# Clipper
## A Low-Latency Online Prediction Serving System

**Daniel Crankshaw**

Xin Wang, Giulio Zhou,
Michael Franklin, Joseph Gonzalez, Ion Stoica

**NSDI 2017**

*March 29, 2017*

# Learning

## TensorFlow: A system for large-scale machine learning

Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey D...
Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isa...
Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Be...
Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xia...

Google Brain

**Abstract**

TensorFlow is a machine learning system that operates at ...

datasets, and moving t...
based TensorFlow on a ...
first-generation system, ...

---

## GraphLab: A New Framework For Parallel Machine Learning

**Yucheng Low**
Carnegie Mellon University
ylow@cs.cmu.edu

**Joseph Gonzalez**
Carnegie Mellon University
jegonzal@cs.cmu.edu

**Aapo Kyrola**
Carnegie Mellon University
akyrola@cs.cmu.edu

**Danny Bickson**
Carnegie Mellon University

**Carlos Guestrin**
Carnegie Mellon University

**Joseph M. Hellerstein**
UC Berkeley

---

## Spark: Cluster Computing with Working Sets

Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, Ion Stoica
*University of California, Berkeley*

MapReduce/Dryad job, each job must reload the ...
from disk, incurring a significant performance pen...

- **Interactive analytics**: Hadoop is often used to ...
  ad-hoc exploratory queries on large datasets, thro...
  SQL interfaces such as Pig [21] and Hive [1]. Ide...
  a user would be able to load a dataset of interest ...
  memory across a number of machines and query i...
  peatedly. However, with Hadoop, each query in...
  significant latency (tens of seconds) because it run...
  a separate MapReduce job and reads data from dis...

This paper presents a new cluster computing fra...
which supports applicati...
viding similar scal...
MapReduce.

...in Spark is that ...
, which represents ...
oned across a set ...
tion is lost. User...
ory across machi...
e-like *parallel ope*...
through a notion ...
lost, the RDD ha...
as derived from ot...
t partition. Altho...
emory abstraction ...
expressivity on th...
ty on the other ha...
for a variety of ap...
d in Scala [5], a...
language for the ...
programming inte...
ddition, Spark ca...
d version of the S...
to define RDDs ...
se them in parallel...
Spark is the first ...
pose programmin...
cess large dataset...
entation of Spark ...
with the system is e...
tperform Hadoop ...
orkloads and can ...
dataset with sub-se...
ed as follows. Sec...

---

## GraphX: Graph Processing in a Distributed Dataflow Framework

Joseph E. Gonzalez[*], Reynold S. Xin[*†], Ankur Dave[*], Daniel Crankshaw[*]
Michael J. Franklin[*], Ion Stoica[*†]
[*]*UC Berkeley AMPLab*   [†]*Databricks*

**Abstract**

---

## Project Adam: Building an Efficient and Scalable Deep Learning Training System

Trishul Chilimbi   Yutaka Suzue   Johnson Apacible   Karthik Kalyanaraman
*Microsoft Research*

**ABSTRACT**

Large deep neural ...
demonstrated state-o...
recognition tasks. ...
extremely time cons...
amount of compute ...
implementation of a...
comprised of commo...
models that exhibits ...
and task accuracy o...
achieves high efficie...
system co-design ...
workload computatio...
asynchrony throug...
performance and sho...
accuracy of trained ...
more efficient and ...
thought possible and ...
a large 2 billion ...
accuracy in compara...
category image class...
previously held the r...
show that task accur...
Our results provid...
distributed systems-...
using current trainin...

**1. INTRODUC...**

Traditional statistical...
table of data and a ...
table correspond to ...
columns correspond ...
underlying data set. ...
algorithms can be a p...

---

## Caffe: Convolutional Architecture for Fast Feature Embedding[*]

Yangqing Jia[*], Evan Shelhamer[*], Jeff Donahue, Sergey Karayev,
Jonathan Long, Ross Girshick, Sergio Guadarrama, Trevor Darrell
SUBMITTED to ACM MULTIMEDIA 2014 OPEN SOURCE SOFTWARE COMPETITION
UC Berkeley EECS, Berkeley, CA 94702
{jiayq,shelhamer,jdonahue,sergeyk,jonlong,rbg,sguada,trevor}@eecs.berkeley.edu

**ABSTRACT**

Caffe provides multimedia scientists and practitioners with
a clean and modifiable framework for state-of-the-art deep
learning algorithms and a collection of reference models.
The framework is a BSD-licensed C++ library with Python
and MATLAB bindings for training and deploying general-
purpose convolutional neural networks and other deep mod-
els efficiently on commodity architectures. Caffe fits indus-
try and internet-scale media needs by CUDA GPU computa-
tion, processing over 40 million images a day on a single K40
or Titan GPU ($\approx$ 2.5 ms per image). By separating model
representation from actual implementation, Caffe allows ex-
perimentation and seamless switching among platforms for
ease of development and deployment from prototyping ma-
chines to cloud environments.

Caffe is maintained and developed by the Berkeley Vi-
sion and Learning Center (BVLC) with the help of an ac-
tive community of contributors on GitHub. It powers on-
going research projects, large-scale industrial applications,

**1. INTRODUCTION**

A key problem in multimedia data analysis is discovery of
effective representations for sensory inputs—images, sound-
waves, haptics, etc. While performance of conventional,
handcrafted features has plateaued in recent years, new de-
velopments in deep compositional architectures have kept
performance levels rising [8]. Deep models have outper-
formed hand-engineered feature representations in many do-
mains, and made learning possible in domains where engi-
neered features were lacking entirely.

We are particularly motivated by large-scale visual recog-
nition, where a specific type of deep architecture has achieved
a commanding lead on the state-of-the-art. These *Con-
volutional Neural Networks*, or CNNs, are discriminatively
trained via back-propagation through layers of convolutional
filters and other operations such as rectification and pooling.
Following the early success of digit classification in the 90's,
these models have recently surpassed all known methods for
large-scale visual recognition, and have been adopted by in-
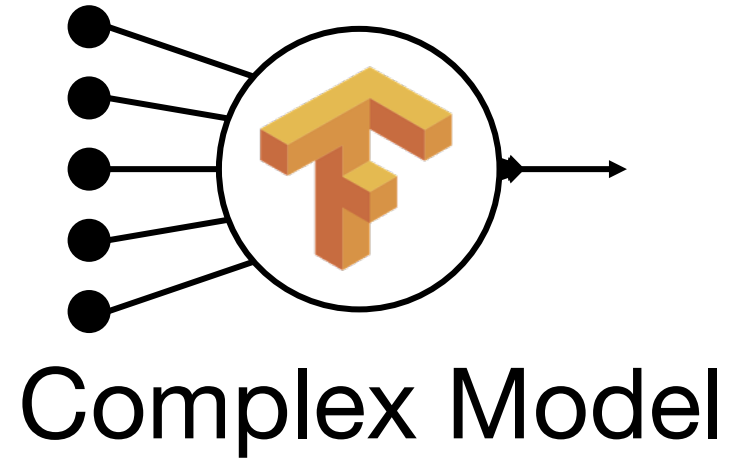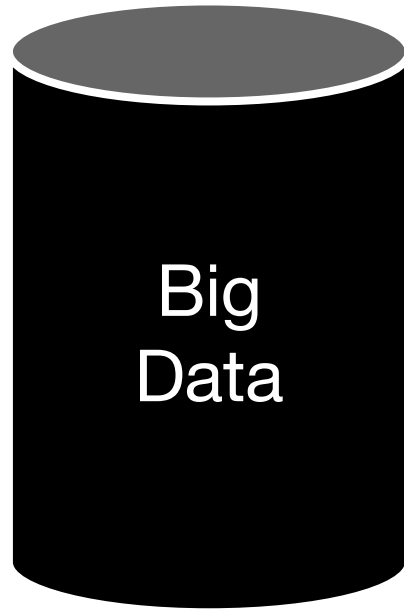
---

## Parameter Server for Distributed Machine Learning

Mu Li[1], Li Zhou[1], Zichao Yang[1], Aaron Li[1], Fei Xia[1],
David G. Andersen[1] and Alexander Smola[1,2]
[1]Carnegie Mellon University
[2]Google Strategic Technologies
{muli, lizhou, zichaoy, aaronli, feixia, dga}@cs.cmu.edu, alex@smola.org

**Abstract**

We propose a parameter server framework to solve distributed machine learning
problems. Both data and workload are distributed into client nodes, while server
nodes maintain globally shared parameters, which are represented as sparse vec-
tors and matrices. The framework manages asynchronous data communications
between clients and servers. Flexible consistency models, elastic scalability and
fault tolerance are supported by this framework. We present algorithms and theo-
retical analysis for challenging nonconvex and nonsmooth problems. To demon-
strate the scalability of the proposed framework, we show experimental results on
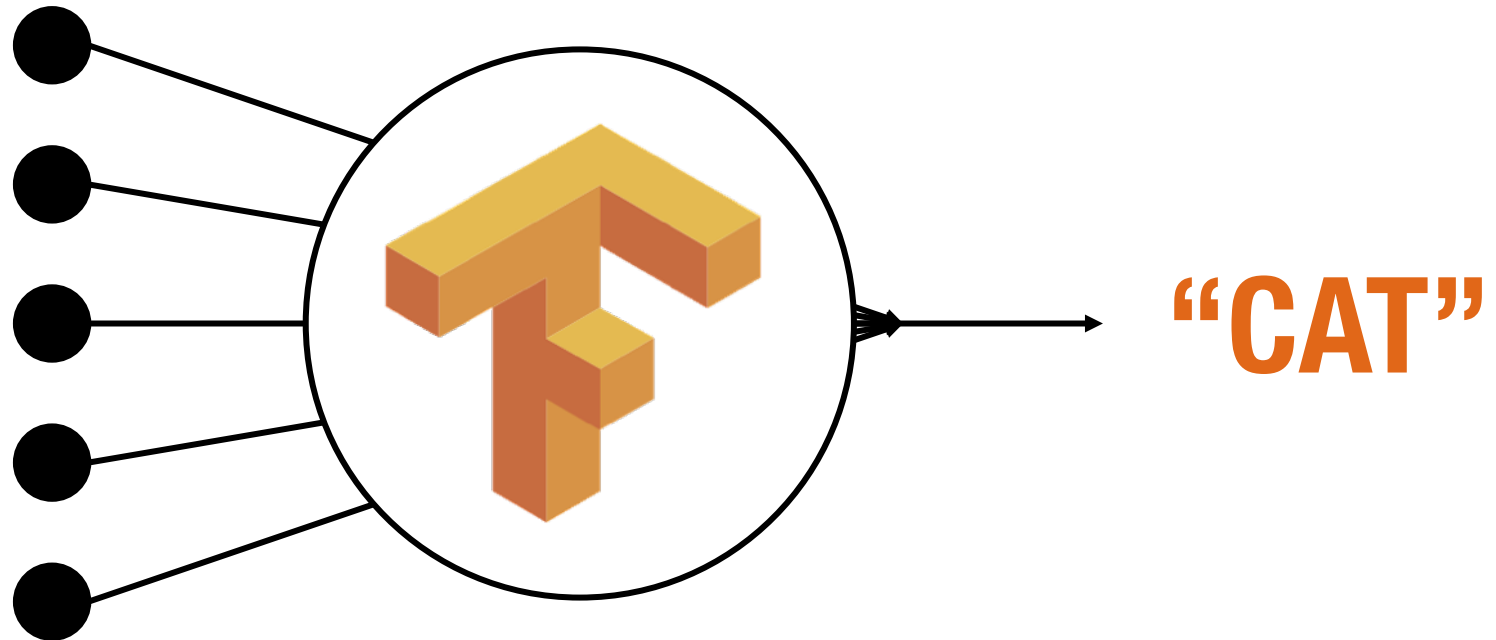real data with billions of parameters.
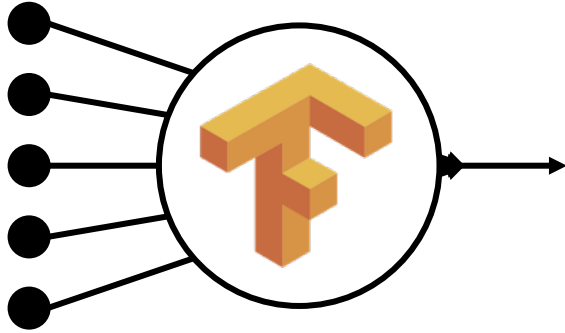
# Learning Produces a Trained Model
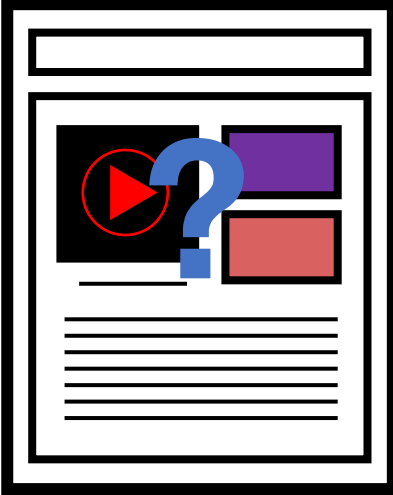
*Query*                                     *Decision*



Model                                     **"CAT"**

**Learning**

**Serving**

Big Data

Training

Model

Query

Decision

Application

# Learning

Big Data

Training

# Serving

Model

Query

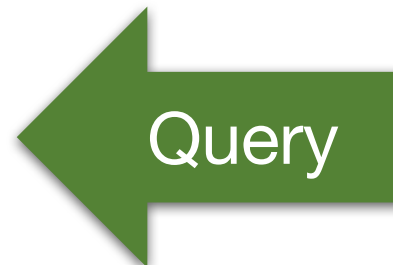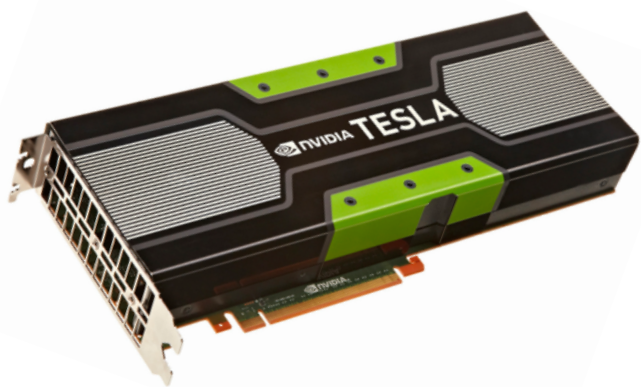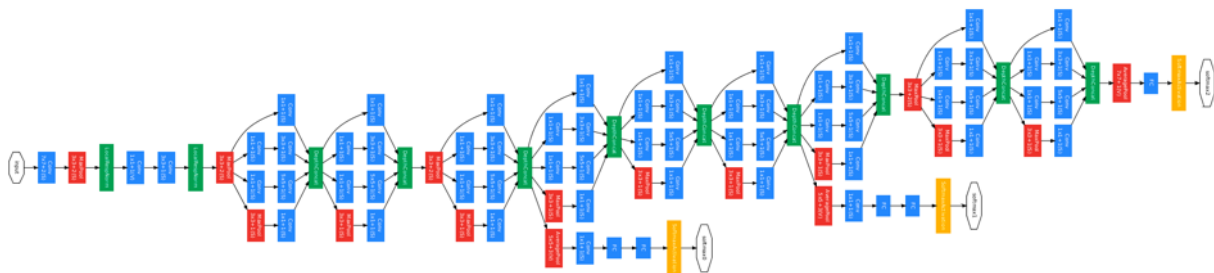Decision

Application

*Prediction-Serving for interactive applications*

*Timescale: ~10s of milliseconds*

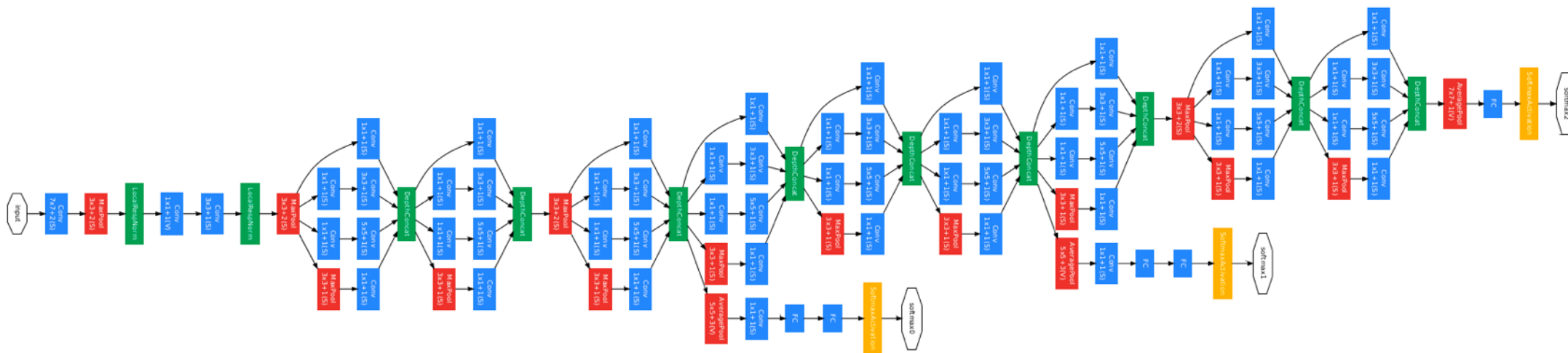# Prediction-Serving Raises New Challenges

# Prediction-Serving Challenges



*Support low-latency, high-throughput serving workloads*



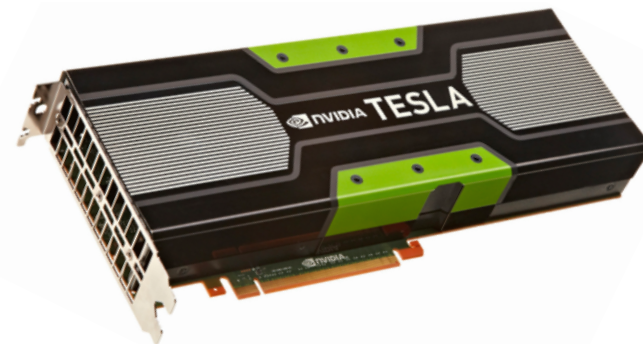*Large and growing ecosystem of ML models and frameworks*

# Support low-latency, high-throughput serving workloads



# Models getting more complex

➤ *10s of GFLOPs [1]*

# Deployed on critical path

➤ *Maintain SLOs under heavy load*

*Using specialized hardware for predictions*

[1] *Deep Residual Learning for Image Recognition.* He et al. CVPR 2015.

# Google Translate

Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi
yonghui,schuster,zhifengc,qvl,mnorouzi@google.com

Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, Jeffrey Dean

## Serving

Google

Translate                                    Turn off instant translation

140 billion words a day[1]

0/5000

82,000 GPUs
running 24/7

**Invented New Hardware!
Tensor Processing Unit
(TPU)**

[1] https://www.nytimes.com/2016/12/14/magazine/the-great-ai-awakening.html
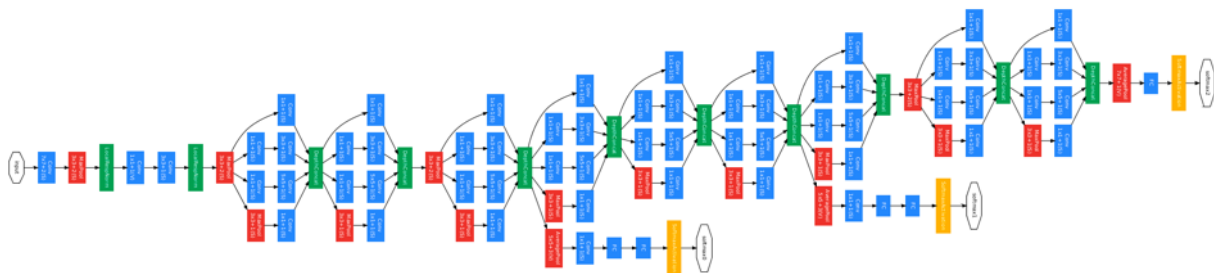
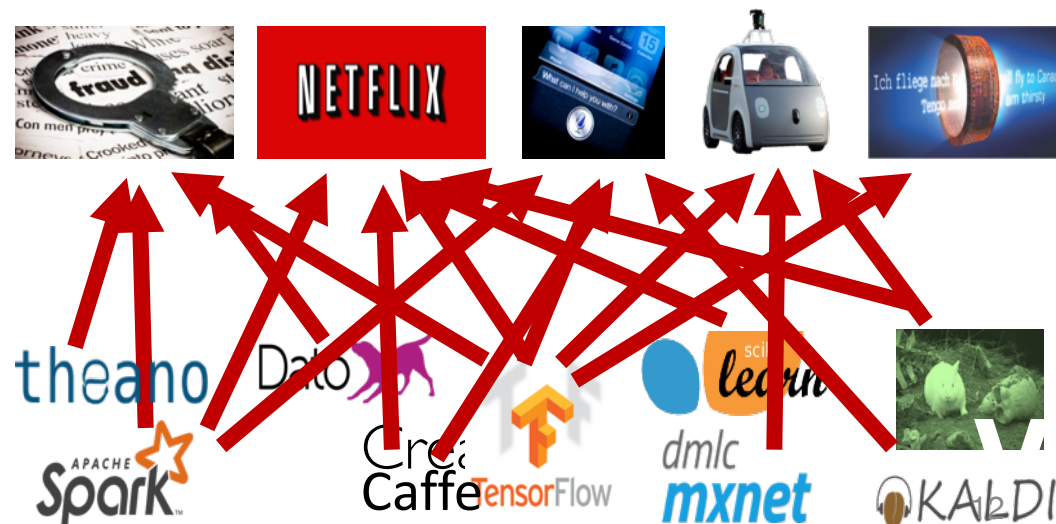# Big Companies Build One-Off Systems

**Problems:**

➢ Expensive to build and maintain

   ➢ Highly specialized and require ML and systems expertise

➢ Tightly-coupled model and application

   ➢ Difficult to change or update model

➢ Only supports single ML framework

# Prediction-Serving Challenges

*Support low-latency, high-throughput serving workloads*

*Large and growing ecosystem of ML models and frameworks*
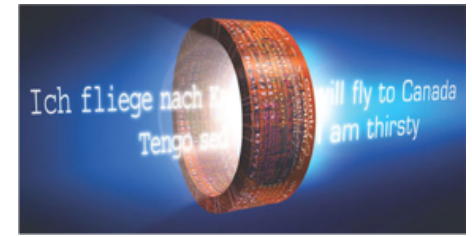
*Large and growing ecosystem of ML models and frameworks*

Fraud Detection

Content Rec.

Personal Asst.

Robotic Control

Machine Translation

theano

Dato Create

Caffe

TensorFlow

scikit learn

dmlc mxnet

VW

KALDI

APACHE Spark

13

Large and growing ecosystem of ML models and frameworks

**Difficult to deploy and brittle to manage**

*Varying physical resource requirements*

*But most companies can't build new serving systems…*

# Use existing systems: Offline Scoring
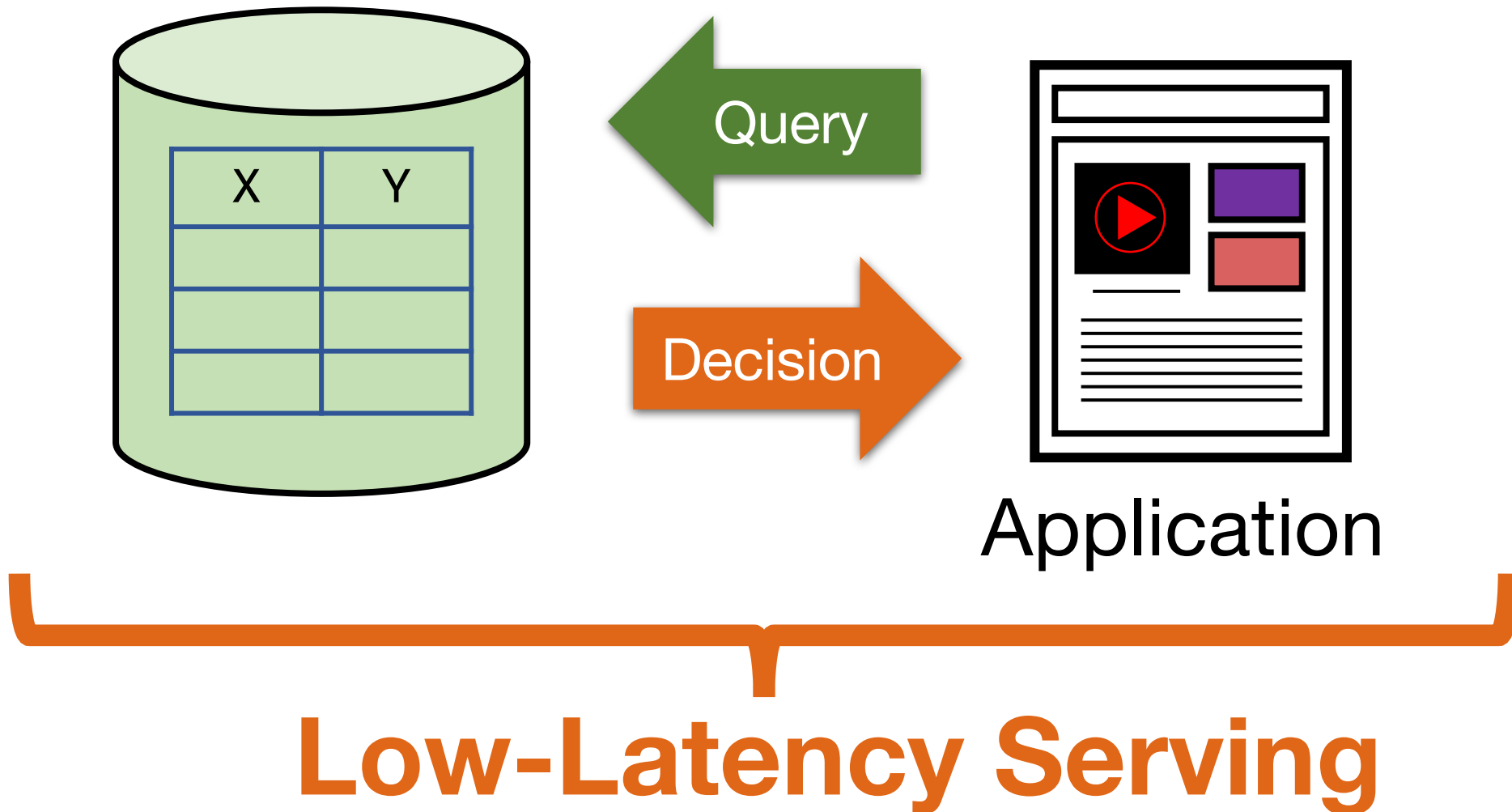
Use existing systems: Offline Scoring

Big Data → Training → Model → Scoring → Datastore (X, Y)

Batch Analytics

# Use existing systems: Offline Scoring

## Look up decision in datastore



Application

**Low-Latency Serving**
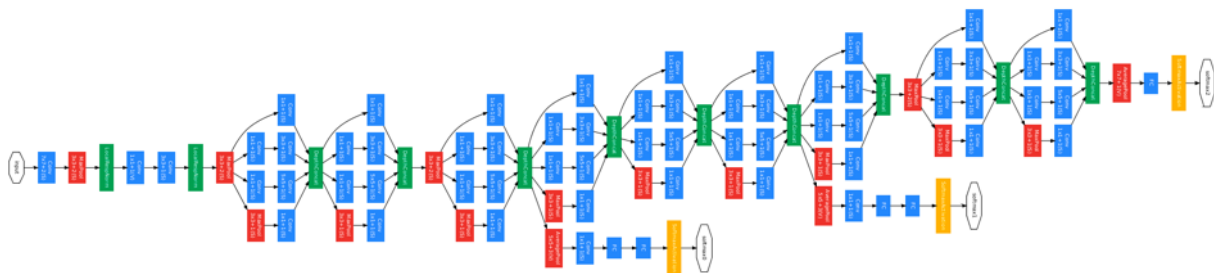
# Use existing systems: Offline Scoring

**Problems:**

➢ Requires full set of queries ahead of time

  ➢ Small and bounded input domain

➢ Wasted computation and space

  ➢ Can render and store unneeded predictions
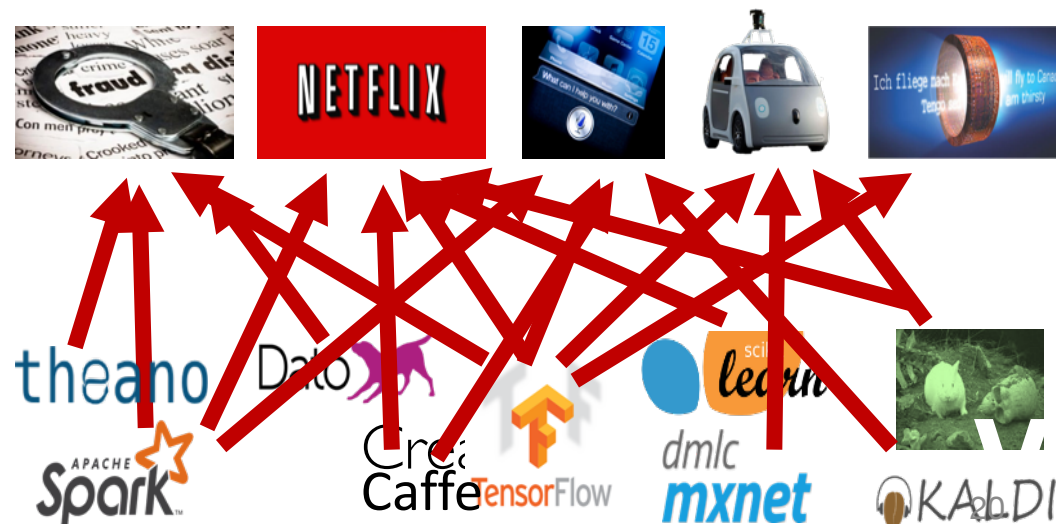
➢ Costly to update

  ➢ Re-run batch job

# Prediction-Serving Challenges



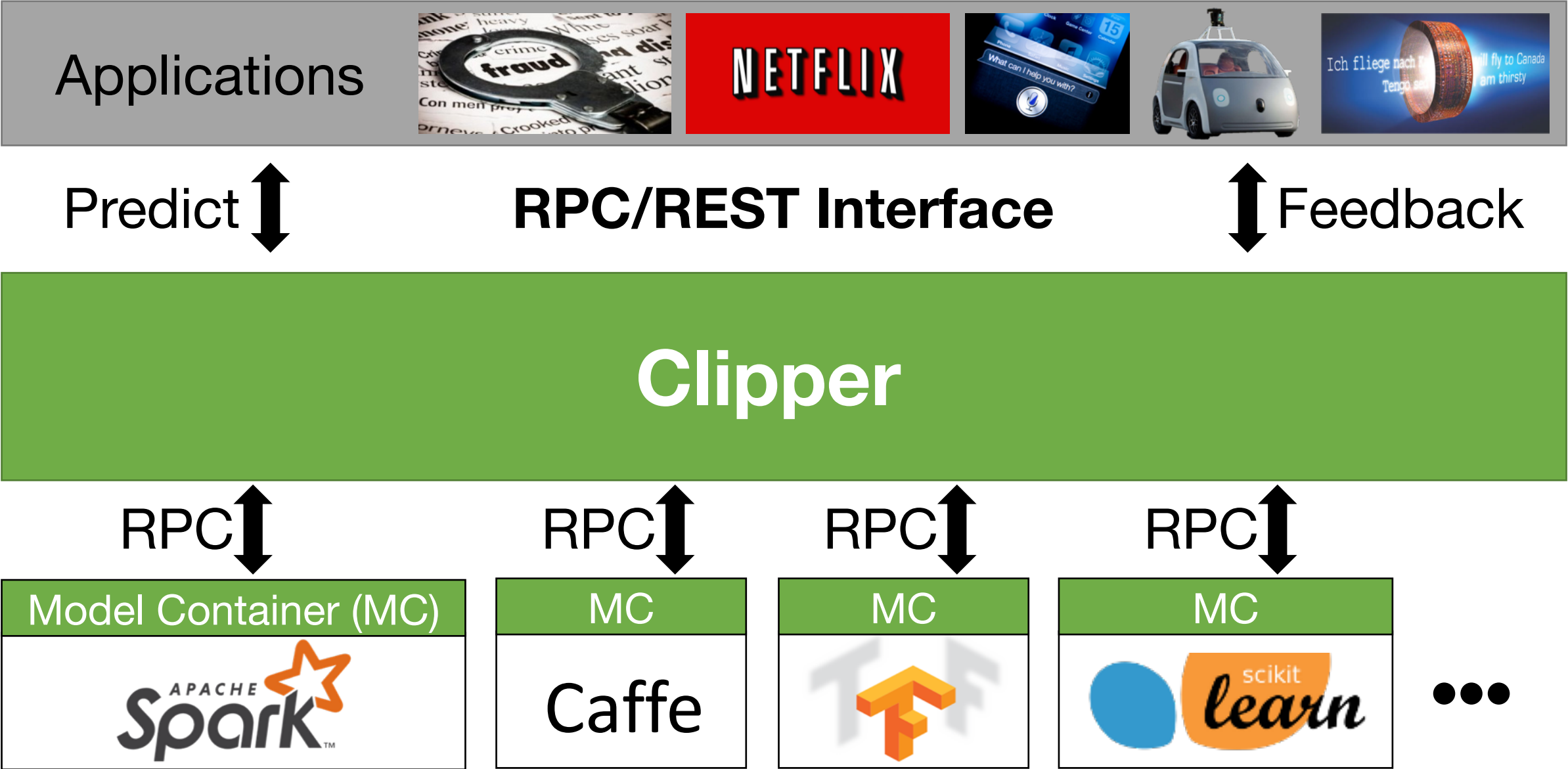*Support low-latency, high-throughput serving workloads*



*Large and growing ecosystem of ML models and frameworks*

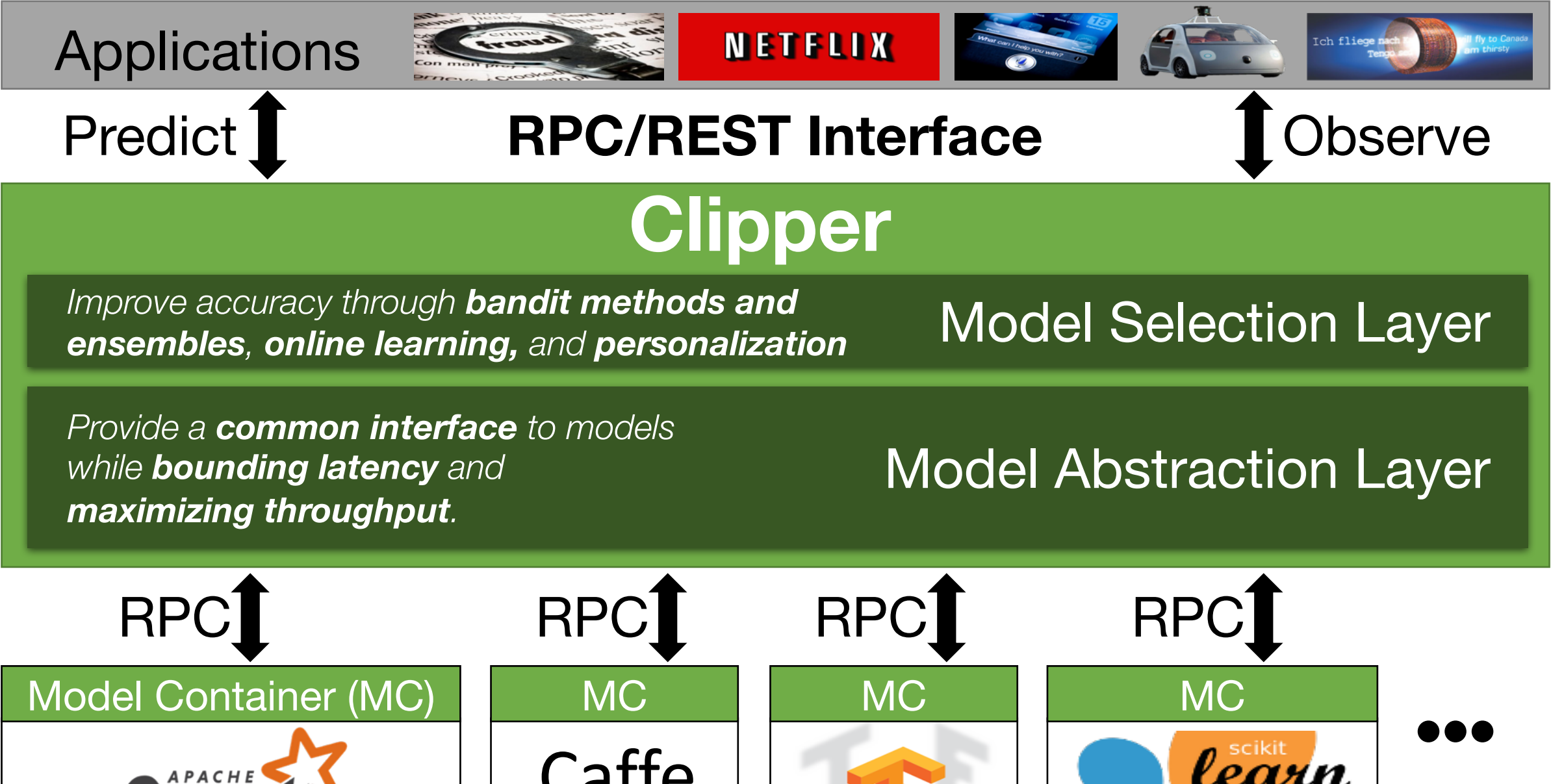*How does Clipper address these challenges?*

# Clipper Solutions

❑ *Simplifies deployment through layered architecture*

❑ *Serves many models across ML frameworks concurrently*

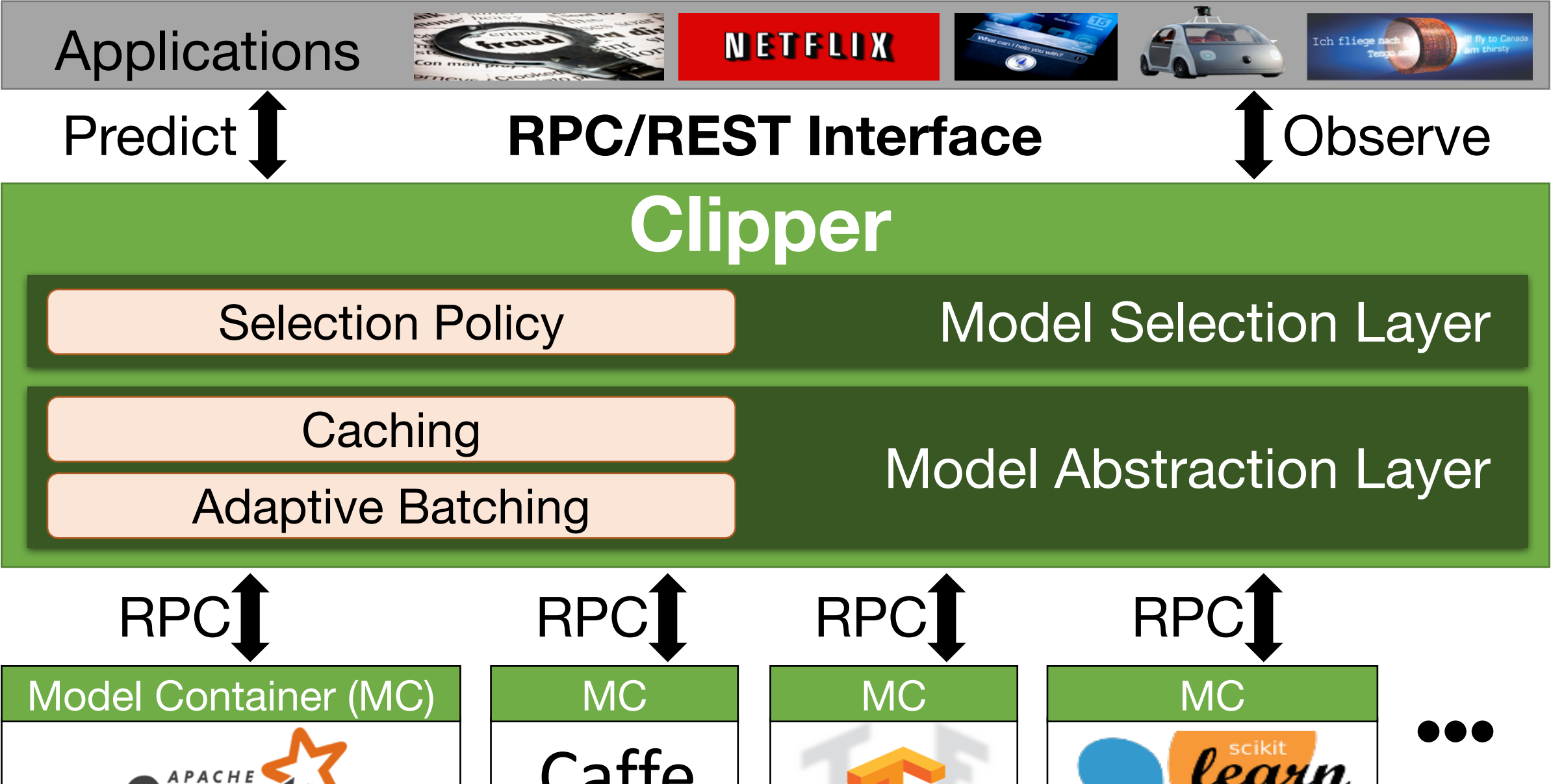❑ *Employs caching, batching, scale-out for high-performance serving*

# Clipper Decouples Applications and Models



**Applications**

↕ **Predict**   **RPC/REST Interface**   **Feedback** ↕

## Clipper

↕ RPC    ↕ RPC    ↕ RPC    ↕ RPC

Model Container (MC)   MC   MC   MC

Spark   Caffe   TF   scikit learn   •••

# Clipper Architecture

Applications

**Predict** ⇕  **RPC/REST Interface**  ⇕ Observe

## Clipper

*Improve accuracy through **bandit methods and ensembles**, **online learning**, and **personalization*** — Model Selection Layer

*Provide a **common interface** to models while **bounding latency** and **maximizing throughput**.* — Model Abstraction Layer

RPC ⇕    RPC ⇕    RPC ⇕    RPC ⇕

| Model Container (MC) | MC | MC | MC |
|---|---|---|---|

APACHE ★   Caffe   scikit learn   ●●●

# Clipper Architecture

Applications

Predict ⇕ **RPC/REST Interface** ⇕ Observe

## Clipper

**Model Selection Layer**
- Selection Policy

**Model Abstraction Layer**
- Caching
- Adaptive Batching

RPC ⇕  RPC ⇕  RPC ⇕  RPC ⇕

Model Container (MC)    MC    MC    MC    ●●●

# Clipper Implementation

Applications  NETFLIX   

Predict ↕ **RPC/REST Interface** ↕ Observe

## Clipper

Core system: 5000 lines of Rust 

RPC:

- 100 lines of Python

- 250 lines of Rust

- 200 lines of C++

Model Abstraction Layer

Caching

Adaptive Batching

RPC ⟷ Model Container (MC) — APACHE Spark

RPC ⟷ MC — Caffe

RPC ⟷ MC — TF

RPC ⟷ MC — scikit learn

●●●

**Common Interface → Simplifies Deployment:**

➢ Evaluate models using original code & systems

# Container-based Model Deployment

*Implement Model API:*

```python
class ModelContainer:
    def __init__(model_data)
    def predict_batch(inputs)
```

# Container-based Model Deployment

*Implement Model API:*

```
class ModelContainer:
    def __init__(model_data)
    def predict_batch(inputs)
```

➢ Implemented in many languages
   ➢ Python
   ➢ Java
   ➢ C/C++

# Container-based Model Deployment

*Model implementation packaged in container*

**Model Container (MC)**

```
class ModelContainer:
    def __init__(model_data)
    def predict_batch(inputs)
```
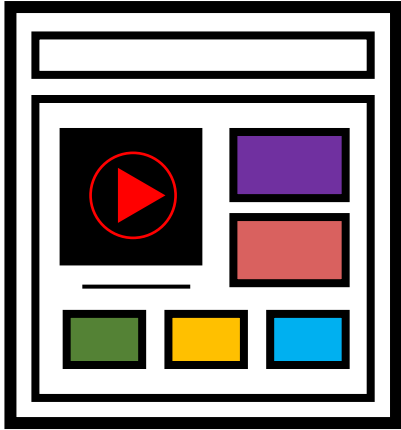
# Container-based Model Deployment

**Common Interface → Simplifies Deployment:**

➢ Evaluate models using original code & systems

➢ Models run in separate processes as Docker containers

    ➢ Resource isolation

Model Abstraction Layer

RPC          RPC          RPC          RPC          RPC          RPC

el Container (MC)    MC          MC          MC          MC          MC

**Common Interface → Simplifies Deployment:**

➢ Evaluate models using original code & systems

➢ Models run in separate processes as Docker containers

  ➢ Resource isolation

  ➢ Scale-out

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

# *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
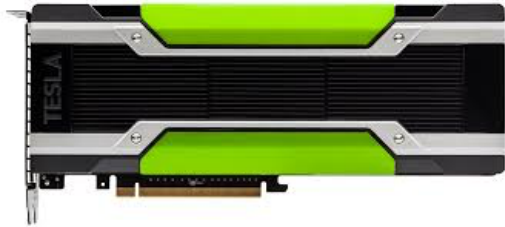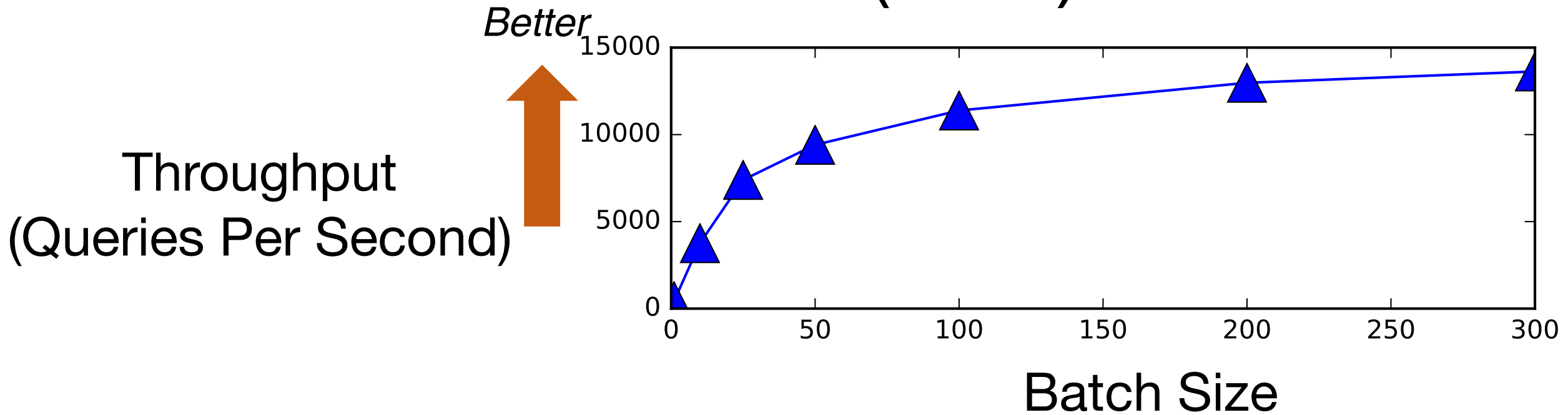
  ➤ model and framework

  ➤ system load

# *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
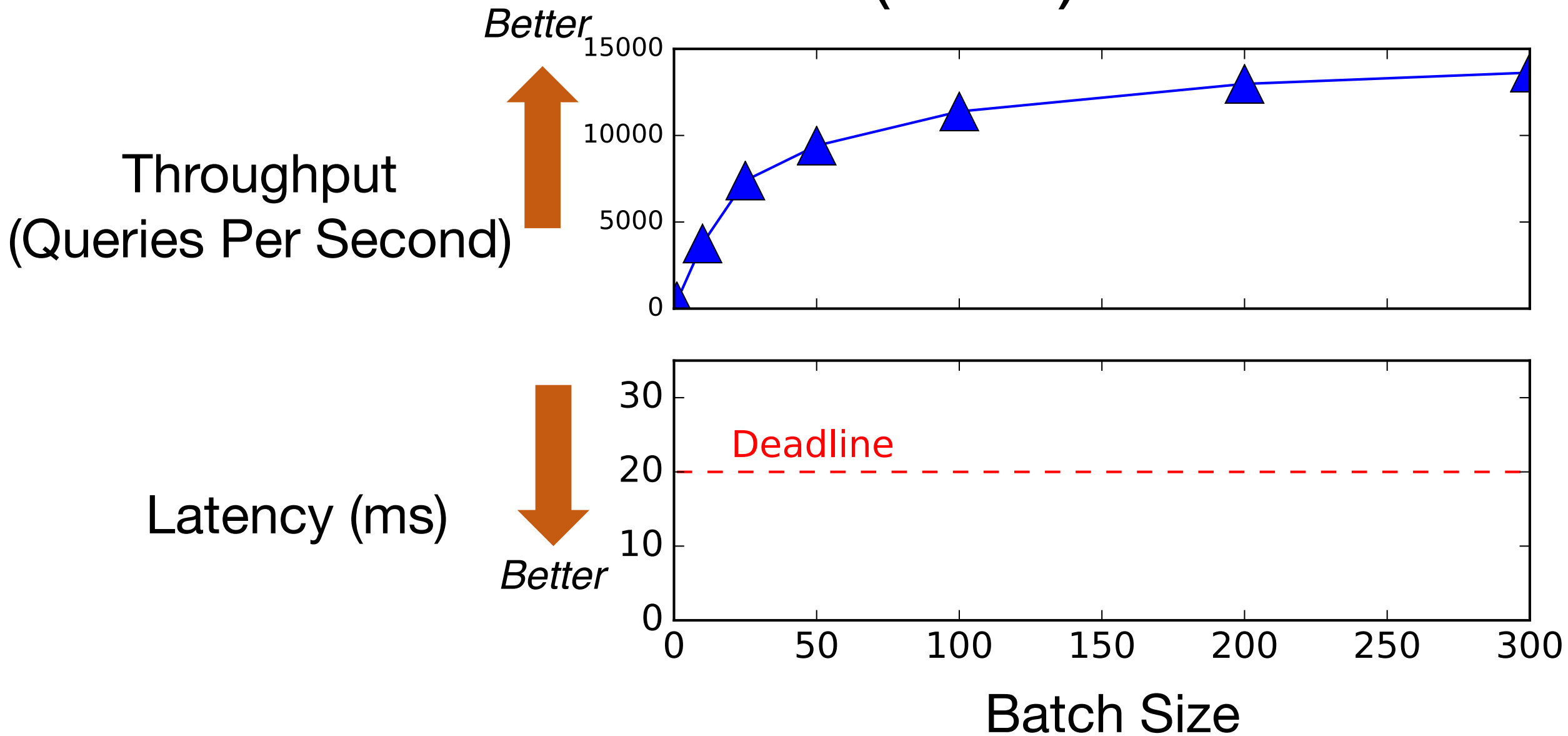  ➢ model and framework
  ➢ system load

**Clipper Solution:**

*Adaptively tradeoff latency and throughput…*

➢ Inc. batch size *until the latency objective is exceeded* (**Additive Increase**)

➢ If latency exceeds SLO cut batch size by a fraction (**Multiplicative Decrease**)
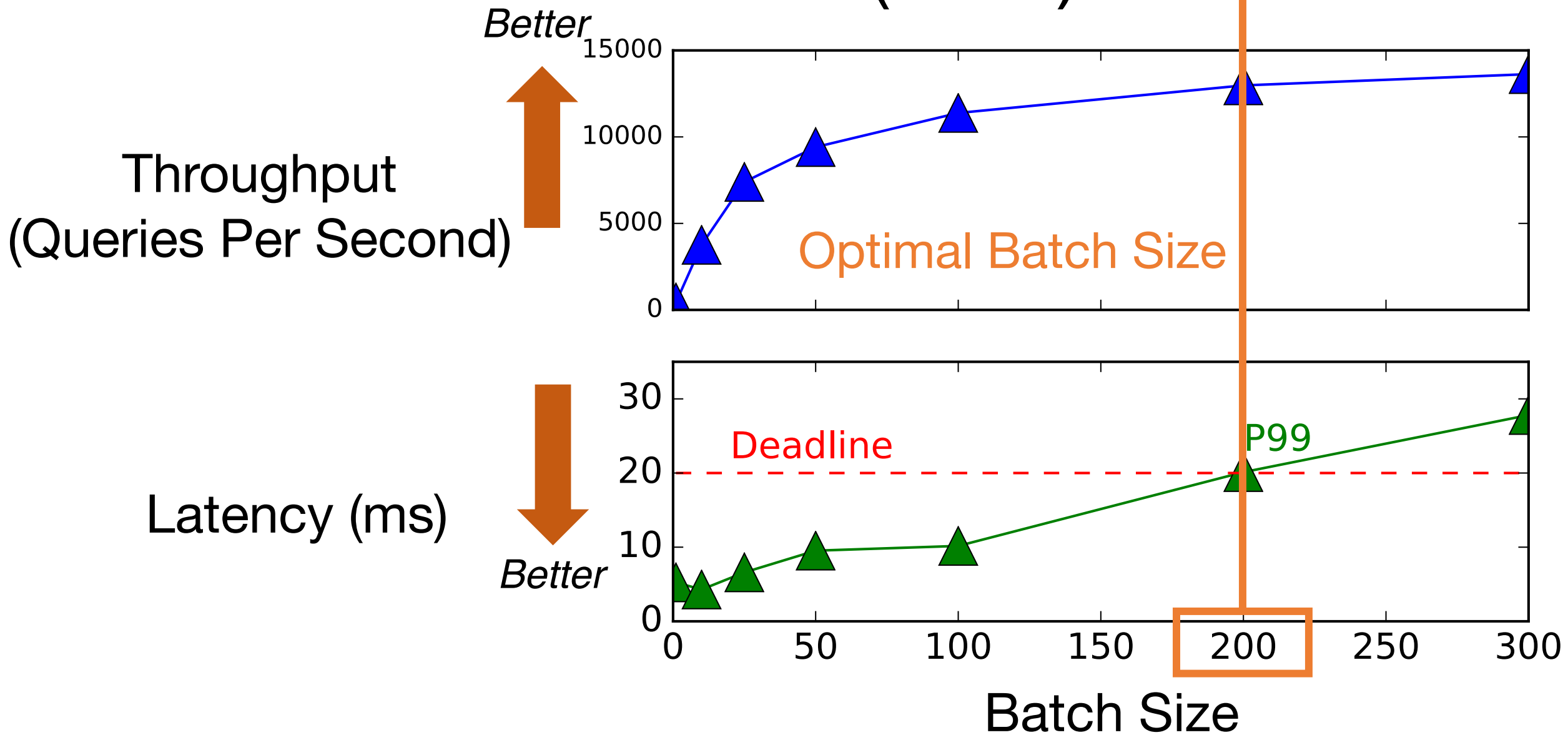
# Tensor Flow Conv. Net (GPU)

*Better*

Throughput
(Queries Per Second)

Batch Size

# Tensor Flow Conv. Net (GPU)

*Better*

Throughput
(Queries Per Second)

15000

10000

5000

0

Latency (ms)

*Better*

30

20

10

0

Deadline

0    50    100    150    200    250    300

Batch Size

# Tensor Flow Conv. Net (GPU)

*Better*

Throughput
(Queries Per Second)

Optimal Batch Size

15000

10000

5000

0

Deadline

P99

30

20

10

0

Latency (ms)

*Better*

0    50    100    150    200    250    300

Batch Size

Figure: Throughput (QPS) and P99 Latency (ms) comparison between Adaptive and No Batching strategies across different models.

**Throughput (QPS)** — Better ↑ (Adaptive vs No Batching):
- No-Op: 48386 vs 8963
- Random Forest (SKLearn): 22859 vs 317
- Linear SVM (PySpark): 29350 vs 7206
- Linear SVM (SKLearn): 48934 vs 1920
- Kernel SVM (SKLearn): 197 vs 203
- Log Regression (SKLearn): 47219 vs 1921

**P99 Latency (ms)** — Better ↓, 20 ms is Fast Enough (Adaptive vs No Batching):
- No-Op: 20 vs 0
- Random Forest (SKLearn): 20 vs 5
- Linear SVM (PySpark): 20 vs 0
- Linear SVM (SKLearn): 20 vs 0
- Kernel SVM (SKLearn): 28 vs 5
- Log Regression (SKLearn): 20 vs 0

# Overhead of decoupled architecture



Applications

**RPC/REST Interface**

Predict ↕          Feedback ↕

**Clipper**

RPC ↕     RPC ↕     RPC ↕     RPC ↕

MC        MC        MC        MC
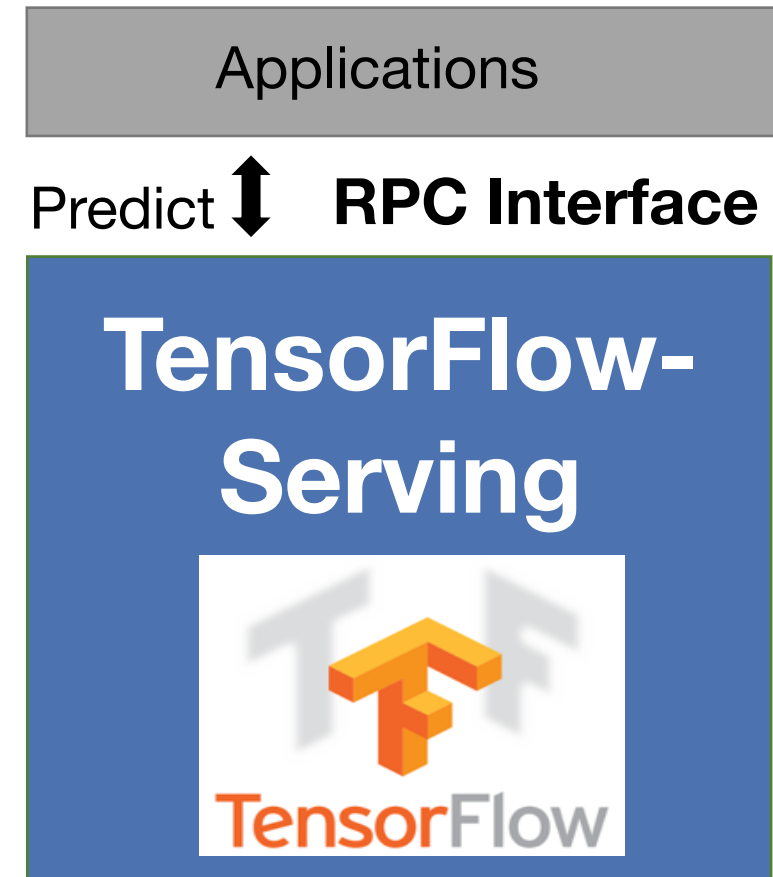
Spark     Caffe     TF        scikit learn     •••

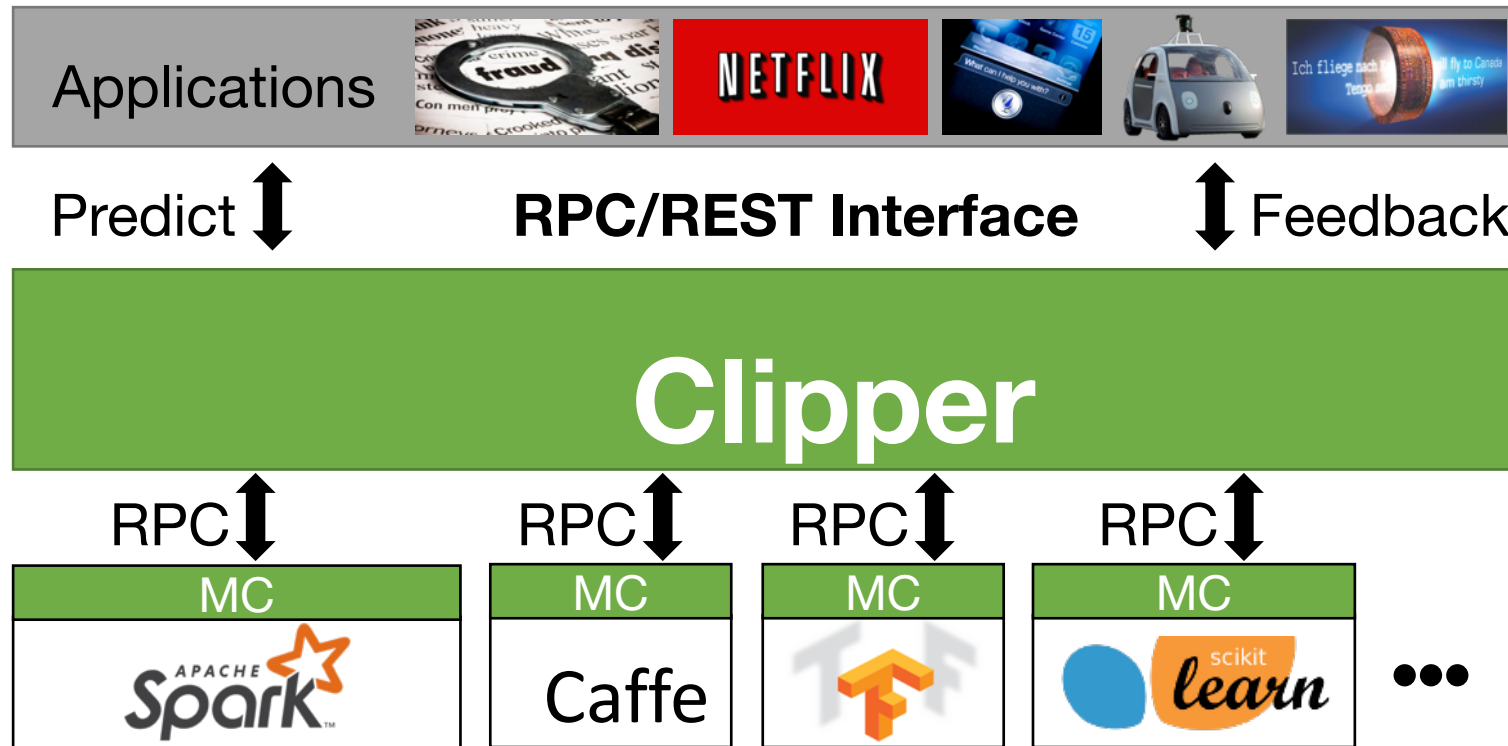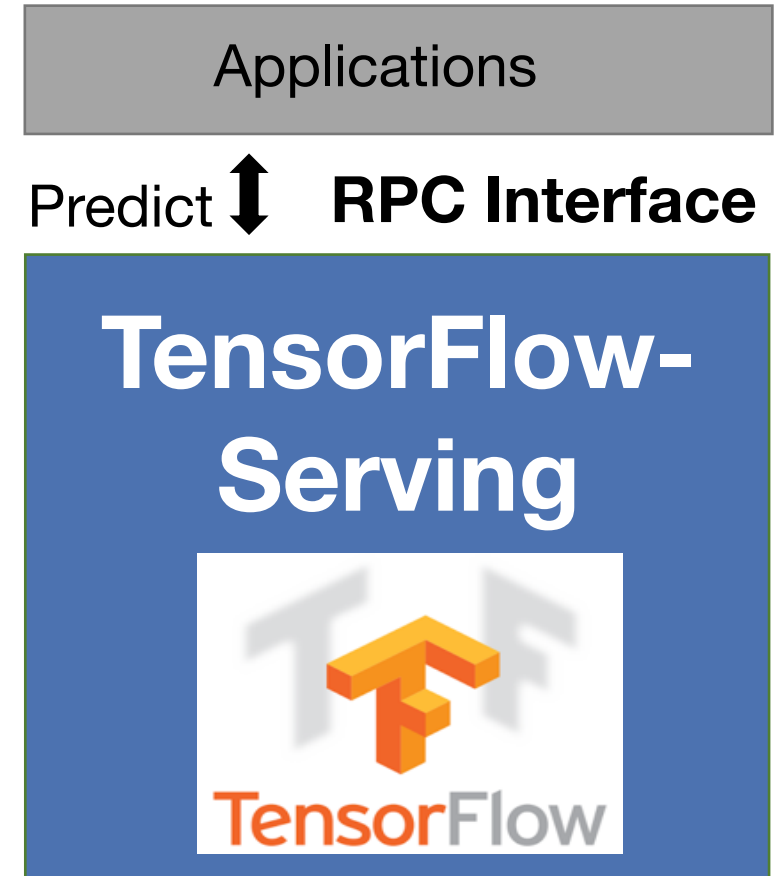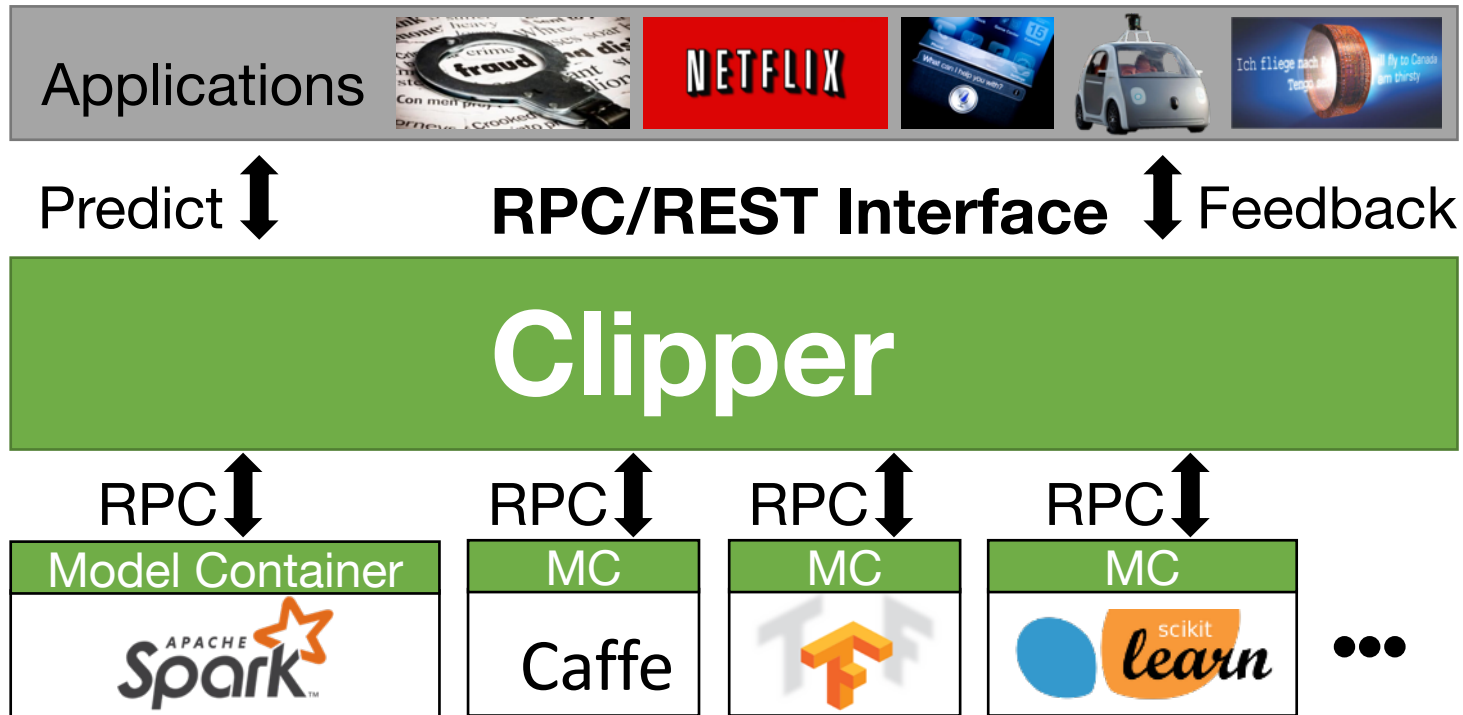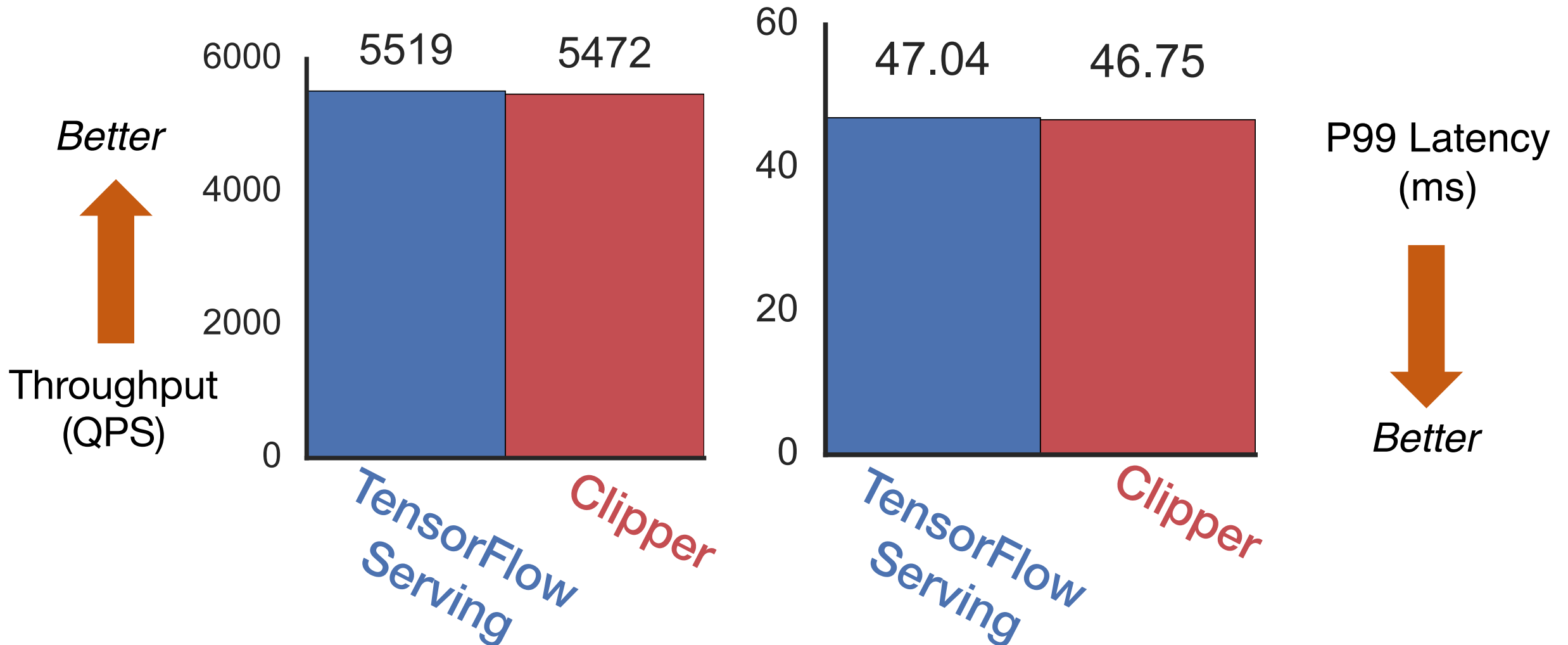# Overhead of decoupled architecture

# Overhead of decoupled architecture

# Overhead of decoupled architecture

Model: AlexNet trained on CIFAR-10
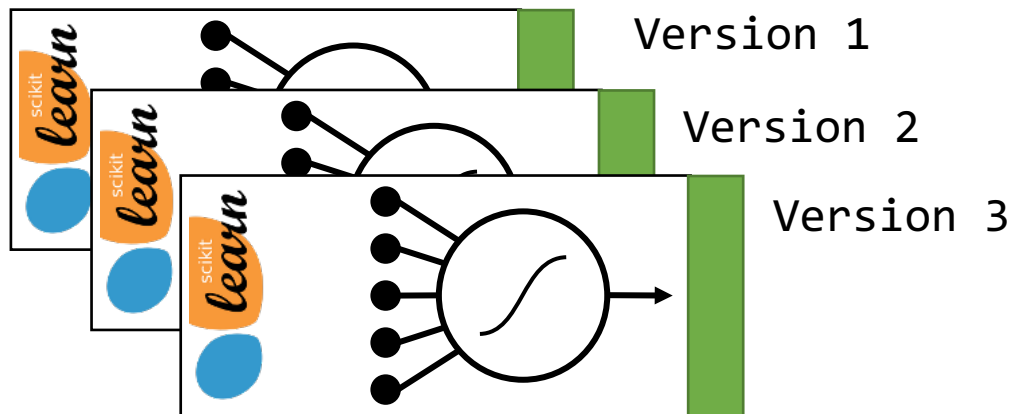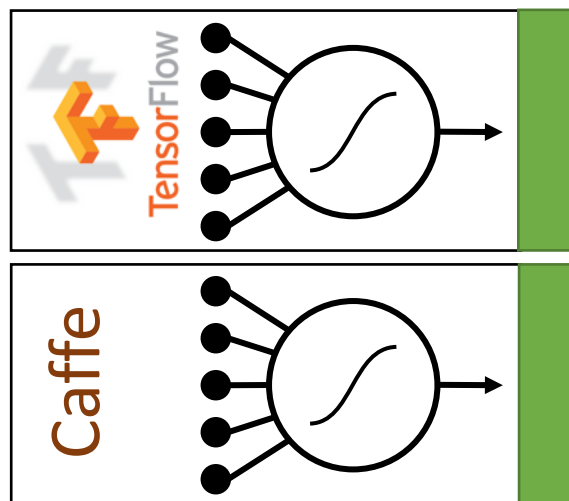
# Clipper Architecture

Applications

Predict ↕ **RPC/REST Interface** ↕ Observe

## Clipper

*Improve accuracy through **bandit methods and ensembles**, **online learning**, and **personalization*** Model Selection Layer

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

RPC ↕        RPC ↕        RPC ↕        RPC ↕

Model Container (MC)    MC    MC    MC    •••

APACHE ⭐    Caffe    [image]    scikit learn

**Clipper**

*Improve accuracy through* **bandit methods and ensembles**, **online learning,** *and* **personalization**

# Model Selection Layer
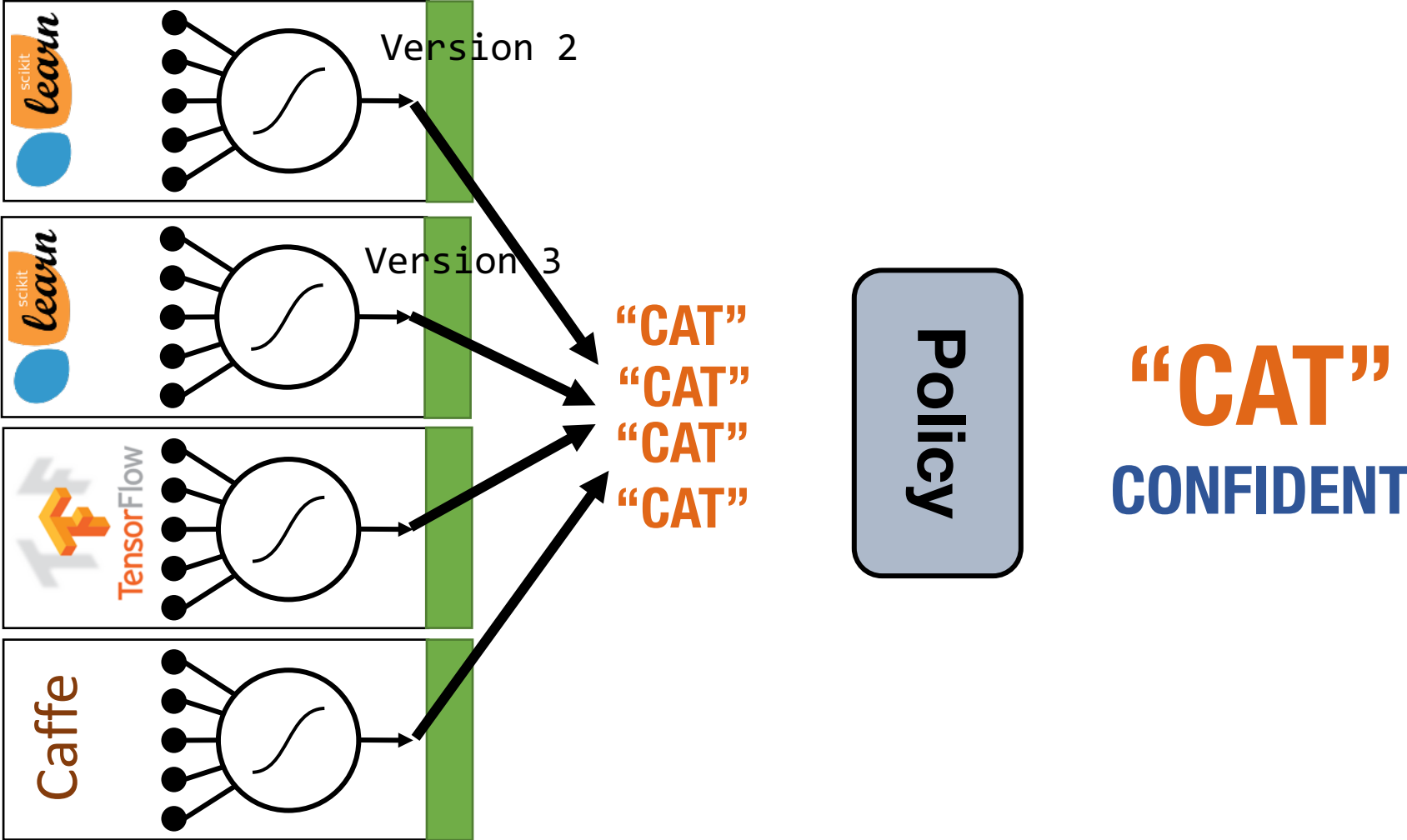
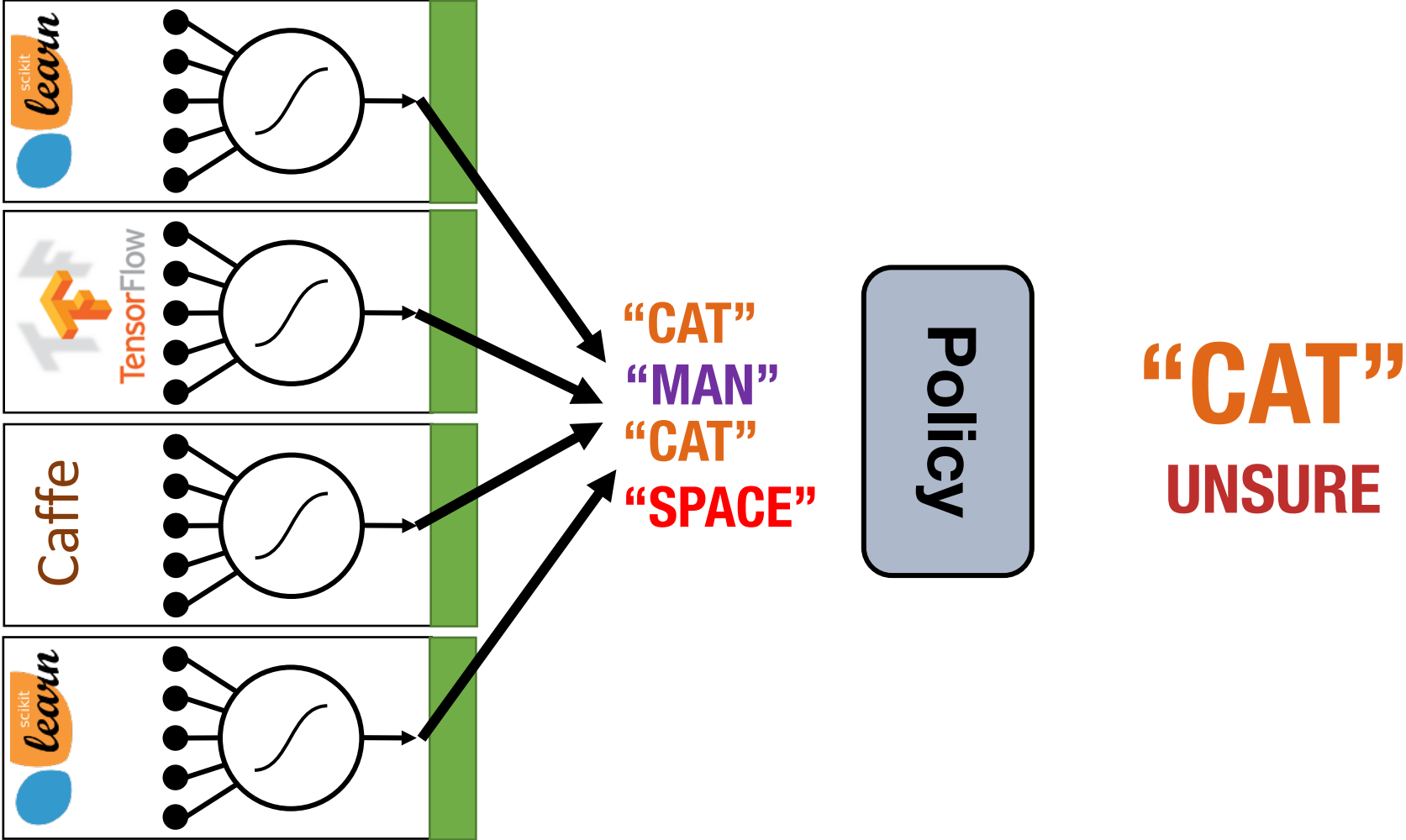Version 1

Version 2

Version 3

*Periodic retraining*

*Experiment with new models and frameworks*

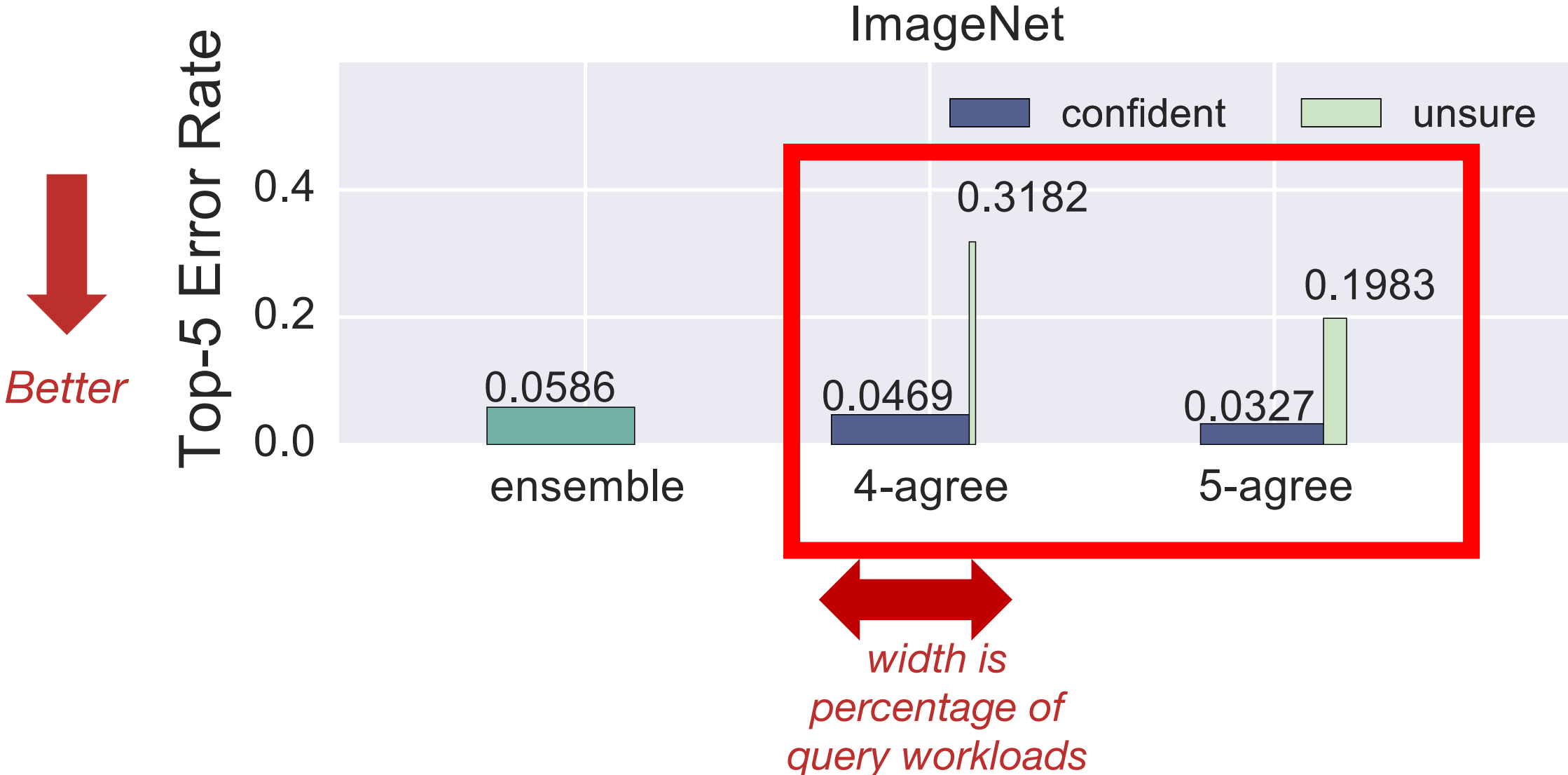# Selection Policy: Estimate confidence

# Selection Policy: Estimate confidence

# Selection Policy: Estimate confidence

# Selection Policy: Estimate confidence

**Clipper**

*Selection policies supported by Clipper*

➤ Exploit multiple models to estimate confidence

➤ Use multi-armed bandit algorithms to learn optimal model-selection online

➤ Online personalization across ML frameworks

*\*See paper for details*

# Conclusion

➤ *Prediction-serving* is an important and *challenging* area for *systems research*

  ➤ Support *low-latency, high-throughput* serving workloads

  ➤ Serve *large* and growing *ecosystem of ML frameworks*

➤ *Clipper* is a *first step* towards addressing these challenges

  ➤ *Simplifies deployment* through layered architecture

  ➤ Serves many models *across ML frameworks* concurrently

  ➤ Employs *caching, adaptive batching, container scale-out* to meet interactive serving workload demands

➤ Beyond academic prototype to build a real, *open-source system*

*https://github.com/ucbrise/clipper*
**crankshaw@cs.berkeley.edu**

# GPU Cluster Scaling