

The New DBfication of ML/AI

Arun Kumar



UC San Diego
JACOBS SCHOOL OF ENGINEERING
Computer Science and Engineering

UC San Diego
HALİCİOĞLU DATA SCIENCE INSTITUTE

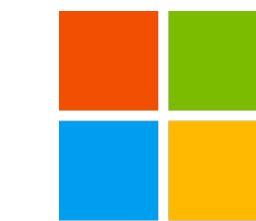
DB & AI Seminar, China

June 26, 2022

Golden Age of ML/AI



FACEBOOK

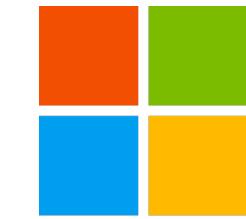


Microsoft

Golden Age of ML/AI



FACEBOOK



Microsoft



Insurance



Retail

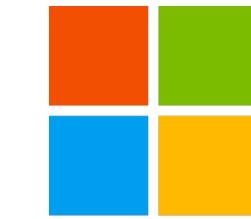


Sciences

Golden Age of ML/AI



FACEBOOK



Microsoft



Healthcare



Insurance



Retail



Sciences

\$ 38 billion
in 2019*



\$ 500 billion!
by 2024*

*International Data Corporation

Golden Age of ML/AI



FACEBOOK



Healthcare



Insurance



Retail



Sciences

\$ 38 billion
in 2019*



\$ 500 billion!
by 2024*

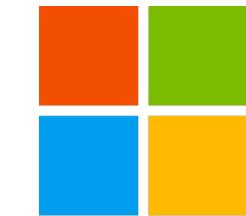
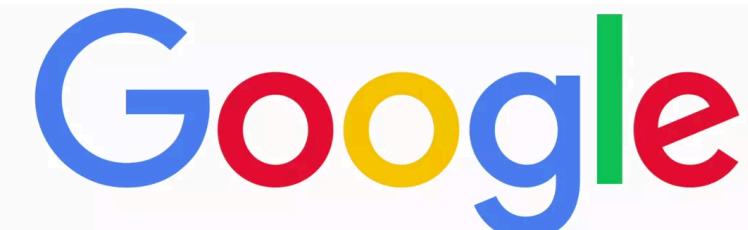


*International Data Corporation

Golden Age of ML/AI



FACEBOOK



Microsoft



Healthcare



Insurance



Retail



Sciences

\$ 38 billion
in 2019*



\$ 500 billion!
by 2024*

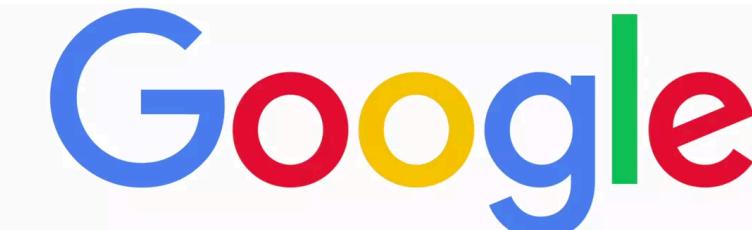


*International Data Corporation

Golden Age of ML/AI



FACEBOOK



Healthcare



Insurance



Retail



Sciences

\$ 38 billion
in 2019*



\$ 500 billion!
by 2024*

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/AI-based data analytics from
a data management/systems standpoint

My Research

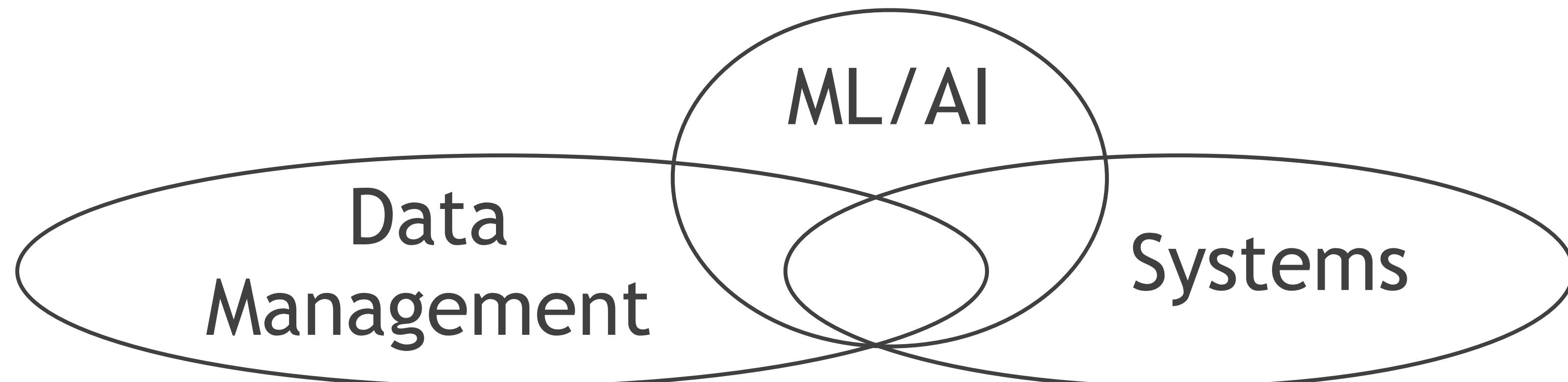
New abstractions, algorithms, and software systems
to “*democratize*” ML/AI-based data analytics from
a data management/systems standpoint

$$\text{Democratization} = \text{System Efficiency} \quad (\text{Reduce costs}) + \text{Human Efficiency} \quad (\text{Improve productivity})$$

My Research

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

$$\text{Democratization} = \text{System Efficiency} \quad (\text{Reduce costs}) + \text{Human Efficiency} \quad (\text{Improve productivity})$$



My Research

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

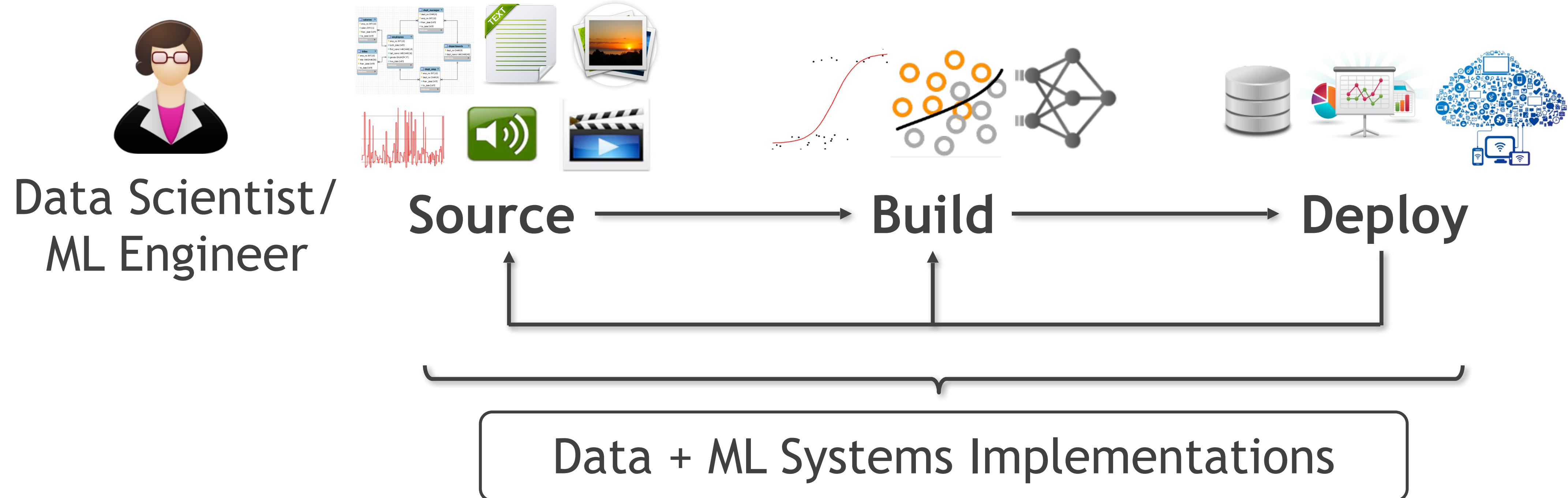
$$\text{Democratization} = \text{System Efficiency} \quad (\text{Reduce costs}) + \text{Human Efficiency} \quad (\text{Improve productivity})$$

Practical and scalable data systems for ML/AI analytics

Inspired by *relational database systems* principles

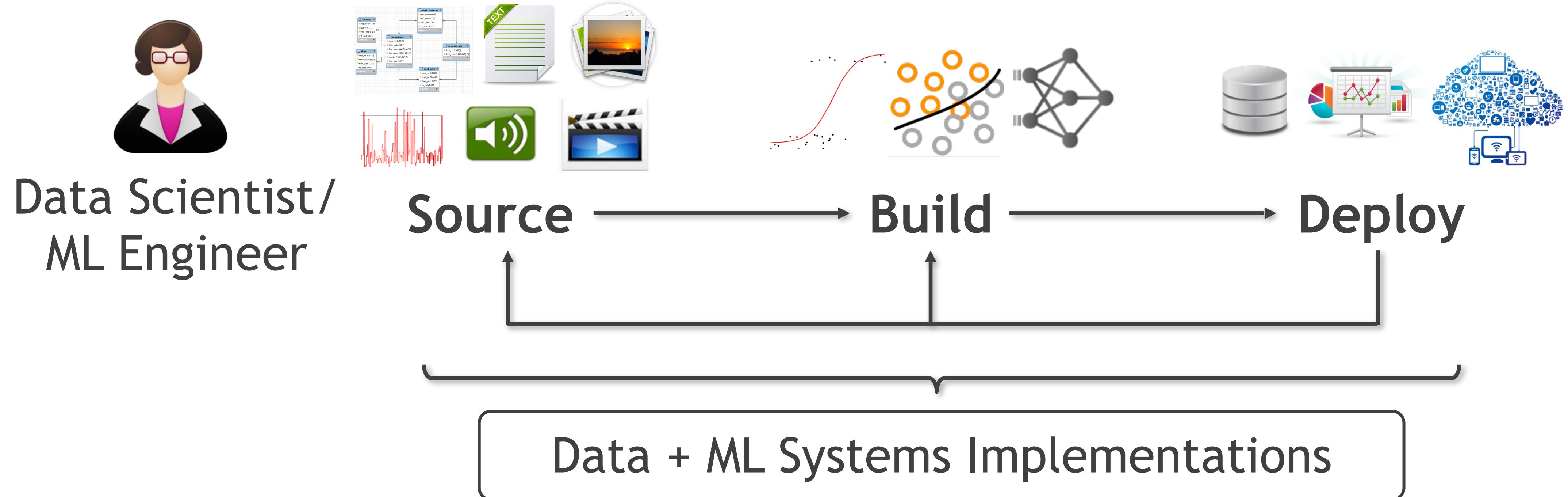
Exploit insights from *learning theory* and *optimization theory*

End-to-End ML Application Lifecycle



<https://ADALabUCSD.github.io>

End-to-End ML Application Lifecycle



Research Approach : *Abstract* key steps + *Formalize* computation + *Automate* grunt work + *Optimize* execution

Outline

- | The New DBfication of ML/AI
- | Two Examples from My Research
- | Accelerating the DBfication of ML/AI

DB-style Practical Concerns in ML

DB-style Practical Concerns in ML

Key concerns in ML:

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

Q: What if the dataset is larger than single-node RAM?

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

Scalability (and efficiency at scale)

Q: What if the dataset is larger than single-node RAM?

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

Scalability (and efficiency at scale)

Q: How are the features and models configured?

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

Scalability (and efficiency at scale)

Usability

Q: How are the features and models configured?

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

Scalability (and efficiency at scale)

Usability

Q: How does it fit within production systems and workflows?

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

Scalability (and efficiency at scale)

Usability

Manageability

Q: How does it fit within production systems and workflows?

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

Scalability (and efficiency at scale)

Usability

Manageability

Q: How to simplify the implementation of such systems?

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

Scalability (and efficiency at scale)

Usability

Manageability

Developability

Q: How to simplify the implementation of such systems?

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

Scalability (and efficiency at scale)

Usability

Manageability

Developability

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

Scalability (and efficiency at scale)

Usability

Manageability

Developability

*Long-standing
concerns in the
DB systems
world!*

DB-style Practical Concerns in ML

Key concerns in ML:

Accuracy

Runtime efficiency (sometimes)

Additional key *practical* concerns in ML Systems:

Scalability (and efficiency at scale)

Usability

Manageability

Developability

*Long-standing
concerns in the
DB systems
world!*

Can often trade off accuracy a bit to gain on the rest!

The Push to “Platformize” ML/AI

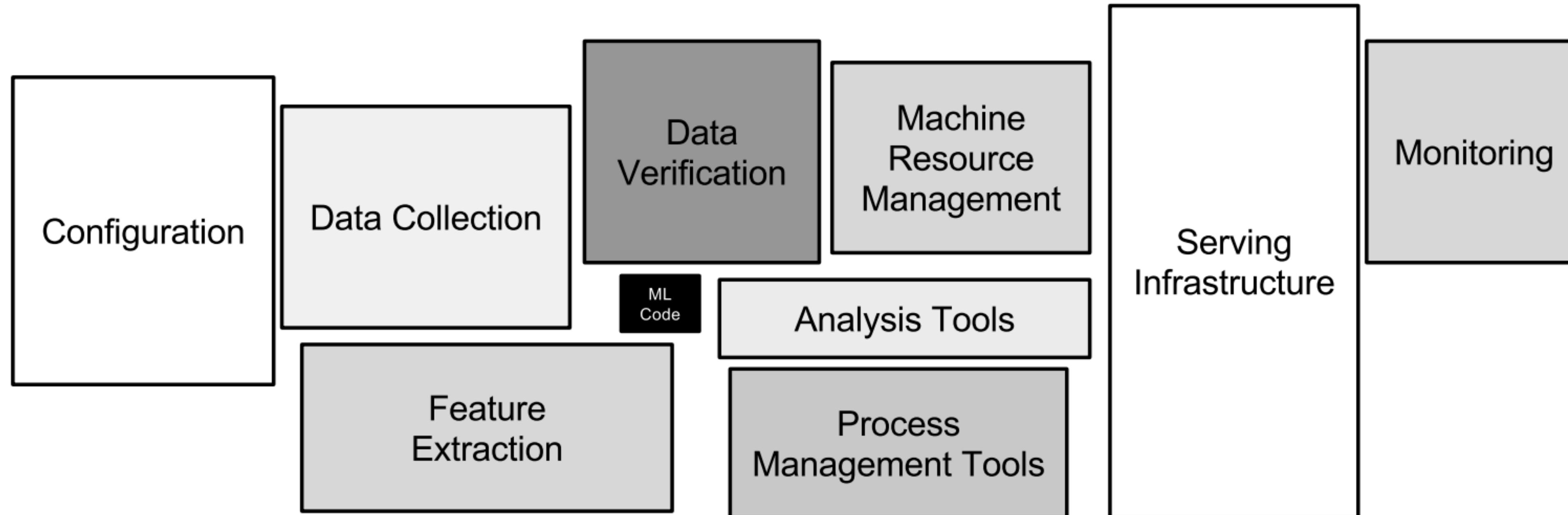


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

Google

A Brief History of “Platformizing” ML

A Brief History of “Platformizing” ML

1980s



S

A Brief History of “Platformizing” ML



Mid
1990s

1980s

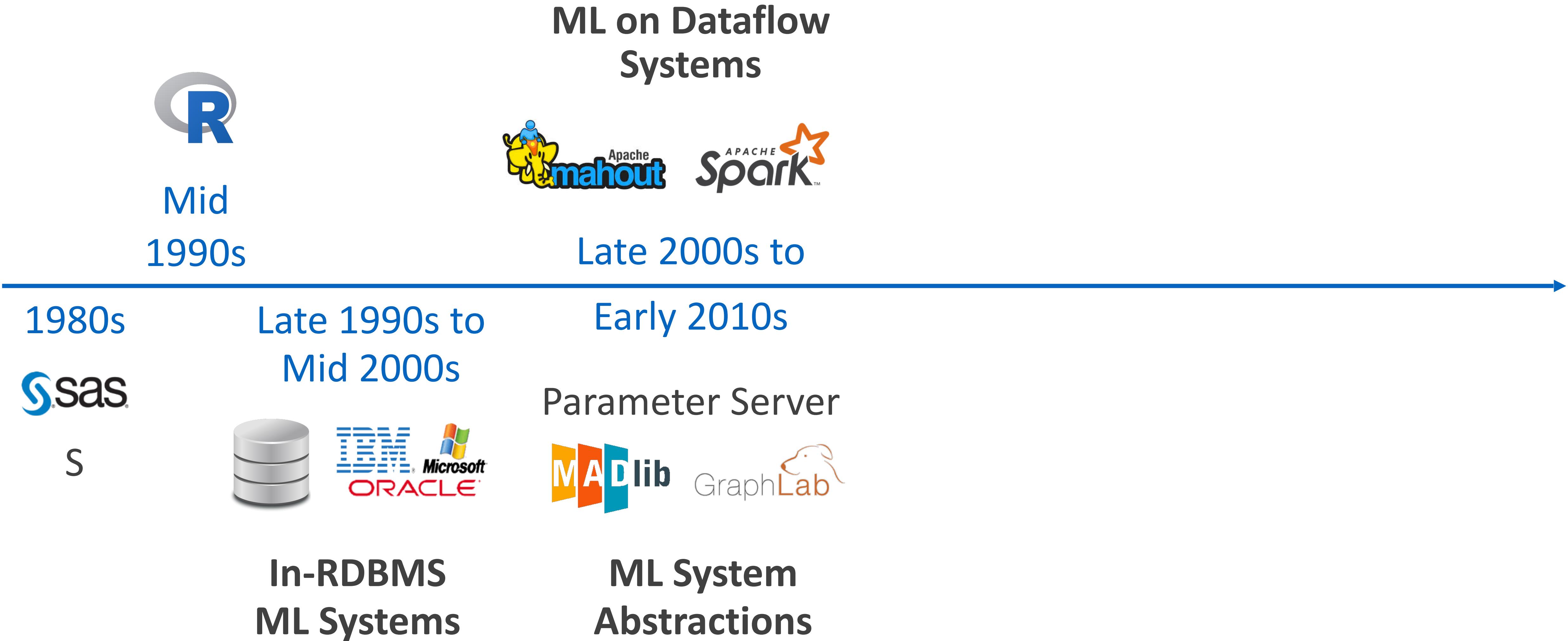


S

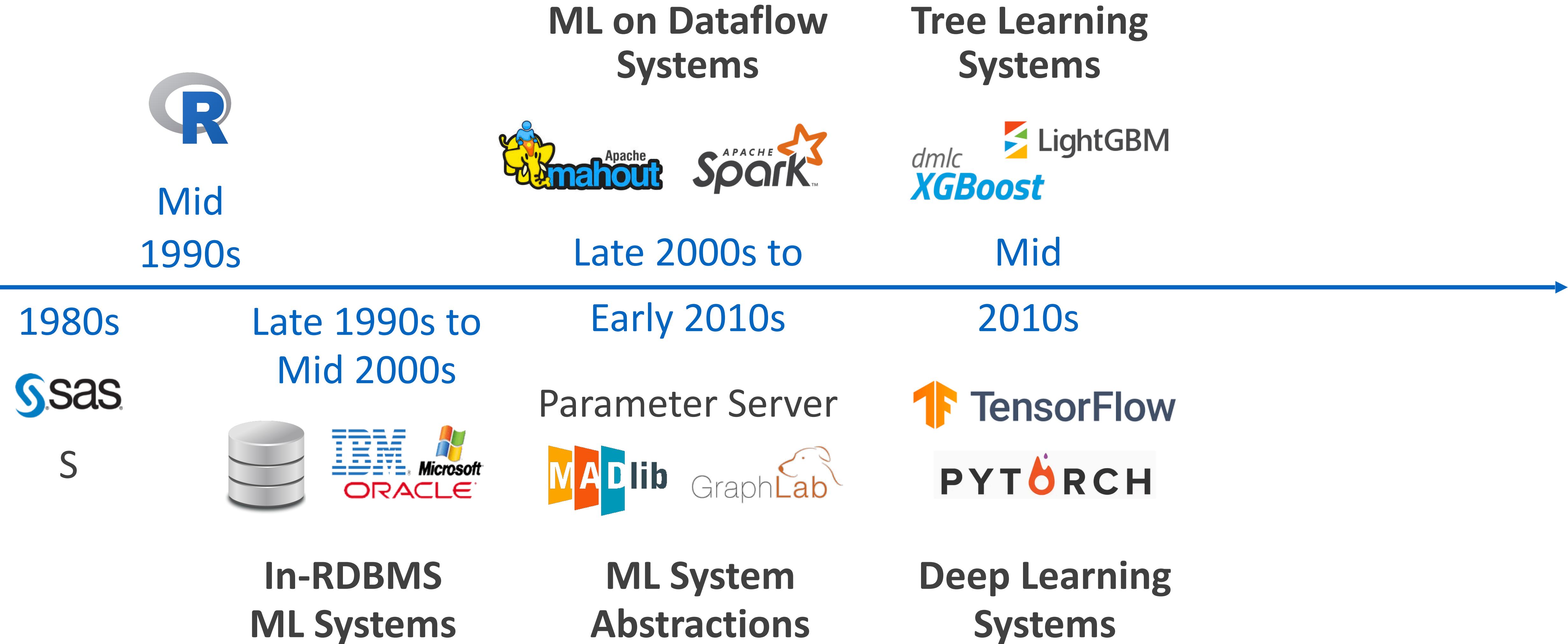
A Brief History of “Platformizing” ML



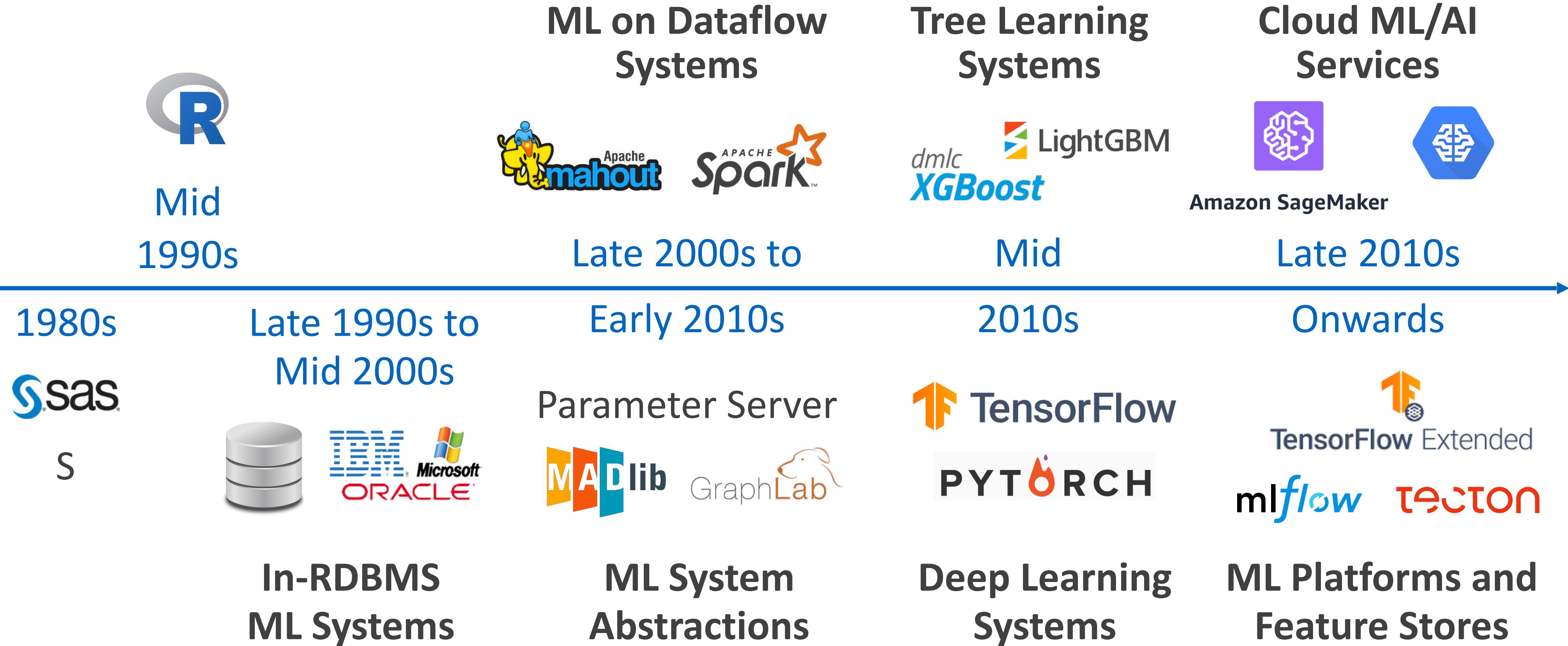
A Brief History of “Platformizing” ML



A Brief History of “Platformizing” ML

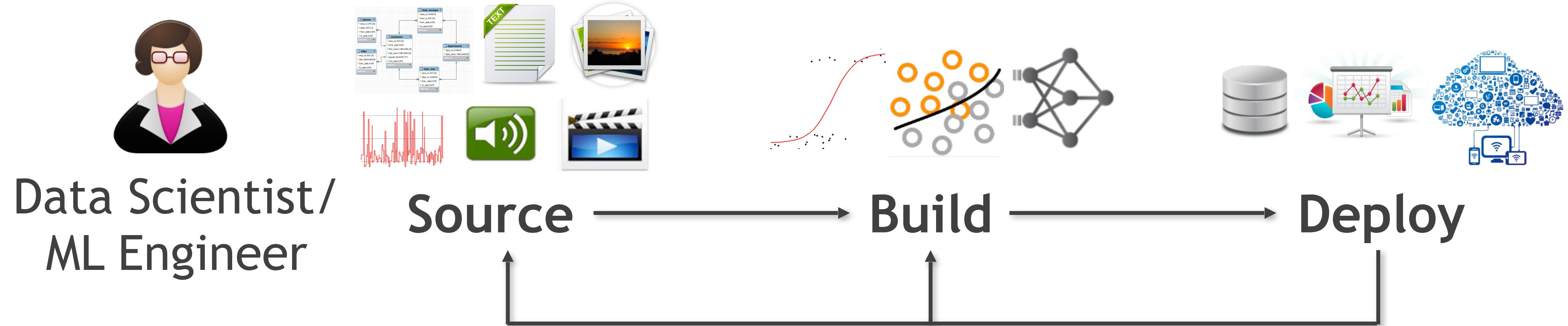


A Brief History of “Platformizing” ML

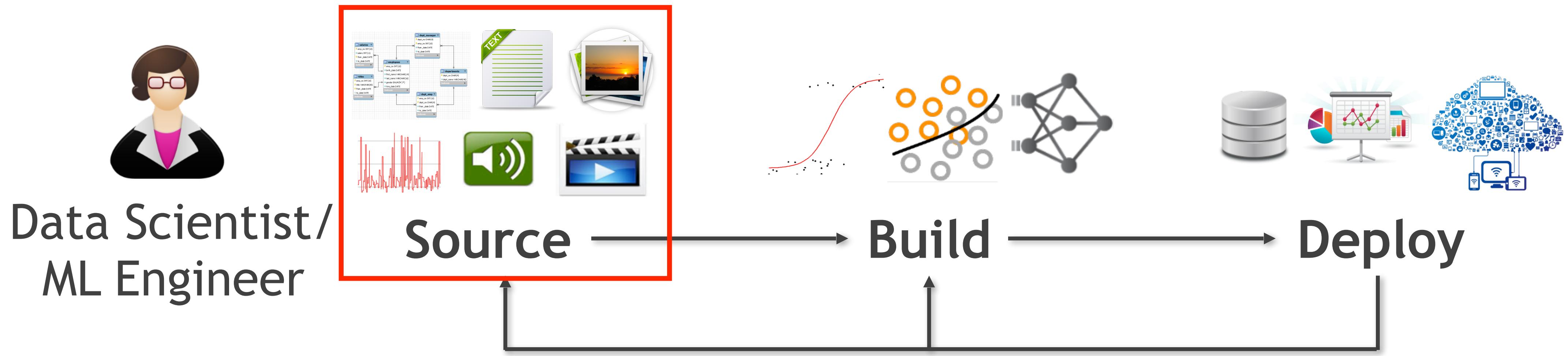


*The way I see it, the rise of ML systems/
platforms today resembles the rise of RDBMSs
circa early 1980s*

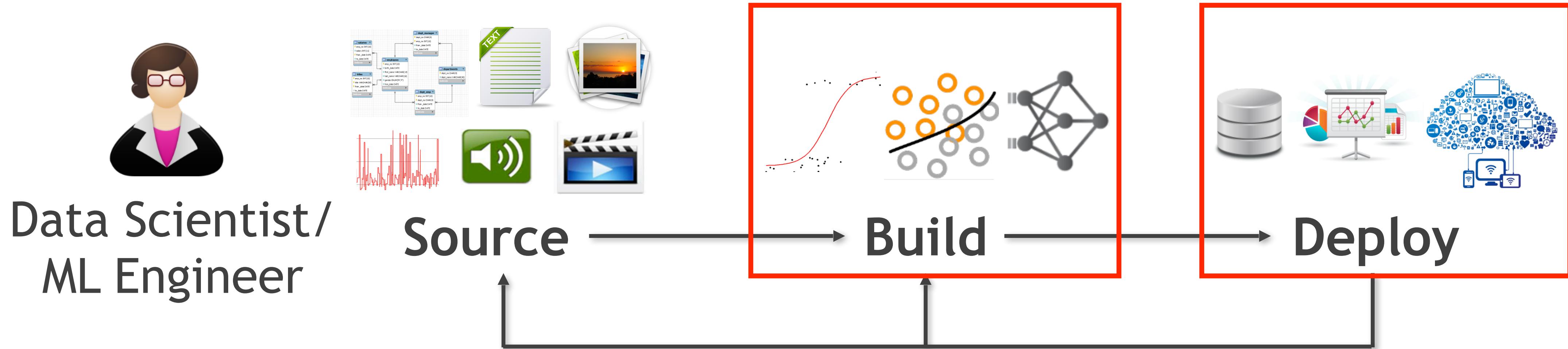
The New DBfication of ML/AI



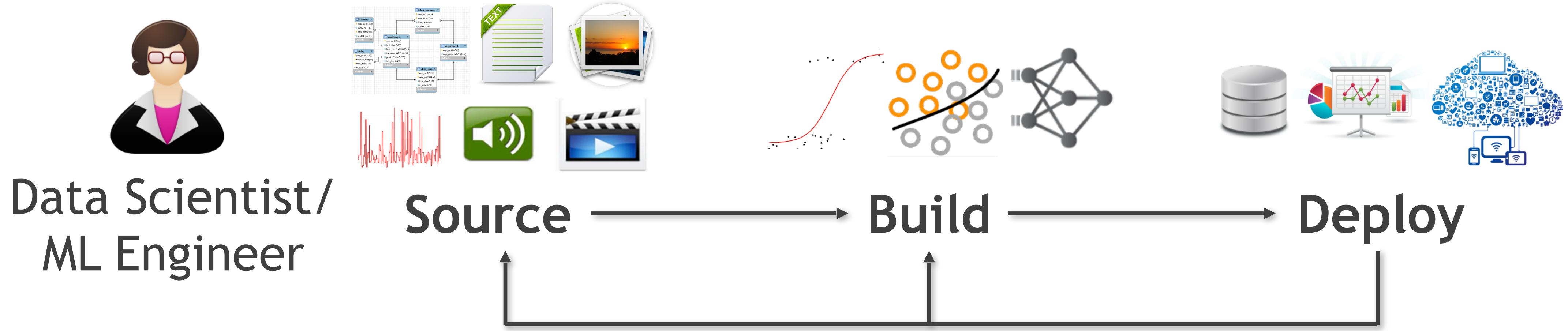
The New DBfication of ML/AI



The New DBfication of ML/AI



The New DBfication of ML/AI



Metadata Management for ML
Data Prep/Cleaning for ML
Multimodal ML Query Models
Data Search, Labeling, etc.

...

Scalable Data Systems for ML
Query Optimization for ML
Cloud and Streaming Infra.
Provenance and Debugging

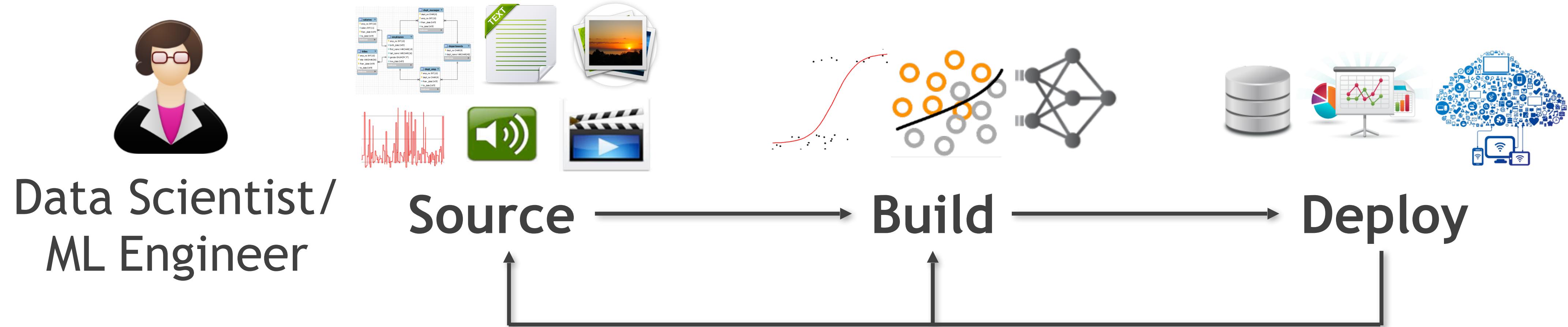
...

Benchmark Frameworks and Data
Fairness, Transparency, Privacy, etc.

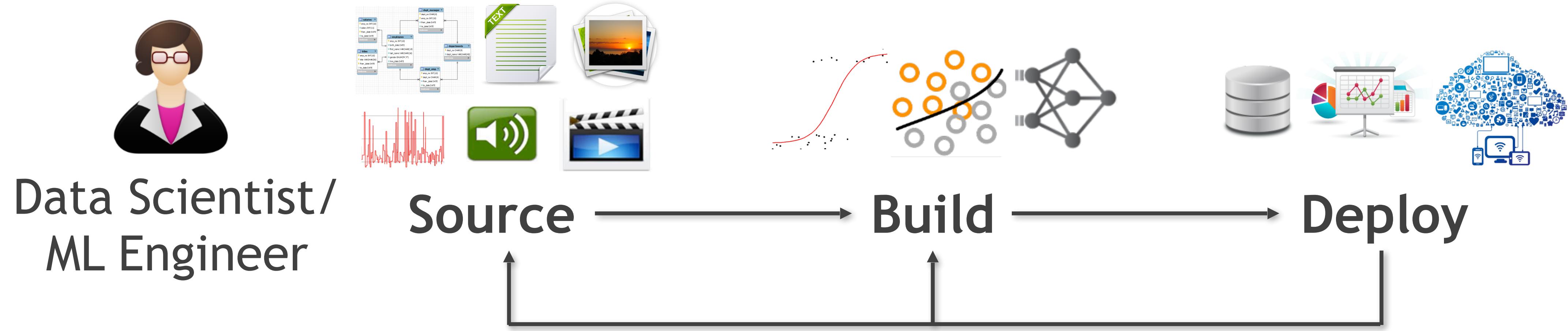
Outline

- | The New DBfication of ML/AI
- | Two Examples from My Research
- | Accelerating the DBfication of ML/AI

The New DBfication of ML/AI



The New DBfication of ML/AI



Metadata Management for ML
Data Prep/Cleaning for ML
Multimodal ML Query Models
Data Search, Labeling, etc.

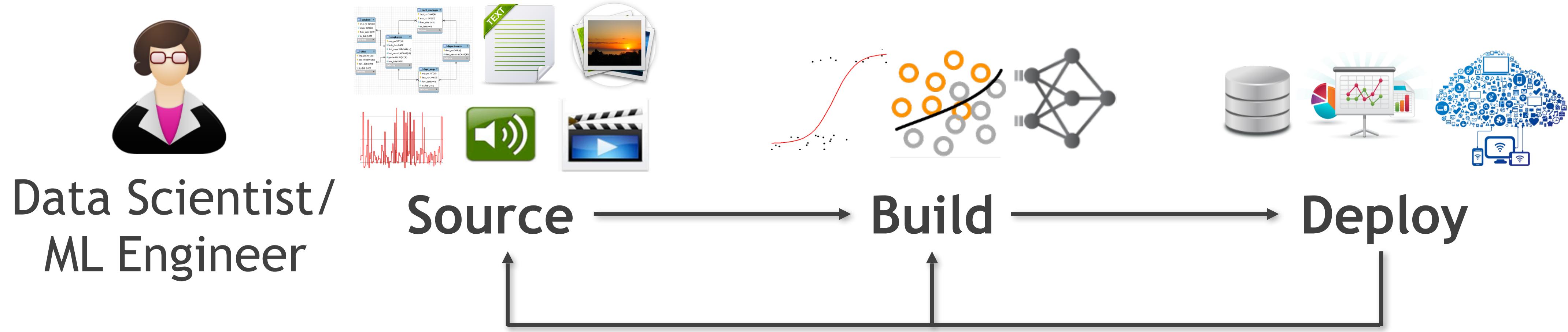
...

Scalable Data Systems for ML
Query Optimization for ML
Cloud and Streaming Infra.
Provenance and Debugging

...

Benchmark Frameworks and Data
Fairness, Transparency, Privacy, etc.

The New DBfication of ML/AI



Metadata Management for ML
Data Prep/Cleaning for ML
Multimodal ML Query Models
Data Search, Labeling, etc.

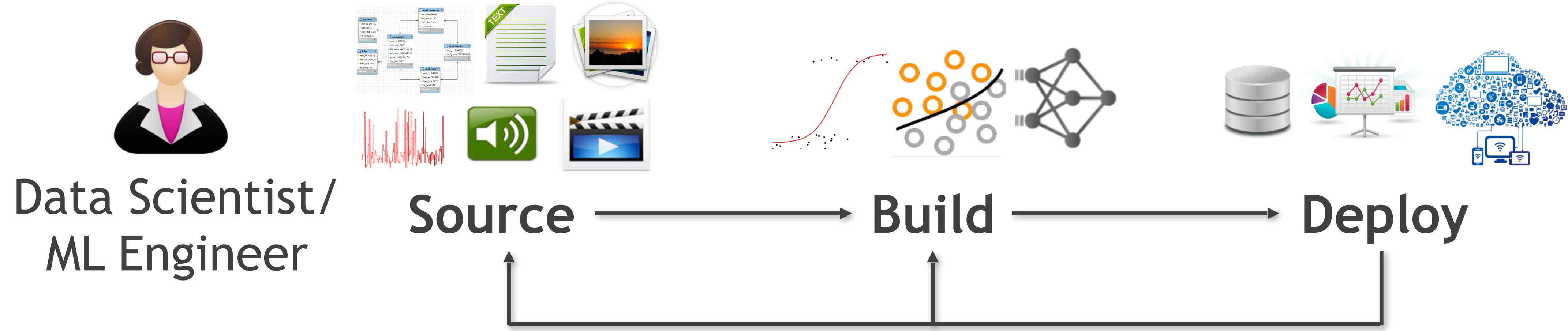
Scalable Data Systems for ML
Query Optimization for ML
Cloud and Streaming Infra.
Provenance and Debugging

...

Benchmark Frameworks and Data
Fairness, Transparency, Privacy, etc.

...

The New DBfication of ML/AI



Metadata Management for ML
Data Prep/Cleaning for ML
Multimodal ML Query Models
Data Search, Labeling, etc.

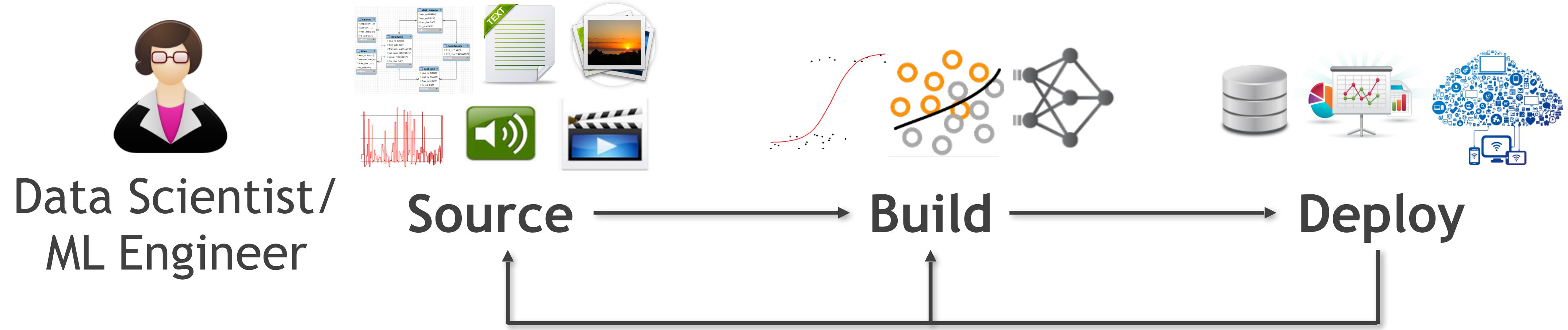
...

Scalable Data Systems for ML
Query Optimization for ML
Cloud and Streaming Infra.
Provenance and Debugging

...

Benchmark Frameworks and Data
Fairness, Transparency, Privacy, etc.

The New DBfication of ML/AI



Metadata Management for ML
Data Prep/Cleaning for ML
Multimodal ML Query Models
Data Search, Labeling, etc.

...

Scalable Data Systems for ML
Query Optimization for ML
Cloud and Streaming Infra.
Provenance and Debugging

...

Benchmark Frameworks and Data
Fairness, Transparency, Privacy, etc.

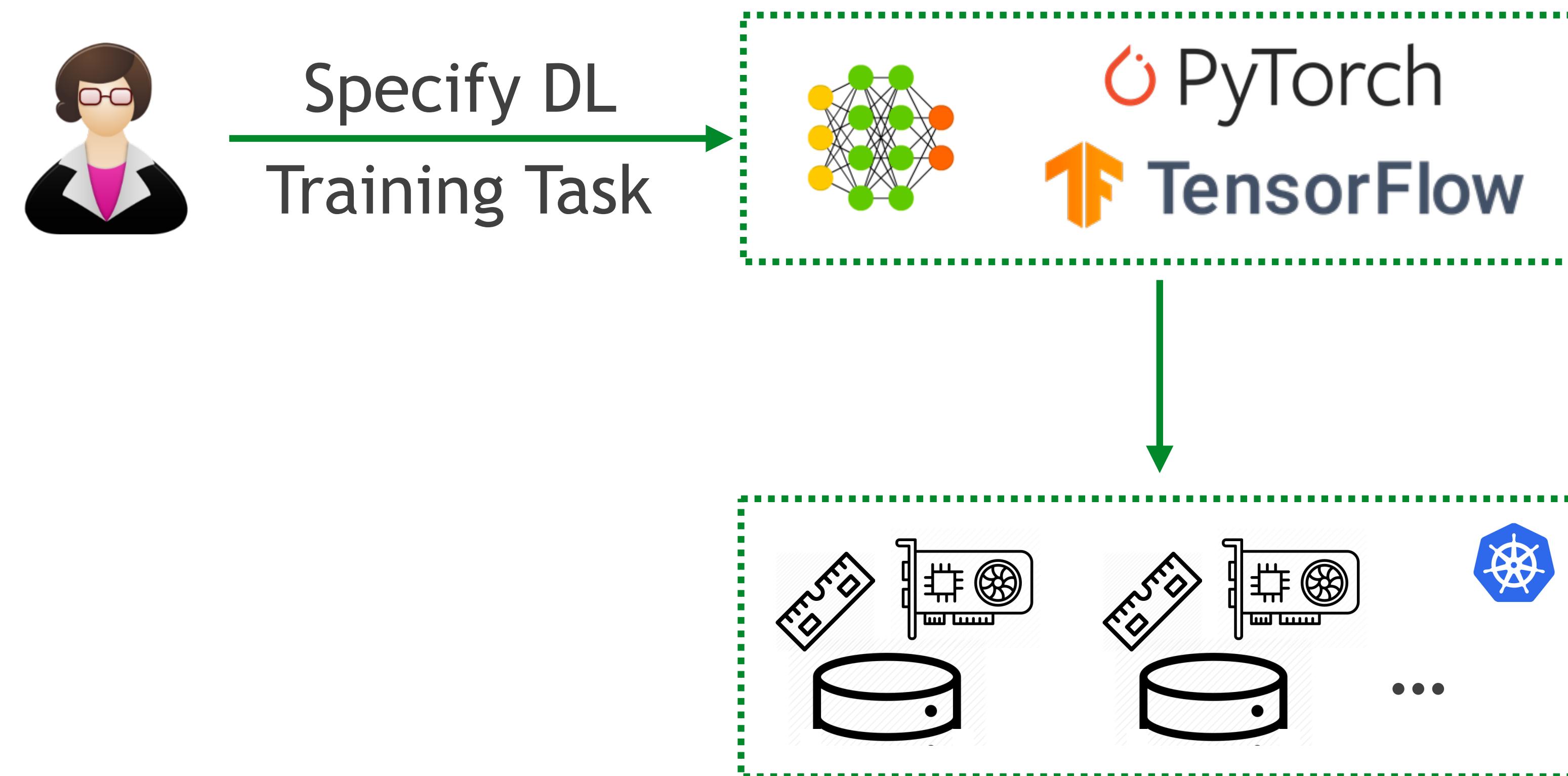
Outline

- | The New DBfication of ML/AI
- | Two Examples from My Research
 - | Example 1: Scalable DL Systems
 - | Example 2: Auto Data Prep for ML
- | Accelerating the DBfication of ML/AI

Outline

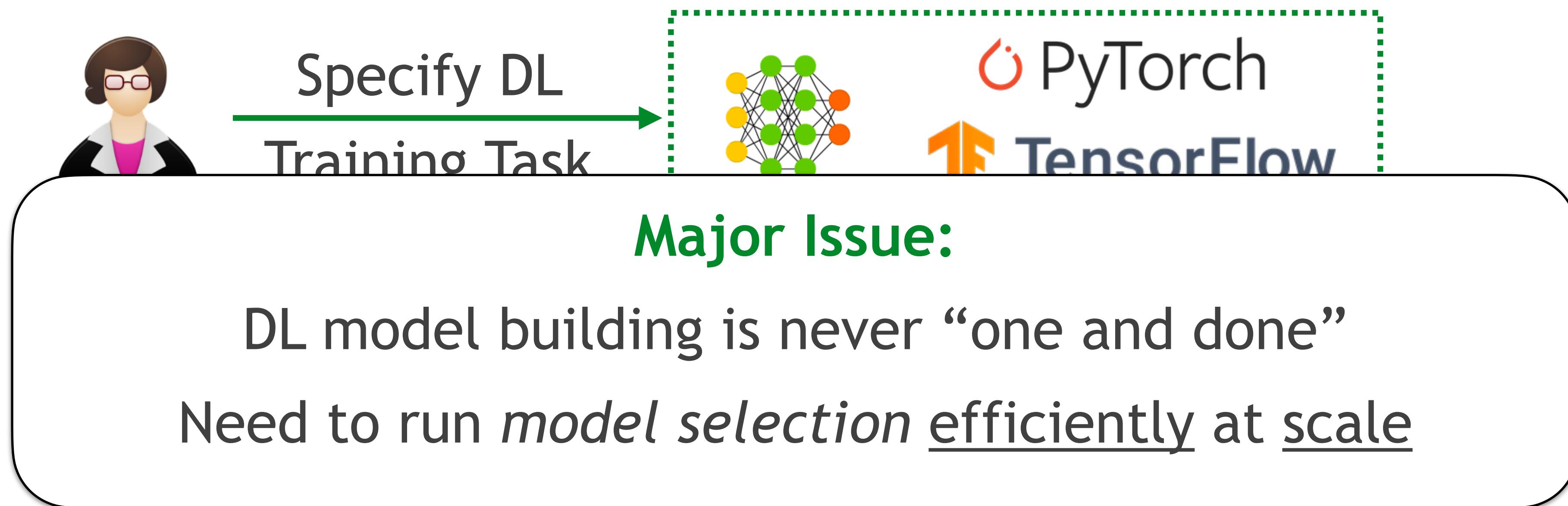
- | The New DBfication of ML/AI
- | Two Examples from My Research
 - | Example 1: Scalable DL Systems
 - | Example 2: Auto Data Prep for ML
- | Accelerating the DBfication of ML/AI

Example 1: DL Systems at Scale



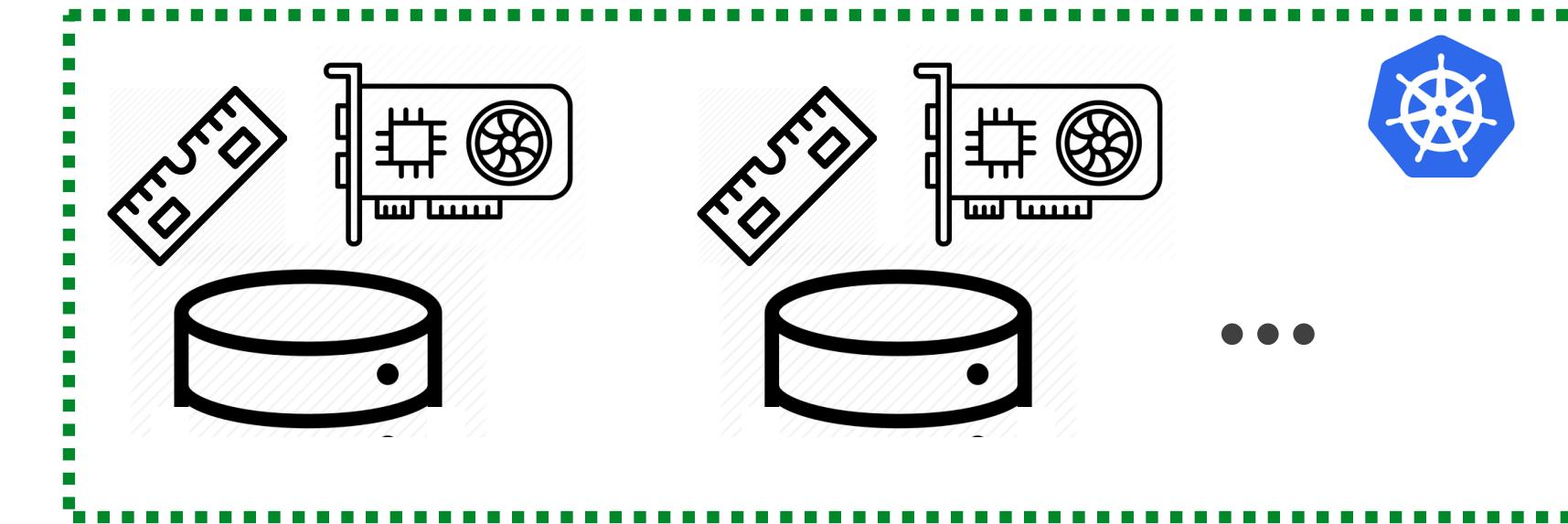
Cluster: GPU+CPU, memory, disks

Example 1: DL Systems at Scale



Cluster: GPU+CPU, memory, disks

Project Cerebro

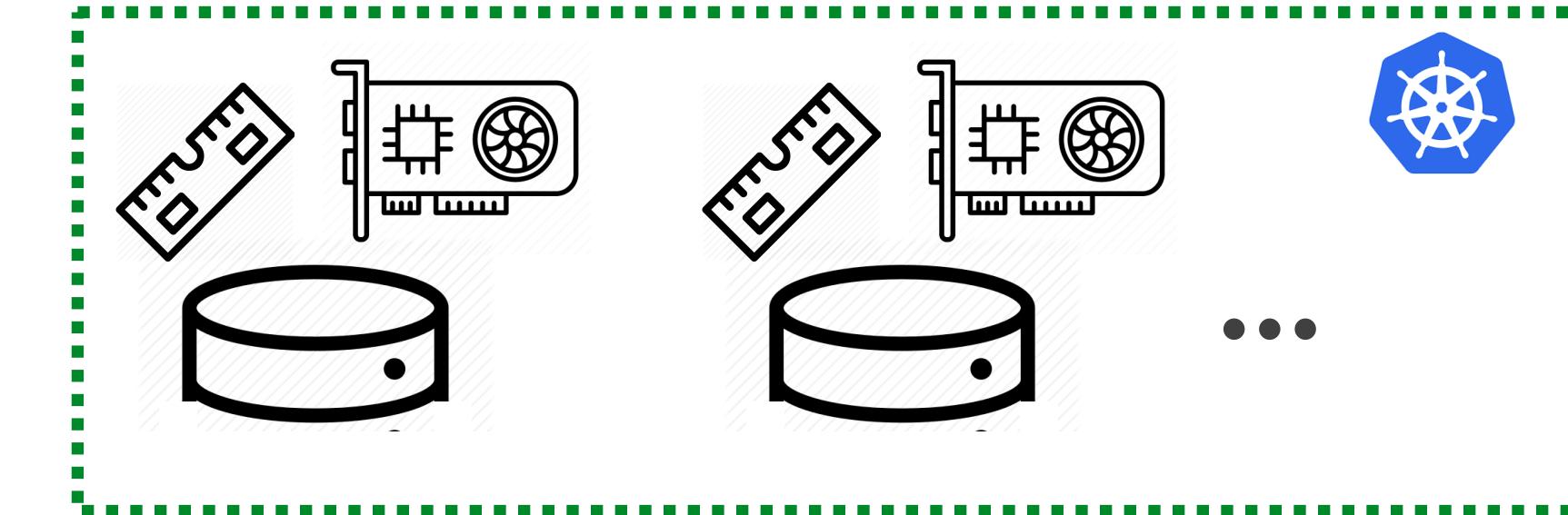


Plain/Kubernetes cluster

Project Cerebro

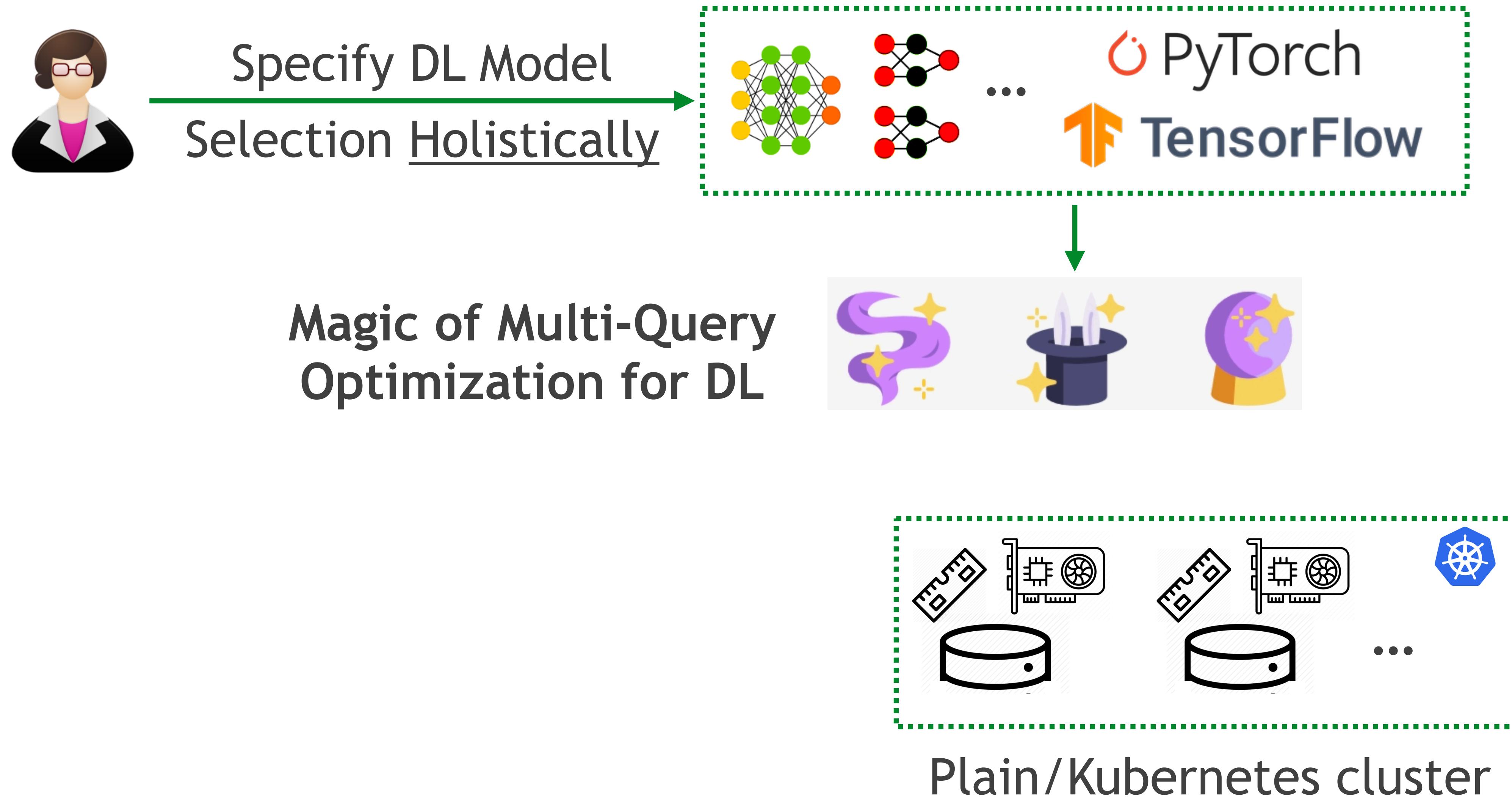


Specify DL Model
Selection Holistically

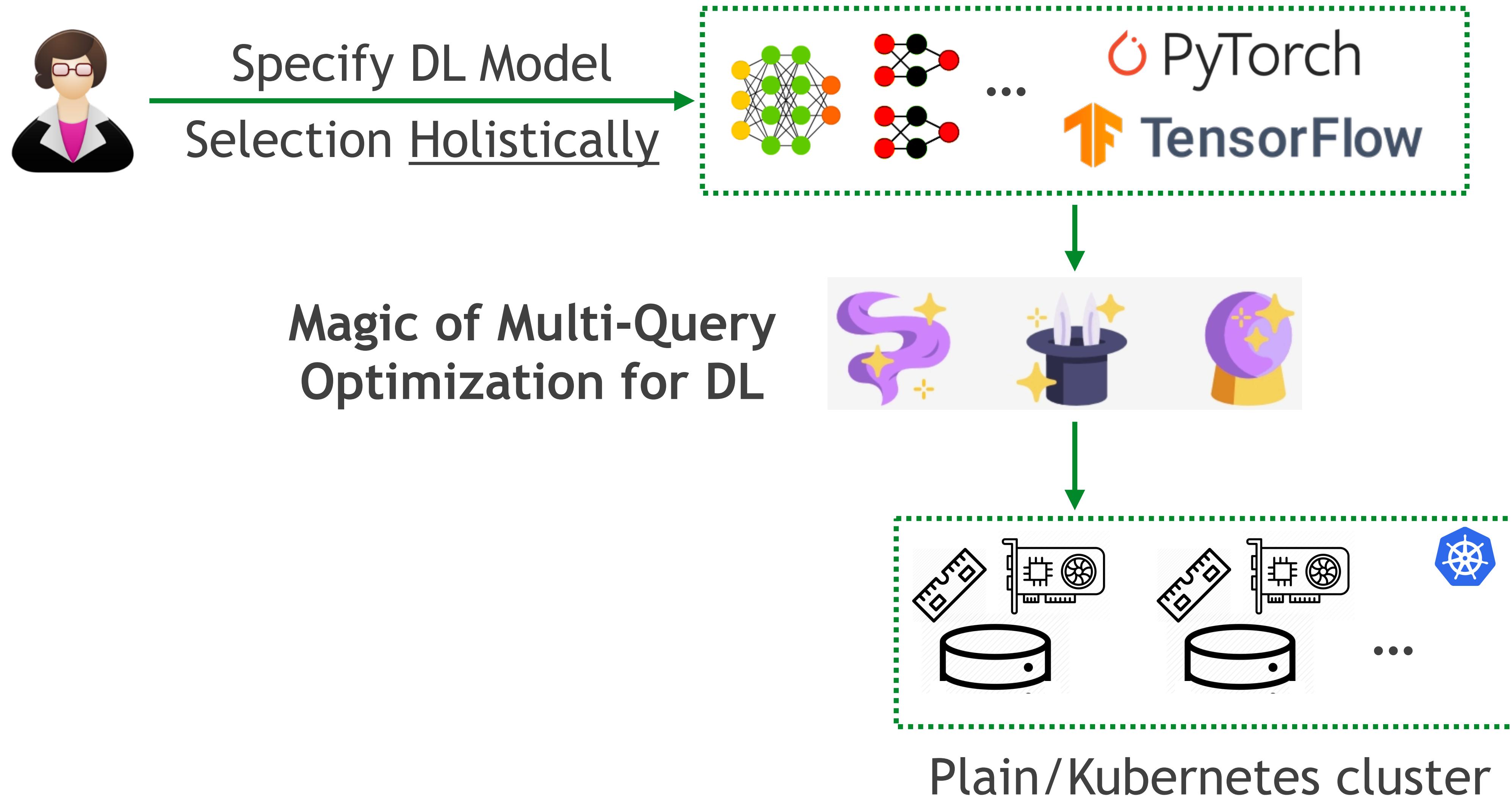


Plain/Kubernetes cluster

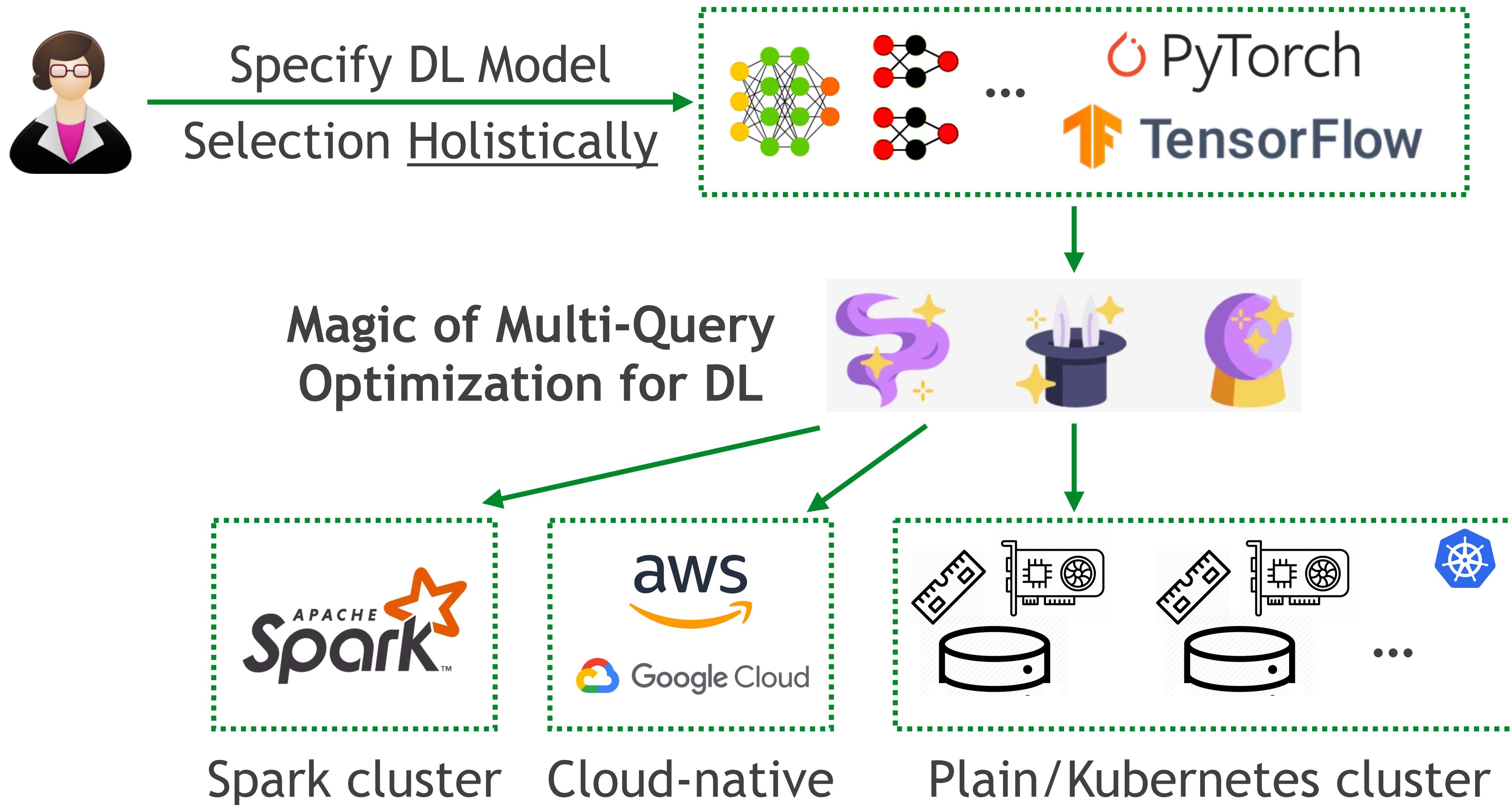
Project Cerebro



Project Cerebro



Project Cerebro



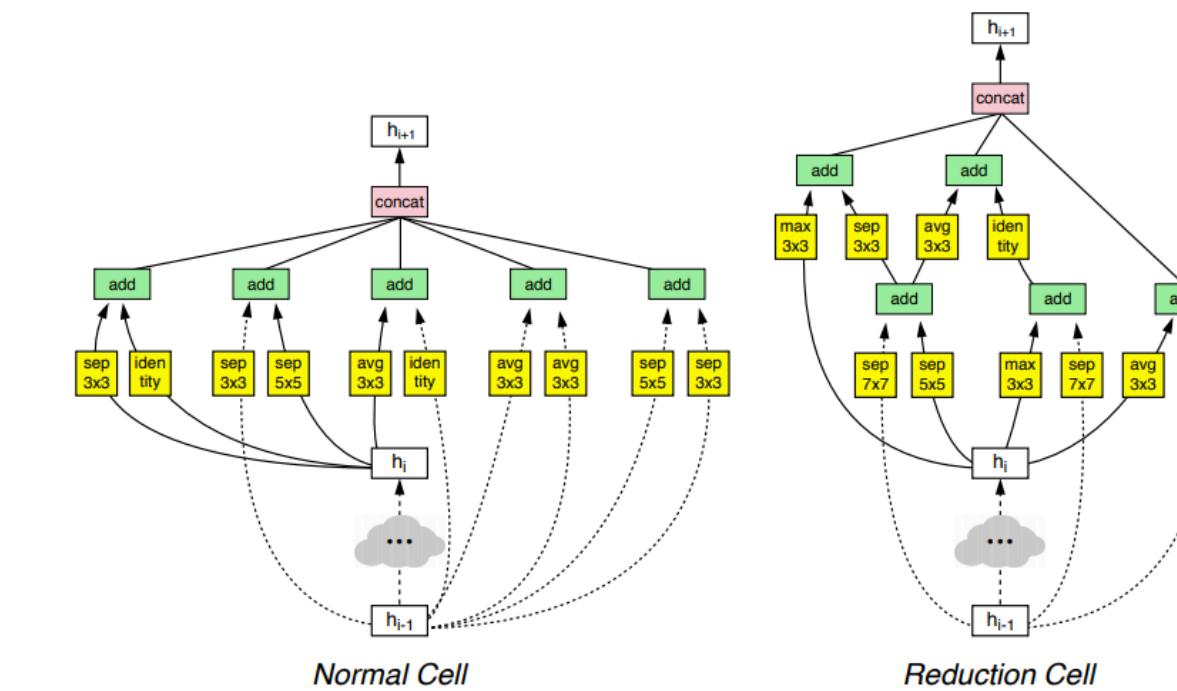
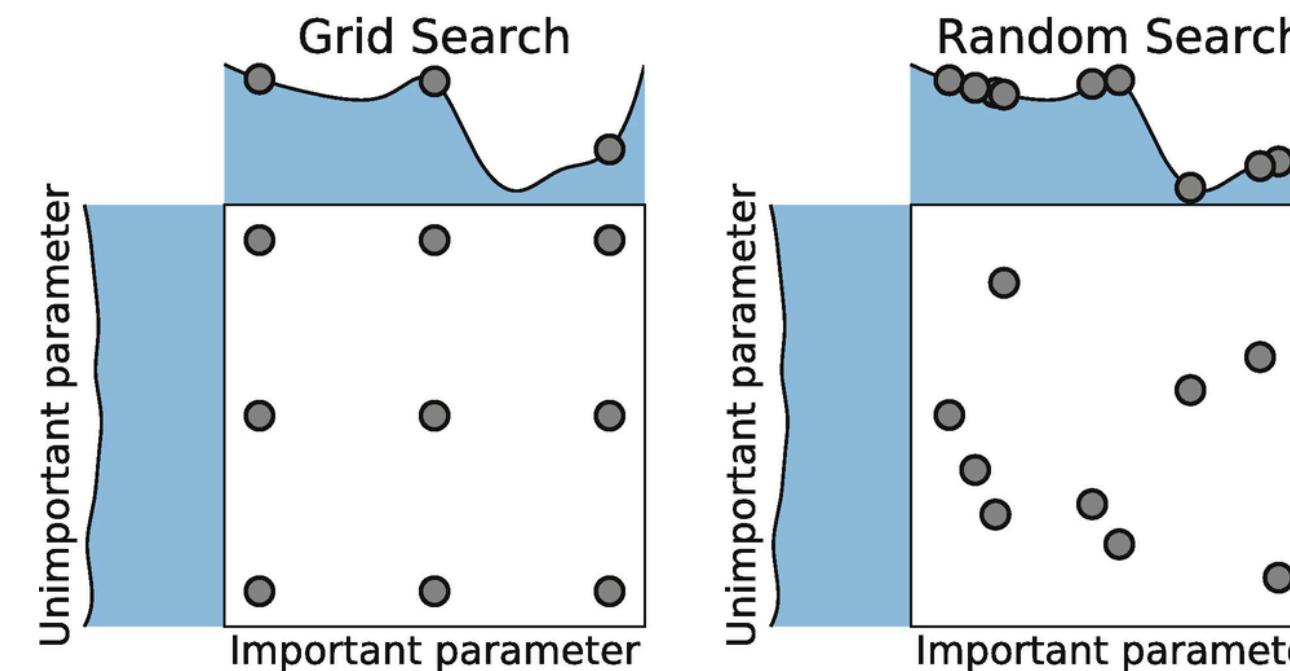
Key Example: MQO for DL Model Selection

Cerebro: Efficient and Reproducible Model Selection on Deep Learning Systems. **SIGMOD DEEM'19**
Cerebro: A Data System for Optimized Deep Learning Model Selection. **VLDB'20**

Key Example: MQO for DL Model Selection

Workload:

Hyper-parameter tuning and neural architecture exploration yield numerous (even 100s) of DL training configs

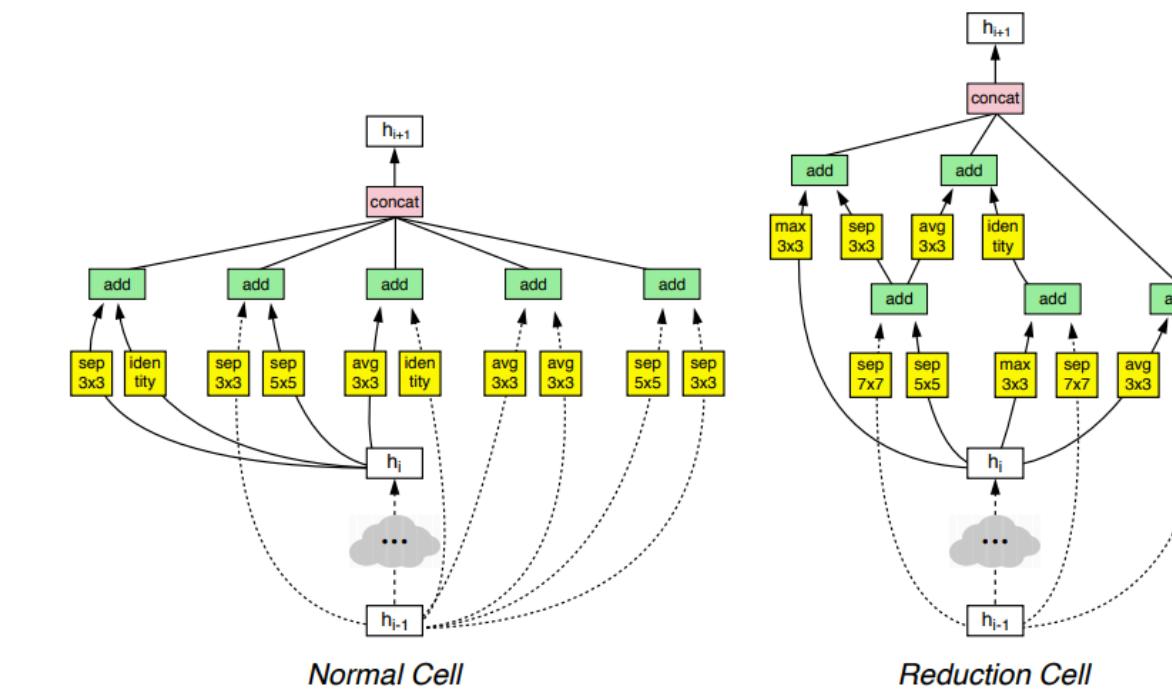
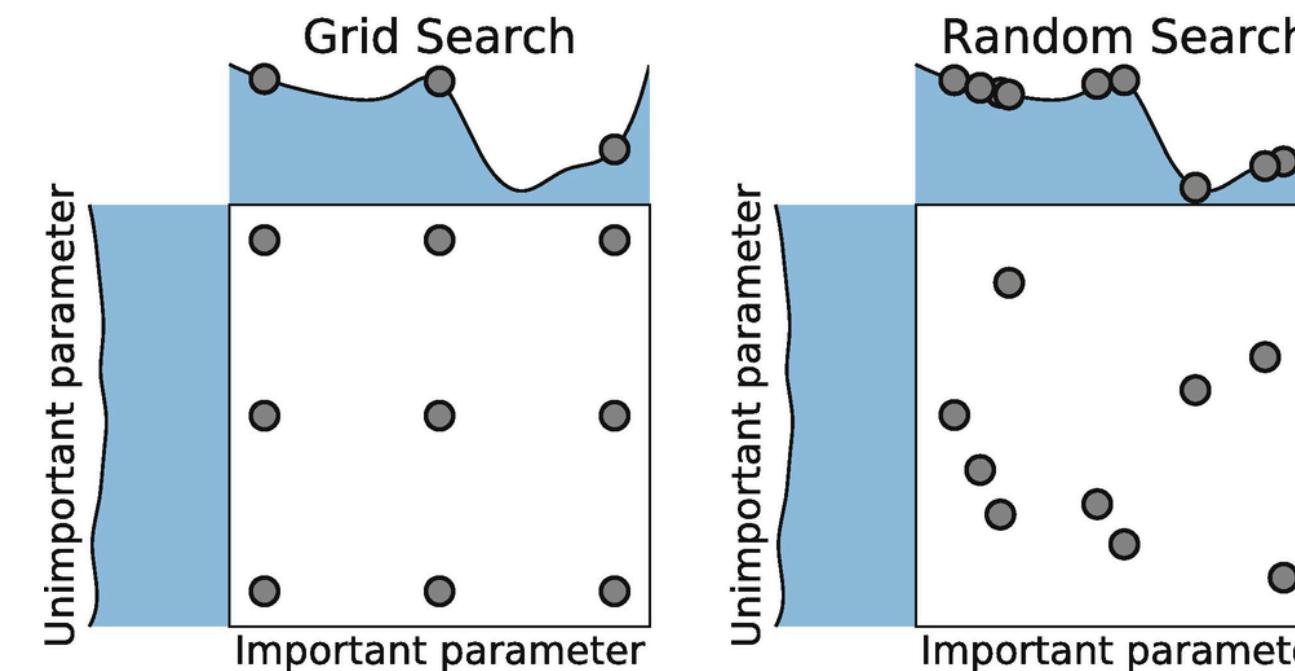


HyperBand,
AutoKeras,
ASHA, PBT,
HyperOpt,
NAS, etc.

Key Example: MQO for DL Model Selection

Workload:

Hyper-parameter tuning and neural architecture exploration yield numerous (even 100s) of DL training configs



HyperBand,
AutoKeras,
ASHA, PBT,
HyperOpt,
NAS, etc.

Multi-Query:

Run multiple configs on large dataset in one go

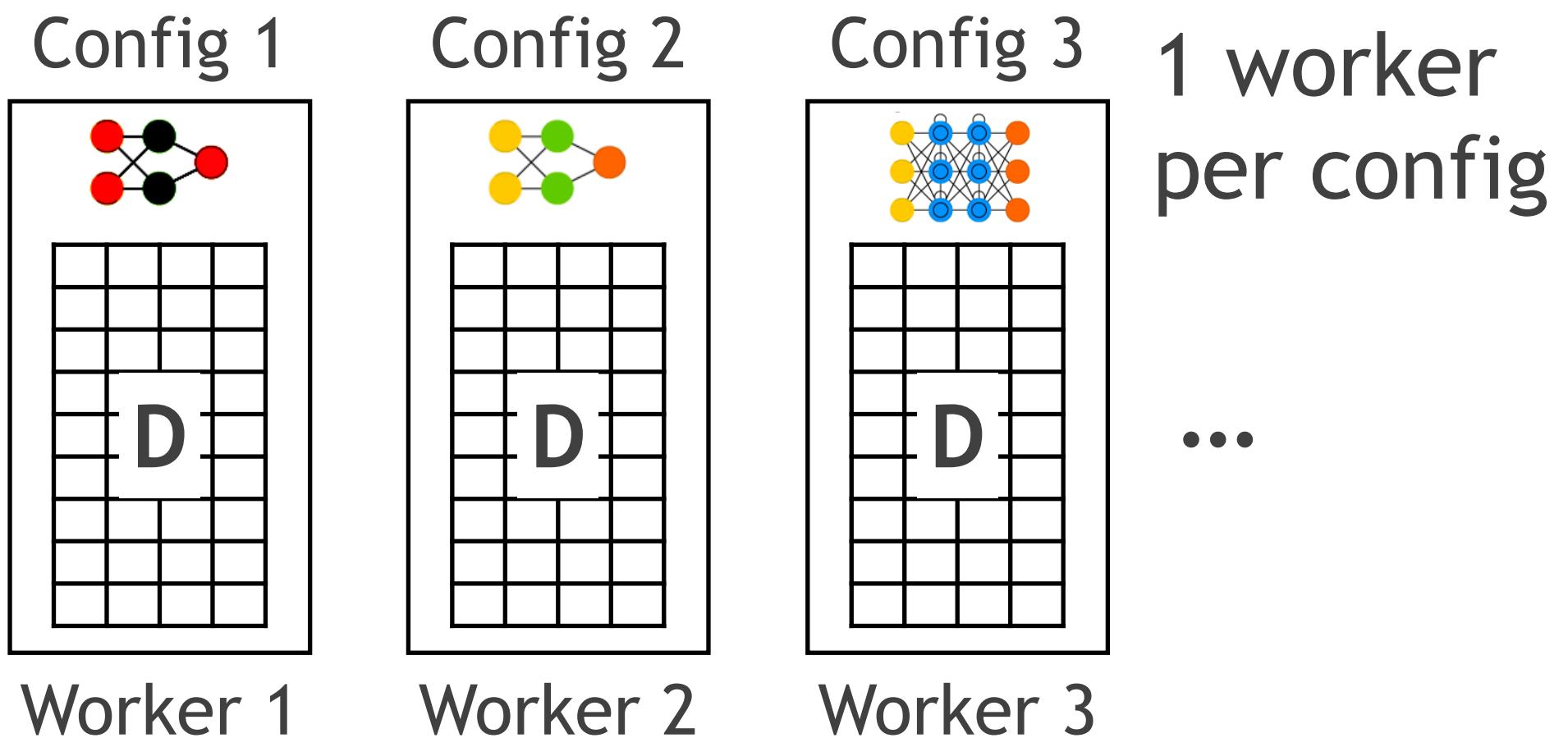
Positioning Cerebro's Technique vs Prior Art

We devise a novel execution strategy called **Model Hopper Parallelism (MOP)**

Positioning Cerebro's Technique vs Prior Art

We devise a novel execution strategy called **Model Hopper Parallelism (MOP)**

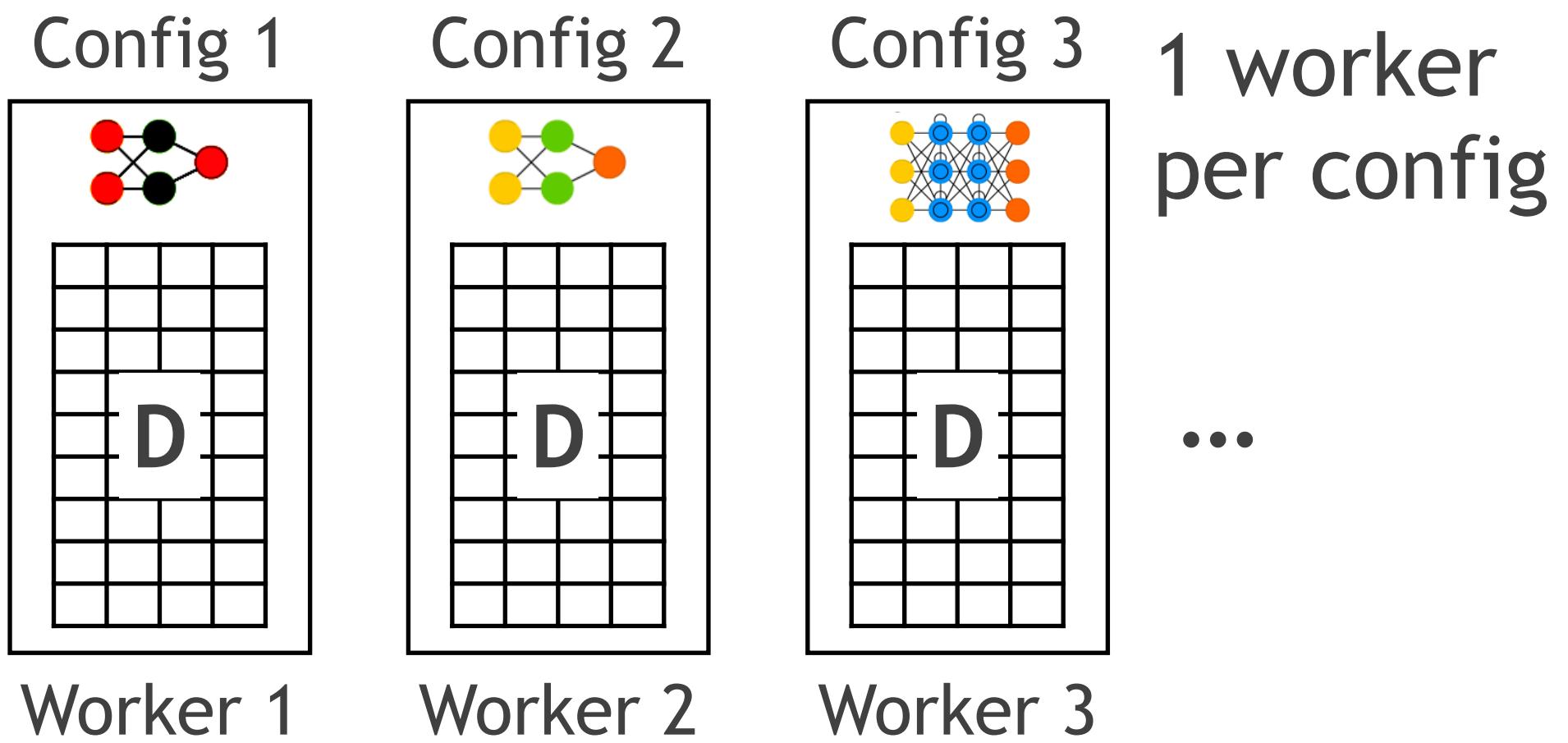
Task Parallelism:



Positioning Cerebro's Technique vs Prior Art

We devise a novel execution strategy called **Model Hopper Parallelism (MOP)**

Task Parallelism:

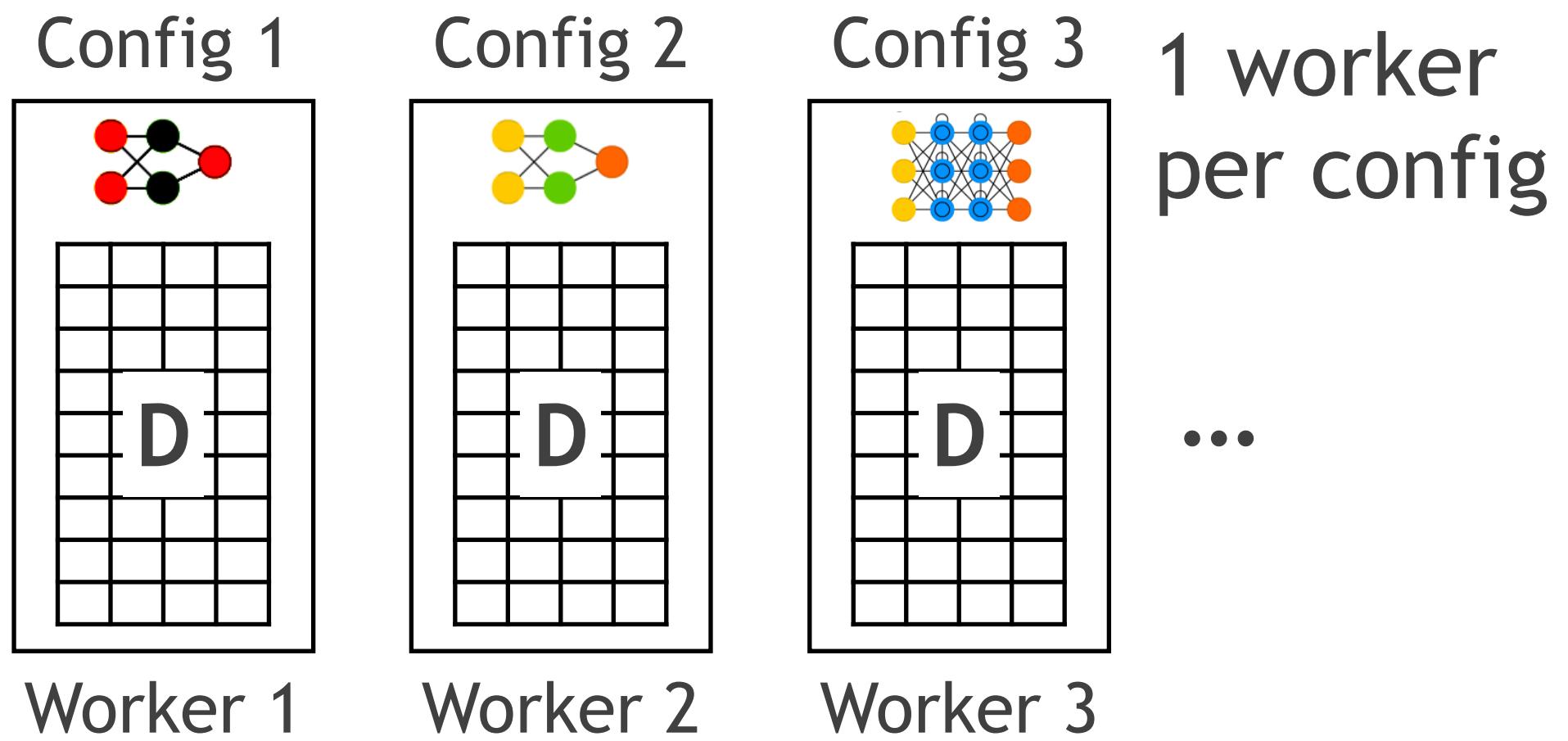


- + High throughput model selection
- + Best accuracy from Sequential SGD
- Low data scalability; wastes space
(copy) or network (remote read)

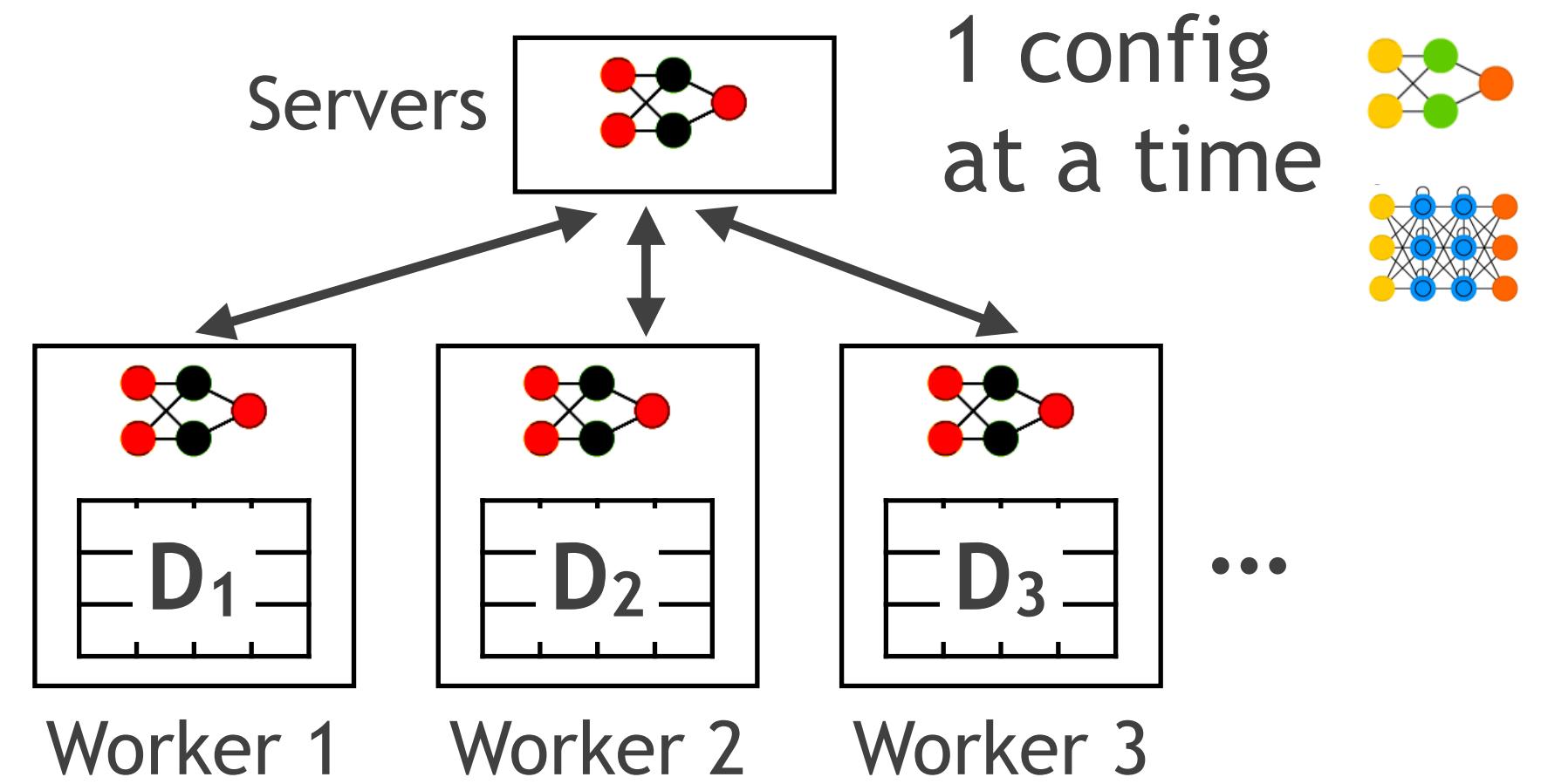
Positioning Cerebro's Technique vs Prior Art

We devise a novel execution strategy called **Model Hopper Parallelism (MOP)**

Task Parallelism:



Data Parallelism:

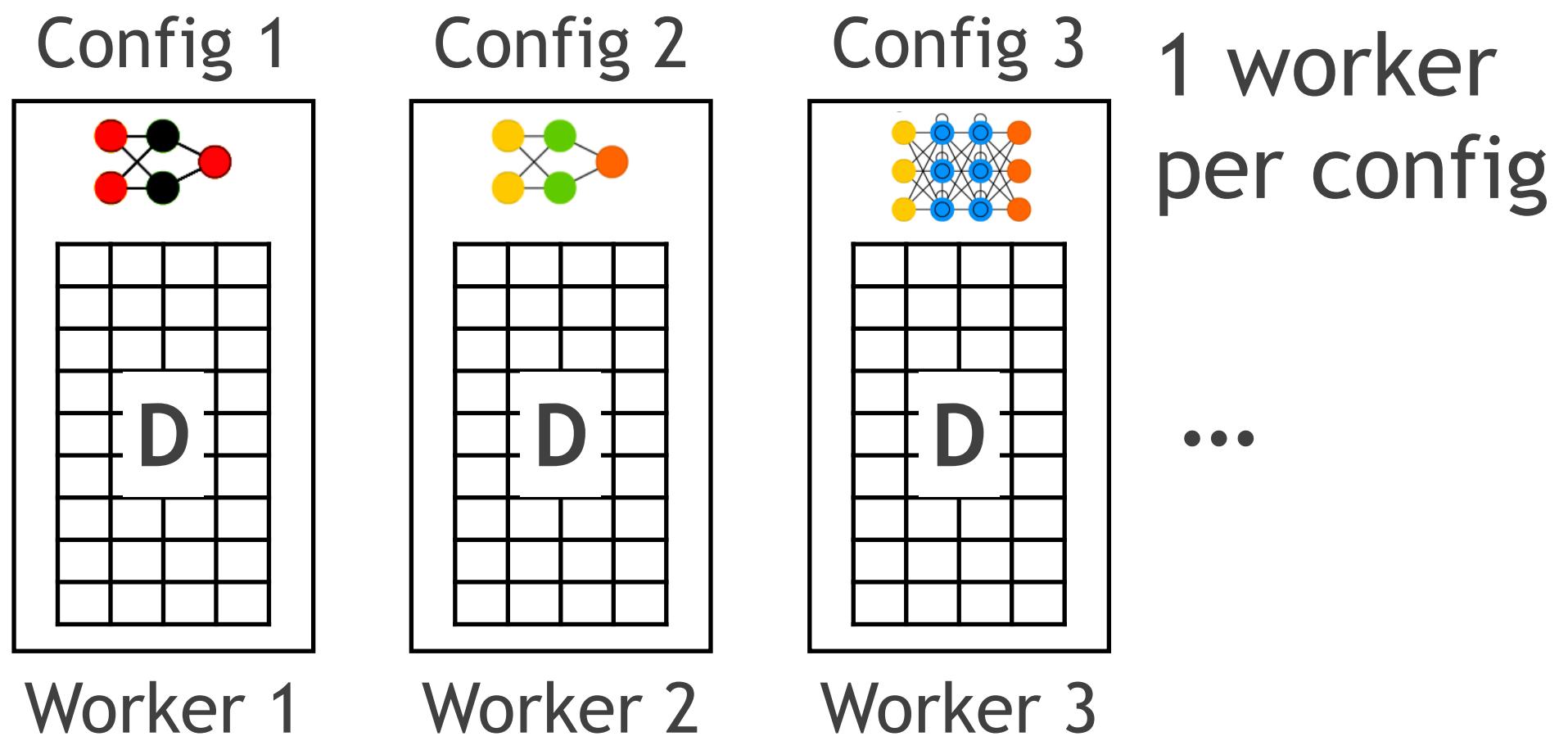


- + High throughput model selection
- + Best accuracy from Sequential SGD
- Low data scalability; wastes space (copy) or network (remote read)

Positioning Cerebro's Technique vs Prior Art

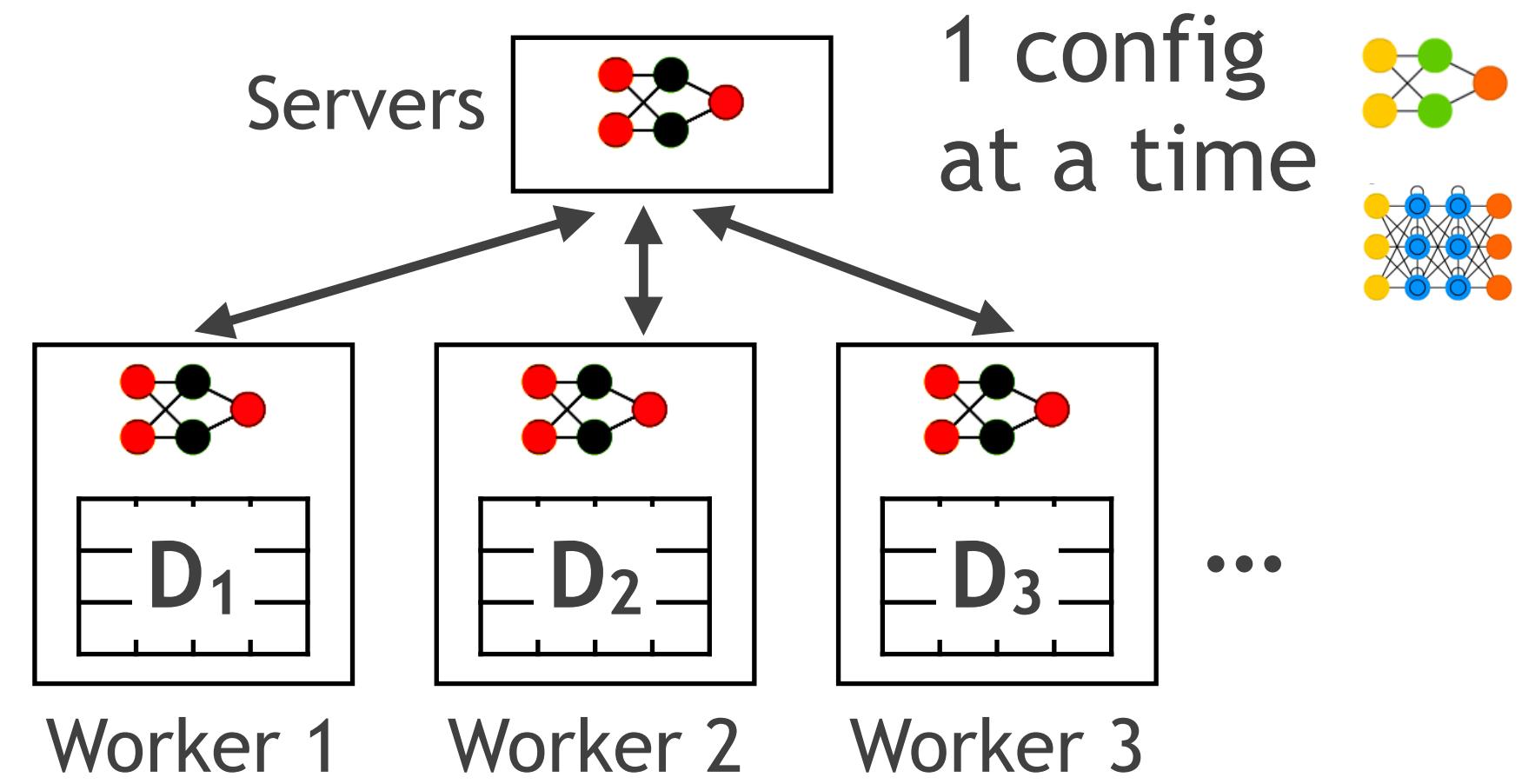
We devise a novel execution strategy called **Model Hopper Parallelism (MOP)**

Task Parallelism:



- + High throughput model selection
- + Best accuracy from Sequential SGD
- Low data scalability; wastes space (copy) or network (remote read)

Data Parallelism:

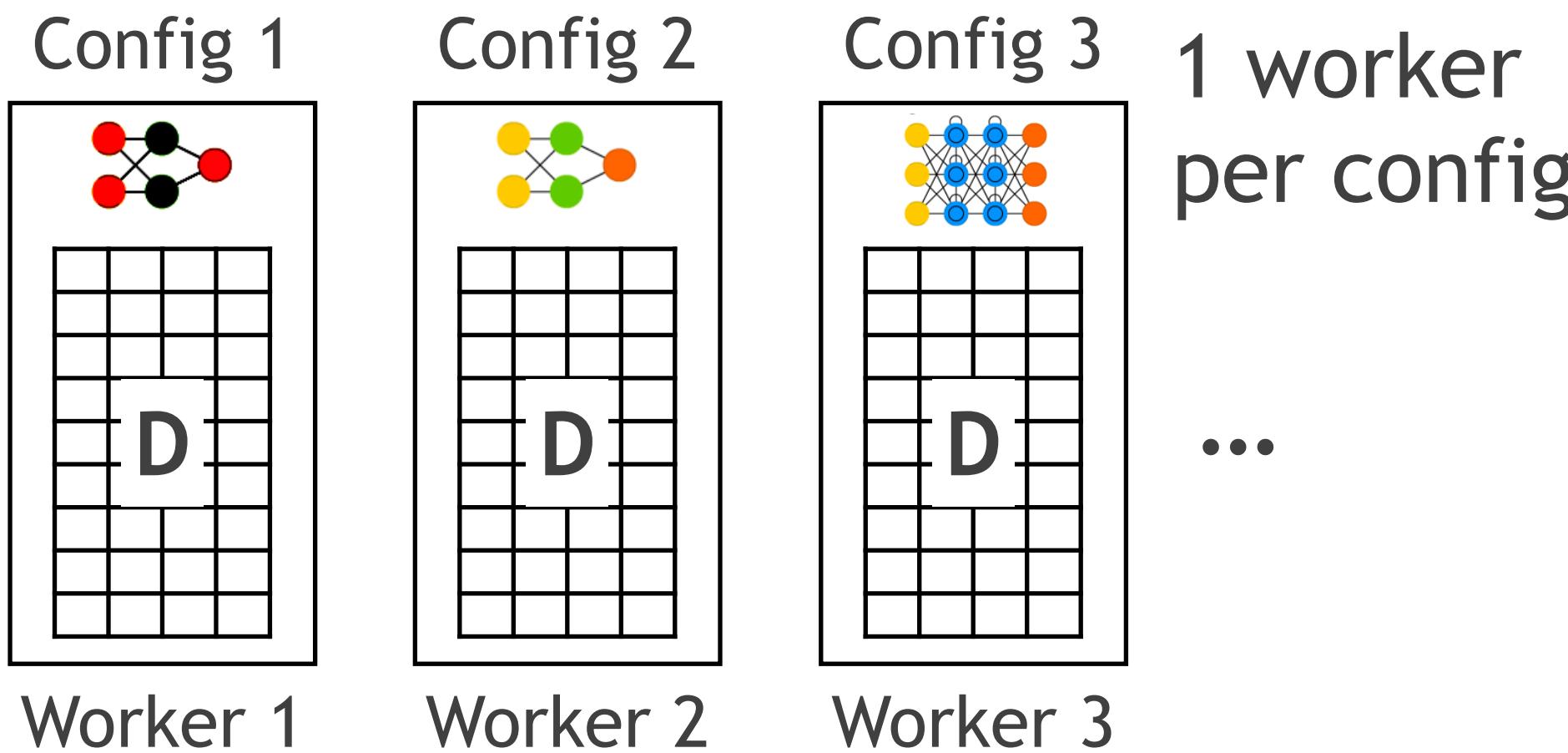


- + High data scalability via sharding
- Shard-level SGD does not converge; mini-batch level is too slow
- Very low throughput overall

Positioning Cerebro's Technique vs Prior Art

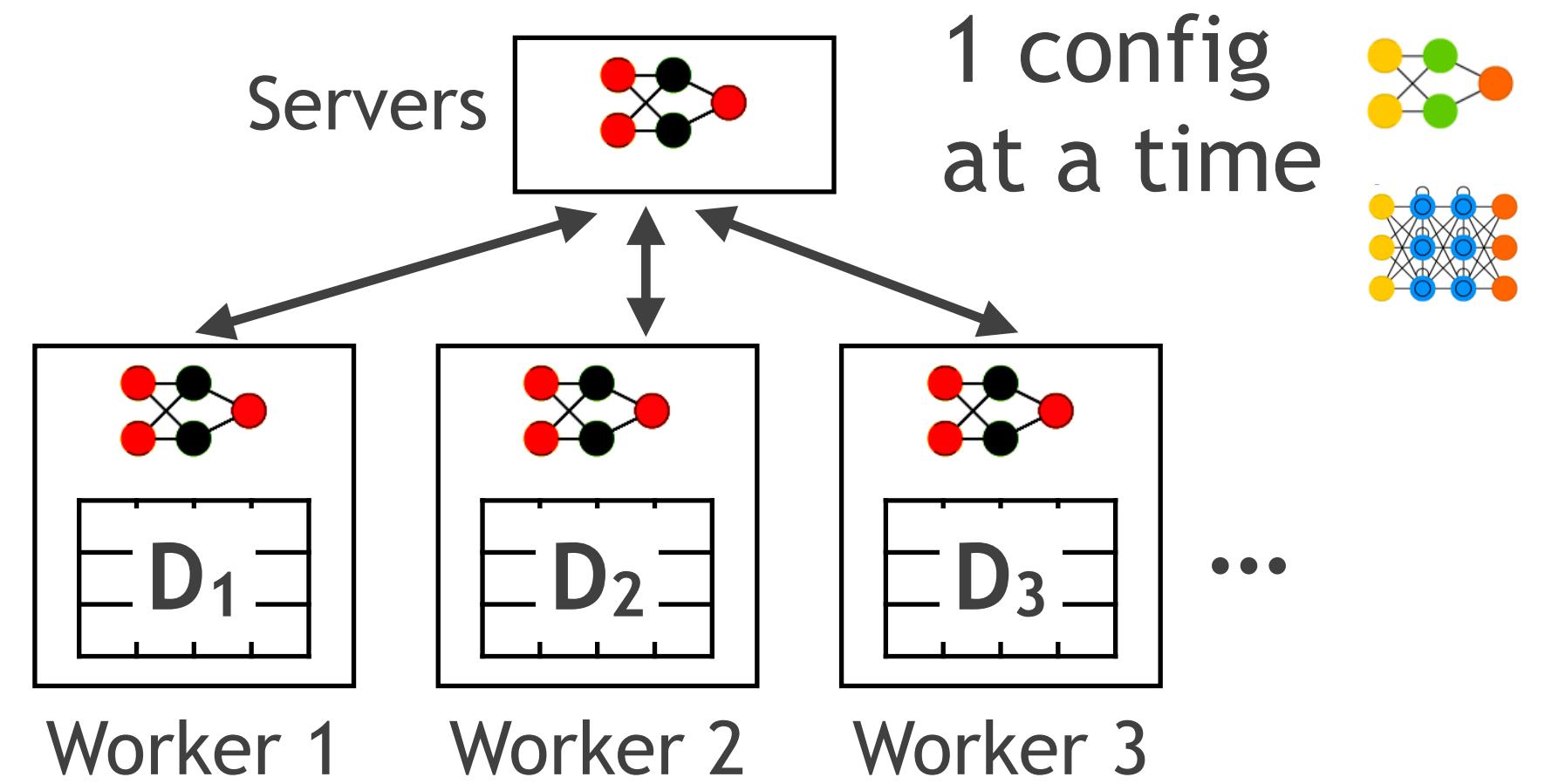
We devise a novel execution strategy called **Model Hopper Parallelism (MOP)**

Task Parallelism:



- + High throughput model selection
- + Best accuracy from Sequential SGD
- Low data scalability; wastes space (copy) or network (remote read)

Data Parallelism:



MOP:

New hybrid of data and task parallelism
Best of both worlds!

- + High data scalability via sharding
- Shard-level SGD does not converge; mini-batch level is too slow
- Very low throughput overall



Model Hopper Parallelism (MOP)

Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

Model Hopper Parallelism (MOP)

Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

Shuffle and partition dataset

Run n DNNs on n workers

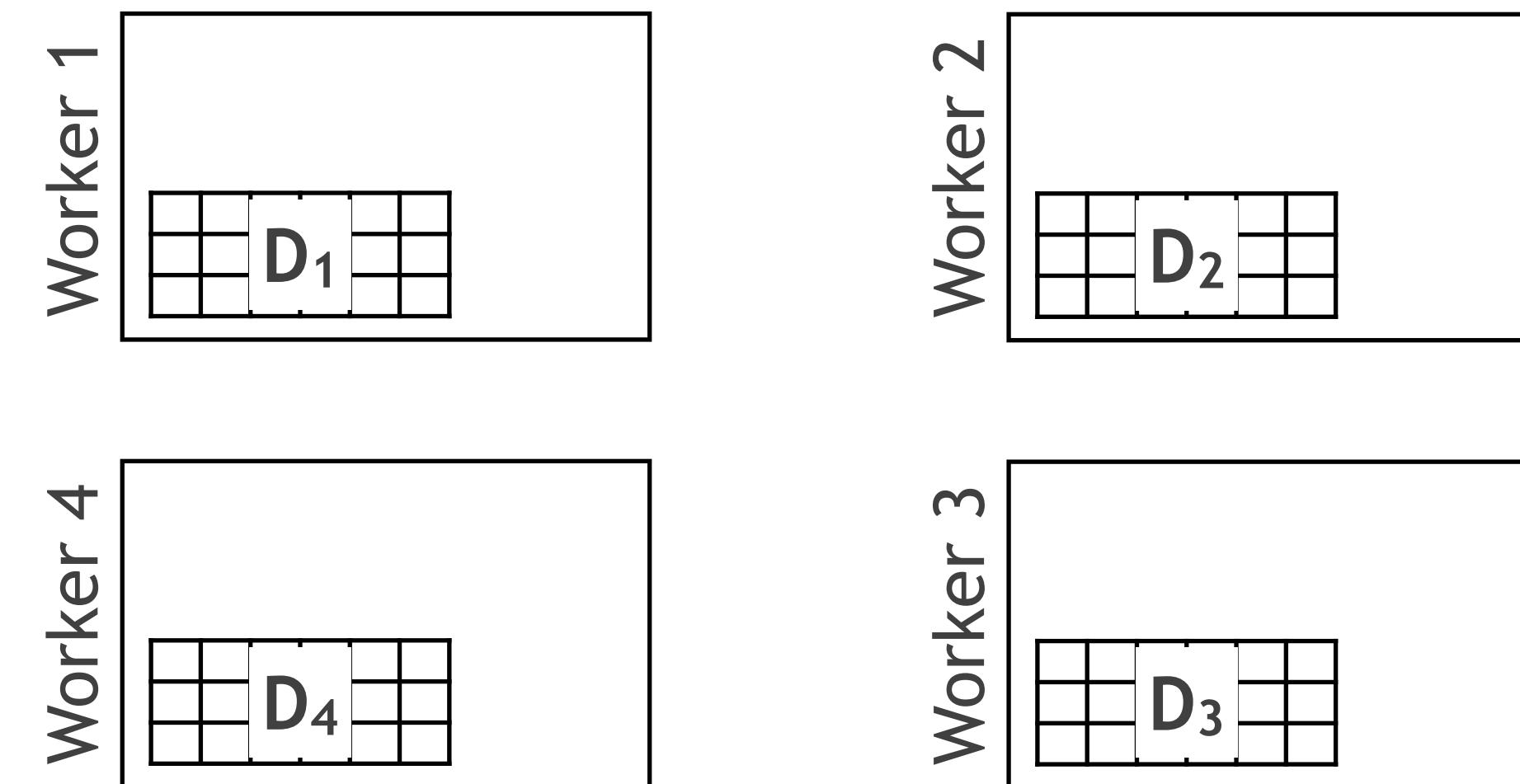
Model Hopper Parallelism (MOP)

Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

Shuffle and partition dataset

Run n DNNs on n workers



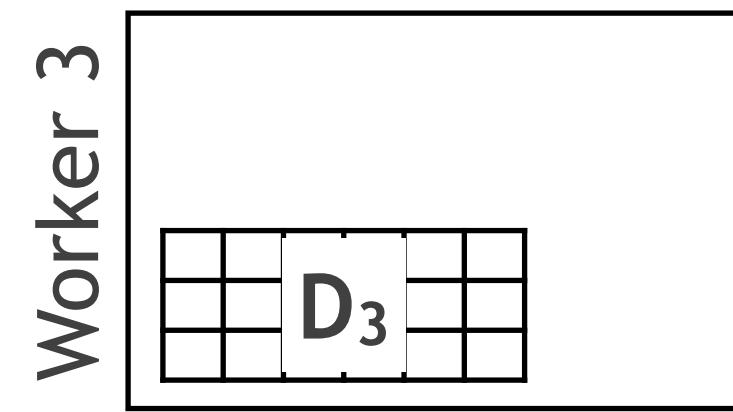
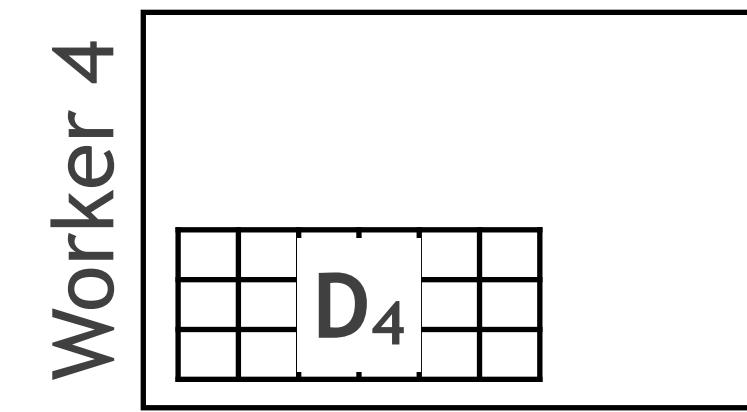
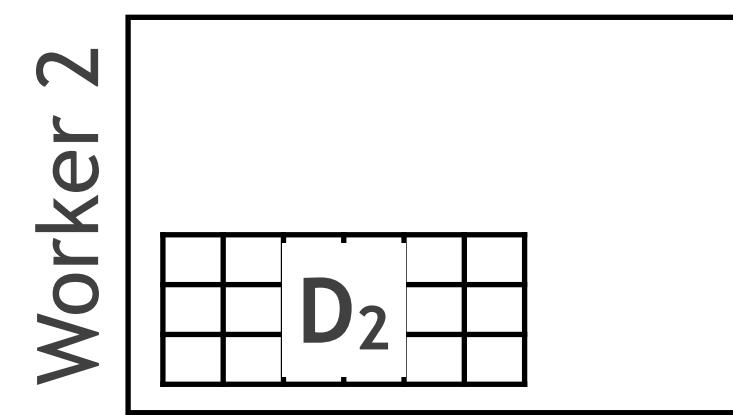
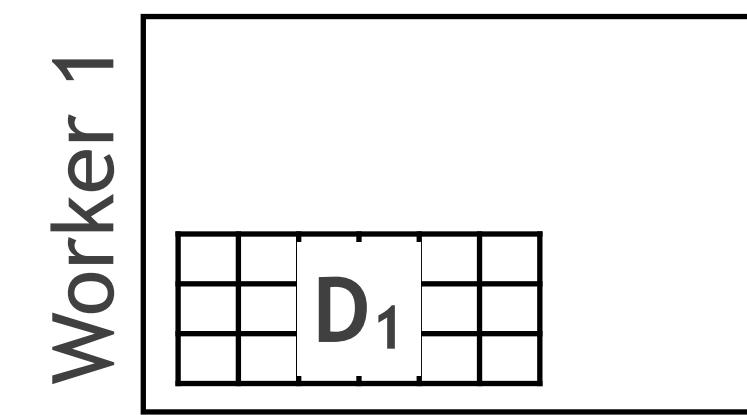
Model Hopper Parallelism (MOP)

Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

Shuffle and partition dataset
Run n DNNs on n workers

Epoch 1.1 starts in parallel



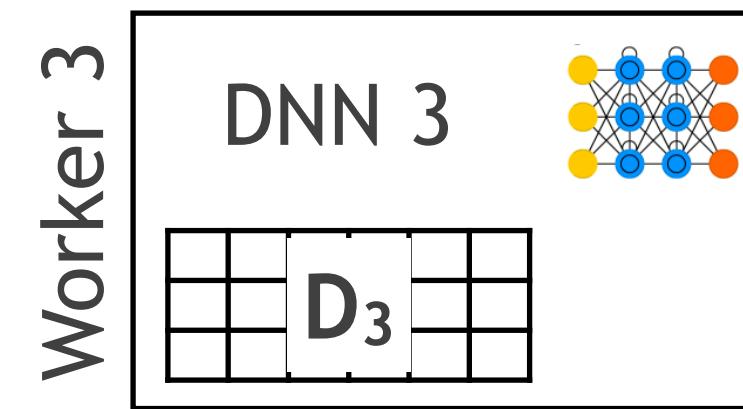
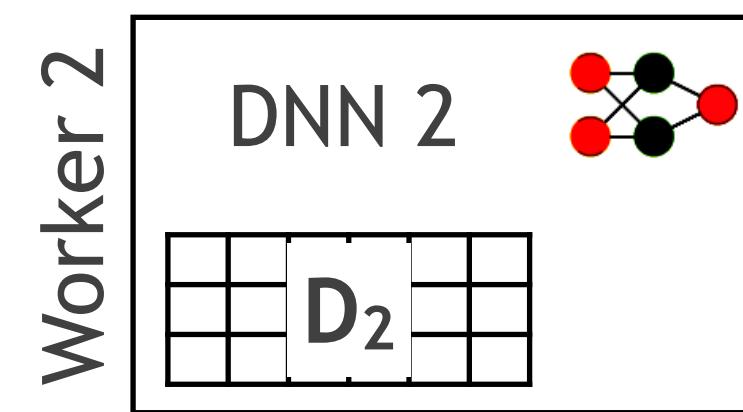
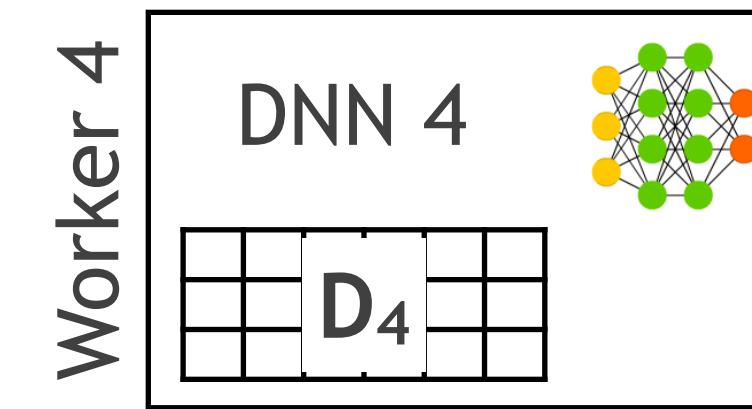
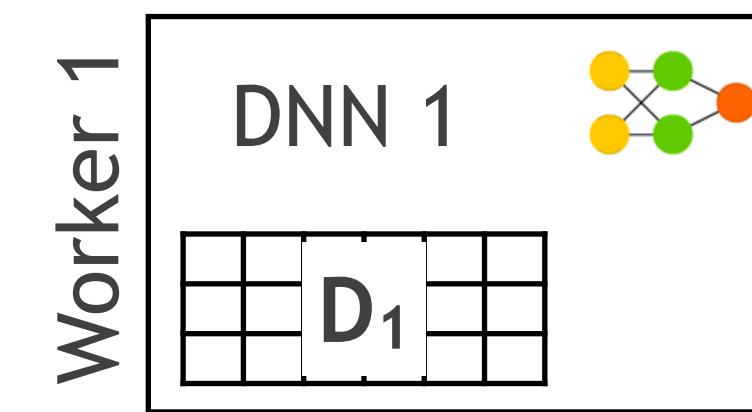
Model Hopper Parallelism (MOP)

Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

Shuffle and partition dataset
Run n DNNs on n workers

Epoch 1.1 starts in parallel



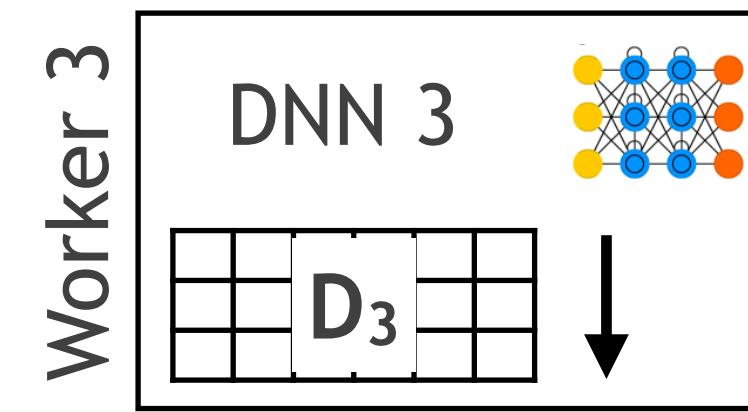
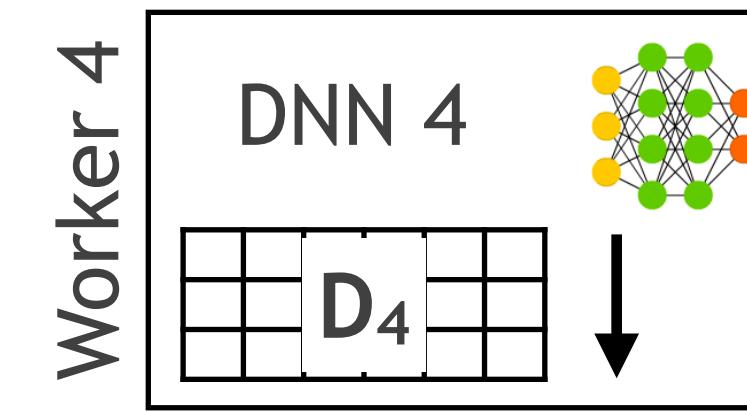
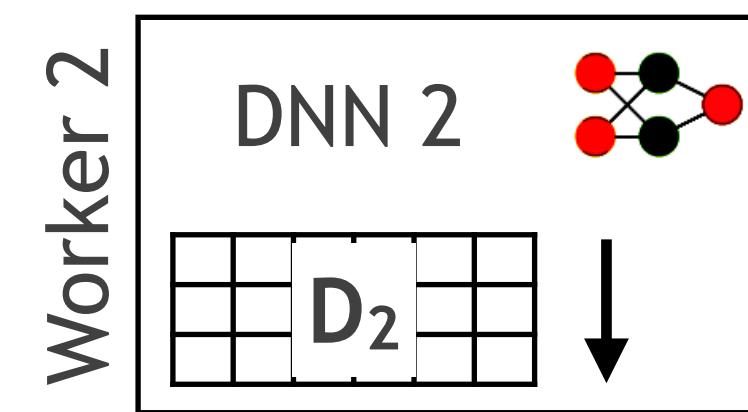
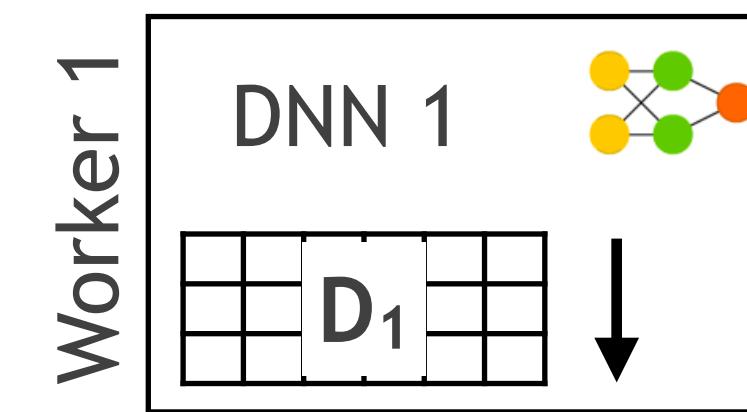
Model Hopper Parallelism (MOP)

Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

Shuffle and partition dataset
Run n DNNs on n workers

Epoch 1.1 starts in parallel



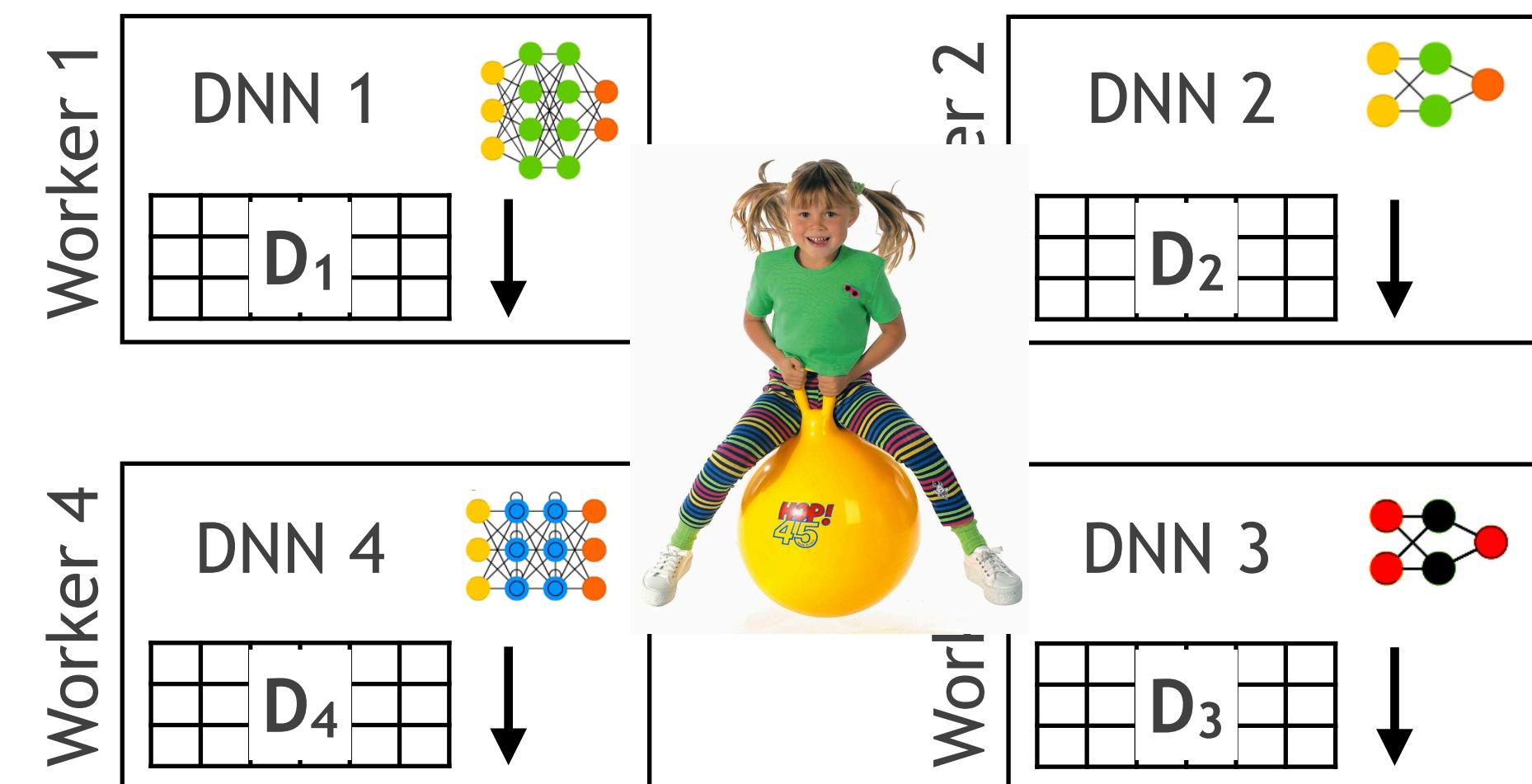
Model Hopper Parallelism (MOP)

Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

Shuffle and partition dataset
Run n DNNs on n workers

Epoch 1.1 starts in parallel



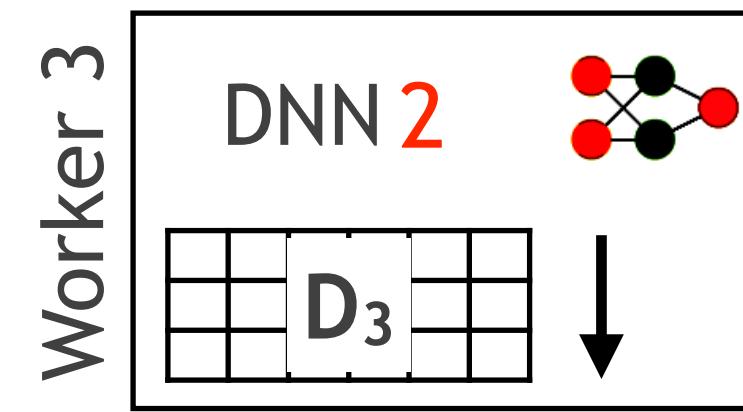
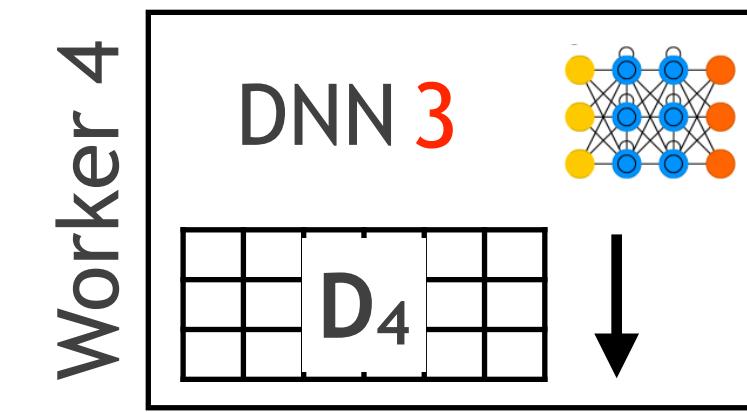
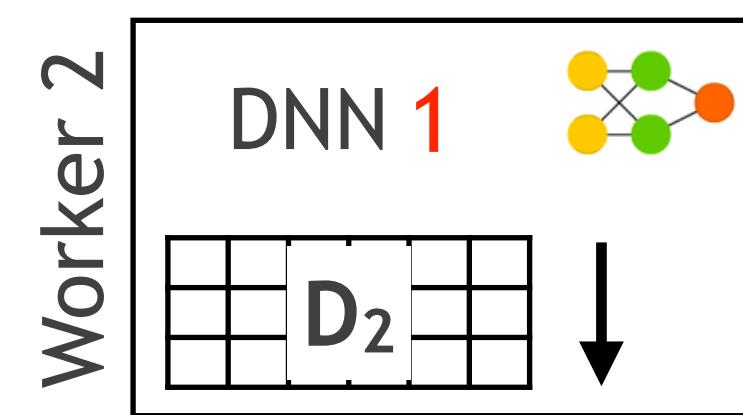
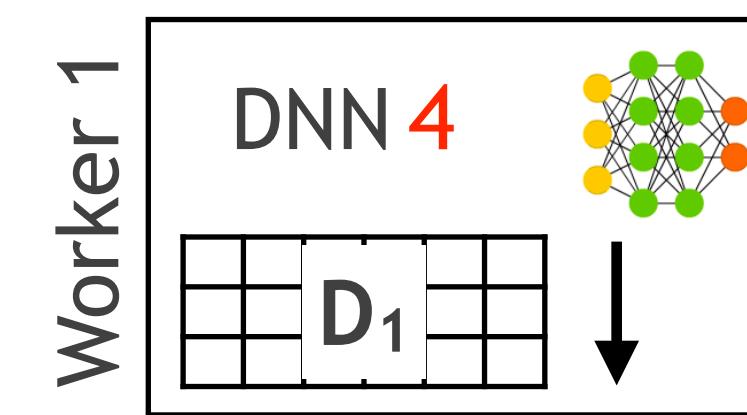
Model Hopper Parallelism (MOP)

Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

Shuffle and partition dataset
Run n DNNs on n workers

Epoch 1.1 starts in parallel



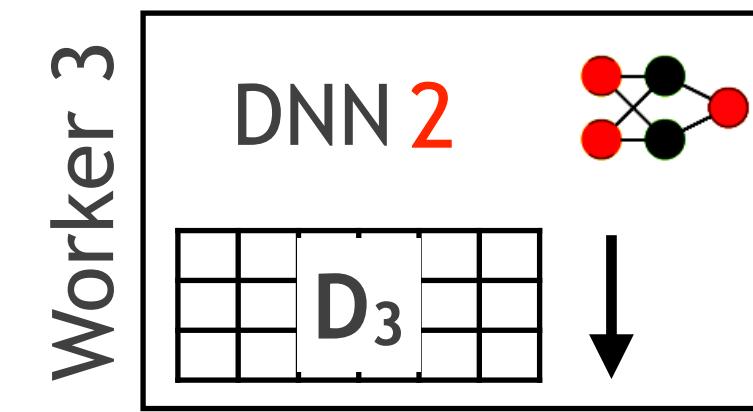
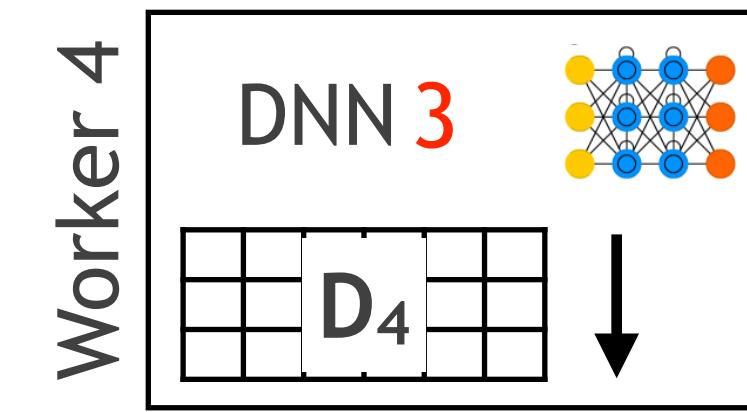
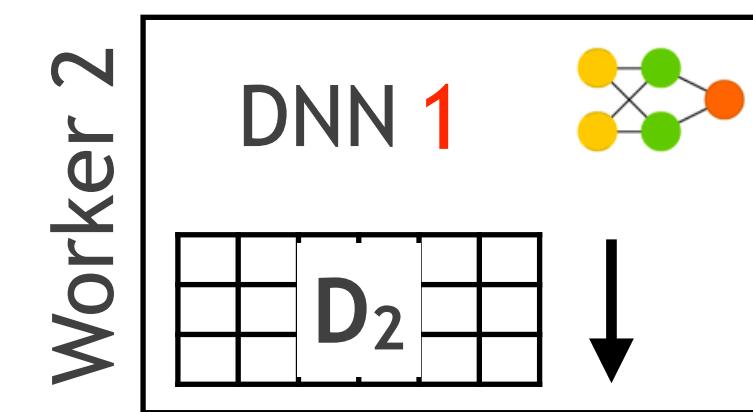
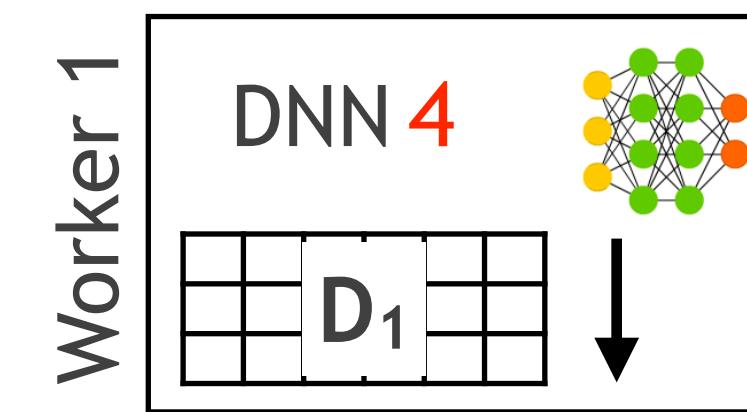
Model Hopper Parallelism (MOP)

Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

Shuffle and partition dataset
Run n DNNs on n workers

Epoch 1.2 starts in parallel



Model Hopper Parallelism (MOP)

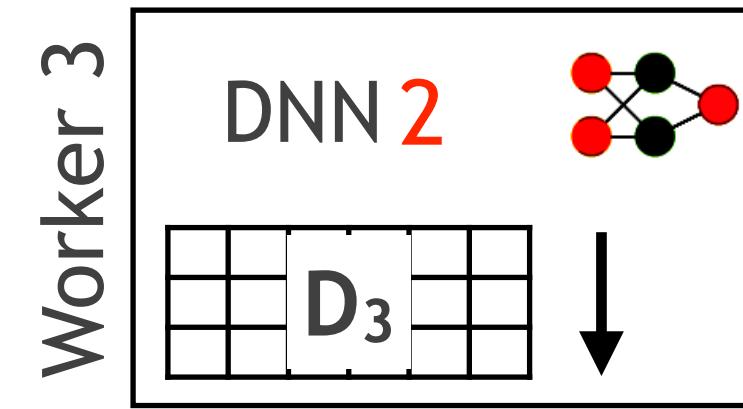
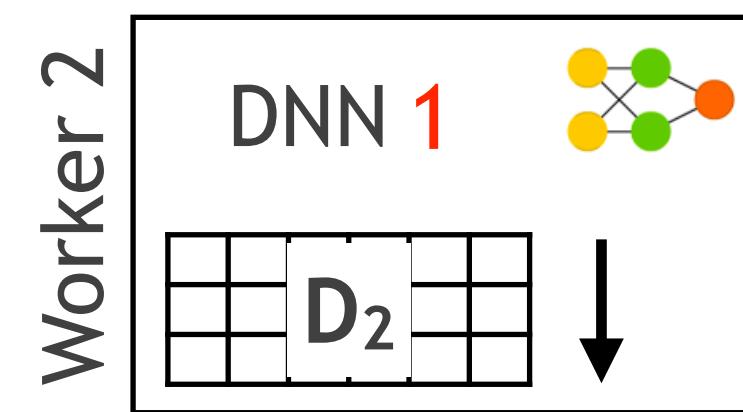
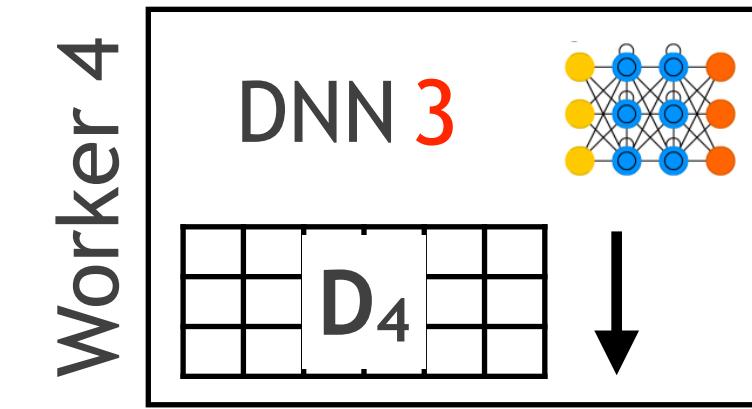
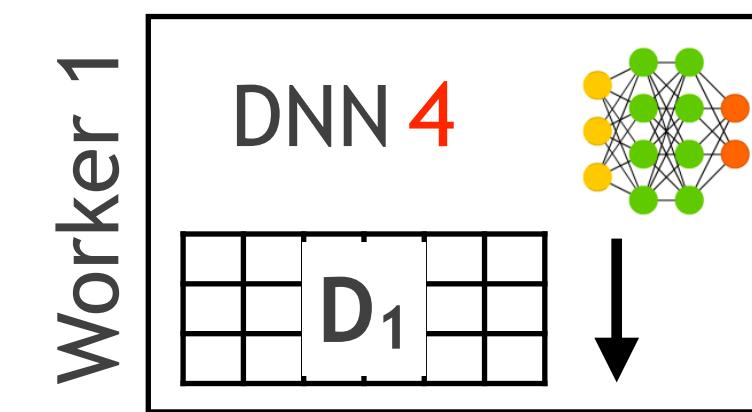
Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

Shuffle and partition dataset
Run n DNNs on n workers

Each model keeps “hopping” until
it sees all of D

Epoch 1.2 starts in parallel



Model Hopper Parallelism (MOP)

Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

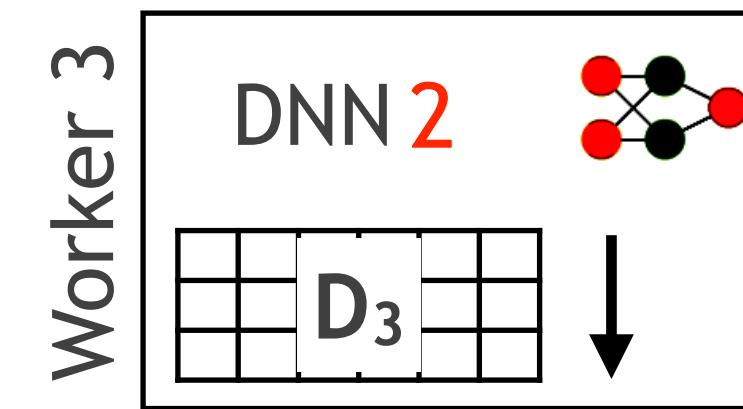
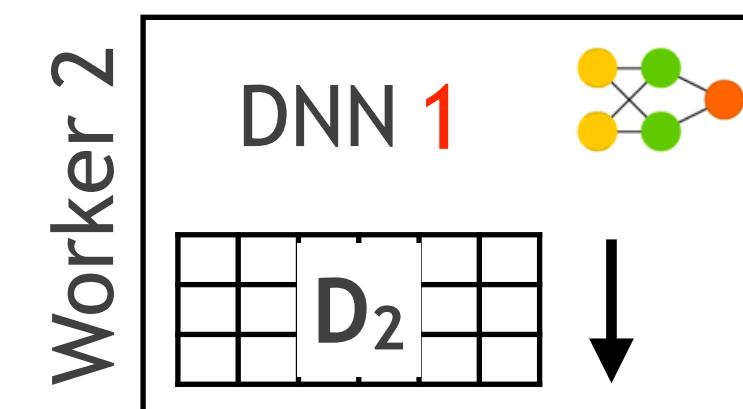
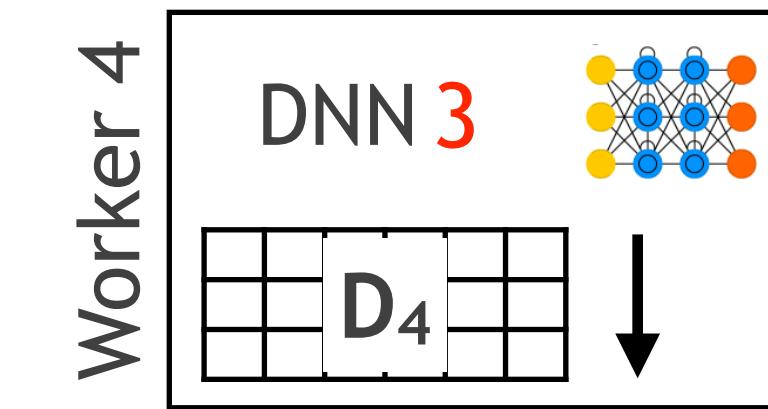
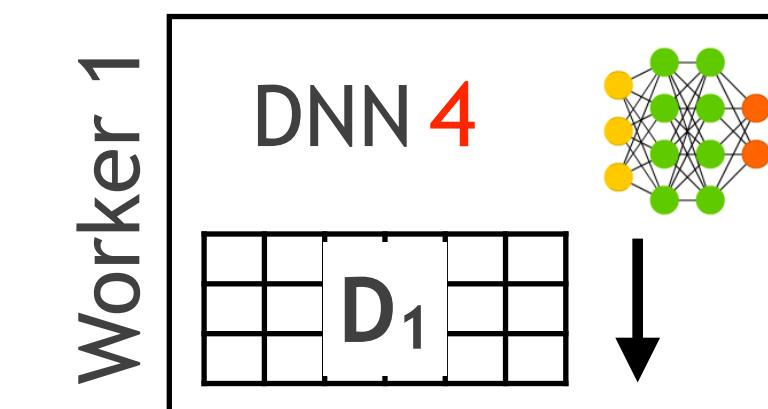
Shuffle and partition dataset
Run n DNNs on n workers

Each model keeps “hopping” until
it sees all of D

Strong theoretical guarantees:

1. *Equivalence* to sequential SGD
2. Hits lower bound on comm. cost

Epoch 1.2 starts in parallel



Model Hopper Parallelism (MOP)

Insight from Optimization Theory:

SGD is robust to *data ordering randomness*

Technical Challenges (See paper for details):

Resource-aware scheduling of evolving model configs' hops

Support data replication, fault tolerance, and elasticity

Non-disruptive integration with existing DL systems

1. *Equivalence* to sequential SGD
2. Hits lower bound on comm. cost

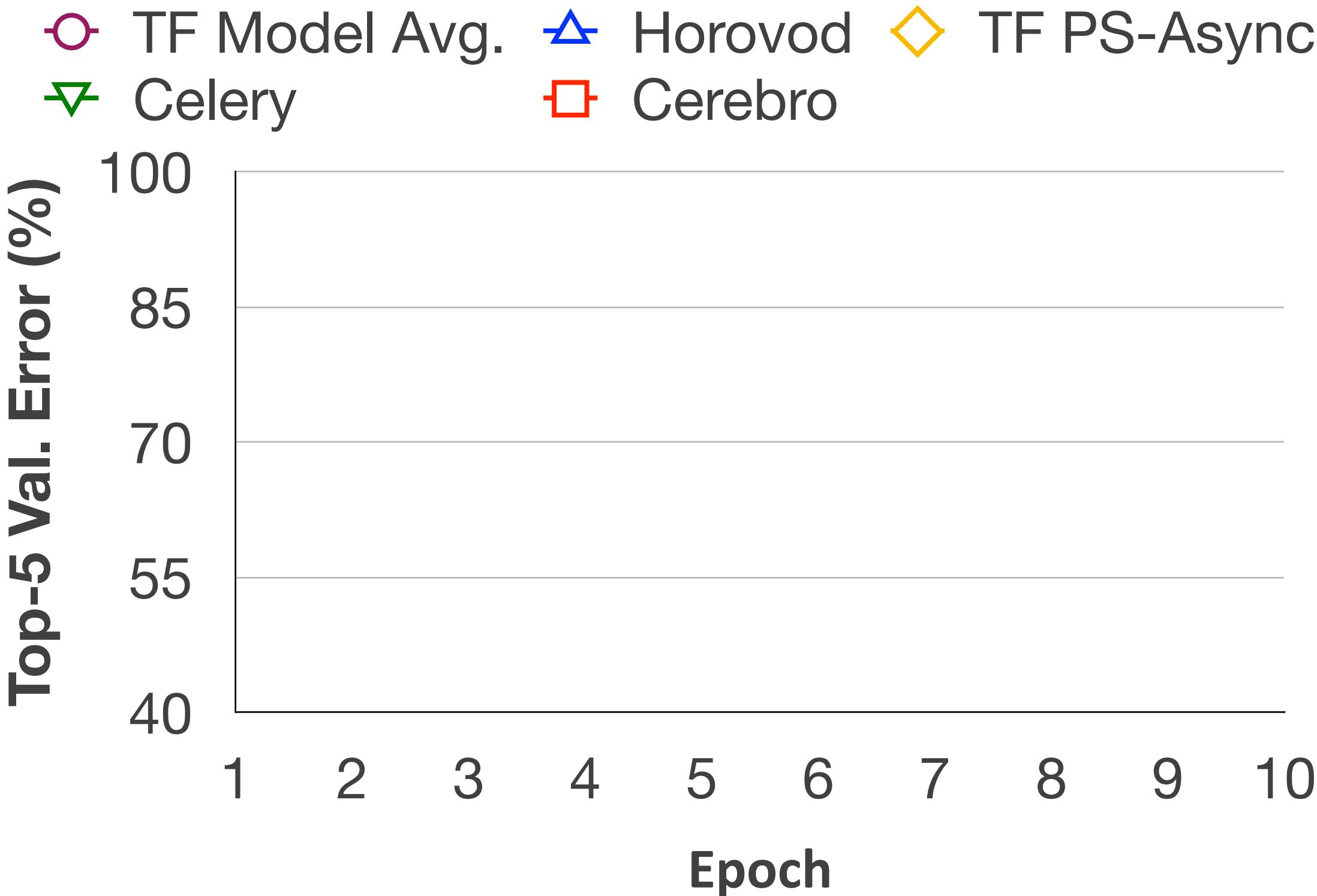


Experimental Evaluation

Setup: ImageNet; 16 CNN configurations; TensorFlow; 8 GPU nodes

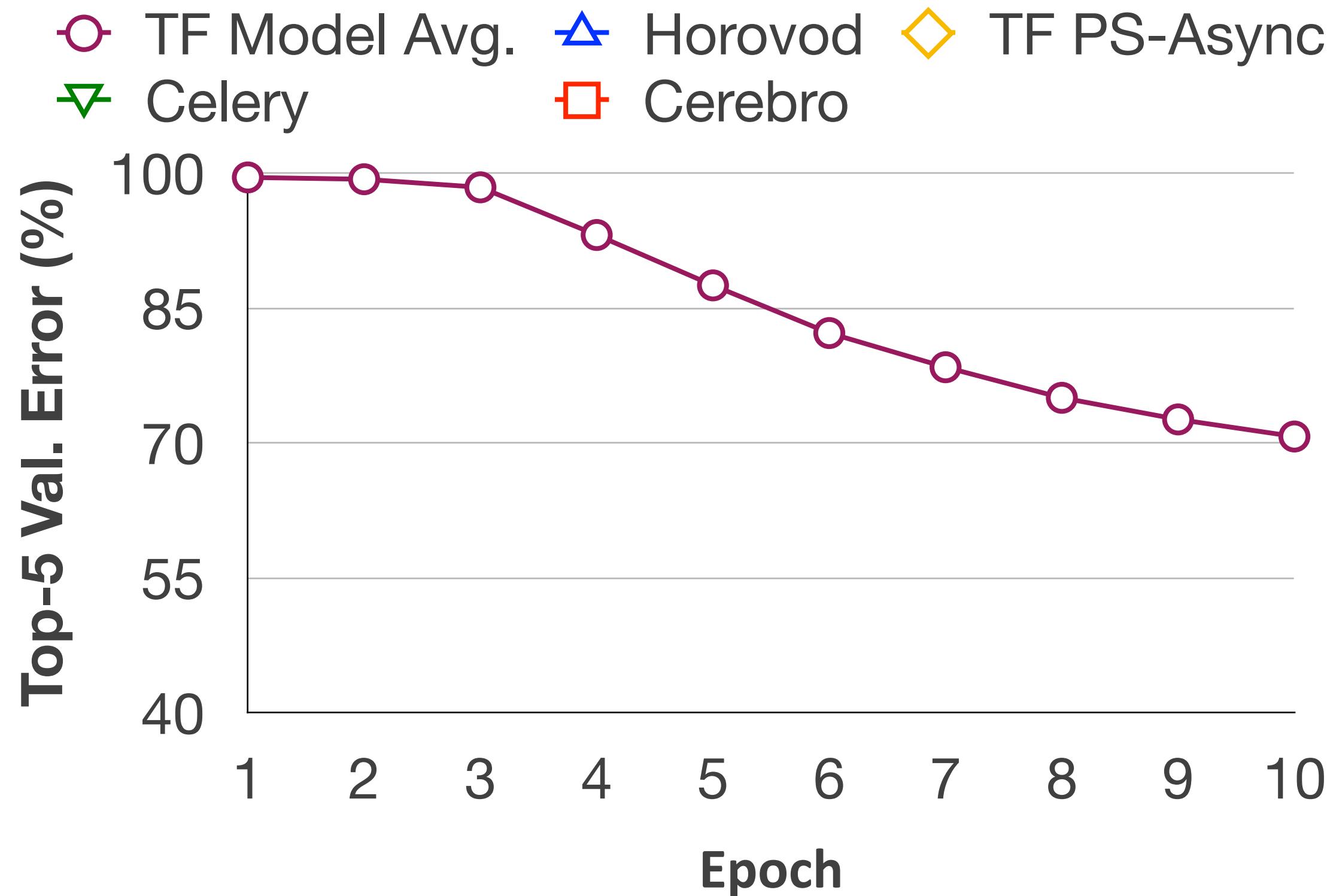
Experimental Evaluation

Setup: ImageNet; 16 CNN configurations; TensorFlow; 8 GPU nodes



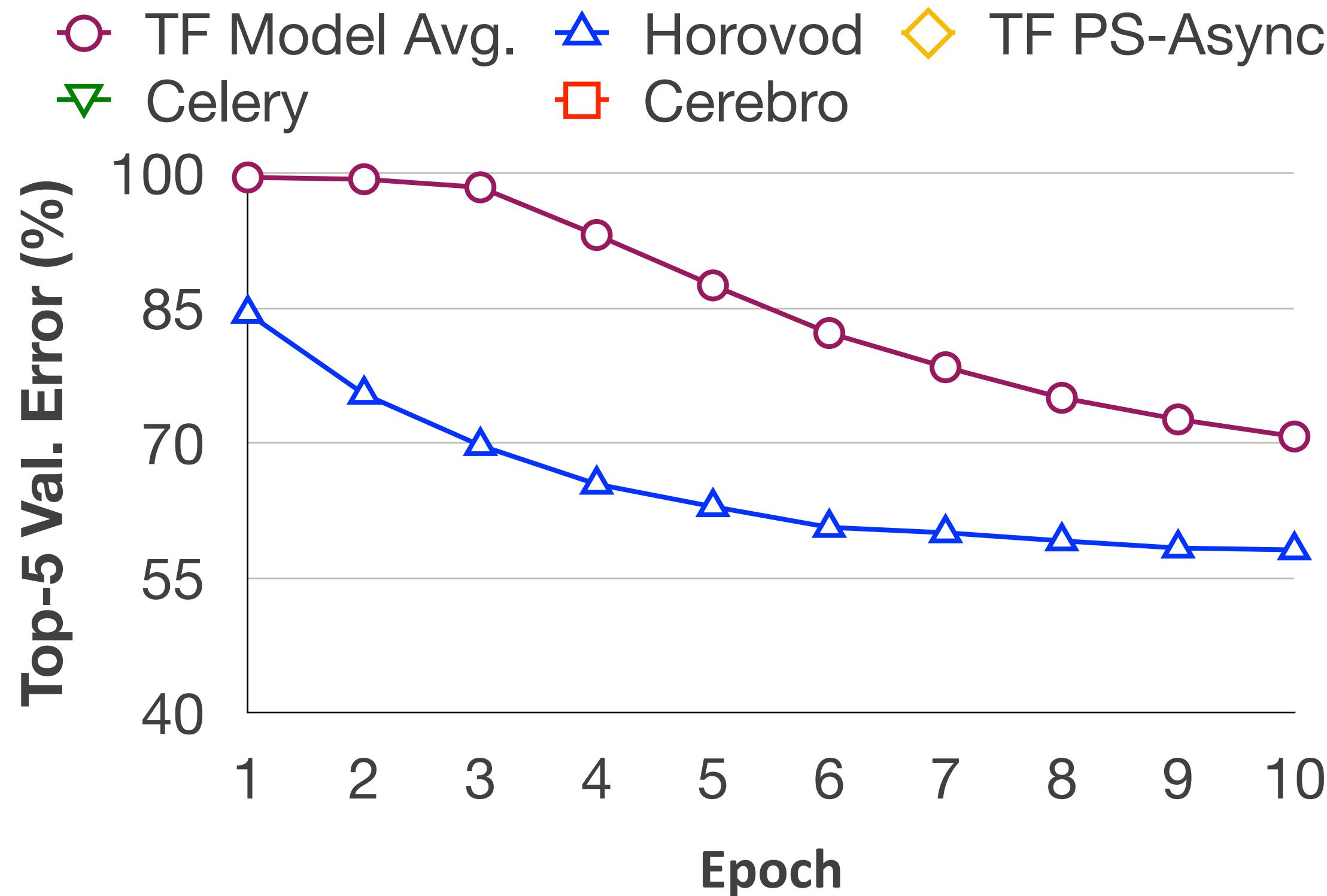
Experimental Evaluation

Setup: ImageNet; 16 CNN configurations; TensorFlow; 8 GPU nodes



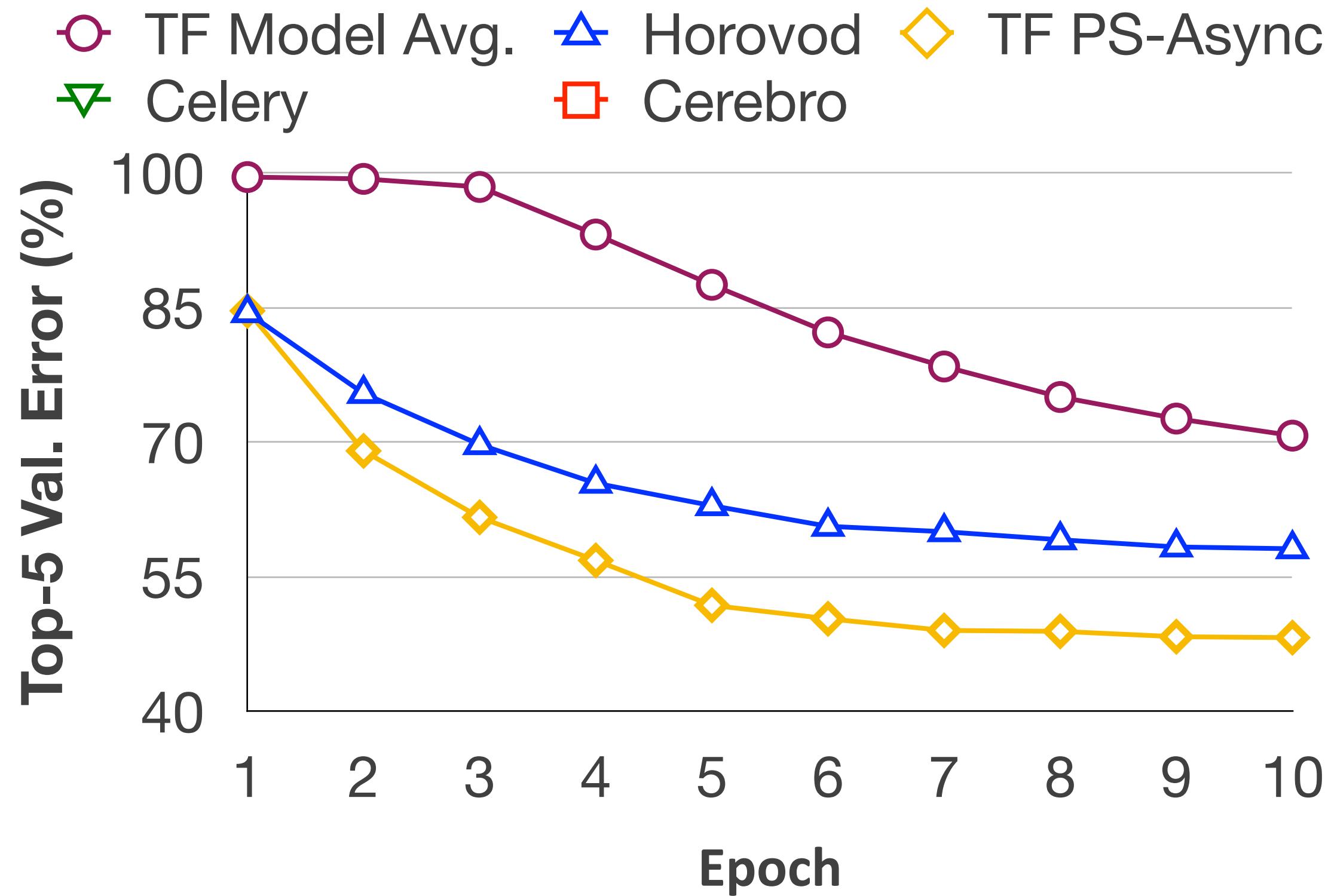
Experimental Evaluation

Setup: ImageNet; 16 CNN configurations; TensorFlow; 8 GPU nodes



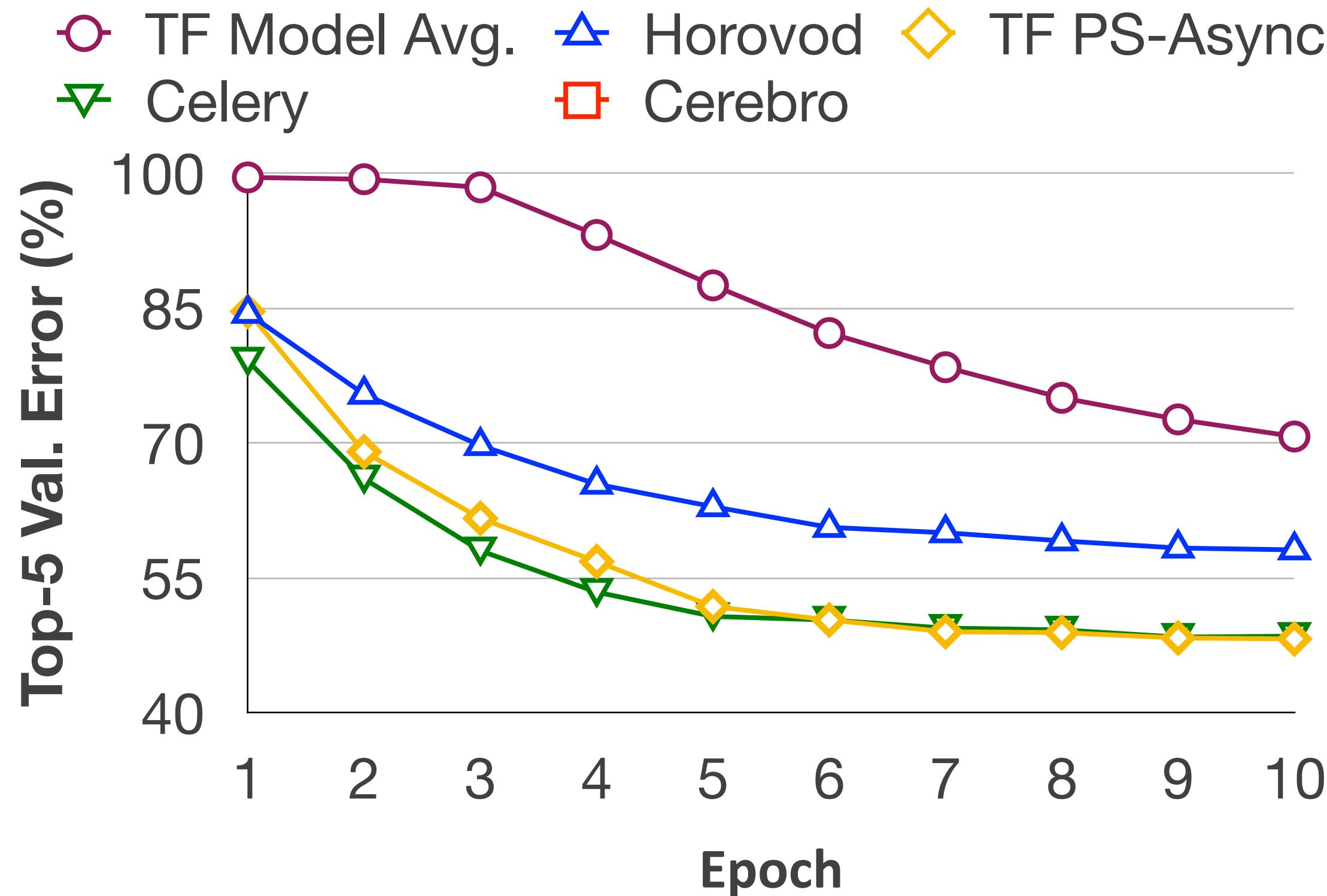
Experimental Evaluation

Setup: ImageNet; 16 CNN configurations; TensorFlow; 8 GPU nodes



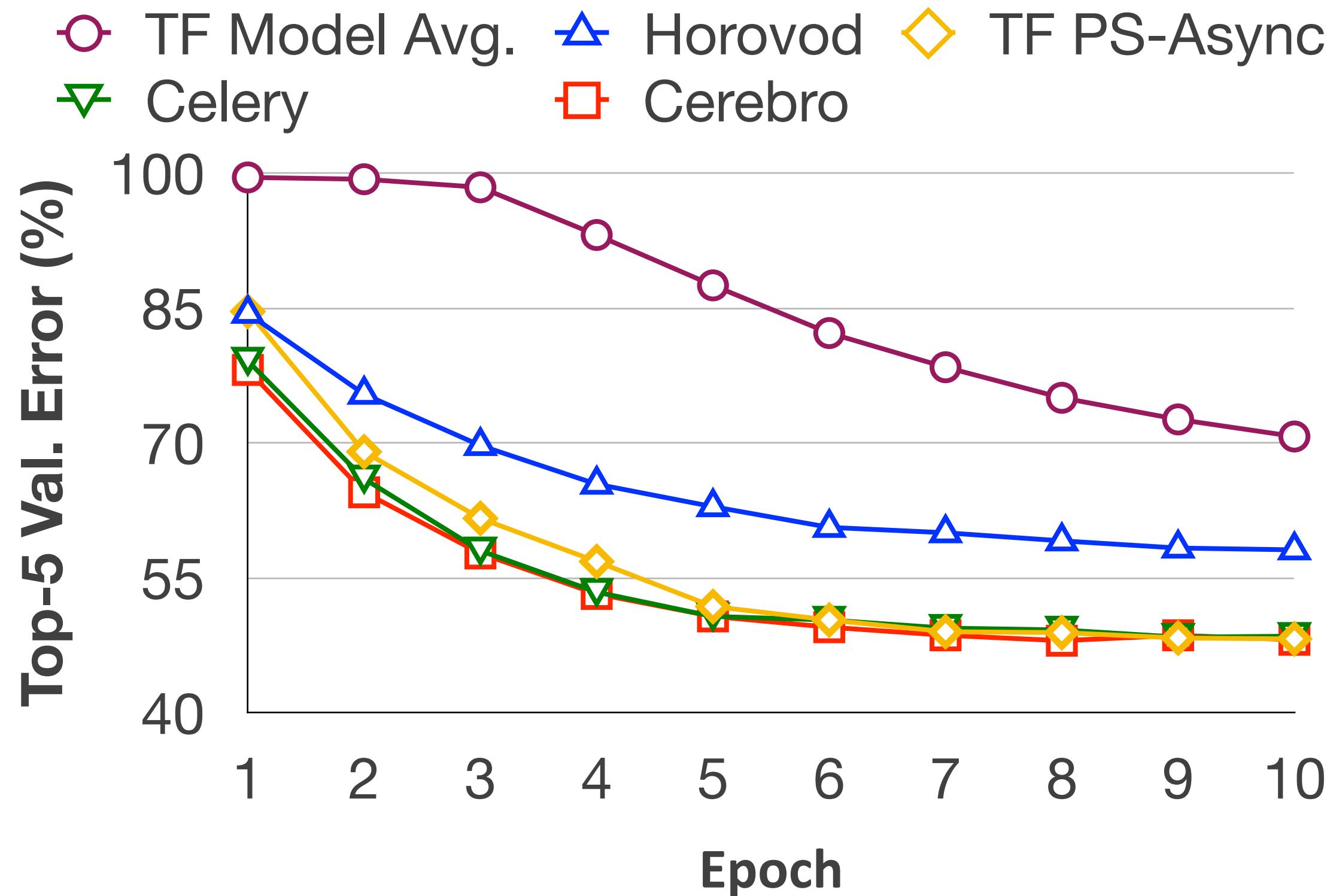
Experimental Evaluation

Setup: ImageNet; 16 CNN configurations; TensorFlow; 8 GPU nodes



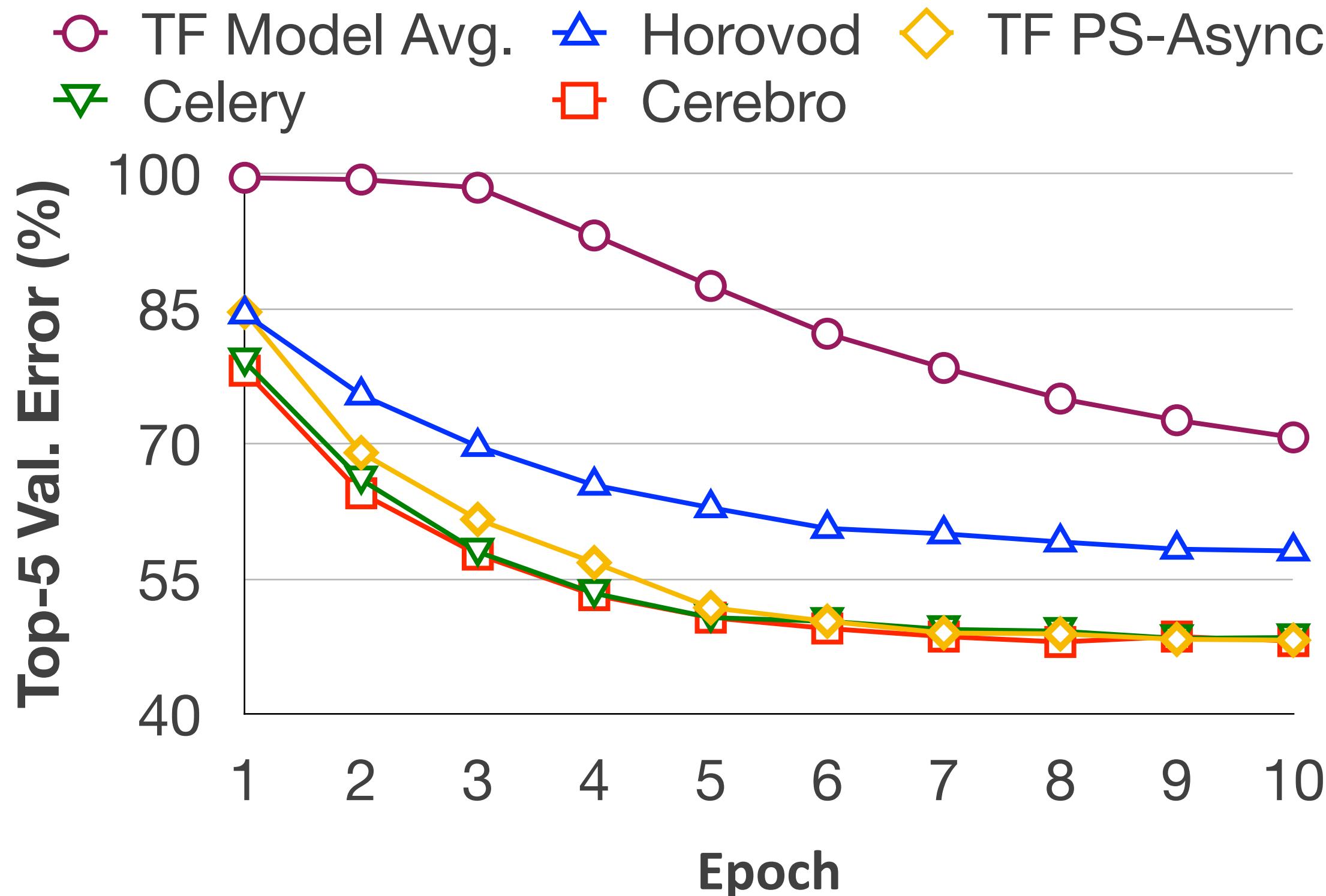
Experimental Evaluation

Setup: ImageNet; 16 CNN configurations; TensorFlow; 8 GPU nodes



Experimental Evaluation

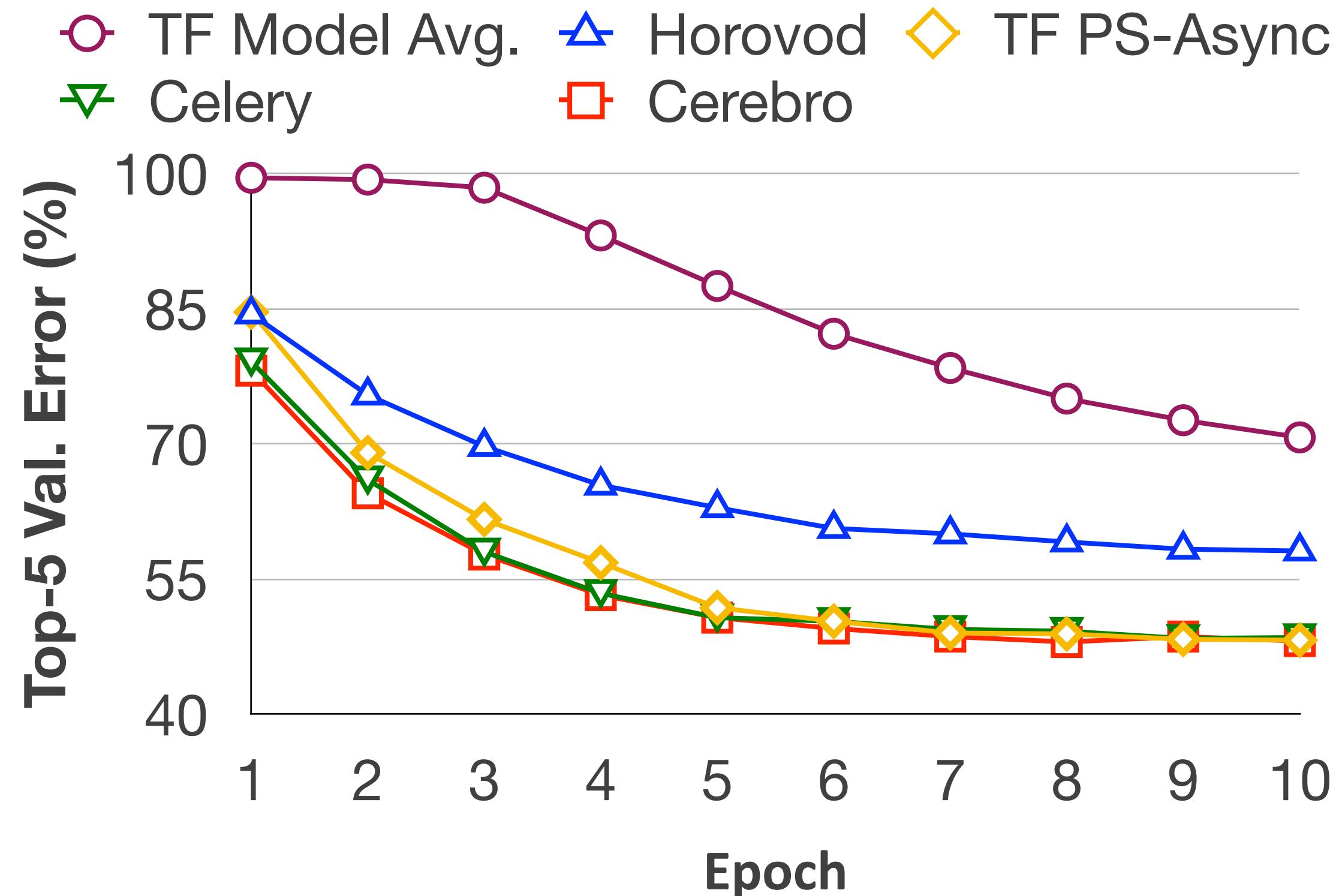
Setup: ImageNet; 16 CNN configurations; TensorFlow; 8 GPU nodes



Method	Runtime (hrs)	Memory
TF PS-Async.	190.0 (!)	200 GB
Horovod	54.2	200 GB
TF Model Averaging	19.7	200 GB
Celery (Task-Parallel)	17.2-19.0*	1600 GB (!)
MOP/Cerebro	17.7	200 GB

Experimental Evaluation

Setup: ImageNet; 16 CNN configurations; TensorFlow; 8 GPU nodes

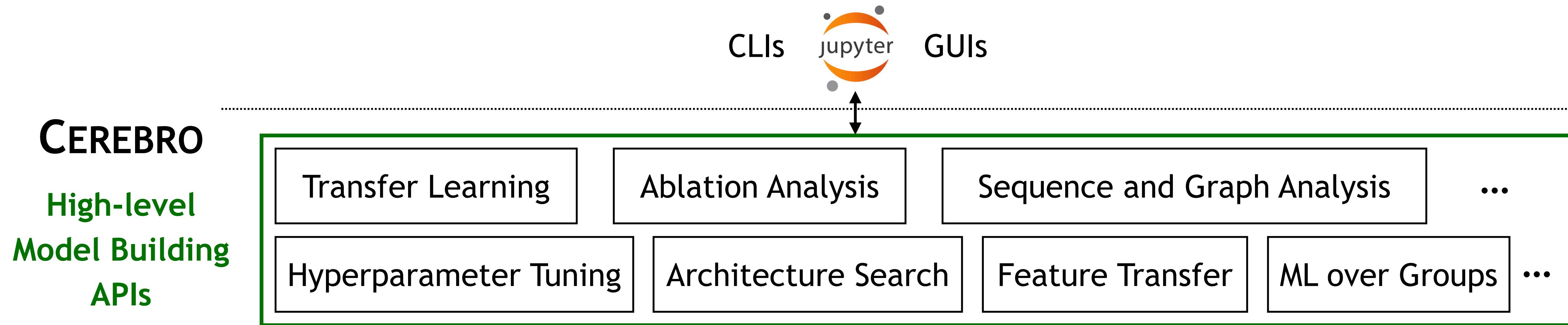


Method	Runtime (hrs)	Memory
TF PS-Async.	190.0 (!)	200 GB
Horovod	54.2	200 GB
TF Model Averaging	19.7	200 GB
Celery (Task-Parallel)	17.2-19.0*	1600 GB (!)
MOP/Cerebro	17.7	200 GB

Cerebro offers overall best combination of resource efficiency and accuracy

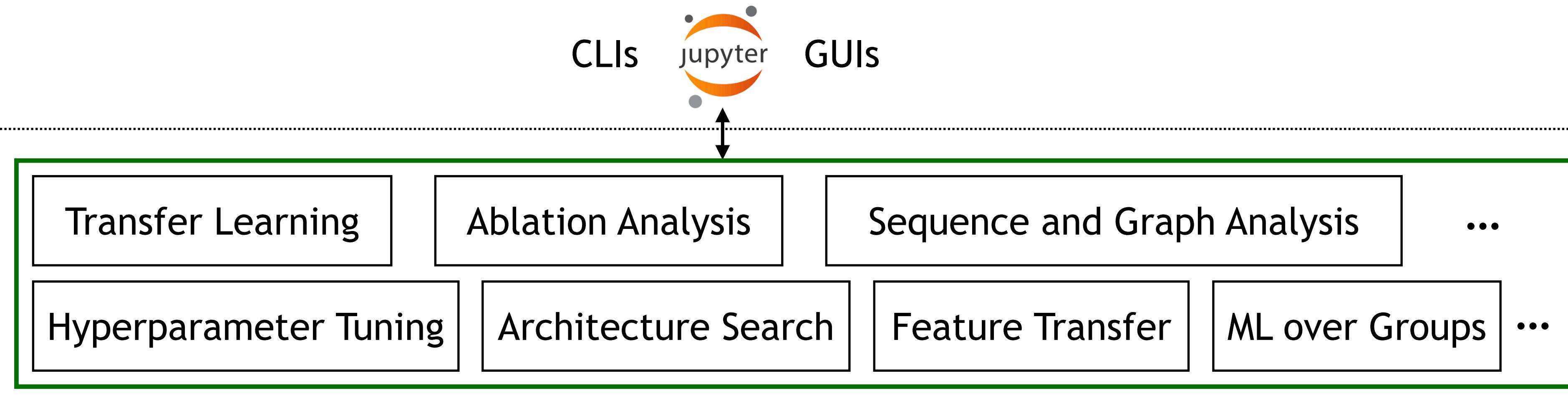
Cerebro: Efficient and Reproducible Model Selection on Deep Learning Systems. SIGMOD DEEM'19
Cerebro: A Data System for Optimized Deep Learning Model Selection. VLDB'20

Full Vision of the Cerebro Platform

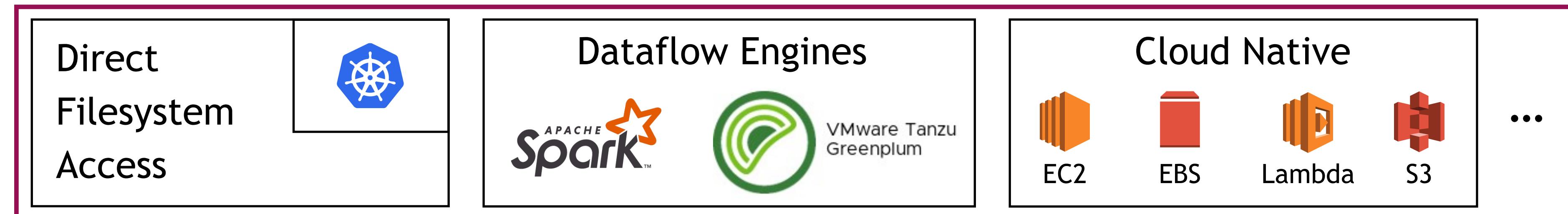


Full Vision of the Cerebro Platform

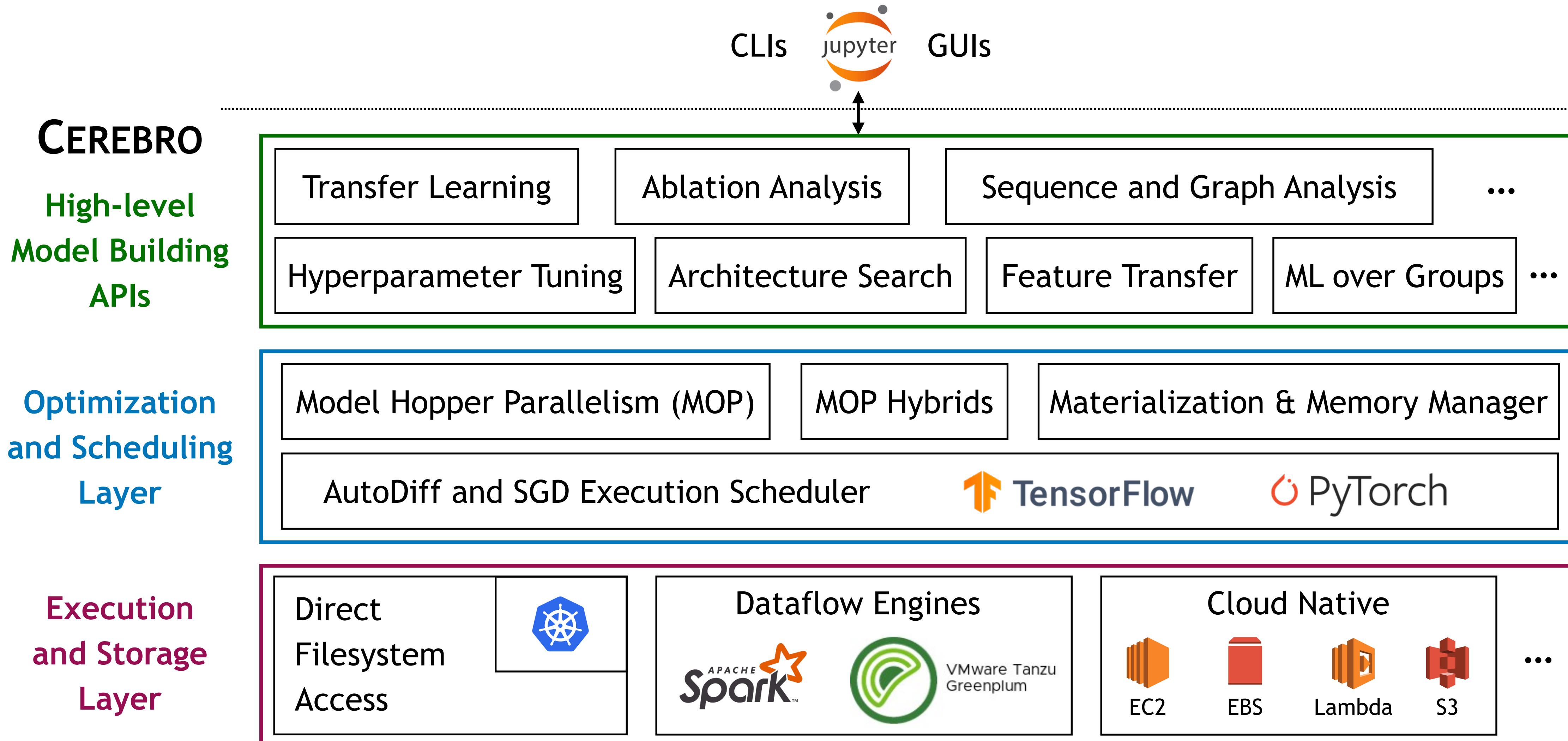
CEREBRO
High-level
Model Building
APIs



Execution
and Storage
Layer



Full Vision of the Cerebro Platform

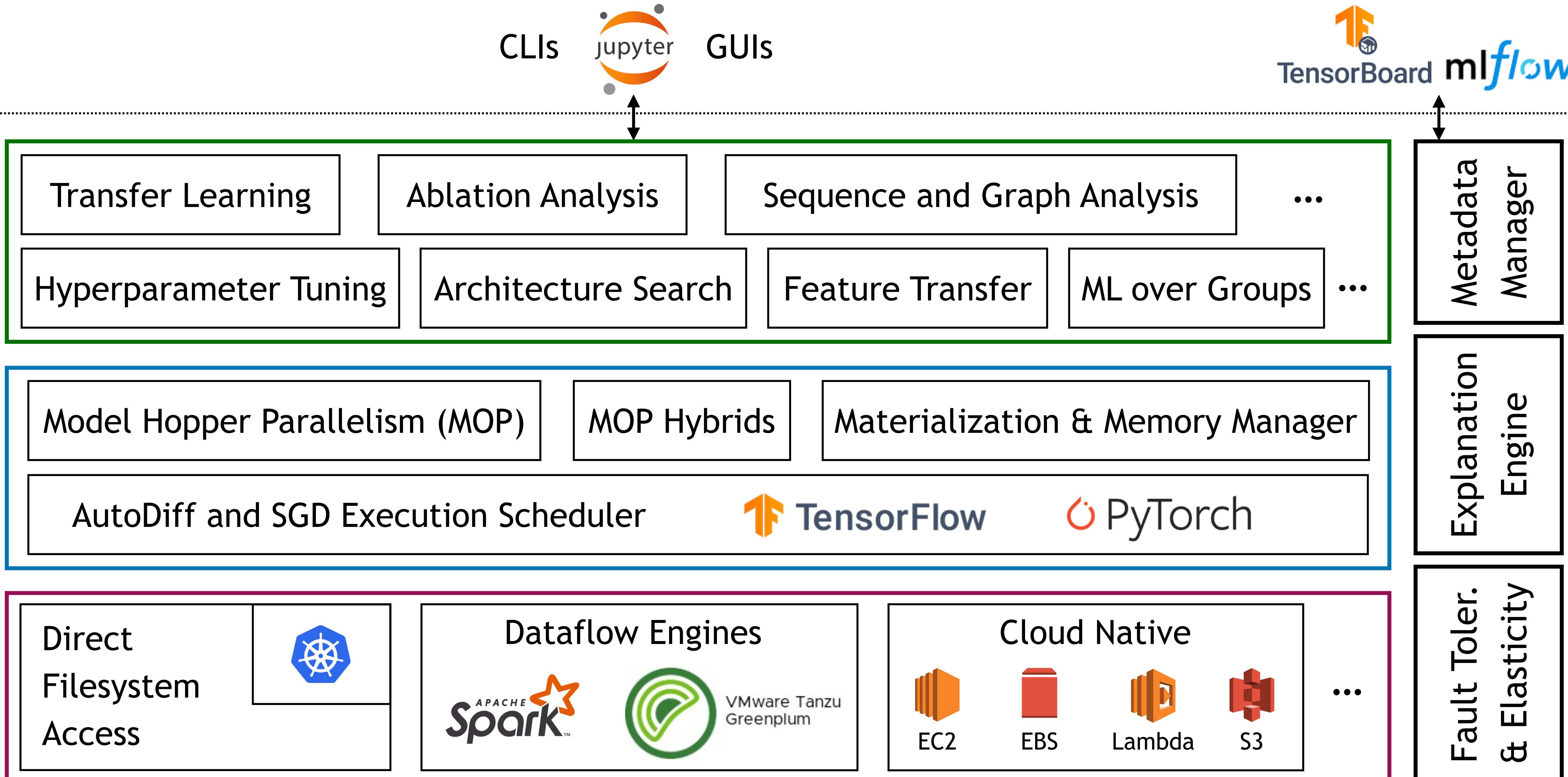
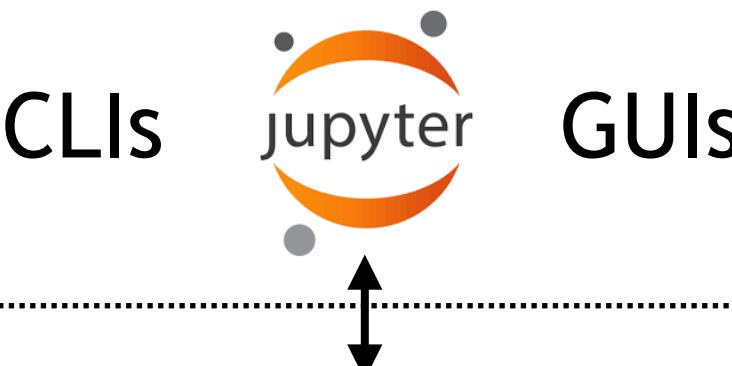


Full Vision of the Cerebro Platform

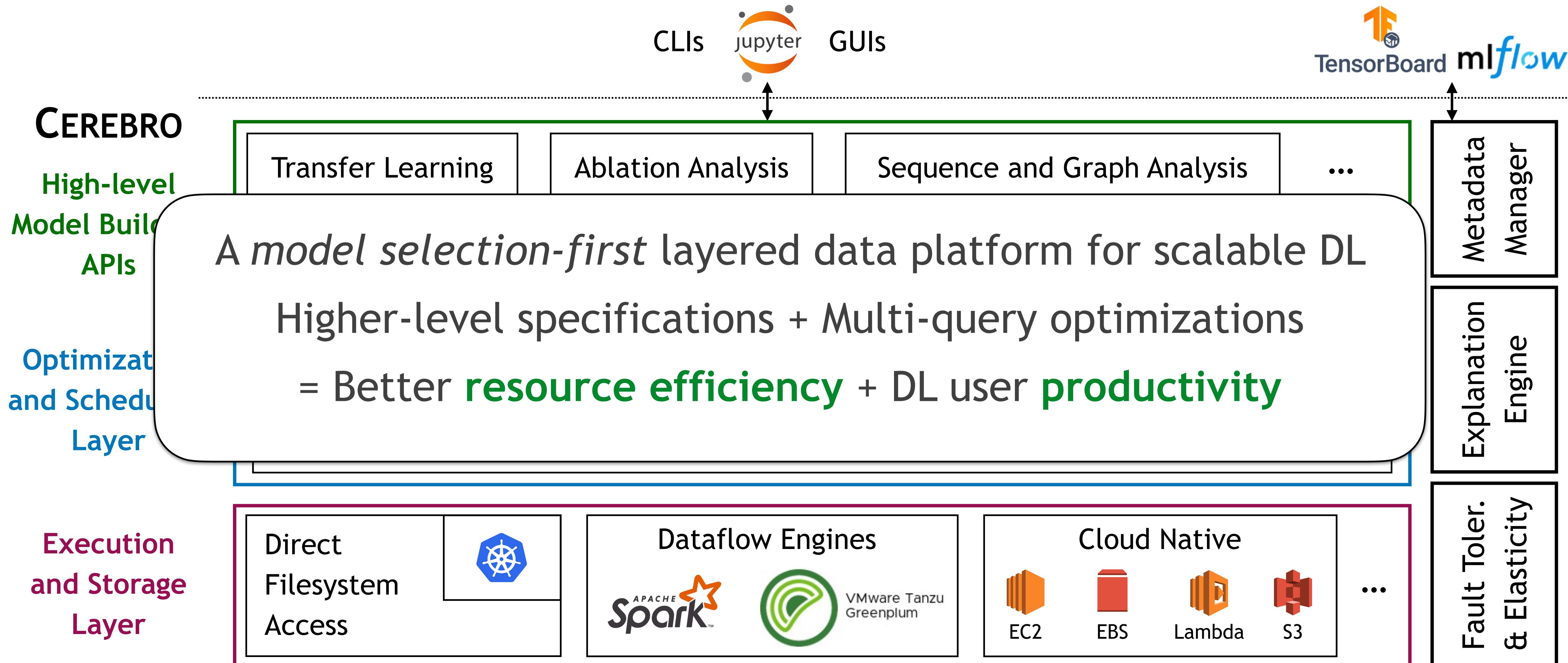
CERE BRO
High-level
Model Building
APIs

Optimization
and Scheduling
Layer

Execution
and Storage
Layer



Full Vision of the Cerebro Platform



Cerebro: Early Impact and Trajectory

Used on TBs by UCSD Public Health scientists; Physics & Materials Eng. next
Shipped with Apache MADlib by Pivotal; VMware collaboration
Cerebro+Spark released; all open sourced under Apache License v2.0
Coming soon: More workloads + optimizations + domain science use cases

Cerebro: Early Impact and Trajectory

Used on TBs by UCSD Public Health scientists; Physics & Materials Eng. next
Shipped with Apache MADlib by Pivotal; VMware collaboration
Cerebro+Spark released; all open sourced under Apache License v2.0
Coming soon: More workloads + optimizations + domain science use cases

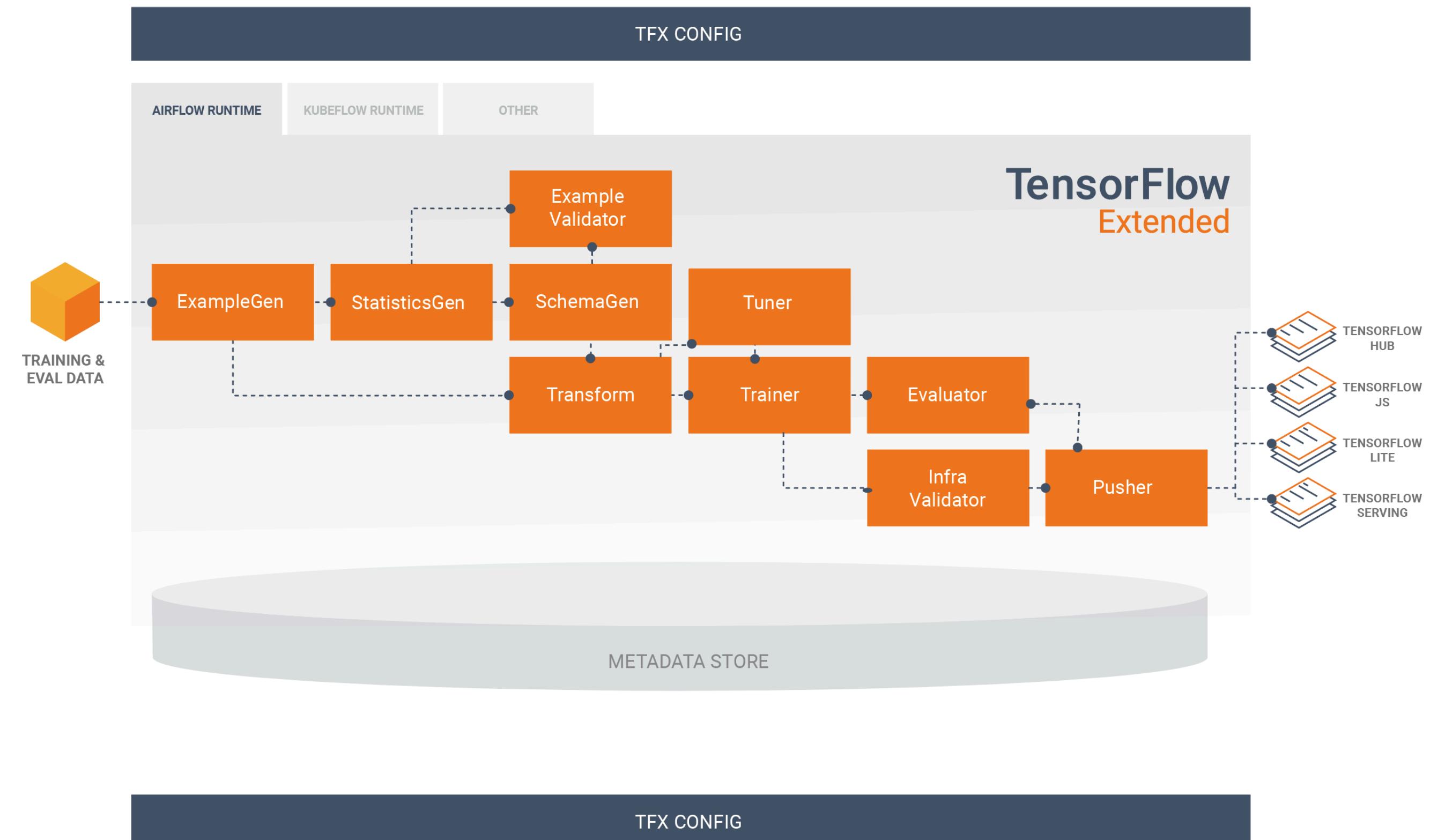
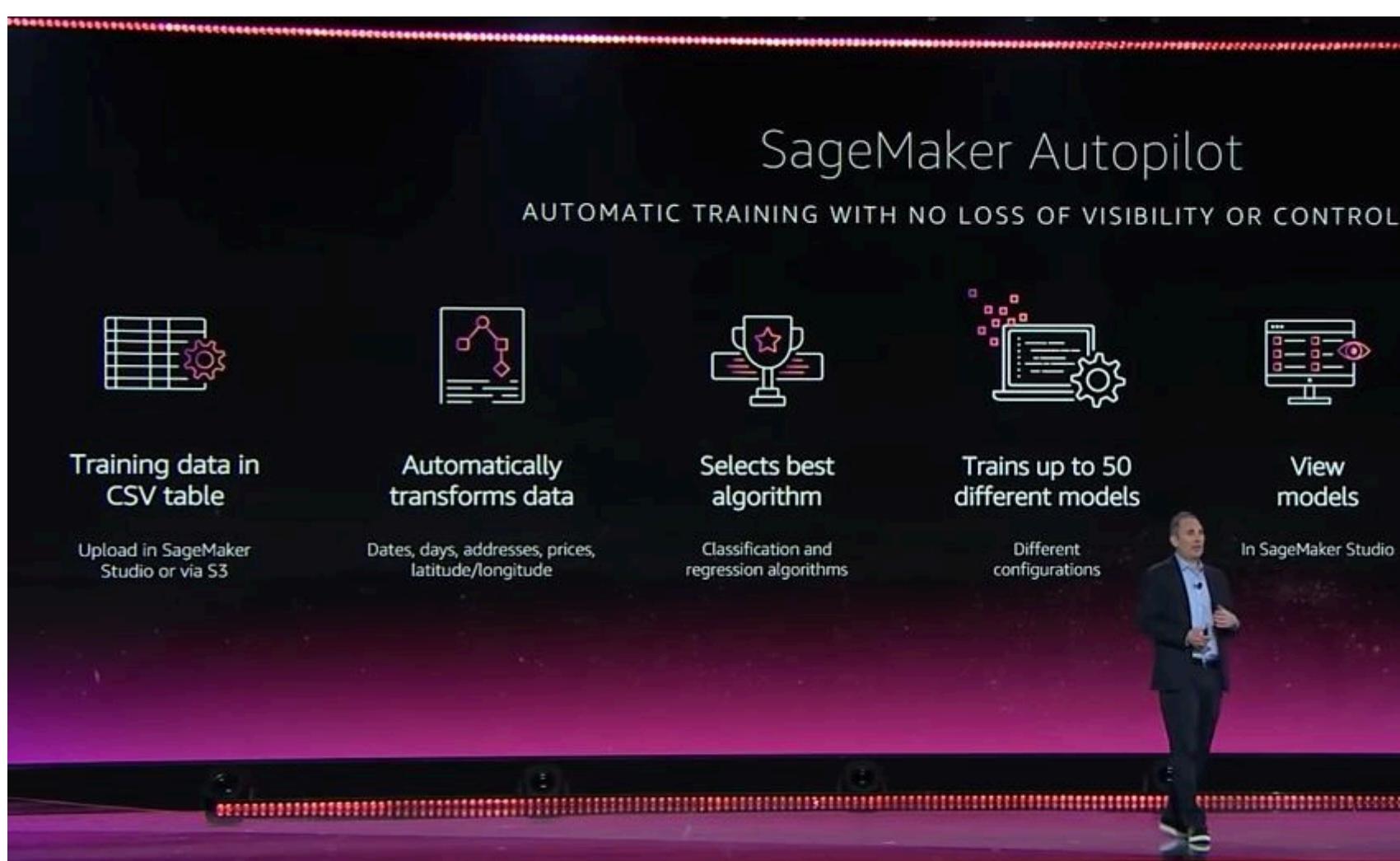
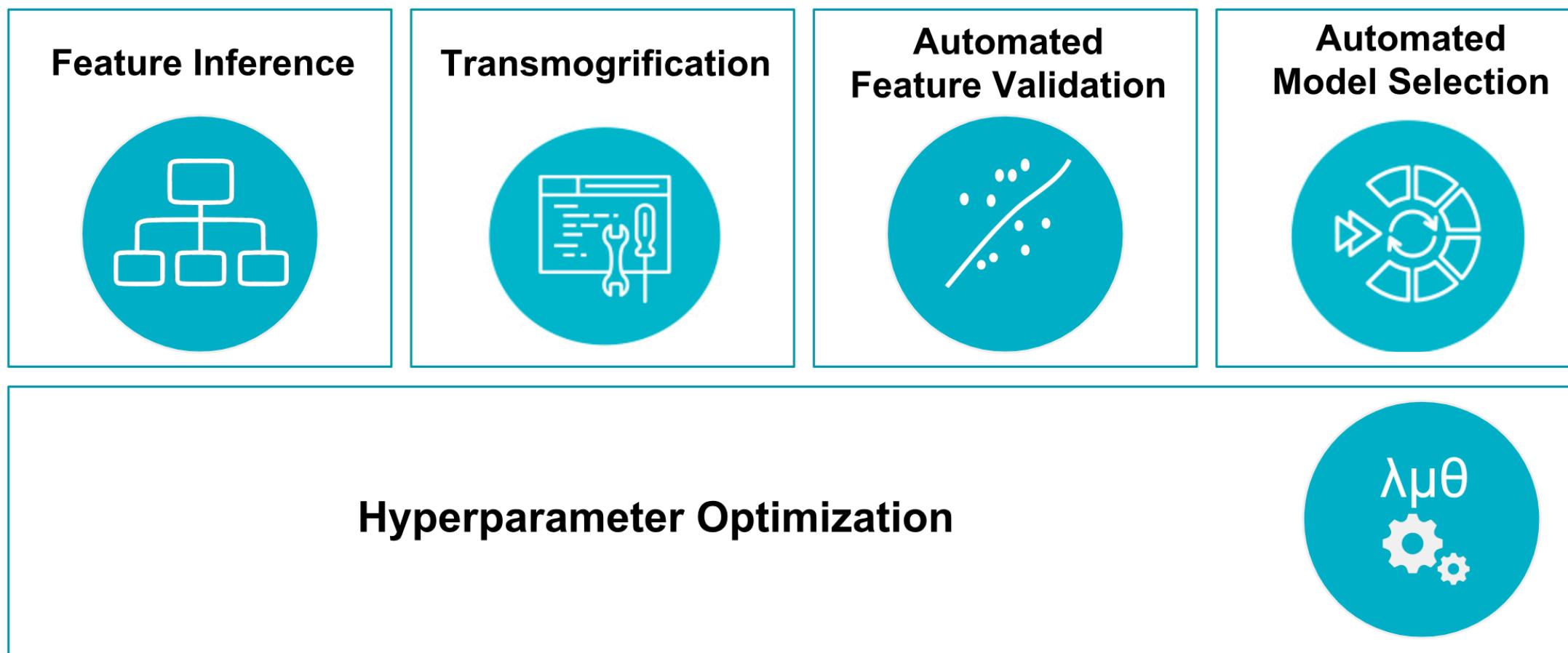


Outline

- | The New DBfication of ML/AI
- | Two Examples from My Research
 - | Example 1: Scalable DL Systems
 - | Example 2: Auto Data Prep for ML
- | Accelerating the DBfication of ML/AI

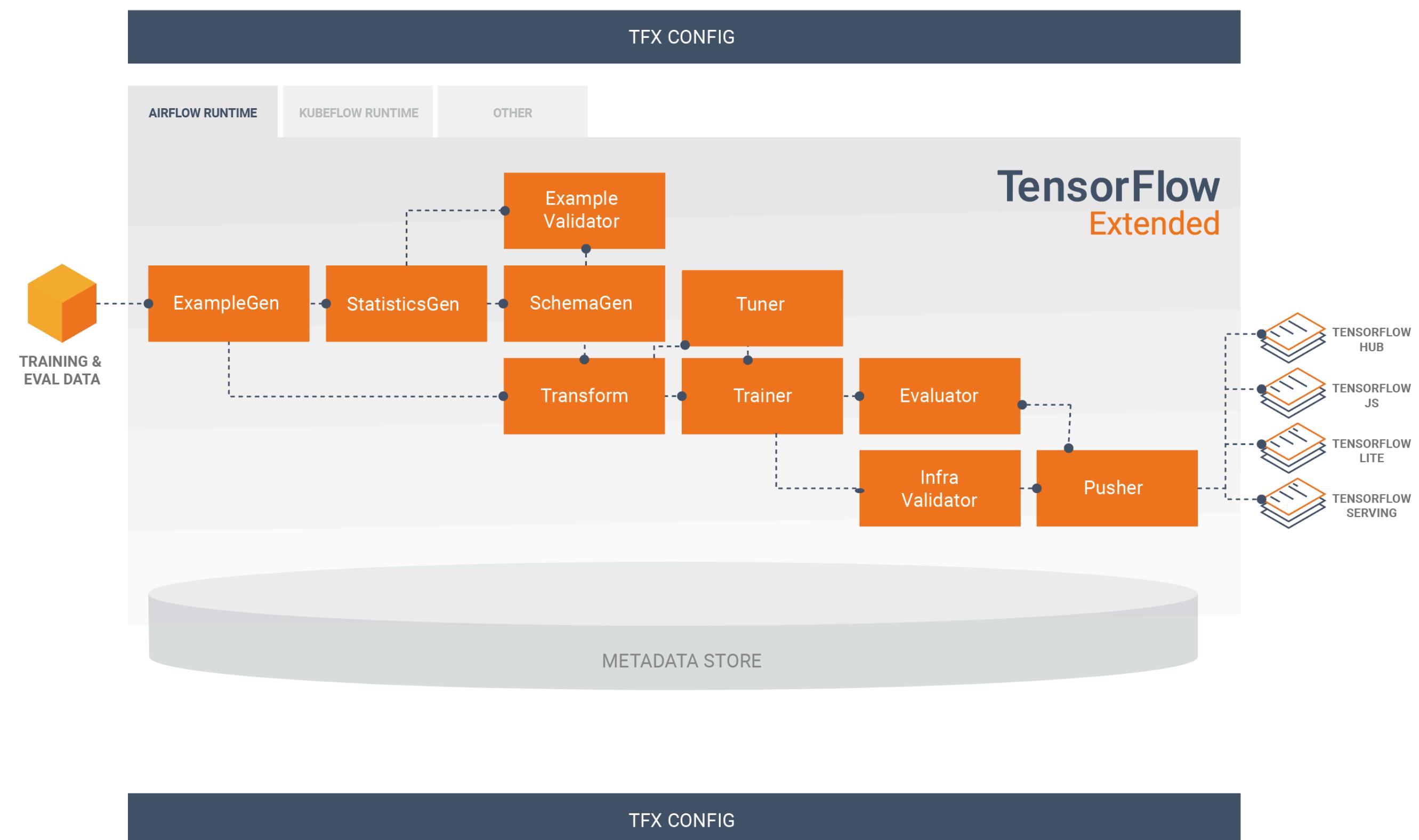
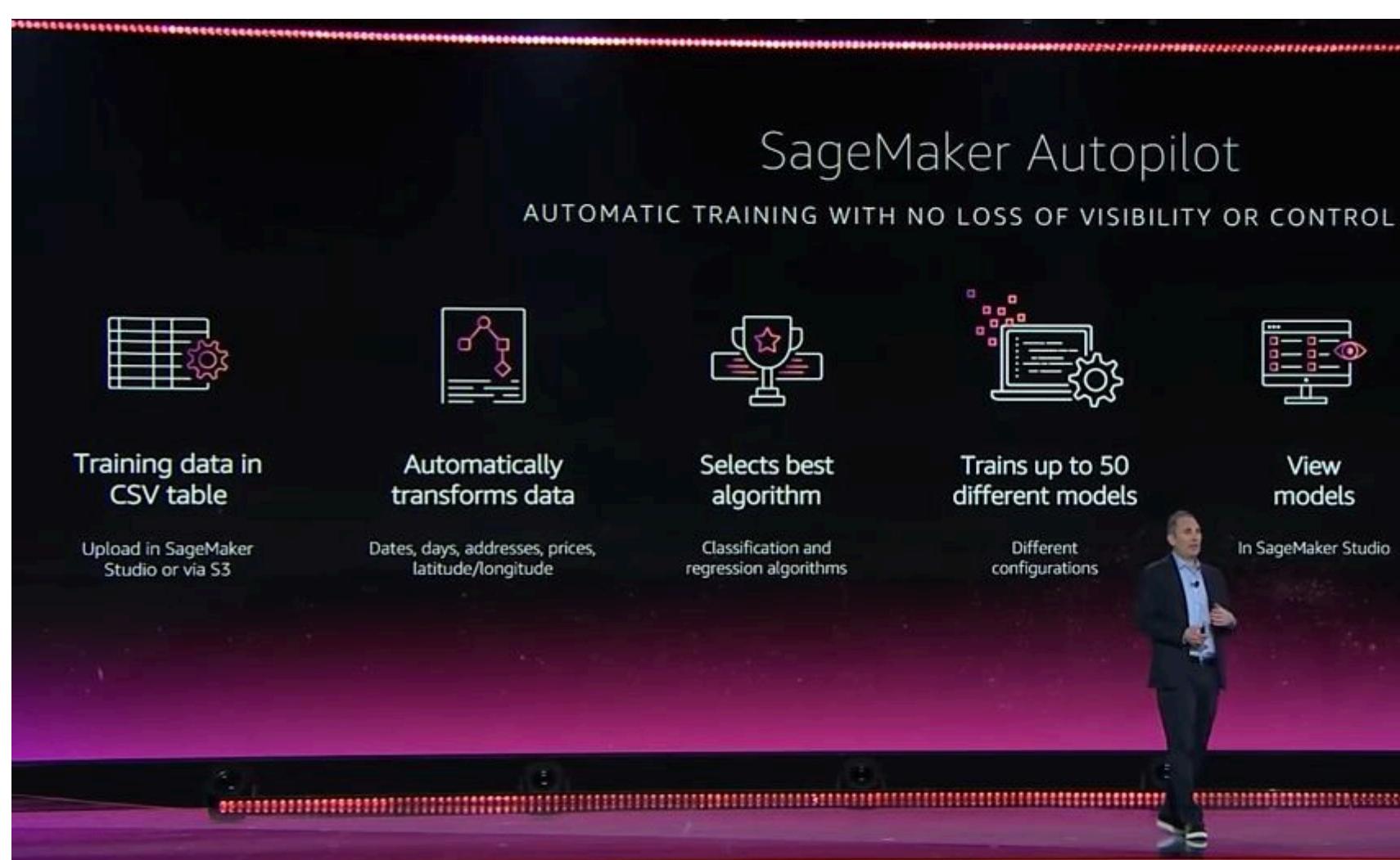
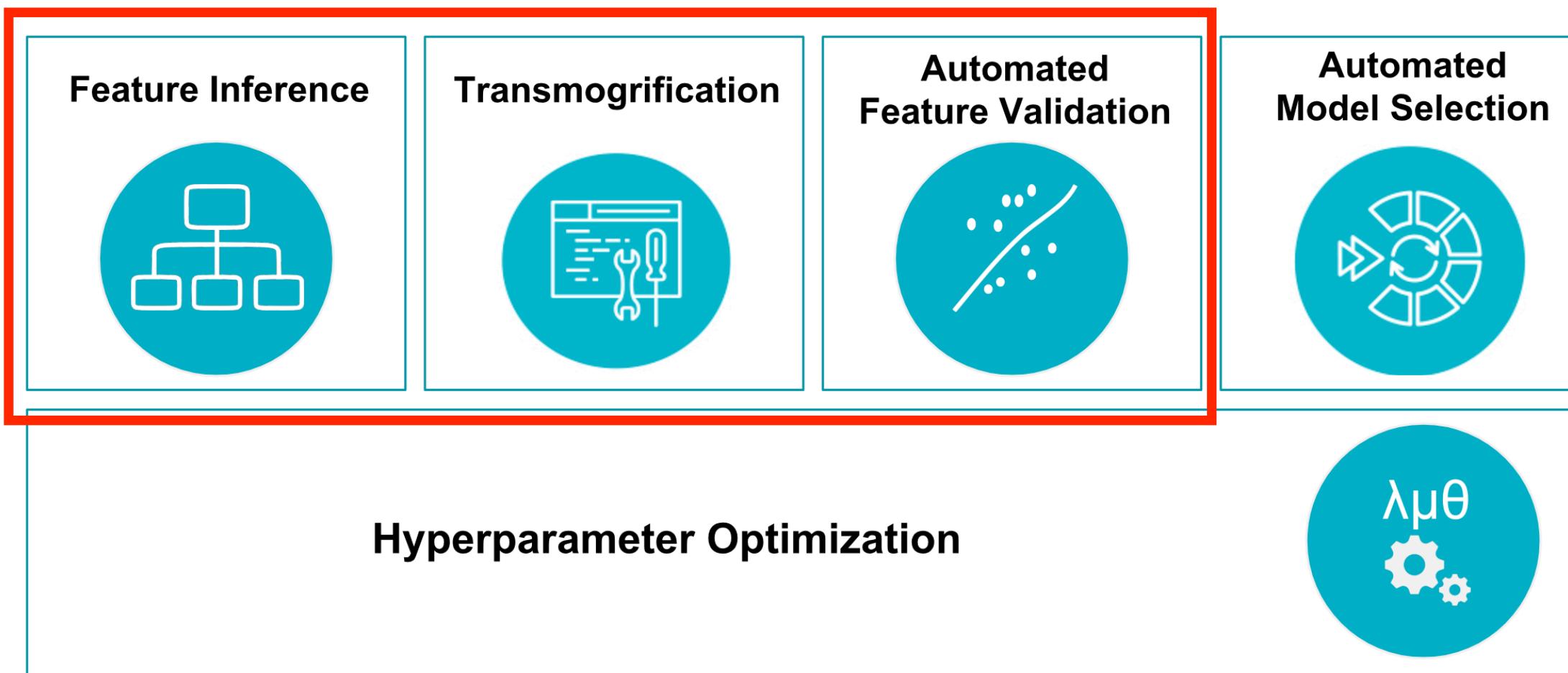
Example 2: Auto Data Prep for ML

TransmogrifAI



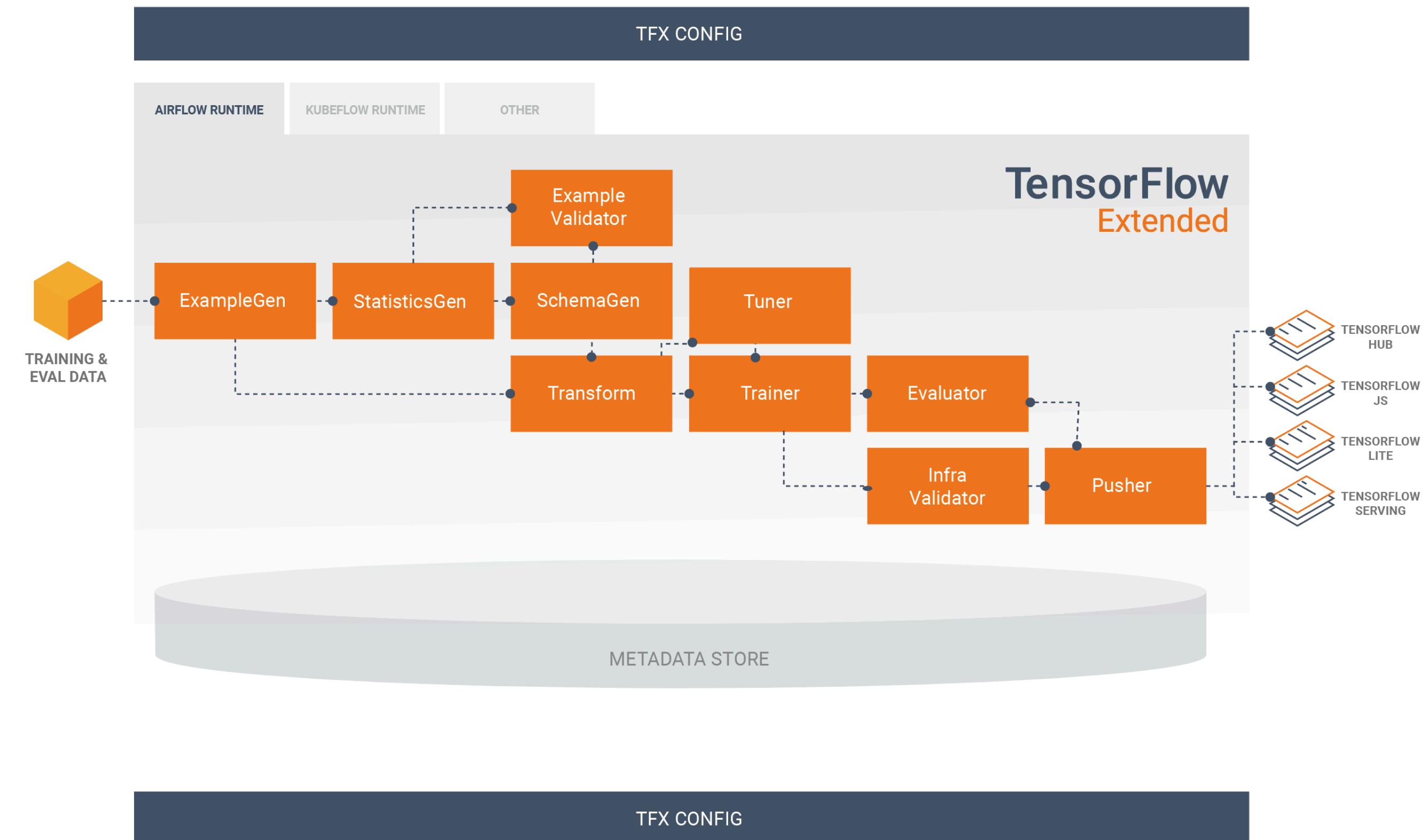
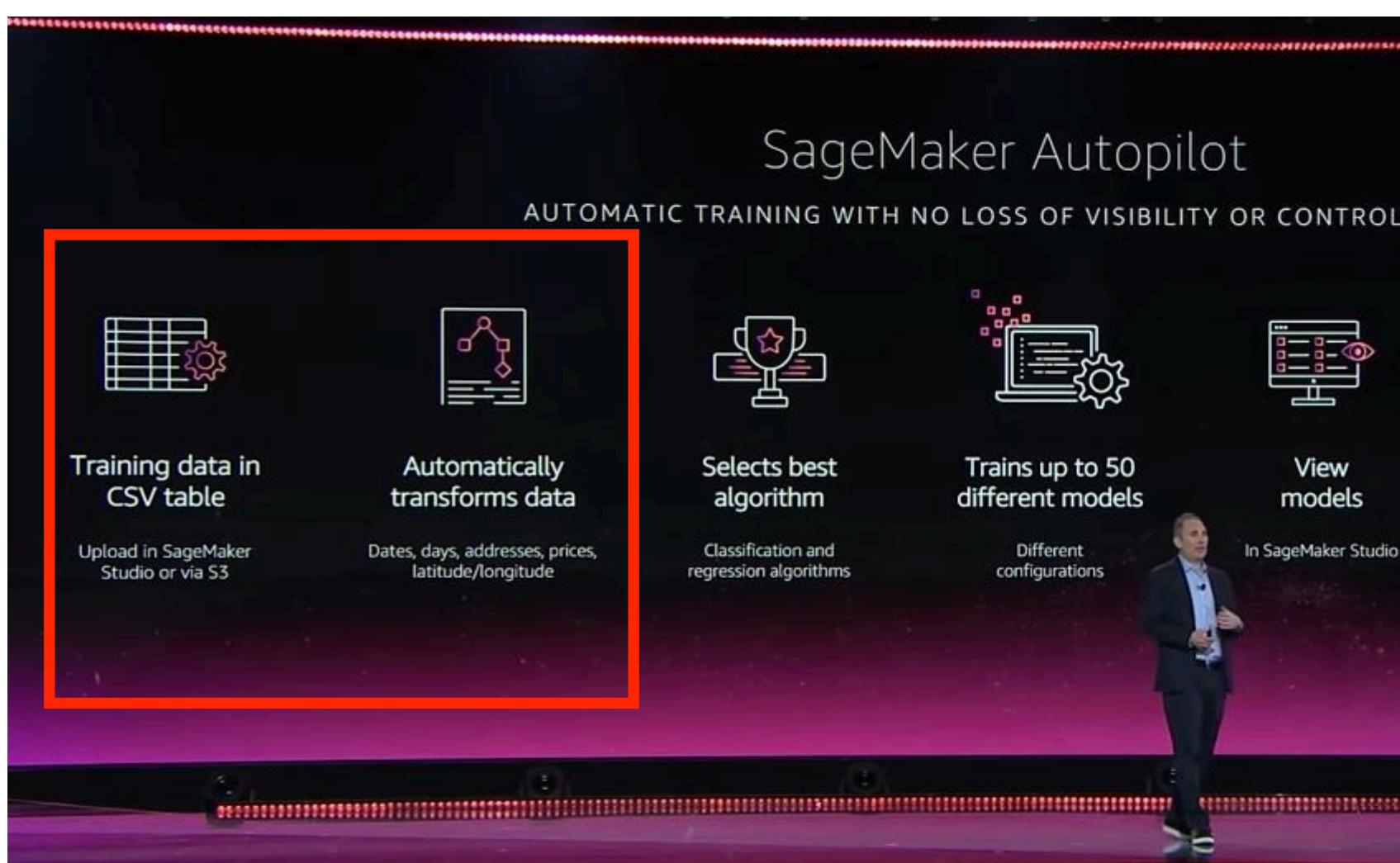
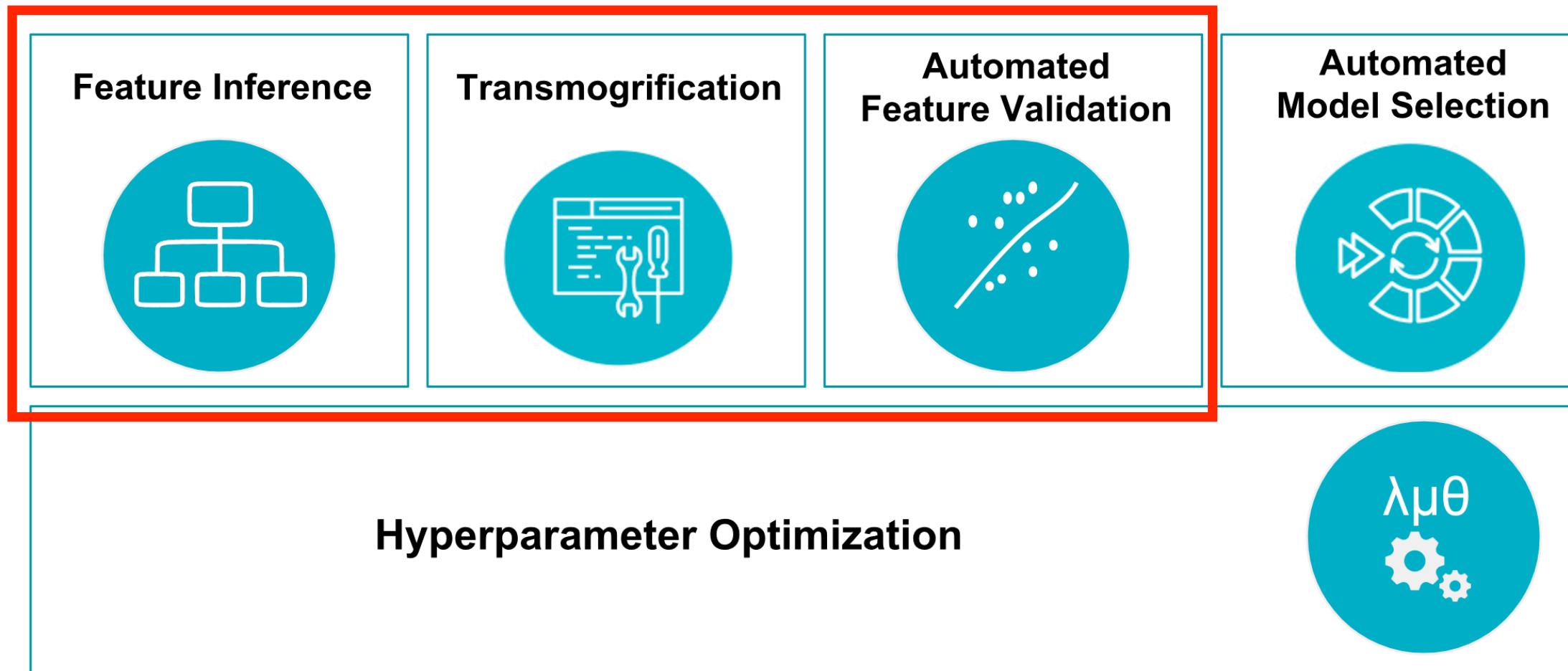
Example 2: Auto Data Prep for ML

TransmogrifAI



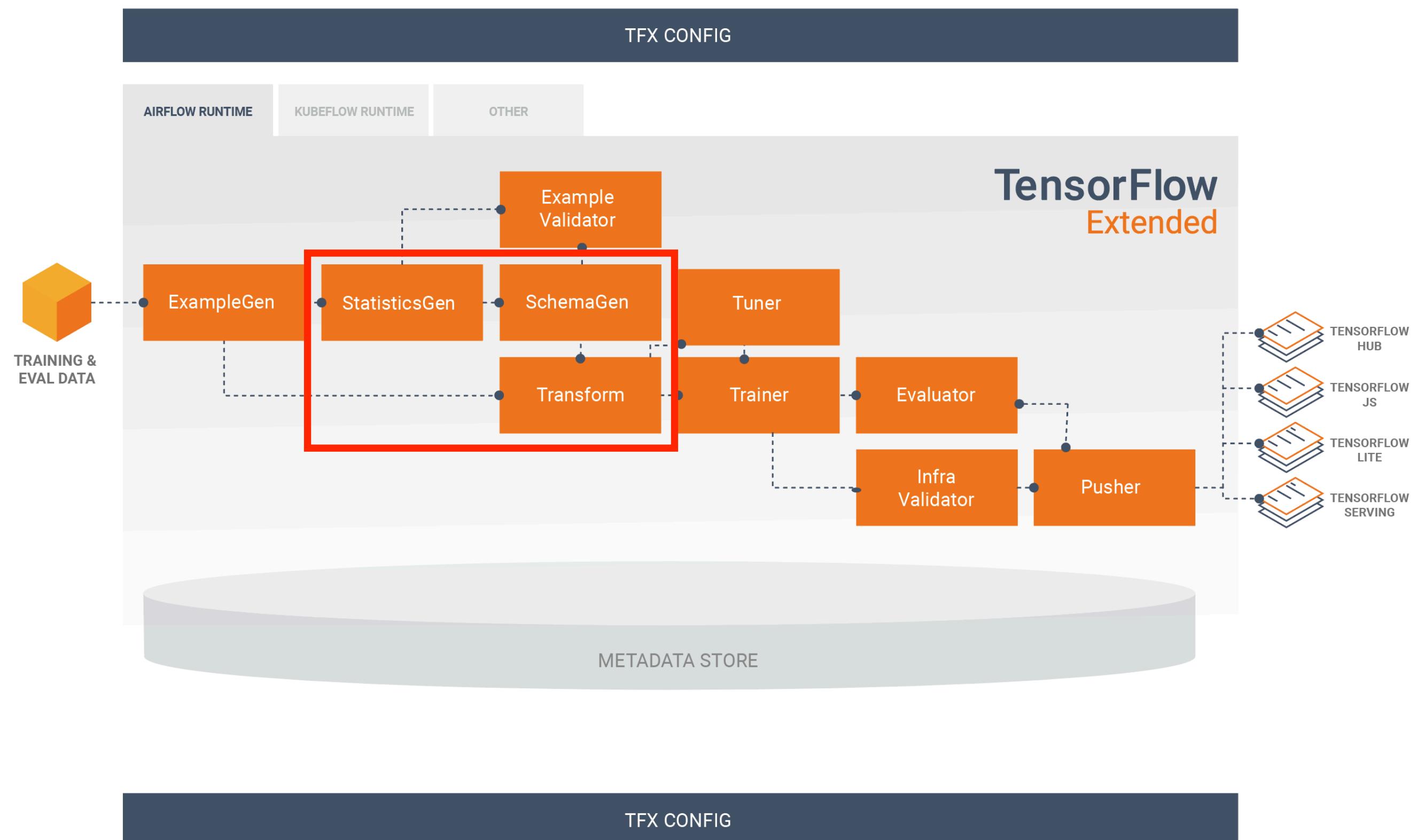
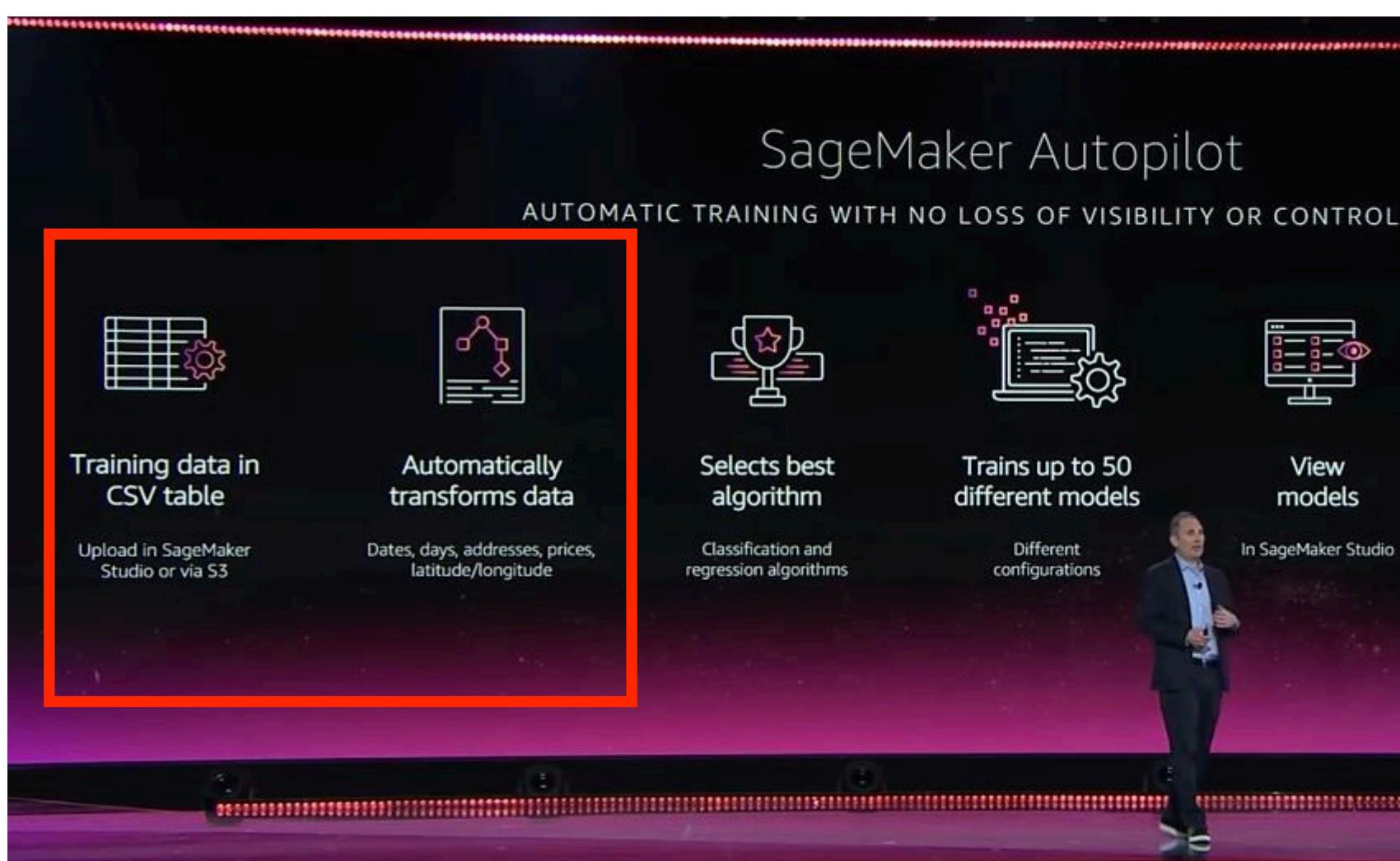
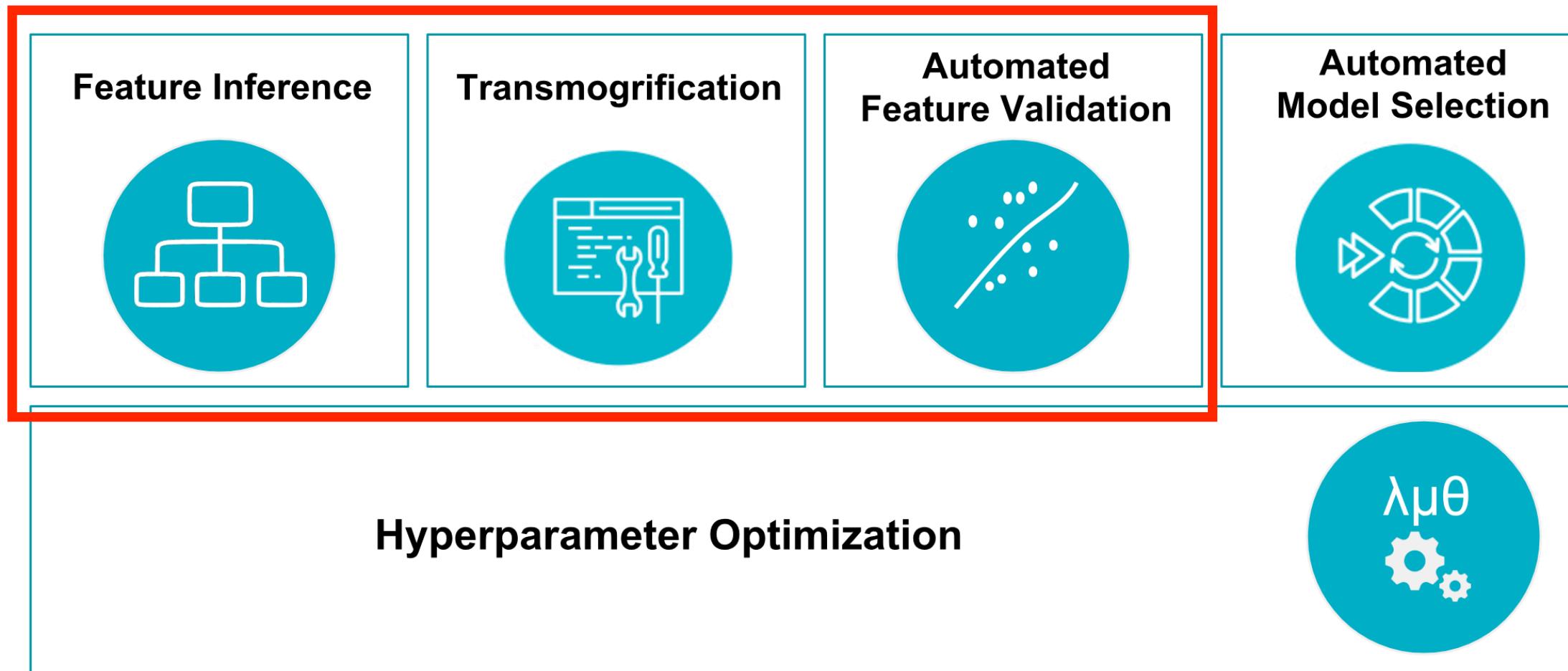
Example 2: Auto Data Prep for ML

TransmogrifAI



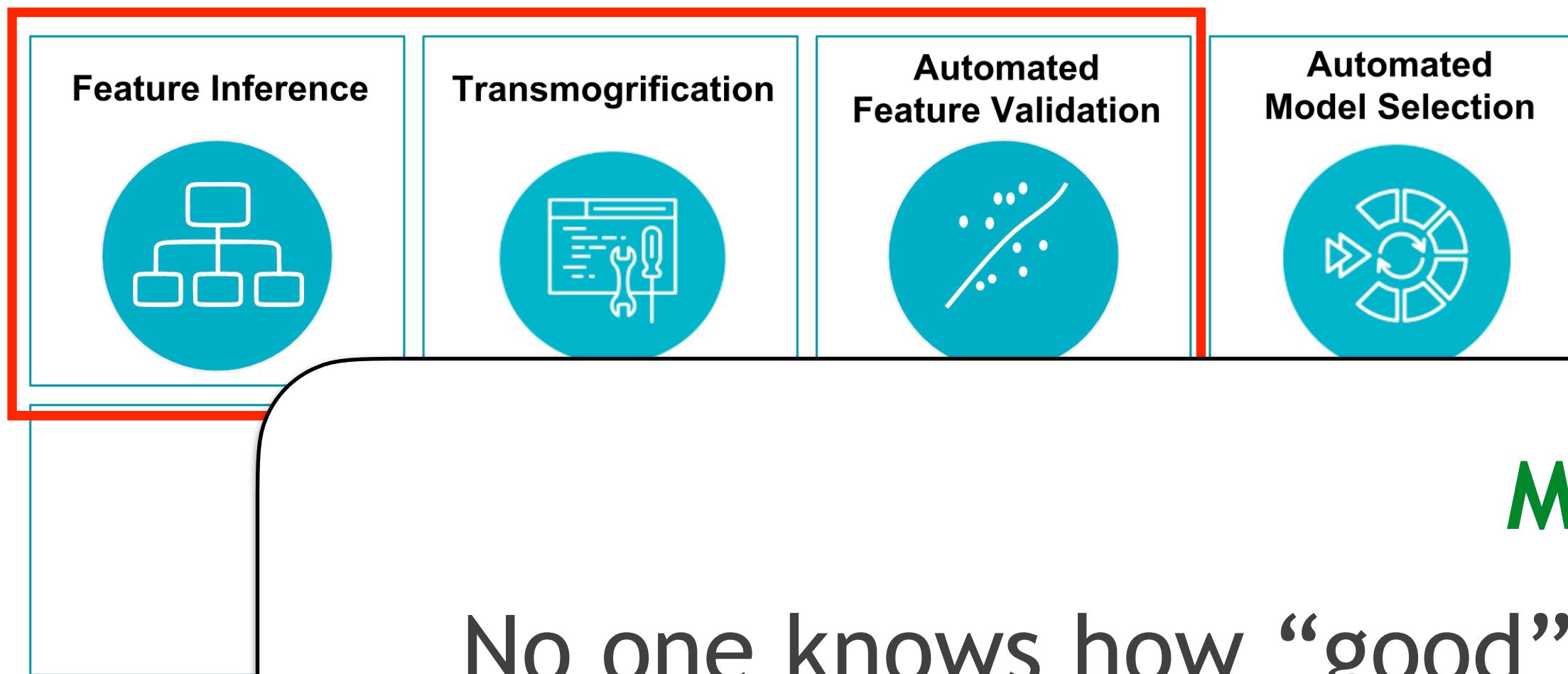
Example 2: Auto Data Prep for ML

TransmogrifAI



Example 2: Auto Data Prep for ML

TransmogrifAI

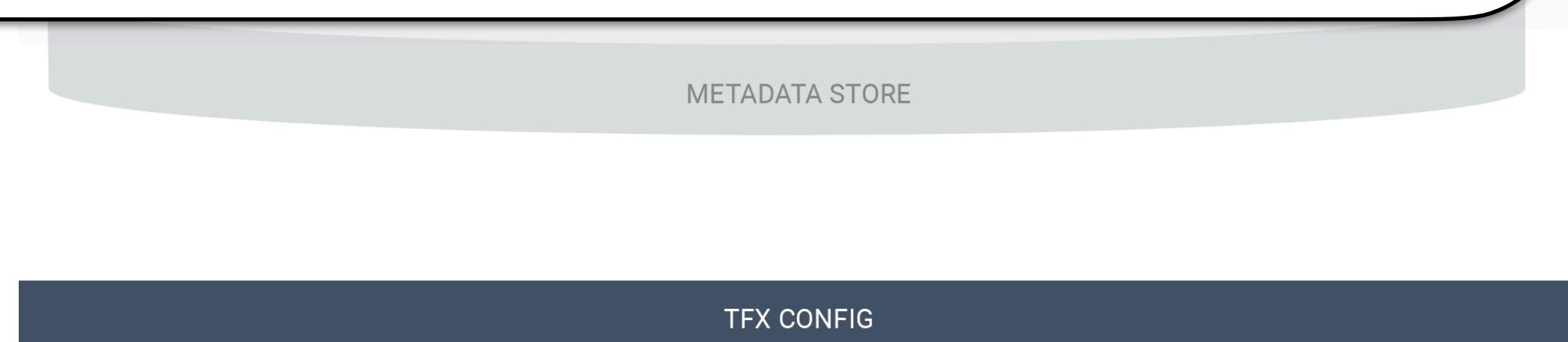
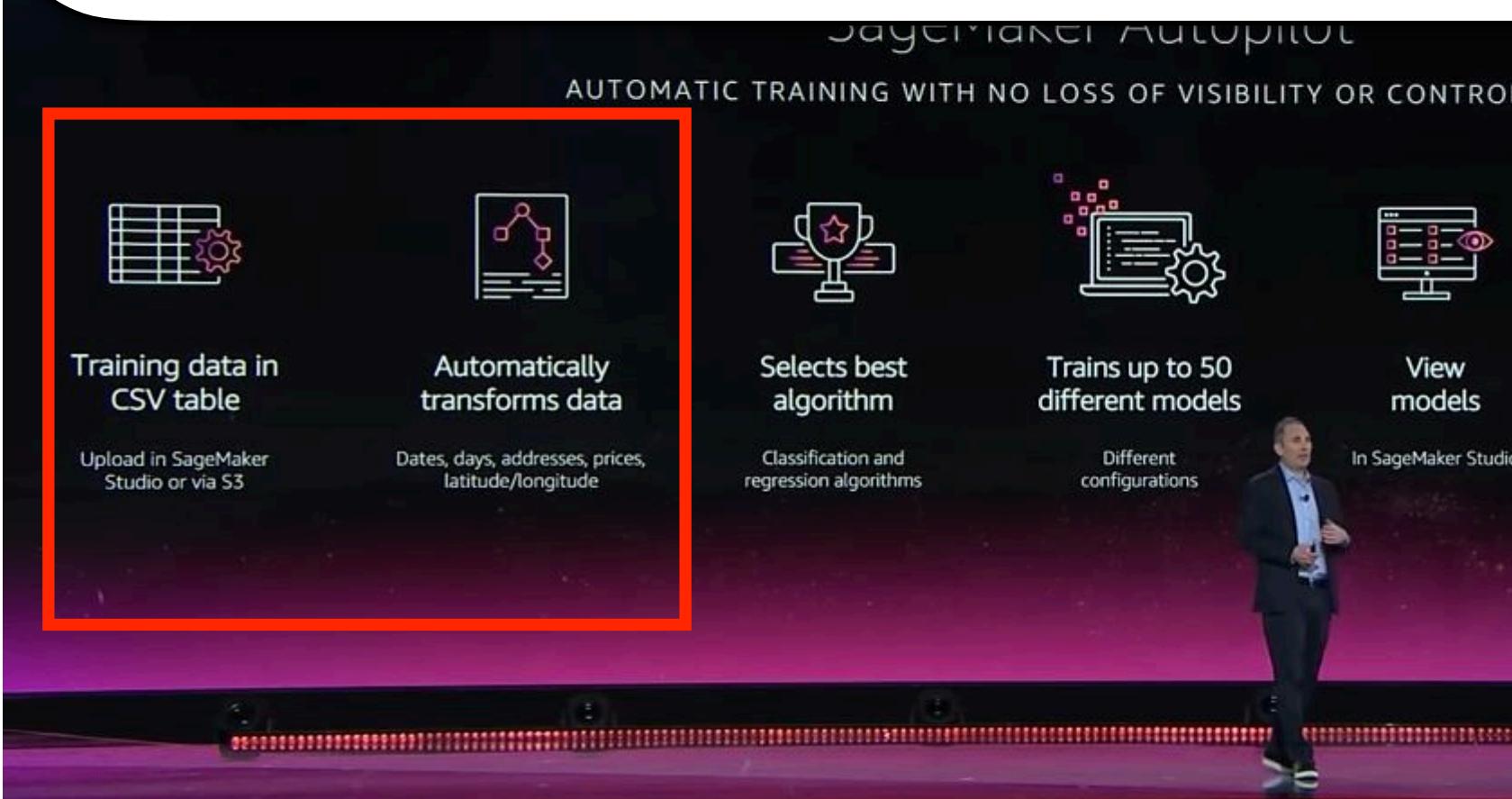


TensorFlow Extended

- TENSORFLOW HUB
- TENSORFLOW JS
- TENSORFLOW LITE
- TENSORFLOW SERVING

Major Issue:

No one knows how “good” objectively such auto data prep is!



Project SortingHat

Benchmarks for
ML Data Prep =

The ML Data Prep Zoo: Towards Semi-Automatic Data Preparation for ML. SIGMOD DEEM'19
Towards Benchmarking Feature Type Inference for AutoML Platforms. SIGMOD'21

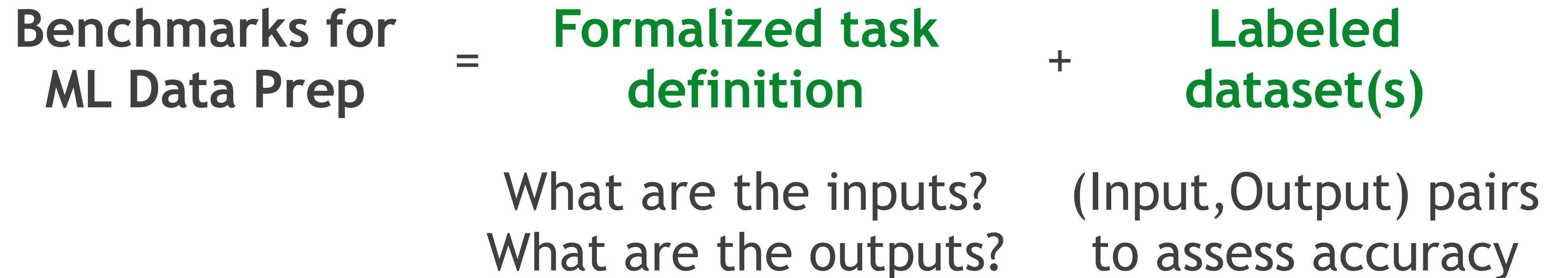
Project SortingHat

Benchmarks for
ML Data Prep = Formalized task
definition

What are the inputs?
What are the outputs?

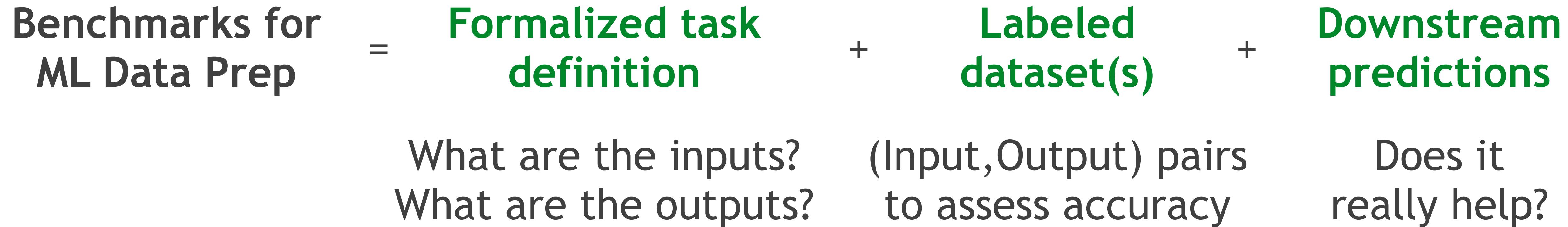
The ML Data Prep Zoo: Towards Semi-Automatic Data Preparation for ML. SIGMOD DEEM'19
Towards Benchmarking Feature Type Inference for AutoML Platforms. SIGMOD'21

Project SortingHat



*The ML Data Prep Zoo: Towards Semi-Automatic Data Preparation for ML. SIGMOD DEEM'19
Towards Benchmarking Feature Type Inference for AutoML Platforms. SIGMOD'21*

Project SortingHat



*The ML Data Prep Zoo: Towards Semi-Automatic Data Preparation for ML. SIGMOD DEEM'19
Towards Benchmarking Feature Type Inference for AutoML Platforms. SIGMOD'21*

Project SortingHat

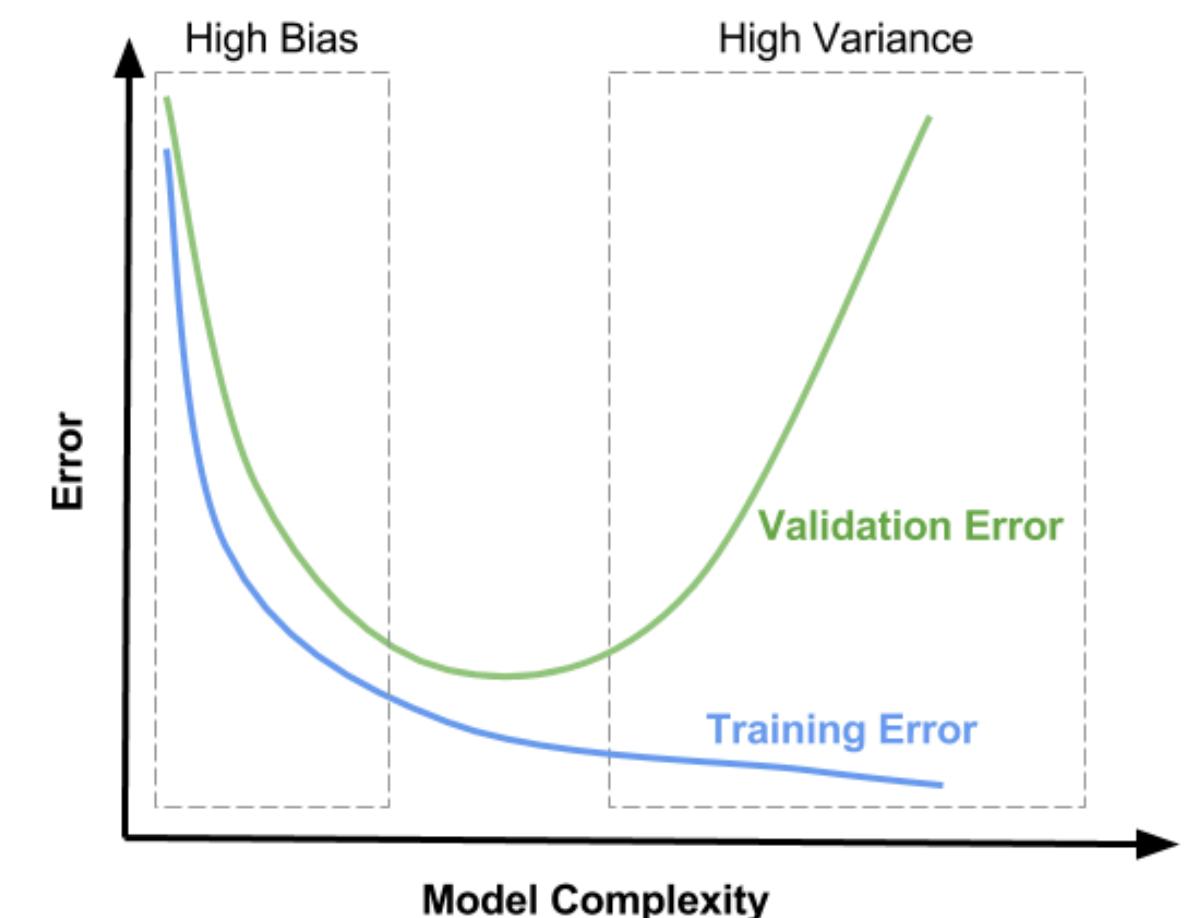
Benchmarks for ML Data Prep = **Formalized task definition** + **Labeled dataset(s)** + **Downstream predictions**

What are the inputs?
What are the outputs?

(Input,Output) pairs
to assess accuracy

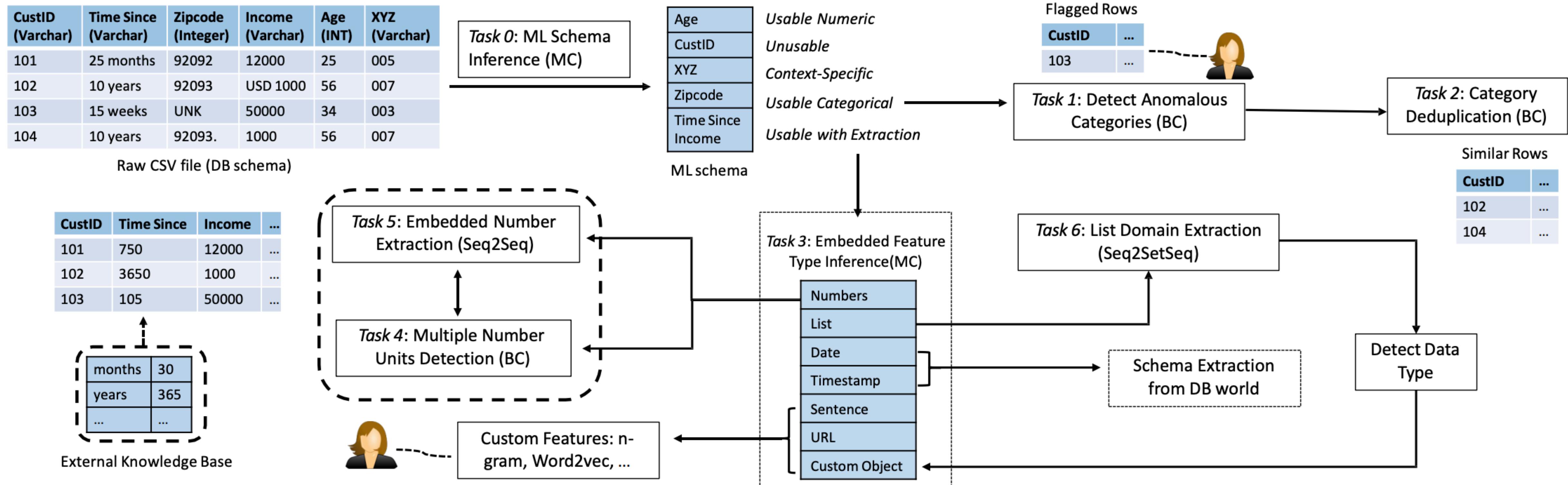
Does it
really help?

How exactly does data prep affect ML accuracy?
Are “failsafes” possible for errors in auto data prep?



Project SortingHat

Current major tasks covered in our ML Data Prep Zoo benchmarks:



The ML Data Prep Zoo: Towards Semi-Automatic Data Preparation for ML. **SIGMOD DEEM'19**
Towards Benchmarking Feature Type Inference for AutoML Platforms. **SIGMOD'21**

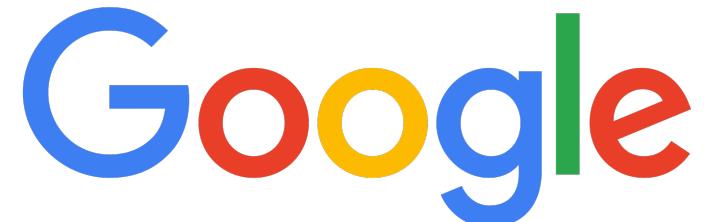
SortingHat: Early Impact and Trajectory

Integrated with TFDV in TFX in collaboration with Google

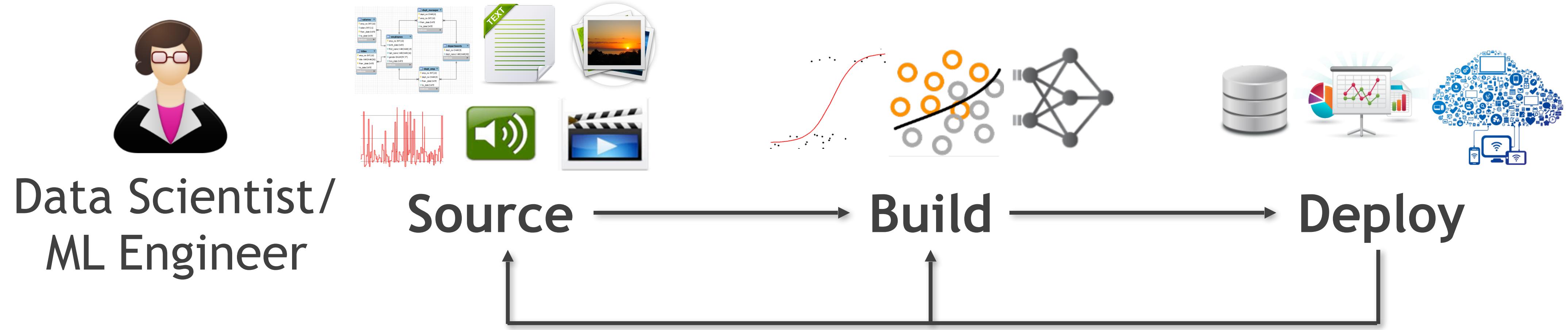
Being tested on Google's internal AutoML benchmarks

Our model used in NSF-funded “AI Maker” model + dataset search tool

Coming soon: AWS DBMS and ML/AI products; OpenML; catalog search



The New DBfication of ML/AI



Metadata Management for ML
Data Prep/Cleaning for ML
Multimodal ML Query Models
Data Search, Labeling, etc.

...

Benchmark Frameworks and Data
Fairness, Transparency, Privacy, etc.

Scalable Data Systems for ML
Query Optimization for ML
Cloud and Streaming Infra.
Provenance and Debugging

...

*Leaving these problems open => Huge waste of time/
effort/money/energy/etc. by users of ML/AI!*

*Leaving these problems open => Huge waste of time/
effort/money/energy/etc. by users of ML/AI!*

*It is high time for DB, ML/AI, and more areas to join
forces to accelerate the DBfication of ML/AI!*

Outline

■ The New DBfication of ML/AI

■ Two Examples from My Research

■ Accelerating the DBfication of ML/AI

Accelerating the DBfication of ML/AI

1. Learn the fundamentals of ML/AI algorithms and theory.

Kinda like learning logic, RA, SQL, etc. for RDBMSs

Review ML/AI algorithms courses in your institution or online

3 key books: Hastie et al. (Stat. ML); Mitchell (ML); Courville et al. (DL)

Need-to-know spectrum: DL for DL sys.; ML theory for accuracy tradeoffs

Accelerating the DBfication of ML/AI

1. Learn the fundamentals of ML/AI algorithms and theory.

Kinda like learning logic, RA, SQL, etc. for RDBMSs

Review ML/AI algorithms courses in your institution or online

3 key books: Hastie et al. (Stat. ML); Mitchell (ML); Courville et al. (DL)

Need-to-know spectrum: DL for DL sys.; ML theory for accuracy tradeoffs

2. Check out my “DB for ML” grad course and research book.

CSE 234/291: *Data Systems for Machine Learning*. UC San Diego. 2020.

Data Management in Machine Learning Systems. Morgan & Claypool Publishers. 2019.

Accelerating the DBfication of ML/AI

3. Check out recent “DB for ML” tutorials and papers.

Data Management in Machine Learning: Challenges, Techniques, and Systems. SIGMOD 2017.

Data Management Challenges in Production Machine Learning. SIGMOD 2017.

Data Integration and Machine Learning: A Natural Synergy. VLDB 2018.

Data Collection and Quality Challenges for Deep Learning. VLDB 2020.

Accelerating the DBfication of ML/AI

3. Check out recent “DB for ML” tutorials and papers.

Data Management in Machine Learning: Challenges, Techniques, and Systems. SIGMOD 2017.

Data Management Challenges in Production Machine Learning. SIGMOD 2017.

Data Integration and Machine Learning: A Natural Synergy. VLDB 2018.

Data Collection and Quality Challenges for Deep Learning. VLDB 2020.

4. Check out SIGMOD DEEM and HILDA Workshops. Check out MLSys.

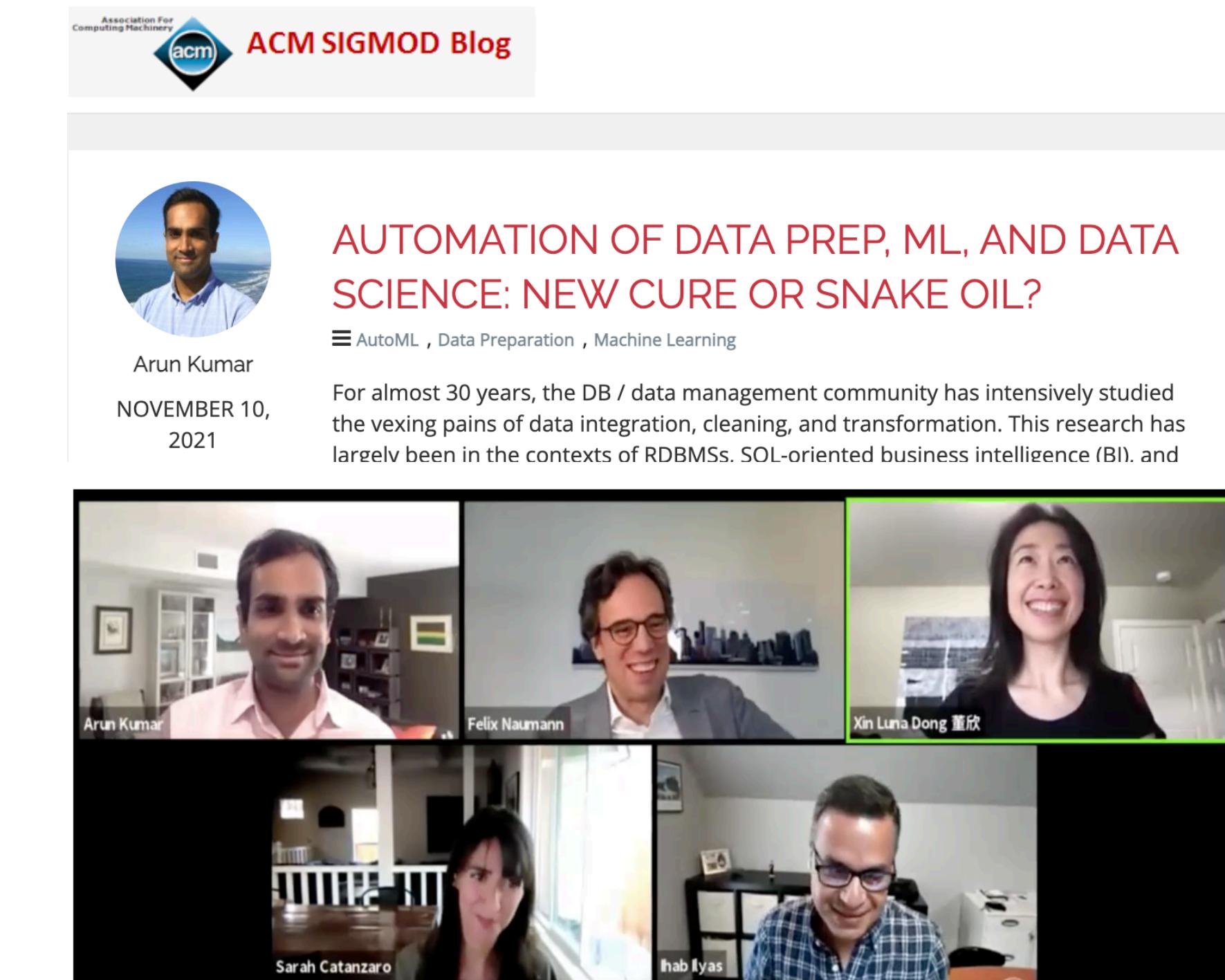
Accelerating the DBfication of ML/AI

5. Check out topical panel discussions on “DB for ML” stuff.



Photos from the workshop. L to R: (1) The DEEM audience. (2) The Panelists: Matei, Joaquin, Jens, Joey, and Manasi (Martin not pictured). (3) Advocatus Diaboli.

SIGMOD DEEM'18



SIGMOD'21

Accelerating the DBfication of ML/AI

6. MOST IMPORTANT: Speak/collaborate with ML/AI users, build REAL stuff, and help transfer research to practice.

Data scientists, Data analysts, ML engineers, MLOps engineers, etc.

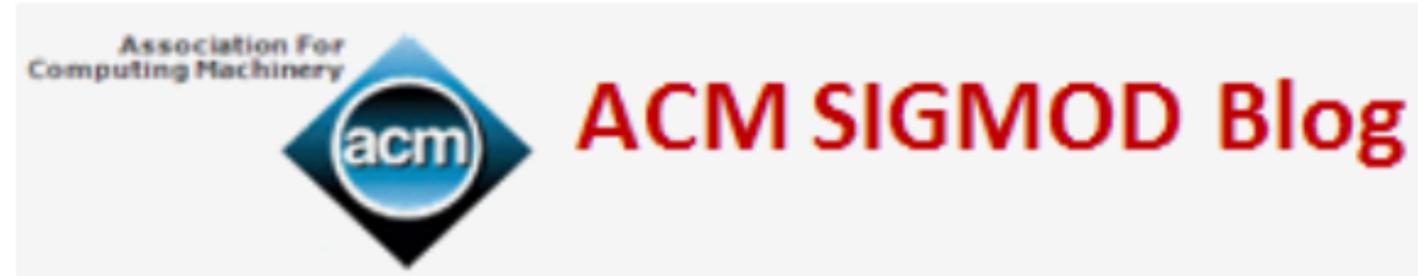
Domain sciences, enterprises, Web companies, cloud vendors, policy, etc.

Create open source artifacts, both software and datasets

Attend/speak at industry venues: Spark/Data+AI Summit, FOSDEM, etc.

...

New Publication Avenues



Alon Halevy, Arun
Kumar and
Nesime Tatbul

FEBRUARY 10,
2020

SCALABLE DATA SCIENCE: A NEW RESEARCH TRACK CATEGORY AT PVLDB VOL 14 / VLDB

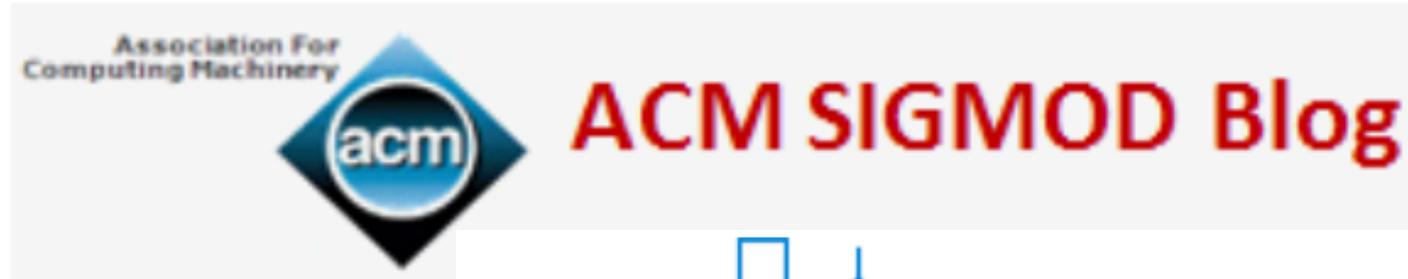
2021

≡ Uncategorized

This post introduces and explains the newly created category of “Scalable Data Science” within the Research Track of PVLDB. This category comes into effect for volume 14, i.e., submissions starting April 1, 2020, which will be evaluated by the Review Board of PVLDB vol 14 for presentation at VLDB 2021.

The Growth of Data Science

New Publication Avenues



Alon Halev
Kumar
Nesime
FEBRUARY
2020

SIGMOD 2022 CALL FOR RESEARCH PAPERS

- **Data Science Track**

We invite the submission of original research in data science targeting the entire data life cycle of real applications. This data life cycle encompasses databases/data management/data systems/data engineering often leveraging statistical, Machine Learning and Artificial Intelligence methods and using massive and heterogeneous collections of potentially messy datasets. Data science papers study phenomena at scales and granularities never before possible. Such papers are expected to focus on data-intensive components of data science pipelines; and solve problems in areas of interest to the

ARCH
LDB

ata
t for
y the

New Publication Avenues

The screenshot shows the homepage of the ACM SIGMOD Blog for the Fifth Conference on Machine Learning and Systems. The page features a sidebar on the left with navigation links for Year (2022) ▾, Help ▾, My Registrations, Profile ▾, Contact Us, Sponsor Info, Conflicts of Interest, Code of Conduct, Proceedings, and MLSys. The main content area includes the conference title, location (Santa Clara Convention Center), and dates (Mon Apr 11th through Thu the 14th, 2022). It also features sections for Registration (with Pricing and Register starting Feb 06 01 PM PST buttons), Conference Overview (describing the intersection of machine learning and systems research), and Schedule (with a Video Library 2021 button). A social media sharing button for Twitter is present.

phenomena at scales and granularities never before possible. Such papers are expected to focus on data-intensive components of data science pipelines; and solve problems in areas of interest to the

New Publication Avenues

The screenshot shows the ACM SIGMOD Blog website. At the top left is the ACM logo and the text "Association For Computing Machinery". The main title "ACM SIGMOD Blog" is in red. Below it, the text "Fifth Conference on Machine Learning and Systems" is displayed. A sidebar on the left has links for "Year (2022) ▾", "Help ▾", and "My Registrations". The main content area shows the location as "Santa Clara Convention Center" and the dates as "Mon Apr 11th through Thu the 14th, 2022". Below this, there are two buttons: "Registration" and "Conference Overview". At the bottom, there is a navigation bar with links for "HOME", "PROGRAM", "ATTENDING", "CALLS", "KDD CUP", "SPONSORS", and "ORGANIZERS". To the left of the main content area, there is a logo for "KDD2021" featuring a blue geometric pattern.

Call for Applied Data Science Track Papers

MLSys

phenomena at scales and granularities never before possible. Such papers are expected to focus on data-intensive components of data science pipelines; and solve problems in areas of interest to the

- applications
- Data preparation, feature selection, and feature extraction

I hope the DB community steps up to the grand challenge of DBfication of ML/AI.

I hope the DB community steps up to the grand challenge of DBfication of ML/AI.

Partner with folks in ML/AI, systems, HCI etc; work with ML/AI users in domain sciences and industry.

My Terrific Advisees



Yuhao Zhang
PhD & MS



Kabir Nagrecha
PhD & BS



Xiuwen Zheng
PhD & MS



Supun Nakandala
PhD'22 -> Databricks



Pradyumna Sridhara
MS



Vignesh Nanda Kumar
MS



Kyle Luoma
PhD



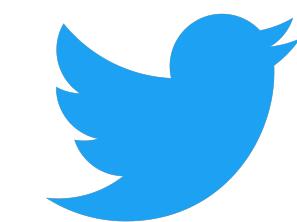
Vraj Shah
PhD'22 -> IBM
Research Almaden

<https://ADALabUCSD.github.io>

arunkk@eng.ucsd.edu



github.com/ADALabUCSD



@TweetAtAKK

ACKS:

