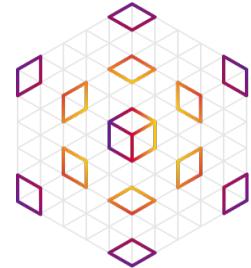
By fdeleo - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=60109157>



By Guillaume Chevalier - File:The\_LSTM\_Cell.svg, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=109362147>



# NLP



## Natural language processing

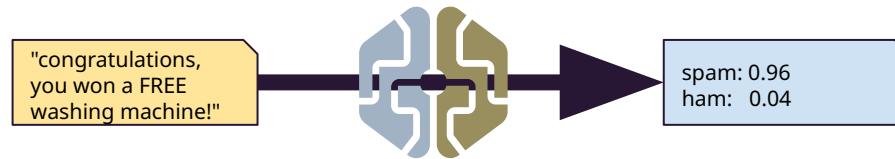
A very wide field of application for ML (since ... forever)

**translation, summarization, classification, clustering, autocomplete, generation, chatbots, sentiment analysis, ...**

Most of the data "out there" is unstructured (e.g. **text**)

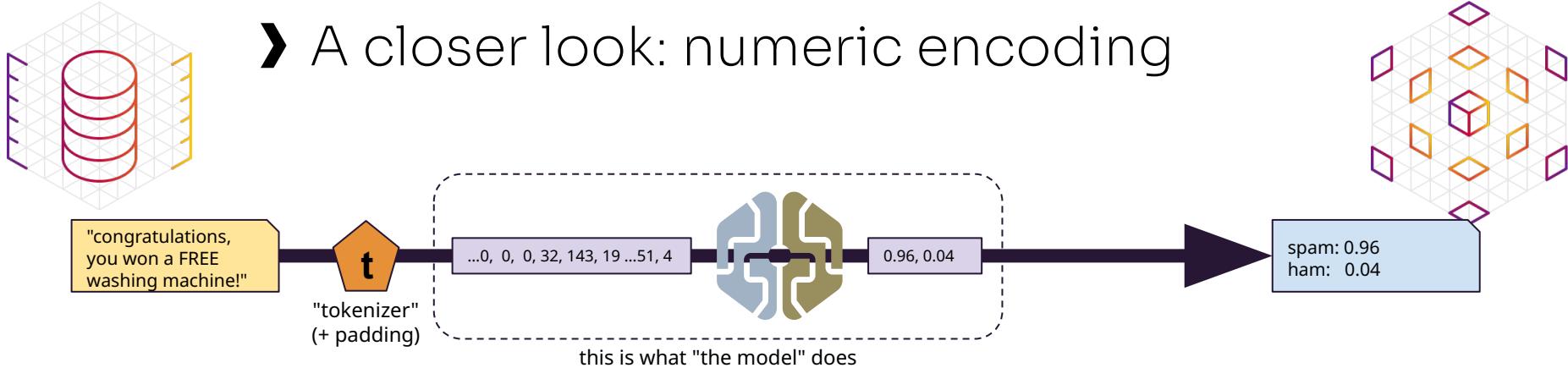


**LSTM very apt at many of these tasks (memory along sequences)**

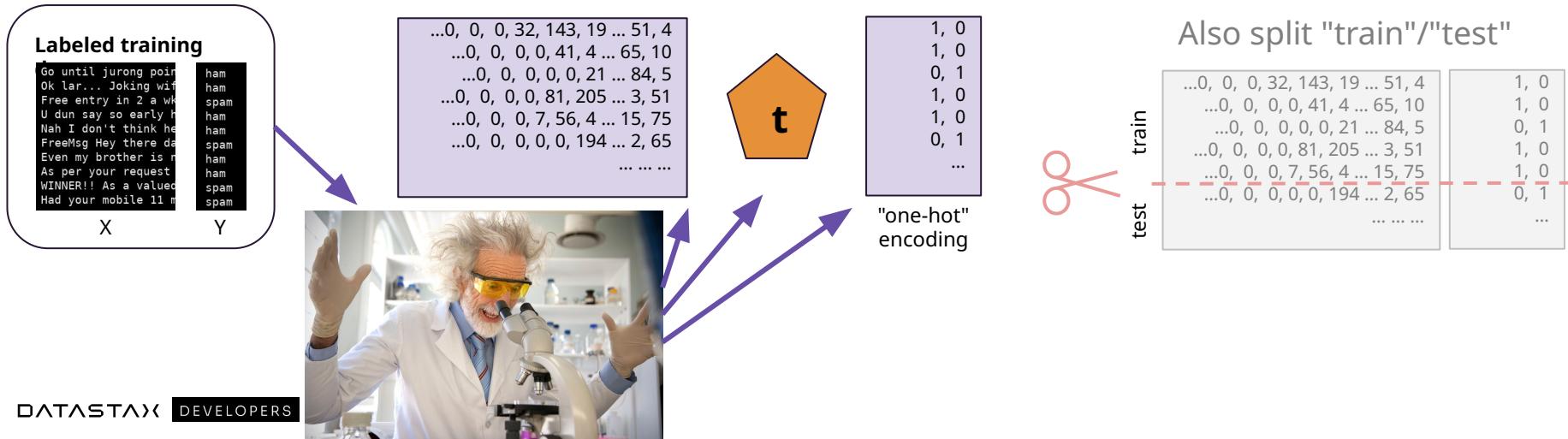


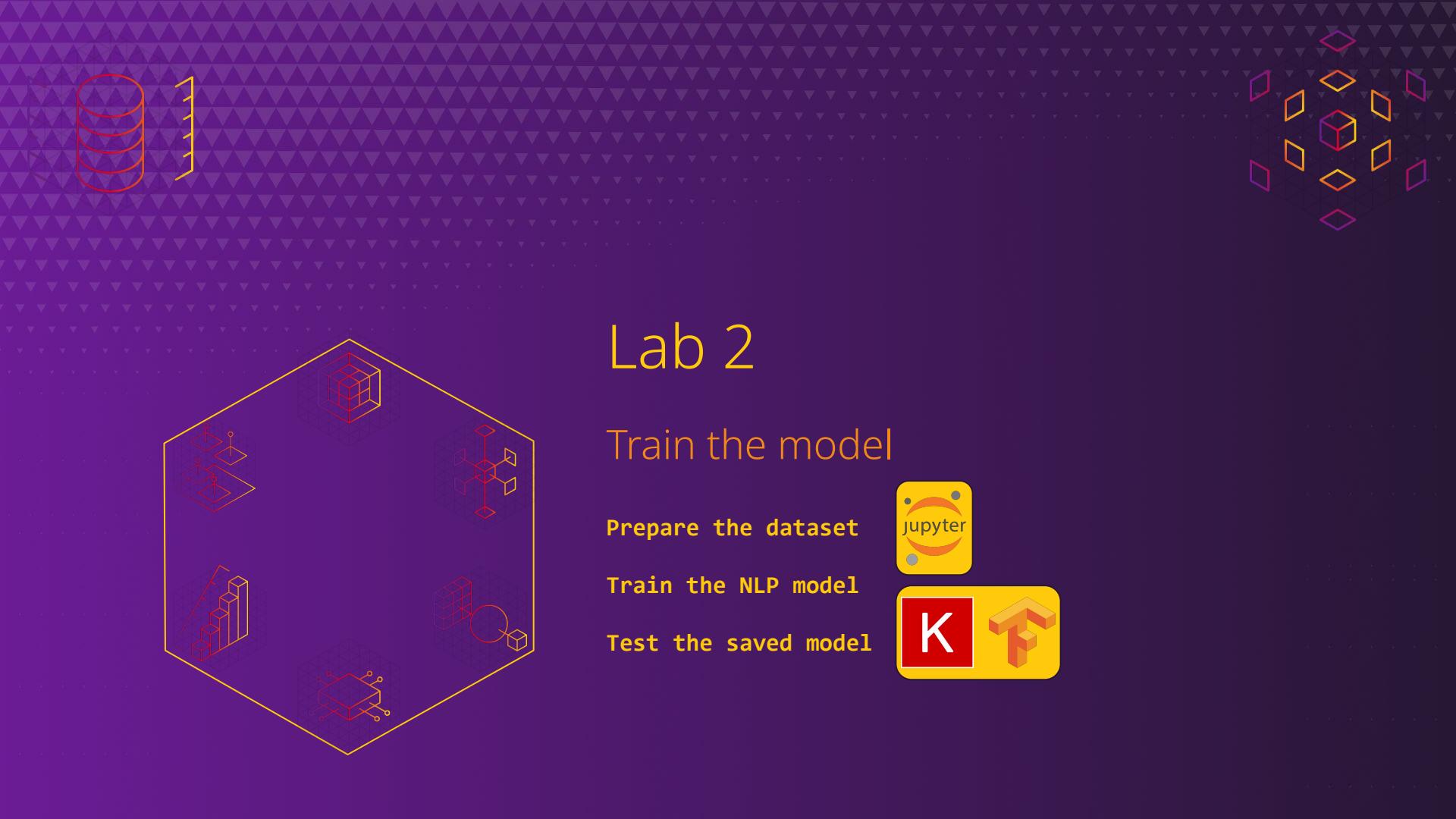
*An evolving field (Word2Vec ~ 2013; transformers e.g. BERT ~ 2018, "foundation models" ~ 2020, and we all know where we're now)*

# ► A closer look: numeric encoding



We must then *prepare the dataset* before training (including building the tokenizer):





# Lab 2

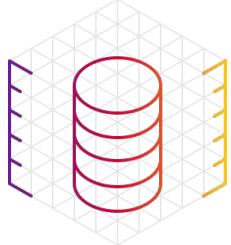
Train the model

Prepare the dataset

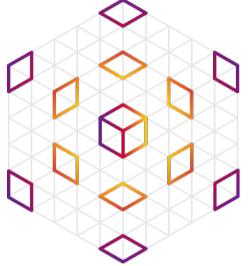
Train the NLP model

Test the saved model





# › The antispam model architecture



```
model = Sequential()  
model.add(Embedding(maxNumWords, embedDim, 1))  
model.add(SpatialDropout1D(0.4))  
model.add(LSTM(LstmOut, dropout=0.3, recurrent_dropout=0.3))  
model.add(Dense(2, activation='softmax'))  
model.compile(loss='categorical_crossentropy', optimizer='adam')
```

***A stack of layers:***

one translates an input number into a vector

one randomly disables pieces of network to enhance training ("dropout")

then the LSTM (recurrent *within the layer*)

and a final ordinary layer reducing it all to a 2-component output (spam/ham)

***"finalize" the model: ready to train!***

inputs

# » During training ...

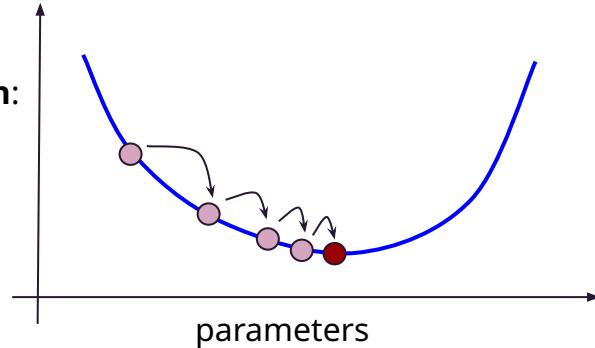
## "Training":

(progressively) tuning model parameters to make predictions as close to labeled data as possible

### loss function:

distance between provided labels and model output for training data

e.g. "categorical cross-entropy"  
(from information theory)



## The reality:

lots of parameters, it's a mess!

```
=====
Total params: 291,034
Trainable params: 291,034
Non-trainable params: 0
```

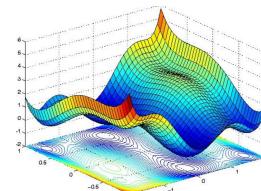
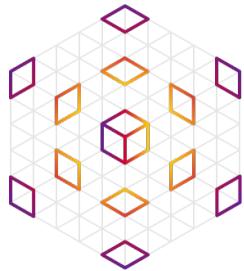
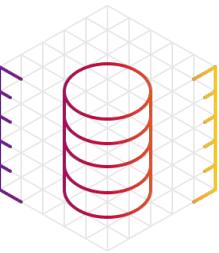
"**Stochastic gradient descent**" (SGD)



**Backpropagation** (glorified chain rule)



**Drop-out**



# Approaching the trained model

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=42)
```

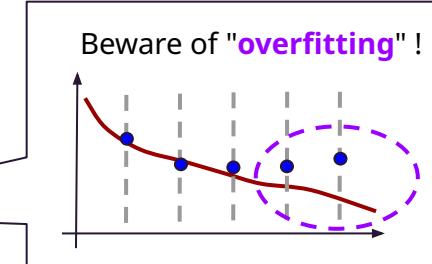
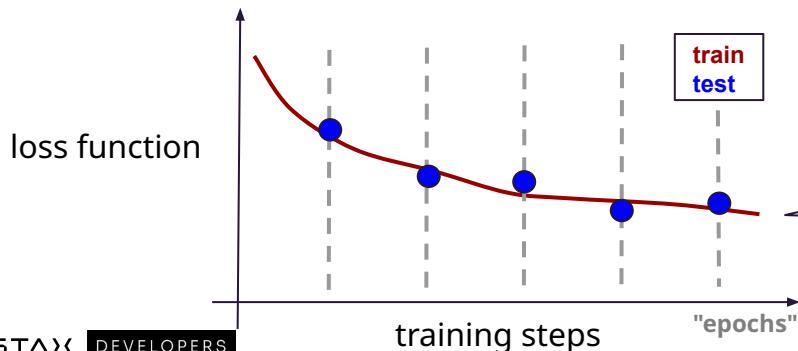
We split "train"/"test"

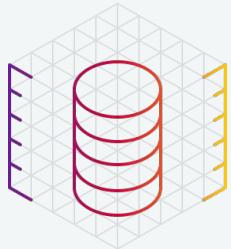


	train	test
...0,	0, 0, 0, 32, 143, 19 ... 51, 4	1, 0
...0,	0, 0, 0, 0, 41, 4 ... 65, 10	1, 0
...0,	0, 0, 0, 0, 21 ... 84, 5	0, 1
...0,	0, 0, 0, 0, 81, 205 ... 3, 51	1, 0
...0,	0, 0, 0, 7, 56, 4 ... 15, 75	1, 0
...0,	0, 0, 0, 0, 194 ... 2, 65	0, 1
	.... ...	...

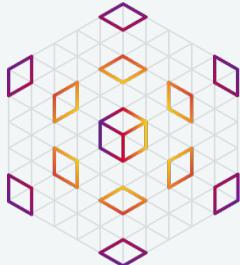
**Train on part of input dataset,  
periodically validate on another portion of it**

*(independently verify we're actually learning the right problem)*





# » Agenda



**01**

**Housekeeping**  
Live and Hands-on

**02**

**Our Goal**  
Use case and Technologies

**03**

**Database Setup**  
Data model & Astra DB

**04**

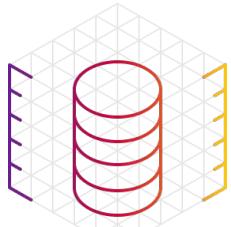
**AI**  
Train a text classifier

**05**

**API**  
Bring to production

**06**

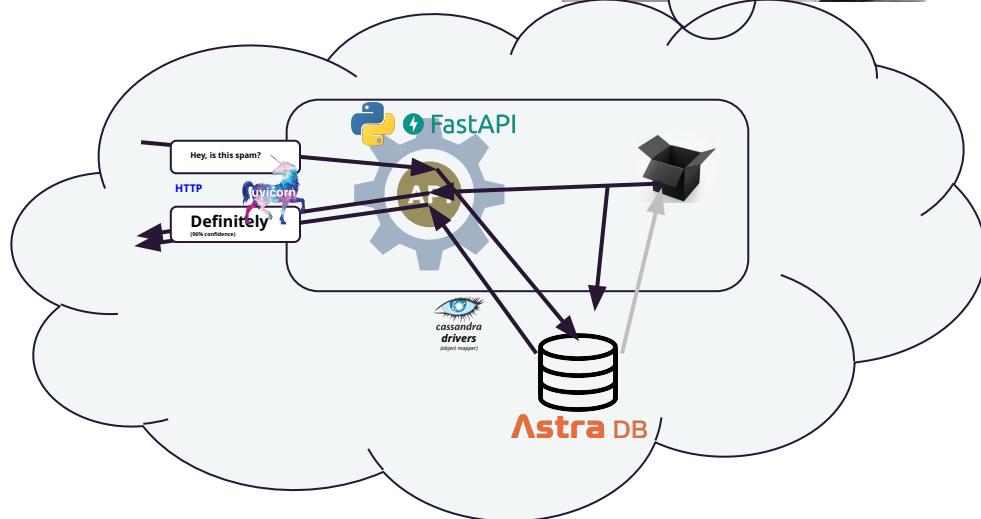
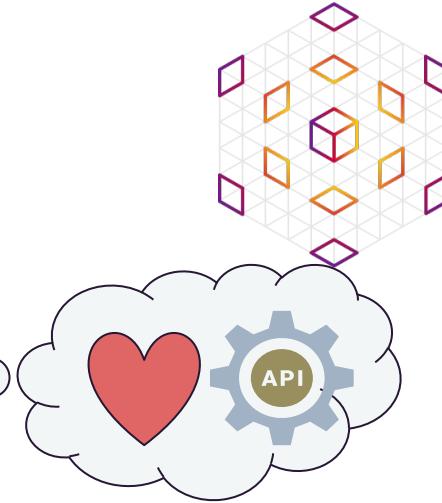
**What's next?**  
Quiz, Homework, Agenda



# ➤ What now?



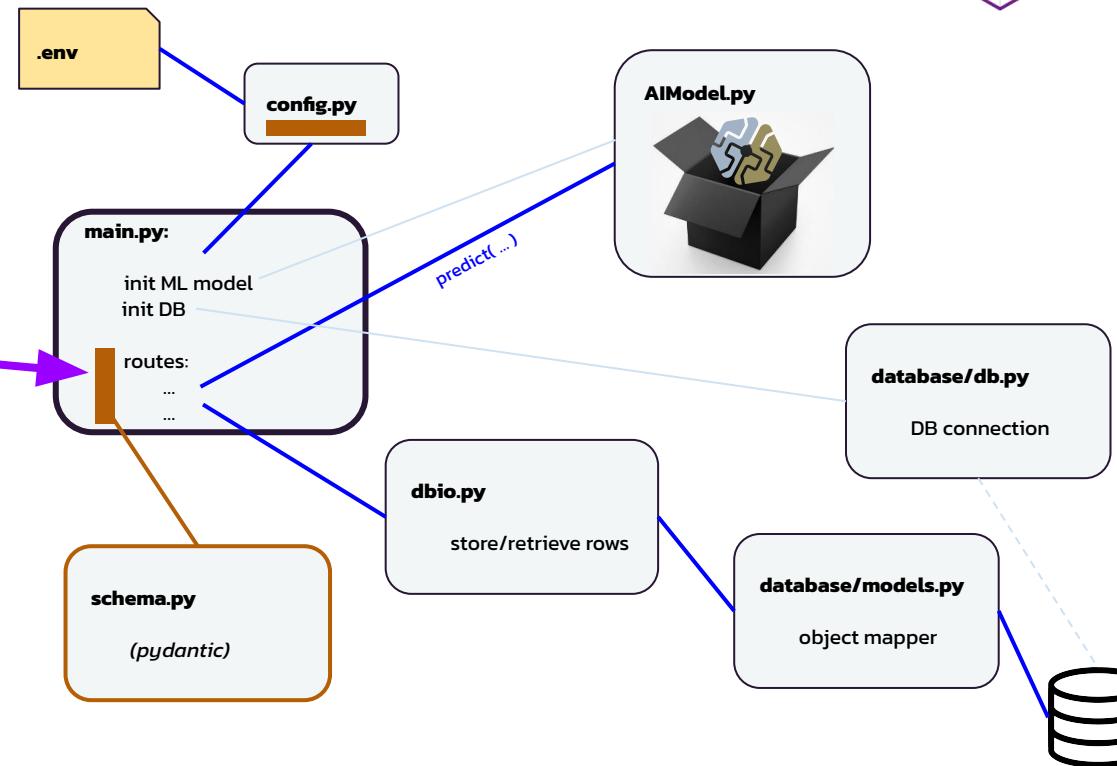
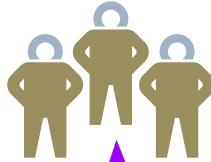
Here, take this



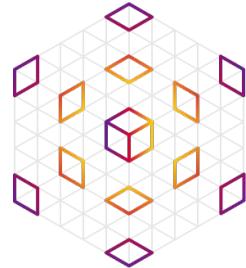
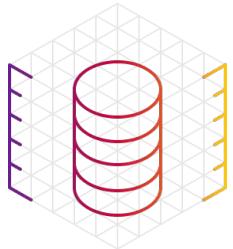
# ➤ API structure

api

```
└── main.py  
└── AIModel.py  
└── schema.py  
└── dbio.py  
└── database  
    └── db.py  
    └── models.py  
└── config.py  
└── minimain.py  
└── MockSpamAIModel.py  
└── tests (...)
```



# ➤ Features & highlights



## ⚡ FastAPI fundamentals: routes, body validation

caching on DB to avoid re-processing the same input

Spam Classifier API 0.1.0 OAS 3.0.0

Swagger API

A sample API exposing a Keras text classifier model.

classification Requests for text classifications.

GET /prediction Single Text Prediction

POST /prediction Single Text Prediction

API docs, Swagger(Open API) UI

call logs & **StreamingResponse**

**GET** and query string

```
@app.get('/prediction', response_model=PredictionResult, tags=['classification'])
def single_text_prediction_get(request: Request, query: SingleTextQuery = Depends()):
```

```
settings = getSettings()
cached = None if query.skip_cache else readCachedPrediction(settings.model_version,
storeCallsToLog([query.text], request.client[0])
if not cached:
    result = spamClassifier.predict([query.text], echoInput=query.echo_input)[0]
    cachePrediction(query.text, result)
    result['from_cache'] = False
    return PredictionResult(**result)
else:
    cached['from_cache'] = True
    return PredictionResult(**cached)
```

```
@app.get('/recent_log', response_model=List[CallerLogEntry], tags=['info'])
def get_recent_calls_log(request: Request):
    caller_id = request.client[0]
    called_hour = getThisHour()
    #
    return StreamingResponse(formatCallerLogJSON(caller_id, called_hour))
```



# Lab 3

API

Serve with FastAPI

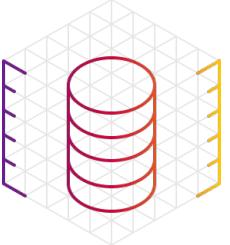


Explore Swagger (Open API) UI

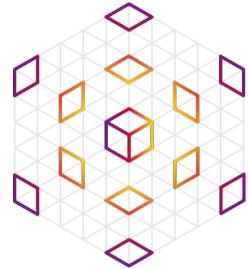


Look at the DB





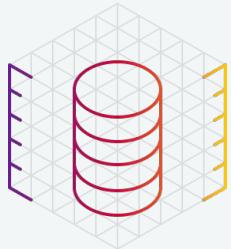
## » More topics



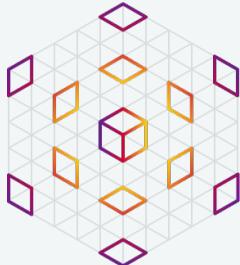
Want to know more on ...

- feature stores
- architectural challenges of ML-based applications
- model versioning with  **FastAPI**
- other MLOps-related issues

Check out **[github.com/datastaxdevs/workshop-ai-as-api](https://github.com/datastaxdevs/workshop-ai-as-api)** (Appendix I, Appendix II)



# › Agenda



**01**

## **Housekeeping**

Live and Hands-on

**02**

## **Our Goal**

Use case and Technologies

**03**

## **Database Setup**

Data model & Astra DB

**04**

## **AI**

Train a text classifier

**05**

## **API**

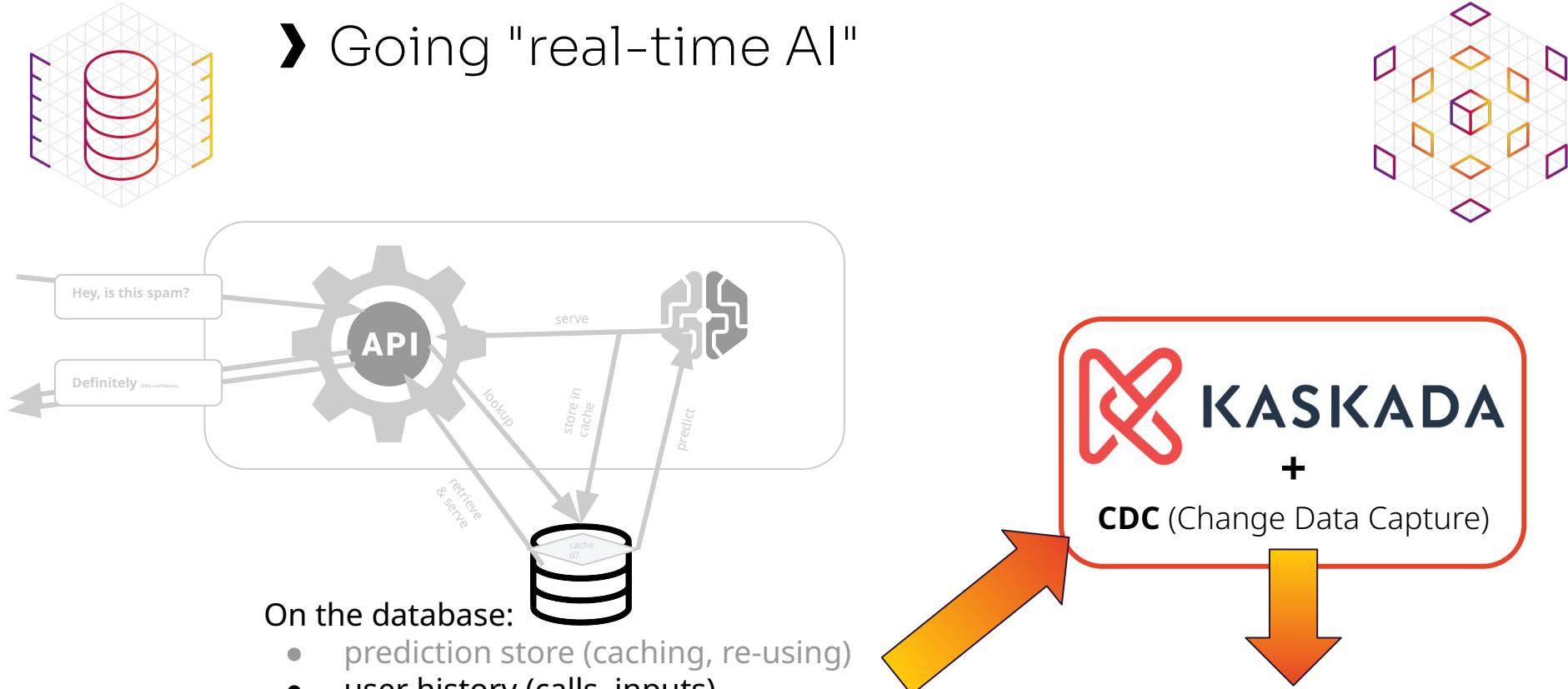
Bring to production

**06**

## **What's next?**

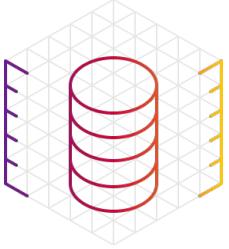
Quiz, Homework, Agenda

# ➤ Going "real-time AI"

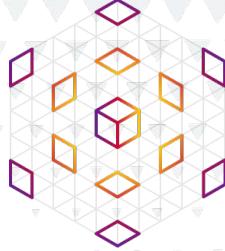


On the database:

- prediction store (caching, re-using)
- user history (calls, inputs)
- per-sender time-series (spam/ham events)



# Planet Cassandra Discord server



!discord

The screenshot shows the Discord interface for the Planet Cassandra server. On the left is a sidebar with a tree view of channels:

- # Browse Channels
- # moderator-only
- WELCOME
  - start-here
  - code-of-conduct
  - FAQ
- # discord-help
- @ PLANET CASSANDRA
  - # cassandra-chat
  - announcements
- CONTRIBUTORS
  - contributor-guide
  - submit-a-link
  - content-submission
  - content-review
  - website-help
- EVENTS
  - # workshop-chat
  - ama-stage
- > STAFF
- > SERVER STATS
- > ARCHIVED

The main window shows a channel named "# workshop-chat". A message from user @FJ asks about a badge for a workshop. User Stefano Lottini responds with instructions for completing an assignment and filling out a Google form. There is a back-and-forth between users @FJ, Stefano Lottini, and cgellersm01 regarding access to the Google form. User @FJ submits a Netflix Clone homework, and Stefano Lottini encourages him to get it reviewed.

```
@FJ Hello, just checking if anyone or the mods can tell me how the badge for yesterday's workshop Netflix clone works?
Stefano Lottini 03/24/2023 1:12 AM
Hello! To receive a badge you are required to complete a small assignment and then fill a google form with some "theory questions" + uploading a screenshot for proof of your assignment. You can find all instructions here: https://github.com/datastaxdevs/workshop-graphql-netflix#homework
A few hours ago I graded all pending submissions and released the badges (you get notified by email when that happens). We usually do that quickly after new submissions come in (1-2 days generally)

@FJ Hello! To receive a badge you are required to complete a small assignment and then fill a google form with some "theory quest"
FJ 03/24/2023 1:13 AM
thank you...
cgellersm01 03/24/2023 2:58 AM
Hey! Sorry but I tried to access the google form both yesterday and today, and I don't seem to have permissions to do it, I'm logged in with the same account I used to inscribe to the workshop. Any ideas on why I can't access the form? Thx in advance
cgellersm01 03/24/2023 3:34 AM
Hey! Sorry but I tried to access the google form both yesterday and today, and I don't seem to have permissions to do it, I'm logged in with the same account I used to inscribe to the workshop. Any ideas on why I can't access the form? Thx in advance
Stefano Lottini 03/24/2023 3:34 AM
Hmm, I think you need a google account. @Aleks could probably check the form settings (yesterday we got several submissions so I think the form is all right)
FJ 03/24/2023 4:27 AM
Just submitted my Netflix Clone - Homework, to get the badge, hope to get it reviewed soon and get the badge 😊 ^ it was a very exciting workshop!!
@FJ Just submitted my Netflix Clone - Homework, to get the badge, hope to get it reviewed soon and get the badge 😊 ^ it was a very exciting workshop!
Stefano Lottini 03/24/2023 4:33 AM
Glad to hear that! I can see your submission - and will get to it shortly 😊
cgellersm01 03/24/2023 10:39 AM
Yup! That was it, I was using my education email with a weird @ so it wasn't supported, thanks a lot!
```

dtsx.io/discord



## DataStax Developers

@DataStaxDevs  
29.9K subscribers

[HOME](#)[VIDEOS](#)[SHORTS](#)[LIVE](#)[PLAYLISTS](#)[COMMUNITY](#)[CHANNELS](#)[ABOUT](#)[Recently uploaded](#)[Popular](#)

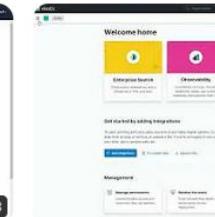
The Legend of DataStax - Cassandra Summit Training Day!



Cassandra Summit Training Day - Sunday March 12th

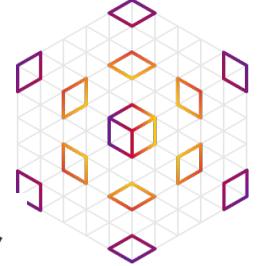
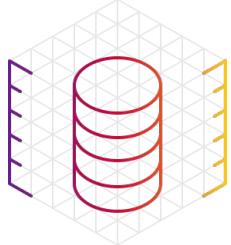


CDC for Astra DB Demo: Sink to Astra DB  
92 views • 2 weeks ago



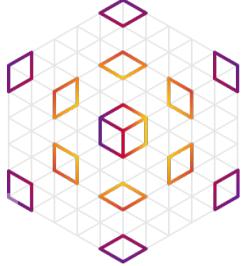
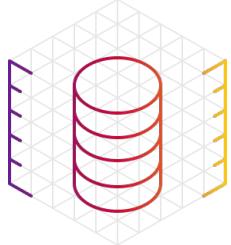
CDC for Astra DB Demo: ElasticSearch





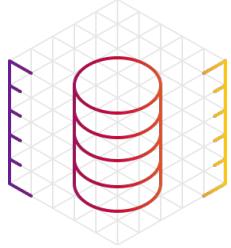
# ➤ Badges



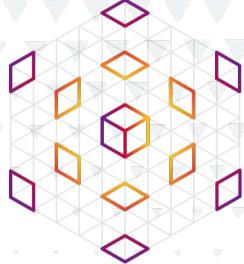


# › Badges





Stay in touch!

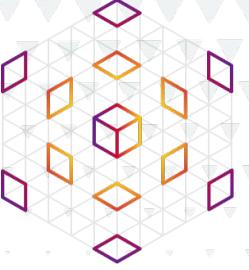
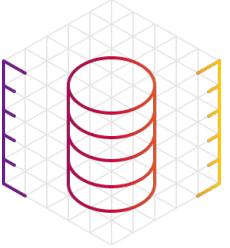


**Discord:** [dtsx.io/discord](https://dtsx.io/discord)

**Academy:** [academy.datastax.com](https://academy.datastax.com)

**Workshops:** [datastax.com/workshops](https://datastax.com/workshops)

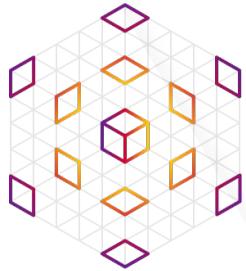
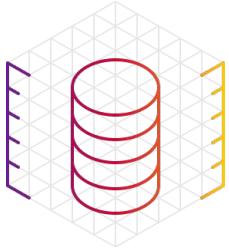
**YouTube:** [@DataStax Developers](#)



› You are not alone in your journey

- Discord:** [dtsx.io/discord](https://dtsx.io/discord)
- StackOverflow<sup>(\*)</sup>:** [stackoverflow.com/questions/tagged/cassandra](https://stackoverflow.com/questions/tagged/cassandra)
- DBA Stack Exchange<sup>(\*)</sup>:** [dba.stackexchange.com/questions/tagged/cassandra](https://dba.stackexchange.com/questions/tagged/cassandra)

*(\*) For best results, follow the cassandra tag*



# Thank You