



# Recs at Reasonable Scale

*NVIDIA RecSys Summit 2022*



Jacopo Tagliabue  
Coveo

# Ciao!



Entrepreneur (co-founder of [Tooso](#)), now [Director of AI](#) (TSX:CVO).

[OS contributor](#) (1000+ Github stars) and [evangelist](#) on “reasonable scale” ML and e-commerce tech.

Researcher ([25+ research papers](#) in top ML venues), [Prof. of MLSys](#) at NYU, Co-organizer of SIGIR eCom.

# Anyone Can Cook



# Anyone Can Cook Build Great RecSys



README.md

## You Don't Need a Bigger Boat

An end-to-end (Metaflow-based) implementation of an intent prediction (and session recommendation) flow for kids who can't MLOps good and [wanna learn to do other stuff good too](#).

This is a WIP - check back often for updates.

### Philosophical Motivations

There are plenty of tutorials and blog posts around the Internet on data pipelines and tooling. However:

- they (for good pedagogical reasons) tend to focus on *one* tool / step at a time, leaving us to wonder how the rest of the pipeline works;
- they (for good pedagogical reasons) tend to work in a toy-world fashion, leaving us to wonder what would happen when a real dataset and a real-world problem enter the scene.

This repository (and soon-to-be-drafted written tutorial) aims to fill these gaps. In particular:



# So, why almost nobody does?

- The majority of innovation in e-commerce tech comes from very few, large, public B2C players (with **one exception**).
- Why? And what can we do about it?

ABSTRACT

You Do Not Need a Bigger Boat: Recommendations at Reasonable Scale in a (Mostly) Serverless and Open Stack

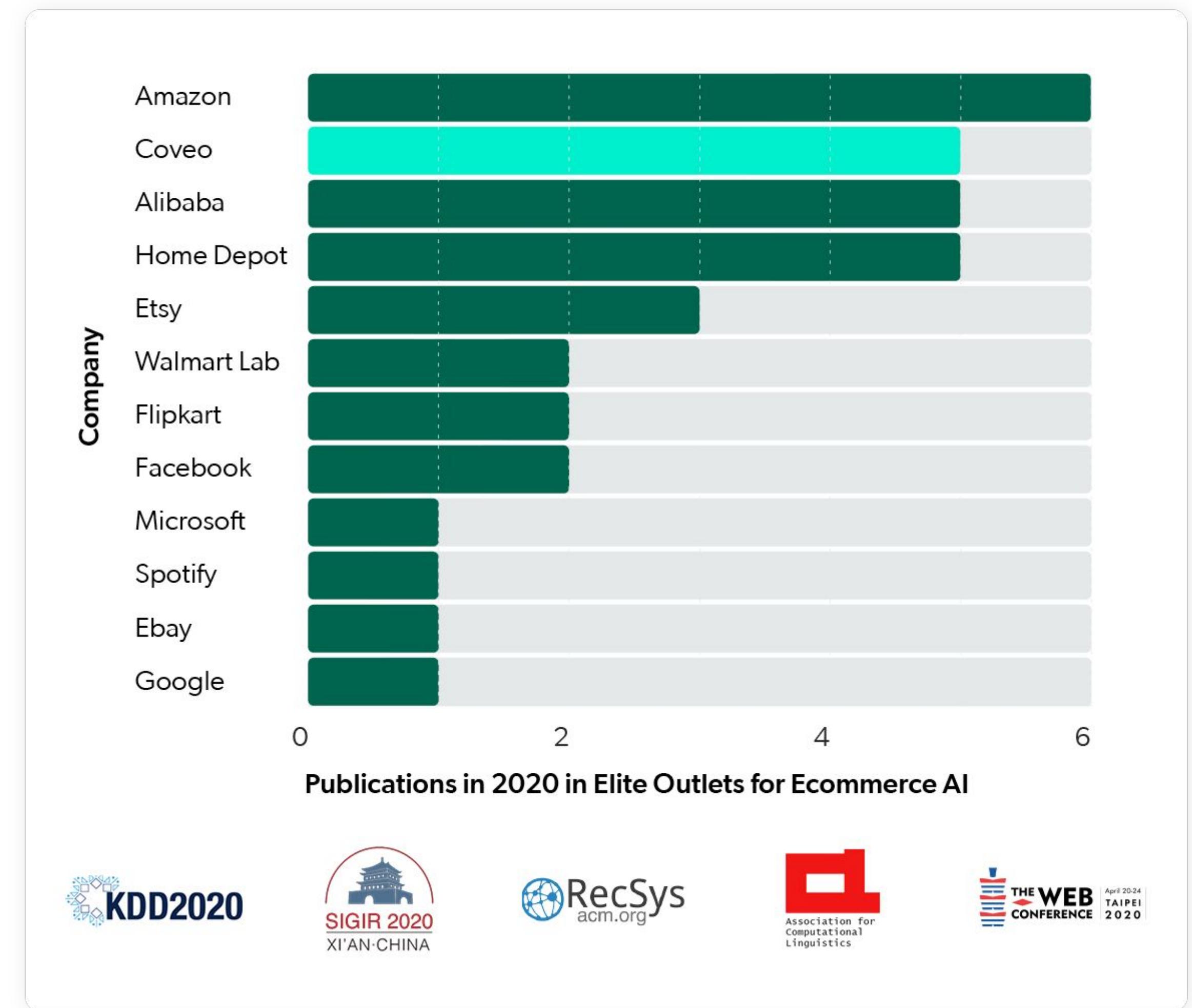
Author:  Jacopo Tagliabue [Authors Info & Claims](#)

RecSys '21: Fifteenth ACM Conference on Recommender Systems • September 2021 • Pages 598-600 • <https://doi.org/10.1145/3460231.3474604>

Online: 13 September 2021 [Publication History](#)

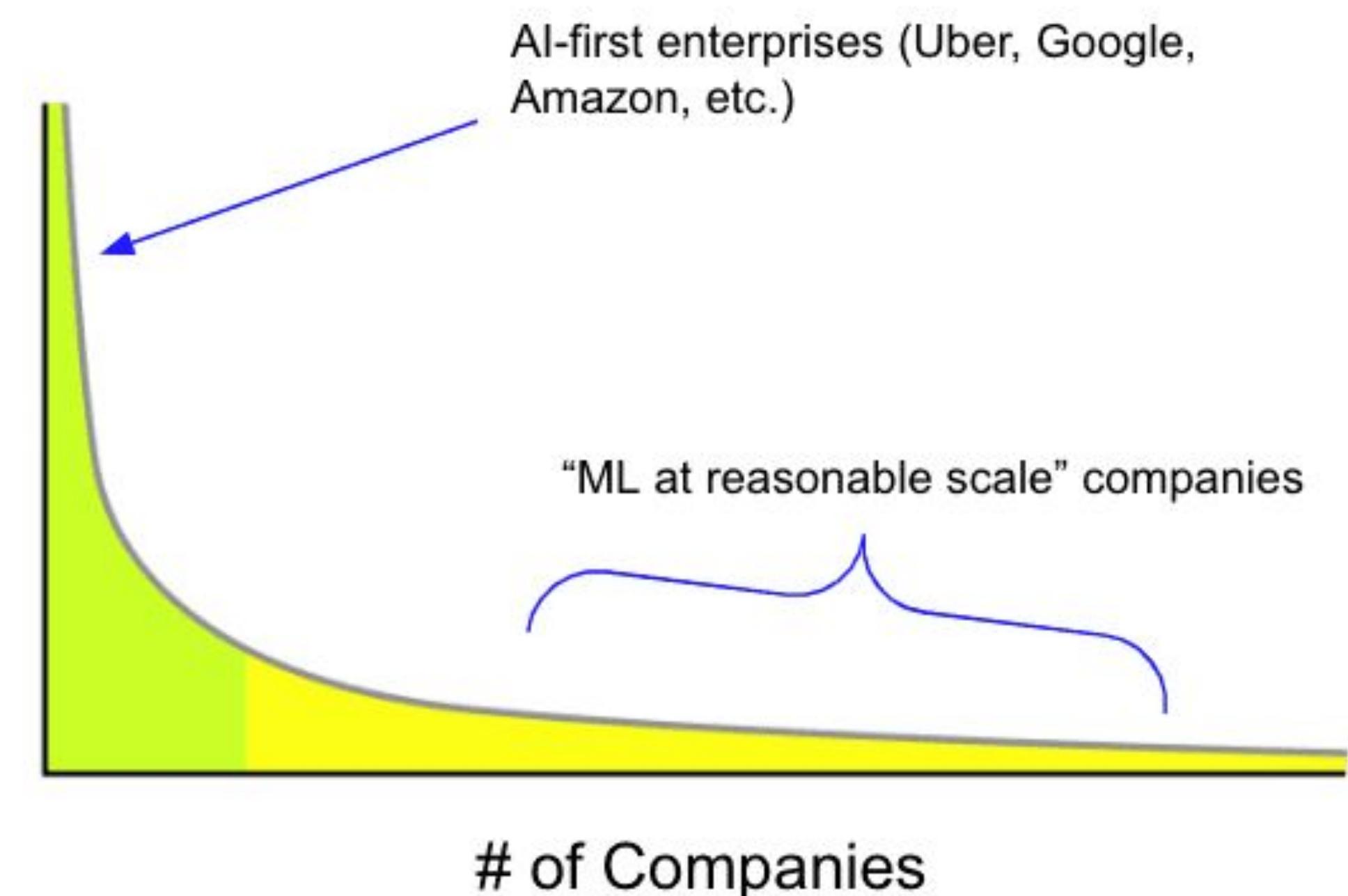
0 547

Get Access



# Barriers to entry are high

- Implementing cutting-edge recSys algorithms is hard.
- Models are only one part of the challenge, MLOps is hard too (GPU provisioning, scheduling, versioning, etc.).
- Most literature is skewed towards the resources, needs, and the problems of few companies.



If you think getting a model into production needs:

- a data engineer for the ETL
- a data scientist for the modelling
- a machine learning engineer for MLOps
- a software engineer for production

Where does it stop?

What's the minimum team size to get started?

# Big Tech fetish is **real**

... and we are not making it easy to start.

 **Edge#176: Meta's New Architecture for Build AI Agents that Can Reason Like Humans and Animals**

The new architecture pr  
AI agents.

Mar 24 

 **Edge#72: Tecton – The Enterprise Feature Store Built by the Creators of Uber's ML Platform**

Mar 18 

*This is an example of what our subscribers receive from us and AI.*

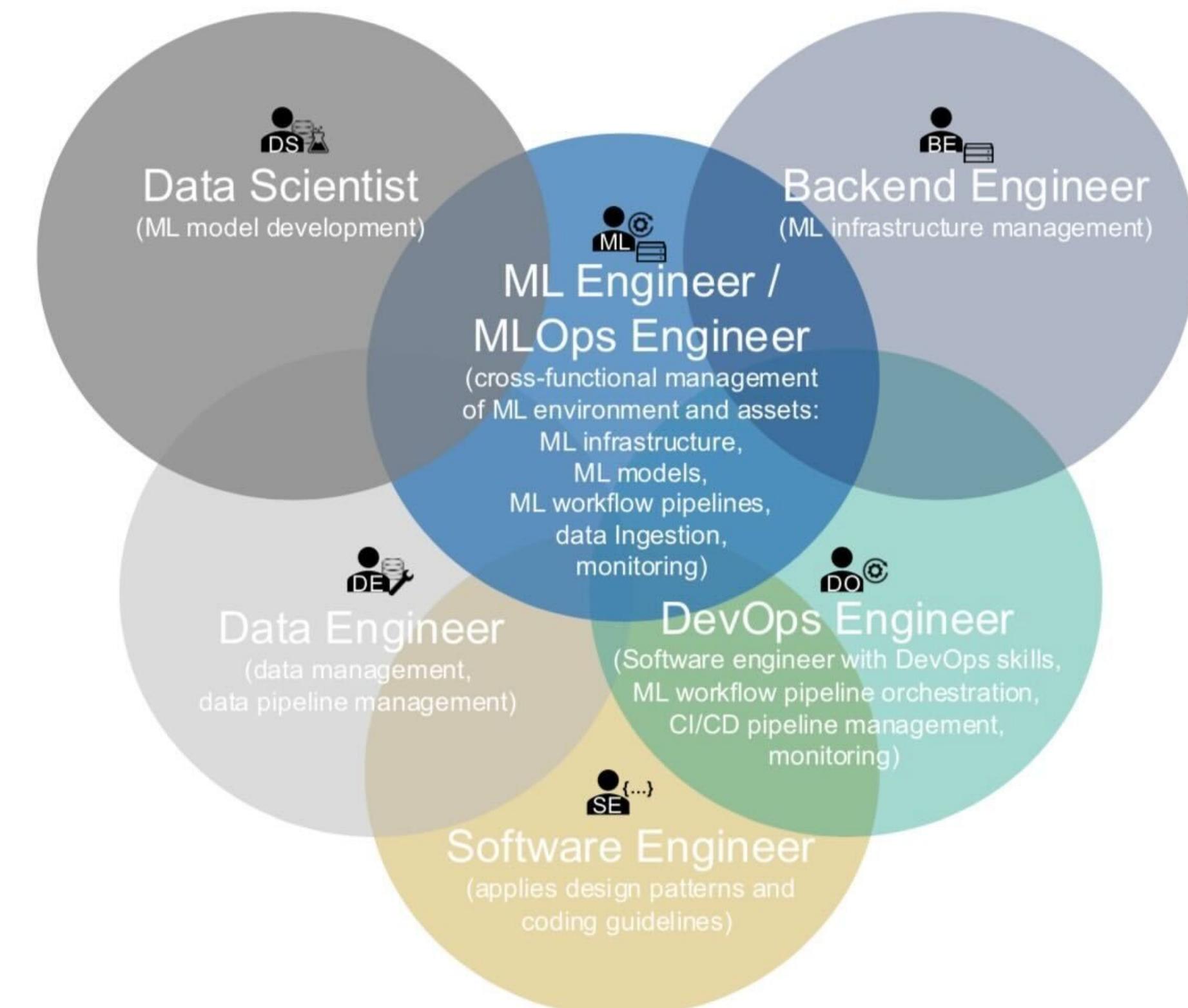
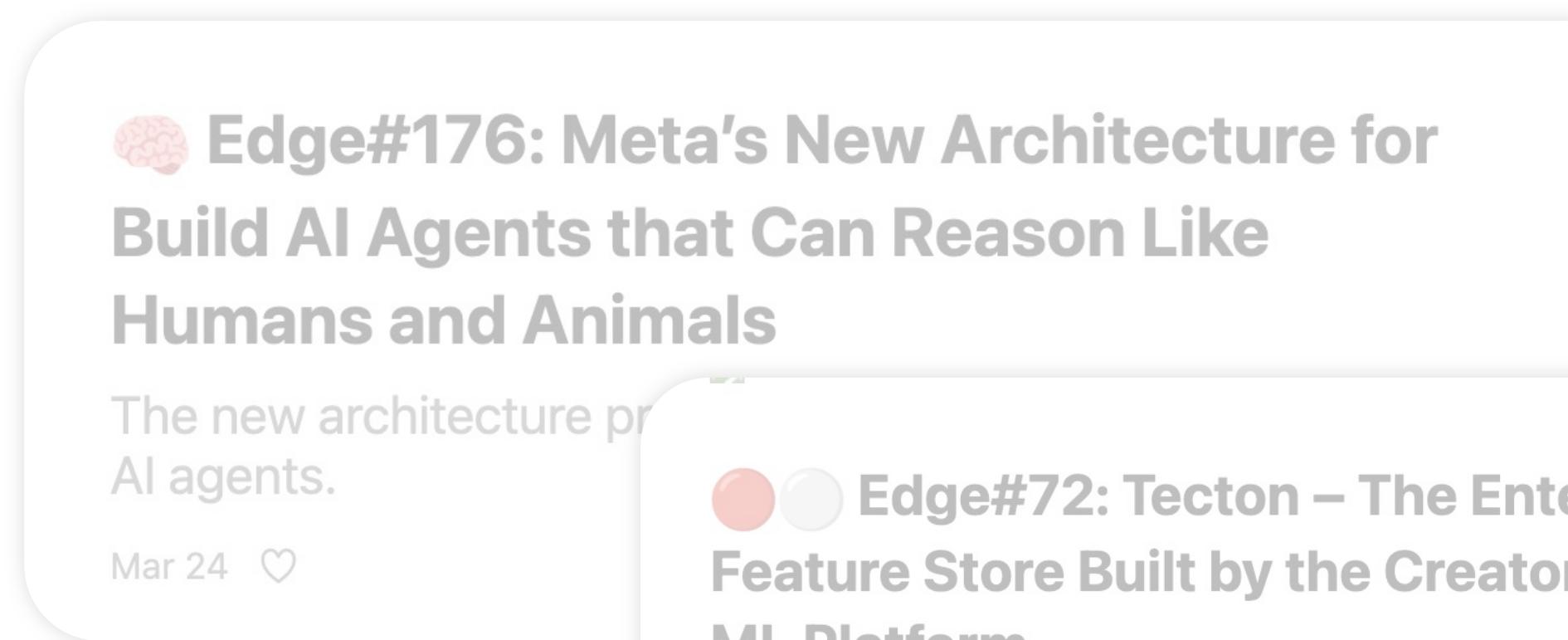
 **New Week, New AI Super Model**

Weekly news digest curated by the industry insiders

Apr 10   

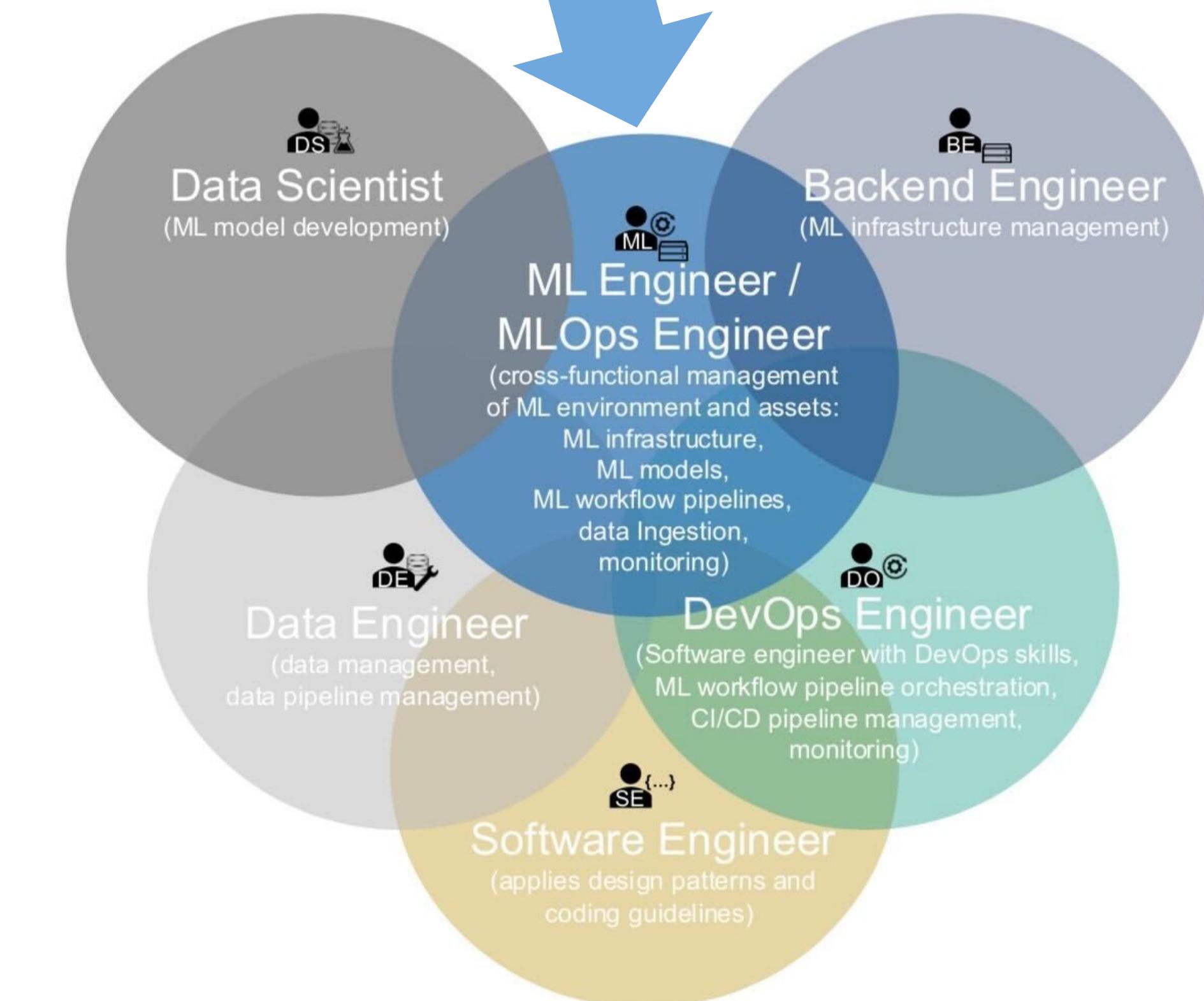
# Big Tech fetish is real

... and we are not making it easy to start.



# Big Tech fetish is real

I'm Batman, an MLOps Engineer



Edge#176: Meta's New Architecture for Build AI Agents that Can Reason Like Humans and Animals

The new architecture pr  
AI agents.

Mar 24 ❤

Edge#72: Tecton – The Enterprise Feature Store Built by the Creators of Uber's ML Platform

Mar 18 ❤

New Week, New AI Super Model

Weekly news digest curated by the industry insiders

Apr 10 ❤ ○ ↗

# Myth-Busting RecSys

- RecSys algorithms are hard.  
Yes, but good abstractions mean we don't start from scratch anymore.
- MLops is hard too.  
Yes, but a new ecosystem means we can buy what we don't want to maintain / build.
- Literature is skewed towards the problems of few companies.  
Yes, but sometimes you and Big Tech have the same problem (e.g. testing).

# Talk is cheap, show me the code

README.md

## recs-at-resonable-scale

Recommendations at "Reasonable Scale": joining dataOps with deep learning recSys with Merlin and Metaflow

### Overview

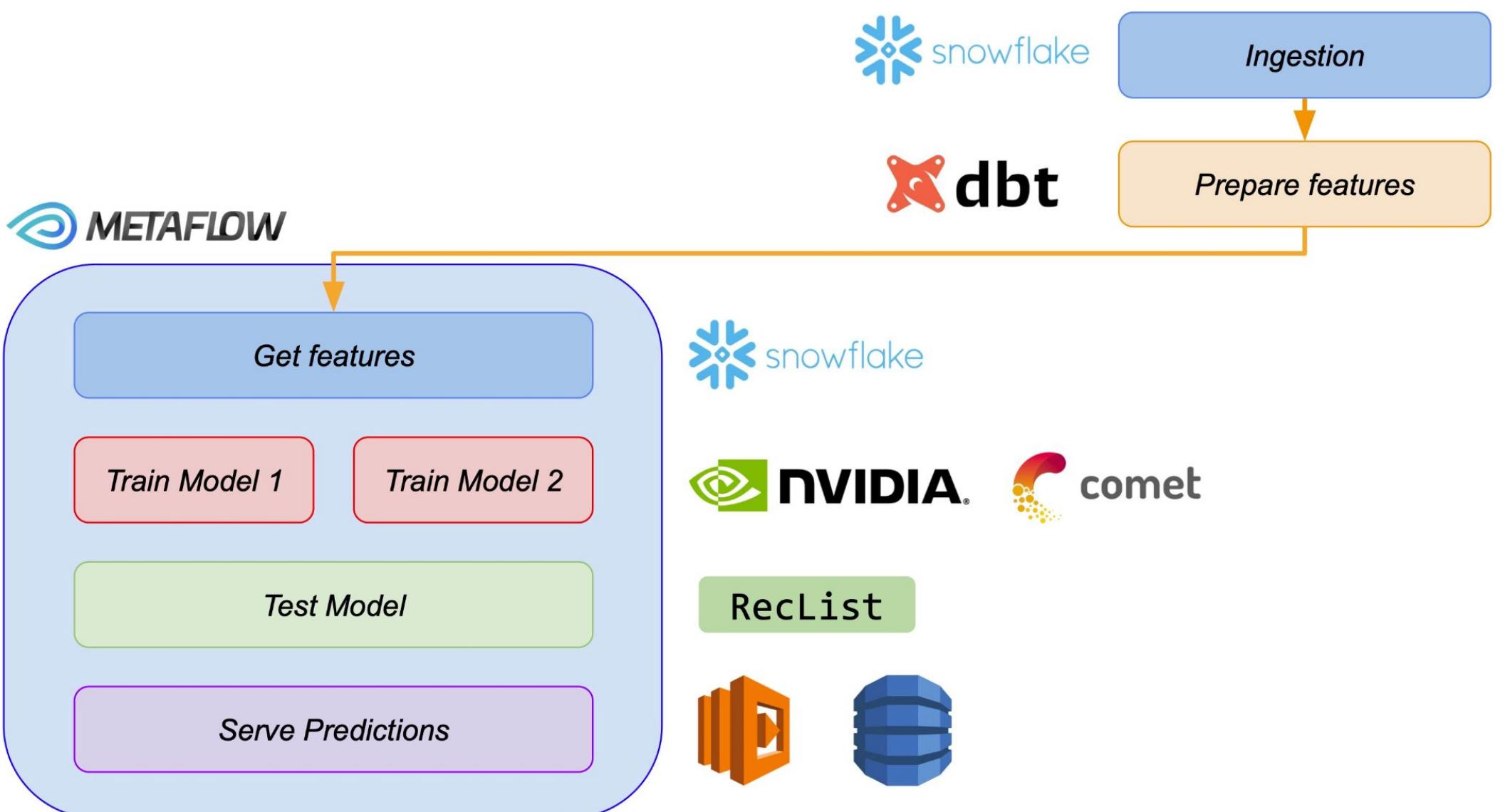
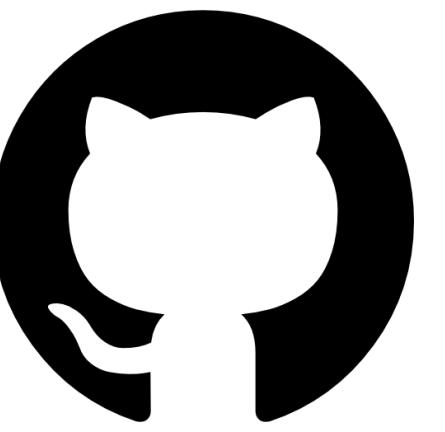
July 2022: this is a WIP, come back often for updates, a blog post and my NVIDIA talk (FORTHCOMING)!

This project is a collaboration with the Outerbounds and NVIDIA Merlin teams, in an effort to release as open source code a realistic data and ML pipeline for cutting edge recommender systems "that just works". Anyone can do great ML, not just Big Tech, if you know how to pick and choose your tools.

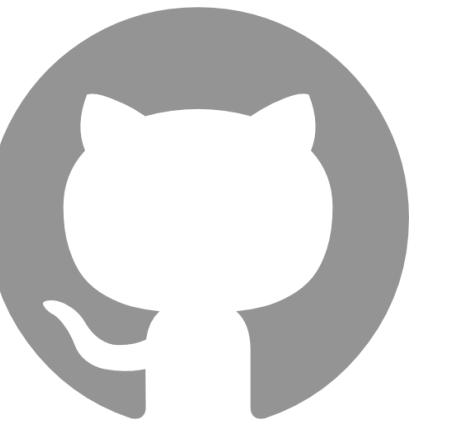
TL;DR: (after setup) a single ML person is able to train a cutting edge deep learning model (actually, several versions of it in parallel), test it and deploy it without any explicit infrastructure work, without talking to any DevOps person, without using anything that is not Python or SQL.

As a use case, we pick a popular RecSys challenge, user-item recommendations for the fashion industry: given the past purchases of a shopper, can we train a model to predict what he/she will buy next? In the current V1.0, we target a typical offline training, cached predictions setup: we prepare in advance the top-k recommendations for our users, and store them in a fast cache to be served when shoppers go online.

Our goal is to build a pipeline with all the necessary real-world ingredients:



# Talk is cheap, show me the code



**“After a 10 minutes setup, one MLE can run an entire RecSys pipeline - from data to API - with no explicit infrastructure work.”**

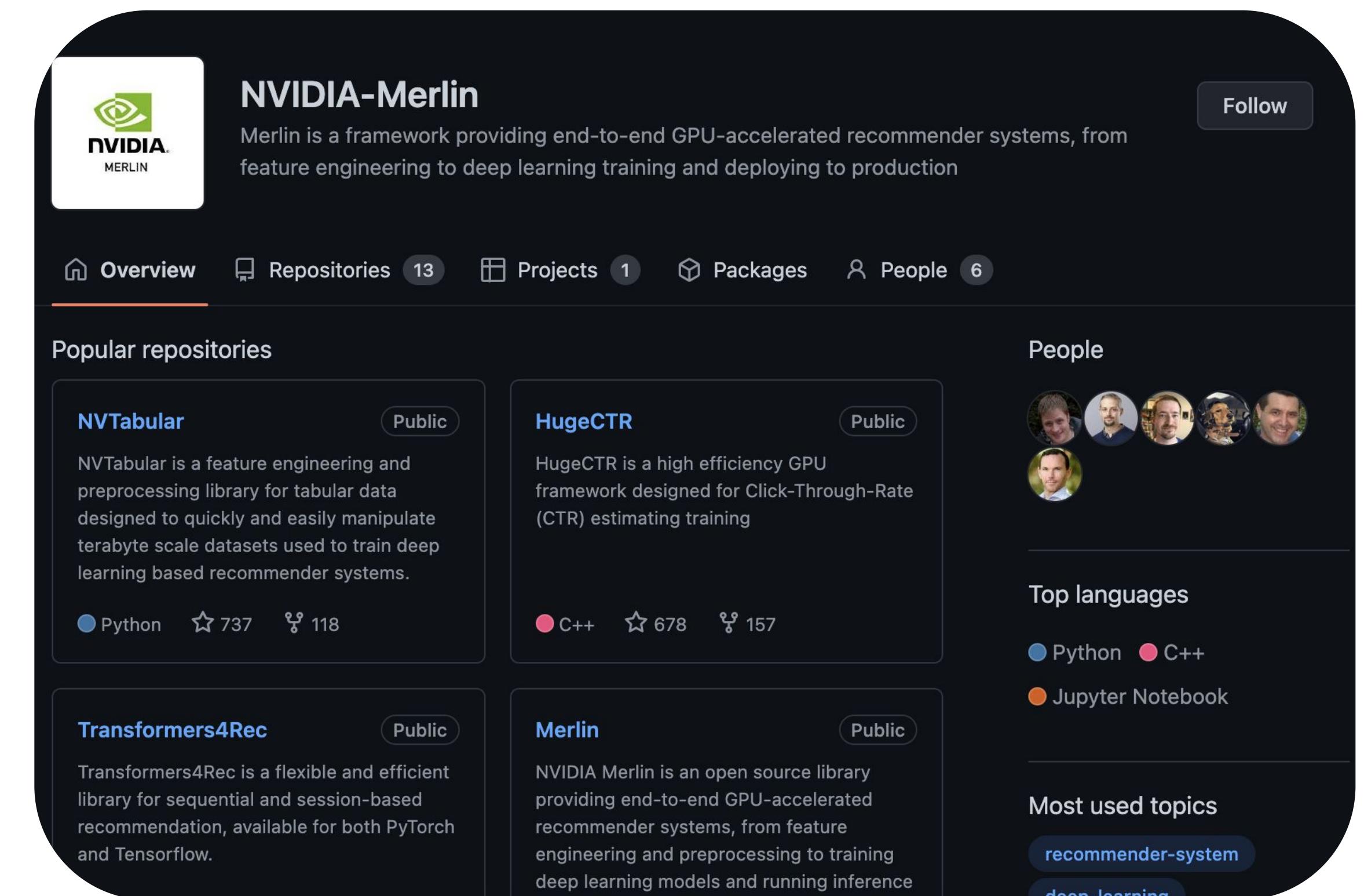


“What one fool **can do**, another can.”

RecSys that **just works**

# RecSys as an algorithmic challenge

- The “HuggingFace” moment of **RecSys**: the field has converged to a set of best practices and building blocks that are effective in most cases.
- **NVIDIA Merlin** is an open source library for DL-based RecSys “with an easy-to-use API” exposing common architectures.  
**How easy?**



# RecSys as an algorithmic challenge

- The “HuggingFace” moment of **RecSys**: the field has converged to a set of best practices and building blocks that are effective in most cases.
- **NVIDIA Merlin** is an open source library for DL-based RecSys “with an easy-to-use API” exposing common architectures.  
**How easy?**

```
# train the model
model = mm.TwoTowerModel(
    train.schema,
    query_tower=mm.MLPBlock([128, 64], no_activation_last_layer=True),
    samplers=[mm.InBatchSampler()],
    embedding_options=mm.EmbeddingOptions(infer_embedding_sizes=True)
)
model.compile(optimizer="adam", run_eagerly=False, metrics=[mm.RecallAt(10), mm.NDCGAt(10)])
model.fit(train, validation_data=valid, batch_size=self.hypers['BATCH_SIZE'], epochs=int(self.N_EPOCHS))
# test the model on validation set and store the results in a MF variable
self.metrics = model.evaluate(valid, batch_size=1024, return_dict=True)
print("\n\n====> Eval results: {}\n\n".format(self.metrics))
```

**VERY EASY!**

MLOps without much Ops

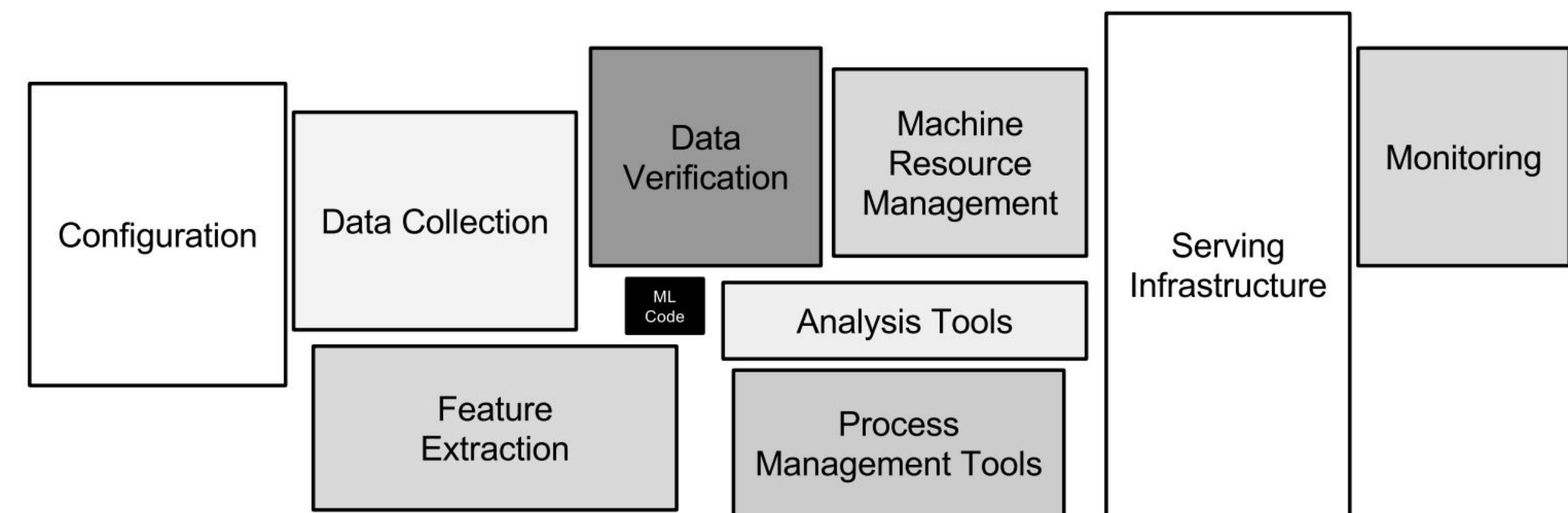
# RecSys as an MLOps challenge

- **DataOps (ingestion, preparation, quality) is as important as MLOps**

- Leverage the Modern Data Stack to provide ultra-scalable data support

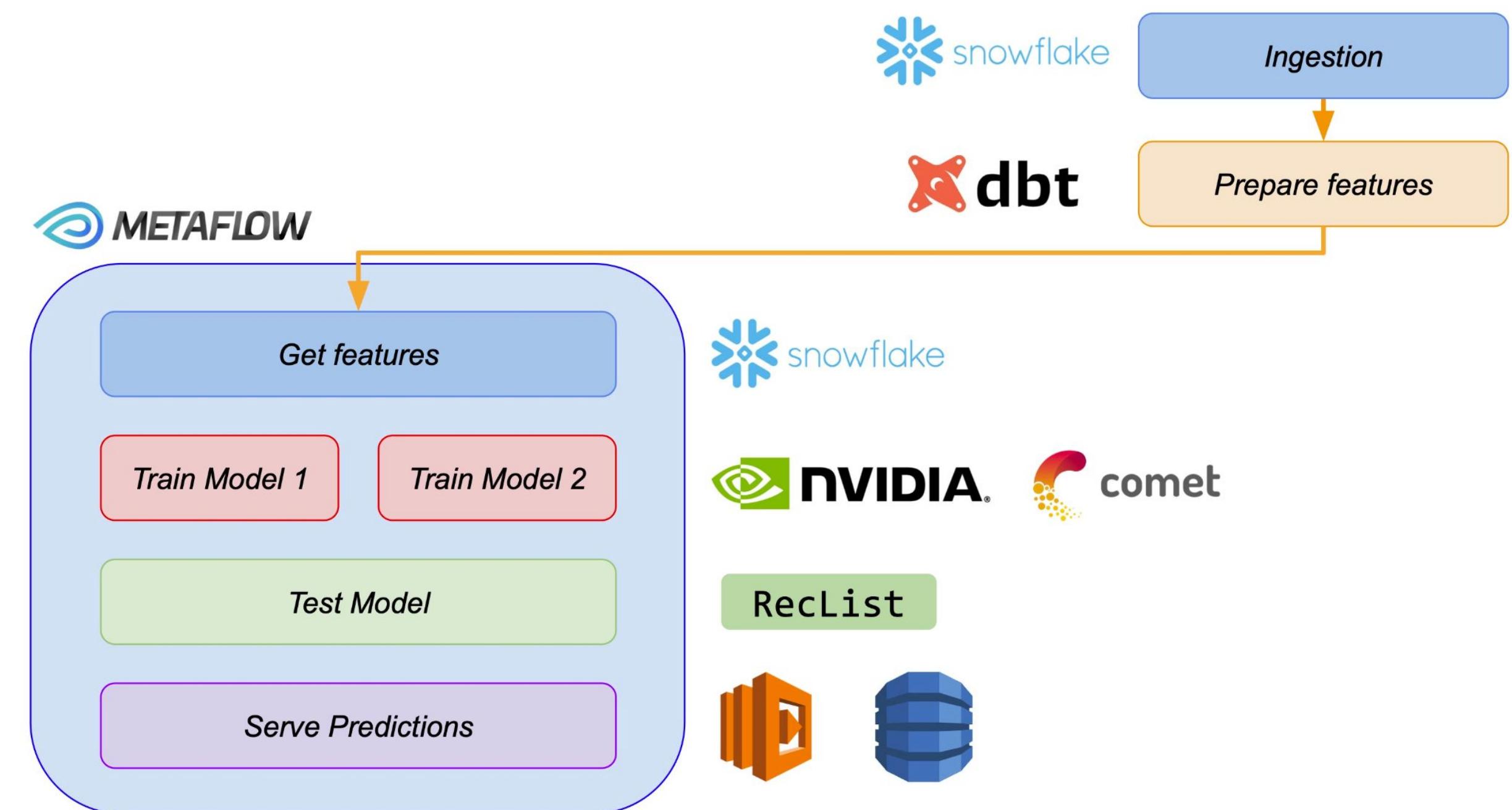
- **Infrastructure management is tedious and error prone**

- Strive for velocity, embrace PaaS / FaaS for non-core activities
  - From “infra as code” to “infra in code”: **Python is all you need.**



# RecSys as an MLOps challenge

1. **Add NVIDIA Merlin**, to train cutting-edge models.
2. **Add the MDS**, to provide scalable DataOps.
3. **Add a Python orchestrator**, to provide Python abstractions over GPU training, scheduling, versioning.
4. **Add a tracker**, to track your progress through experiments.
5. **Add a FaaS+SaaS endpoint**, to serve predictions to shoppers.



Start with the problem (and the stack)  
**you have**, not the one you wish you had.

Wild Wild Tests

# RecSys as a testing challenge

- **RecSys evaluation** is typically a point-wise metric over held-out data points, say HIT RATE, or MRR.
- RecSys are some of the most ubiquitous form of Machine Learning systems, and *even one mistake can have dire consequences*.
- Typical testing protocols fall short of measuring real-world generalization performance.



maxwell ✅

@maxwellstrachan

...

absolutely horrifying: amazon's algorithm identified a product that people were using to commit suicide.

it then started to recommend additional products that made the suicide easier to execute.

[nytimes.com/2022/02/04/tec...](https://nytimes.com/2022/02/04/tech/amazon-suicide-recommendations.html)

# RecSys as a testing challenge

The ~~heart~~ power-law is deceitful above all things: two models may have the same exact “accuracy” while behaving very differently.

HR:  $57 / 110 = 0.52$



VS

HR:  $57 / 110 = 0.52$



#: 100, H: 50, HR: 0.5

#: 10, H: 7, HR: 0.7

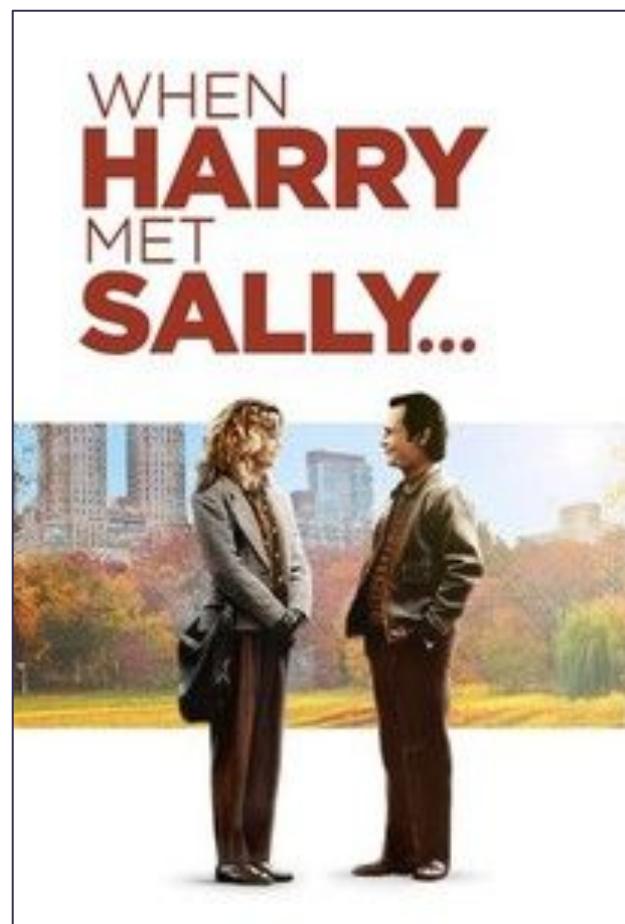
#: 100, H: 56, HR: 0.56

#: 10, H: 1, HR: 0.1

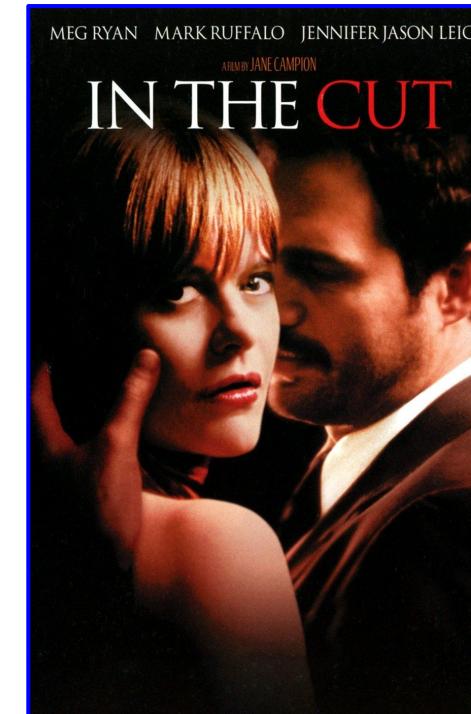
# RecSys as a testing challenge

Not all predictions are wrong in the same way: in particular, some mistakes are worse than others.

Test Truth



Model #1



Model #2



*This is a reasonable mistake!*

# RecSys as a testing challenge

- **Introducing RecList**, an open source package for behavioral testing.
- RecList provides abstractions over datasets and best practices to help with testing and monitoring!
- RecList is a community project supported by forward-looking companies: **check them out!**

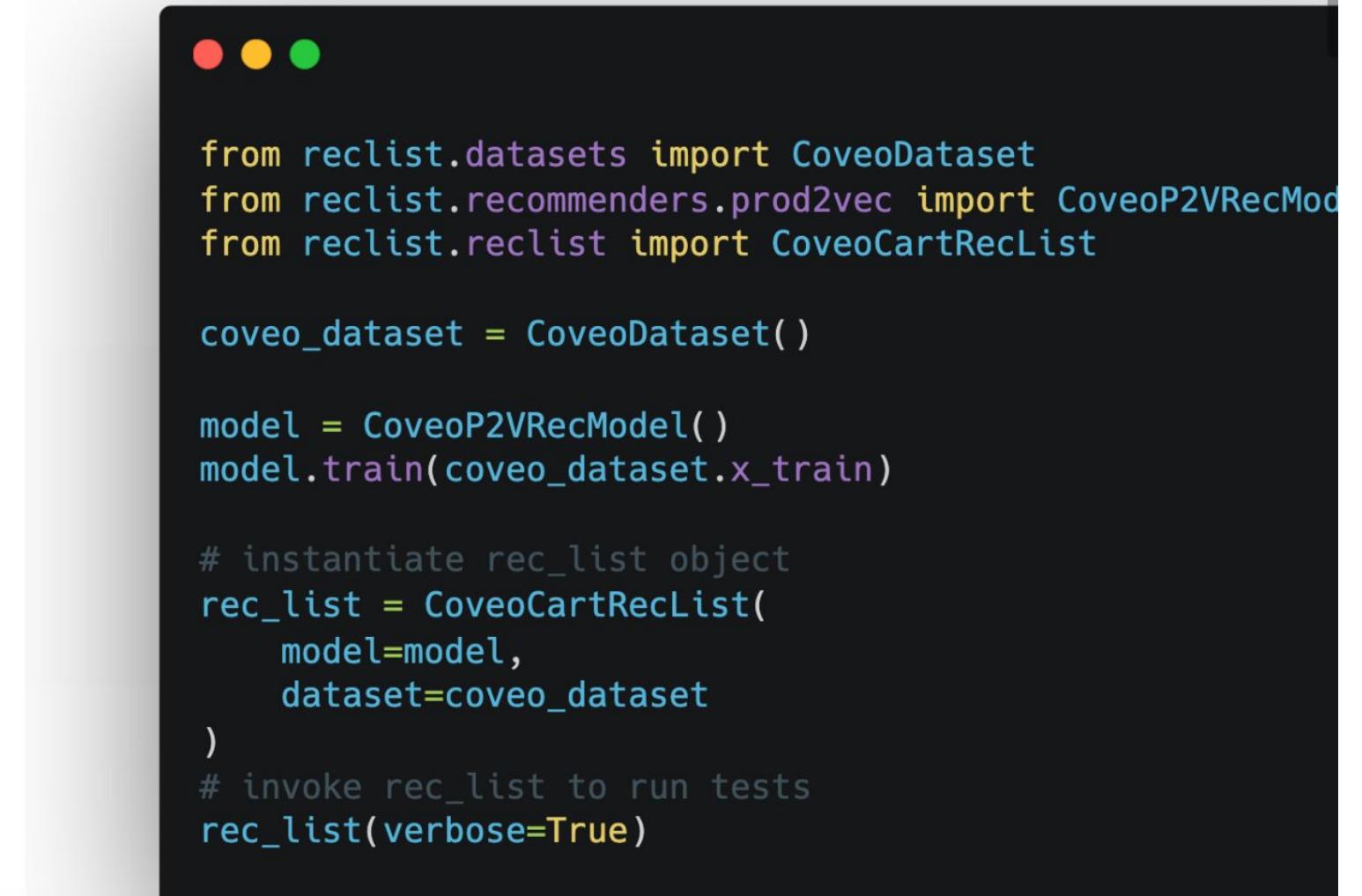


## RecList

RecList is an open source library providing behavioral, “black-box” testing for recommender systems. Inspired by the pioneering work of [Ribeiro et al. 2020](#) in NLP, we introduce a general plug-and-play procedure to scale up behavioral testing, with an easy-to-extend interface for custom use cases.

To streamline comparisons among existing models, RecList ships with popular datasets and ready-made behavioral tests: read the [our TDS blog post](#) as a gentle introduction to the main use cases, and try out our [colab](#) to get started with the code.

We are actively working towards our beta, with new capabilities and improved documentation and tutorials: check back



```
from reclist.datasets import CoveoDataset
from reclist.recommenders.prod2vec import CoveoP2VRecModel
from reclist.reclist import CoveoCartRecList

coveo_dataset = CoveoDataset()

model = CoveoP2VRecModel()
model.train(coveo_dataset.x_train)

# instantiate rec_list object
rec_list = CoveoCartRecList(
    model=model,
    dataset=coveo_dataset
)
# invoke rec_list to run tests
rec_list(verbose=True)
```

<https://reclist.io/>

# Bonus: Join us at CIKM!

**CIKM 22 Data Challenge:**  
<https://reclist.io/cikm2022-cup/>



The image shows a dark blue rectangular banner with rounded corners. In the top left corner, the word "EVALRS" is written in yellow capital letters. In the top right corner, there are three white links: "CHALLENGE", "ORGANIZERS", and "IMPORTANT DATES". The background of the banner is a night-time city skyline with illuminated buildings, including a prominent skyscraper with a golden spire. Overlaid on the center of the image is the word "EvalRS" in a white, italicized serif font. At the bottom, the words "A ROUNDED EVALUATION OF RECOMMENDER SYSTEMS" are written in large, bold, white capital letters.

**EVALRS**

CHALLENGE   ORGANIZERS   IMPORTANT DATES

*EvalRS*

A ROUNDED EVALUATION OF  
RECOMMENDER SYSTEMS

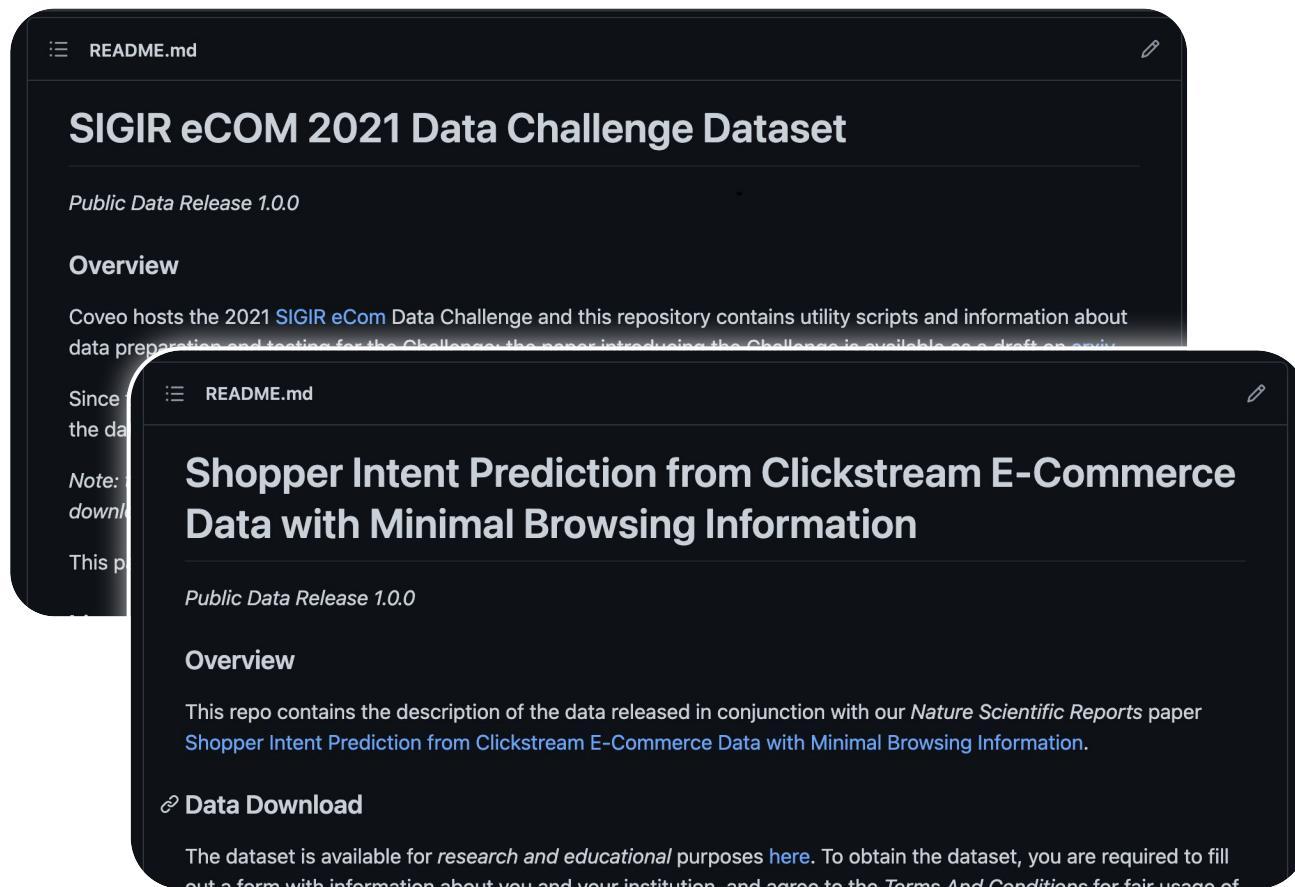
The **future** at Reasonable Scale

# A roadmap between RecSys and MLOps

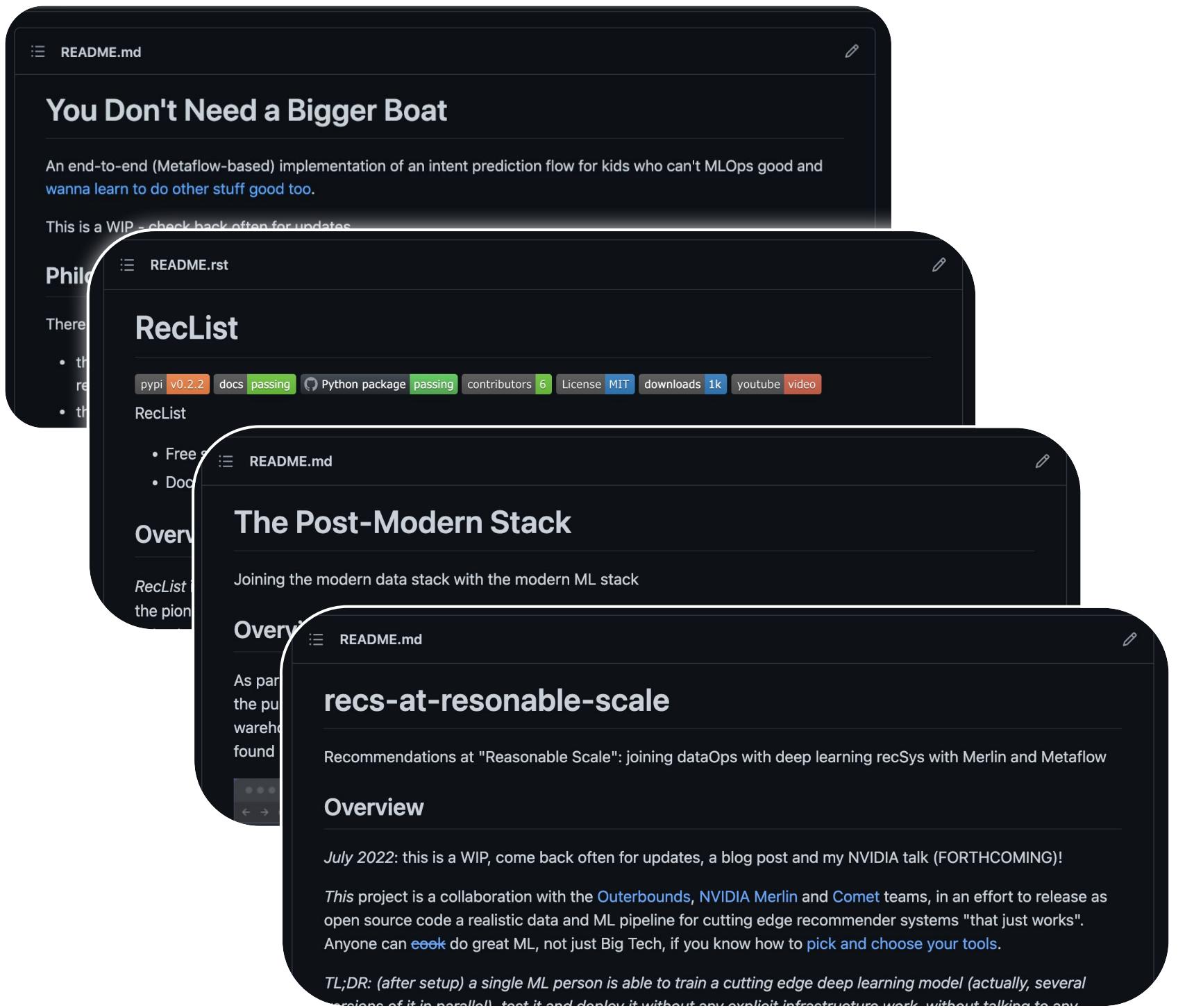
1. **Training and Optimizing Offline:** recommendations are often more similar to *interventions* than predictions.
2. **Deployment:** while “simple” deployment is a commodity, more complex models are still hard to handle.
3. **Testing and Monitoring:** performance on test sets is a small part of a rounded evaluation process.

# Ask *not* what the RS can do for you...

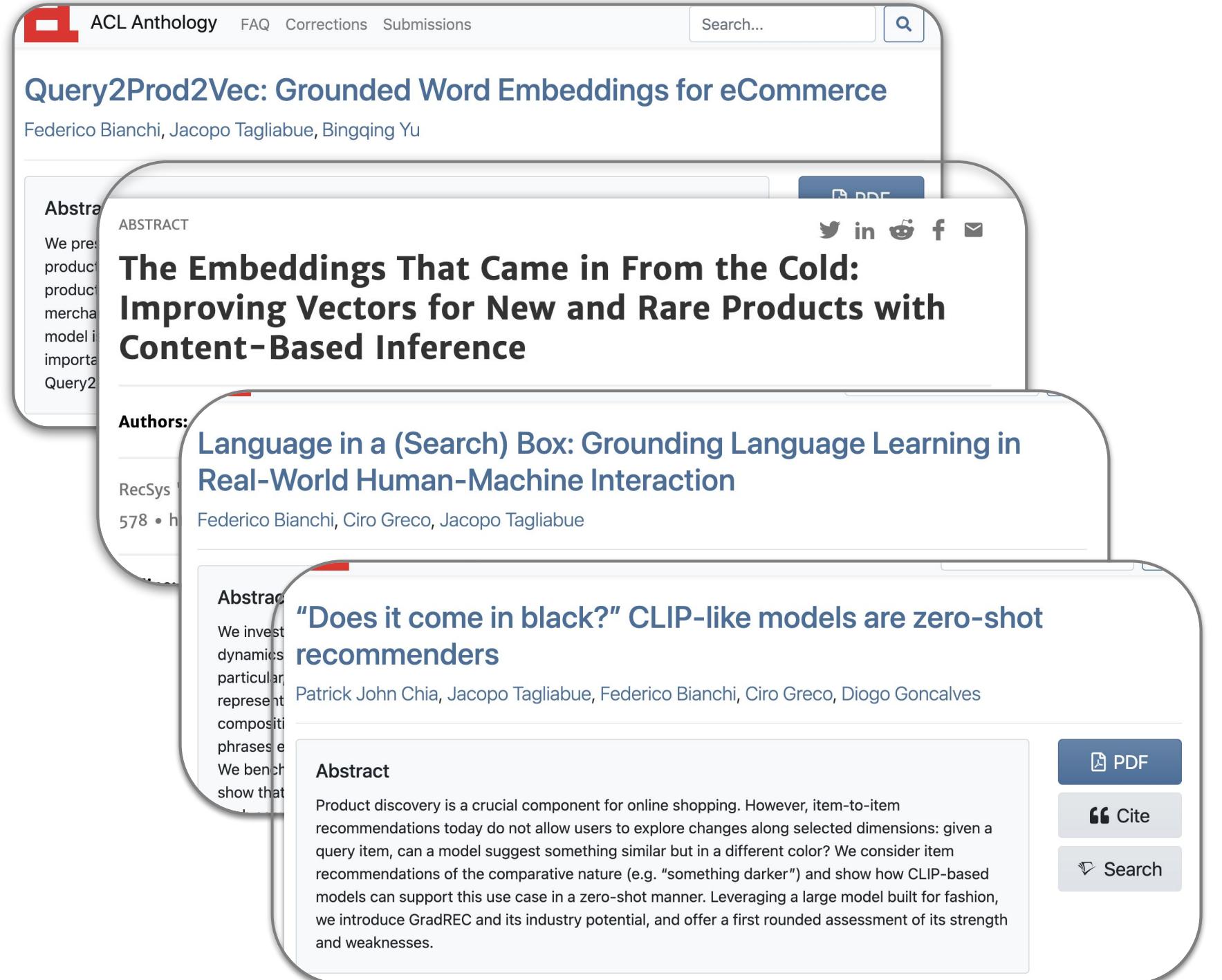
## Open Data



## Open Source

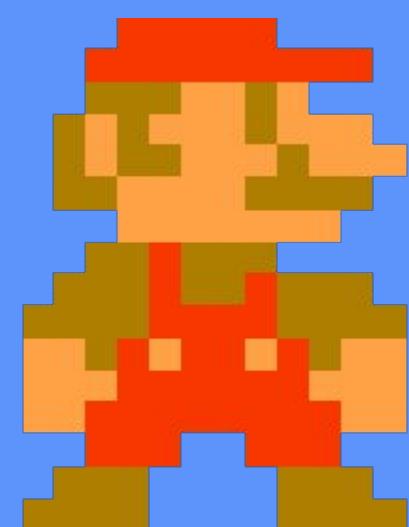
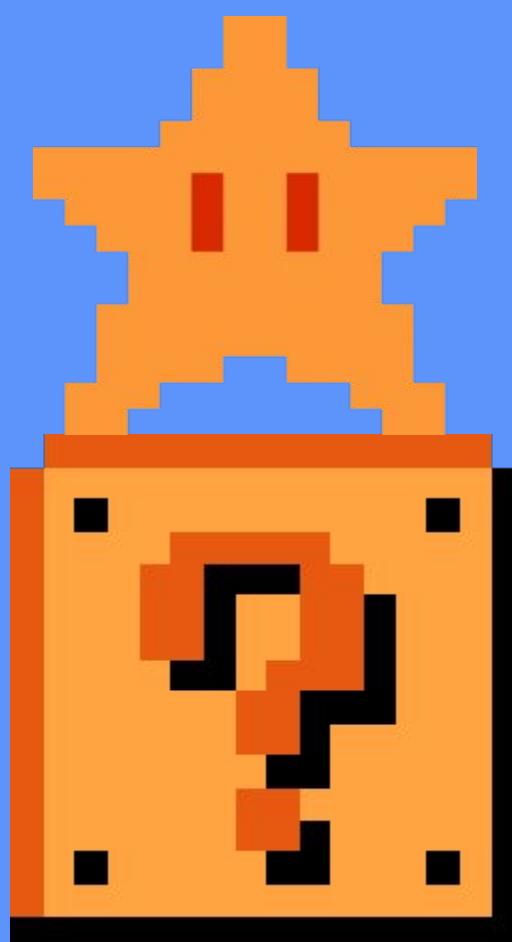


## Open Science



Check out / share / star our open  
source **repositories**!

Wanna work with us? Get in touch!



“We can only see a short distance ahead, but we  
can see plenty there that needs to be done.”

See you, RecSys cowboys