

# Prediction Serving

*what happens after learning?*

**Joseph E. Gonzalez**

Asst. Professor, UC Berkeley

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

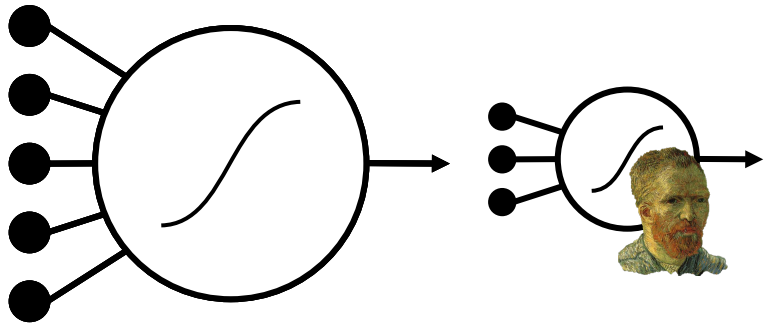
Co-founder, Dato Inc.

[joseph@dato.com](mailto:joseph@dato.com)



# Outline

 **VELOX**



**Clipper**

theano

Dato



Caffe

amc



learn

scikit



vw

Keystone

ML

Cre

TensorFlow

mxnet



KALDI

Daniel Crankshaw, Xin Wang  
Michael Franklin, & Ion Stoica

# Learning



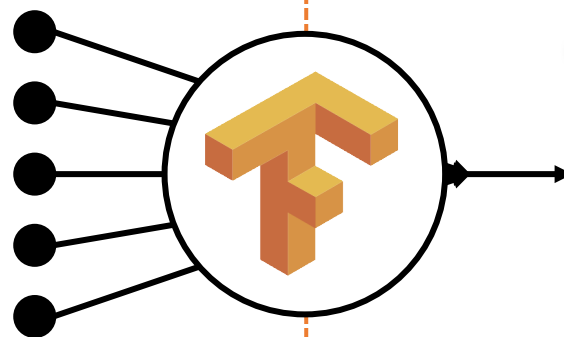
**Timescale:** minutes to days

**Systems:** offline and batch optimized

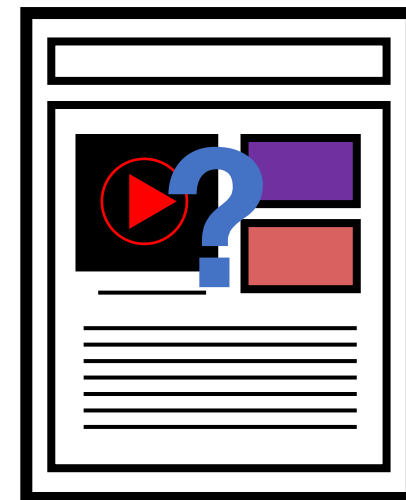
*Heavily studied ... major focus of the **AMPLab***

# Learning

# Inference



Big Model



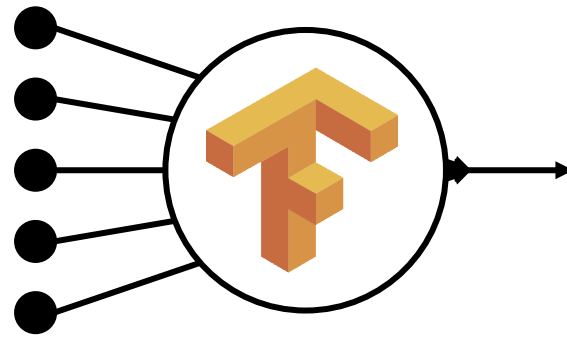
Application



# Learning

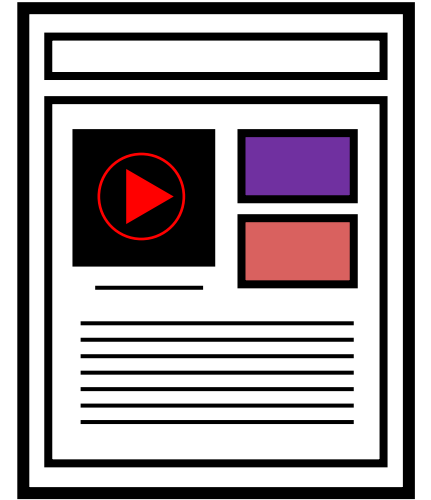


Training



Big Model

# Inference



Application

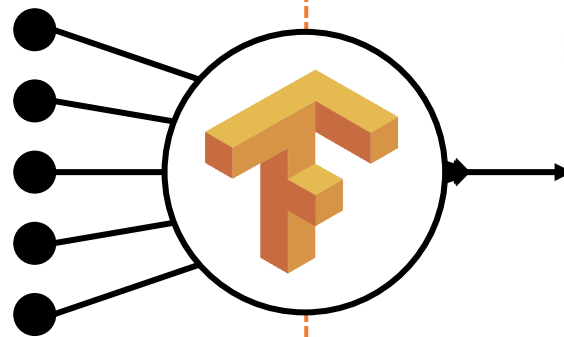
**Timescale:** ~10 milliseconds

**Systems:** *online* and *latency* optimized

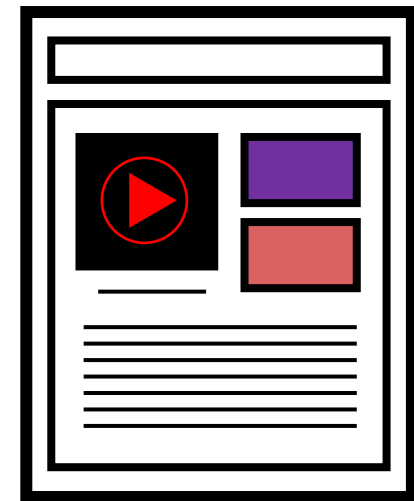
*Less studied ...*

# Learning

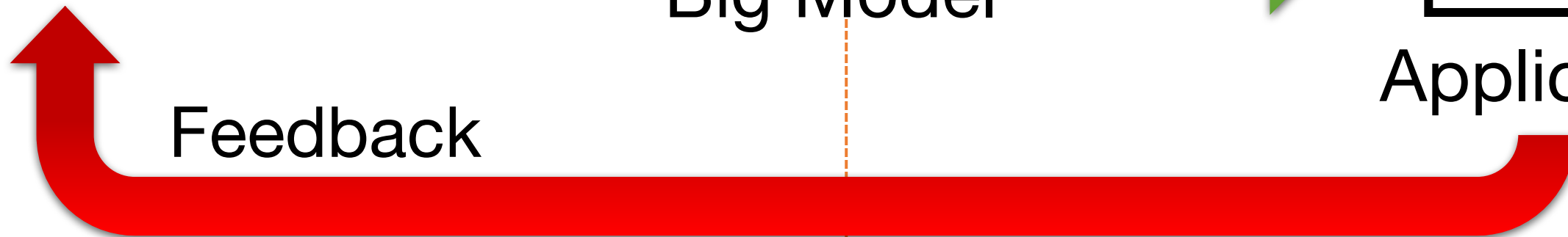
# Inference



Big Model

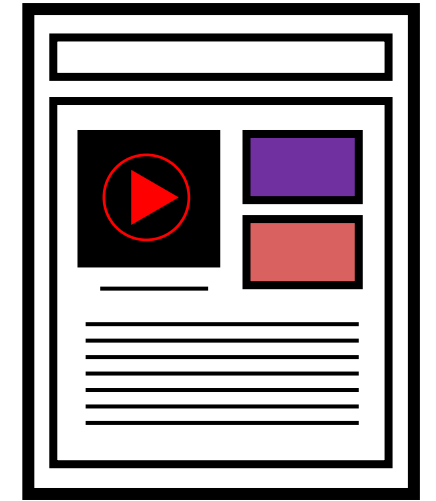
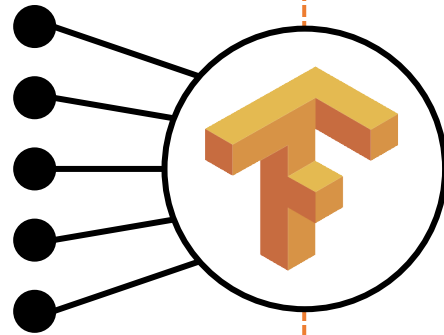


Application



# Learning

# Inference



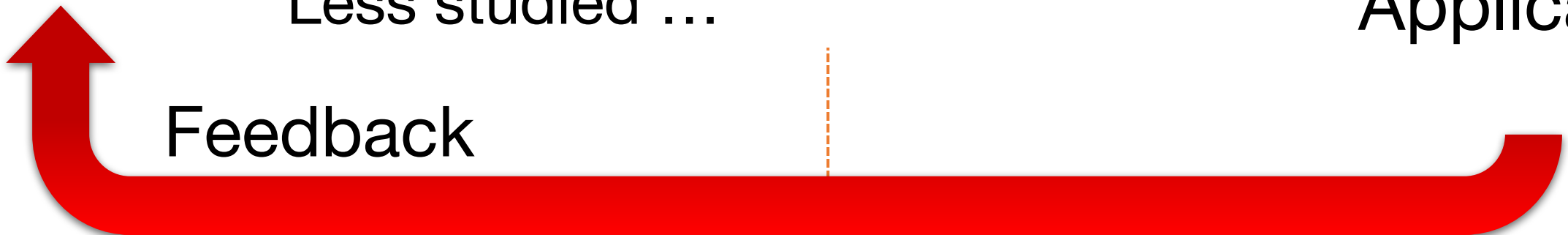
Application

**Timescale:** hours to weeks

**Systems:** combination of systems

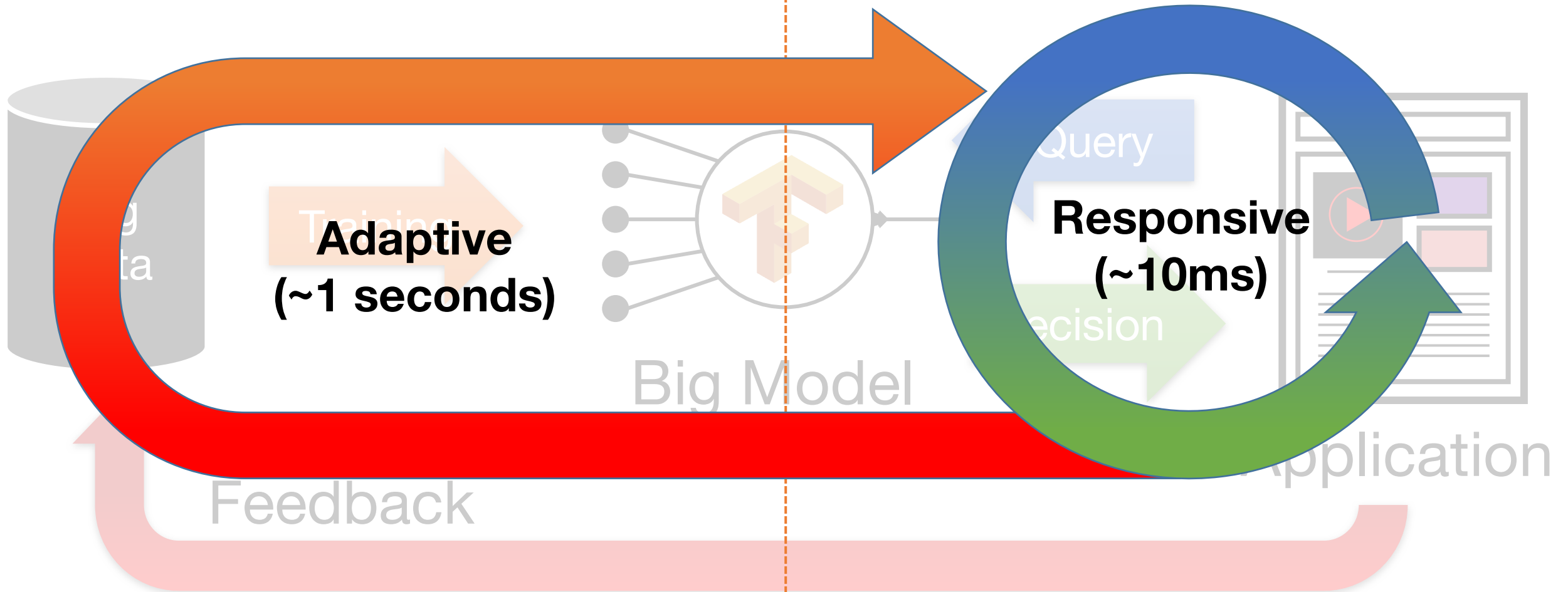
Less studied ...

Feedback



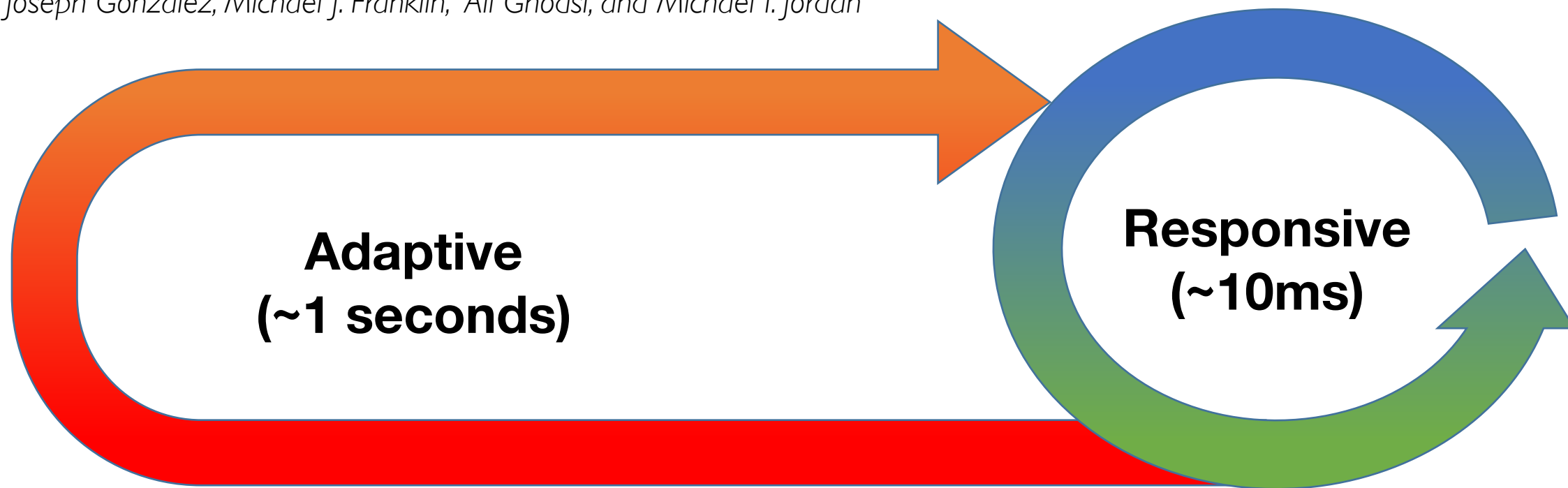
# Learning

# Inference



# **VELOX** Model Serving System [CIDR'15]

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

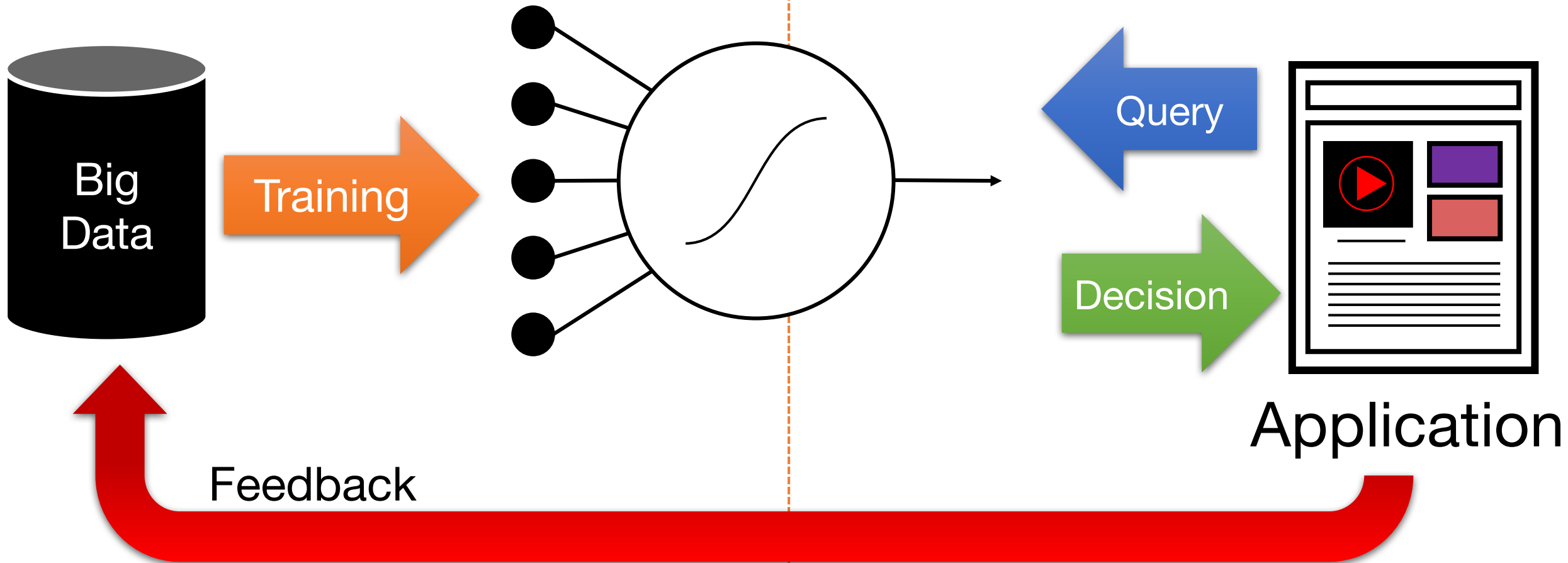


## **Key Insight:**

*Decompose models into fast and slow changing components*

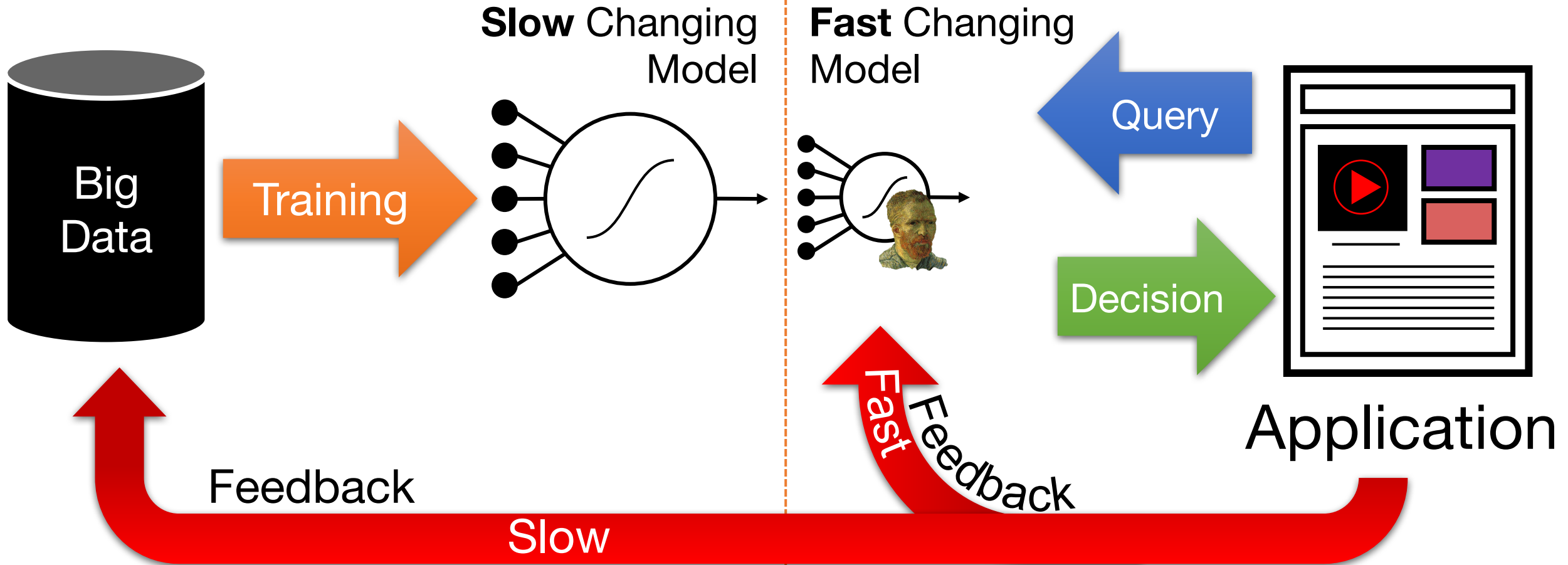
# Learning

# Inference



# Learning

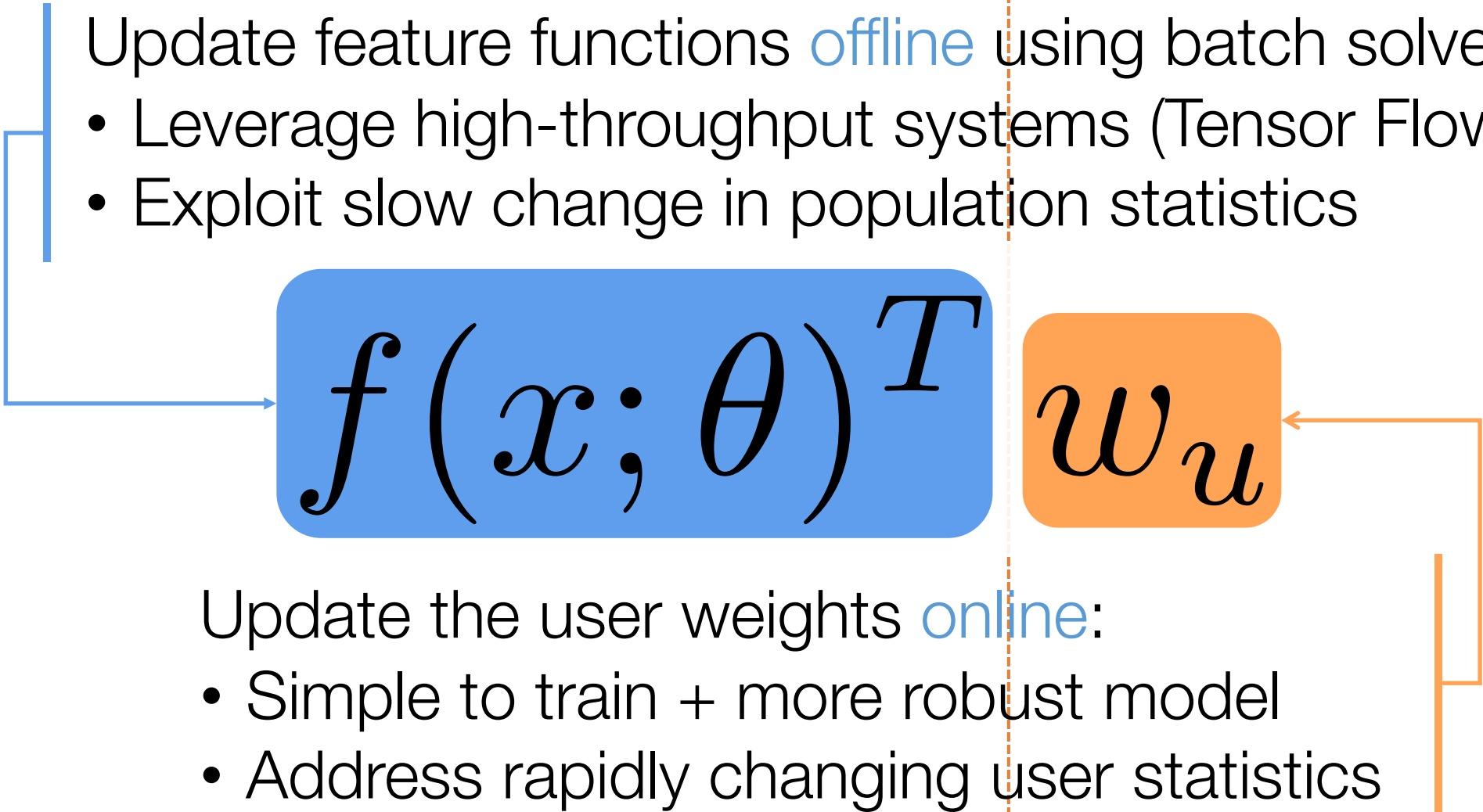
# Inference



# Hybrid Offline + Online Learning

Update feature functions **offline** using batch solvers

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



The diagram illustrates the components of a hybrid learning model. A vertical dashed orange line separates the offline and online learning processes. On the left, a blue rounded rectangle contains the feature function  $f(x; \theta)^T$ . A blue arrow points from the text 'Update feature functions offline' to this box. On the right, an orange rounded rectangle contains the user weight  $w_u$ . An orange arrow points from the text 'Update the user weights online:' to this box. A blue line connects the top of the feature function box to the top of the user weight box, and an orange line connects the bottom of the user weight box to the bottom of the feature function box, indicating a relationship or interaction between the two.

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

Update the user weights **online**:

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



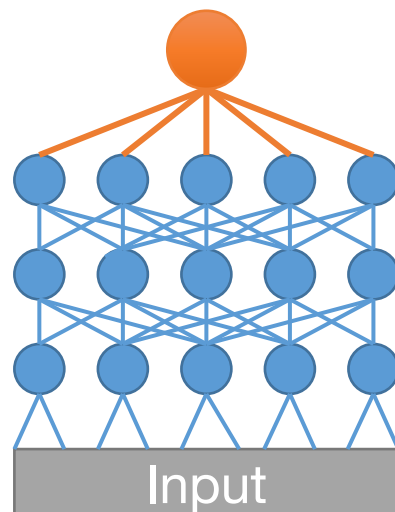
# Common modeling structure

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

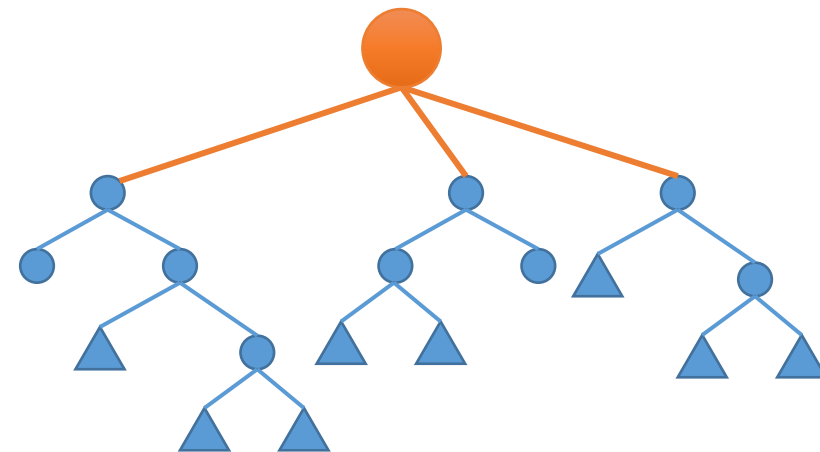
Matrix  
Factorization



Deep  
Learning

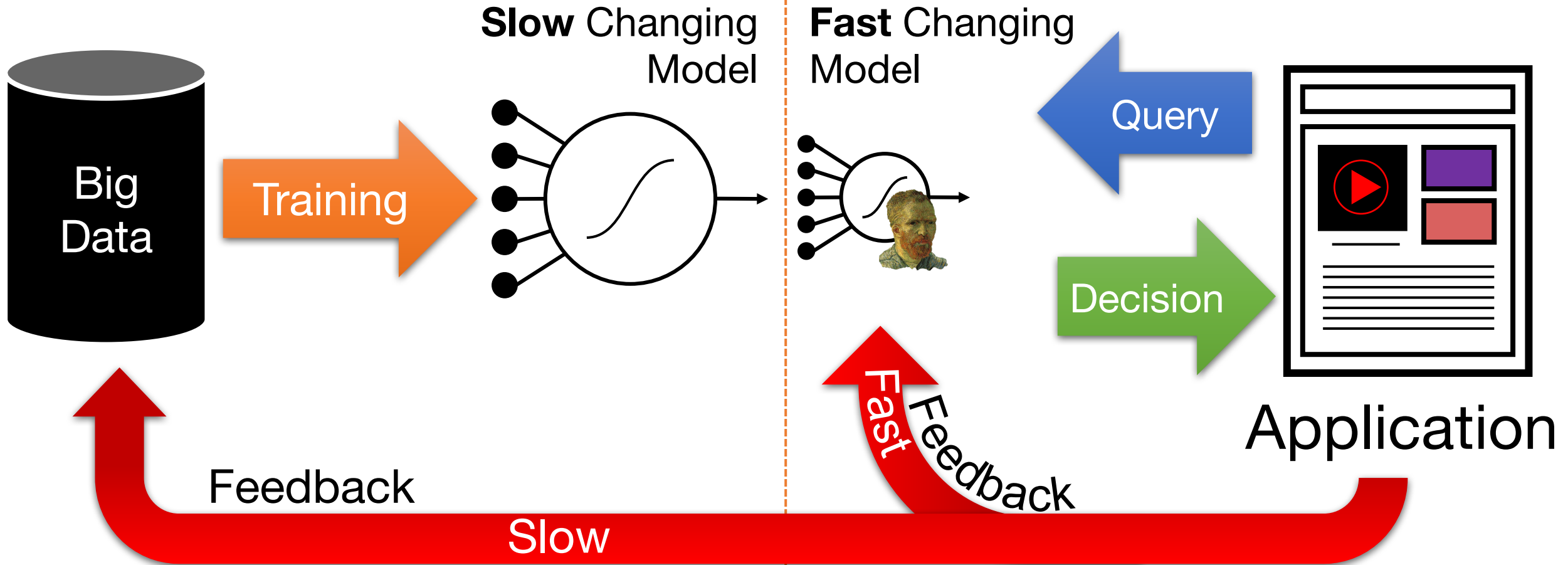


Ensemble  
Methods



# Learning

# Inference



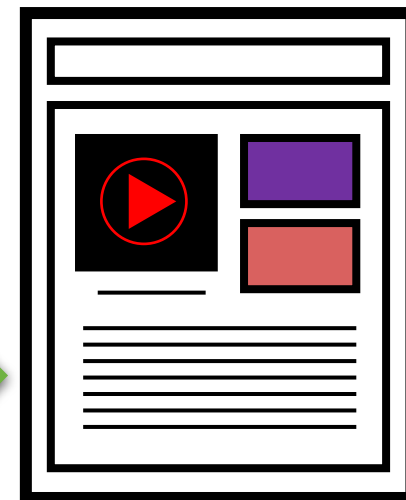
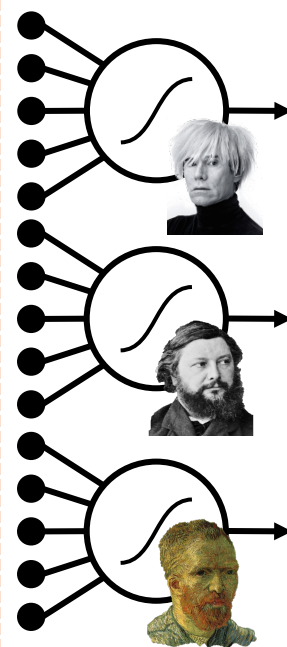
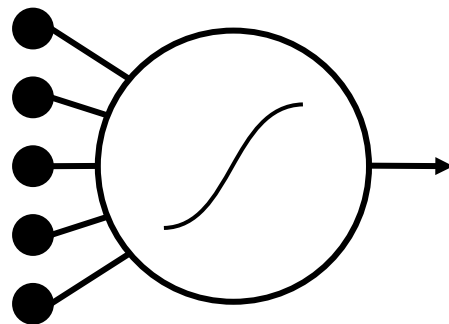
# Learning

Fast Changing  
Model per user

# Inference



Slow Changing  
Model



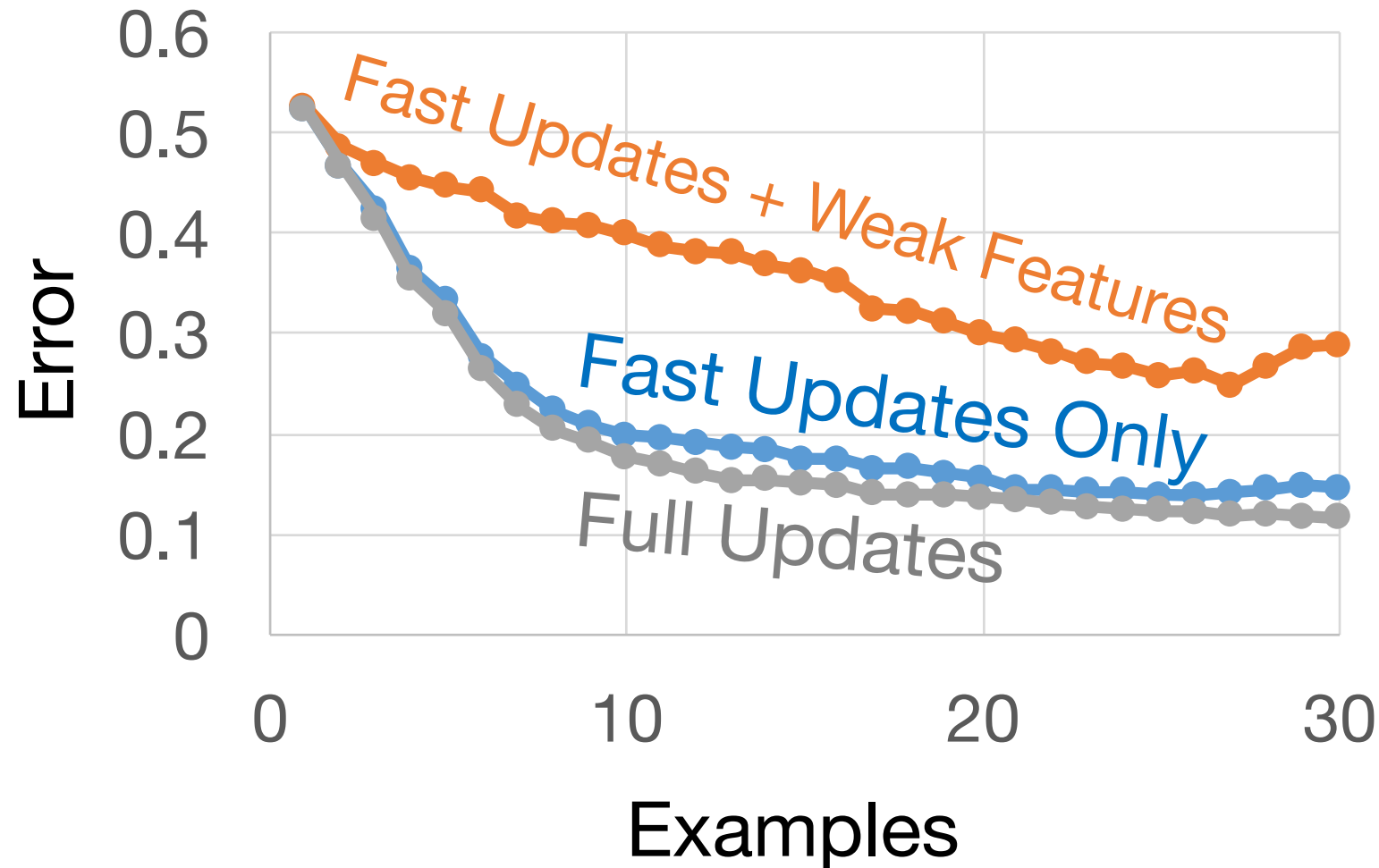
Application



Feedback

Slow

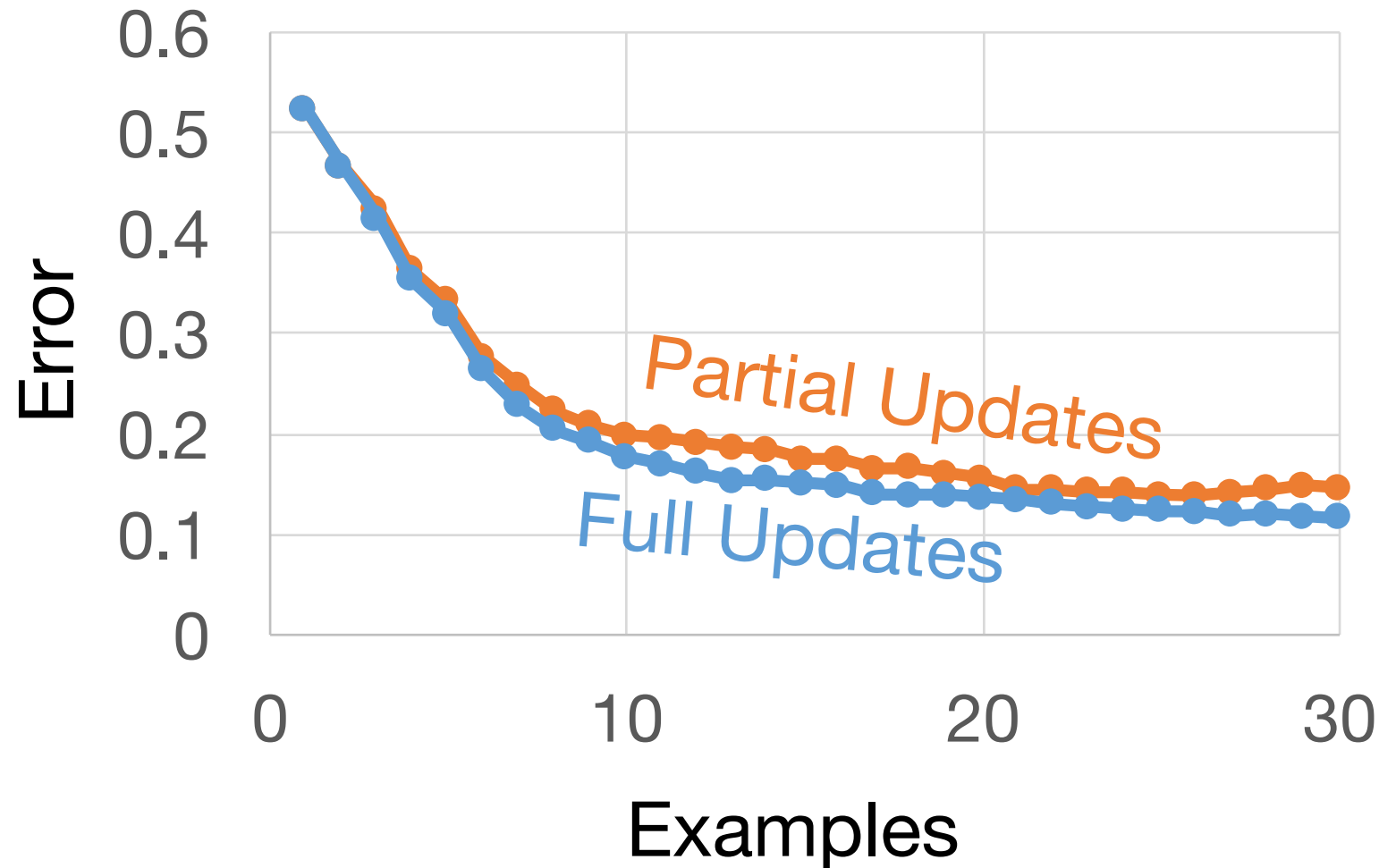
# Velox Online Learning for Recommendations (20-News Groups)



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

*>4 orders-of-magnitude  
**faster adaptation**  
given sufficient offline  
training data*

# Velox Online Learning for Recommendations (20-News Groups)



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

*>4 orders-of-magnitude  
faster adaptation*

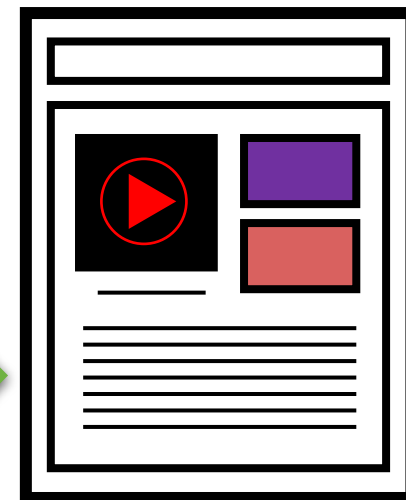
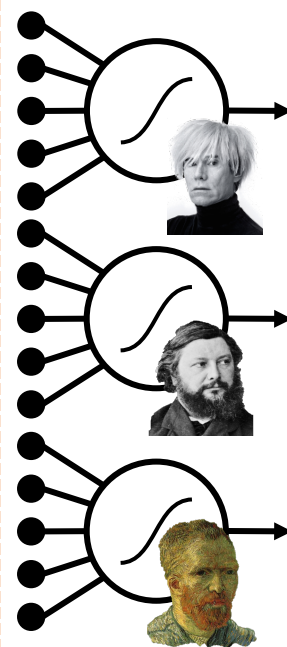
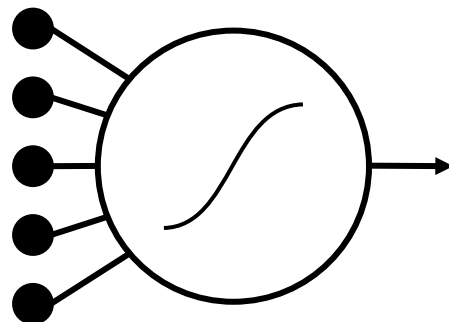
# Learning

Fast Changing  
Model per user

# Inference

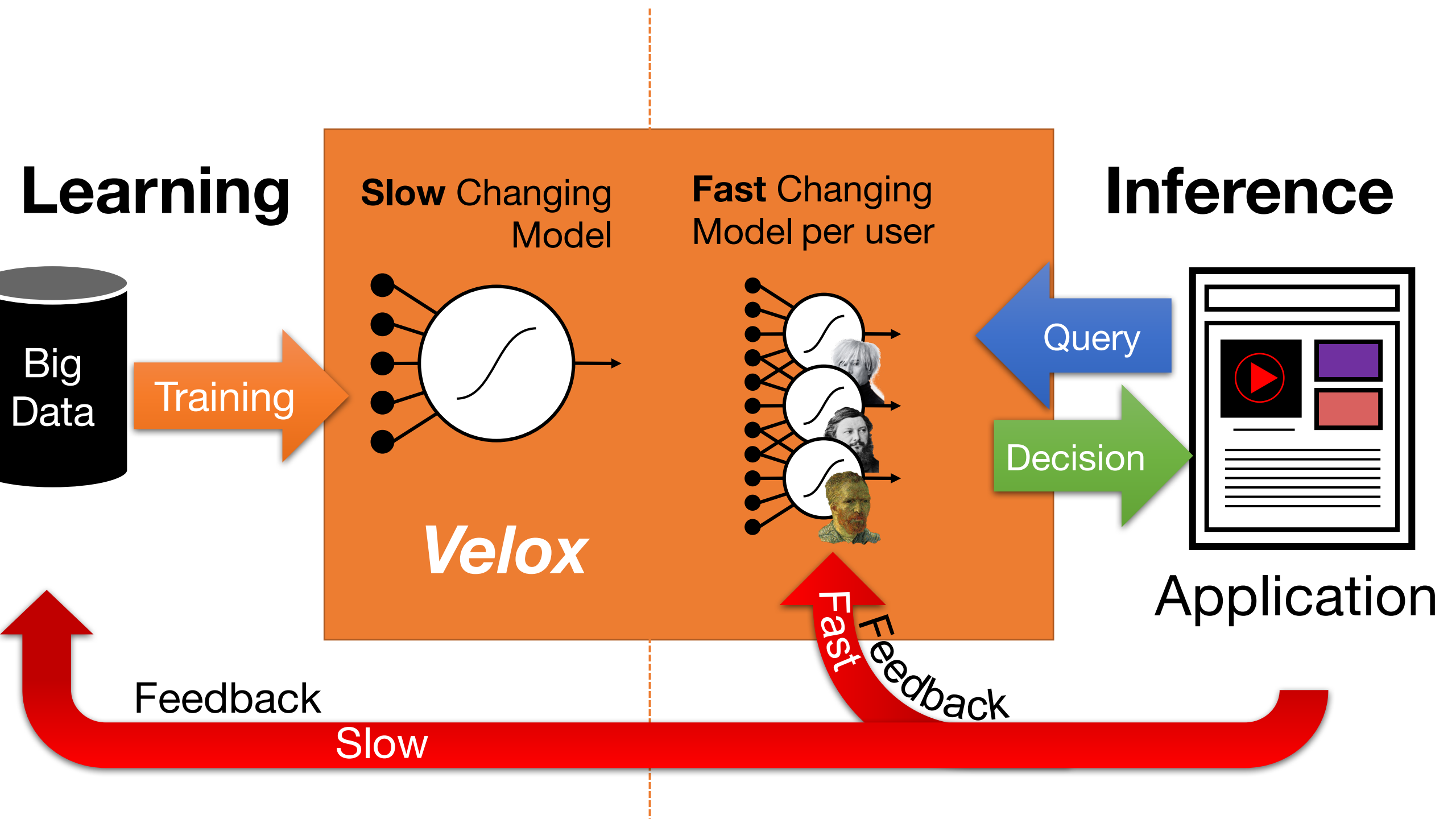


Slow Changing  
Model



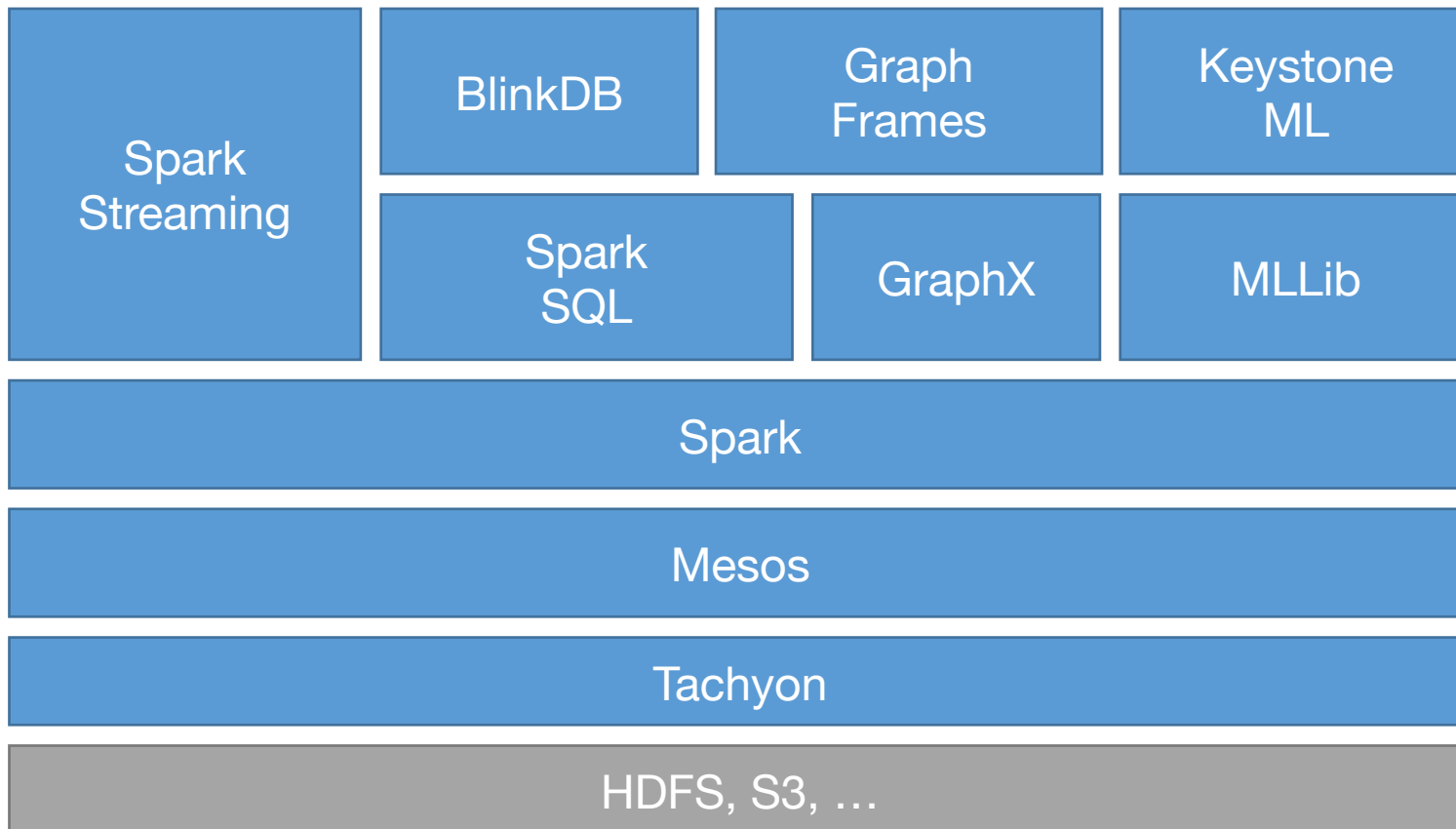
Application





# **VELOX**: the Missing Piece of BDAS

## Learning

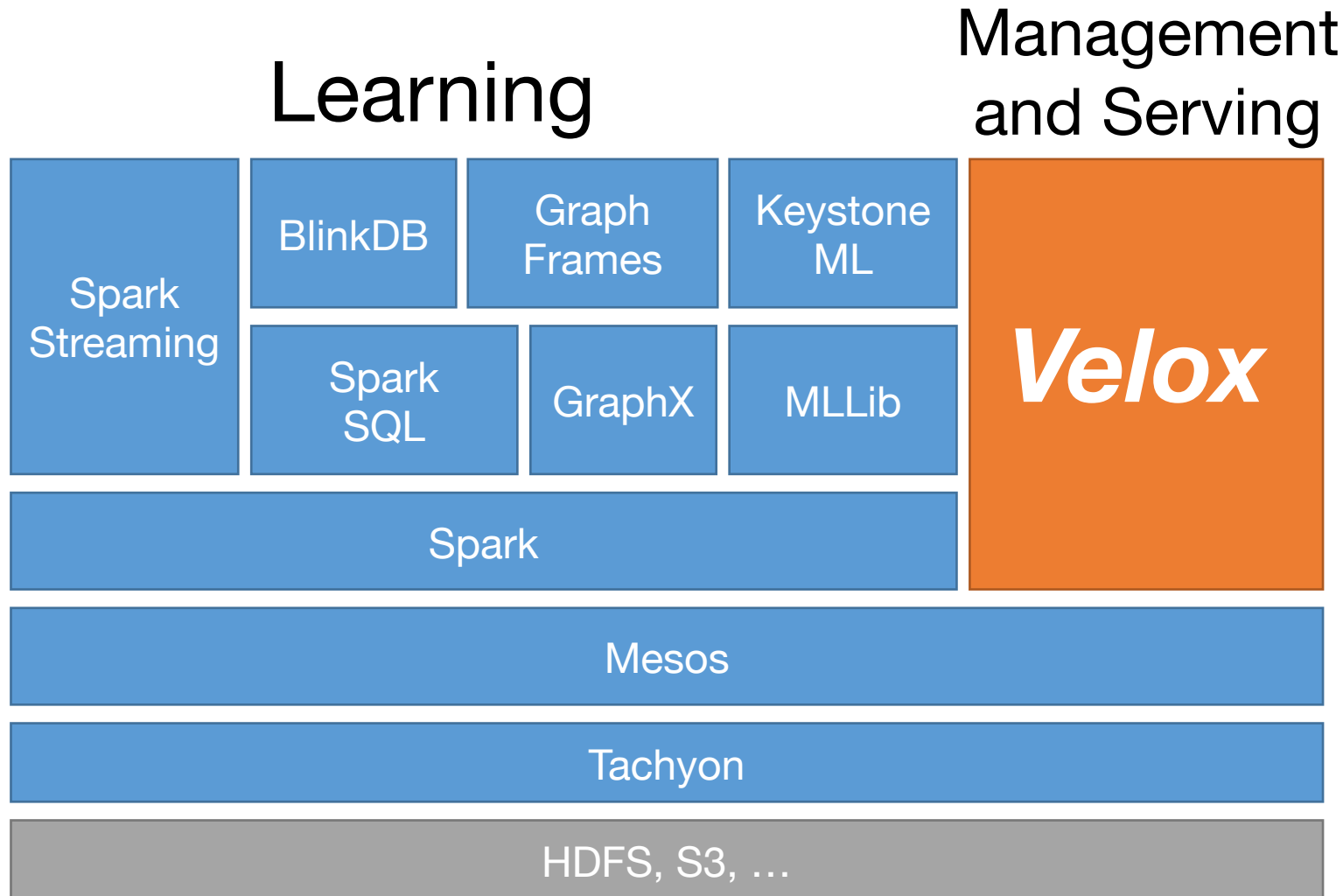


— **amplab** 

**B**erkeley  
**D**ata  
**A**nalytics  
**S**tack



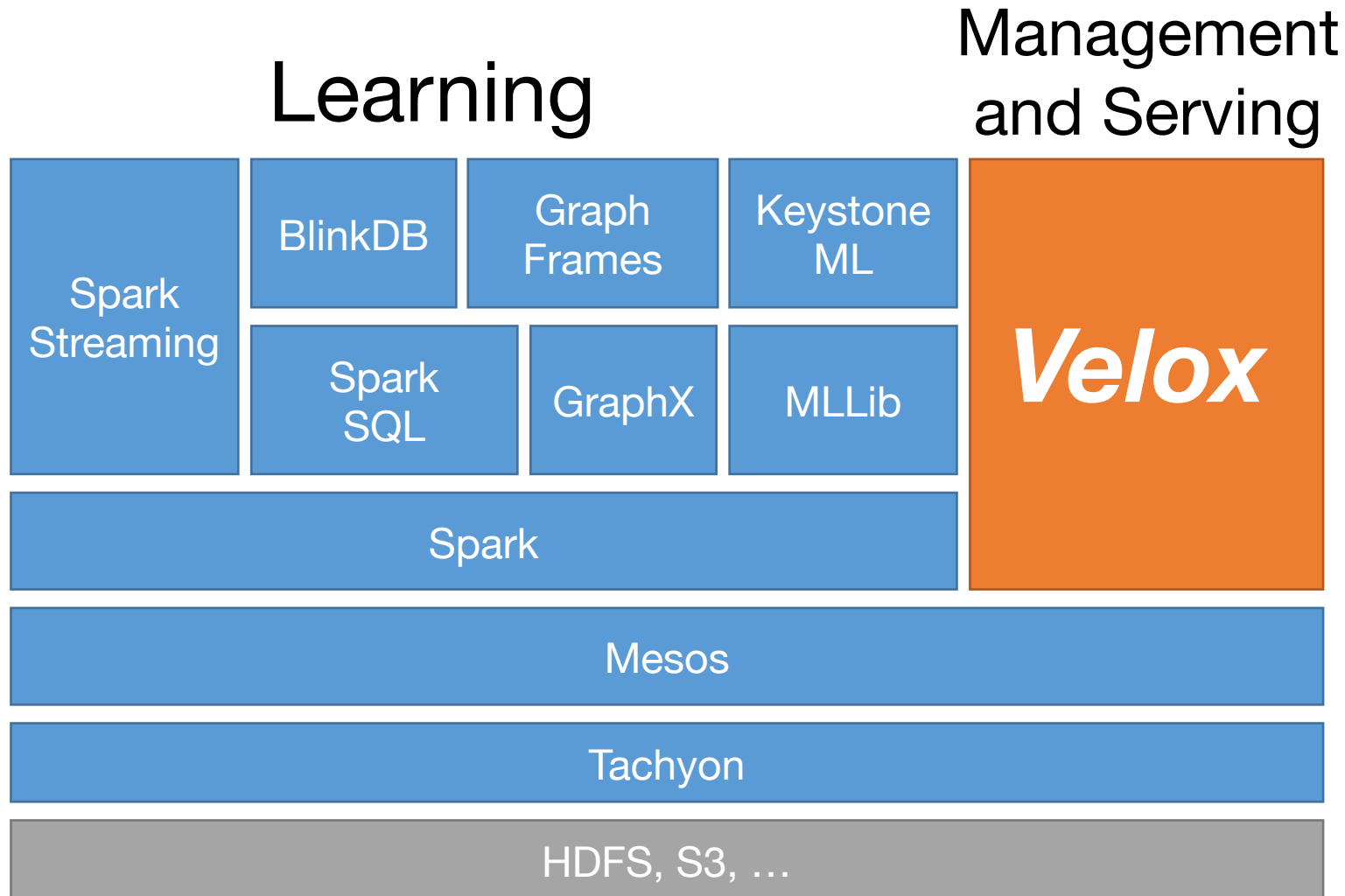
# **VELOX**: the Missing Piece of BDAS



— **amplab** 

**B**erkeley  
**D**ata  
**A**nalytics  
**S**tack

# **VELOX**: the Missing Piece of BDAS



— **amplab** 

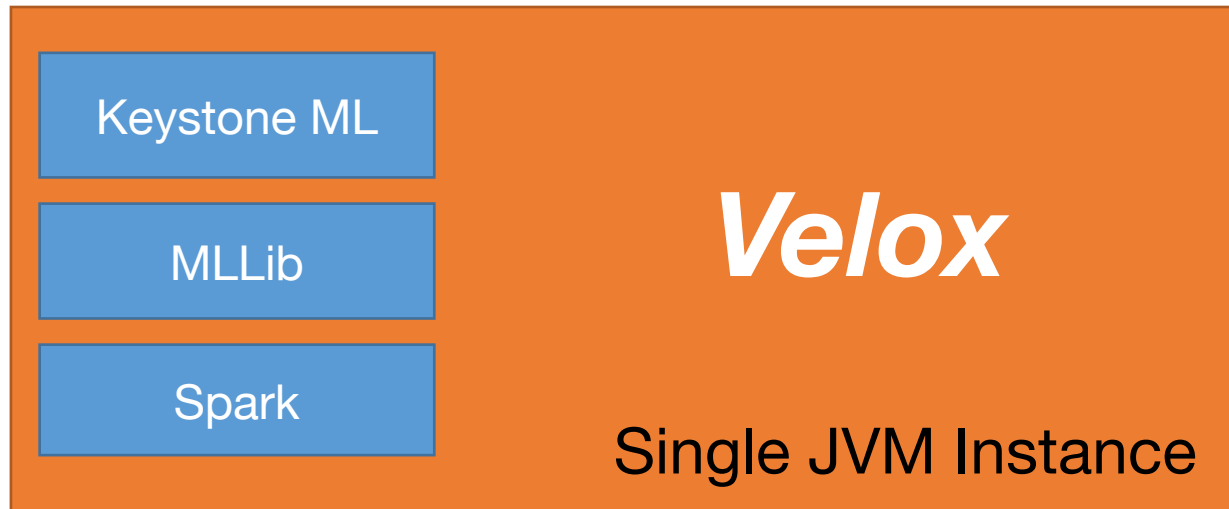
**B**erkeley  
**D**ata  
**A**nalytics  
**S**tack

# **VELOX** Architecture

Fraud  
Detection



Content  
Rec.



# **VELOX** Architecture

Fraud  
Detection



Content  
Rec.



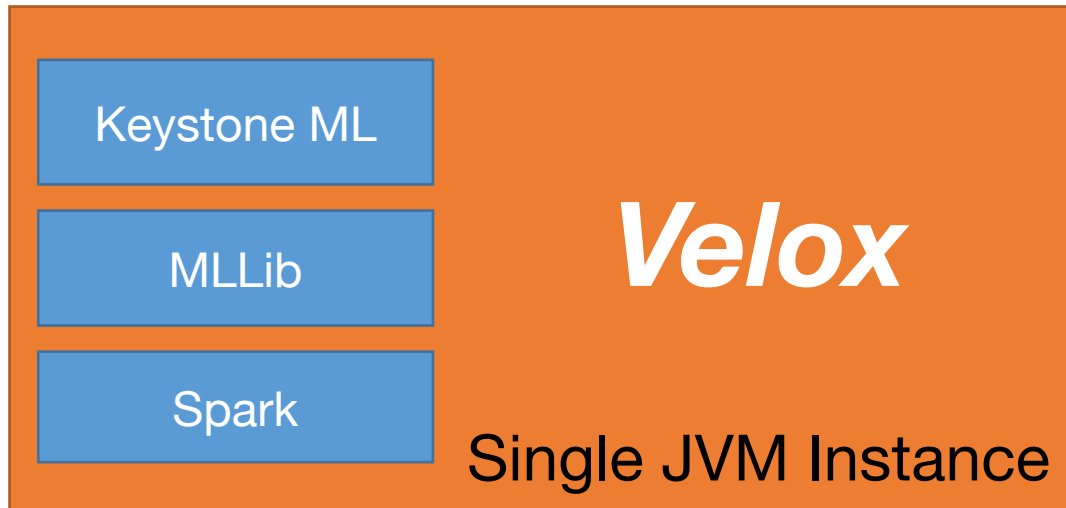
Personal  
Asst.



Robotic  
Control



Machine  
Translation



Caffe

# **VELOX** as a Middle Layer Arch?

Fraud  
Detection



Content  
Rec.



Personal  
Asst.



Robotic  
Control



Machine  
Translation



Generalize *Velox*?

theano

Dato



Create

Caffe



TensorFlow



dmlc

mxnet



KALDI

KeystoneML

# Clipper

## A Low-Latency Online Prediction Serving System

Daniel Crankshaw

Xin Wang

Michael Franklin

**Joseph E. Gonzalez**

Ion Stoica





# Clipper Generalizes Velox Across ML Frameworks

Fraud  
Detection



Content  
Rec.



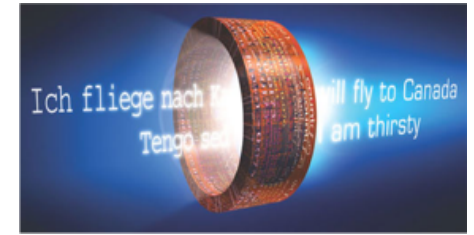
Personal  
Asst.



Robotic  
Control



Machine  
Translation



## Clipper

theano

Dato



Create

Caffe



TensorFlow



dmlc  
mxnet



KALDI

KeystoneML



# 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 prediction interface*
- **bound latency** and **maximize throughput**
  - through *approximate caching* and *adaptive batching*
- enable *robust* **online learning** and **personalization**
  - through generalized *split-model correction policies*

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



# Clipper Design Goals

Low and **bounded** latency predictions

- interactive applications need reliable latency objectives

Up-to-date and personalized predictions **across models** and **frameworks**

- generalize the split model decomposition

Optimize **throughput** for performance under heavy load

- single query can trigger many predictions

**Simplify** deployment

- serve models using the original code and systems

# Clipper Architecture

Fraud  
Detection



Content  
Rec.



Personal  
Asst.



Robotic  
Control



Machine  
Translation



## Clipper

theano

Dato



Create

Caffe



TensorFlow



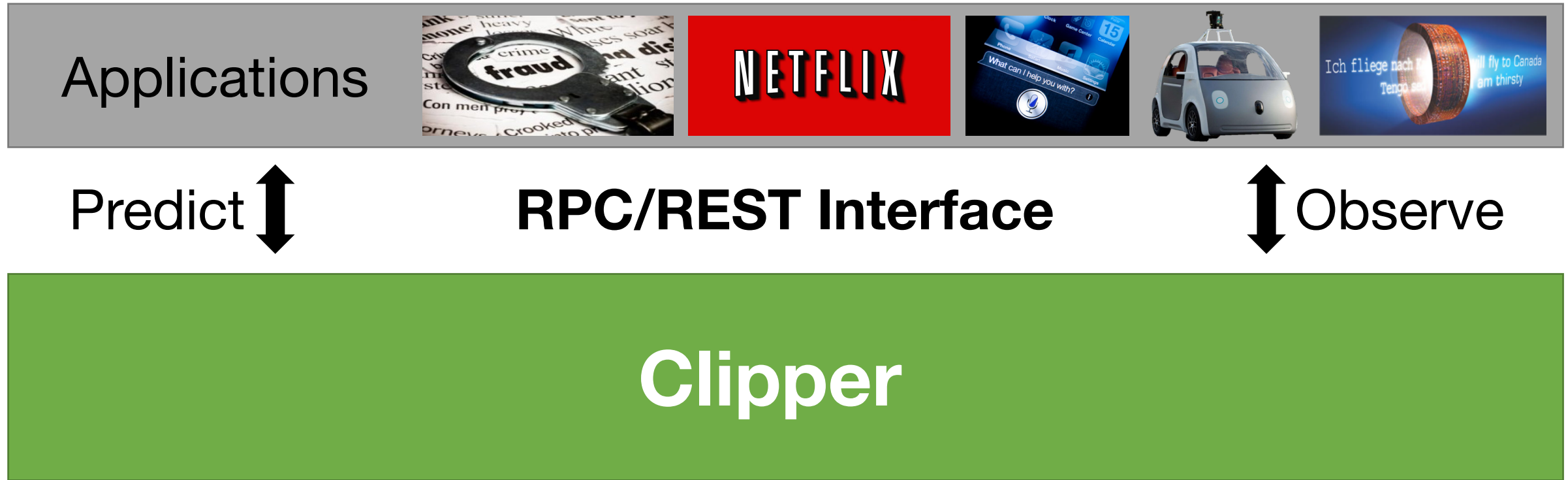
dmlc  
mxnet



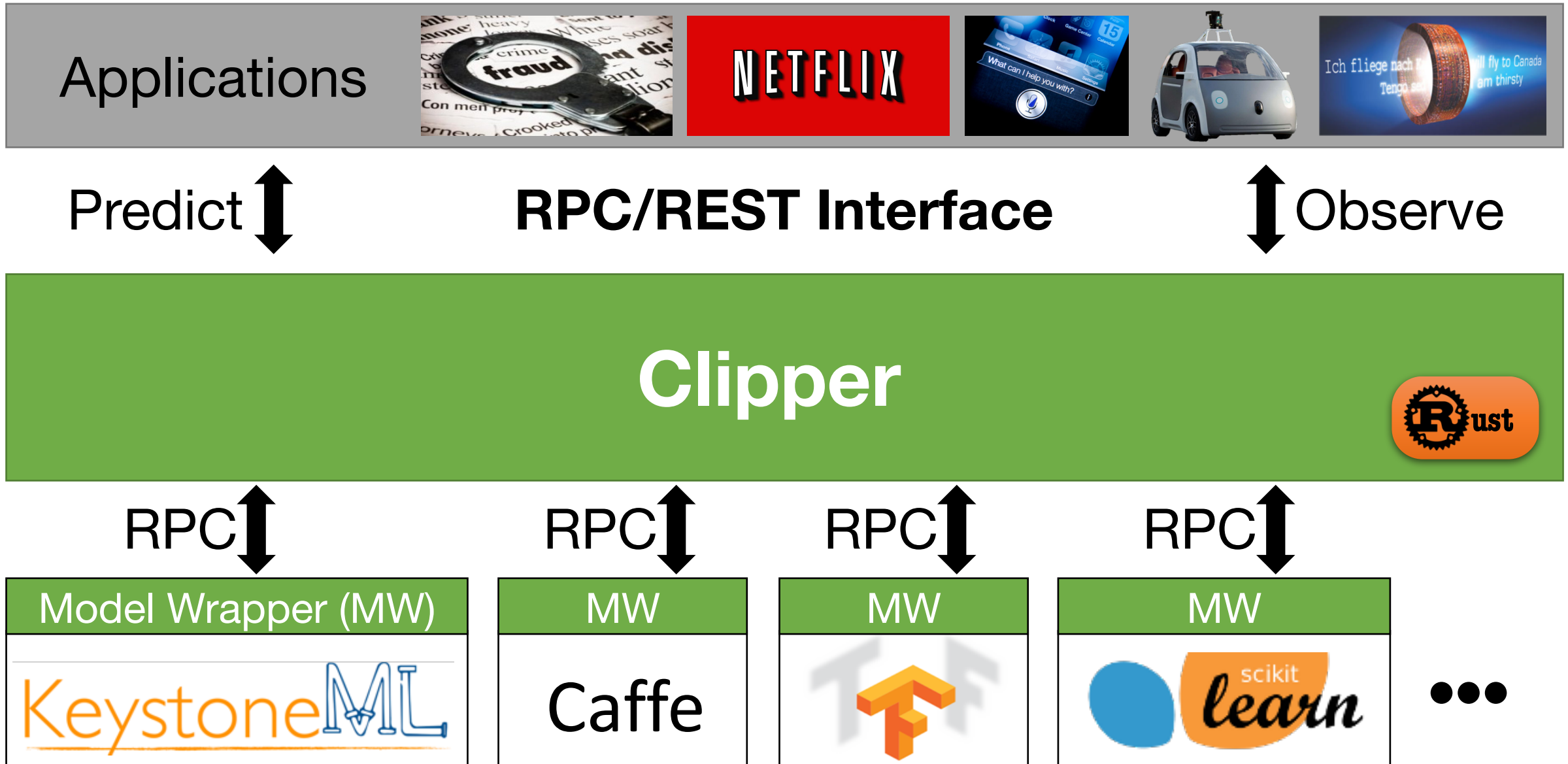
KALDI

KeystoneML

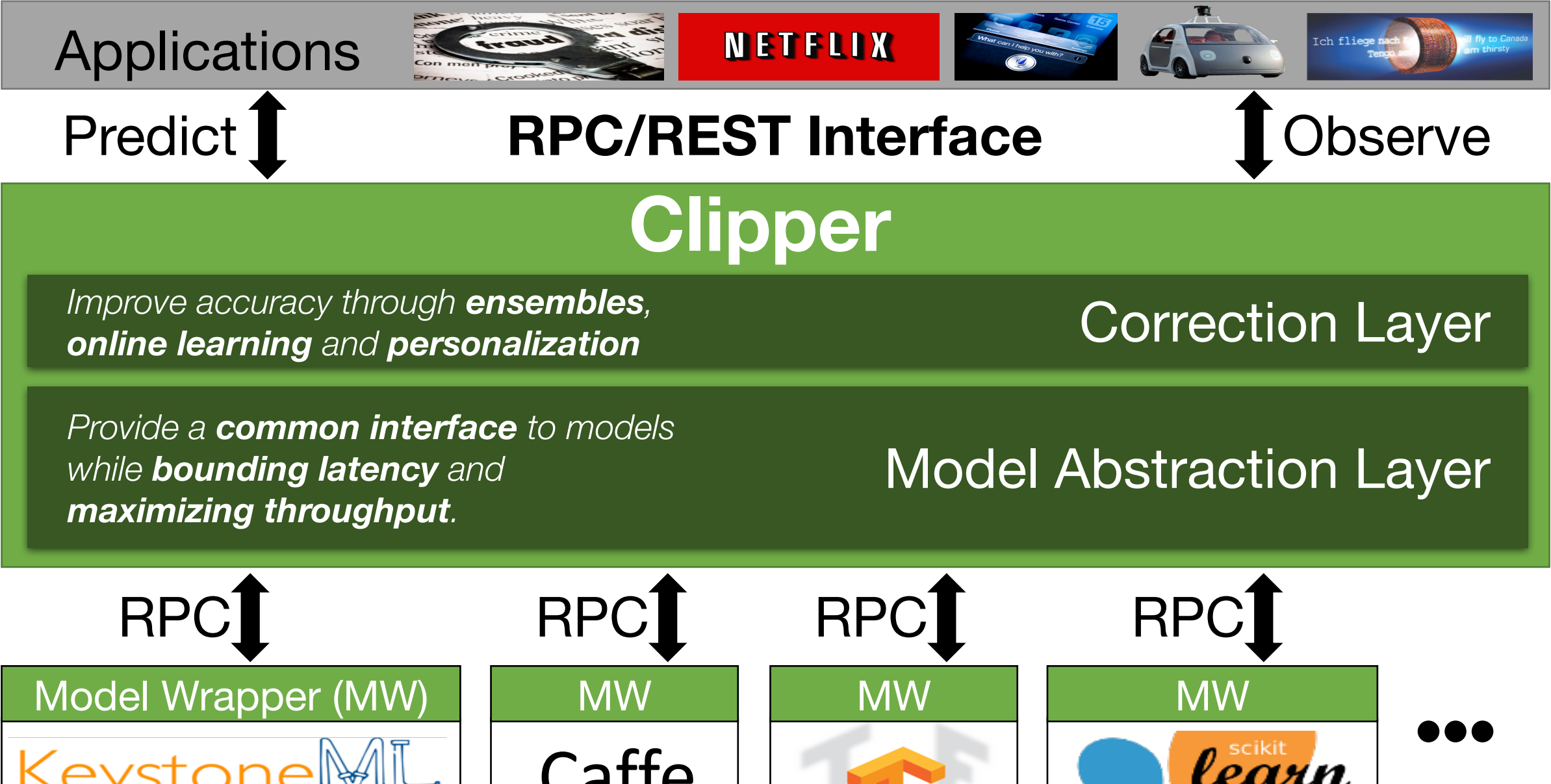
# Clipper Architecture



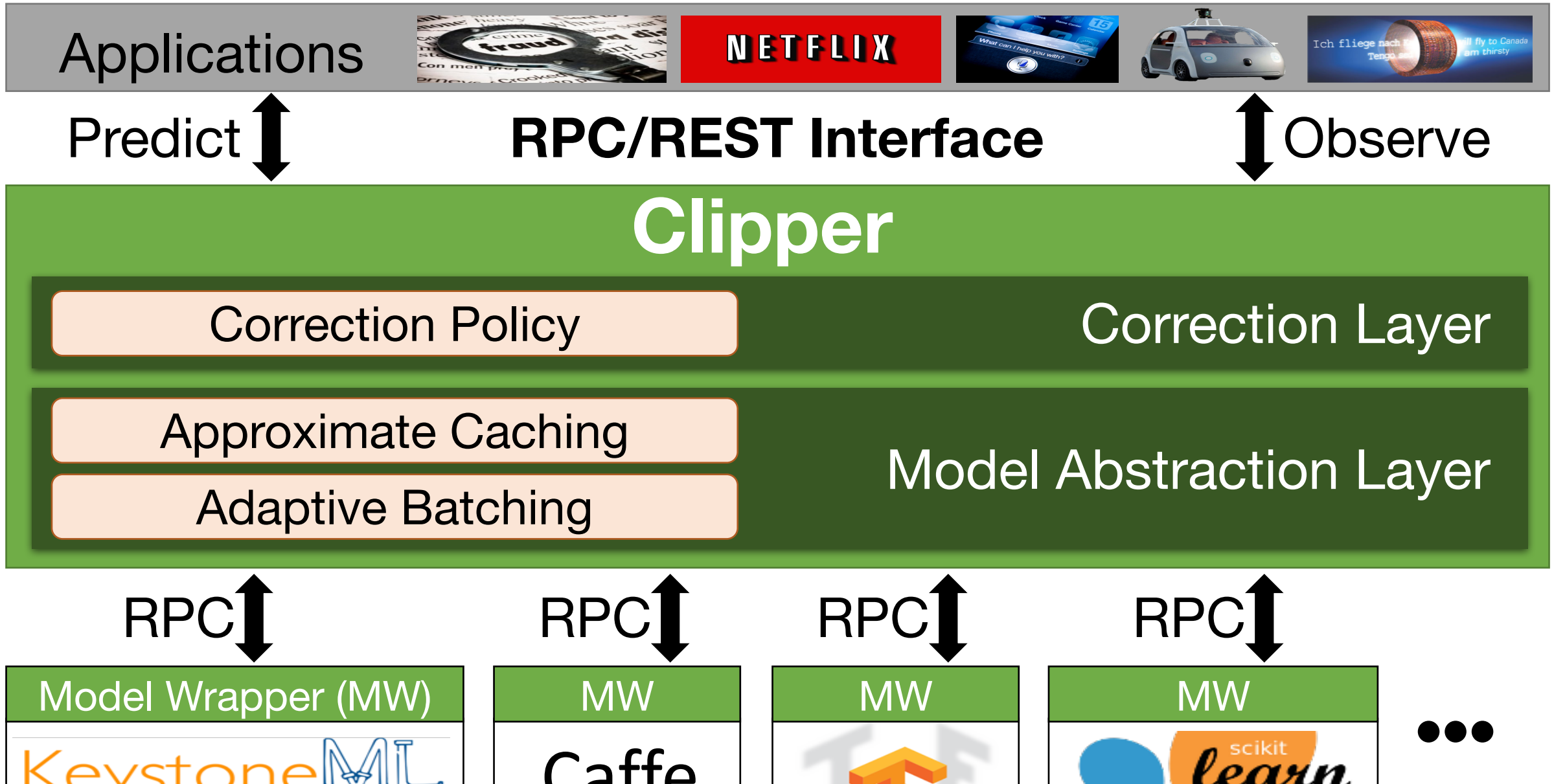
# Clipper Architecture



# Clipper Architecture



# Clipper Architecture



Approximate Caching

Adaptive Batching

Model Abstraction Layer

RPC

Model Wrapper (MW)

KeystoneML

RPC

MW

Caffe

RPC

MW



RPC

MW



...

Provides a unified generic prediction API across **frameworks**

- **Reduce Latency** → Approximate Caching
- **Increase Throughput** → Adaptive Batching
- **Simplify Deployment** → RPC + Model Wrapper

Approximate Caching

Adaptive Batching

Model Abstraction Layer

RPC↕

Model Wrapper (MW)

KeystoneML

RPC↕

MW

Caffe

RPC↕

MW



RPC↕

MW



...



Approximate Caching

Adaptive Batching

Model Abstraction Layer

RPC

RPC

RPC

RPC

Model Wrapper (MW)

MW

MW

MW

KeystoneML

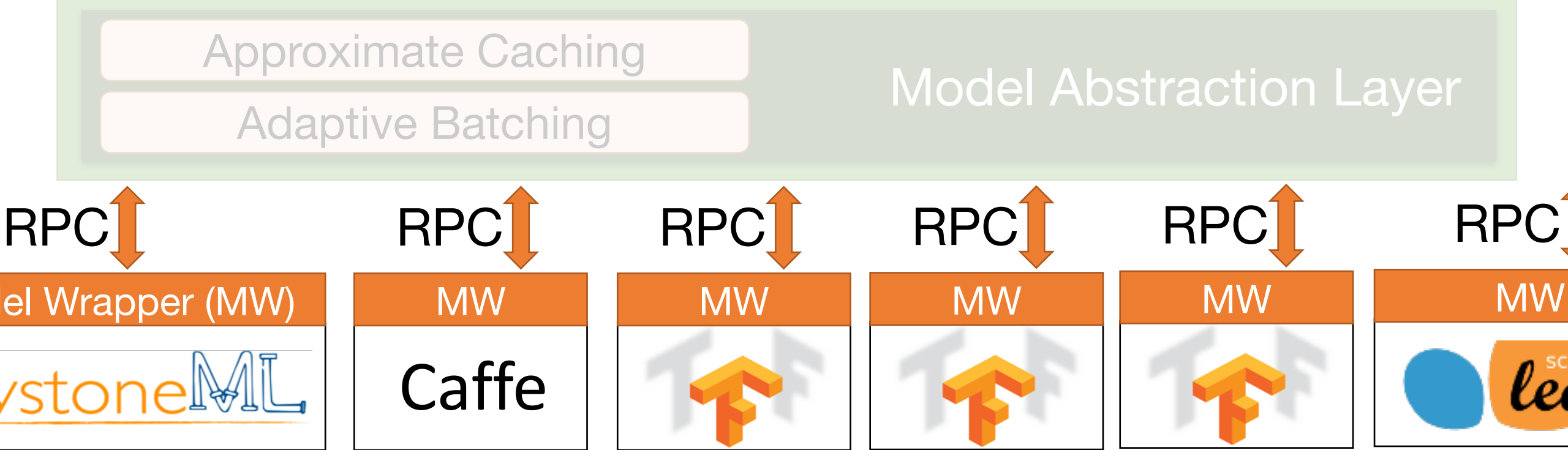
Caffe



...

## Common Interface → Simplifies Deployment:

- Evaluate models using original code & systems
- Models run in separate processes
  - Resource isolation



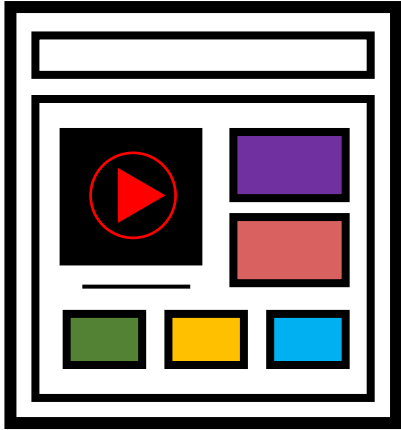
## Common Interface → Simplifies Deployment:

- Evaluate models using original code & systems
- Models run in separate processes
  - Resource isolation
  - Scale-out

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

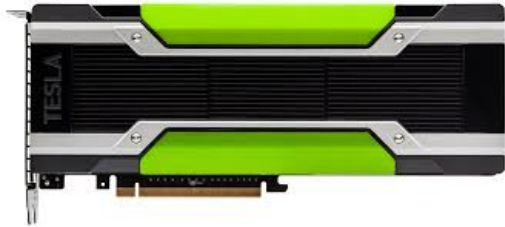
# Adaptive Batching to Improve Throughput

- Why batching helps:



A single page load may generate many queries

Hardware Acceleration



Helps amortize system overhead

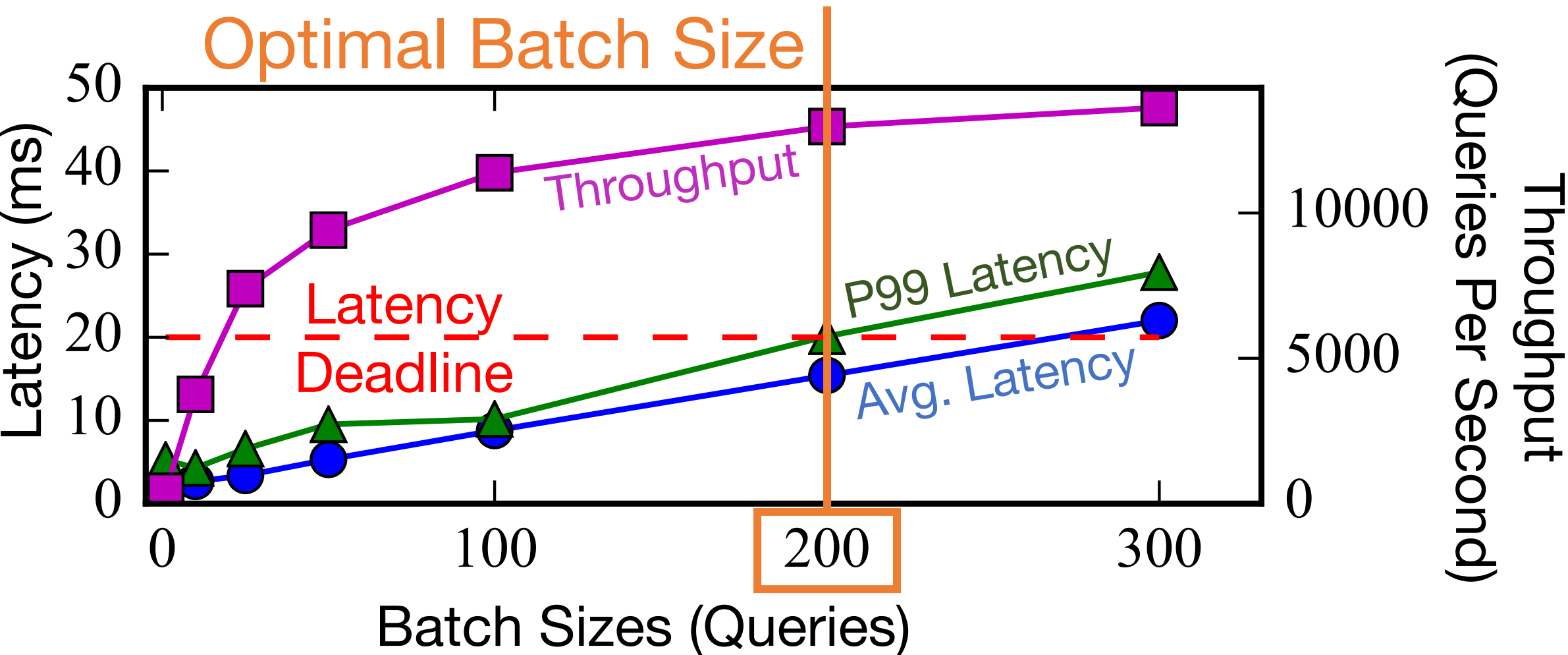
- Optimal batch depends on:
  - hardware configuration
  - model and framework
  - system load

## Clipper Solution:

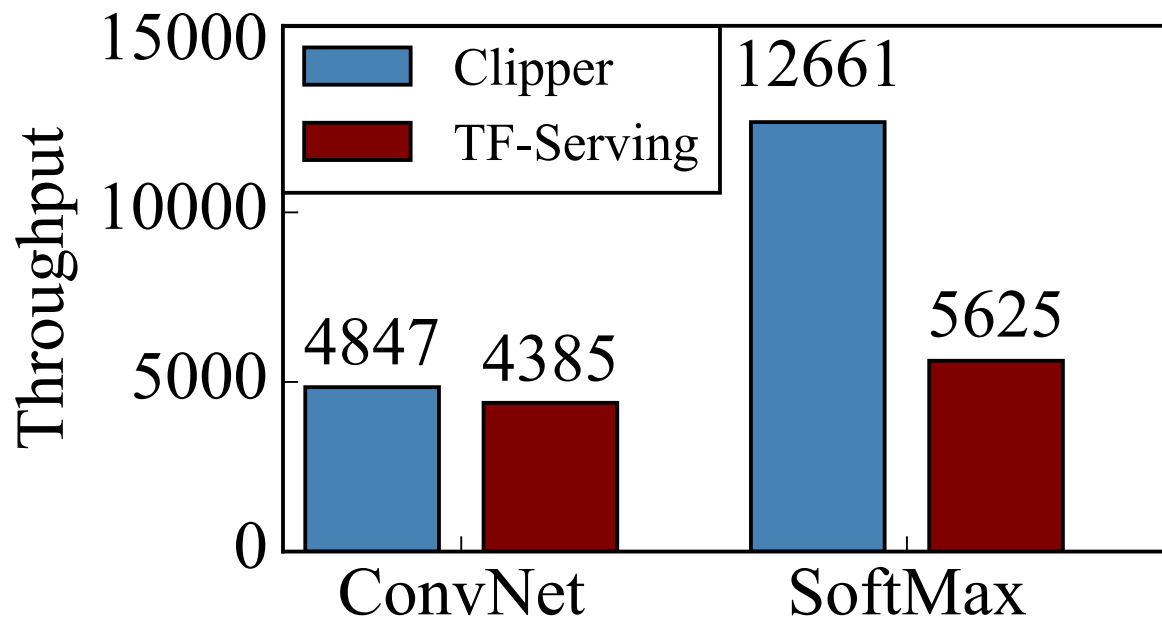
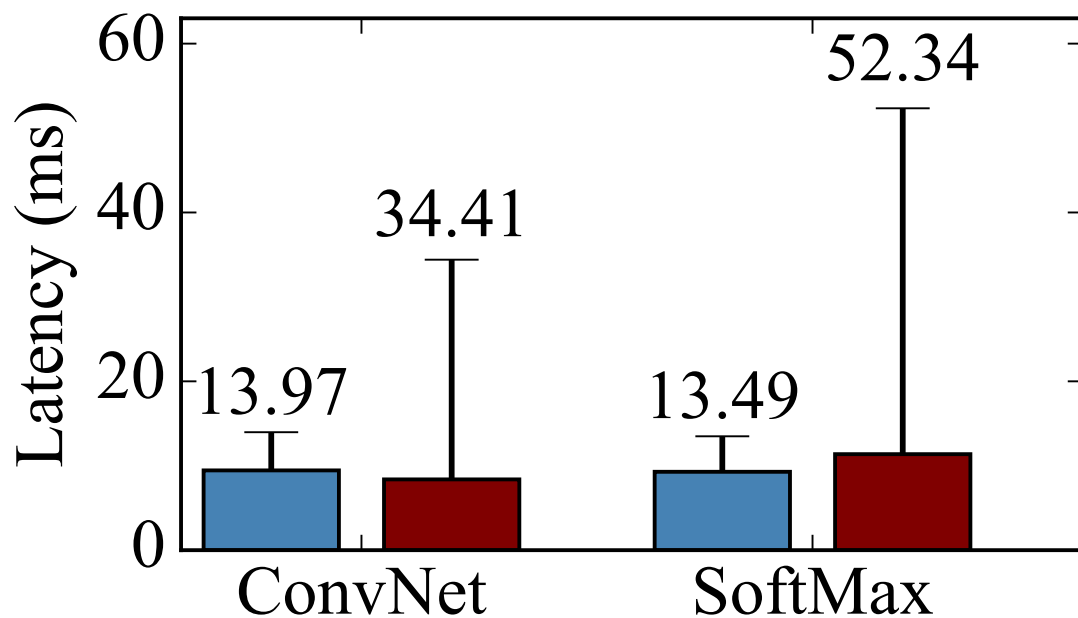
be as **slow** as **allowed...**

- Inc. batch size *until the latency objective is exceeded* (**Additive Increase**)
- If latency exceeds SLO cut batch size by a fraction (**Multiplicative Decrease**)

# Tensor Flow Conv. Net (GPU)



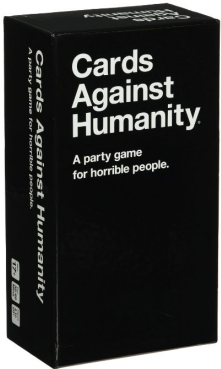
# Comparison to TensorFlow Serving



**Takeaway:** *Clipper is able to **match the average latency** of TensorFlow Serving while reducing **tail latency (2x)** and **improving throughput (2x)***

# Approximate Caching to Reduce Latency

- Opportunity for caching



Popular items may be evaluated frequently

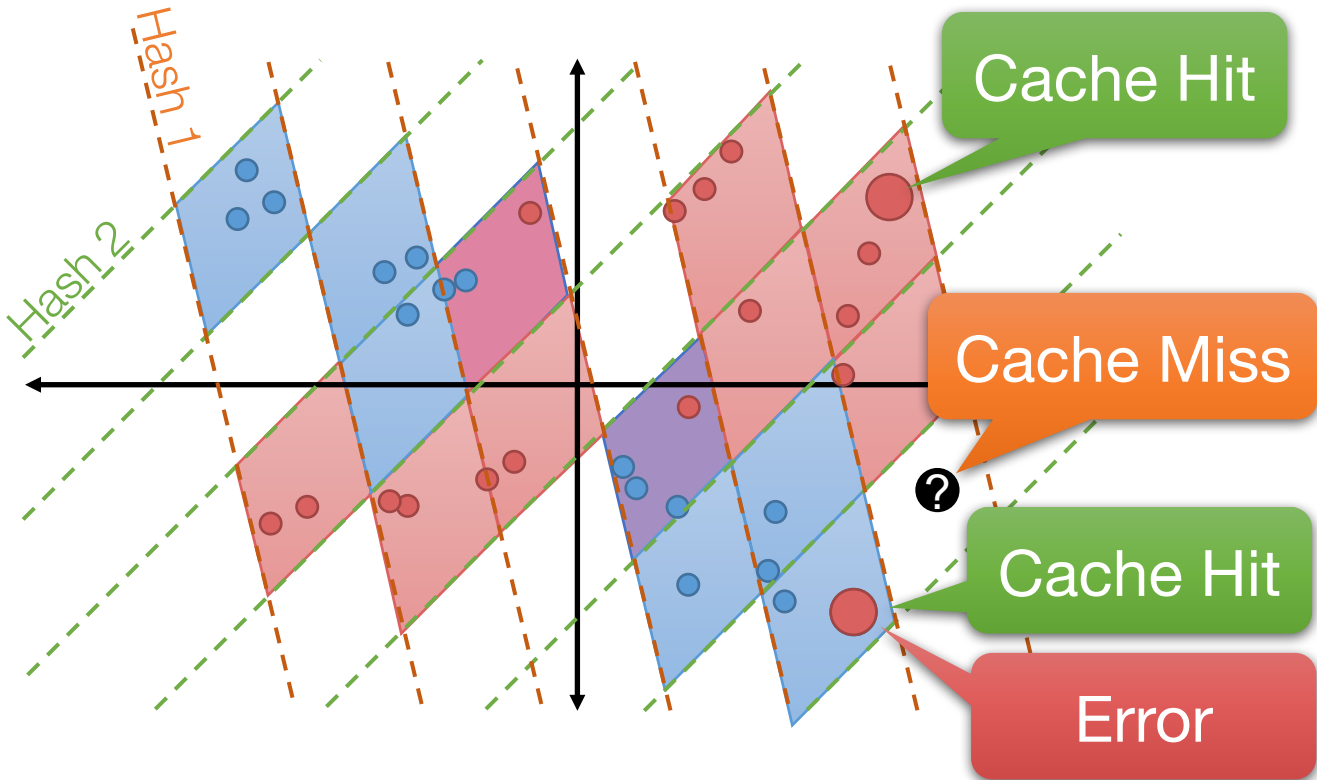
- Need for **approximation**



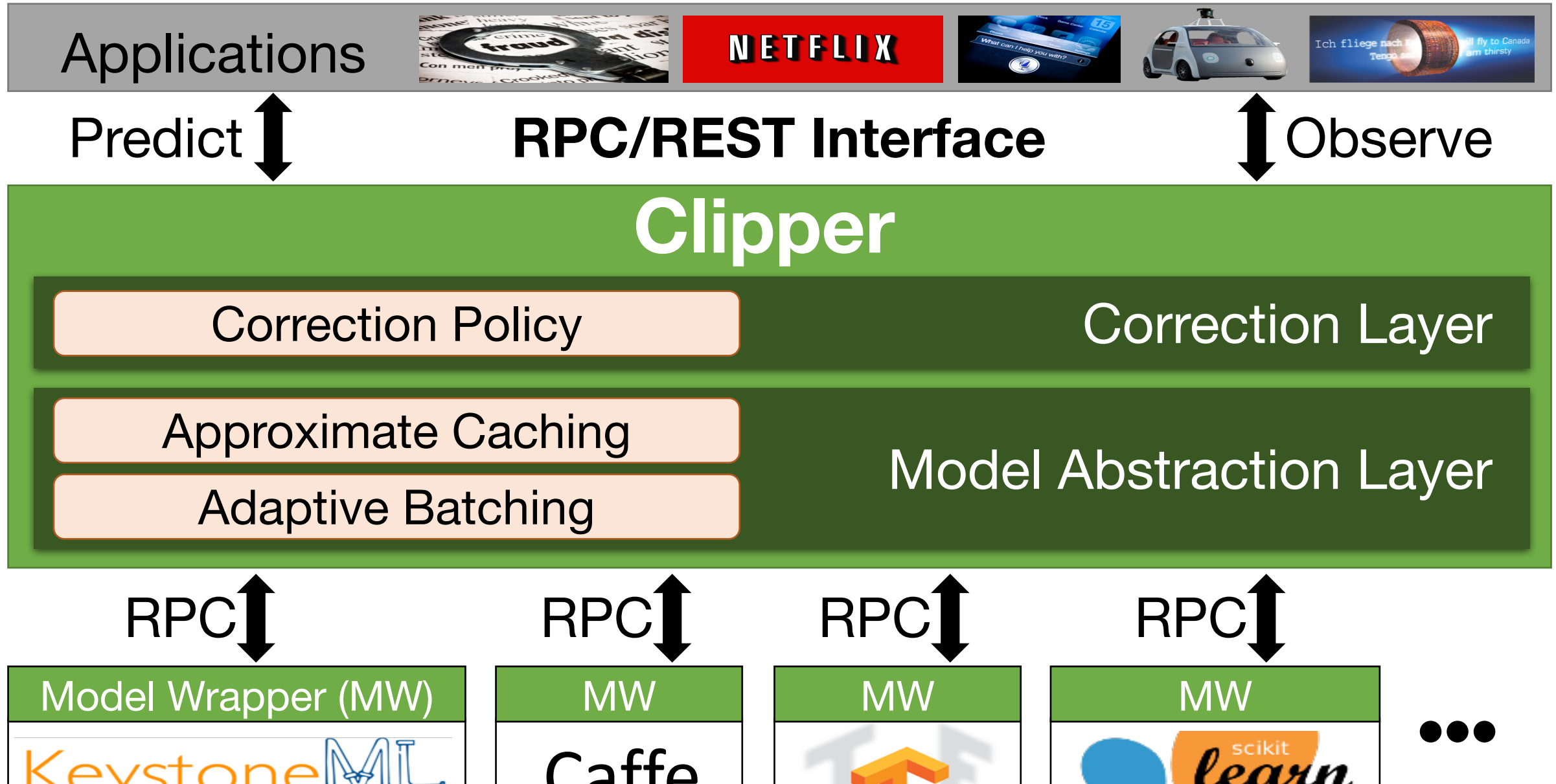
High Dimensional and continuous valued queries have low cache hit rate.

## Clipper Solution: **Approximate Caching**

apply *locality sensitive hash functions*



# Clipper Architecture



## Goal:

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

Generalize the **split-model** insight from Velox to achieve:

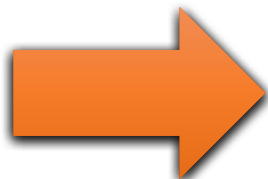
- **robust predictions** by combining multiple models & frameworks
- **online learning** and **personalization** by correcting and personalizing **predictions** in response to feedback



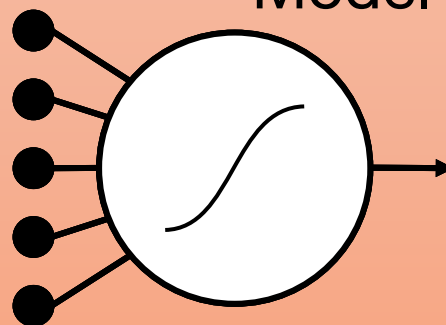
# Learning

# Inference

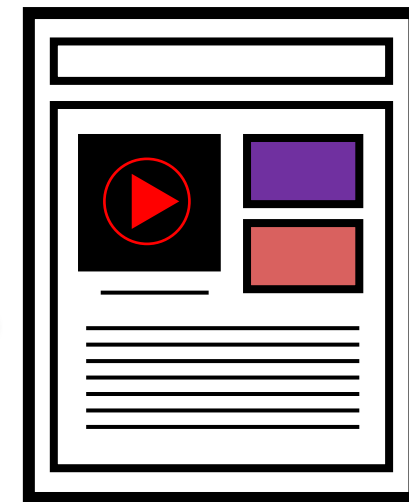
## Velox



**Slow** Changing Model



**Fast** Changing User Model



# Application



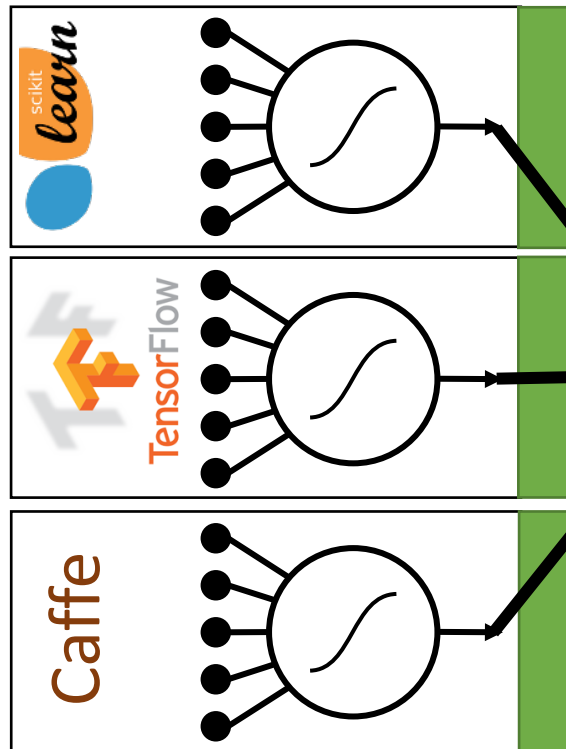
Feedback

Slow



# Learning

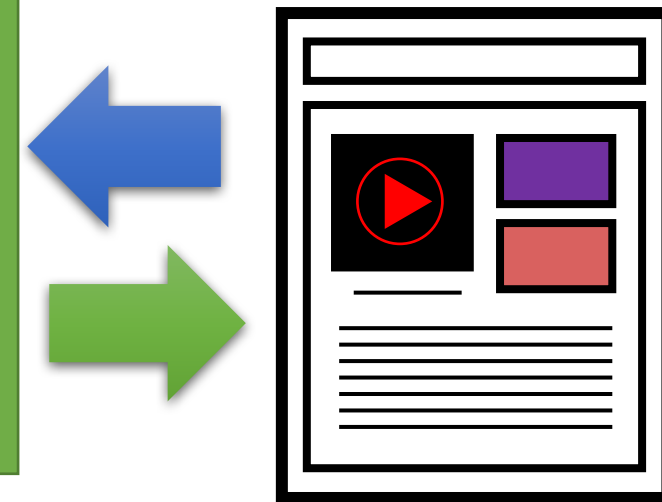
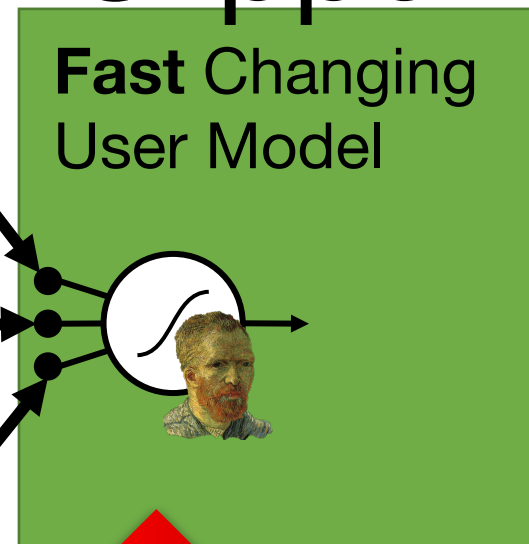
Slow Changing Model



# Inference

## Clipper

Fast Changing User Model



Application

Feedback

Slow

Fast Feedback

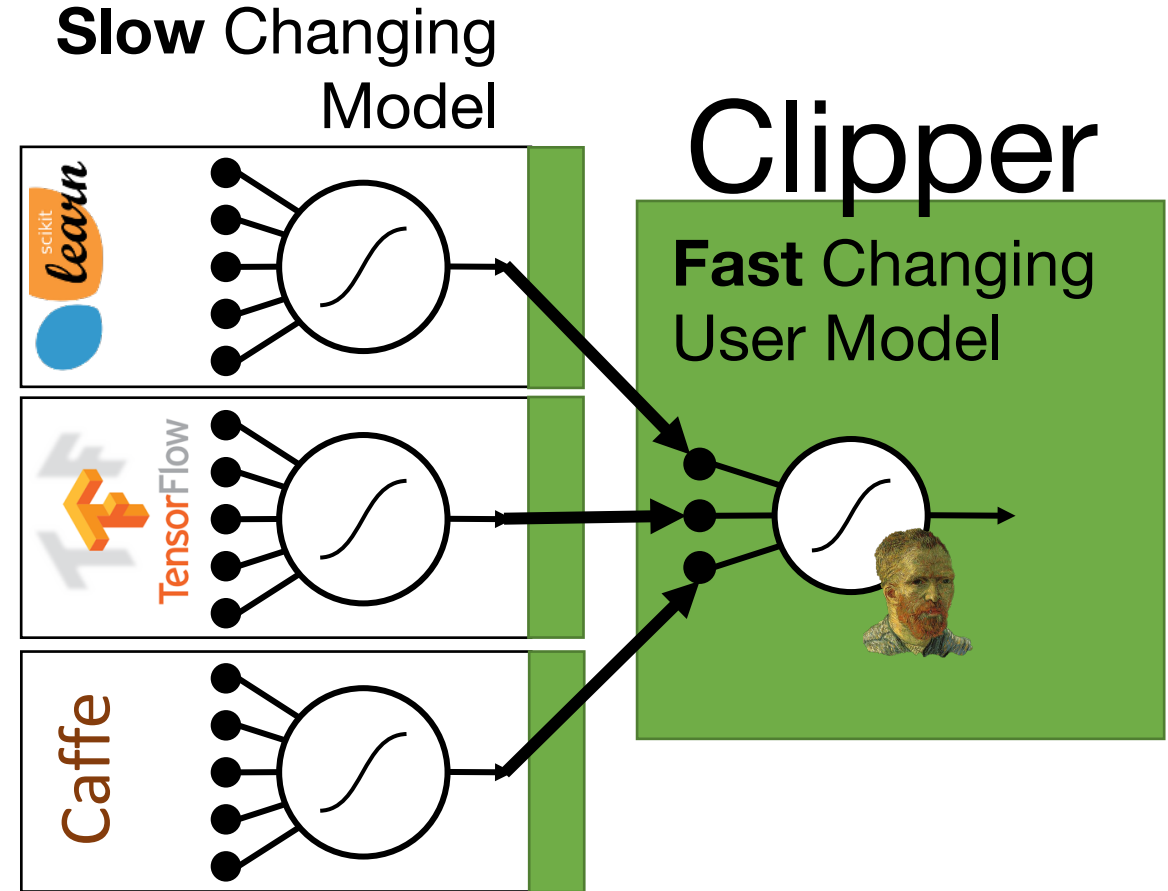


Big Data

# Correction Policy

Improves prediction **accuracy** by:

- Incorporating real-time **feedback**
- Managing **personalization**
- **Combine** models & **frameworks**
  - enables frameworks to **compete**



# Improved 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 state-of-the-art models

# Improved Prediction Accuracy

| System     | Model        | Relative Error | # Errors |
|------------|--------------|----------------|----------|
| Caffe      | ResNet       | 9.02%          | 6525     |
| Caffe      | Inception v3 | 6.18%          | 5760     |
| Caffe      | Ensemble     | 5.86%          | 4512     |
| TensorFlow | Inception v3 | 6.18%          | 3088     |
| Clipper    | Ensemble     | 5.86%          | 2930     |

5.2% relative improvement  
in prediction accuracy!

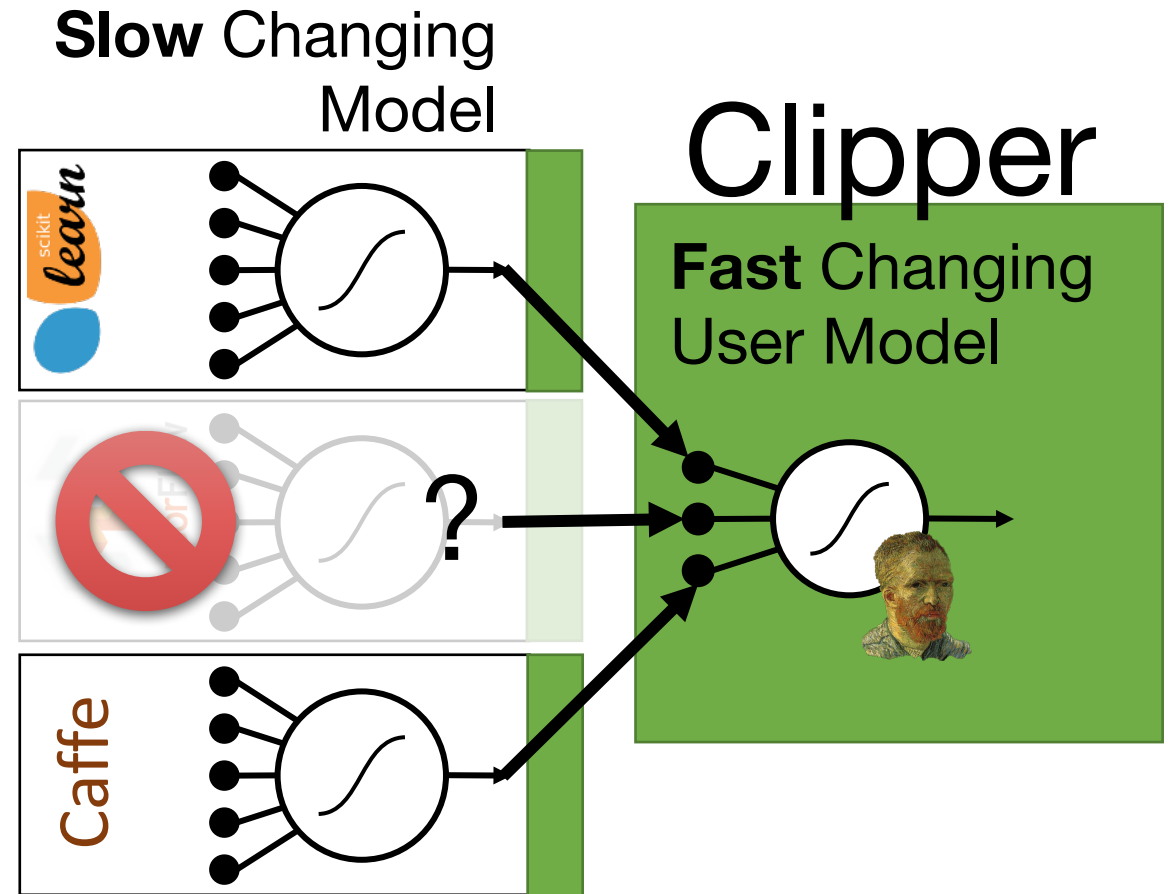
# Cost of Ensembles

## Increased Load

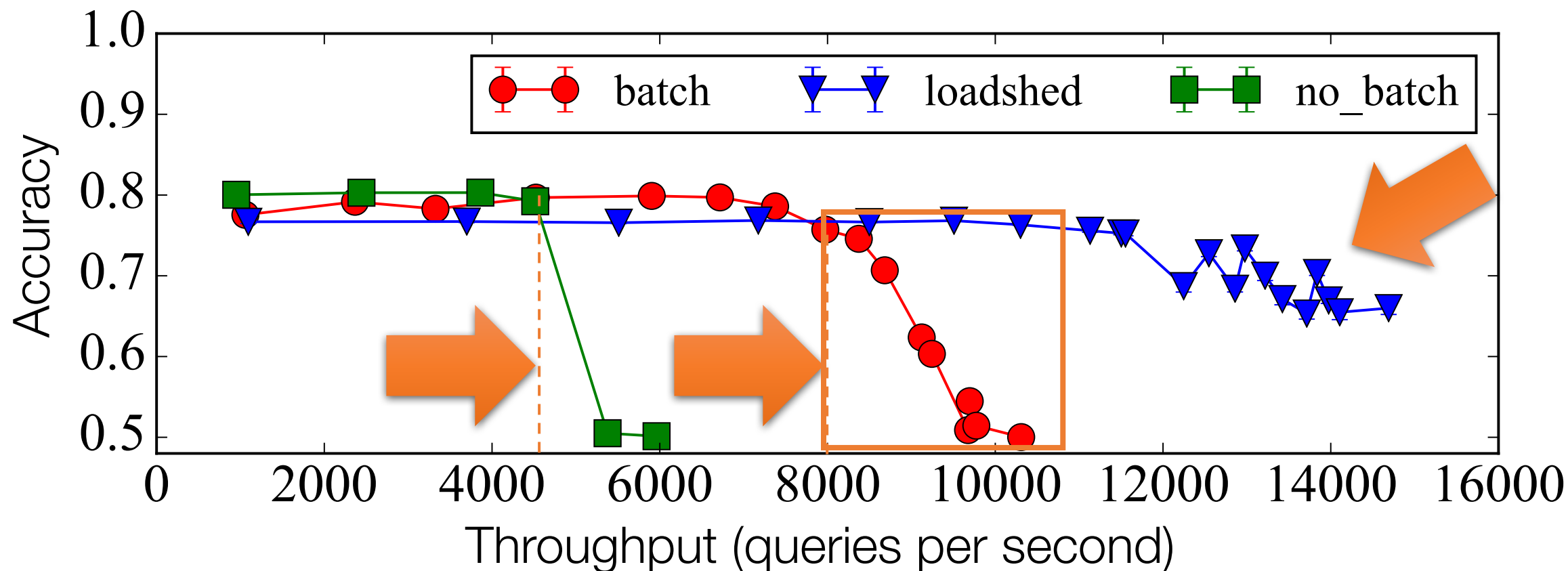
- *Solutions:*
  - **Caching** and **Batching**
  - **Load-shedding** correction policy can prioritize frameworks

## Stragglers

- e.g., framework fails to meet SLO
- *Solution:* **Anytime** predictions
  - Correction policy must render predictions with missing inputs
  - e.g., built-in correction policies **substitute expected value**



# Evaluation of Throughput Under Heavy Load



**Takeaway:** Clipper is able to **gracefully degrade accuracy** to maintain availability under heavy load.

# Conclusion

Clipper sits between applications and ML frameworks to

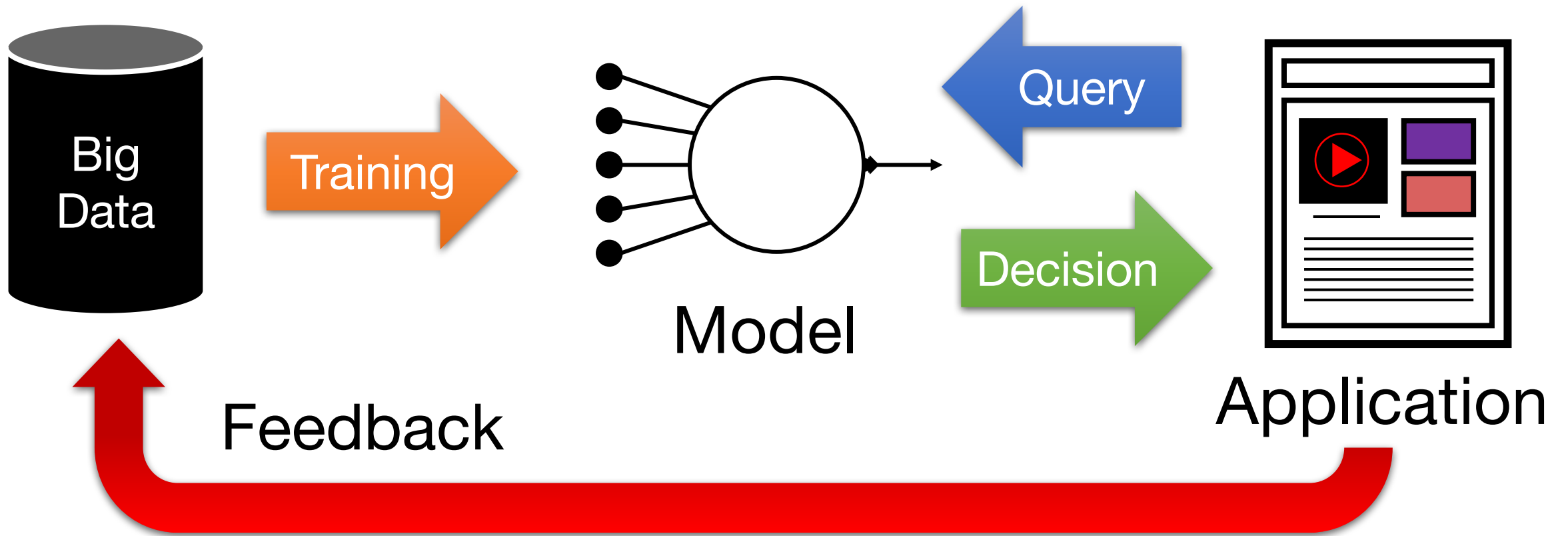


Clipper

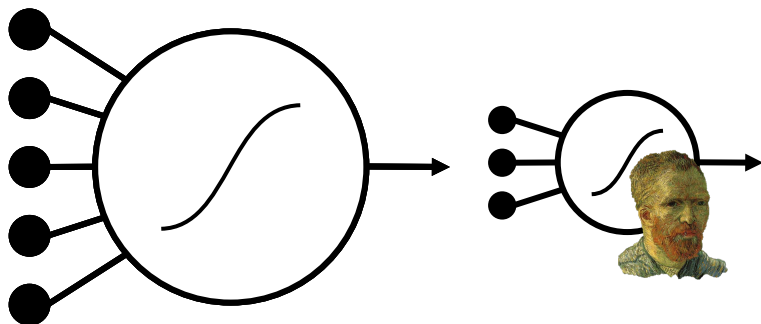


- to **simplifying deployment**
  - **bound latency** and **increase throughput**
  - and enable **real-time learning** and **personalization**
- across **machine learning frameworks**





 **VELOX**



# Ongoing & Future Research Directions

- Serving and updating RL models
- Bandit techniques in correction policies
- Splitting inference across the cloud and the client to reduce latency and bandwidth requirements
- Secure model evaluation on the client (model DRM)

# Coarsening + Anytime Predictions

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

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

