



# Research at the intersection of AI + Systems

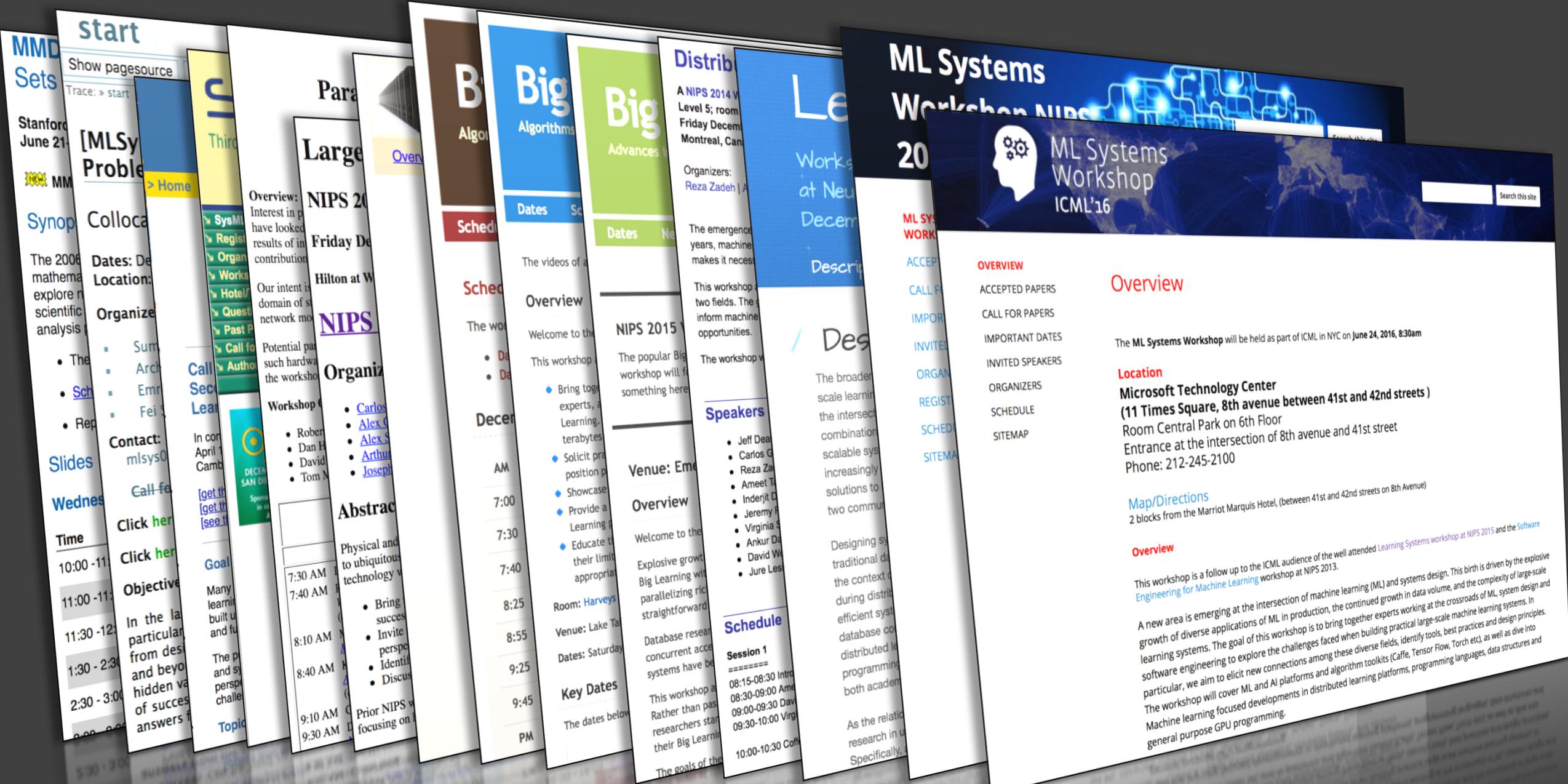
Joseph E. Gonzalez

Assistant Professor, UC Berkeley

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

# Looking Back on AI Systems

Going back to when I started graduate school ...



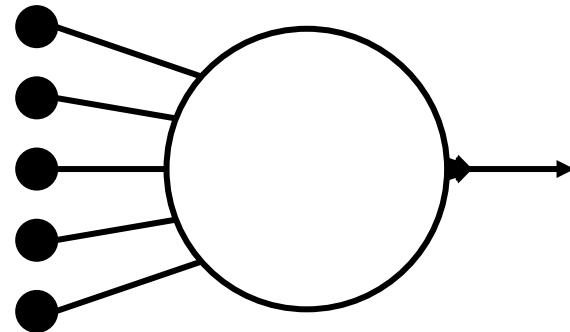
# Machine learning community has had an evolving focus on AI Systems



# Learning



The focus of AI Systems research has been on model training.



Enabling **Machine Learning** and **Systems** Innovations

**Stochastic  
Optimization**

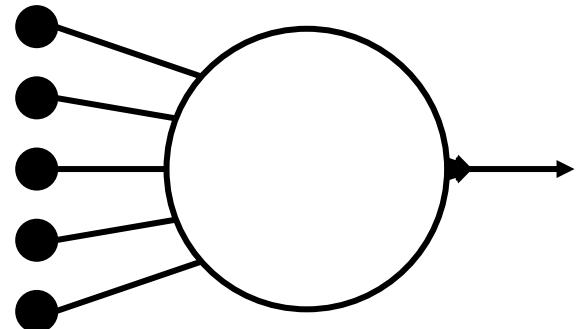
**Distributed  
Dataflow Systems**

**Deep Learning  
(CNN/RNN)**

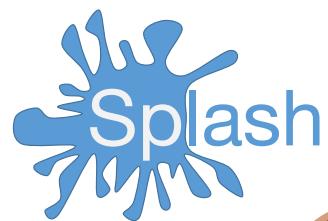
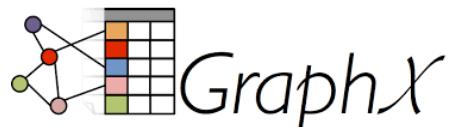
**Domain Specific  
Languages (TensorFlow)**

**Symbolic  
Methods**

**GPU / TPU  
Acceleration**



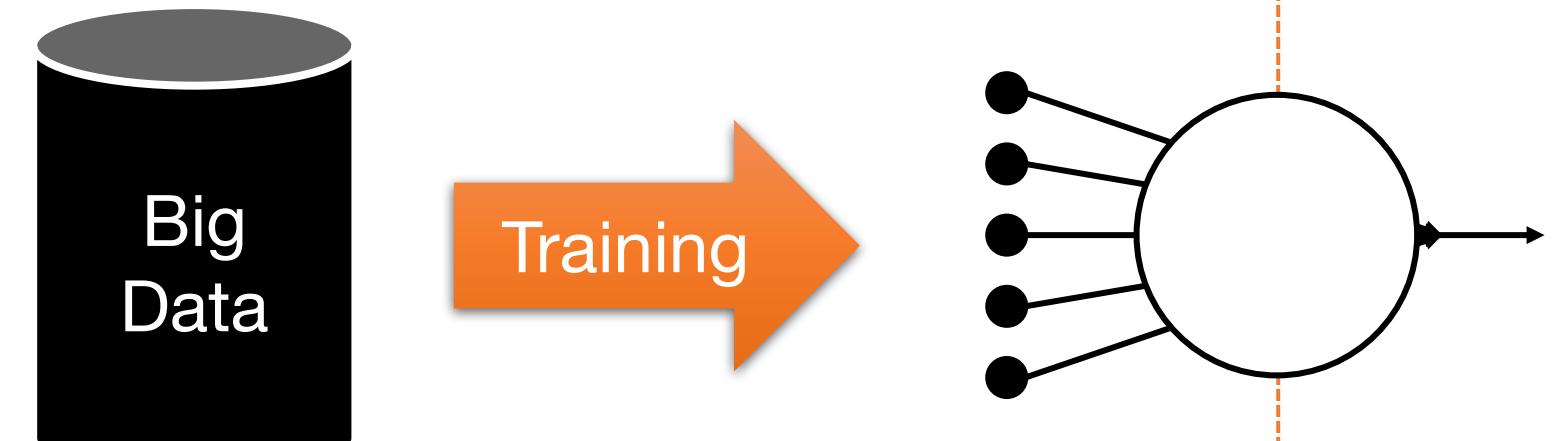
KeystoneML



rllab



# Learning

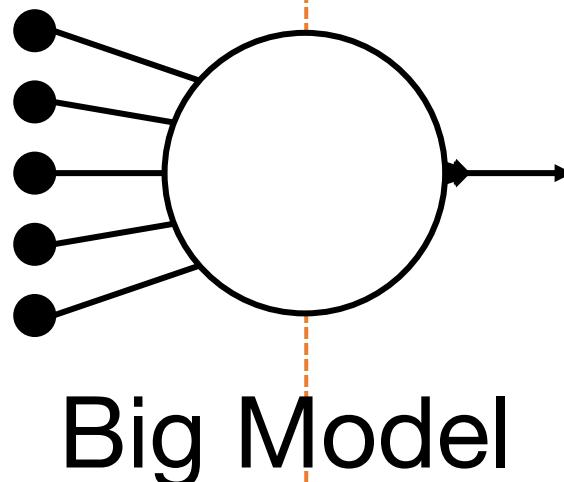


?

# Learning



Training



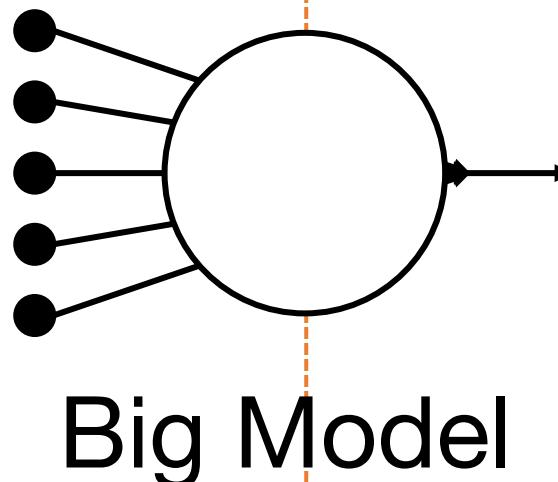
# Drive Actions



# Learning



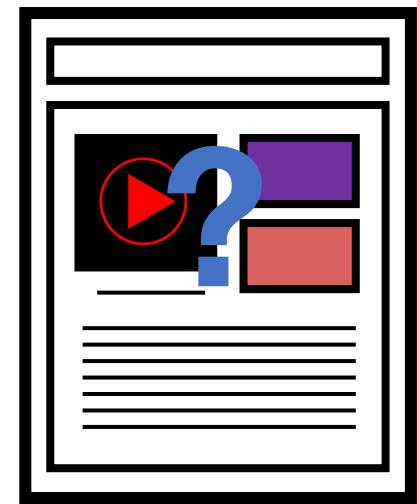
Training



# Prediction

Query

Decision

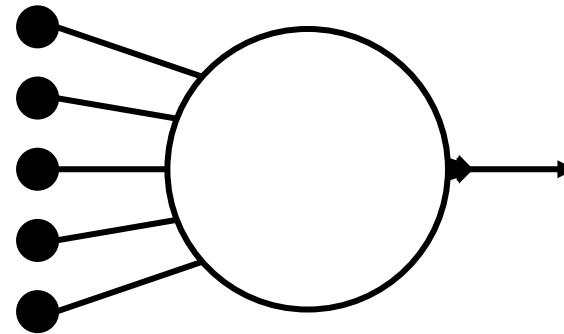


Application

# Learning

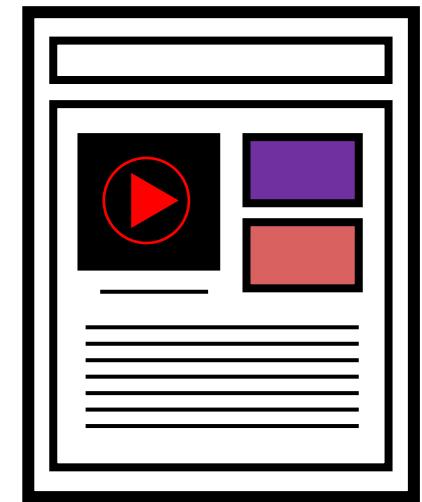
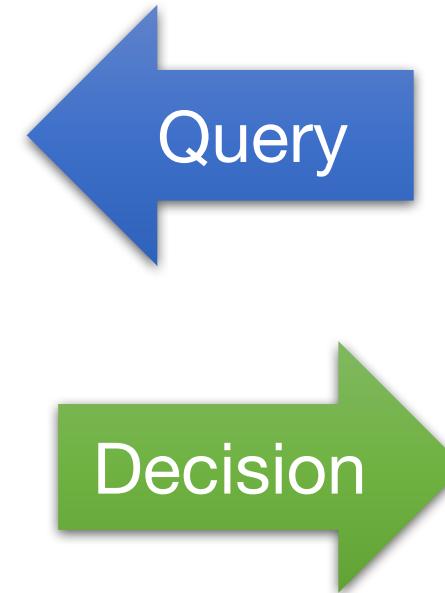


Training



Big Model

# Prediction



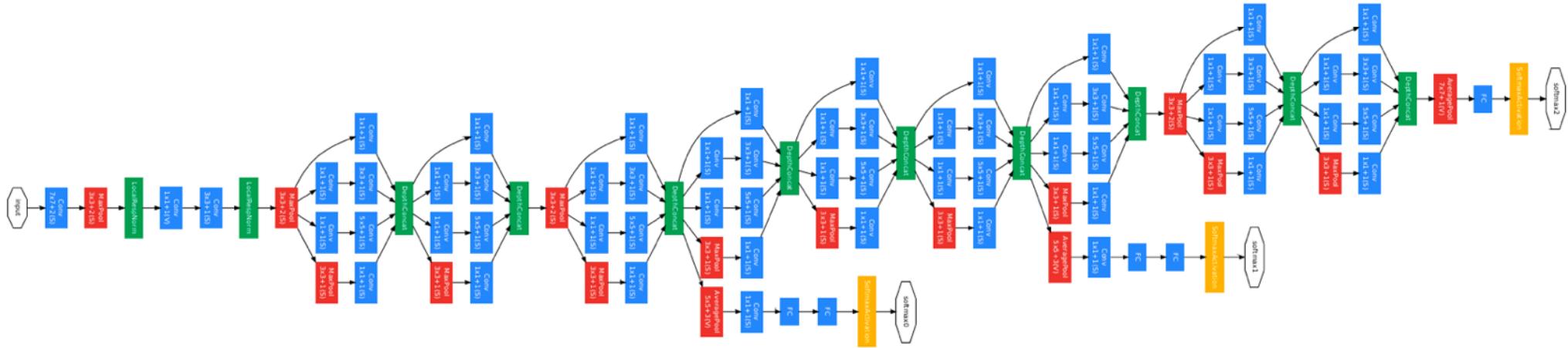
Application

**Goal:** ~10 ms under heavy load

*Complicated by Deep Learning*

→ New **ML Algorithms** and **Systems**

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



# *Models getting more complex*

- 10s of GFLOPs [1]
  - Recurrent nets

*Deployed on critical path*

  - Maintain latency goals under heavy load

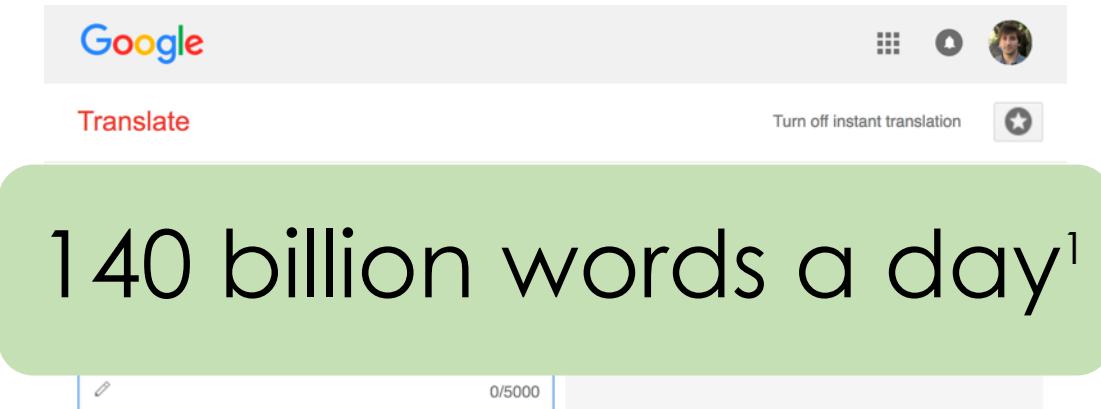
[1] Deep Residual Learning for Image Recognition. He et al. CVPR 2015.



# Using specialized hardware for predictions

# Google Translate

## Serving



82,000 GPUs  
running 24/7

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

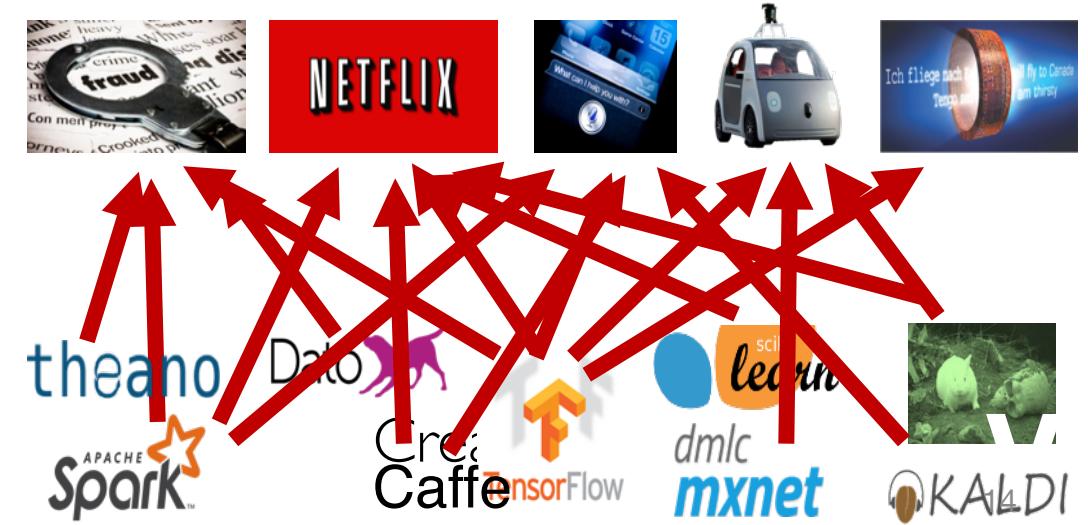
*"If each of the **world's Android phones** used the new Google voice search for just **three minutes a day**, these engineers realized, the company would **need twice as many data centers.**"*  
– Wired

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

# Prediction-Serving Challenges



Support low-latency, high-throughput serving workloads

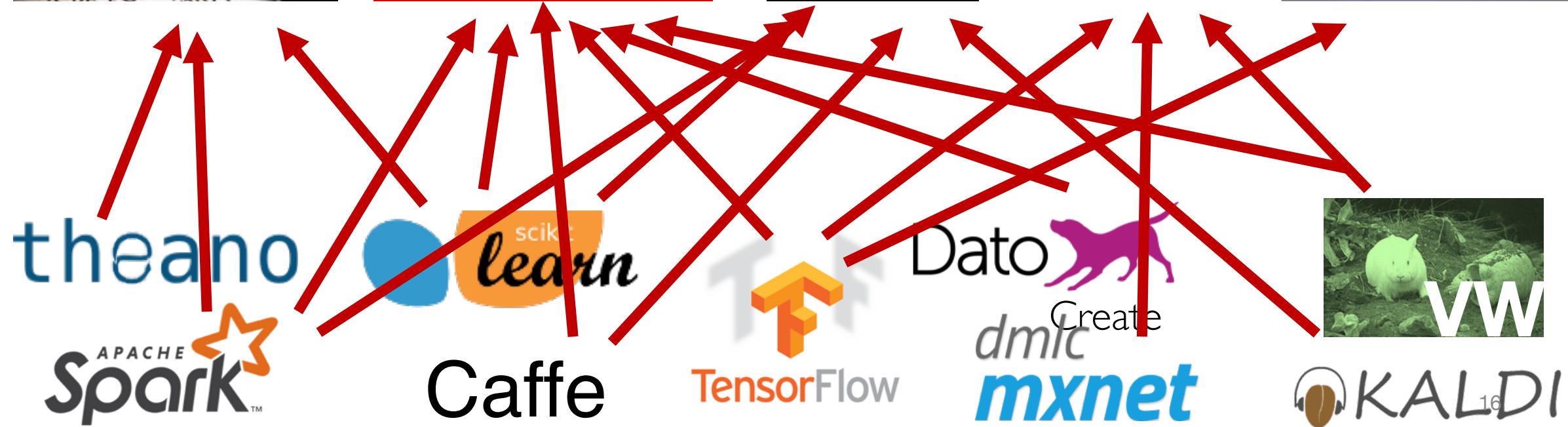


Large and growing ecosystem of ML models and frameworks

# Wide range of application and frameworks



# Wide range of application and frameworks



# One-Off Systems for High-Value Tasks

## Problems:

Expensive to build and maintain

- Requires **AI + Systems expertise**

Tightly-coupled model, framework, and application

- Difficult to **update models** and **add new frameworks**

# Prediction Serving is an Open Problem

- Computationally challenging
  - Need low latency & high throughput
- No standard technology or abstractions for serving models



Low Latency  
Prediction Serving System  
[NSDI'17]

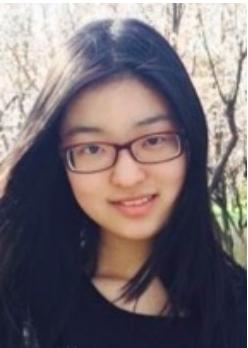
**IDK** Prediction  
Cascades  
Learning to make fast predictions  
[Work in Progress]

# Clipper

## Low Latency Prediction Serving System



Daniel  
Crankshaw



Xin  
Wang



Yika  
Luo



Giulio  
Zhou



Corey  
Zumar

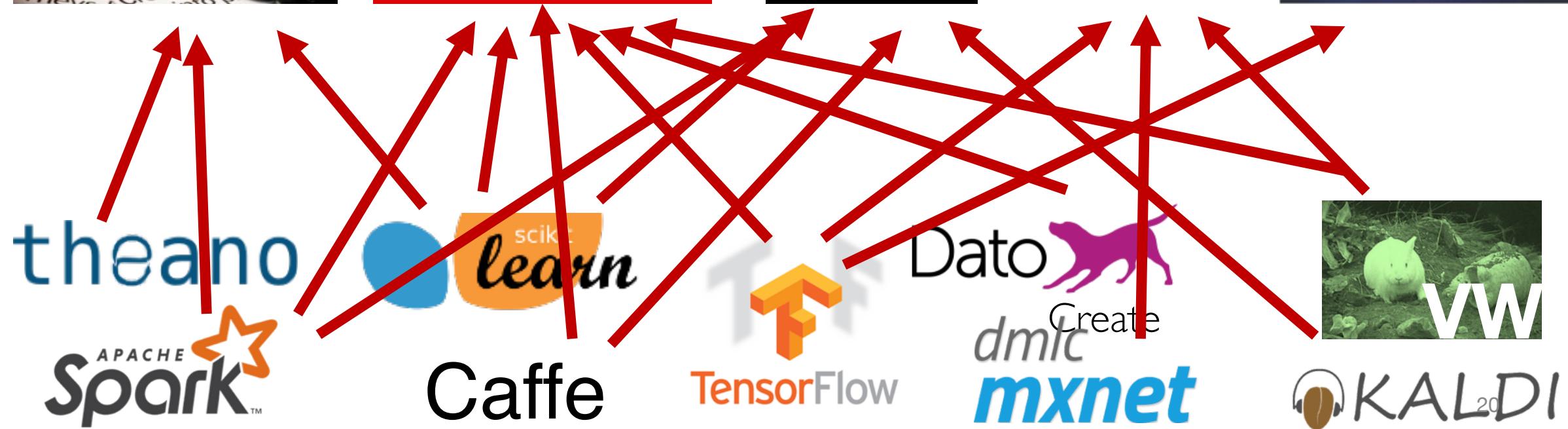


Alexey  
Tumanov



Ion  
Stoica

# Wide range of application and frameworks





---

# Middle layer for prediction serving.

Common  
Abstraction

System  
Optimizations

---

theano  
APACHE  
Spark™

The scikit-learn logo, which consists of two overlapping circles, one blue and one orange, with the word "scikit" in white on the blue circle and "learn" in black on the orange circle.

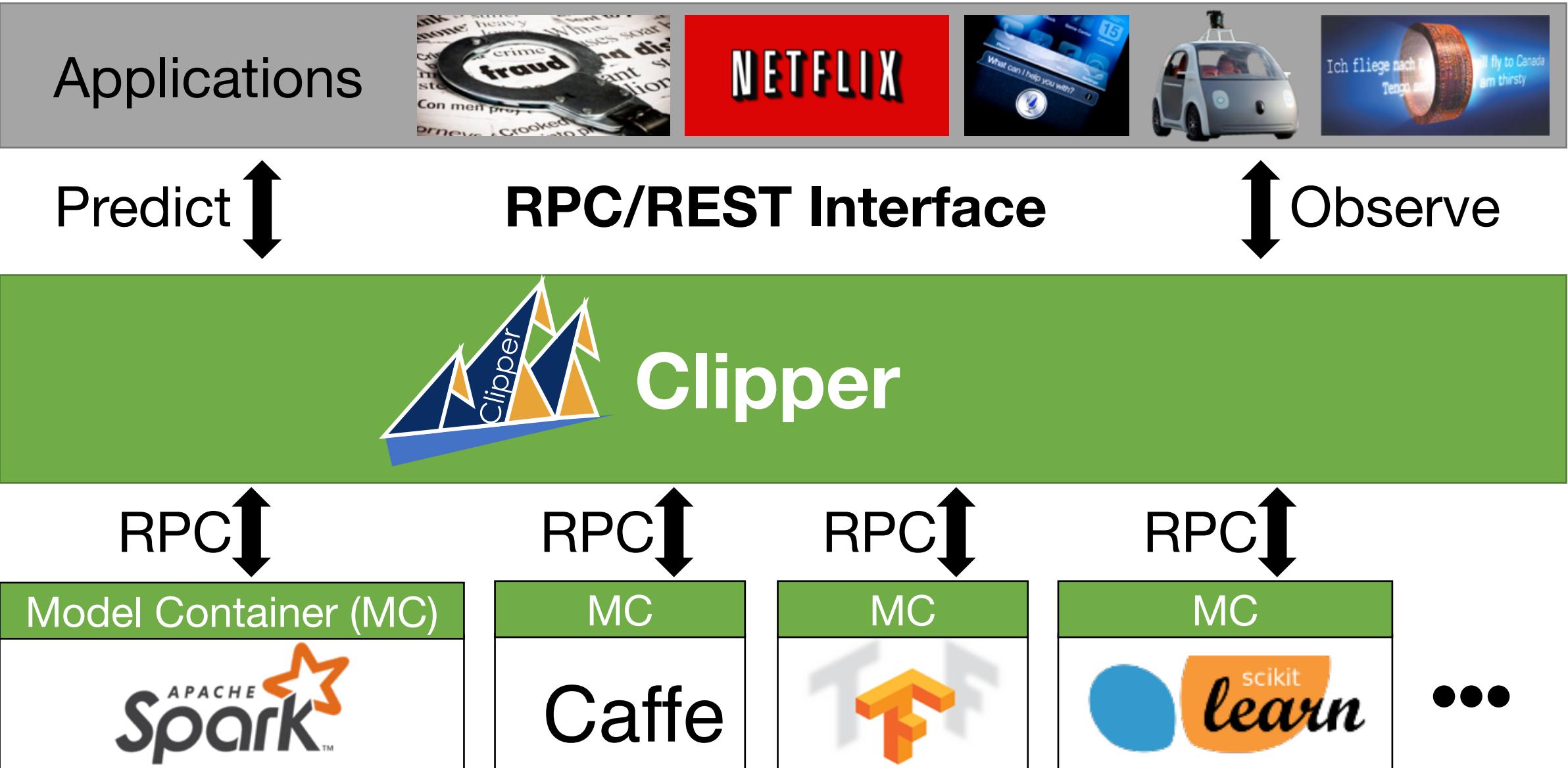
Caffe

The TensorFlow logo, featuring a stylized orange "F" shape composed of many smaller triangles, with the word "TensorFlow" in a smaller, sans-serif font below it.

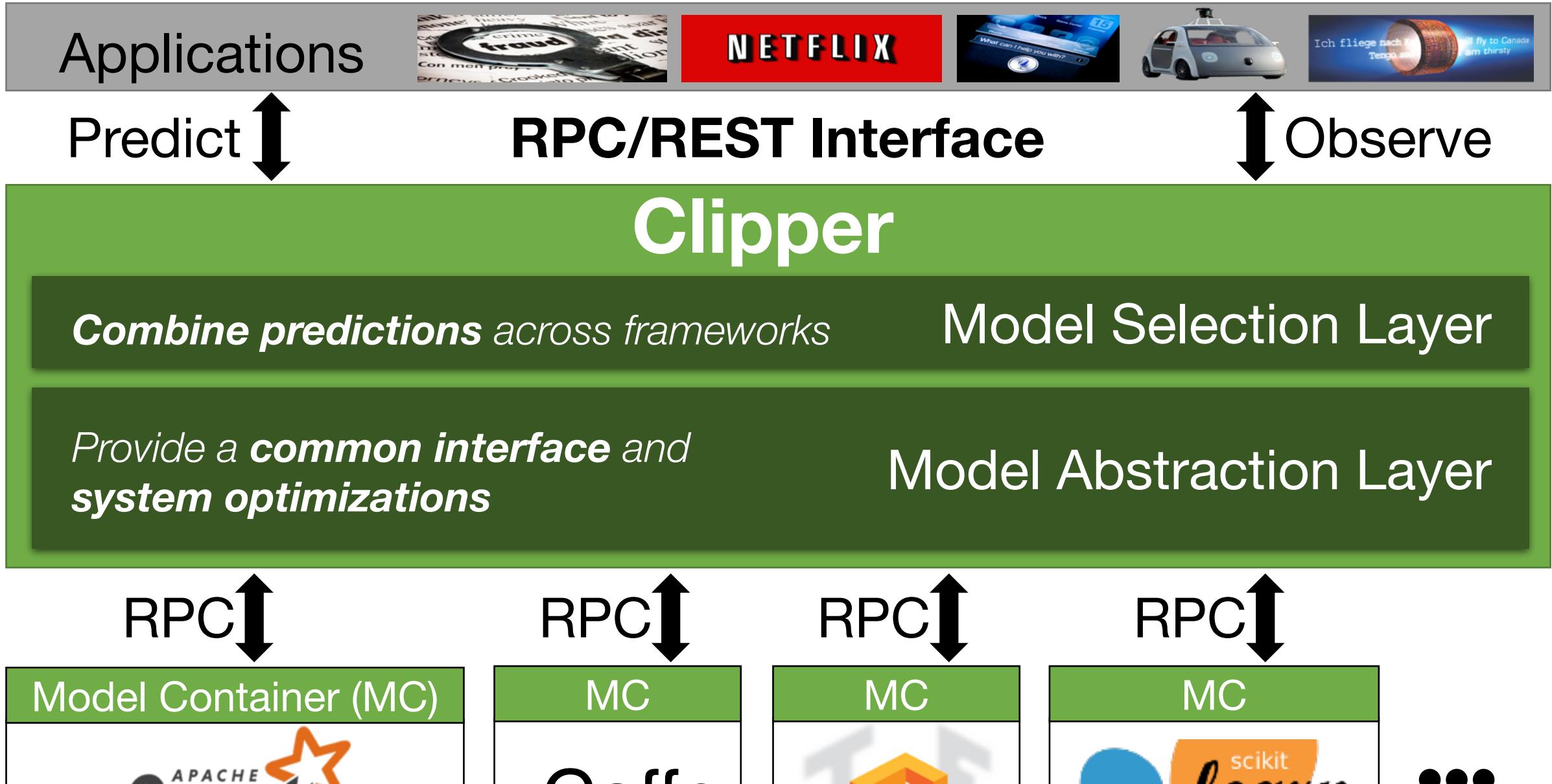
Dato  
Create  
dmrc  
mxnet

The Volkswagen logo, featuring a white rabbit standing on a grassy hill next to a large, stylized green "VW" monogram.The OKALDI logo, featuring a brown coffee bean with a white "OK" and "ALDI" monogram above it, with a small "21" at the bottom right.

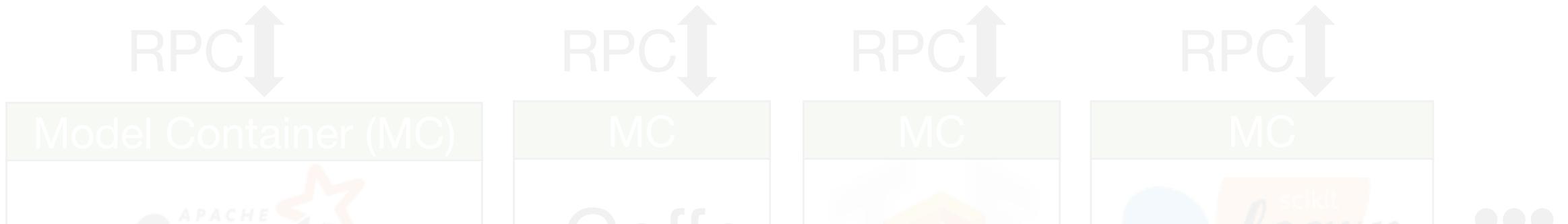
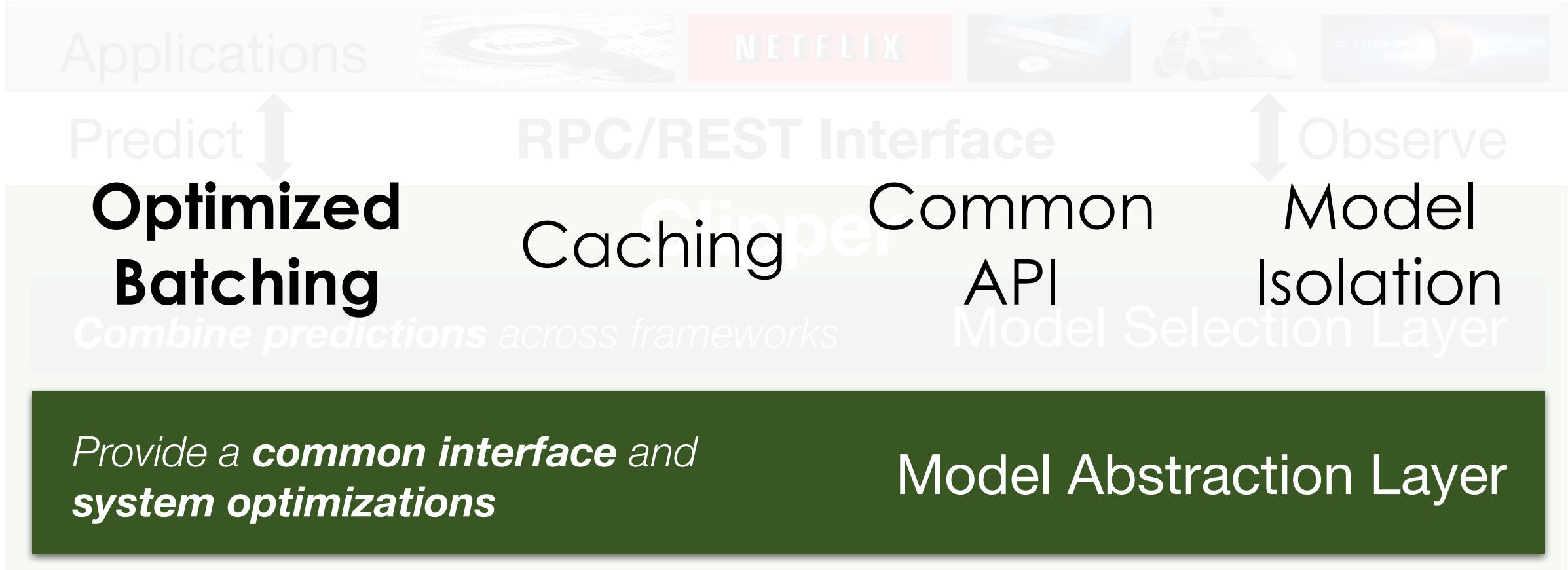
# Clipper Decouples Applications and Models



# Clipper Architecture

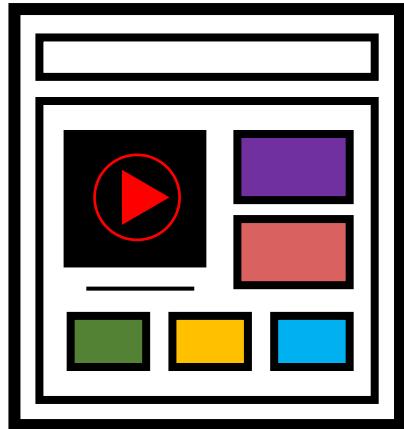


# Clipper Architecture



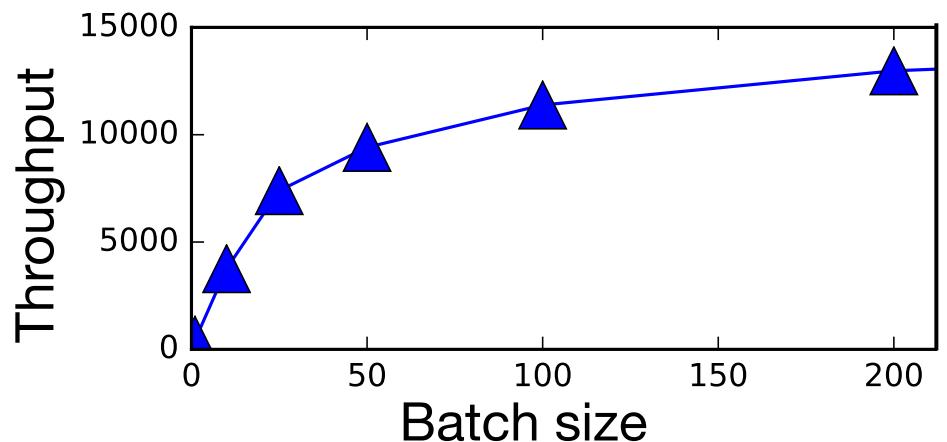
# Batching to Improve Throughput

- Why batching helps:



A single page load may generate many queries

Throughput-optimized frameworks



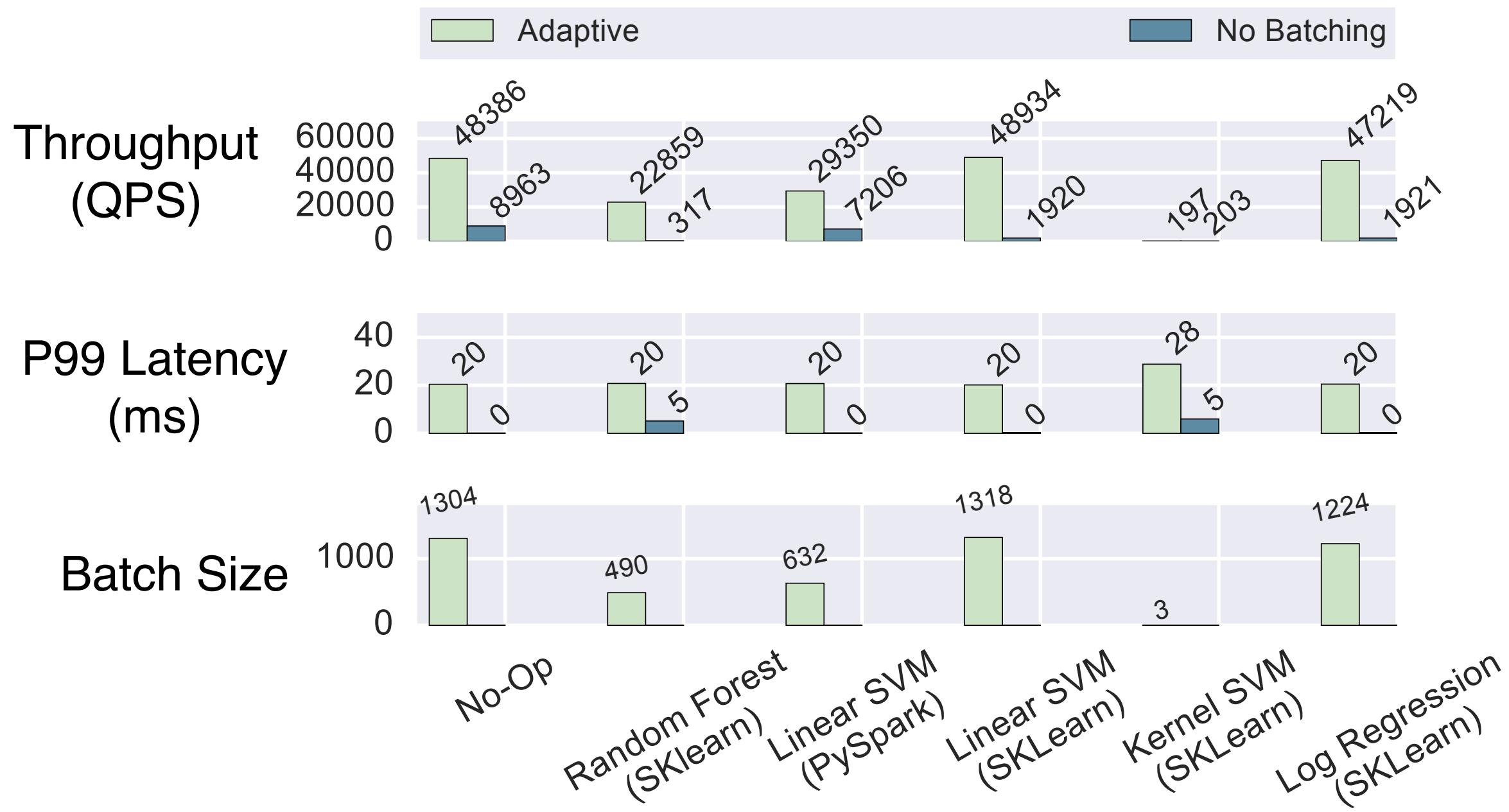
- Optimal batch depends on:

- hardware configuration
- model and framework
- system load

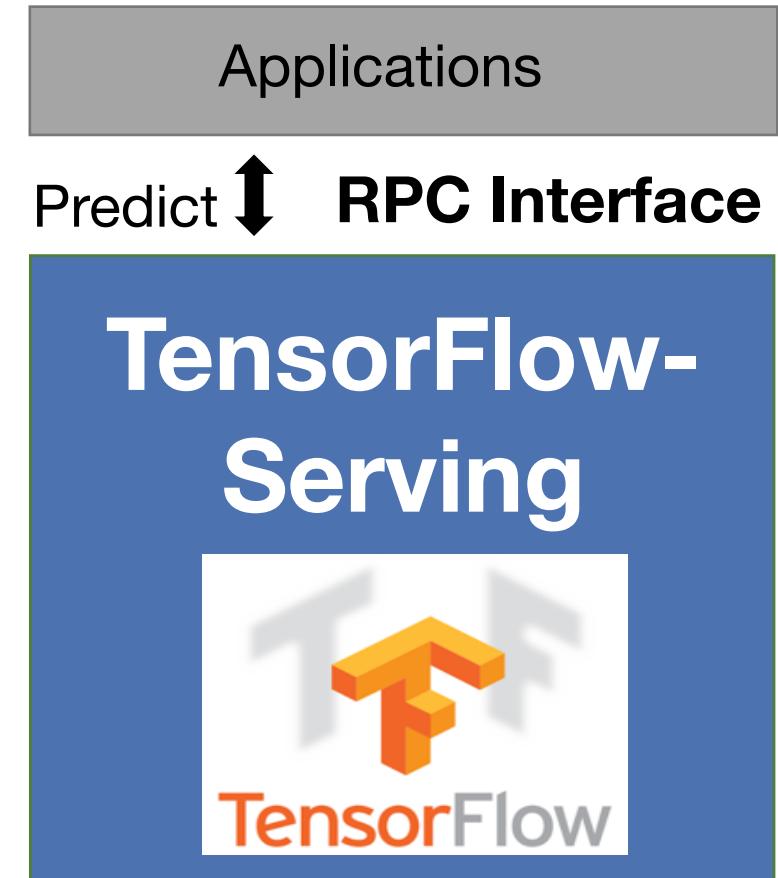
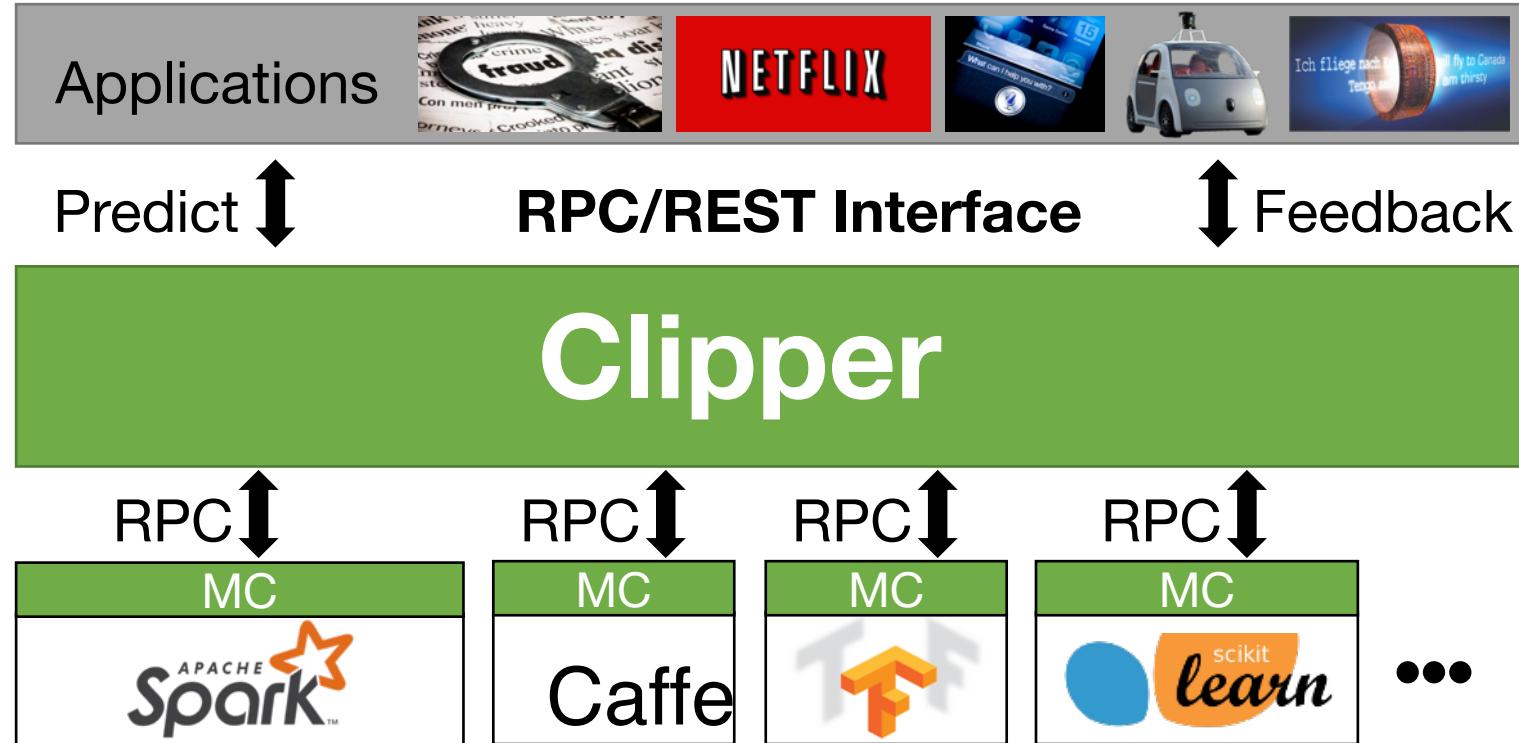
## Clipper Solution:

*Adaptively tradeoff latency and throughput...*

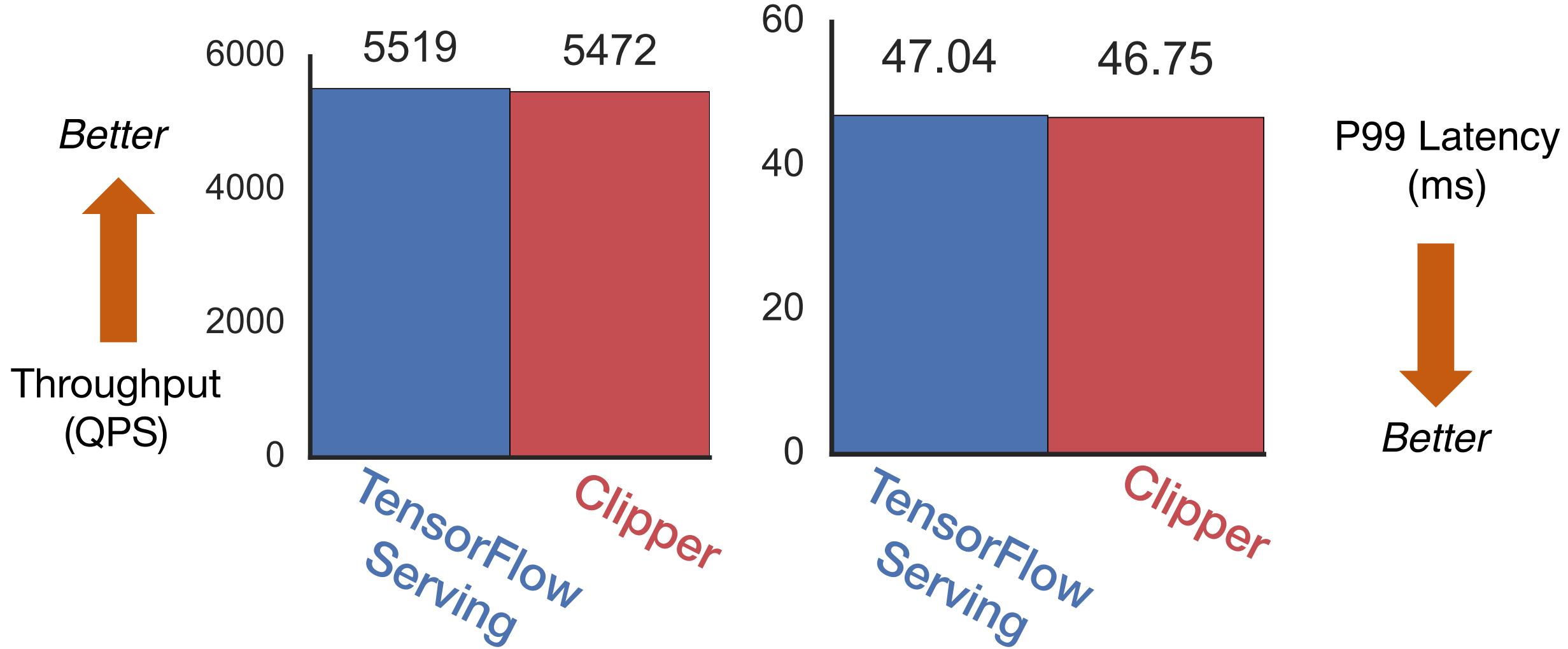
- Inc. batch size *until the latency objective is exceeded (Additive Increase)*
- If latency exceeds SLO cut batch size by a fraction (**Multiplicative Decrease**)



# Overhead of decoupled architecture

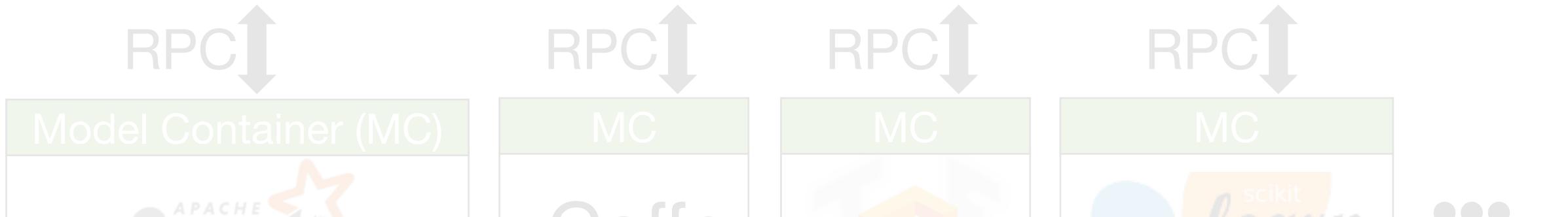
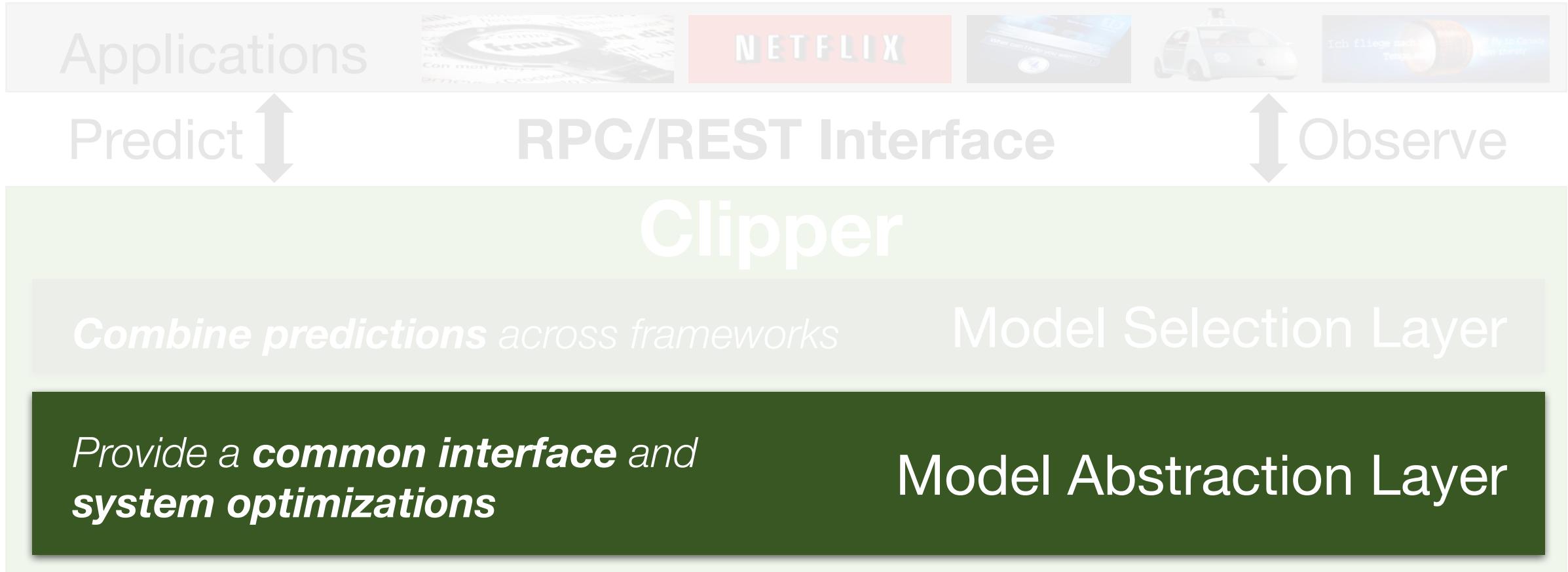


# Overhead of decoupled architecture

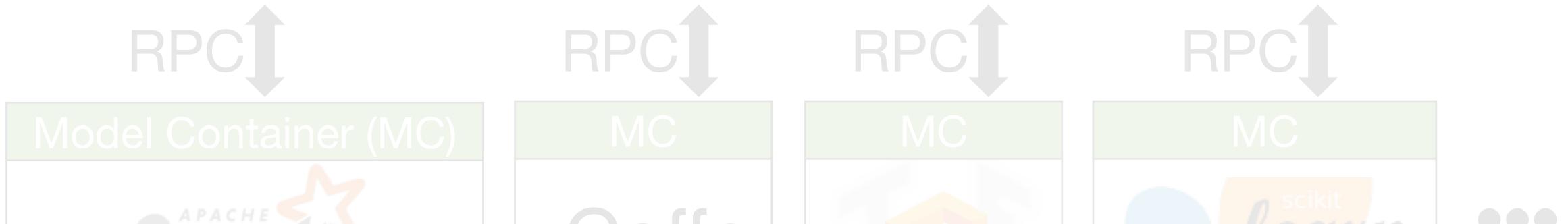
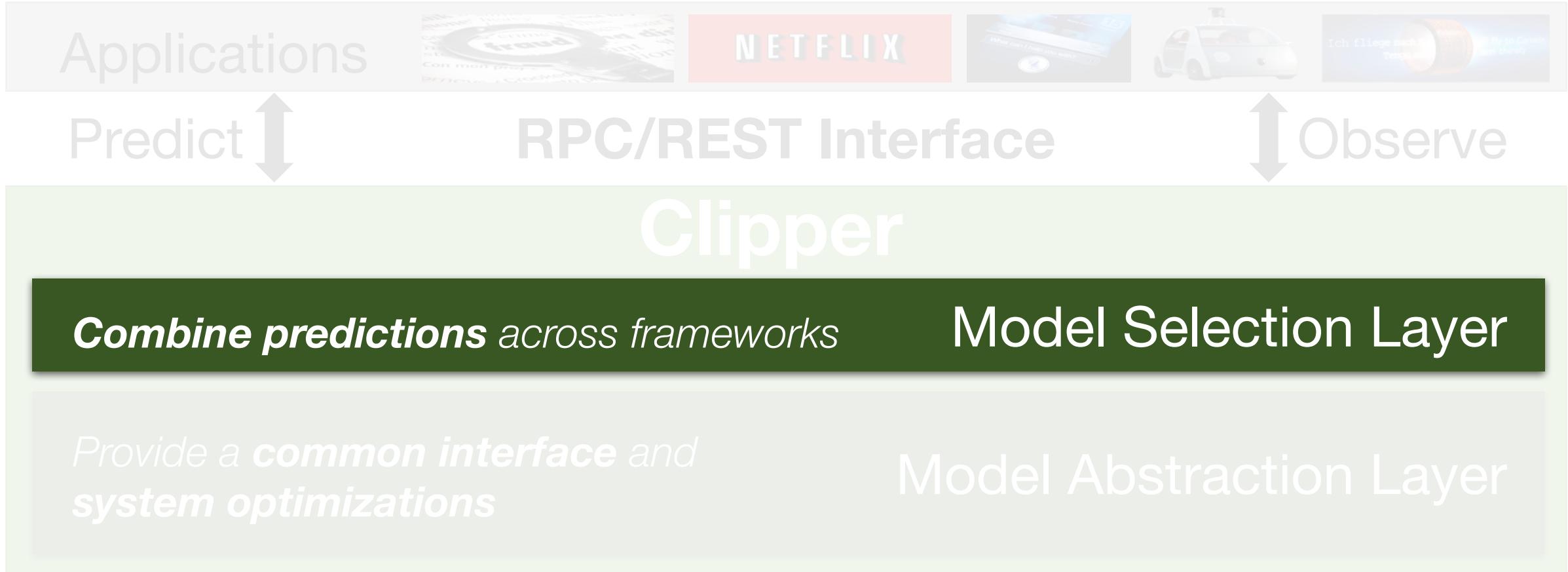


Decentralized system matches performance of centralized design.

# Clipper Architecture



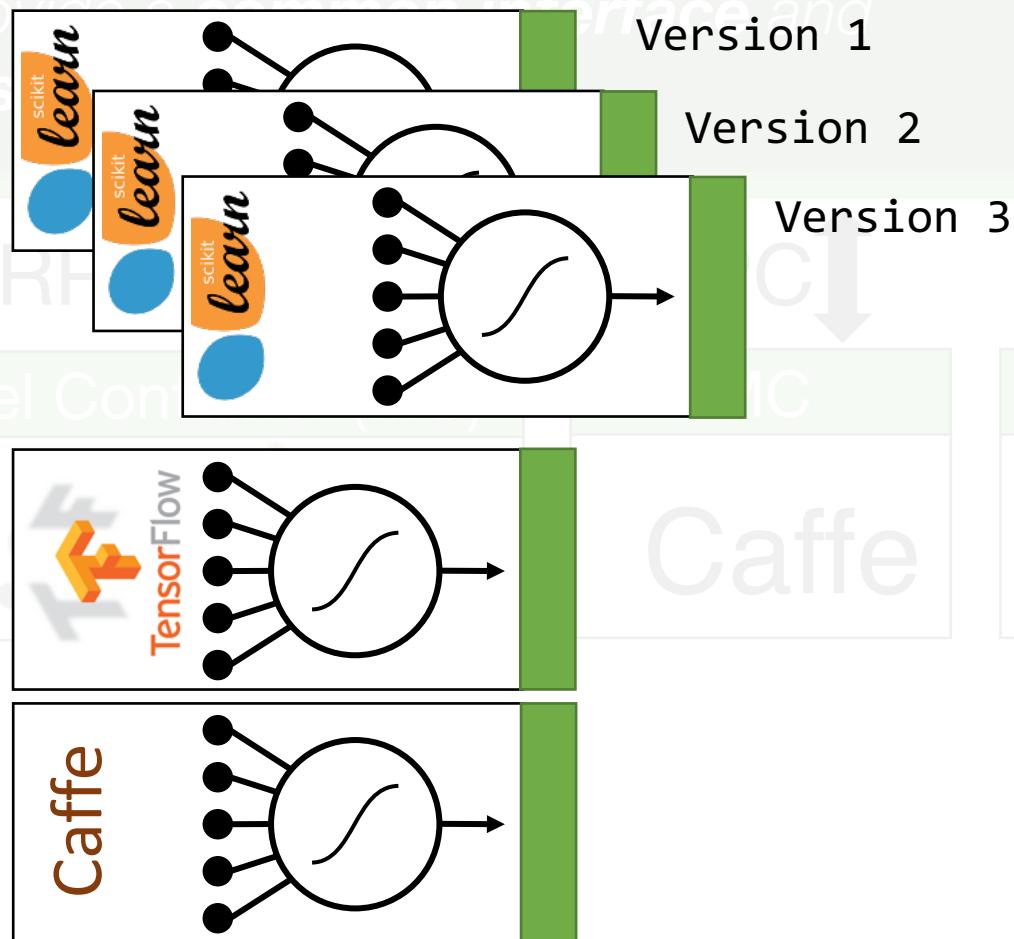
# Clipper Architecture



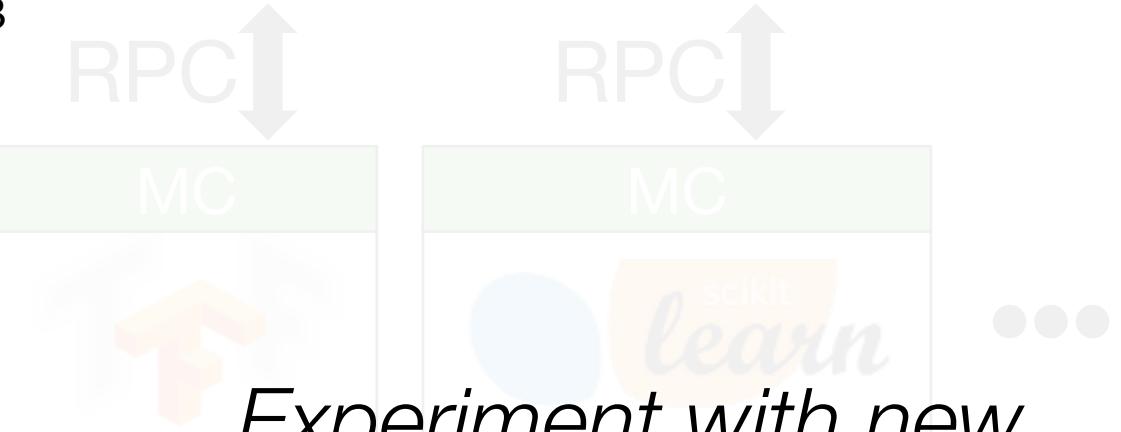
# Clipper

**Combine predictions** across frameworks

Model Selection Layer

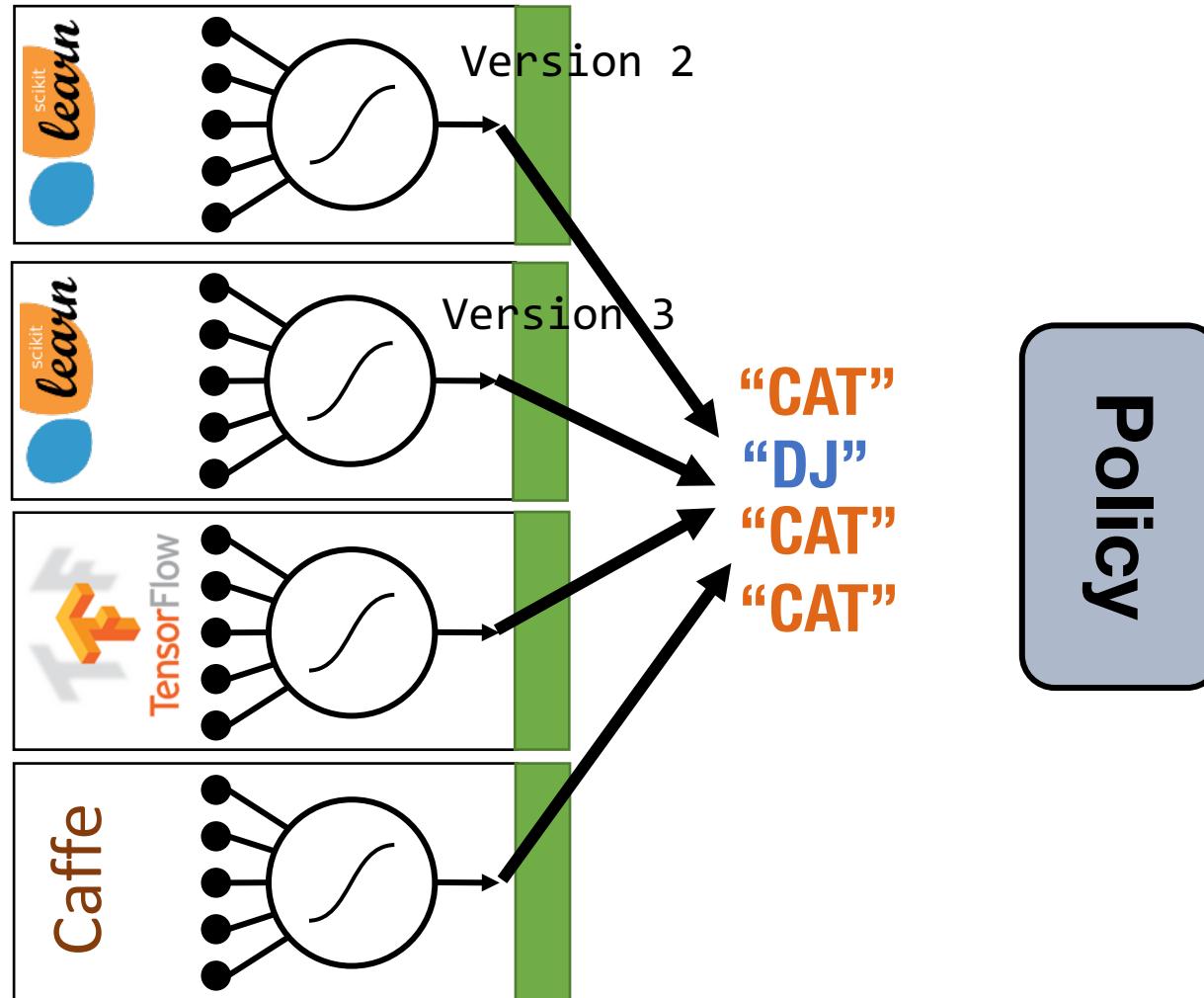


Model Abstraction Layer  
Periodic retraining



Experiment with new  
models and frameworks

# Selection Policy can Calibrate Confidence

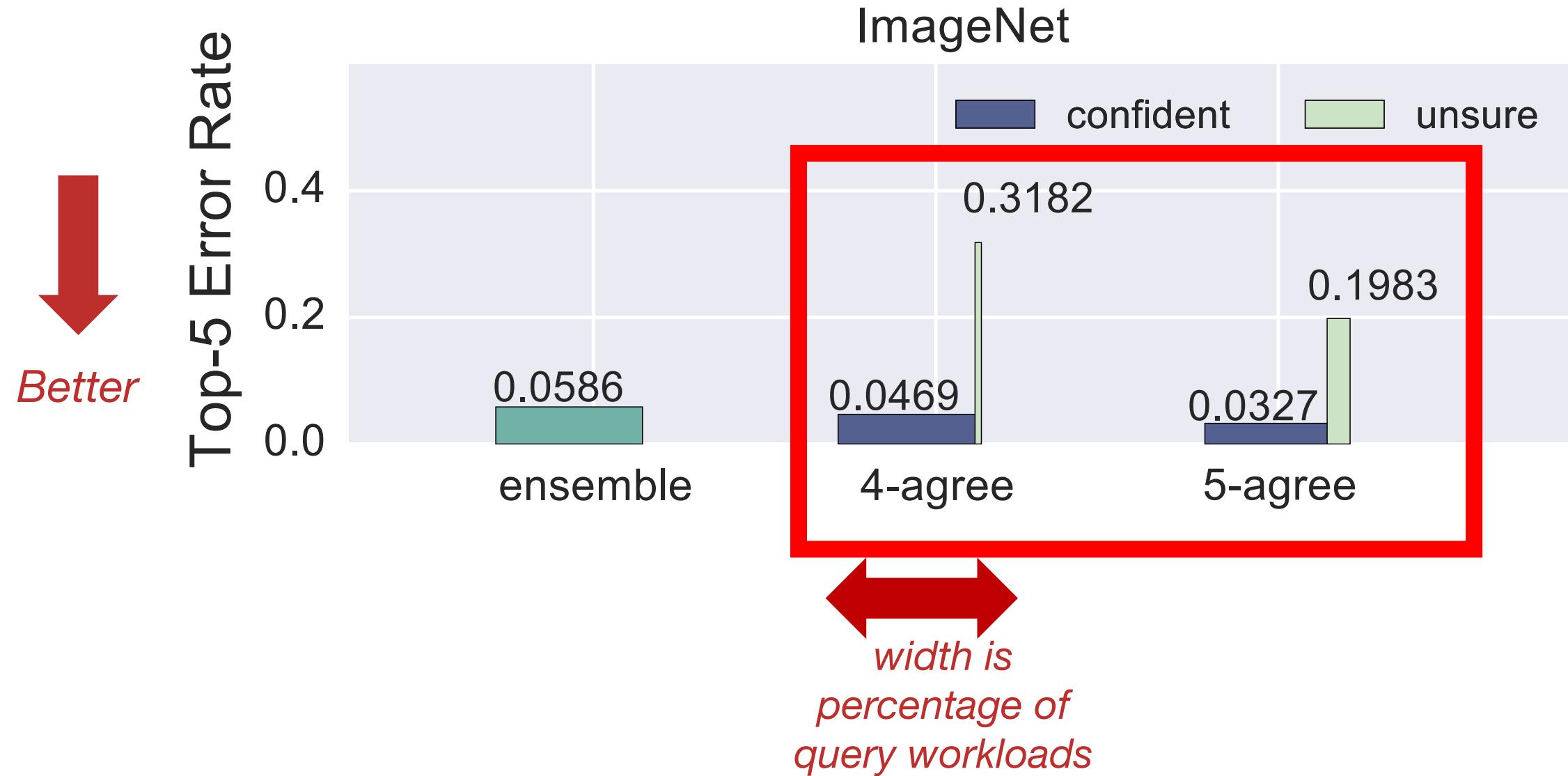


**"CAT"**  
**UNSURE**

# Selection Policy: Estimate confidence



# Selection Policy: Estimate confidence





# Open Research Questions

- Efficient execution of **complex model compositions**
  - Optimal batching to achieve end-to-end latency goals
- Automatic **model failure identification and correction**
  - Use anomaly detection techniques to identify model failures
- Prediction **serving on the edge**
  - Allowing models to span cloud and edge infrastructure



Low Latency  
Prediction Serving System  
[NSDI'17]

IDK Prediction  
Cascades

Learning to make fast predictions.  
[Work in Progress]

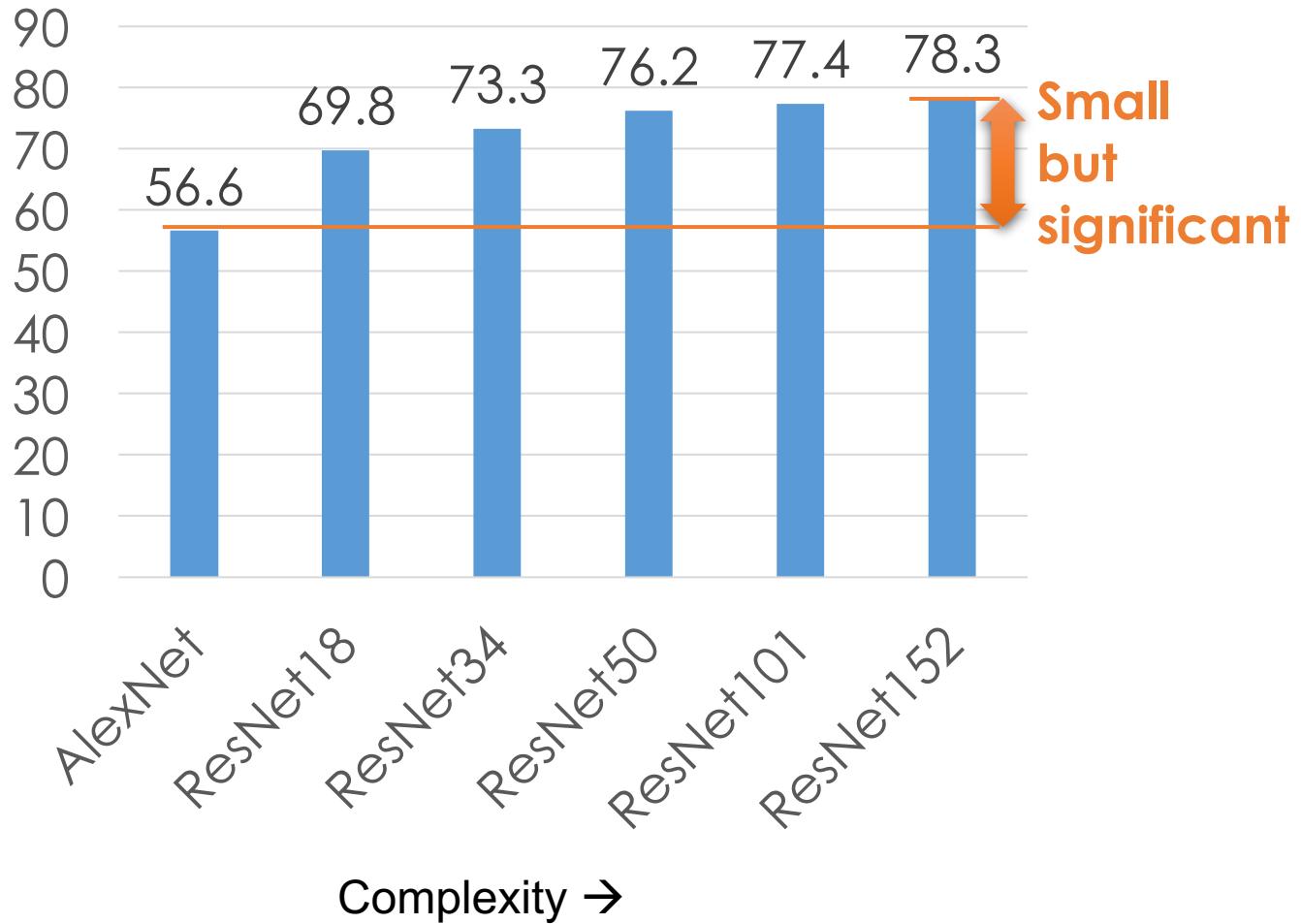


Low Latency  
Prediction Serving System  
[NSDI'17]

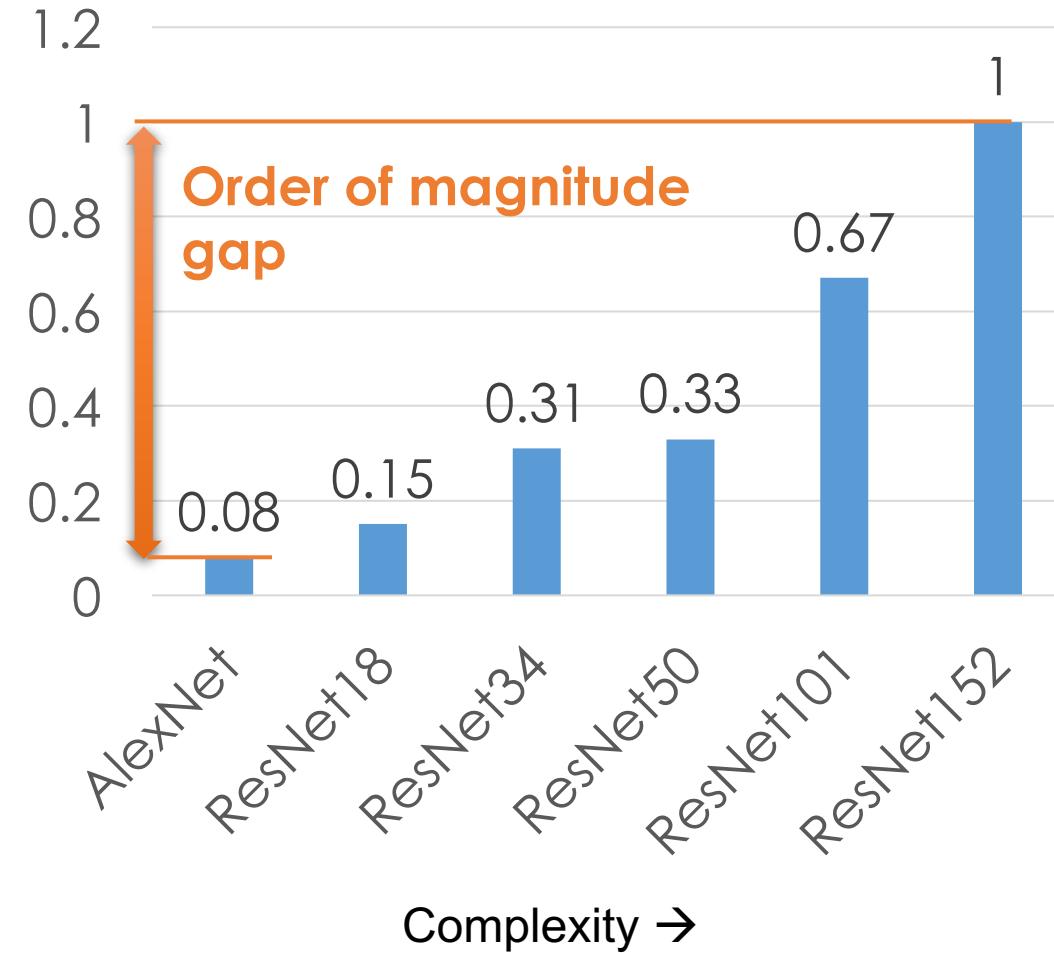
# IDK Prediction Cascades

Learning to make fast predictions.  
[Work in Progress]

### Accuracy



### Relative Cost



Model **costs are increasing** much  
**faster than** gains in **accuracy**.

# IDK Prediction Cascades

Simple models for simple tasks

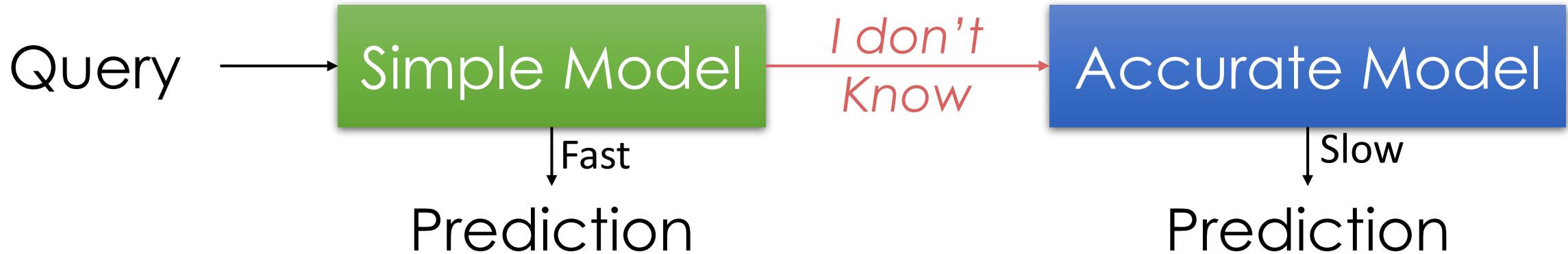
<https://arxiv.org/abs/1706.00885>



Xin  
Wang

Yika  
Luo

Daniel  
Crankshaw    Alexey  
Tumanov

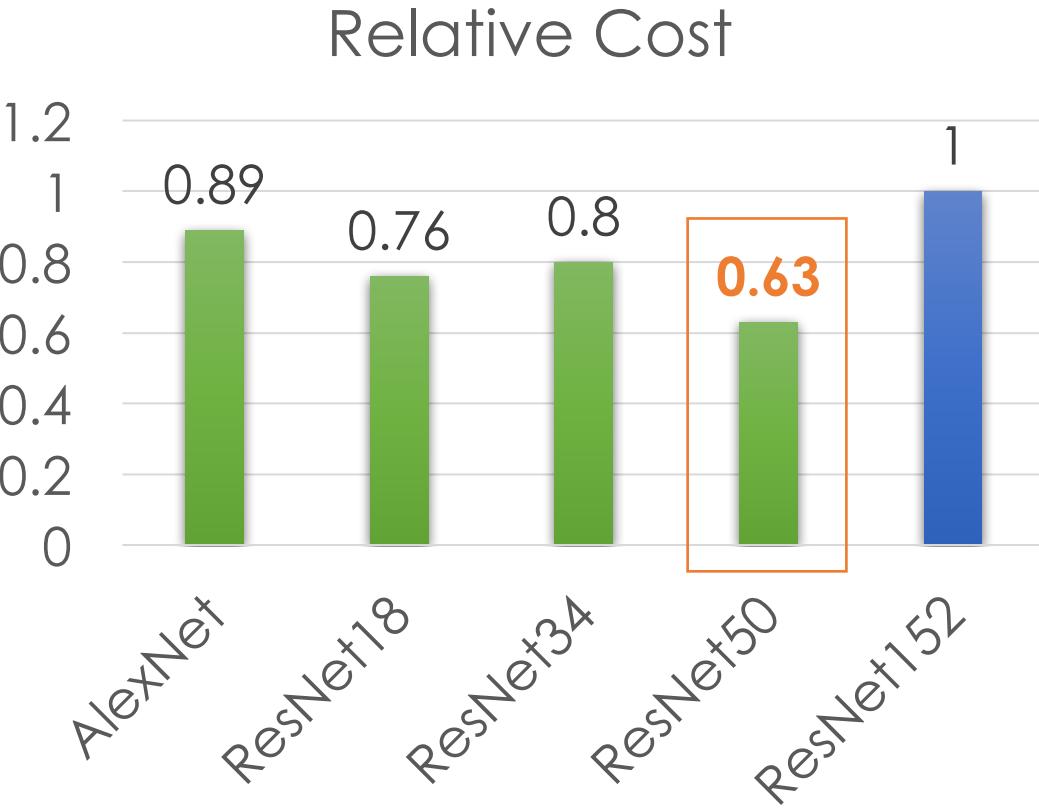
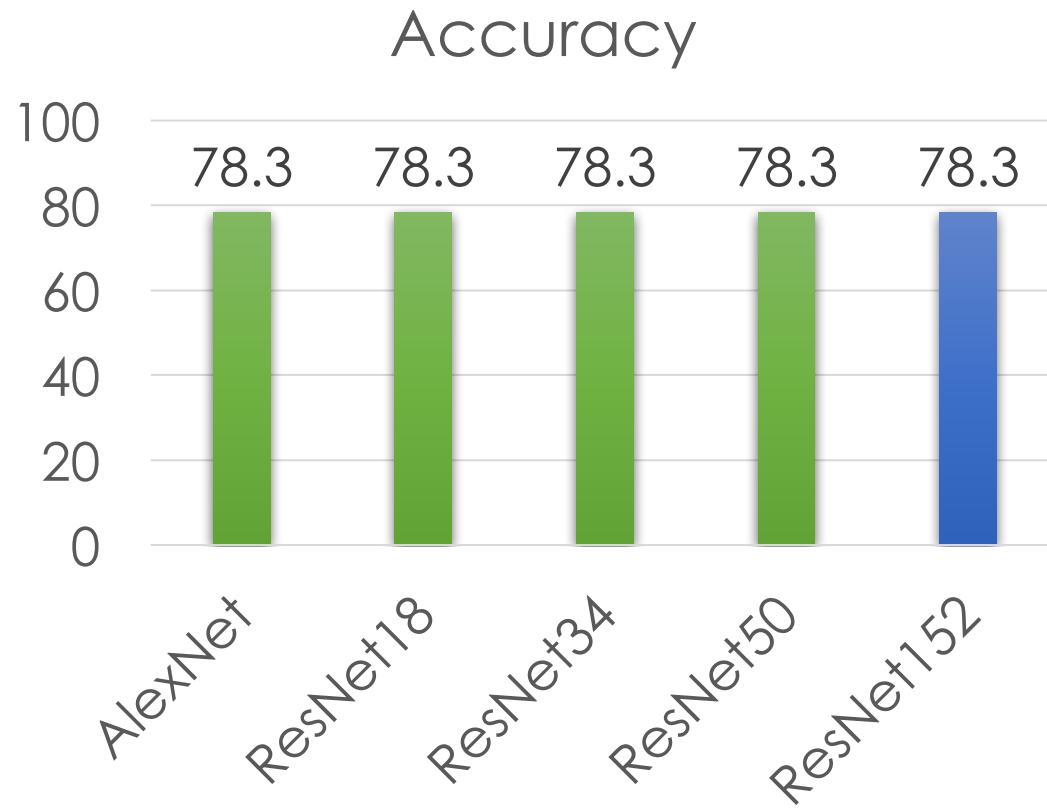


Combine **fast (inaccurate) models** with **slow (accurate) models** to maximize accuracy while reducing computational costs.

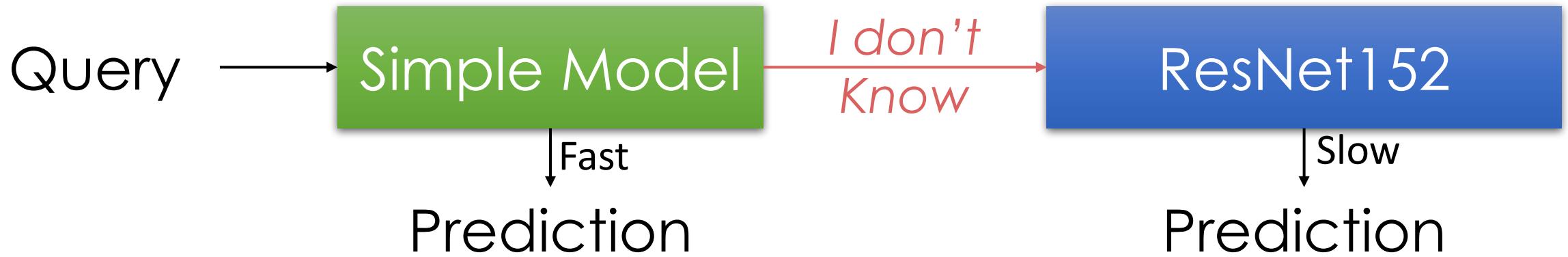
Query



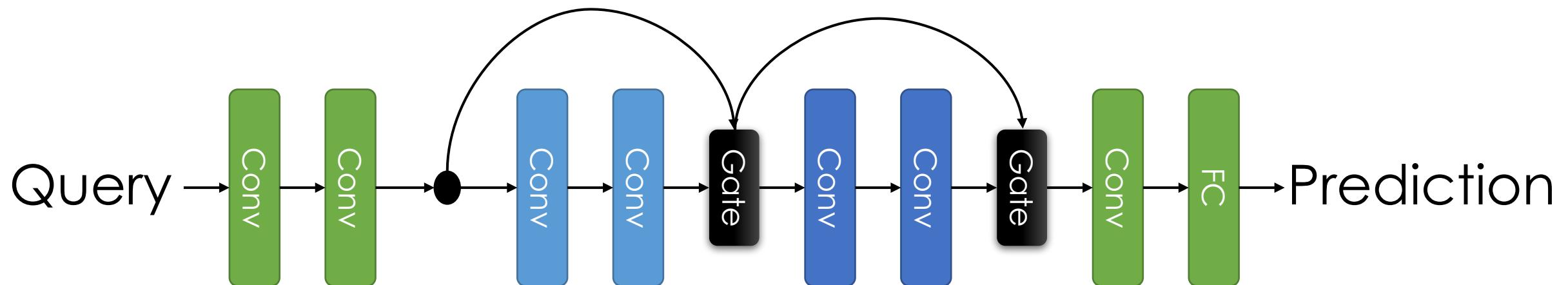
ResNet152

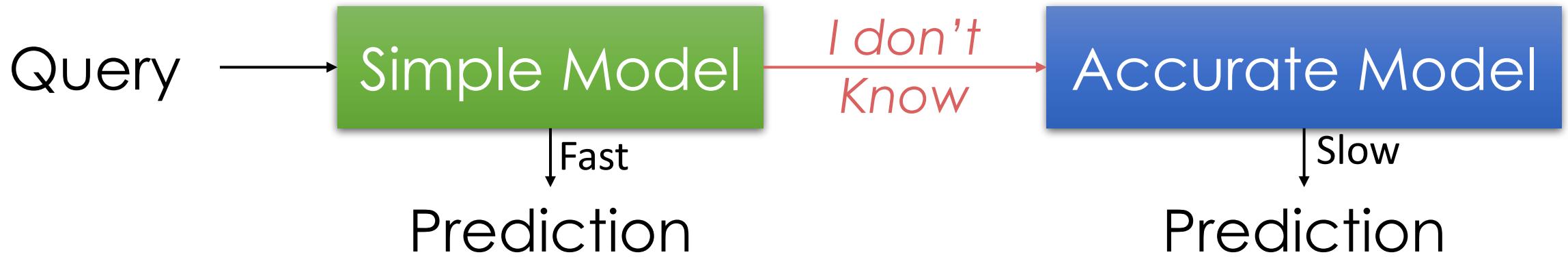


**37% reduction in runtime  
@ no loss in accuracy**

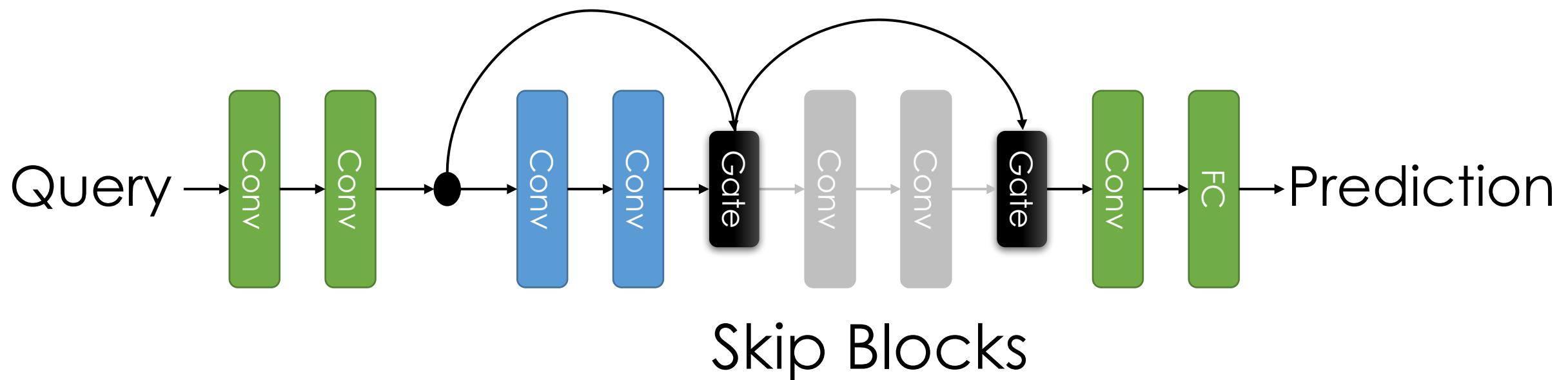


## ➤ Cascades within a Model

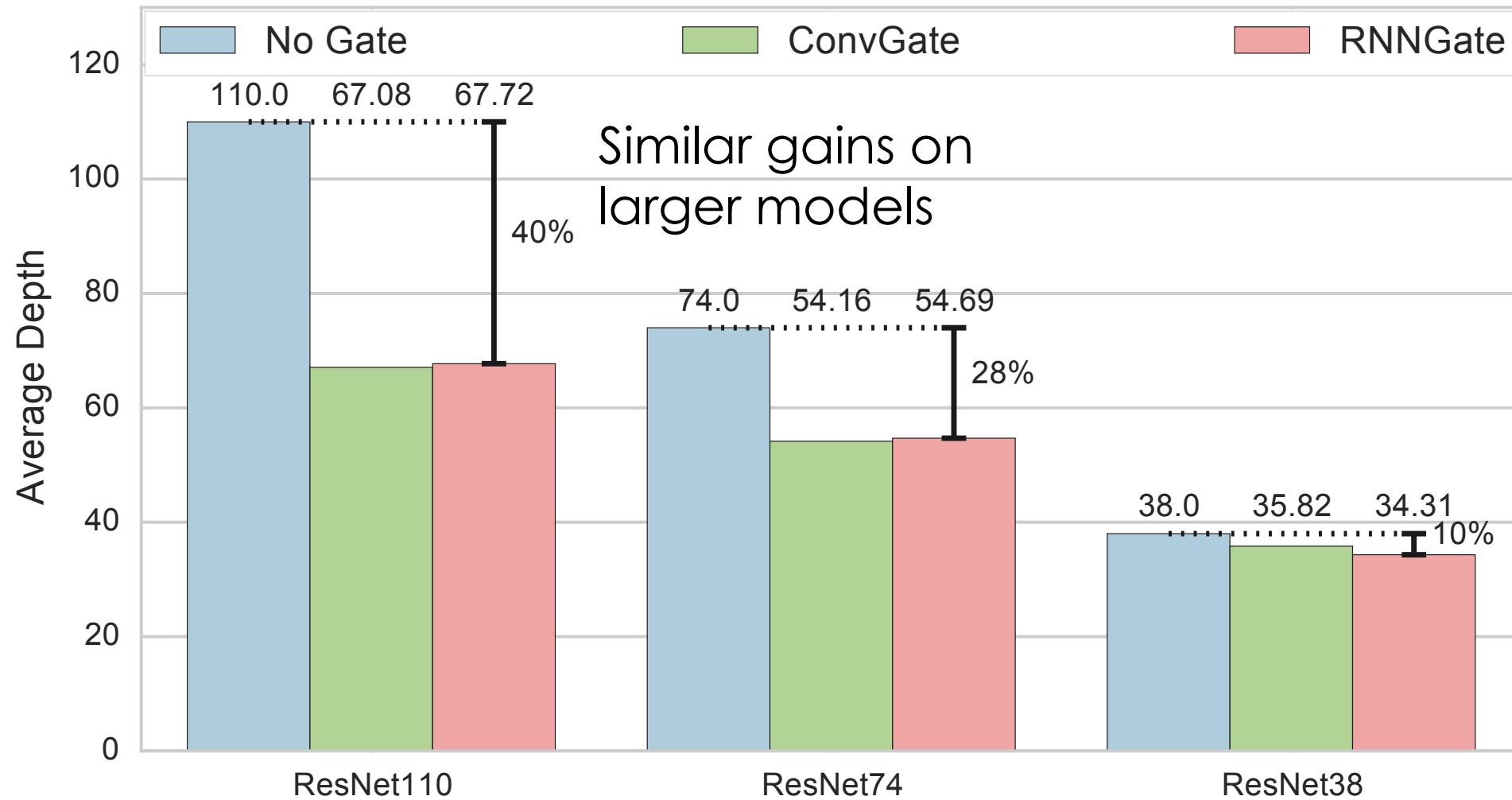




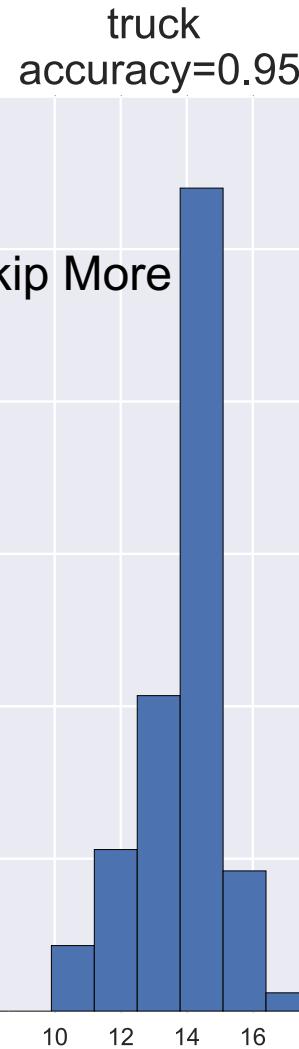
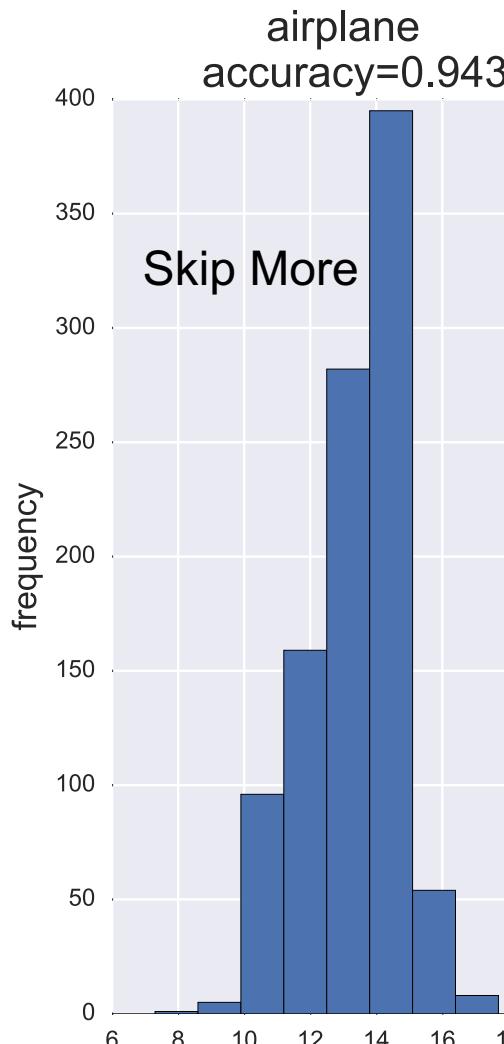
## ➤ Cascades within a Model



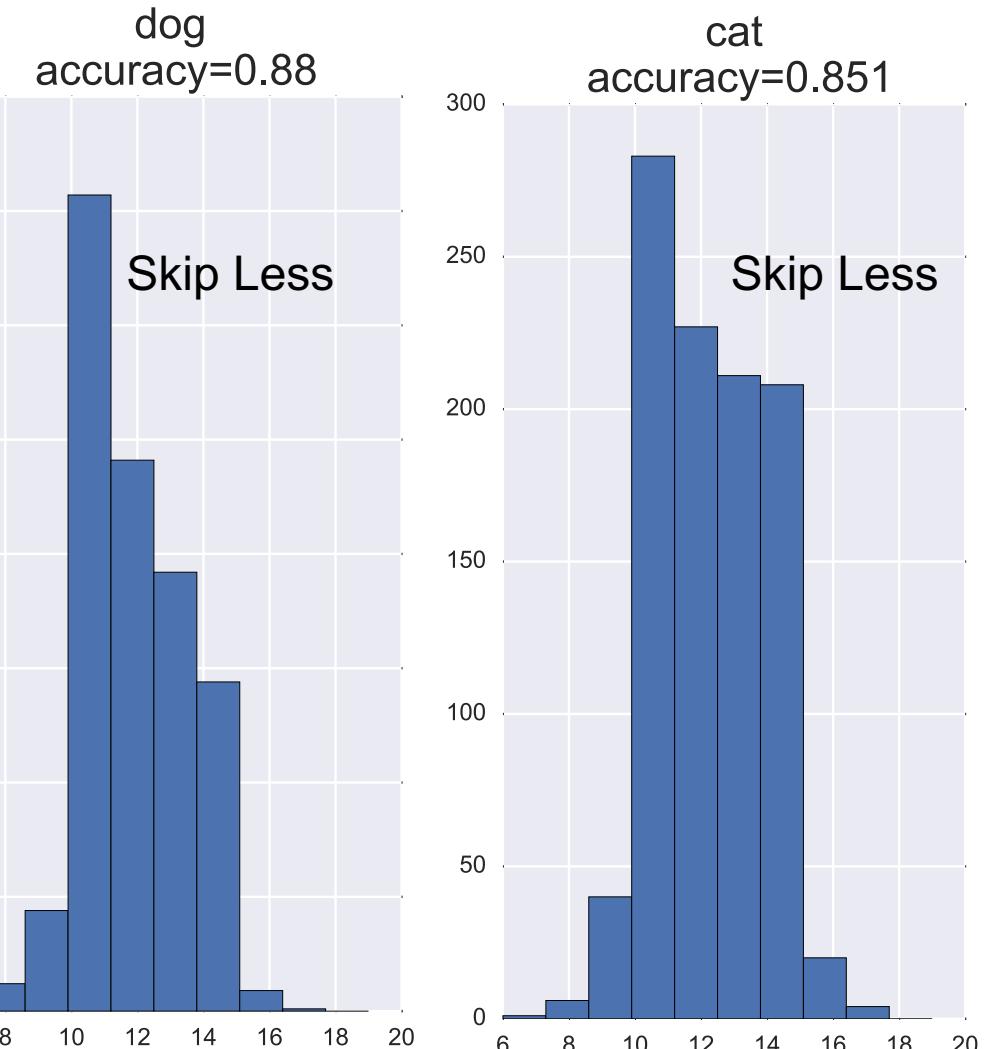
# Cascading reduces computational cost



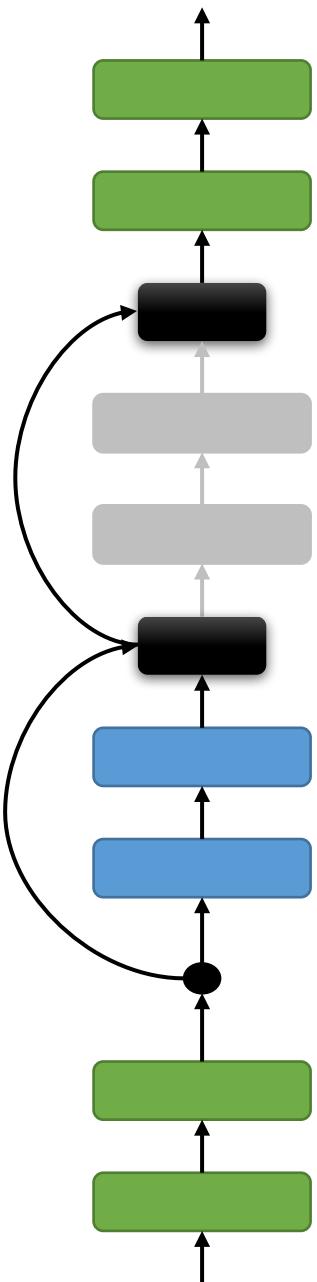
# Easy Images



# Difficult Images



Number of Layers Skipped



# Future Directions for Cascades

- Using **reinforcement learning** techniques to reduce gating costs
- **Query triage** during **load spikes** → forcing fractions of the network to go dark
- **Irregular execution** →
  - complicates batching
  - Issues for parallel execution

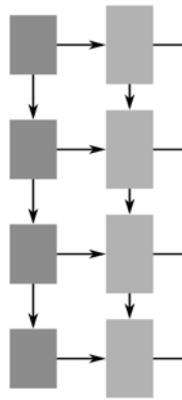


Low Latency  
Prediction Serving System  
[NSDI'17]

**IDK** Prediction Cascades  
Simple models for simple tasks  
[Work in Progress]

Other AI Systems Projects in RISE

**Jarvis**  
Managing the Machine  
Learning Lifecycle

 **Ray**  
Distributed Python for  
Reinforcement Learning



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

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