

Wayfair's Mercury Platform: Scaling ML Applications via Programmatic Feature Definition, Build, and Maintenance

Alexander Hristov, Staff ML Scientist, Wayfair



Key Idea for Mercury:

- Prioritize **maintainability** and **efficiency** over flexibility in defining new features.
- Improve model performance by increasing access to signals
- **Define, build, consume, and maintain** features as a **group**, rather than individually.



Key Idea for Mercury:

- Prioritize **maintainability** and **efficiency** over flexibility in defining new features.
- Improve model performance by increasing access to signals
- **Define, build, consume, and maintain** features as a **group**, rather than individually.

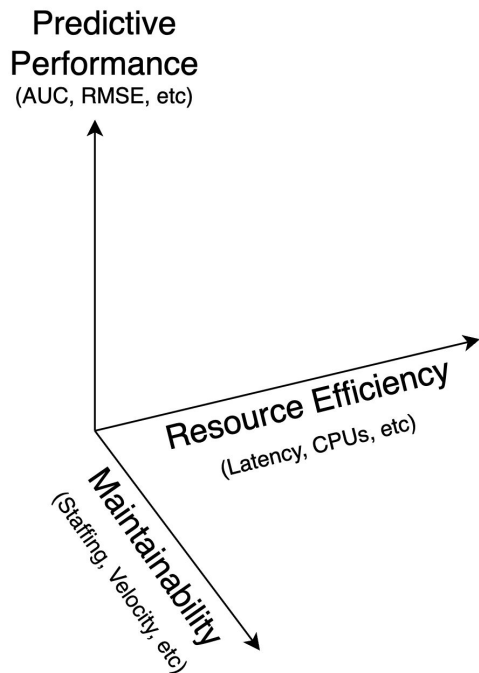


Part 1: Motivation

Why Trade Flexibility for Efficiency in Feature Platform Design

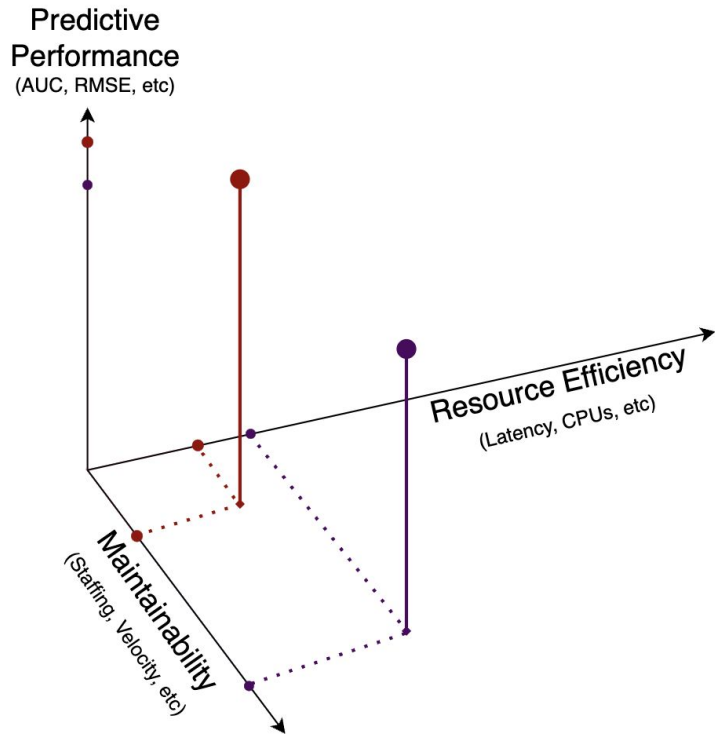


All ML systems make tradeoffs.

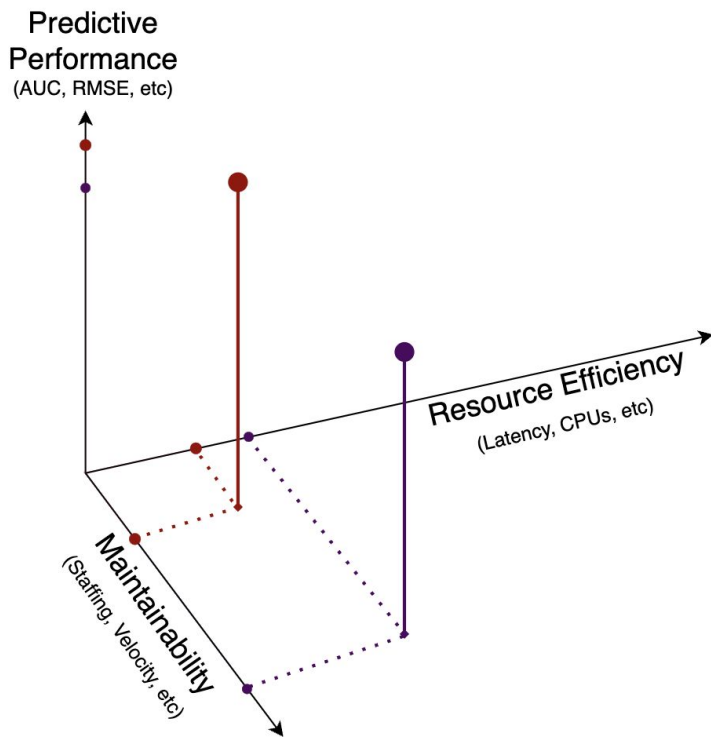




Which ML system would you prefer to deploy?



Which ML system would you prefer to deploy?



It depends on ...

- Team: size and maturity
- Resources: budget and elasticity
- Problem value, i.e. $d\$/dRMSE$
- What else could you be doing?
 - Wide teams: $0 \rightarrow 1$, then onto another task
 - Tall teams: Iterate on the same problem

In other words,

**The more predictive model is not always better.
Success comes from managing tradeoffs.**

Feature platform flexibility is an easily overlooked tradeoff.

Nonetheless, determining this tradeoff at platform level can help teams scale.

Very Flexible:

Make any feature in any way

- More focus on “best case” performance
- Hardest to provide cost, latency guarantees at platform level
- Some changes hard to automate

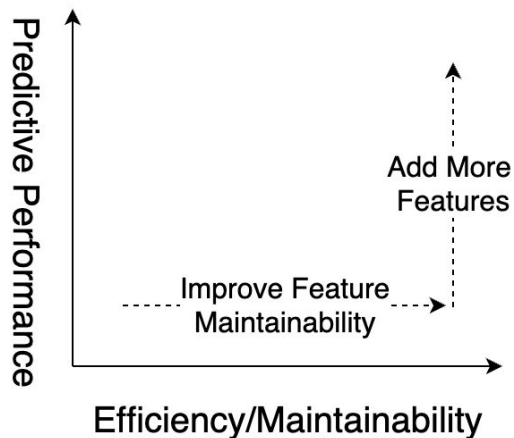
Very Opinionated:

Features must follow specific approach

- More focus on “out of the box” performance
- Stronger platform guarantees for speed, cost
- Easier to automate

Be more opinionated, because ML team time is bounded.

1. Improve predictive performance obtained in bounded time, not the hypothetical best-case performance.
2. Improve not just velocity to deploy new things, but also **carrying capacity**, the ability to support more features and models well.



Mercury is strongly opinionated that users should define, build, consume and maintain features as groups, not one at a time.

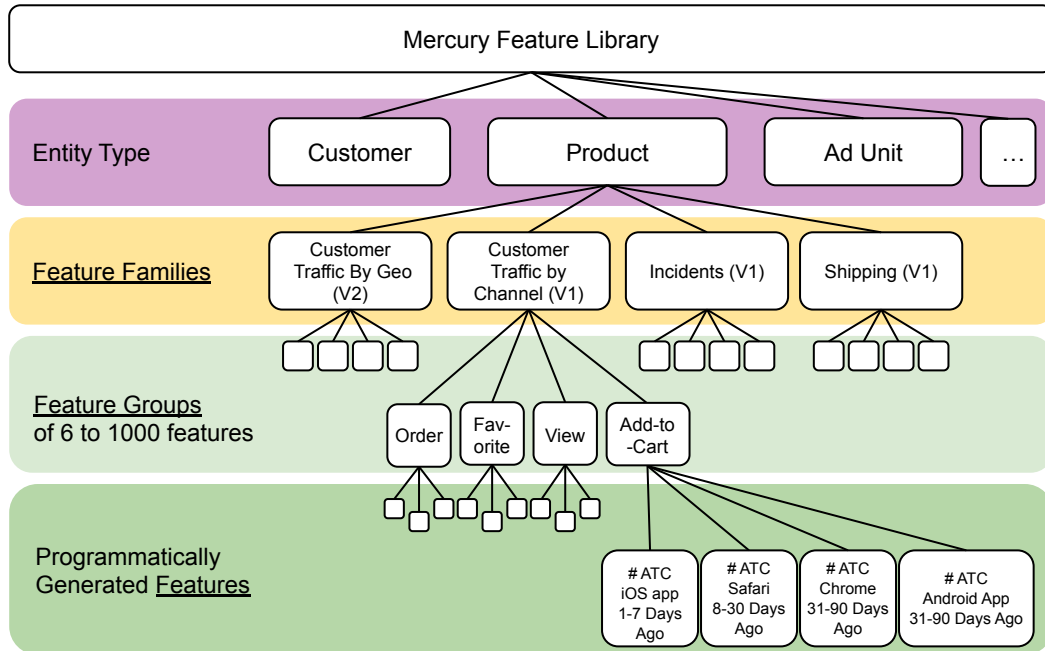
Many tabular features can share a pattern, creating a game of “fill in the blanks”

_____ for a product from sessions on _____ between _____ days ago		
# of add-to-cart	iOS app	1 and 7
avg. # per order	desktop chrome	8 and 30
highest review	mobile safari	31 and 90
...	...	91 and 365
		...

Days since a customer last _____ a _____		
	browsed for	Sofa
	added to cart	Desk

Features with the same **template** and source **event** can build together

Mercury groups features hierarchically for performance and maintainability



For discovery and versioning, group features sharing entity of analysis (Geo, Channel)

For performance, group features together that share actions (atc, clicks, etc)



Programmatic features improve scale and performance in the long run.

1000s

Features Defined

Trillions

Feature Values
Computed Per Day

20+

ML Applications using
Mercury

Billions

Daily Model Predictions

100s

Model objects deployed
use Mercury

User Experience Improvements:

Predictive Performance: 2-9% for 3 preexisting applications that migrated to Mercury

Deploy Speed: 2-5x across multiple teams that used Mercury for new projects

Maintainability: One team reported saving >4 hrs per week per model pipeline

Disclaimer:

Scale isn't free, comes with restrictions on latency and flexibility that may not make sense in all cases.

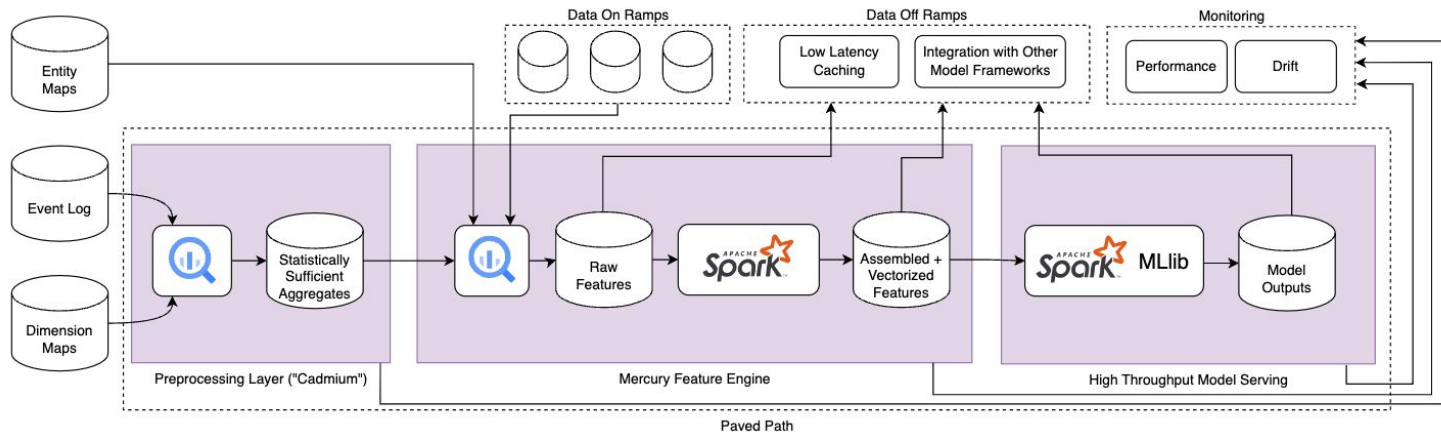


Part 2: Key Architecture Enablers

Defining and building thousands of features

High Level Logical Architecture

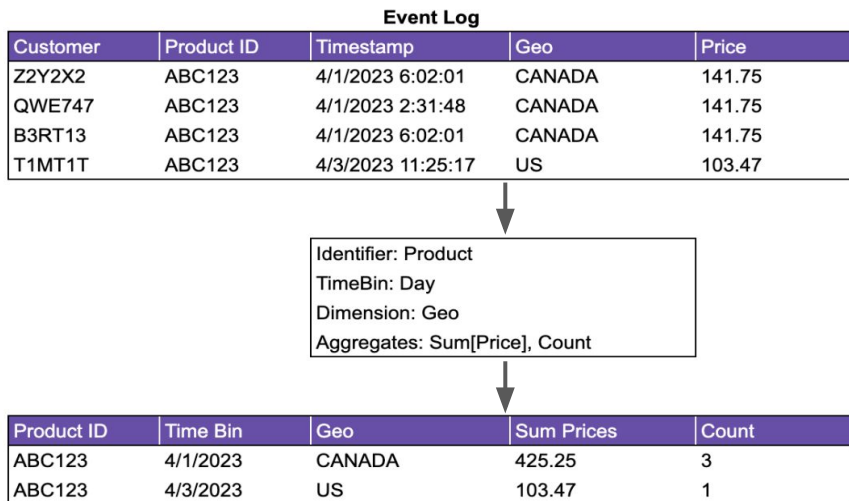
The paved path is highly opinionated, still able to support many applications



Key Lessons:

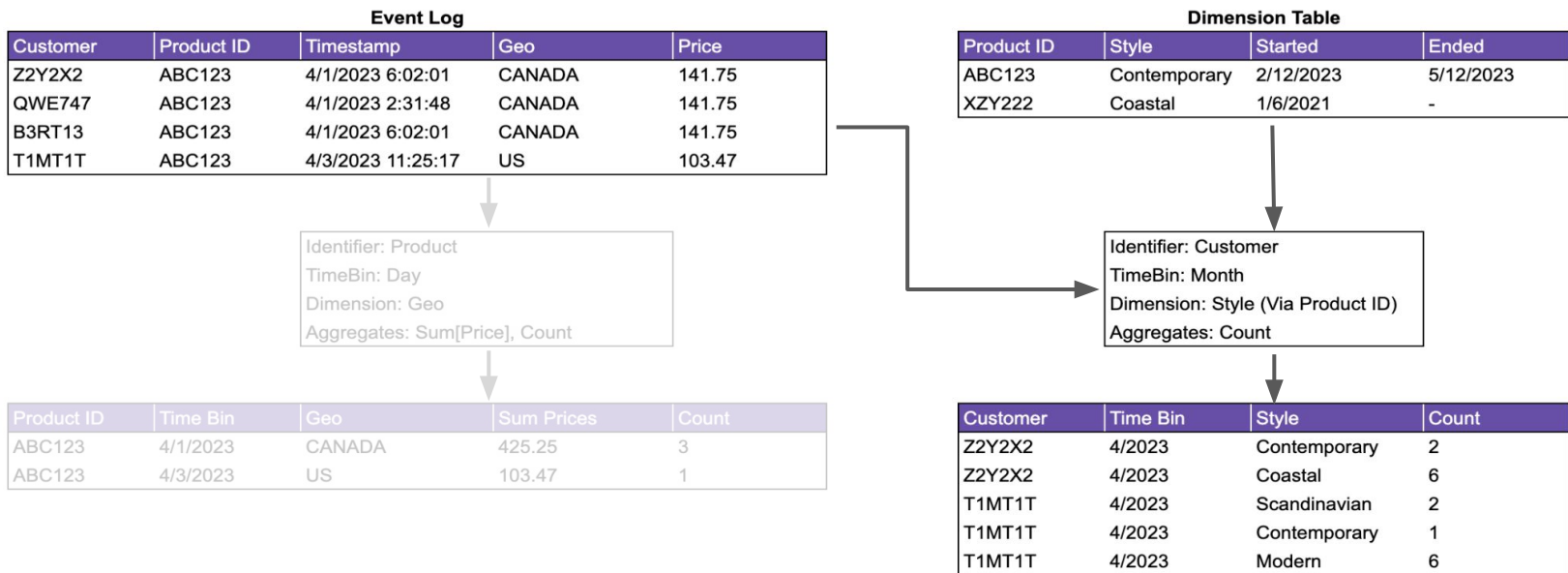
1. Pre-aggregation is an important step for improving maintainability and performance
2. Build features together and package as vectors
3. Emphasize reproducibility of feature computation
4. Drive reuse in multiple ways

Mercury preprocesses event logs to improve maintainability and performance



Try to offload as much compute to this stage: schema standardization, analytic transforms, etc

Mercury preprocesses event logs to improve maintainability and performance



Try to offload as much compute to this stage: schema standardization, analytic transforms, etc

Simple aggregation patterns are flexible enough to support many uses

Supports statistically sufficient aggregates for



Counts & Totals:

store a daily count or sum



Averages & Quantiles:

Store daily sum or %-tile and count



Ratios:

Store daily sums for the two values

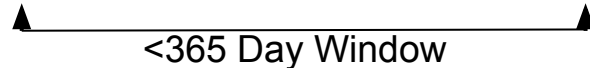


Min & Max:

Store daily min and max

... and different windows

Daily - Used for features with up to a 1 year lookback window



Monthly - Used for features with a lifetime-lookback window



Being opinionated does however obstruct some features, like

- distinct counts
- sequence graph transforms

Features build all at once and are packaged into a sparse vector

Example: Avg. price & total count of orders per supplier in (geo) over previous # to # days

Requested Entities

Supplier ID	Date
AA	2023-05-01 00:00:00
YY	2023-05-01 00:00:00

Entity Map Snapshot

Supplier ID	Date	Product ID
AA	2023-05-01 00:00:00	ABC123
YY	2023-05-01 00:00:00	DEF456
YY	2023-05-01 00:00:00	TUV789

Preaggregated Data

Product ID	Time Bin	Geo	Sum Prices	Count
ABC123	2023-04-15	CA	500	3
ABC123	2023-04-12	CA	100	1
DEF456	2023-03-01	DE	160	4
DEF456	2023-03-01	US	181	6
TUV789	2023-04-01	US	144	7

- Reduces repeated scans
- All features within the group are “deployed”
- Vectors from multiple groups can be concatenated

Supplier ID	Date	Time Diff	Geo	Count	Average Price
AA	2023-05-01	(7, 30]	CA	4	150
YY	2023-05-01	(30, 90]	US	4	40
YY	2023-05-01	(30, 90]	DE	13	25

Relative Position

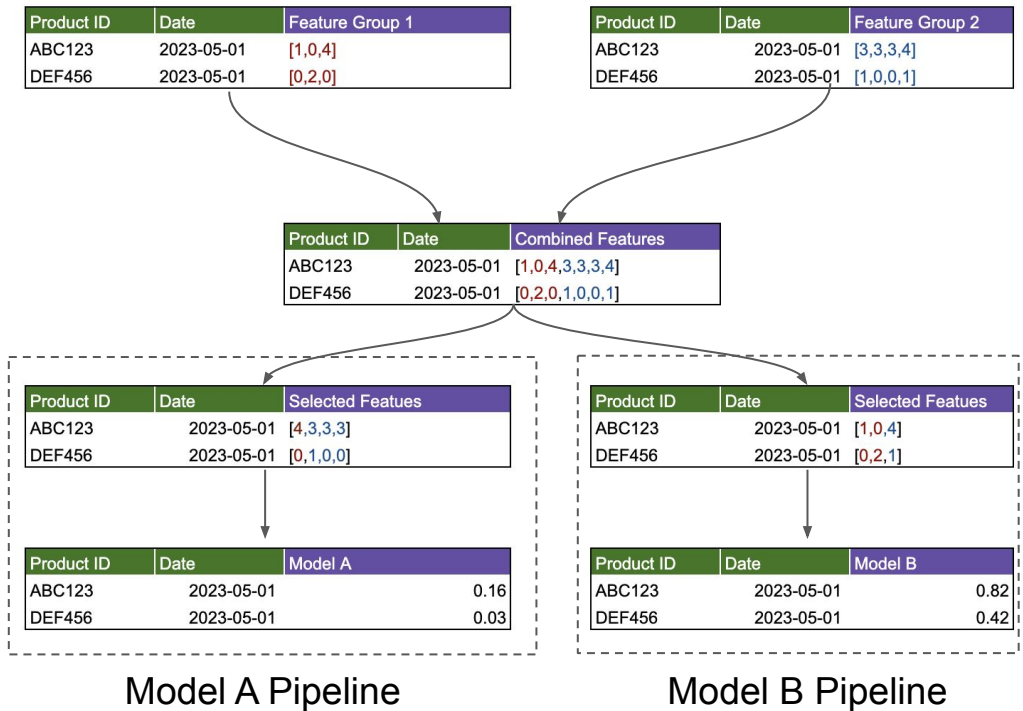
Values

Supplier ID	Date	Features
AA	2023-05-01	[0,.....,4,150,.....]
YY	2023-05-01	[0,.....,13,25,.....,4,40,.....]



**All identifiers and numbers for illustrative purposes only and do not reflect real data.

Despite building features as large vectors, we are still have flexibility in how models consume the features through SparkML



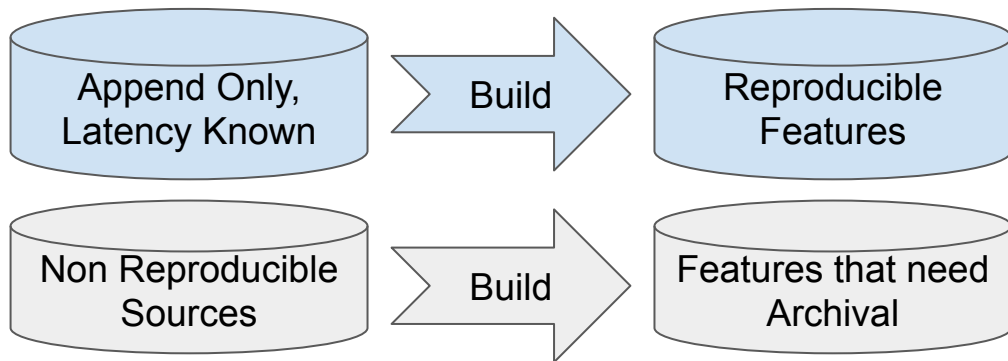
Feature Group Concatenation

Model Specific Feature Selection
via spark VectorSlicer

Inference

Write-only data sources help with reproducible compute, tunable retention

Path towards online-offline consistency often depends on the contract with upstream data.

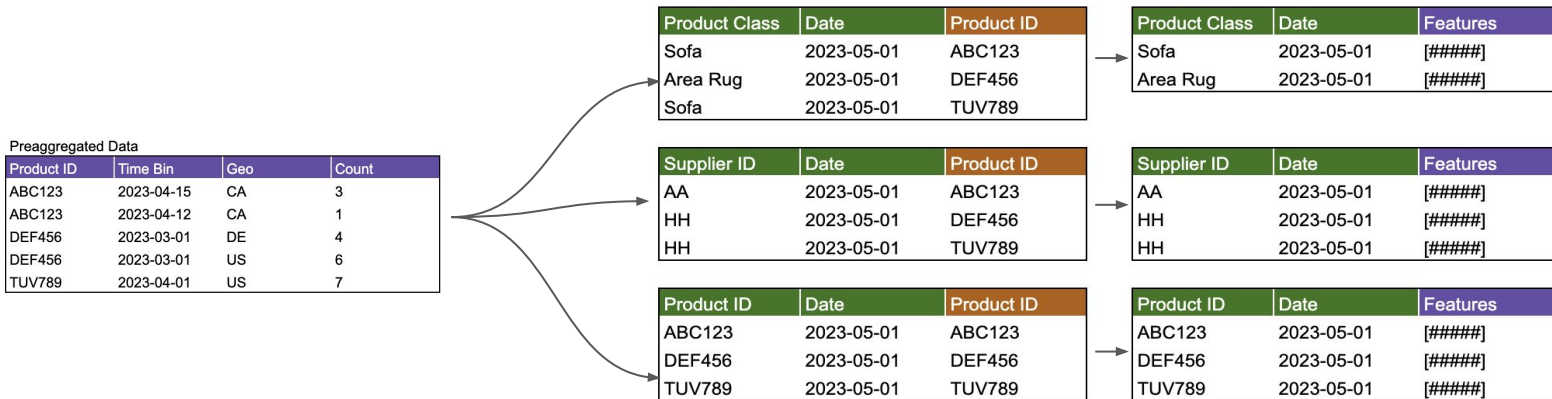


This then determines the requirements and capabilities of a feature platform, such as

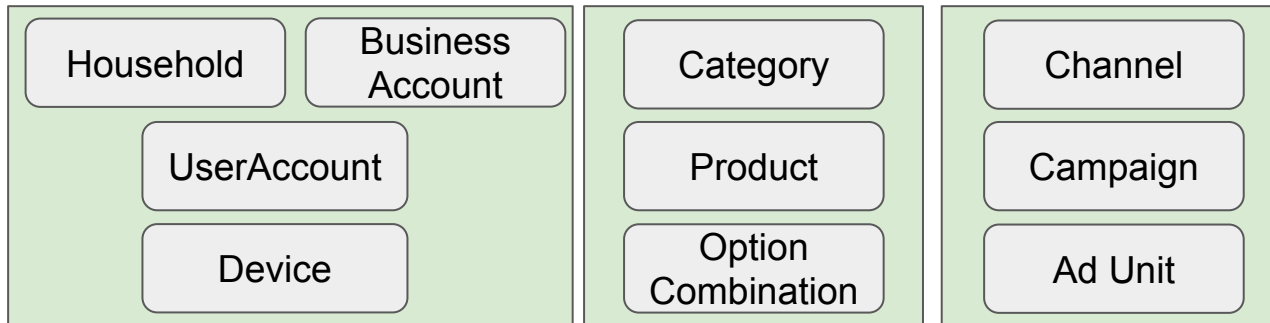
- **Library:** Recipes to make features
- **Engine:** Turns recipes from the library into data
- **Archival:** Store feature values for training and offline eval
- **Cache:** Low latency retrieval.

Features can also be reused across entities

The feature compute steps are independent of entity map used



There are many entity hierarchies to which this can apply:



**All identifiers and numbers for illustrative purposes only and do not reflect real data.



Succinctness, automated discovery, and lifecycle management are also important drivers of reuse



300 lines of config can produce thousands of features, compare that to 50 lines of sql for 1 feature.

It's actually possible to review the code for features you use.

Not that you should actually have to.

We use scalable feature selection methods to help find predictive signals from the library of thousands.

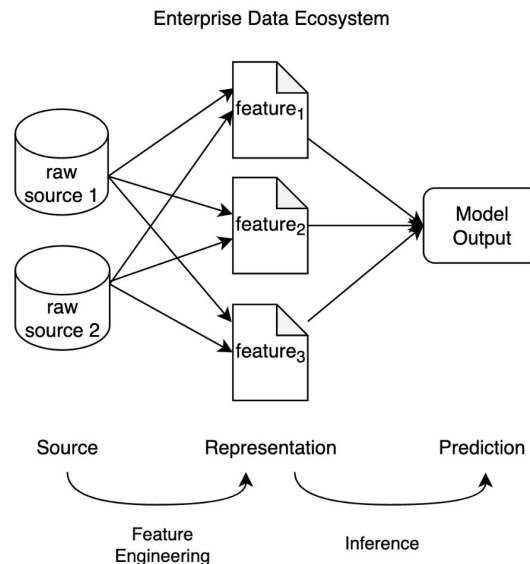
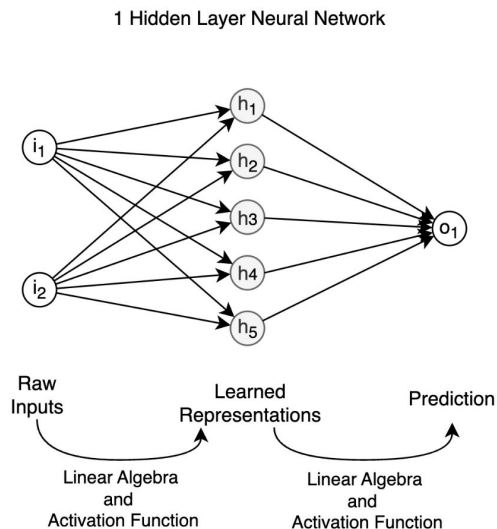


But as helpful as all that is, it's expectations on maintenance of the features that most influence users commitments to reuse features.

Part 3: Intuition

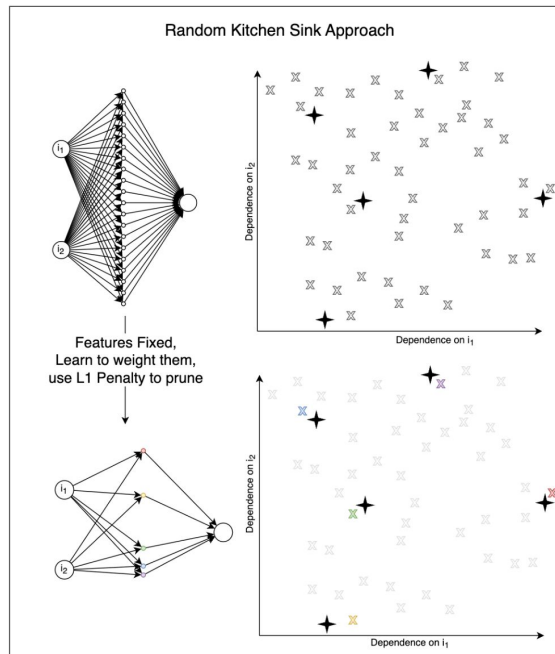
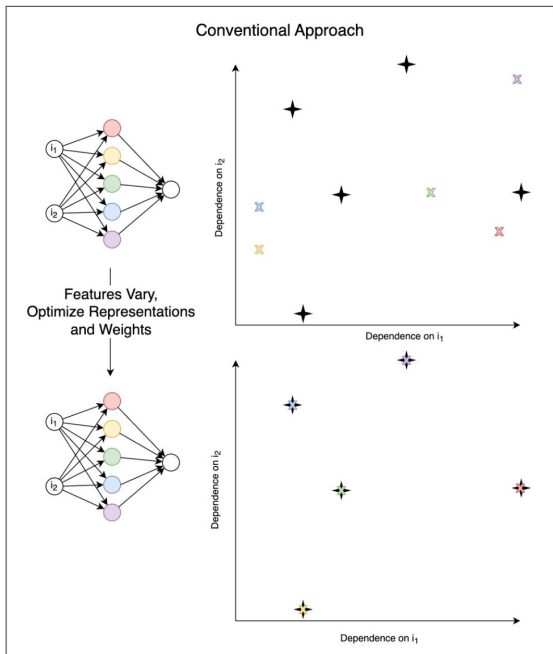
Not-So-Random Kitchen Sinks

To provide intuition for why Mercury is effective in practice we will rely on a metaphor between a neural network and an enterprise data ecosystem.



“Random Kitchen Sinks” provide intuition for the effectiveness of programmatic feature definitions.

Weighted Sums of Random Kitchen Sinks: Replacing minimization with randomization in Learning. A. Rahimi and B. Recht. NIPS, 2008.

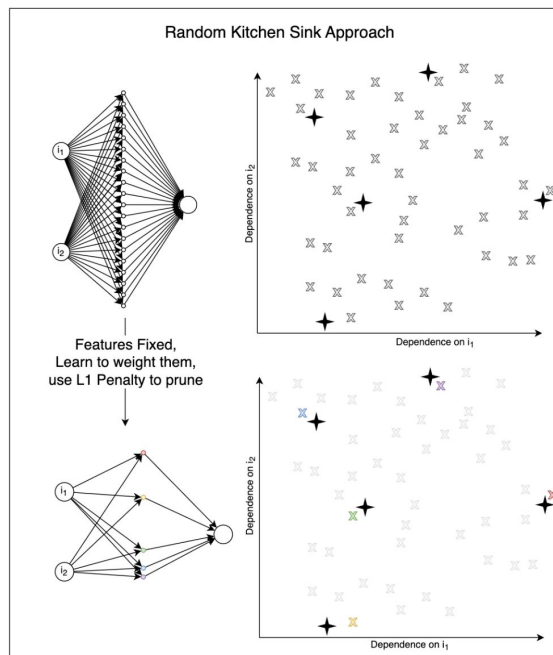
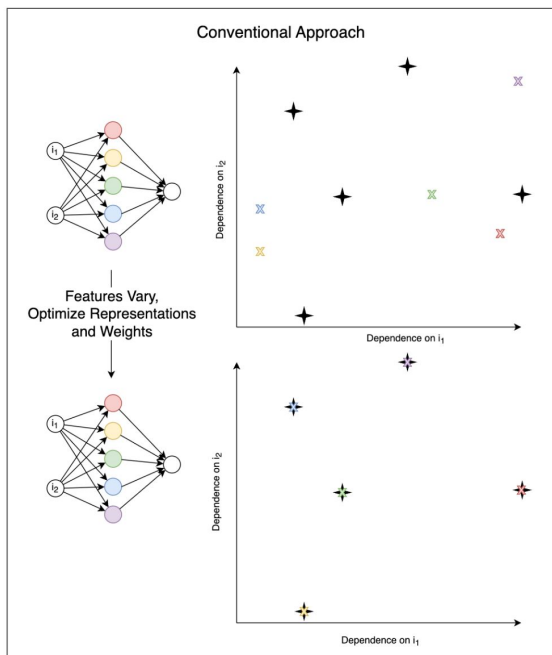


The paper shows random kitchen sinks

1. Approximate best case performance up to proven bound
2. Train up to ~10x faster

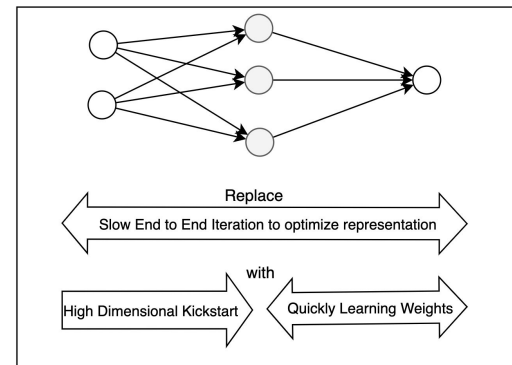
“Random Kitchen Sinks” provide intuition for the effectiveness of programmatic feature definitions.

Weighted Sums of Random Kitchen Sinks: Replacing minimization with randomization in Learning. A. Rahimi and B. Recht. NIPS, 2008.

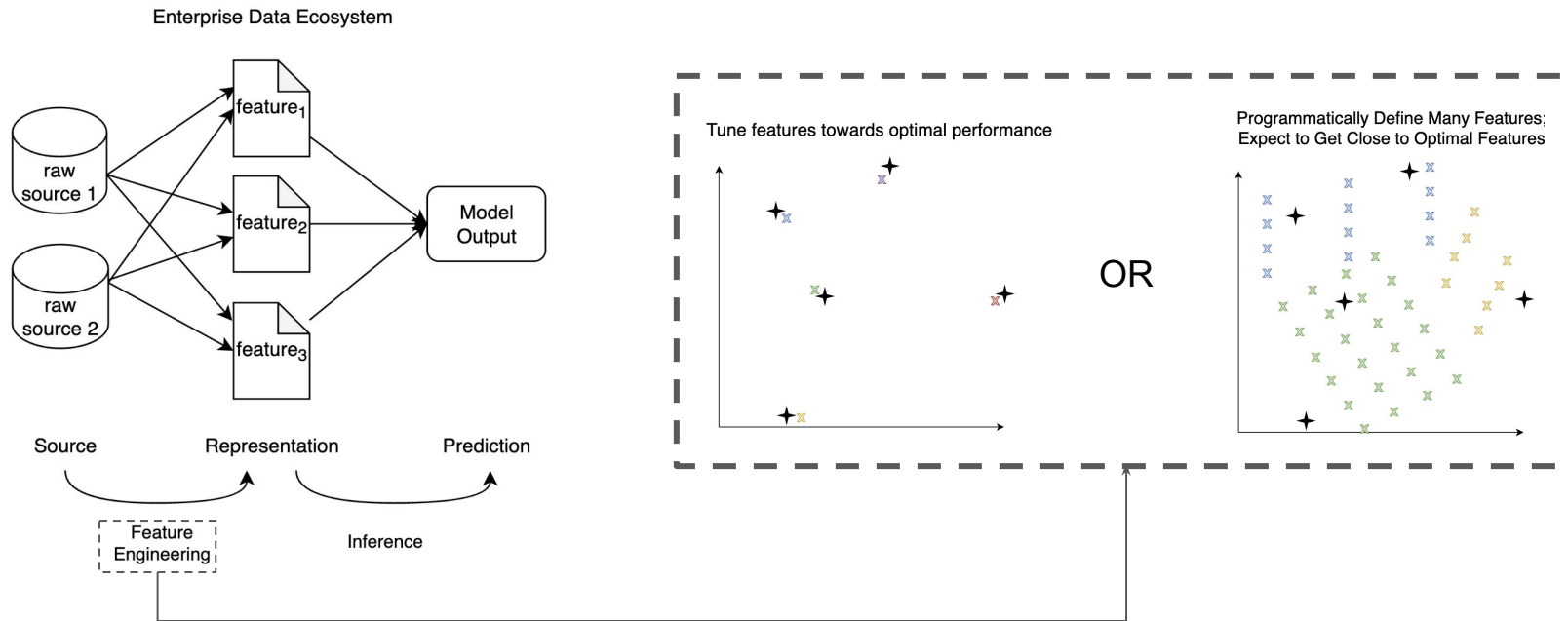


The paper shows random kitchen sinks

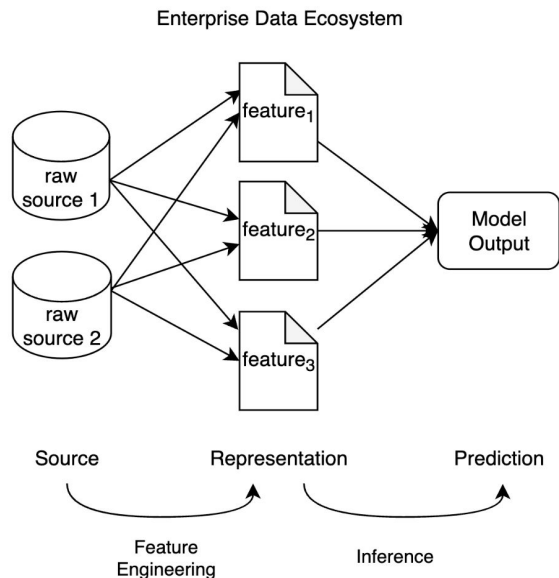
1. Approximate best case performance up to proven bound
2. Train up to ~10x faster



Programmatic feature definitions can be thought of as a deterministic analog of the random kitchen sink in the context of an enterprise data warehouse



Currently excited to explore this approach as a complement to AutoML



“AutoML” frameworks have been powerful ways of helping ML practitioners be more productive.

They provide tools to accelerate:

- Architecture / Hyperparameter Tuning
- Ensembling
- Feature Postprocessing
 - Winsorization
 - Imputation
 - Dimensionality Reduction
 - Polynomial Transformations

... but need to be given a dataframe.

Thank you for your attention. Questions?

And thank you to Mercury contributors:

- Jacob Baron
- Jesse Fredrickson
- Patrick Andruszkiewicz
- Vipul Dalsukrai
- Kurt Zimmer
- Minjoo Kim
- Masoum Mosmer