



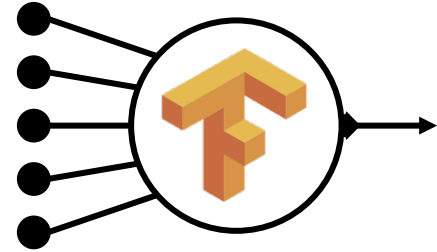
# RISE

## to the Challenges of AI Systems

Joseph E. Gonzalez

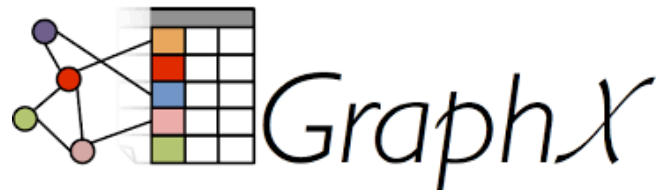
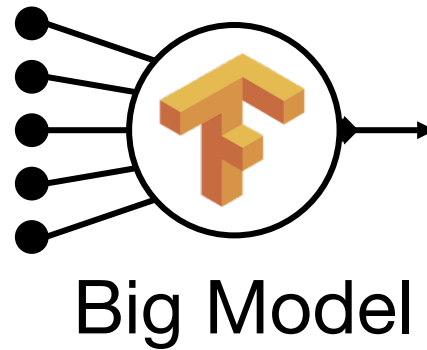
Assistant Professor, UC Berkeley

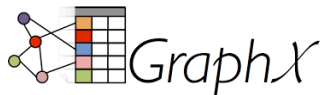
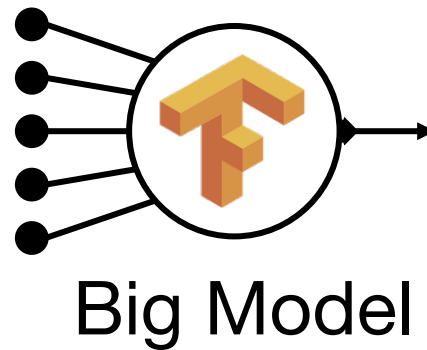
[jegonzal@cs.berkeley.edu](mailto:jegonzal@cs.berkeley.edu)



Big Model

Large-Scale parallel and distributed systems







# How to do Research in AI Systems

- Manage **Complexity**

- seek parsimony in system design
- great systems research is often about what features are taken away
- Do a few things well and be composable

- Identify **Tradeoffs**

- With each design decision what do you gain and lose?
- What trade-offs are fundamental?

- **Evaluate** your System

- **Positive:** How fast and scalable is it and *why*?
- **Negative:** When does it fail and what are its *limitations*?

# Hemingway\*

## Modeling Throughput and Convergence for ML Workloads



Shivaram  
Venkataraman

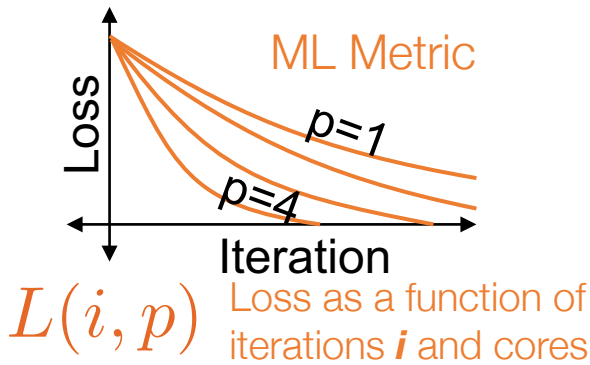
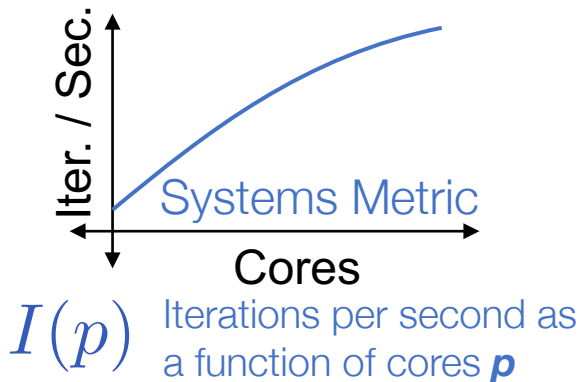


Xinghao  
Pan



Zi  
Zheng

- What is the best algorithm and level of parallelism for an ML task?
  - **Trade-off:** Parallelism, Coordination, & Convergence
- **Research challenge:** Can we model this trade-off explicitly?



We can estimate  $I$  from data on many systems

We can estimate  $L$  from data for our problem

# Hemingway\*

## Modeling Throughput and Convergence for ML Workloads



Shivaram  
Venkataraman



Xinghao  
Pan



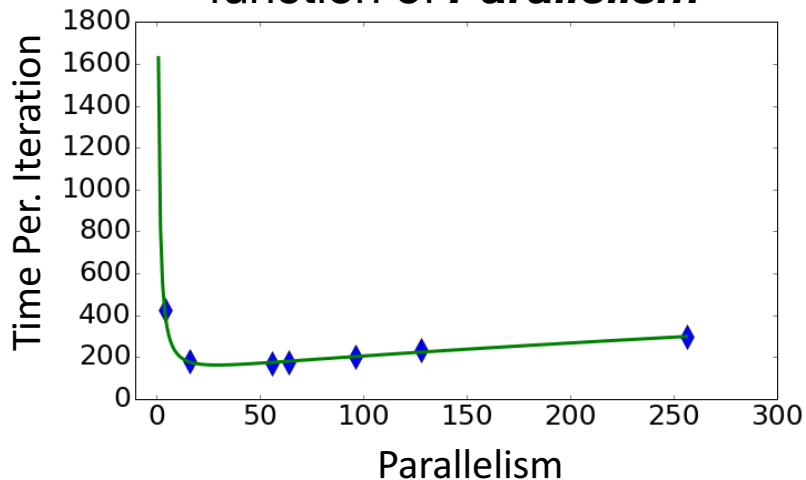
Zi  
Zheng

- What is the best algorithm and level of parallelism for an ML task?
  - **Trade-off:** Parallelism, Coordination, & Convergence
- **Research challenge:** Can we model this trade-off explicitly?

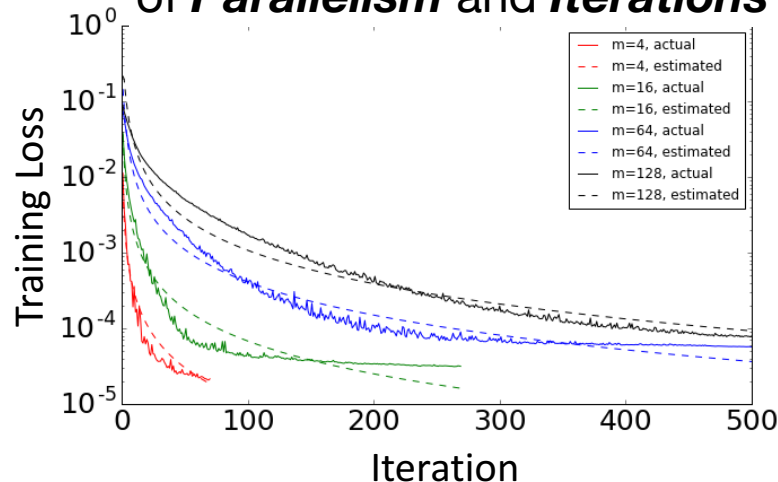
$$\left. \begin{array}{l} L(i, p) \text{ Loss as a function of} \\ \text{iterations } i \text{ and cores } p \\ I(p) \text{ Iterations per second as} \\ \text{a function of cores } p \end{array} \right\} \text{loss}(t, p) = L(t * I(p), p)$$

- How long does it take to get to a given loss?
- Given a time budget and number of cores which algorithm will give the best result?

## System Performance as a function of *Parallelism*

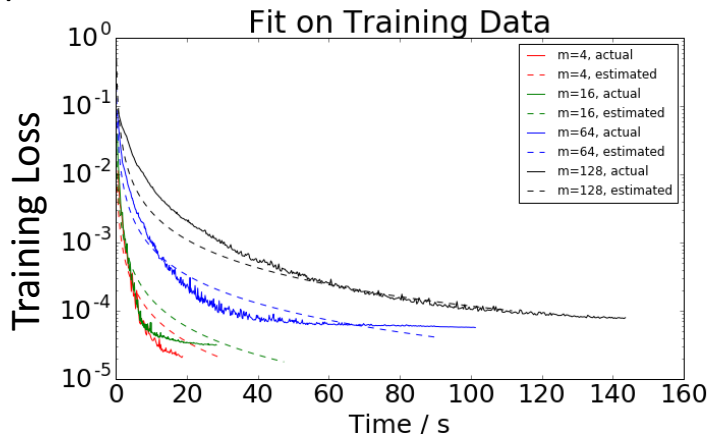


## Convergence as a function of *Parallelism* and *Iterations*



### **Hemingway: Modeling Distributed Optimization Algorithms.**

Xinghao Pan, Shivaram Venkataraman,  
Zizheng Tai, Joseph Gonzalez.  
NIPS'16 ML-Sys Workshop.



## Convergence as a fn. of *Time* and *Parallelism*

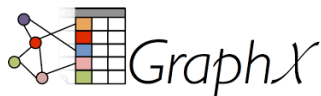
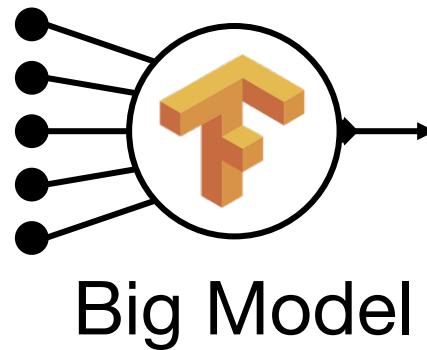
# Take away ...

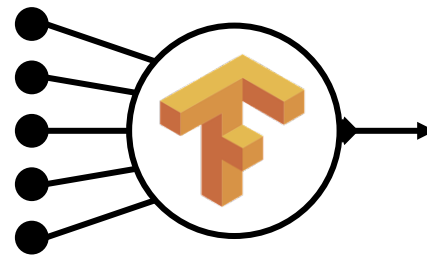
try to decouple

**System  
Improvements**

**Algorithm  
Improvements**

use data collection + sparse modeling  
to understand your system

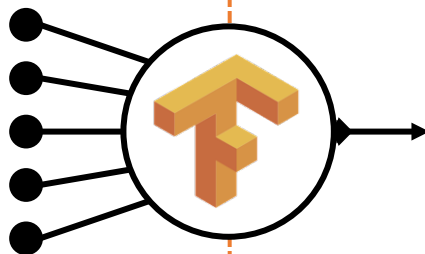




Big Model

— amplab 

# Learning

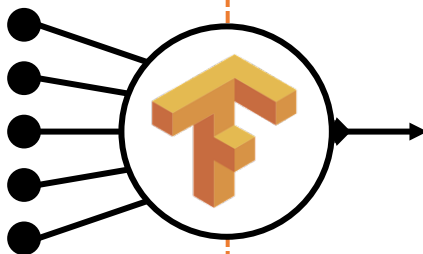


Big Model





# Learning

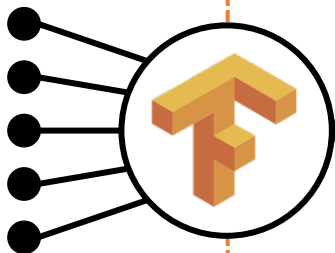


Big Model

# Conference Papers



# Learning



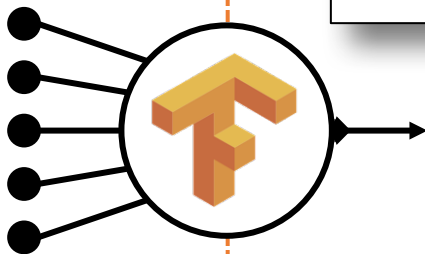
## Big Model

## Conference Papers



## Dashboards and Reports

# Learning



## Big Model

## Conference Papers



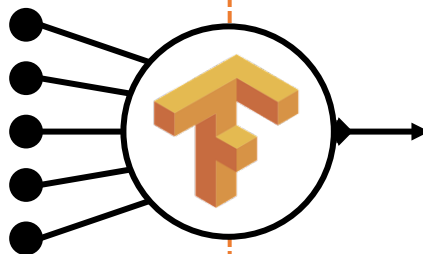
## Dashboards and Reports



## Drive Actions



# Learning

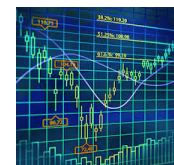


Big Model

# Drive Actions

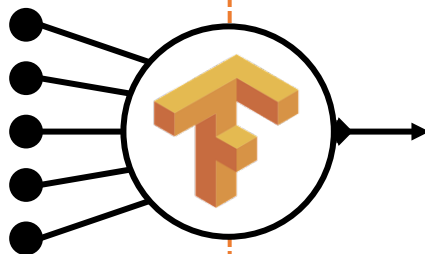


Hi, I'm Cortana.



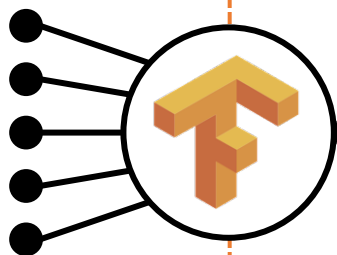
# Learning

# Inference



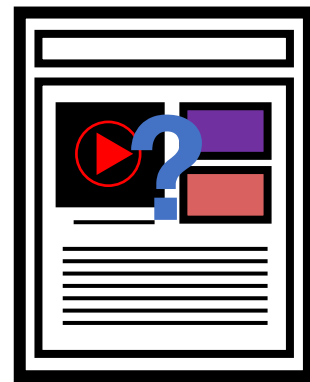
Big Model

# Learning



Big Model

# Inference

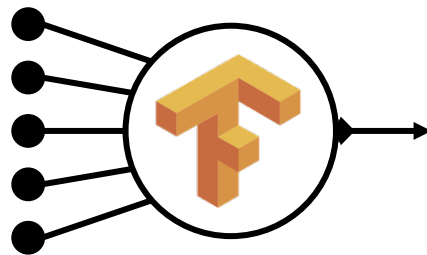


Application

## Learning



Training

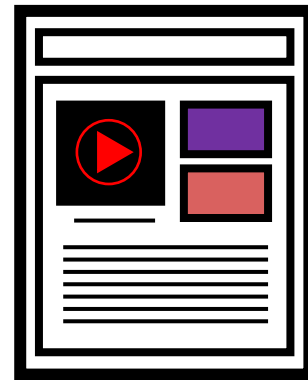


Big Model

## Inference

Query

Decision



Application

Often **overlooked**

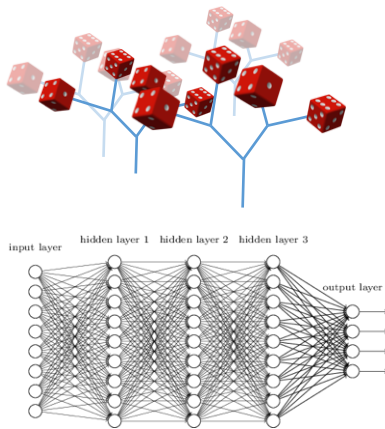
Timescale: ~10 milliseconds

**Billions of Queries a Day → Costly**

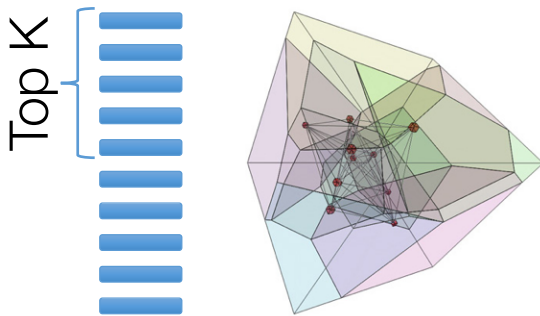
# why is **Inference** challenging?

Need to render **low latency** ( $< 10\text{ms}$ ) predictions for **complex**

## Models



## Queries



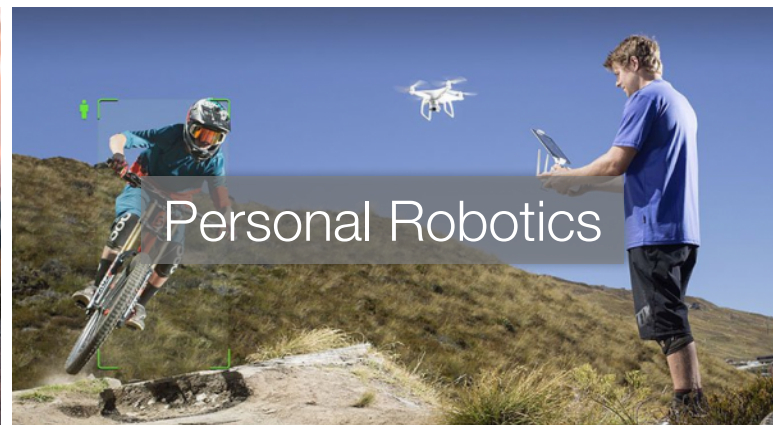
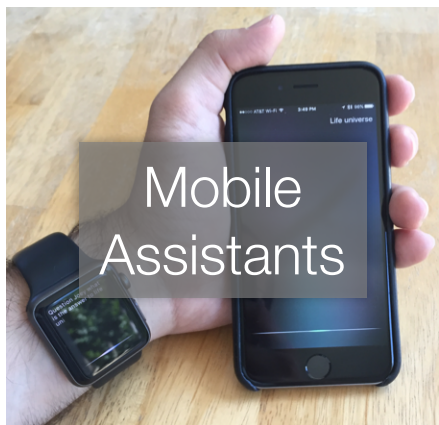
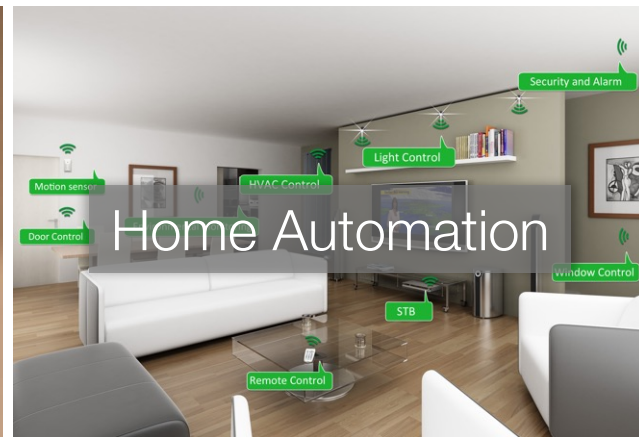
## Features

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

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



# Inference is moving beyond the cloud



# Inference is moving beyond the cloud



## Opportunities

- Reduce latency and improve privacy
- Address network partitions

## Research Challenges

- Minimize **power consumption**
- **Limited hardware** & long life-cycles
- Develop new **hybrid models** to leverage the cloud and edge devices



# Robust Inference is critical

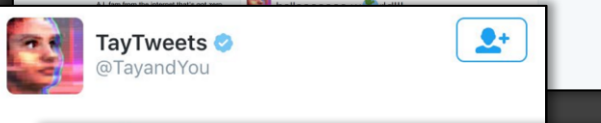
Self “*Parking*” Cars



Self “*Driving*” Cars

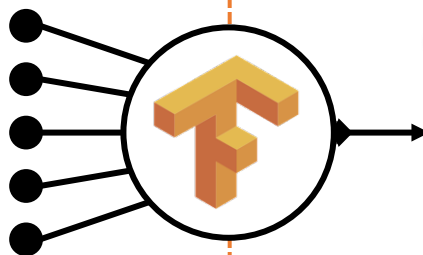


Chat AIs

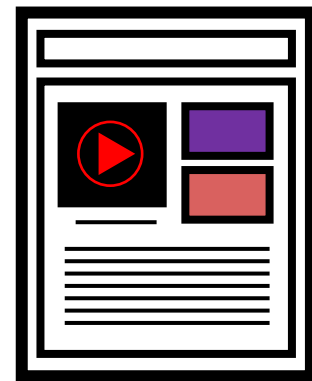


# Learning

# Inference



Big Model

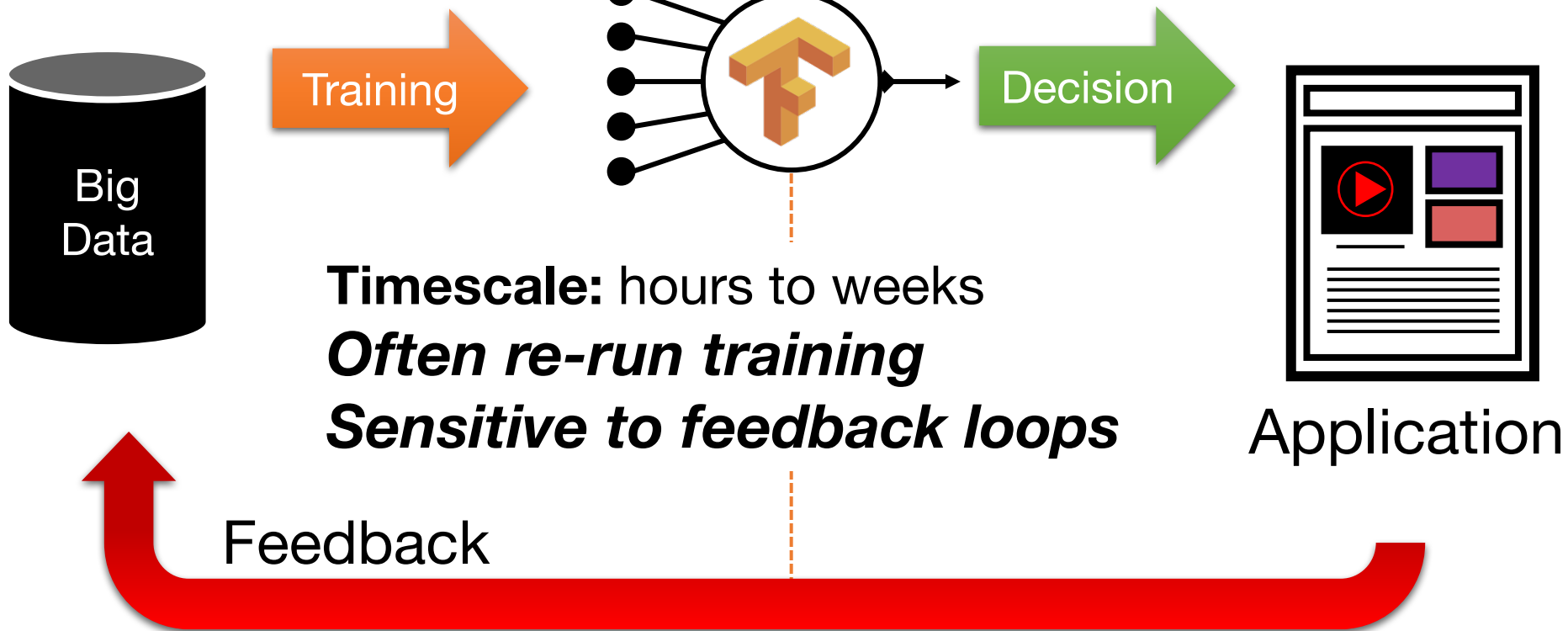


Application



# Learning

# Inference



# Why is **Closing the Loop** challenging?



**Implicit and Delayed**  
Feedback



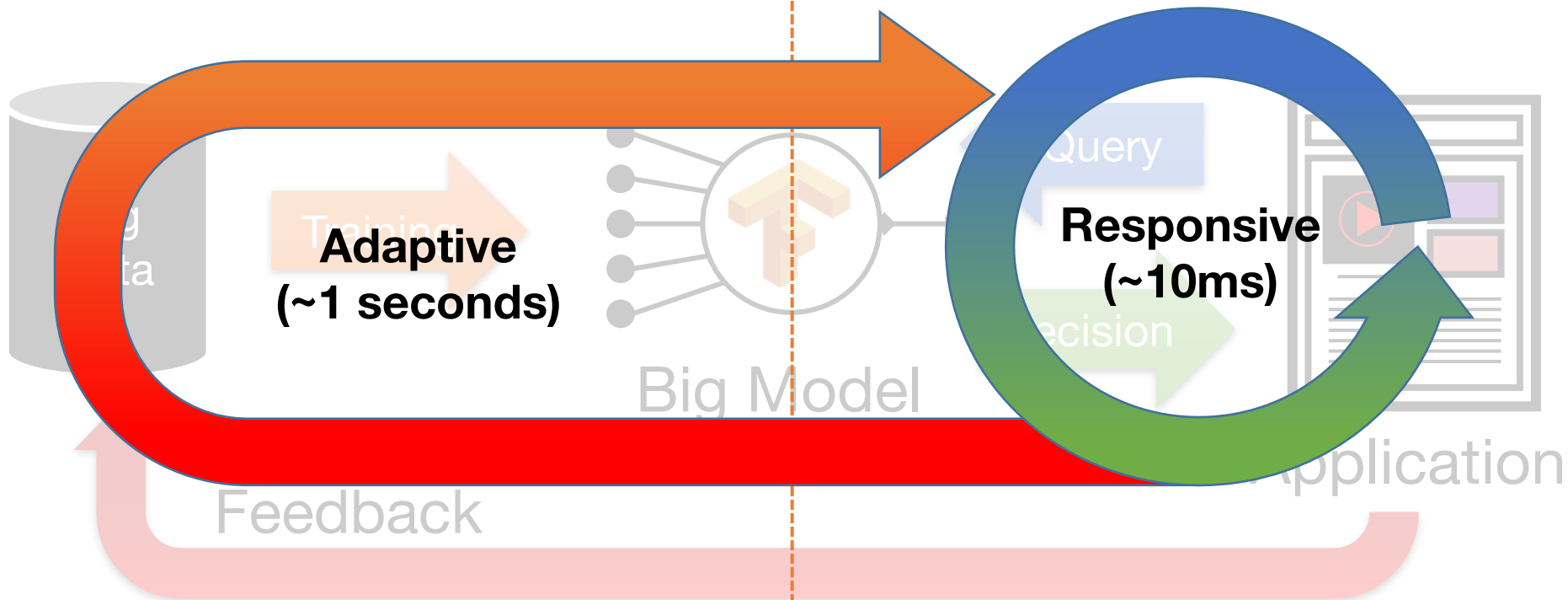
Self Reinforcing  
**Feedback Loops**

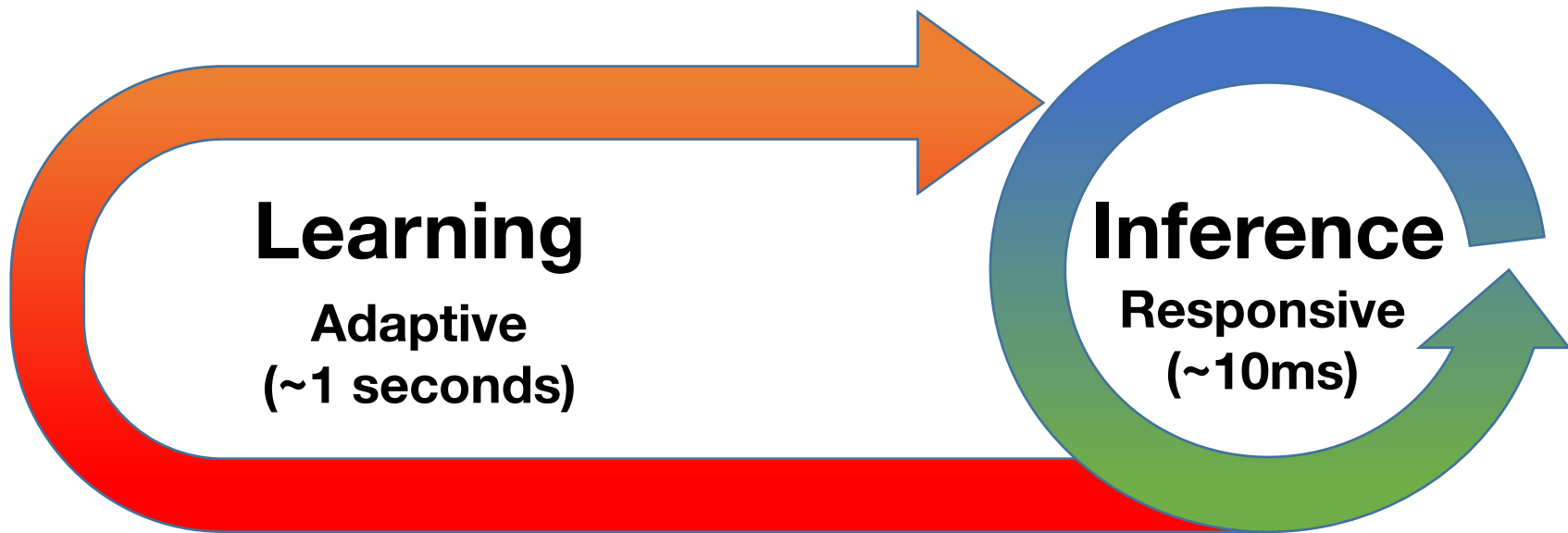


World Changes  
**at varying rates**

# Learning

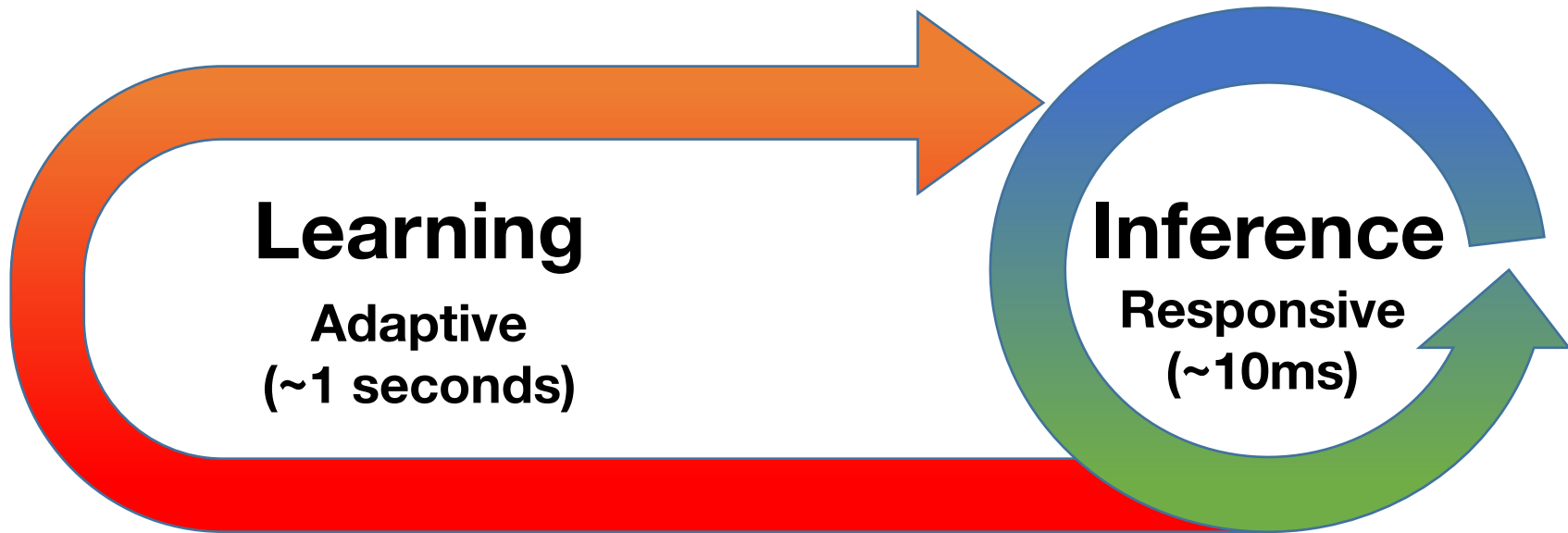
# Inference





?

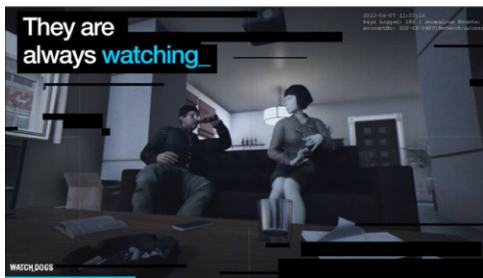




**Secure**

# Intelligence in **Sensitive Contexts**

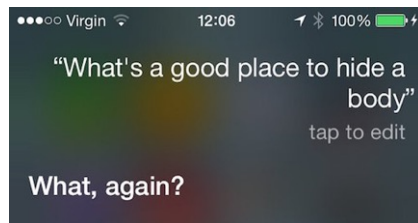
Augmented Reality



Home Monitoring



Voice Technologies



Medical Imaging



Protect the **data**, the **model**, and the **query**

# Protect the **data**, the **model**, and the **query**

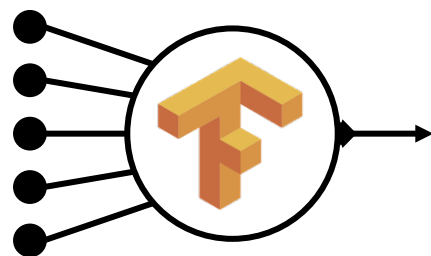
High-Value **Data is Sensitive**



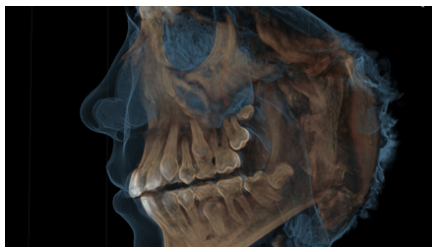
- Medical Info.
- Home video
- Finance

**Models** capture **value** in data

- Core Asset
- Sensitive



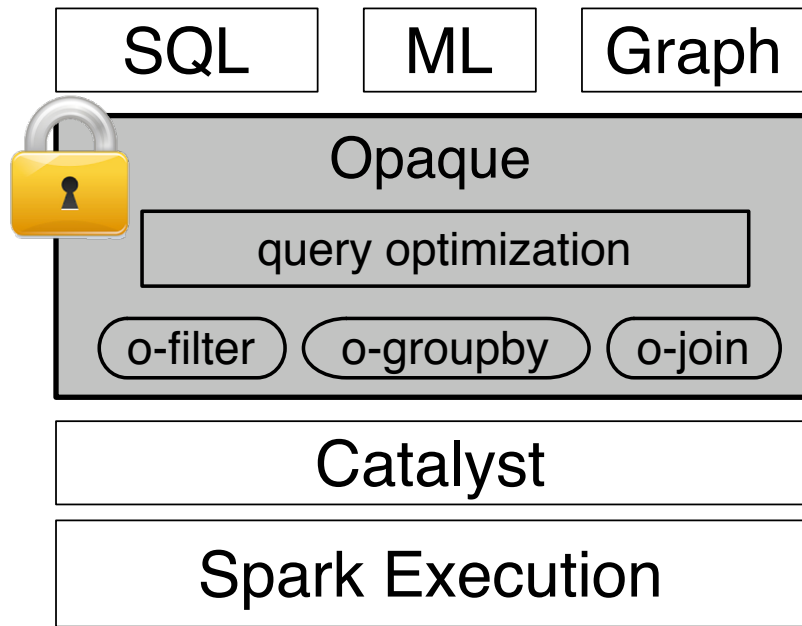
**Queries** can be as sensitive as the data

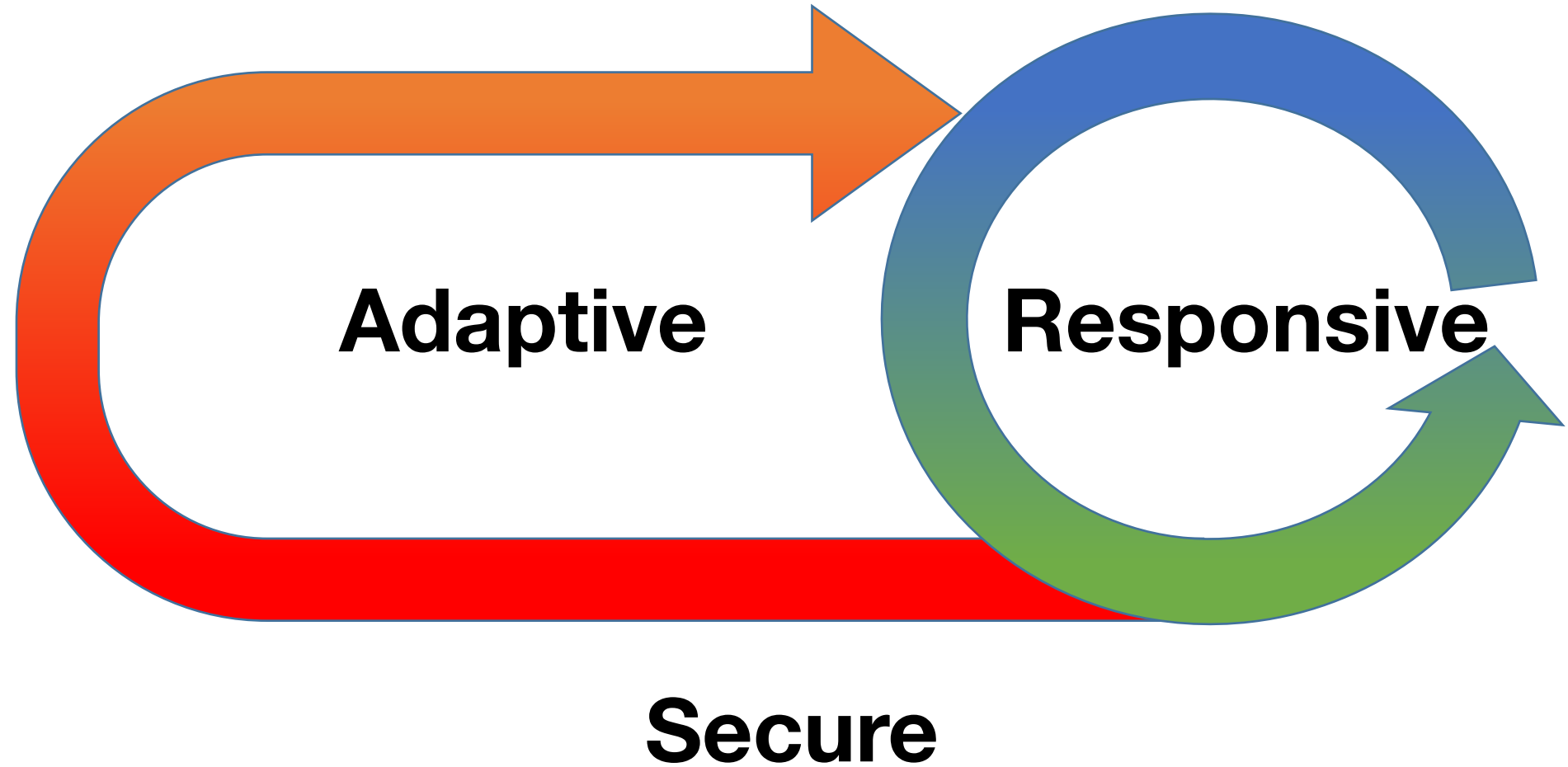


# Opaque: Analytics on Secure Enclaves

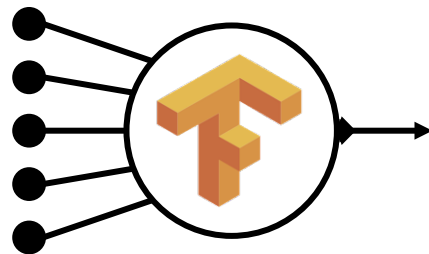
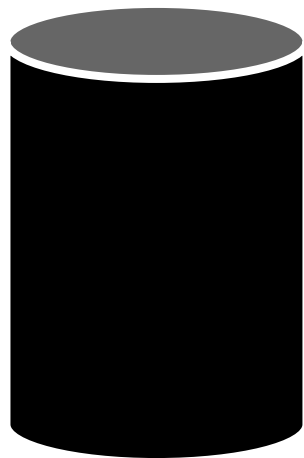
*Exploit hardware support to enable computing on encrypted data*

- **Today:** prototype system running in Apache Spark
  - support SQL queries in untrusted cloud
  - ~50% reduction in perf.
- **Future:** enable prediction serving on enc. queries



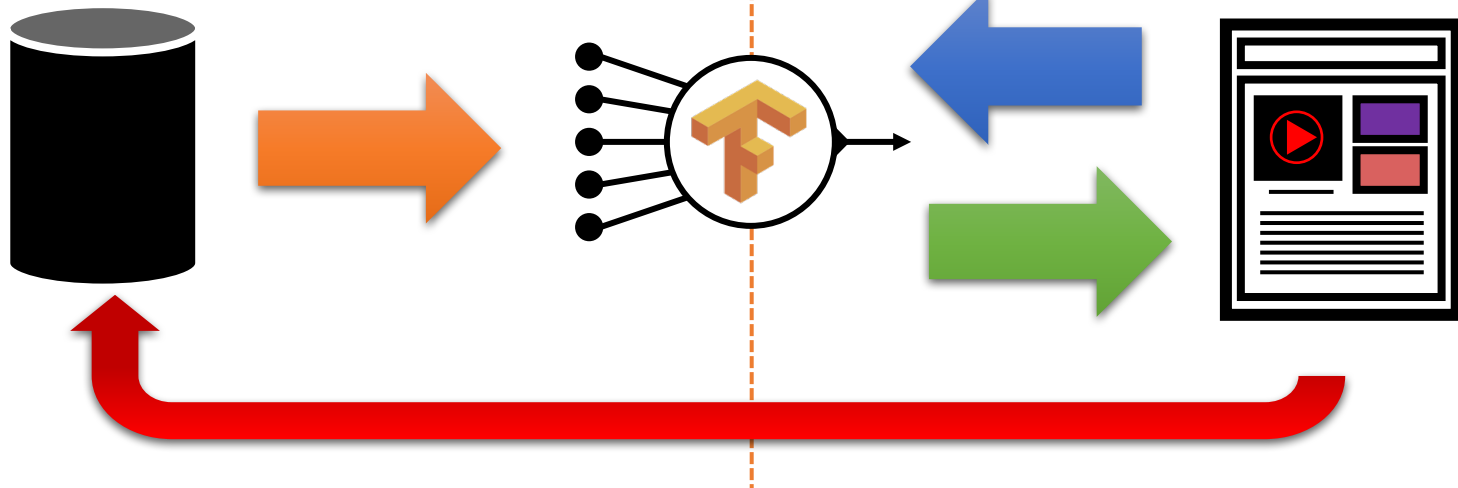


- amplab 



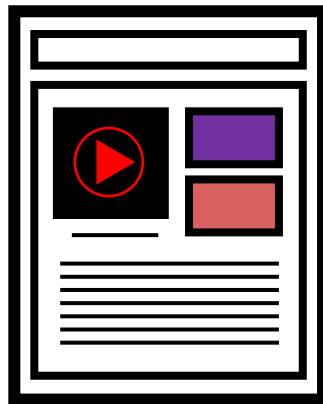
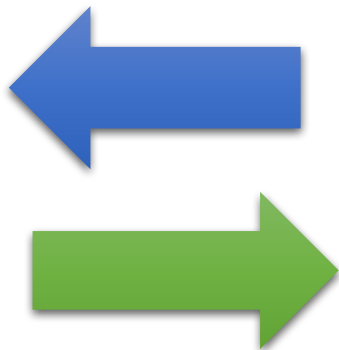
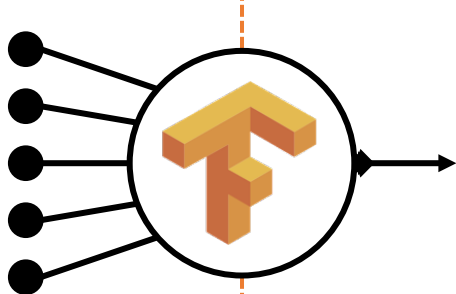


-amplab



# Clipper

A Low-Latency Online Prediction  
Serving System



**NSDI'17**

**Daniel Crankshaw**

Xin Wang

Giulio Zhou

Michael J. Franklin

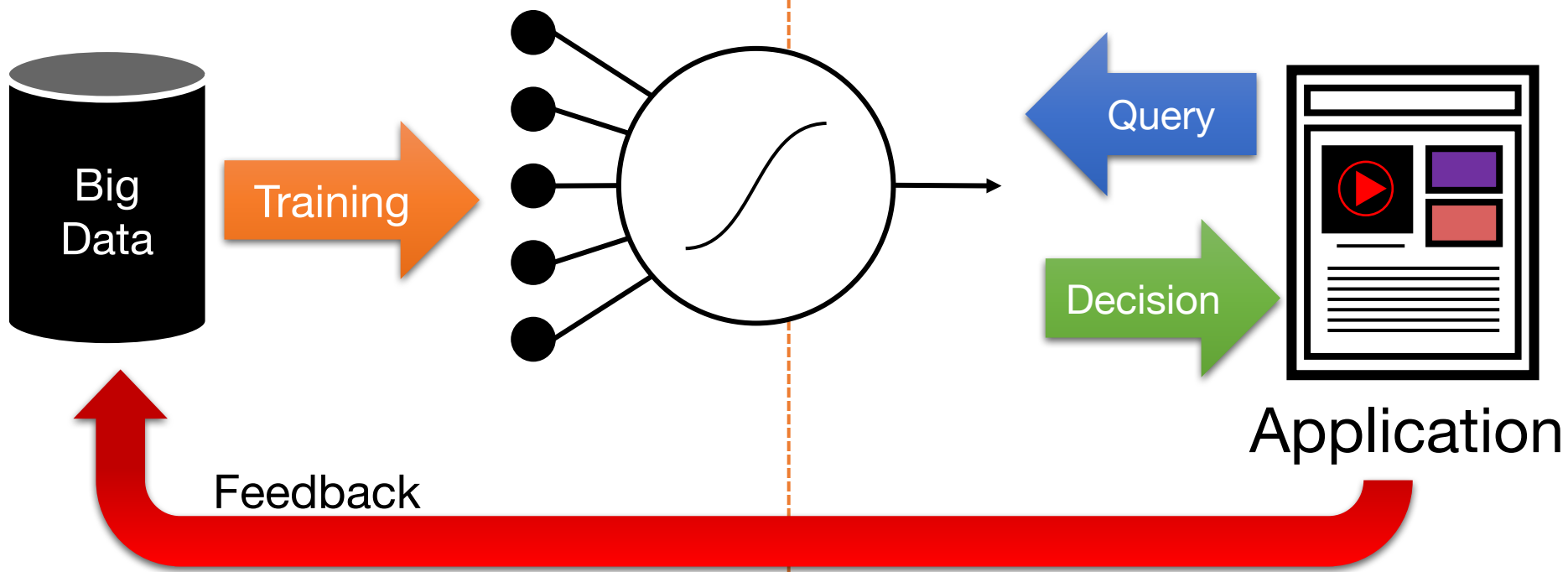
Joseph E. Gonzalez

Ion Stoica



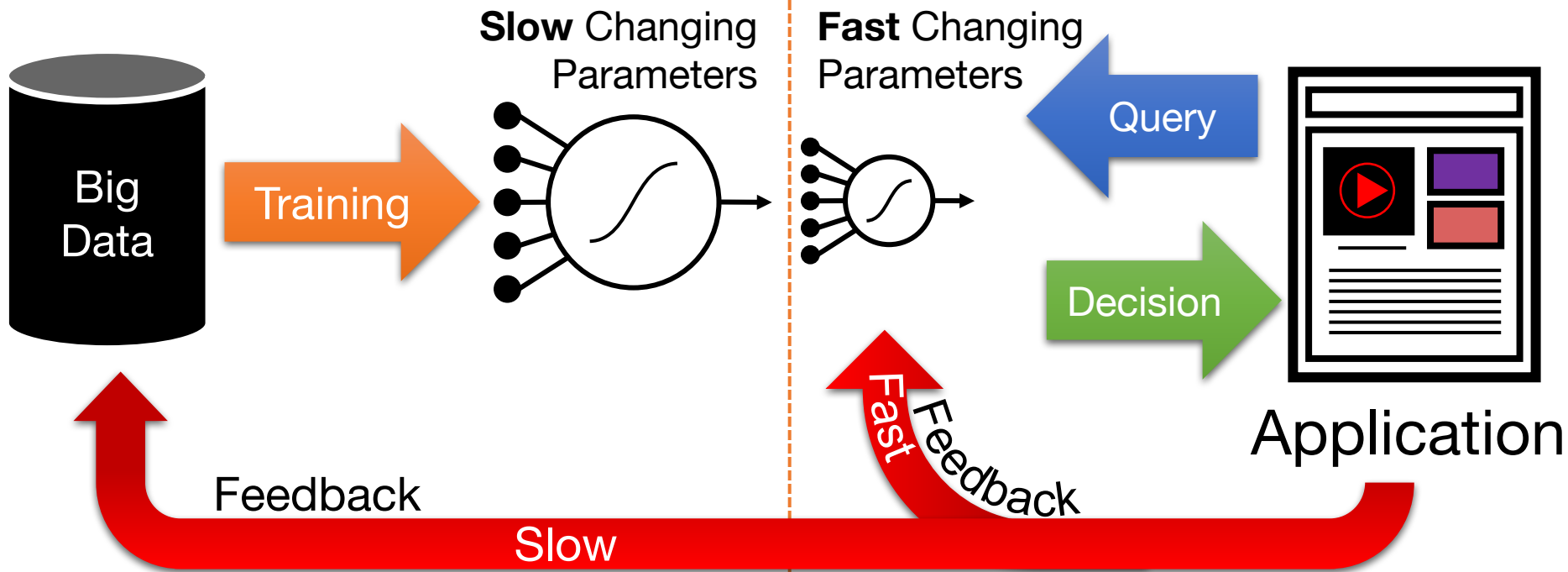
# Learning

# Inference



# Learning

# Inference



# Hybrid Offline + Online Learning

Update “feature” functions **offline** using batch solvers

- Leverage high-throughput systems (Tensor Flow)
- Exploit slow change in population statistics

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

Update the user weights **online**:

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

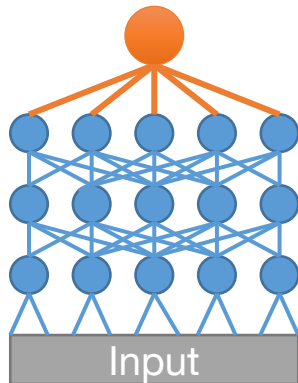
# Common modeling structure

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

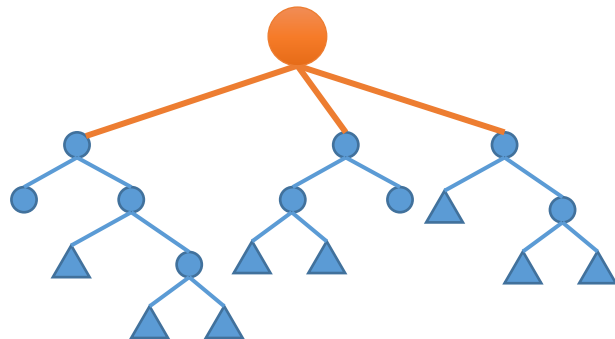
Matrix  
Factorization



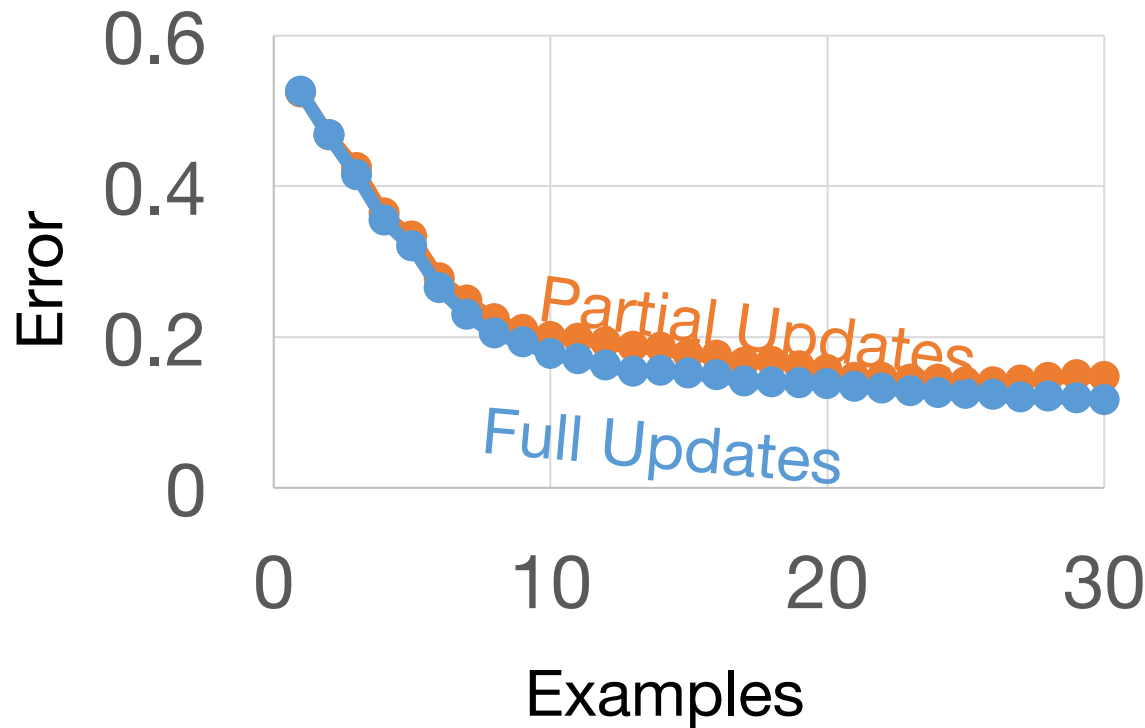
Deep  
Learning



Ensemble  
Methods



# Clipper Online Learning for Recommendations (Simulated News Rec.)

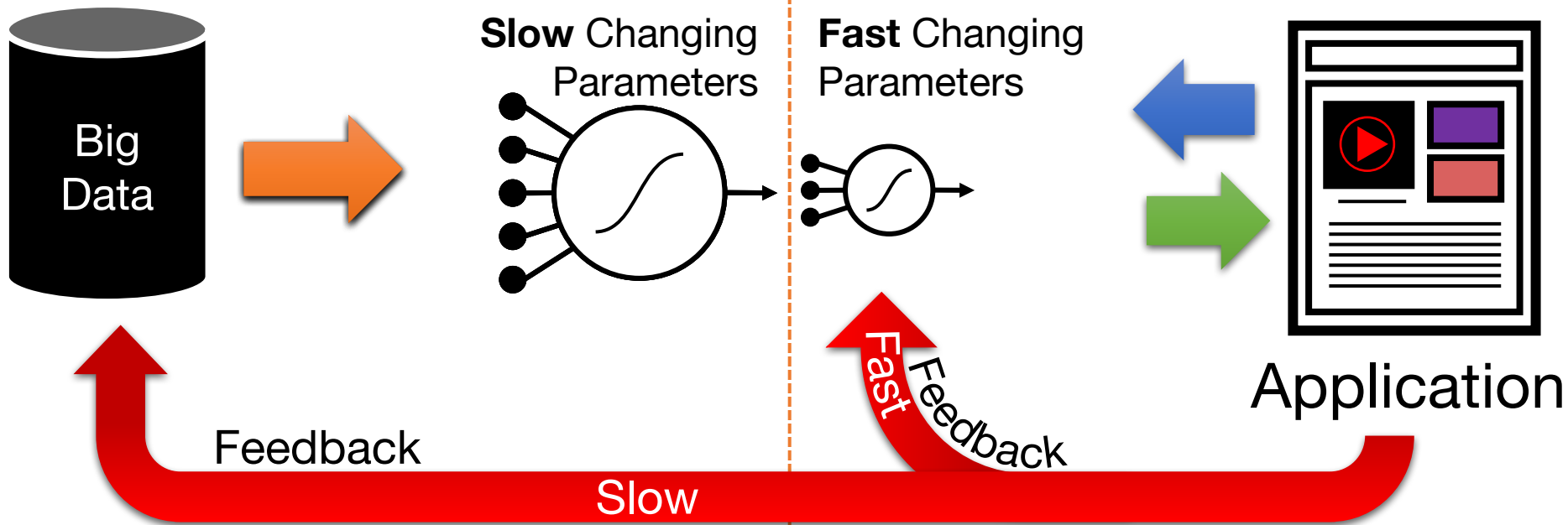


**Partial Updates:** 0.4 ms  
**Retraining:** 7.1 seconds

>4 orders-of-magnitude **faster adaptation**

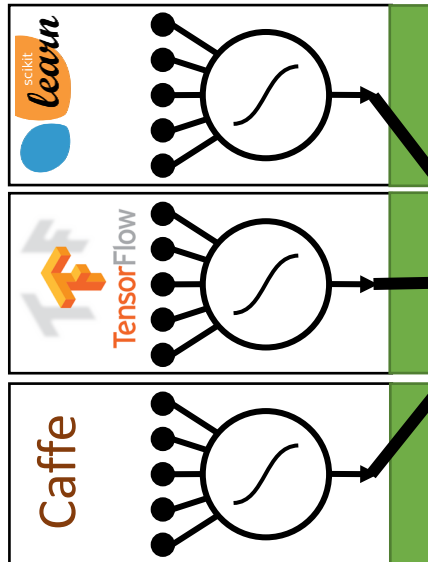
# Learning

# Inference



# Learning

Slow Changing  
Parameters



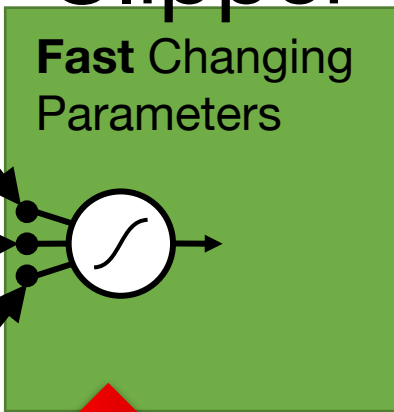
Feedback

Slow

# Inference

Clipper

Fast Changing  
Parameters



Application

Fast  
Feedback

# Clipper Serves Predictions across ML Frameworks

Fraud  
Detection



Content  
Rec.



Personal  
Asst.



Robotic  
Control



Machine  
Translation



## Clipper

theano

Dato



Create

KeystoneML

Caffe



TensorFlow



dmlc  
mxnet



KALDI





# Clipper

## Key Insight:

*The challenges of prediction serving can be addressed between end-user applications and machine learning frameworks*

As a result, Clipper is able to:

- **hide complexity** by
  - providing a common interface to applications
- **bound latency** and **maximize throughput**
  - through *caching, adaptive batching, model replication*
- enable *robust* **online learning** and **personalization**
  - through *model selection and ensemble algorithms*

***without modifying*** machine learning **frameworks** or front-end **applications**

# Clipper Architecture

Fraud  
Detection



Content  
Rec.



Personal  
Asst.



Robotic  
Control



Machine  
Translation



## Clipper

theano

Dato



Create

KeystoneML

Caffe



TensorFlow

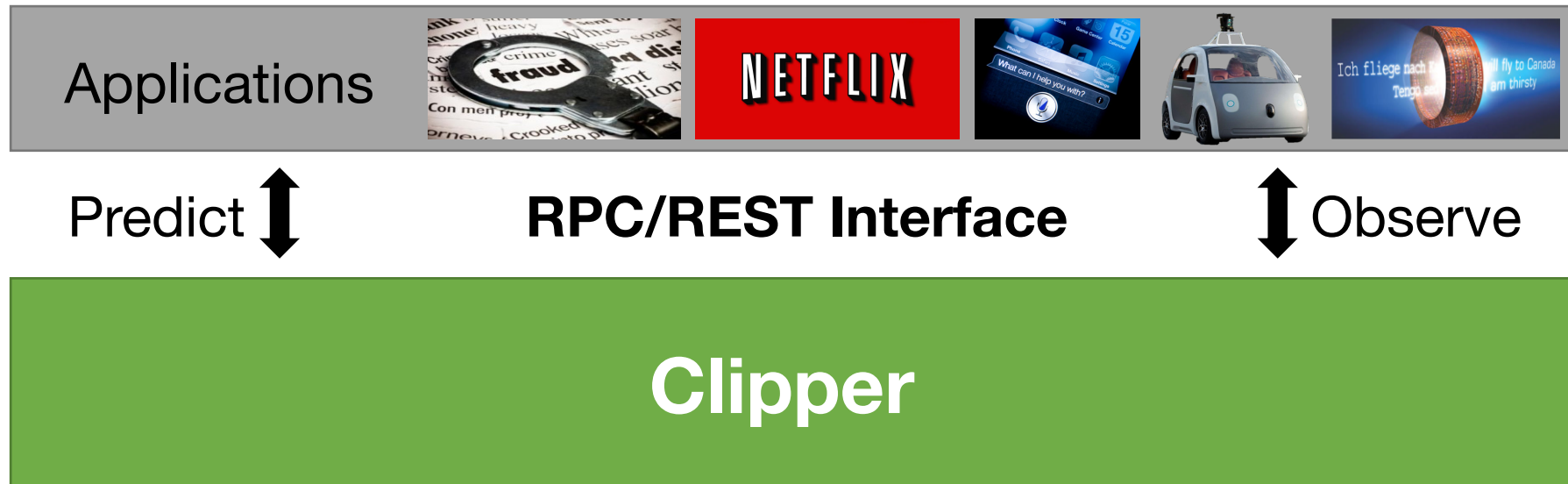


dmlc  
mxnet

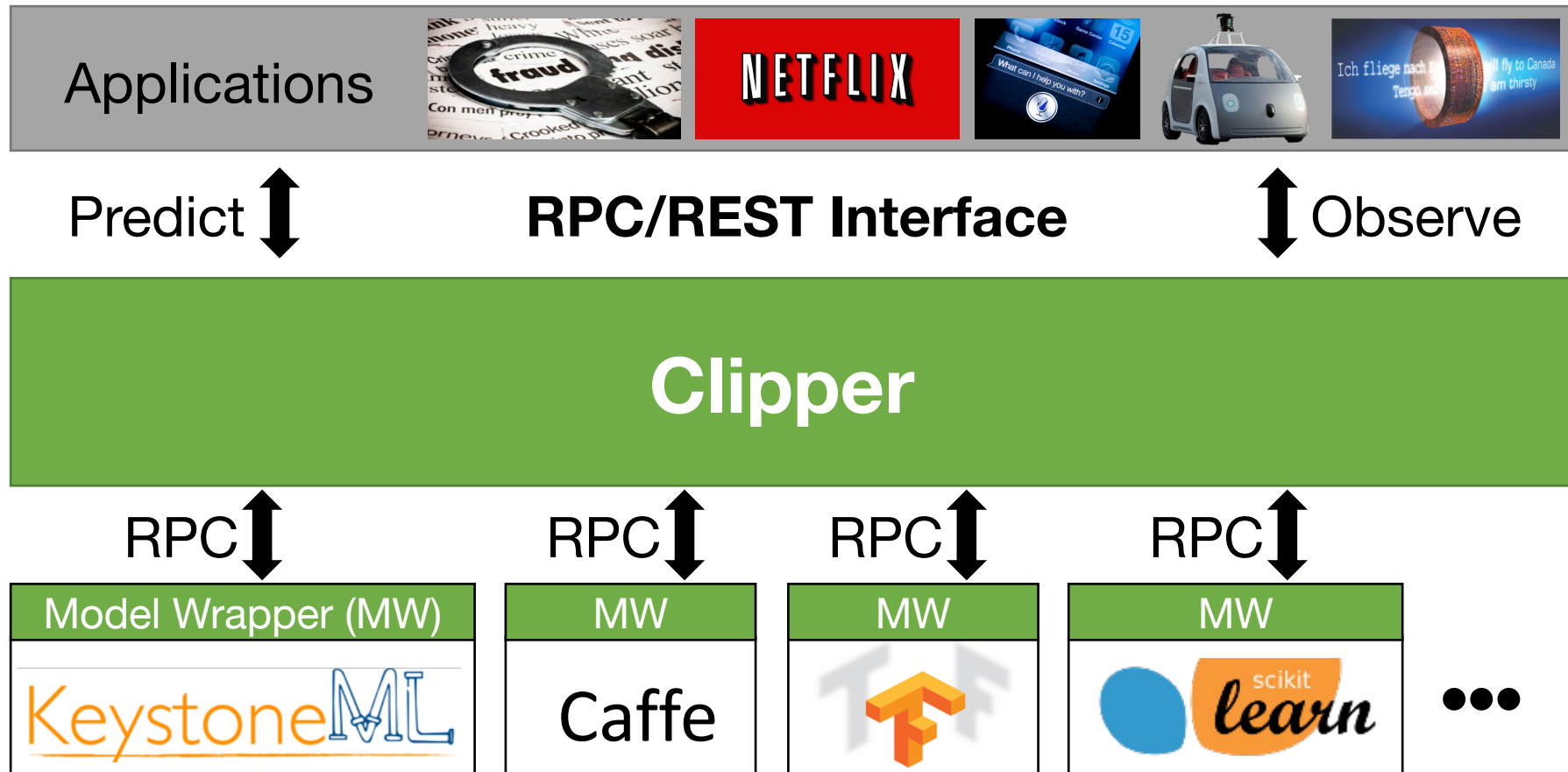


KALDI

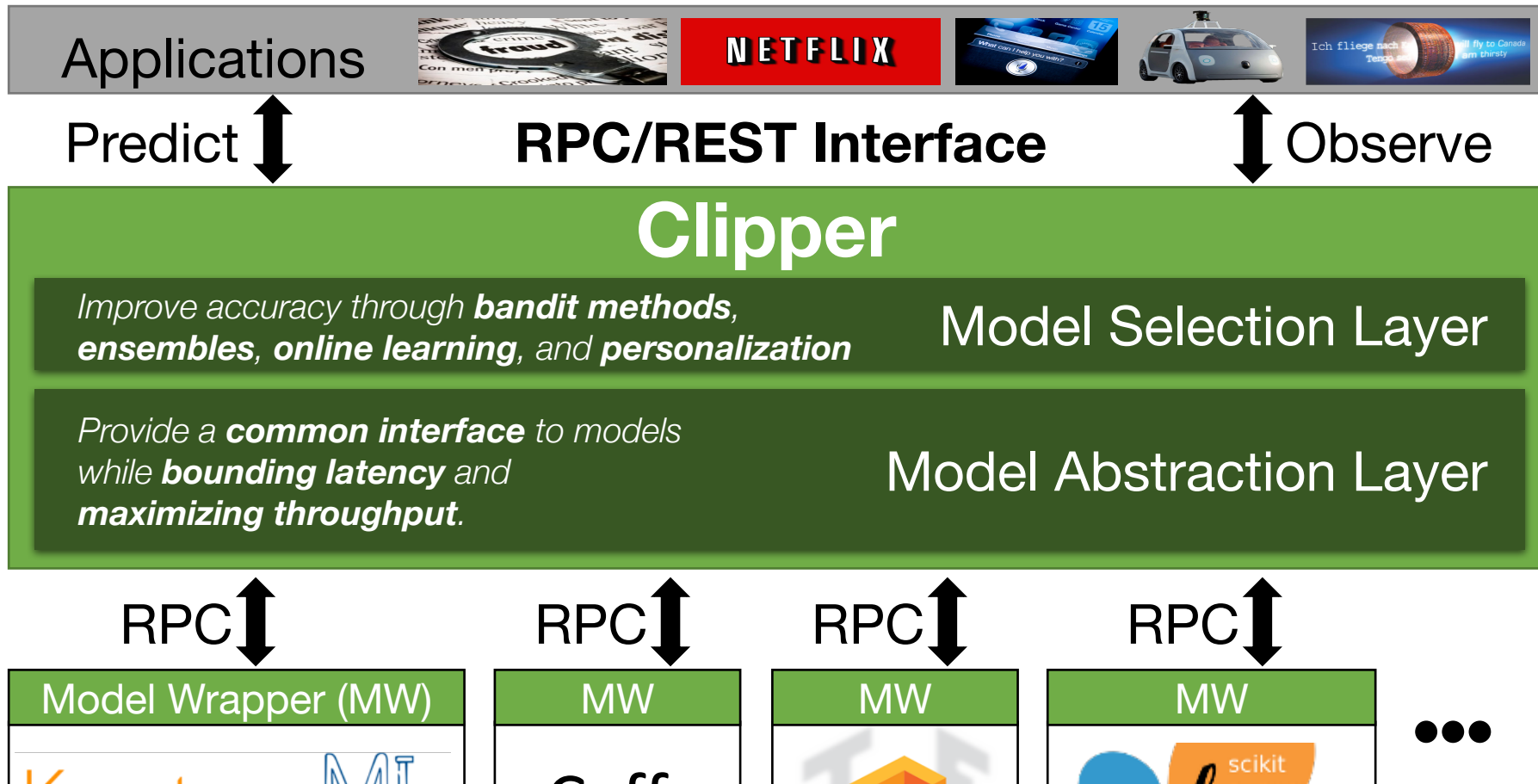
# Clipper Architecture



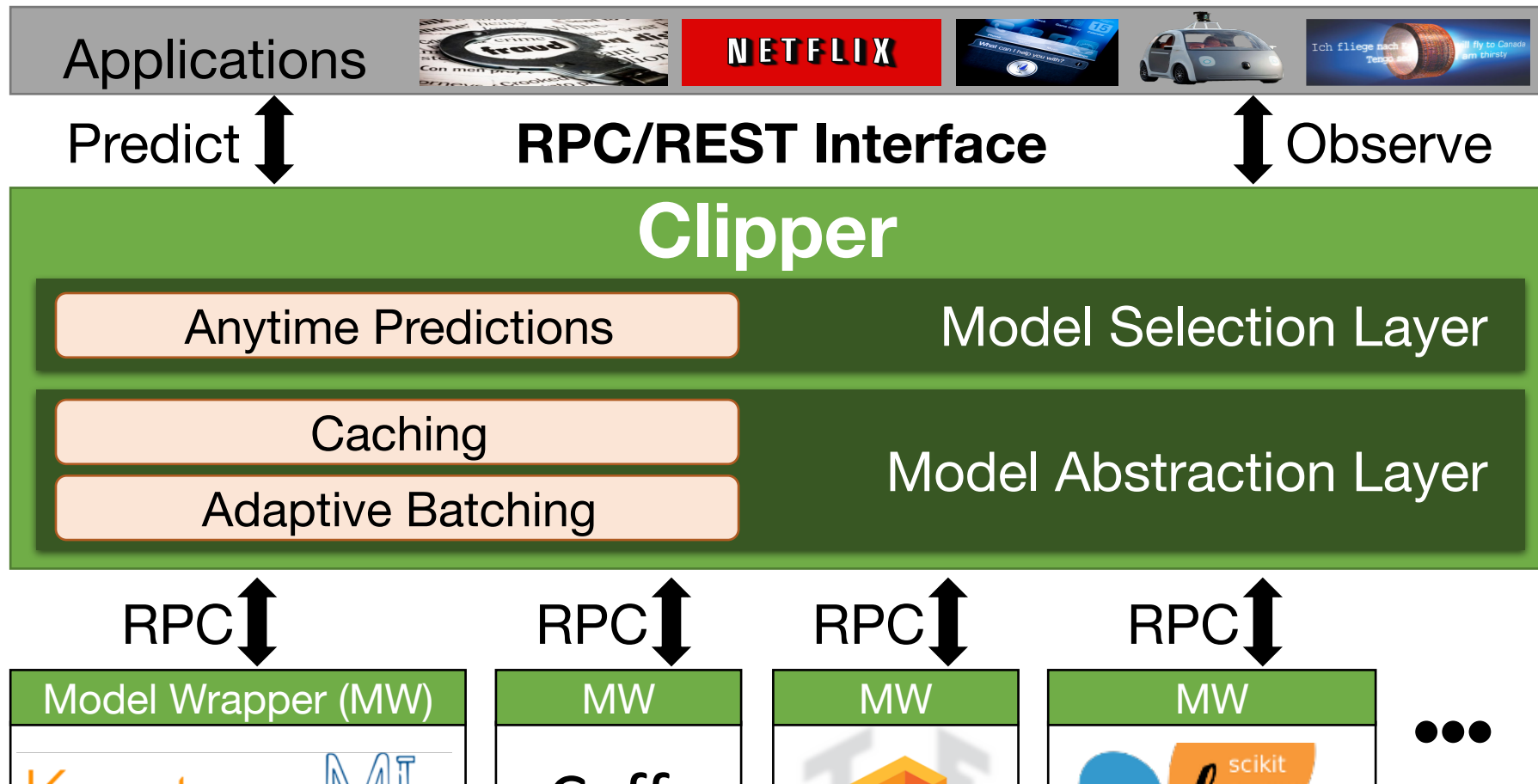
# Clipper Architecture



# Clipper Architecture



# Clipper Architecture



Caching

Adaptive Batching

Model Abstraction Layer

RPC↕

Model Wrapper (MW)

KeystoneML

RPC↕

MW

Caffe

RPC↕

MW



RPC↕

MW



...

Caching

Adaptive Batching

Model Abstraction Layer

RPC↕



RPC↕



RPC↕



RPC↕

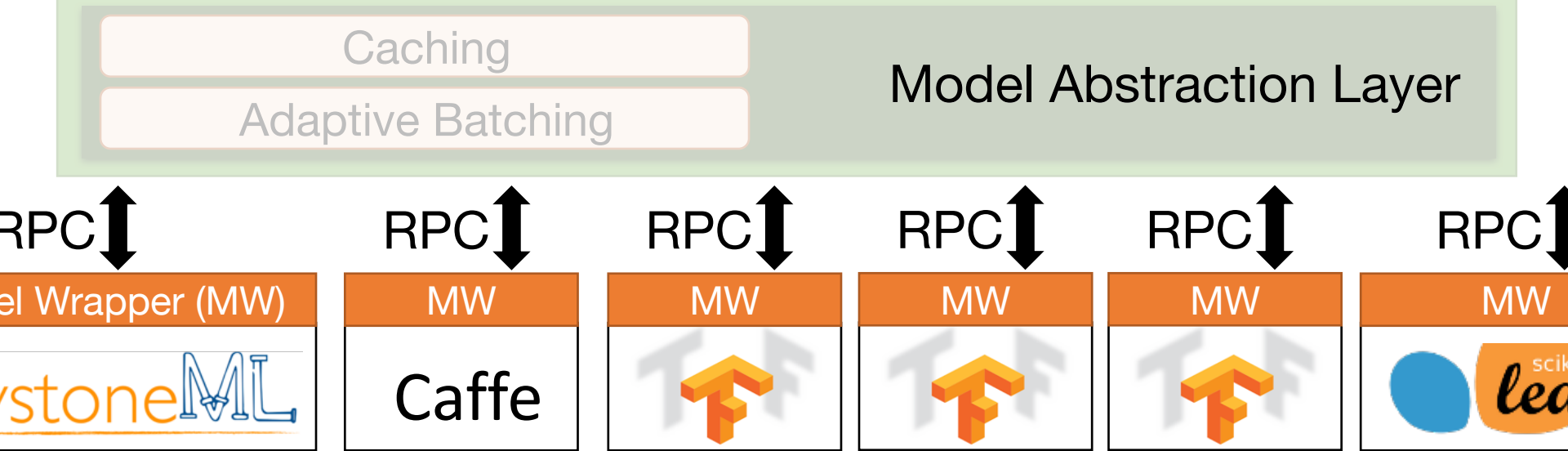


...

Provide a **common interface** to models while **bounding latency** and **maximizing throughput**.

- Models run in separate processes as Docker containers
  - Resource isolation





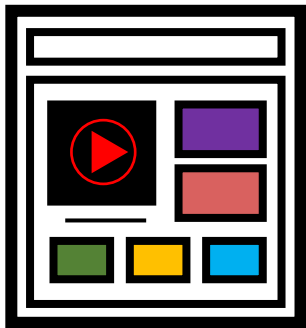
Provide a **common interface** to models while **bounding latency** and **maximizing throughput**.

- Models run in separate processes as Docker containers
  - Resource isolation
- Scaling under heavy load

**Problem:** frameworks optimized for **batch processing** not **latency**

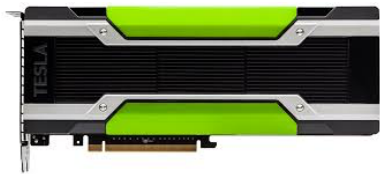
# *Adaptive Batching* to Improve Throughput

- Why batching helps:



A single page load may generate many queries

Hardware Acceleration



Helps amortize system overhead

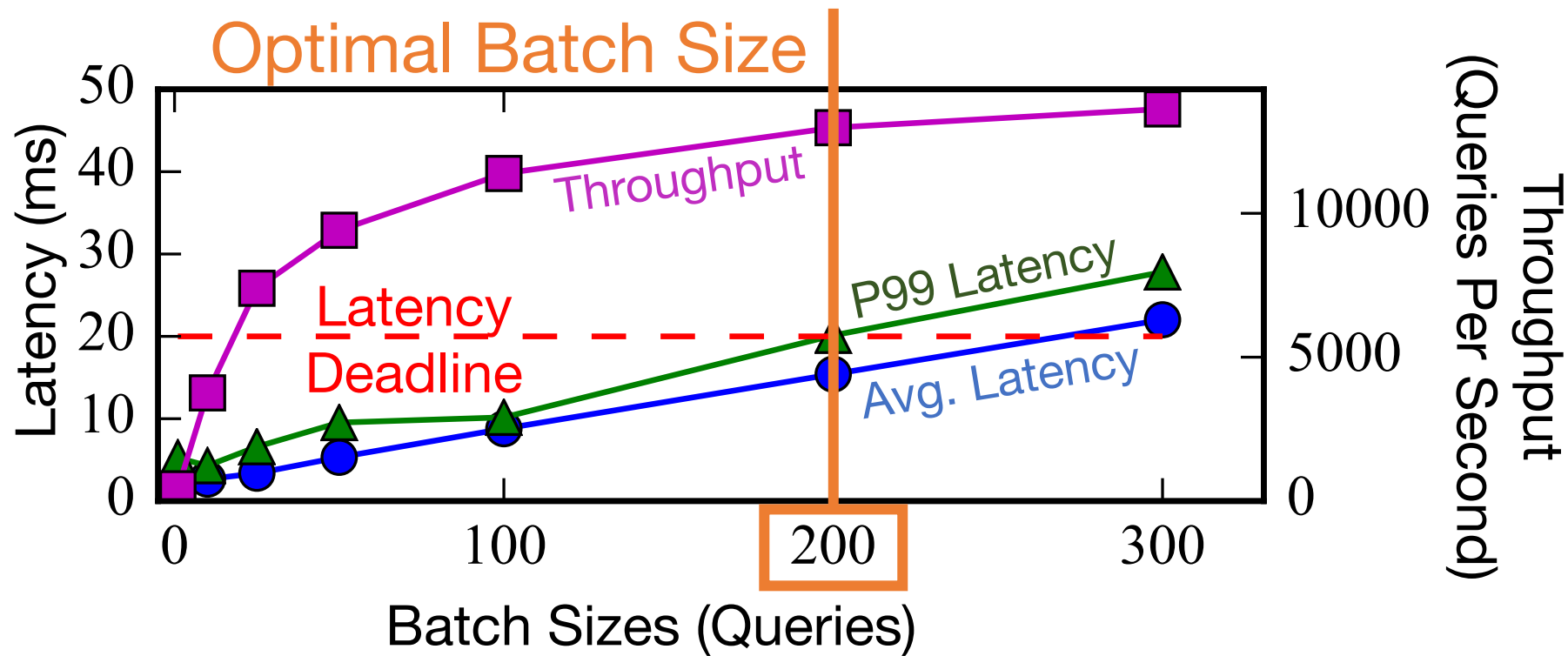
- Optimal batch depends on:
  - hardware configuration
  - model and framework
  - system load

## Clipper Solution:

be as ***slow*** as ***allowed***...

- Application specifies latency objective
- Clipper uses TCP-like tuning algorithm to **increase latency** up to the objective

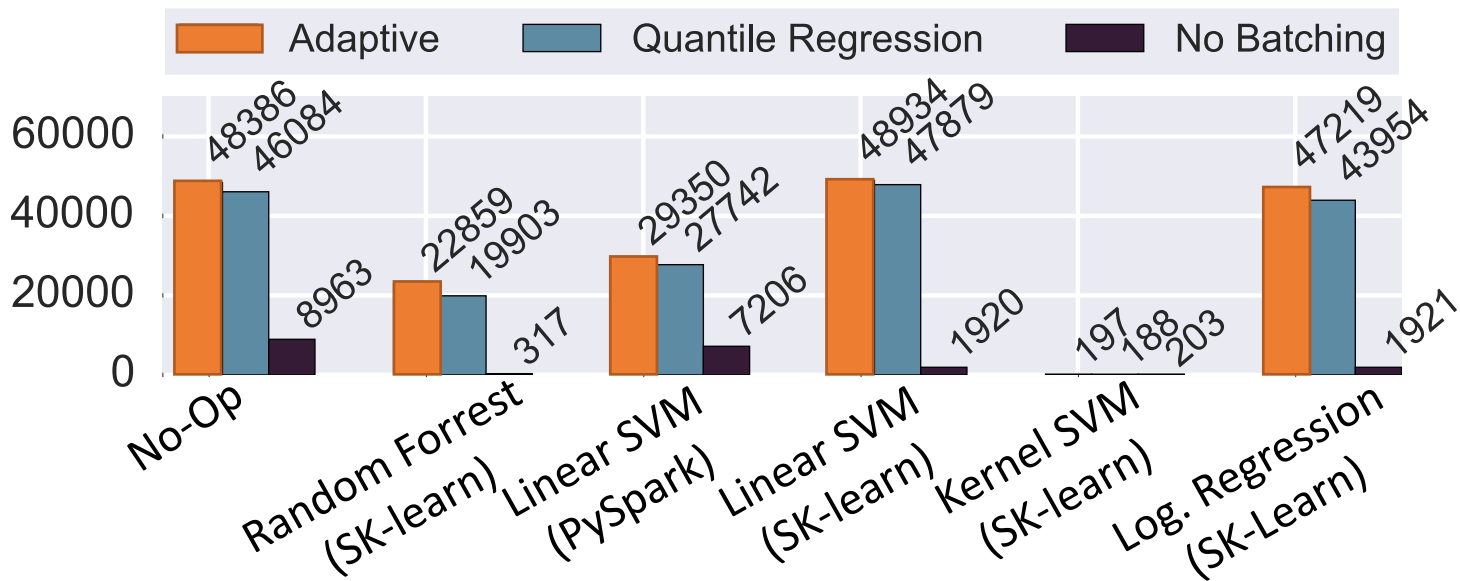
# Tensor Flow Conv. Net (GPU)

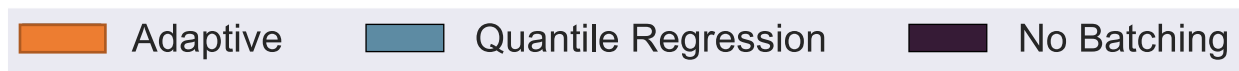


Throughput  
(QPS)



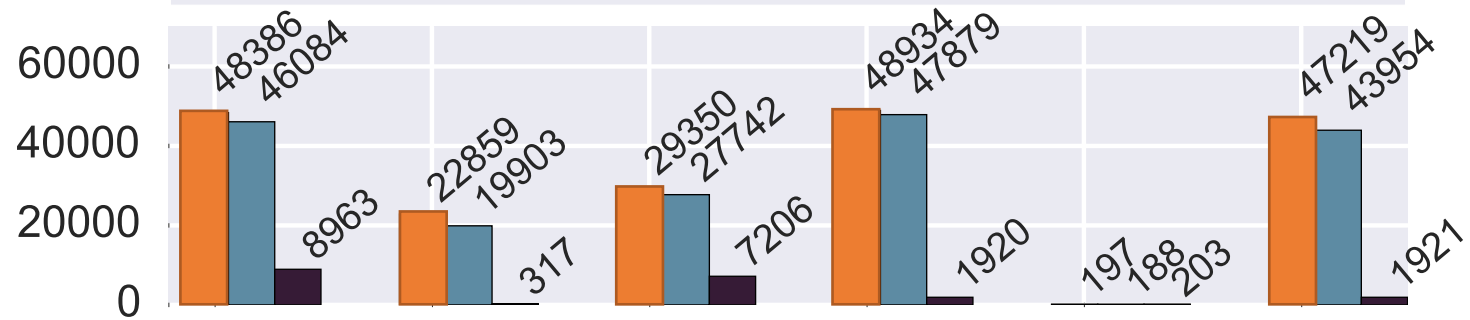
Better





Throughput  
(QPS)

↑  
Better



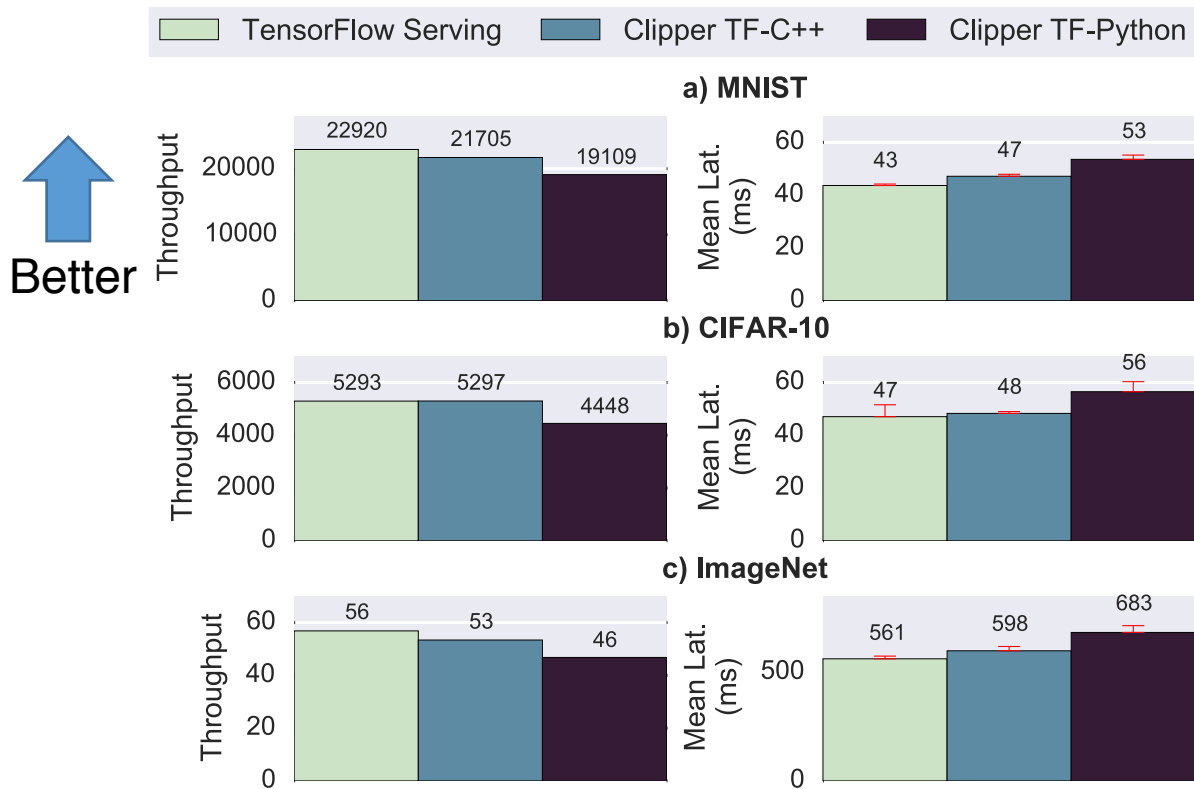
P99 Latency

↓  
Better

20000 is  
Good Enough



# Overhead of modularity?



Better

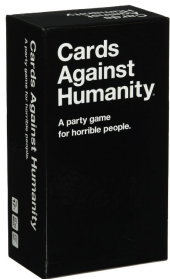


40000 is  
Good Enough

*The decoupled Clipper architecture  
can be as fast as the  
in-process approach adopted by  
TensorFlow-Serving*

# Approximate Caching to Reduce Latency

- Opportunity for caching



Popular items may be evaluated frequently

- Need for **approximation**

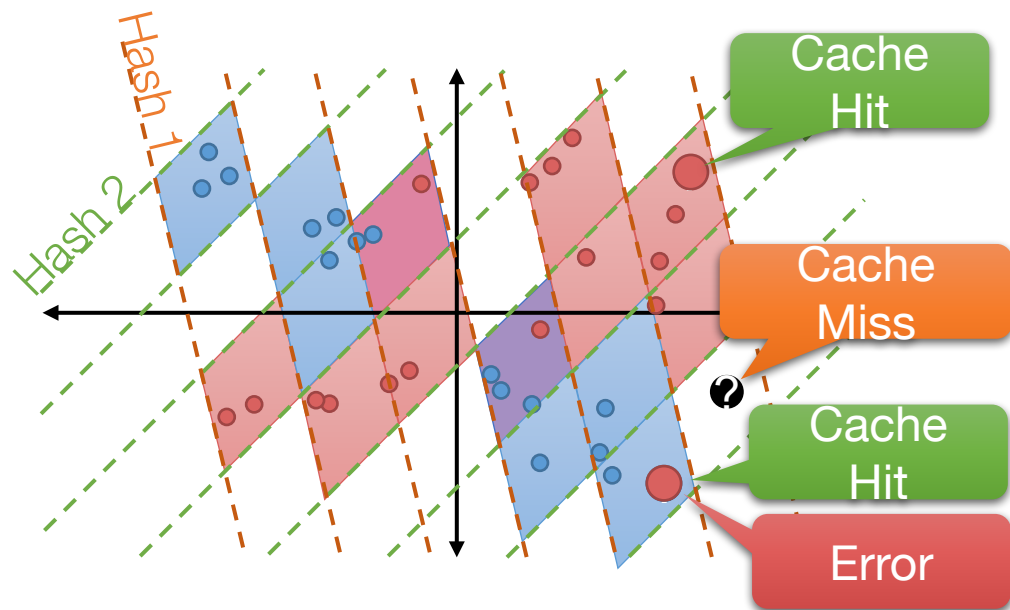


**Images**

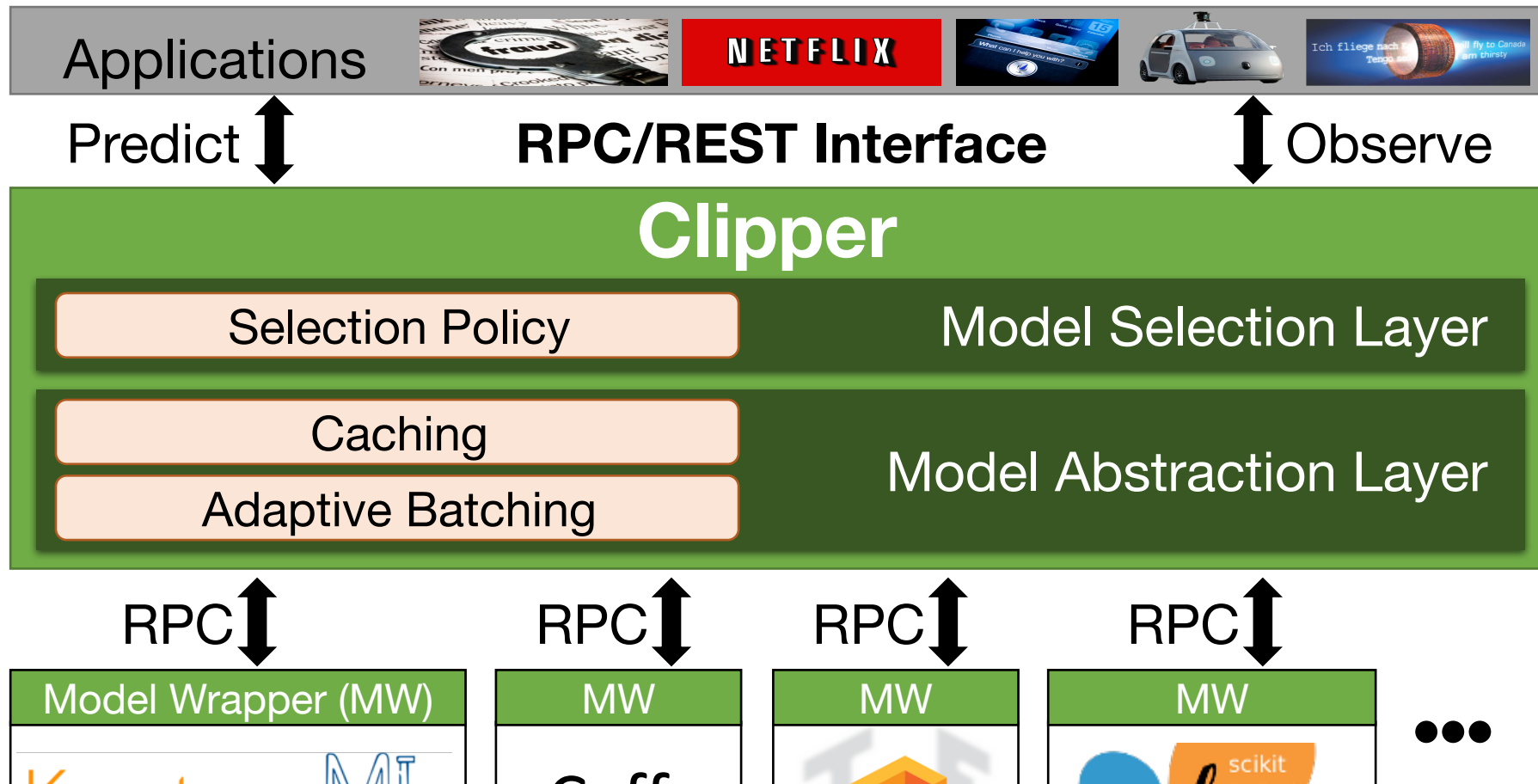
High Dimensional and continuous valued queries have low cache hit rate.

## Clipper Solution: **Approximate Caching**

apply *locality sensitive hash functions*



# Clipper Architecture





## Goal:

*Maximize **accuracy** through **bandits**, **ensembles**, **online learning**, and **personalization***

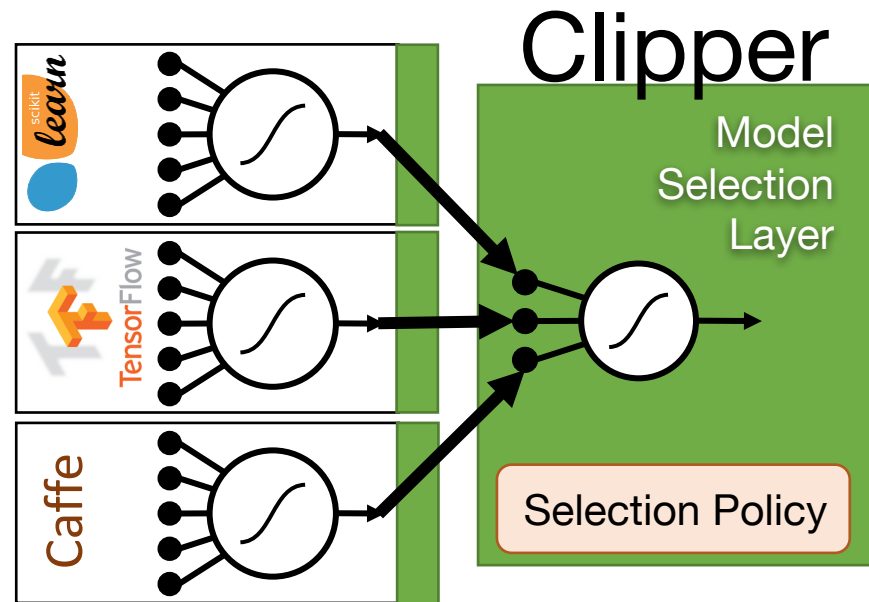
Incorporate feedback in real-time to achieve:

- **robust predictions** by adaptively combining predictions from multiple models and frameworks
- **online learning** and **personalization** by selecting and personalizing **predictions** in response to feedback

# Model Selection Policy

Improves prediction **accuracy** by:

- **Combining predictions** from multiple frameworks
  - Ensemble methods
- Incorporate real-time **feedback**
  - Personalized ensembles
  - Bandit algorithms
- Estimates **confidence** of predictions
  - Agreement between models



# Ensemble Prediction **Accuracy** (ImageNet)

System	Model	Error Rate	#Errors
Caffe	VGG	13.05%	6525
Caffe	LeNet	11.52%	5760
Caffe	ResNet	9.02%	4512
TensorFlow	Inception v3	6.18%	3088

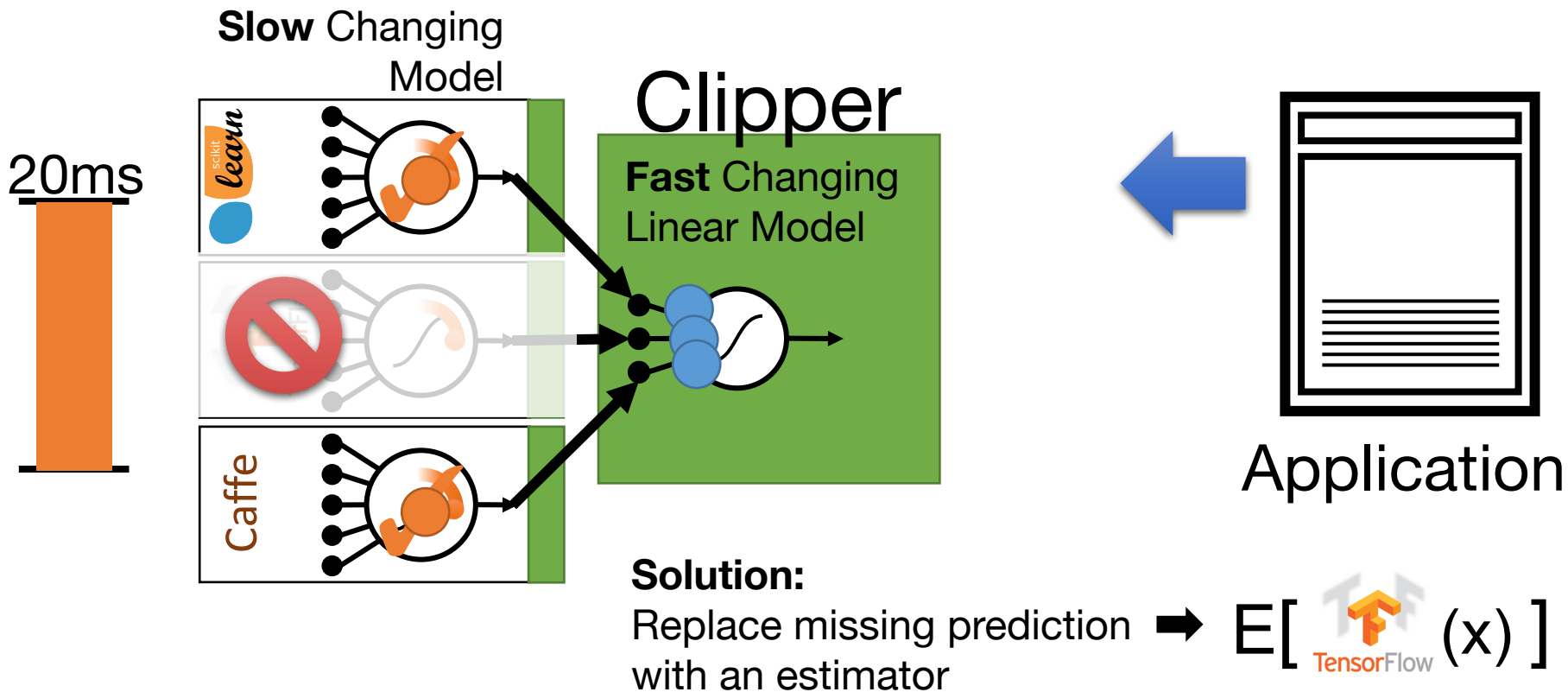
sequence of pre-trained models

# Ensemble Prediction **Accuracy** (ImageNet)

System	Model	Top-5 Error	Errors
Caffe	ResNet	9.02%	6525
Caffe	Inception v3	6.18%	5760
Caffe	Ensemble	5.86%	4512
TensorFlow	Inception v3	6.18%	3088
<b>Clipper</b>	<b>Ensemble</b>	<b>5.86%</b>	<b>2930</b>

**5.2% relative improvement  
in prediction accuracy!**

# Ensemble Methods Create Stragglers

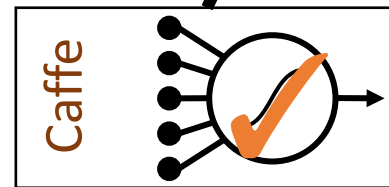
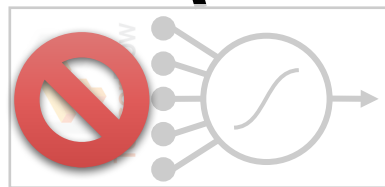
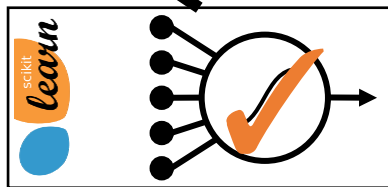


# Anytime Predictions



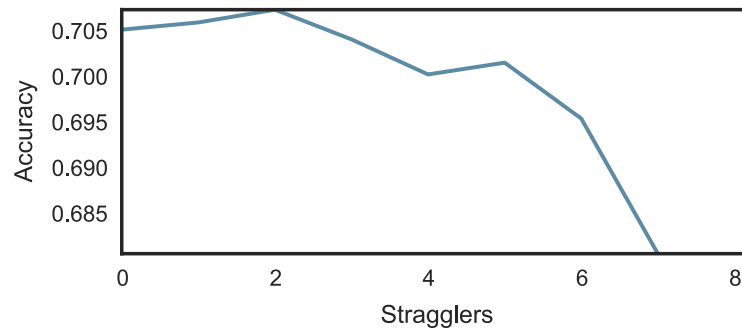
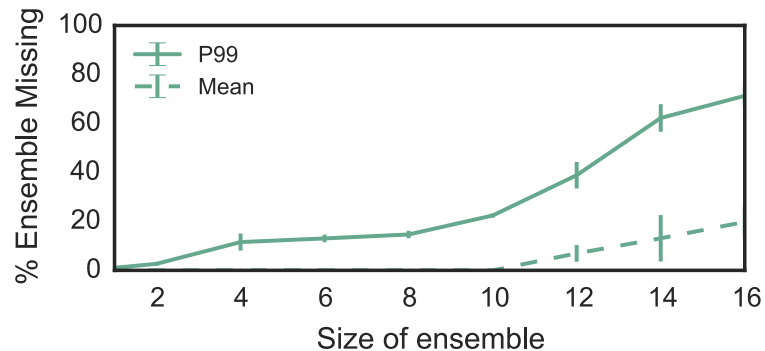
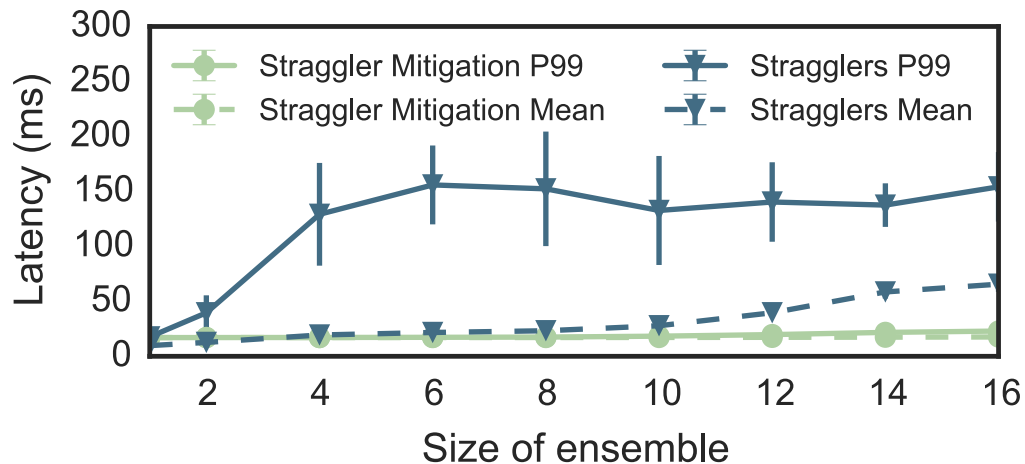
$$\underbrace{W_{\text{scikit}} f_{\text{scikit}}(x)} + \underbrace{W_{\text{TF}} \mathbb{E}_X [f_{\text{TF}}(X)]}_{\text{Fast Changing Model}} + \underbrace{W_{\text{Caffe}} f_{\text{Caffe}}(x)}$$

**Slow** Changing Model

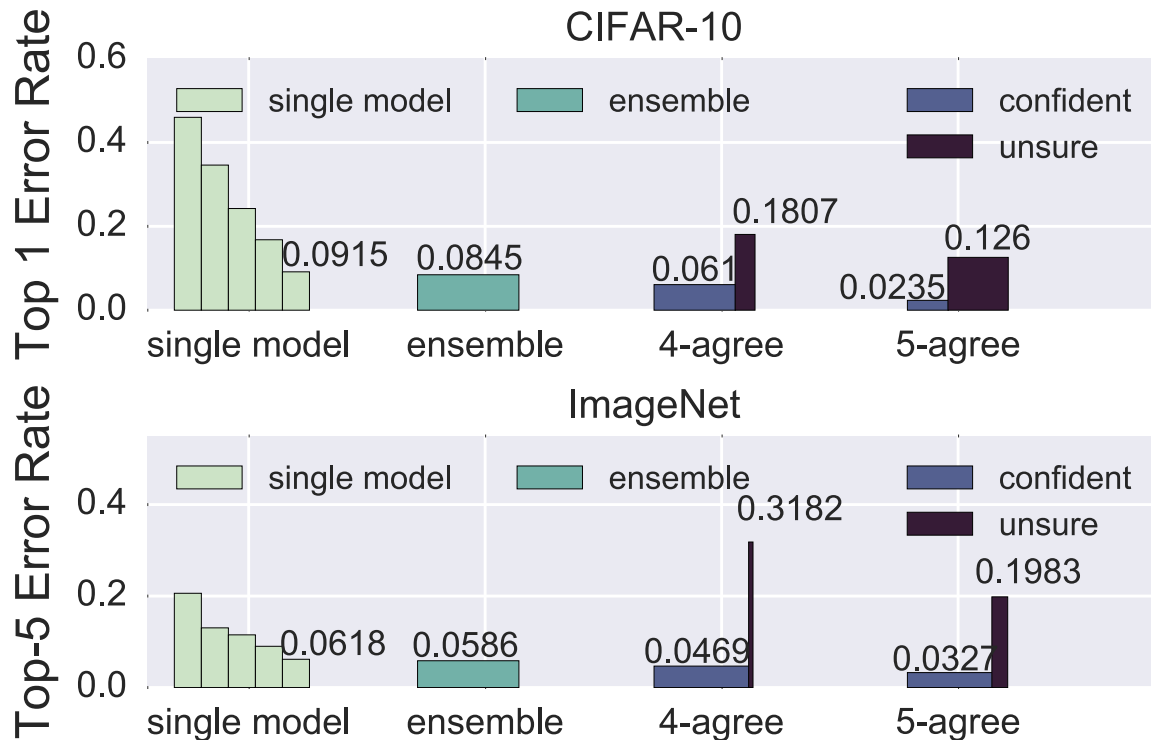


# Anytime Predictions

- Tolerates some loss of models
- Depends heavily on ensemble




# Ensemble's to Estimate Confidence







# Clipper



Clipper is a prediction serving system that spans multiple ML Frameworks and is designed to

- to **simplifying** model serving
- **bound latency** and **increase throughput**
- and enable **real-time learning** and **personalization across** machine learning frameworks

*“Clipper: A Low-Latency Online Prediction Serving System”*

<https://github.com/ucbrise/clipper> (open source)

# Ongoing Clipper Subprojects

- **Adaptive Batching for Prediction**

- Leverage internal data-parallelism and hardware acceleration

- **Approximate Caching**

- Detect “similar” queries and re-use cached predictions

- **Prediction Cascades**

- Automatically deriving cascades of increasingly GPU intensive models

- **RL/Control**

- Serving and updating RL policies based on feedback

- **Scheduling and resource allocation**

- Reduce the need to over-provision for bursty workloads



We are developing new technologies that will enable applications to make low-latency intelligent decisions on live data with strong security guarantees.

