

Intelligent Services

Serving Machine Learning

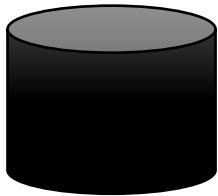
Joseph E. Gonzalez

jegonzal@cs.berkeley.edu; Assistant Professor @ UC Berkeley

joseph@dato.com; Co-Founder @ Dato Inc.

Contemporary Learning Systems

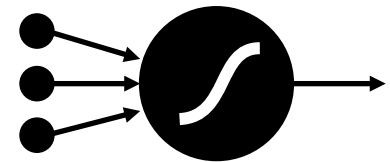
Big
Data



Training



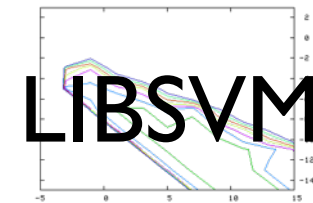
Big
Models



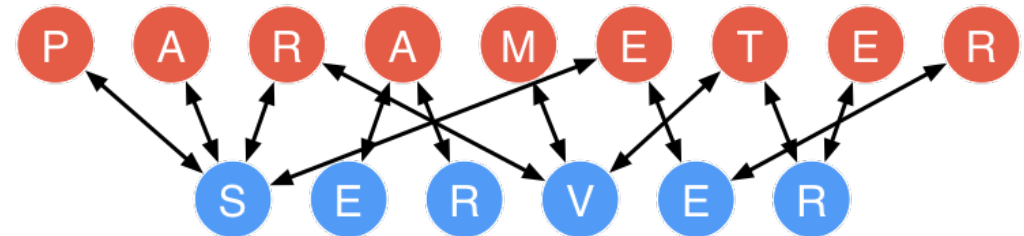
Contemporary Learning Systems



BIDMach

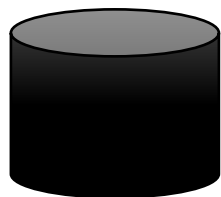


H₂O



What happens *after* we train a model?

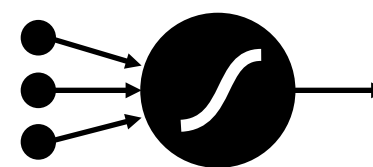
Data



Training



Model



Conference
Papers



Dashboards and
Reports



Drive Actions



What happens *after* we train a model?

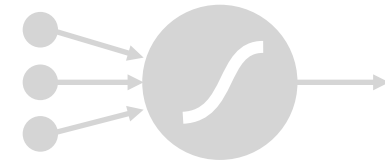
Data



Training



Model



Conference
Papers



Dashboards and
Reports



Drive Actions



Suggesting Items
at Checkout



Fraud
Detection



Cognitive
Assistance



Internet of
Things



Low-Latency



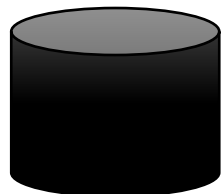
Personalized



Rapidly Changing



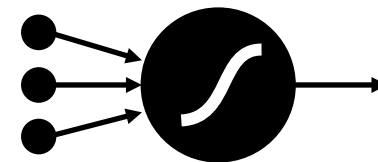
Data

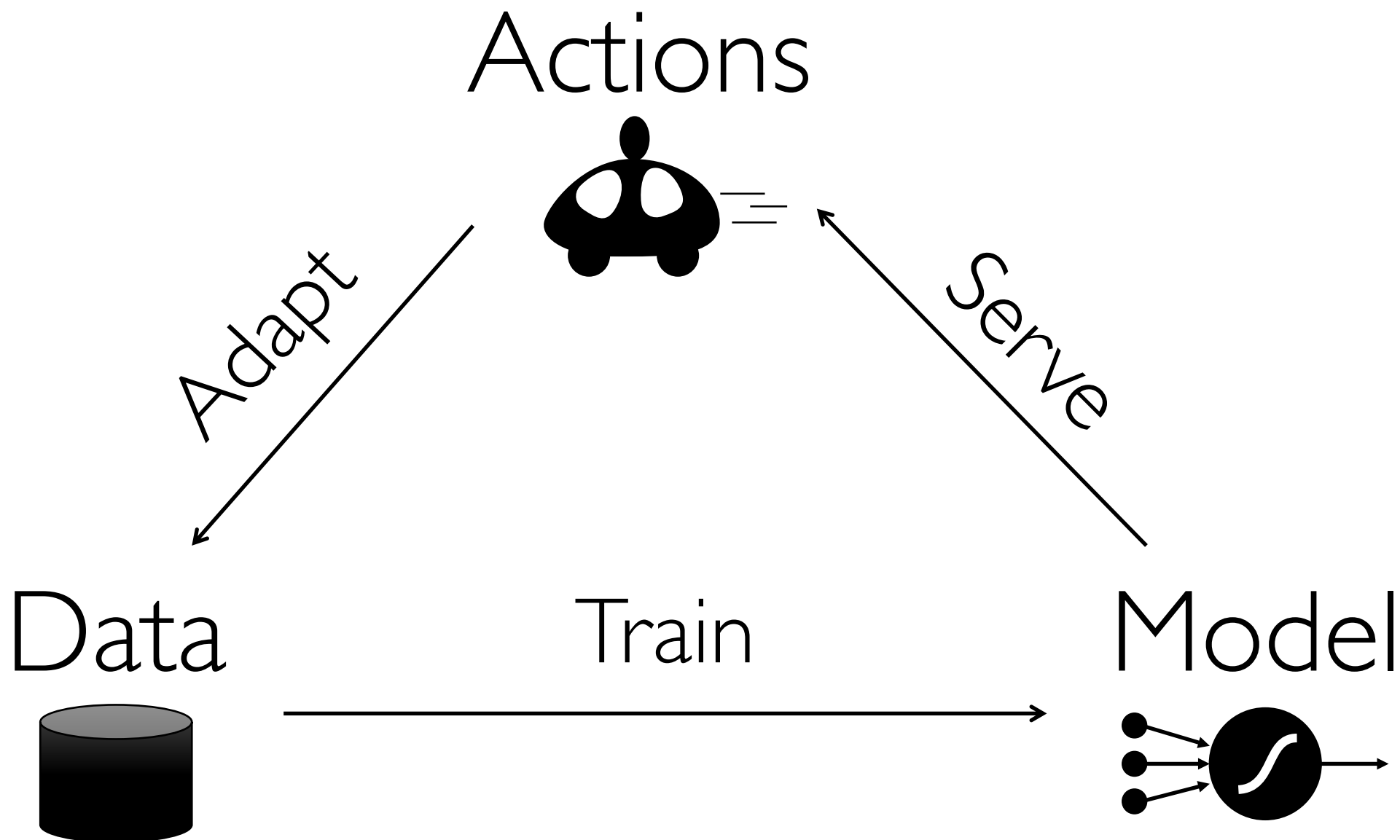


Train



Model



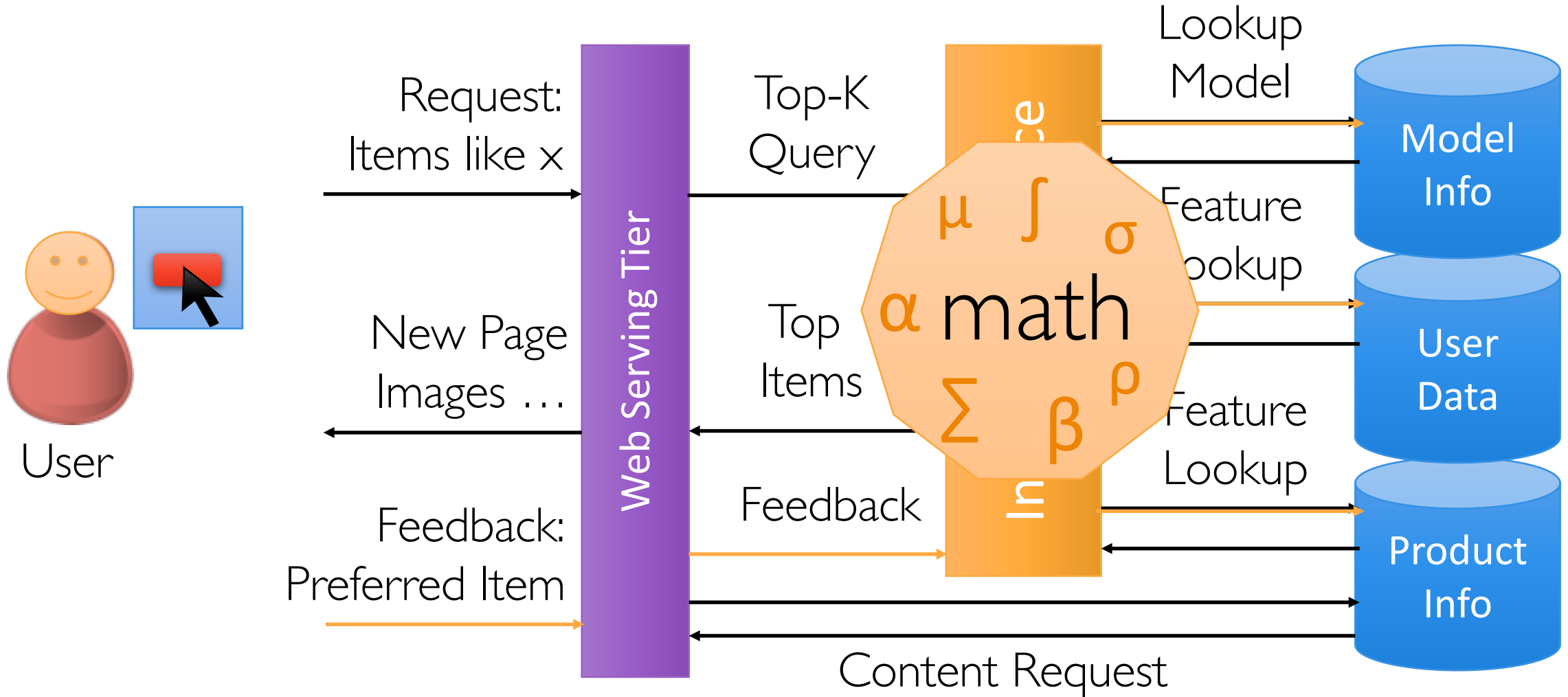


Machine
Learning



Intelligent
Services

The Life of a Query in an Intelligent Service



Essential Attributes of Intelligent Services

Responsive

Intelligent applications
are interactive

Adaptive

ML models out-of-date the
moment learning is done

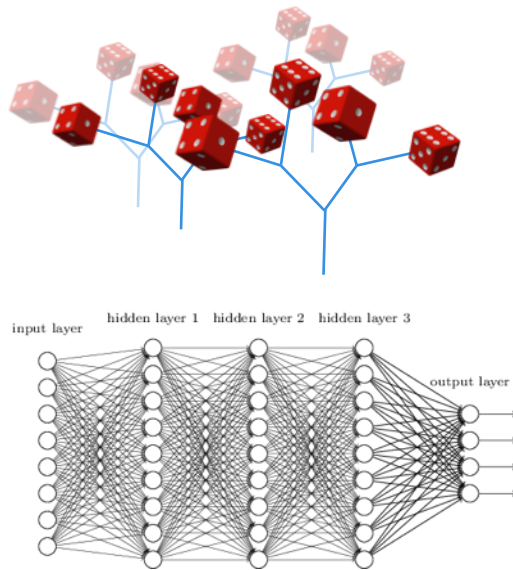
Manageable

Many models
created by multiple people

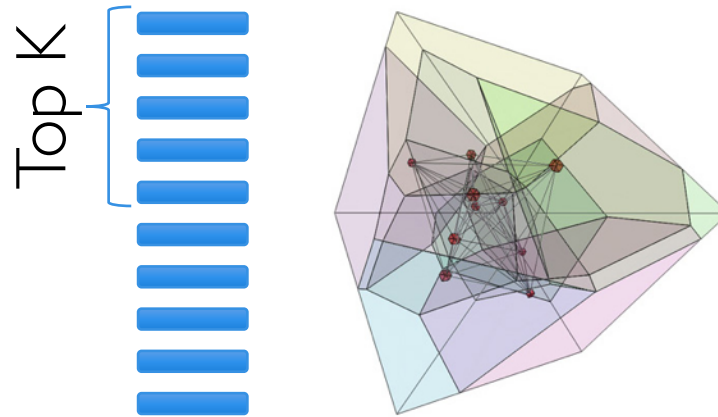
Responsive: *Now and Always*

Compute predictions in $< 20\text{ms}$ for complex

Models



Queries

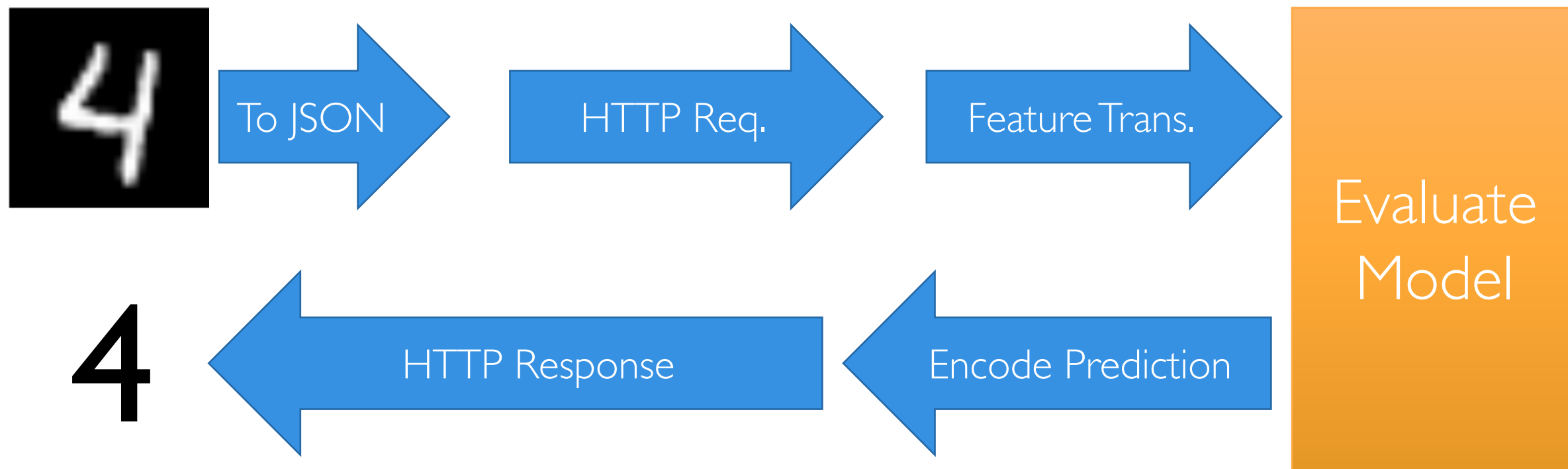


Features

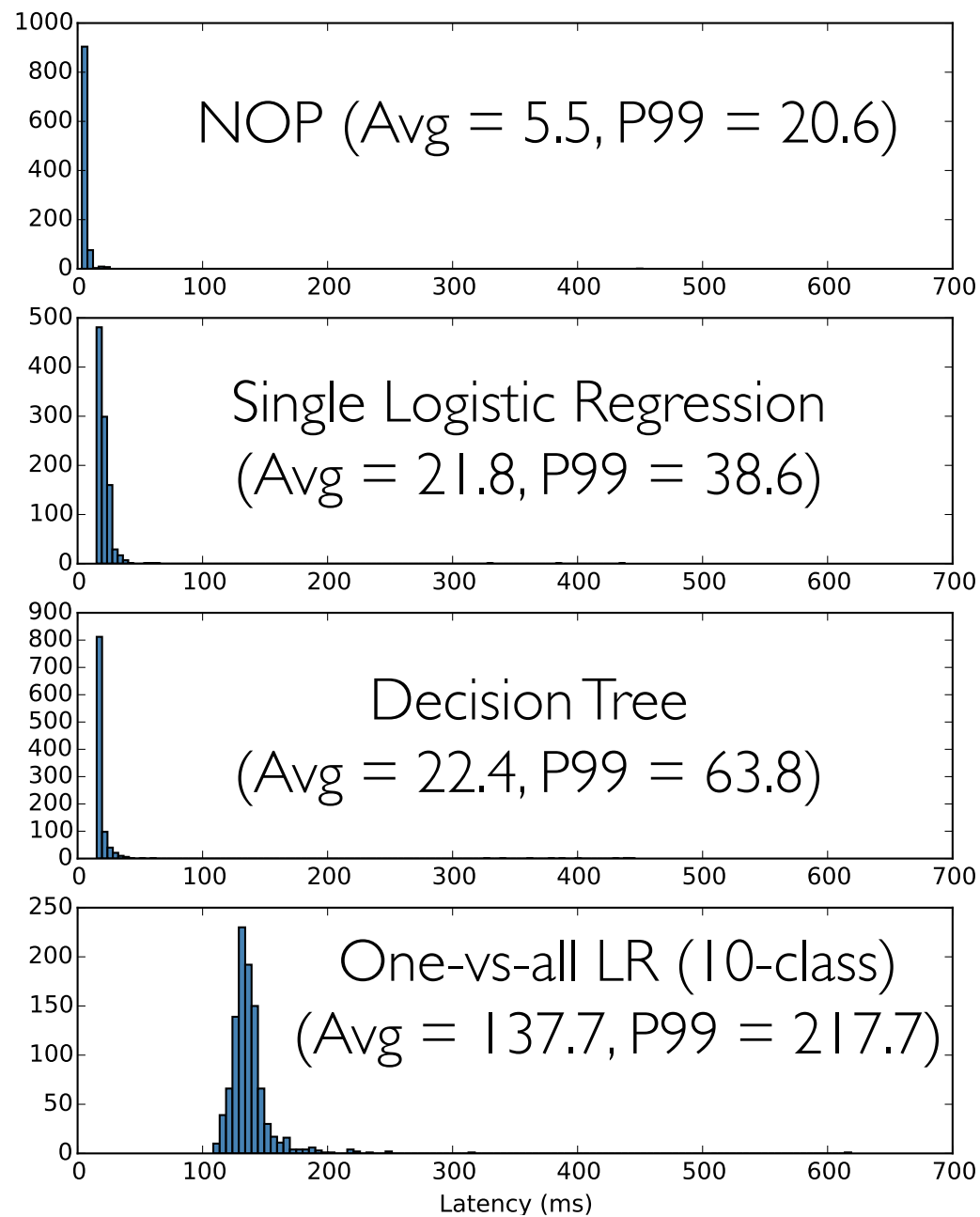
```
SELECT * FROM  
users JOIN items,  
click_logs, pages  
WHERE ...
```

under heavy *query load* with system *failures*.

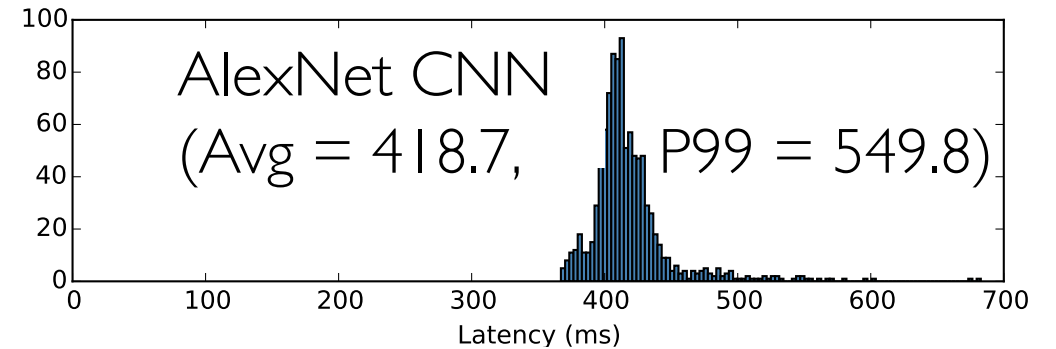
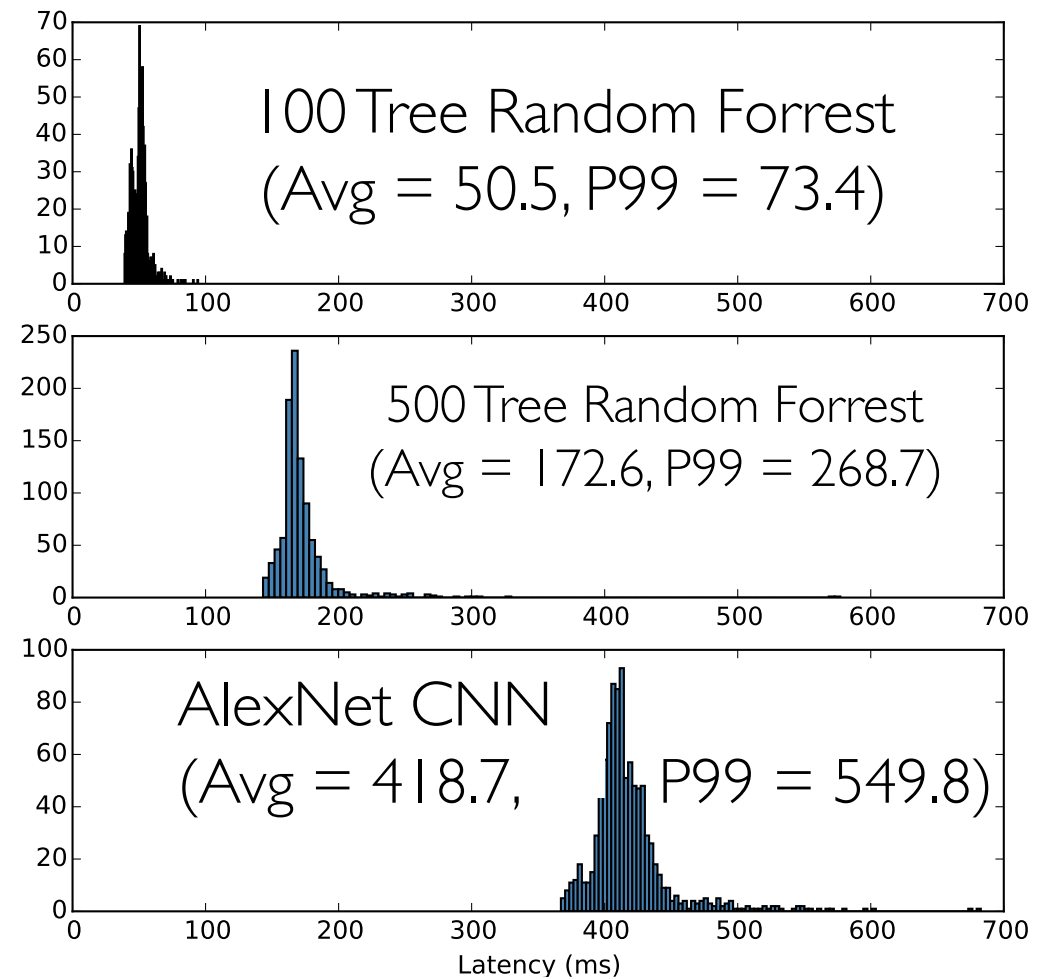
Experiment: End-to-end Latency in Spark MLlib



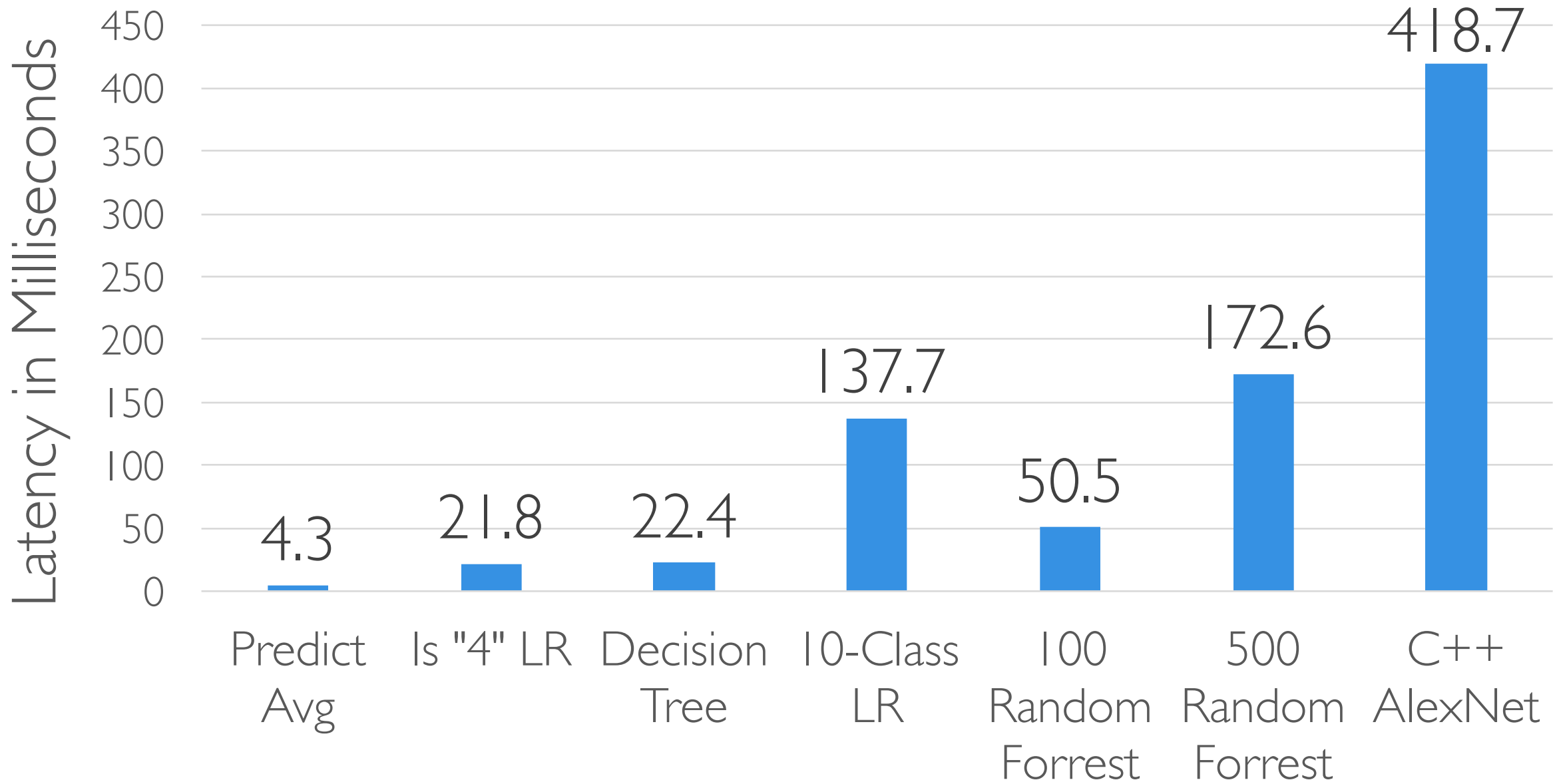
Count out of 1000



End-to-end Latency for Digits Classification
784 dimension input
Served using MLib and Dato Inc.



Latency measured in milliseconds



Adaptive to Change at All Scales

Population

Granularity of Data

Session



Shopping
for Mom

Shopping
for Me



Months

Rate of Change

Minutes

Adaptive to Change at All Scales

Population

Granularity of Data

Session



Law of Large Numbers
→ Change Slow

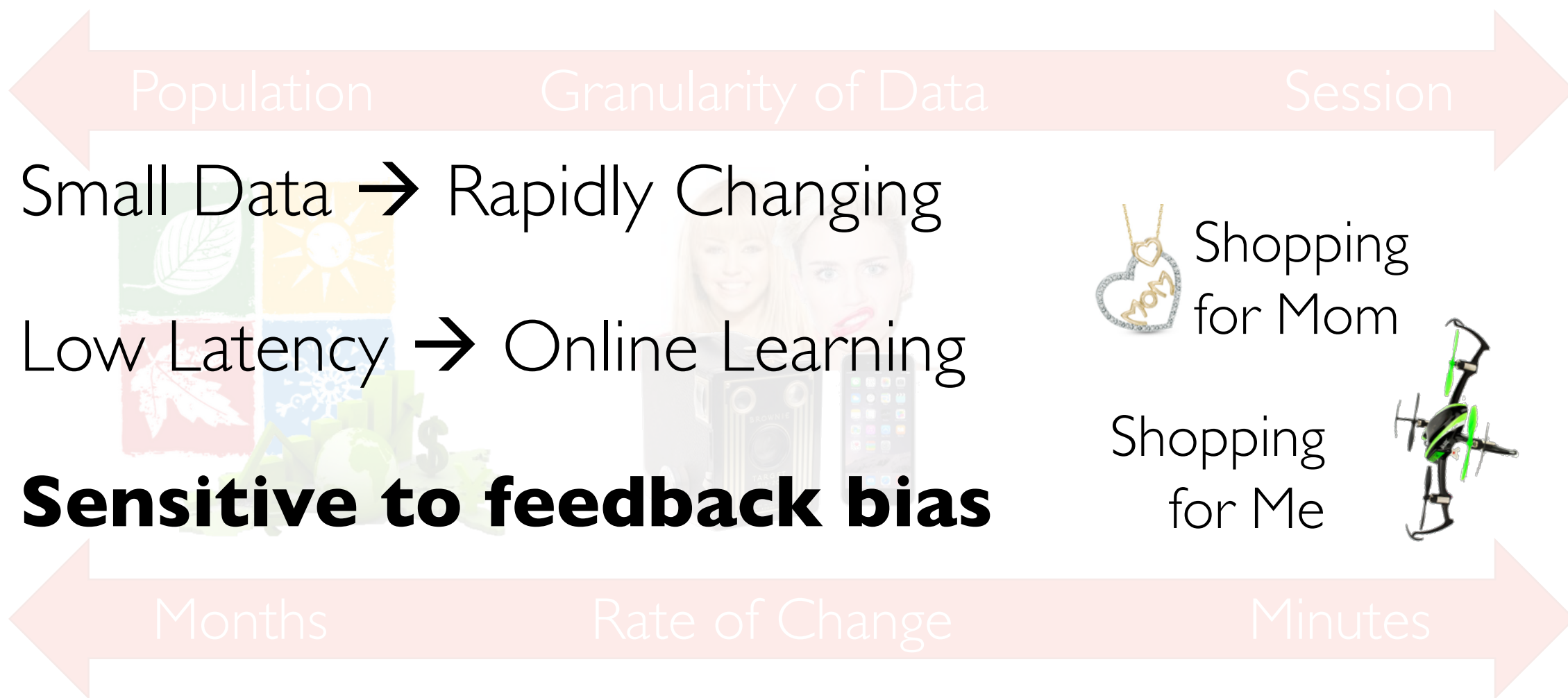
Rely on efficient offline retraining
→ High-throughput Systems

Months

Rate of Change

Minutes

Adaptive to Change at All Scales



The Feedback Loop

I once looked at cameras on Amazon ...

Opportunity for
Bandit Algorithms

Bandits present new challenges:

- computation overhead
- complicates caching + indexing

Similar cameras
and accessories



My Amazon Homepage

Exploration / Exploitation Tradeoff

*Systems that can take **actions** can
adversely bias future **data**.*

*Opportunity for **Bandits**!*

Bandits present new challenges:

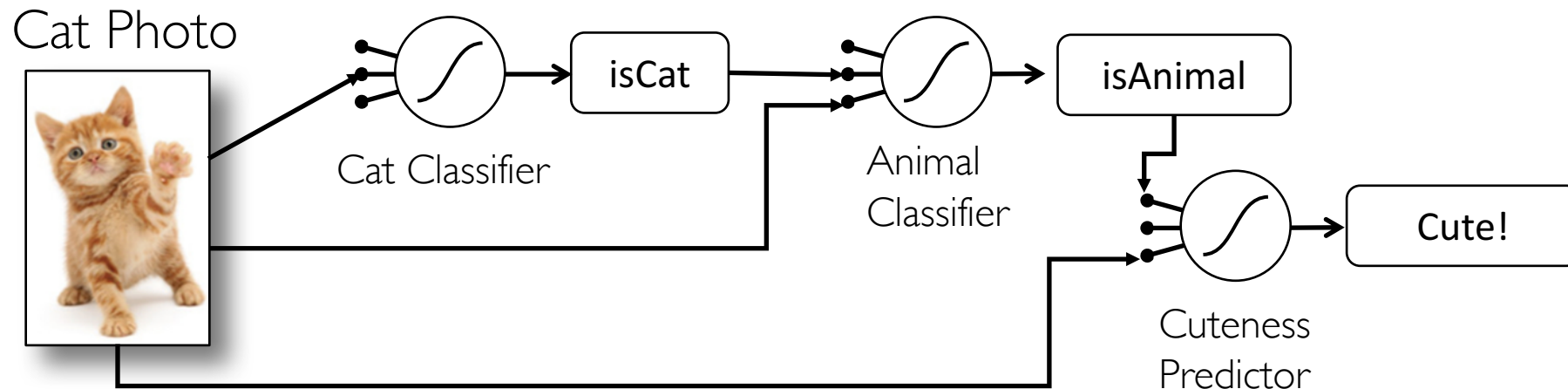
- Complicates caching + indexing
- tuning + counterfactual reasoning

Management: Collaborative Development

Teams of data-scientists working on similar tasks

➤ “*competing*” features and models

Complex model dependencies:





UC Berkeley AMPLab

*Daniel Crankshaw, Xin Wang, Joseph Gonzalez
Peter Bailis, Haoyuan, Zhao Zhang,
Michael J. Franklin, Ali Ghodsi,
and Michael I. Jordan*





UC Berkeley AMPLab

*Daniel Crankshaw, Xin Wang, Joseph Gonzalez
Peter Bailis, Haoyuan, Zhao Zhang,
Michael J. Franklin, Ali Ghodsi,
and Michael I. Jordan*

Active Research Project



Velox Model Serving System

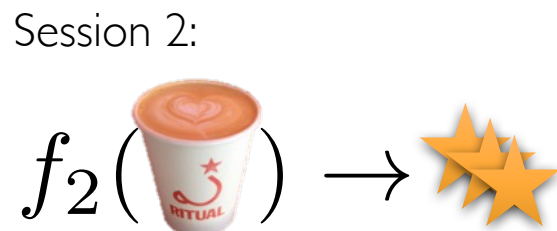
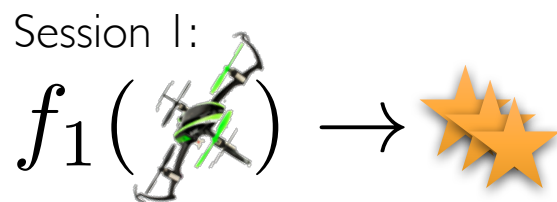
[CIDR'15, LearningSys'15]

Focuses on the multi-task learning (MTL) domain

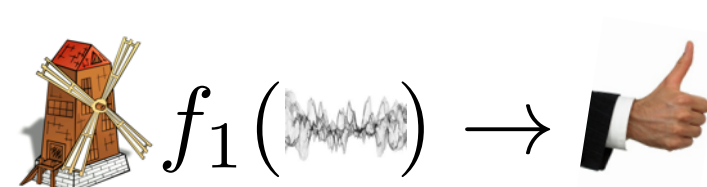
Spam
Classification



Content Rec.
Scoring



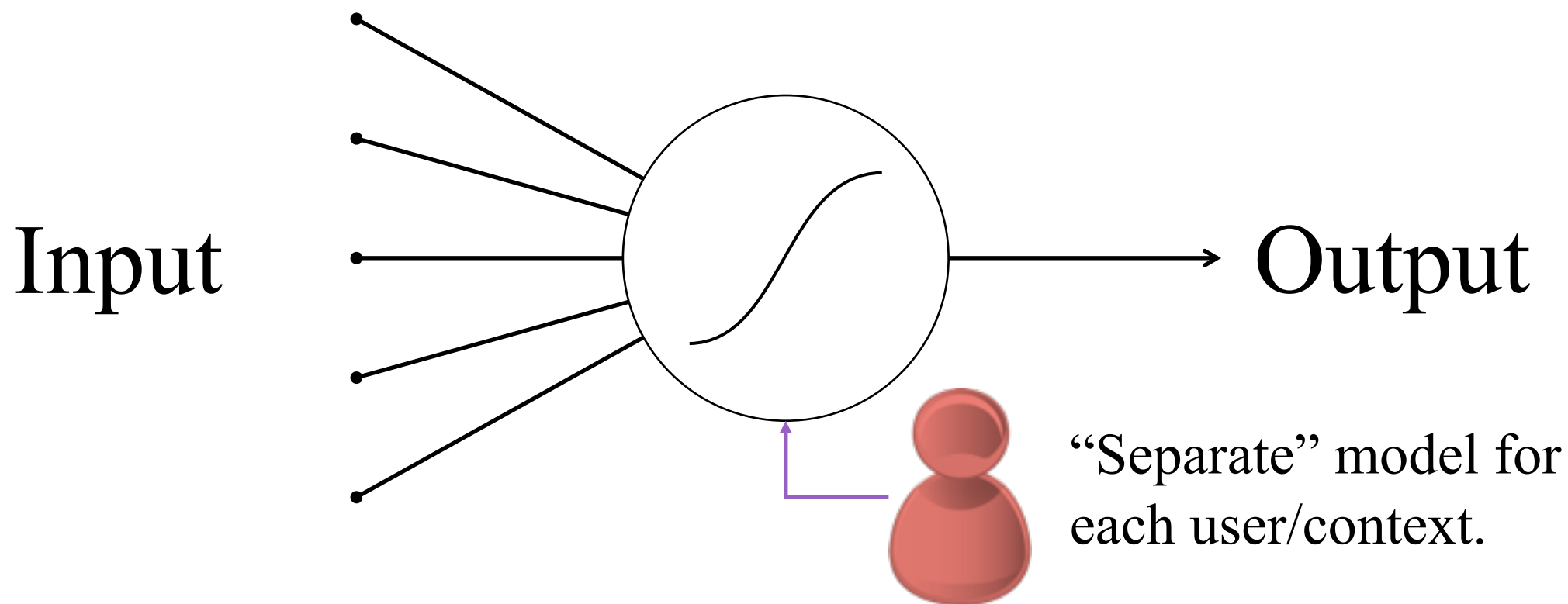
Localized
Anomaly Detection



Velox Model Serving System

[CIDR'15, [LearningSys'15](#)]

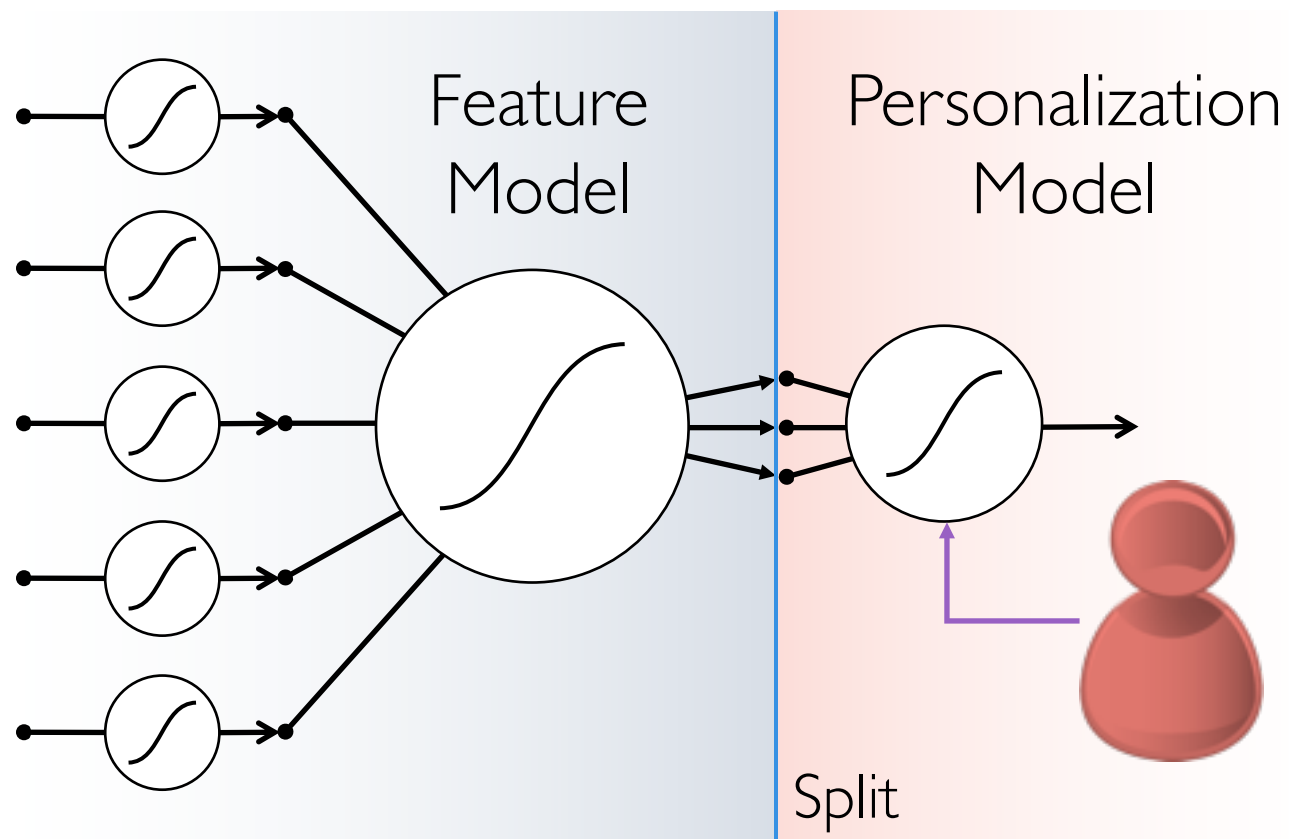
Personalized Models (Multitask Learning)



Velox Model Serving System

[CIDR'15, [LearningSys'15](#)]

Personalized Models (Multitask Learning)



Hybrid Offline + Online Learning

Update feature functions *offline* using batch solvers

- Leverage high-throughput systems (Apache Spark)
- Exploit slow change in population statistics

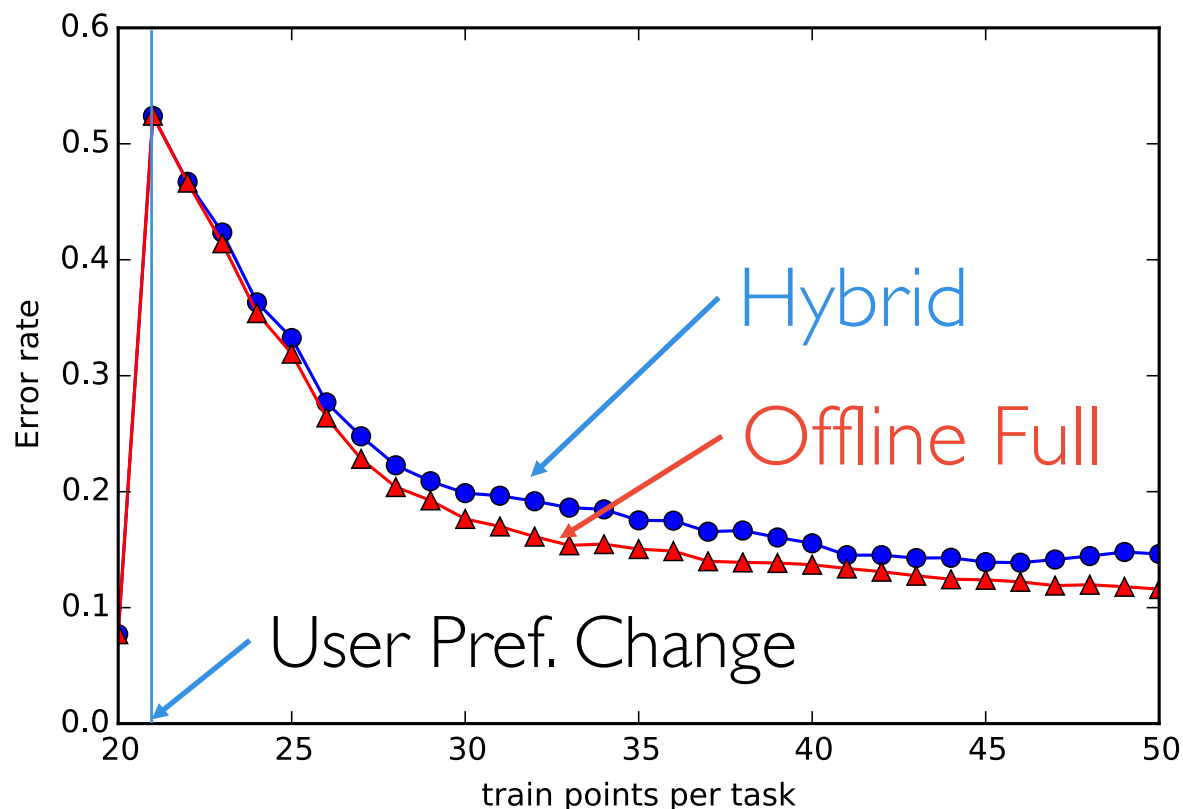
$$f(x; \theta)^T w_u$$

Update the user weights *online*:

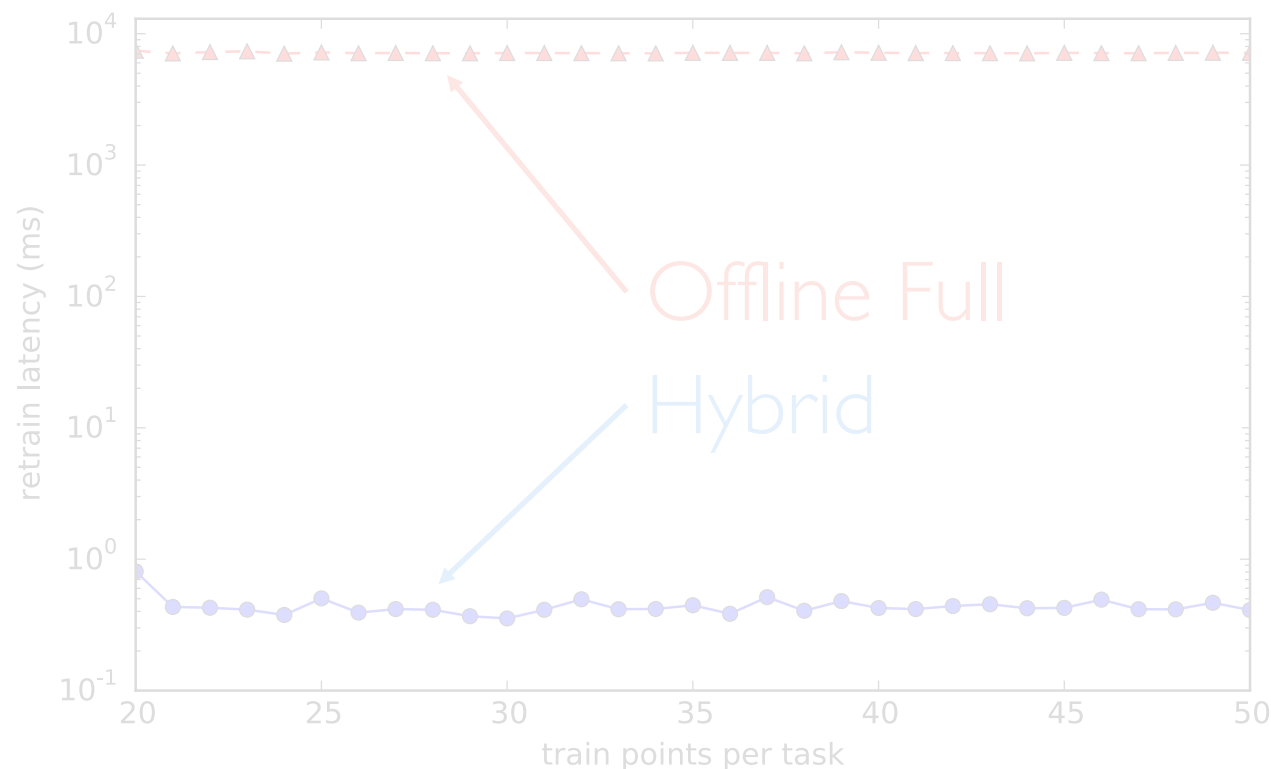
- Simple to train + more robust model
- Address rapidly changing user statistics

Hybrid Online + Offline Learning Results

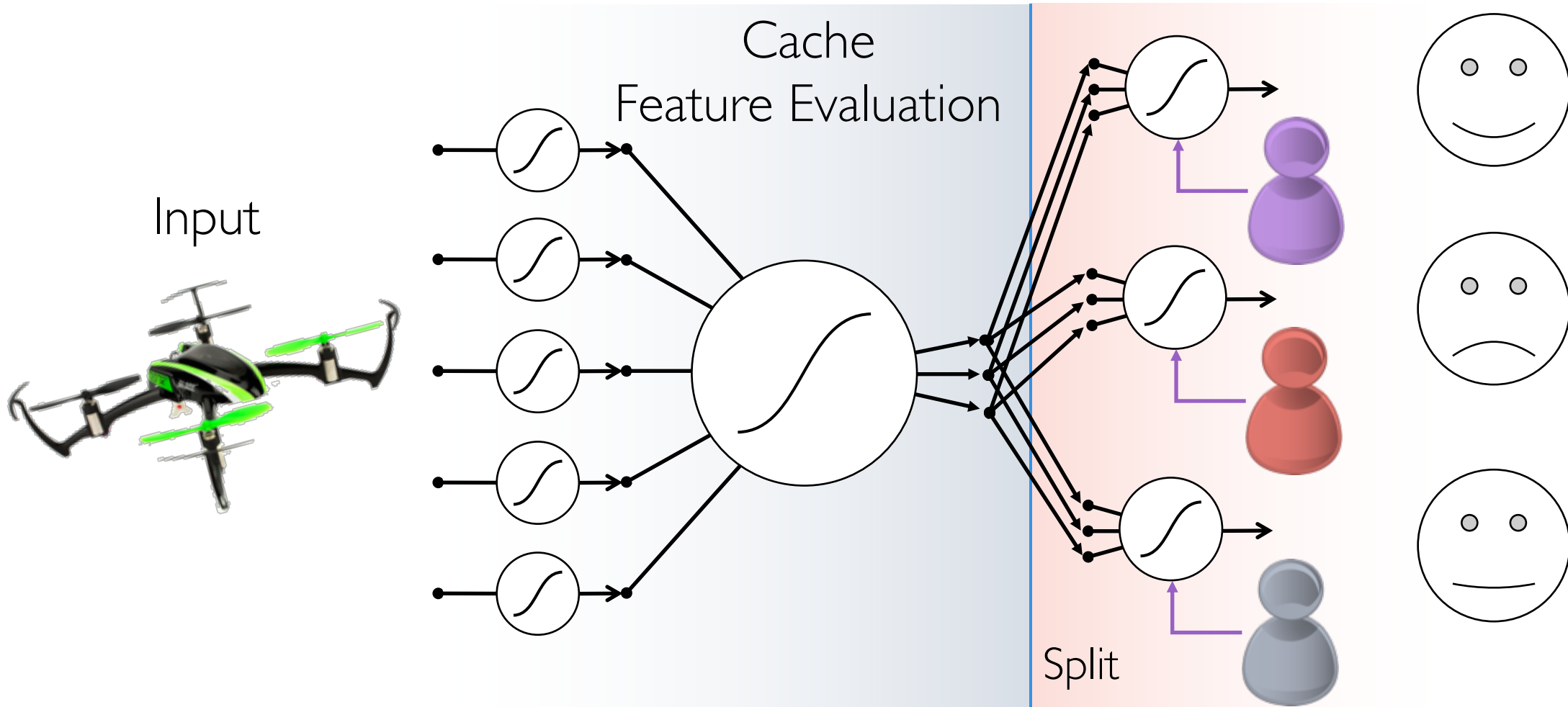
Similar Test Error



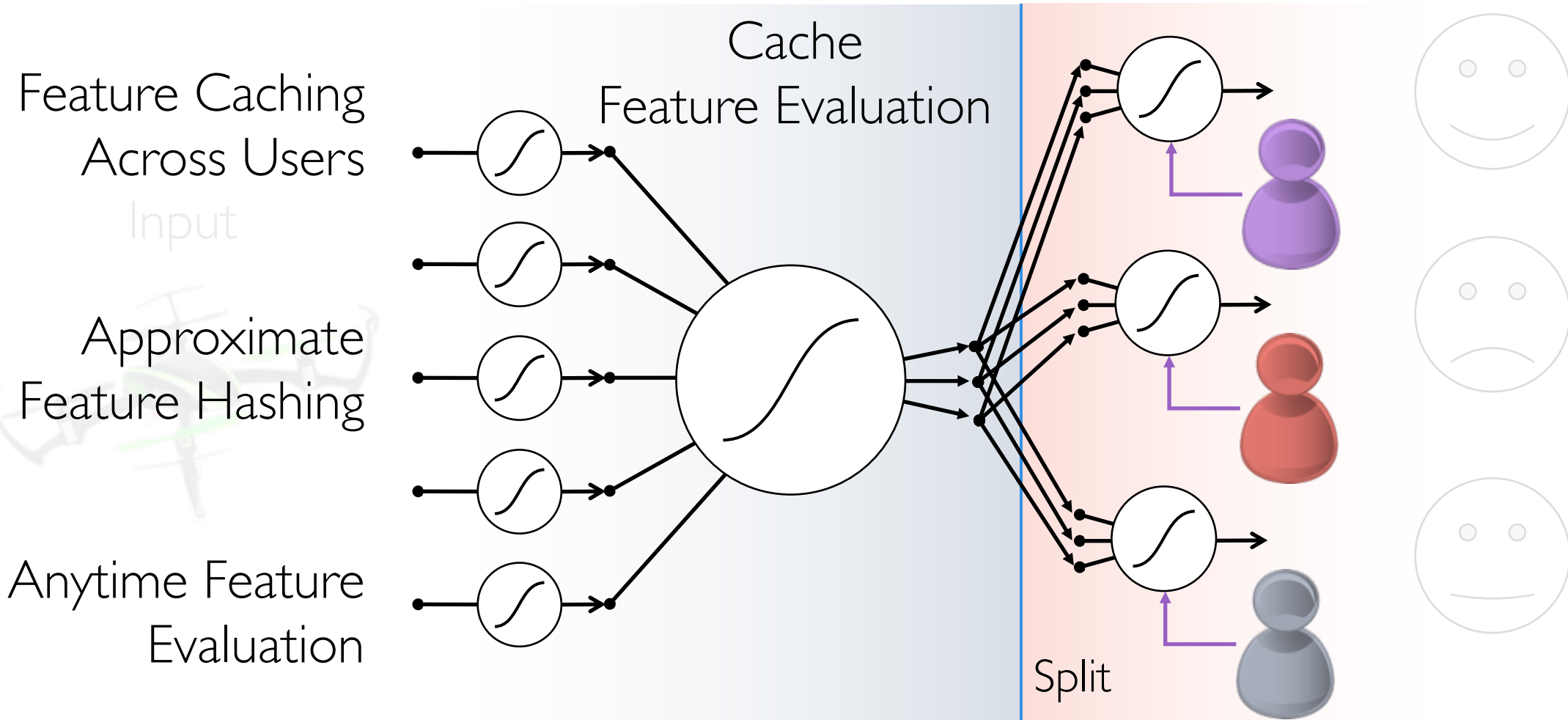
Substantially Faster Training



Evaluating the Model



Evaluating the Model



Feature Caching

New input: x

Compute feature: $f(x; \theta)$

Hash input: $h(x)$ — Store result in table

			$f(x; \theta)$	
--	--	--	----------------	--

Feature Hash Table

LSH Cache Coarsening

New input $z \neq x$

Hash new input: $h(z)$

False cache collision

Use Wrong Value!

→ LSH hash fn.

			$f(x; \theta)$	
--	--	--	----------------	--

Feature Hash Table

LSH Cache Coarsening

Locality-Sensitive Hashing:

$$x \approx z \quad \Rightarrow \quad h(x) = h(z)$$

Locality-Sensitive Caching:


$f(x; \theta) \approx f(z; \theta)$	\Rightarrow	$f(h(x)) = h(z)$
-------------------------------------	---------------	------------------

Feature Hash Table

Use Value Anyways!
→ Req. LSH

Anytime Predictions

Compute features asynchronously:

$$\text{---} w_{u1} + \text{---} w_{u2} + \text{---} w_{u3}$$


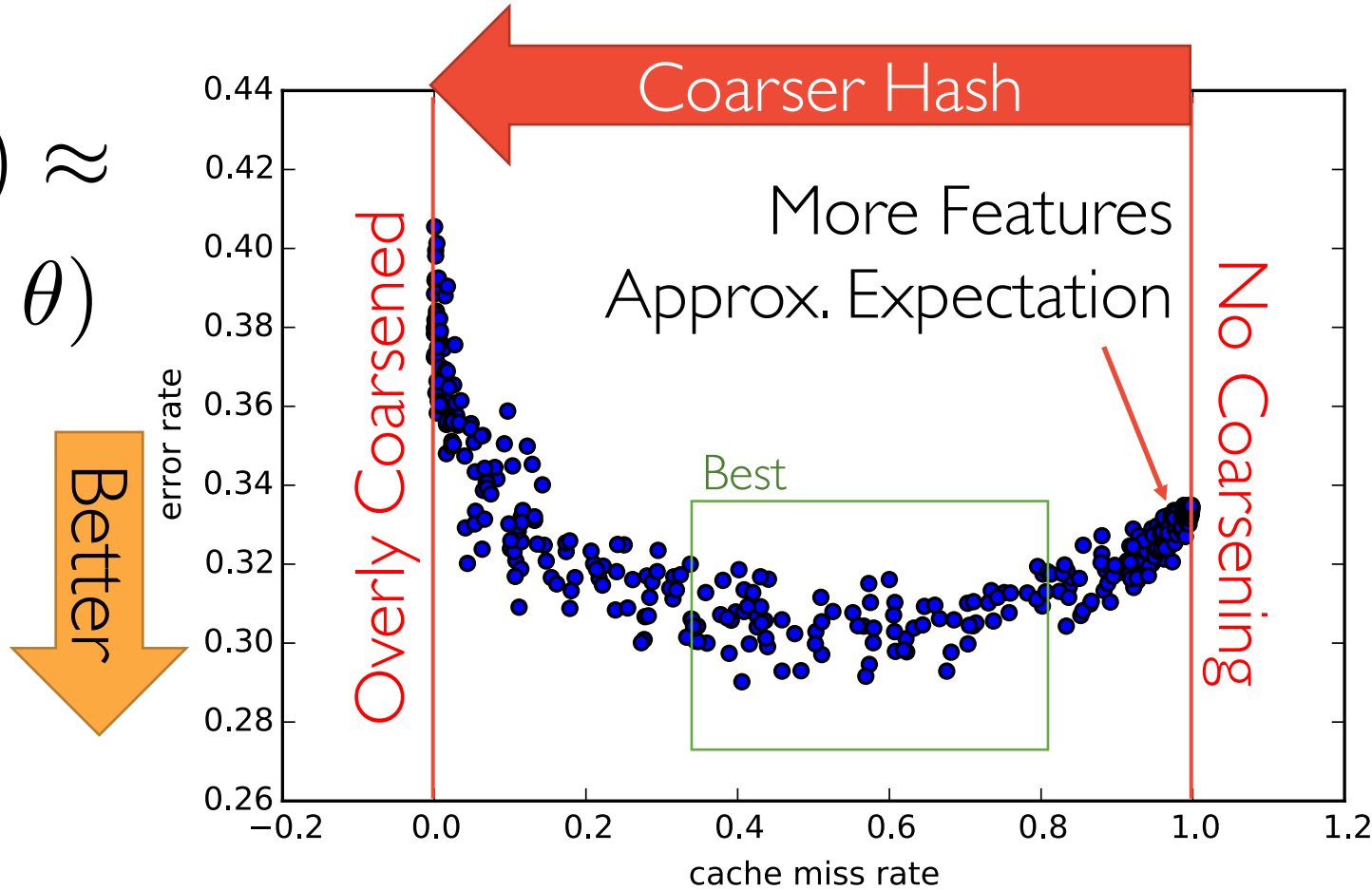
if a particular element does not arrive use estimator instead

Always able to render a prediction by the latency deadline

Coarsening + Anytime Predictions

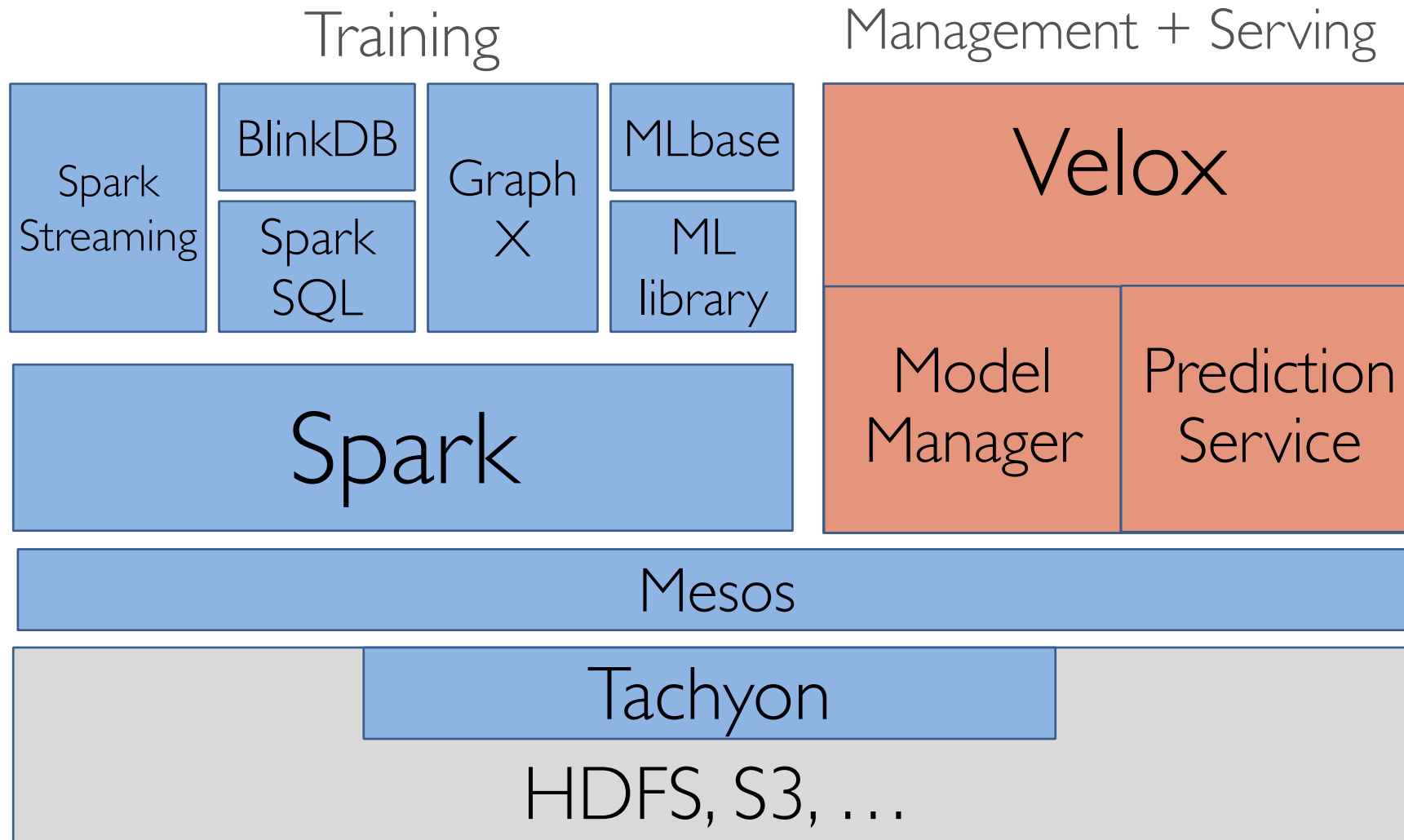
$$f_i(x; \theta) \approx f_i(z; \theta)$$

$$f_i(x; \theta) \approx \mathbb{E} [f_i(x; \theta)]$$



Checkout our poster!

Part of Berkeley Data Analytics Stack



Dato Predictive Services

Production ready model serving and management system

- Elastic scaling and load balancing of docker.io containers
- AWS Cloudwatch Metrics and Reporting
- Serves Dato Create models, scikit-learn, and custom python
- Distributed shared caching: *scale-out to address latency*
- REST management API: Demo?



UC Berkeley AMPLab

*Daniel Crankshaw, Xin Wang, Joseph Gonzalez
Peter Bailis, Haoyuan, Zhao Zhang,
Michael J. Franklin, Ali Ghodsi,
and Michael I. Jordan*



Predictive Services

Responsive

Adaptive

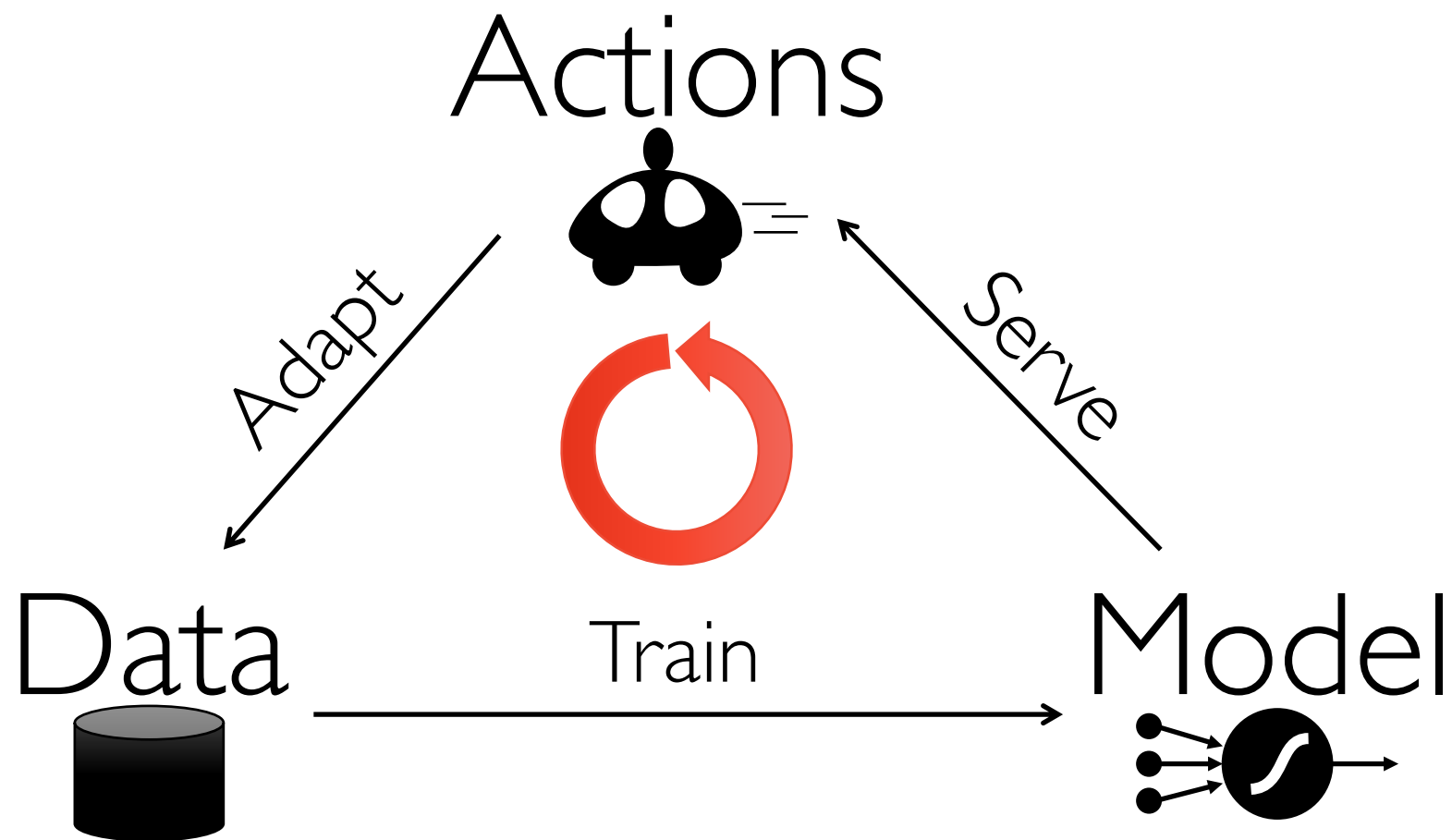
Manageable

Key Insights:

Caching, Bandits, &
Management

Online/Offline Learning
Latency vs. Accuracy

Future of Learning Systems



Thank You

Joseph E. Gonzalez

jegonzal@cs.berkeley.edu, Assistant Professor @ UC Berkeley

joseph@dato.com, Co-Founder @ Dato