

The Relational Data Borg is Learning

fdbresearch.github.io

relational.ai

Dan Olteanu

University of Zurich

VLDB 2020 Keynote
Virtual Tokyo, Sept 1, 2020



Acknowledgments

FDB team, in particular:



Ahmet



Amir



Haozhe

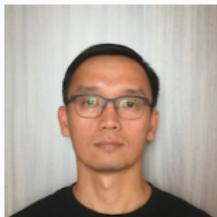


Max



Milos

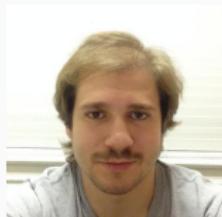
RelationalAI team, in particular:



Hung



Long



Mahmoud



Molham

Database Research In Data Science Era

Reasons for DB research community to feel empowered:

1. **Pervasiveness of relational data** in data science
 - Hard fact
2. **Widespread need for efficient data processing**
 - Core to our community's *raison d'être*
3. **New processing challenges** posed by data science workloads
 - DB's approach reminiscent of Star Trek's Borg Collective

These reasons also serve as motivation for our work.

Star Trek Borg

Co-opt technology and knowledge of alien species
to become ever more powerful and versatile



Relational Data Borg

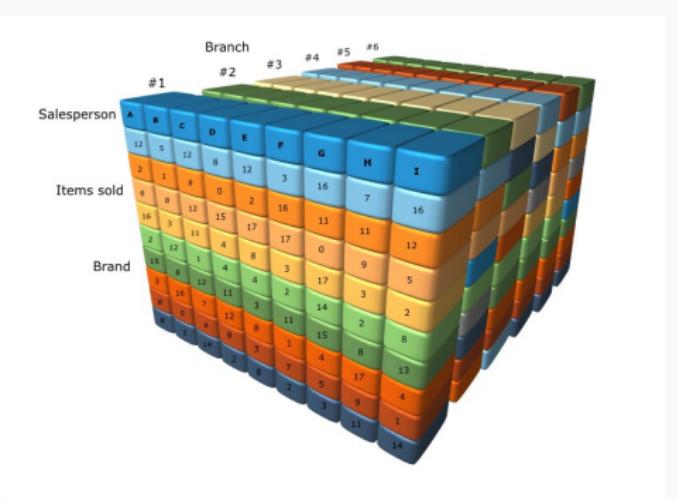
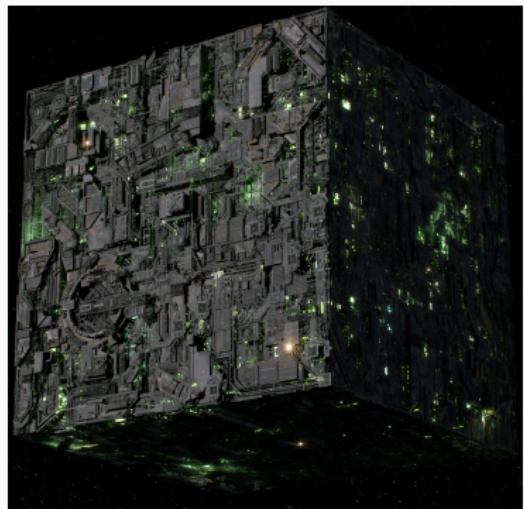


Assimilate ideas and applications of related fields
to adapt to new requirements and
become ever more powerful and versatile

Unlike in Star Trek, the Relational Data Borg

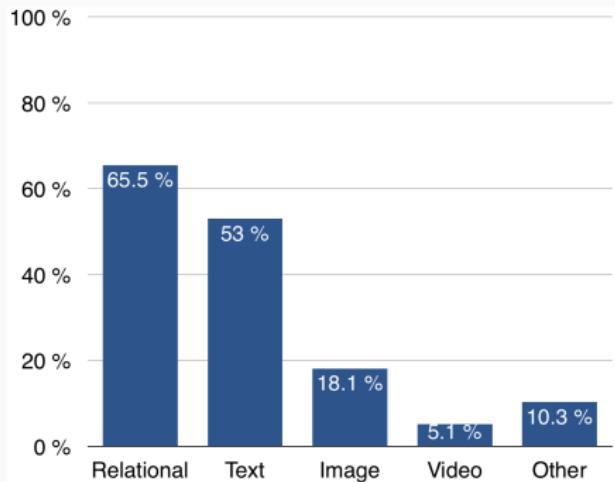
- moves fast
- has great skin complexion and
- is reasonably happy

Borg Cube vs Data Cube

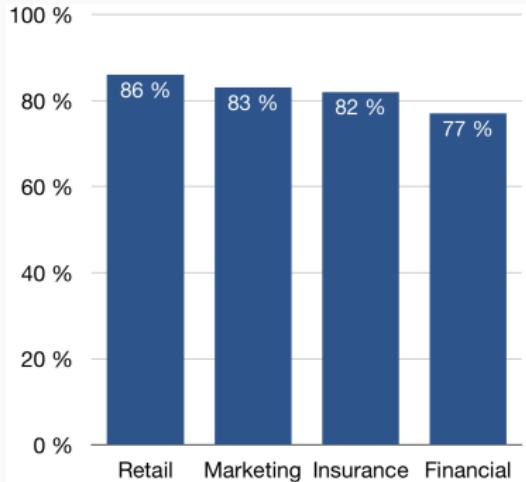


Relational Data is Ubiquitous

Kaggle Survey: Most Data Scientists use Relational Data at Work!



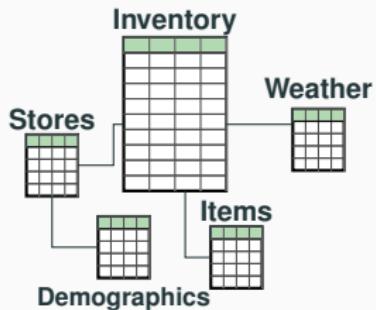
Overall



By Industry

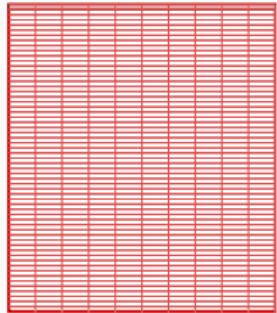
Source: The State of Data Science & Machine Learning 2017, Kaggle, October 2017
(based on 2017 Kaggle survey of 16,000 ML practitioners)

State of Affairs in Learning over Relational Data



Feature Extraction Query →
Inventory \bowtie Stores \bowtie Items
 \bowtie Weather \bowtie Demographics

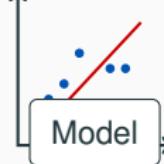
10,000s of Features



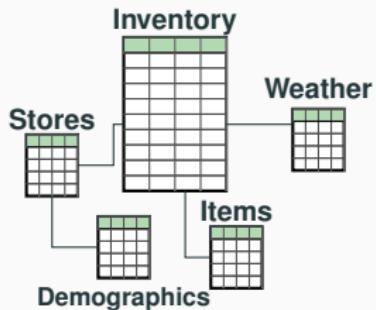
Training Dataset

ML Tool

Model



State of Affairs in Learning over Relational Data

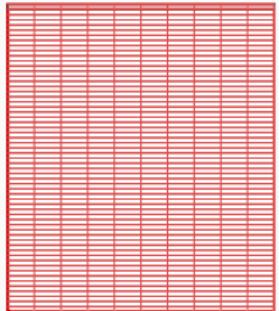


Feature Extraction Query

Inventory \bowtie Stores \bowtie Items
 \bowtie Weather \bowtie Demographics

Relational Data

10,000s of Features

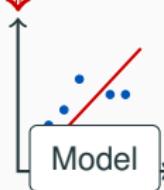


Training Dataset

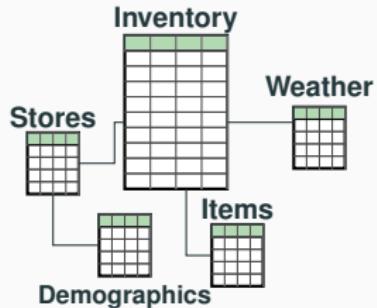
Structure-Agnostic Learning:

ML Tool

Model



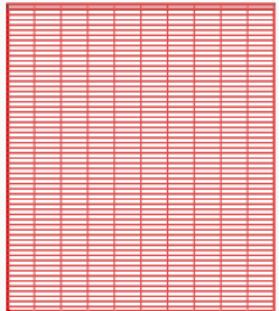
State of Affairs in Learning over Relational Data



Relational Data

Feature Extraction Query →
Inventory \bowtie Stores \bowtie Items
 \bowtie Weather \bowtie Demographics

10,000s of Features



Training Dataset

Structure-Agnostic Learning:

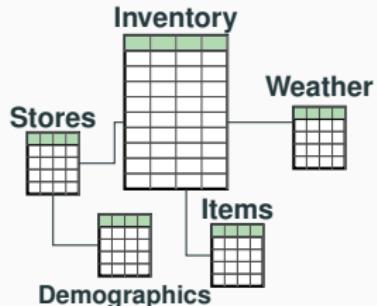
1. Unnecessary data matrix materialization

Relational structure carefully crafted by domain experts thrown away

ML Tool

Model

State of Affairs in Learning over Relational Data

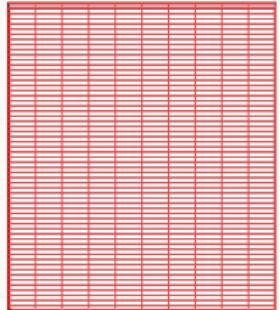


Feature Extraction Query

Inventory \bowtie Stores \bowtie Items
 \bowtie Weather \bowtie Demographics

Relational Data

10,000s of Features



Training Dataset

Structure-Agnostic Learning:

1. **Unnecessary** data matrix materialization

Relational structure carefully crafted by domain experts thrown away

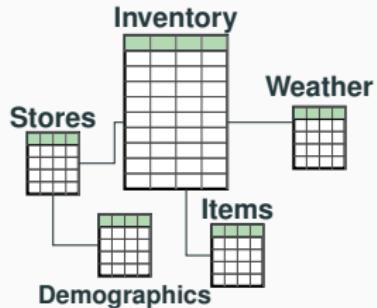
2. **Expensive** data move

Training dataset can be **order-of-magnitude larger** than the input DB

ML Tool

Model

State of Affairs in Learning over Relational Data

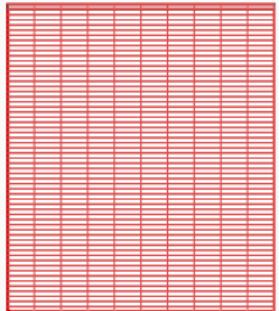


Feature Extraction Query

Inventory \bowtie Stores \bowtie Items
 \bowtie Weather \bowtie Demographics

Relational Data

10,000s of Features



Training Dataset

Structure-Agnostic Learning:

1. **Unnecessary** data matrix materialization

Relational structure carefully crafted by domain experts thrown away

2. **Expensive** data move

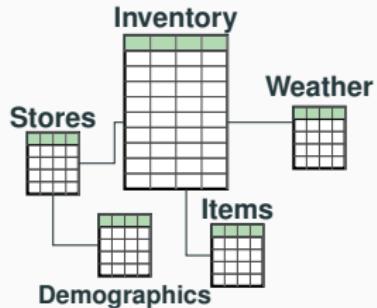
Training dataset can be **order-of-magnitude larger** than the input DB

3. **Bloated** one-hot encoding

ML Tool

Model

State of Affairs in Learning over Relational Data

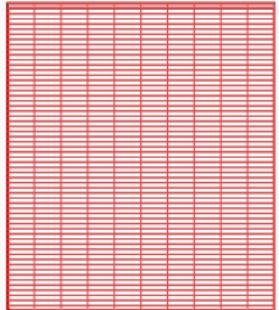


Feature Extraction Query

Inventory \bowtie Stores \bowtie Items
 \bowtie Weather \bowtie Demographics

Relational Data

10,000s of Features



Training Dataset

Structure-Agnostic Learning:

1. **Unnecessary** data matrix materialization

Relational structure carefully crafted by domain experts thrown away

2. **Expensive** data move

Training dataset can be **order-of-magnitude larger** than the input DB

3. **Bloated** one-hot encoding

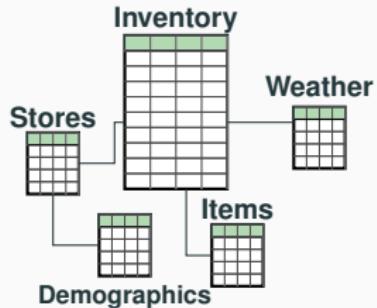
4. **High** maintenance cost

Recomputation from scratch after updates

ML Tool

Model

State of Affairs in Learning over Relational Data

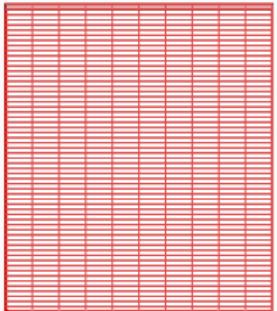


Feature Extraction Query

Inventory \bowtie Stores \bowtie Items
 \bowtie Weather \bowtie Demographics

Relational Data

10,000s of Features



Training Dataset

Structure-Agnostic Learning:

1. **Unnecessary** data matrix materialization

Relational structure carefully crafted by domain experts thrown away

2. **Expensive** data move

Training dataset can be **order-of-magnitude larger** than the input DB

3. **Bloated** one-hot encoding

4. **High** maintenance cost

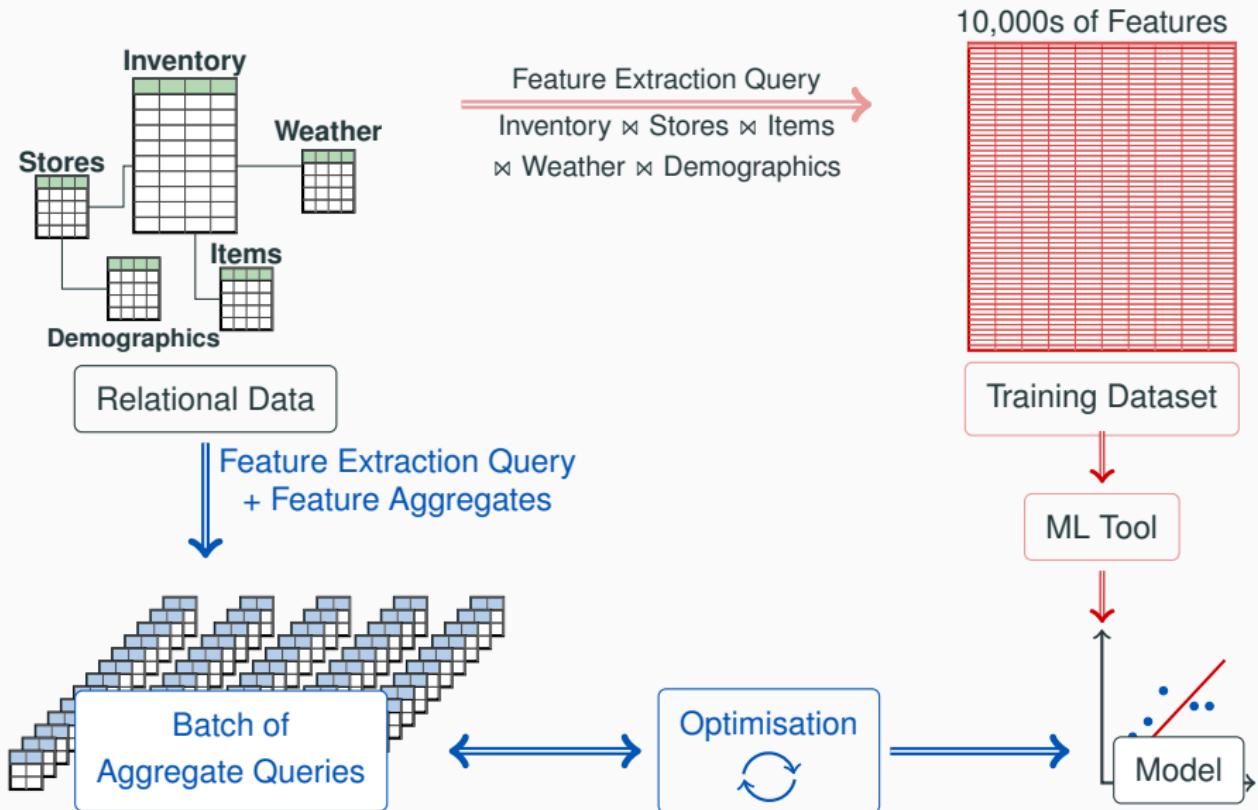
Recomputation from scratch after updates

5. **Limitations** inherited from both DB and ML tools

ML Tool

Model

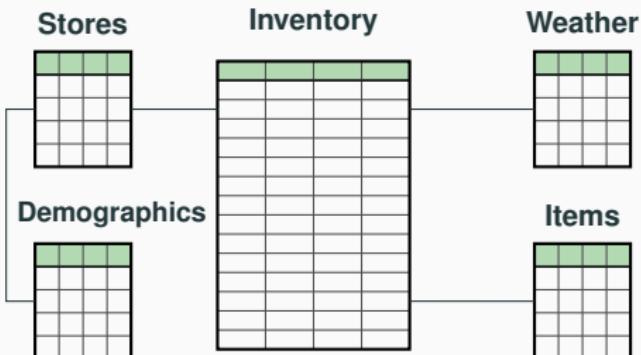
Structure-Aware Learning over Relational Data



Conjecture

The learning time and accuracy of the model can be drastically improved by exploiting the structure and semantics of the underlying multi-relational database.

Structure-aware Learning FASTER than Feature Extraction Query alone



Relation	Cardinality	Arity (Keys+Values)	File Size (CSV)
Inventory	84,055,817	3 + 1	2 GB
Items	5,618	1 + 4	129 KB
Stores	1,317	1 + 14	139 KB
Demographics	1,302	1 + 15	161 KB
Weather	1,159,457	2 + 6	33 MB
Join	84,055,817	3 + 41	23GB

Structure-aware versus Structure-agnostic Learning

Train a linear regression model to predict *inventory* given all features

PostgreSQL+TensorFlow

	Time	Size (CSV)
Database	—	2.1 GB
Join	152.06 secs	23 GB
Export	351.76 secs	23 GB
Shuffling	5,488.73 secs	23 GB
Query batch	—	—
Grad Descent	7,249.58 secs	—
Total time	13,242.13 secs	

Structure-aware versus Structure-agnostic Learning

Train a linear regression model to predict *inventory* given all features

	PostgreSQL+TensorFlow		Our approach (SIGMOD'19)	
	Time	Size (CSV)	Time	Size (CSV)
Database	–	2.1 GB	–	2.1 GB
Join	152.06 secs	23 GB	–	–
Export	351.76 secs	23 GB	–	–
Shuffling	5,488.73 secs	23 GB	–	–
Query batch	–	–	6.08 secs	37 KB
Grad Descent	7,249.58 secs	–	0.05 secs	–
Total time	13,242.13 secs		6.13 secs	

2,160× faster while being more accurate (RMSE on 2% test data)

Structure-aware versus Structure-agnostic Learning

Train a linear regression model to predict *inventory* given all features

	PostgreSQL+TensorFlow		Our approach (SIGMOD'19)	
	Time	Size (CSV)	Time	Size (CSV)
Database	–	2.1 GB	–	2.1 GB
Join	152.06 secs	23 GB	–	–
Export	351.76 secs	23 GB	–	–
Shuffling	5,488.73 secs	23 GB	–	–
Query batch	–	–	6.08 secs	37 KB
Grad Descent	7,249.58 secs	–	0.05 secs	–
Total time	13,242.13 secs		6.13 secs	

2,160× faster while being more accurate (RMSE on 2% test data)

TensorFlow trains one model. Our approach takes < 0.1 sec for any extra model over a subset of the given feature set.

TensorFlow's Behaviour is the Rule, not the Exception!

Similar behaviour (or outright failure) for more:

- **datasets**: Favorita, TPC-DS, Yelp, Housing
- **systems**:
 - used in industry: R, scikit-learn, Python StatsModels, mpack, XGBoost, MADlib
 - academic prototypes: Morpheus, libFM
- **models**: decision trees, factorisation machines, k -means, ..

This is to be contrasted with the scalability of DBMSs!

**How to achieve this performance
improvement?**

Idea 1: Turn the ML Problem into a DB Problem



Through DB Glasses, Everything is a Batch of Queries

Workload	Query Batch
Linear Regression	$\text{SUM}(X_i * X_j)$
Covariance Matrix	$\text{SUM}(X_i) \text{ GROUP BY } X_j$ $\text{SUM}(1) \text{ GROUP BY } X_i, X_j$
Decision Tree Node	$\text{VARIANCE}(Y) \text{ WHERE } X_j = c_j$
Mutual Information	$\text{SUM}(1) \text{ GROUP BY } X_i$
Rk-means	$\text{SUM}(1) \text{ GROUP BY } X_j$ $\text{SUM}(1) \text{ GROUP BY Center}_1, \dots, \text{Center}_k$

Through DB Glasses, Everything is a Batch of Queries

Workload	Query Batch
Linear Regression	$\text{SUM}(X_i * X_j)$ [WHERE $\sum_k X_k * w_k < c$]
Covariance Matrix	$\text{SUM}(X_i)$ GROUP BY X_j [WHERE ...]
(Non)poly. loss	$\text{SUM}(1)$ GROUP BY X_i, X_j [WHERE ...]
Decision Tree Node	$\text{VARIANCE}(Y)$ WHERE $X_j = c_j$
Mutual Information	$\text{SUM}(1)$ GROUP BY X_i
Rk-means	$\text{SUM}(1)$ GROUP BY X_j $\text{SUM}(1)$ GROUP BY Center ₁ , ..., Center _k

Through DB Glasses, Everything is a Batch of Queries

Workload	Query Batch	# Queries
Linear Regression	$\text{SUM}(X_i * X_j)$ [WHERE $\sum_k X_k * w_k < c$]	814
Covariance Matrix	$\text{SUM}(X_i)$ GROUP BY X_j [WHERE ...]	
(Non)poly. loss	$\text{SUM}(1)$ GROUP BY X_i, X_j [WHERE ...]	
Decision Tree Node	$\text{VARIANCE}(Y)$ WHERE $X_j = c_j$	3,141
Mutual Information	$\text{SUM}(1)$ GROUP BY X_i	56
Rk-means	$\text{SUM}(1)$ GROUP BY X_j $\text{SUM}(1)$ GROUP BY Center ₁ , ..., Center _k	41

(# Queries shown for Retailer dataset with 39 attributes)

Queries in a batch:

- Same aggregates but over different attributes
- Expressed over the same join of the database relations

AMPLE opportunities for sharing computation in a batch.

Models under Consideration

So far:

- Polynomial regression
- Factorisation machines
- Classification/regression trees
- Mutual information
- Chow Liu trees
- k -means clustering
- k -nearest neighbours
- (robust, ordinal) PCA
- SVM

On-going:

- Boosting regression trees
- AdaBoost
- Sum-product networks
- Random forests
- Logistic regression
- Linear algebra:
 - QR decomposition
 - SVD
 - low-rank matrix factorisation

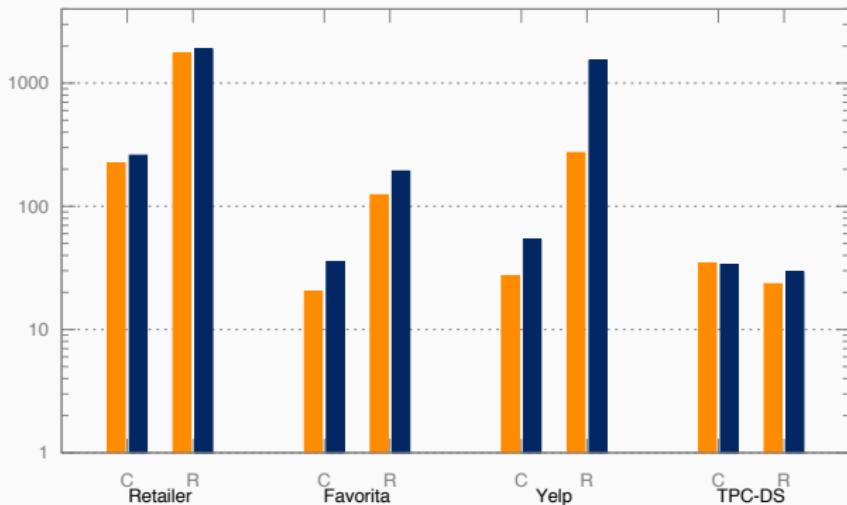
All these cases can benefit from **structure-aware computation**

Natural Attempt:

Use Existing DB System to Compute Query Batch

Existing DBMSs are **NOT** Designed for Query Batches

Relative Speedup for **Our Approach** over **DBX** and **MonetDB**



C = Covariance Matrix; R = Regression Tree Node; AWS d2.xlarge (4 vCPUs, 32GB)

Existing DSMSs are **NOT** Designed for Query Batches

Task: Maintain the covariance matrix over Retailer

- Round-robin insertions in all relations
- All maintenance strategies implemented in DBToaster



Azure DS14, Intel Xeon, 2.40GHz, 112GB, 1 thread; one hour timeout