

Factorized Machine Learning: Paths and Roadblocks

Arun Kumar



UC San Diego
JACOBS SCHOOL OF ENGINEERING
Computer Science and Engineering

UC San Diego
HALICIOĞLU DATA SCIENCE INSTITUTE

Factorized Databases Workshop

Aug 3, 2022

Golden Age of ML Analytics



FACEBOOK



Golden Age of ML Analytics



FACEBOOK



Golden Age of ML Analytics



FACEBOOK



MAYO CLINIC

Healthcare



Insurance



Retail



Sciences

\$ 38 billion
in 2019*



\$ 500 billion
by 2025*

*International Data Corporation

Golden Age of ML Analytics



FACEBOOK



MAYO CLINIC

Healthcare



Insurance



Retail



Sciences

\$ 38 billion
in 2019*



\$ 500 billion
by 2025*



*International Data Corporation

Golden Age of ML Analytics



FACEBOOK



MAYO CLINIC

Healthcare



Insurance



Retail



Sciences

\$ 38 billion
in 2019*



\$ 500 billion
by 2025*



*International Data Corporation

Golden Age of ML Analytics



FACEBOOK



MAYO CLINIC

Healthcare



Insurance



Retail



Sciences

\$ 38 billion
in 2019*



\$ 500 billion
by 2025*

Still, fundamental efficiency and usability bottlenecks in the end-to-end process of building and deploying ML applications

*International Data Corporation

My Research

New abstractions, algorithms, and software systems
to “*democratize*” ML-based data analytics from
a data management/systems standpoint

My Research

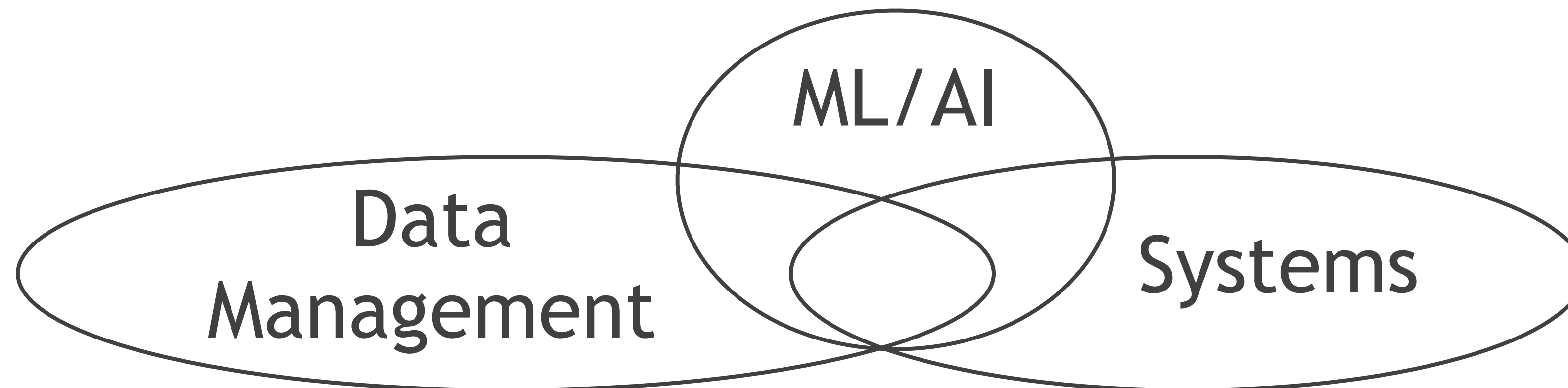
New abstractions, algorithms, and software systems
to “*democratize*” ML-based data analytics from
a data management/systems standpoint

Democratization = System Efficiency
(Reduce costs) + Human Efficiency
(Improve productivity)

My Research

New abstractions, algorithms, and software systems
to “*democratize*” ML-based data analytics from
a data management/systems standpoint

Democratization = System Efficiency (Reduce costs) + Human Efficiency (Improve productivity)



My Research

New abstractions, algorithms, and software systems
to “*democratize*” ML-based data analytics from
a data management/systems standpoint

Democratization = System Efficiency (Reduce costs) + Human Efficiency (Improve productivity)

Practical and scalable data systems for ML analytics

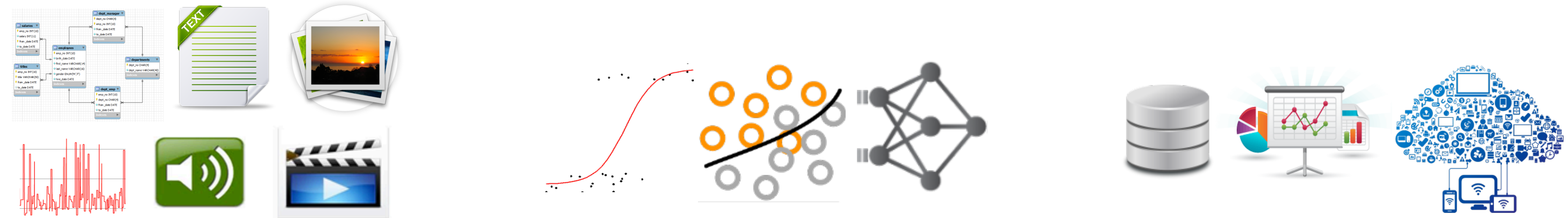
Inspired by *relational database systems* principles

Exploit insights from *learning theory* and *optimization theory*

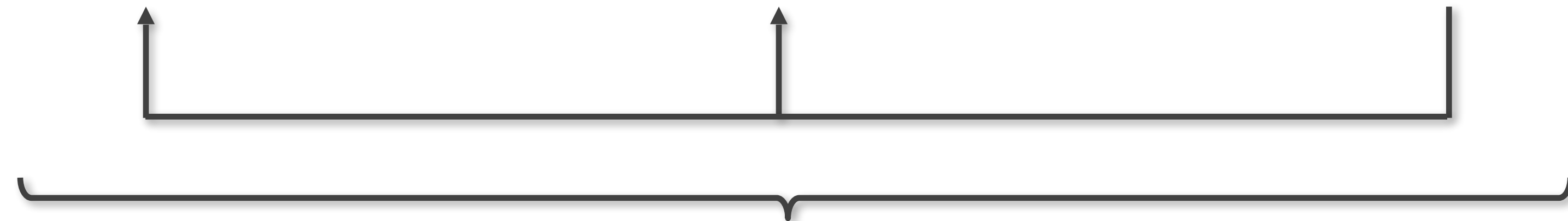
End-to-End ML Application Lifecycle



Data Scientist/
ML Engineer



Source → Build → Deploy



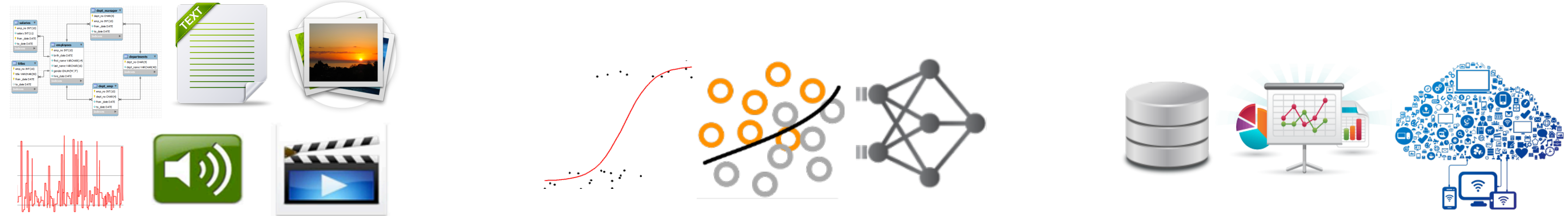
Data + ML Systems Implementations

<https://ADALabUCSD.github.io>

End-to-End ML Application Lifecycle



Data Scientist/
ML Engineer



Source → Build → Deploy

Data + ML Systems Implementations

Research
Approach :

Abstract
key steps

+

Formalize
computation

+

Automate
grunt work

+





Optimize
execution

<https://ADALabUCSD.github.io>

Outline

- Introducing ML over Joins
- Orion: Factorized ML
- Morpheus and Extensions
- Roadblocks and Musings

Outline

- 4m  Introducing ML over Joins
- 4m  Orion: Factorized ML
- 10m  Morpheus and Extensions
- 4m  Roadblocks and Musings

ML after Joins: The Problem

ML after Joins: The Problem

A fundamental bottleneck in feature engineering on structured data:

Many datasets
are multi-table



ML toolkits assume
single-table inputs

ML after Joins: The Problem

A fundamental bottleneck in feature engineering on structured data:

Many datasets
are multi-table



ML toolkits assume
single-table inputs



Materialize
join output

ML after Joins: The Problem

A fundamental bottleneck in feature engineering on structured data:

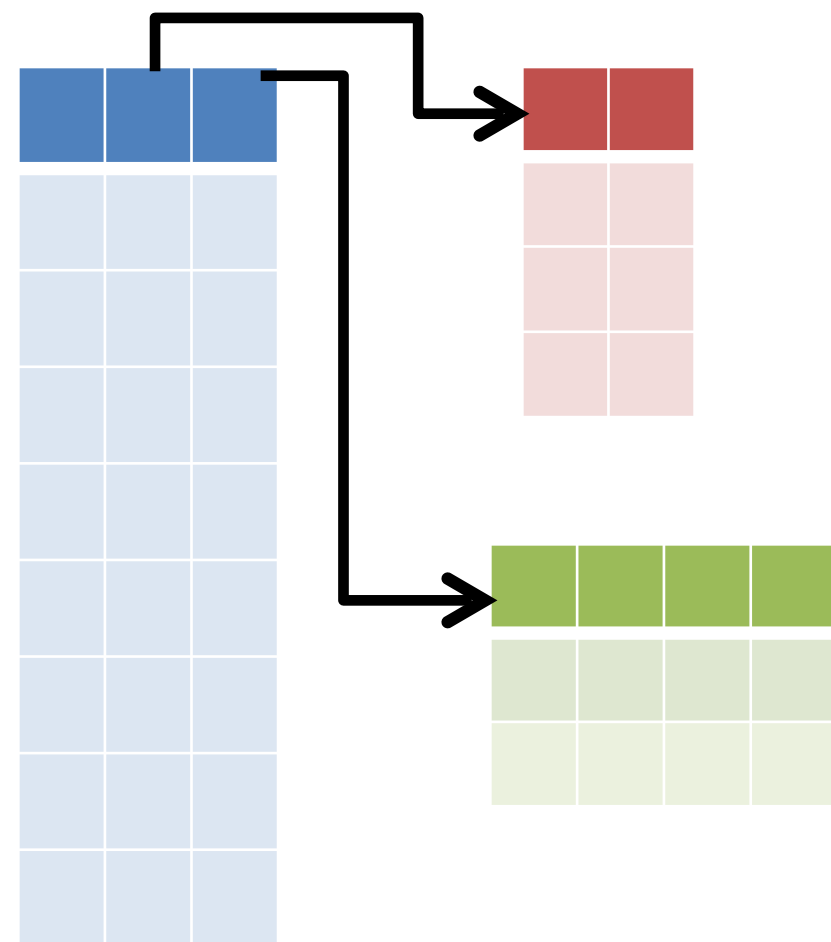
Many datasets are multi-table



ML toolkits assume single-table inputs



Materialize
join output



ML after Joins: The Problem

A fundamental bottleneck in feature engineering on structured data:

Many datasets are multi-table



ML toolkits assume single-table inputs



Materialize join output



ML after Joins: The Problem

A fundamental bottleneck in feature engineering on structured data:

Many datasets are multi-table



ML toolkits assume single-table inputs



Materialize join output



ML after Joins: The Problem

A fundamental bottleneck in feature engineering on structured data:

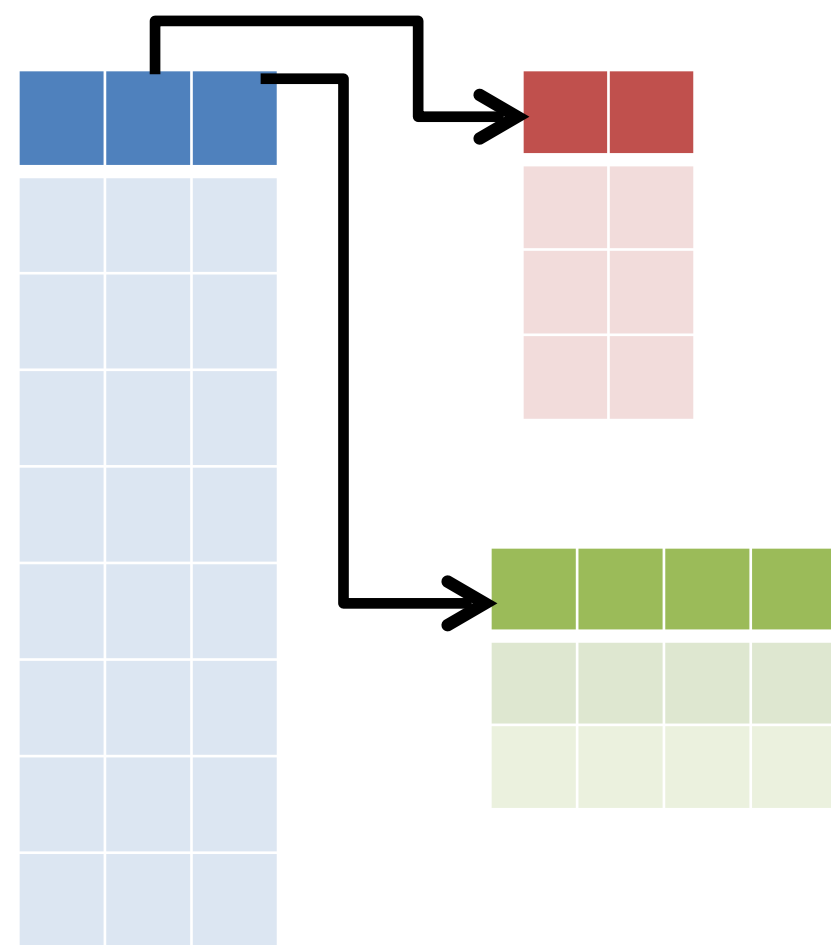
Many datasets are multi-table



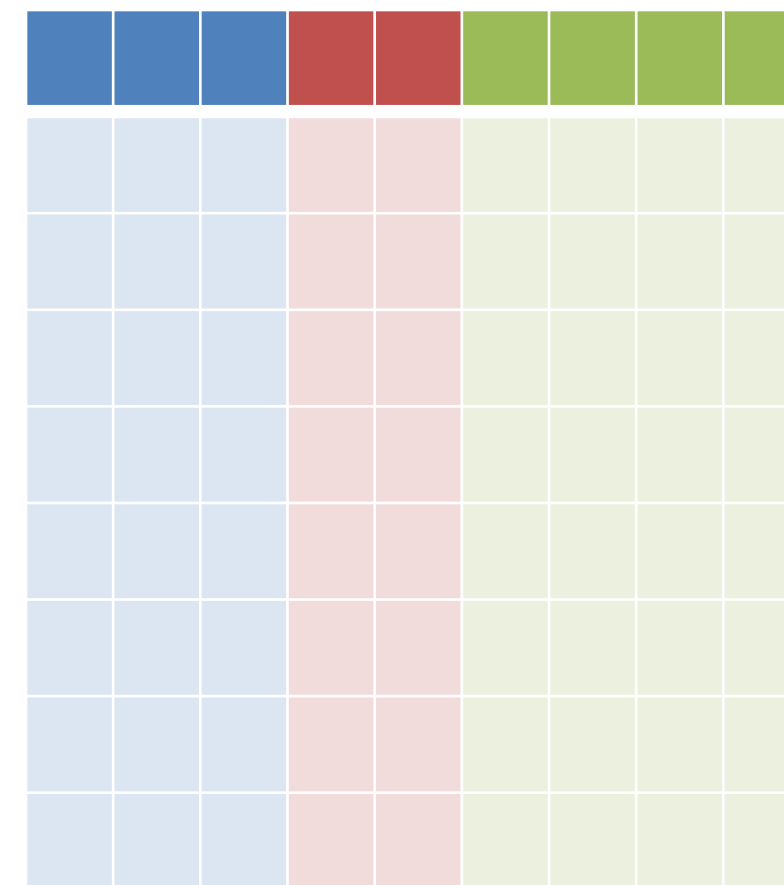
ML toolkits assume single-table inputs



Materialize join output



Key-Foreign
Key (KFK) Joins



✘ System efficiency



ML after Joins: The Problem

A fundamental bottleneck in feature engineering on structured data:

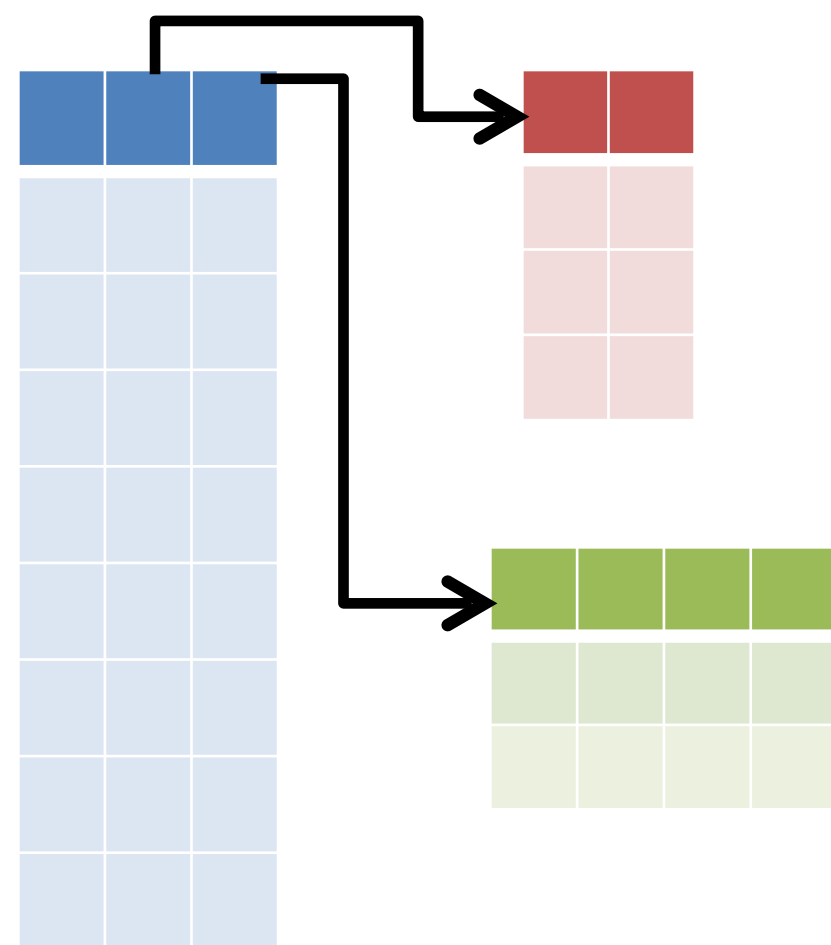
Many datasets are multi-table



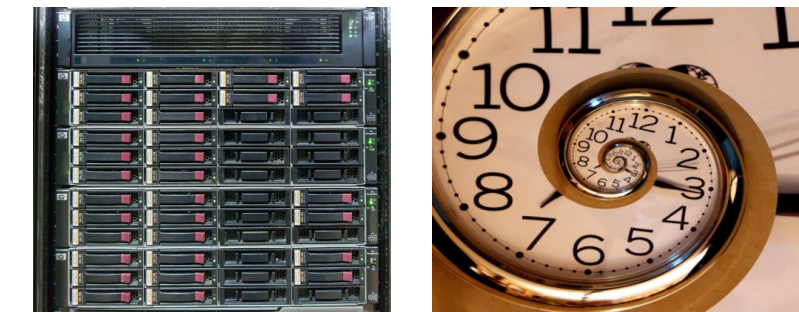
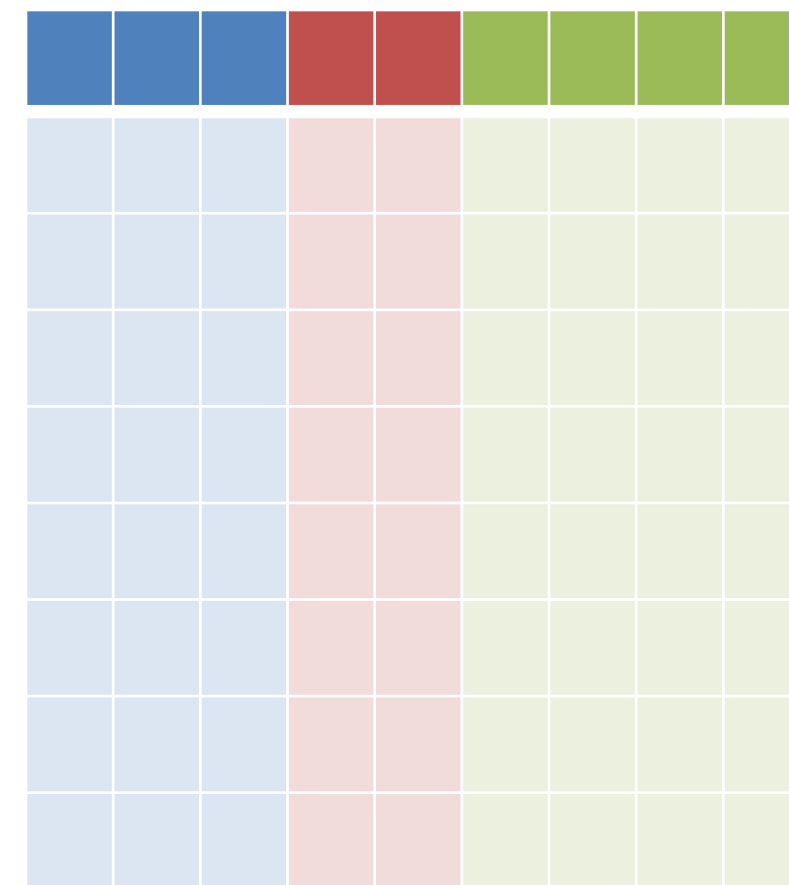
ML toolkits assume single-table inputs



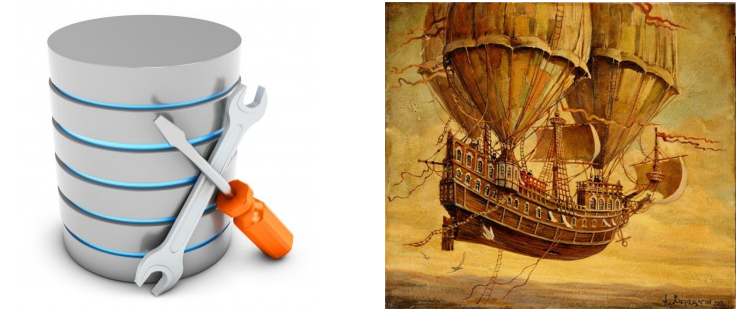
Materialize join output



Key-Foreign
Key (KFK) Joins



✘ System efficiency



✘ Human efficiency



ML after Joins: The Problem

A fundamental bottleneck in feature engineering on structured data:

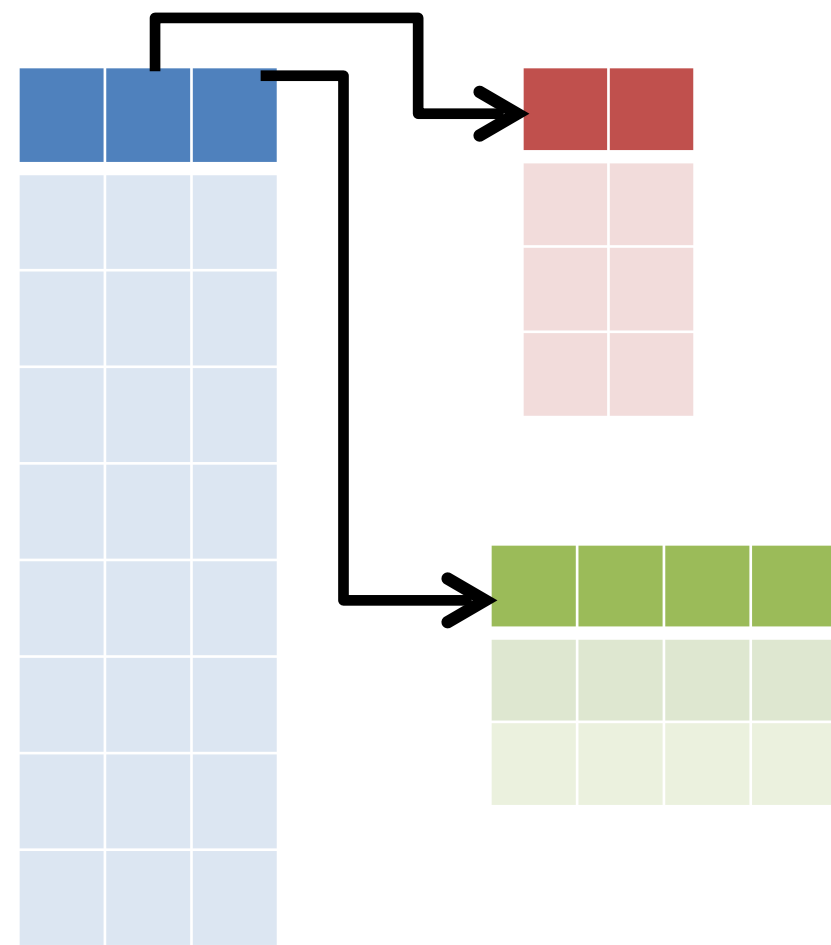
Many datasets are multi-table



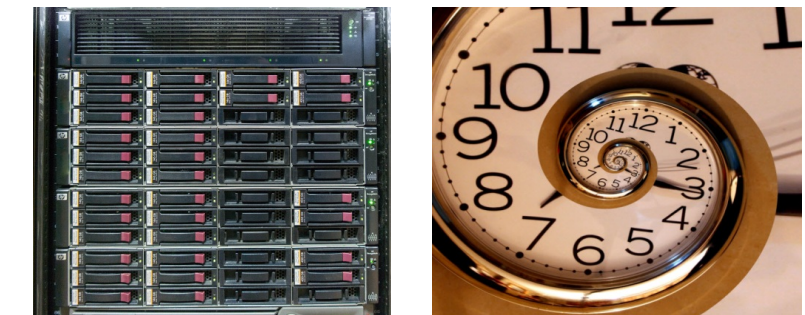
ML toolkits assume single-table inputs



Materialize join output



Key-Foreign
Key (KFK) Joins



✘ System efficiency



✘ Human efficiency



ML over Joins: Overview



✘ System efficiency

✘ Human efficiency

ML over Joins: Overview

Avoid Joins Physically

ORION, MORPHEUS



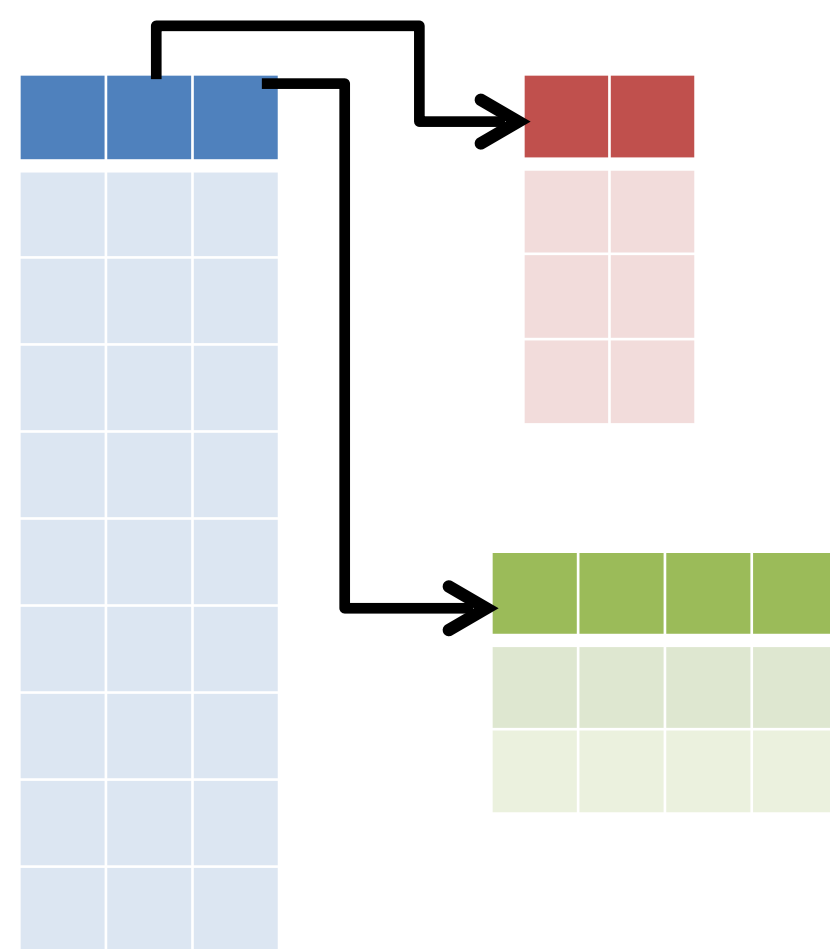
✘ System efficiency

✘ Human efficiency

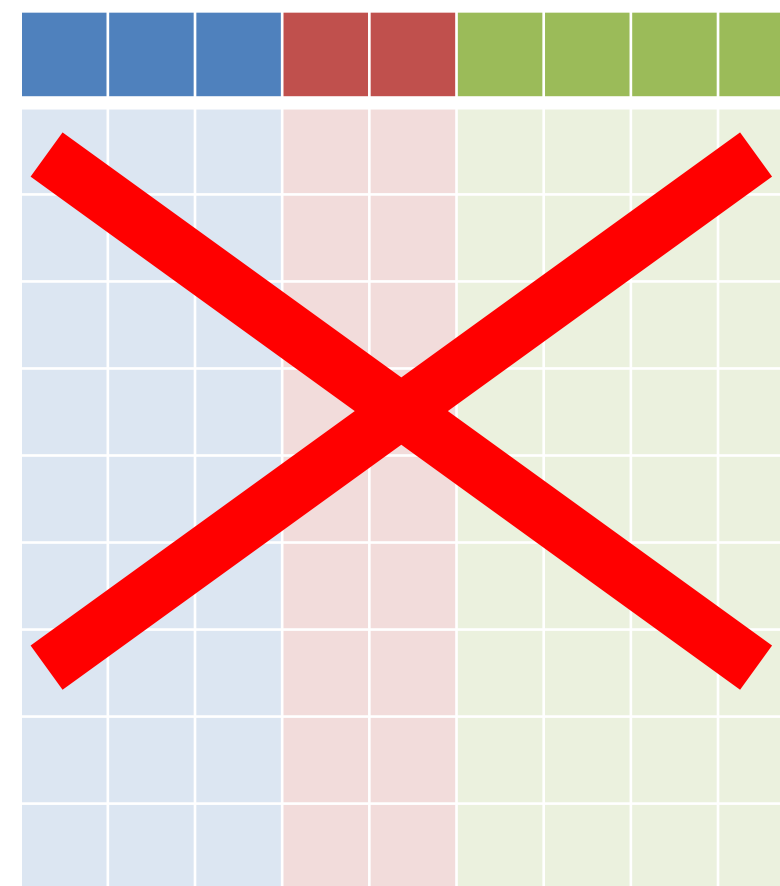
ML over Joins: Overview

Avoid Joins Physically

ORION, MORPHEUS



Key-Foreign
Key (KFK) Joins



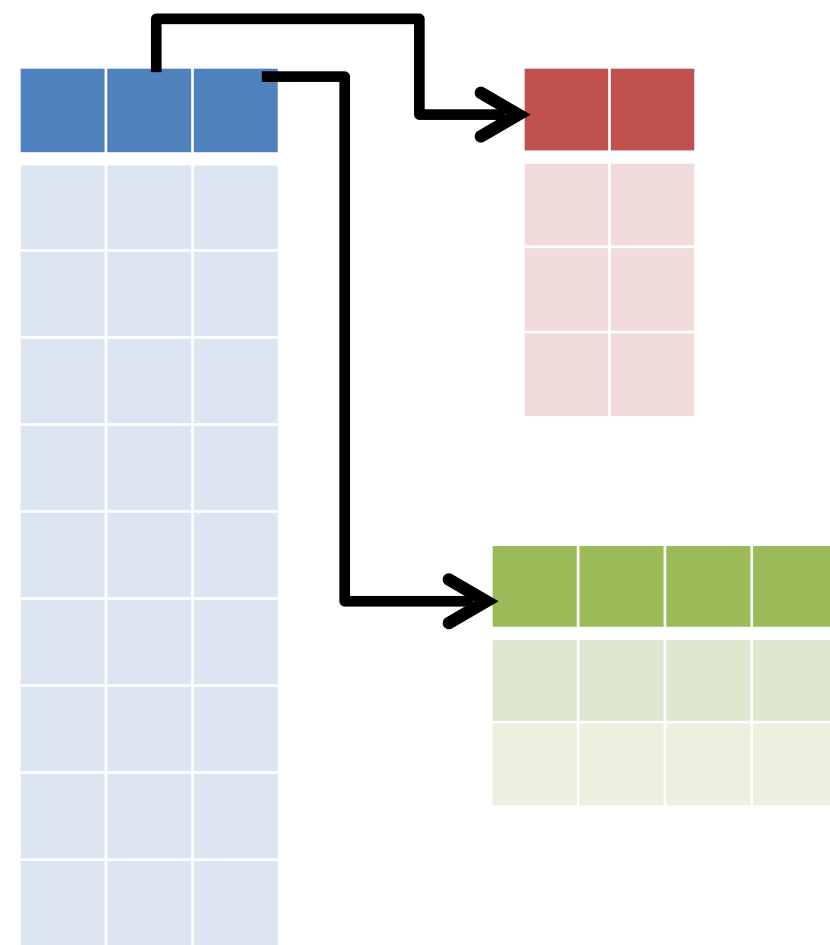
✘ System efficiency

✘ Human efficiency

ML over Joins: Overview

Avoid Joins Physically

ORION, MORPHEUS



✘ System efficiency

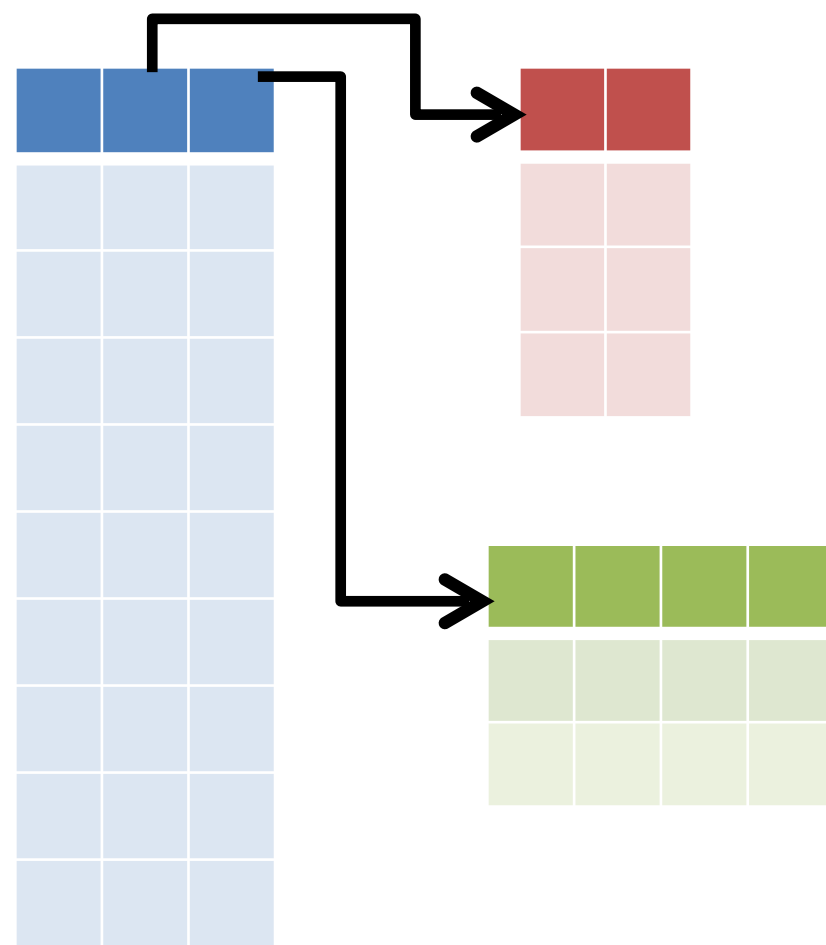
✘ Human efficiency

ML over Joins: Overview

Avoid Joins Physically

ORION, MORPHEUS

*Runs faster,
same accuracy*



✘ System efficiency

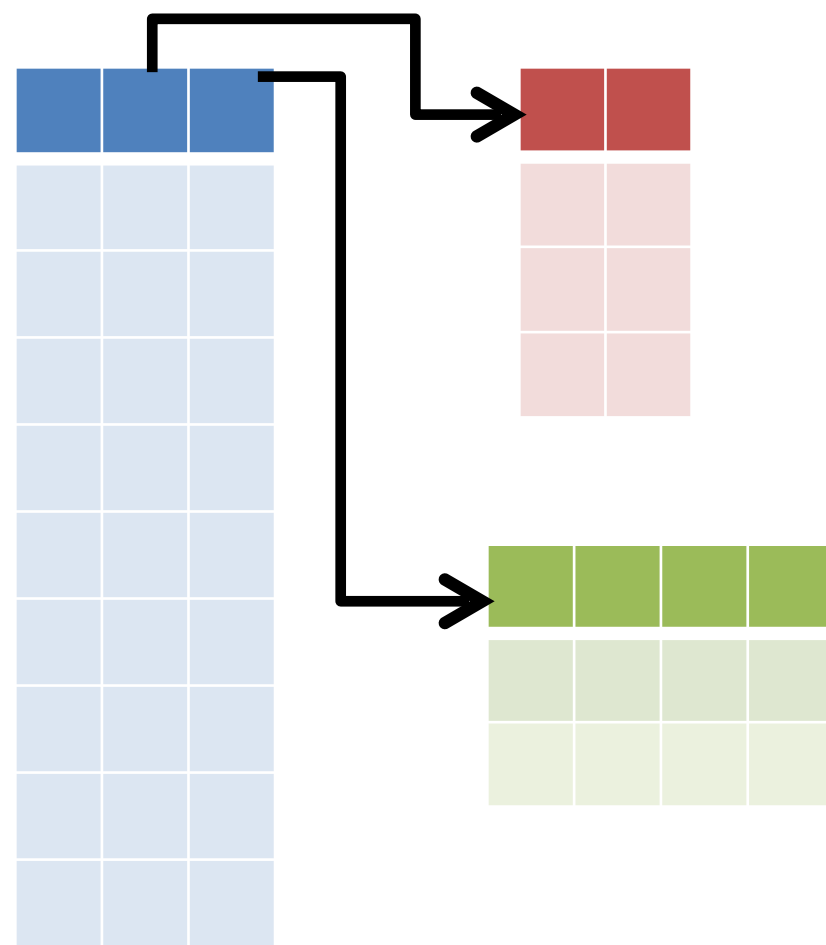
✘ Human efficiency

ML over Joins: Overview

Avoid Joins Physically

ORION, MORPHEUS

*Runs faster,
same accuracy*



✘ System efficiency

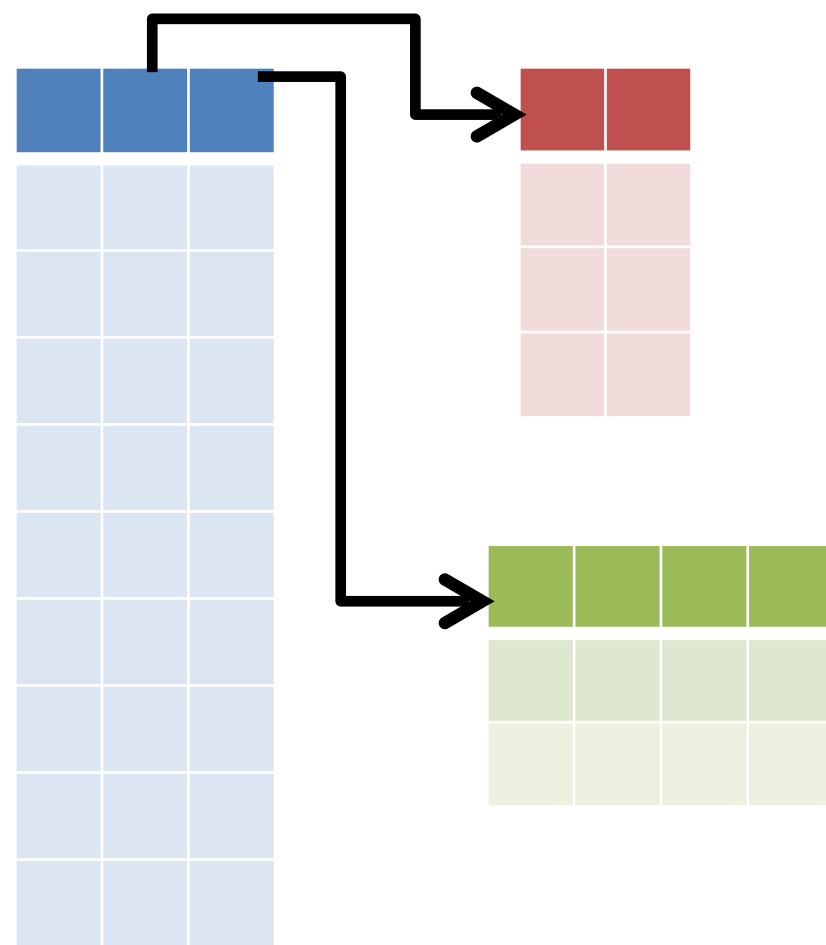
+ Human efficiency

ML over Joins: Overview

Avoid Joins Physically

ORION, MORPHEUS

*Runs faster,
same accuracy*



+ System efficiency

+ Human efficiency

ML over Joins: Overview

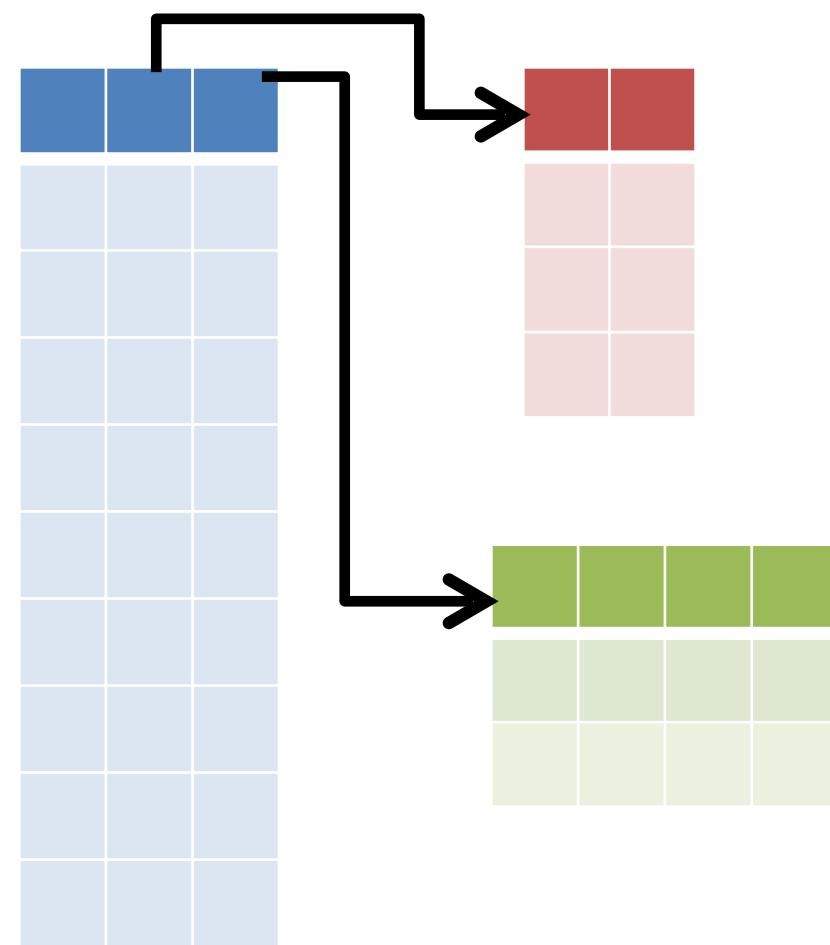
Avoid Joins Physically

ORION, MORPHEUS

*Runs faster,
same accuracy*

Avoid Joins Logically

HAMLET, HAMLET++



+ System efficiency

+ Human efficiency

ML over Joins: Overview

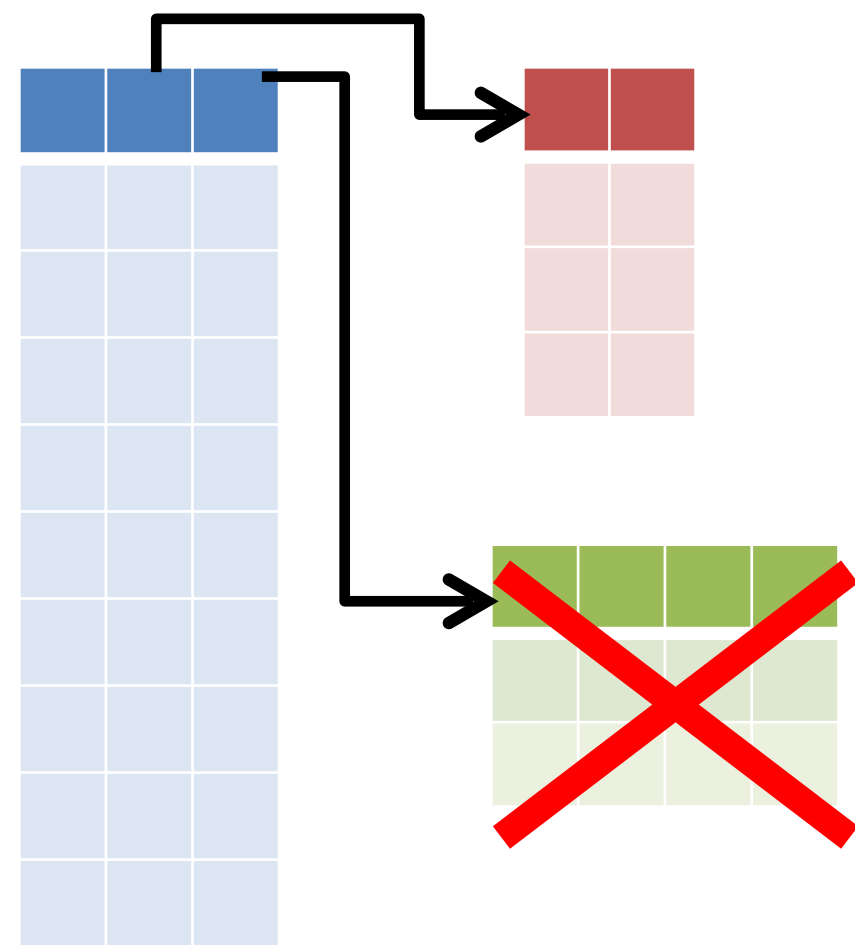
Avoid Joins Physically

ORION, MORPHEUS

*Runs faster,
same accuracy*

Avoid Joins Logically

HAMLET, HAMLET++



+ System efficiency

+ Human efficiency

ML over Joins: Overview

Avoid Joins **Physically**

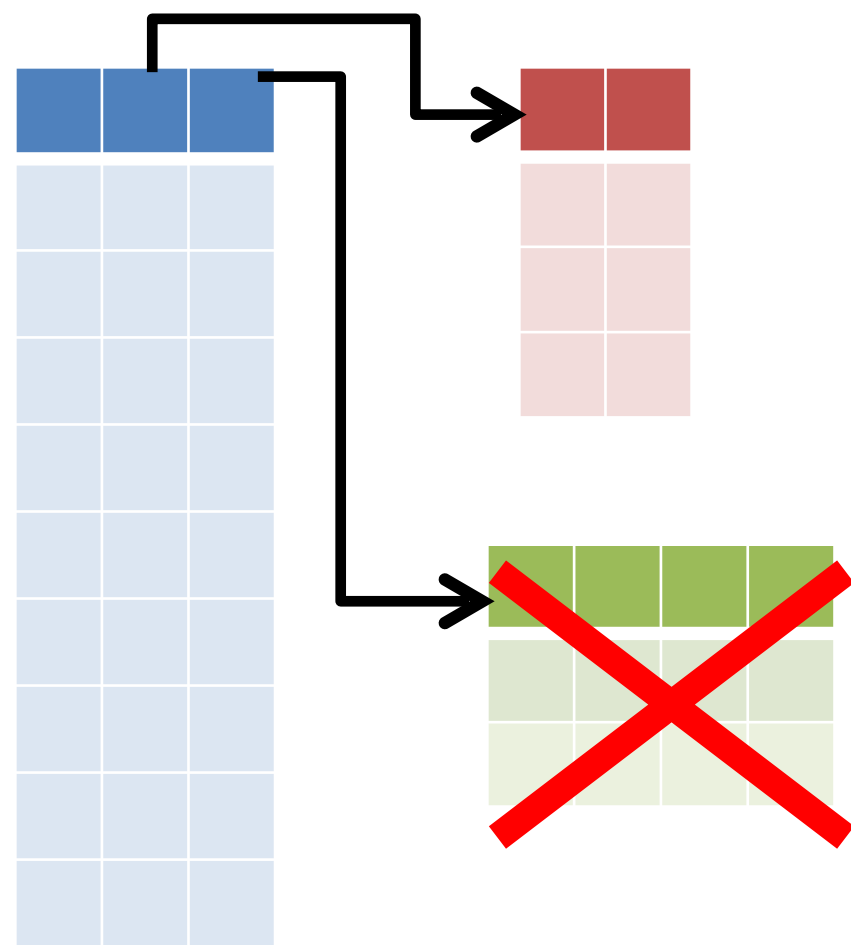
ORION, MORPHEUS

*Runs faster,
same accuracy*

Avoid Joins **Logically**

HAMLET, HAMLET++

*Even faster,
similar accuracy*



+ System efficiency

+ Human efficiency

ML over Joins: Overview

Avoid Joins Physically

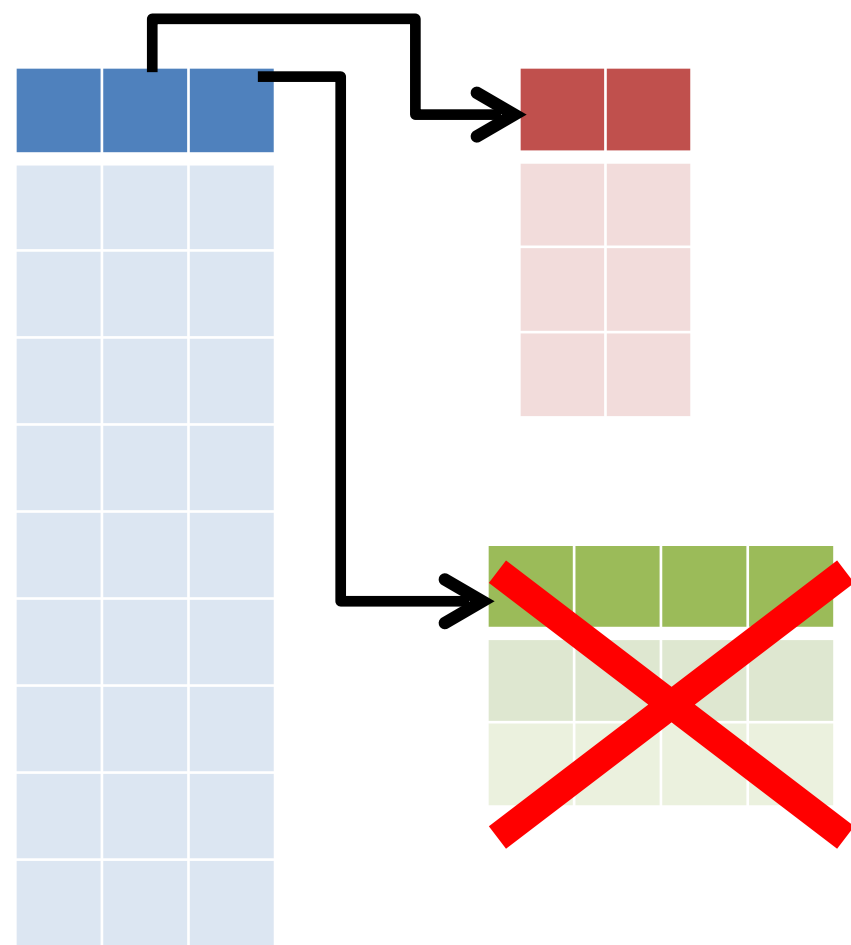
ORION, MORPHEUS

*Runs faster,
same accuracy*

Avoid Joins Logically

HAMLET, HAMLET++

*Even faster,
similar accuracy*



++ System efficiency

++ Human efficiency

ML over Joins: Overview

Avoid Joins Physically

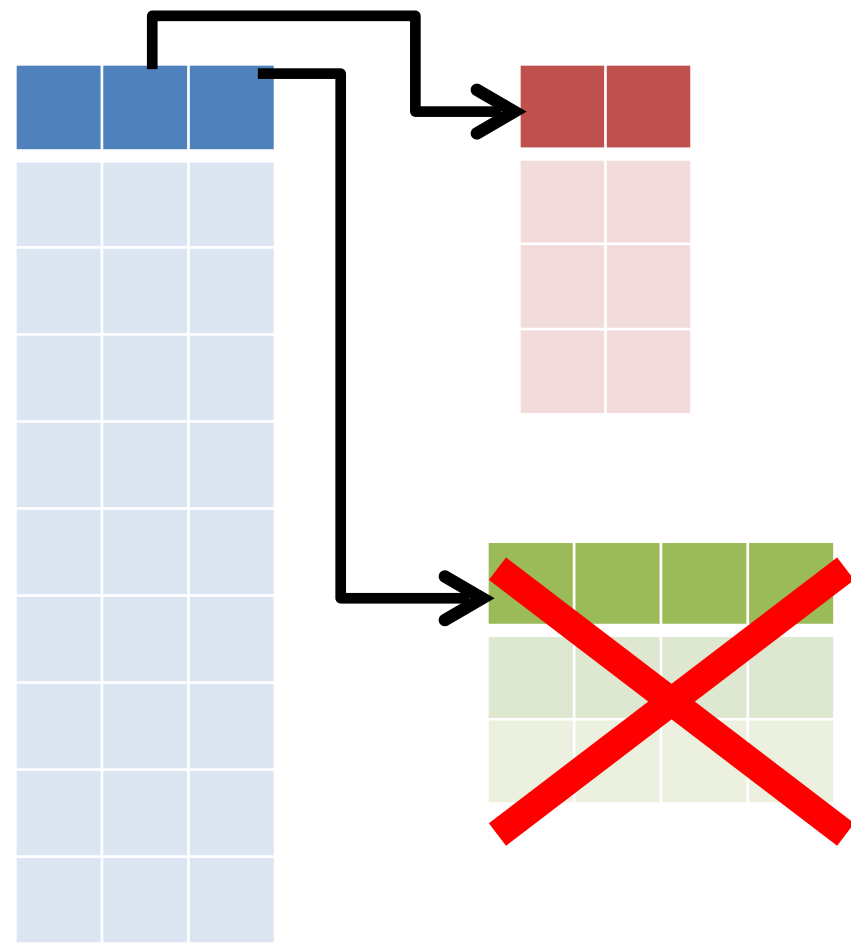
ORION, MORPHEUS

*Runs faster,
same accuracy*

Avoid Joins Logically

HAMLET, HAMLET++

*Even faster,
similar accuracy*



++ System efficiency

++ Human efficiency

Running Example for ML over Joins

ML Task: Classify if a customer will *churn* or not



Customers

Foreign Key

Employers

CID	Churn?	Gender	Age	EmpID	EmpID	State	Revenue
1	Yes	Female	33	AMZN	AMZN	WA	136b
2	No	Male	51	GOOG	GOOG	CA	89b
3	Yes	Other	46	GOOG	MSFT	WA	85b
4	No	Female	27	MSFT
...

Running Example for ML over Joins

ML Task: Classify if a customer will *churn* or not



Customers

Foreign Key

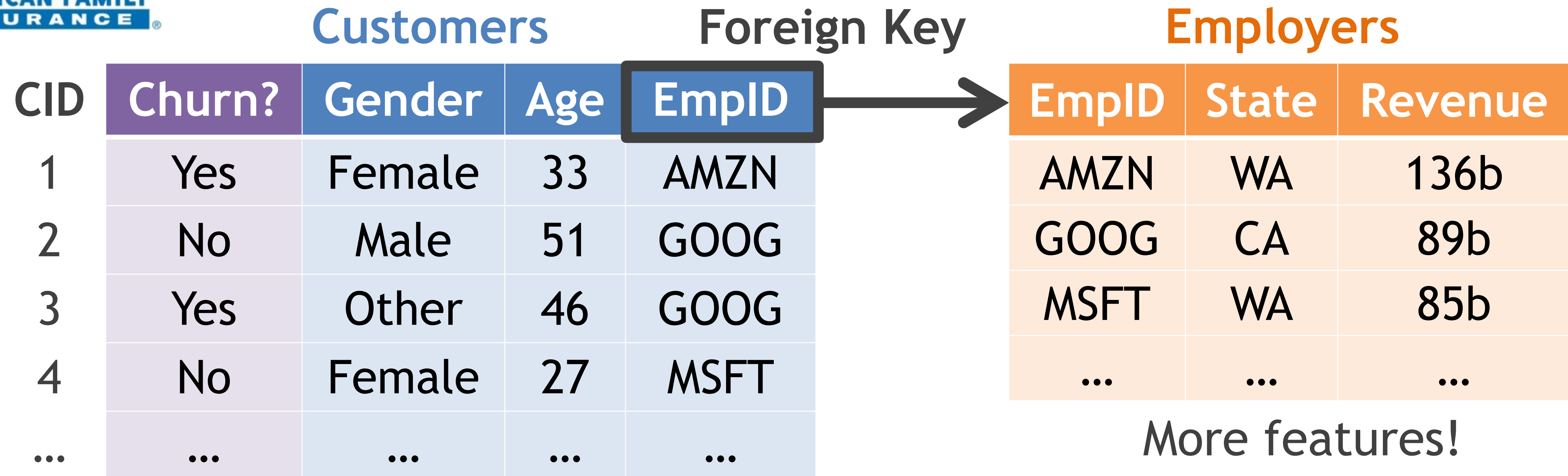
Employers

CID	Churn?	Gender	Age	EmpID	EmpID	State	Revenue
1	Yes	Female	33	AMZN	AMZN	WA	136b
2	No	Male	51	GOOG	GOOG	CA	89b
3	Yes	Other	46	GOOG	MSFT	WA	85b
4	No	Female	27	MSFT
...

More features!

Running Example for ML over Joins

ML Task: Classify if a customer will *churn* or not



More joins possible, e.g., with neighborhood data, weather data, etc.
Materializing such joins can blow up the data, even by over 10x!

Outline

- 4m ■ Introducing ML over Joins
- 4m ■ **Orion: Factorized ML**
- 10m ■ Morpheus and Extensions
- 4m ■ Roadblocks and Musings

ORION: Factorized ML

ORION: Factorized ML

Insight: Decompose ML computations and push them down through joins

ORION: Factorized ML

Insight: Decompose ML computations and push them down through joins

Focus: Generalized Linear Models (GLMs) solved using (batch) gradient descent methods

$$\mathbf{X} \equiv [\mathbf{X}_C \mathbf{X}_E]$$

ORION: Factorized ML

Insight: Decompose ML computations and push them down through joins

Focus: Generalized Linear Models (GLMs) solved using (batch) gradient descent methods

$$\nabla L(\mathbf{w}) = \sum_{i=1}^N g(\mathbf{w}^T \mathbf{x}_i, y_i) \mathbf{x}_i \quad \mathbf{X} \equiv [\mathbf{X}_C \mathbf{X}_E]$$

ORION: Factorized ML

Insight: Decompose ML computations and push them down through joins

Focus: Generalized Linear Models (GLMs) solved using (batch) gradient descent methods

$$\nabla L(\mathbf{w}) = \sum_{i=1}^N g(\mathbf{w}^T \mathbf{x}_i, y_i) \mathbf{x}_i \quad \mathbf{X} \equiv [\mathbf{X}_C \ \mathbf{X}_E]$$

ORION: Factorized ML

Insight: Decompose ML computations and push them down through joins

Focus: Generalized Linear Models (GLMs) solved using (batch) gradient descent methods

$$\nabla L(\mathbf{w}) = \sum_{i=1}^N g(\mathbf{w}^T \mathbf{x}_i, y_i) \mathbf{x}_i \quad \mathbf{X} \equiv [\mathbf{X}_C \ \mathbf{X}_E]$$

$$\mathbf{w}^T \mathbf{X} = [\mathbf{w}_C^T \ \mathbf{w}_E^T] \begin{bmatrix} \mathbf{X}_C \\ \mathbf{X}_E \end{bmatrix} = \mathbf{w}_C^T \mathbf{X}_C + \mathbf{w}_E^T \mathbf{X}_E$$

ORION: Factorized ML

Insight: Decompose ML computations and push them down through joins

Focus: Generalized Linear Models (GLMs) solved using (batch) gradient descent methods

$$\nabla L(\mathbf{w}) = \sum_{i=1}^N g(\mathbf{w}^T \mathbf{x}_i, y_i) \mathbf{x}_i \quad \mathbf{X} \equiv [\mathbf{X}_C \ \mathbf{X}_E]$$

$$\mathbf{w}^T \mathbf{x} = [\mathbf{w}_C^T \ \mathbf{w}_E^T] \begin{bmatrix} \mathbf{x}_C \\ \mathbf{x}_E \end{bmatrix} = \mathbf{w}_C^T \mathbf{x}_C + \boxed{\mathbf{w}_E^T \mathbf{x}_E}$$

ORION: Factorized ML

Insight: Decompose ML computations and push them down through joins

Focus: Generalized Linear Models (GLMs) solved using (batch) gradient descent methods

$$\nabla L(\mathbf{w}) = \sum_{i=1}^N g(\mathbf{w}^T \mathbf{x}_i, y_i) \mathbf{x}_i \quad \mathbf{X} \equiv [\mathbf{X}_C \ \mathbf{X}_E]$$

$$\mathbf{w}^T \mathbf{x} = [\mathbf{w}_C^T \ \mathbf{w}_E^T] \begin{bmatrix} \mathbf{x}_C \\ \mathbf{x}_E \end{bmatrix} = \mathbf{w}_C^T \mathbf{x}_C + \mathbf{w}_E^T \mathbf{x}_E$$

ORION: Factorized ML

Insight: Decompose ML computations and push them down through joins

Focus: Generalized Linear Models (GLMs) solved using (batch) gradient descent methods

$$\nabla L(\mathbf{w}) = \sum_{i=1}^N g(\mathbf{w}^T \mathbf{x}_i, y_i) \mathbf{x}_i \quad \mathbf{X} \equiv [\mathbf{X}_C \ \mathbf{X}_E]$$

$$\mathbf{w}^T \mathbf{X} = [\mathbf{w}_C^T \ \mathbf{w}_E^T] \begin{bmatrix} \mathbf{X}_C \\ \mathbf{X}_E \end{bmatrix} = \mathbf{w}_C^T \mathbf{X}_C + \mathbf{w}_E^T \mathbf{X}_E$$

1 full iteration requires 2 scans of Employers, 1 scan of Customers

ORION: Factorized ML

Insight: Decompose ML computations and push them down through joins

Focus: Generalized Linear Models (GLMs) solved using (batch) gradient descent methods

$$\nabla L(\mathbf{w}) = \sum_{i=1}^N g(\mathbf{w}^T \mathbf{x}_i, y_i) \mathbf{x}_i \quad \mathbf{X} \equiv [\mathbf{X}_C \ \mathbf{X}_E]$$

$$\mathbf{w}^T \mathbf{X} = [\mathbf{w}_C^T \ \mathbf{w}_E^T] \begin{bmatrix} \mathbf{X}_C \\ \mathbf{X}_E \end{bmatrix} = \mathbf{w}_C^T \mathbf{X}_C + \mathbf{w}_E^T \mathbf{X}_E$$

1 full iteration requires 2 scans of Employers, 1 scan of Customers

Challenges Tackled: Scalability; developability

ORION: Implementations

Learning Generalized Linear Models over Normalized Data. **SIGMOD 2015**
Demonstration of Santoku: Optimizing Machine Learning over Normalized Data. **VLDB 2015**

ORION: Implementations

Prototyped on PostgreSQL with UDAFs (MADlib style)

Distributed prototype with MapReduce on Hive & Spark (MLlib style)

Extended to Naive Bayes, k-means clustering, decision trees as R package

Learning Generalized Linear Models over Normalized Data. **SIGMOD 2015**

Demonstration of Santoku: Optimizing Machine Learning over Normalized Data. **VLDB 2015**

ORION: Implementations

Prototyped on PostgreSQL with UDAFs (MADlib style)

Distributed prototype with MapReduce on Hive & Spark (MLlib style)

Extended to Naive Bayes, k-means clustering, decision trees as R package

Explored for
production use cases:



(Web security)



(Retail)



(Ads)

Learning Generalized Linear Models over Normalized Data. **SIGMOD 2015**

Demonstration of Santoku: Optimizing Machine Learning over Normalized Data. **VLDB 2015**

Q: *Can we avoid manual rewriting of each ML algorithm and “automate” factorized ML on top of ML tools?*

Outline

- 4m ■ Introducing ML over Joins
- 4m ■ Orion: Factorized ML
- 10m ■ **Morpheus and Extensions**
- 4m ■ Roadblocks and Musings

MORPHEUS: Generalizing ORION

MORPHEUS: Generalizing ORION

Goal: *Automate* factorized ML to many ML algorithms in a unified way

MORPHEUS: Generalizing ORION

Goal: Automate factorized ML to many ML algorithms in a unified way

Idea: Many ML algorithms are bulk *linear algebra* (LA) programs
Create a framework for rewrite rules for LA ops

MORPHEUS: Generalizing ORION

Goal: Automate factorized ML to many ML algorithms in a unified way

Idea: Many ML algorithms are bulk *linear algebra* (LA) programs
Create a framework for rewrite rules for LA ops

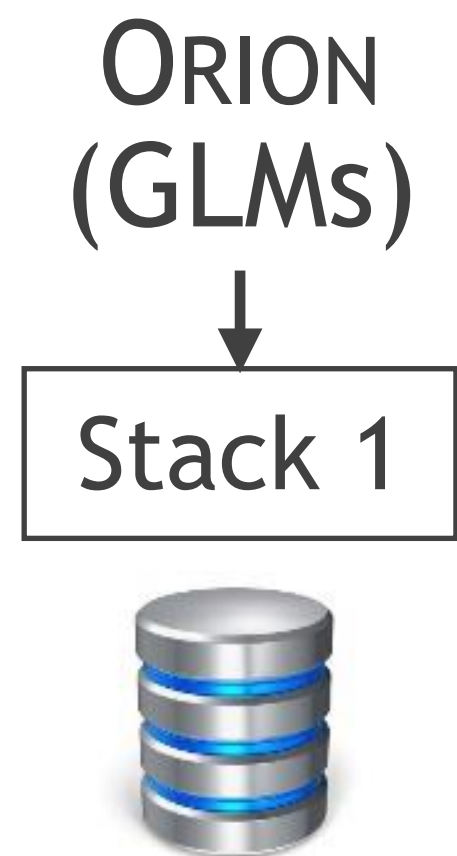
Factorized ML: Prior Work

MORPHEUS: Generalizing ORION

Goal: Automate factorized ML to many ML algorithms in a unified way

Idea: Many ML algorithms are bulk *linear algebra* (LA) programs
Create a framework for rewrite rules for LA ops

Factorized ML: Prior Work

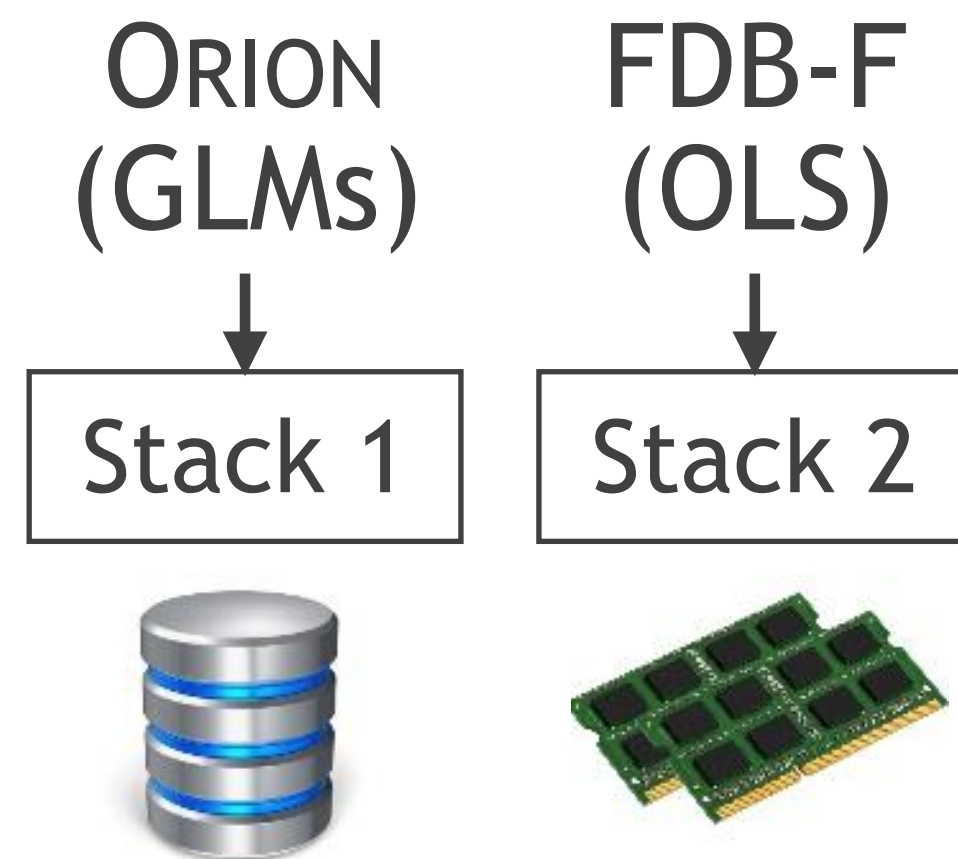


MORPHEUS: Generalizing ORION

Goal: Automate factorized ML to many ML algorithms in a unified way

Idea: Many ML algorithms are bulk *linear algebra* (LA) programs
Create a framework for rewrite rules for LA ops

Factorized ML: Prior Work

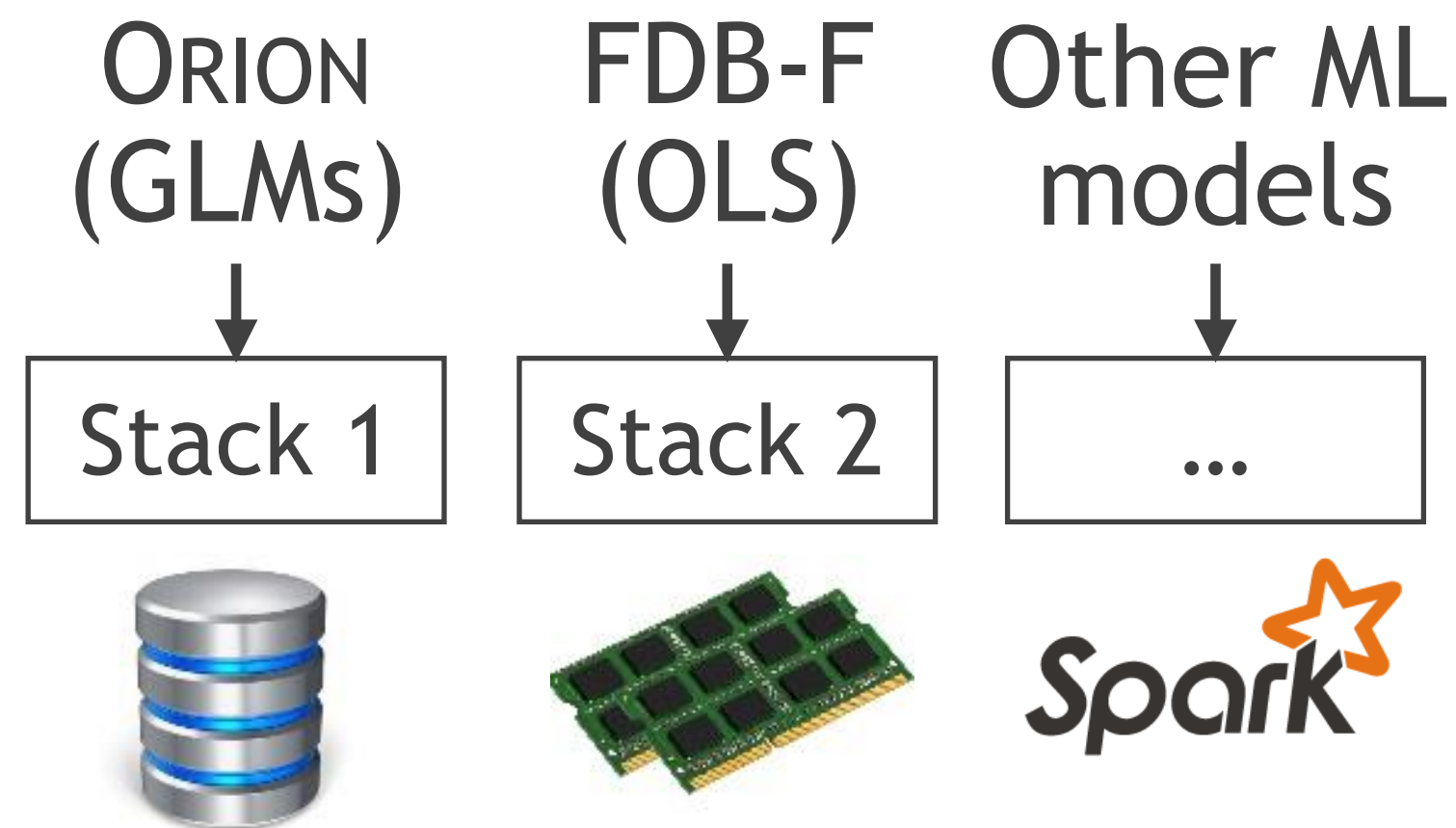


MORPHEUS: Generalizing ORION

Goal: Automate factorized ML to many ML algorithms in a unified way

Idea: Many ML algorithms are bulk *linear algebra* (LA) programs
Create a framework for rewrite rules for LA ops

Factorized ML: Prior Work

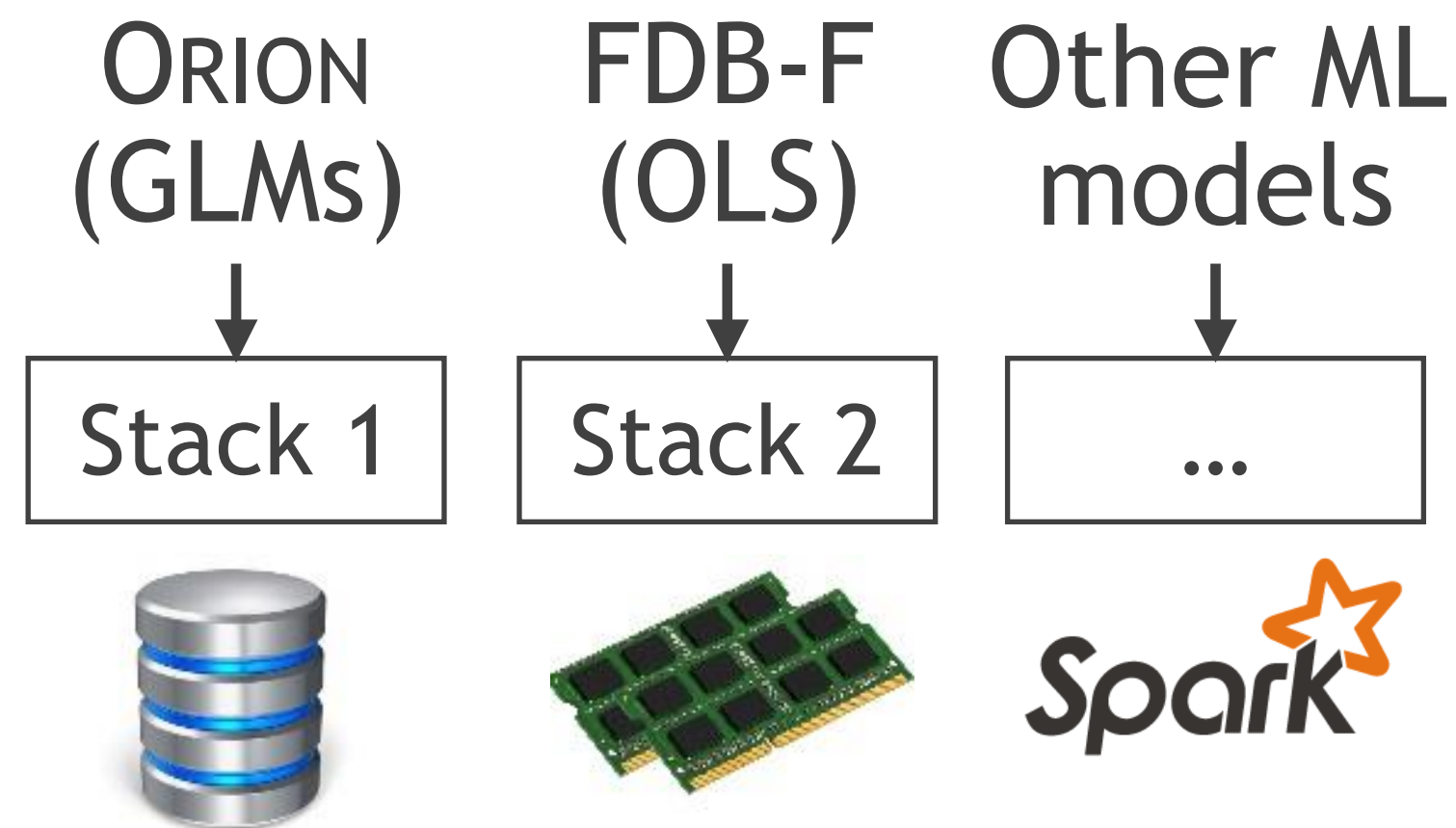


MORPHEUS: Generalizing ORION

Goal: Automate factorized ML to many ML algorithms in a unified way

Idea: Many ML algorithms are bulk *linear algebra* (LA) programs
Create a framework for rewrite rules for LA ops

Factorized ML: Prior Work



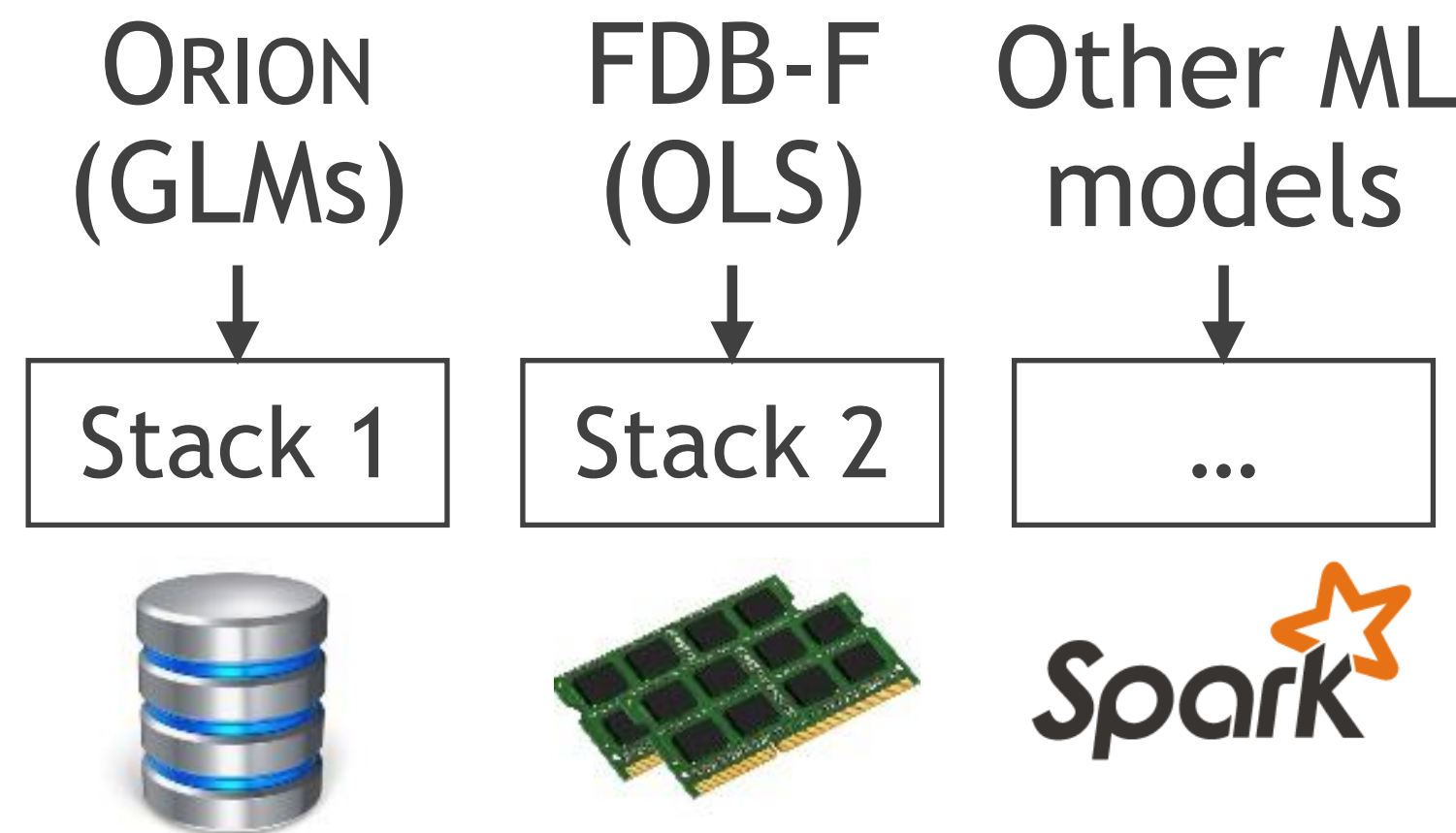
The MORPHEUS Approach

MORPHEUS: Generalizing ORION

Goal: Automate factorized ML to many ML algorithms in a unified way

Idea: Many ML algorithms are bulk *linear algebra* (LA) programs
Create a framework for rewrite rules for LA ops

Factorized ML: Prior Work



The MORPHEUS Approach

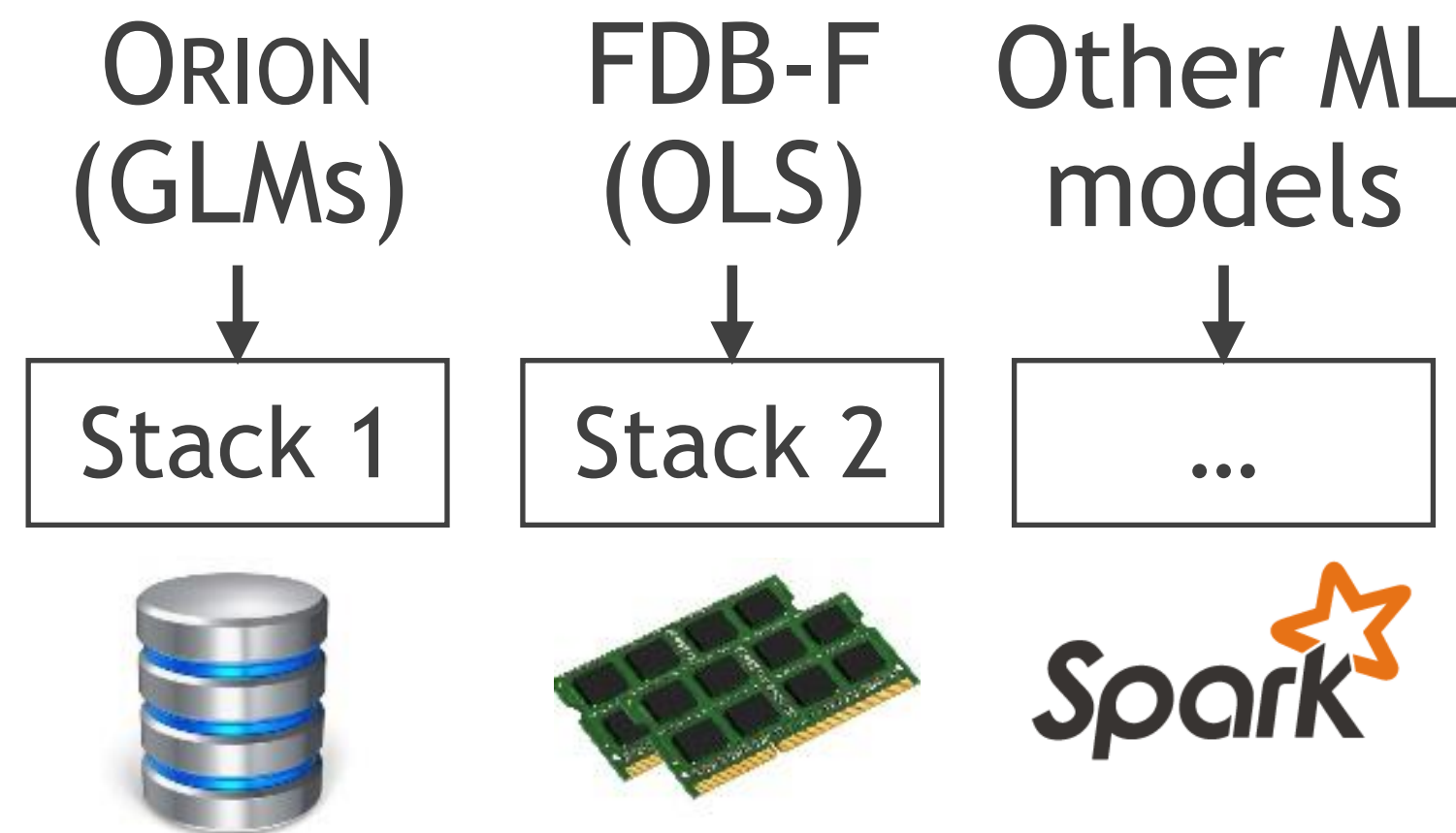
GLMs OLS K-Means NMF ...

MORPHEUS: Generalizing ORION

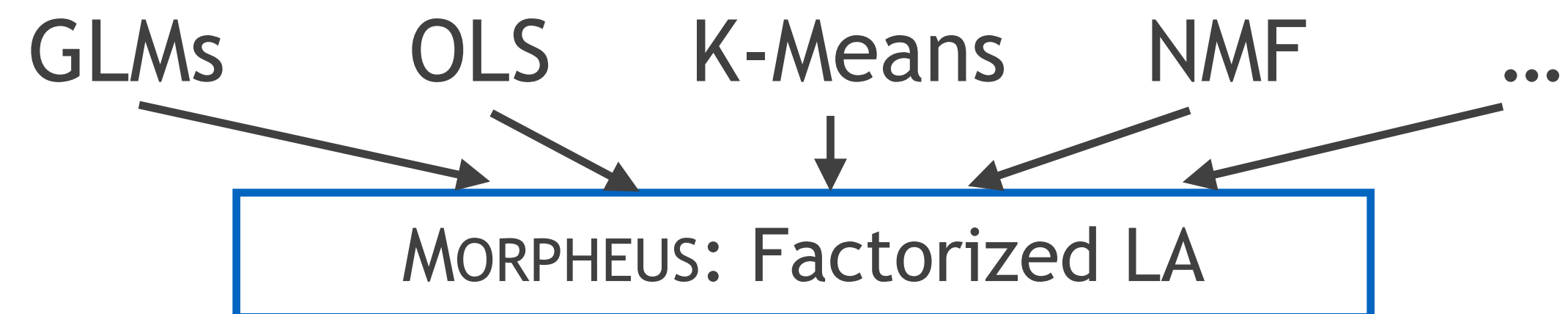
Goal: Automate factorized ML to many ML algorithms in a unified way

Idea: Many ML algorithms are bulk *linear algebra* (LA) programs
Create a framework for rewrite rules for LA ops

Factorized ML: Prior Work



The MORPHEUS Approach

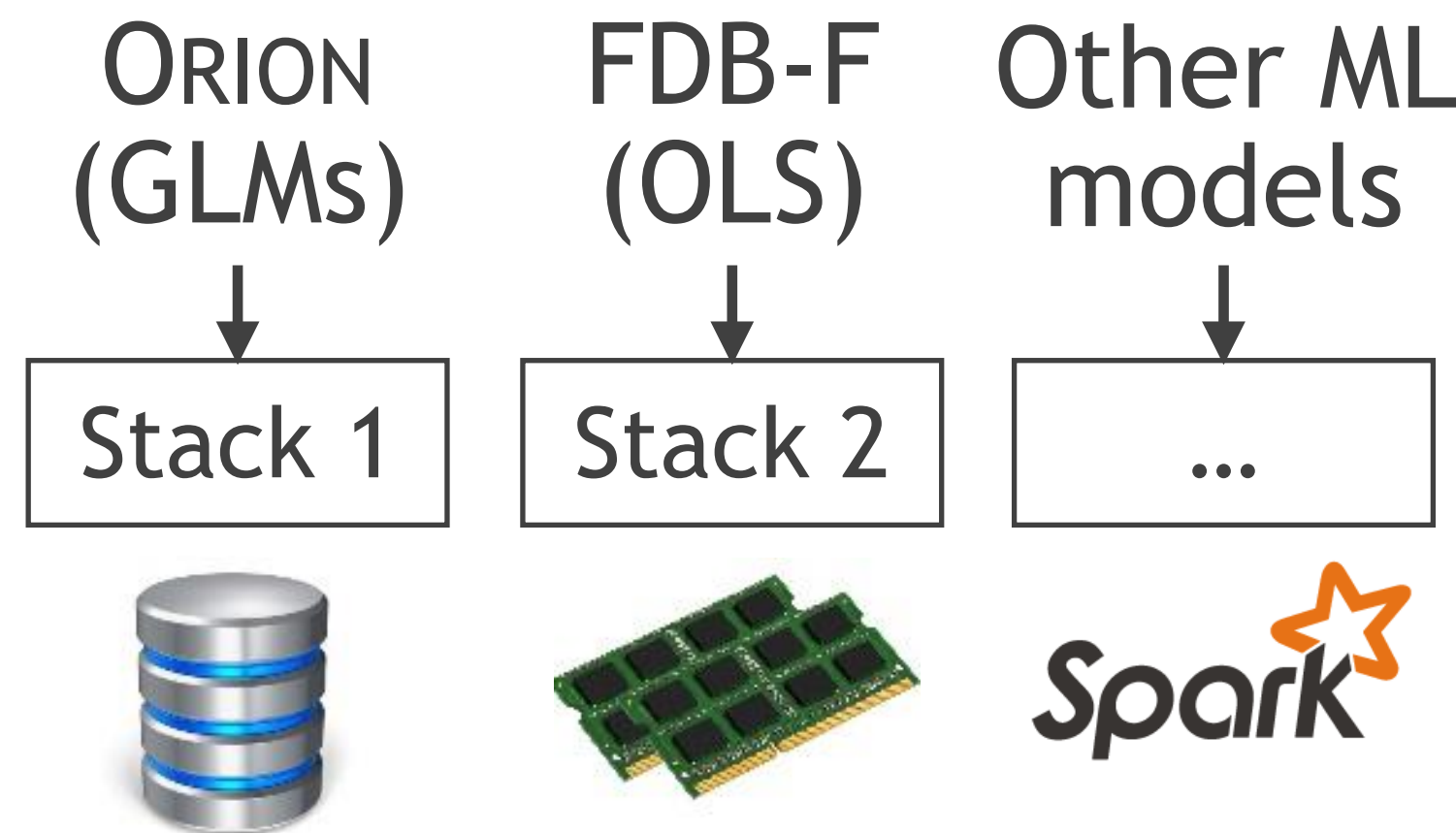


MORPHEUS: Generalizing ORION

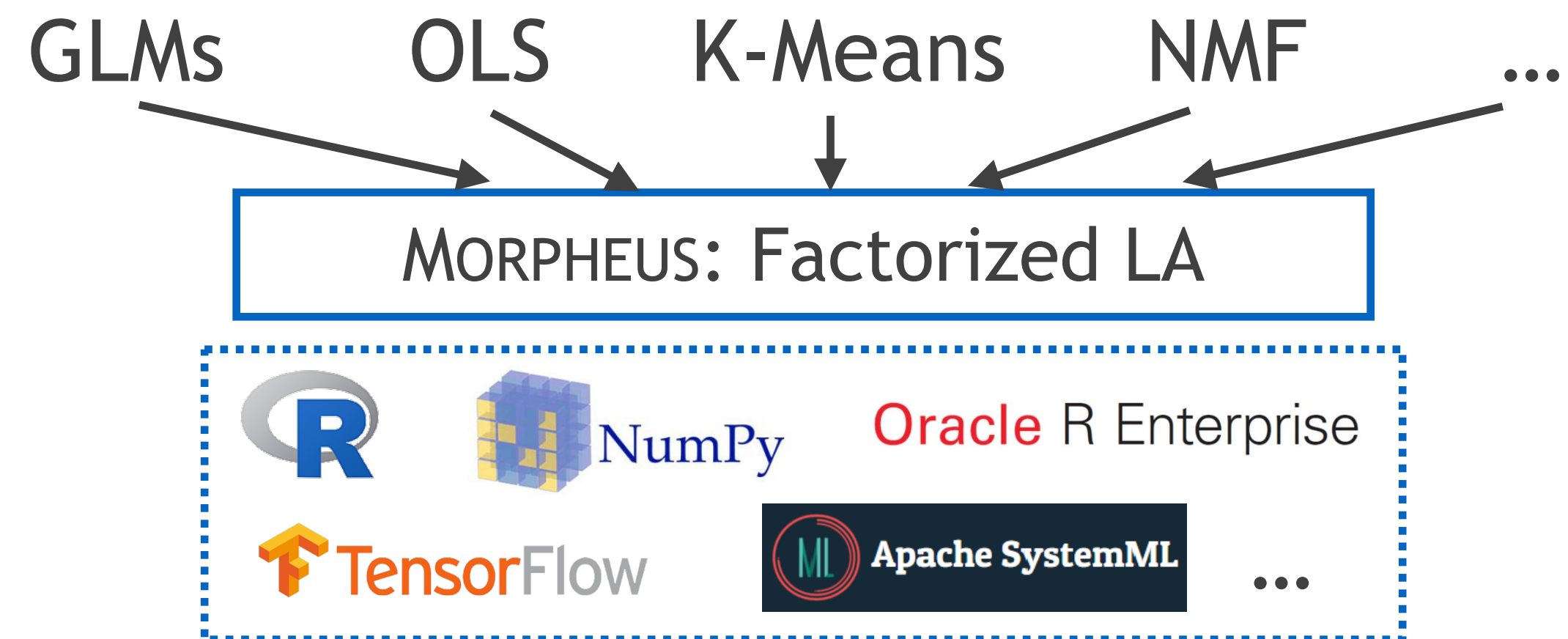
Goal: Automate factorized ML to many ML algorithms in a unified way

Idea: Many ML algorithms are bulk *linear algebra* (LA) programs
Create a framework for rewrite rules for LA ops

Factorized ML: Prior Work



The MORPHEUS Approach



Bulk LA-based ML Algorithms

Bulk LA-based ML Algorithms

Ordinary Least Squares linear regression with normal equations

Input: Regular matrix T, Y, w
 $w = \text{ginv}(\text{crossprod}(T))(T^\top Y)$

Bulk LA-based ML Algorithms

Ordinary Least Squares linear regression with normal equations

```
Input: Regular matrix  $T, Y, w$   
 $w = \text{ginv}(\text{crossprod}(T))(T^\top Y)$ 
```

Logistic regression with BGD; works for L-BFGS and Conjugate Gradient too

```
Input: Regular matrix  $T, Y, w, \alpha$   
for  $i$  in  $1 : \text{max\_iter}$  do  
  |  $w = w + \alpha * (T^\top (Y / (1 + \exp(Tw))))$   
end
```

MORPHEUS: High-level Architecture

MORPHEUS: High-level Architecture

MORPHEUS

Rewrite Rules for
Factorized LA ops
on an LA tool



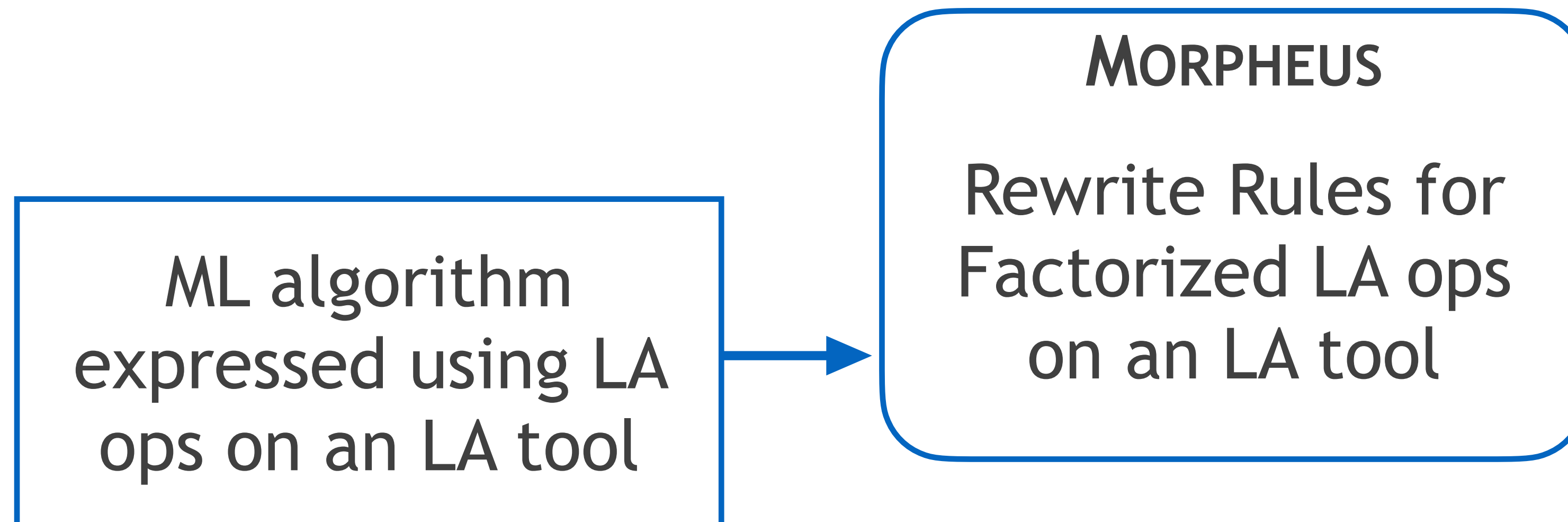
Oracle R Enterprise



NumPy



MORPHEUS: High-level Architecture



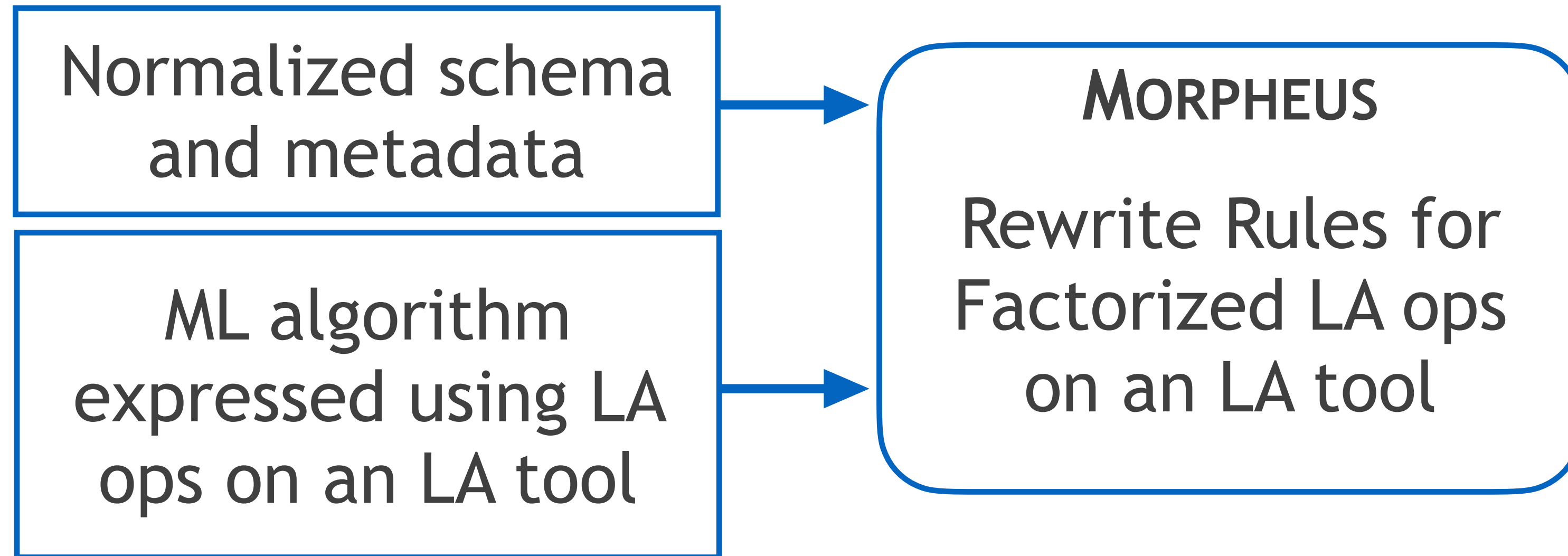
Oracle R Enterprise



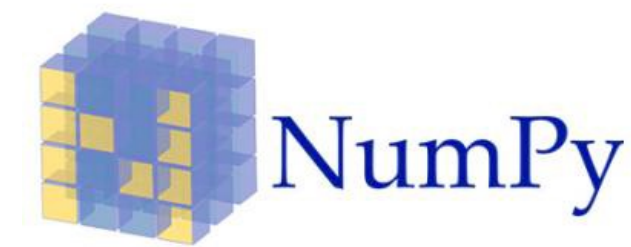
NumPy



MORPHEUS: High-level Architecture



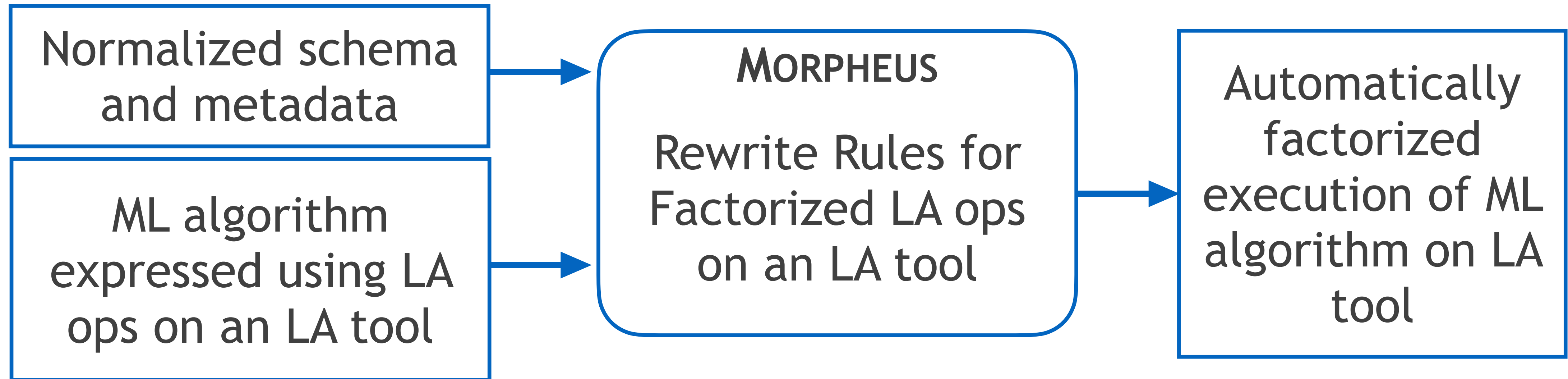
Oracle R Enterprise



NumPy



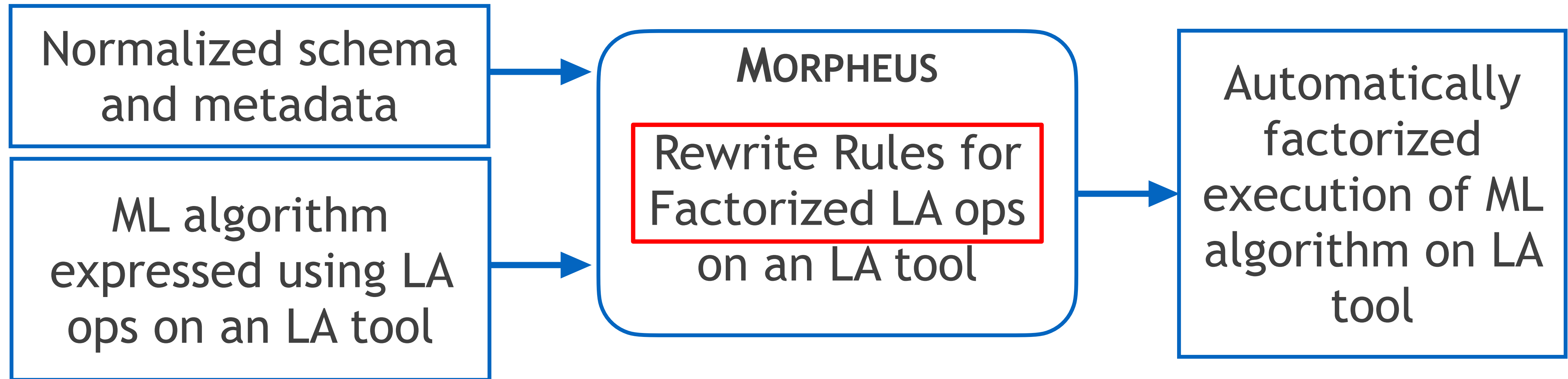
MORPHEUS: High-level Architecture



Oracle R Enterprise



MORPHEUS: High-level Architecture



Oracle R Enterprise



MORPHEUS: Factorized LA Rewrite Rules

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

Employers

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

Customers Employers

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

$S_{n \times d_S}$

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

$S_{n \times d_S}$ $R_{n_R \times d_R}$

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, \boxed{X}) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S \ X_R]$ Customers Employers

$T_{n \times d}$ $S_{n \times d_S}$ $R_{n_R \times d_R}$

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

$T_{n \times d}$ $S_{n \times d_S}$ $R_{n_R \times d_R}$

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

$T_{n \times d}$ $S_{n \times d_S}$ $K_{n \times n_R}$ $R_{n_R \times d_R}$

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

$T_{n \times d}$ $S_{n \times d_S}$ $K_{n \times n_R}$ $R_{n_R \times d_R}$

$$K[i, j] = \begin{cases} 1, & \text{if } \mathbf{S}[i].FK = j \\ 0, & \text{o/w} \end{cases}$$

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

$T_{n \times d}$ $S_{n \times d_S}$ $K_{n \times n_R}$ $R_{n_R \times d_R}$

$$T = [S \quad KR] \quad K[i, j] = \begin{cases} 1, & \text{if } \mathbf{S}[i].FK = j \\ 0, & \text{o/w} \end{cases}$$

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

$T_{n \times d}$ $S_{n \times d_S}$ $K_{n \times n_R}$ $R_{n_R \times d_R}$

$$T = [S \quad KR] \quad K[i, j] = \begin{cases} 1, & \text{if } \mathbf{S}[i].FK = j \\ 0, & \text{o/w} \end{cases}$$

Framework of algebraic rewrite rules for many LA operations

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

$T_{n \times d}$ $S_{n \times d_S}$ $K_{n \times n_R}$ $R_{n_R \times d_R}$

$$T = [S \quad KR] \quad K[i, j] = \begin{cases} 1, & \text{if } \mathbf{S}[i].FK = j \\ 0, & \text{o/w} \end{cases}$$

Framework of algebraic rewrite rules for many LA operations

Left Matrix Multiplication:

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

$T_{n \times d}$ $S_{n \times d_S}$ $K_{n \times n_R}$ $R_{n_R \times d_R}$

$$T = [S \quad KR] \quad K[i, j] = \begin{cases} 1, & \text{if } \mathbf{S}[i].FK = j \\ 0, & \text{o/w} \end{cases}$$

Framework of algebraic rewrite rules for many LA operations

Left Matrix Multiplication: $Tw \rightarrow Sw_S + K(Rw_R)$

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

$T_{n \times d}$ $S_{n \times d_S}$ $K_{n \times n_R}$ $R_{n_R \times d_R}$

$$T = [S \quad KR] \quad K[i, j] = \begin{cases} 1, & \text{if } \mathbf{S}[i].FK = j \\ 0, & \text{o/w} \end{cases}$$

Framework of algebraic rewrite rules for many LA operations

Left Matrix Multiplication: $Tw \rightarrow Sw_S + K(Rw_R)$

MORPHEUS: Factorized LA Rewrite Rules

New Abstraction: “Normalized Matrix” to represent join in LA

$$\mathbf{T}(ID, X) \leftarrow \pi(\mathbf{S}(ID, X_S, FK) \bowtie_{FK=RID} \mathbf{R}(RID, X_R))$$

$X \equiv [X_S X_R]$ Customers Employers

$T_{n \times d}$ $S_{n \times d_S}$ $K_{n \times n_R}$ $R_{n_R \times d_R}$

$$T = [S \quad KR] \quad K[i, j] = \begin{cases} 1, & \text{if } \mathbf{S}[i].FK = j \\ 0, & \text{o/w} \end{cases}$$

Framework of algebraic rewrite rules for many LA operations

Left Matrix Multiplication: $Tw \rightarrow Sw_S + K(Rw_R)$

GLMs, K-means clustering, NMF, etc. *automatically* factorized

Automatically Factorized ML in MORPHEUS

Automatically Factorized ML in MORPHEUS

```
Input: Regular matrix  $T, Y, w, \alpha$   
for  $i$  in  $1 : max\_iter$  do  
  |  $w = w + \alpha * (T^T (Y / (1 + \exp(Tw))))$   
end
```

Automatically Factorized ML in MORPHEUS

```
Input: Regular matrix  $T, Y, w, \alpha$   
for  $i$  in  $1 : max\_iter$  do  
  |  $w = w + \alpha * (T^T (Y / (1 + \exp(Tw))))$   
end
```

$T \equiv (S, K, R)$ ↓ MORPHEUS

Automatically Factorized ML in MORPHEUS

```
Input: Regular matrix  $T, Y, w, \alpha$   
for  $i$  in  $1 : max\_iter$  do  
|  $w = w + \alpha * (T^\top (Y / (1 + \exp(Tw))))$   
end
```

$T \equiv (S, K, R)$ ↓ MORPHEUS

```
Input: Normalized matrix  $(S, K, R), Y, w, \alpha$   
for  $i$  in  $1 : max\_iter$  do  
|  $P = (Y / (1 + \exp(Sw[1 : d_S, ] +$   
|  $\qquad\qquad\qquad K(Rw[d_S + 1 : d_S + d_R, ]))))^\top$   
|  $w = w + \alpha * [PS, (PK)R]^\top$   
end
```

LA Operations Factorized in MORPHEUS

Table 1: Operators and functions of linear algebra handled in this paper over a normalized matrix T .

Op Type	Name	Expression	Output Type	Parameter X or x	Factorizable
Element-wise Scalar Op	Arithmetic Op ($\otimes = +, -, *, /, \hat{\cdot}$, etc)	$T \otimes x$ or $x \otimes T$	Normalized Matrix	A scalar	Yes
	Transpose	T^\top		N/A	
	Scalar Function f (e.g., log, exp, sin)	$f(T)$		Parameters for f	
Aggregation	Row Summation	rowSums(T)	Column Vector	N/A	
	Column Summation	colSums(T)	Row Vector		
	Summation	sum(T)	Scalar		
Multiplication	Left Multiplication	TX	Regular Matrix	$(d_S + d_R) \times d_X$ matrix	
	Right Multiplication	XT		$n_X \times n_S$ matrix	
	Cross-product	crossprod(T)		N/A	
Inversion	Pseudoinverse	ginv(T)			
Element-wise Matrix Op	Arithmetic Op ($\otimes = +, -, *, /, \hat{\cdot}$, etc)	$X \otimes T$ or $T \otimes X$		$n_S \times (d_S + d_R)$ matrix	No

Snapshot of Empirical Results

Snapshot of Empirical Results

Prototype in R (and Python) for listed LA ops; ~800 LOC; commodity machine

Snapshot of Empirical Results

Prototype in R (and Python) for listed LA ops; ~800 LOC; commodity machine



S: Ratings

R_1 : Users

R_2 : Businesses

Snapshot of Empirical Results

Prototype in R (and Python) for listed LA ops; ~800 LOC; commodity machine



S: Ratings

R₁: Users

R₂: Businesses



S: Listings

R₁: Hotels

R₂: Search details

Snapshot of Empirical Results

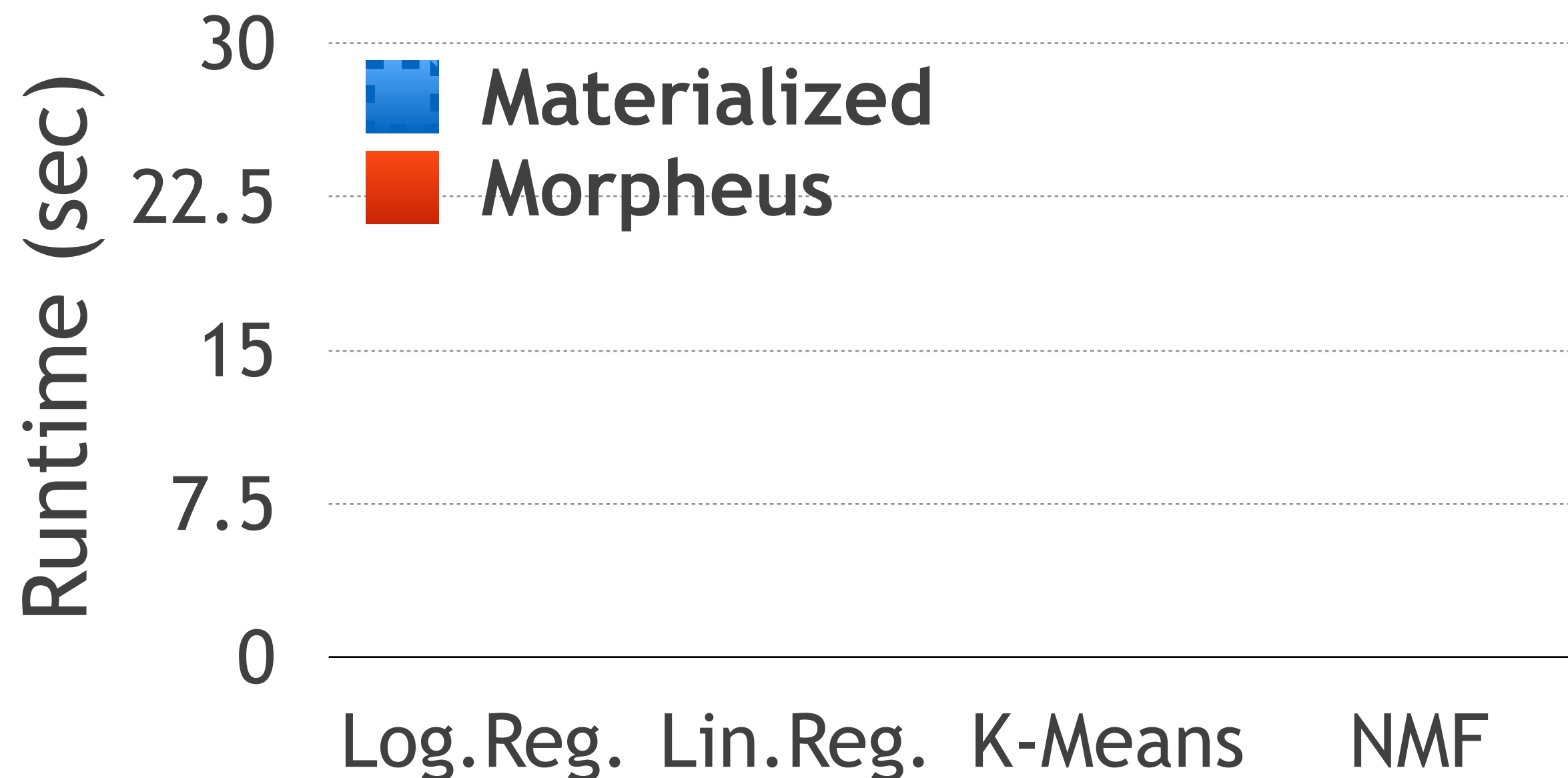
Prototype in R (and Python) for listed LA ops; ~800 LOC; commodity machine



S: Ratings
R₁: Users
R₂: Businesses



S: Listings
R₁: Hotels
R₂: Search details



Snapshot of Empirical Results

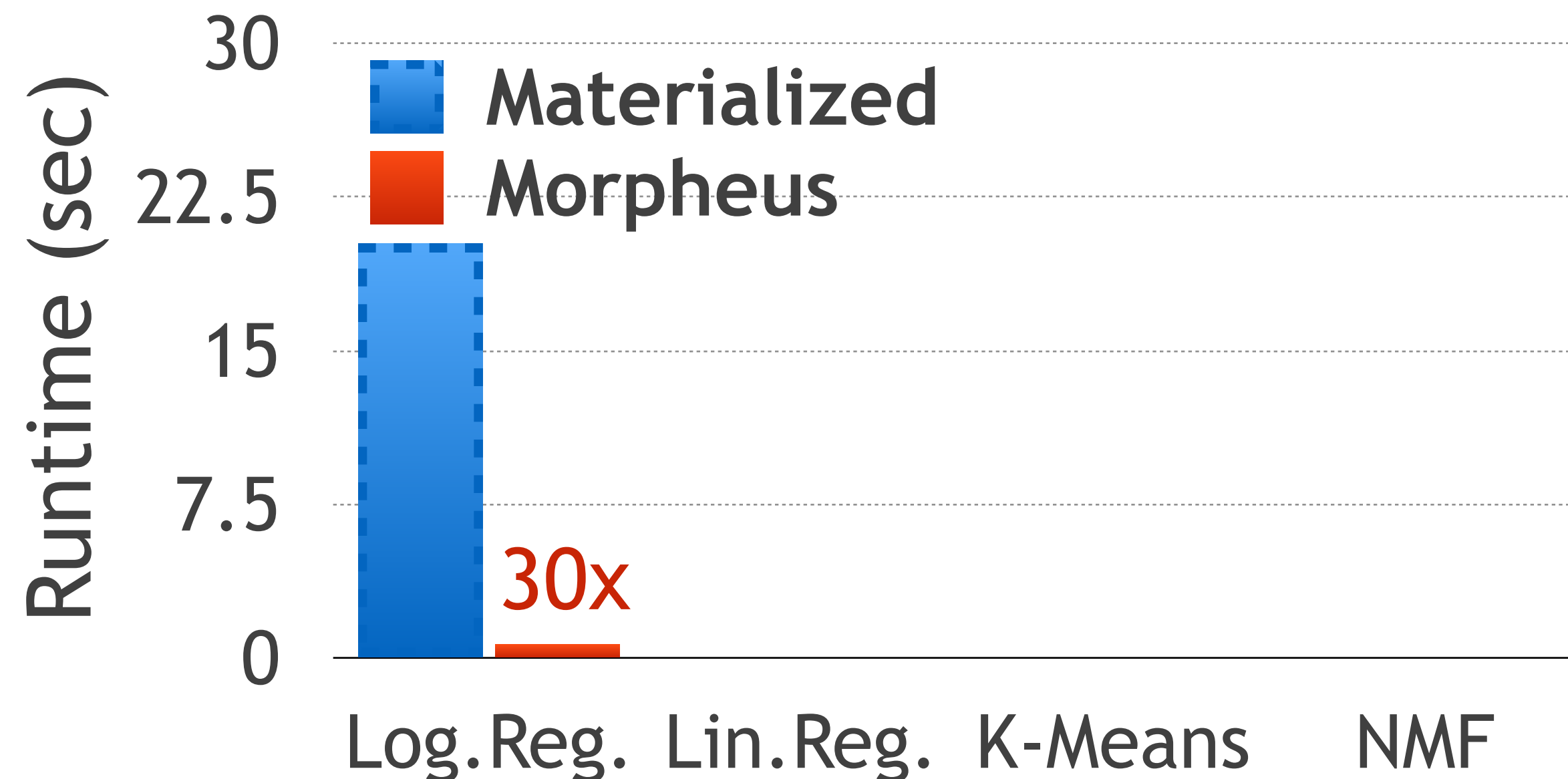
Prototype in R (and Python) for listed LA ops; ~800 LOC; commodity machine



S: Ratings
R₁: Users
R₂: Businesses



S: Listings
R₁: Hotels
R₂: Search details



Snapshot of Empirical Results

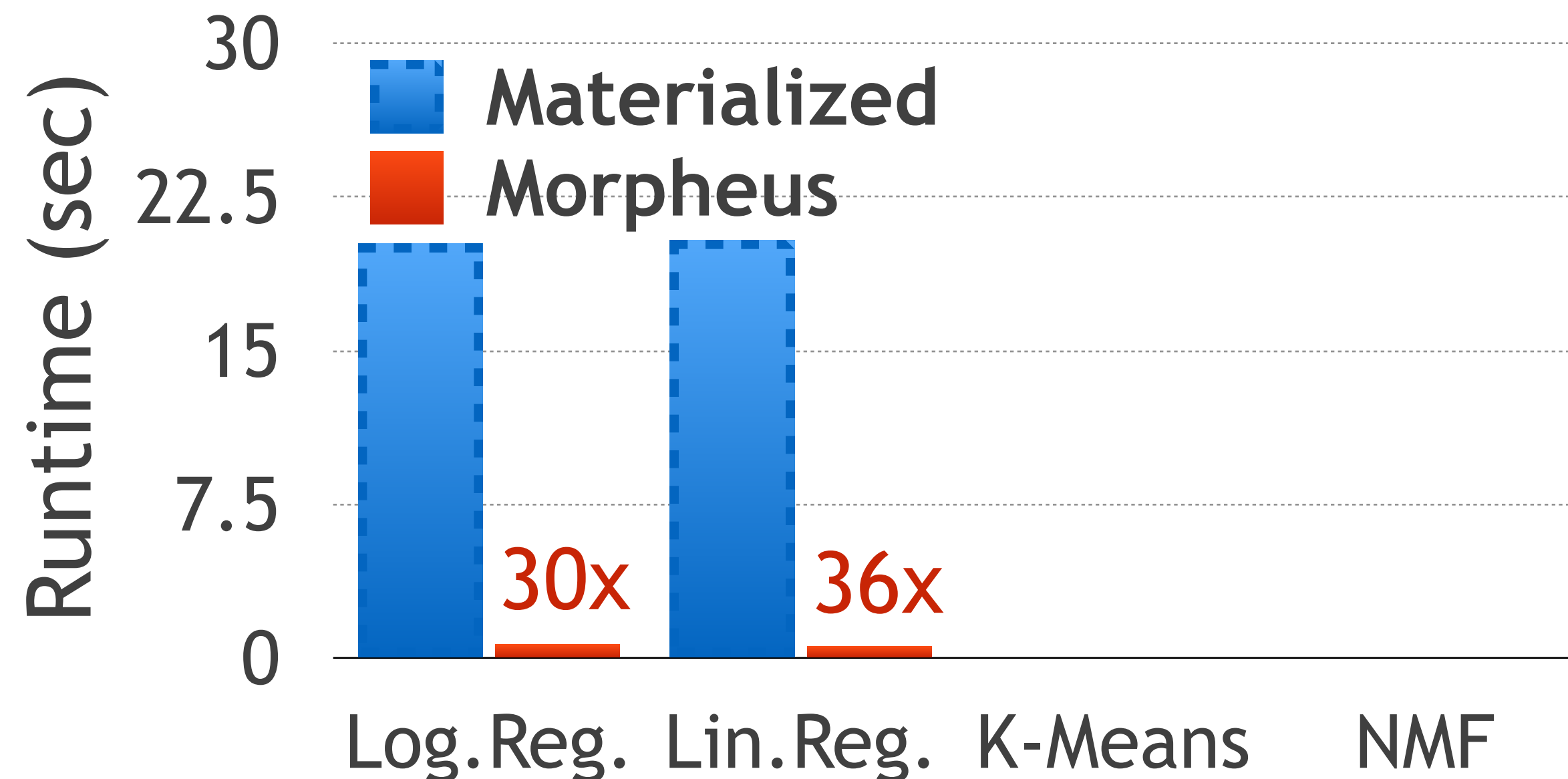
Prototype in R (and Python) for listed LA ops; ~800 LOC; commodity machine



S: Ratings
R₁: Users
R₂: Businesses



S: Listings
R₁: Hotels
R₂: Search details



Snapshot of Empirical Results

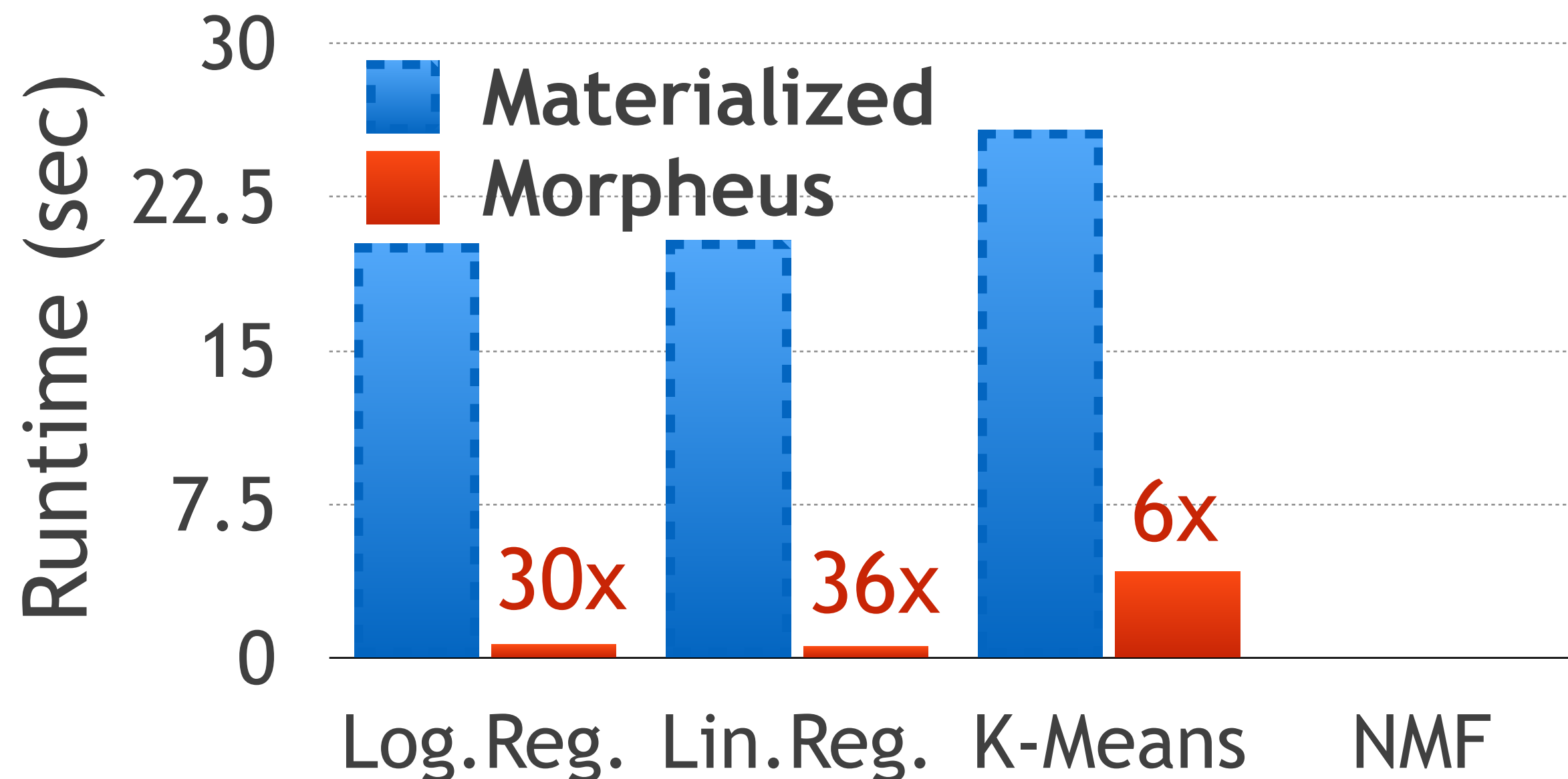
Prototype in R (and Python) for listed LA ops; ~800 LOC; commodity machine



S: Ratings
R₁: Users
R₂: Businesses



S: Listings
R₁: Hotels
R₂: Search details



Snapshot of Empirical Results

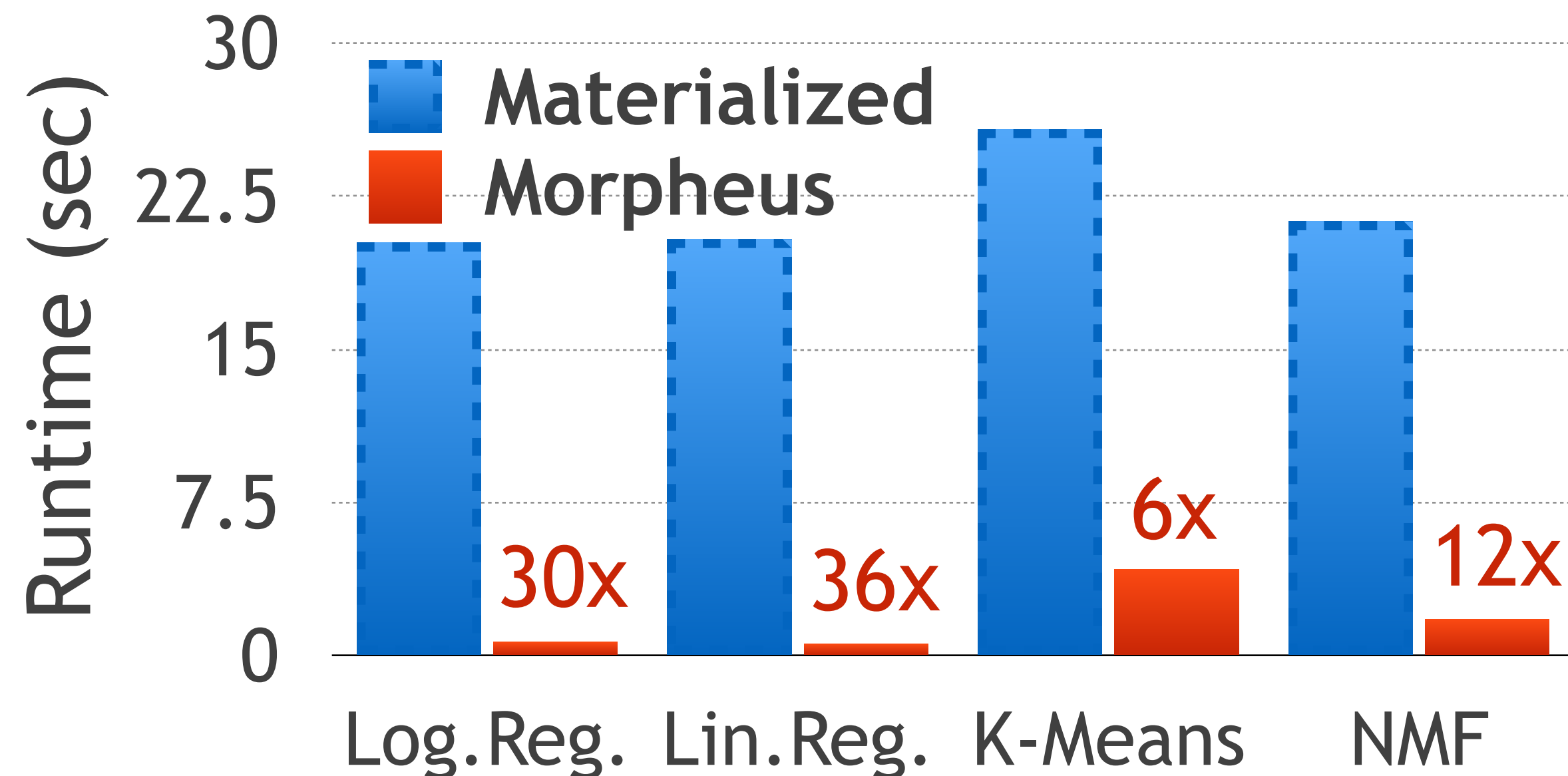
Prototype in R (and Python) for listed LA ops; ~800 LOC; commodity machine



S: Ratings
R₁: Users
R₂: Businesses



S: Listings
R₁: Hotels
R₂: Search details

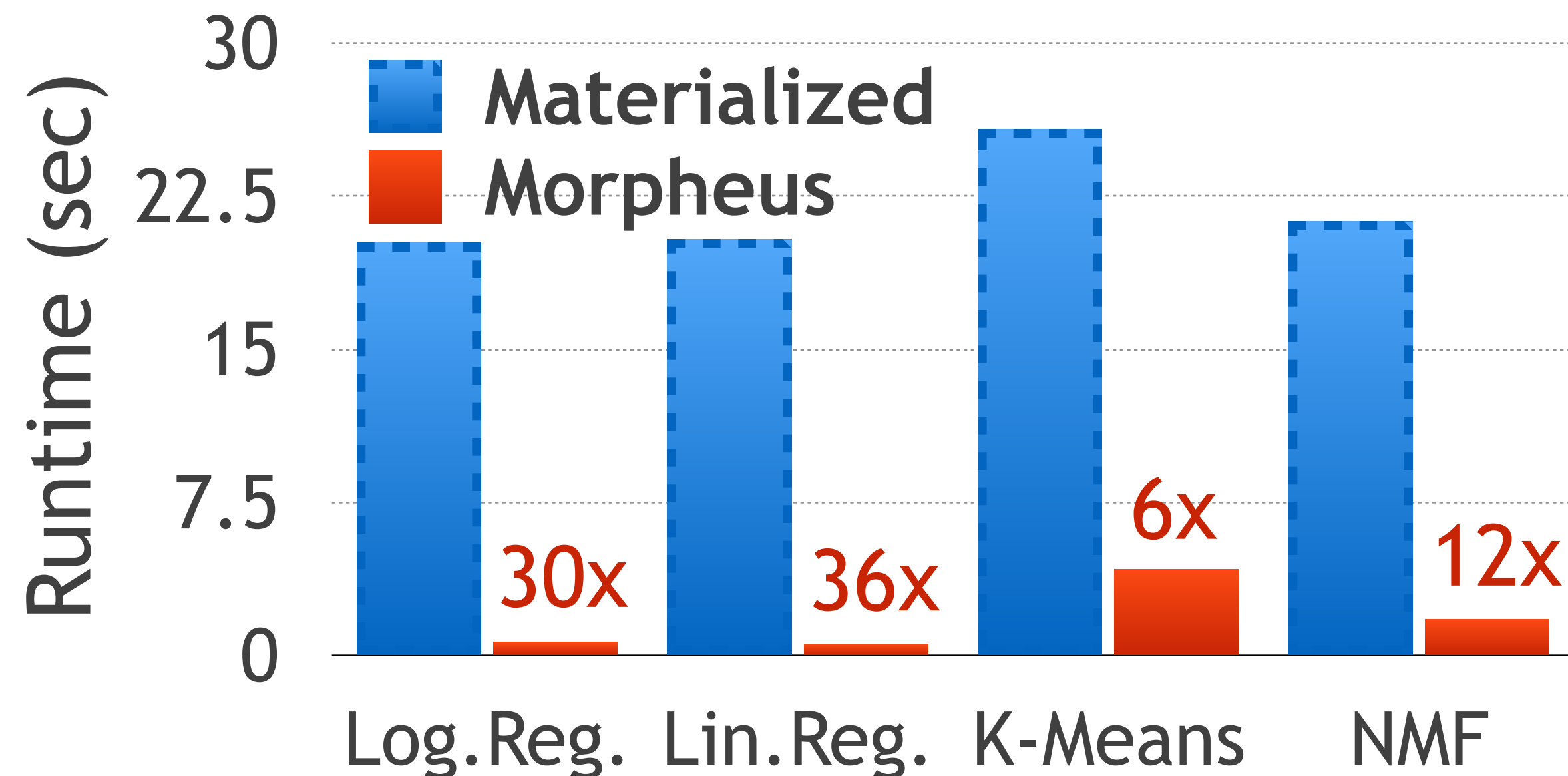


Snapshot of Empirical Results

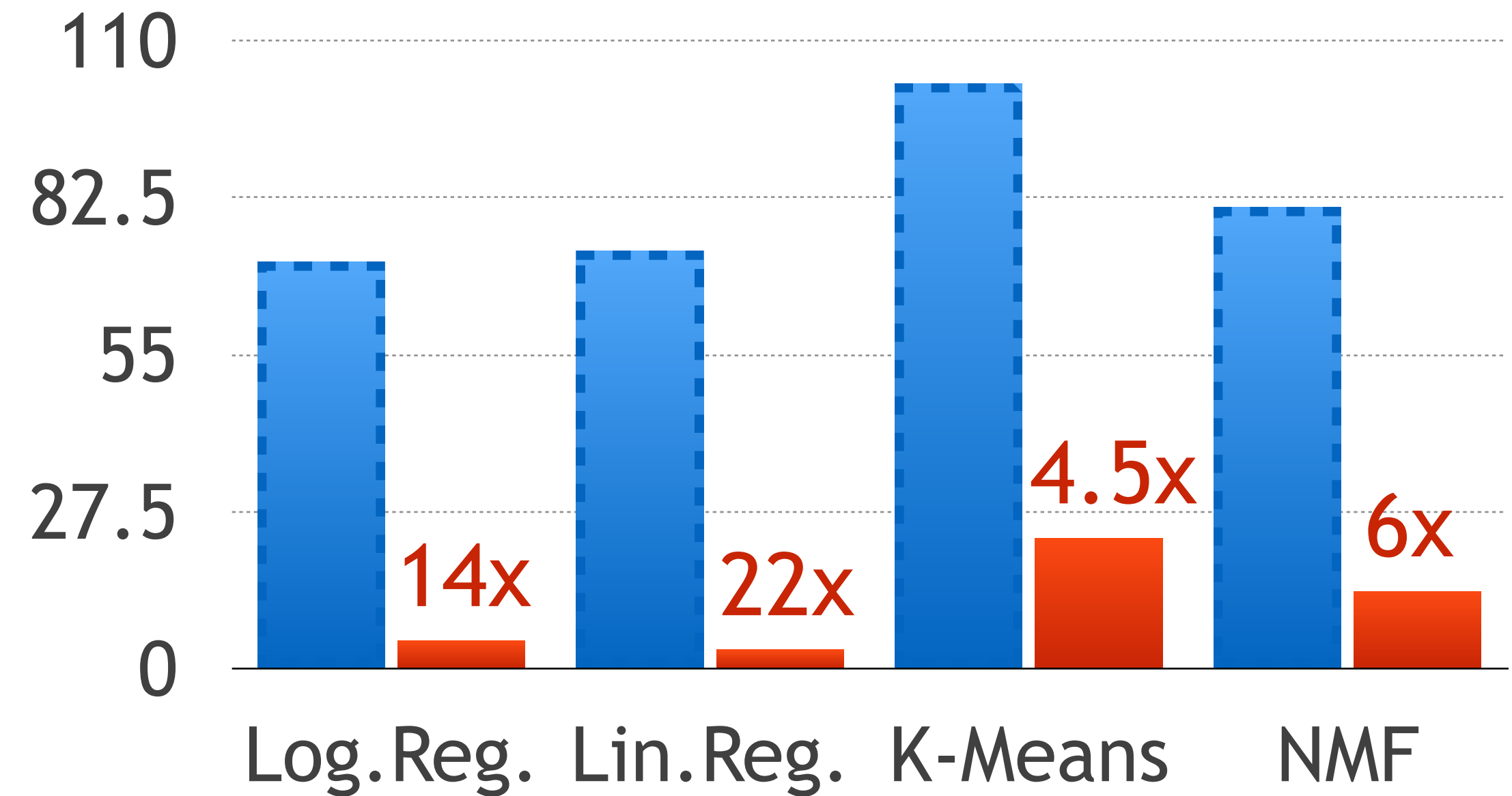
Prototype in R (and Python) for listed LA ops; ~800 LOC; commodity machine



S: Ratings
R₁: Users
R₂: Businesses



S: Listings
R₁: Hotels
R₂: Search details



When is MORPHEUS not likely to be beneficial?

When is MORPHEUS not likely to be beneficial?

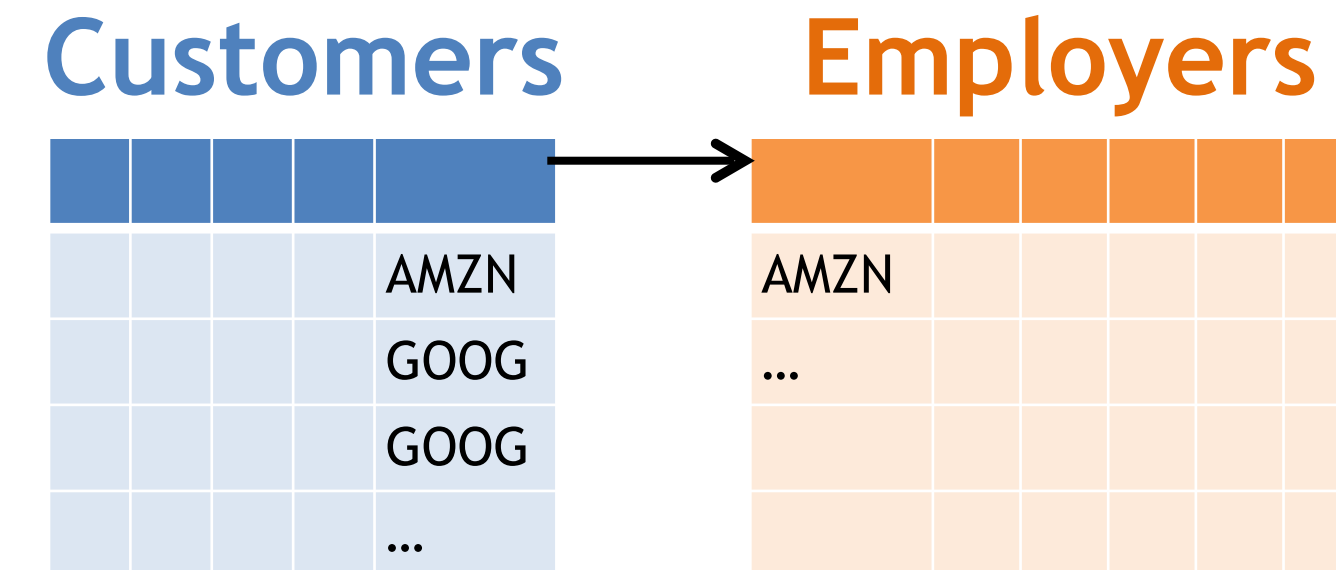
Short Answer: When the join(s) do *not* introduce much redundancy

When is MORPHEUS not likely to be beneficial?

Short Answer: When the join(s) do *not* introduce much redundancy

Case 1:

Fact table is *not much taller* than dimension table(s)

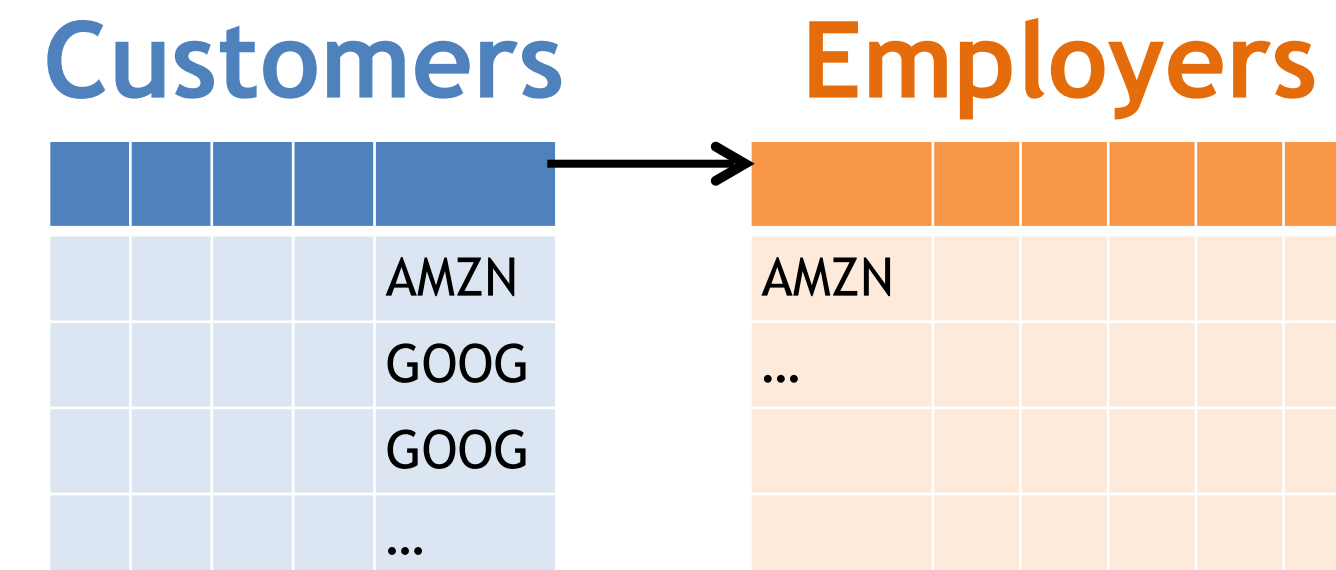


When is MORPHEUS not likely to be beneficial?

Short Answer: When the join(s) do *not* introduce much redundancy

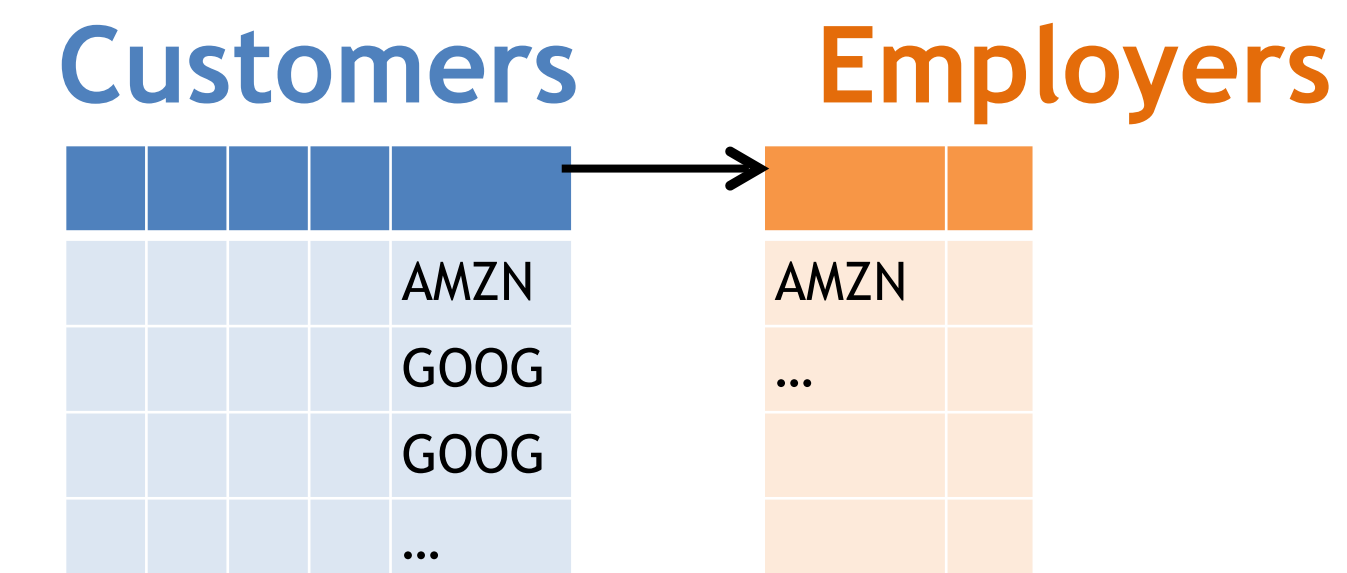
Case 1:

Fact table is *not much taller* than dimension table(s)



Case 2:

Dimension table has *much fewer* features than fact table

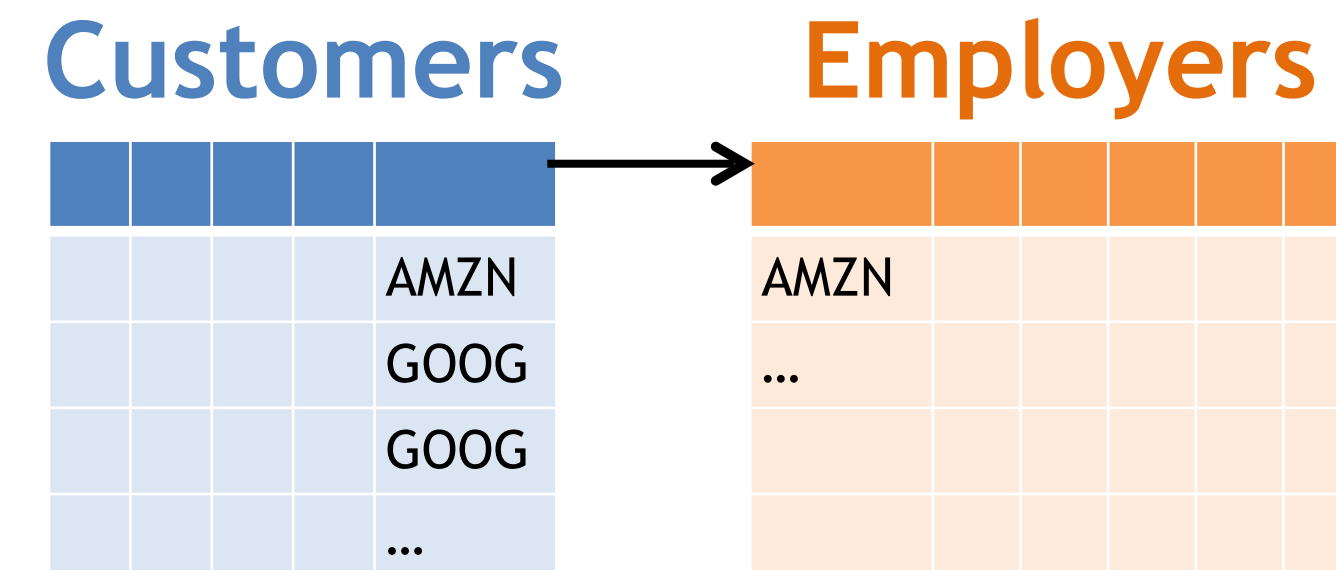


When is MORPHEUS not likely to be beneficial?

Short Answer: When the join(s) do *not* introduce much redundancy

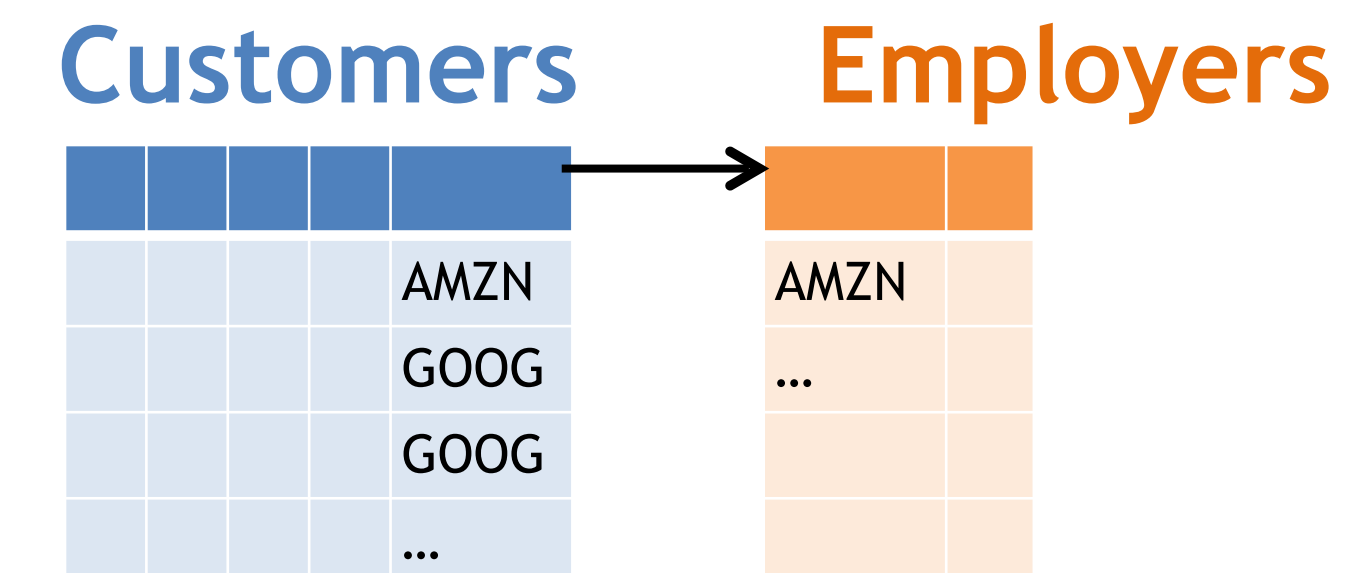
Case 1:

Fact table is *not much taller* than dimension table(s)



Case 2:

Dimension table has *much fewer* features than fact table



Case 3: MLPs do *not* have much computational redundancy (anyway)

MORPHEUS: Implementations and Extensions

Towards Linear Algebra over Normalized Data. **VLDB 2017**

Enabling and Optimizing Non-linear Feature Interactions in Factorized Linear Algebra. **SIGMOD 2019**

Tuple-oriented Compression for Large-scale Mini-batch Stochastic Gradient Descent. **SIGMOD 2019**

MORPHEUS: Implementations and Extensions

Library released for both R and Python NumPy

Towards Linear Algebra over Normalized Data. [VLDB 2017](#)

Enabling and Optimizing Non-linear Feature Interactions in Factorized Linear Algebra. [SIGMOD 2019](#)

Tuple-oriented Compression for Large-scale Mini-batch Stochastic Gradient Descent. [SIGMOD 2019](#)

MORPHEUS: Implementations and Extensions

Library released for both R and Python NumPy

Supports star schemas for many LA ops; snowflakes can be reduced to star

Towards Linear Algebra over Normalized Data. [VLDB 2017](#)

Enabling and Optimizing Non-linear Feature Interactions in Factorized Linear Algebra. [SIGMOD 2019](#)

Tuple-oriented Compression for Large-scale Mini-batch Stochastic Gradient Descent. [SIGMOD 2019](#)

MORPHEUS: Implementations and Extensions

Library released for both R and Python NumPy

Supports star schemas for many LA ops; snowflakes can be reduced to star

Some data cleaning/prep ops also factorized

Towards Linear Algebra over Normalized Data. [VLDB 2017](#)

Enabling and Optimizing Non-linear Feature Interactions in Factorized Linear Algebra. [SIGMOD 2019](#)

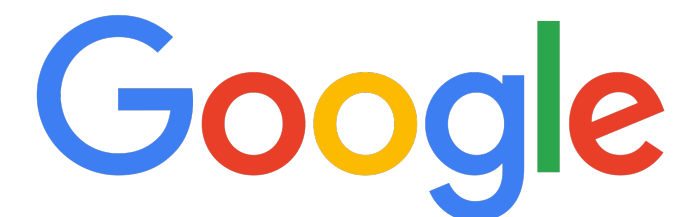
Tuple-oriented Compression for Large-scale Mini-batch Stochastic Gradient Descent. [SIGMOD 2019](#)

MORPHEUS: Implementations and Extensions

Library released for both R and Python NumPy

Supports star schemas for many LA ops; snowflakes can be reduced to star

Some data cleaning/prep ops also factorized



Towards Linear Algebra over Normalized Data. [VLDB 2017](#)

Enabling and Optimizing Non-linear Feature Interactions in Factorized Linear Algebra. [SIGMOD 2019](#)

Tuple-oriented Compression for Large-scale Mini-batch Stochastic Gradient Descent. [SIGMOD 2019](#)

MORPHEUS: Implementations and Extensions

Library released for both R and Python NumPy

Supports star schemas for many LA ops; snowflakes can be reduced to star

Some data cleaning/prep ops also factorized



MORPHEUSFI: Second-order feature interactions in Morpheus

Towards Linear Algebra over Normalized Data. [VLDB 2017](#)

Enabling and Optimizing Non-linear Feature Interactions in Factorized Linear Algebra. [SIGMOD 2019](#)

Tuple-oriented Compression for Large-scale Mini-batch Stochastic Gradient Descent. [SIGMOD 2019](#)

MORPHEUS: Implementations and Extensions

Library released for both R and Python NumPy

Supports star schemas for many LA ops; snowflakes can be reduced to star

Some data cleaning/prep ops also factorized



MORPHEUSFI: Second-order feature interactions in Morpheus

MORPHEUSFLOW: “Lazy join” for SGD in TensorFlow

Towards Linear Algebra over Normalized Data. [VLDB 2017](#)

Enabling and Optimizing Non-linear Feature Interactions in Factorized Linear Algebra. [SIGMOD 2019](#)

Tuple-oriented Compression for Large-scale Mini-batch Stochastic Gradient Descent. [SIGMOD 2019](#)

MORPHEUS: Implementations and Extensions

Library released for both R and Python NumPy

Supports star schemas for many LA ops; snowflakes can be reduced to star

Some data cleaning/prep ops also factorized



MORPHEUSFI: Second-order feature interactions in Morpheus

MORPHEUSFLOW: “Lazy join” for SGD in TensorFlow

TOC: Generalized data compression for SGD

Towards Linear Algebra over Normalized Data. [VLDB 2017](#)

Enabling and Optimizing Non-linear Feature Interactions in Factorized Linear Algebra. [SIGMOD 2019](#)

Tuple-oriented Compression for Large-scale Mini-batch Stochastic Gradient Descent. [SIGMOD 2019](#)

TRINITY: MORPHEUS Meets Oracle GraalVM

TRINITY: MORPHEUS Meets Oracle GraalVM

Goal: *Automate Morpheus itself to many PLs in a unified way*

TRINITY: MORPHEUS Meets Oracle GraalVM

Goal: *Automate* Morpheus itself to many PLs in a unified way

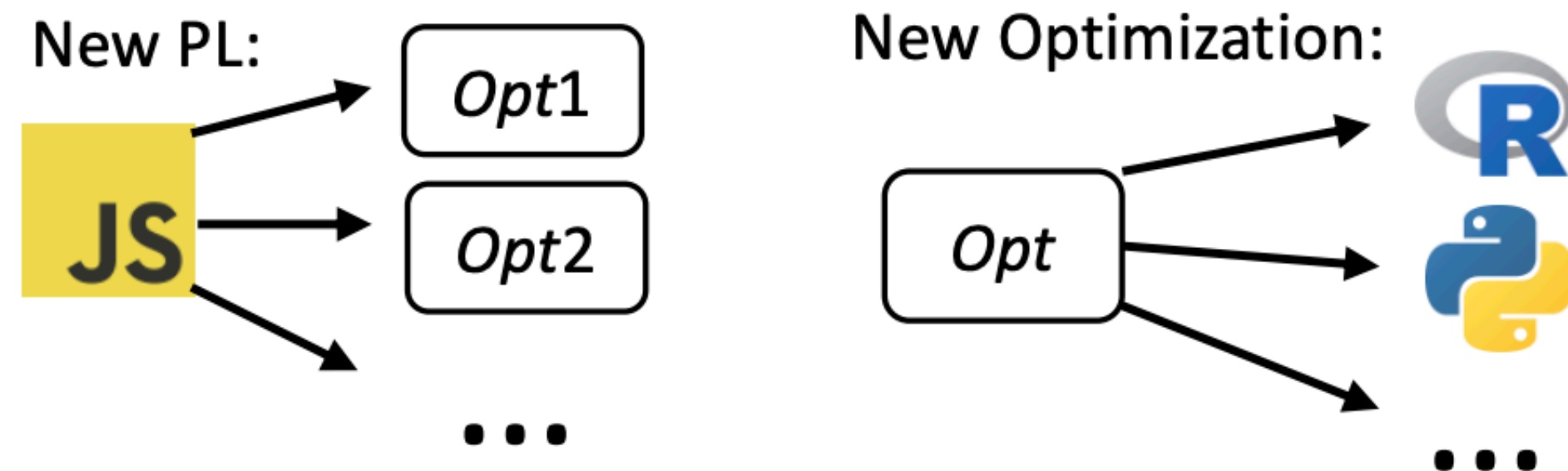
Idea: Exploit GraalVM, an industrial-strength *polyglot* compiler + runtime for data science workloads (R, Py, Javascript, etc.)

TRINITY: MORPHEUS Meets Oracle GraalVM

Goal: Automate Morpheus itself to many PLs in a unified way

Idea: Exploit GraalVM, an industrial-strength *polyglot* compiler + runtime for data science workloads (R, Py, Javascript, etc.)

World without Trinity

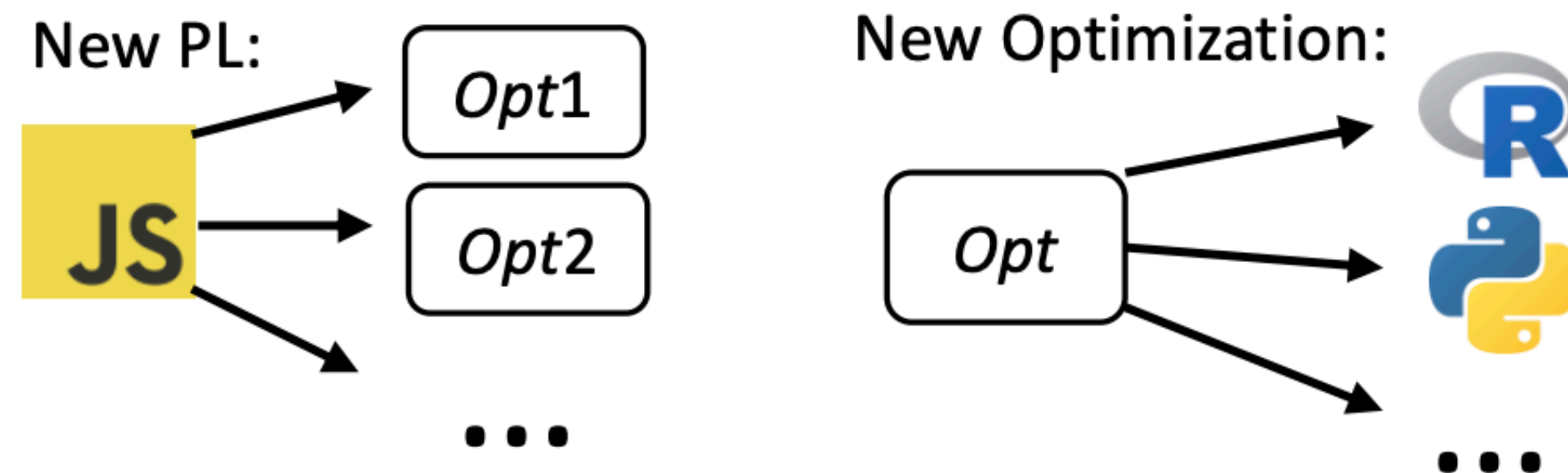


TRINITY: MORPHEUS Meets Oracle GraalVM

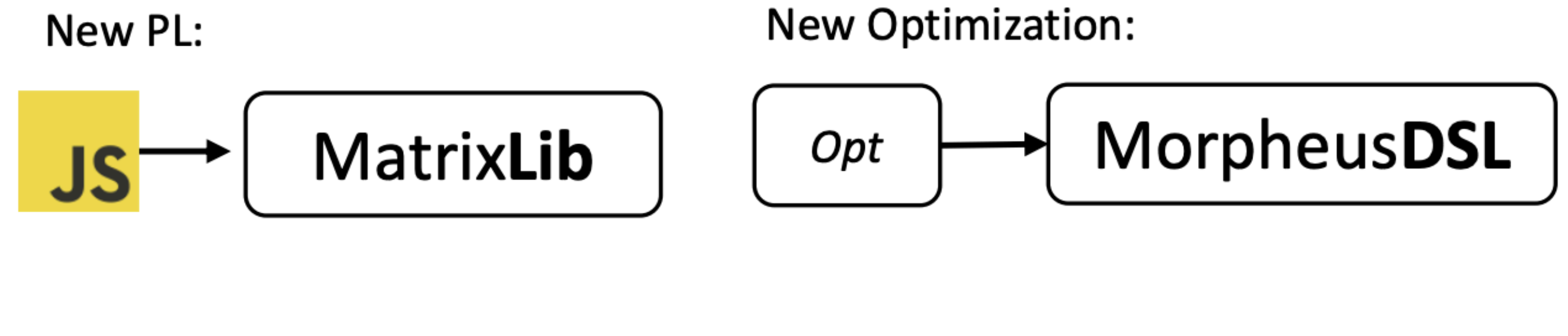
Goal: Automate Morpheus itself to many PLs in a unified way

Idea: Exploit GraalVM, an industrial-strength *polyglot* compiler + runtime for data science workloads (R, Py, Javascript, etc.)

World without Trinity



World with Trinity

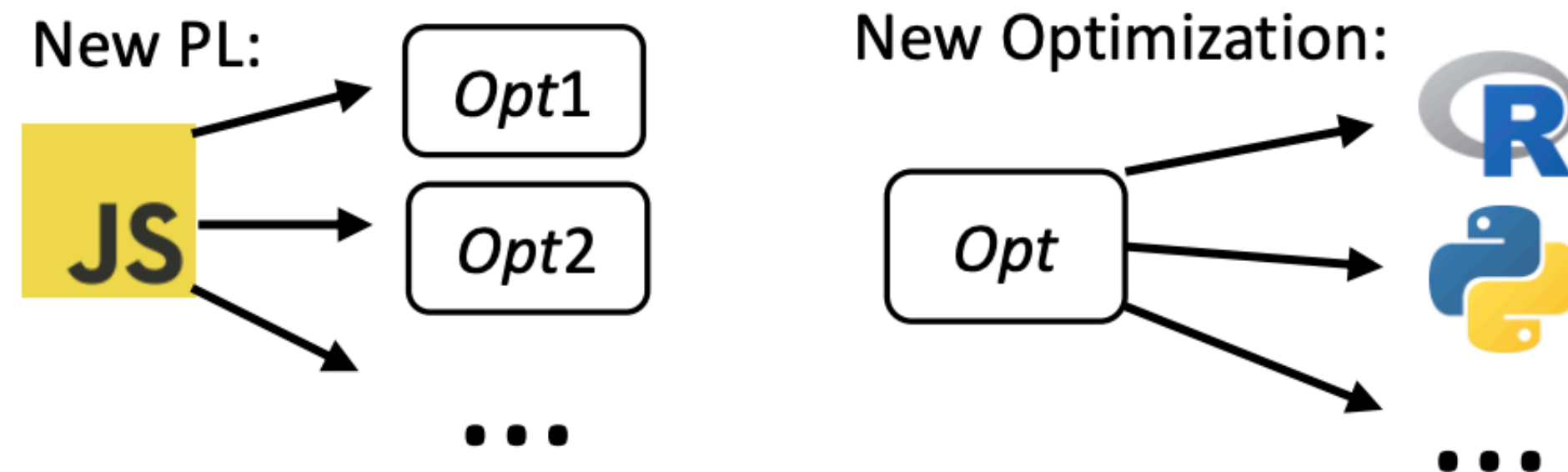


TRINITY: MORPHEUS Meets Oracle GraalVM

Goal: Automate Morpheus itself to many PLs in a unified way

Idea: Exploit GraalVM, an industrial-strength *polyglot* compiler + runtime for data science workloads (R, Py, Javascript, etc.)

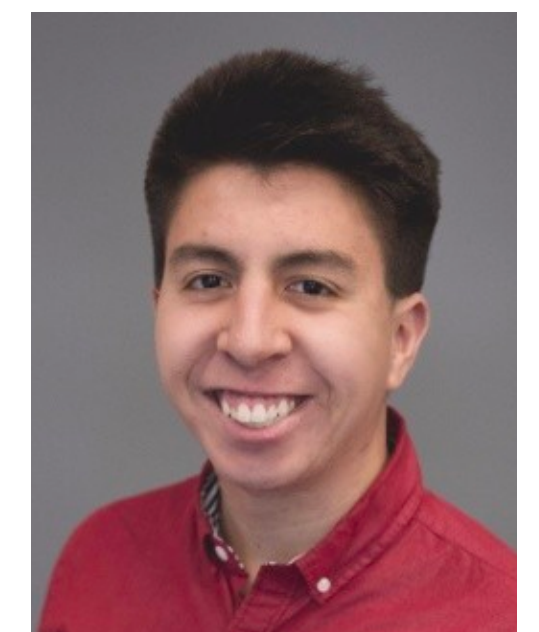
World without Trinity



World with Trinity



Learn more about Trinity from David Justo at 3:30pm today!



Towards A Polyglot Framework for Factorized ML. [VLDB 2021](#)

Outline

- 4m ■ Introducing ML over Joins
- 4m ■ Orion: Factorized ML
- 10m ■ Morpheus and Extensions
- 4m ■ Roadblocks and Musings

Roadblocks for Factorized ML

Roadblocks for Factorized ML

Observation: Factorized ML yet to have big practical impact on any path. :-/

Roadblocks for Factorized ML

Observation: Factorized ML yet to have big practical impact on any path. :-/

Reason 1: Applicability to *business-critical* ML algorithms limited

Tree ensembles rule tabular data; factorized ML gains marginal there

GLMs, clustering, etc. often not big bottleneck in real-world pipelines

Roadblocks for Factorized ML

Observation: Factorized ML yet to have big practical impact on any path. :-/

Reason 1: Applicability to *business-critical* ML algorithms limited

Tree ensembles rule tabular data; factorized ML gains marginal there

GLMs, clustering, etc. often not big bottleneck in real-world pipelines

Roadblocks for Factorized ML

Observation: Factorized ML yet to have big practical impact on any path. :-/

Reason 1: Applicability to *business-critical* ML algorithms limited

Tree ensembles rule tabular data; factorized ML gains marginal there
GLMs, clustering, etc. often not big bottleneck in real-world pipelines

Reason 2: *Implementation* effort to make it practical still non-trivial

Orion-style: UDFs too complex to implement/maintain on RDBMS/Spark

Morpheus-style: ML not always written as LA scripts; hidden C++ callouts

Trinity-style: Likely promising; over the wall at Oracle now! :)

Plug: First Textbook on ML Systems

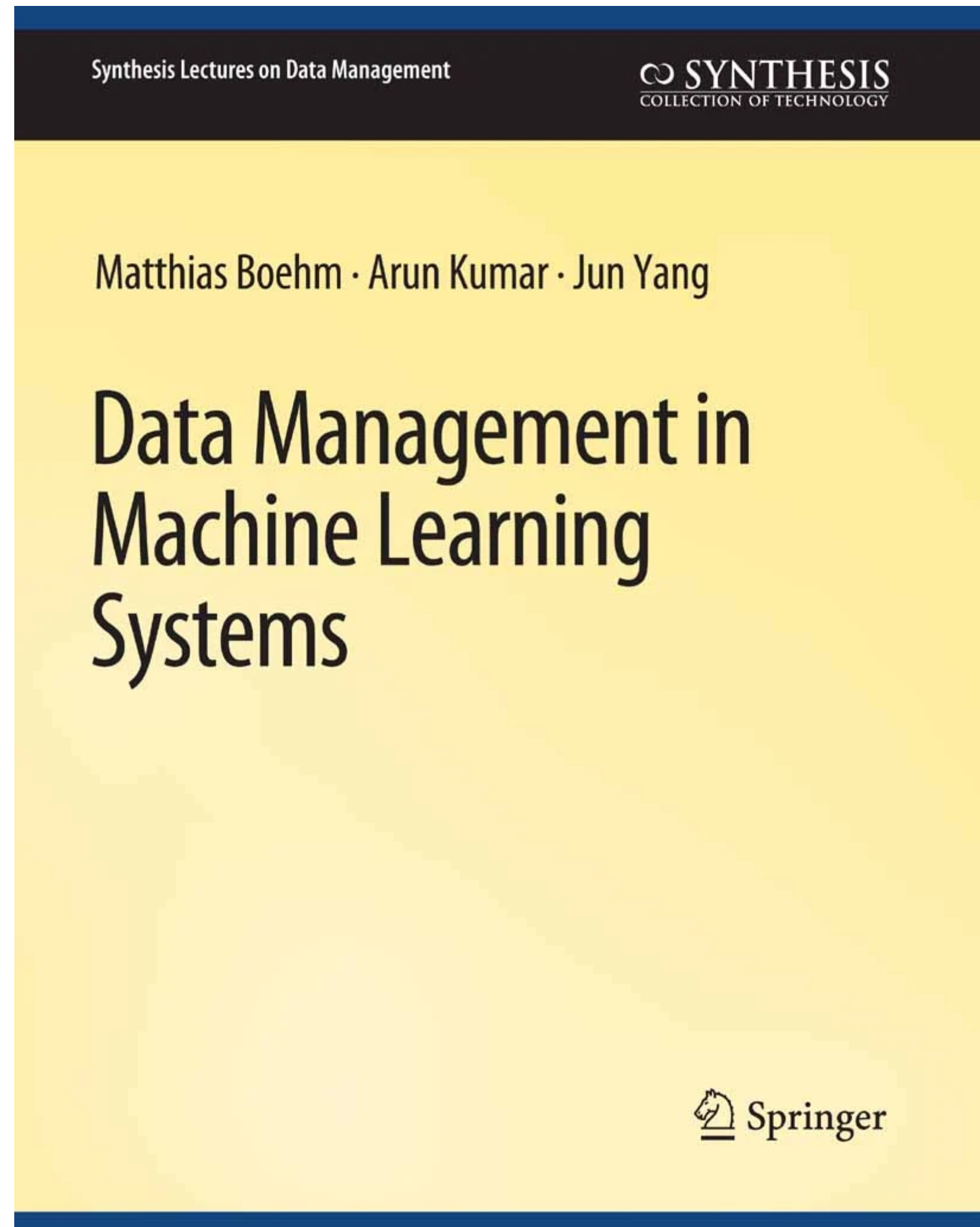


Table of Contents

Introduction

ML Through Database Queries and UDFs

Multi-Table ML and Deep Systems Integration

Rewrites and Optimization

Execution Strategies

Data Access Methods

Resource Heterogeneity and Elasticity

Systems for ML Lifecycle Tasks

Conclusions

Bibliography

Authors' Biographies

<https://tinyurl.com/MLSystemsBook>

<https://adalabucsd.github.io>

arunkk@eng.ucsd.edu



github.com/ADALabUCSD



@TweetAtAKK

ACKS:

