

# Managing ML Pipelines: Feature Stores and the Coming Wave of Embedding Ecosystems

VLDB 2021

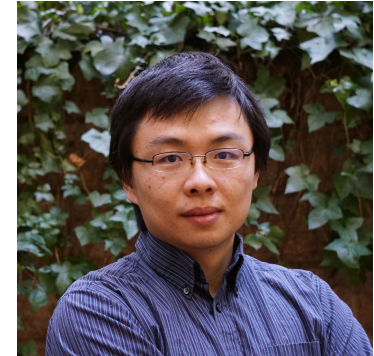
# Speakers



Laurel Orr  
Stanford



Atindriyo Sanyal  
Uber AI



Xiao Ling  
Apple



Megan Leszczynski  
Stanford



Karan Goel  
Stanford

# Modern ML Pipelines

ML pipelines help engineers build and deploy models

**Standardization**

**Reproducibility**

**Easier to Maintain**

## H2O AutoML

Scalable AutoML in H2O-3 Open Source



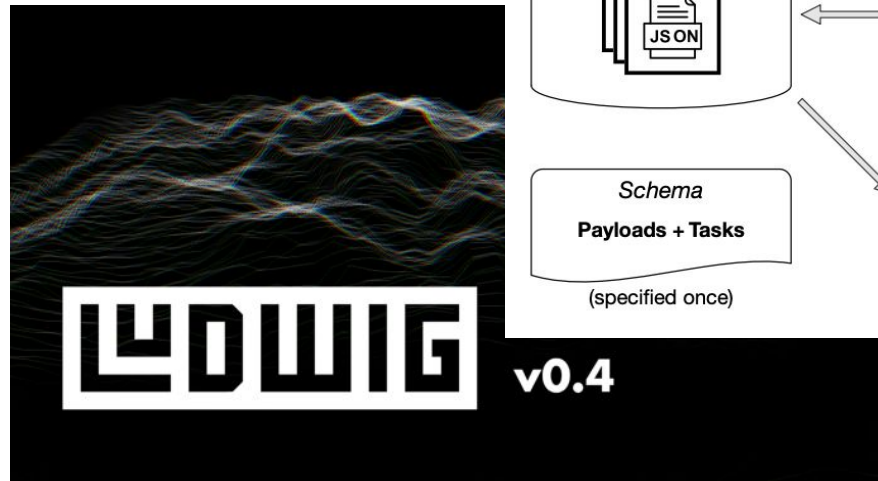
Data Preprocessing



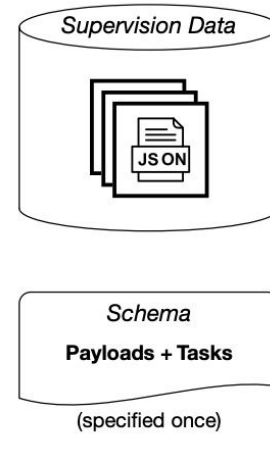
Model Generation



Ensembles



## Overton



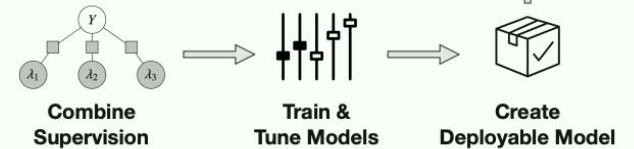
### Actions

- Add/augment slices
- Add labeling functions
- Add synthetic examples

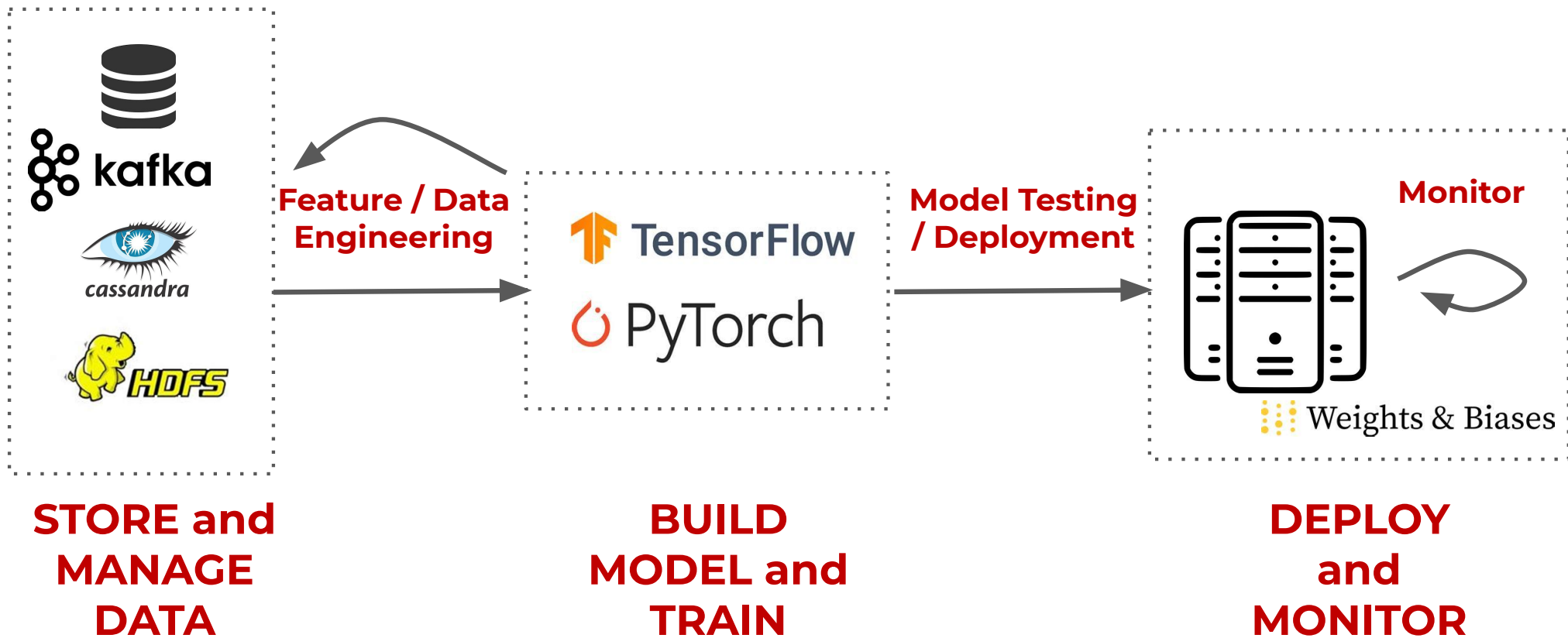
### Fine-grained quality reports

	task 1	task 2
slice 1	✓	✓
slice 2	✓	✗
slice 3	✓	✓
slice 4	✗	✓

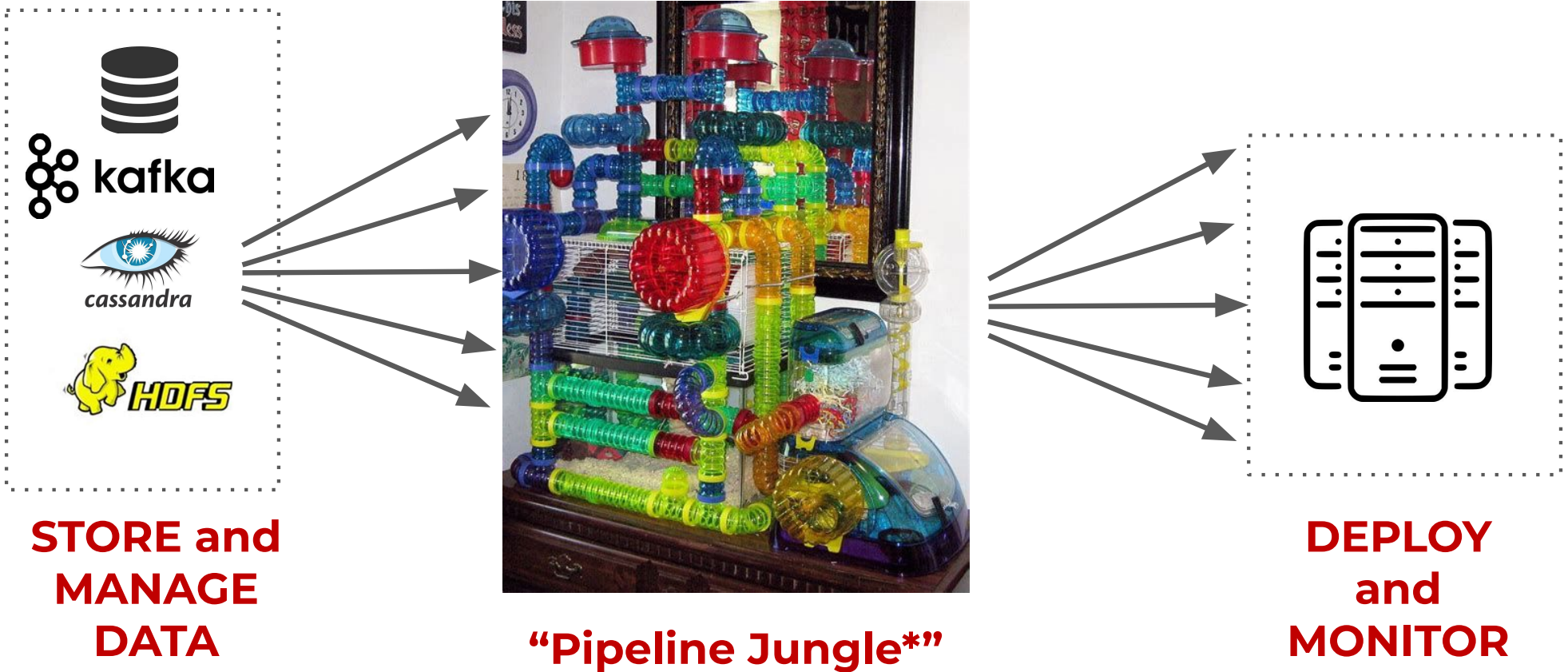
### Overton



# Engineer Workflow Today



# Engineer Workflow of Yesteryear (< 2017-8)



\*Sculley, David, et al. "Hidden technical debt in machine learning systems." *Neurips* (2015)

# The “Pipeline Jungle” Experience

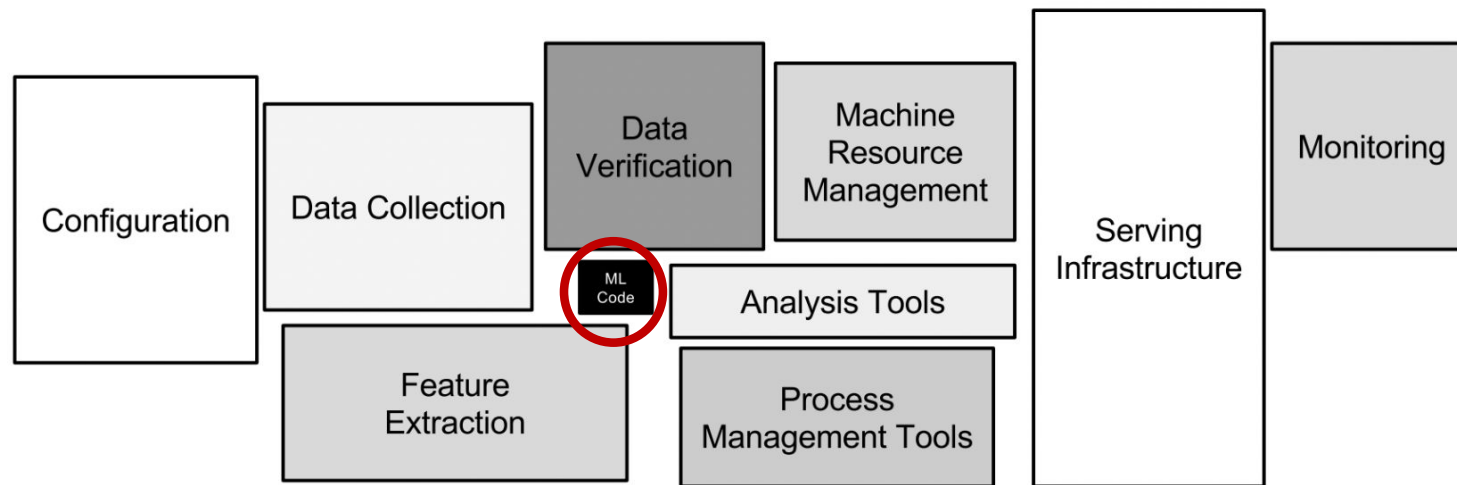
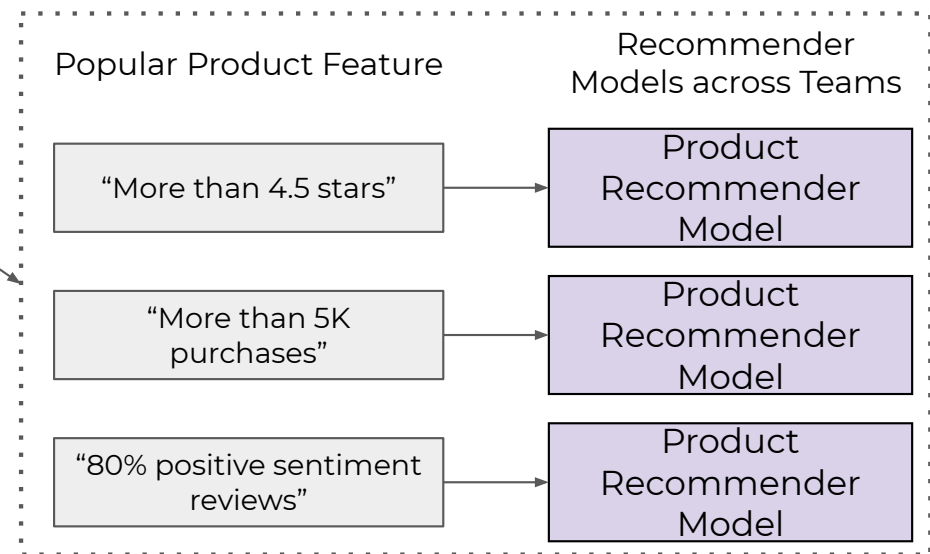


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

# The “Pipeline Jungle” Experience

The challenges to deploying a model:

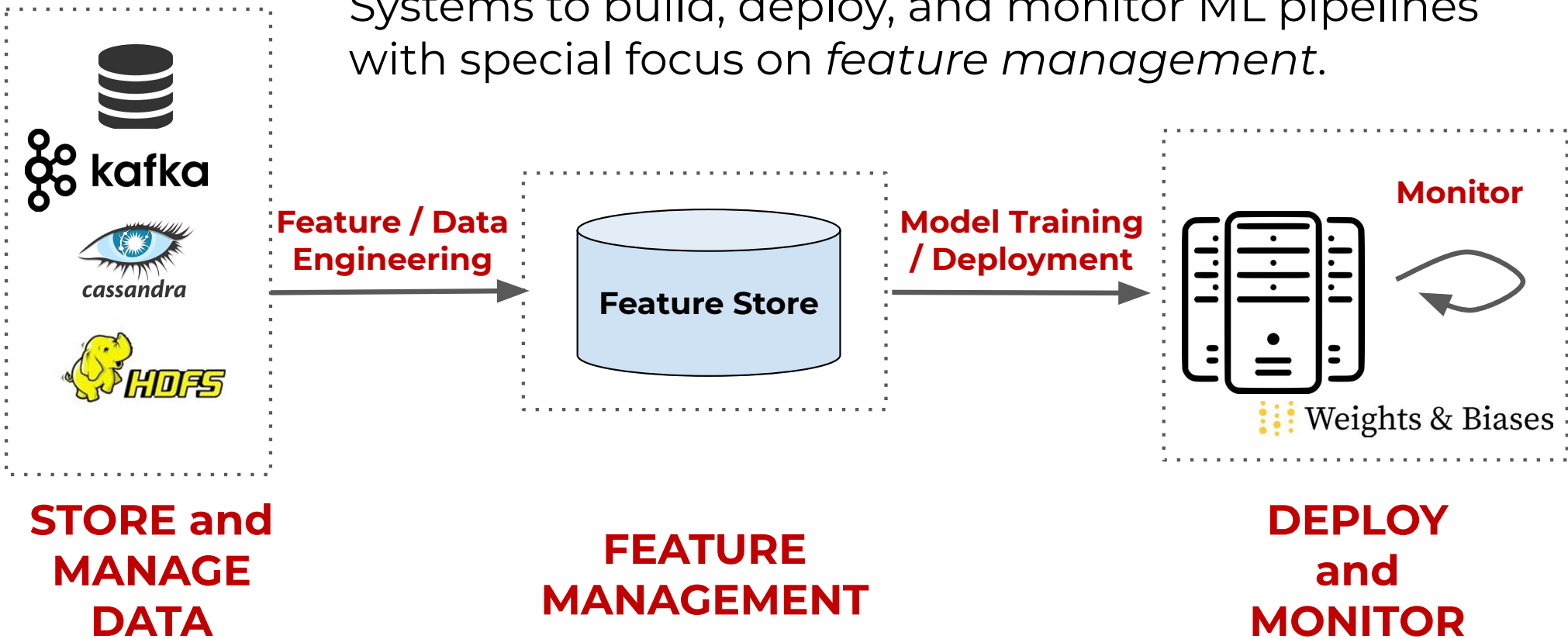
- One-off feature definitions
- Lack of reproducibility
- Inconsistent storage
- No standard evaluations and testing
- Difficult to detect and recover from errors
- ...



PART 1

# Feature Store Solution

Systems to build, deploy, and monitor ML pipelines with special focus on *feature management*.



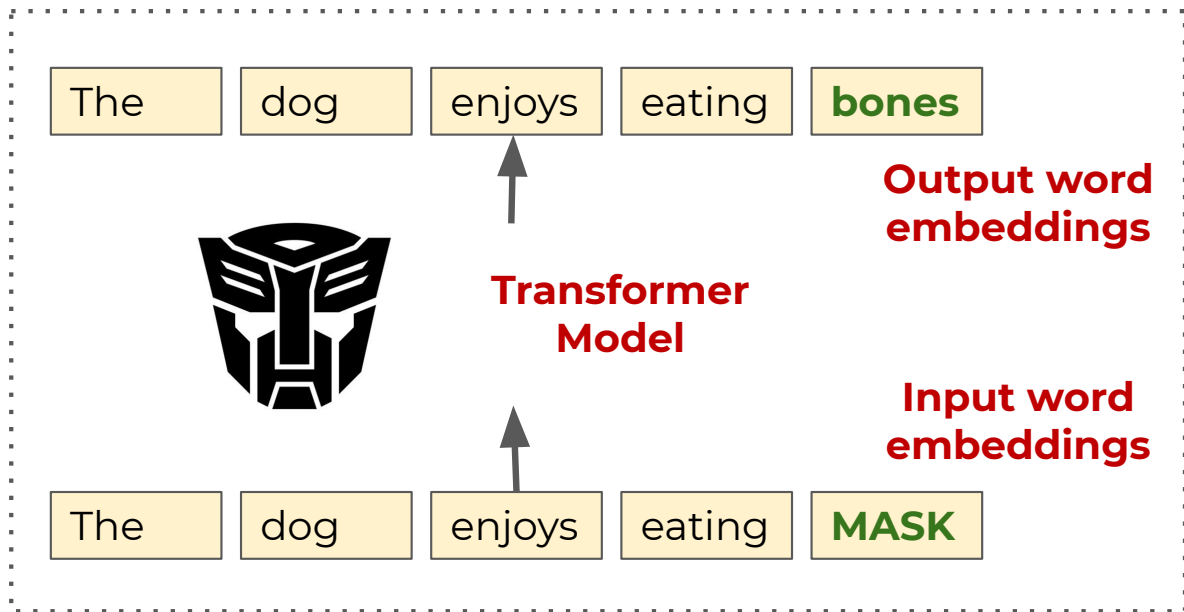


# Enter Self-Supervision

Paradigm where models learn embedding representations of underlying training data *without* needed manually labels.

# Self-Supervision Example: Transformers and MLM

Learn word embeddings by train a language model to predict a masked word in a given context.

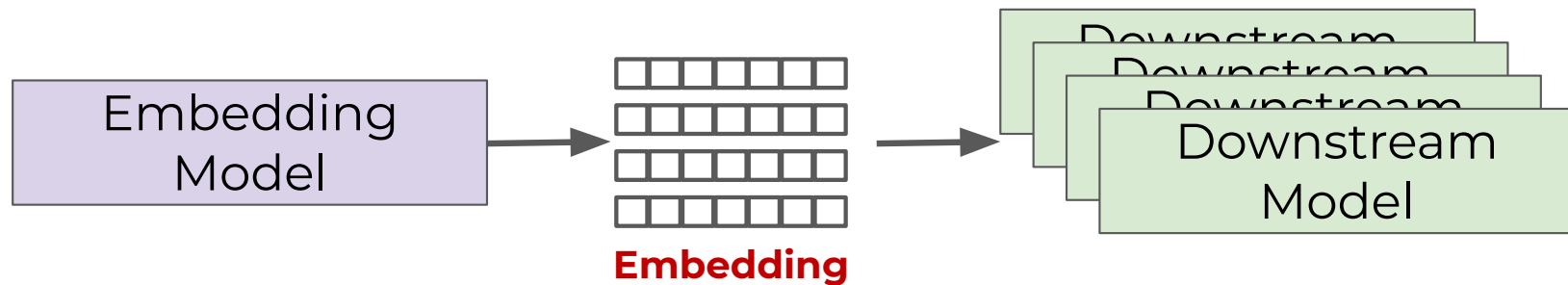


Word embeddings encode contextual information.

# Enter Self-Supervision

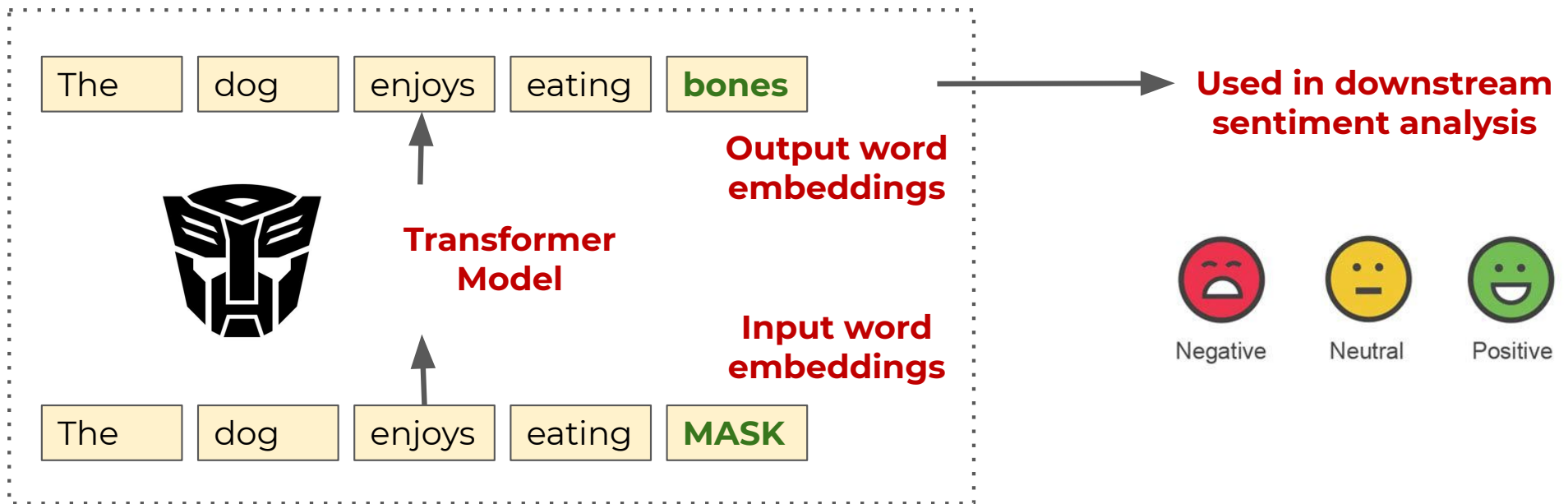
Paradigm where models learn embedding representations of the underlying training data *without* manual labels.

Embeddings are then used in downstream models.



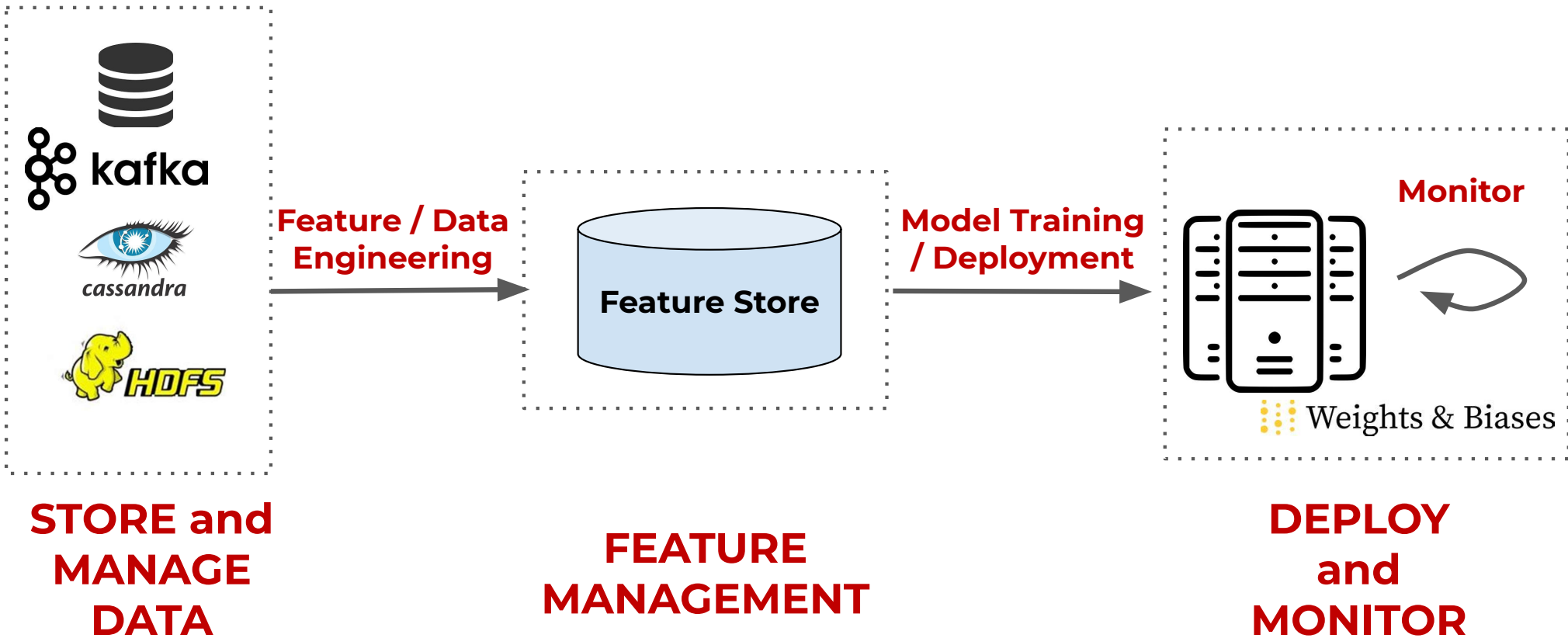
# Self-Supervision Example: Transformers

Learn word embeddings by train a language model to predict a masked word in a given context.



Word embeddings encode contextual information.

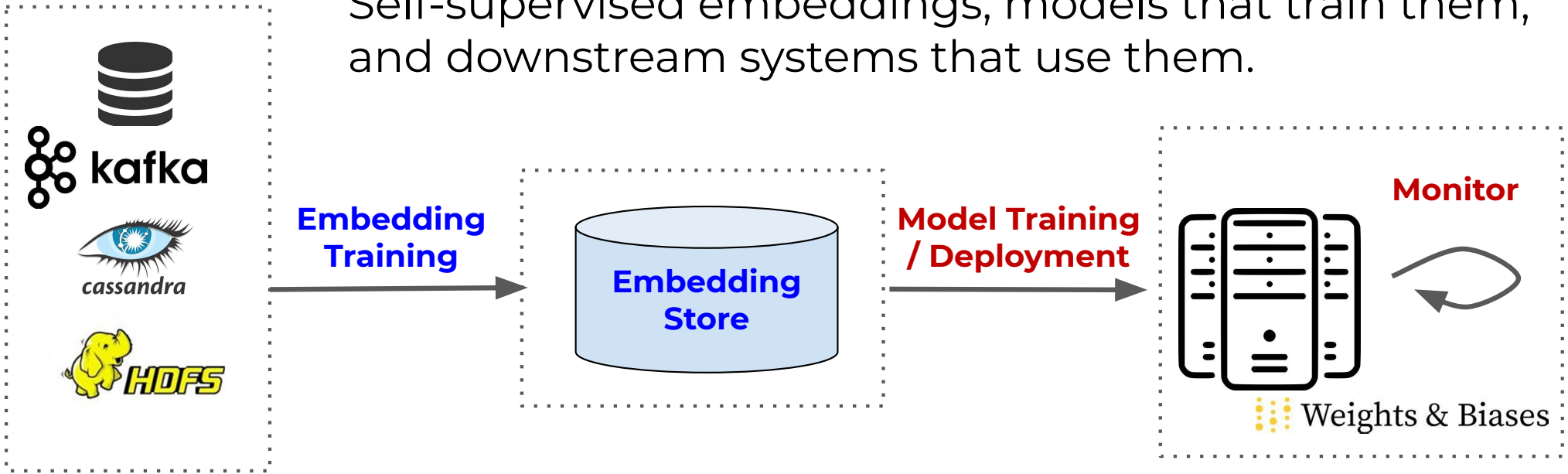
# Recall Feature Store Solution



PART 2

# Embedding Ecosystems

Self-supervised embeddings, models that train them, and downstream systems that use them.



**STORE and  
MANAGE  
DATA**

**EMBEDDING  
MANAGEMENT**

**DEPLOY  
and  
MONITOR**

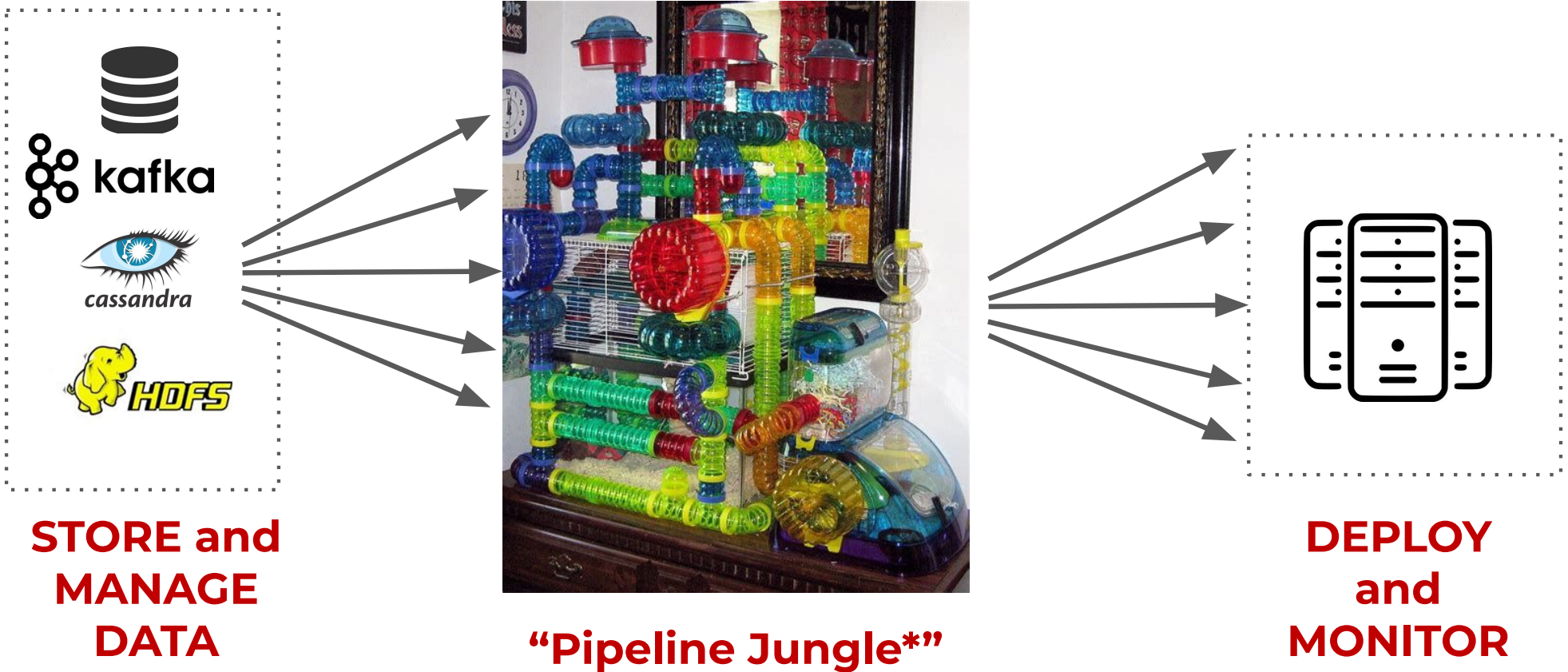
Reduction in engineer effort

Higher quality downstream systems

One embedding for multiple tasks

# Feature Stores

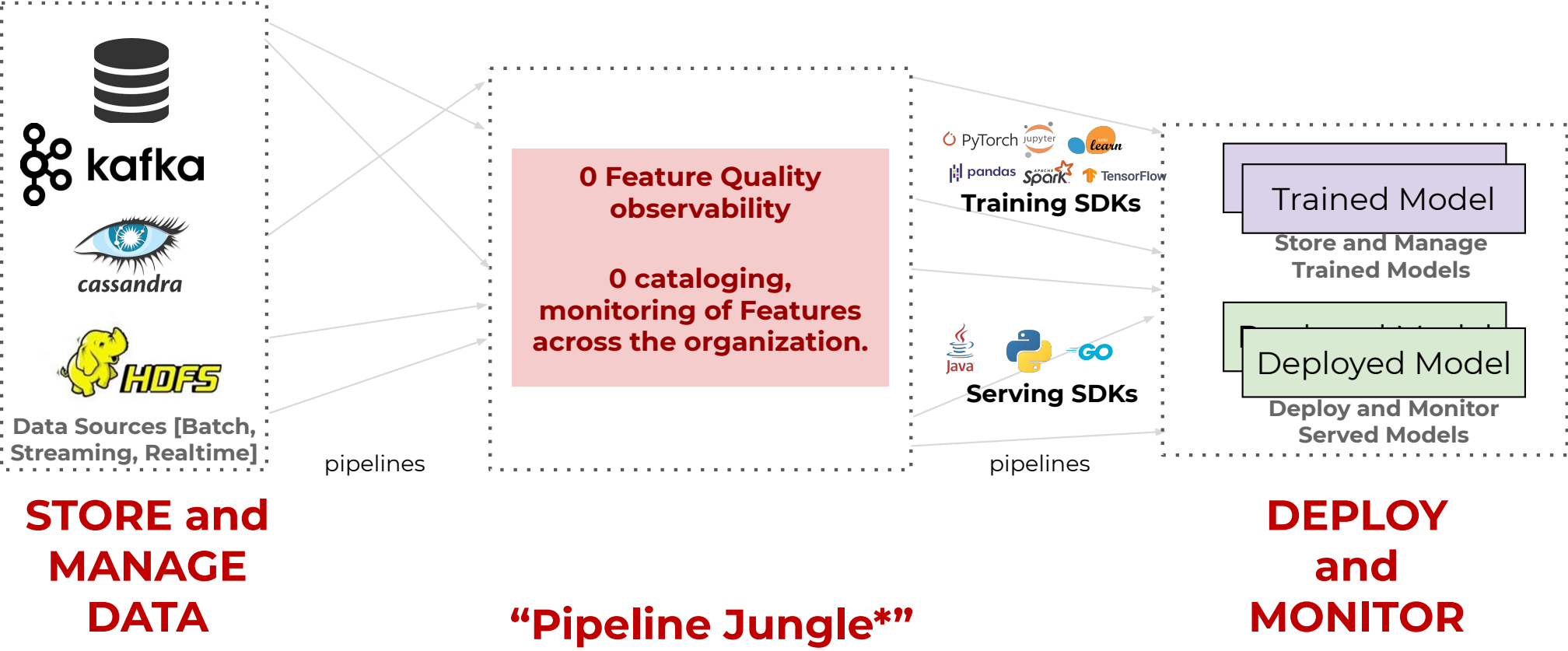
## Engineer Workflow of Yesteryear (< 2017-8)



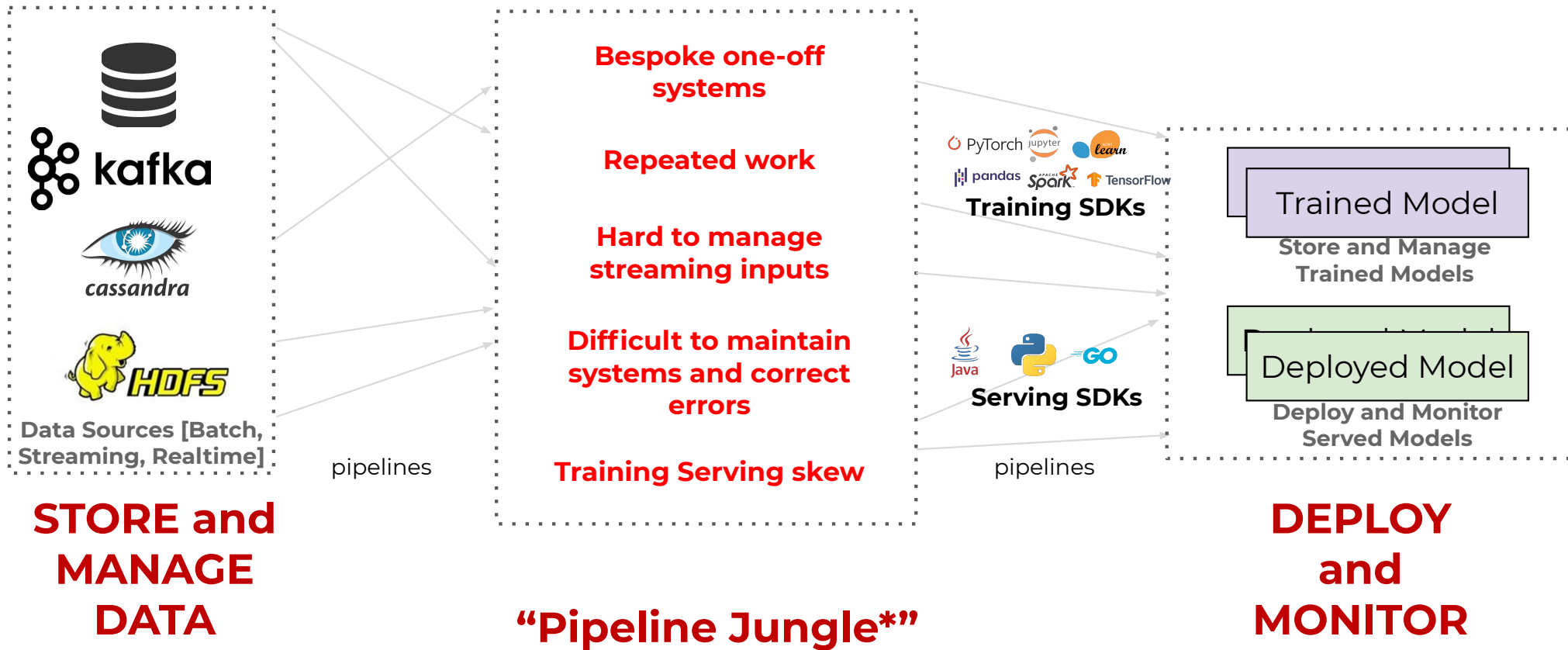
\*Sculley, David, et al. "Hidden technical debt in machine learning systems." *Neurips* (2015)



# Lack of Feature Management



## Lack of Feature Management



## Lack of Feature Management



### Days/Weeks to make data ML ready

- Materializing Features from **various data sources**.
- **Duplicating code while materializing** in training & serving
- **No guarantees** of training-serving parity



### Near 0 monitoring of Features

- **No Feature health metrics** out of the box (due to the various sources problem)
- **No online-offline parity monitoring**, leading to models performing poorly
- **No feature drift** monitoring
- **No idea about Feature impact** on a model



### High latency, unreliable Feature serving in production models at scale

- **Poor Model latencies** leading to bad user experience.
- **No dedicated dynamic resource allocation** for feature engineering
- **Multiple RPC calls at high throughputs to fetch features** dramatically increasing latencies

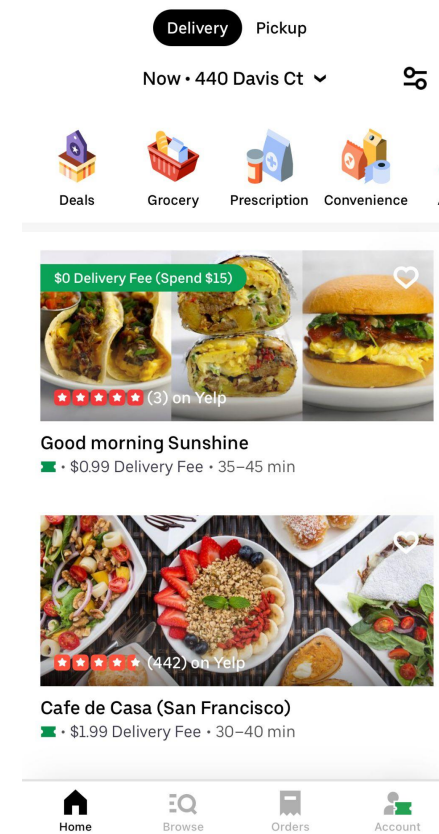
## Feature Stores



## Use case - ETA of an Uber EATS Order

### Key ML Features

- How large is the order? (*order\_size*)
- How busy is the restaurant? (*n\_meal*)
- How quick is the restaurant? (*meal\_preptime*)
- How busy is the traffic? (*n\_busy*)

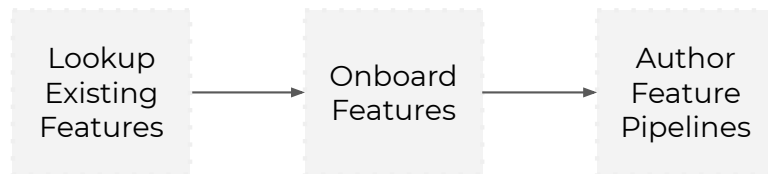


# Palette Feature Store: Workflow

Lookup  
Existing  
Features

- Search for features
  - by `feature_name`
  - by entity (e.g. `eater_features`)
  - by type (e.g. `categorical_features`)
  - by models (e.g. features used in `eta_prediction_model`)  
or any combination ...

## Palette Feature Store: Workflow



- Onboard Features and Author Pipelines
  - Metadata driven onboarding process
  - Feature Pipelines automatically created
  - Immediately available for consumption during Training & Inferencing

# Palette Feature Store: Workflow



```
dsl1 = DSLEstimator(lambdas={
  "region_id": regionId("@palette:restaurant:realtime_feature:lat:r_id", "@palette:restaurant:realtime_feature:lat:r_id")
})

dsl2 = DSLEstimator(lambdas={
  "prep_time": nFill(nVal("@palette:restaurant:batch_features:prep_time:r_id")),
  "n_meal": nFill(nVal("@palette:restaurant:realtime_features:n_nean:r_id")),
  "n_order": nFill(nVal("@basis:n_order")),
  "n_busy": nFill(nVal("@palette:restaurant:service_feature:n_busy:region_id"))
})

ml_pipeline = MLPipeline(dsl1, dsl2)
model = ml_pipeline.fit(basis_dataframe)
michelangelo_api.deploy(model)
```



## Palette Feature Store: Workflow

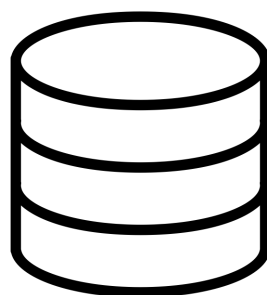


- Monitor Feature Metrics
  - Training-Serving **Skew**
  - Feature **Drift**
  - Feature **Importance**<>**Drift** correlation
  - Feature Quality (**Freshness, Consistency, Null Rate**)

## Feature Store (Palette) Lifecycle

### Feature Preparation

Batch &  
Streaming ETLs



### Feature Storage

Historical & Near Real-Time

Curated & Crowd-sourced

Metadata

Scalable offline access

Scalable online access

Online/Offline data parity

### Feature Discovery

Sharing across Models  
Automatic feature selection



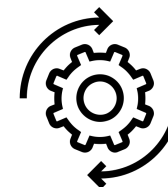
### Feature Monitoring

Data Quality reporting

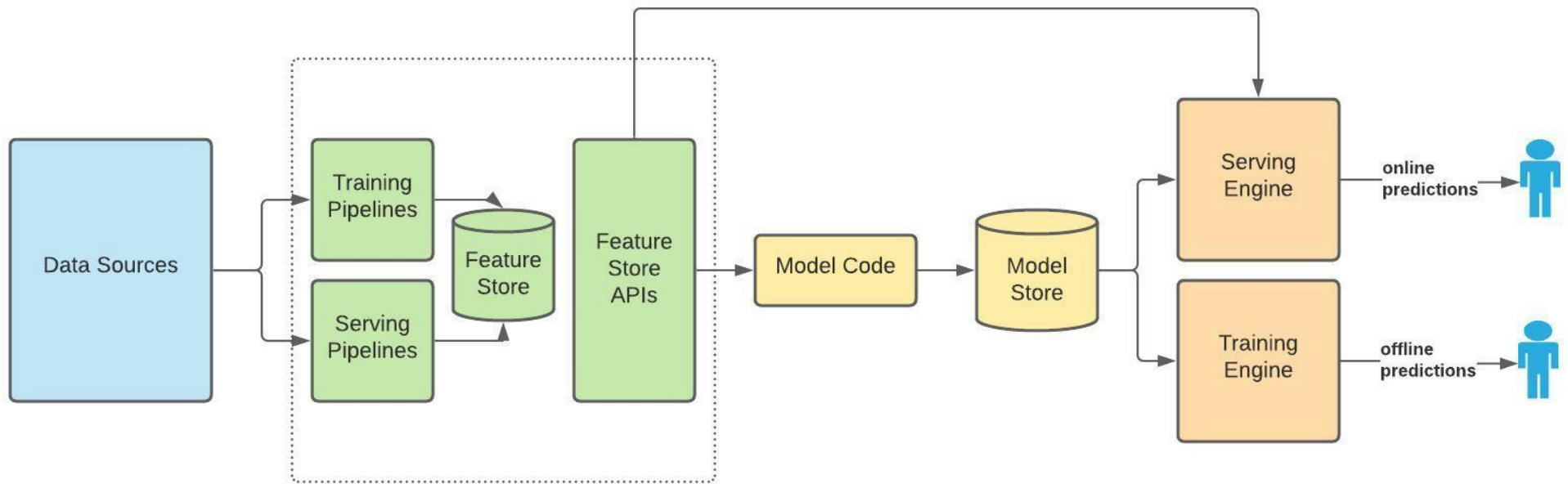


### Feature Transforms

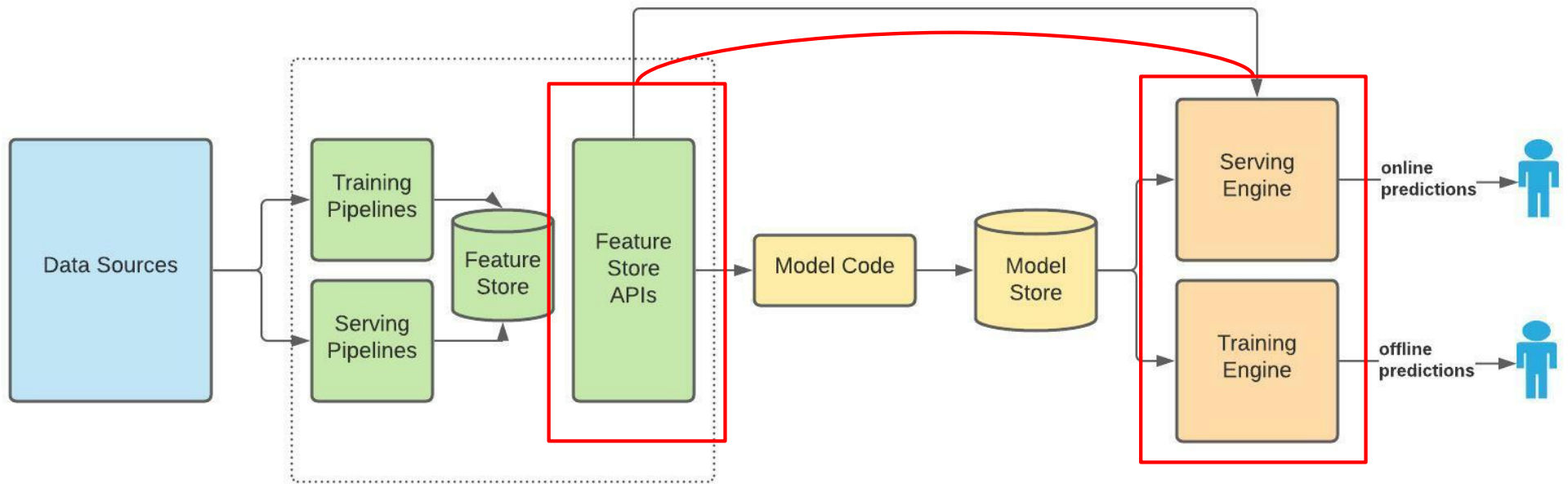
Model specific transforms



# Feature Stores in an End to End ML Platform



## Feature Stores in an End to End ML Platform




# Palette Feature Store Organization

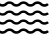
Organized as **entity : feature\_group : feature : join\_key**


*e.g. restaurant : order\_history : meal\_preptime : restaurant\_uuid*

- Feature Store Abstractions:
  - **Entity**: A Top Level Business Unit (e.g. **eater, courier, restaurant**)
  - **Feature Group**: Group of Features commonly used together (e.g. **order\_history**)
  - **Feature**: The Feature (e.g. **meal\_preptime, n\_meal, sum\_orders\_1week**)
  - **Entity Key**: The UUIDs of the entities (e.g. **eater\_uuid, restaurant\_uuid**)
- **Bring** your join keys or UUIDs
- **Join** together cross-entity Feature sets with minimal code
- **Train** on historical Feature values
- **Serve** the latest, most accurate values of Features at Low Latency
- Backed by a dual datastore system (training & serving)
- Get Training Serving parity out of the box

# Feature Types in Palette (Michelangelo)

-  **Batch Features:**
  - Features calibrated on **historical data**
  - Generated via offline **batch jobs**
  - Auto dispersed** for model inferencing
  - E.g. meal\_preptime (average prep time of historical orders)

---
-  **Near Real Time Features:**
  - Features calibrated on **streaming data**
  - Generated via **near real-time streaming jobs** (Flink, AthenaX)
  - Auto dispersed for model training
  - E.g. n\_meal (how busy is the restaurant)

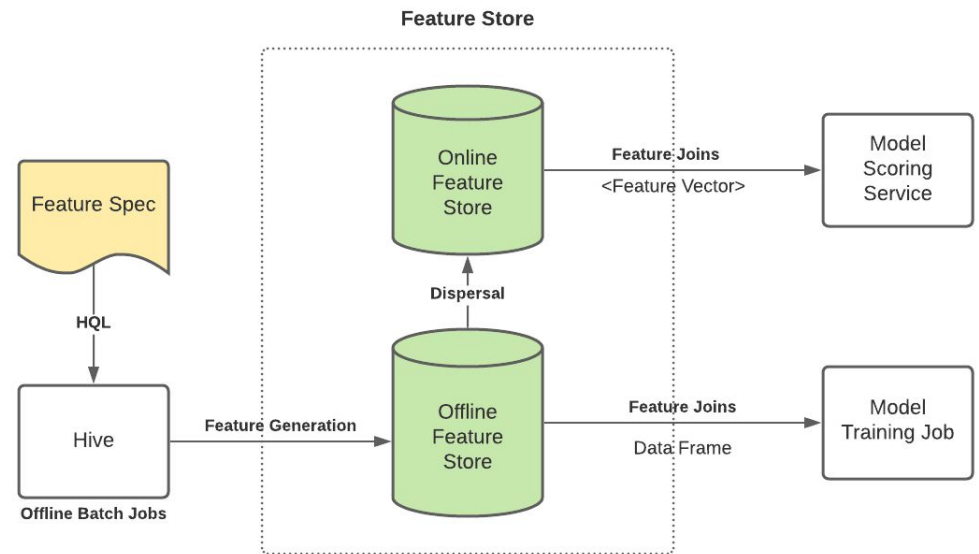
---
-  **RPC Features:** Features retrieved via 3P APIs
  - Features calibrated on 3P API calls
  - Calculated at run time and served to models directly
  - Auto dispersed for model training
  - E.g. location\_geohash (current geohash location of the courier)

---



# Computing Batch Features

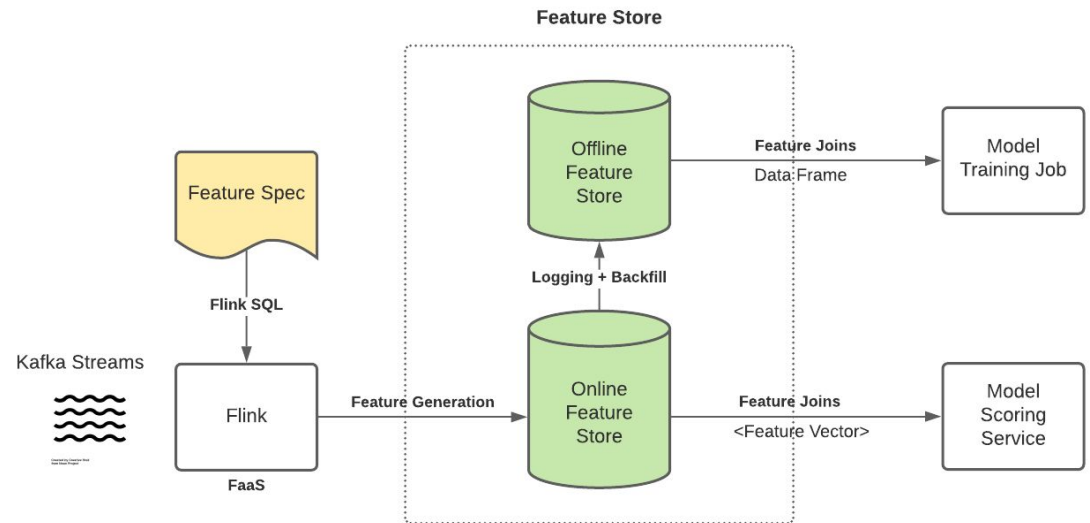
- Computed using **Historical Data**
- Not time sensitive
- Ingested from Hive Queries or Spark Jobs
- Aggregates over days/weeks
- E.g. *meal\_preptime*





# Computing Near Real Time Features

- Signals generated **seconds ago**
- Write Flink SQL to perform real time aggregations
- Materialize to the online store
- Auto ETL and Backfill to the offline store
- E.g. n\_busy (How busy is the restaurant)
  - **Kafka** event streams
  - Perform **Real-Time aggregations**

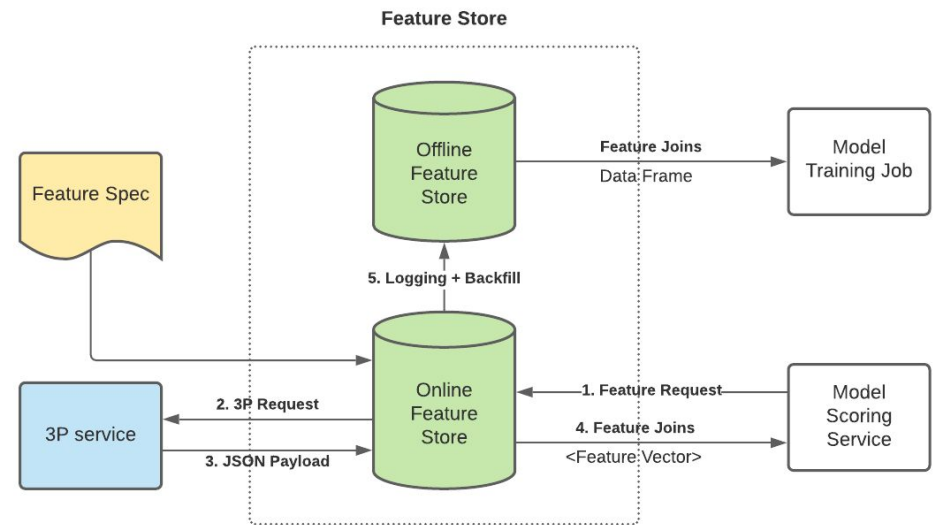






# Computing RPC Features

- Signals generated **in real-time**
- Make RPC calls to Fetch Features behind the scenes
- Auto ETL and Backfill to the offline store
- E.g. *lat/long*:
  - Fetched via HTTP calls





# Feature Extraction & Transformation

- Michelangelo Transformer
  - `transform()` and `scoreInstance()`
  - ML Readable / Writable
  - Extension of the Spark Transformer Framework
  - Parity across Spark and Spark-less environments
    - UDFs / DSLs
    - In-house unit testing framework for parity
- Feature Store APIs as Transformers
  - Feature Engineering as an integral part of the ML Pipeline

## Michelangelo Feature Store APIs as Spark Transformers

```
palette_tx1 = PaletteTransformer( {  
  "nMeal": "@palette:restaurant:realtime_feature:nMeal:r_id",  
  "prep_time": "@palette:restaurant:batch_feature:prep_time:r_id",  
  "lat": "@palette:restaurant:realtime_feature:lat:r_id",  
  "long": "@palette:restaurant:realtime_feature:long:r_id",  
})
```

- Instantiate Palette Transformer with Feature expressions
- Create a pipeline with one or more stages of estimators and transformers
- model = **pipeline.fit()**
- Evaluate your model via **transform()**
- Score your model via **scoreInstance()**

## ↔ DSLs: Feature Manipulation / Imputation

- Write expressions to define Transformations
- Pre-compiled Scala code execution at runtime
- Example Michelangelo code:

```
dsl_est1 = DSLEstimator(lambdas={
  "region_id": regionId("restaurant:fg:lat:r_id", "restaurant:fg:long:r_id")
})

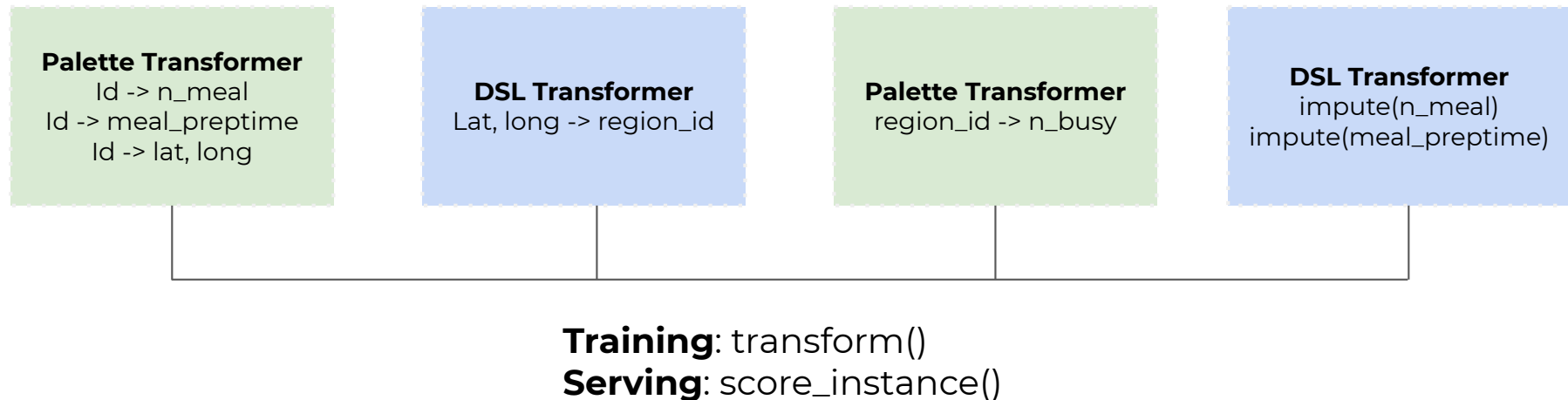
dsl_est2 = DSLEstimator(lambdas={
  "prep_time": nFill(nVal("restaurant:batch_fg:prep_time:r_id"), -1),
  "n_meal":     nFill(nVal("restaurant:realtime_fg:n_meal:r_id"), -1),
  "order_size": nVal("basis:order_size"),
  "n_busy":     nFill(nVal("restaurant:service_fg:n_busy:region_id"), -1)
})
```

**\*fg: feature\_group**

## Uber EATS Transformation Example

### Computation Order

- `n_meal:restaurant_id -> n_meal`
- `meal_preptime:restaurant_id -> meal_preptime -> DSL`
- `busy_scale: restaurant_id -> lat, long -> regionId(lat, long) -> busy_scale`



# Feature Store Results & Takeaways

- **Democratized Usage:** 20K+ Features used across 8K+ production models
- **Model development times** reduced from days to hours
- **Multi Modality Support:** Batch, Realtime and RPC Features with online and offline parity
- **Offline scalability:** Joins across billions of rows
- **Online serving latency:** Parallel IO, fast storage with caching
- **Feature Transformers:** Setup chains of transformations at training/serving time

# Embedding Ecosystems



# Managing ML Pipelines: Feature Stores and the Coming Wave of Embedding Ecosystems

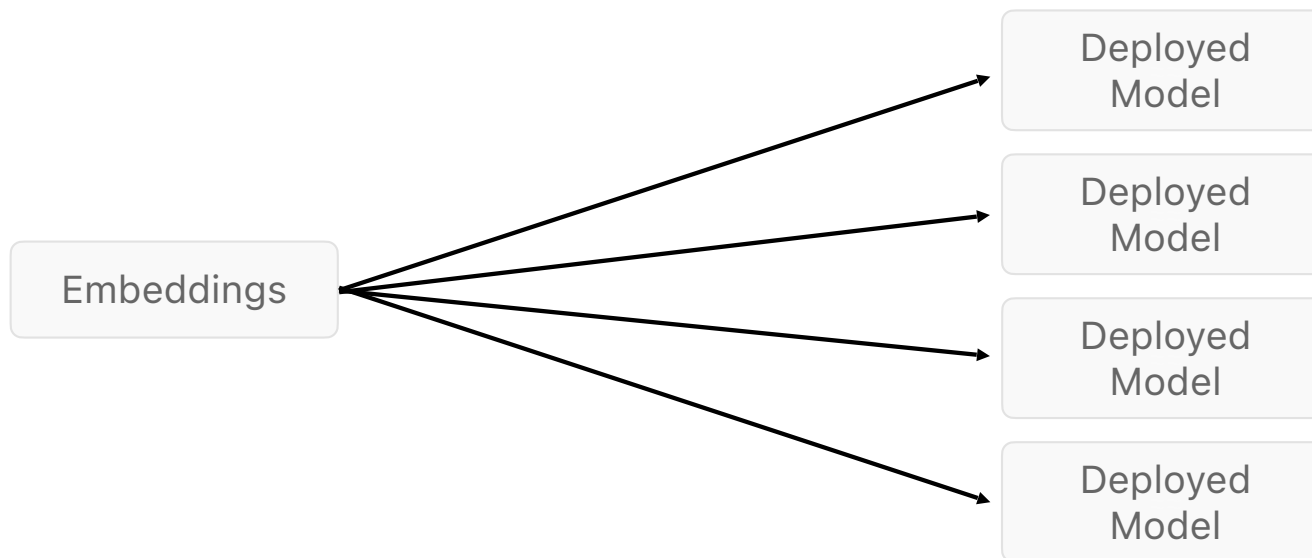
Xiao Ling | VLDB Tutorial 2021





# Recap: Self-Supervision Embeddings

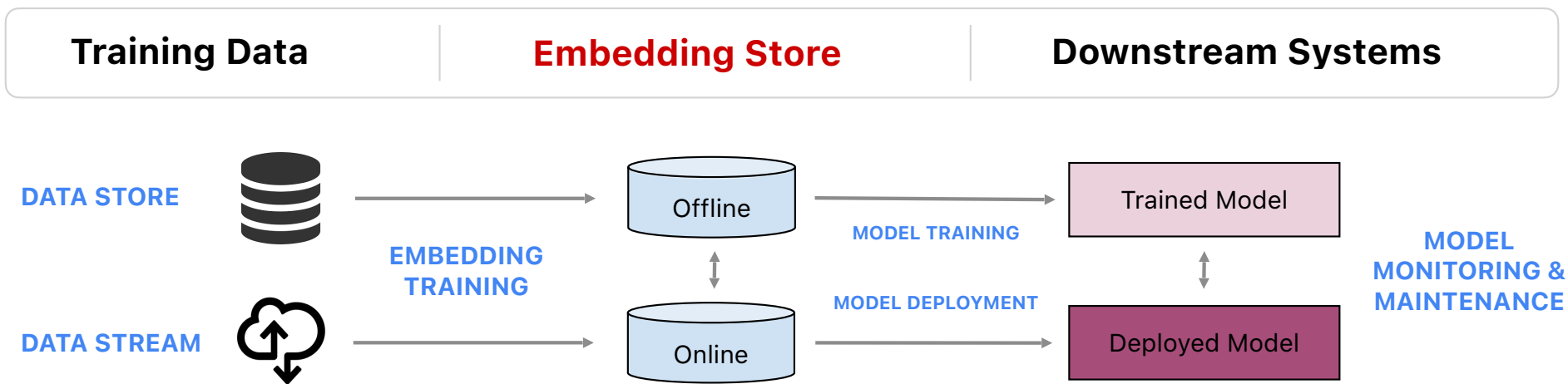
Used in many different downstream systems



Downstream systems require **less supervised data** and provide a **quality lift** compared to hand-tuned predecessors.

# Recap: Embedding Ecosystems

New age of feature store systems manage pretrained embeddings downstream systems use them as inputs.



# Grounding Use Case: Named Entity Disambiguation

Map “strings to things” in a knowledge base.

Key part of assistant, search, and information extraction

How tall is Lincoln?



# The Long Tail of Entities

HEAD



Washington, DC

Q61

>1k examples in training

TORSO



Chevrolet, Corvette

Q56166

(10, 1000]

TAIL

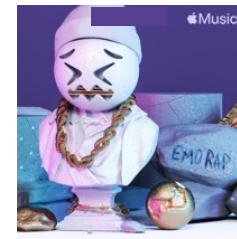


Reddish Potato Beetle

Q14934552

[1, 10]

UNSEEN



Sauce! by XXXTENTACION

Q???

0 example

Popular

IR and BERT/IR work great for the HEAD.

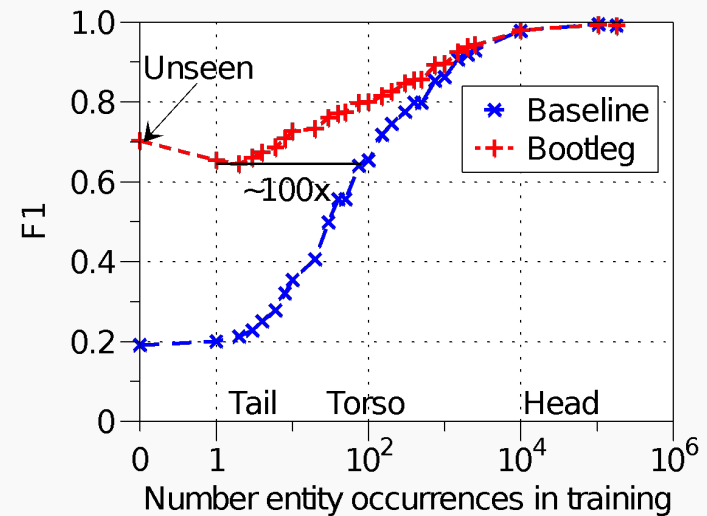
Rare

The **majority** of entities are **rare**!  
13% entities have Wikipedia page.  
< 1% of songs in Wikidata!

# #1 Tail Scalability Challenge

Large number of patterns needed to resolve the tail, making it difficult to scale a system that can learn the patterns.

90 million entities in Wikidata ->  
90\*100 million examples for 60 F1



Subtle reasoning clues are needed for the tail!  
(+40 F1 points by encoding these reasoning patterns)

# Entity Embeddings in Downstream Applications

## Experiment on the entity linking task in an existing Q&A system

- With and without Bootleg-learned entity embeddings
- The entity embeddings significantly improve F1 by a relative 8%

Also, a relative 8% improvement on tail entities!

## #2 Memory Usage

**Embeddings linearly grow per number of entities**

- 128d float32 x 5M  $\approx$  2.4 Gb (English Wikipedia)
- 128d float32 x 96M  $\approx$  46.08 Gb (Wikidata)

**It requires larger and larger servers over time** 

- More computation affects service latency 

**Hard to fit *on device!*** 

# Memory Usage

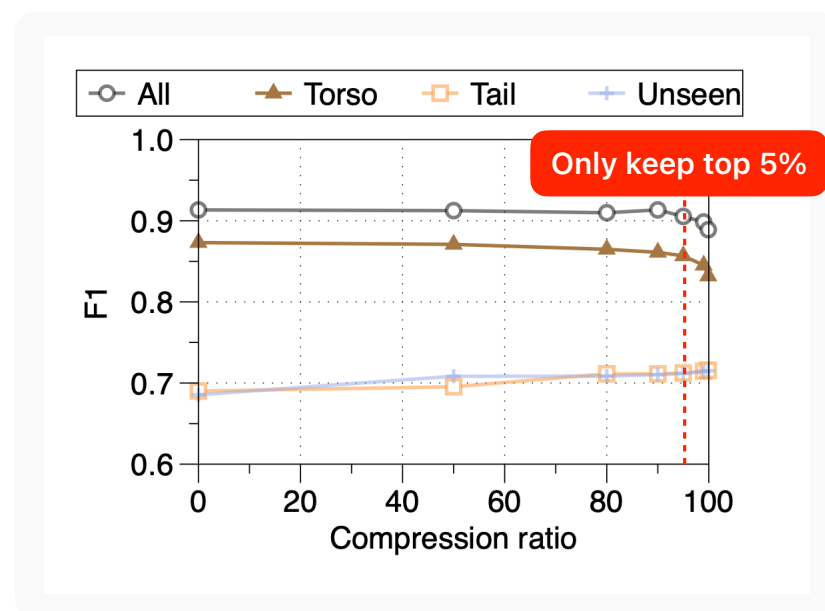
Memory can be saved w/o a big quality sacrifice

**Only keep the top k entity embeddings (i.e., compression ratio 100 - k)**

- Uses a random UNK entity embedding
- Less memory-heavy signals remain

**F1 only drops by 0.8 overall**

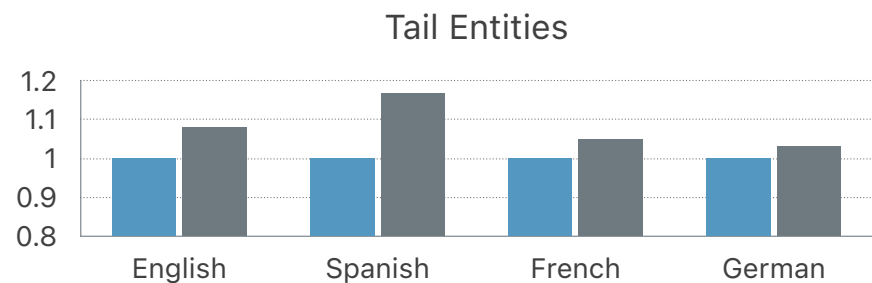
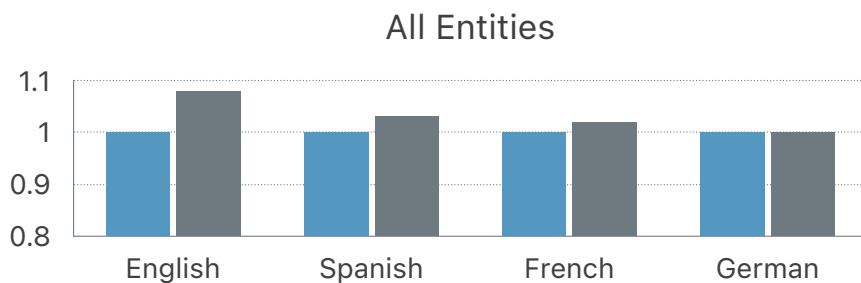
- Memory drops from 5.2GB to 0.3GB!





# #3 Embeddings in i18n languages

## Embeddings work on other languages



## Challenges

- Lack of equally abundant resources in English
- Memory usage increases the size of embeddings by the num of languages

# Multilingual Entity Embeddings

Botha et al. 2020 proposed to train multilingual entity embeddings

- Memory usage doesn't grow with the number of languages
- Entity embeddings trained from resources across languages
- Enabled by a multilingual language model



# #4 Embedding Stability

## Entity embeddings are self-supervised from Wikipedia

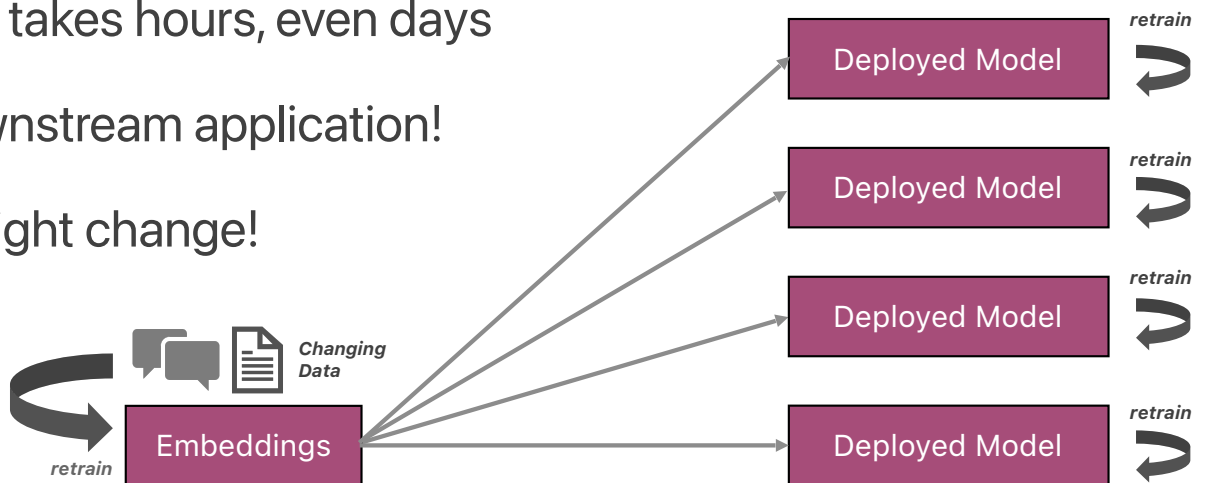
- 20k new articles / month

## Updating the model is hard

- Retraining entity embeddings takes hours, even days

Also, need to retrain *each* downstream application!

- Previous correct prediction might change!



# #5 Model Evaluation and Monitoring

## Are the embeddings

- sensitive to questions?
- vulnerable to attacks?
- biased to entities popular in one country?, Etc...

## Is the downstream application affected by updated embeddings?

- What about 10s or 100s downstream apps?
- How to enable safe regular model updates?



# Summary of Challenges

**#1 Long-tail of entities**

#2 Memory usage

#3 Multi-lingual embeddings

**#4 Embedding stability**

**#5 Model monitoring**



TM and © 2021 Apple Inc. All rights reserved.

# Self-Supervised Training Data: The Challenge of the Long Tail

# Grounding Use Case: Named Entity Disambiguation

Map “strings to things” in a knowledge base.

How tall is *Lincoln*?



Q91

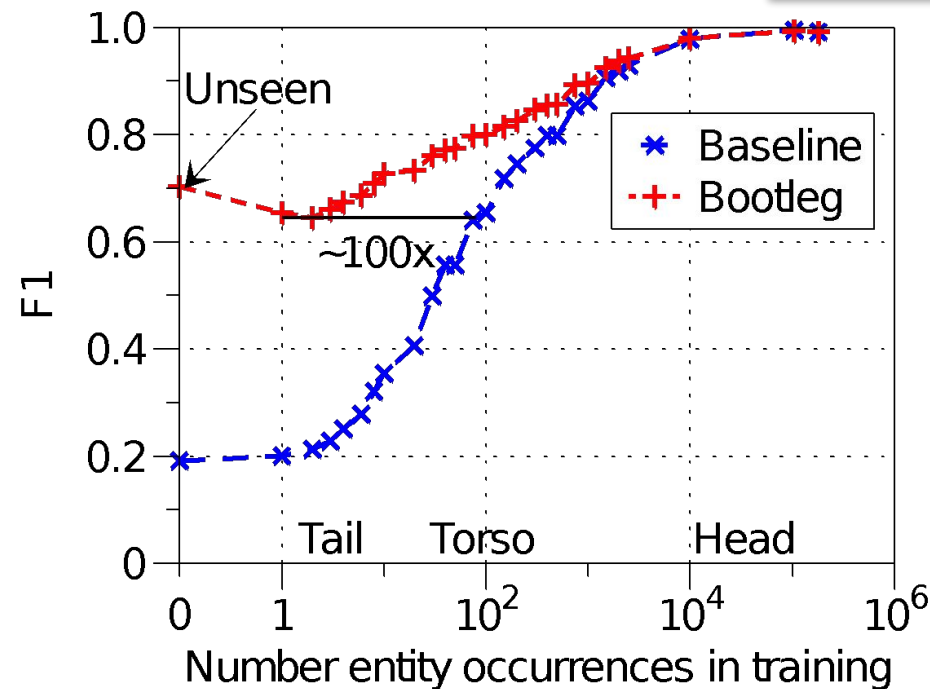


Key part of assistant, search, and information extraction



# Tail Challenge

*Impossible to scale the data to memorize all patterns needed for rare entities*



Subtle reasoning clues are needed for the tail!  
(+40 F1 points by encoding these reasoning patterns)

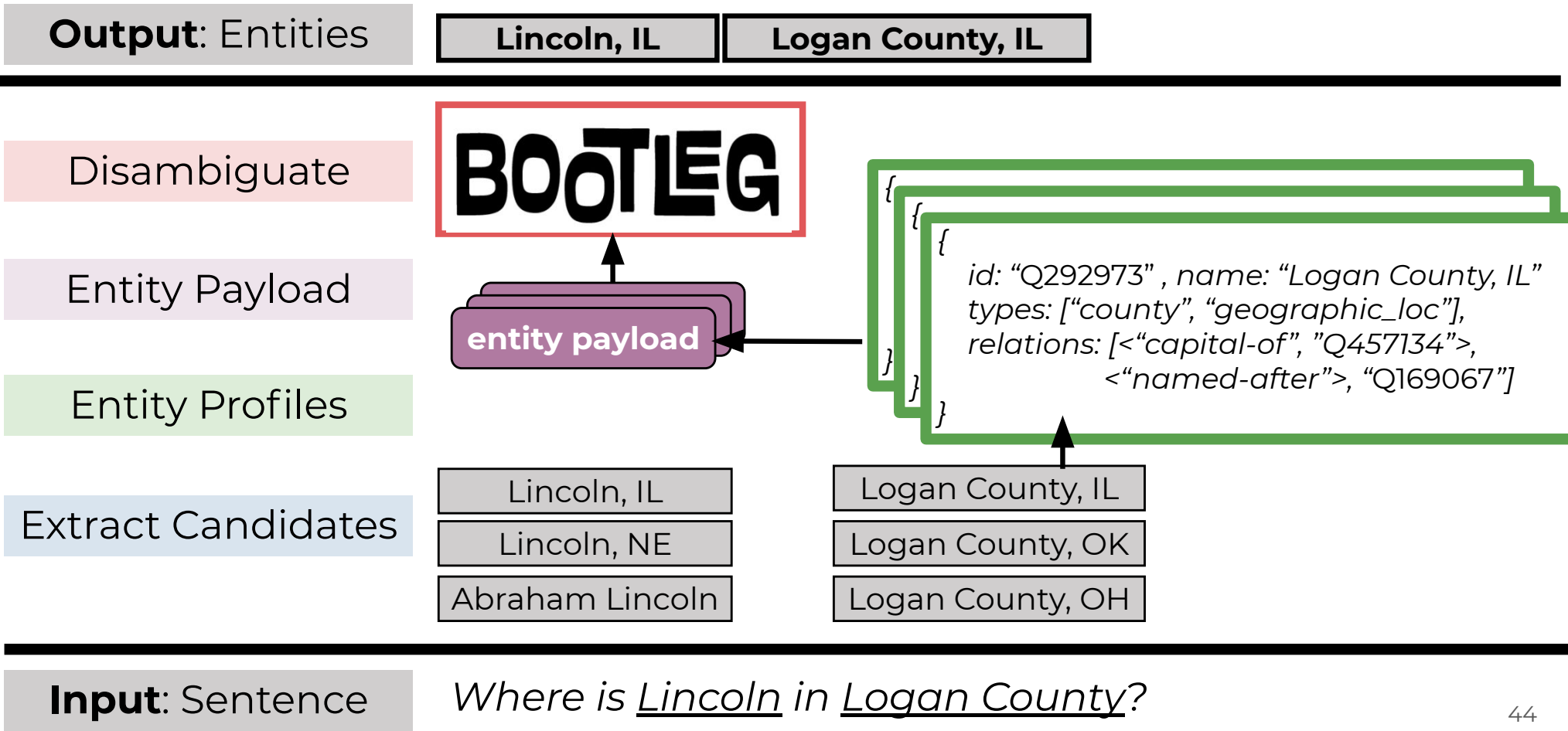
# Bootleg: Tackles the Tail with Structural Knowledge



## BOOTLEG

**Key Idea:** reasoning over *type* and *relationship* signals can resolve unseen entities.

# Disambiguation Inputs and Outputs

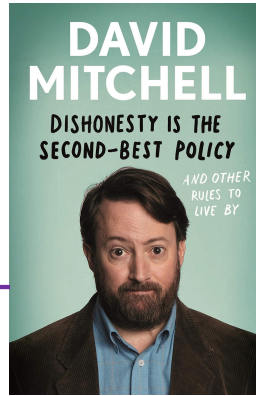


# Reasoning over Relationships



Victoria Mitchell  
(poker player, writer)

spouses

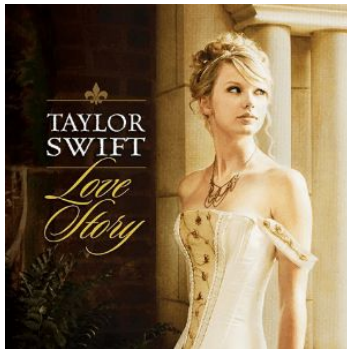


David Mitchell



Victoria Mitchell  
(runner)

David and Victoria Mitchell added spice to their marriage



Love Story by Taylor Swift



Love Story by Andy Williams

Play Love Story by Williams

# Reasoning over Types

How tall is Lincoln?



*People have heights,  
not places or brands*

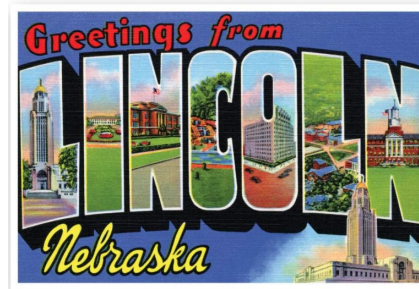
What is the  
cheapest Lincoln?



L I N C O L N

*Brands have prices,  
not places or people*

How many people  
are in Lincoln?



*Places have populations,  
not people or brands*

# Bootleg: Tackles the Tail with Structural Knowledge

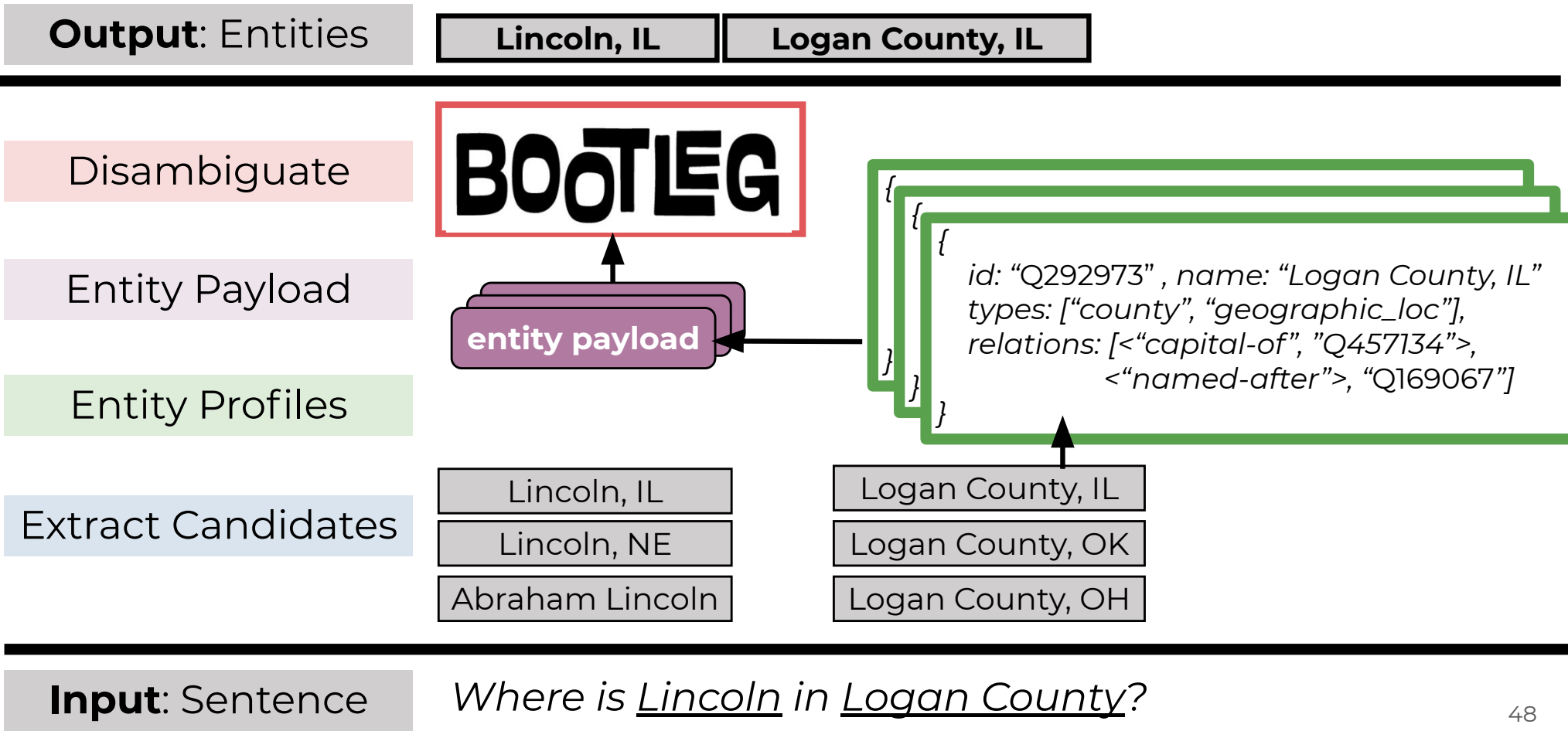


## BOOTLEG

**Key Idea:** reasoning over *type* and *relationship* signals can resolve unseen entities.

**Implementation:** use *embeddings* to teach a model to reason over types and relationships.

# Disambiguation Inputs and Outputs



# Using Embeddings to Encode Signals

## Entity Embedding

<i>key</i>	<i>value</i>
...	
Q3452	
Q36897	
Q12	
Q292973	
Q903278	
Q328475	
...	

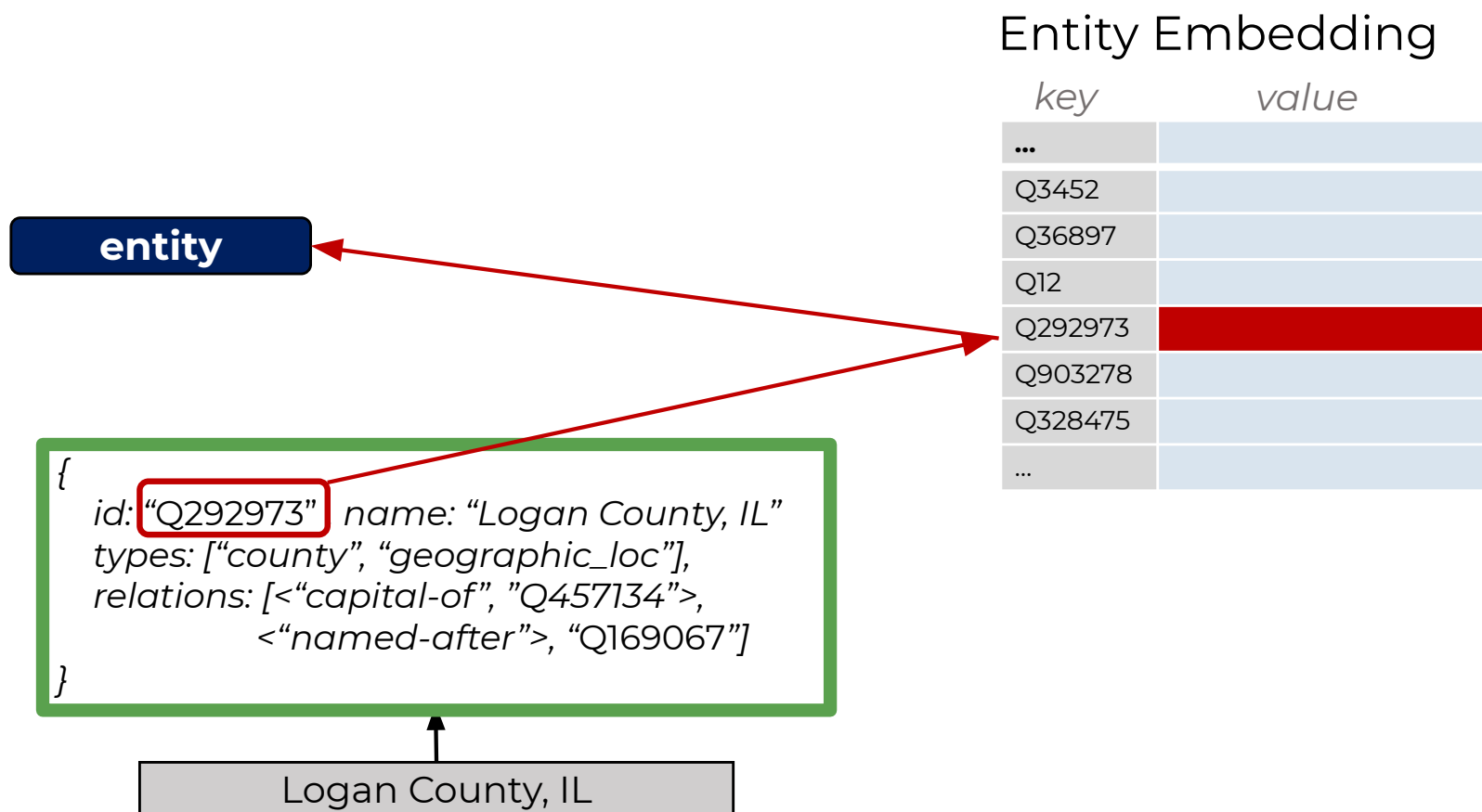
```
{  
  id: "Q292973" , name: "Logan County, IL"  
  types: ["county", "geographic_loc"],  
  relations: [<"capital-of", "Q457134">,  
             <"named-after">, "Q169067"]  
}
```

Logan County, IL

*For each candidate, we use the entity profile to extract (learned) embeddings.*

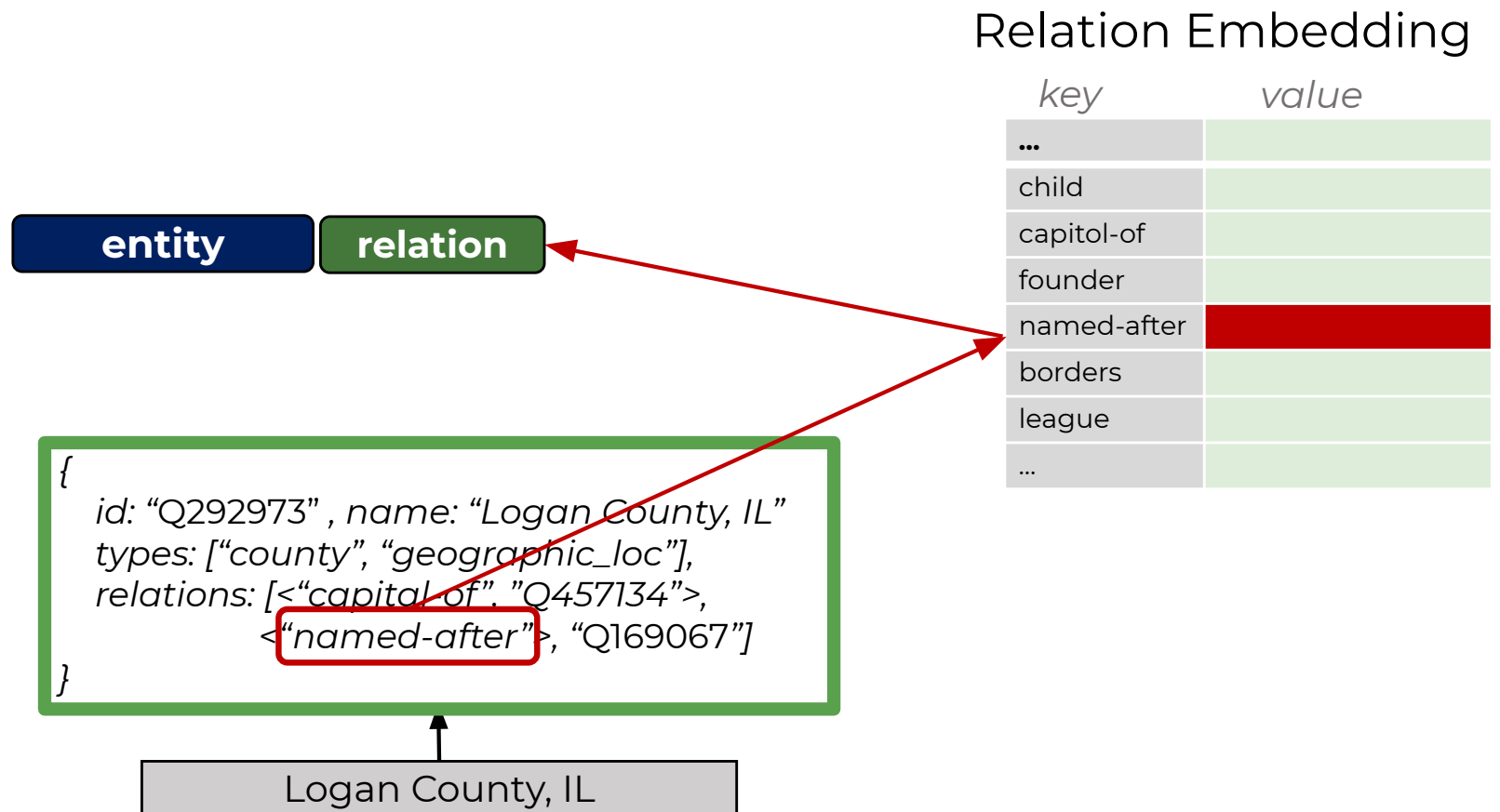


# Using Embeddings to Encode Signals

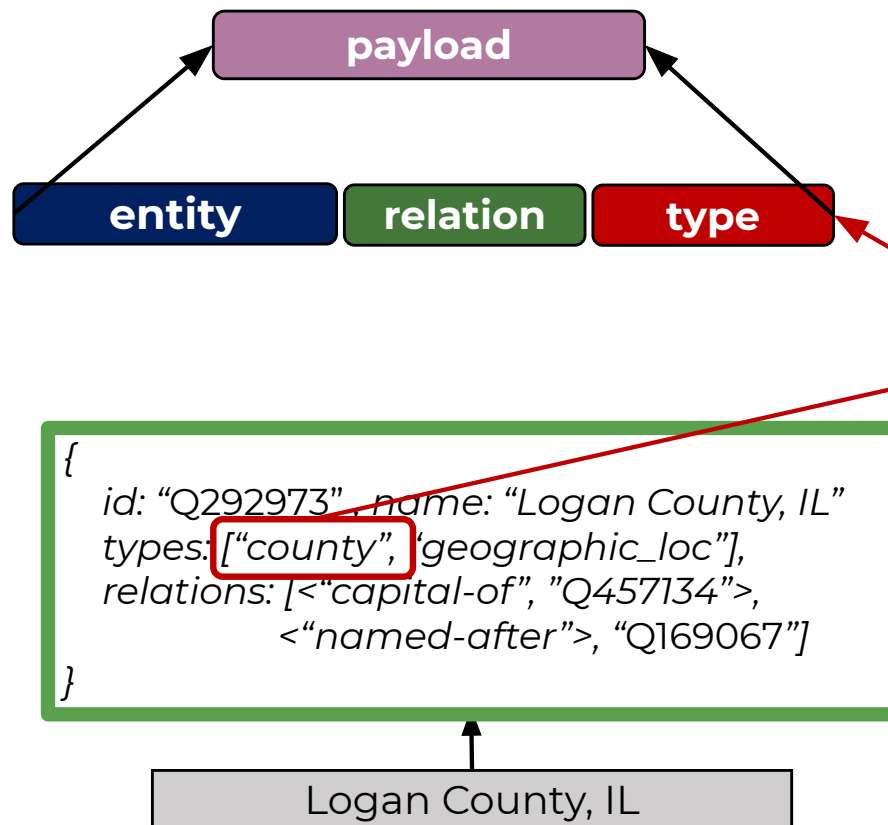


*For each candidate, we use the entity profile to extract (learned) embeddings.*

# Using Embeddings to Encode Signals



# Using Embeddings to Encode Signals

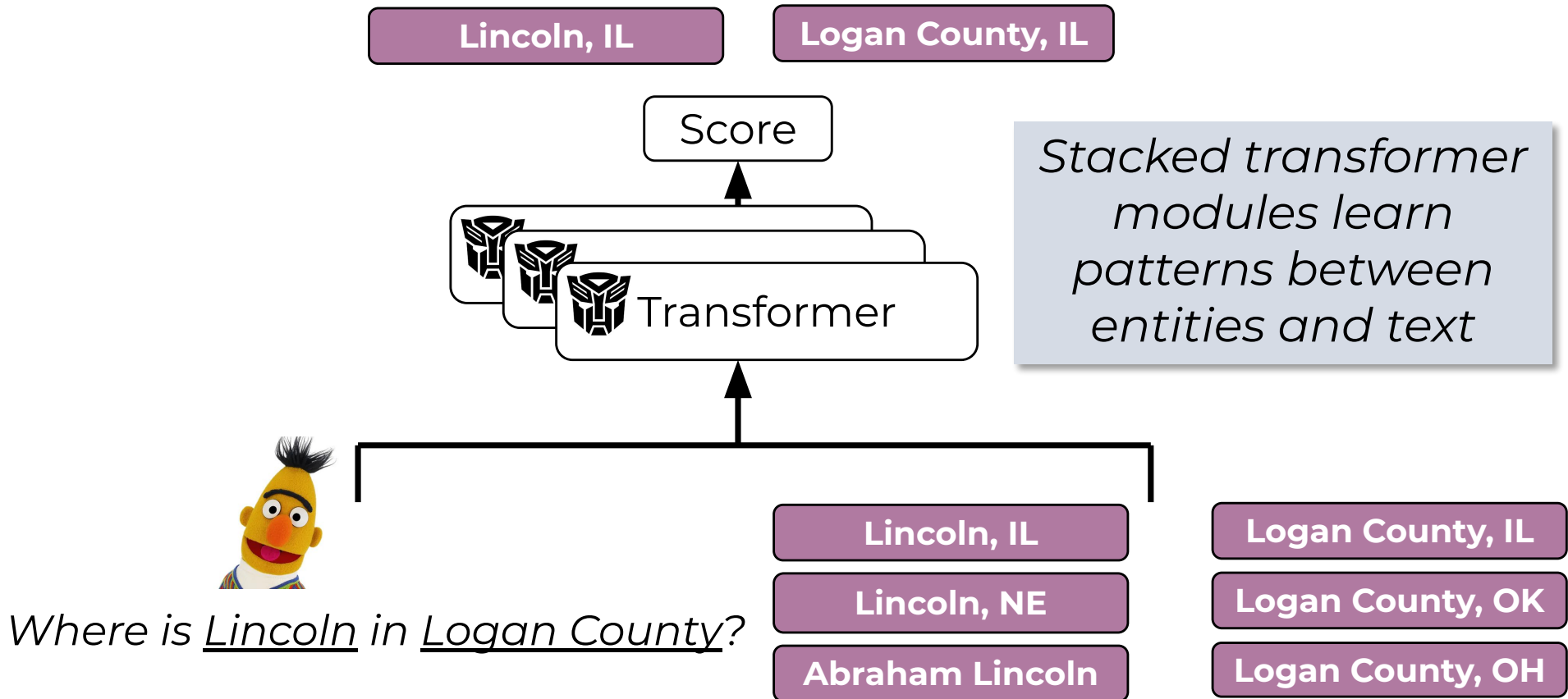


## Type Embedding

key	value
...	
child	
soccer team	
crime	
fruit	
county	
geo-loc	
...	

*The entity payload has embeddings mapping for each structural resource.*

# Bootleg Architecture



*Simplest architecture that supports reasoning over types and relations.*

## Bootleg: Tail Performance

*On the head, BERT-based baseline performs ~ 5 F1 points of Bootleg.*

*On the tail, Bootleg outperforms baseline by > 40 F1 points!*

Evaluation Set	BERT NED Baseline	Bootleg	# Examples
All	85.9	<b>91.3</b>	4,066K
Torso Entities	79.3	<b>87.3</b>	1,912K
Tail Entities	27.8	<b>69.0</b>	163K
<b>Unseen Entities</b>	18.5	<b>68.5</b>	10K

Performance results on Wikipedia dataset.

## Bootleg: Industrial Performance

Included Bootleg embeddings into an Overton production task answering millions of users' factoid queries. We report relative lift.

Evaluation Set	English	Spanish	French	German
All Entities	1.08	1.03	1.02	1.00
Tail Entities	1.08	1.17	1.05	1.03

# Using Bootleg Downstream: SoTA on the TACRED Benchmark

**Goal: extract the relationship between a subject and object pair.**

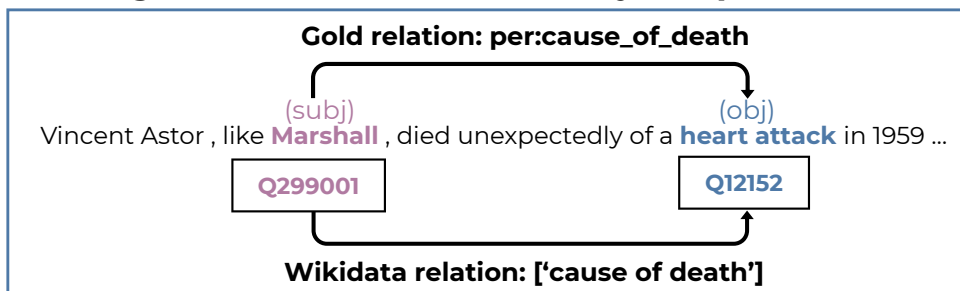


**Micro-Avg. F1 on TACRED Revised test dataset:**

Model	Test F1 Score
SpanBERT	78.0
KnowBERT	79.3
Bootleg+SpanBERT	80.2 (SoTA)

Zhang et al., 2017 and Hennig et al., 2020.

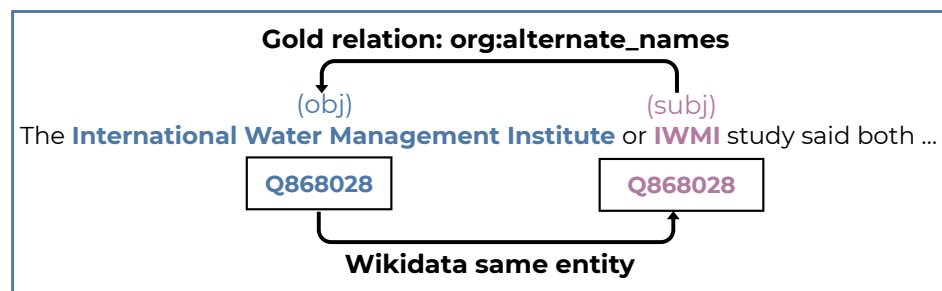
**Bootleg resolves errors made in by the prior SoTA model.**



**SpanBERT** no\_relation ❌

**BOOTLEG** per:cause\_of\_death ✅

*Leveraging type and relation information downstream*



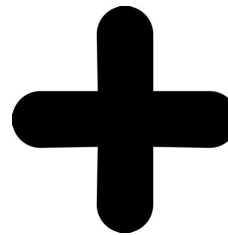
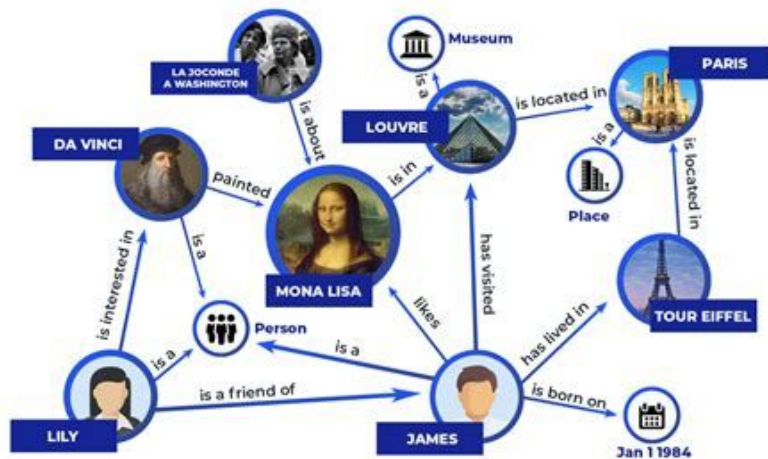
**SpanBERT** no\_relation ❌

**BOOTLEG** org:alternate\_names ✅

*Understand that sub-strings relate to the same entity*

# Self-Supervised Data Take Away

Self-supervised data does not well represent tail distributions -> embeddings may not be high quality for rare entities

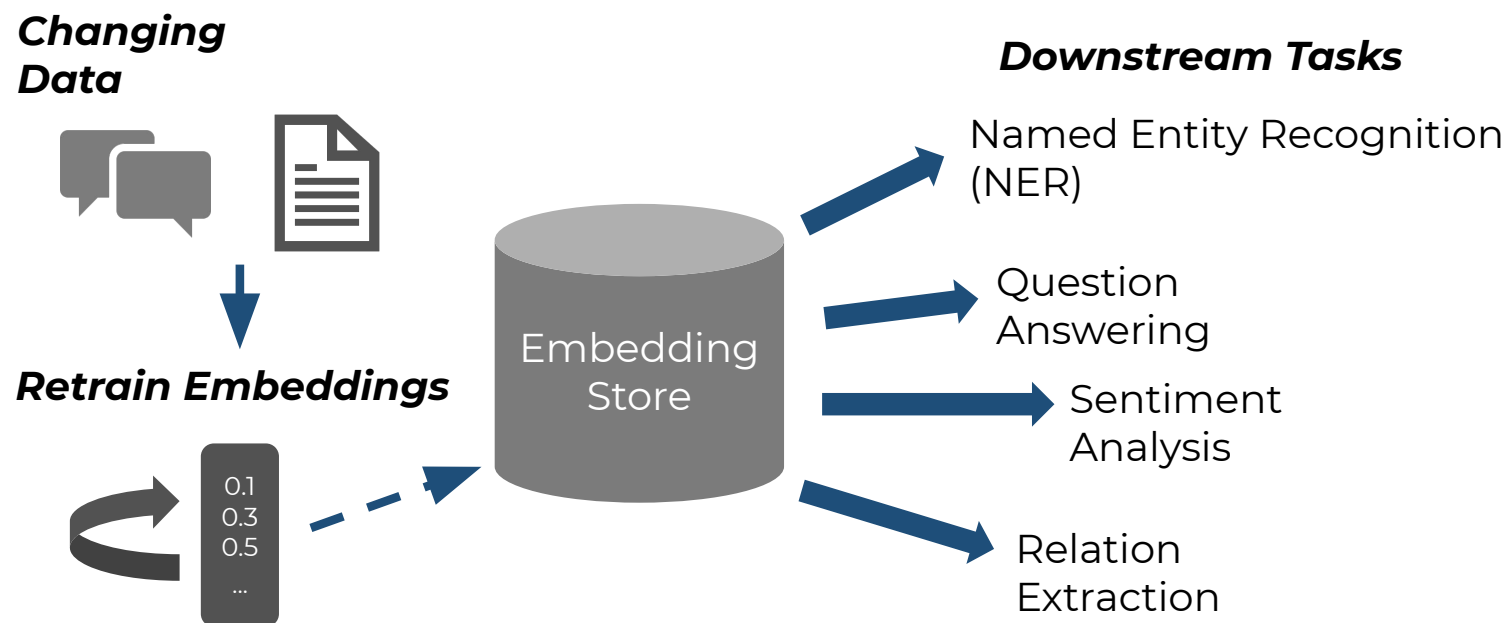


**Solution:** merged unstructured data with structured knowledge that can generalize to the tail.



# Embedding Management: Stability

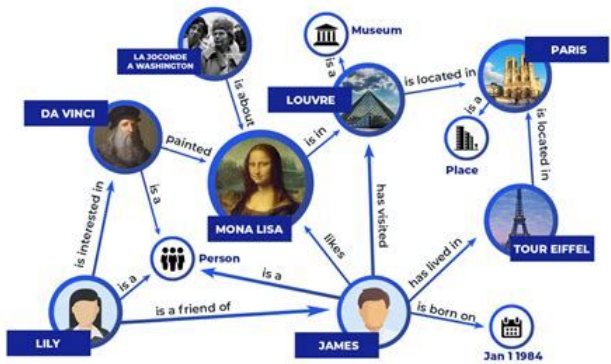
# Problem Setting: Embedding Store



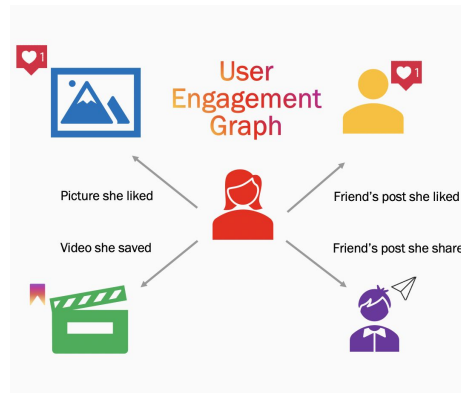
**New embeddings require downstream tasks to be retrained!**

# Why do embeddings need to be retrained?

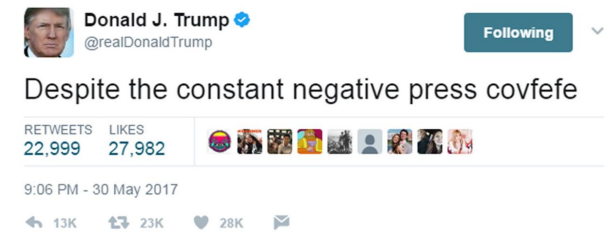
Learn new entities



Leverage new activity<sup>1</sup>



Understand new words



**Model freshness is necessary for user satisfaction in many products.**

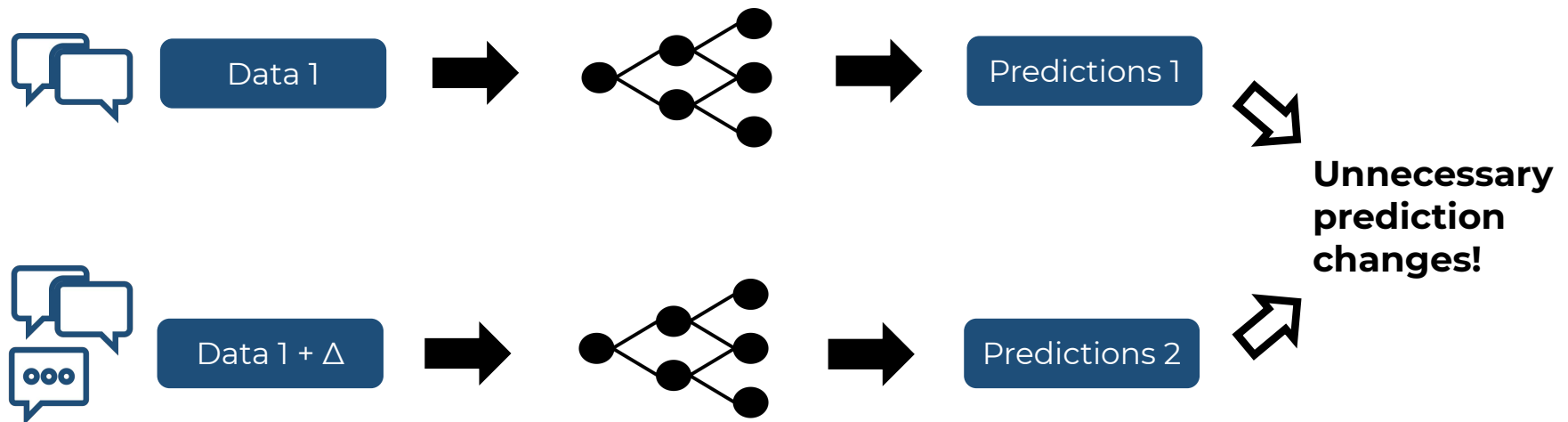
[1] <https://about.instagram.com/blog/engineering/designing-a-constrained-exploration-system>

Google retrains their app store Google Play models *every day*, and Facebook retrains search models *every hour*.

[1] Baylor et al. TFX: A TensorFlow-Based Production-Scale Machine Learning Platform. KDD, 2017.

[2] Hazelwood et al. Applied Machine Learning at Facebook: A Datacenter Infrastructure Perspective. HPCA, 2018.

But model training can be unstable...



***Prediction churn***

[1] Cormier et al. Launch and Iterate: Reducing Prediction Churn. NeurIPS, 2016.

# Challenges of Instability

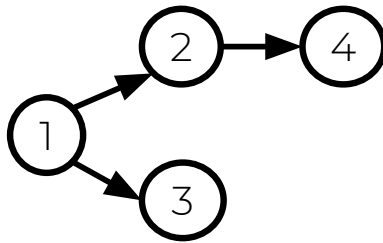
**Debugging**



**Consistent user-experience**



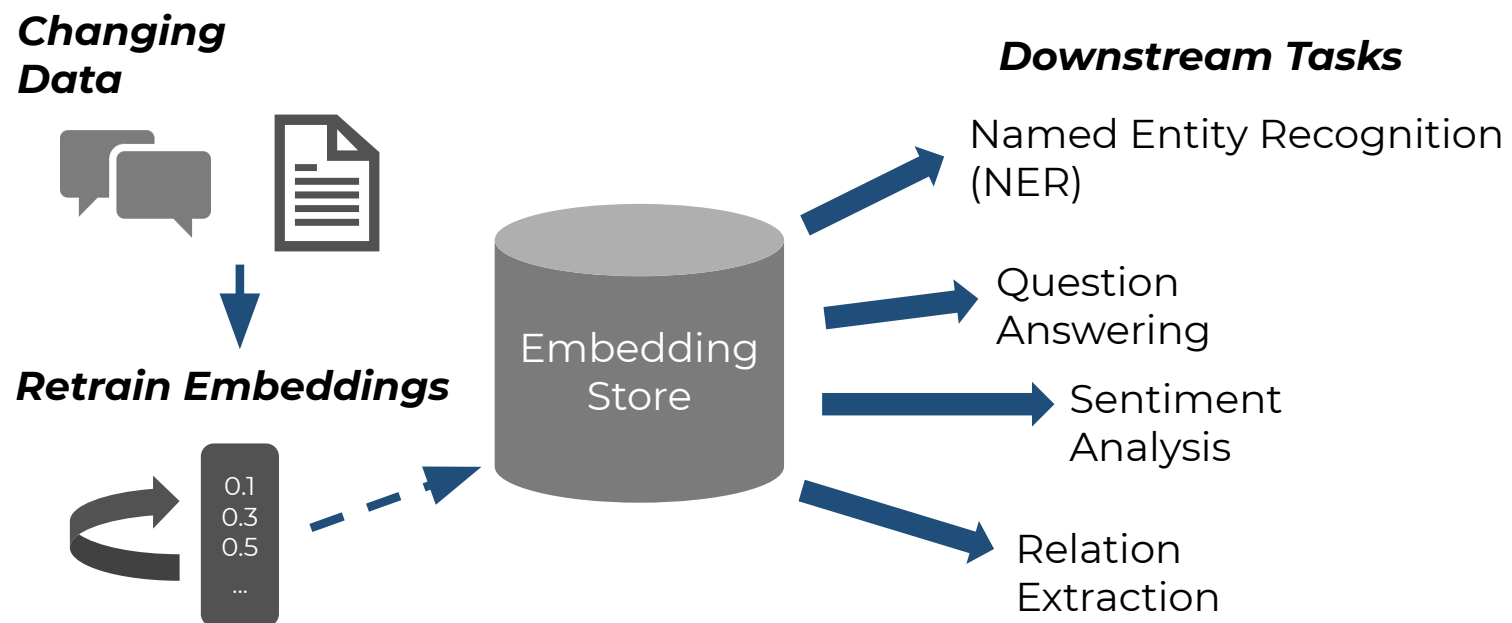
**Model dependencies**



**Research reliability**



# Problem Setting: Embedding Store



**How does the embedding instability propagate to downstream tasks?**

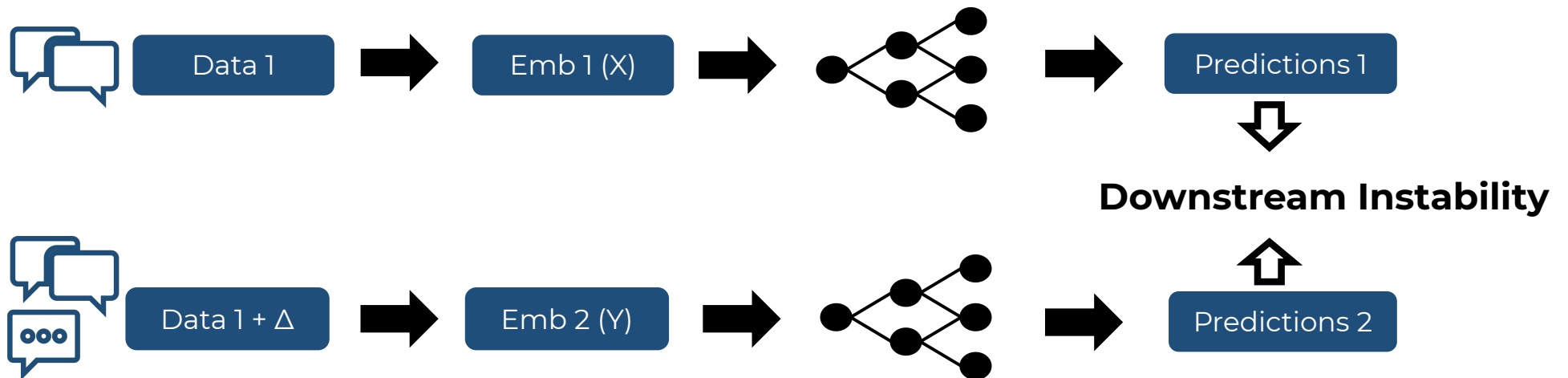
# Outline

- Downstream instability definition
- Stability-space tradeoff
- Measuring embedding quality with distance measures





## Definition: Downstream Instability




Downstream instability = % **prediction disagreement** between models trained on a pair of embeddings

# Embedding Hyperparameters that Impact Storage



<i>Uniform Quantization</i>	<u>32-bit</u>	<u>1-bit</u>
	<b>0.04</b>	<b>0.1</b>
Interval: [-0.1, 0.1]	<b>-0.03</b>	<b>-0.1</b>
	<b>-0.08</b>	<b>-0.1</b>

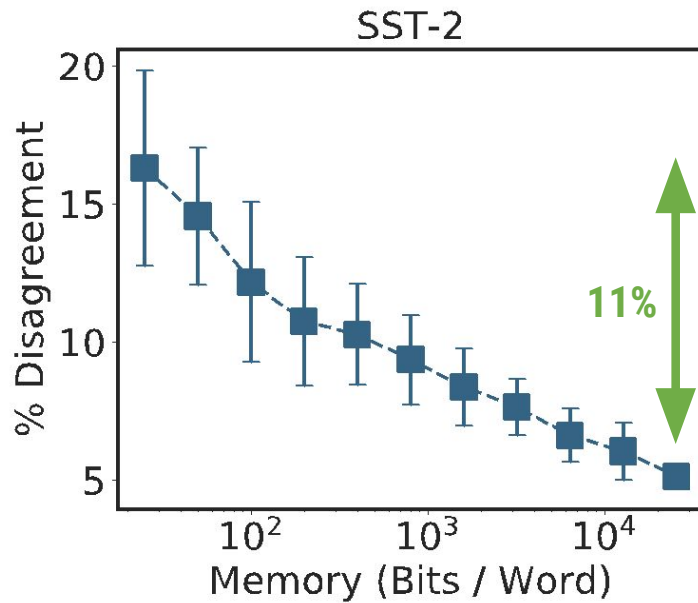


 **Downstream Instability**

[1] May et al. On the downstream performance of compressed word embeddings. NeurIPS, 2019.

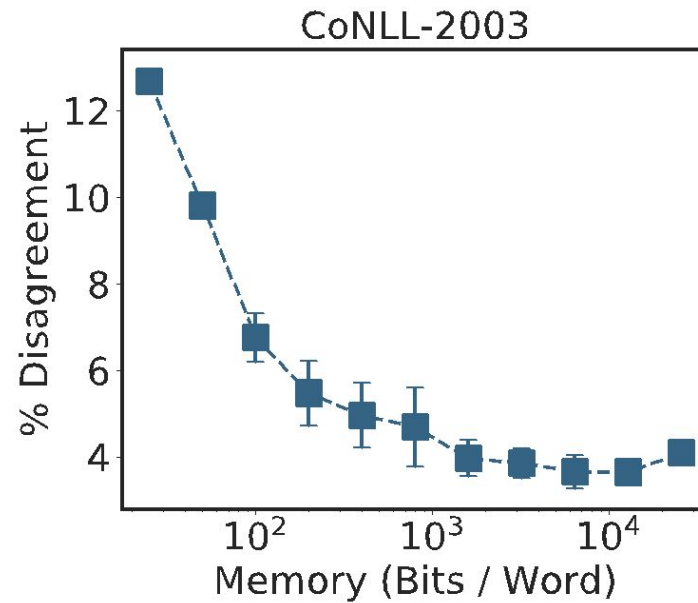
# Stability-Space Tradeoff for Word Embeddings

Sentiment Analysis



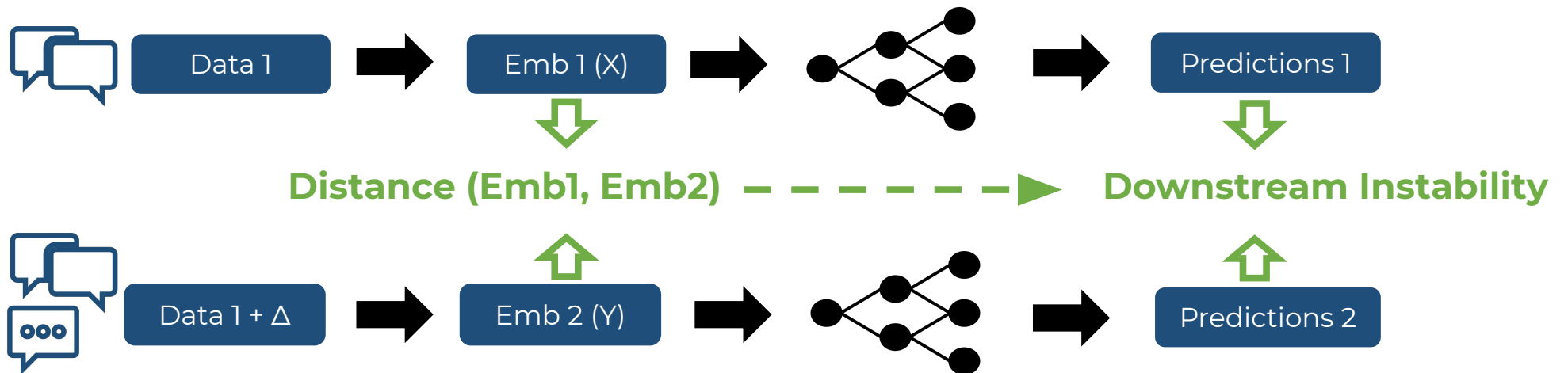
 **Embedding Size**

NER



 **Downstream Instability**

# Goal: Embedding Distance Measure for Instability



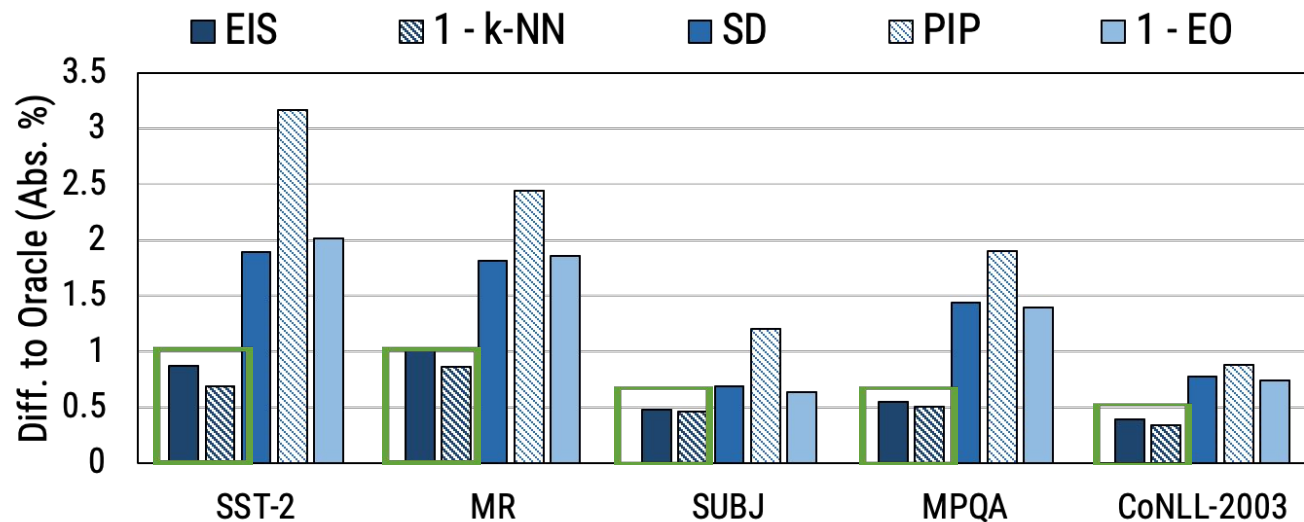
The measure should relate the **distance between the embeddings** to the **downstream instability**.

# Embedding Distance Measures

- k-NN measure [1,2,3]
- Semantic displacement (SD) [4]
- PIP loss [5]
- Eigenspace overlap (EO) [6]
- Eigenspace instability measure (EIS) [7]

[1] Hellrich & Hahn, COLING, 2016; [2] Antoniak & Mimno, TACL, 2018; [3] Wendlandt et al., NAACL-HLT, 2018; [4] Hamilton et al., ACL, 2016; [5] Yin & Shen, NeurIPS, 2018; [6] May et al., NeurIPS, 2019; [7] Leszczynski et al., MLSys 2020

# Using Embedding Distance Measures to Minimize Downstream Instability



**k-NN measure** and **theoretically grounded EIS measure** outperform other measures for selecting embeddings to **minimize downstream instability**.

# Stability Takeaways

- Defined downstream instability with respect to embeddings
- Stability-space tradeoff (precision, dimension)



**Embedding Size**



**Downstream Instability**

- Measuring embedding quality with embedding distance measures
  - EIS and k-NN measures select embeddings with lower downstream instability

# Closing the Loop of Model Development: Monitoring and Patching



# Monitoring and Patching

**Embeddings need to be updated:** distribution shift, changing needs



## **Monitor (when to update)**

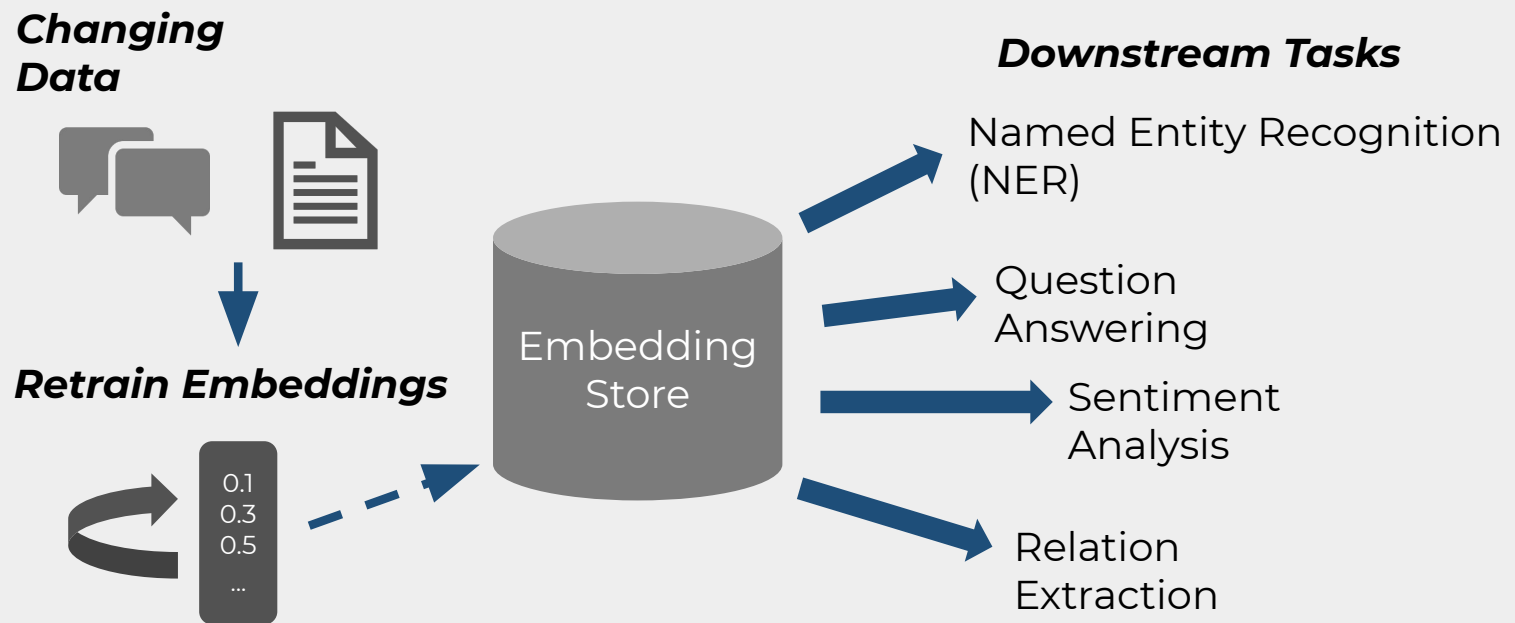
Evaluate and track distribution shift



## **Patch (how to update)**

Fix bugs and improve performance

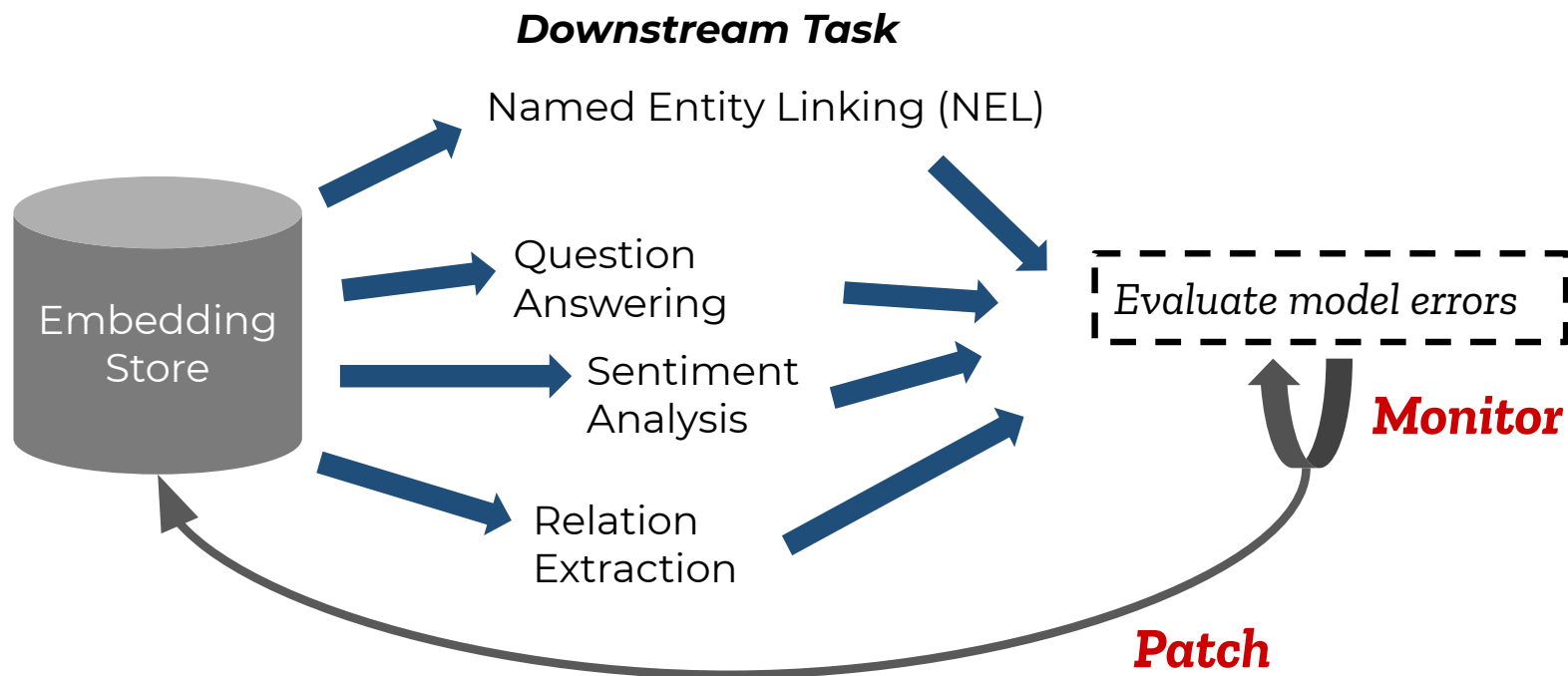
## Remember: Embedding Store



**Important:** update **embeddings not downstream models** → changes propagate down to models!

# Crucial Bottleneck: Evaluation

Can't monitor and patch embeddings without evaluation



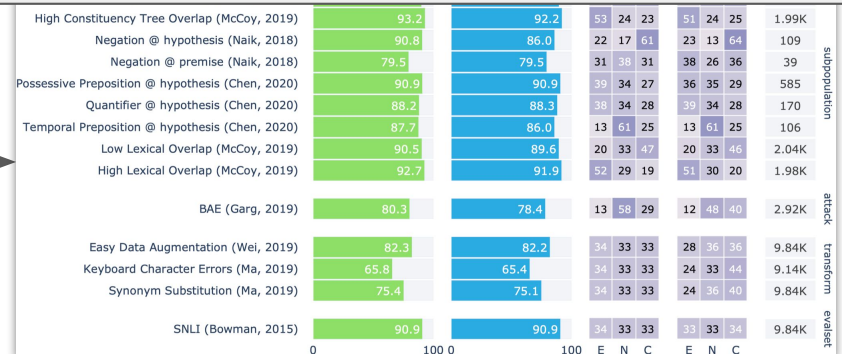
# Crucial Bottleneck: Evaluation

Can't monitor and patch embeddings without evaluation

## Many Evaluation Strategies

- Critical data slices
- Bias / fairness concerns
- Sensitivity to perturbations
- Invariance to transformations
- and more!

## Fine-Grained Evaluation Metrics

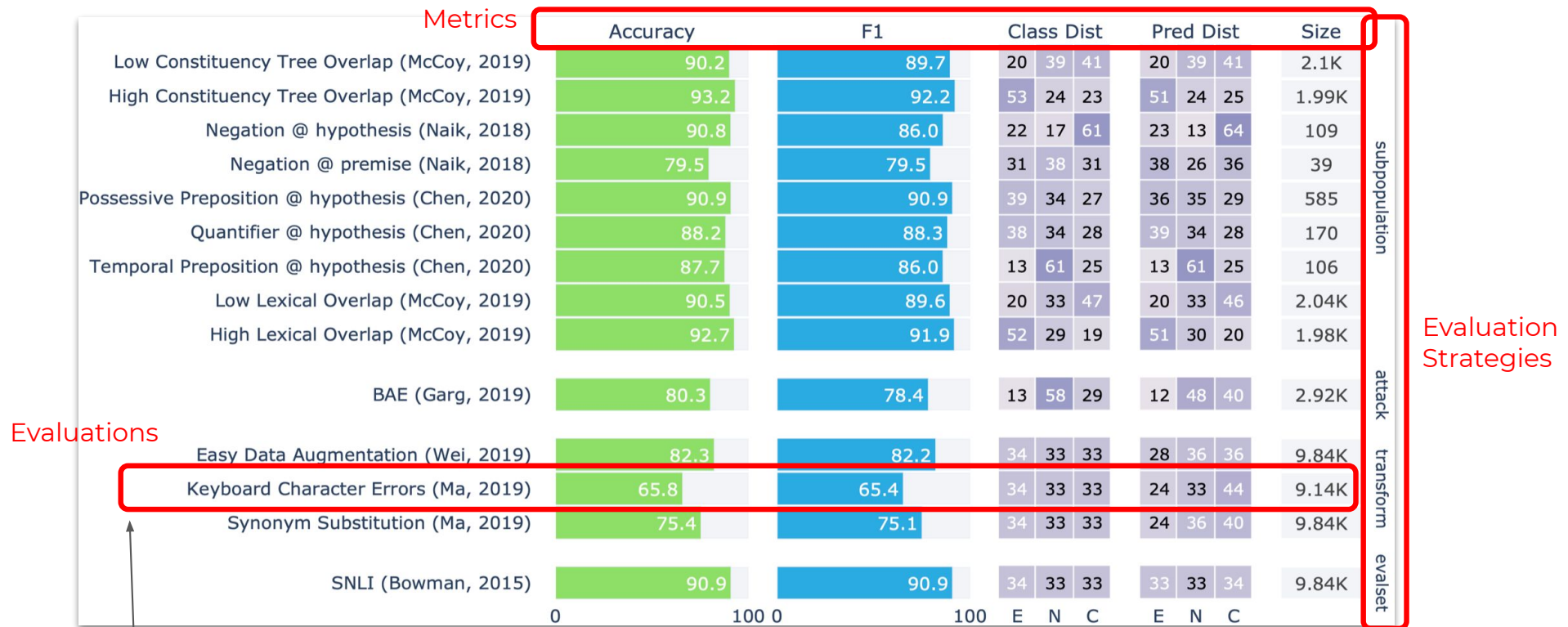


Shift towards **fine-grained evaluation** with new tools (e.g. Robustness Gym, Dynabench)

# Tool: Robustness Gym



Consolidates different evaluation strategies (slices, transformations) and metrics



**Example:** BERT embeddings are sensitive to character errors

# Tool: Robustness Gym

Consolidates many different evaluation types (subpopulations, transformations) and metrics

	Metrics											
	Accuracy		F1		Class Dist			Pred Dist			Size	
Low Constituency Tree Overlap (McCoy, 2019)	90.2		89.7		20	39	41	20	39	41	2.1K	subpopulation
High Constituency Tree Overlap (McCoy, 2019)	93.2		92.2		53	24	23	51	24	25	1.99K	
Negation @ hypothesis (Naik, 2018)	90.8		86.0		22	17	61	23	13	64	109	
Negation @ premise (Naik, 2018)	79.5		79.5		31	38	31	38	26	36	39	
Possessive Preposition @ hypothesis (Chen, 2020)	90.9		90.9		39	34	27	36	35	29	585	
Quantifier @ hypothesis (Chen, 2020)	88.2		88.3		38	34	28	39	34	28	170	
Unlikely Gender (McCoy, 2019)	90.9		90.9		13	61	25	13	61	25	106	
Unlikely Gender (McCoy, 2019)	90.9		90.9		33	33	33	20	33	46	2.04K	
High Lexical Overlap (McCoy, 2019)	92.7		91.9		52	29	19	51	30	20	1.98K	
BAE (Garg, 2019)	80.3		78.4		13	58	29	12	48	40	2.92K	
Easy Data Augmentation (Wei, 2019)	82.3		82.2		34	33	33	28	36	36	9.84K	transform
Keyboard Character Errors (Ma, 2019)	65.8		65.4		34	33	33	24	33	44	9.14K	
Synonym Substitution (Ma, 2019)	75.4		75.1		34	33	33	24	36	40	9.84K	
SNLI (Bowman, 2015)	90.9		90.9		34	33	33	33	33	34	9.84K	
	0	100	0	100	E	N	C	E	N	C		evalset

## Emerging questions

- Discovering important failure modes automatically
- Understanding knowledge captured by an embedding

Evaluations

**Example:** BERT embeddings are sensitive to character errors

# Important Evaluation Strategy: Slice-Based Evaluation

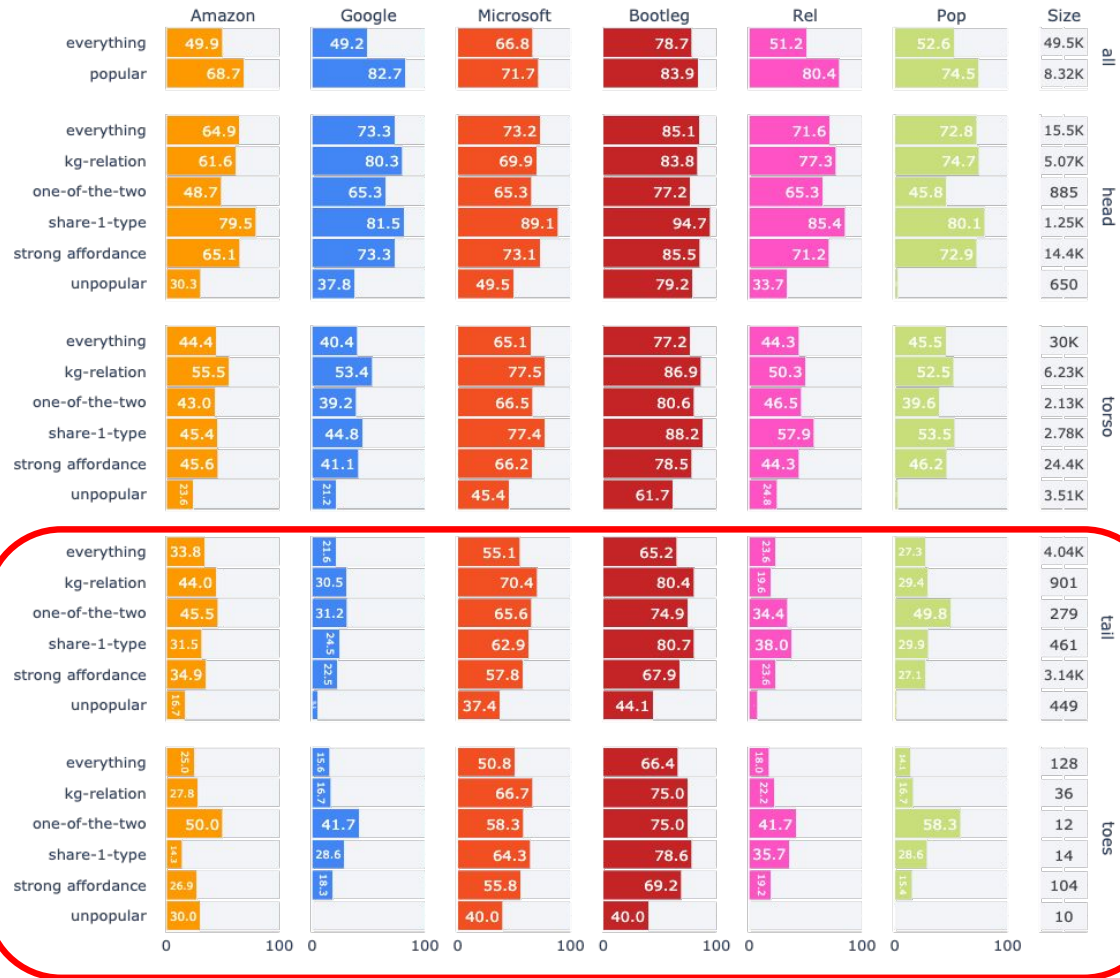
A type of fine-grained evaluation

→ Measure fine-grained performance on critical subpopulations (filtering)

## **Example:**

short passages (< 50 words) in a text dataset

# Example: Named Entity Linking



*Most Named Entity Linking systems are poor on rare entities*

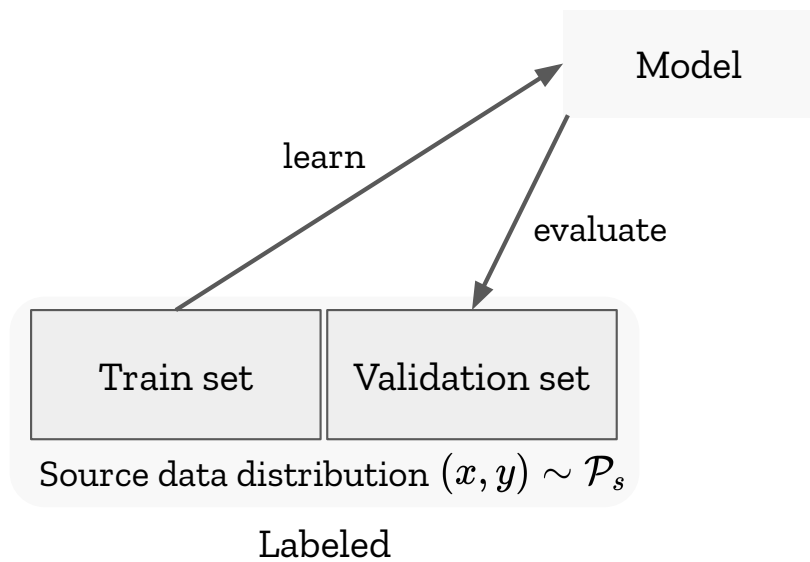
Robustness Gym: Unifying the NLP Evaluation Landscape. NAACL Demo 2021.

Goodwill Hunting: Analyzing and Repurposing Off-the-Shelf Named Entity Linking Systems. NAACL Industry 2021.



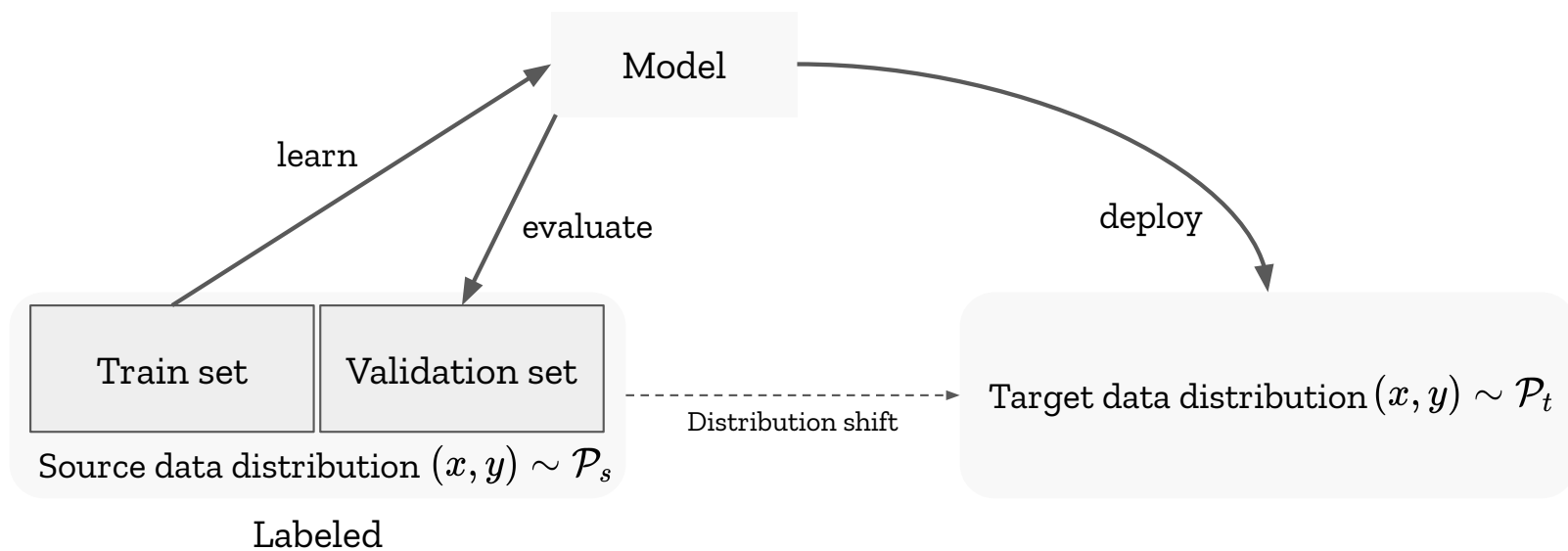
# Evaluation over Time: **Monitoring**

Continually evaluate as the world changes



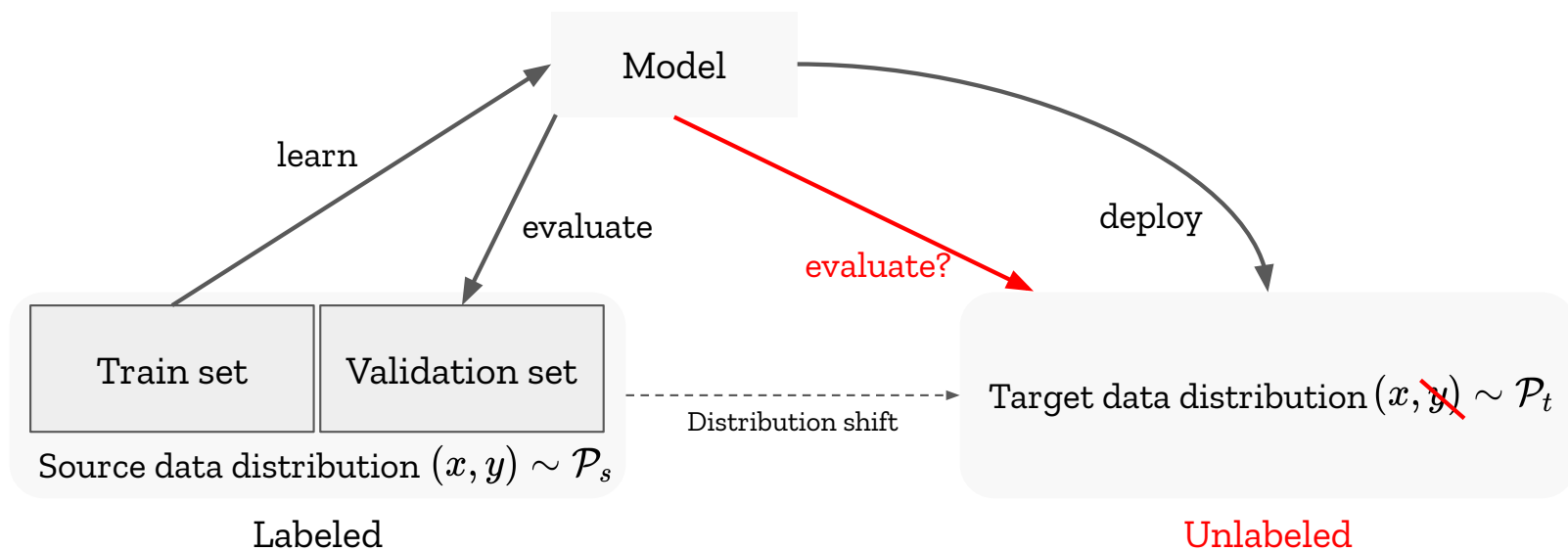
# Evaluation over Time: **Monitoring**

Continually evaluate as the world changes



# Evaluation over Time: **Monitoring**

Continually evaluate as the world changes



Need to monitor model performance on unlabeled data

# Approach: Importance Weighting

## **Estimate metrics on incoming data**

Upweight examples in our dev set more likely to be seen in the future

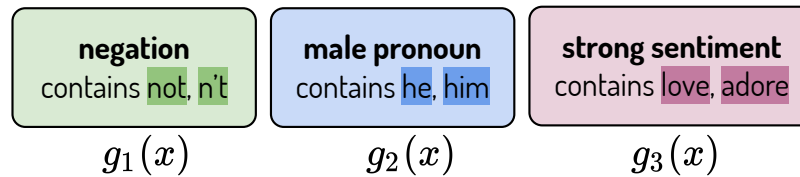
## **Theoretical Foundations**

Density ratio estimation (Sugiyama, 2012)

**Recent work:** accurately estimate performance with *slice-based evaluation + importance weighting*

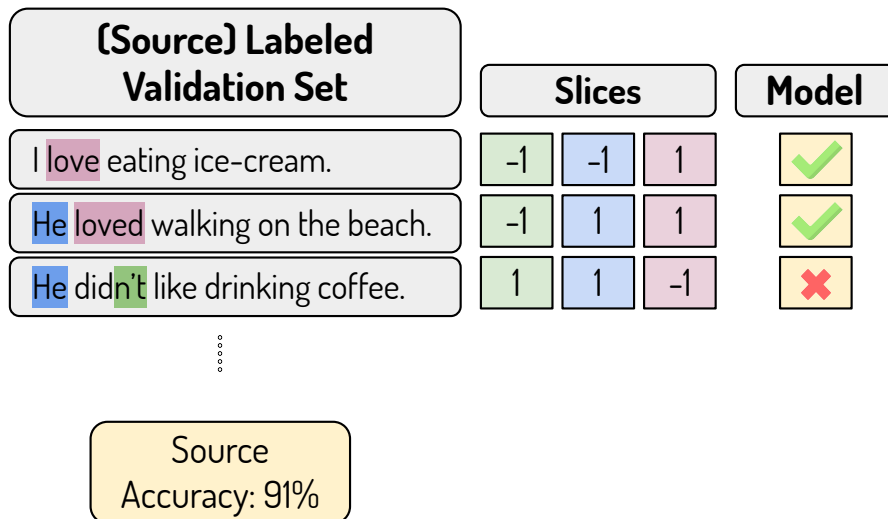
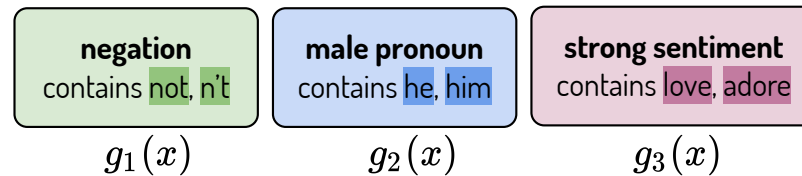
# Mandoline: Slice-based reweighting framework

**Slice:** user-defined grouping of data  $g(x) \in \{-1, 1\}$



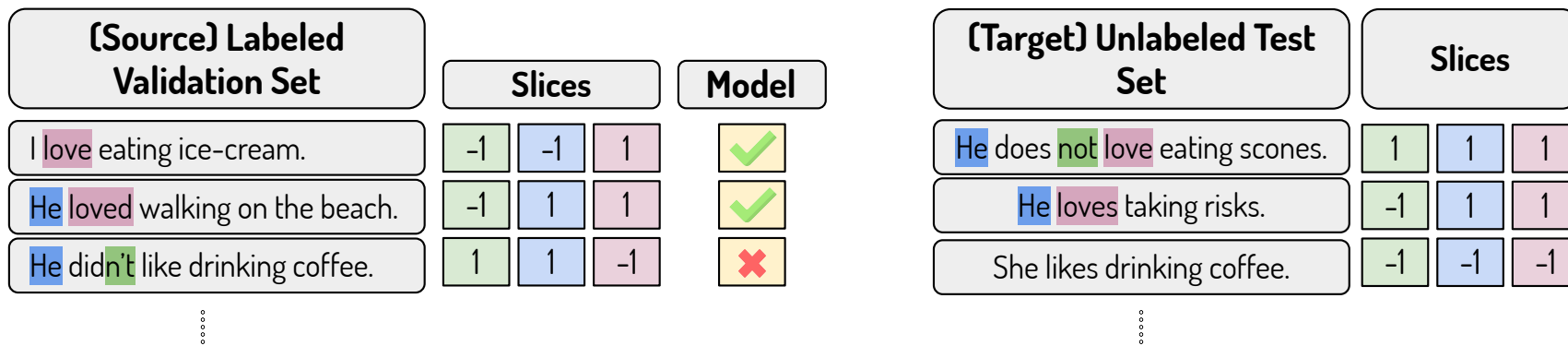
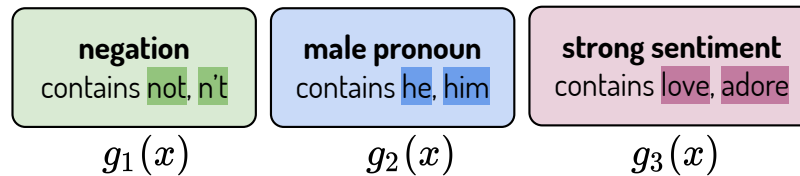
# Mandoline: Slice-based reweighting framework

**Slice:** user-defined grouping of data  $g(x) \in \{-1, 1\}$



# Mandoline: Slice-based reweighting framework

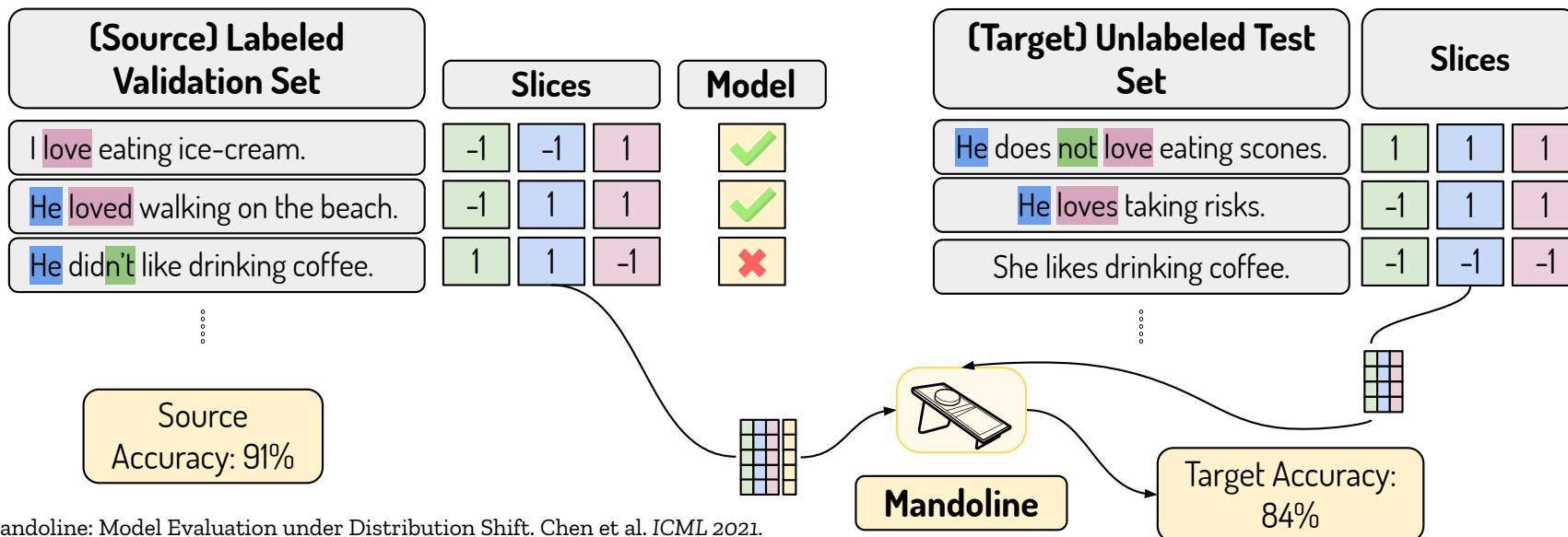
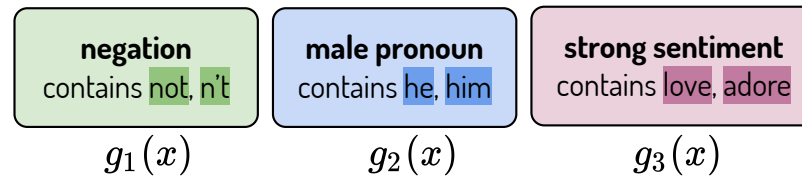
**Slice:** user-defined grouping of data  $g(x) \in \{-1, 1\}$



Source Accuracy: 91%

# Mandoline: Slice-based reweighting framework

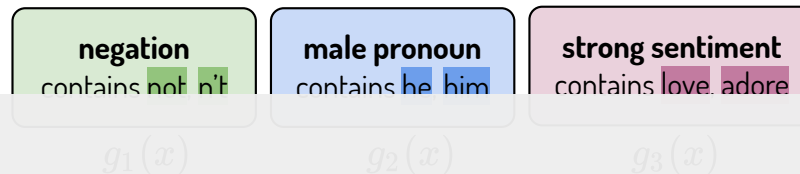
**Slice:** user-defined grouping of data  $g(x) \in \{-1, 1\}$





# Mandoline: Slice-based reweighting framework

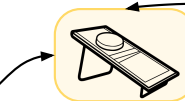
**Slice:** user-defined grouping of data  $g(x) \in \{-1, 1\}$



## Takeaways

- Monitor any model: importance weighting
- Add domain knowledge (slices) to improve monitoring

Source  
Accuracy: 91%



**Mandoline**

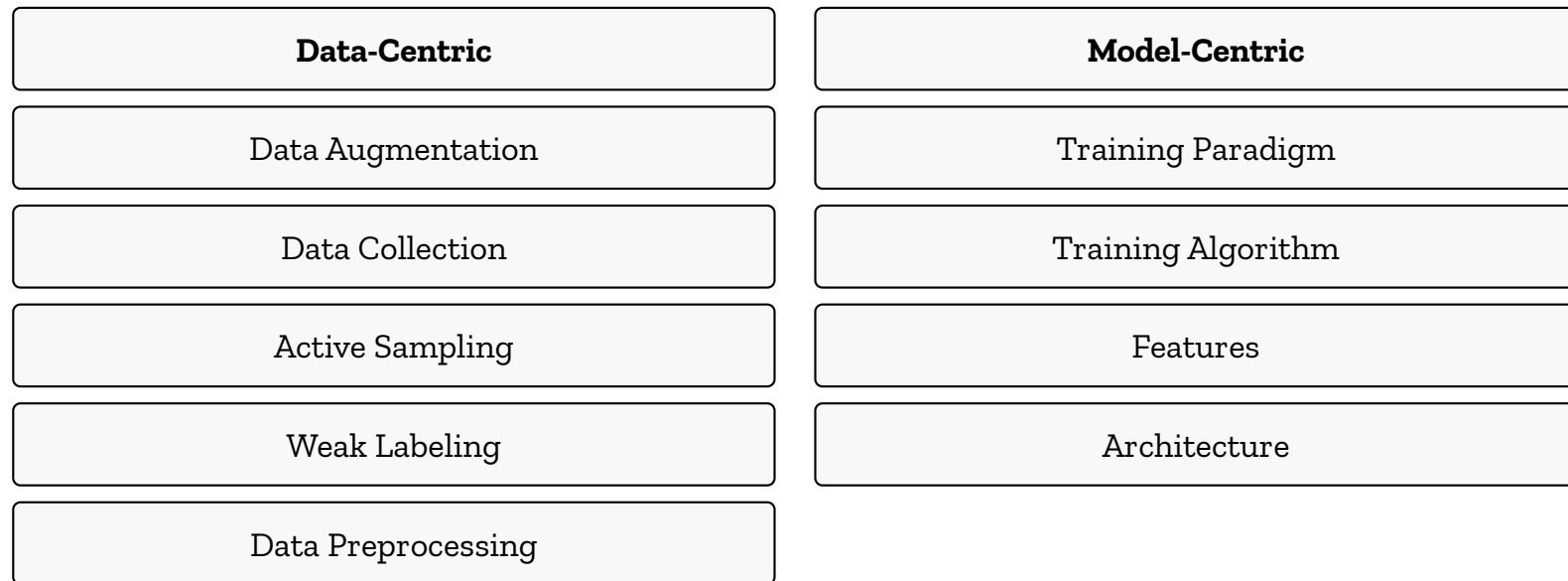
Target Accuracy:  
84%



# Embedding Model **Patching**

Once errors are identified, need to retrain or update embeddings

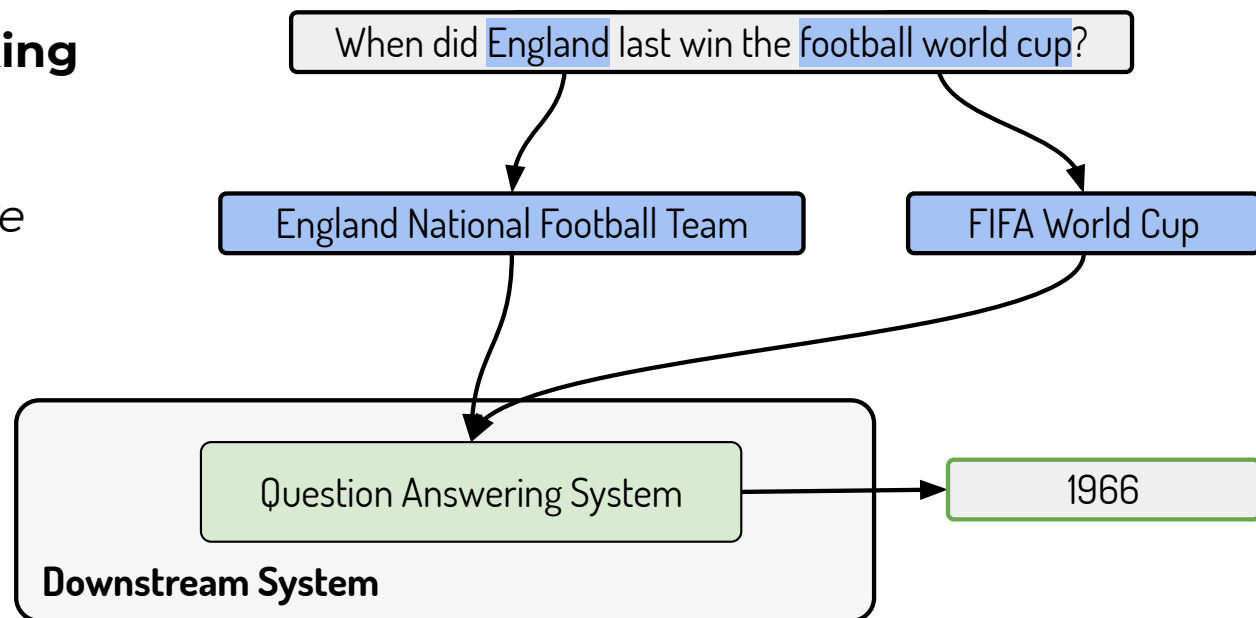
## Many Approaches



# Named Entity Linking

## Named Entity Linking

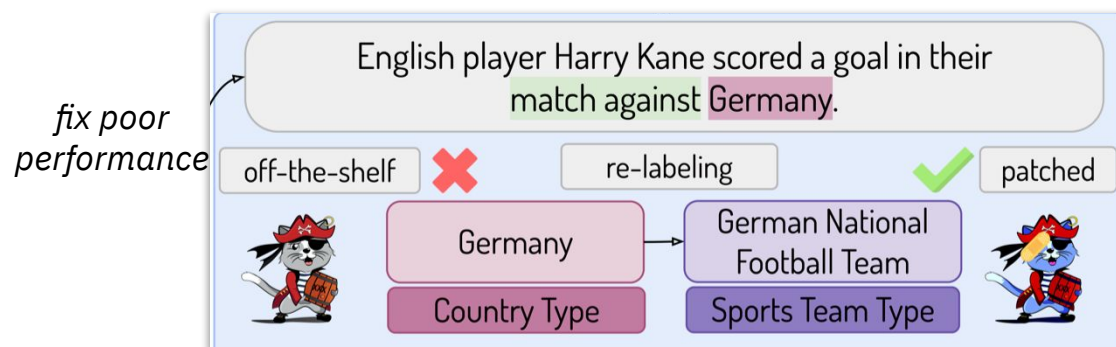
*map "strings" to  
"things" in a  
knowledge base like  
Wikipedia*



**A correct NEL is required for the downstream system!**

# End to End Example: Named Entity Linking

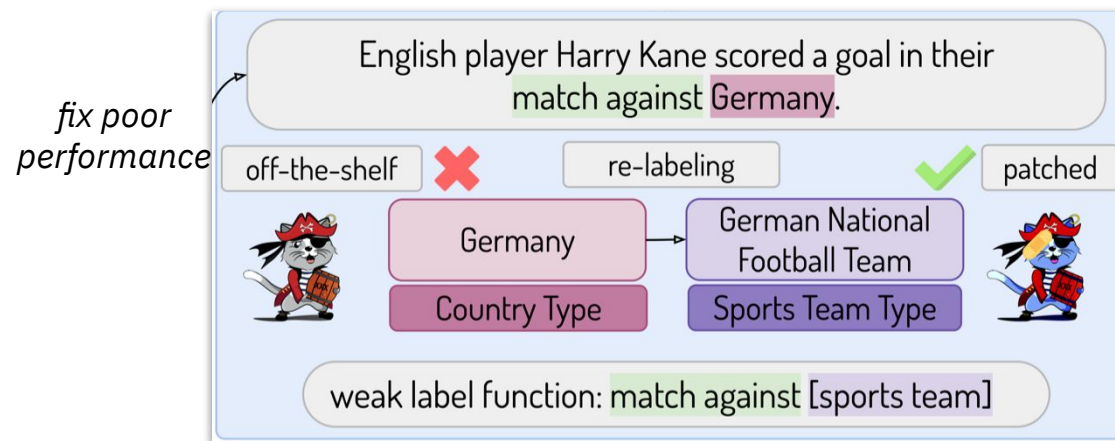
**Repurposing** Bootleg NEL system to patch errors for sports QA



Sports QA: prefer if the model predicted the national sports team instead of the country!

# End to End Example: Named Entity Linking

## Repurposing Bootleg NEL system to patch errors for sports QA



Relabel the training set with a simple heuristic

# End to End Example: Named Entity Linking

**Repurposing** Bootleg NEL system to patch errors for sports QA

*Wikipedia  
examples with  
mentions of  
countries and  
sports teams*

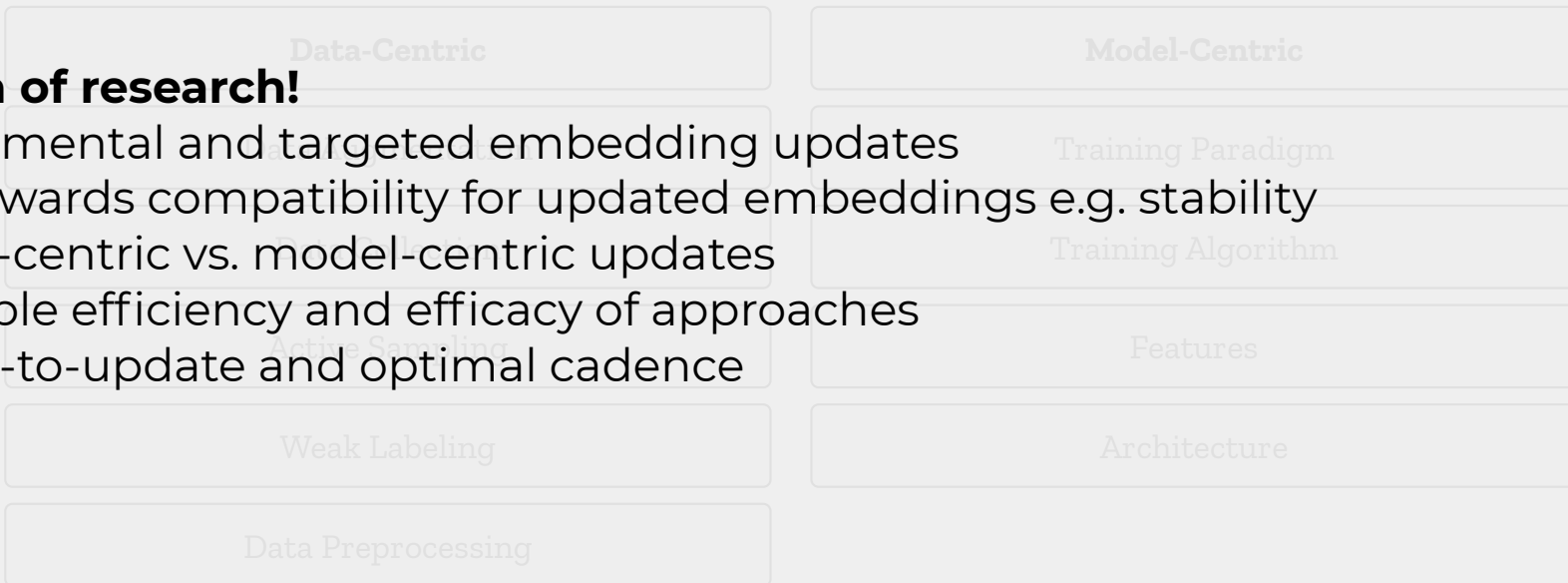
Subpop.	Gold Label	Pred. Label	Size (Off-The-Shelf → Patched)	
	Team	Country	151 → 106	(↓)
		Team	3393 → 3447	(↑)

25% absolute accuracy improvement in sports-related errors

# Embedding Model **Patching**

Once errors are identified, need to retrain or update embeddings

## Many Approaches



### **New area of research!**

- Incremental and targeted embedding updates
- Backwards compatibility for updated embeddings e.g. stability
- Data-centric vs. model-centric updates
- Sample efficiency and efficacy of approaches
- Time-to-update and optimal cadence

# Future Directions



# Embeddings as First Class Citizens

What is the right system for embedding management in ML pipelines?



**Embedding A**



**Embedding B**



## Search

What set of embeddings are best for a specific task?

$(x_1, y_1)$

...

$(x_n, y_n)$



## Provenance

What data had the most "impact" on these embeddings?

*tail performance?*      *syntactic information?*

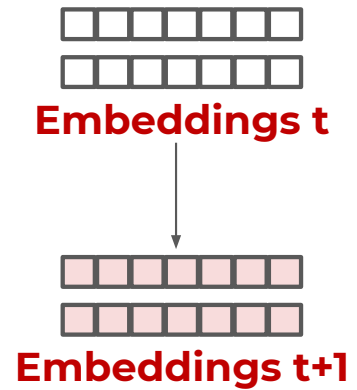
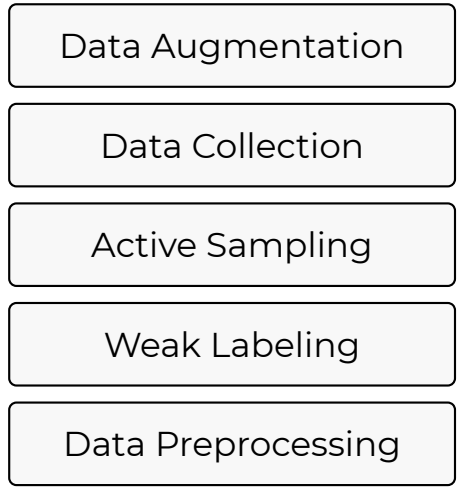
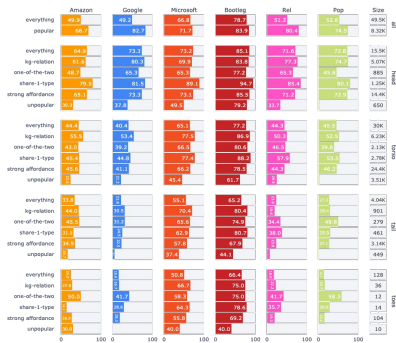


## Quality

What are the right metrics/probes for embedding quality?

# End-to-End Model Patching

How can we automate and provide guidance for embedding patching?



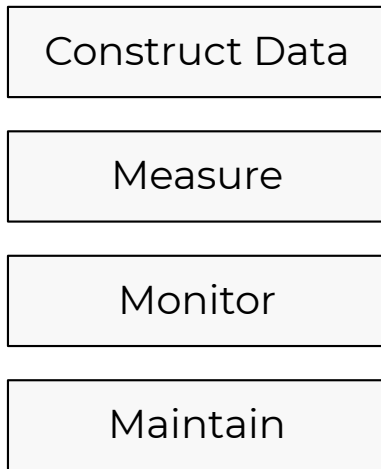
What are the current failure modes?

What data engineering strategy to use?

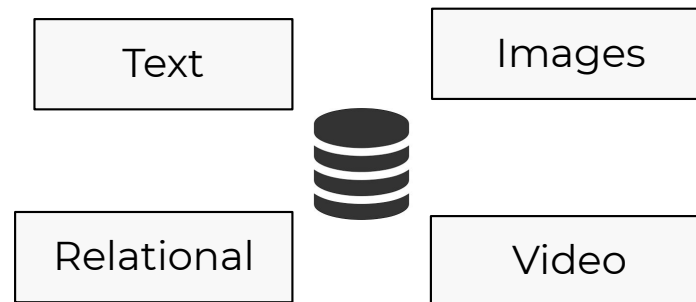
How do I update my models efficiently?

# Interactive Machine Learning

How can we facilitate human interaction with model training and evaluation data?



How can we support the entire lifecycle?



How to integrate multiple modalities?



Data Panels for ML  
Video, Image, Model  
Outputs, ...

Systems to store and manage models and data?