



# Automated *Model-less* Inference Serving

**Francisco Romero\***, **Qian Li\***,

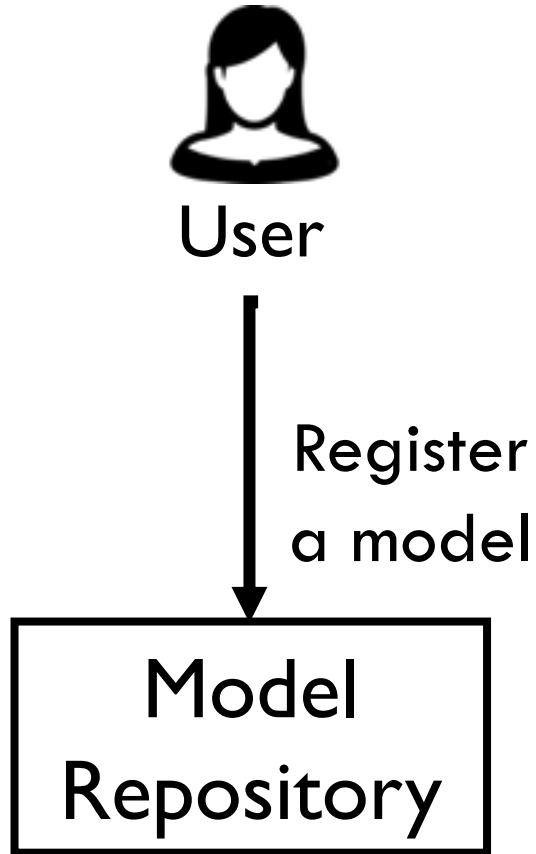
Neeraja J. Yadwadkar, and Christos Kozyrakis



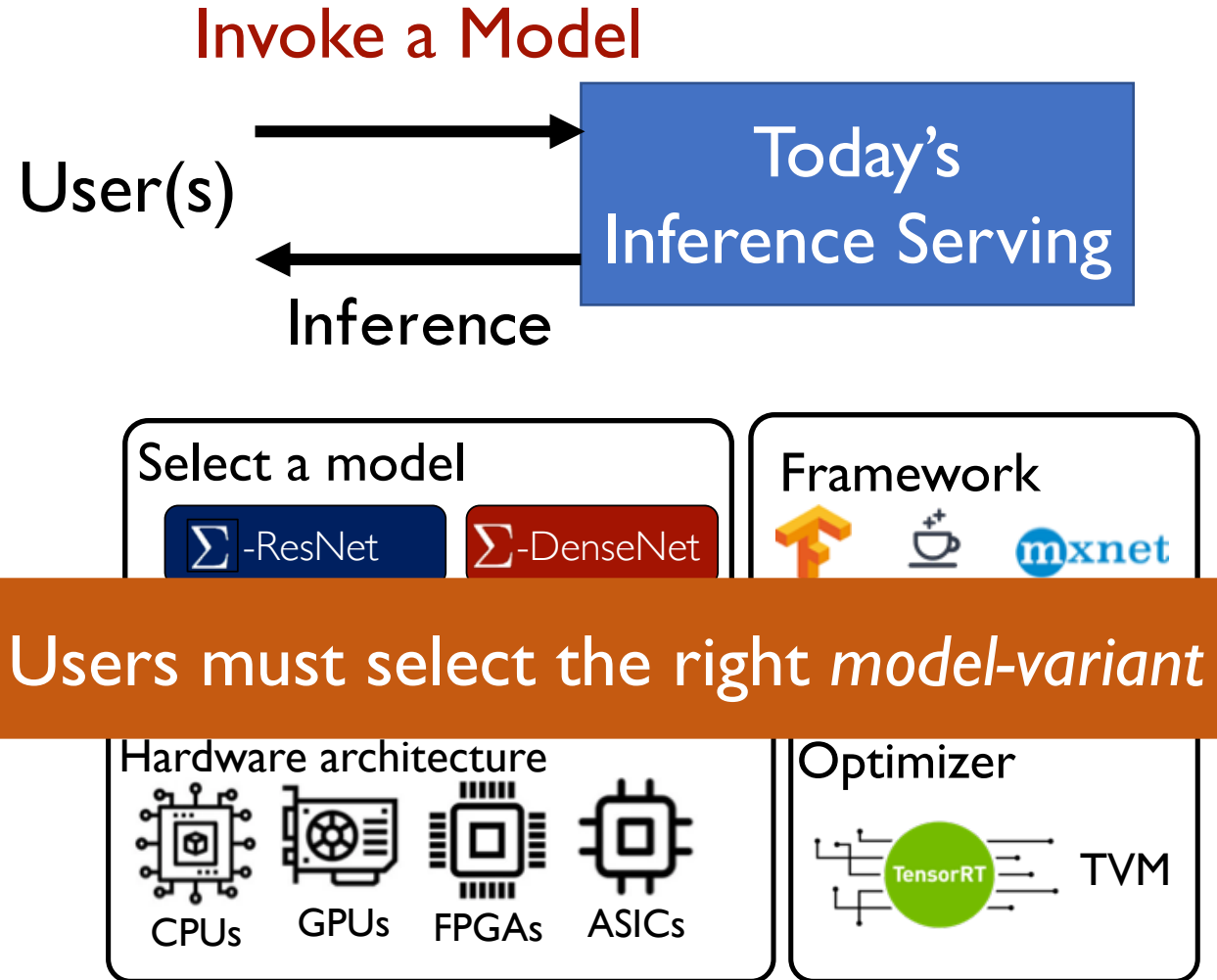
PLATFORMLAB



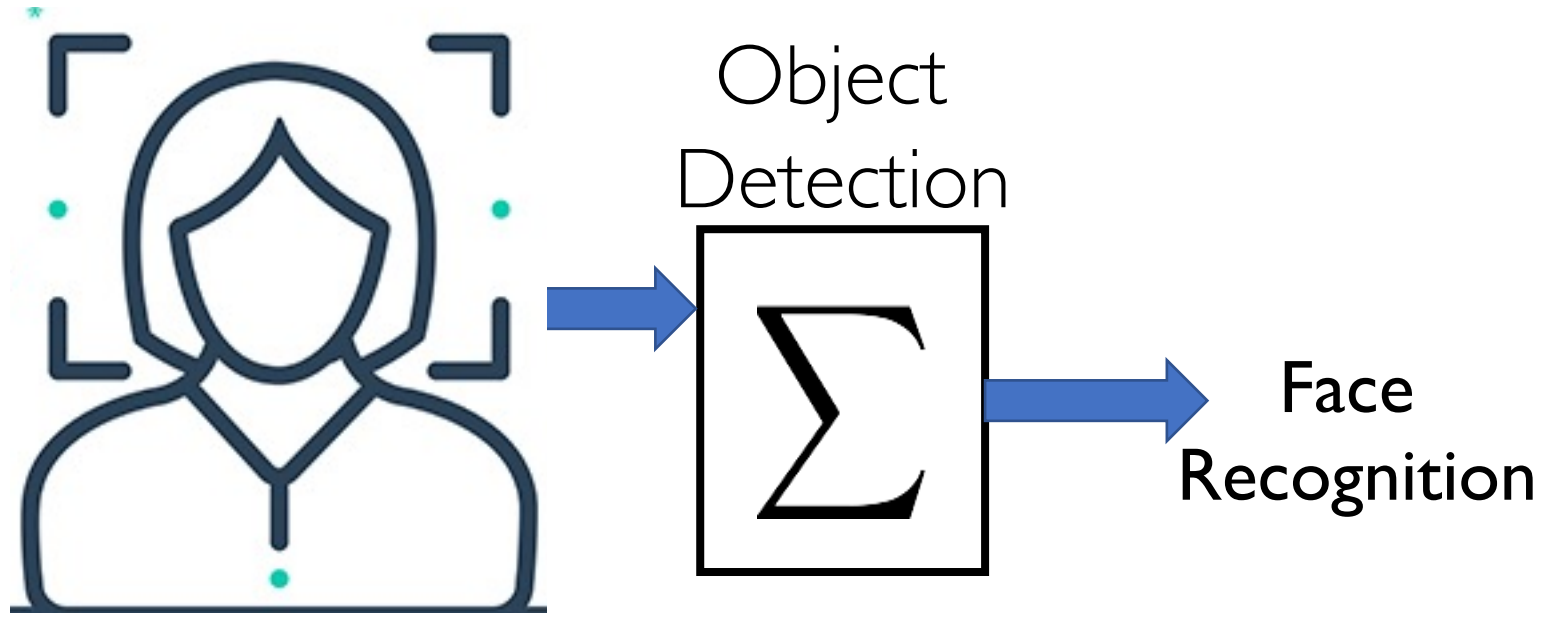
# Model Registration





# Today's Inference Serving





# Diverse application requirements



Example: Face Recognition

 **Accuracy**       **Latency**

Social Media

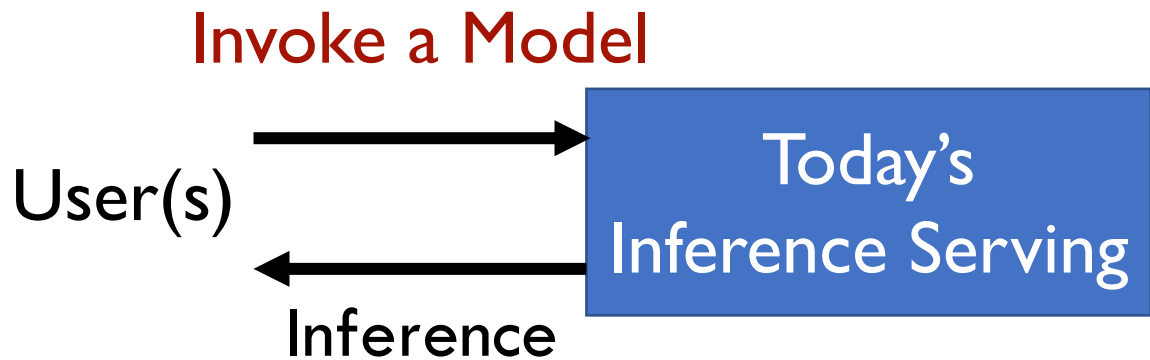
 **Latency**  
**Accuracy**       **Cost**

Navigation for visually impaired person

Applications evolve over time

The performance penalty can be up to 100x

## Today's Inference serving



Select a model

$\Sigma$ -ResNet

$\Sigma$ -DenseNet

$\Sigma$ -SqueezeNet

$\Sigma$ -Inception

Framework



PYTORCH

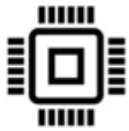
Hardware architecture



CPUs



GPUs



FPGAs



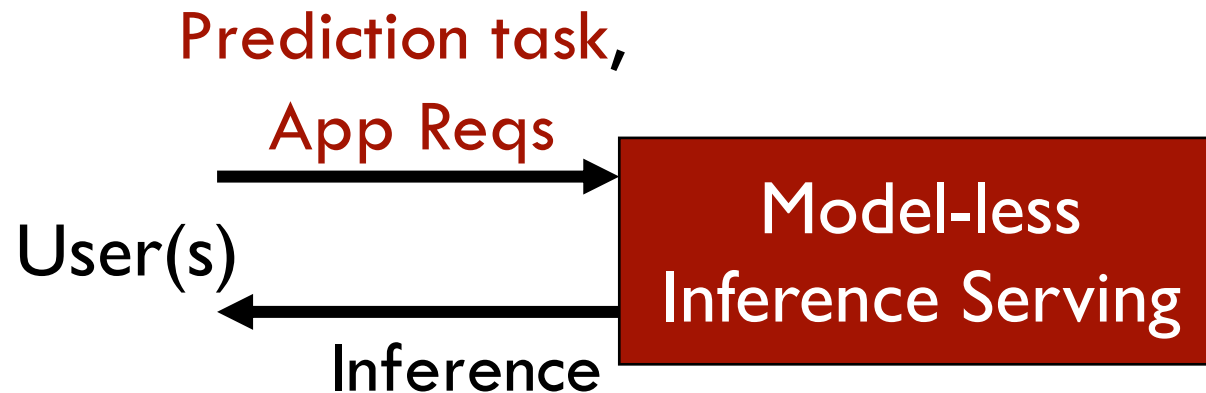
ASICs

Optimizer



TVM

## Model-less Inference serving



Prediction task:

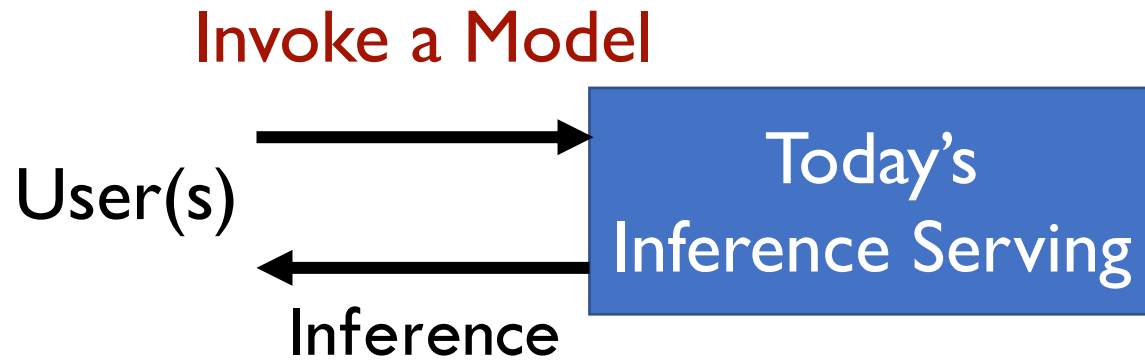
“Face Recognition”

App Reqs:

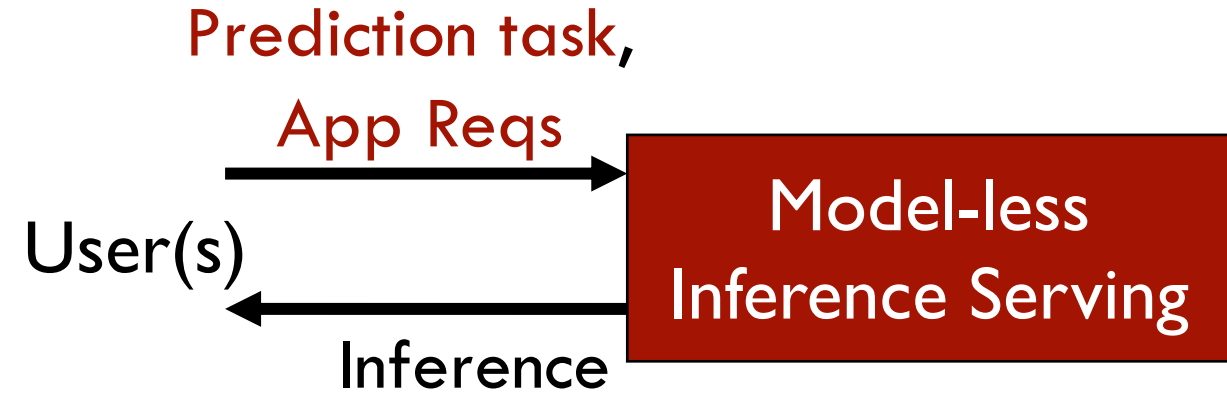
Target accuracy

Target latency

## Today's Inference serving



## Model-less Inference serving



**Easy-to-use:** Automatically and efficiently select a model and hardware

**Cost Efficient:** Share the hardware as well as models across users



## INFerence-as-a-Service system

INFaaS provides a *model-less API* to inference queries that abstracts (a) Model Selection and (b) Resource Provisioning from users.

INFaaS is open-source!

<https://stanford-mast.github.io/INFaaS/>

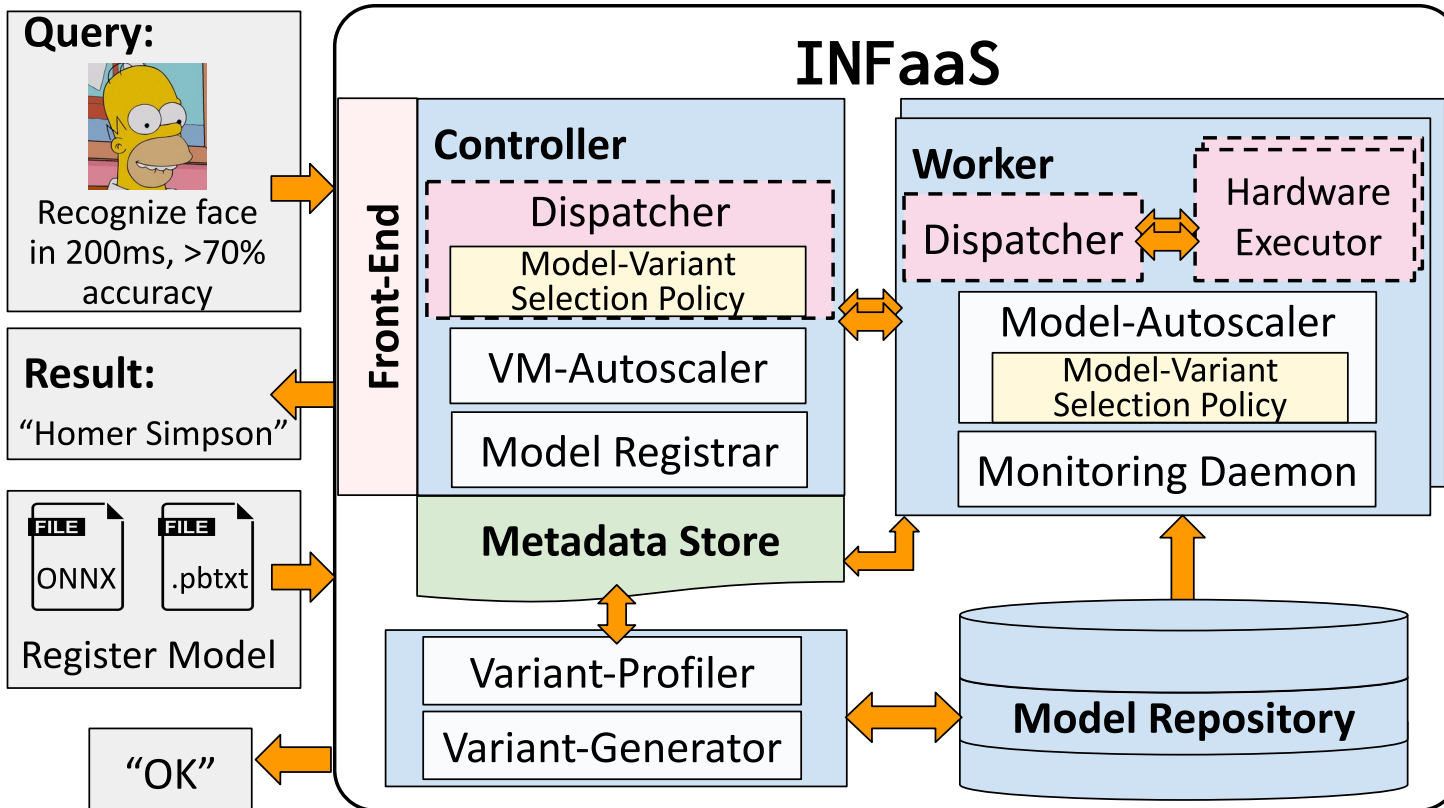
## Goals & Requirements

**Ease-of-use:** Automatically select a model and hardware

## Challenges

- Novice and expert users
- Diverse user requirements
- Large search space
- Decision overhead

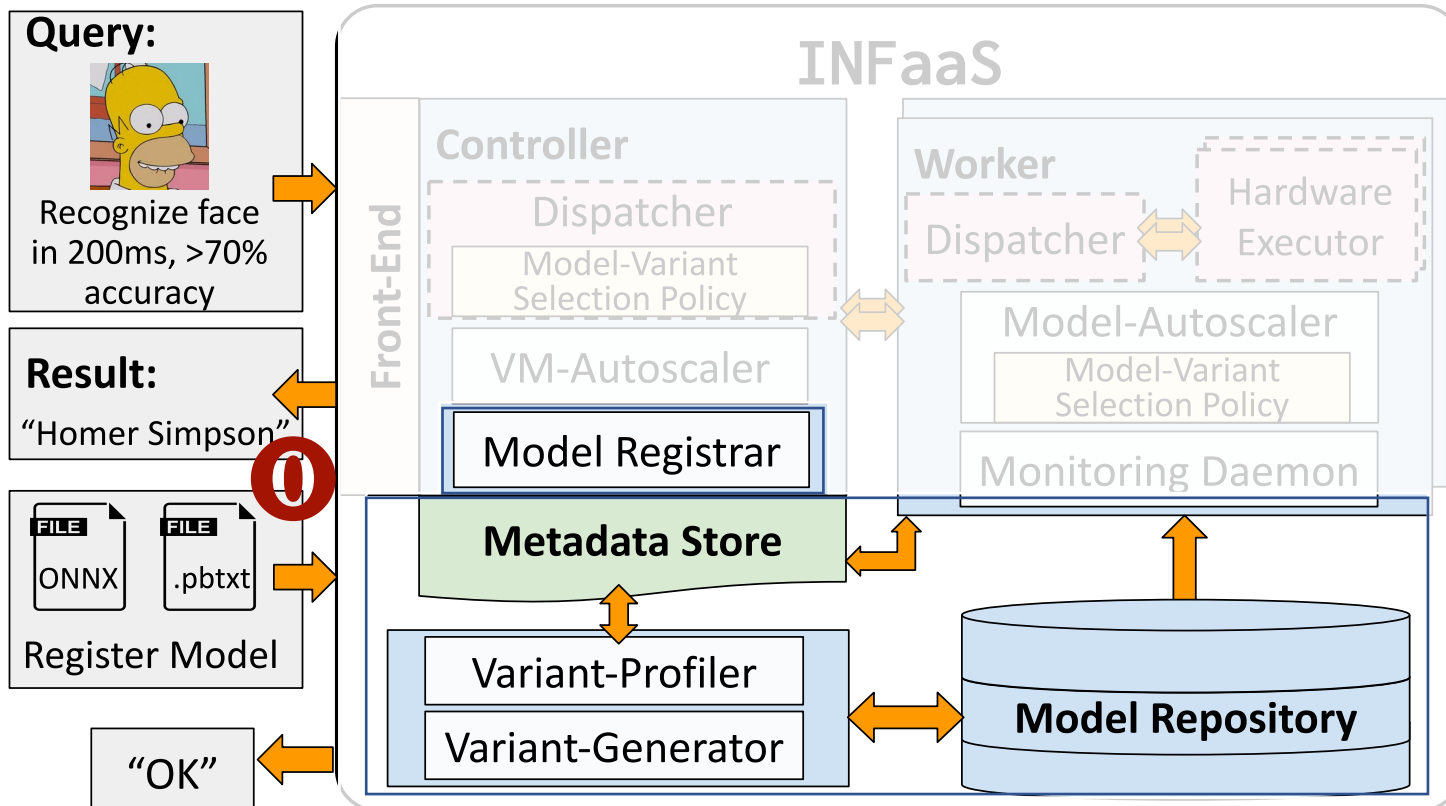
# INFaaS overview



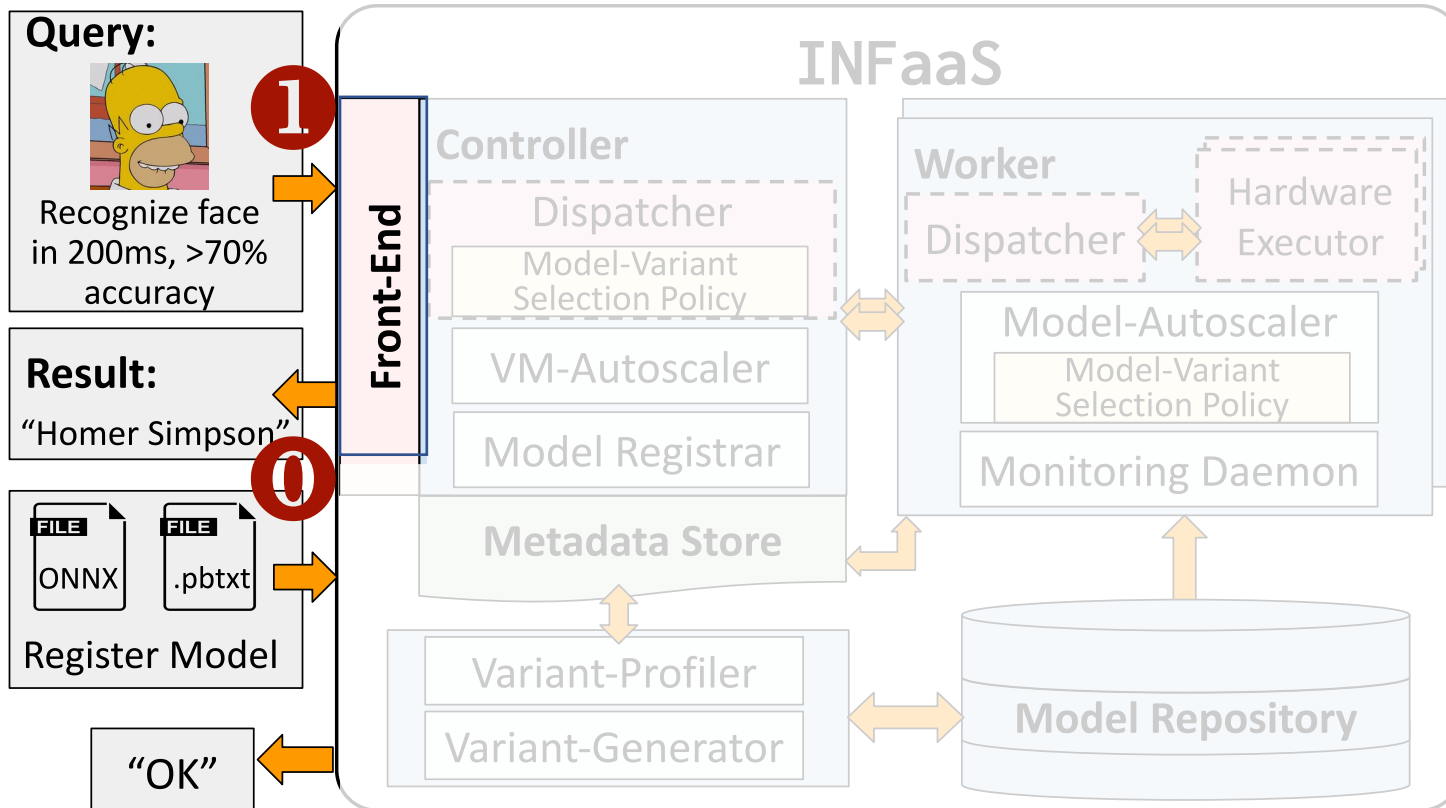


# INFaaS overview

① Users register models



# INFaaS overview



0 Users register models

1 The user submits a query using INFaaS' user API

# Front-end

- Goal: need to map query requirements to models and resources
  - Affects user API, metadata organization, model-variant selection, and autoscaling
- Challenge: needs to be intuitive

# The *model-less* abstraction

Registered model



bert-pytorch-cpu

resnet50-tf-cpu

resnet50-caffe2-gpu

# The *model-less* abstraction

Prediction task



translation

face-detection

Registered model



bert-pytorch-cpu

resnet50-tf-cpu

resnet50-caffe2-gpu

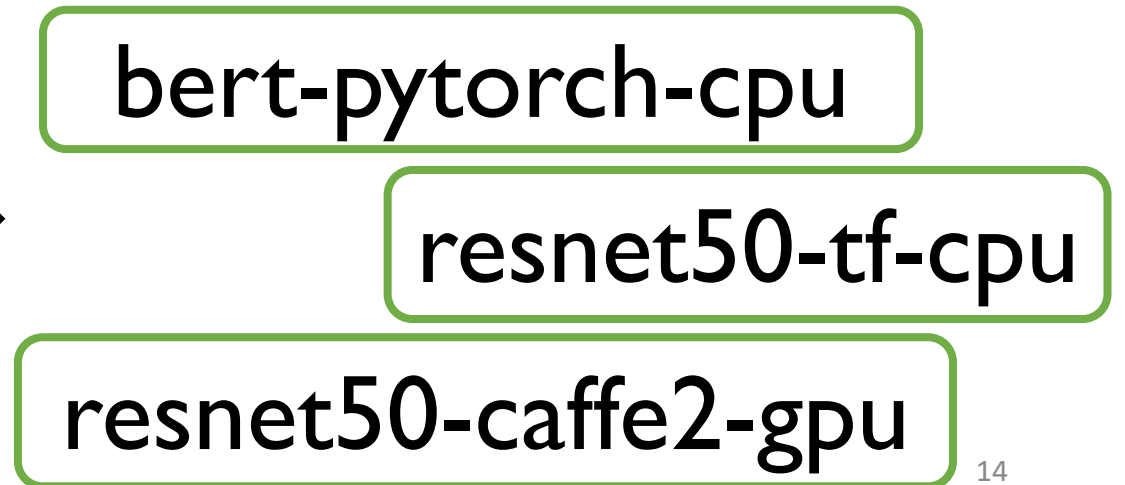
# The *model-less* abstraction

- L** Latency target
- A** Accuracy target

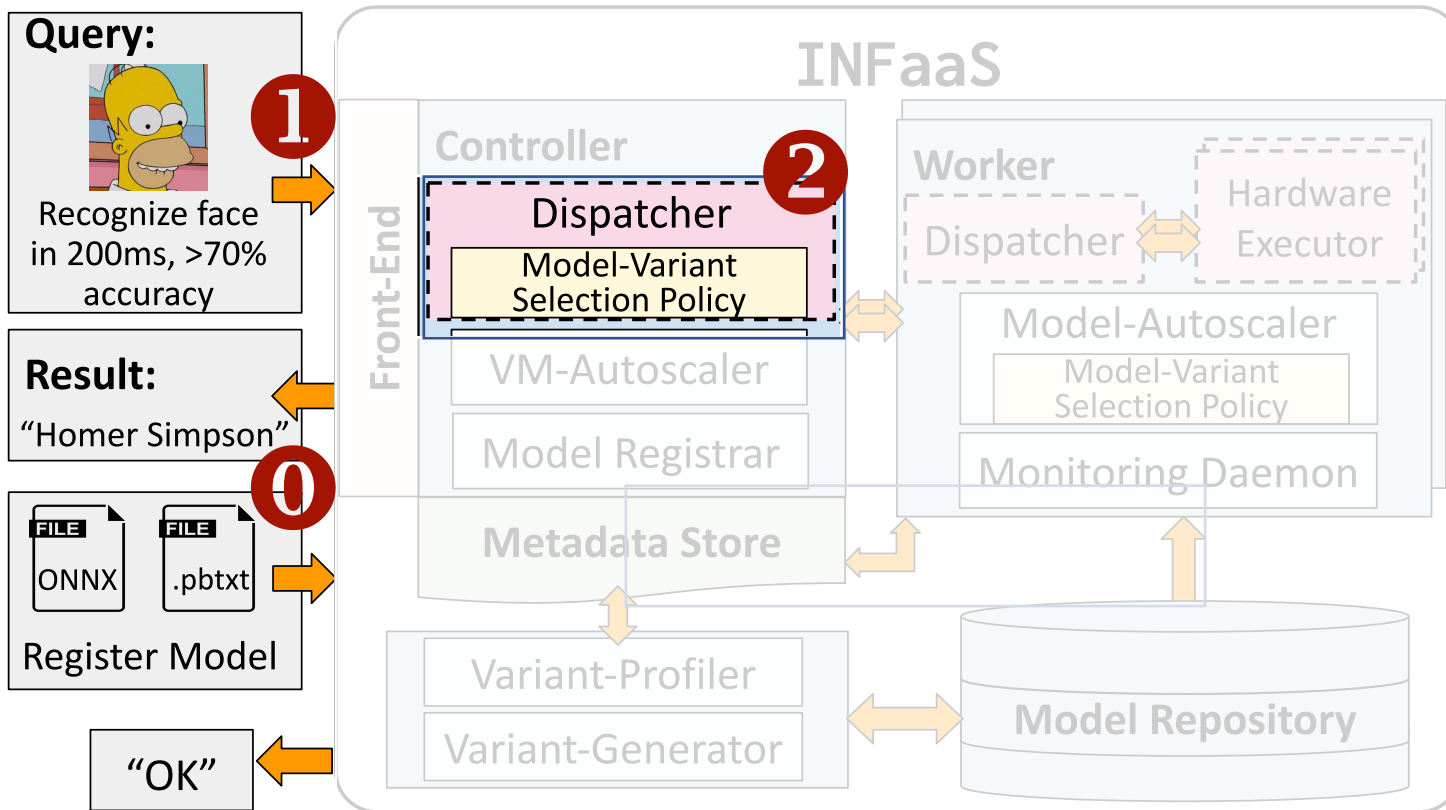
Prediction task



Registered model



# INFaaS overview

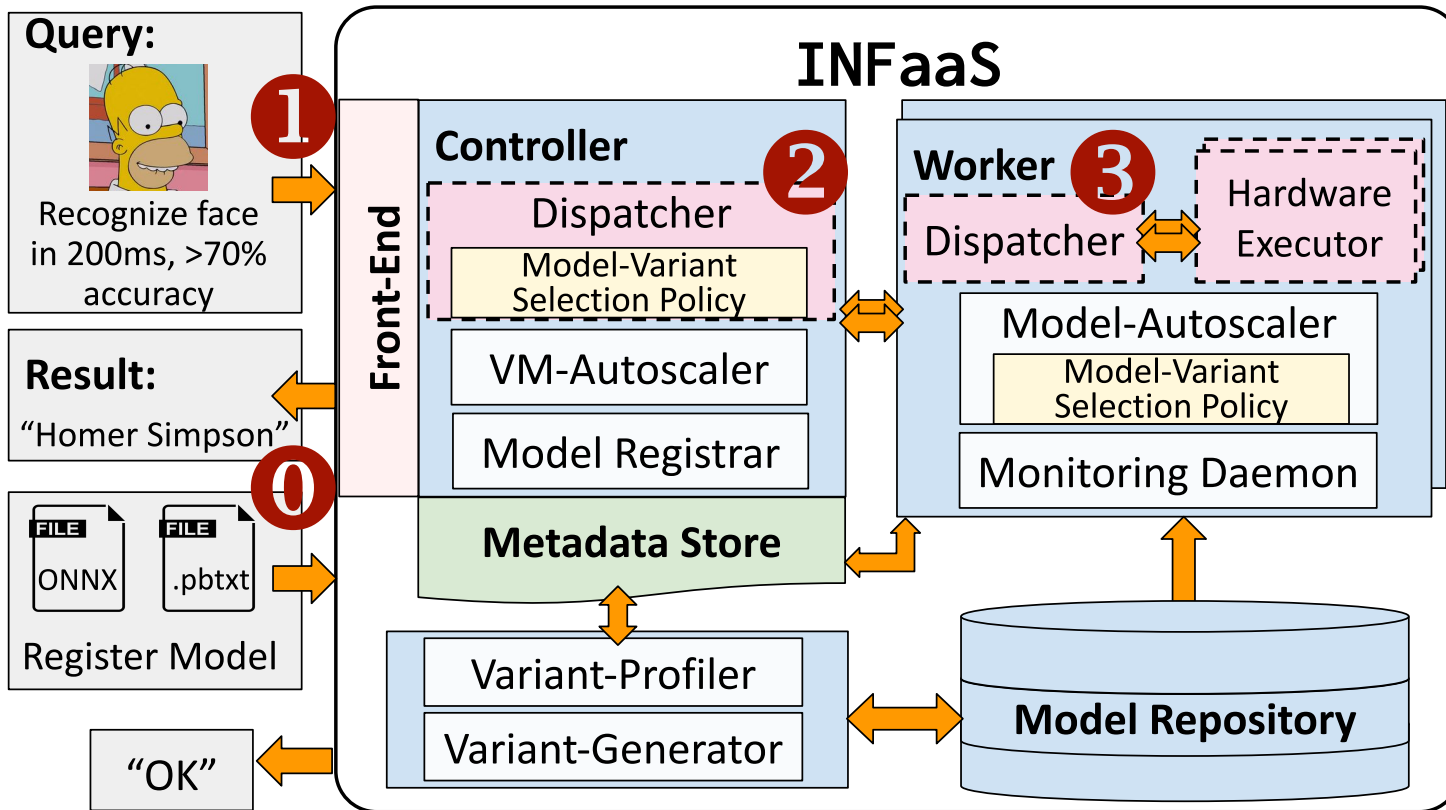


0 Users register models

1 The user submits a query using INFaaS' user API

2 The Controller selects a model-variant, then selects a worker to process the query

# INFaaS overview



0 Users register models

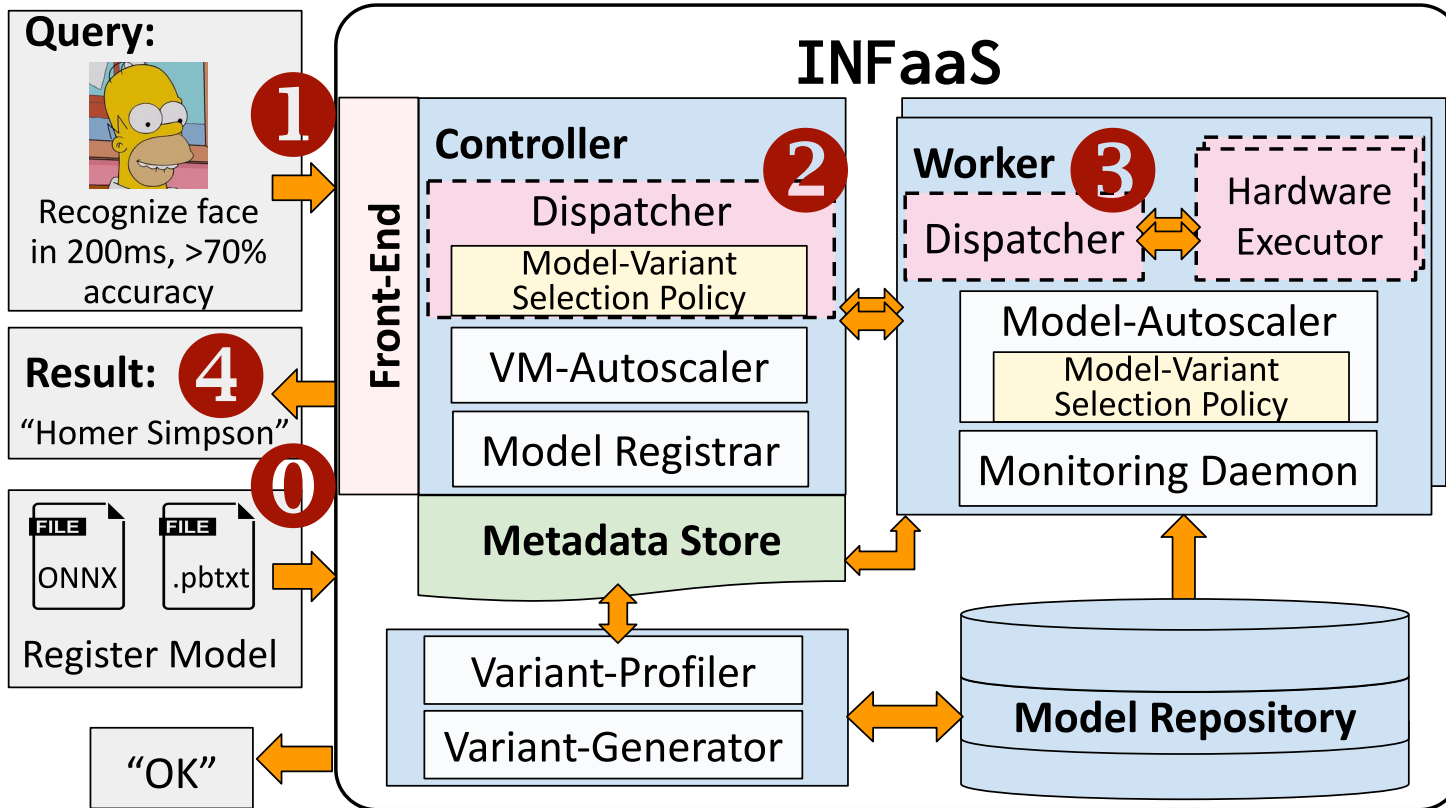
1 The user submits a query using INFaaS' user API

2 The Controller selects a model-variant, then selects a worker to process the query

3 The query proceeds to run on the variant's target hardware platform



# INFaaS overview



0 Users register models

1 The user submits a query using INFaaS' user API

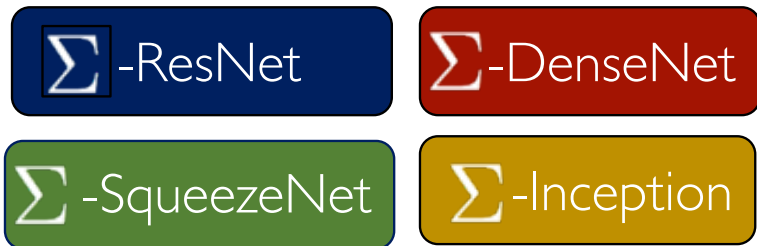
2 The Controller selects a model-variant, then selects a worker to process the query

3 The query proceeds to run on the variant's target hardware platform

4 Upon completion, the result is returned to the user

# Ease-of-use and cost efficiency

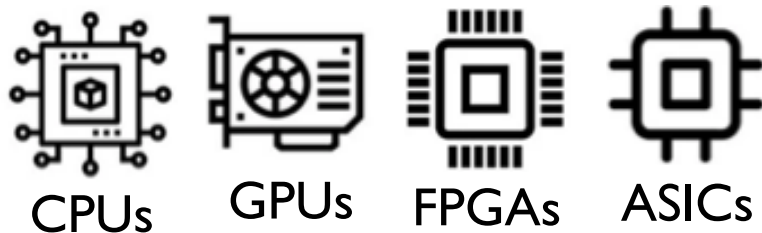
Select a model



Framework



Hardware architecture



Optimizer



INFaaS removes the system configuration burden and improves ease-of-use

## Goals & Requirements

**Ease-of-use:** Automatically select a model and hardware

**Autoscaling:** allocate just enough resources, meet SLOs, minimize cost

## Challenges

- Novice and expert users
- Diverse user requirements
- Large search space
- Decision overhead

- Query load and pattern changes
- Heterogeneous hardware & models
- Scalability

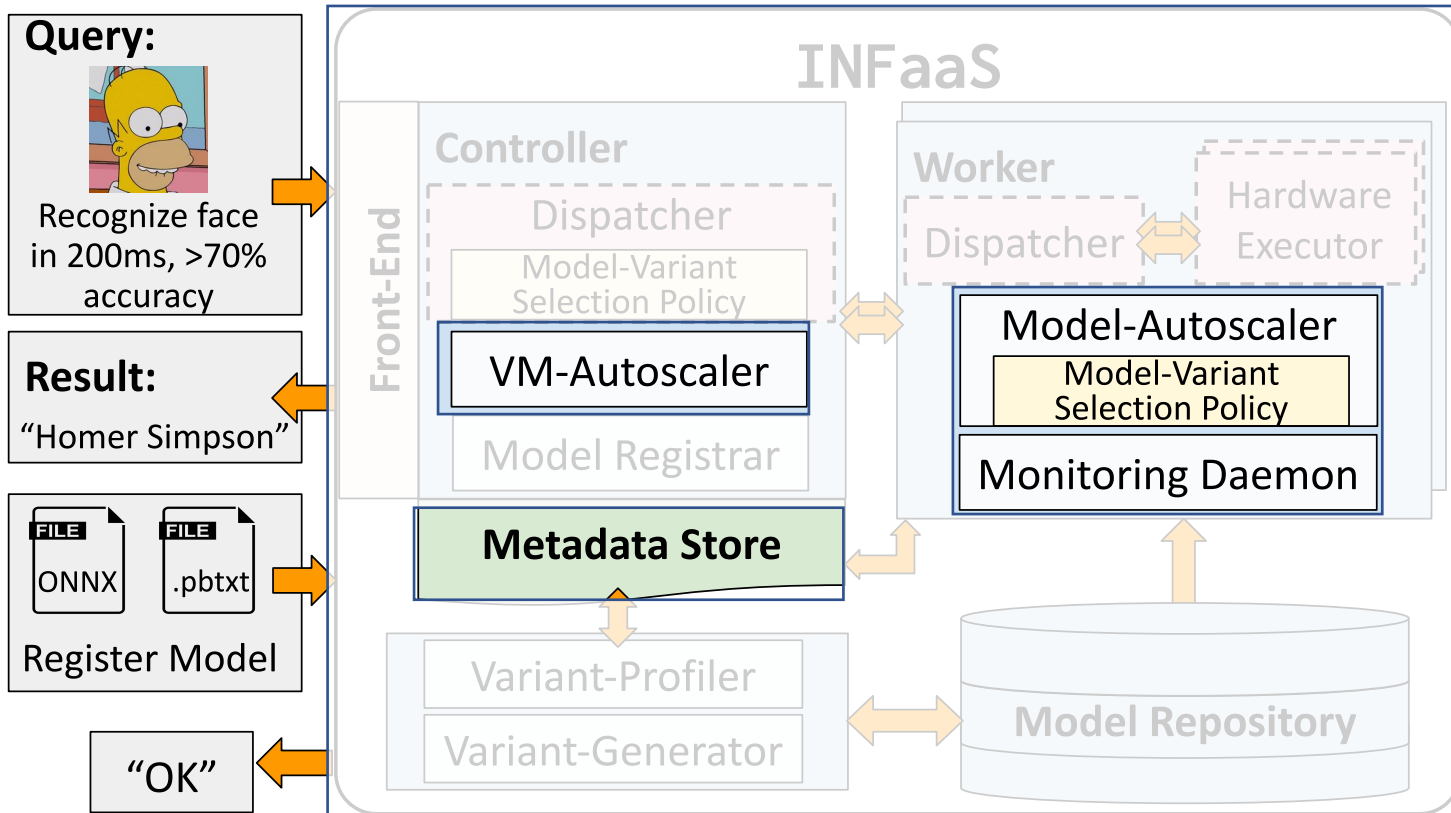
# Existing systems

- Static provisioning (TensorFlow Serving, Triton Inference Server)
  - based on peak load
  - Meet SLOs but expensive
  - Waste resources at low load

# Existing systems

- Static provisioning (TensorFlow Serving, Triton Inference Server)
  - based on peak load
  - Meet SLOs but expensive
  - Waste resources at low load
- Replica-only (Clipper, SageMaker, AI Platform) - replicate individual variants
  - Lower costs but high start-up latency
  - Fail to leverage heterogeneous resources / variants

# Autoscaling

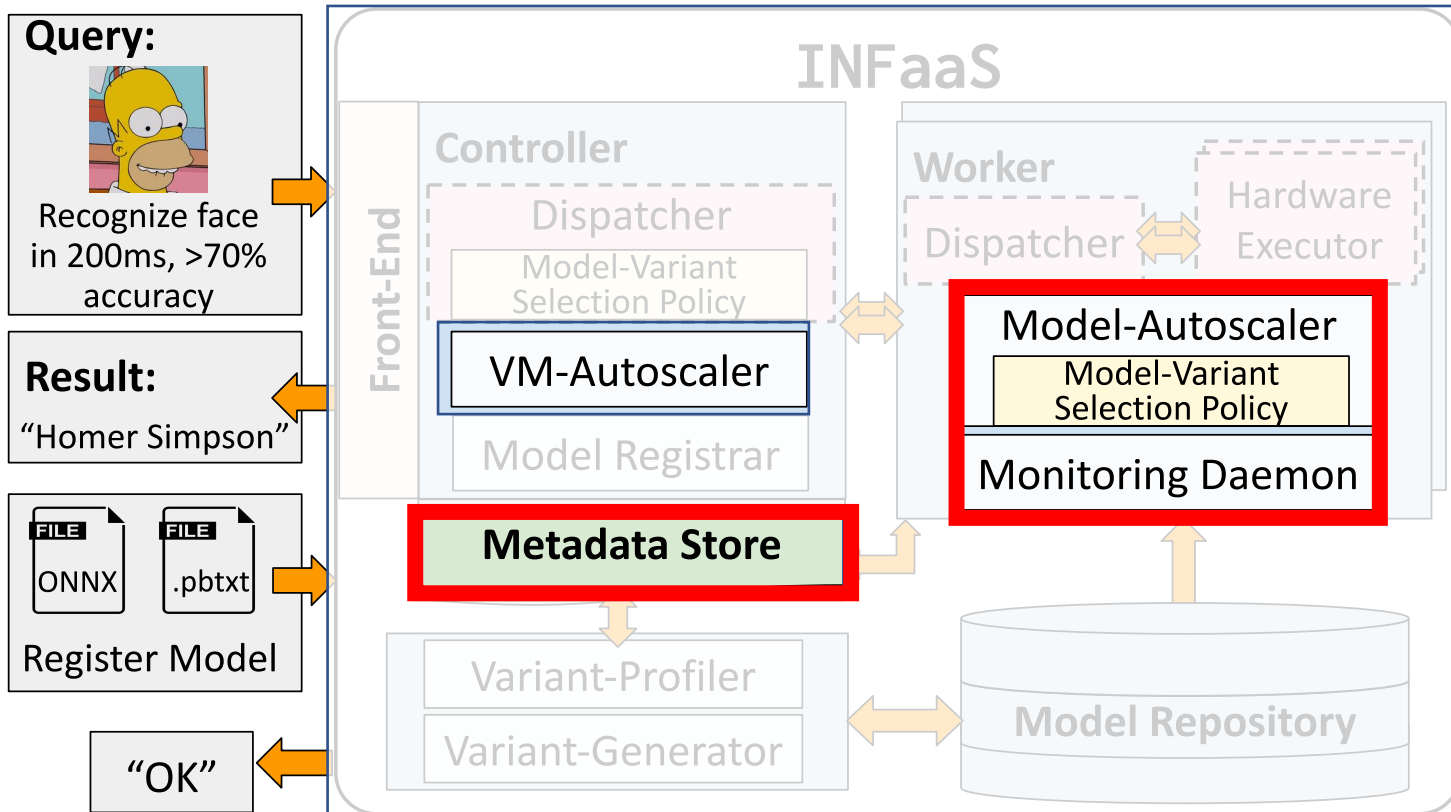


## 3 types of scaling

- Model-horizontal scaling
- Model-vertical scaling  
-> (*Our contribution*)
- VM-autoscaling

## Division of responsibility

# Autoscaling



## 3 types of scaling

- **Model-horizontal scaling**
- **Model-vertical scaling**  
-> *(Our contribution)*
- **VM-autoscaling**

## Division of responsibility

# Model-Autoscaler at each worker

- Goal: Decide the *type* and *number* of model-variants to meet the load and requirements, while minimizing cost
- Formulate as an Integer Programming problem

$$\min \text{Cost}(\text{action}) = \text{Hardware Cost} + \lambda \text{Loading Latency}$$

  
{load, unload} variant instances



# Model-Autoscaler at each worker

- Goal: Decide the *type* and *number* of model-variants to meet the load and requirements, while minimizing cost
- Formulate as an Integer Programming problem

$$\min \text{Cost}(\text{action}) = \text{Hardware Cost} + \lambda \text{Loading Latency}$$

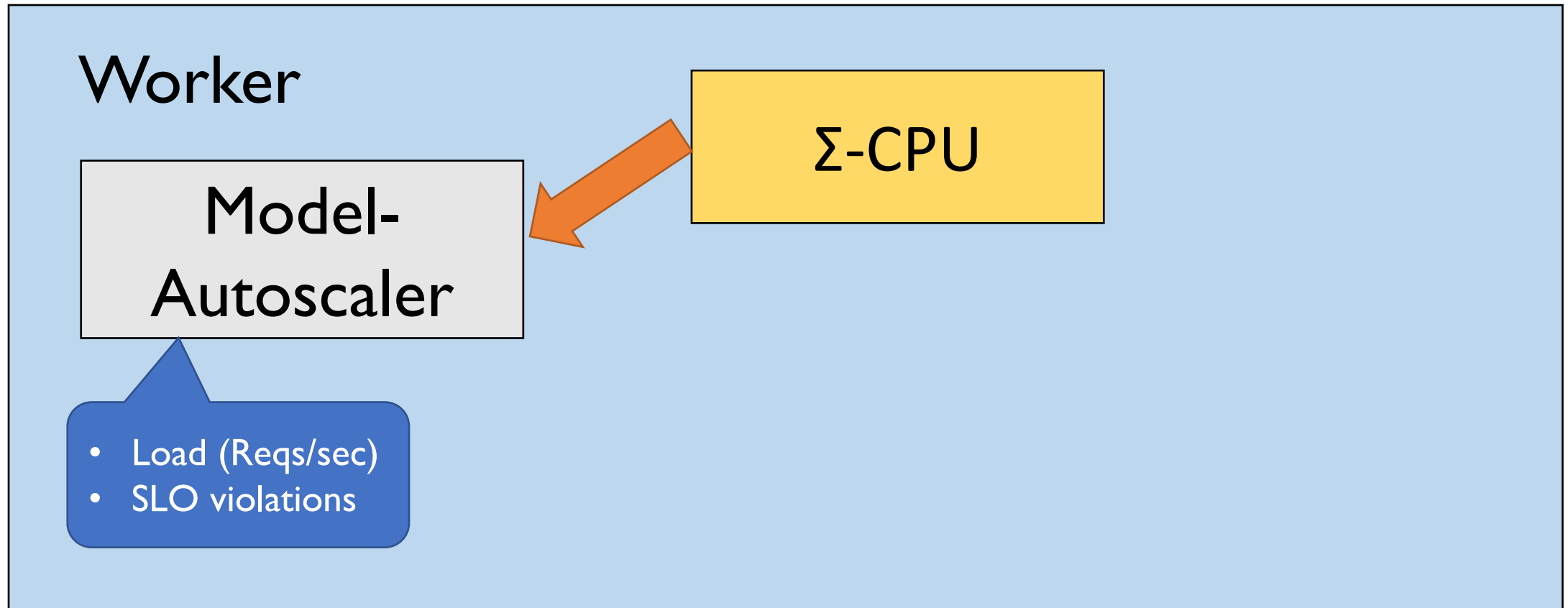
  
{load, unload} variant instances

## Constraints:

- (1) With the chosen scaling action, INFaaS supports the incoming query load.
- (2) The newly-loaded instances satisfy applications' SLOs.
- (3) Do not exceed the total system resources.
- (4) The number of running instances is non-negative.

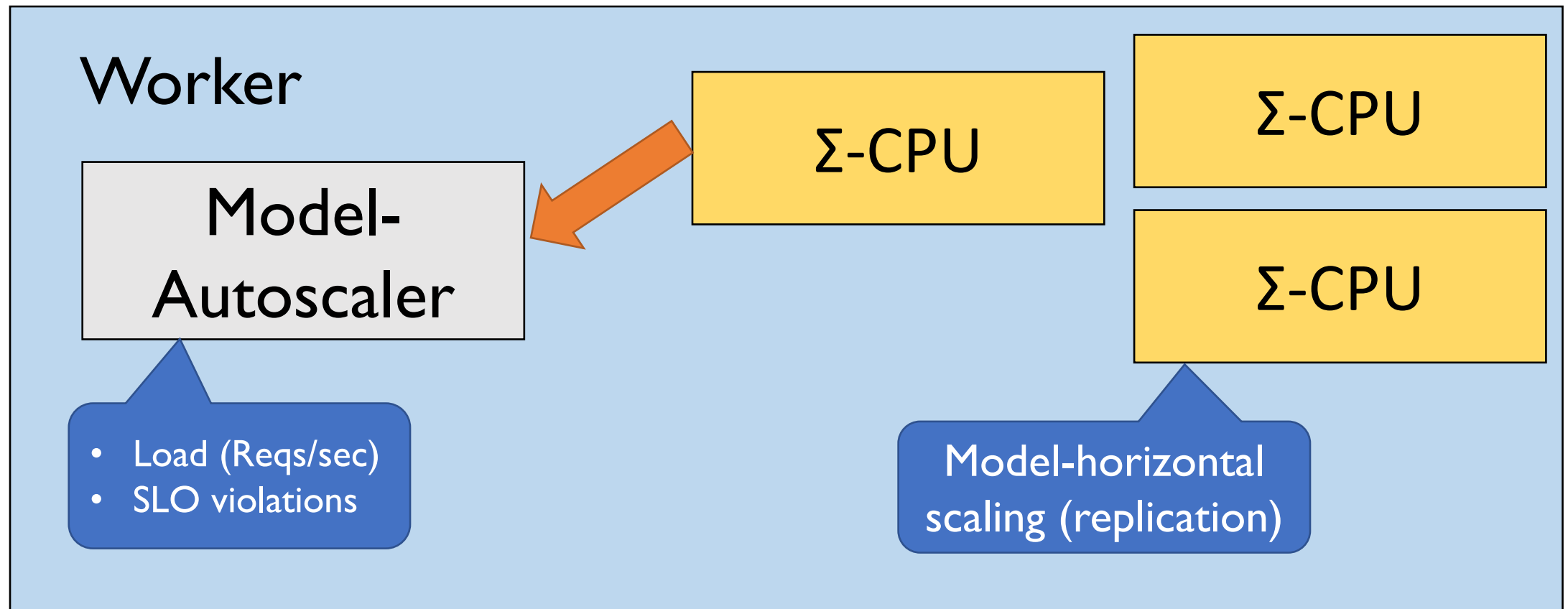
# Model-Autoscaler at each worker

- Respond to changes in load and meet SLOs by: 1) **model-horizontal scaling** and 2) **model-vertical scaling**



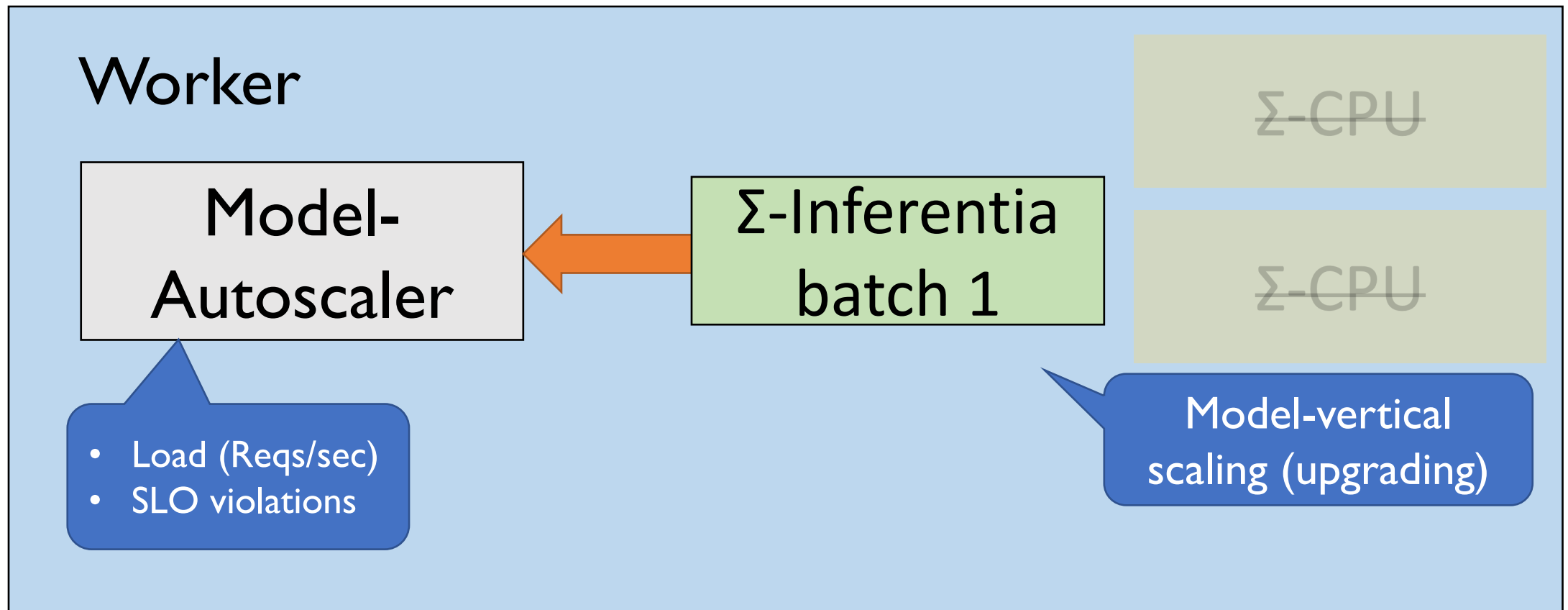
# Model-Autoscaler at each worker

- Respond to changes in load and meet SLOs by: 1) **model-horizontal scaling** and 2) **model-vertical scaling**



# Model-Autoscaler at each worker

- Respond to changes in load and meet SLOs by: 1) **model-horizontal scaling** and 2) **model-vertical scaling**



# Evaluation

- Baselines:
  - **CLIPPER<sup>+</sup>** (Clipper, TIS and TFS): preloaded and persisted *beefy* variants
    - CLIPPER<sup>+</sup><sub>GPU</sub> and CLIPPER<sup>+</sup><sub>CPU</sub>
  - **SM<sup>+</sup>** (InferLine, SageMaker, and AI Platform): model-horizontal scaling, replicated *light-weight* variants on/across worker
    - SM<sup>+</sup><sub>GPU</sub> and SM<sup>+</sup><sub>CPU</sub>

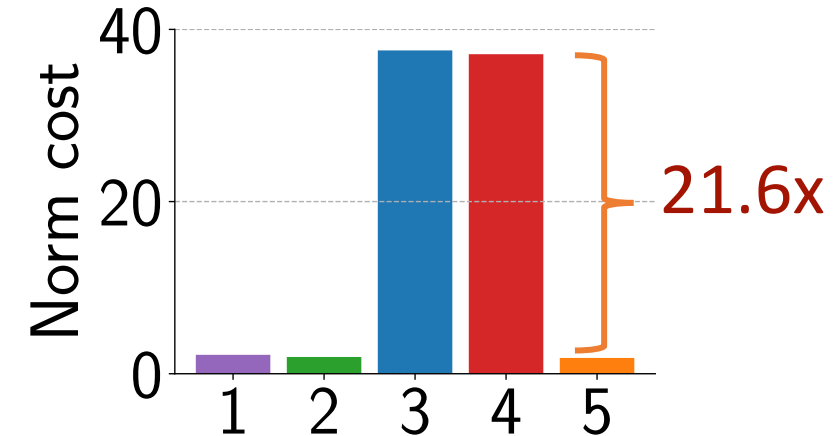
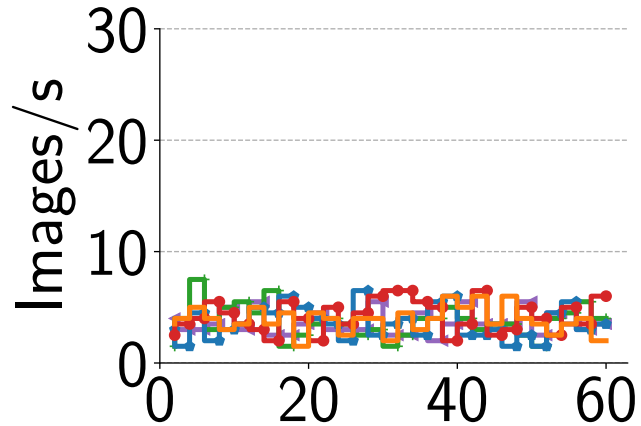
# Evaluation

- We deployed INFaaS on AWS EC2
  - GPU worker has 1 NVIDIA V100 GPU
  - Inference worker has 1 AWS Inference accelerator
  - Controller / CPU worker / client are CPU-only machines

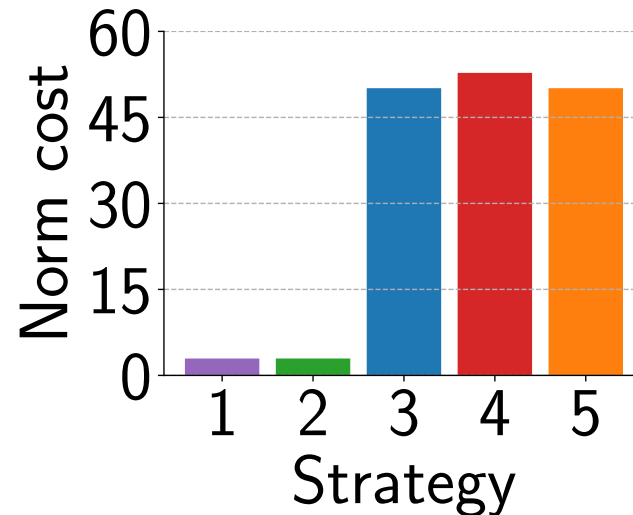
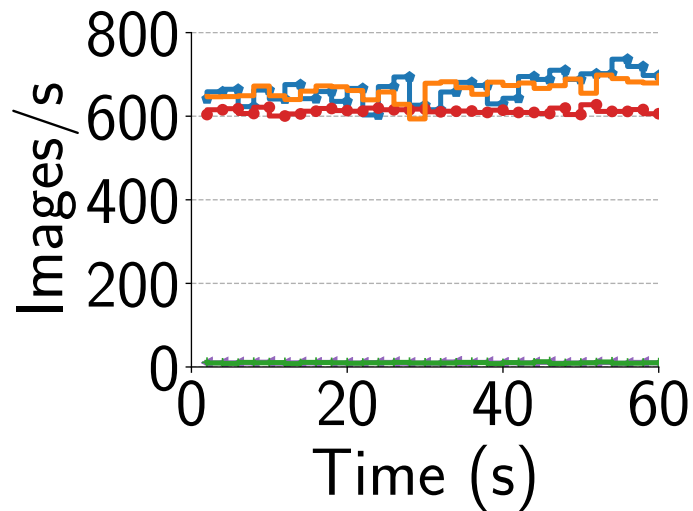
# How well does INFaaS scale with load?

Clipper<sub>CPU</sub><sup>+</sup> SM<sub>CPU</sub><sup>+</sup> Clipper<sub>GPU</sub><sup>+</sup> SM<sub>GPU</sub><sup>+</sup> INFaaS

low load



high load

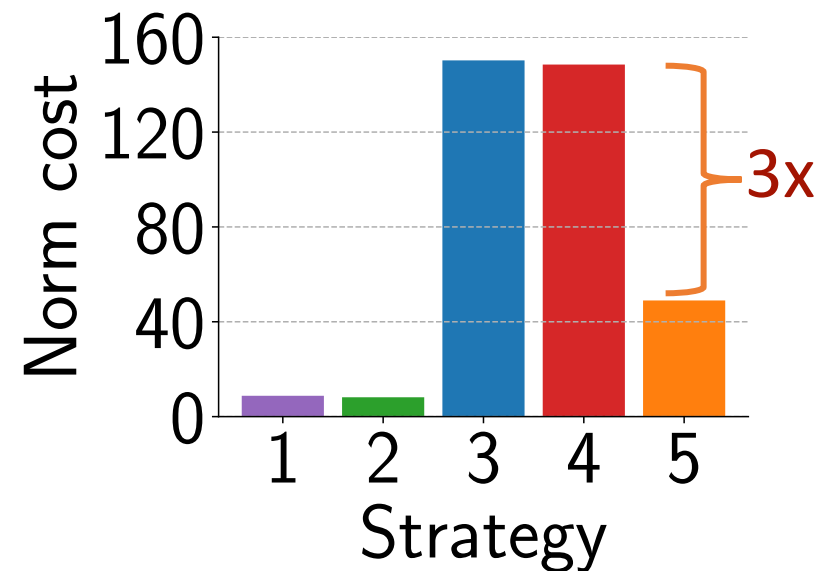
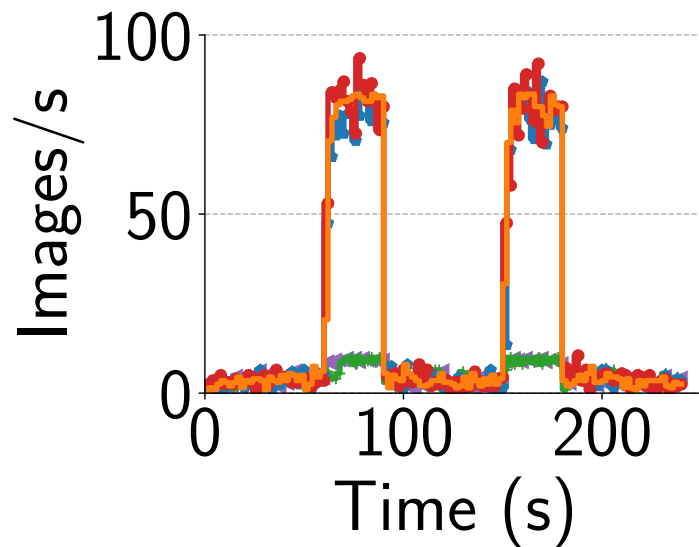


INFaaS achieved load while minimizing cost

# How well does INFaaS scale with load?

Clipper<sub>CPU</sub><sup>+</sup> SM<sub>CPU</sub><sup>+</sup> Clipper<sub>GPU</sub><sup>+</sup> SM<sub>GPU</sub><sup>+</sup> INFaaS

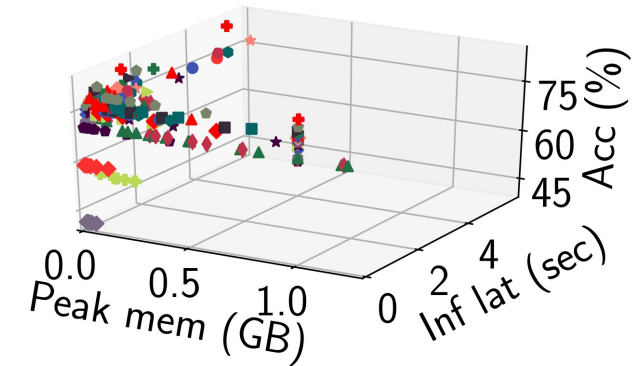
Fluctuating, spiky load



INFaaS reduced cost by 3x by leveraging CPU/Inferentia variant;  
if limited to CPU/GPU variants, still 1.7x cheaper



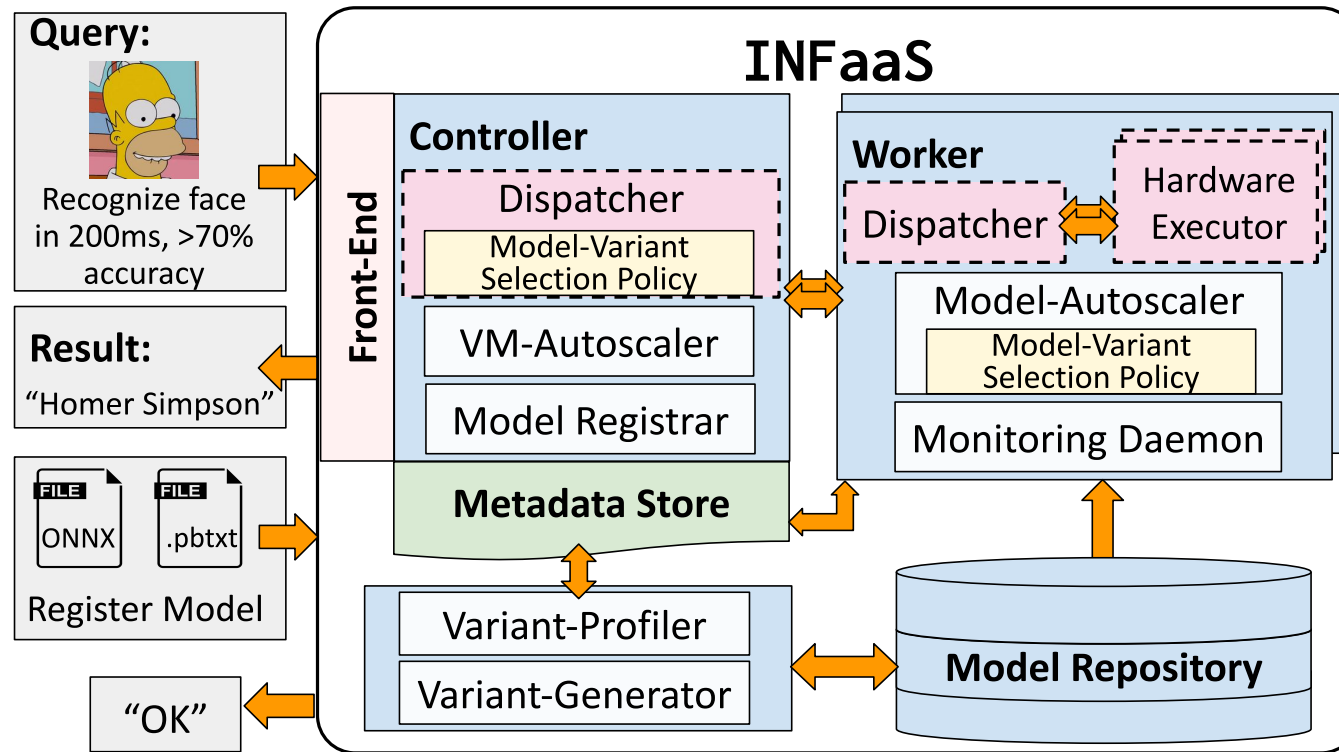
# Putting it all together



- Real workload: Twitter trace (diurnal pattern + spikes) **175** model-variants
- Compared to CLIPPER<sup>+</sup> and SM<sup>+</sup>:
  - **1.1x, 1.3x** higher throughput versus CLIPPER<sup>+</sup>, SM<sup>+</sup>
  - **1.6x, 2.5x** fewer SLO violations compared to CLIPPER<sup>+</sup>, SM<sup>+</sup>
  - **1.23x** lower cost by leveraging CPU, GPU, Inferentia machines

INFaaS achieved high performance, better resource utilization, lower SLO violations, and reduced cost

# Conclusion



<https://stanford-mast.github.io/INFaaS/>

**Contact us:**

{faromero,qianl15,neerajay,kozyraki}@stanford.edu