

The background features a vibrant, multi-colored gradient. The top left is a deep blue, transitioning through purple and magenta to a bright orange and yellow in the center, and finally fading into a light blue and white on the right. A diagonal line separates the darker blue on the left from the lighter blue on the right.

AWS  
re:Invent

**AIM 329**

# Using deep learning to track wildfires and air quality

## **Sanjay Padhi**

Head of AWS Research, US Education  
Amazon Web Services

## **Feng Yan**

Assistant Professor  
University of Nevada, Reno



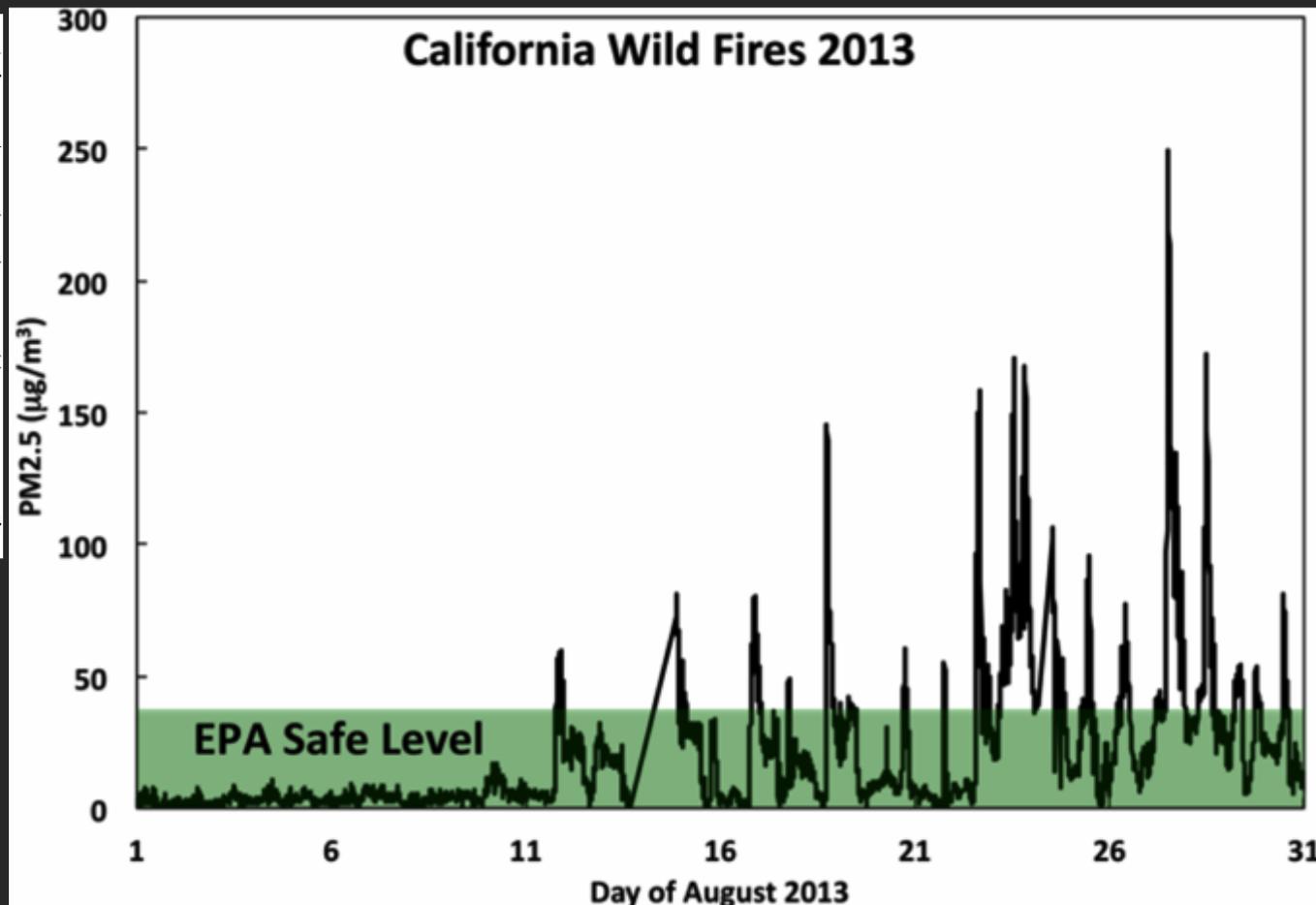
**NSF BIGDATA: IA: Collaborative Research: Protecting Yourself from Wildfire Smoke: Big Data Driven Adaptive Air Quality Prediction Methodologies, 2019–2022**



Feng Yan, Evgenia Smirni, Lei Yang, Heather Holmes,  
Sponsored by NSF and AWS, supported by ALERTWildfire, UNR OIT

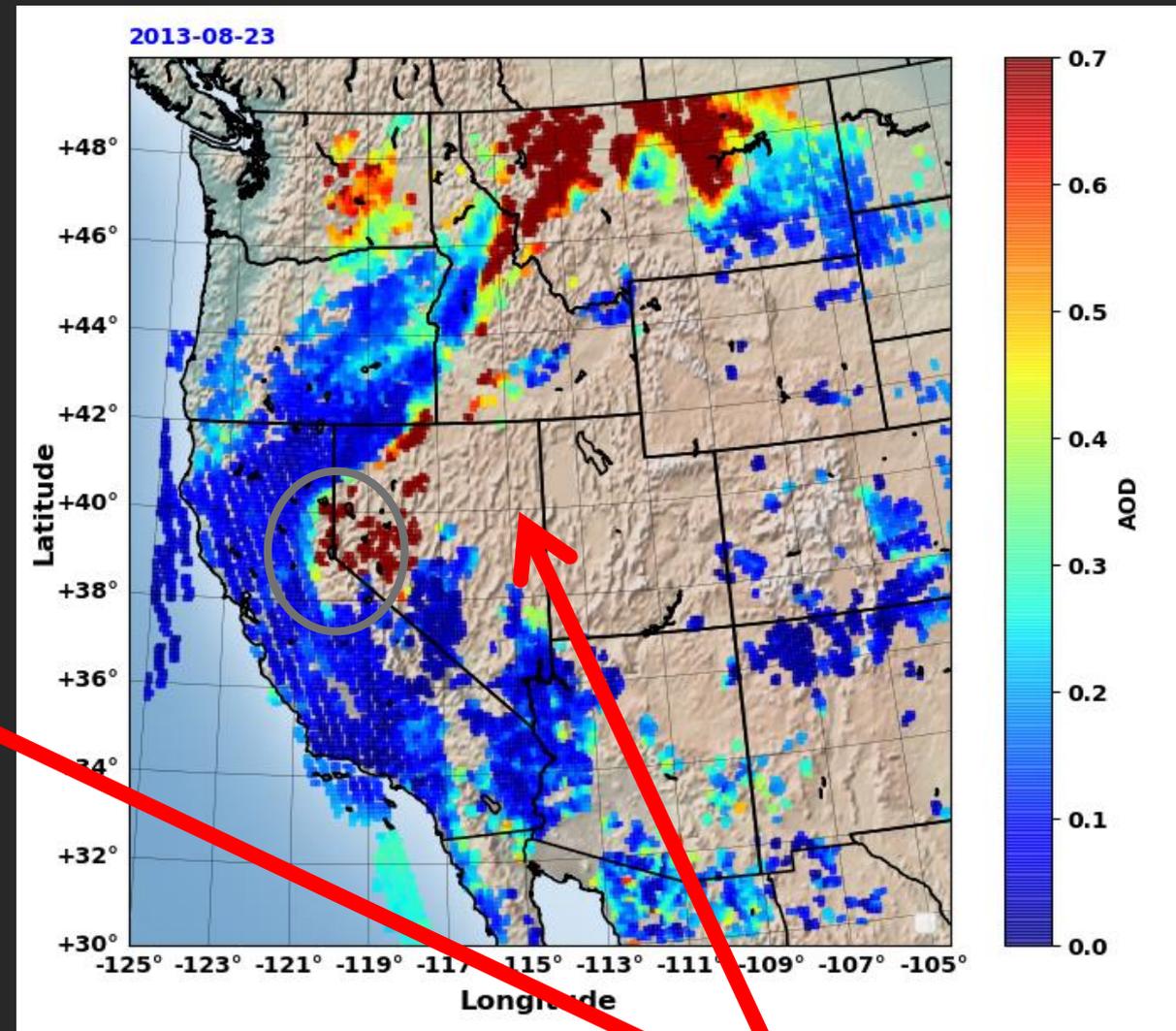
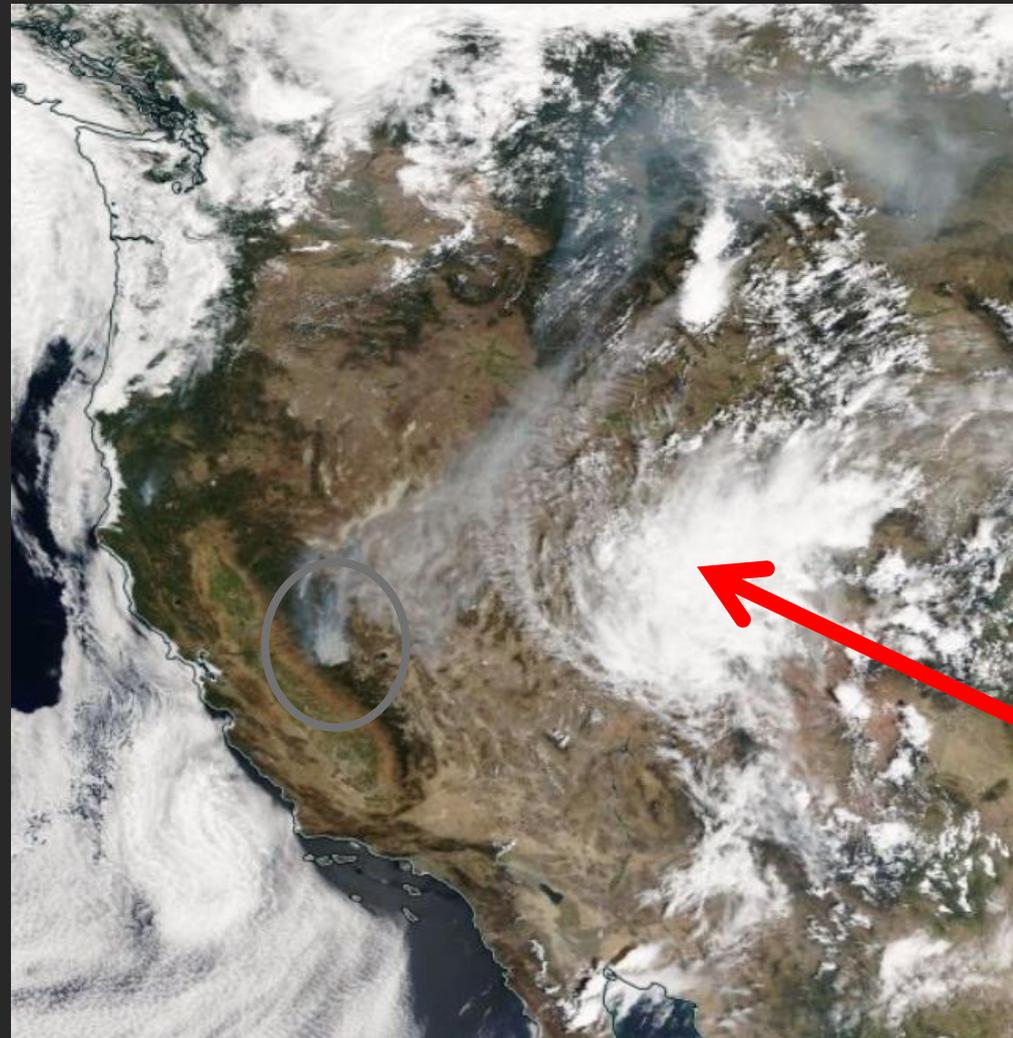
# Why is wildfire smoke dangerous?

- Human health impacts result from wildfire smoke exposure
- Visibility and radiative forcing impacts for climate
- Increasing drought conditions in western US can lead to more fires



# Fire/smoke and clouds: Missing data

August 23, 2013, Yosemite Rim fire in California

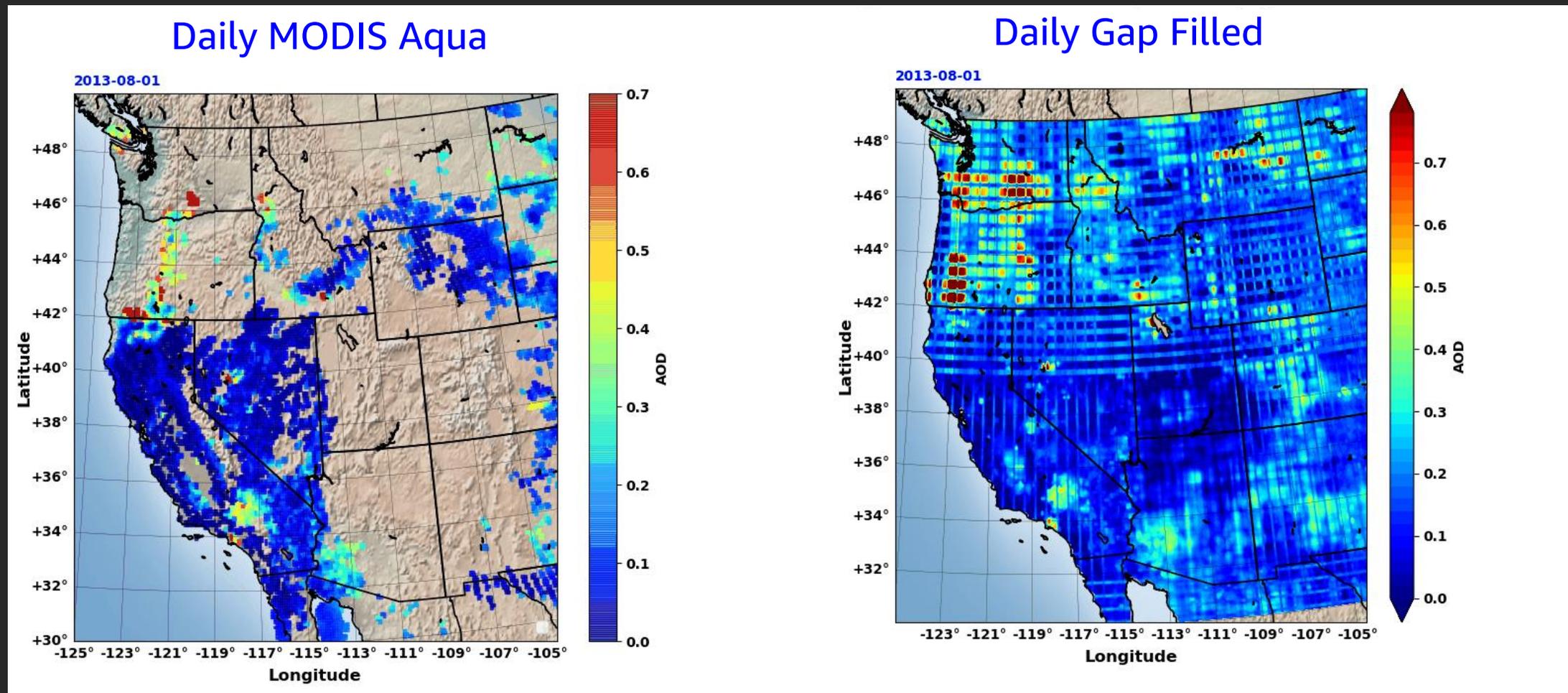


Satellites cannot measure the air pollution below clouds

Missing predictions at a fine-grained level

# MODIS satellite aerosol optical depth

- Example California fire in August 2013
- Daily, 10-km resolution, air pollution estimates from satellites
- Use machine learning to recover missing satellite data



# Smoke spreading quickly

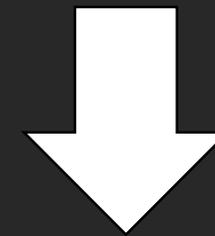
## Reno smoke on August 18, 2013

The smoke was from the American River fire near Sacramento, CA

**Reno, Nevada 24 Hour streaming video, looking to south.**

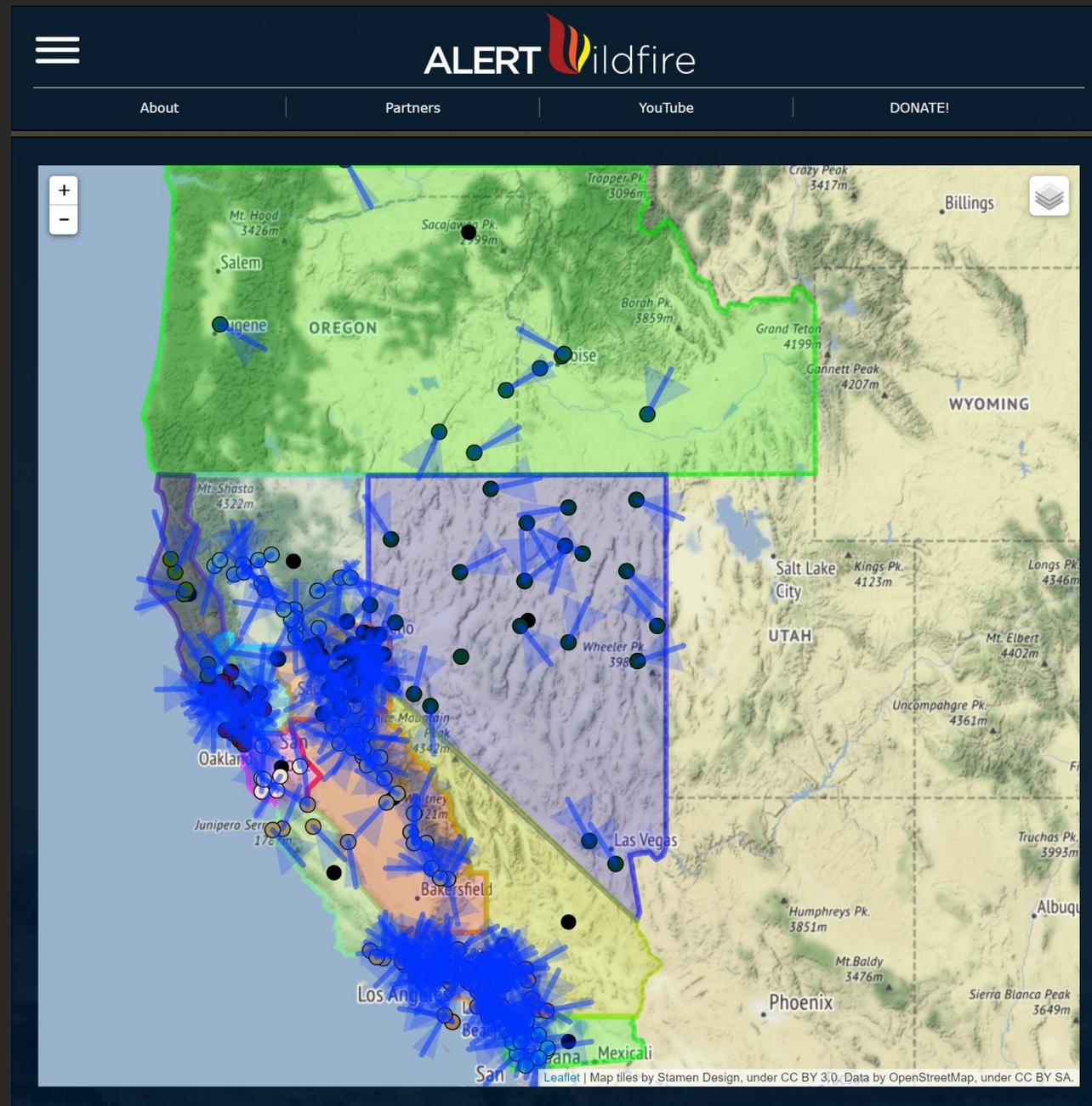


Wildfire smoke can travel  
very fast!



Real-time, fine-grained  
monitoring and prediction

# ALERTWildfire camera network



ALERTWildfire is a fire camera network and associated tools

Use pan-tilt-zoom (PTZ) cameras

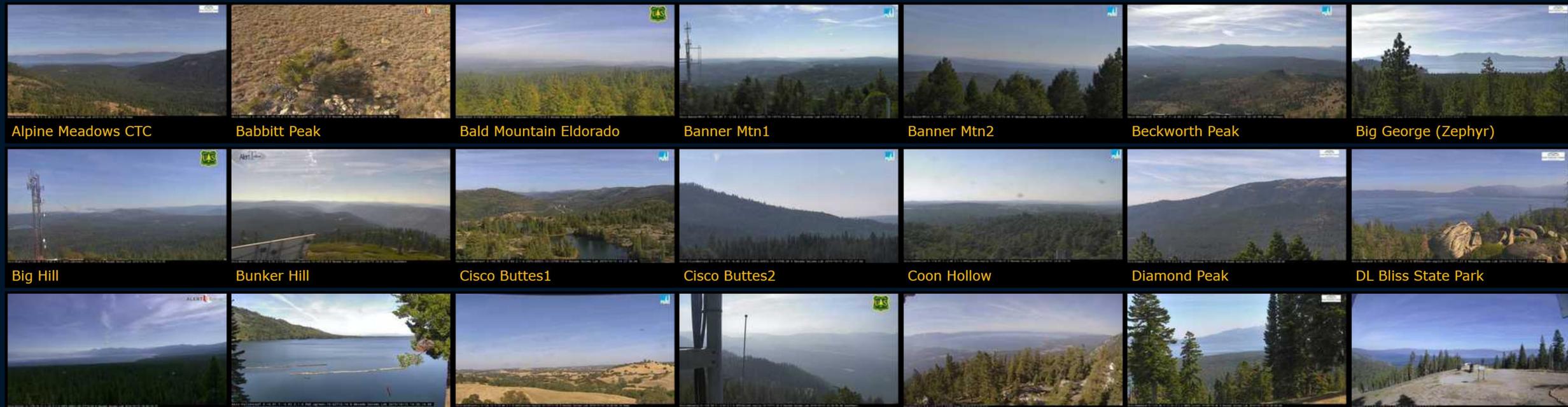
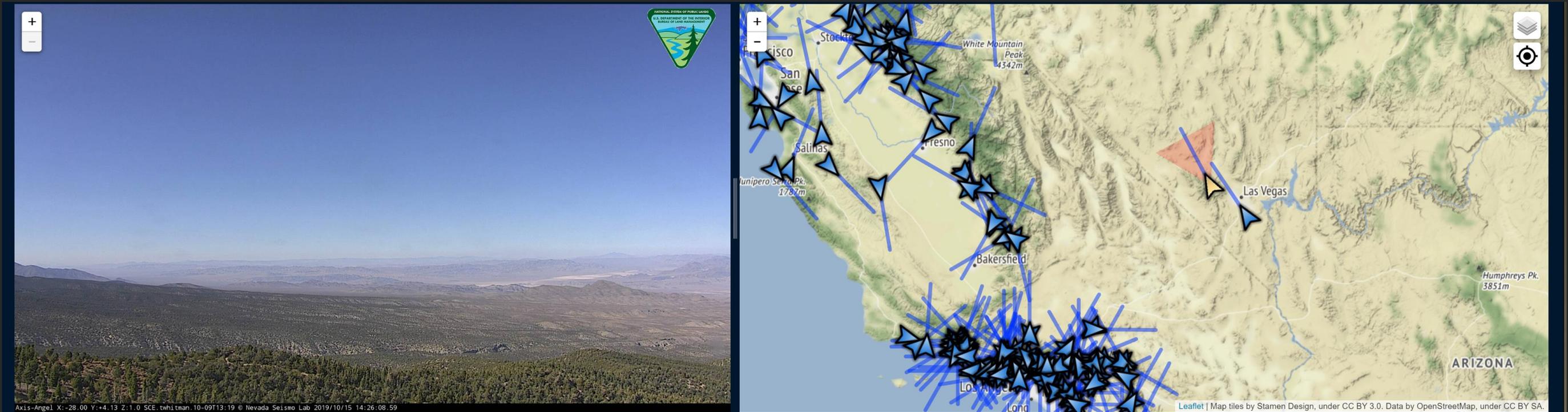
Help firefighters and first responders

Consortium of University of Nevada, Reno (UNR), University of California, San Diego (UCSD), and University of Oregon (UO)



<http://www.alertwildfire.org/>

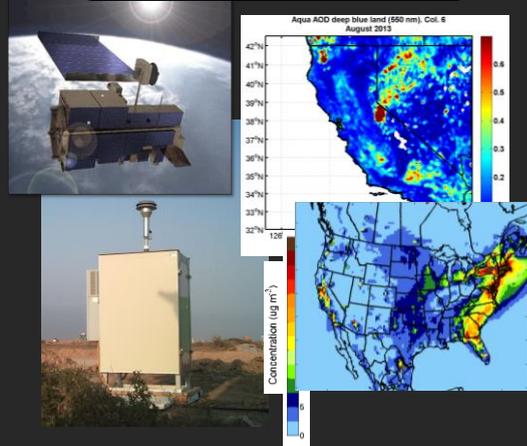
# ALERTWildfire camera network



# Architecture

## Big data

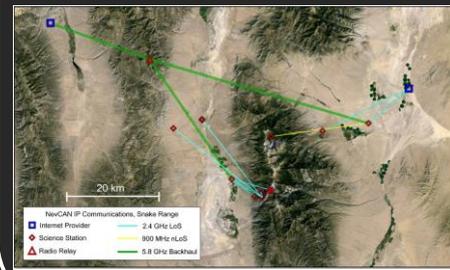
### Existing data sources



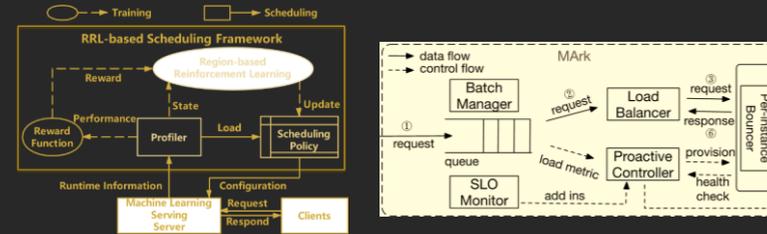
### New camera network data



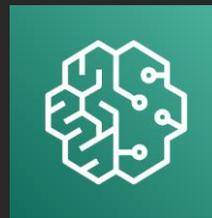
## Edge computing & data transfer



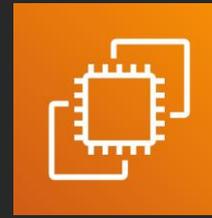
## Machine learning as a service (MLaaS)



### Latency, scaling, efficiency



Amazon SageMaker

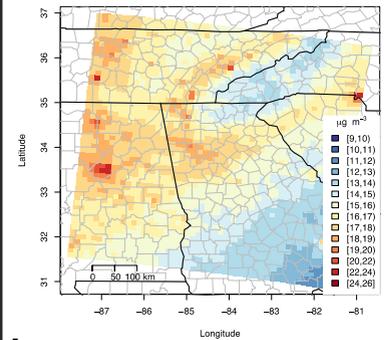


Amazon EC2



AWS Lambda

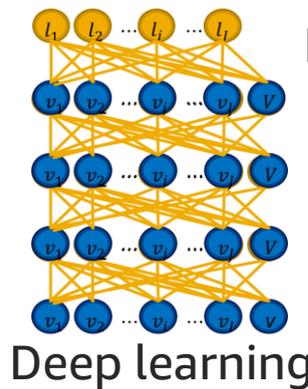
## Prediction methodology



Long-term, coarse-grained model

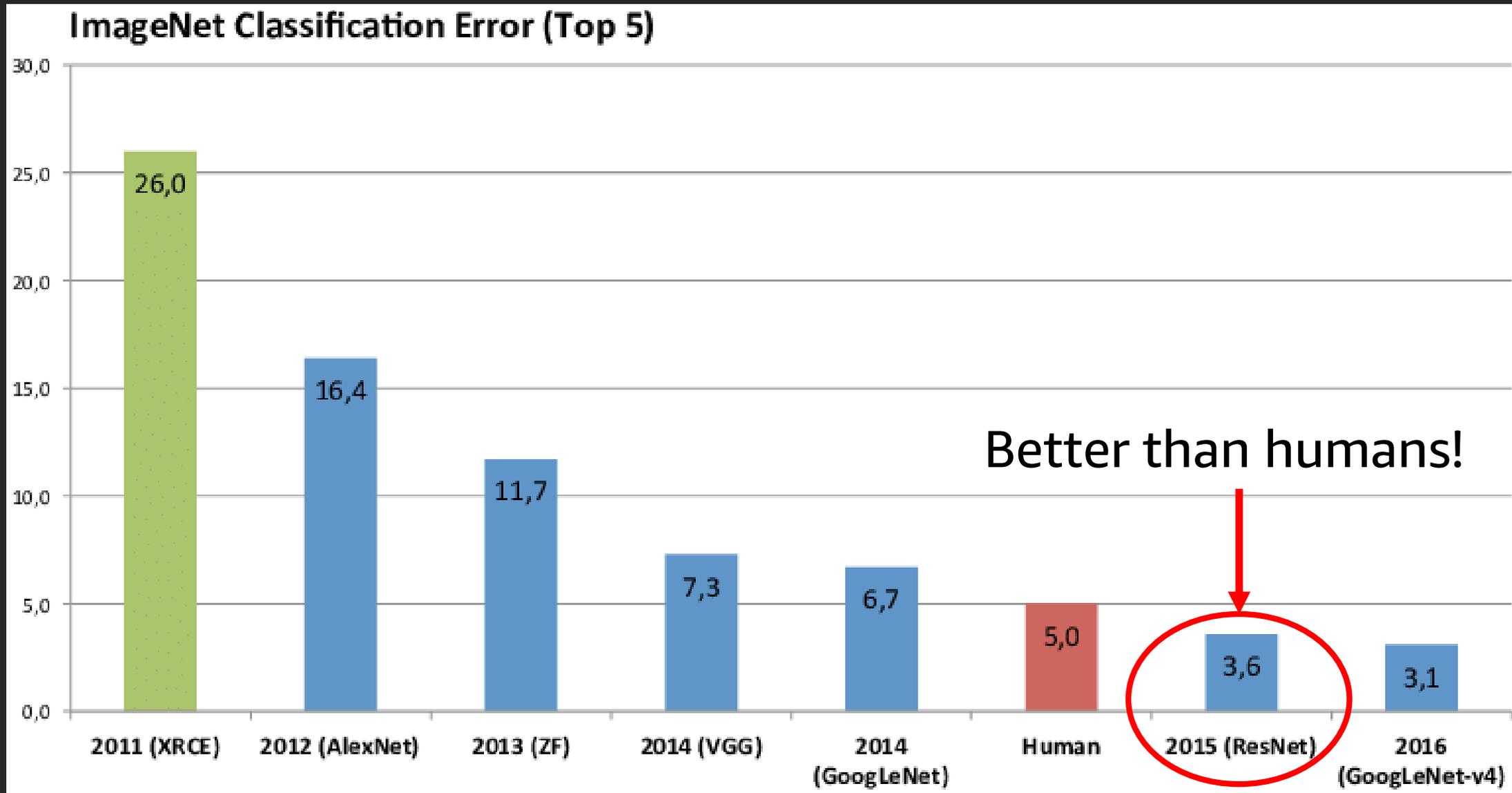


Real-time, fine-grained model



Deep learning

# Deep neural network for image classification

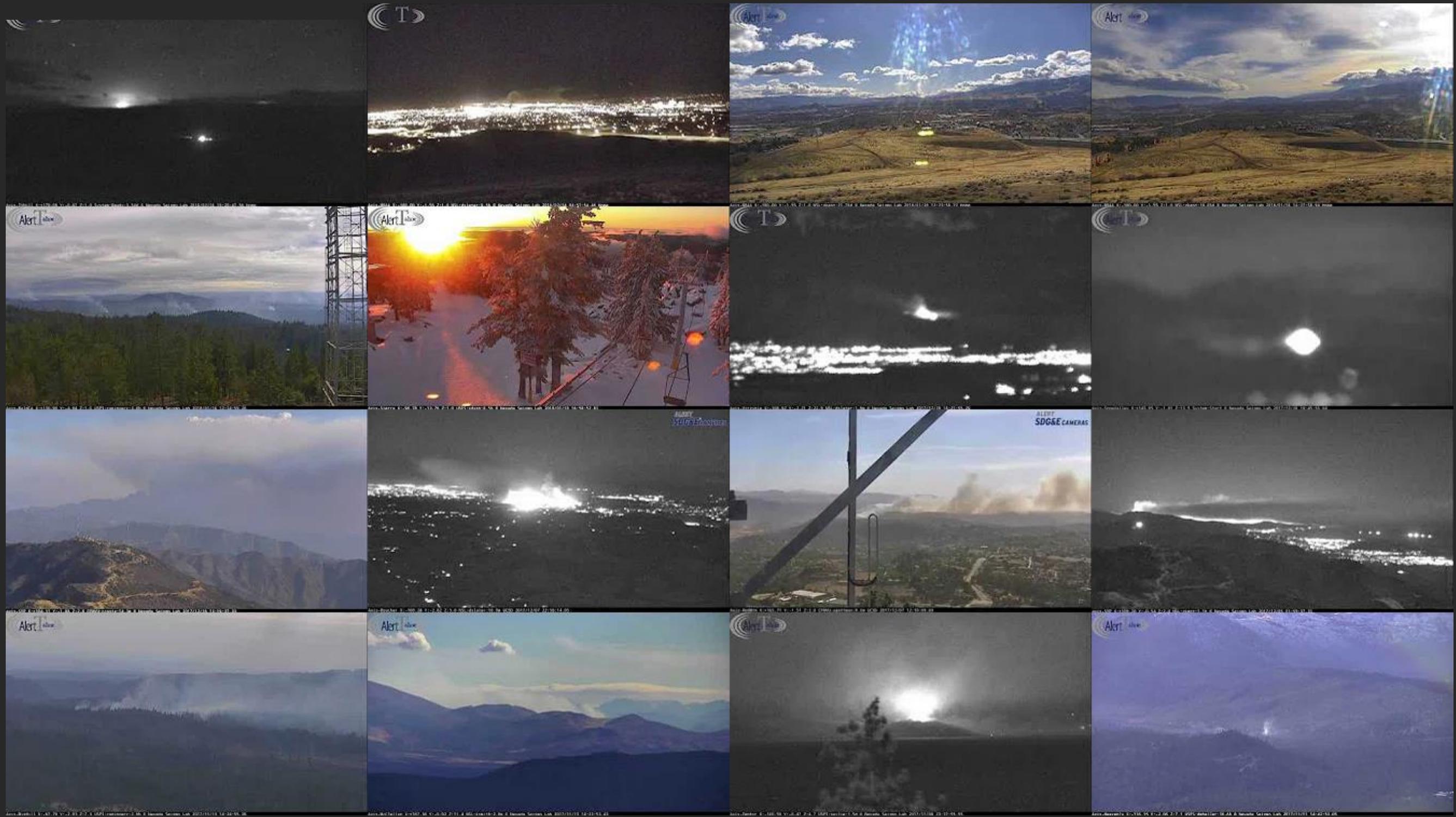


# Data used in machine learning benchmarks



Pictures from MNIST and ILSVRC2012 datasets

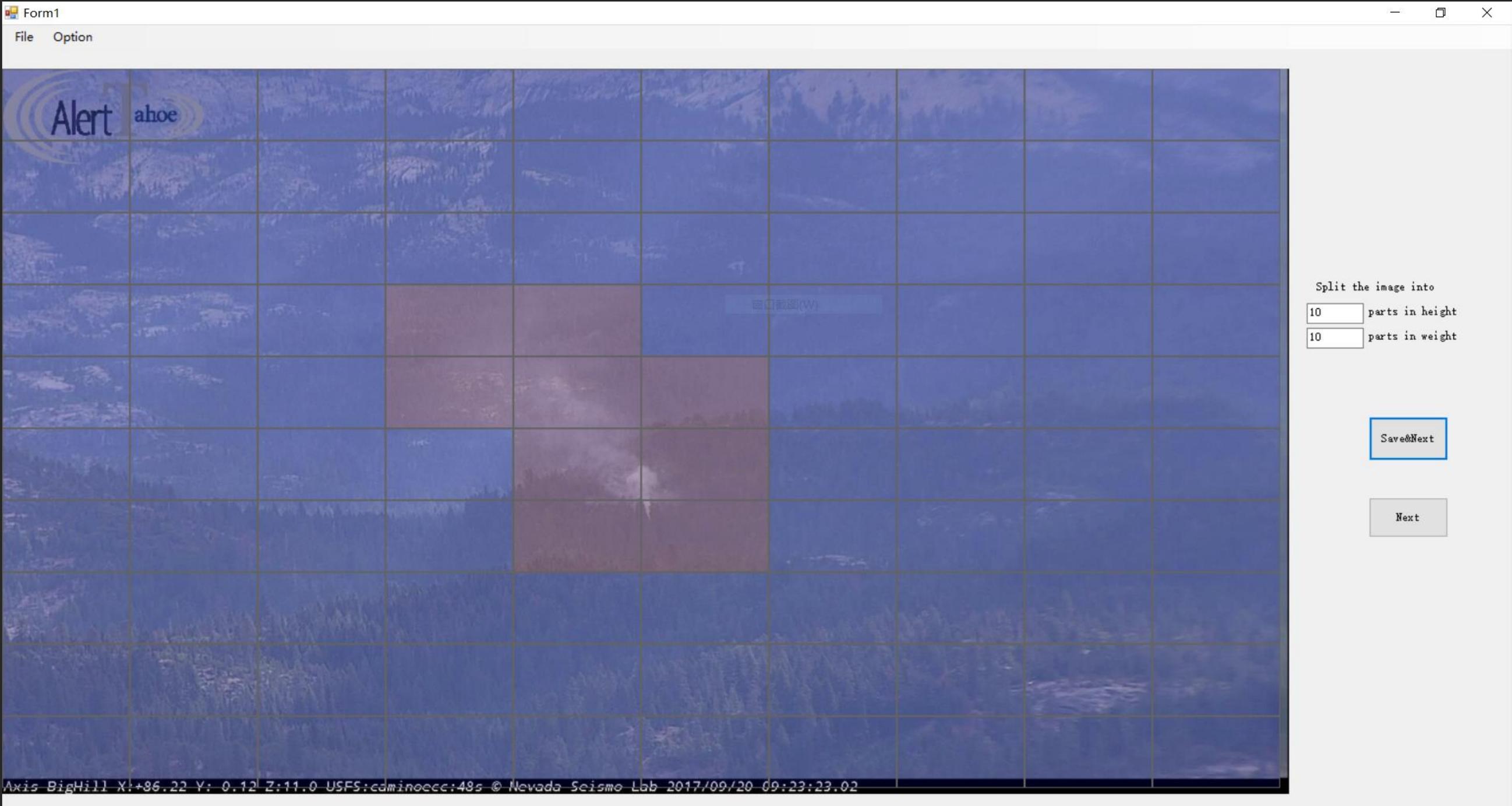
# The camera data we need to deal with

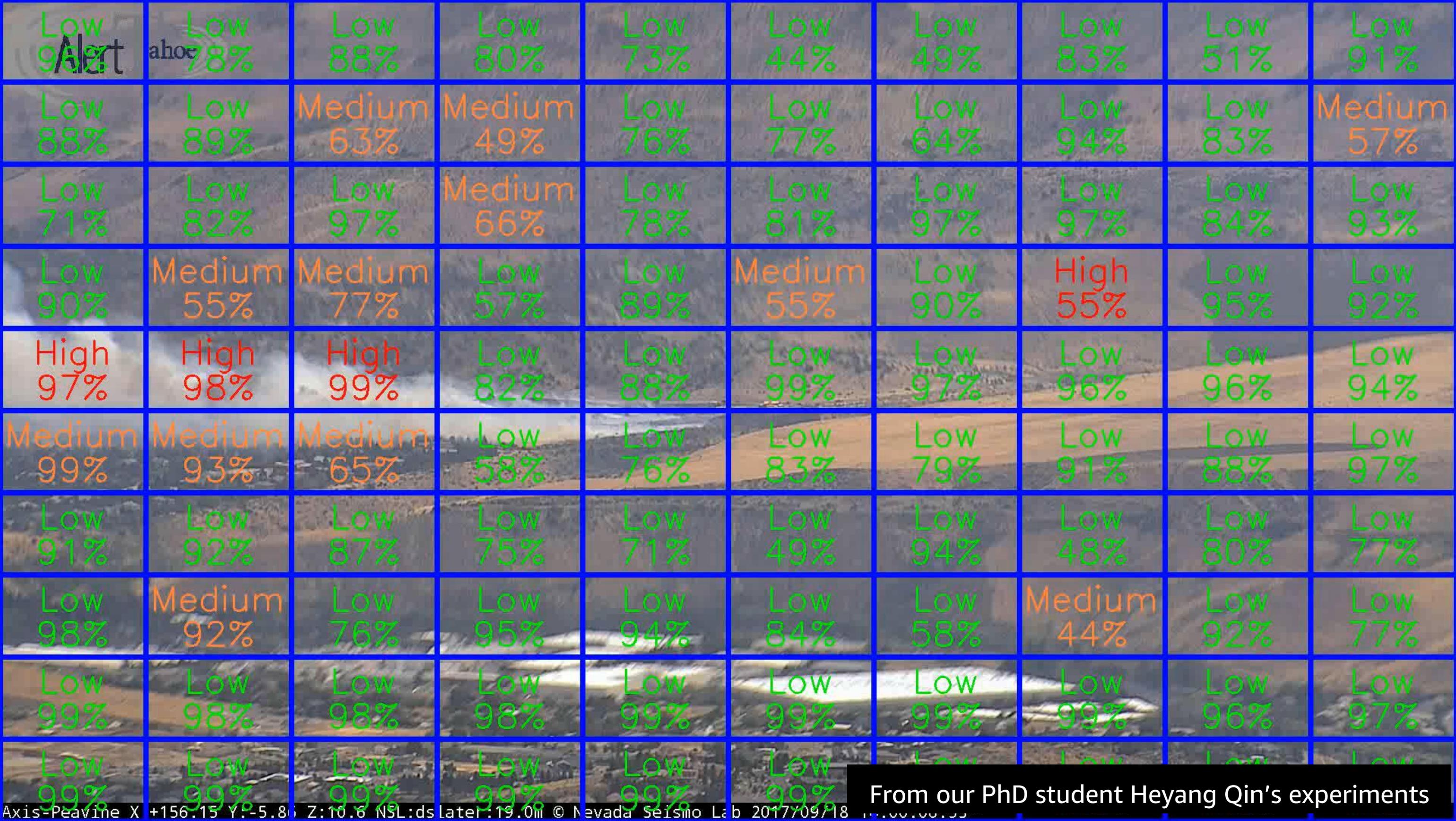


Many different shapes, changing all the time

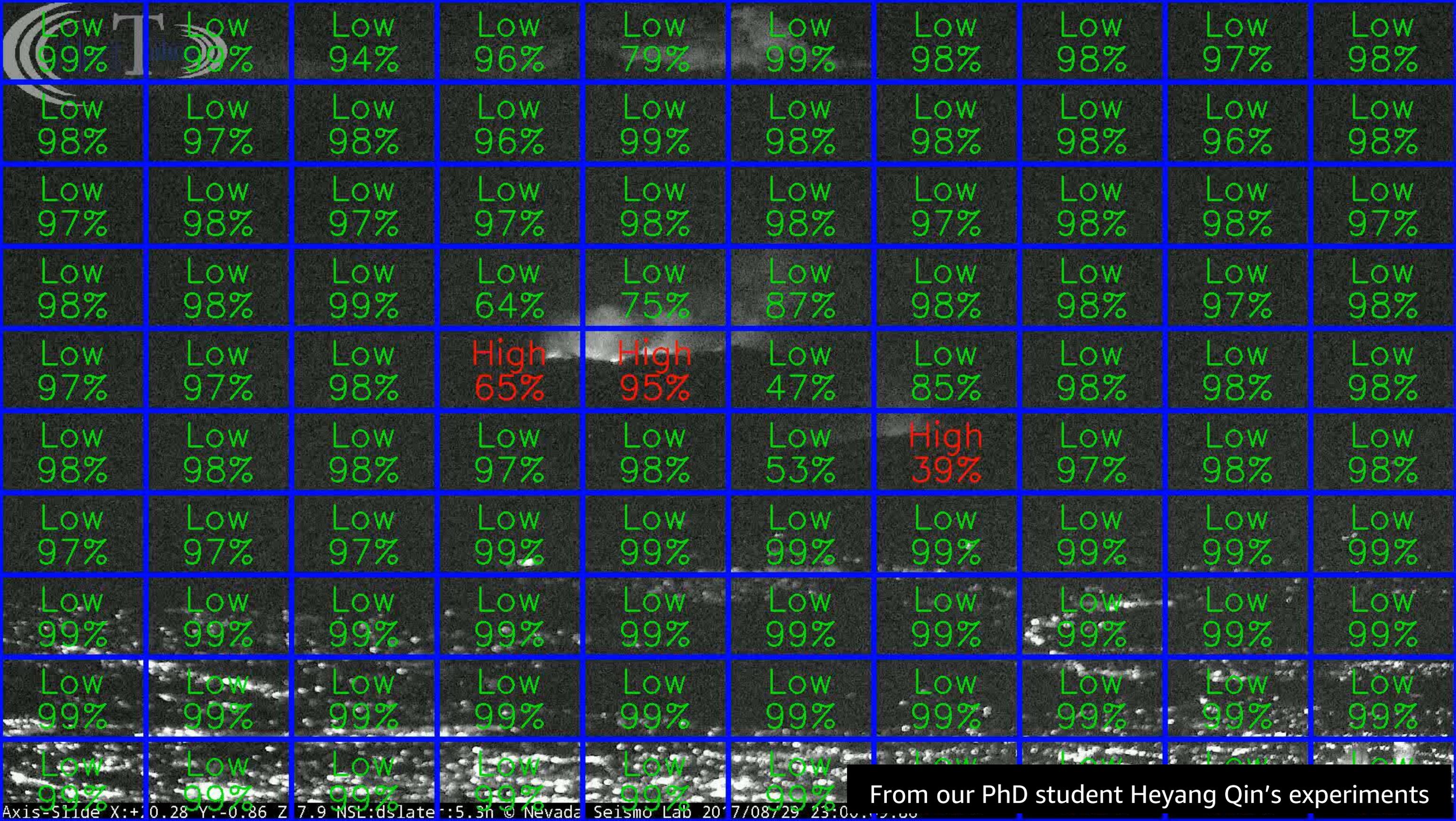


# Identify the smoke regions on images and add labels





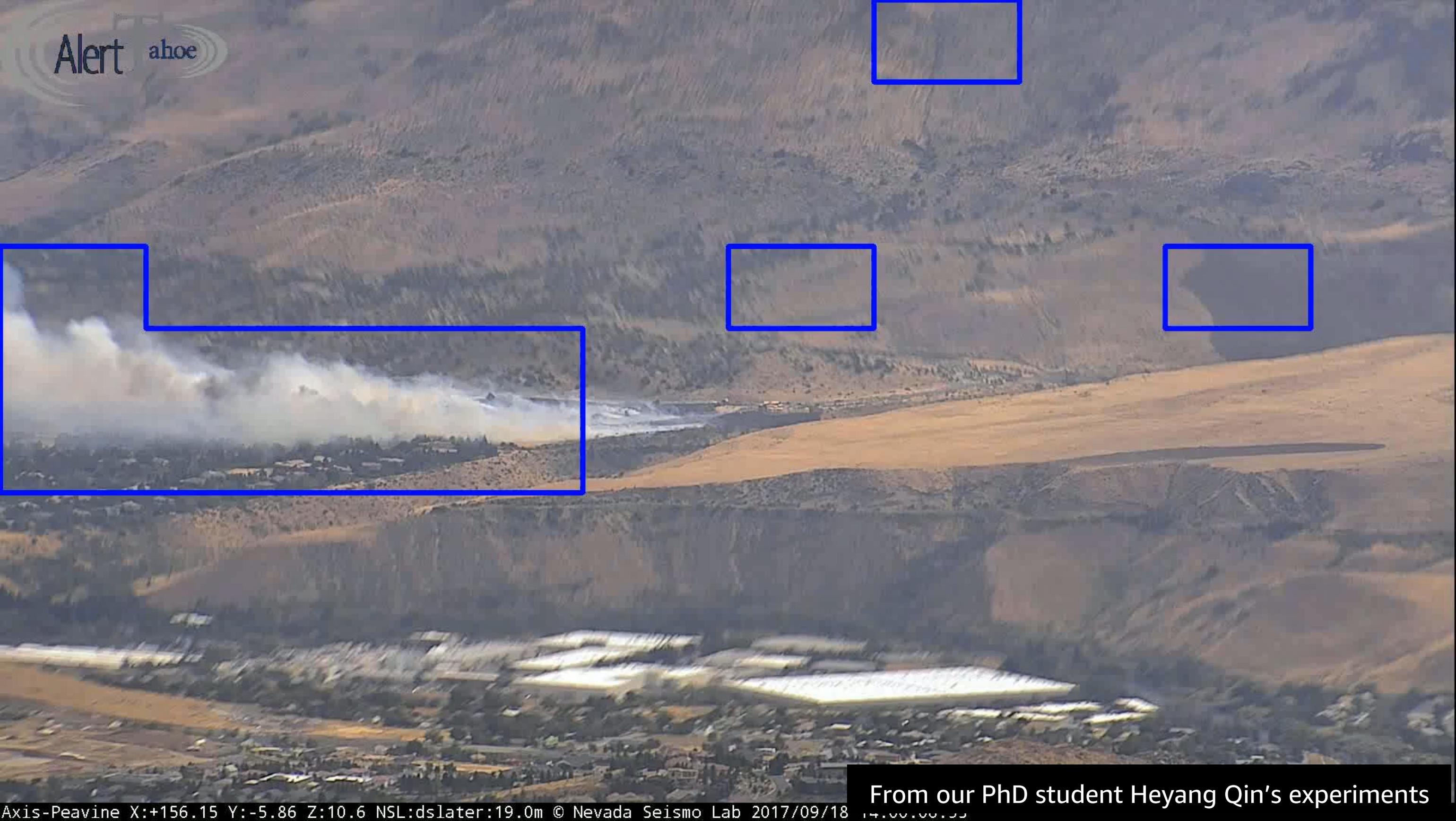
From our PhD student Heyang Qin's experiments



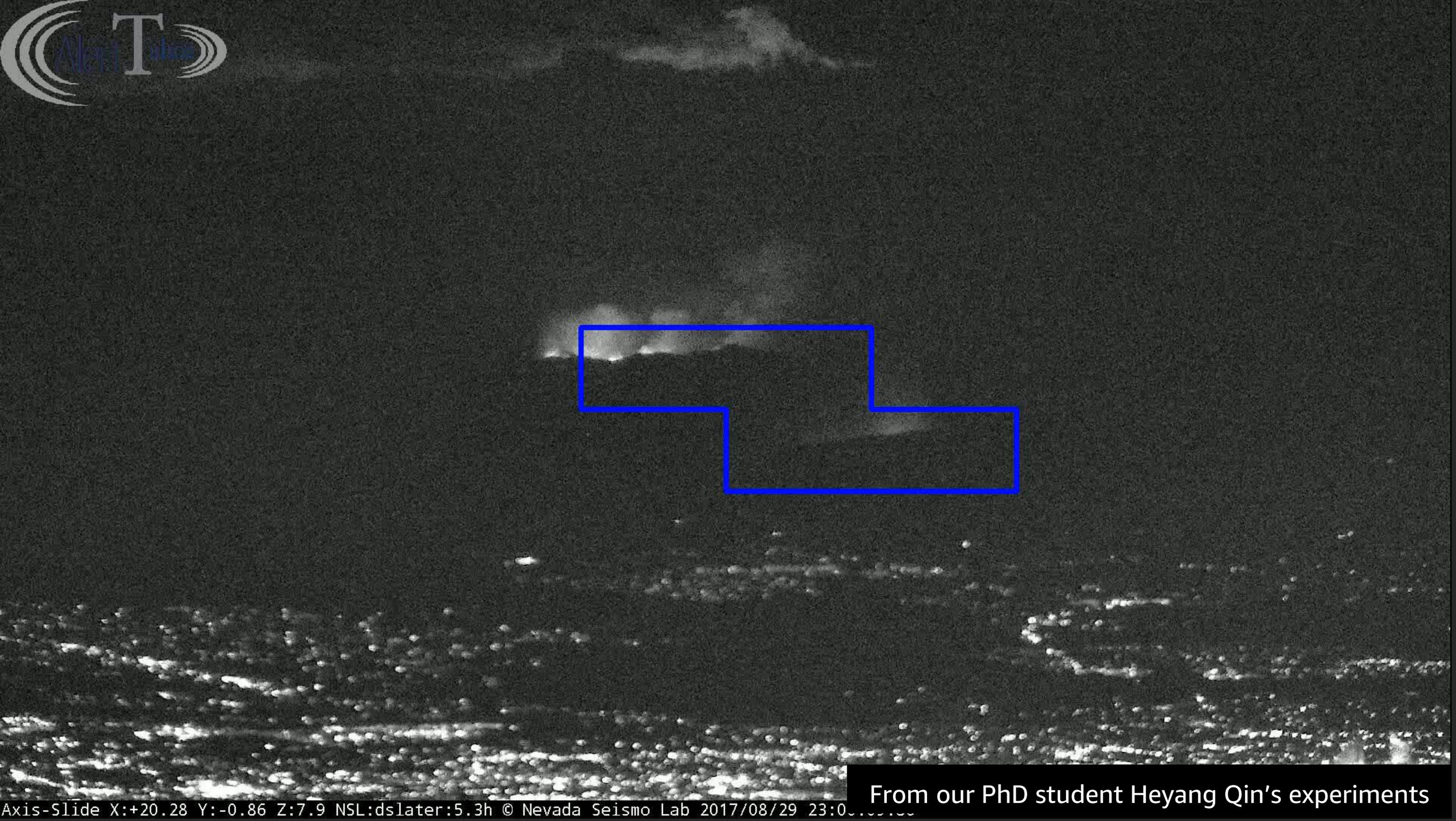
From our PhD student Heyang Qin's experiments

# Use bounding box for detection



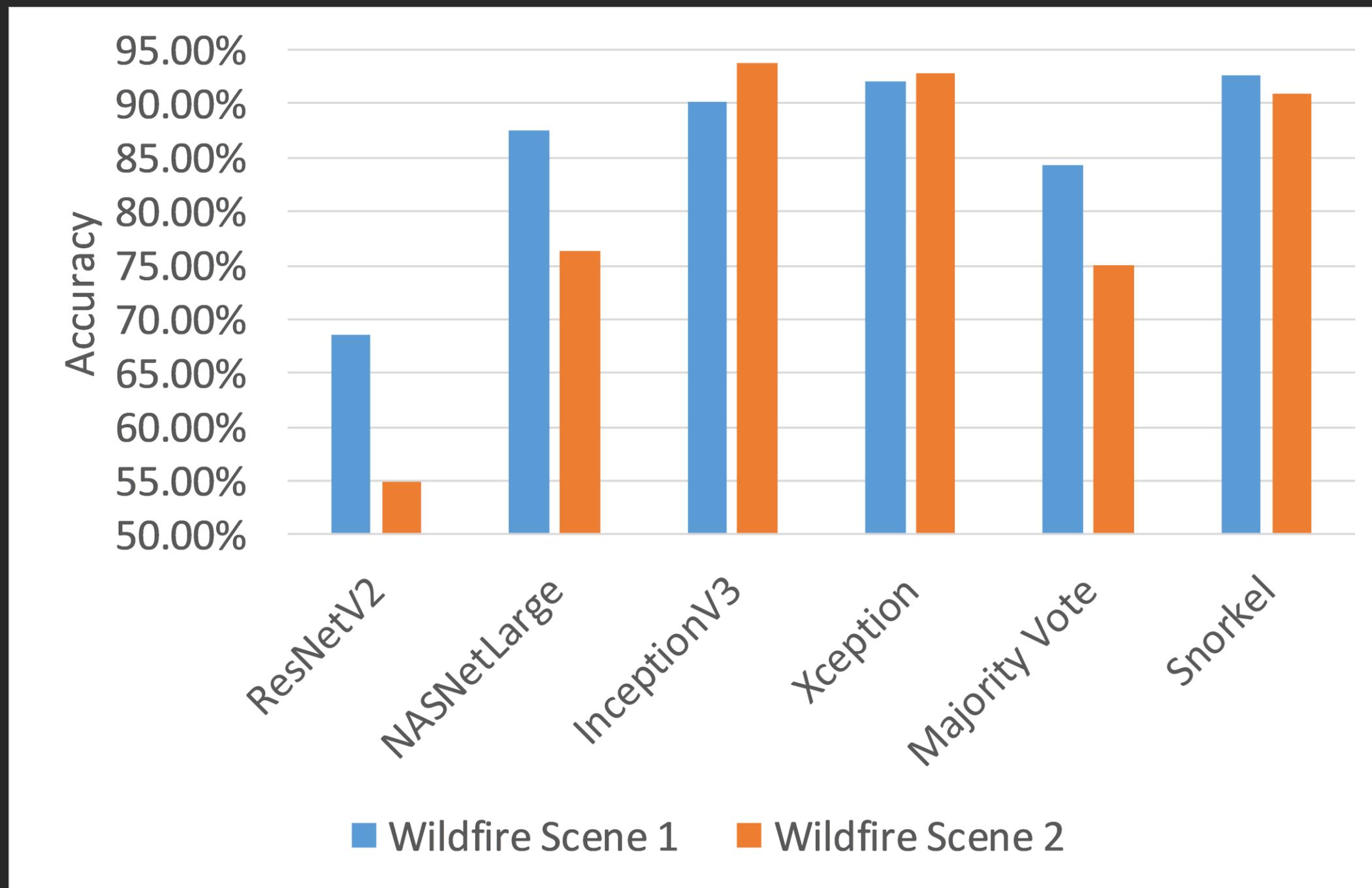


From our PhD student Heyang Qin's experiments

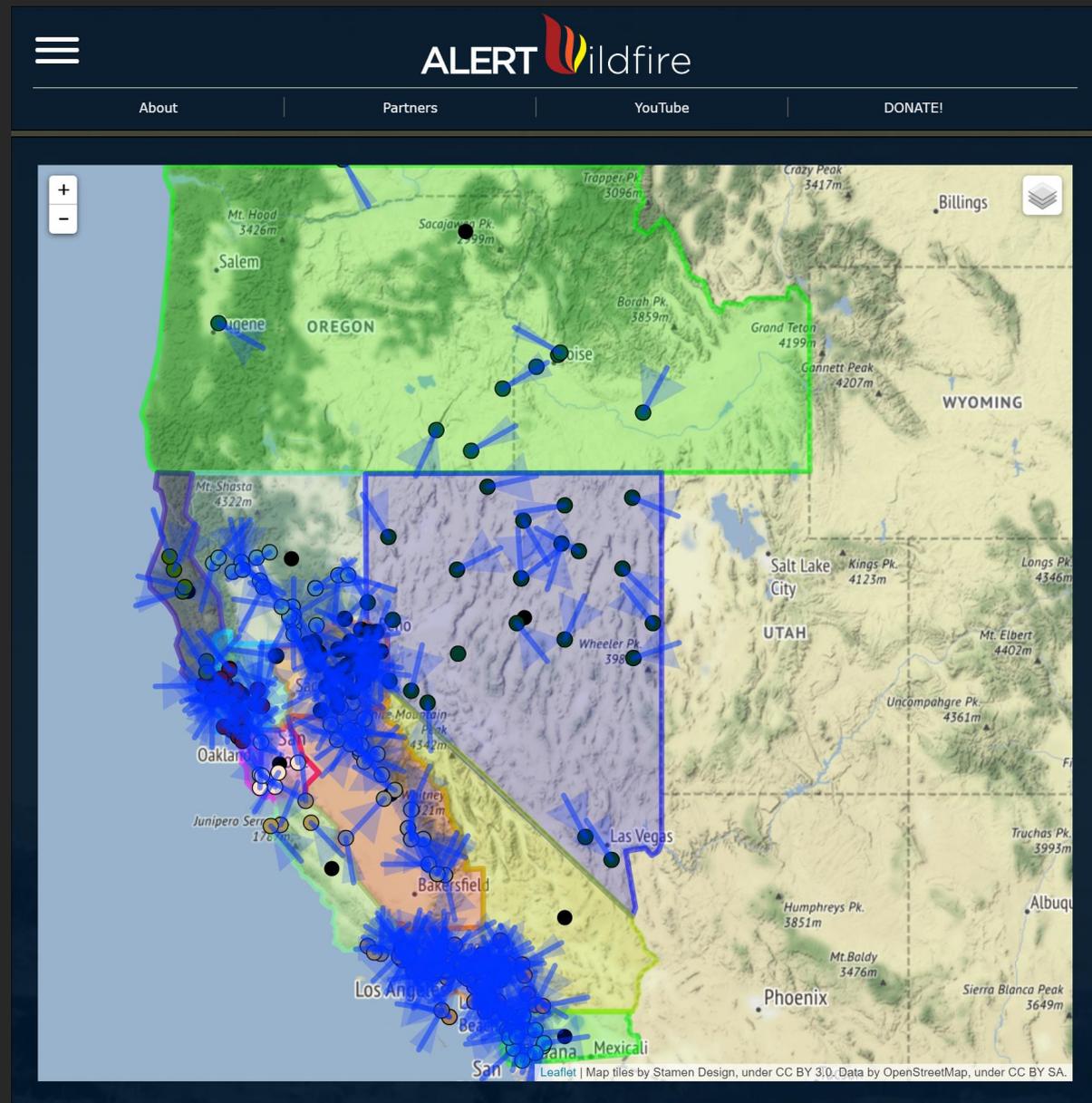


From our PhD student Heyang Qin's experiments

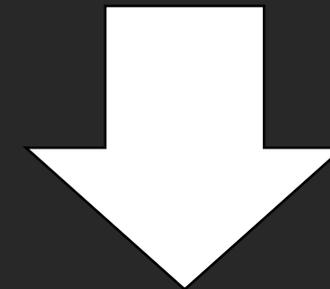
# Prediction accuracy using different models



# Real-time camera image classification

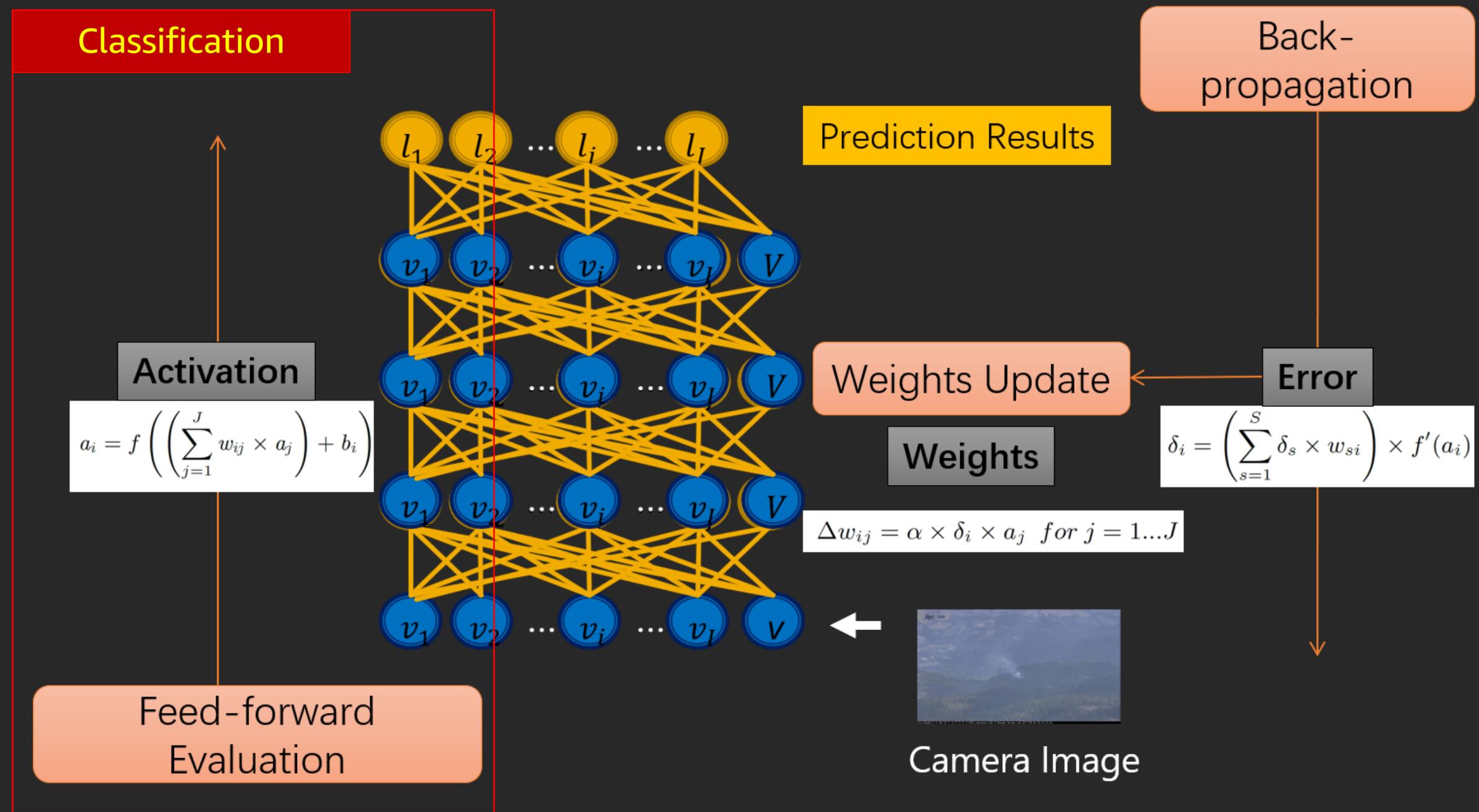


Many cameras dynamically join and leave during wildfire



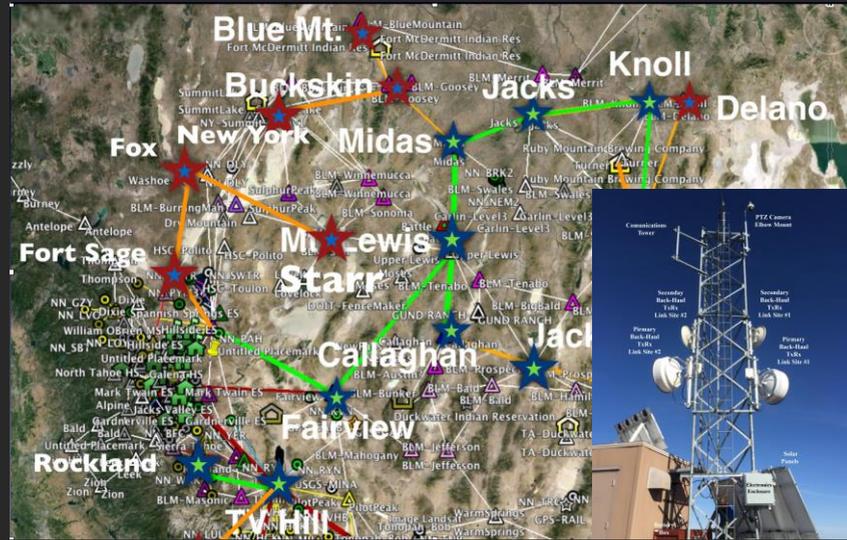
How to provide scalable and real-time image/video classification?

# Classification vs. training



# Use machine learning as a service (MLaaS)

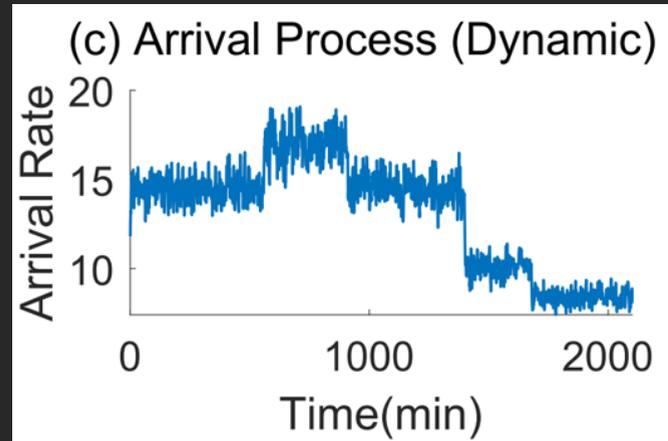
## Camera network



Smoke/wildfire detection

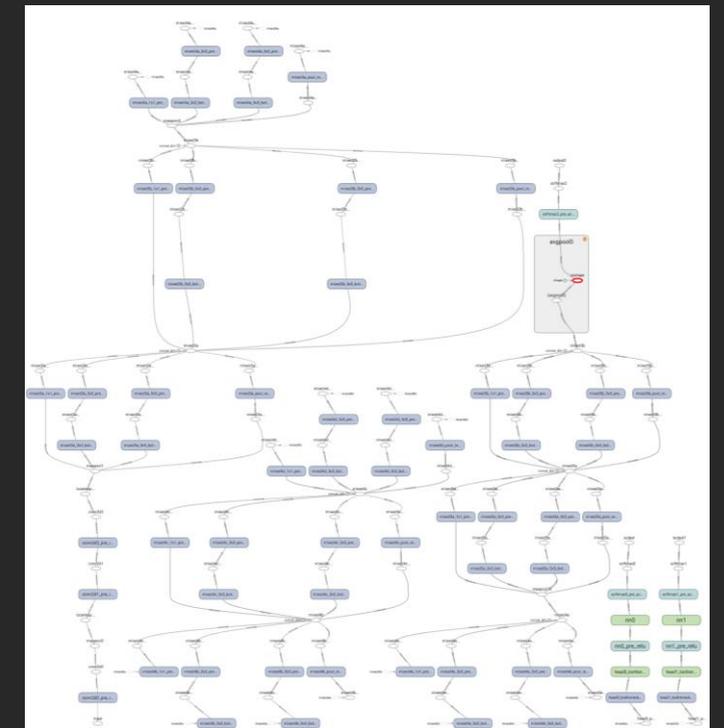
API

Machine learning model is a callable API



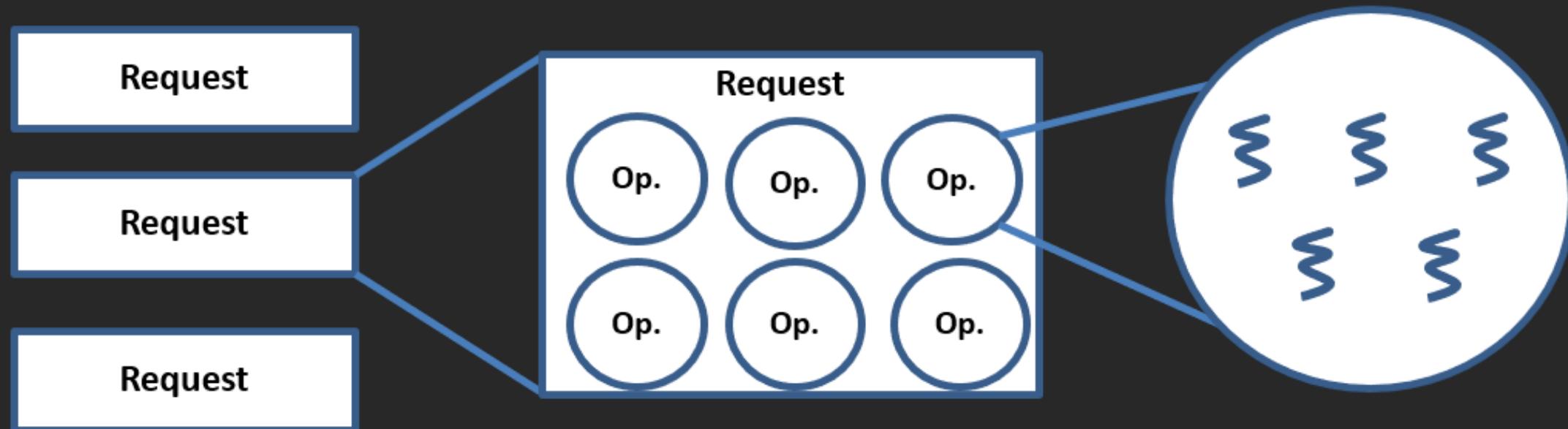
Two major challenges

1. Sequential execution is too slow
2. Dynamic classification workload



# Use parallelism to speed up classification

**Request parallelism**    **Inter-op parallelism**    **Intra-op parallelism**



 **Request**     **Operation**     **Thread**

# Use batching to speed up classification

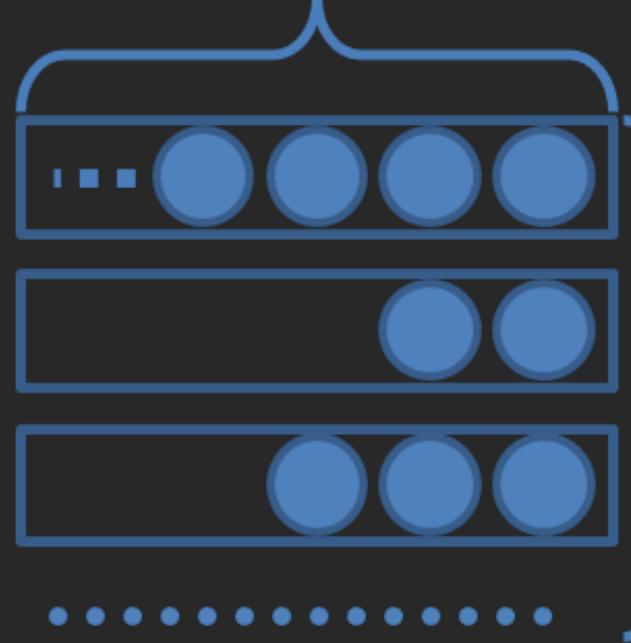
Camera network



Requests

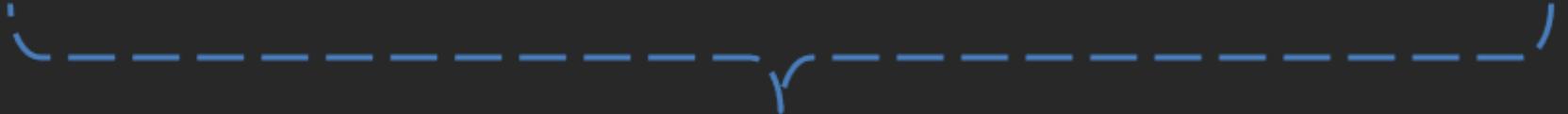


Batch Size



Parallel Batch Threads

Batch timeout

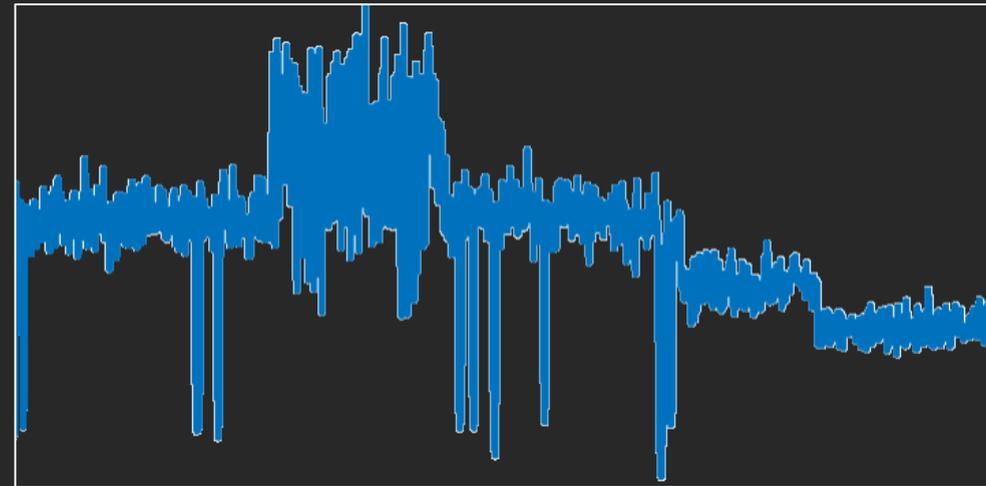


# Problem

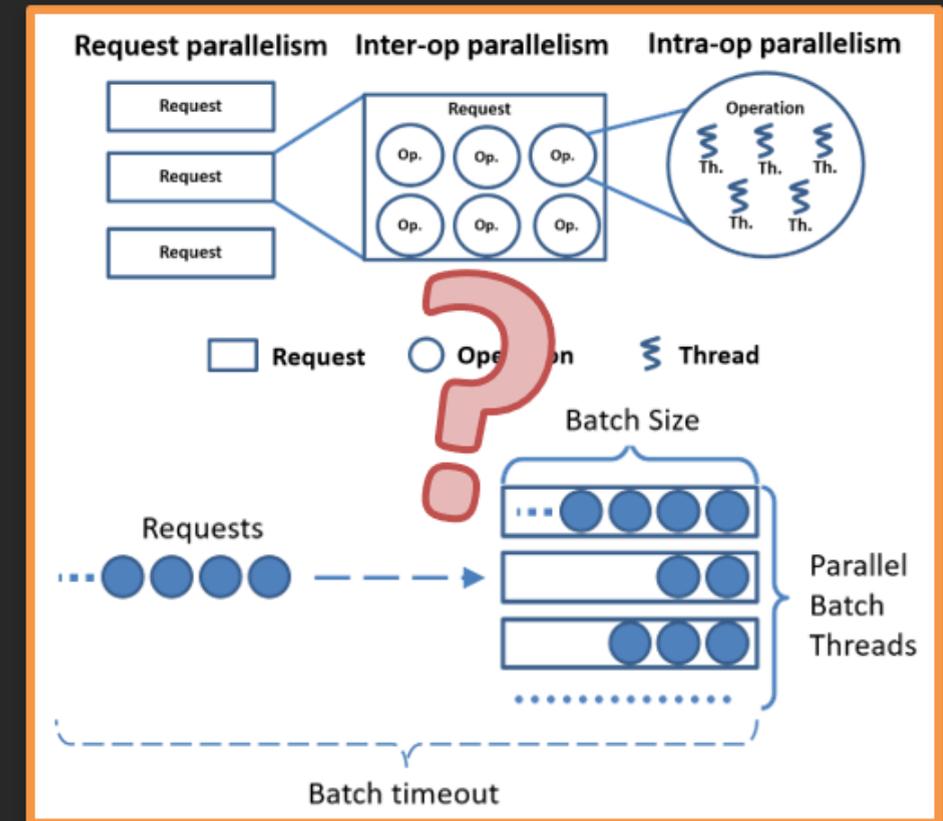
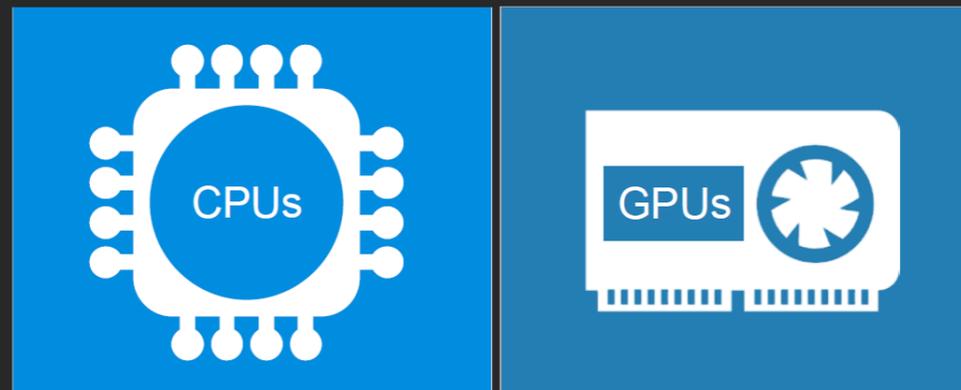
Classification requests form  
Dynamic Arrival Processes

How to tune the system  
configurations?

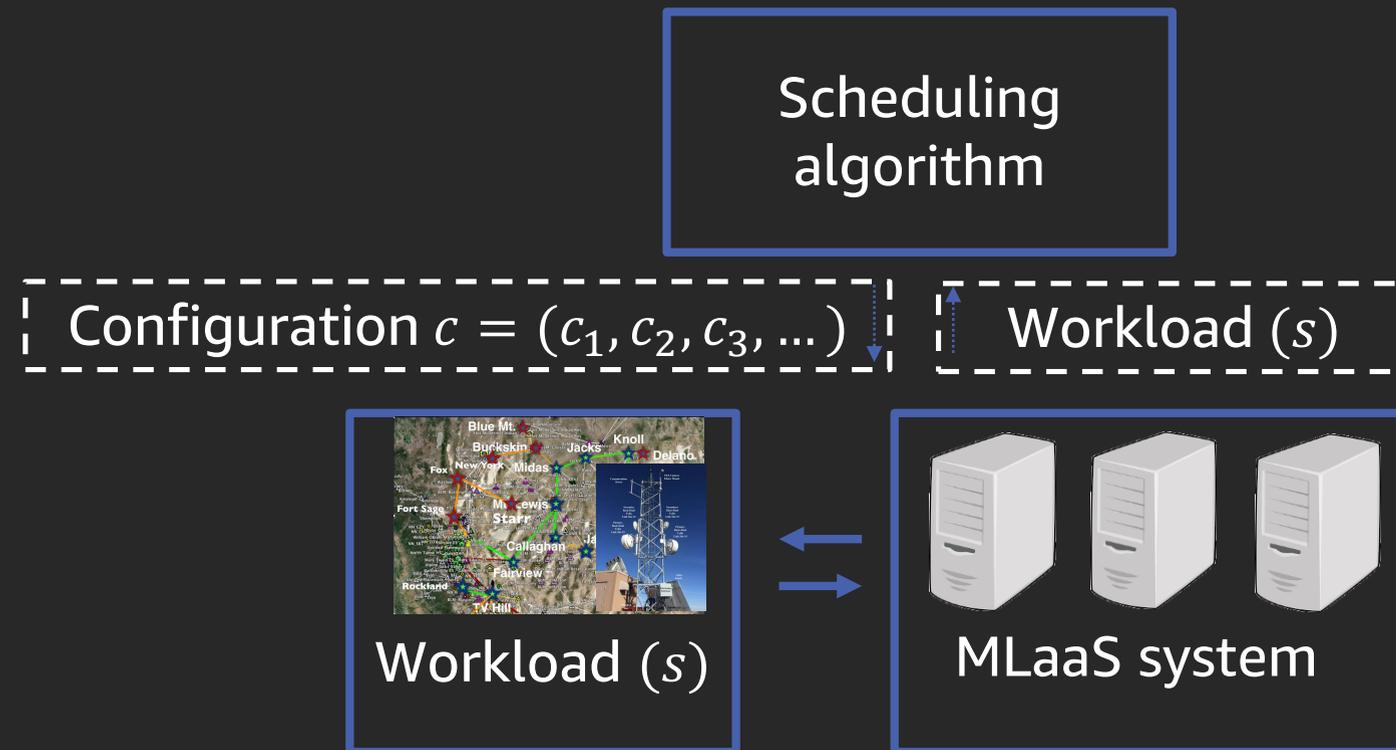
Camera network



Different infrastructures



# Problem definition



Minimize  $r(s, c)$

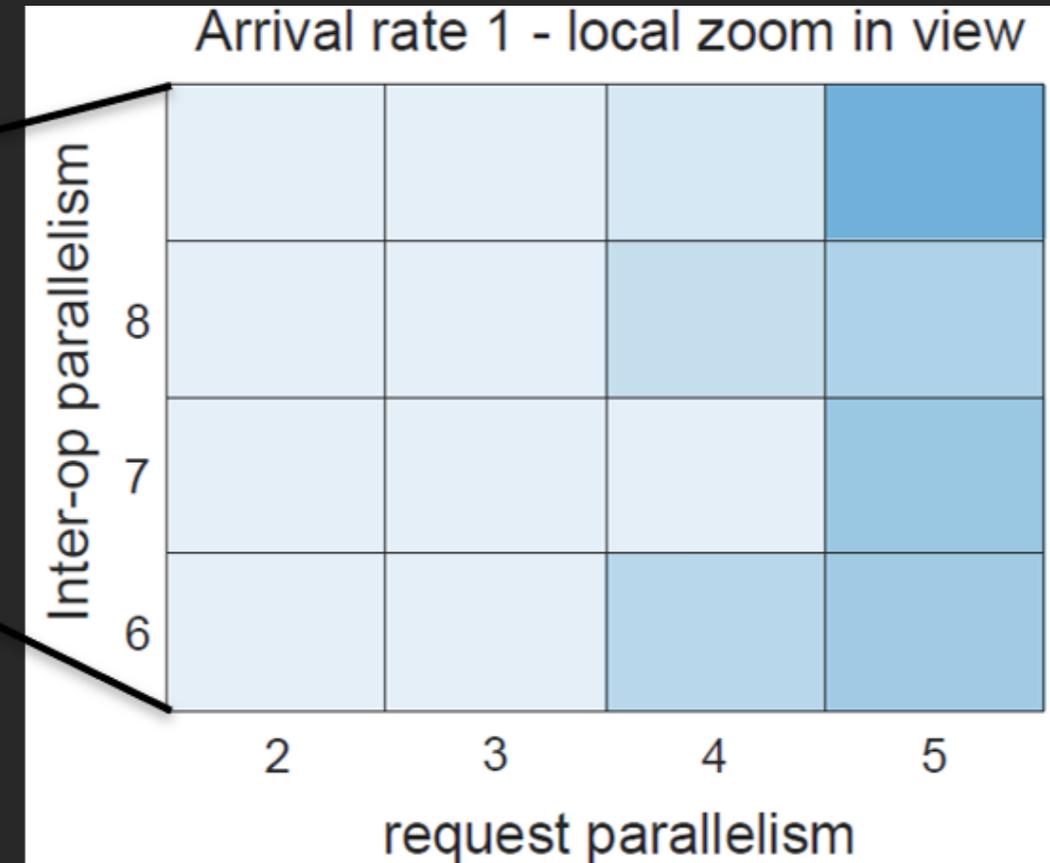
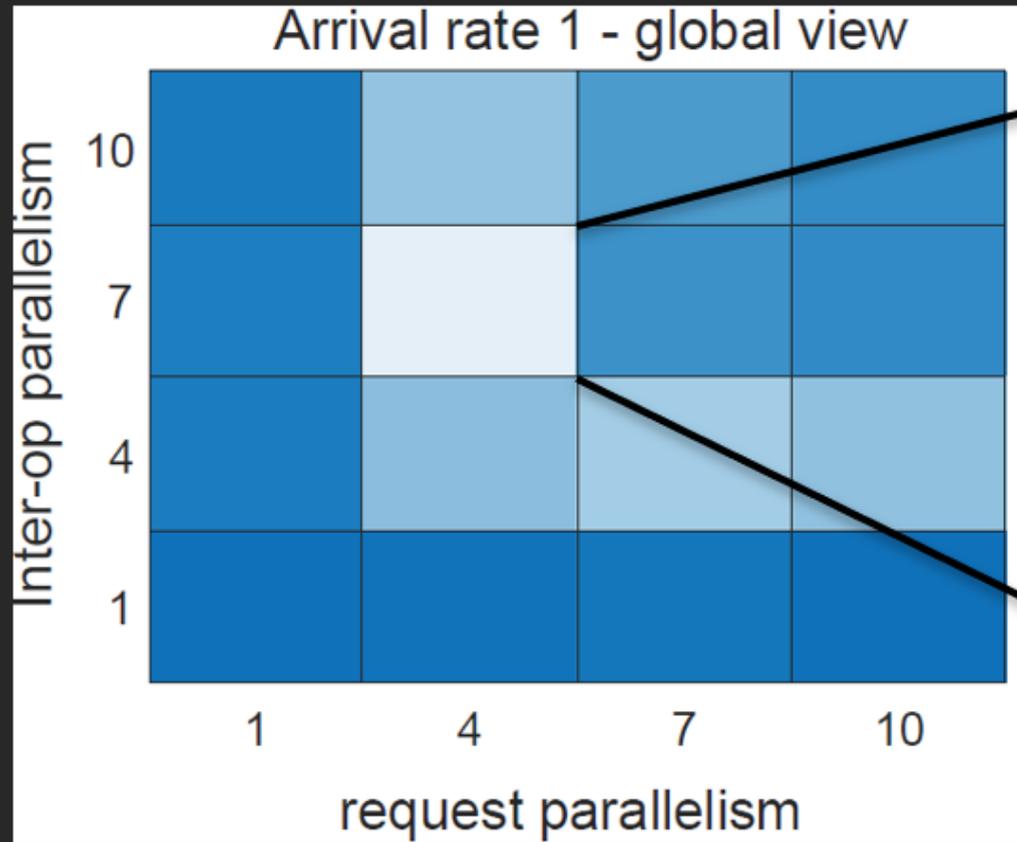
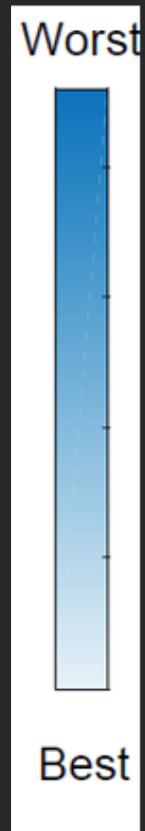
Subject to  $s \in \mathcal{S}, c \in \mathcal{C}$

$r(s, c)$  - Average request latency

$\mathcal{S}$  - Possible workloads

$\mathcal{C}$  - Possible system configs

# Key observations

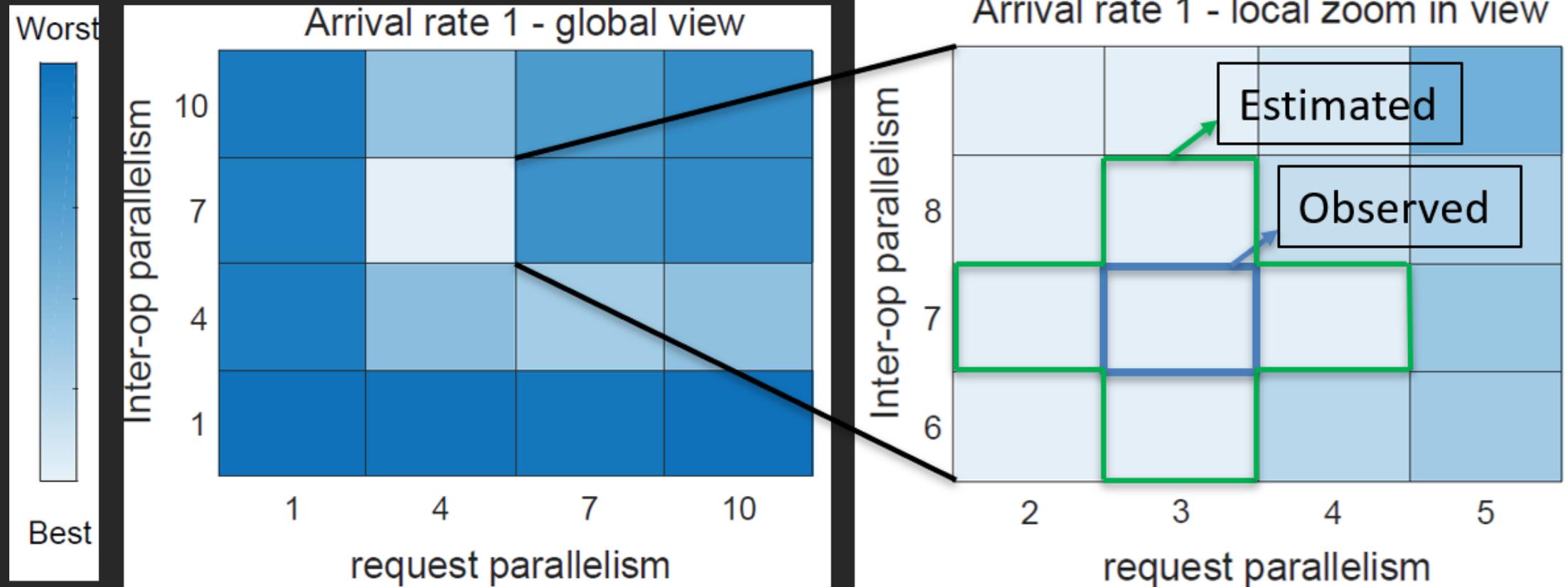


versatile globally

vs

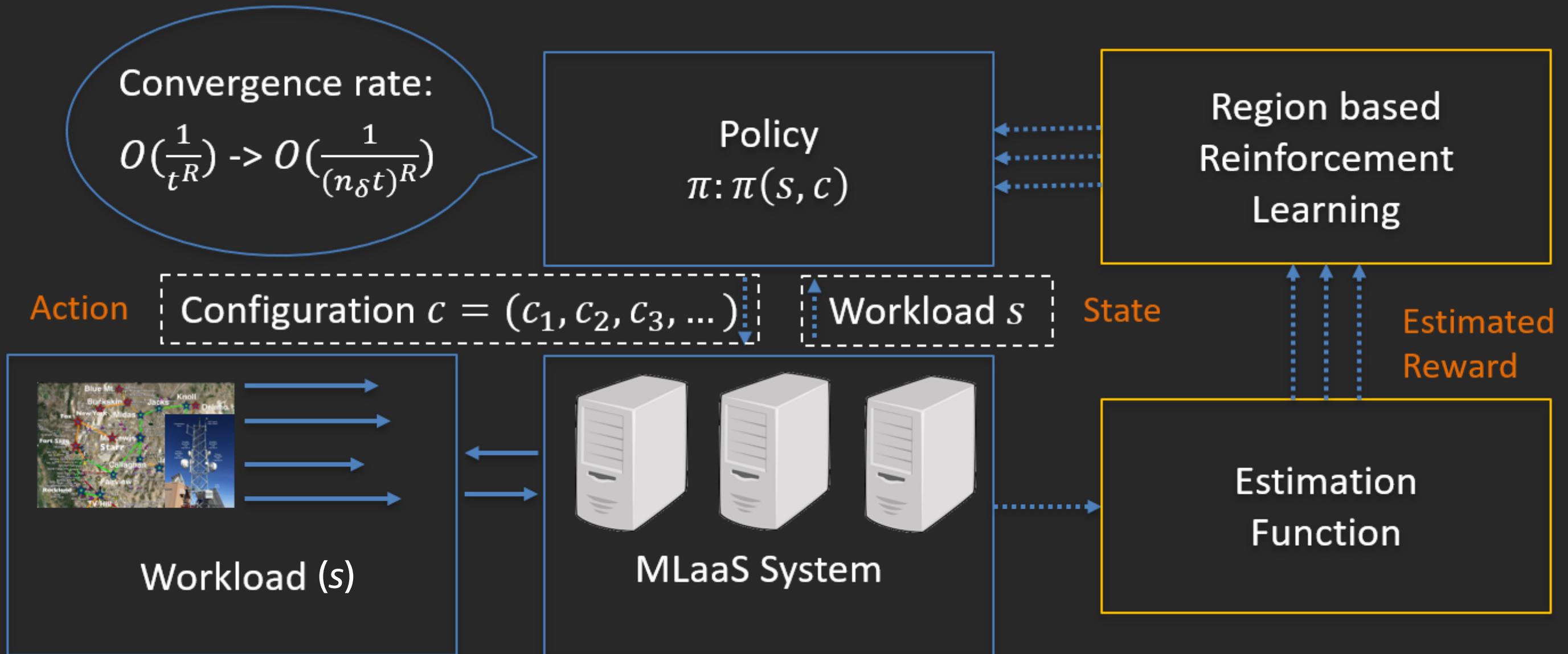
smooth locally

# Key idea

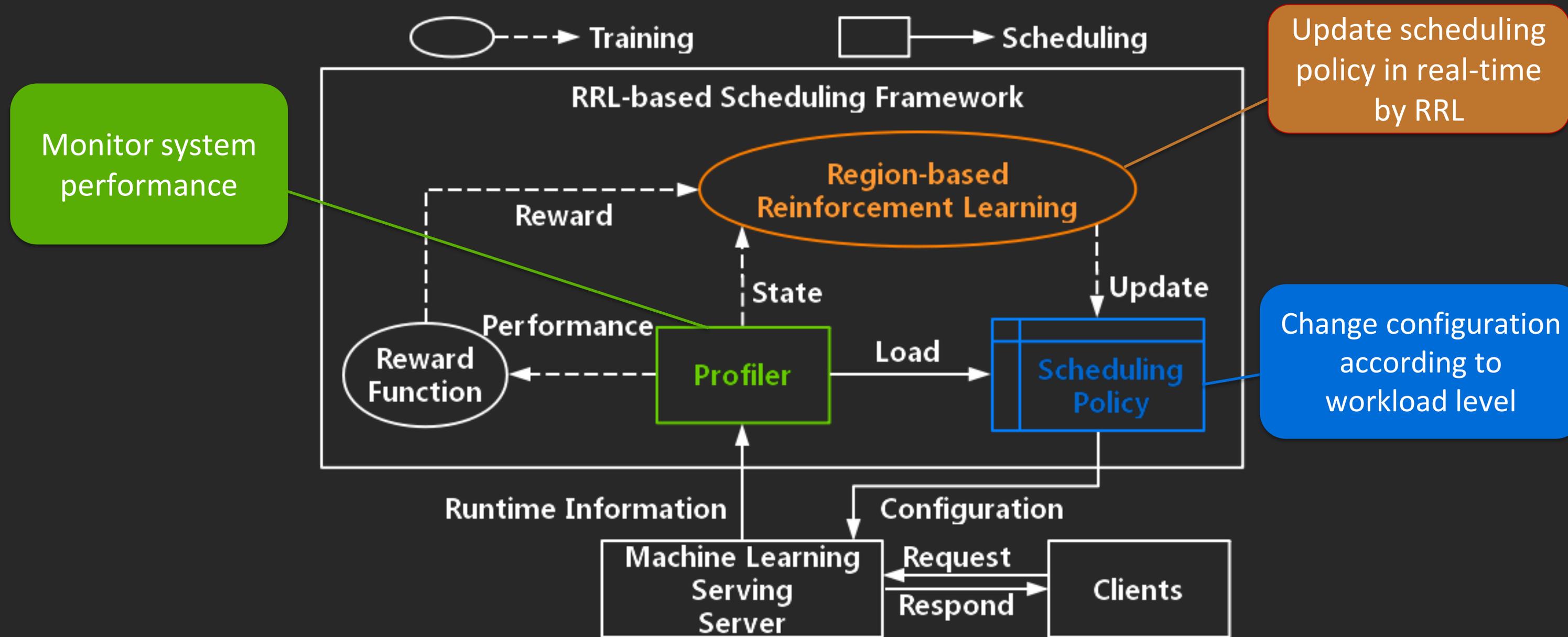


Use a single observation to update the configurations in a nearby region

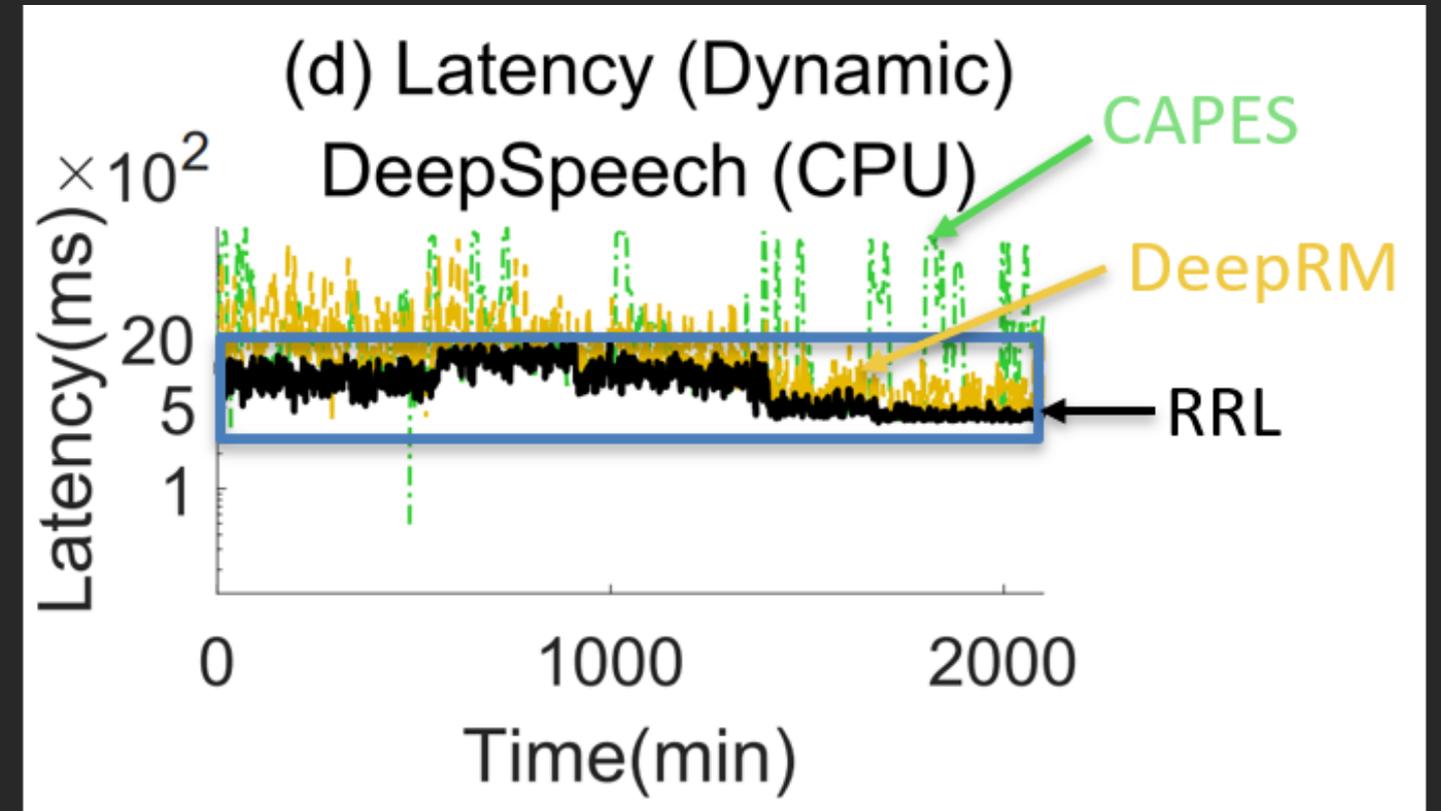
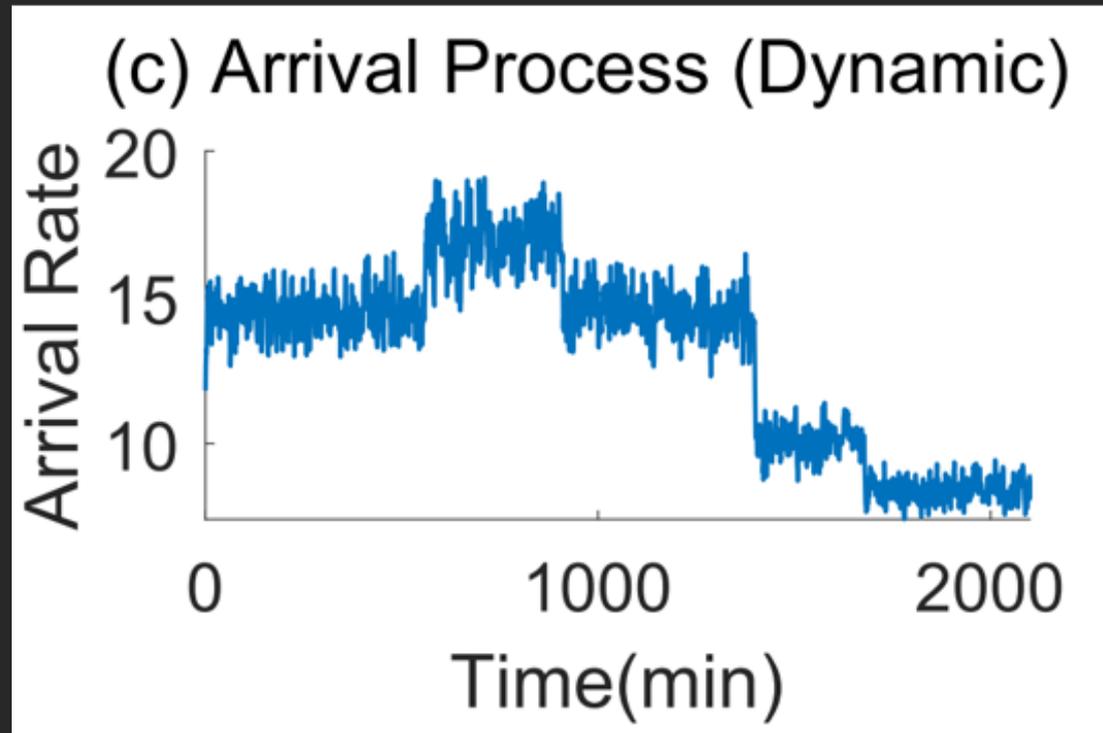
# Region based Reinforcement Learning (RRL)



# System design



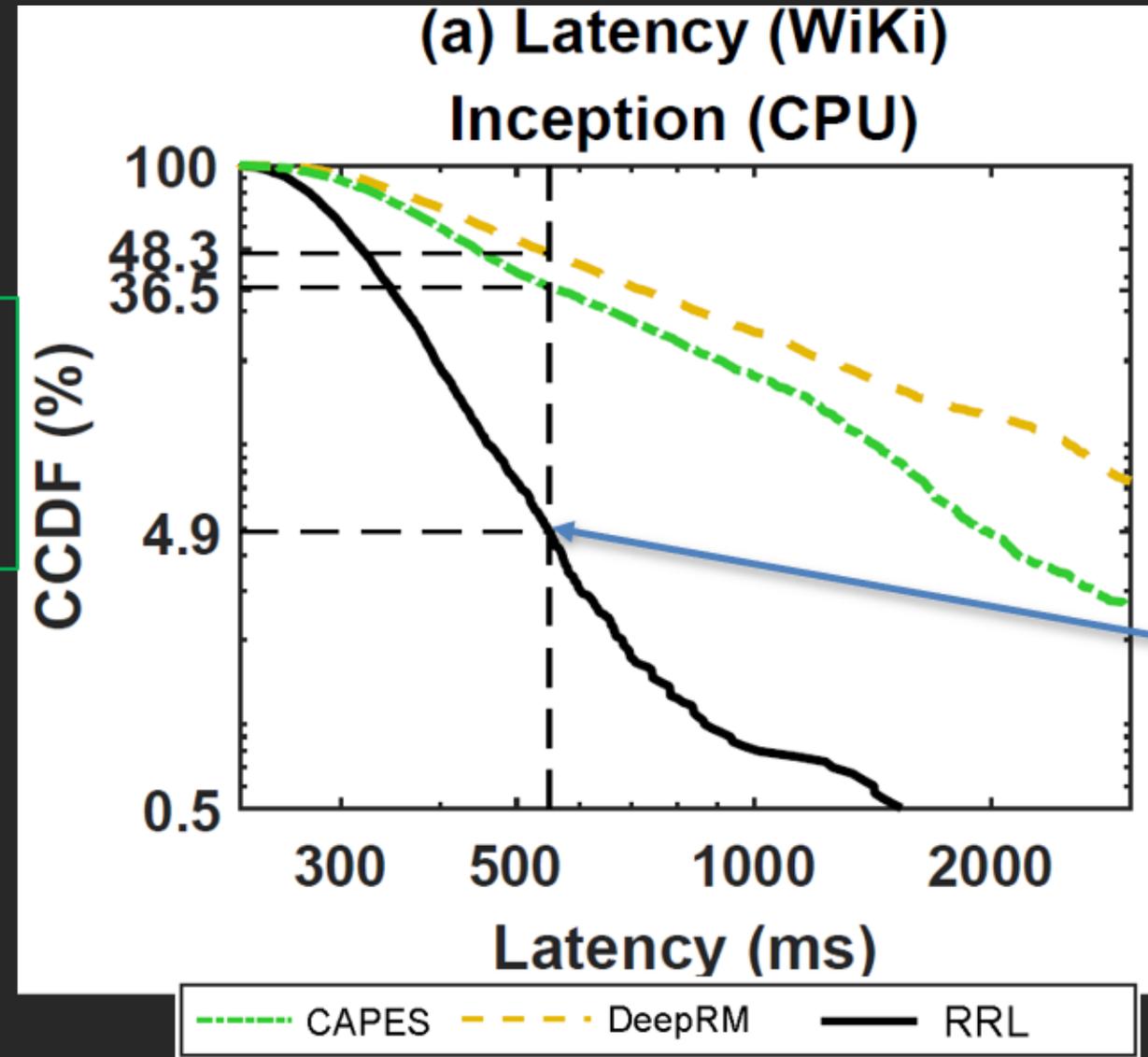
# Results



RRL has 63% less average latency than CAPES and 49% than DeepRM.

# Results

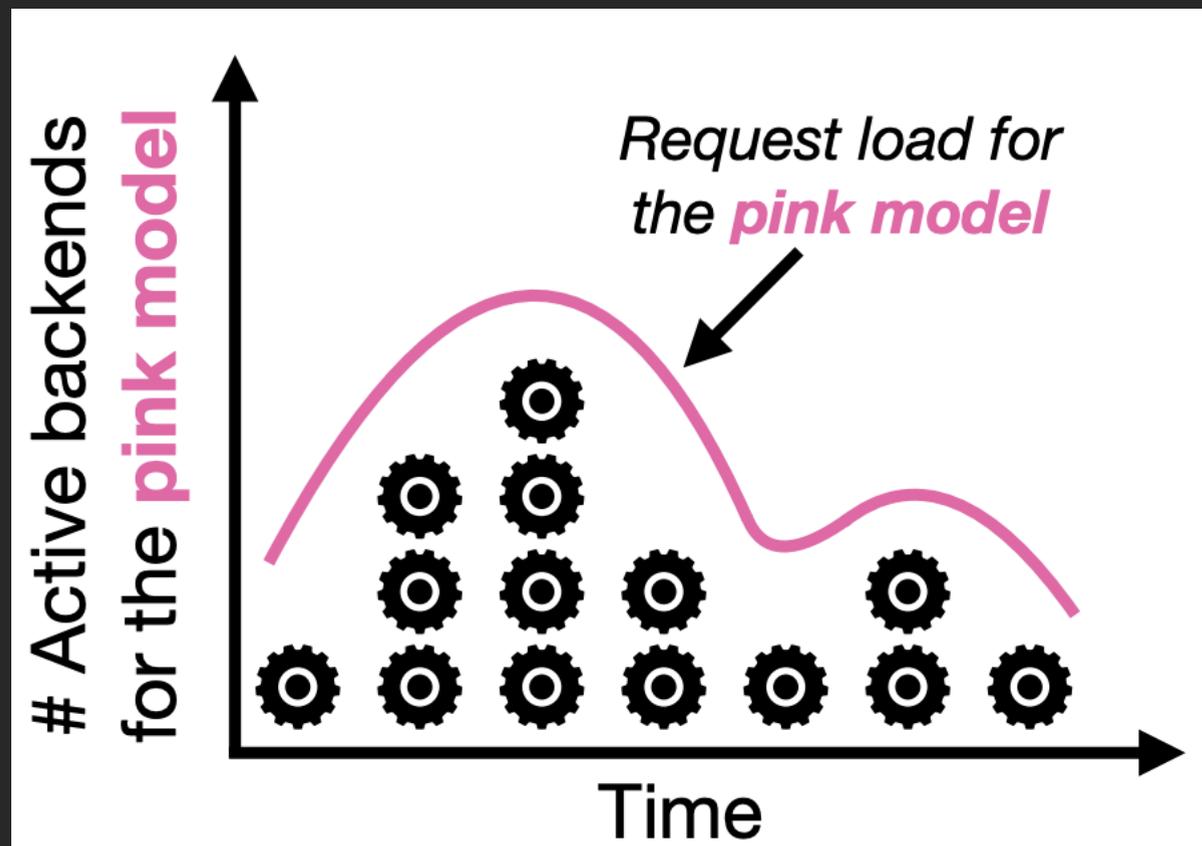
## Complementary cumulative distribution function



Please note that the Y-axis is log scale

The 95th percentile of RRL is 550 ms

# Conventional autoscaling: Amazon SageMaker



Reactive scaling: based on current load

Provisioning  
time  
(minutes)

>>

Execution  
time  
(< 1s)

Hide provisioning time → overprovisioning  
How much overprovisioning?

[1] [https://people.mpi-sws.org/~arpanbg/pdfs/middleware2017\\_slides.pdf](https://people.mpi-sws.org/~arpanbg/pdfs/middleware2017_slides.pdf)

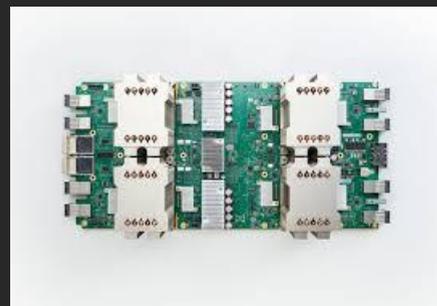
# ML accelerators: GPU, TPU, FPGA

- Mass parallel support
- Essential for training complex models
- **Expensive**

**CPU:** m5.xlarge: \$0.192 per hour

**GPU:** p2.xlarge: \$0.9 per hour

**TPU v2:** \$4.5 per hour



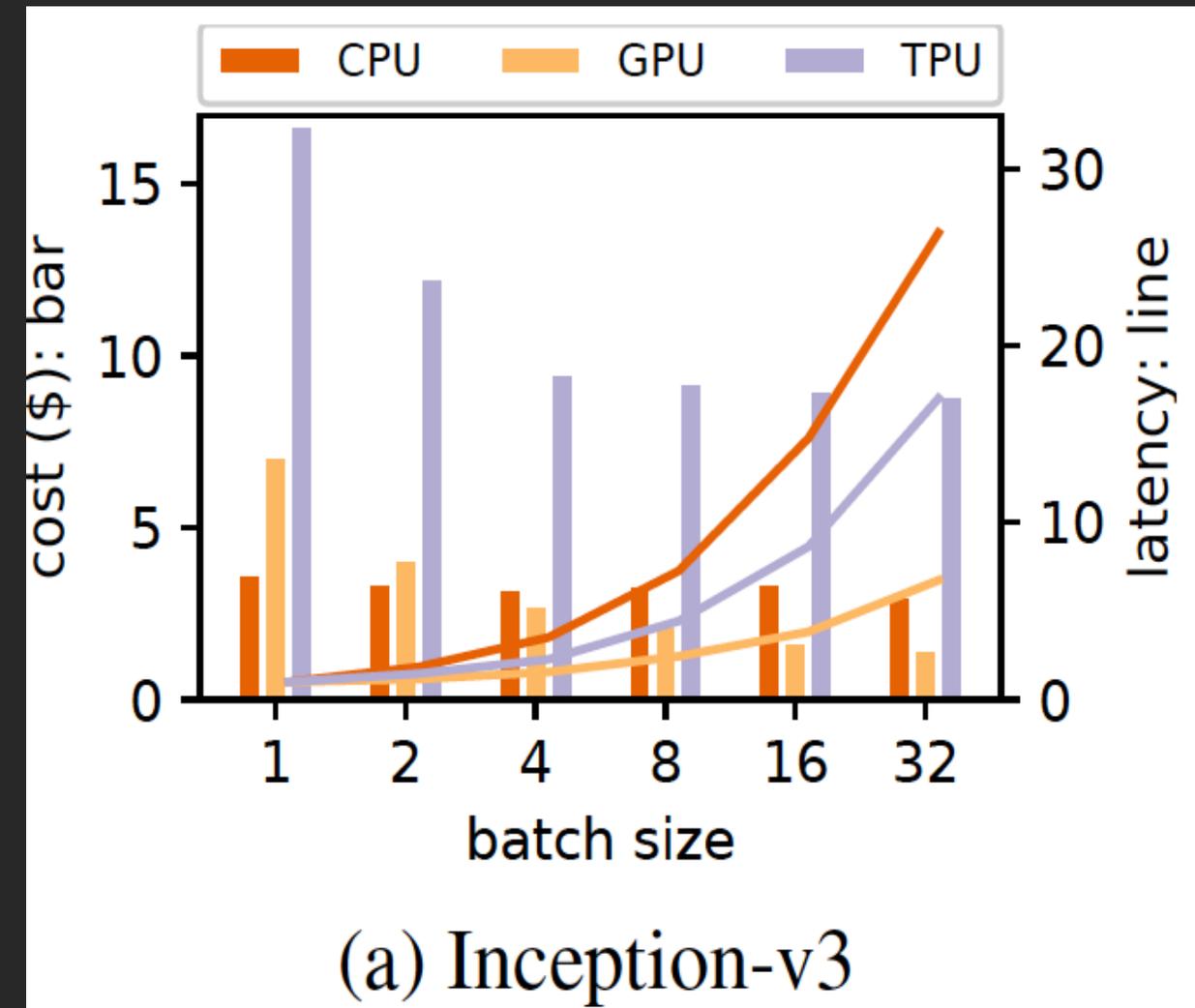
## Inference

- Run comfortably without them
- Way less parallelism

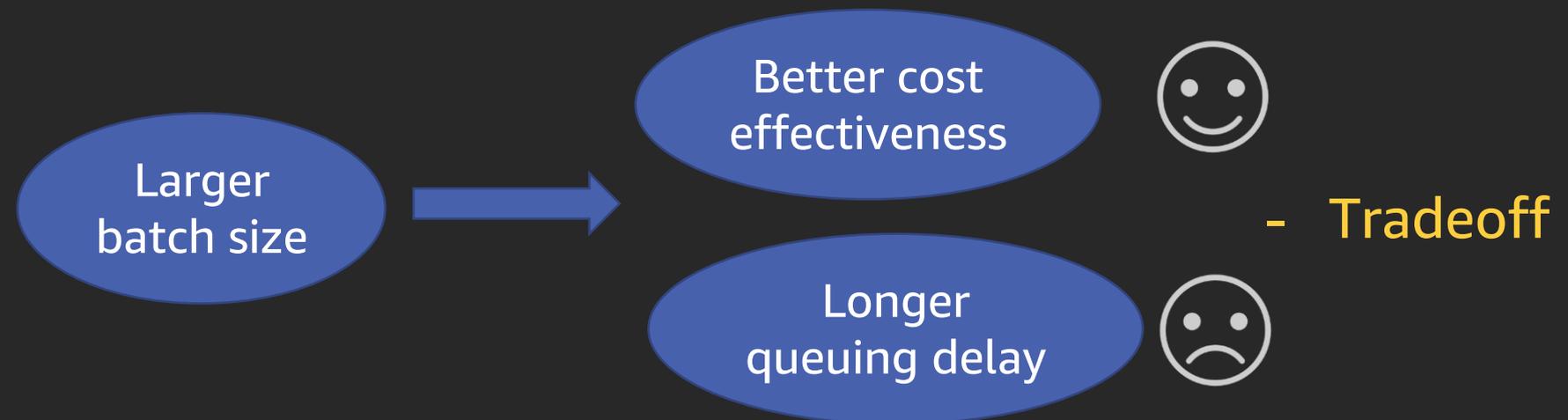
- Choose between CPU and accelerators
- Justify the price tag

# Characterization: CPU vs. GPU

CPU: 1 vCPU, 2 GB mem; GPU: K80; TPU: TPU-v2



- CPU: no significant benefits for small instances
- GPU: substantial benefit
- GPUs can be cheaper, but only with batching and high utilization



# Design consideration: Cloud services for model serving

Infrastructure  
as a Service

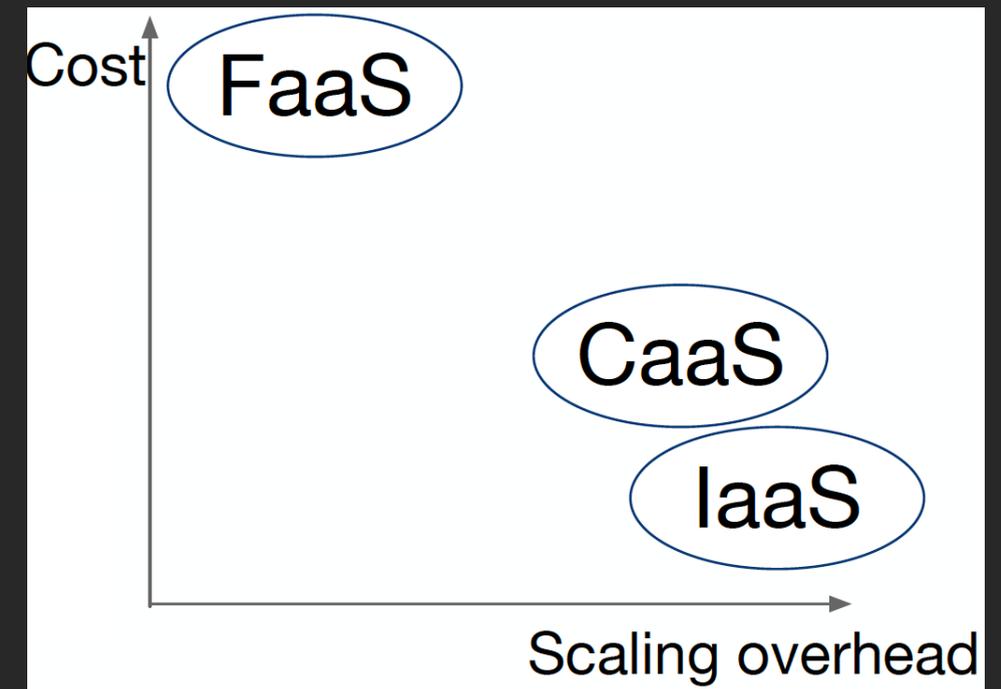
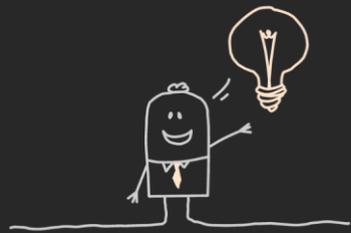
Container  
as a Service

Function  
as a Service  
(FaaS, serverless comp.)

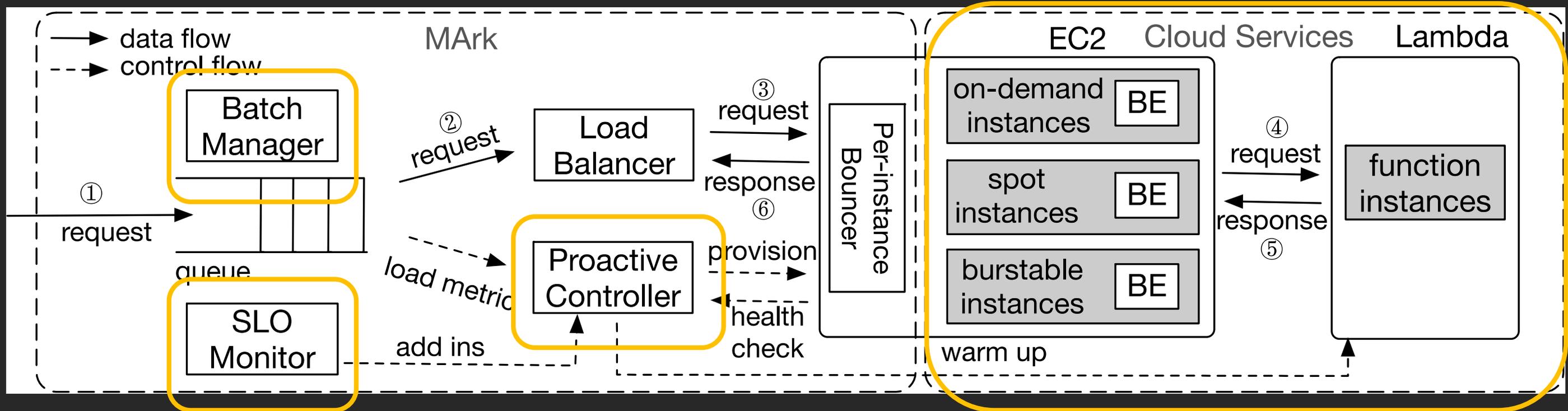
ML Model	EC2		ECS		Lambda	
	\$	$t$ (ms)	\$	$t$ (ms)	\$	$t$ (ms)
Inception-v3	5.0	210	9.17	217	19.0	380
Inception-ResNet	9.3	398	16.4	411	39.3	785
OpenNMT-ende	51.5	2180	96.3	2280	155	3100

EC2: c5.large; ECS: 2vCPU, 4GB mem; Lambda: 3008MB mem

- Combine IaaS's cost advantage with FaaS's scalability
- Instead of overprovisioning IaaS, use FaaS to handle demand surge and spikes



# MARk (Model Ark) using AWS Lambda

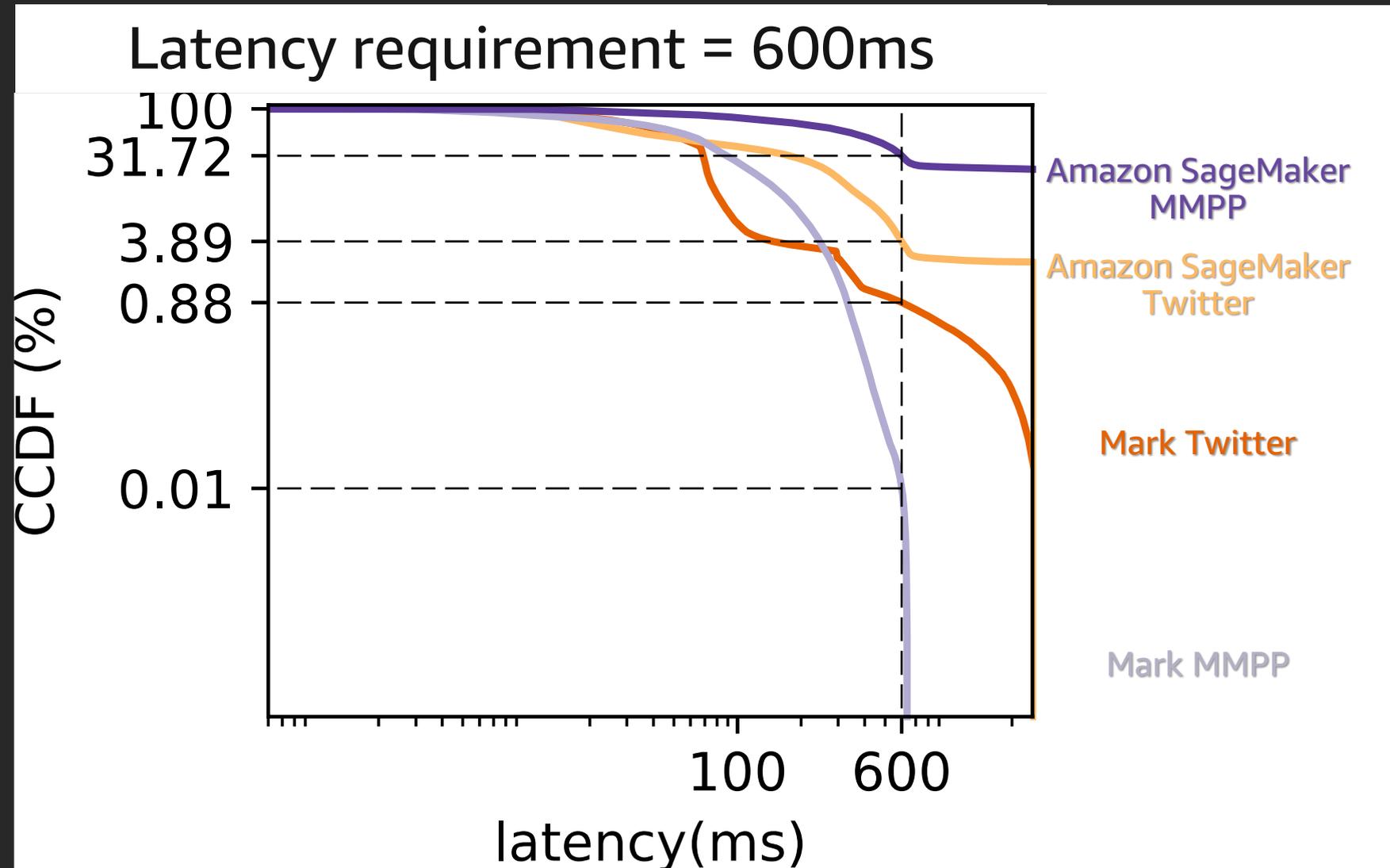


- Weighted round robin for load balancing
- Server front implemented with Sanic framework
- Support different serving frameworks

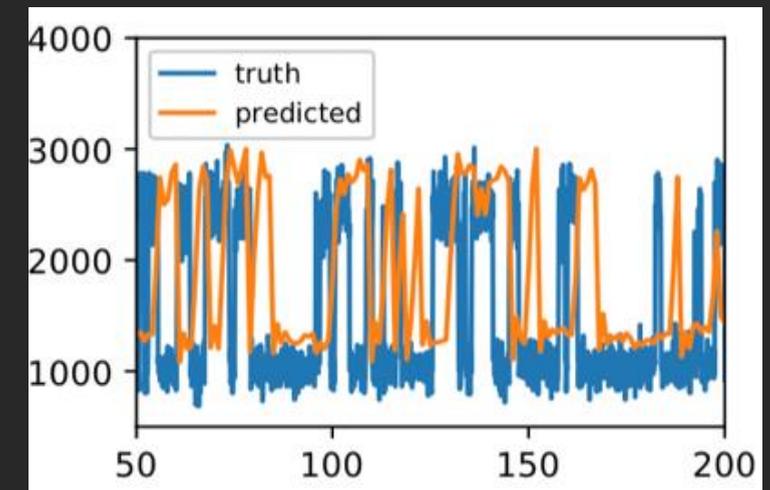
- Nginx and Gunicorn for admission and parallelism control
- Support for Amazon EC2 Spot Instances

# Results

Latency complementary cumulative distribution function



What if workload is unpredictable?



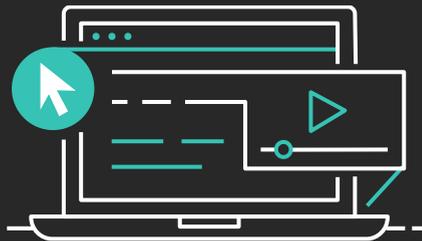
Markov-modulated Poisson Arrivals (MMPP)  
MMPP: unpredictable, highly dynamic workload

# Learn ML with AWS Training and Certification

The same training that our own developers use, now available on demand



Role-based ML learning paths for developers, data scientists, data platform engineers, and business decision makers



70+ free digital ML courses from AWS experts let you learn from real-world challenges tackled at AWS



Validate expertise with the  
**AWS Certified Machine Learning - Specialty** exam

Visit <https://aws.training/machinelearning>

# Thank you!

**Sanjay Padhi**

sanpadhi@amazon.com

**Feng Yan**

fyan@unr.edu



Please complete the session survey in the mobile app.