

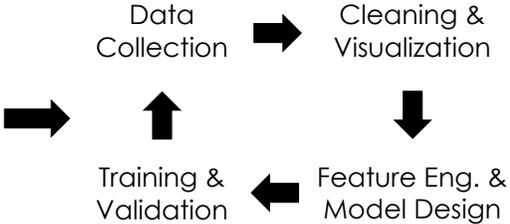
AI-Systems Prediction Serving

Joseph E. Gonzalez
Co-director of the RISE Lab
jegonzal@cs.berkeley.edu



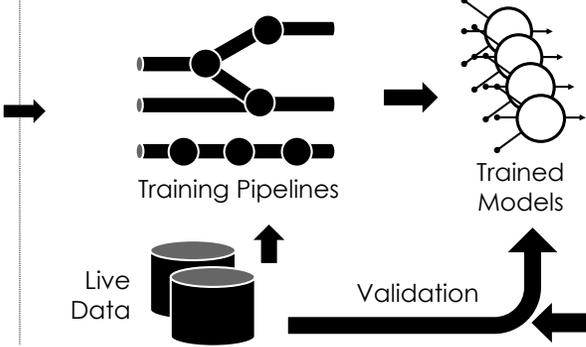
Machine Learning Lifecycle

Model Development



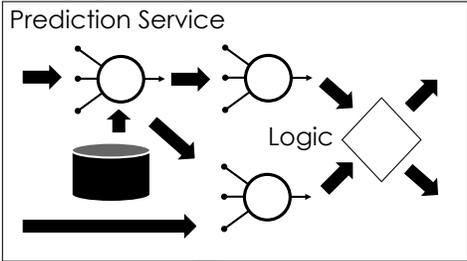
Data Scientist

Training

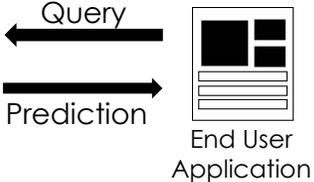


Data Engineer

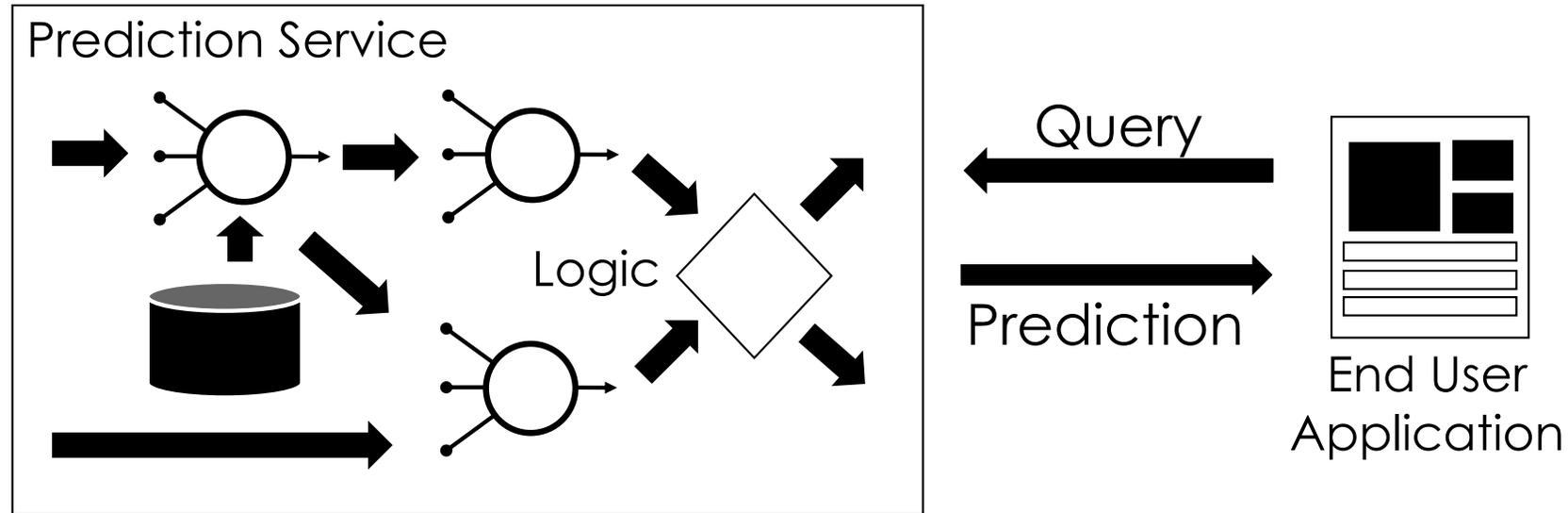
Inference



Data Engineer

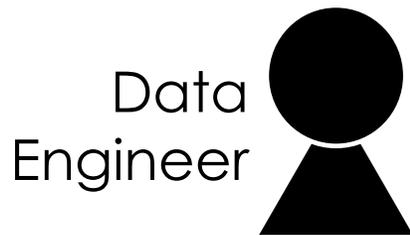


Inference



Feedback

Goal: make predictions in
~10ms under **bursty** load

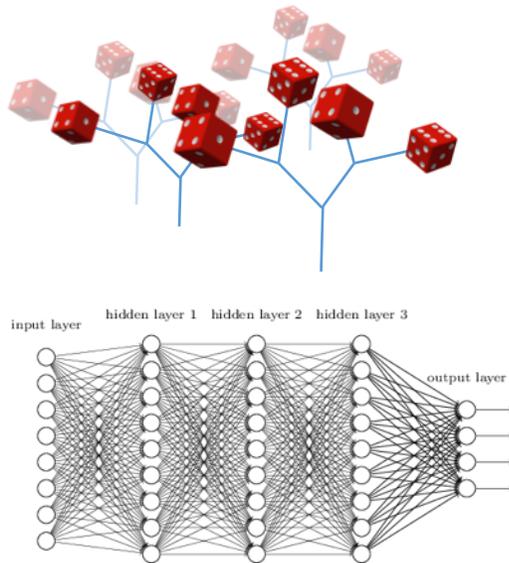


Complicated by **Deep Neural Networks**
→ New **ML Algorithms** and **Systems**

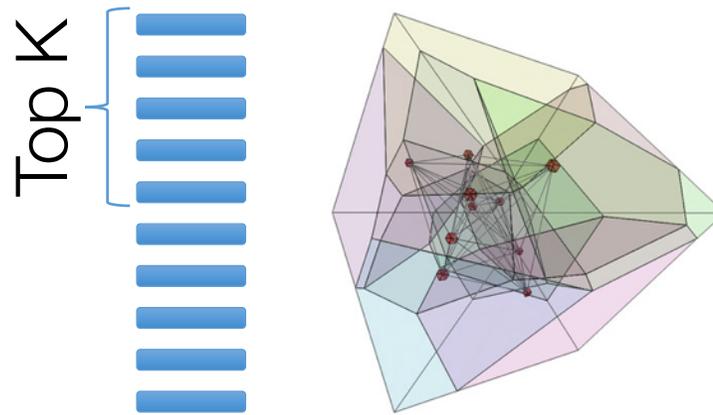
why is Inference challenging?

Need to render **low latency** (< 10ms) predictions for **complex**

Models



Queries



Features

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

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

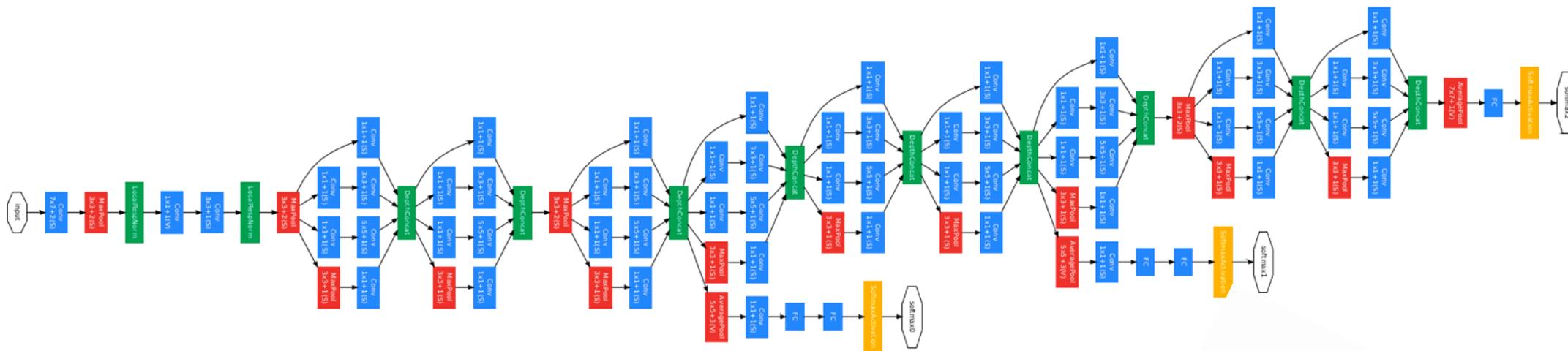
Basic Linear Models (Often High Dimensional)

- Common for **click prediction** and **text filter models** (spam)
- Query x encoded in sparse Bag-of-Words:
 - $x = \text{"The quick brown"} = \{(\text{"brown"}, 1), (\text{"the"}, 1), (\text{"quick"}, 1)\}$

- Rendering a prediction: $\mathbf{Predict}(x) = \sigma \left(\sum_{(w,c) \in x} \theta_w c \right)$

- θ is a large vector of weights for each possible word
 - or word combination (n-gram models) ...
- Optimizations?

Support low-latency, high-throughput serving workloads



Models getting more complex

- 10s of GFLOPs [1]

Deployed on critical path

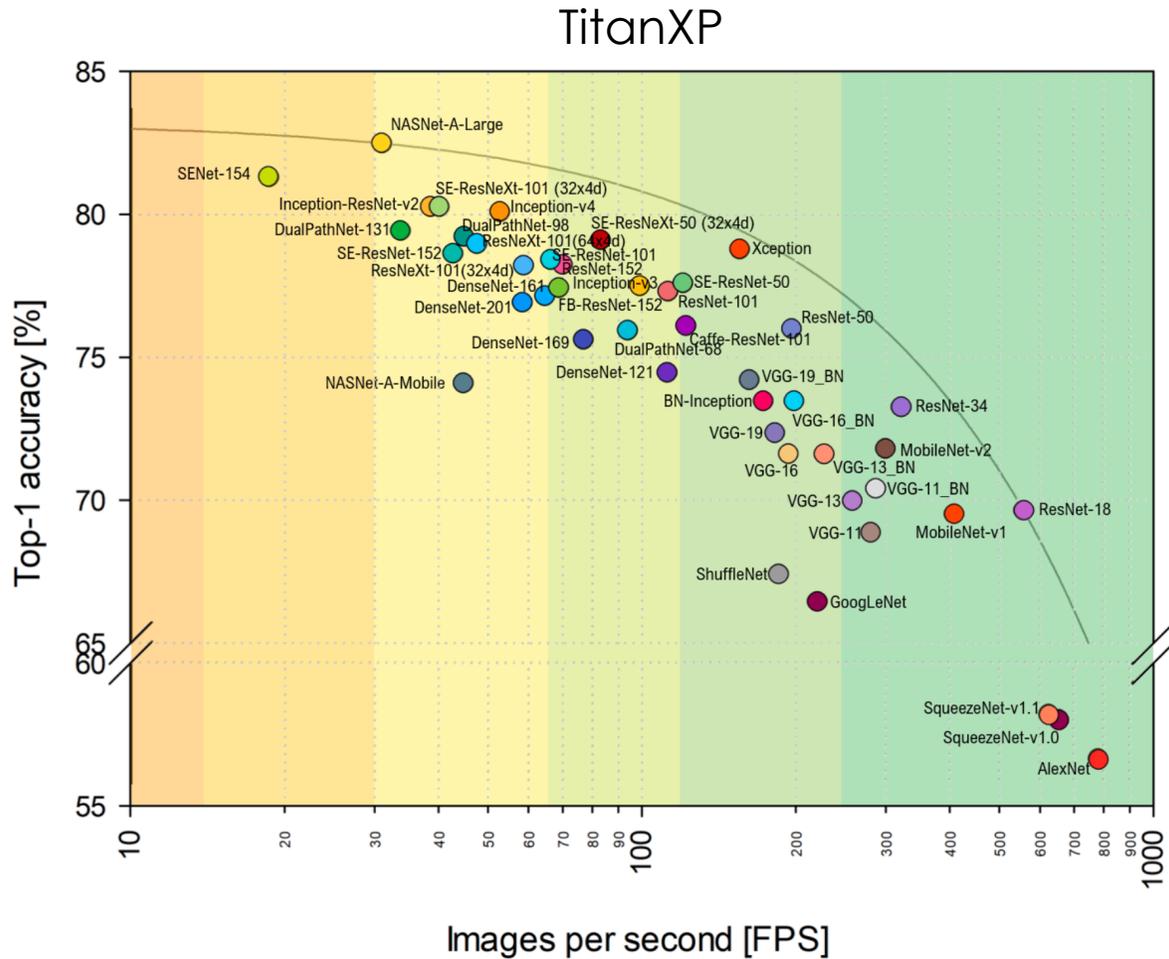
- Maintain SLOs under heavy load



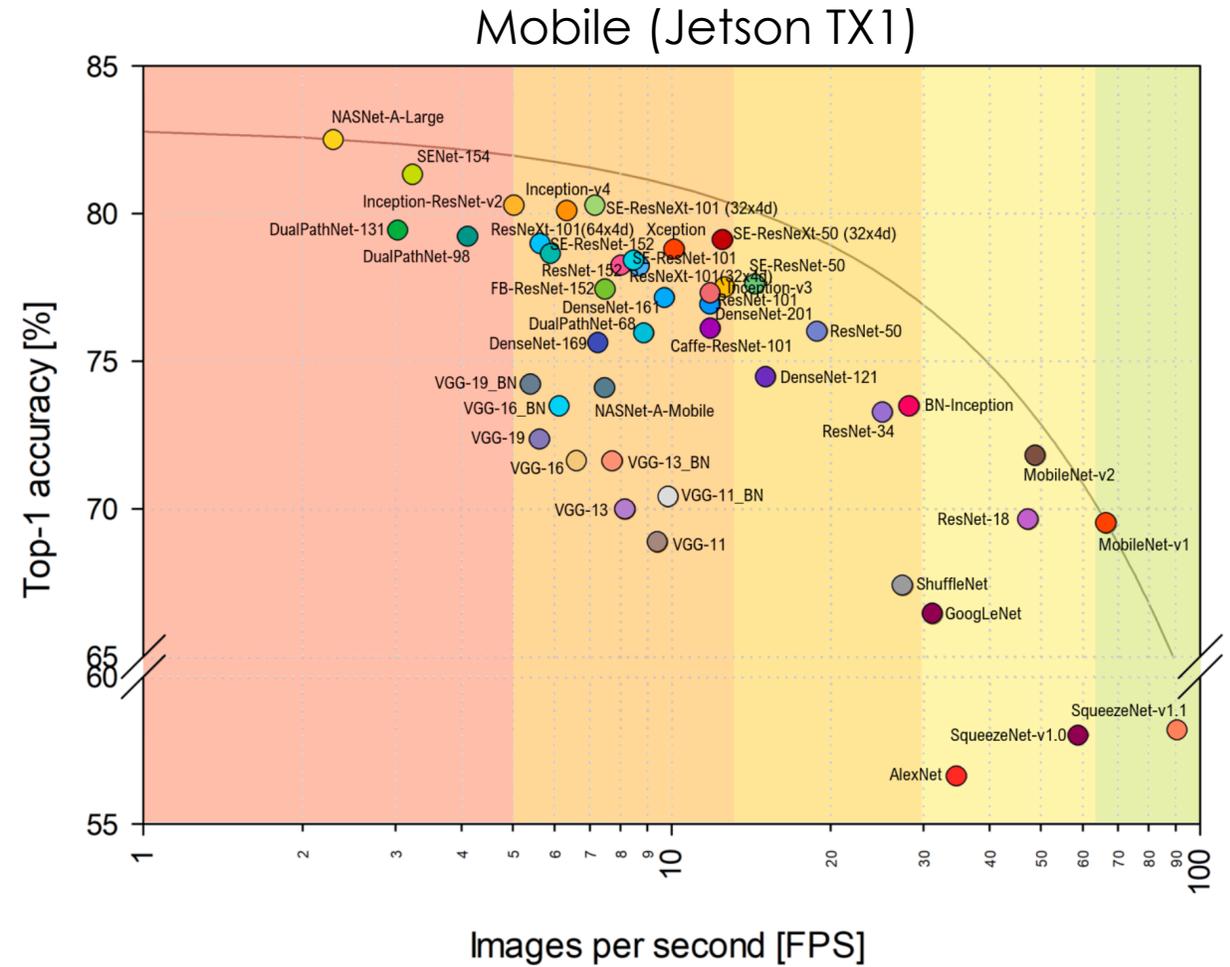
Using specialized hardware for predictions

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

Benchmark Analysis of Representative Deep Neural Network Architectures



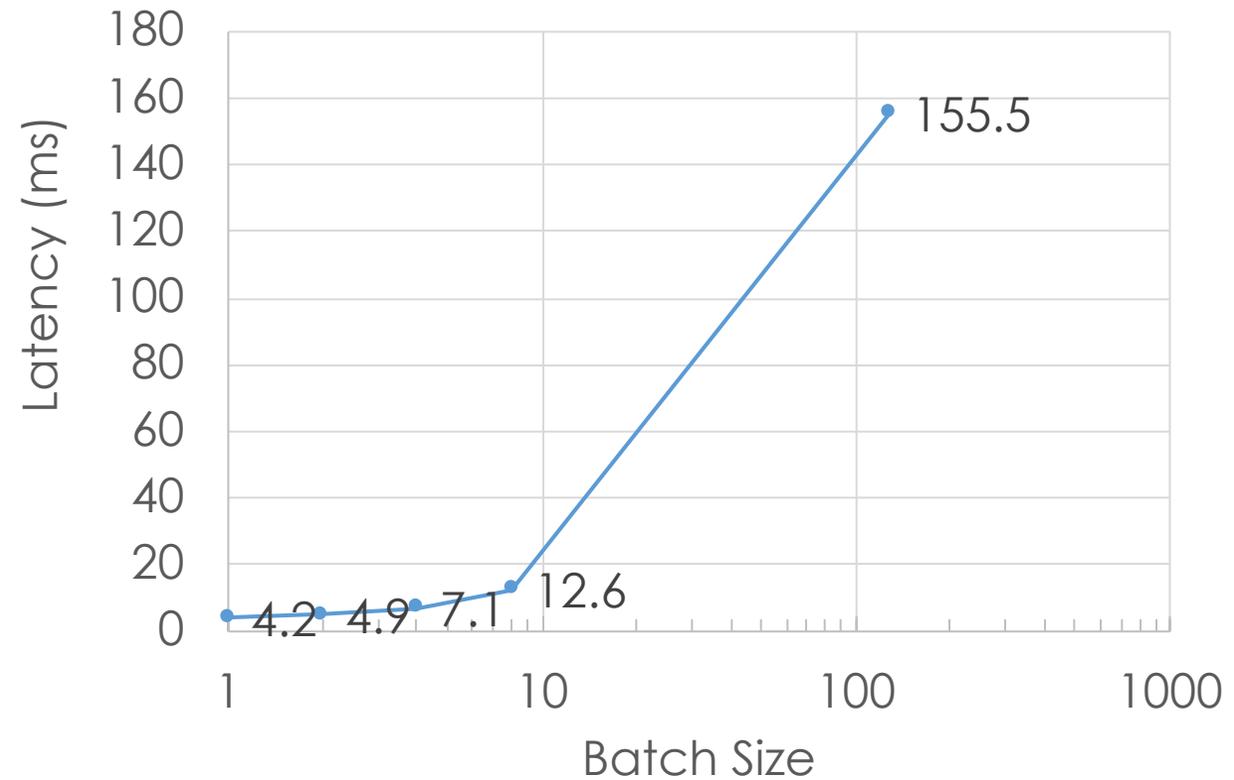
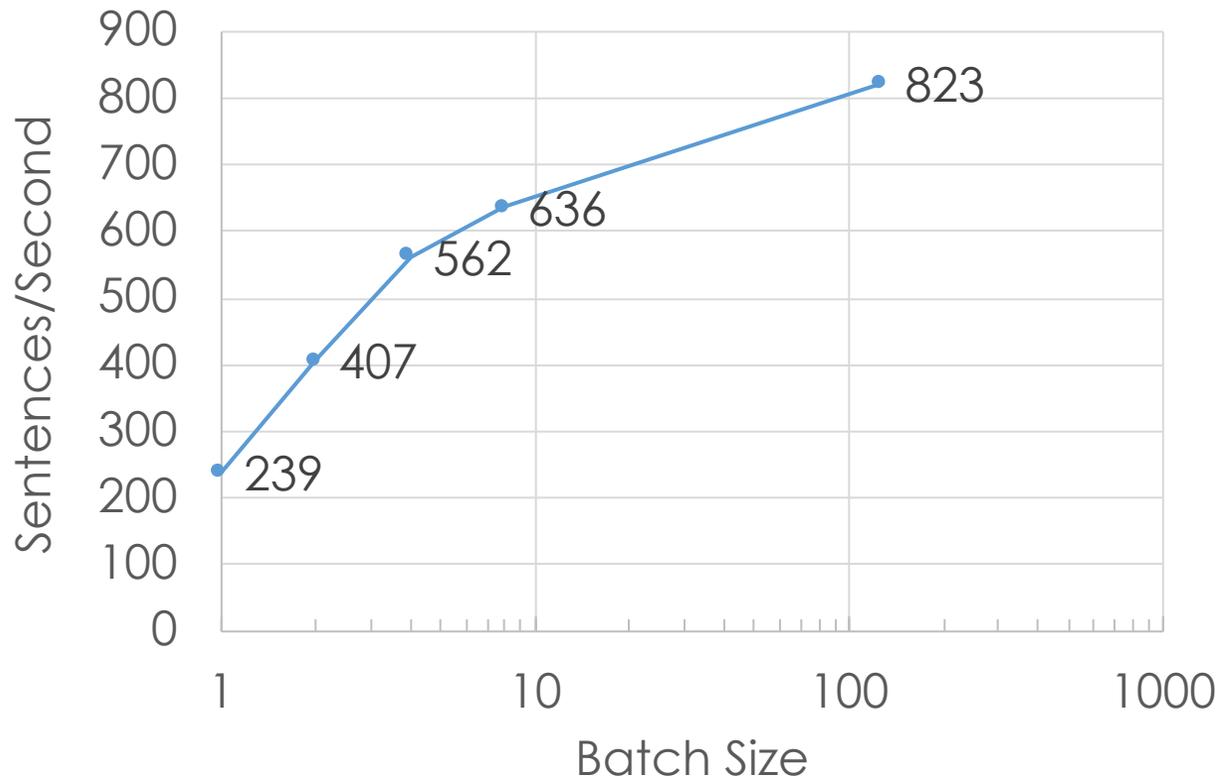
(a)



(b)

FIGURE 3: Top-1 accuracy vs. number of images processed per second (with batch size 1) using the Titan Xp (a) and Jetson TX1 (b).

BERT-Large on a V100 (~\$10K)

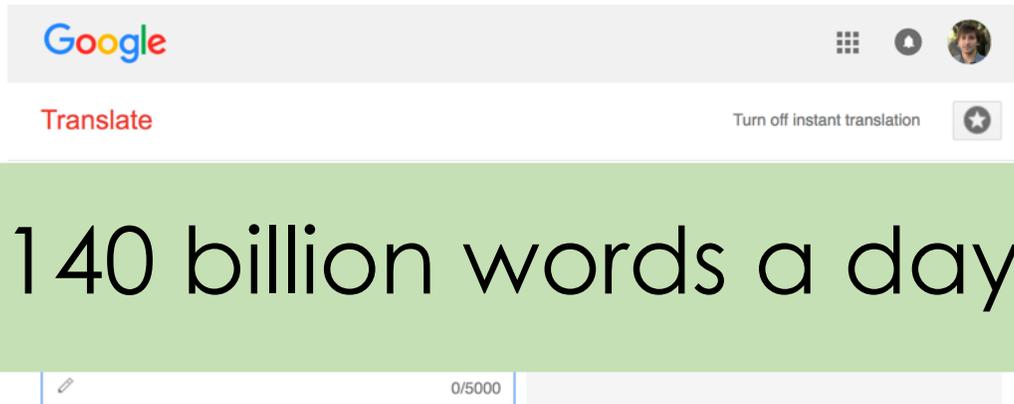


Results included Mixed precision optimizations!

Numbers obtained from: <https://developer.nvidia.com/deep-learning-performance-training-inference>

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

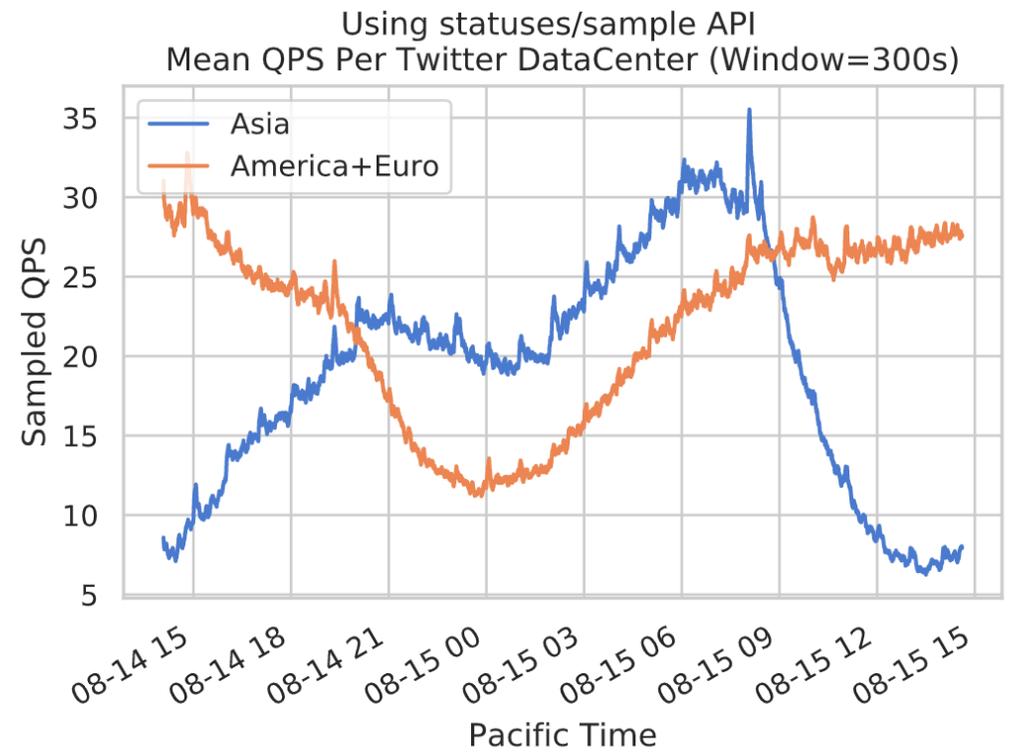
*"If each of the **world's Android phones** used the new Google voice search for just **three minutes a day**, these engineers realized, the company would **need twice as many data centers.**"*
– Wired

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

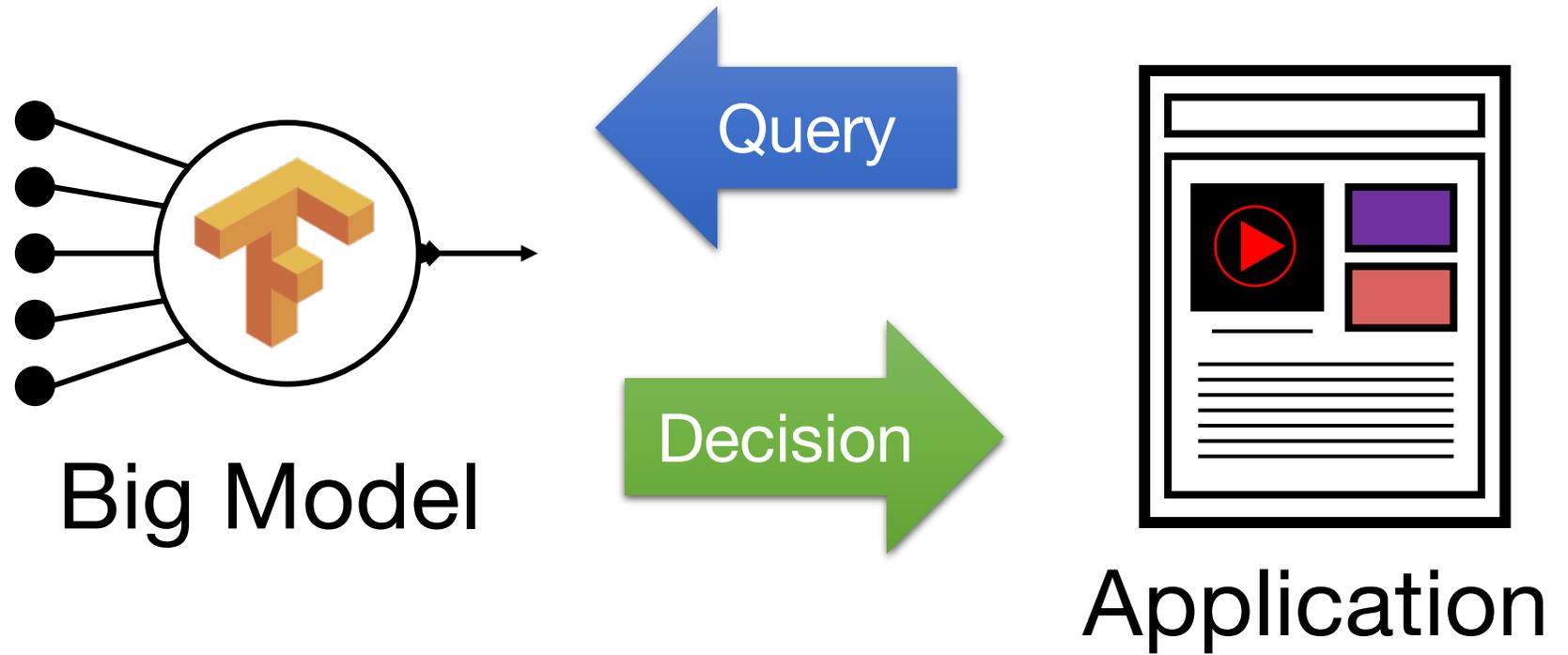
[1] <https://www.nytimes.com/2016/12/14/magazine/the-great-ai-awakening.html>

Other Challenges?

- **Bursty load** →
 - overprovision resources →
 - **expensive**
 - TPU reports 28% utilization of vector units in production
 - Solutions
 - statistical multiplexing → difficult → why?
 - could try to predict arrival process → difficult (impossible?)!
- **Versioning** and **testing models**
- **Prediction pipelines** → more on this soon



Inference



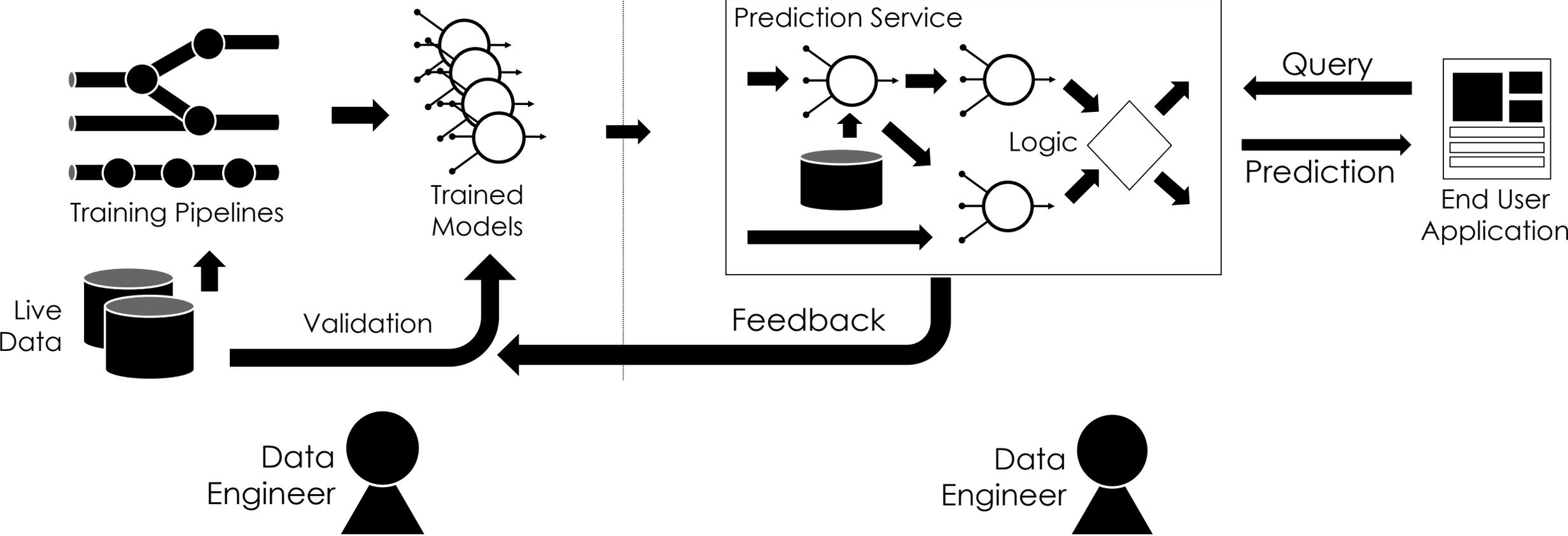
Two Approaches

- **Offline:** Pre-Materialize Predictions
- **Online:** Compute Predictions on the fly

Pre-materialized Predictions

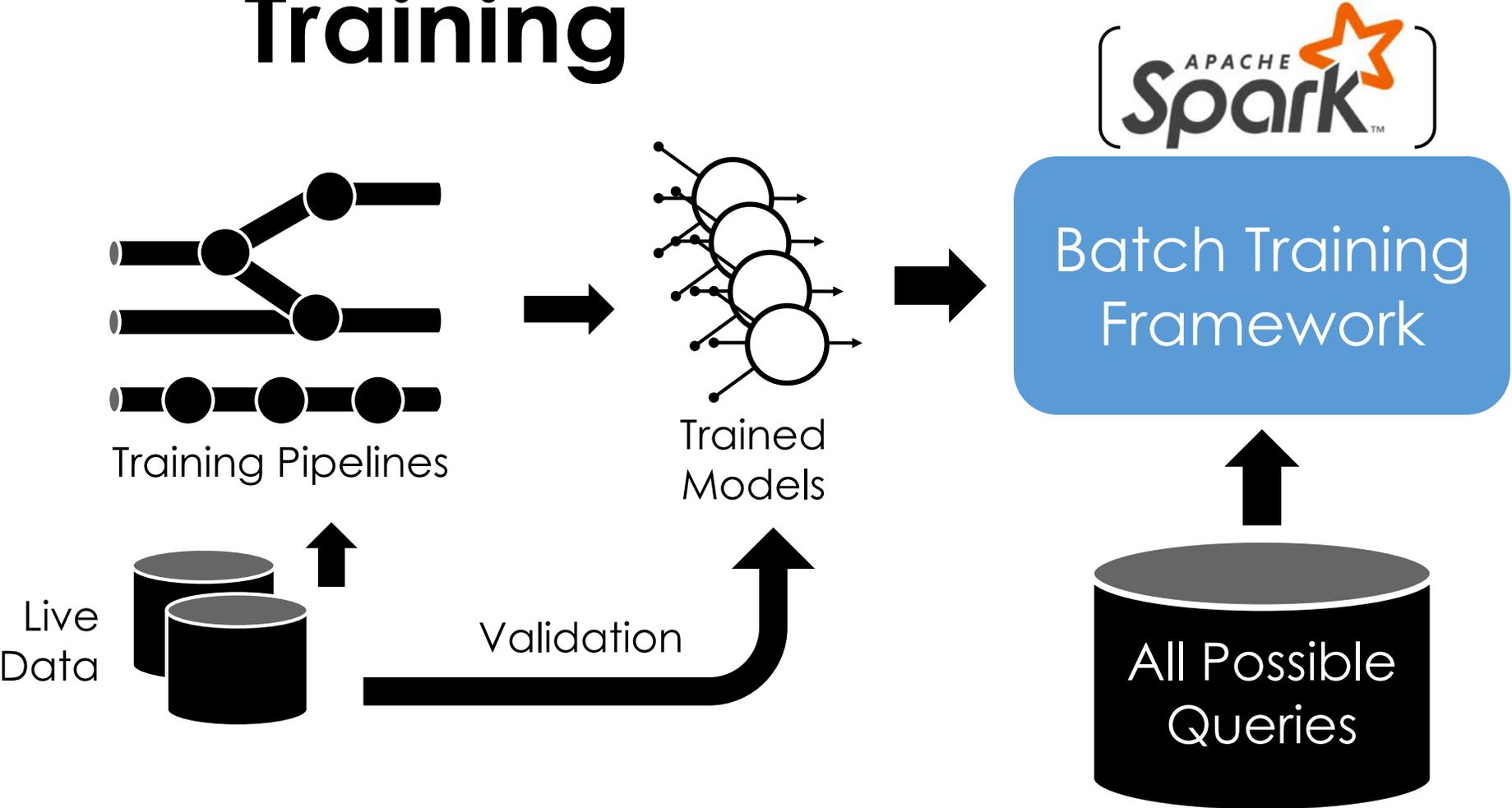
Training

Inference



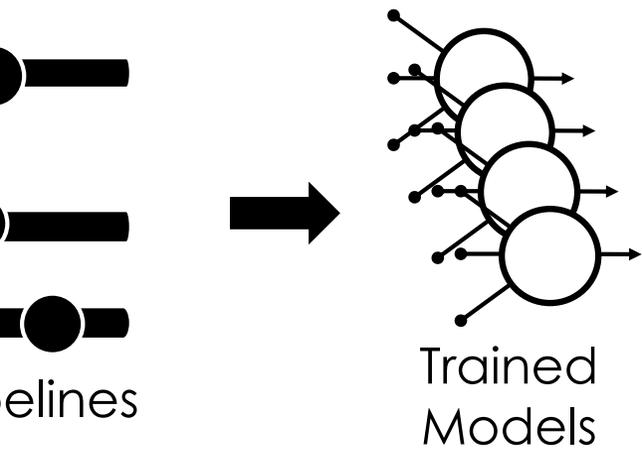
Pre-materialized Predictions

Training



Pre-materialized Predictions

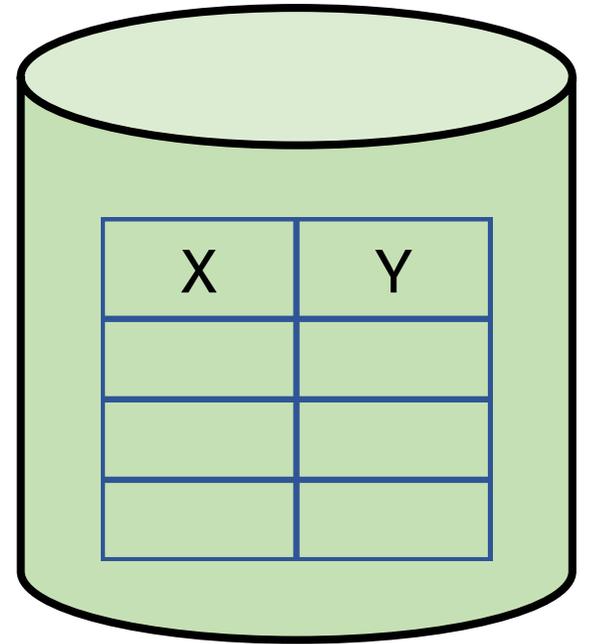
aining



Batch Training Framework

(Scoring)

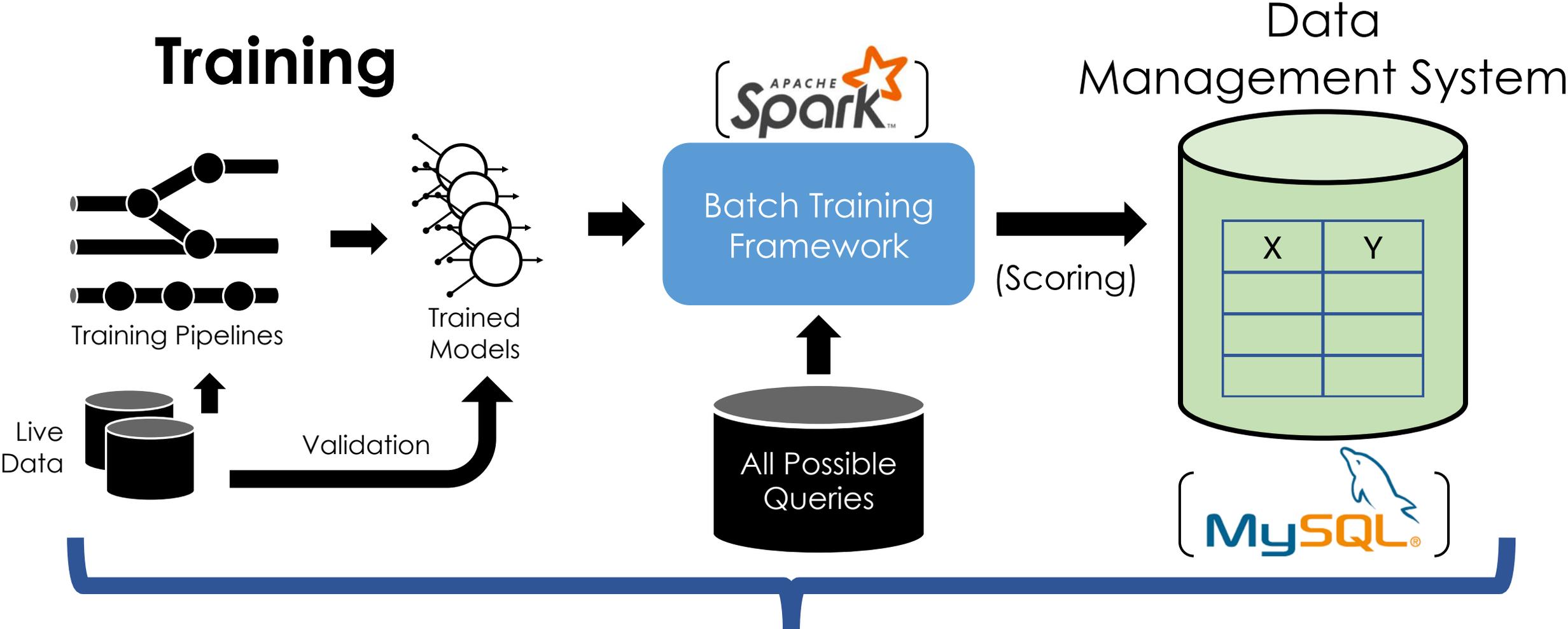
Data Management System



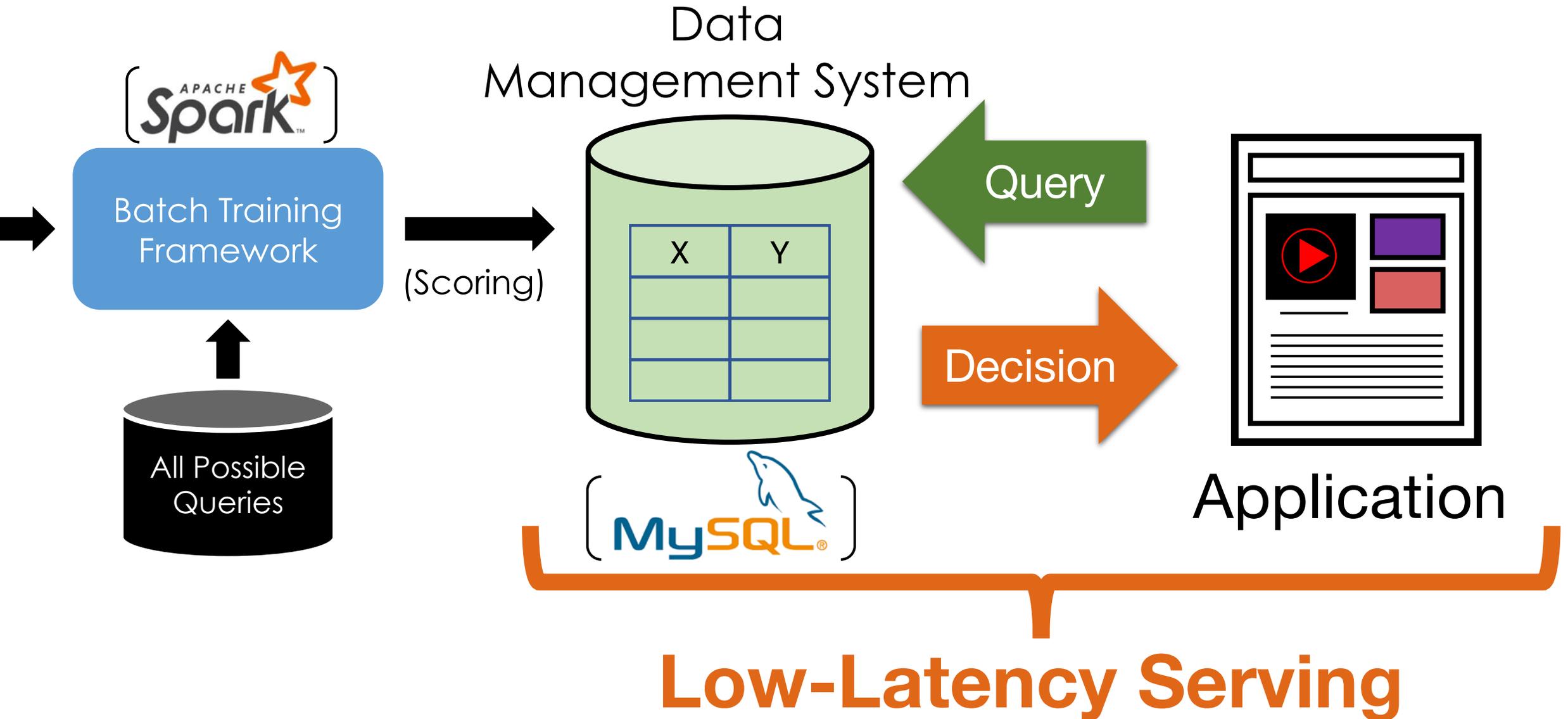
Validation

Pre-materialized Predictions

Training



Serving Pre-materialized Predictions



Serving Pre-materialized Predictions

Advantages:

- Leverage existing **data serving** and **model training** infrastructure
- **Batch processing** improves hardware perf.
- **Indexing** support for complex queries
 - Find all *Pr("cute")* dresses **where price < \$20**
- More **predictable** performance

Low-Latency Serving

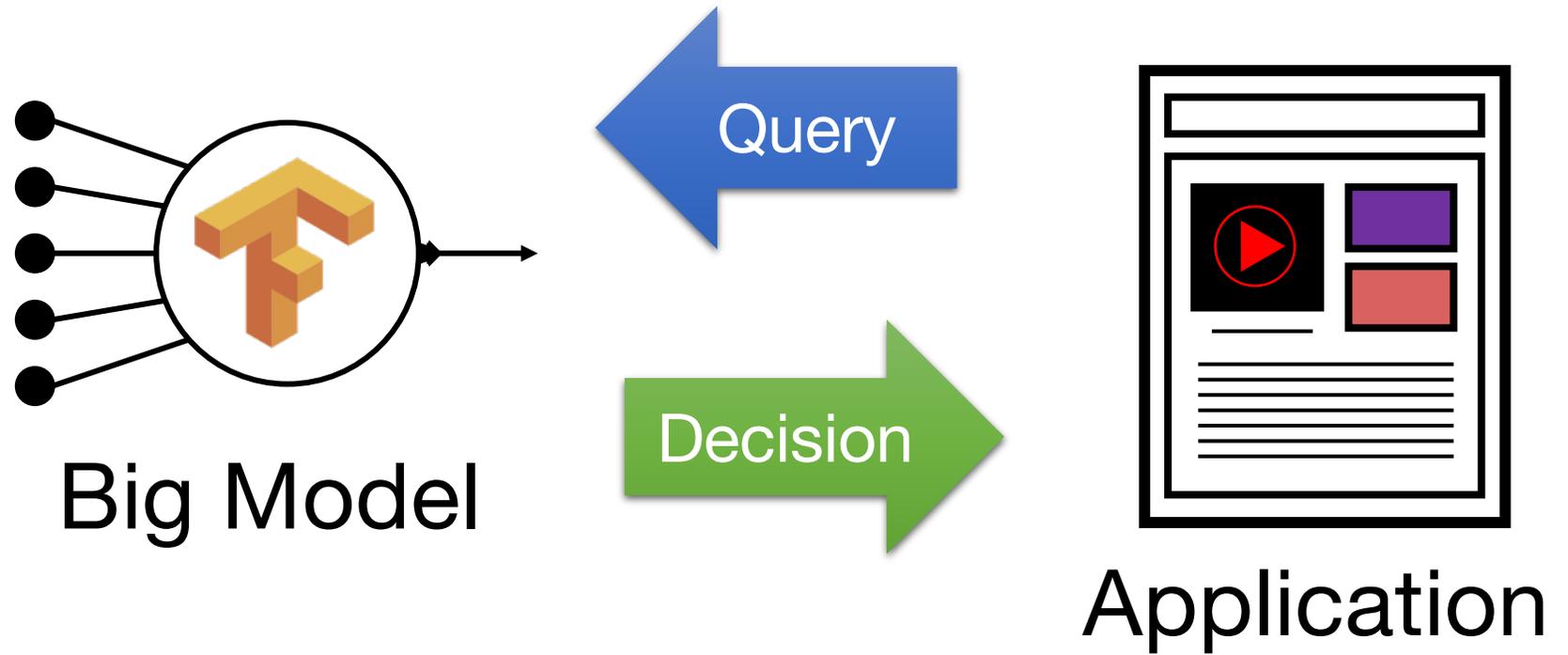
Serving Pre-materialized Predictions

Problems:

- Requires full set of **queries ahead of time**
- Small and **bounded input domain**
- Requires substantial **computation** and **space**
 - Example: *scoring all content for all customers!*
- Costly update → rescore everything!

Low-Latency Serving

Inference

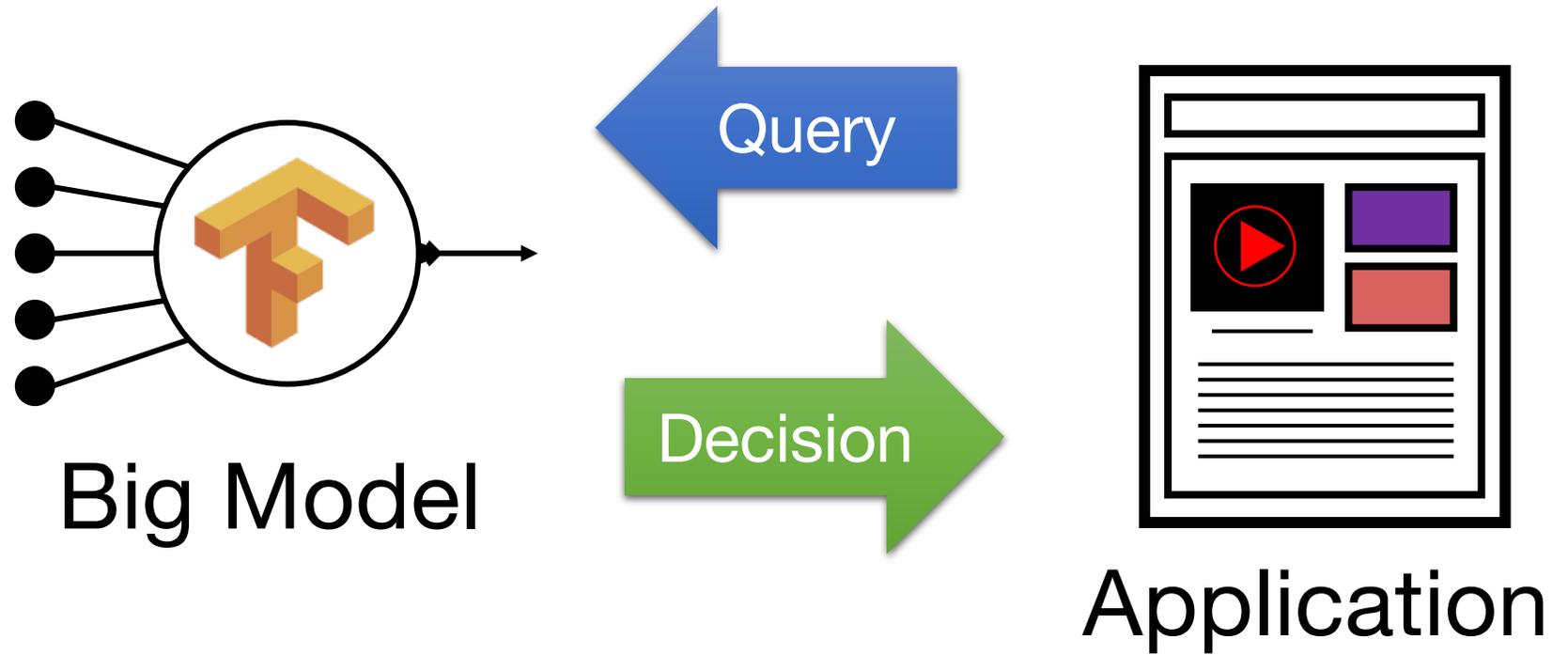


Two Approaches

➤ **Offline:** Pre-Materialize Predictions

➤ **Online:** Compute Predictions on the fly

Inference

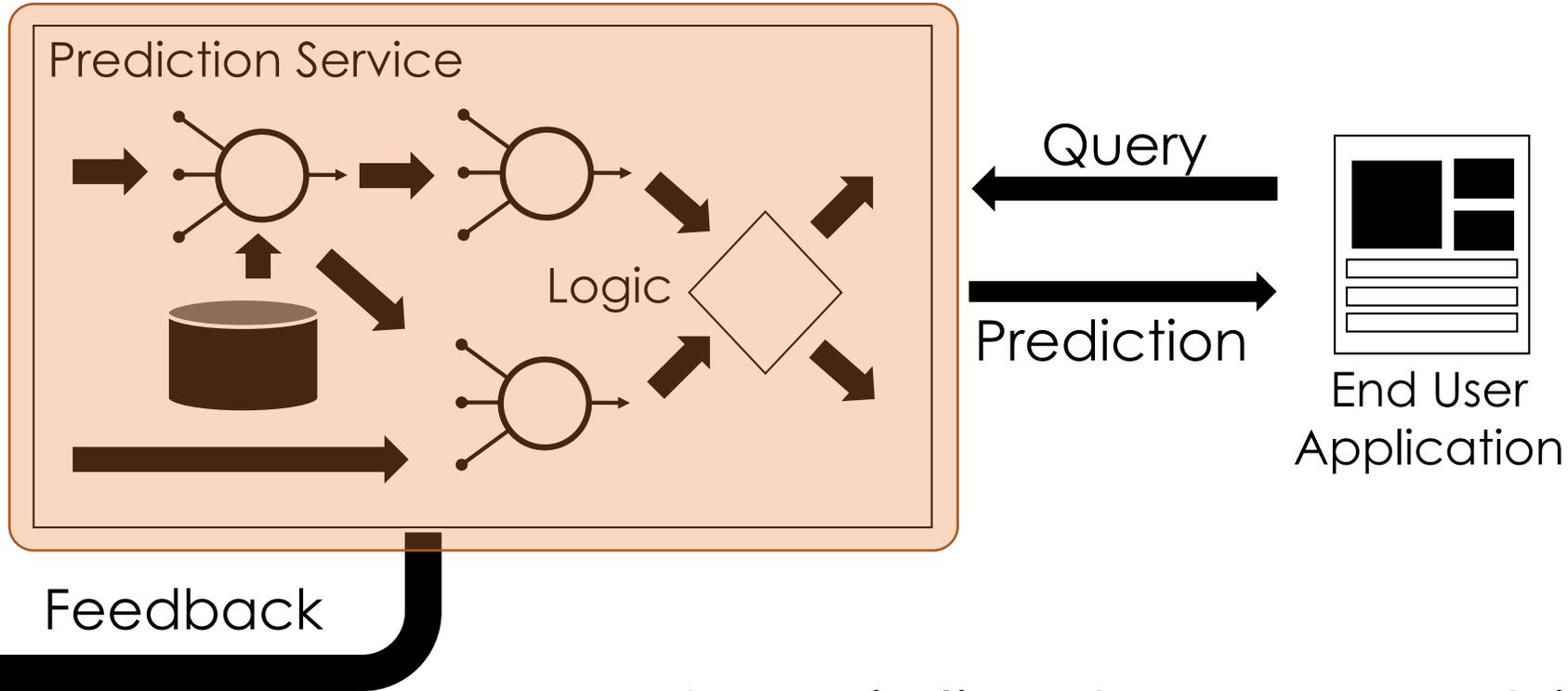


Two Approaches

➤ **Offline:** Pre-Materialize Predictions

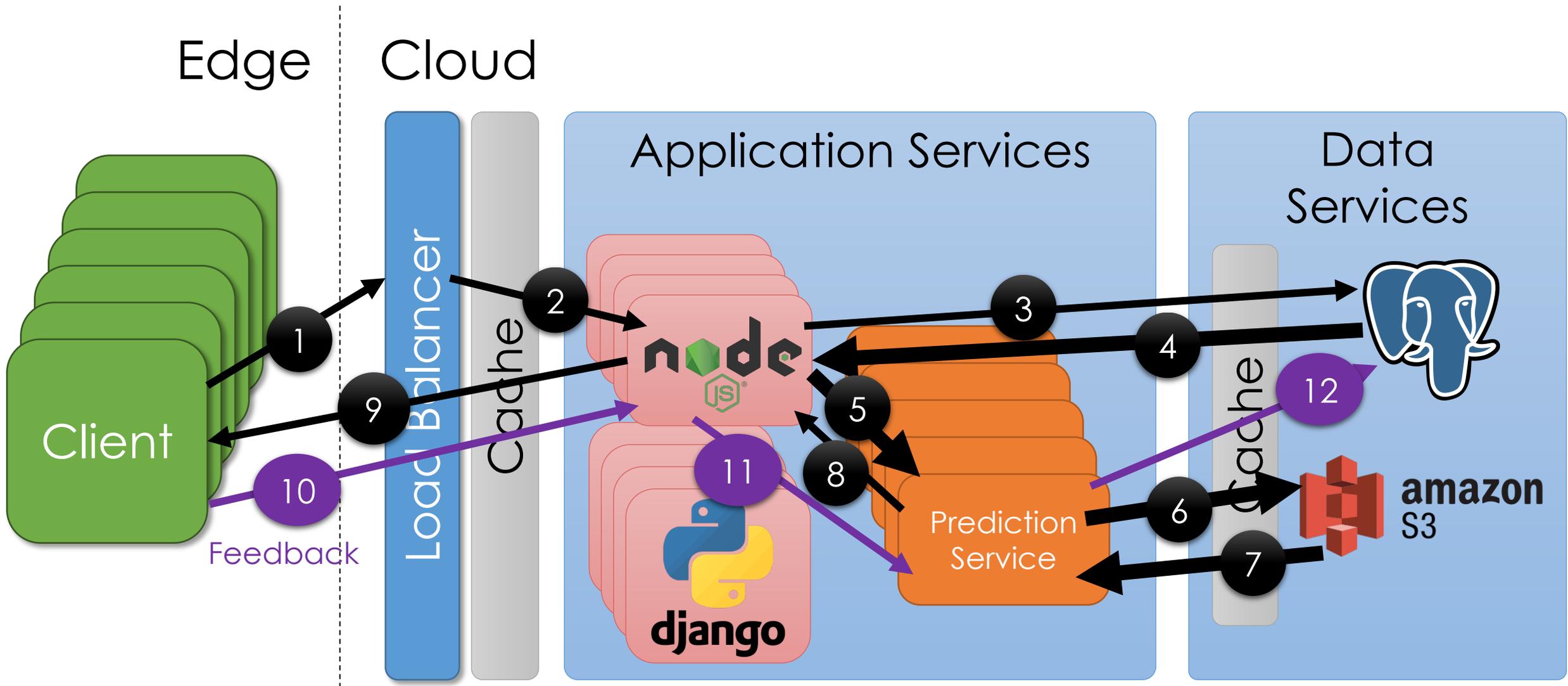
➤ **Online:** Compute Predictions on the fly

Prediction Services



Specialized systems which render predictions at **query time**.

Architecture of a Prediction Service



Architecture of a Prediction Service

Simple Prediction Service Design

**Prediction
Service**



Flask

PYT ORCH

Use existing web technologies.

➤ **Strengths**

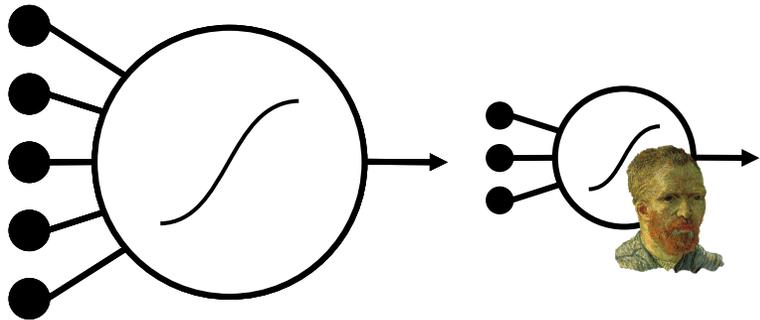
- Leverages existing technologies
- Easy to setup

➤ **Limitations**

- Need to address common issues
 - batching, monitoring, etc...
- Limited isolation between models
- Missed opportunity for common abstraction

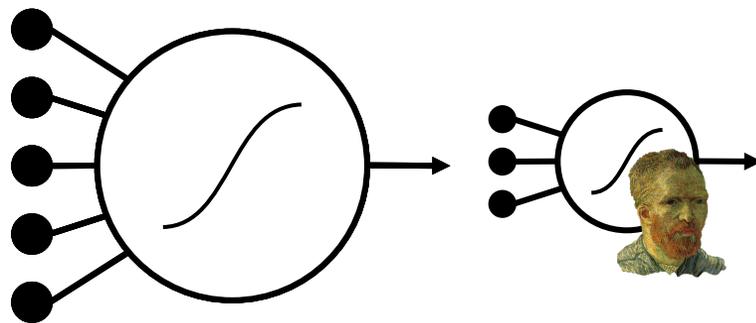
Two Approaches to Prediction Service Design

 **VELOX**



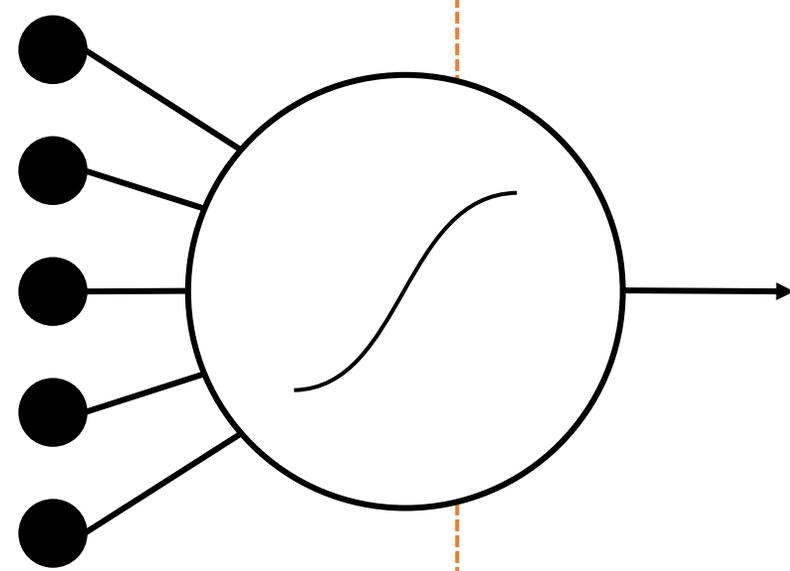
Addressing Feedback by **Learning at Different Speeds**

 **VELOX**



Learning

Inference



Application



Feedback

Hybrid Offline + Online Learning

Update feature functions **offline** using batch solvers

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

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

$$w_u$$

Update the user weights **online**:

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

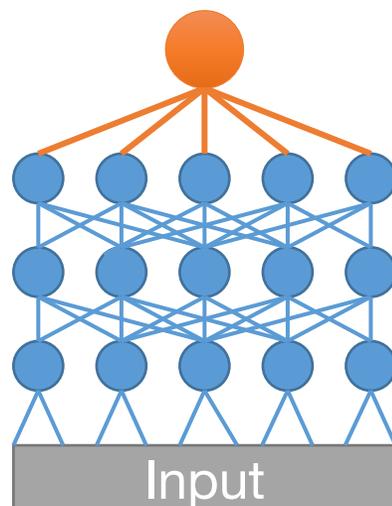
Common modeling structure

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

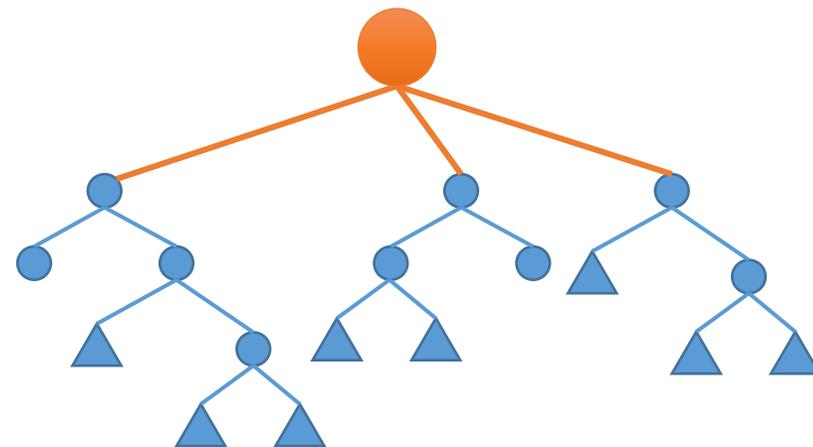
Matrix
Factorization



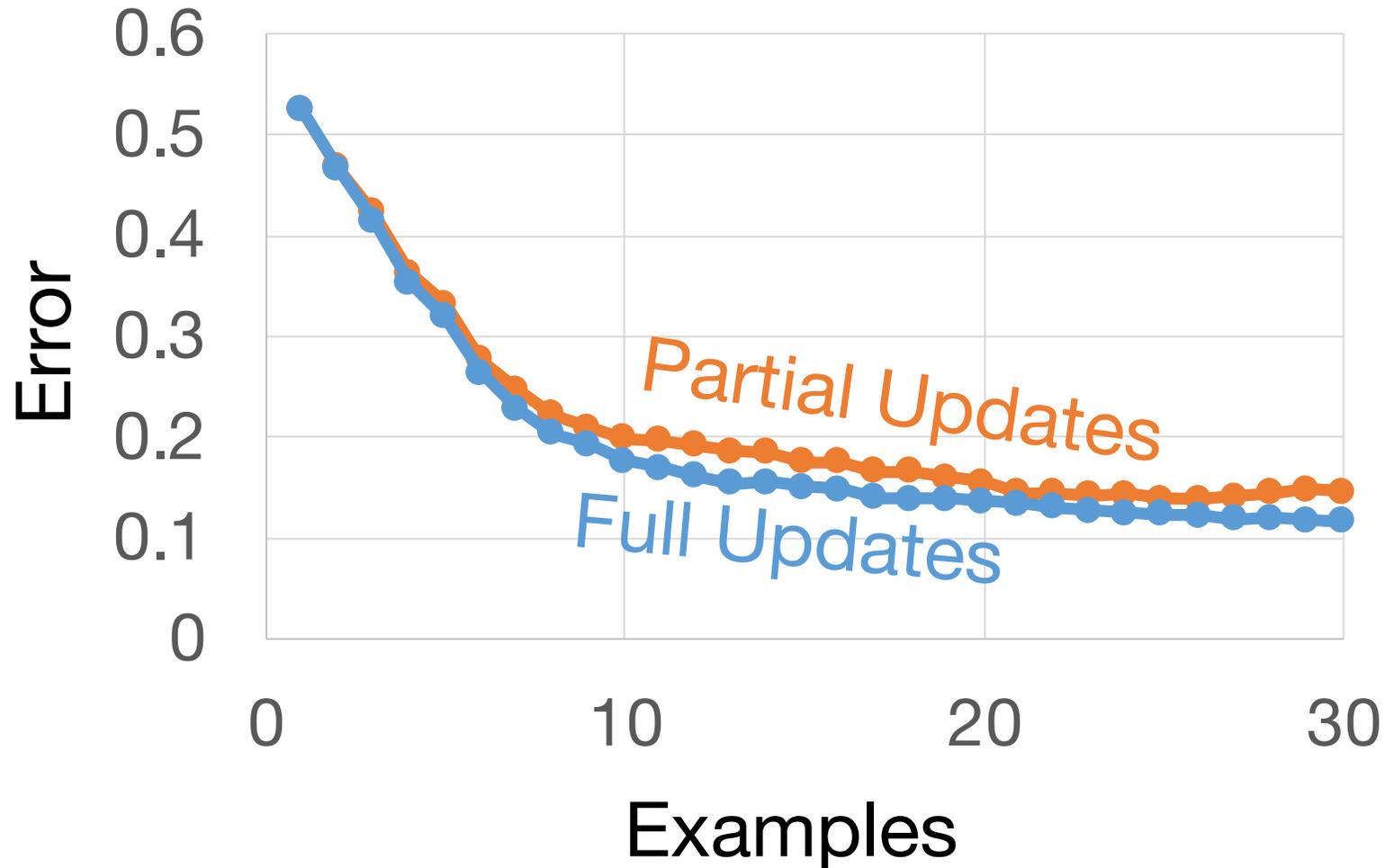
Deep
Learning



Ensemble
Methods



Velox Online Learning for Recommendations (20-News Groups)



Partial Updates: 0.4 ms
Retraining: 7.1 seconds

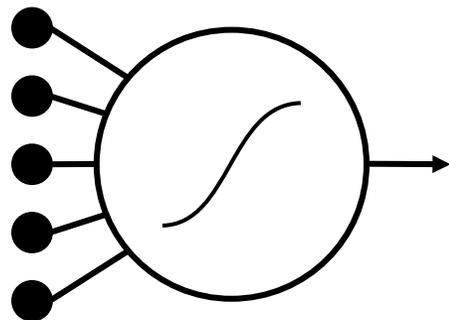
*>4 orders-of-magnitude
faster adaptation*

Learning

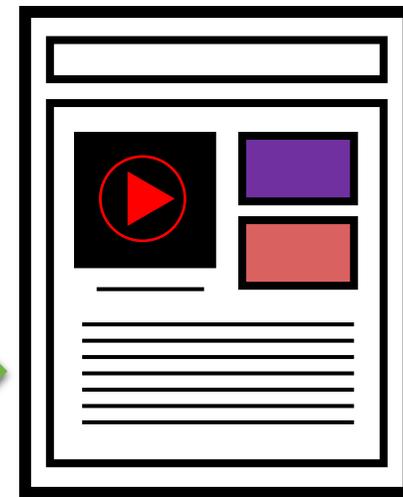
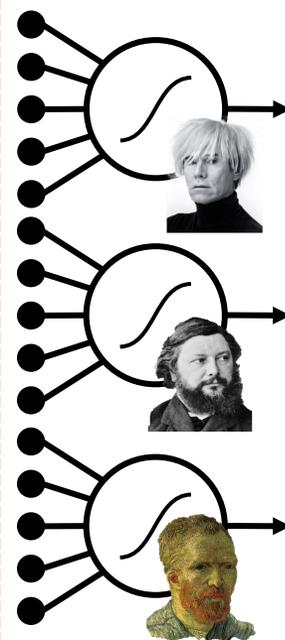
Inference



Slow Changing Model



Fast Changing Model per user

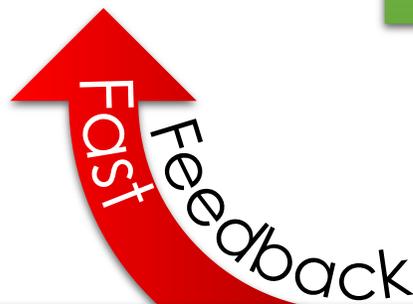


Application



Feedback

Slow

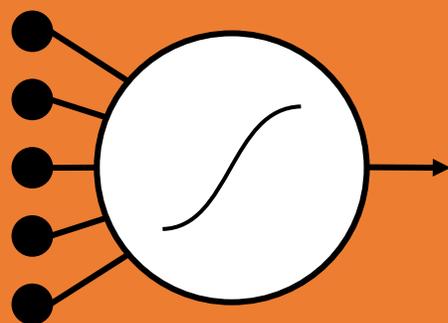


Fast Feedback

Learning



Slow Changing Model

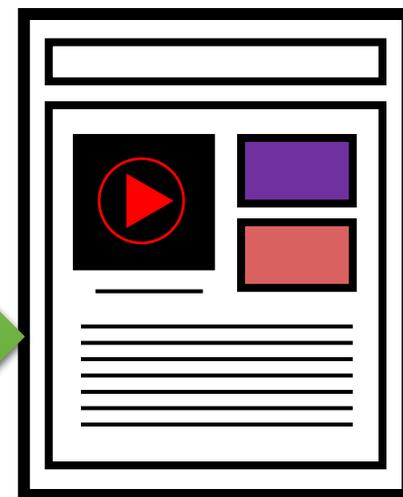
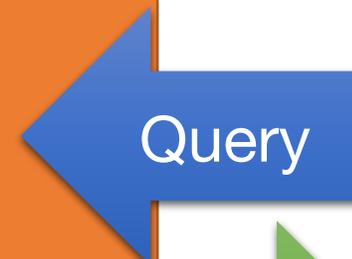


Fast Changing Model per user



Velox

Inference



Application



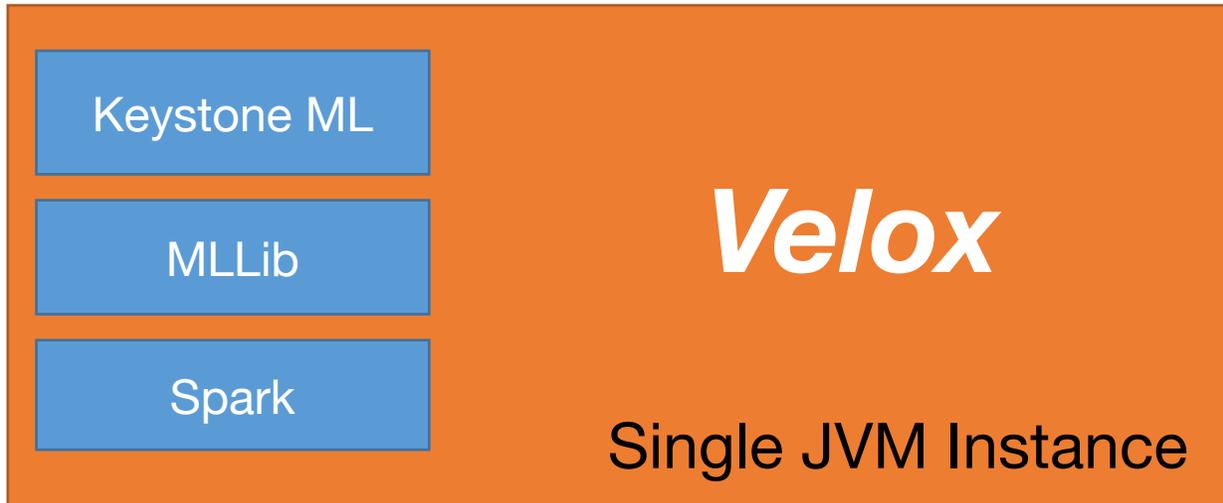
Slow

VELOX Architecture

Fraud
Detection



Content
Rec.



VELOX Architecture

Fraud
Detection



Content
Rec.



Personal
Asst.



Robotic
Control



Machine
Translation



Keystone ML

MLLib

Spark

Velox

Single JVM Instance

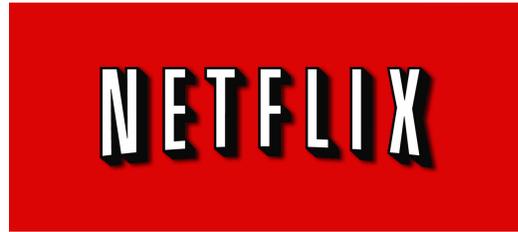


VELOX as a Middle Layer Arch?

Fraud
Detection



Content
Rec.



Personal
Asst.



Robotic
Control



Machine
Translation



Generalize *Velox*?

theano

Dato



Create

Caffe



TensorFlow



dmlc

mxnet



KALDI

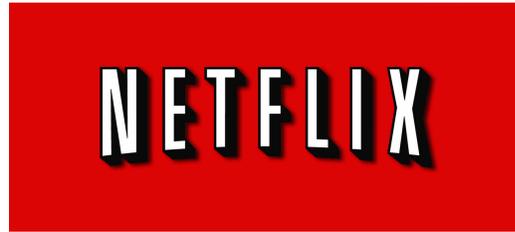
KeystoneML

Clipper **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



Middle layer for prediction serving.

**Common
Abstraction**

**System
Optimizations**

theano

APACHE
Spark

scikit
learn

Caffe

TensorFlow

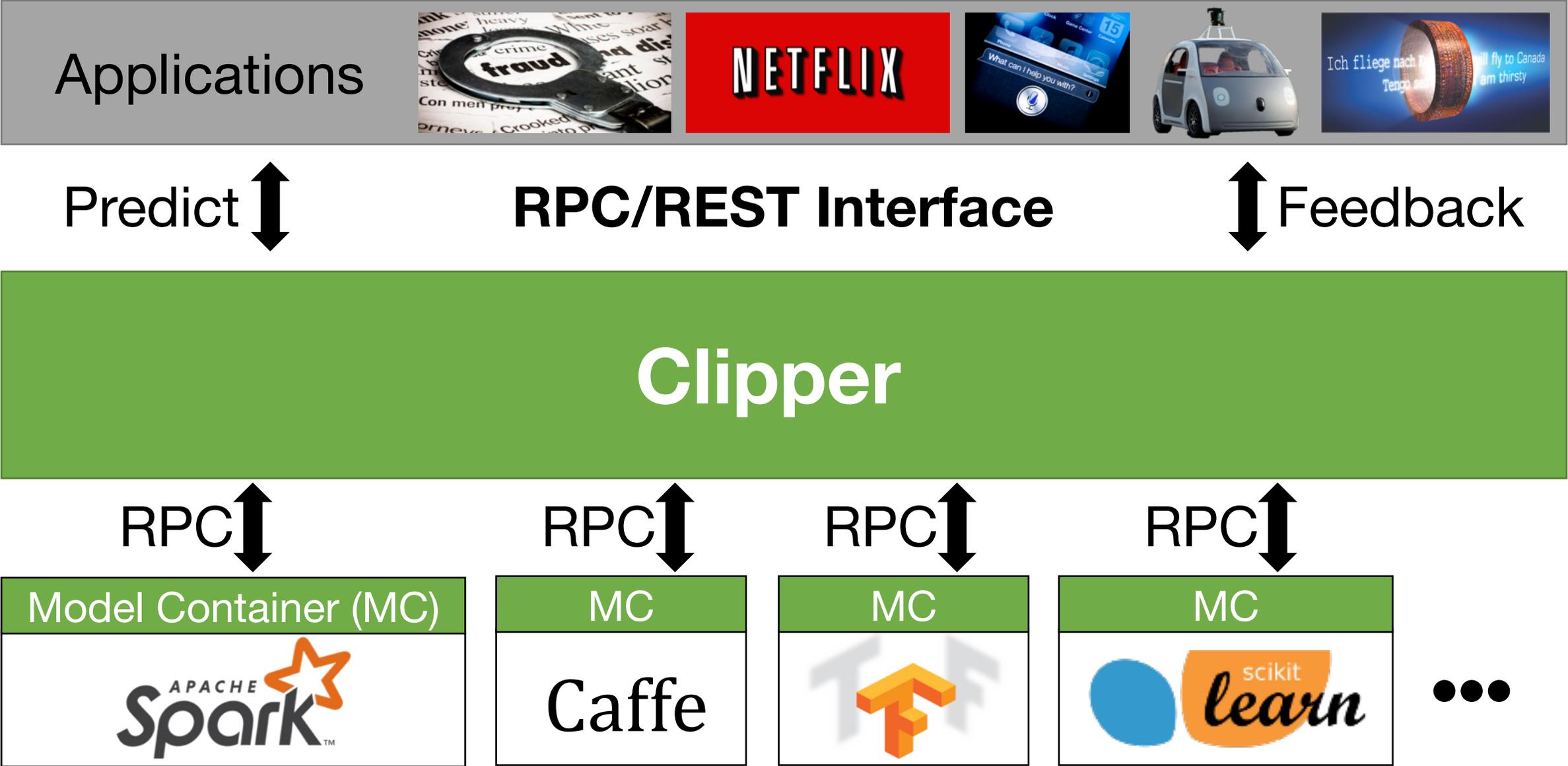
Dato 

Create
dmlc
mxnet

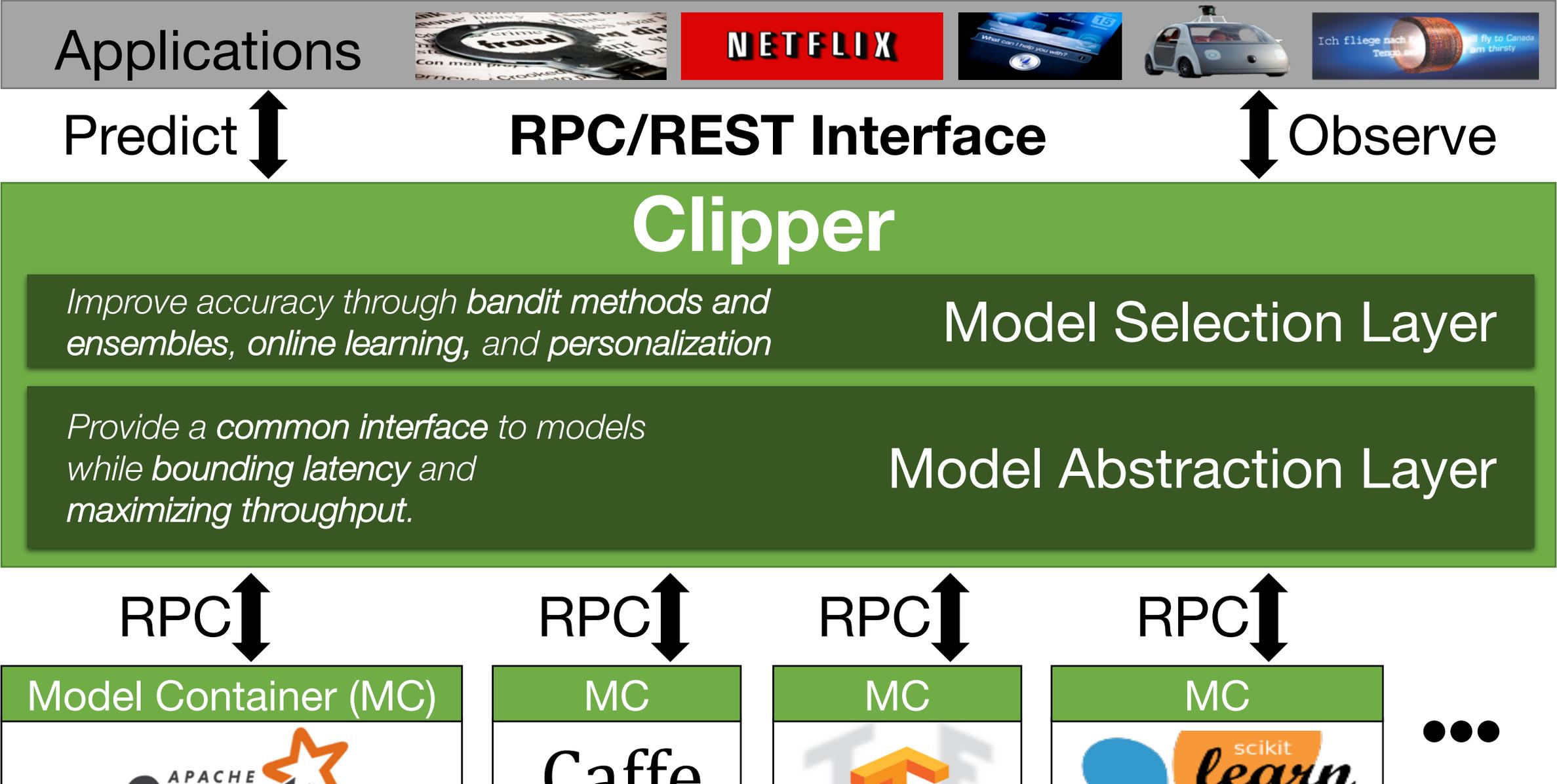


KALDI 39

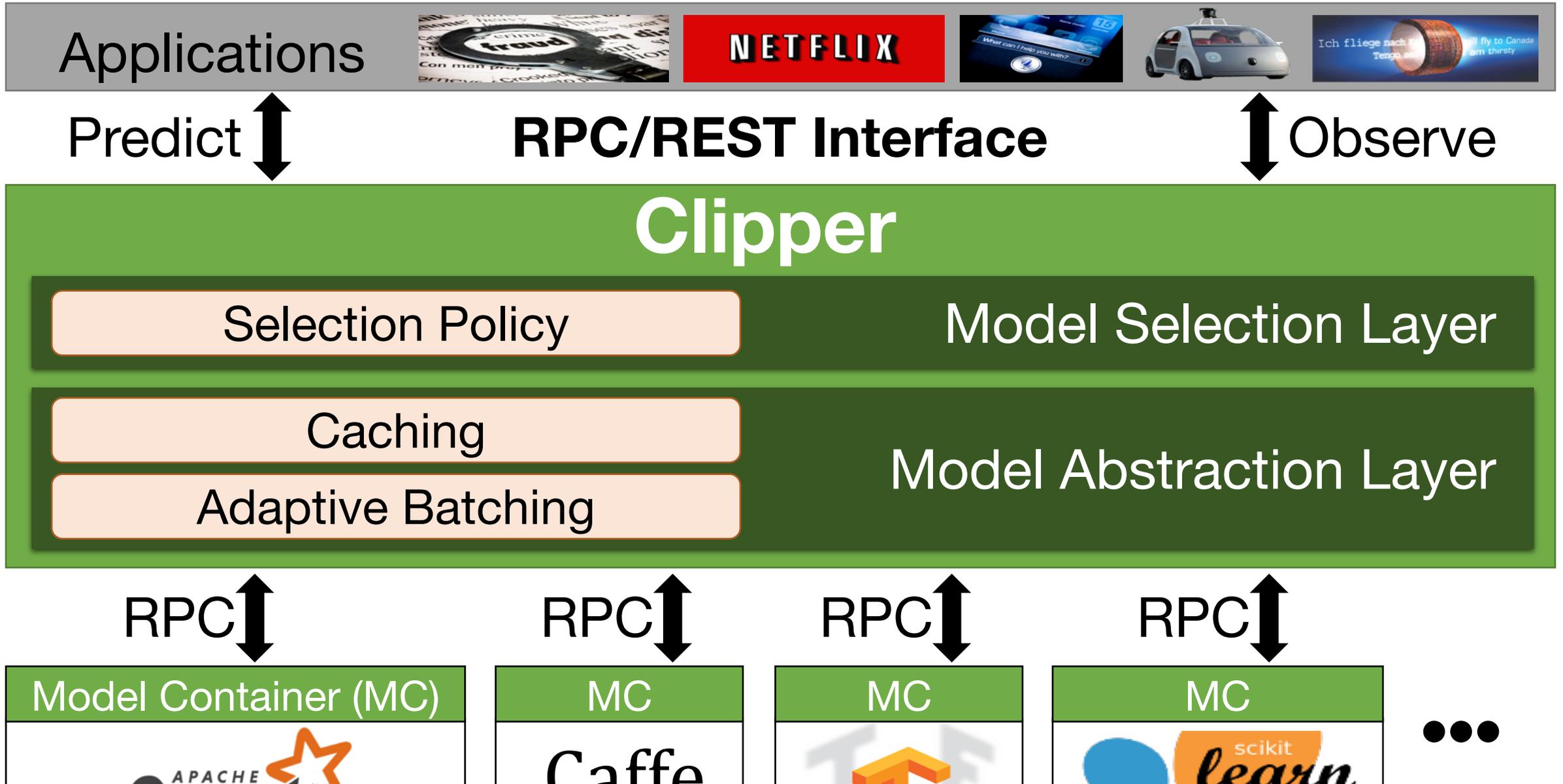
Clipper Decouples Applications and Models



Clipper Architecture



Clipper Architecture



Approximate Caching

Adaptive Batching

Model Abstraction Layer

RPC

RPC

RPC

RPC

Model Wrapper (MW)

MW

MW

MW

KeystoneML

Caffe



Caching

Adaptive Batching

Model Abstraction Layer

RPC

RPC

RPC

RPC

Model Container (MC)

MC

MC

MC



Caffe



Common Interface → Simplifies Deployment:

- Evaluate models using original code & systems

Container-based Model Deployment

Implement Model API:

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

Container-based Model Deployment

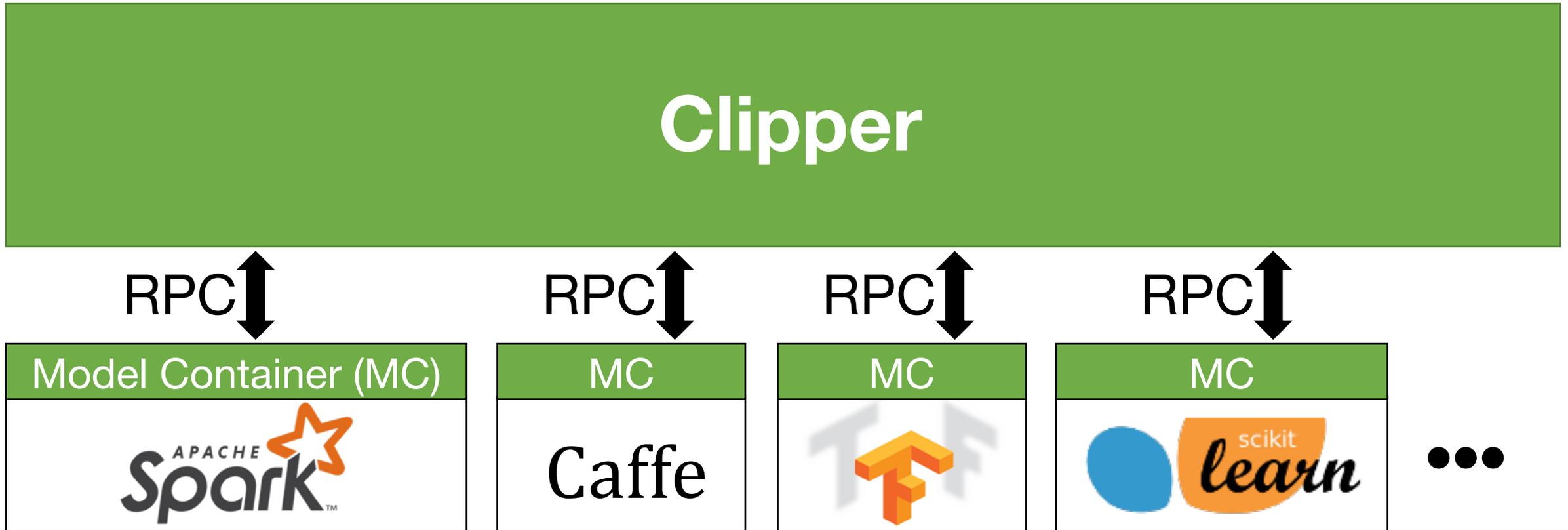
Model implementation packaged in container

Model Container (MC)

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



Container-based Model Deployment



Caching

Adaptive Batching

Model Abstraction Layer

RPC ↑↓

RPC ↑↓

RPC ↑↓

RPC ↑↓

Model Container (MC)

MC

MC

MC

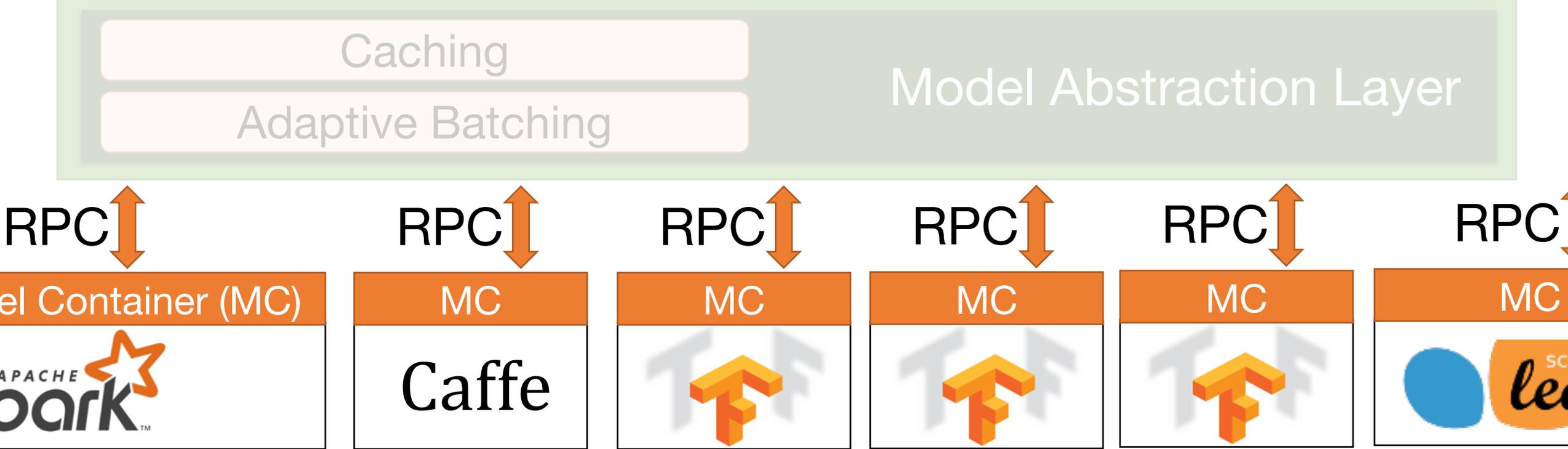


Caffe



Common Interface → Simplifies Deployment:

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



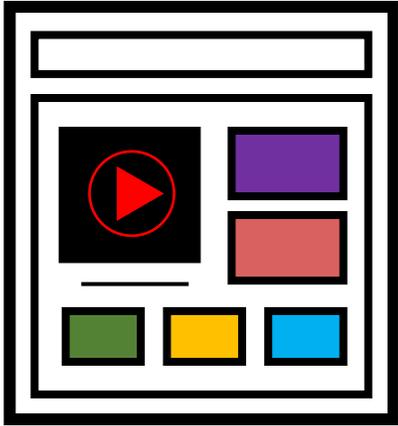
Common Interface → Simplifies Deployment:

- 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**

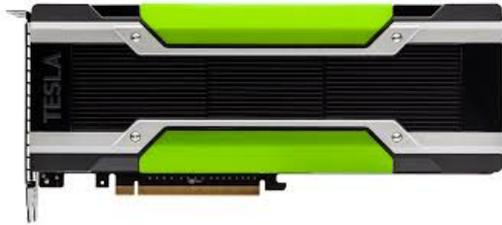
Batching to Improve Throughput

- Why batching helps:



A single page load may generate many queries

Hardware Acceleration

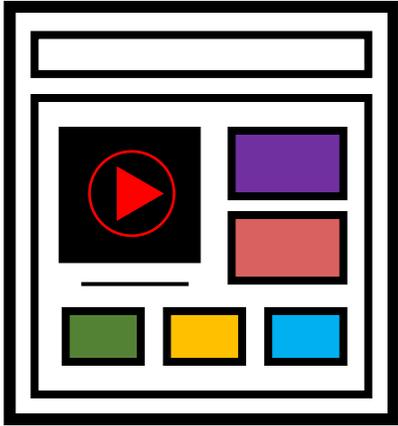


Helps amortize system overhead

- 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

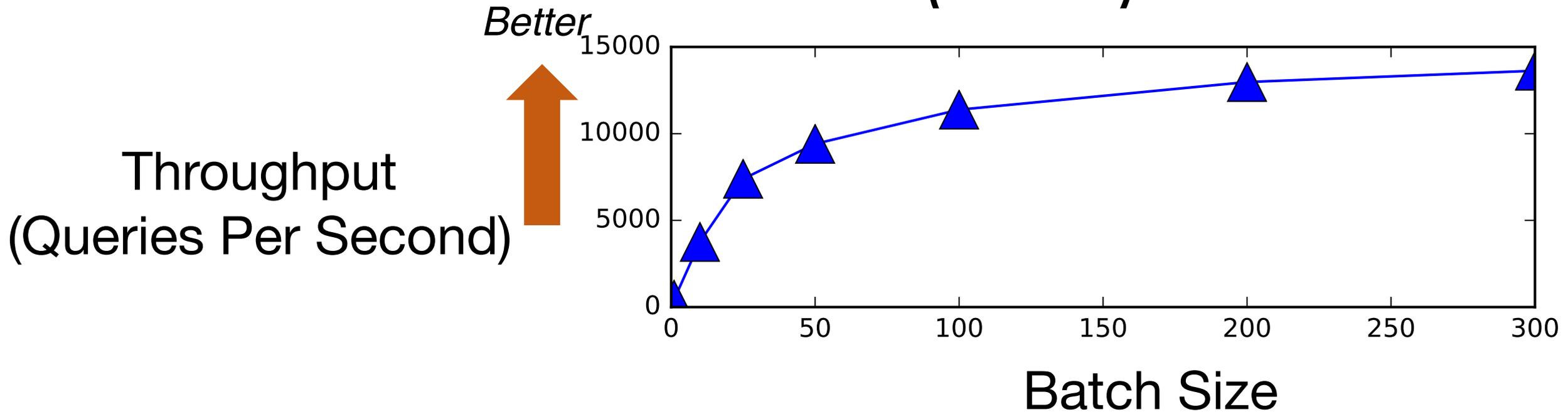
- Optimal batch depends on:
 - hardware configuration
 - model and framework
 - system load

Clipper Solution:

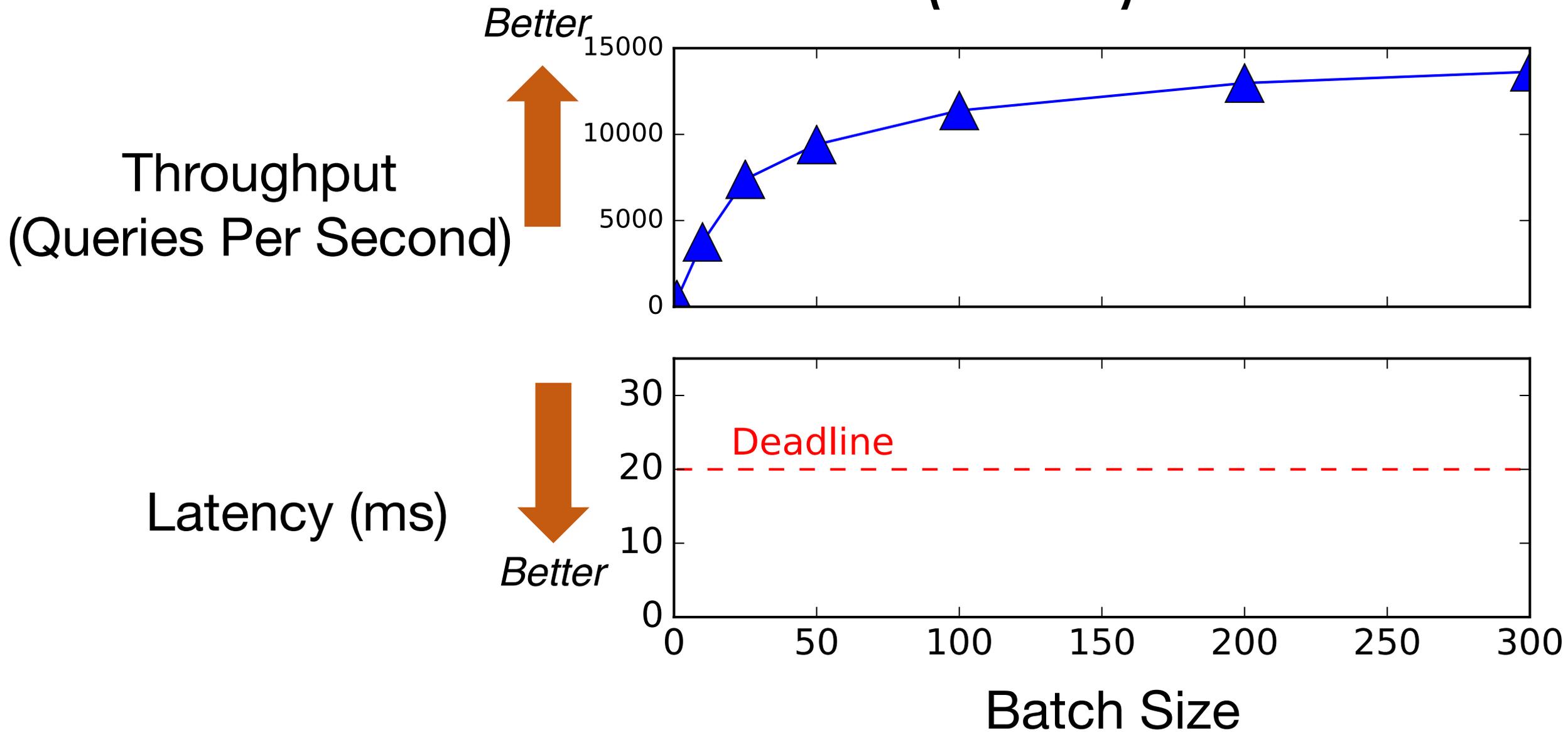
Adaptively tradeoff latency and throughput...

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

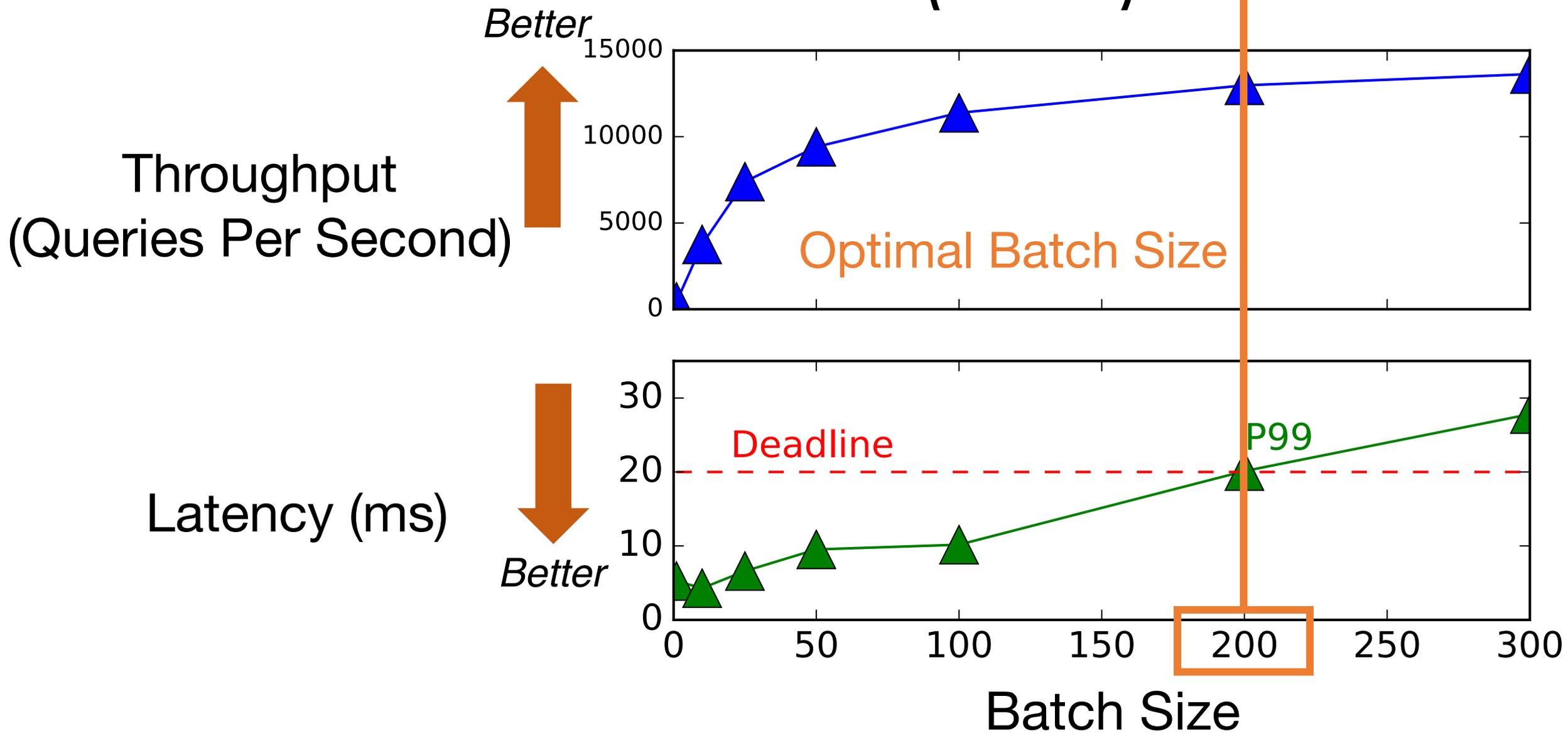
Tensor Flow Conv. Net (GPU)



Tensor Flow Conv. Net (GPU)

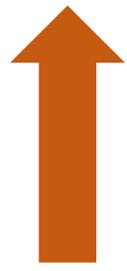


Tensor Flow Conv. Net (GPU)



Throughput
(QPS)

Better



60000
40000
20000
0

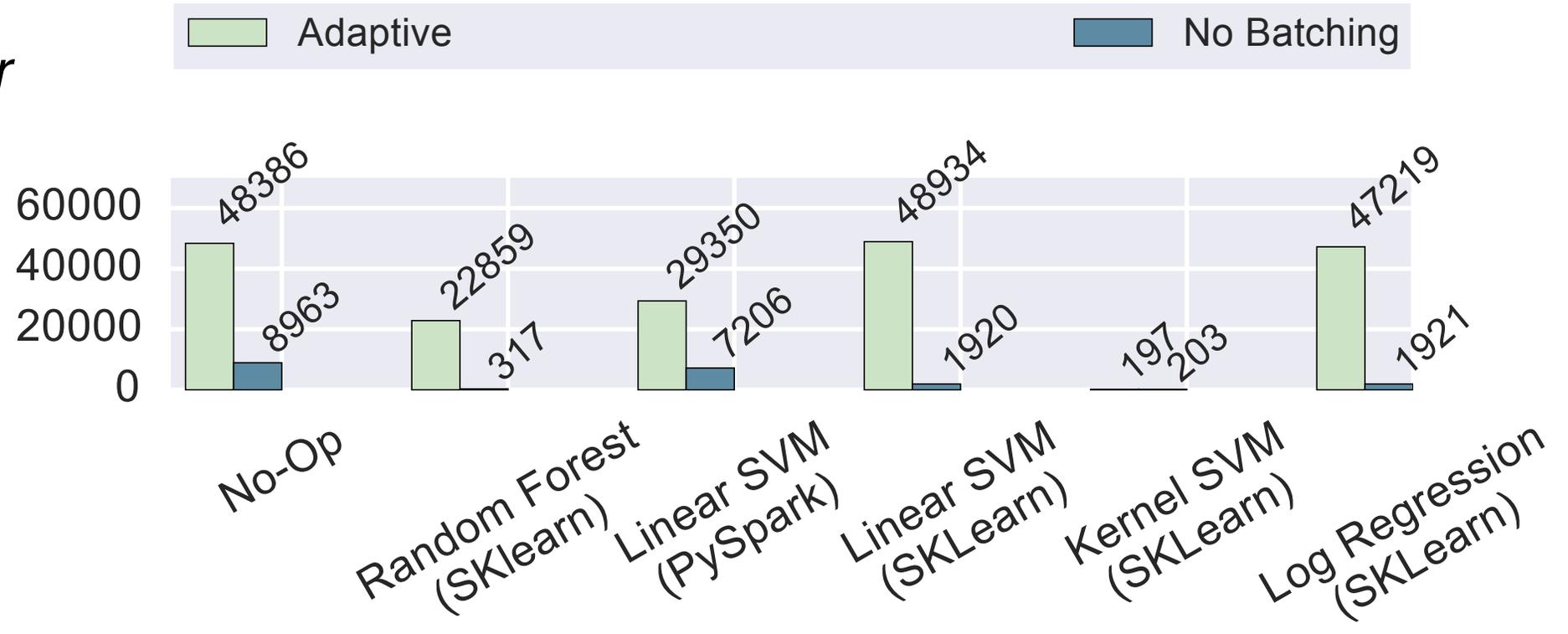
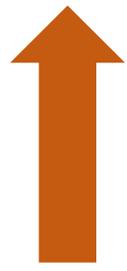
No Batching



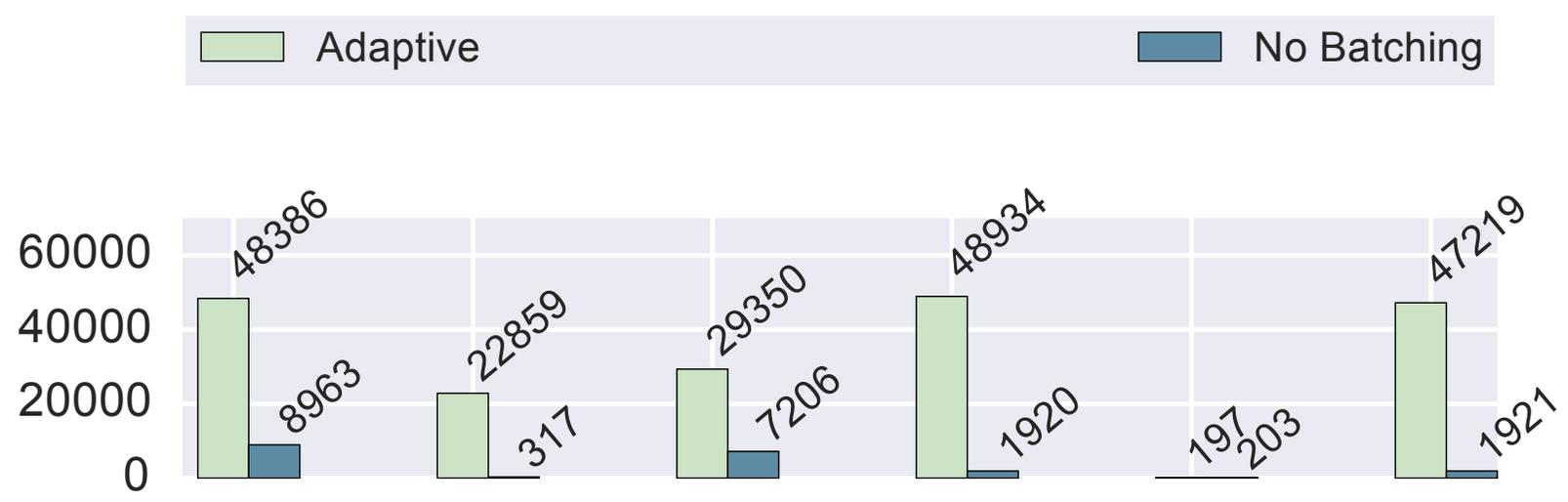
No-Op
Random Forest (SKlearn)
Linear SVM (PySpark)
Linear SVM (SKLearn)
Kernel SVM (SKLearn)
Log Regression (SKLearn)

Throughput
(QPS)

Better

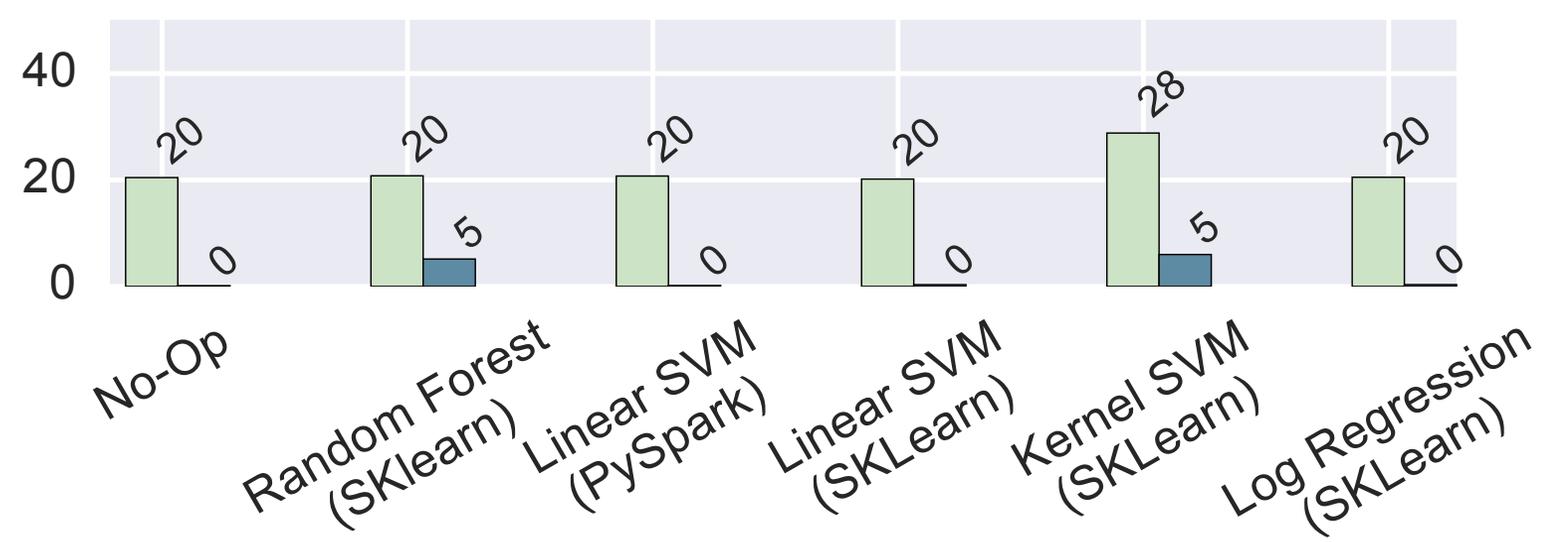
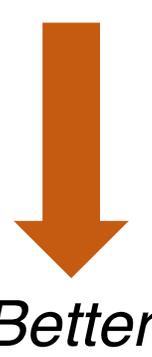


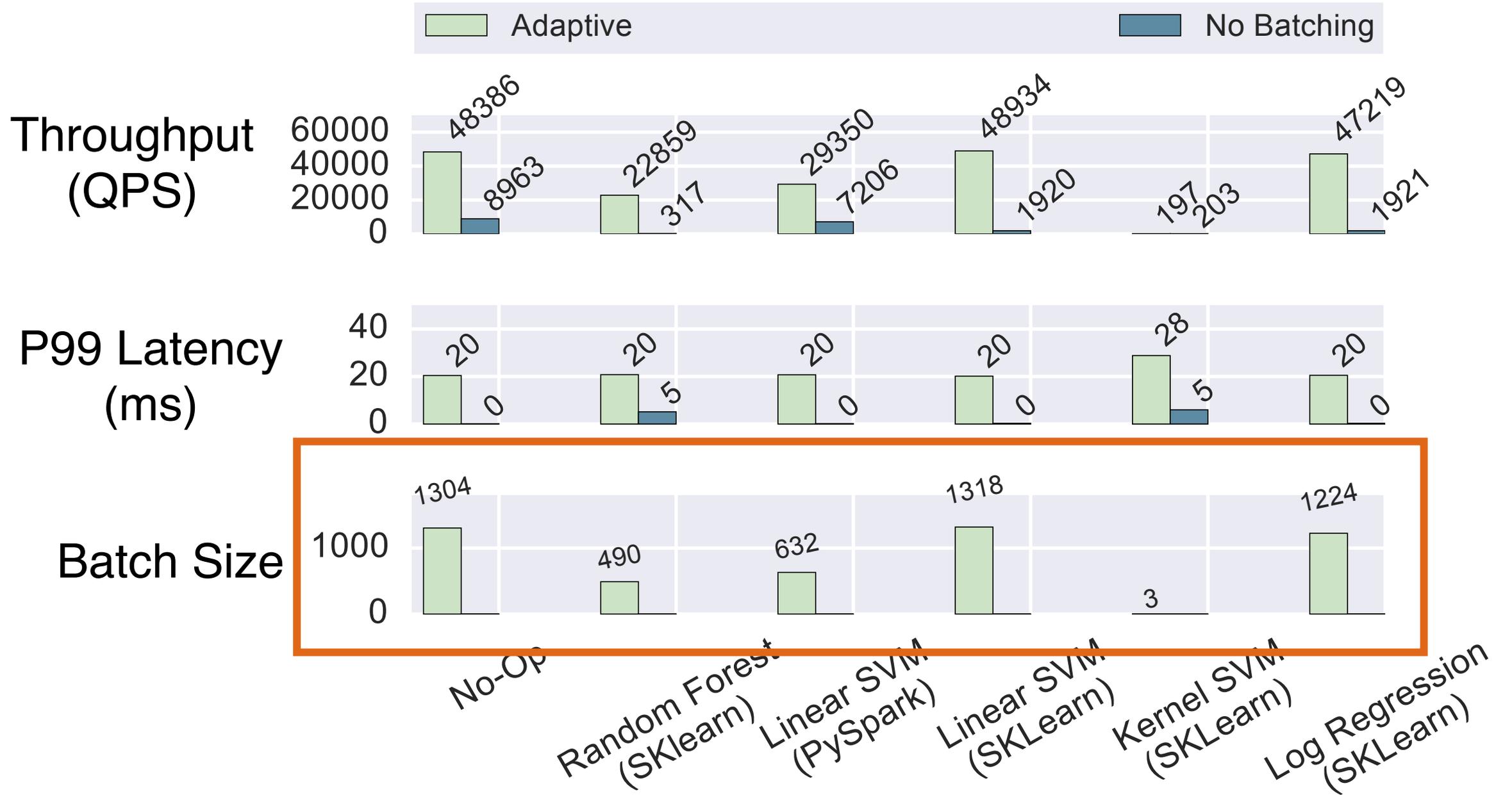
Throughput
(QPS)



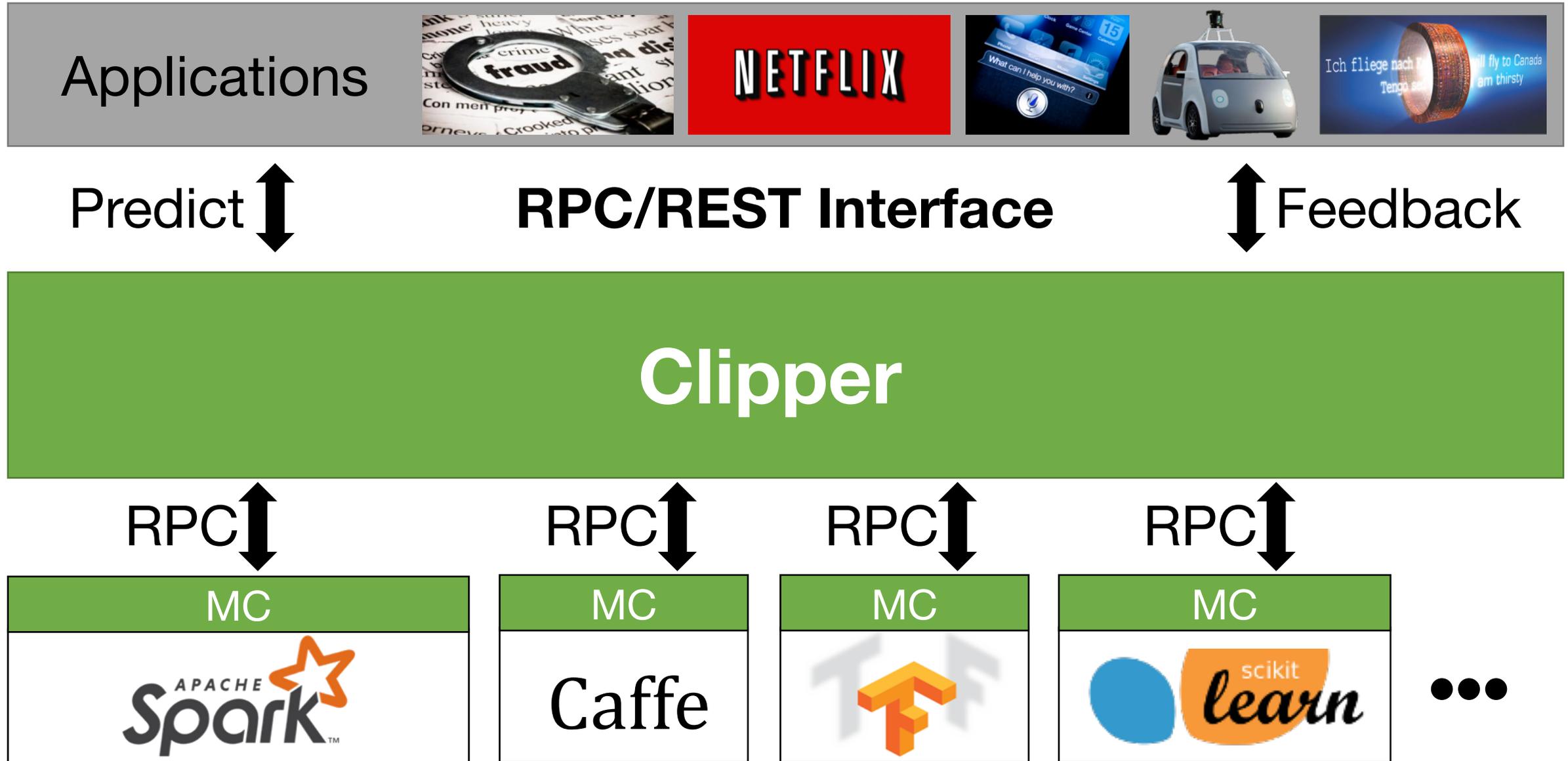
P99 Latency
(ms)

20 ms is
Fast Enough

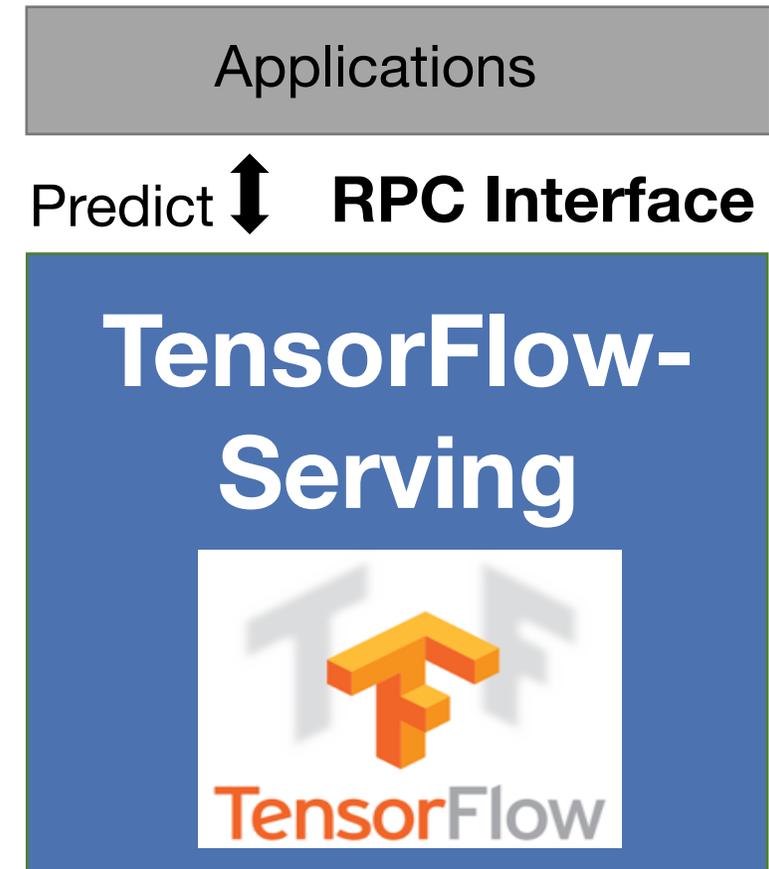
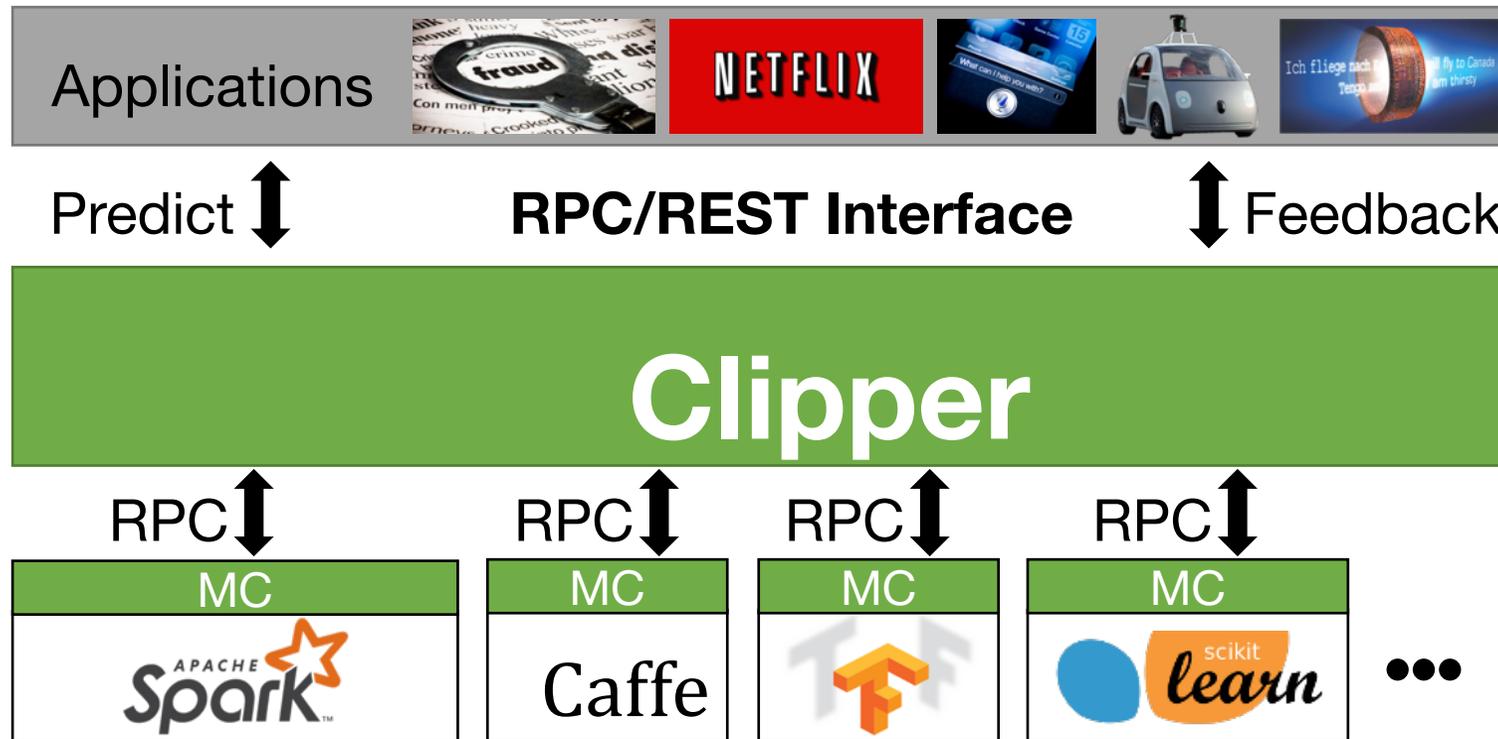




Overhead of decoupled architecture

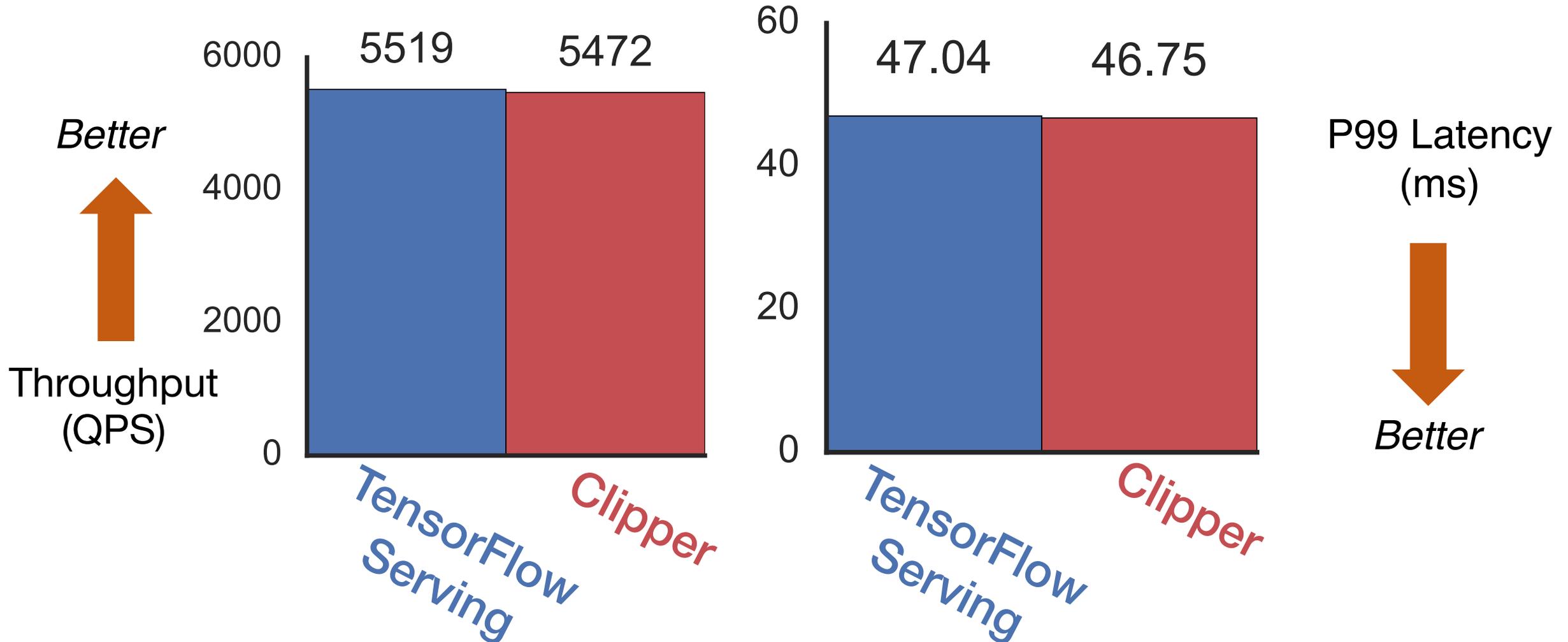


Overhead of decoupled architecture

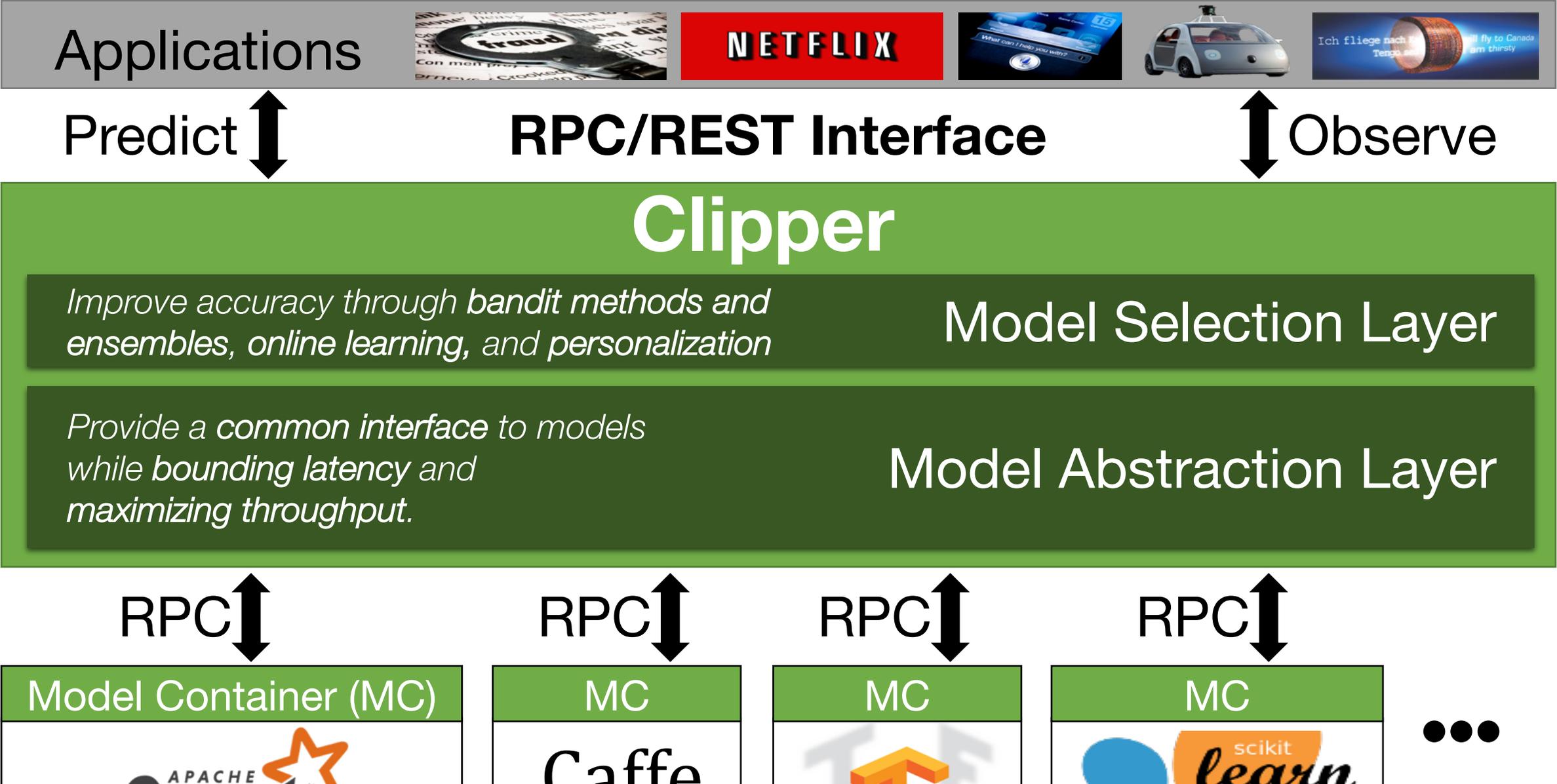


Overhead of decoupled architecture

Model: AlexNet trained on CIFAR-10

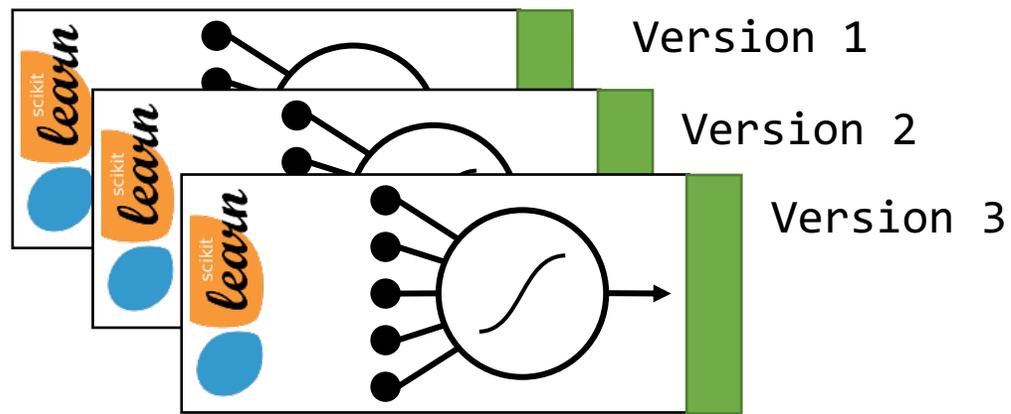


Clipper Architecture

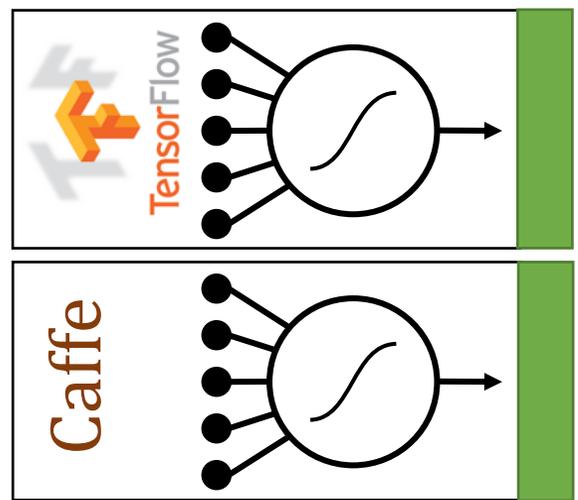


Improve accuracy through bandit methods and ensembles, online learning, and personalization

Model Selection Layer

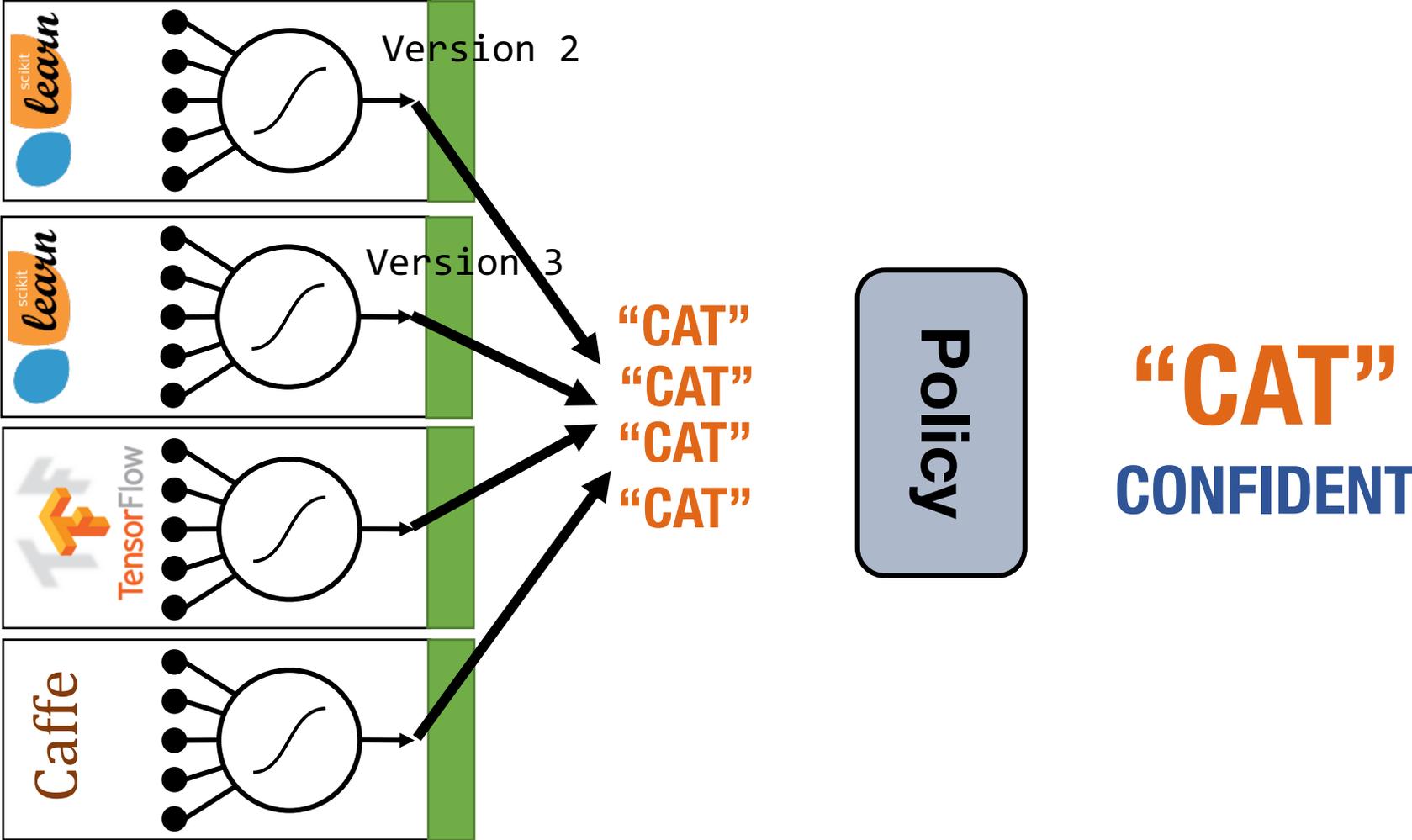


Periodic retraining

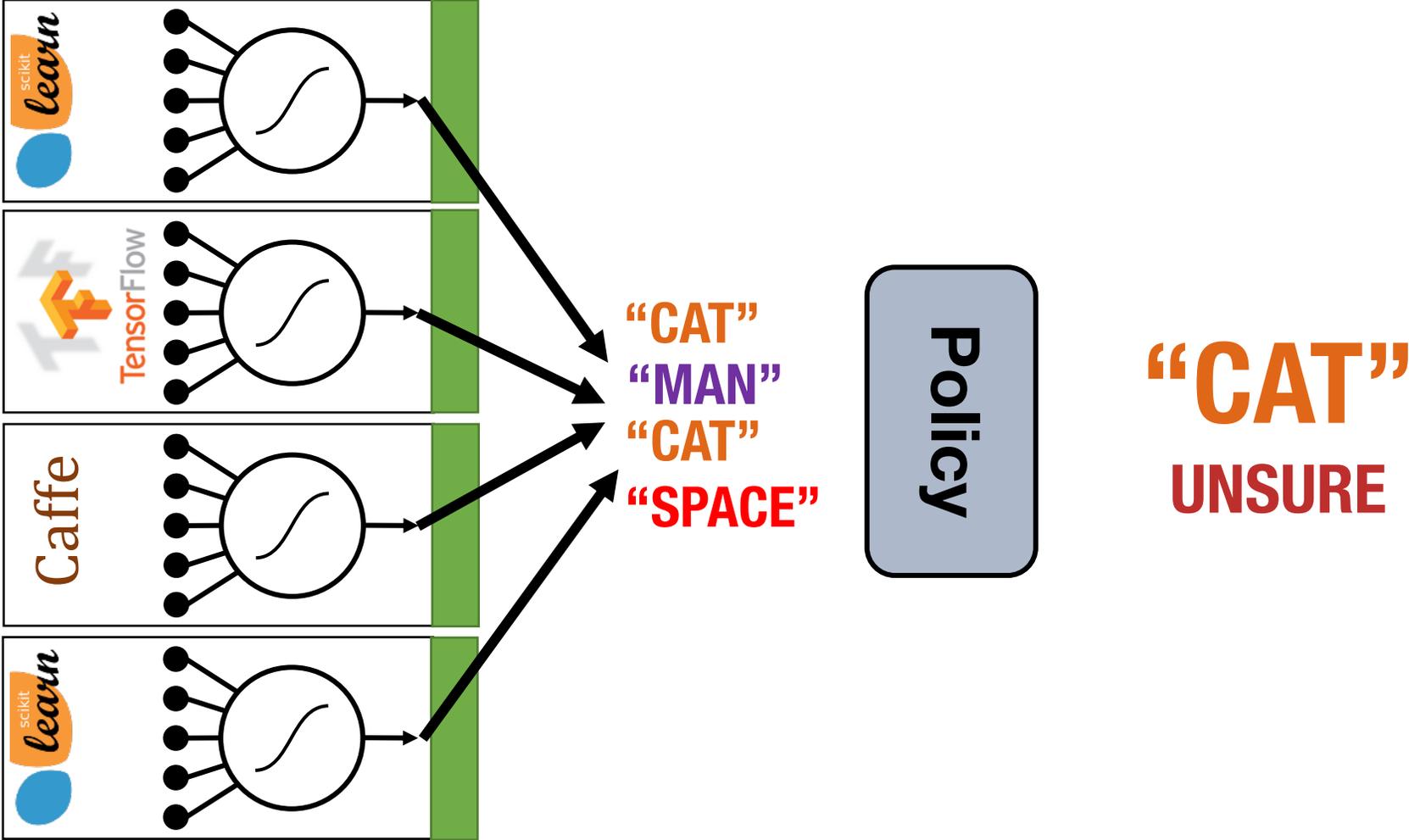


Experiment with new models and frameworks

Selection Policy: Estimate confidence



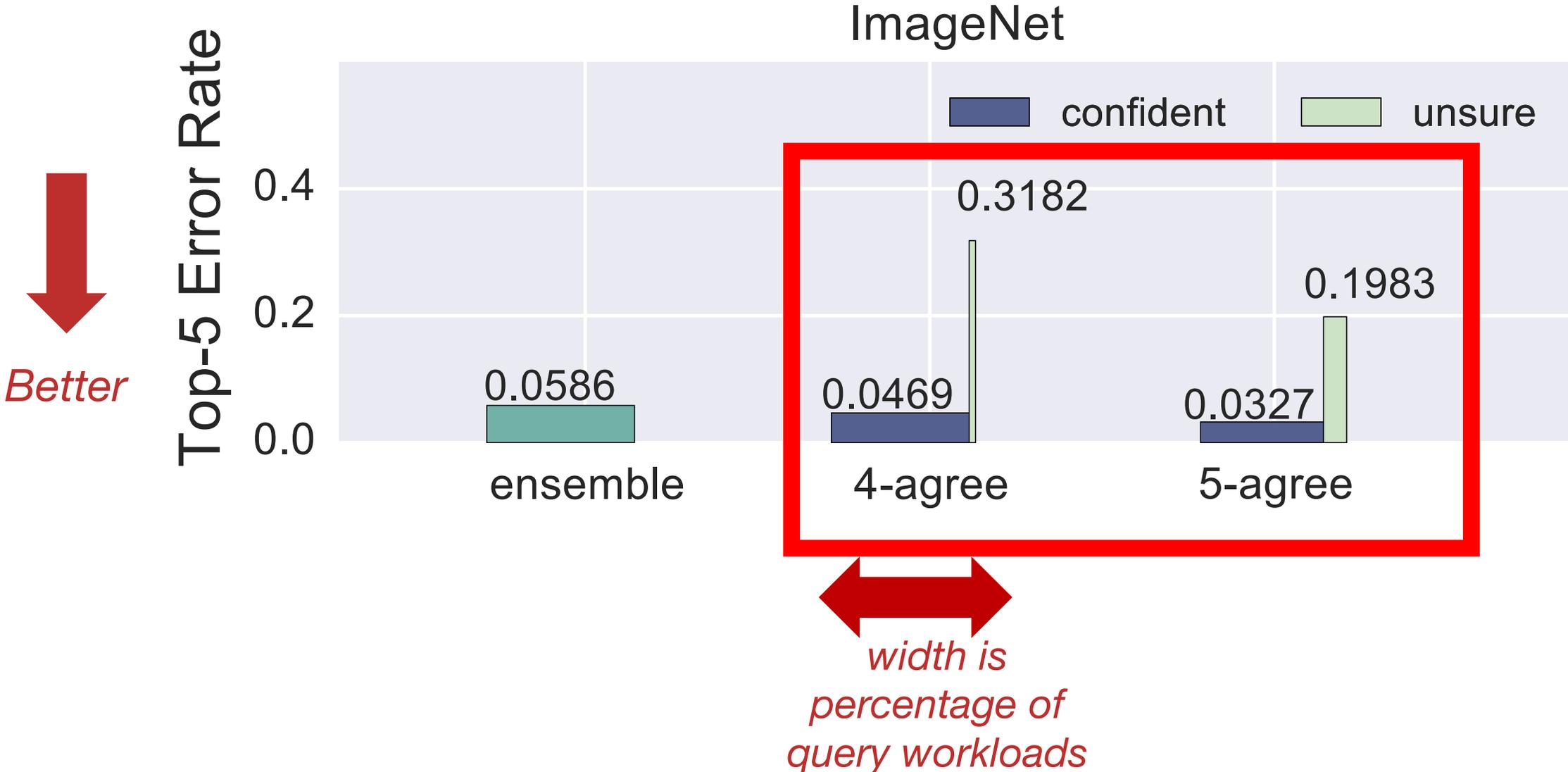
Selection Policy: Estimate confidence



Selection Policy: Estimate confidence



Selection Policy: Estimate confidence



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

Online: Compute Predictions at Query Time

➤ **Examples**

- Speech recognition, image tagging
- Ad-targeting based on search terms, available ads, user features

➤ **Advantages**

- Compute only necessary queries
- Enables models to be changed rapidly and bandit exploration
- Queries do not need to be from small ground set

➤ **Disadvantages**

- Increases complexity and computation overhead of serving system
- Requires low and predictable latency from models

Prediction Pipelines

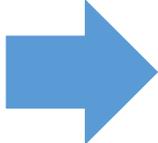
Example

Query Image

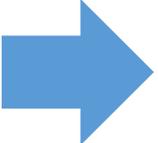
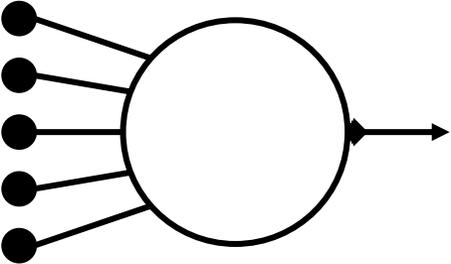


This is my daughter!

Query Image



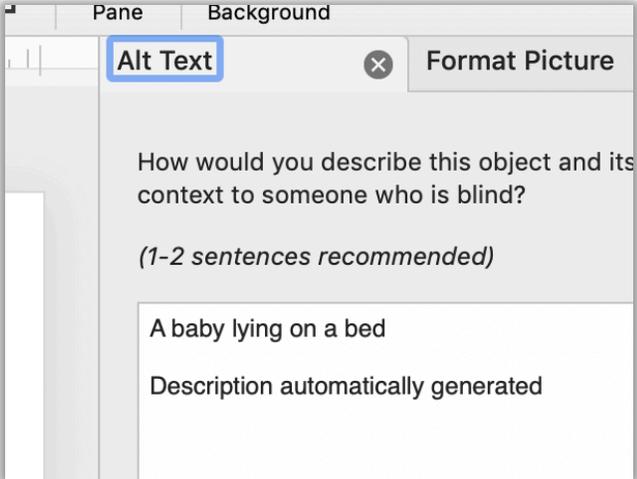
Machine Learning Model



Prediction

“A baby lying on a bed”

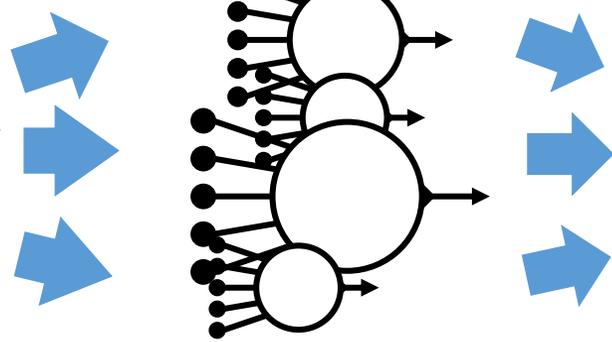
This caption was generated automatically ***in the cloud*** by Microsoft PowerPoint



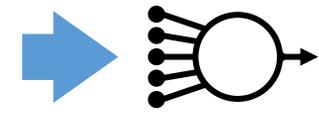
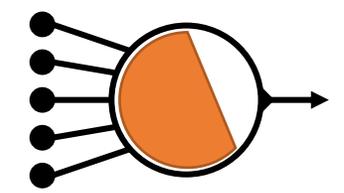
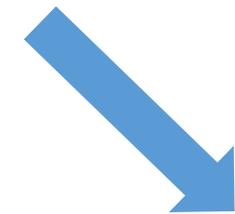
More realistic → Prediction Pipeline

Machine Learning
Model
Ensemble

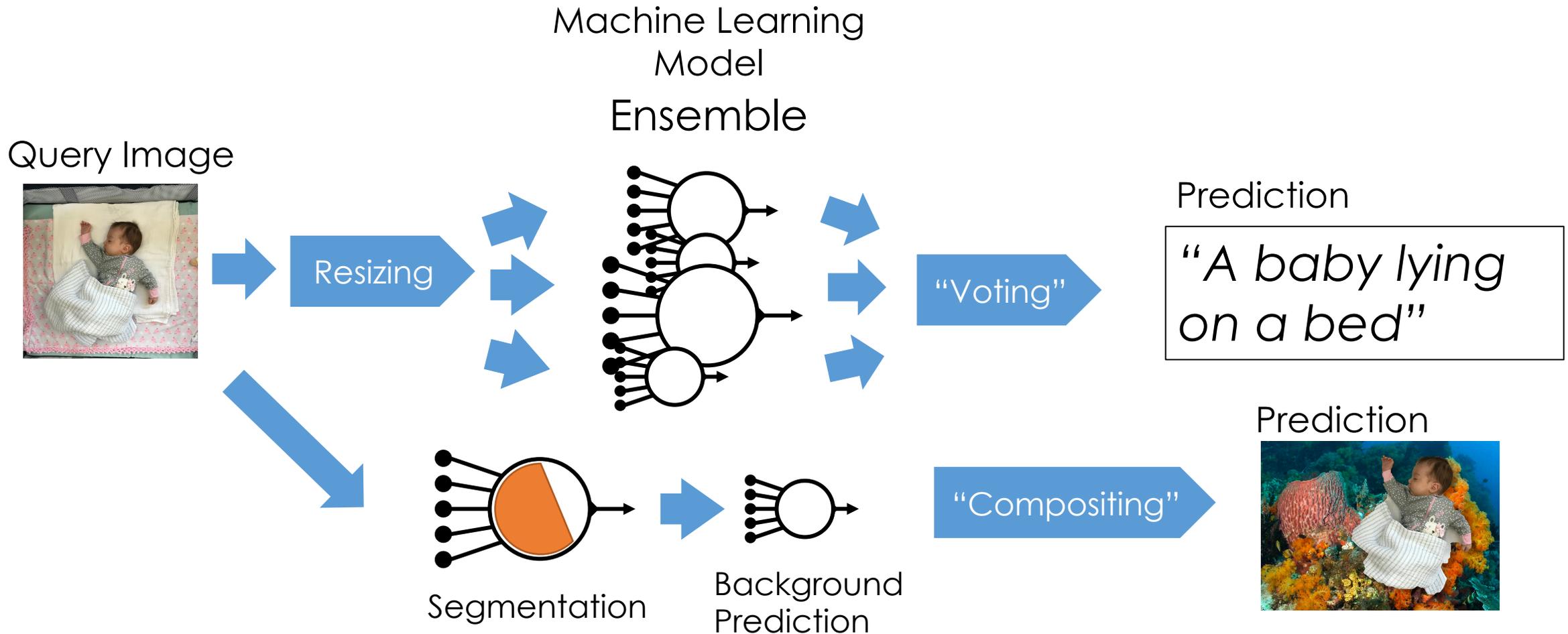
Query Image



Prediction
“A baby lying on a bed”



Prediction

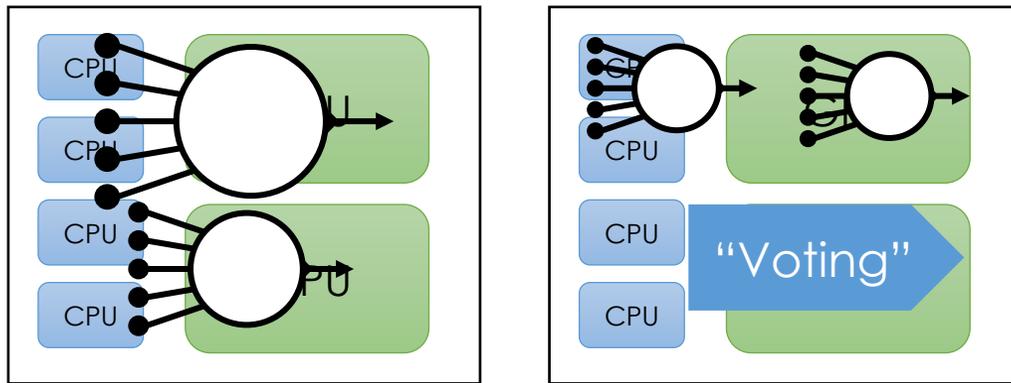



How do we provision resources for these pipelines?

Latency vs. **Throughput** vs. **Cost**

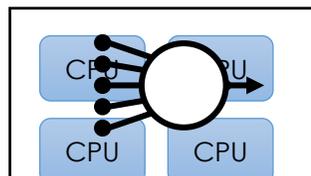
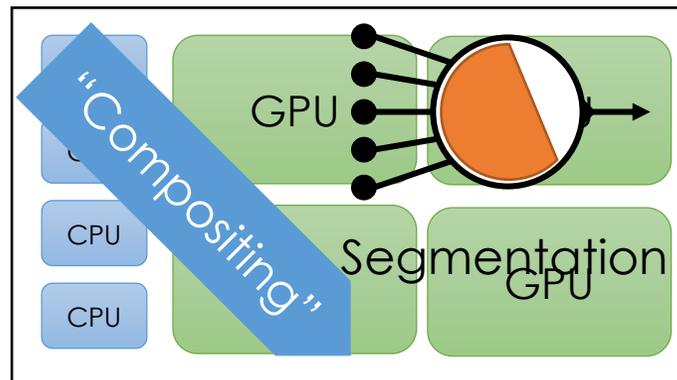
How do we provision resources for these pipelines?

Latency vs. Throughput vs. Cost



Two readings this week will address this problem.

- Pretzel
- InferLine



Background
Prediction

Cloud -- Edge

Example

KUNA

Home video security systems



Technology

- AC Powered Lamp
- Commodity ARM proc.
- 720HD Video
- Microphone & Speaker
- Infrared Motion Sensors

Goals:

- Detect, identify, and record people
- Notify homeowner and open channel of comm.



How does **KUNA** work?



Fast onboard pixel-level filter identifies suspicious change



Key frames are sent to EC2 for further processing



More sophisticated processing to reduce false positives (**costly GPU time**)

KUNA technology challenges



- Splitting classification across **device** and **cloud**.
- **Shared learning** to identify common patterns
 - e.g., traffic in urban environments
- More **efficient prediction rendering** on cloud + edge
 - Running full CV pipeline on all images is very costly



Desired Capabilities

- Event characterization: *“Package delivery at 1:33 PM”*
- Automatic user interaction: *“I would be happy to digitally sign for the package ...”*

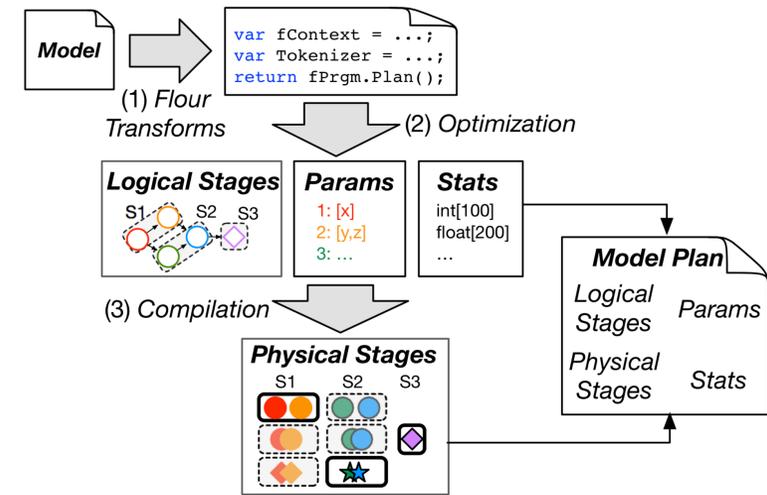


Reading This Week

Reading for the Week

- [Pretzel: Opening the Black Box of Machine Learning Prediction Serving Systems](#) (OSDI'18)
 - Optimizing prediction serving pipeline using compiler and database system techniques
- [InferLine: ML Inference Pipeline Composition Framework](#) (pre-print)
 - Optimizing prediction serving pipeline configurations for deep learning on heterogenous hardware
- [Focus: Querying Large Video Datasets with Low Latency and Low Cost](#) (OSDI'18)
 - Enabling real-time queries on video data with offline pre-processing

Pretzel: Opening the Black Box of Machine Learning Prediction Serving Systems



- Addresses a range of **practical issues**:
 - Dealing with infrequently used models
 - Need to “**page-out**” infrequently used models → **Cold starts**
 - Need to **pack** many models in same machine
 - **Sharing model stages** across prediction pipelines
 - Eliminate **redundant computation** and **memory requirements**
 - Pushing **computation (reuse)** through feature concatenation
 - Generating efficient **binary executables** from high-level DSLs
- **Setting**: Focused on **non-deep learning** pipelines
 - CPU Intensive
 - ML.Net Infrastructure and production workloads

- **Sharing model stages** across prediction pipelines
 - Eliminate **redundant computation** and **memory requirements**
- Pushing **computation (reuse)** through feature concatenation
- Generating efficient **binary executables** from high-level DSLs

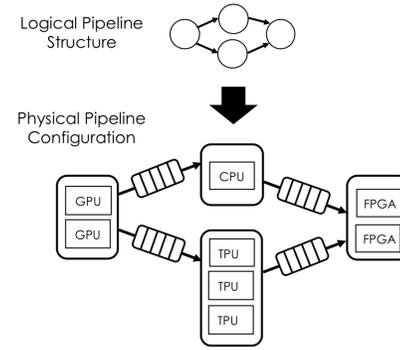
- **Setting:** Focused on **non-deep learning** pipelines
 - CPU Intensive
 - ML.Net Infrastructure and production workloads

- **Big Idea:** Leverage visibility into pipeline achieved by high-level pipeline DSL to optimize execution across pipelines and stages within a pipeline

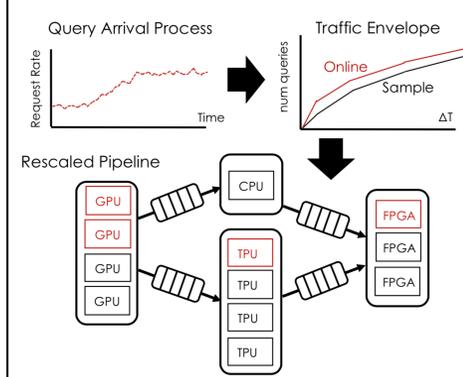
- **What to look for in reading**
 - Motivations driven by real-world **workload profiling**
 - Combination of **offline** and **online** optimization
 - Decomposition of problem into **logical** and physical plans and **runtime scheduling**

InferLine: Prediction Pipeline Provisioning and Management for Tight Latency Objectives

Offline Proactive Planner



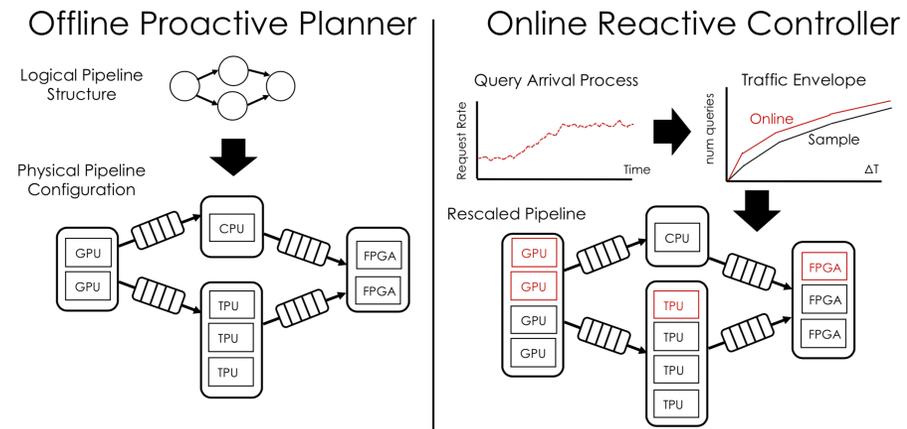
Online Reactive Controller



- **Context:** This is a pre-print paper from **my group**
 - Submitted OSDI'18 and SOSP'19 → rejected ☹️
 - Issues with presentation and contributions
 - How can we review a professor's paper?
 - Your honest feedback is very helpful!
- Why choose this paper?
 - Discusses a range of challenges in **black-box pipeline** management
 - Presents interesting **configuration space**
 - **We need feedback!** (*It is ok to be negative.*)

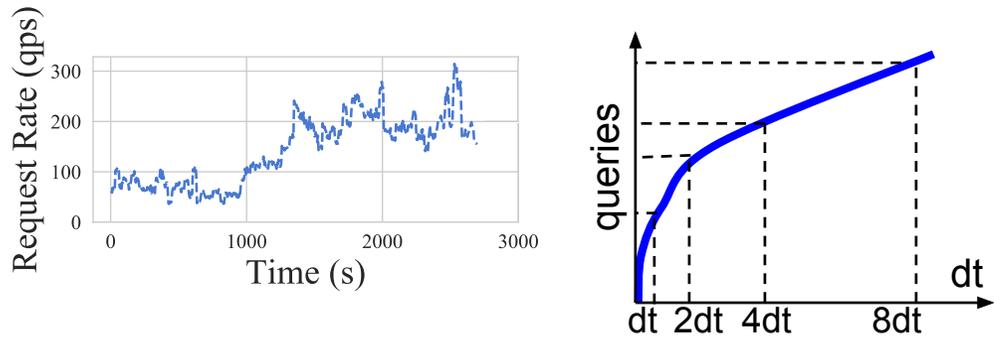
InferLine: Prediction Pipeline Provisioning and Management for Tight Latency Objectives

- **Big Idea:** Optimally configure per-model parameters in a prediction pipeline to achieve probabilistically bounded tail latency at minimal cost.

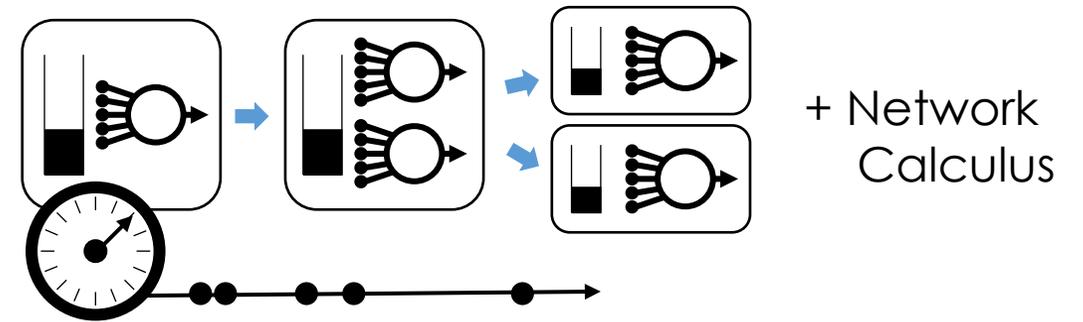


Technical Ideas in the InferLine Project

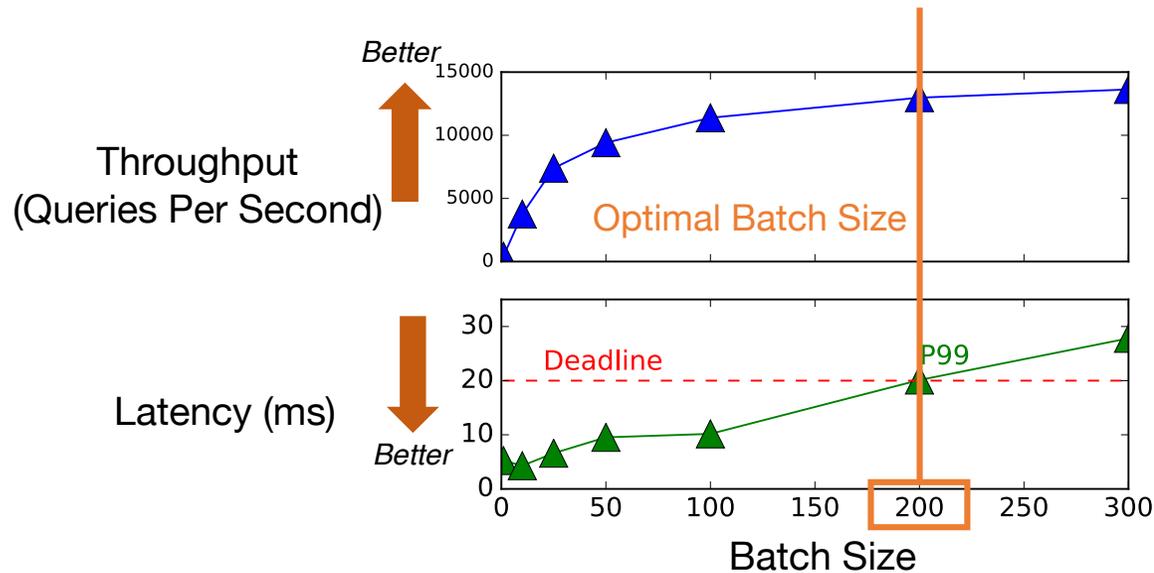
Arrival Process Characterization



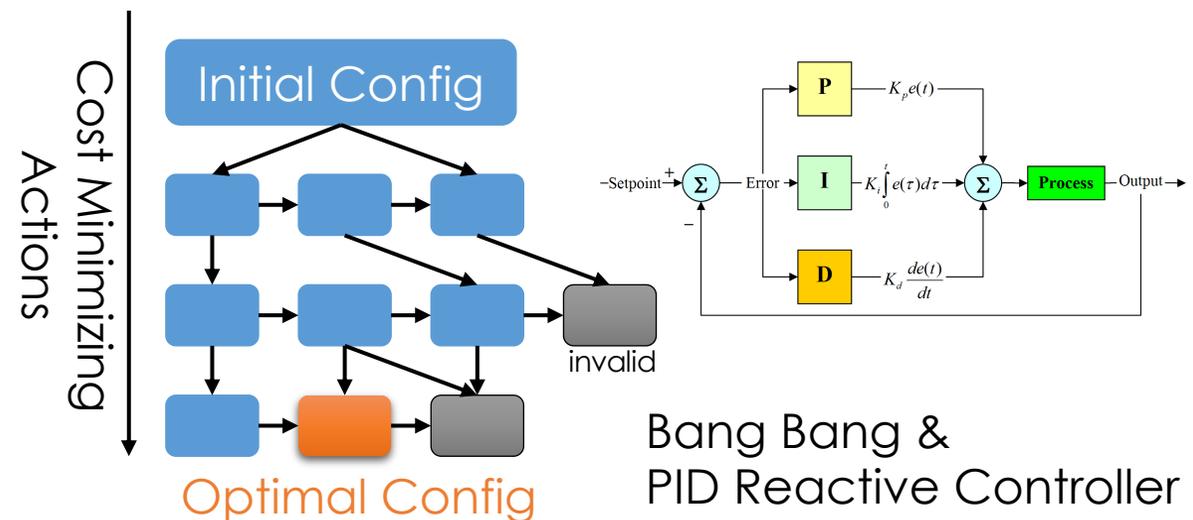
Discrete Event Continuous Time Simulator



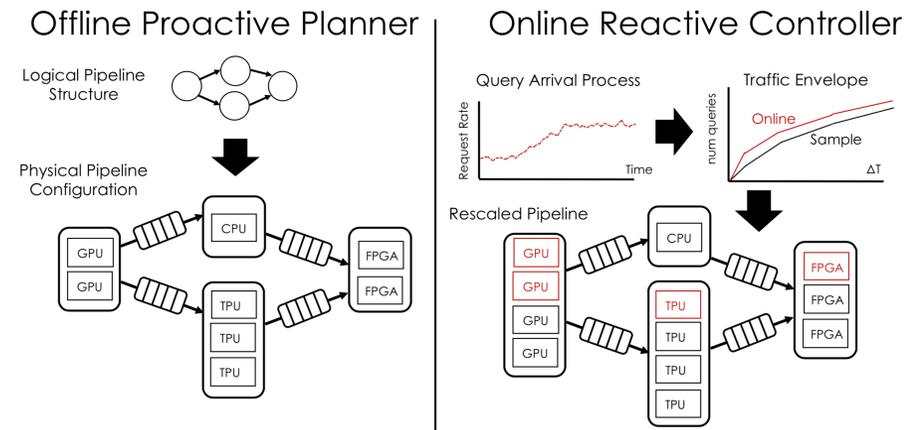
Individual Model Profiles



Proactive and Reactive Optimizer



InferLine: Prediction Pipeline Provisioning and Management for Tight Latency Objectives



➤ Technical Ideas (summary)

- Individual model **performance profiles** + **discrete event simulation** → reason about **end-to-end latency** in the presence of complex **queuing behavior**
- Simple **greedy search heuristic** to configure for each model:
 - Hardware type, number of copies, and batching parameters
- Online re-provisioning using **network calculus** to **optimally resize**

➤ What to look for:

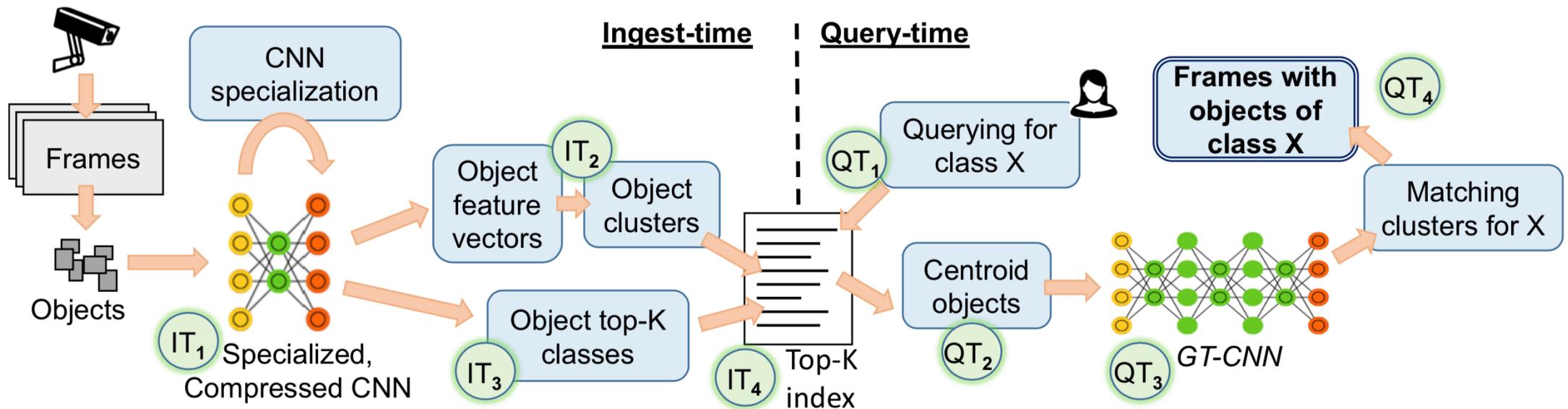
- **Too many ideas not enough contribution?**
- Clarity of presentation

Focus: Querying Large Video Datasets with Low Latency and Low Cost

- **Context:** builds on a line of earlier work
 - [Live Video Analytics at Scale with Approximation and Delay-Tolerance](#) (NSDI'17)
 - [Chameleon: Scalable Adaptation of Video Analytics](#) (SIGCOMM'18)
- **Big Ideas in the Line of Work:**
 - **Trade-off accuracy and latency** in video processing tasks
 - Schedule resources according to **acc.** and **latency requirements**
- **Different “Serving Model”:**
 - Large queries on historical or video streams:
 - *Find all the times where a car and a bike are in a frame.*

Focus: Querying Large Video Datasets with Low Latency and Low Cost

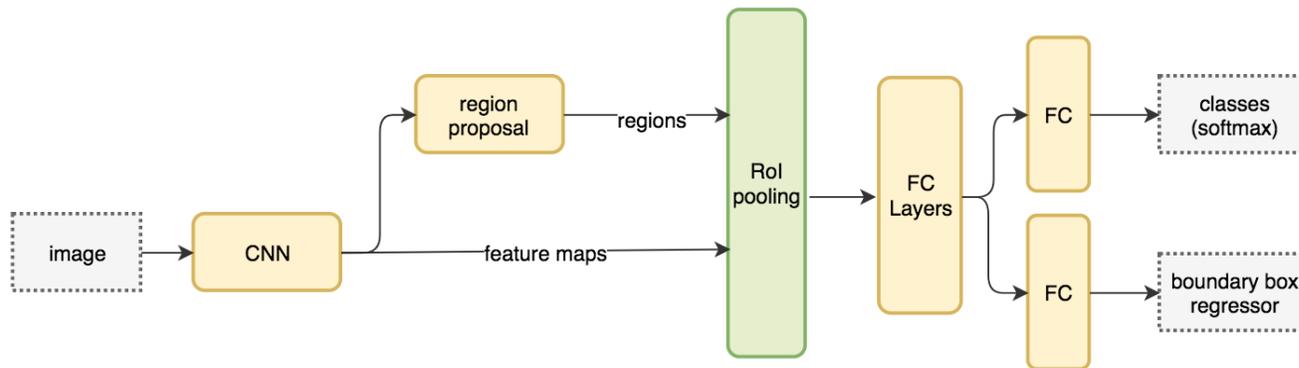
➤ Big Idea(s)



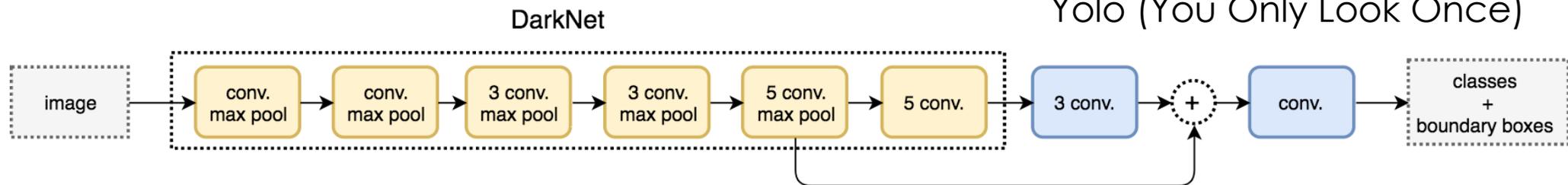
Focus: Querying Large Video Datasets with Low Latency and Low Cost

➤ What to look for:

- Framing of relationships between object detection models and image classification



Faster R-CNN

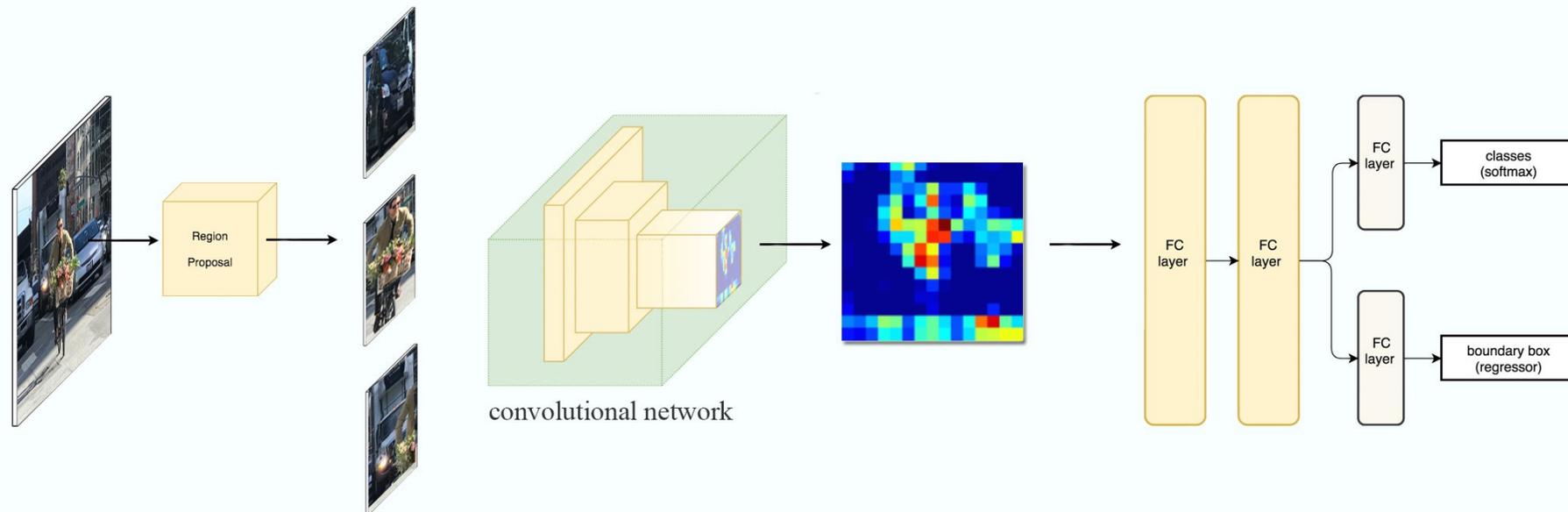


Yolo (You Only Look Once)

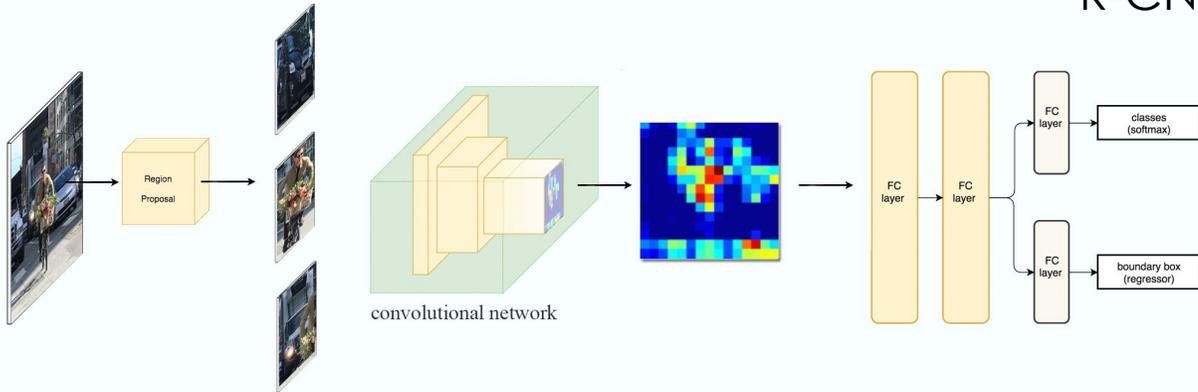
Focus: Querying Large Video Datasets with Low Latency and Low Cost

➤ What to look for:

- Framing of relationships between object detection models and image classification

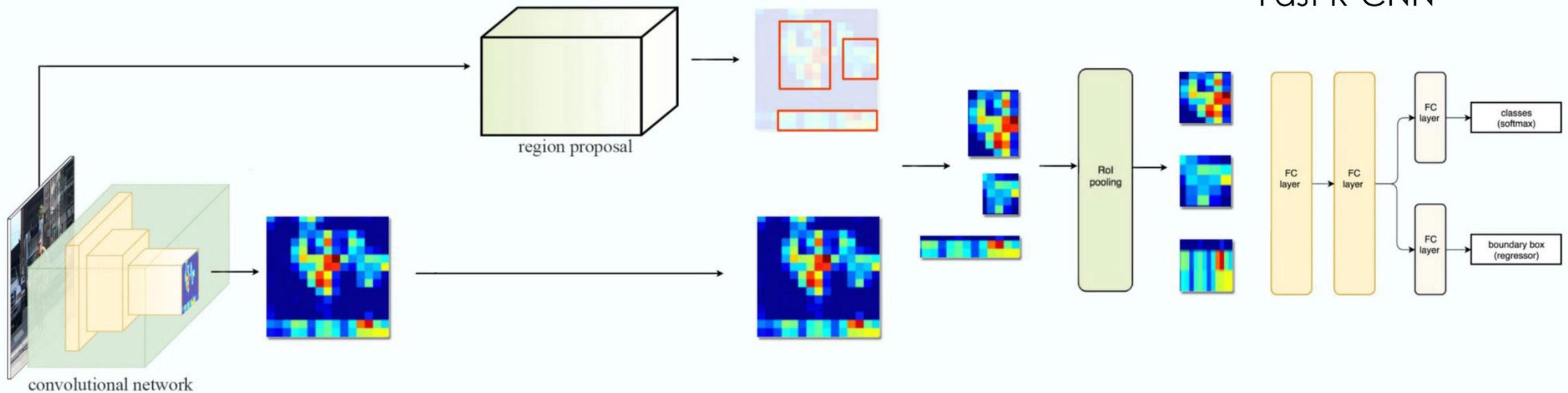


R-CNN



Identify regions of interest and run feature network on each.

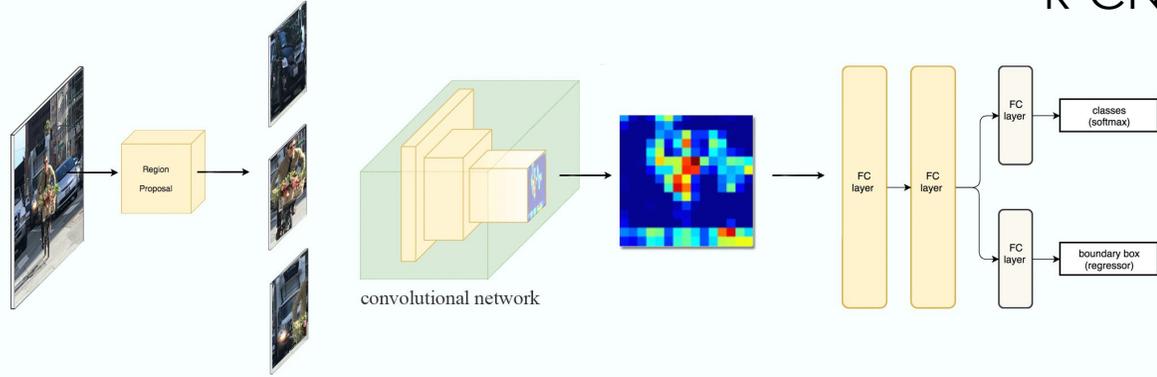
Fast R-CNN



Reuse feature outputs for all externally proposed regions.

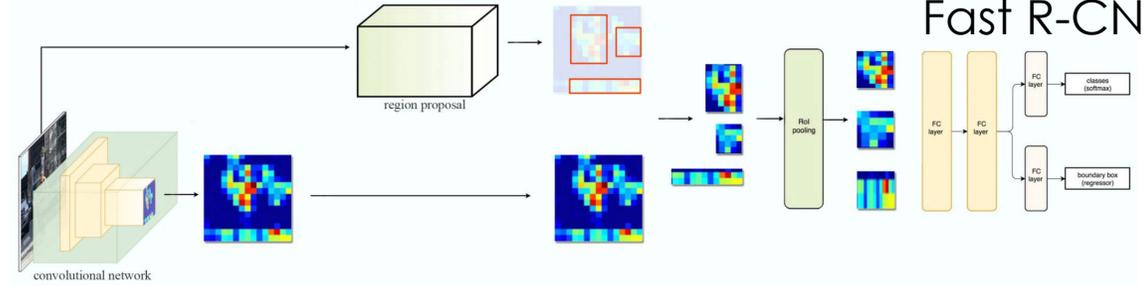
Identify regions of interest and run feature network on each.

R-CNN



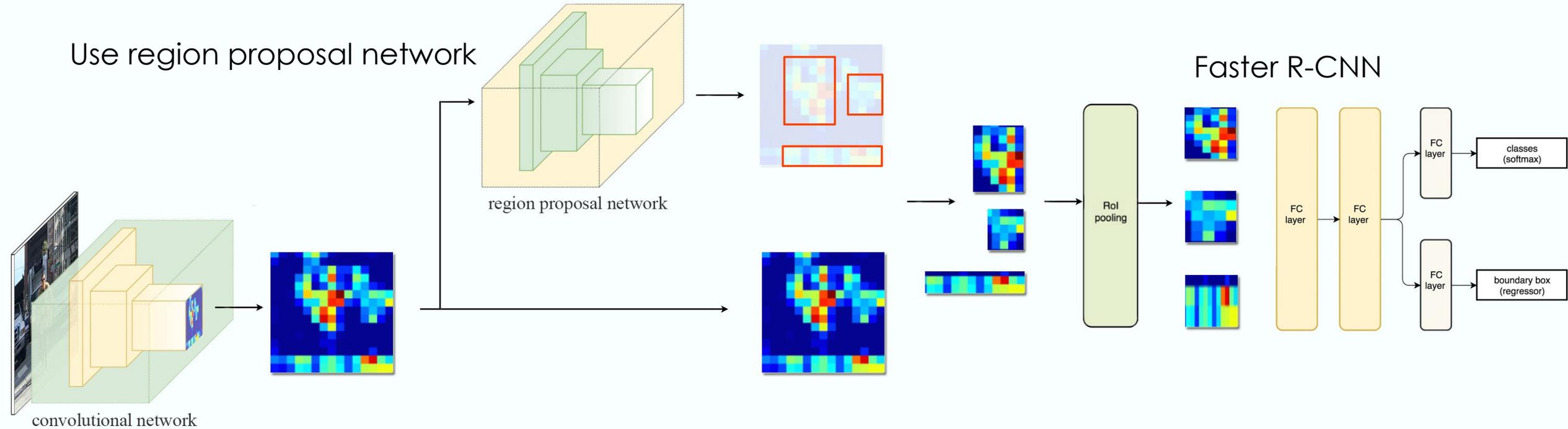
Reuse feature outputs for all externally proposed regions.

Fast R-CNN

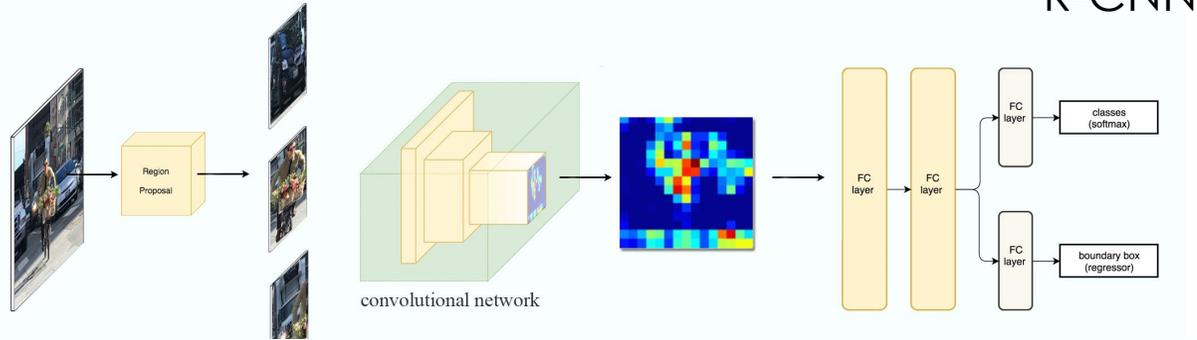


Use region proposal network

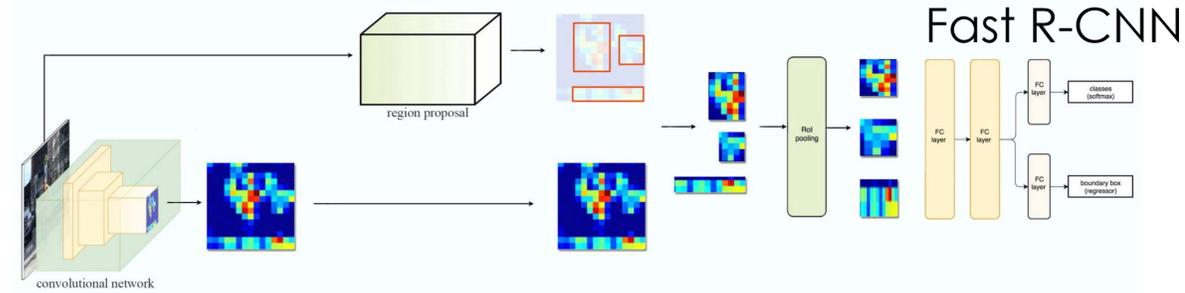
Faster R-CNN



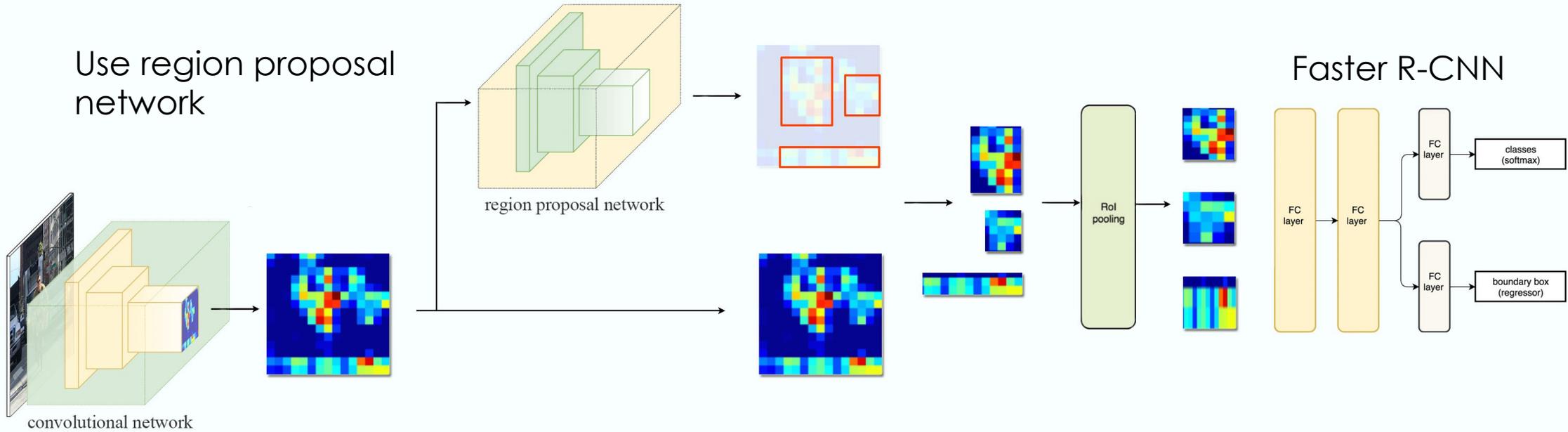
Identify regions of interest and run feature network on each.



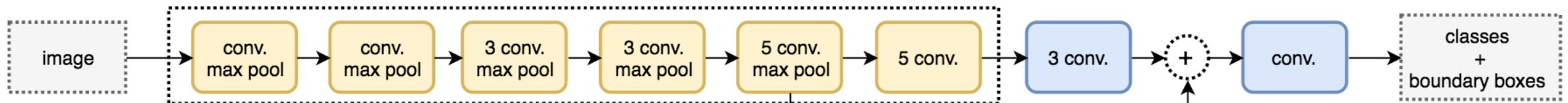
Reuse feature outputs for all externally proposed regions.



Use region proposal network



DarkNet



Yolo (you only look once)

Focus: Querying Large Video Datasets with Low Latency and Low Cost

➤ **What to look for:**

- Framing of relationships between object detection models and image classification
 - Do they leverage the structure of object detection models?
- Split between ingest and query time computation
- Tradeoff between accuracy and latency → how is it evaluated?
- Presentation → many optimizations explored, do they fit together?

Done!

Cascaded Predictions

IDK Prediction Cascades

Simple models for simple tasks



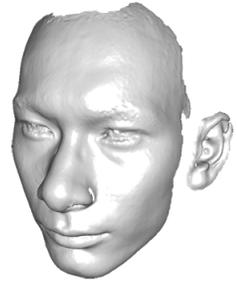
Xin Wang



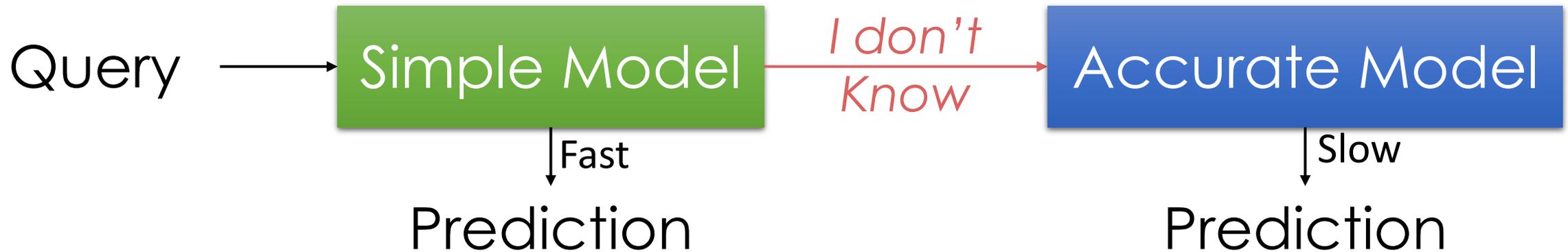
Yika Luo



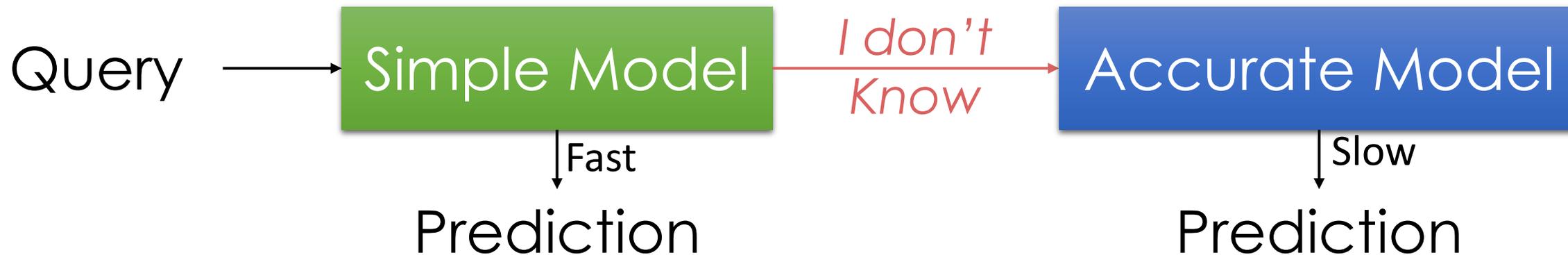
Zi-Yi Duo



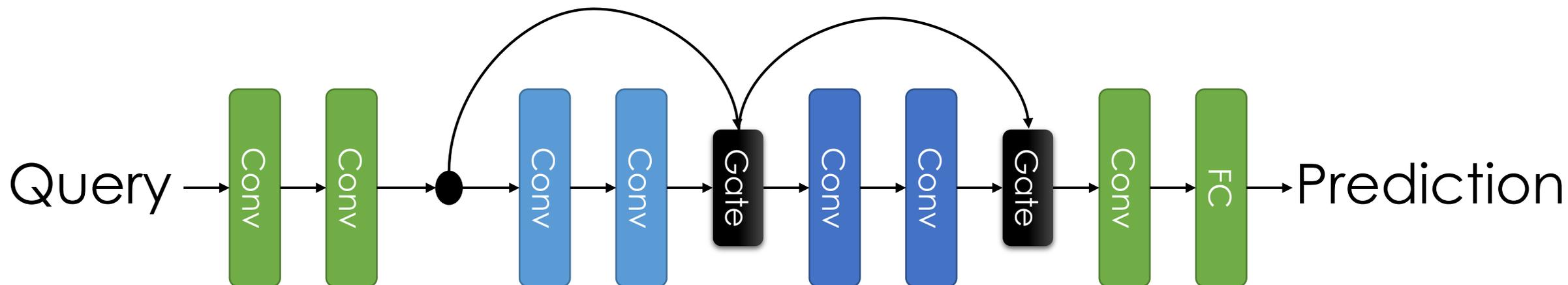
Fisher Yu

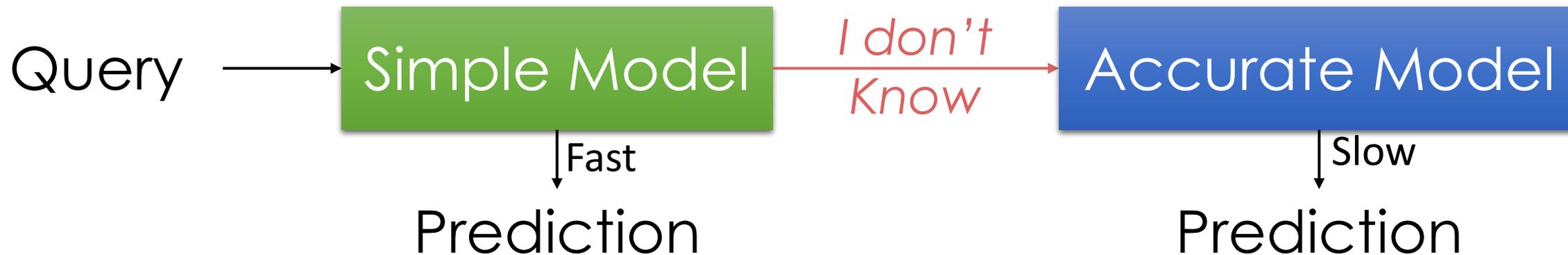


Learn to combine **fast (inaccurate) models** with **slow (accurate) models** to maximize accuracy while reducing computational costs.

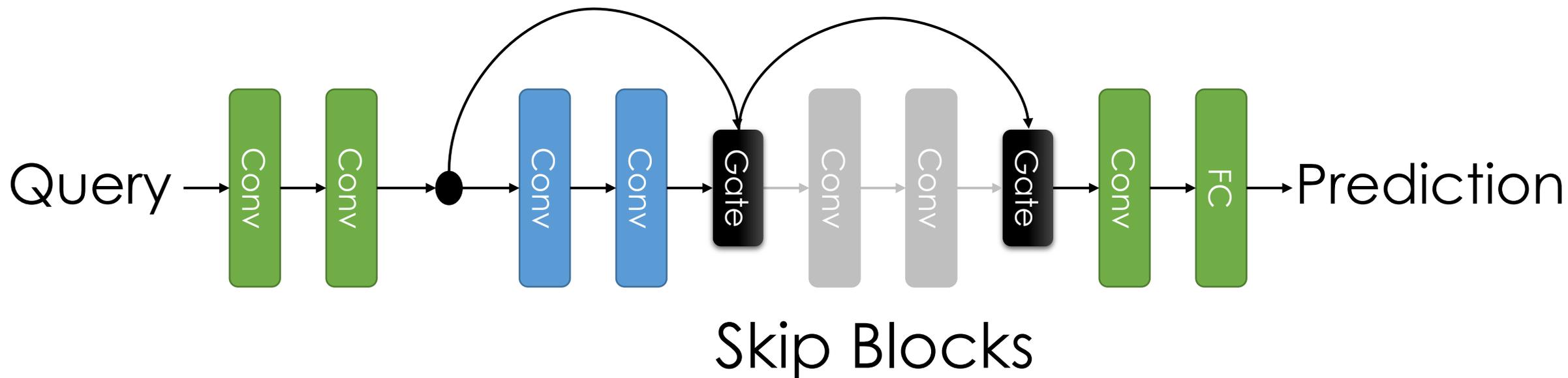


SkipNet: dynamic execution within a model

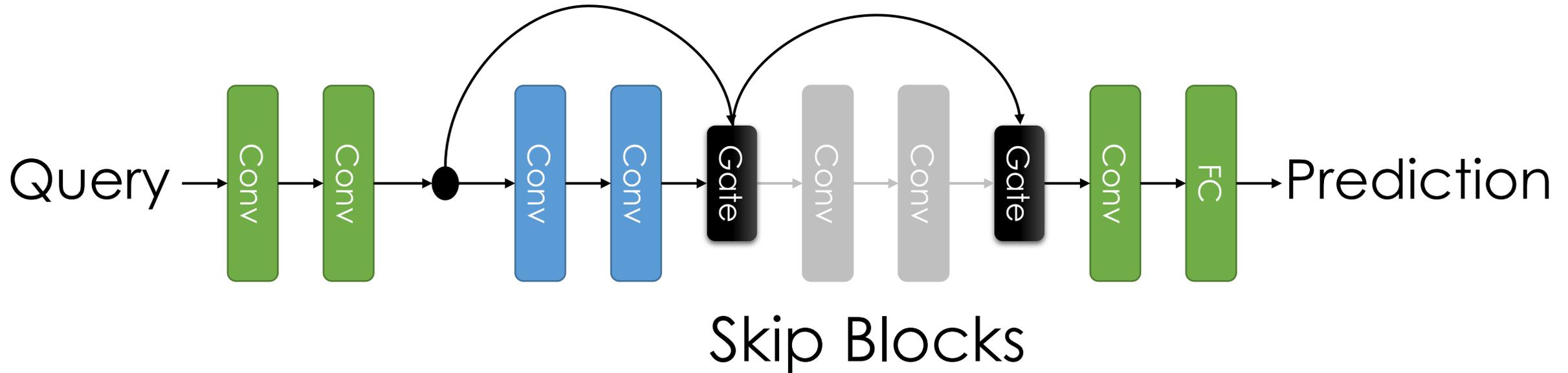




SkipNet: dynamic execution within a model

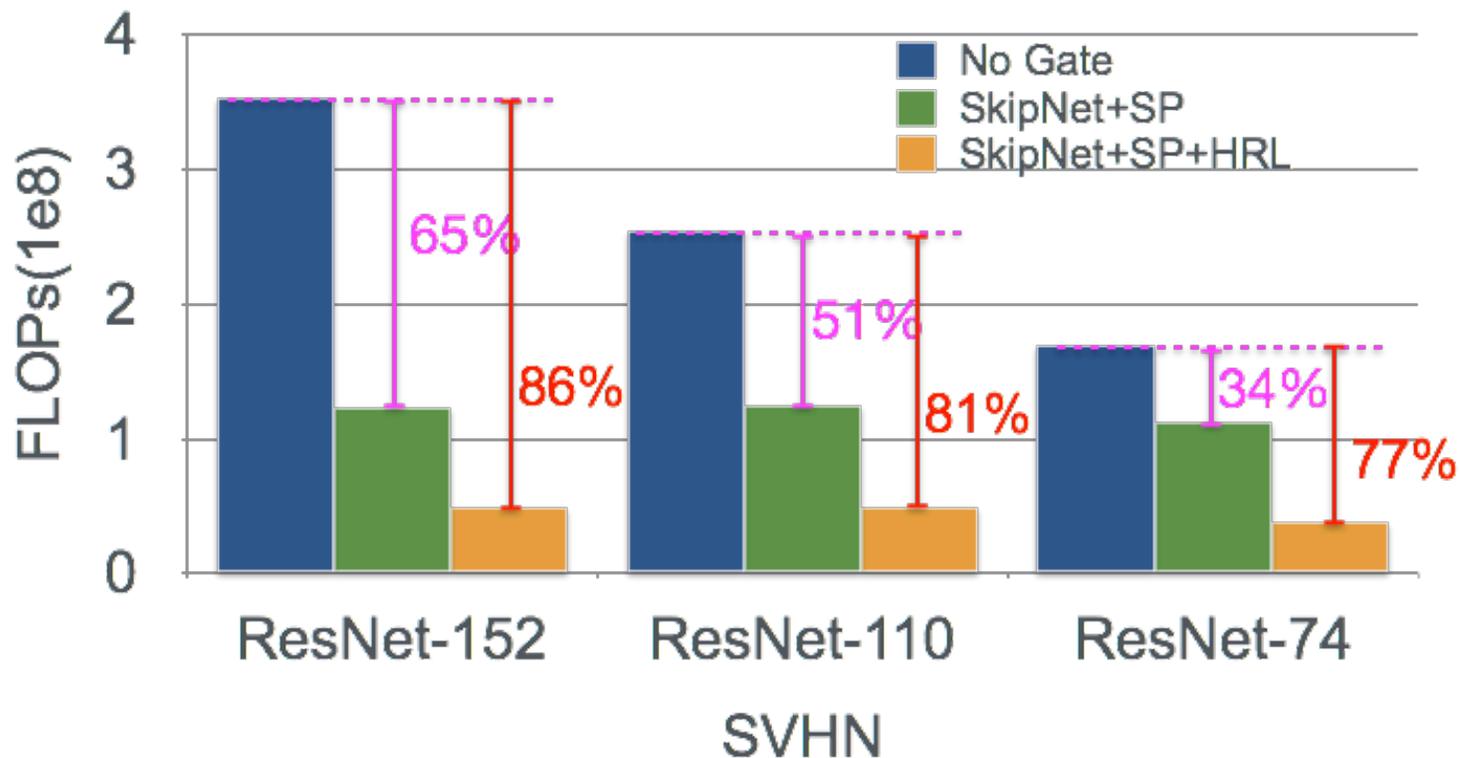


SkipNet: dynamic execution within a model



- Combine **reinforcement learning** with **supervised pre-training** to learn a gating policy

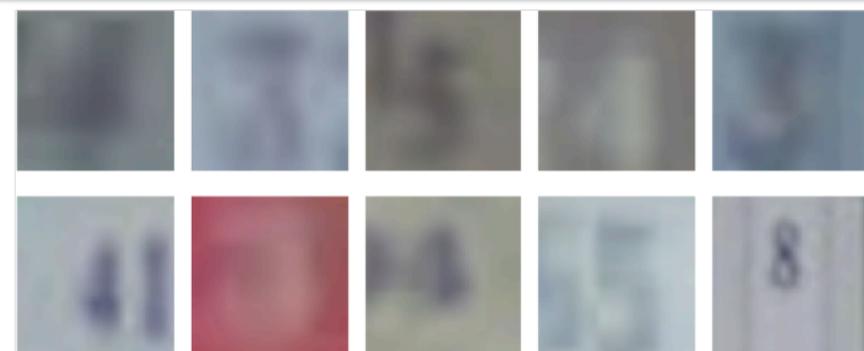
SkipNet Performance



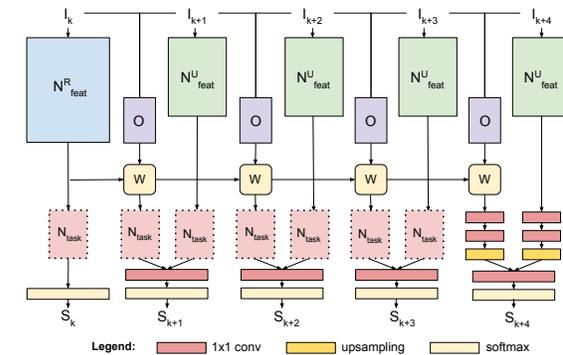
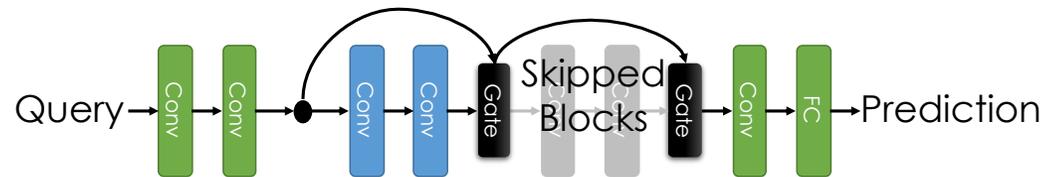
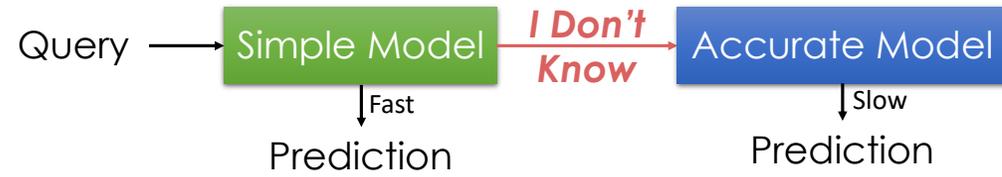
Easy Images
Skip **Many** Layers



Hard Images
Skip **Few** Layers

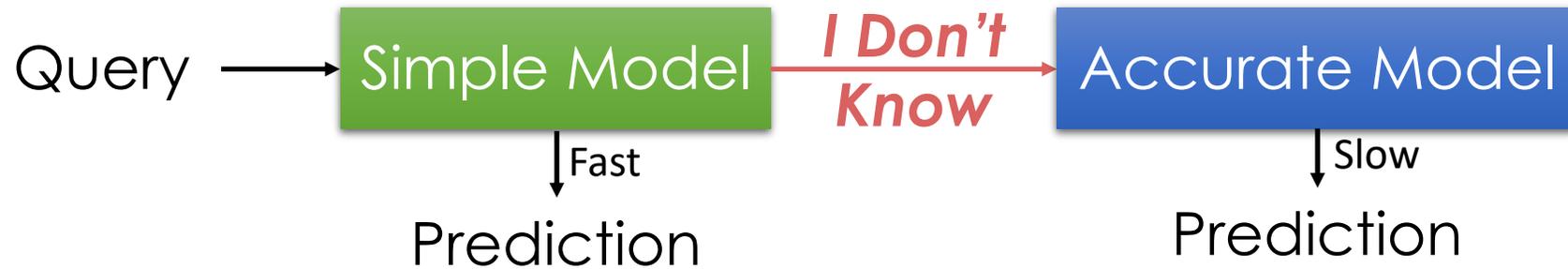


Efficient Neural Networks

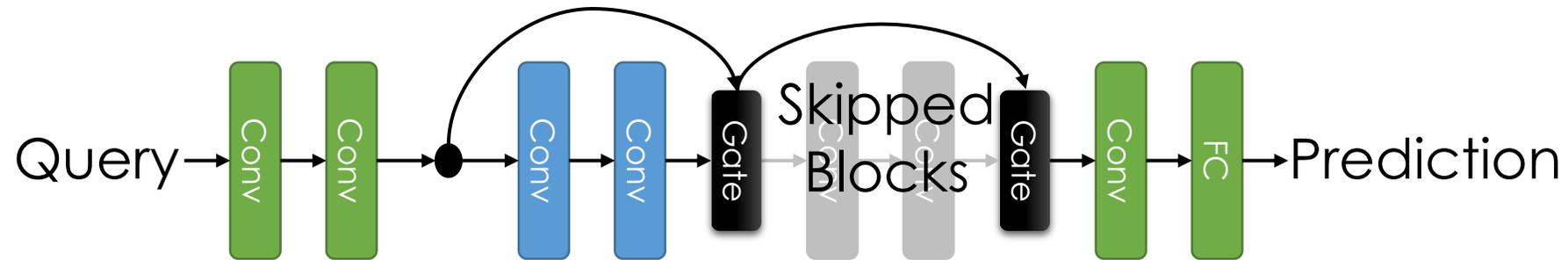


Dynamic Networks for **fast** and **accurate** inference

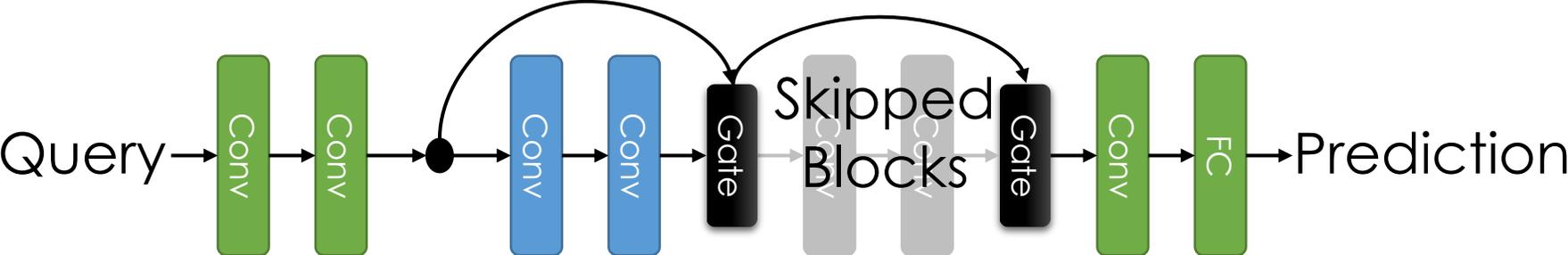
IDK Cascades: Using the fastest model possible [UAI'18]



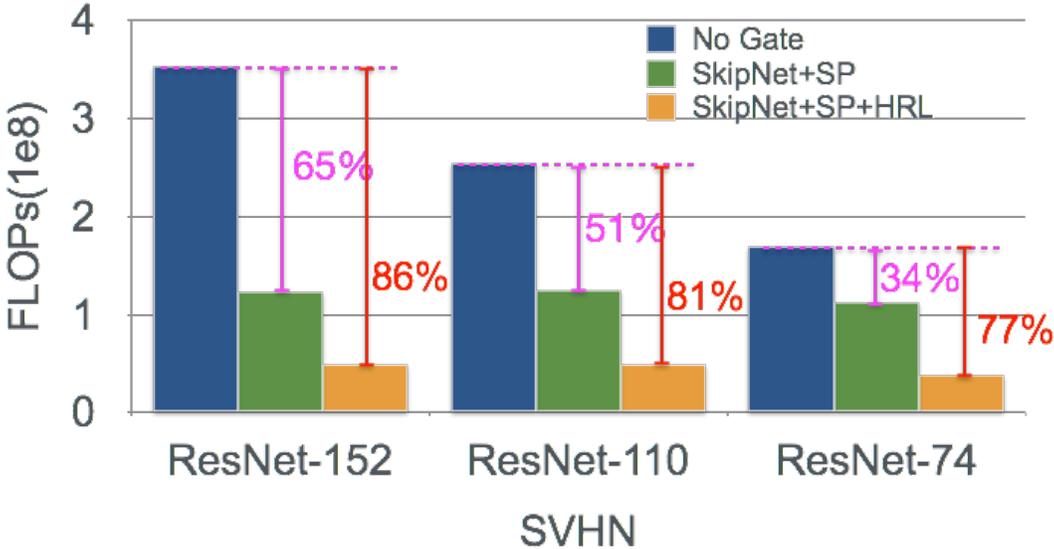
SkipNet: dynamic execution within a model [ECCV'18]



SkipNet: dynamic execution within a model [ECCV'18]



Large Reductions in FLOPs

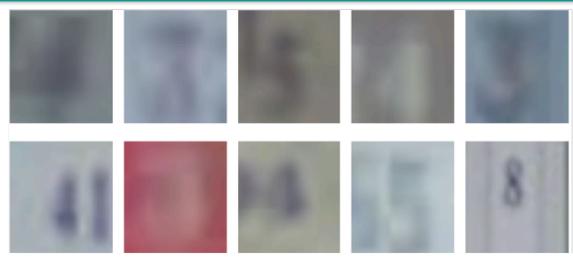


Skip more layers on clear images

Easy Images
Skip **Many** Layers

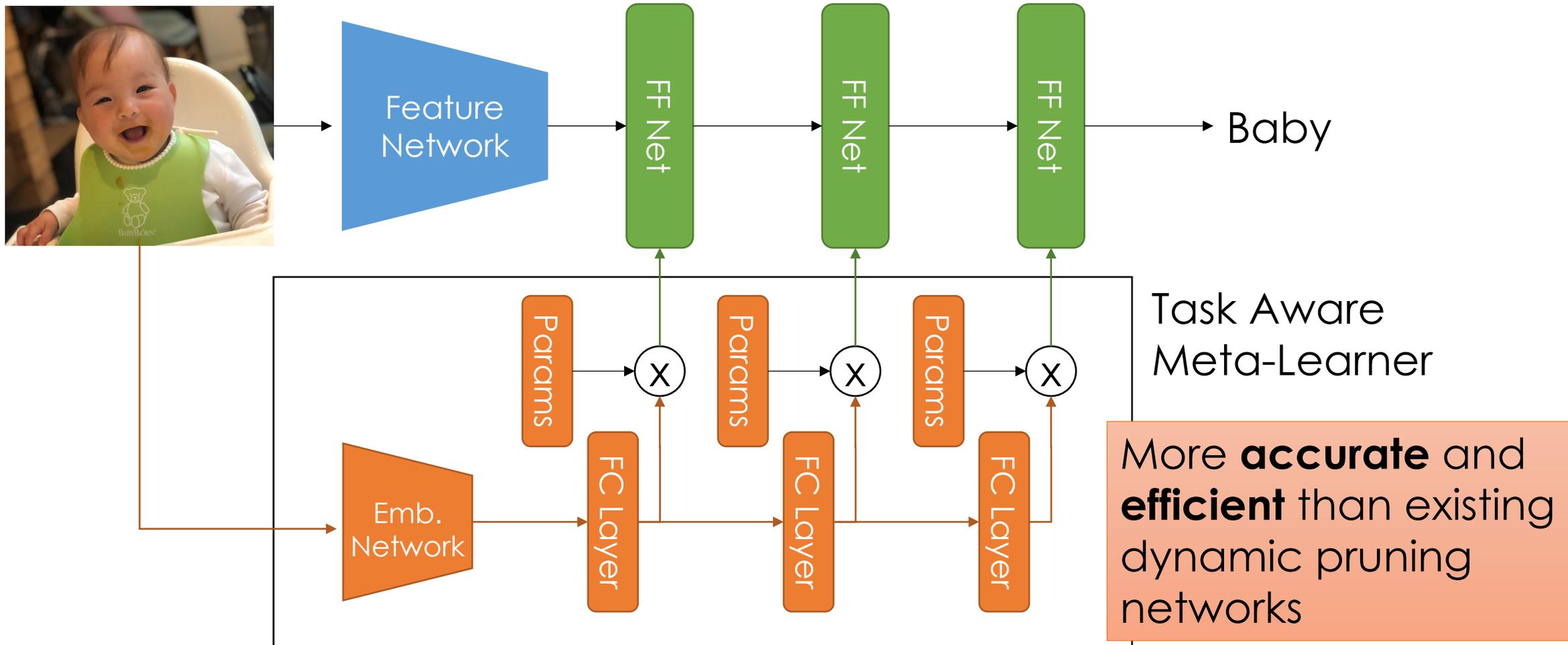


Hard Images
Skip **Few** Layers



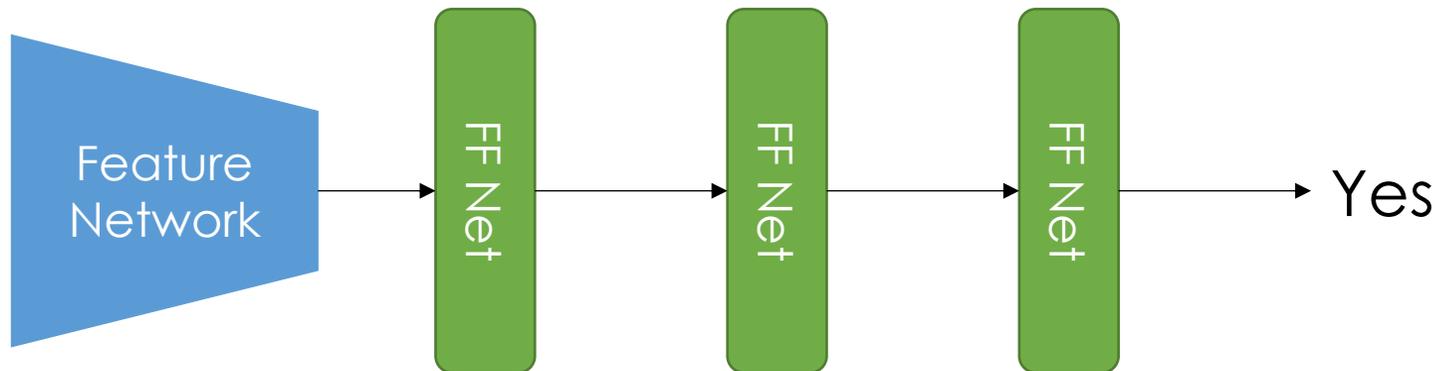
Task Aware Feature Embeddings

[CVPR'19]



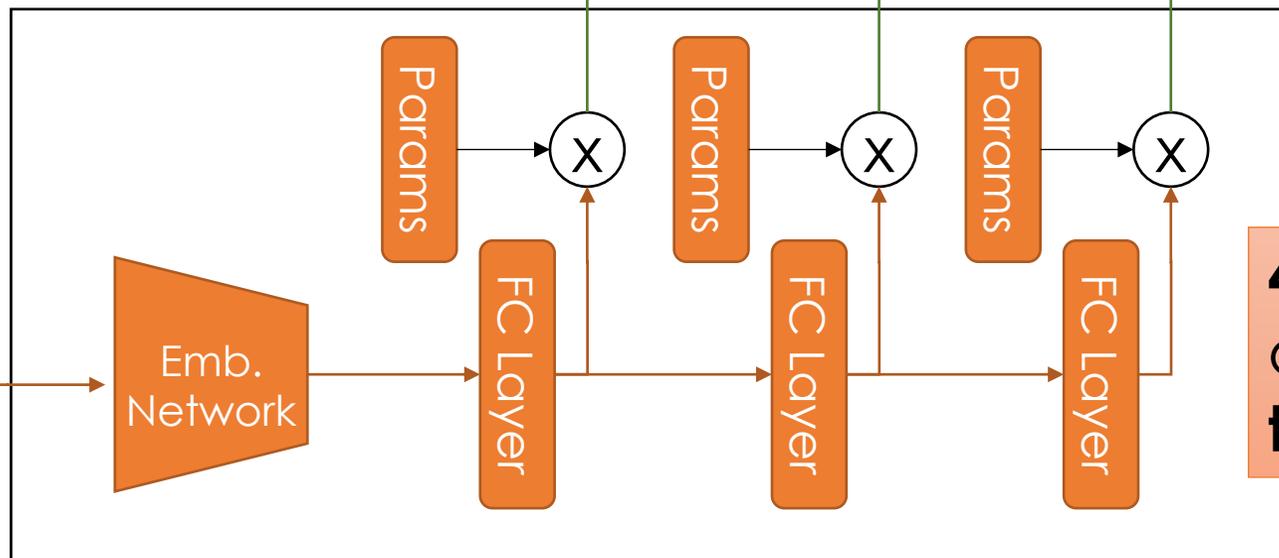
Task Aware Feature Embeddings

[CVPR'19]



Task Description:

"Smiling Baby"

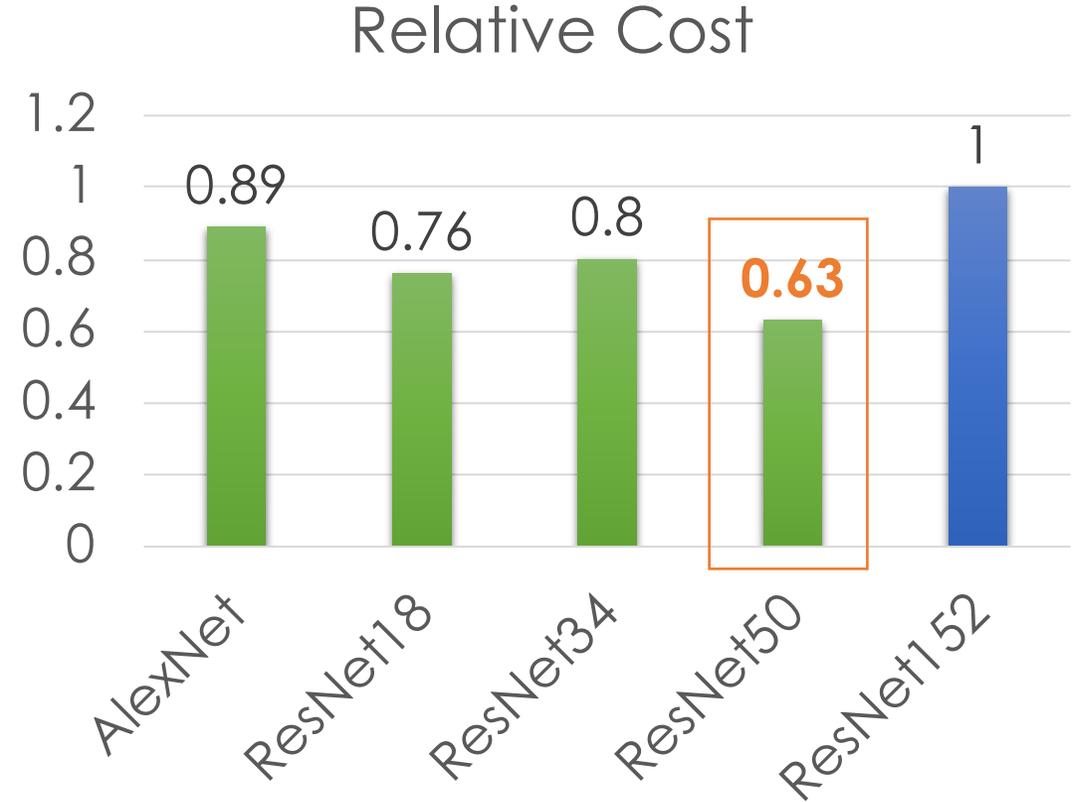
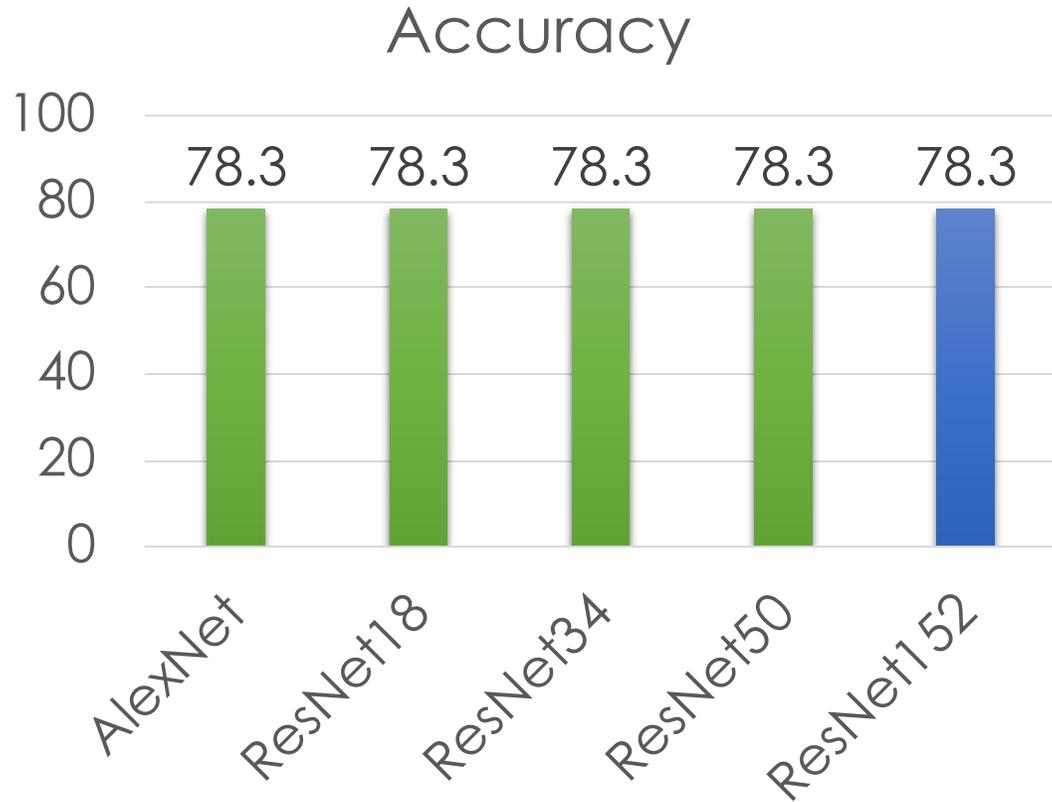
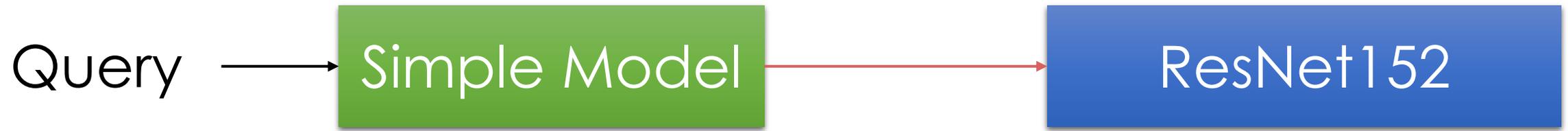


Task Aware
Meta-Learner

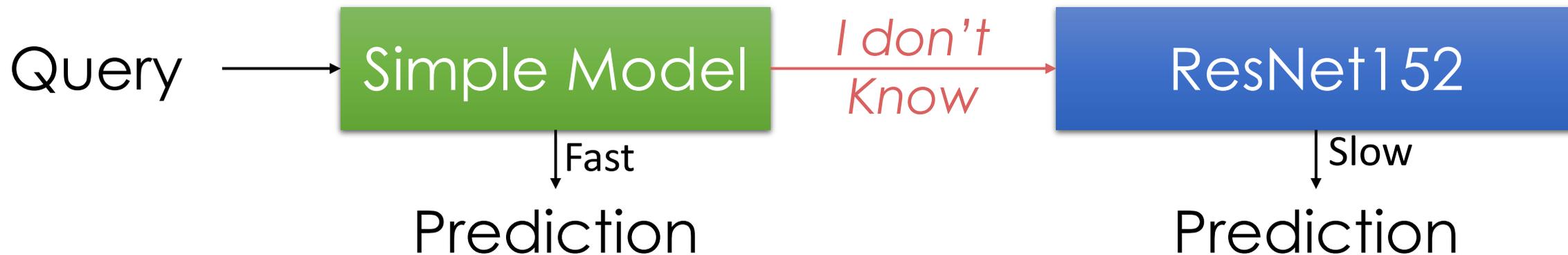
4 - 15% improvement
on **attribute-object**
tasks

Leverage **motion** to improve the **speed** and **accuracy** of **semantic segmentation**

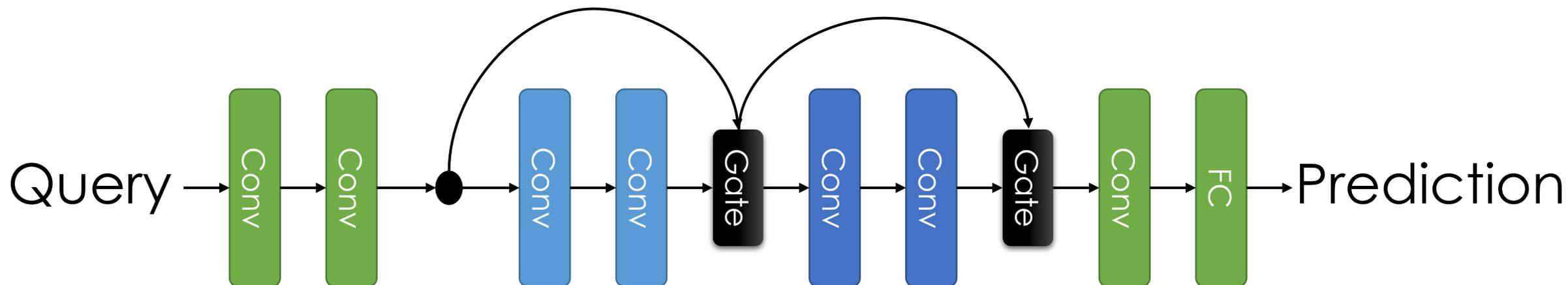


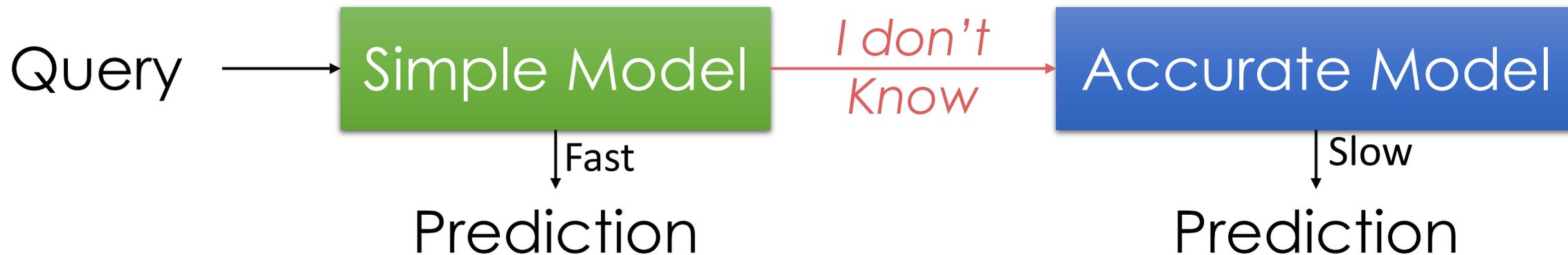


37% reduction in runtime
@ no loss in accuracy

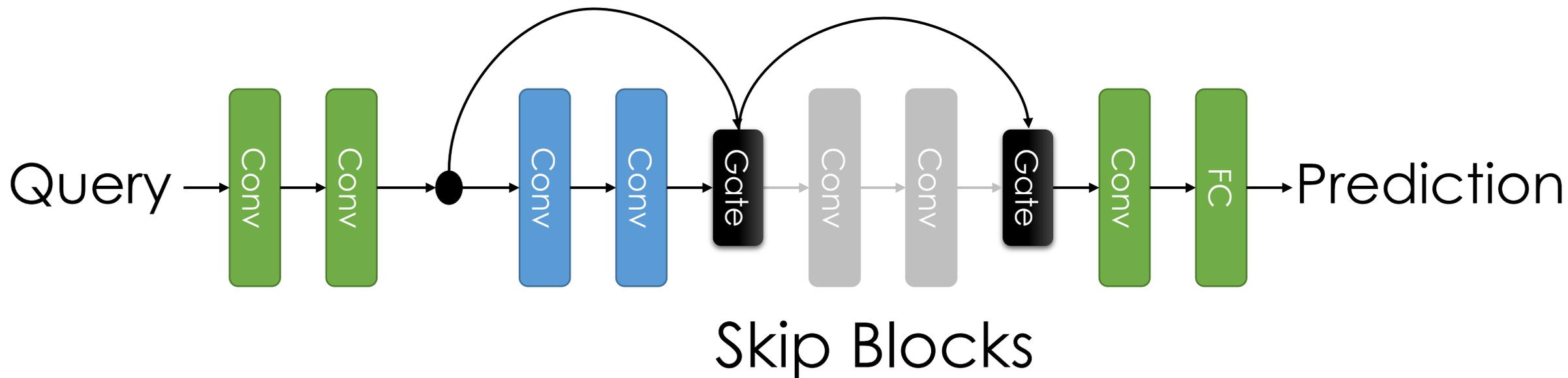


➤ Cascades within a Model

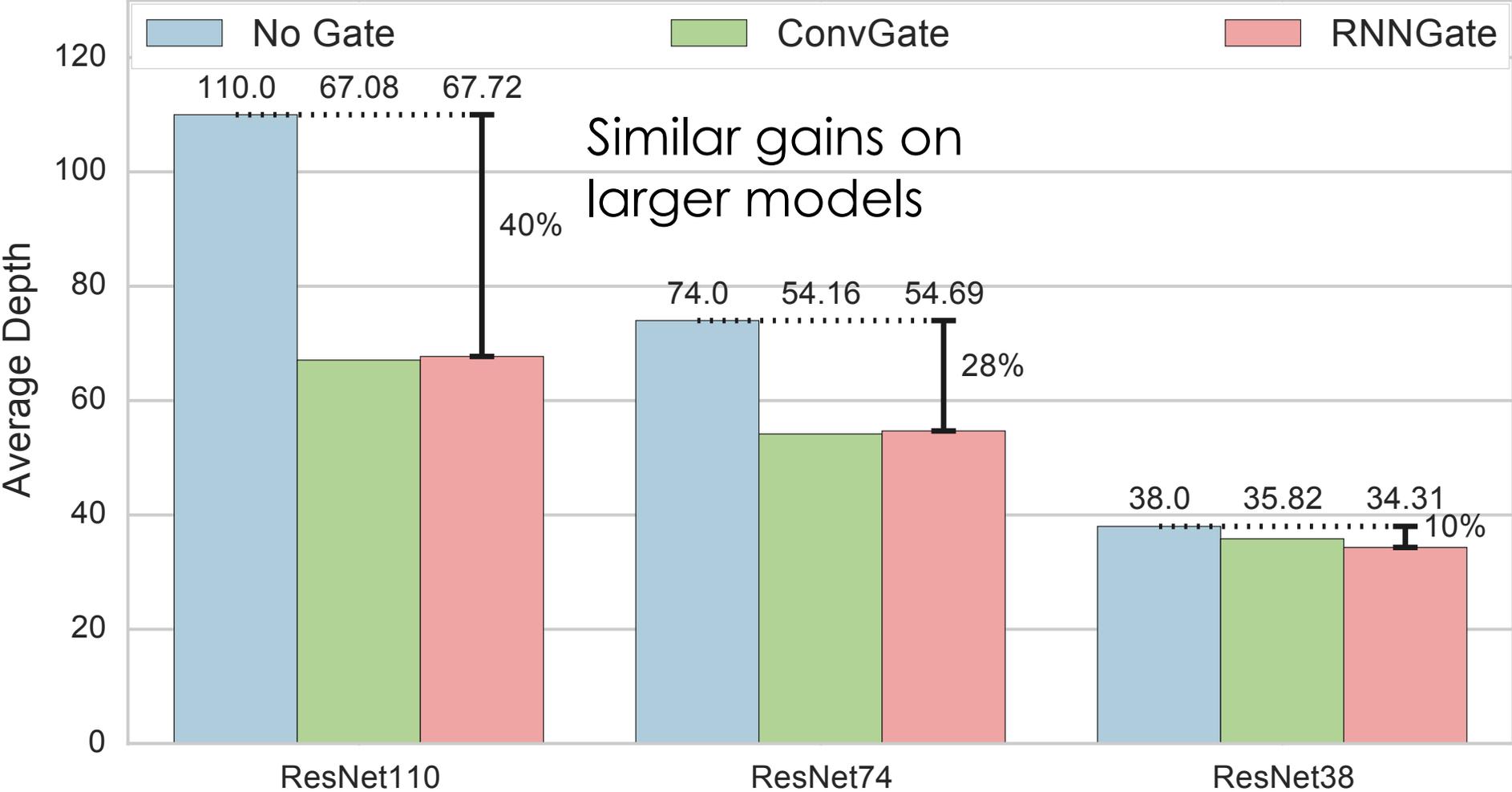




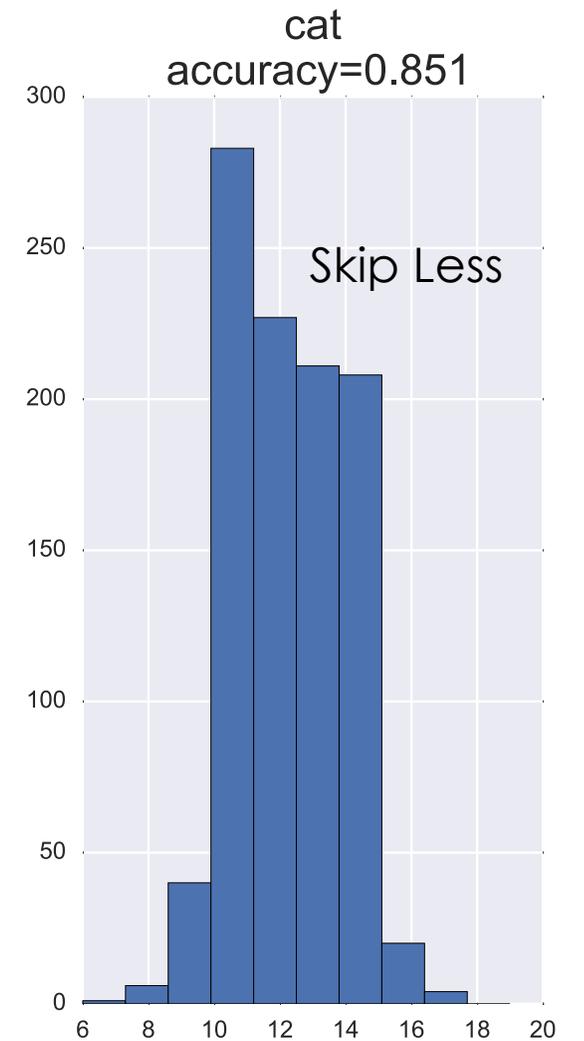
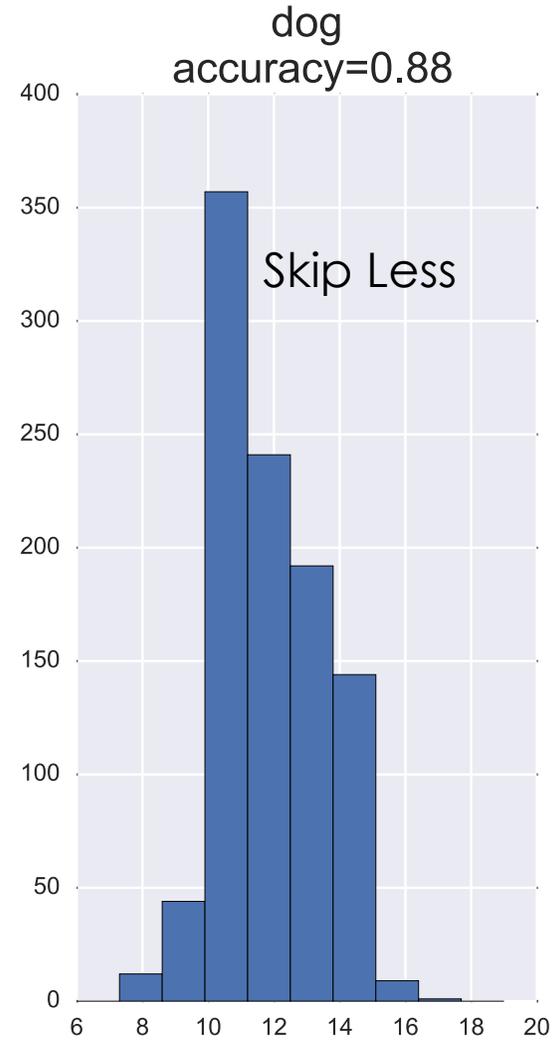
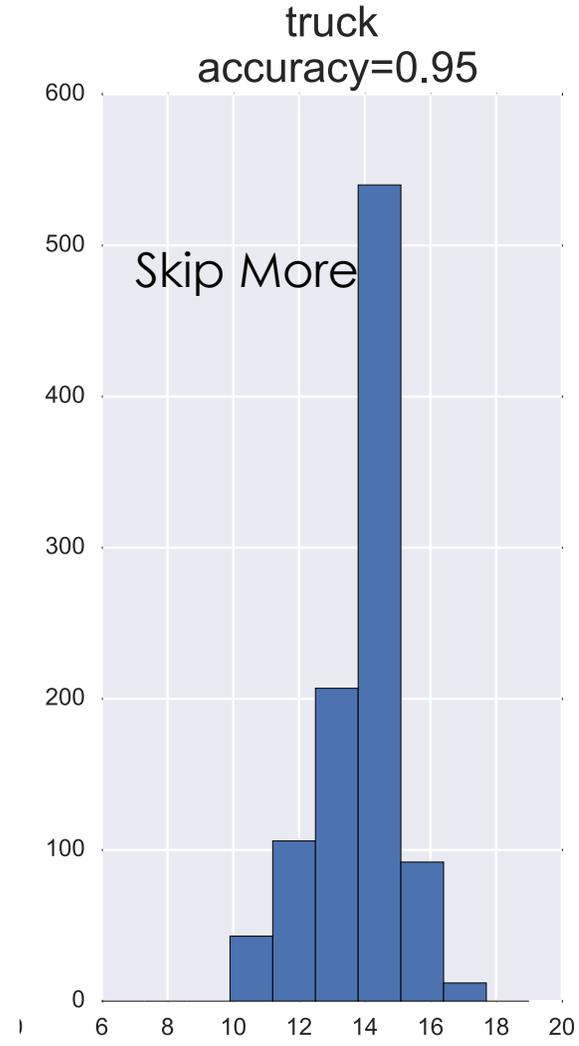
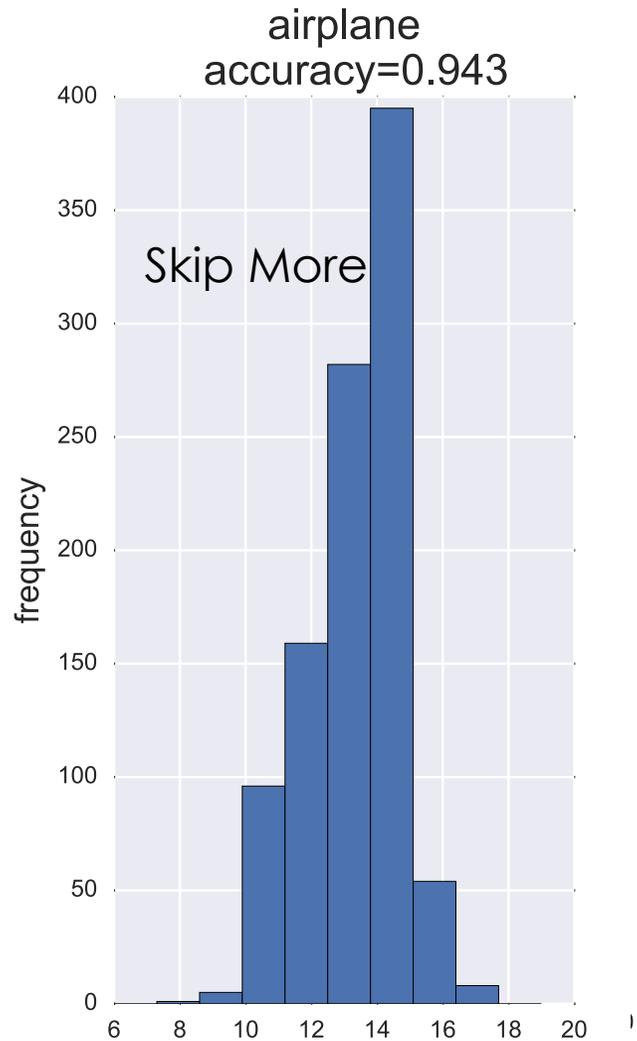
➤ Cascades within a Model



Cascading reduces computational cost



Easy 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

