# Serving DNNs like <u>Clockwork</u>
## Performance Predictability from the Bottom Up

# Serving DNNs like Clockwork
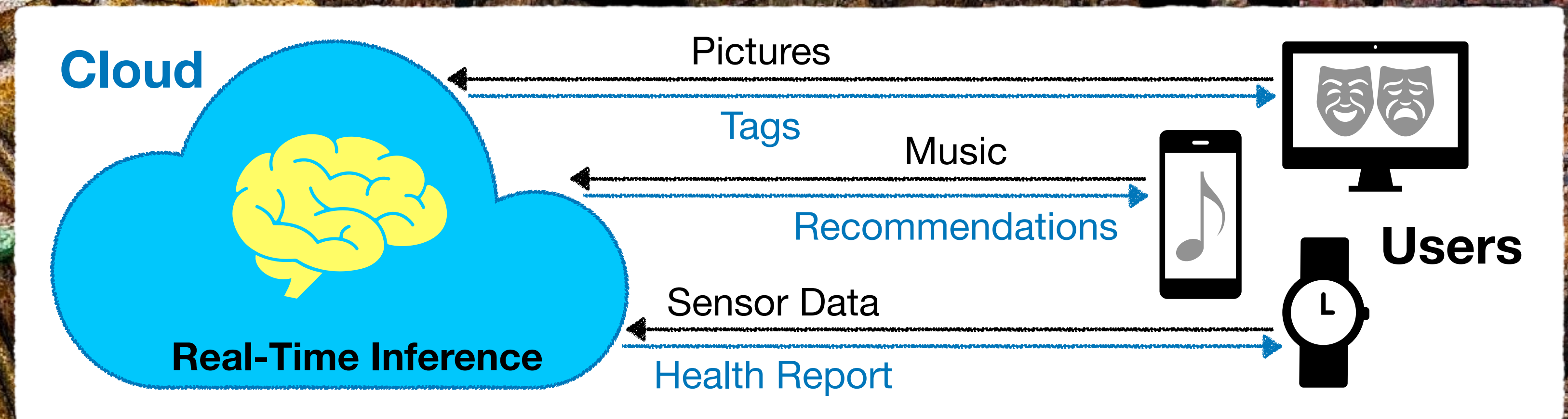## Performance Predictability from the Bottom Up



**DNN inference has a very predictable execution time!**

DNN

**Clockwork**
**End-to-end predictable DNN serving platform for the Cloud**

Cloud

Real-Time Inference

Pictures
Tags
Music
Recommendations
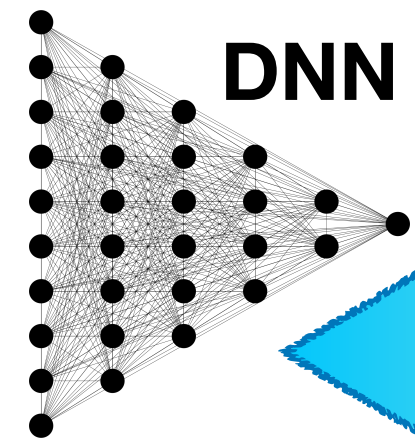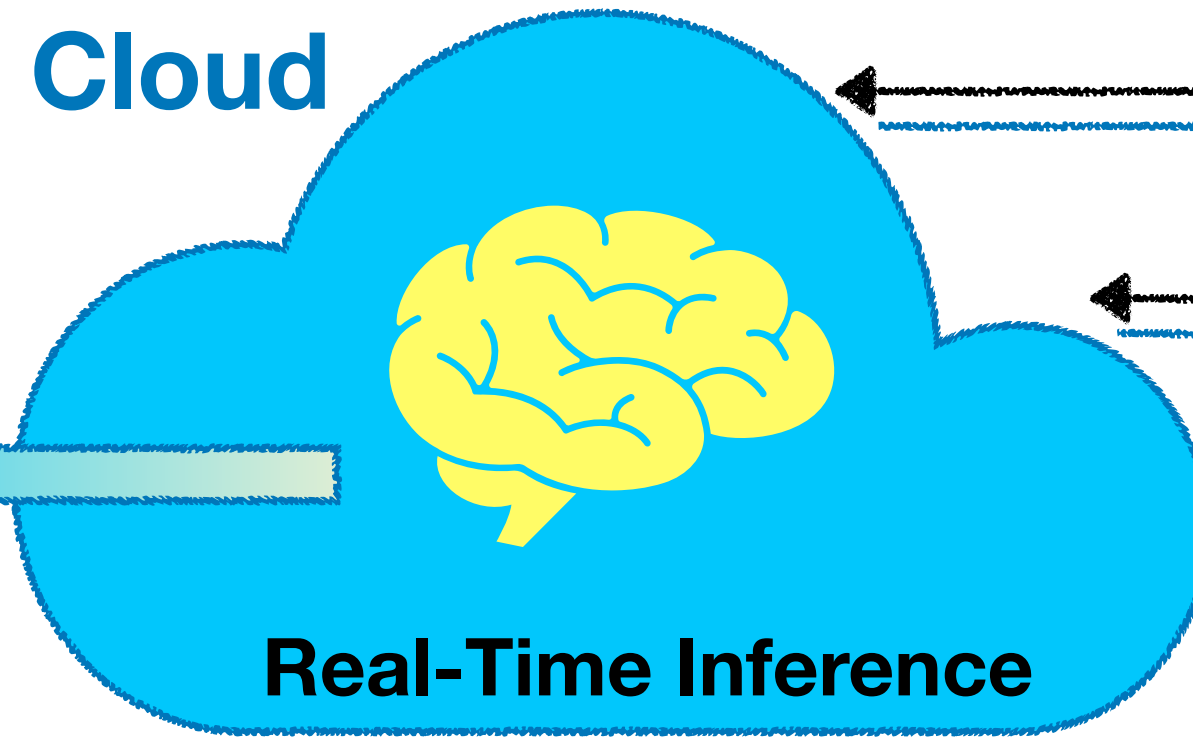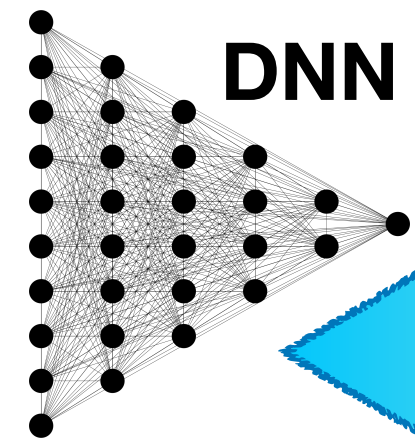Sensor Data
Health Report

Users

# Serving DNNs like Clockwork
## Performance Predictability from the Bottom Up

**DNN inference has a very predictable execution time!**

DNN

**Cloud**

Real-Time Inference

Pictures
Tags
Music
Recommendations
Sensor Data
Health Report

**Users**

**Clockwork**

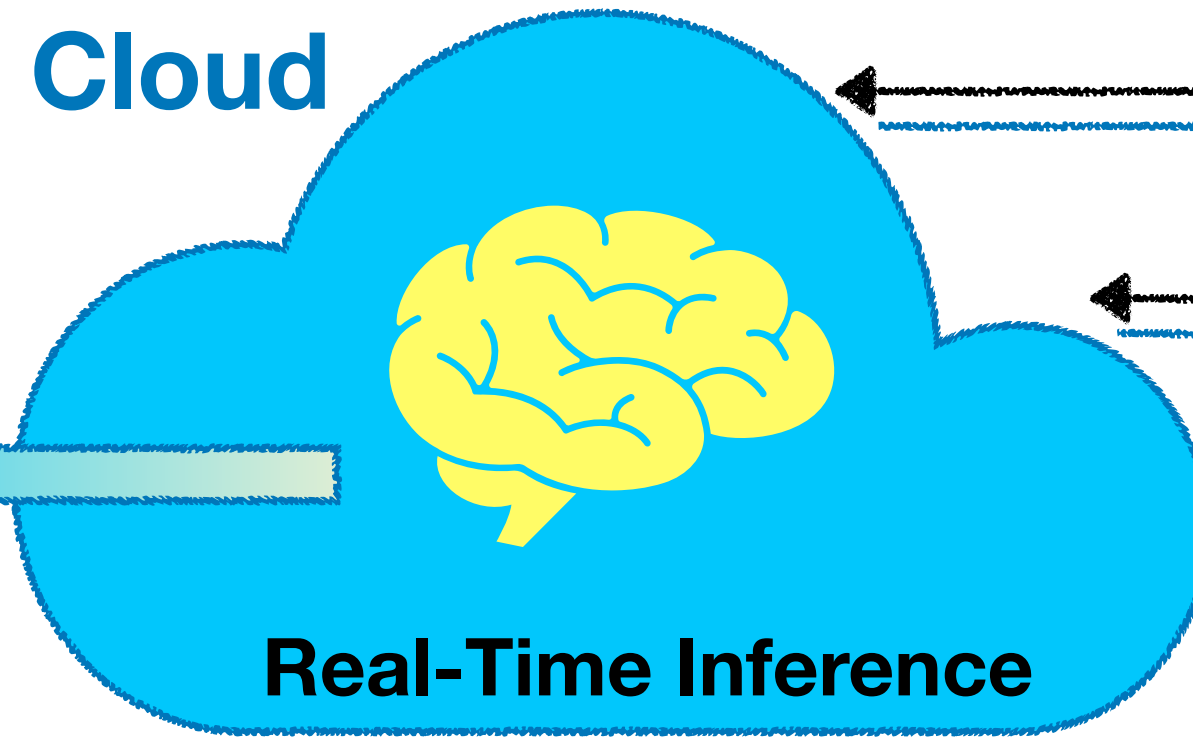**End-to-end predictable DNN serving platform for the Cloud**

✓ Supports 1000s of models concurrently per GPU

✓ Mitigates tail latency, supporting tight latency SLOs (10—100 ms)

✓ Close to ideal goodput under overload, contention, and bursts

# Background

# Inference Serving at the Cloud Scale is Difficult

# Inference Serving at the Cloud Scale is Difficult

**1000s of trained models of different types and resource requirements**

# Inference Serving at the Cloud Scale is Difficult



**1000s of trained models of different types and resource requirements**



**Requests arrive at different rates and regularity**

# Inference Serving at the Cloud Scale is Difficult

### 1000s of trained models of different types and resource requirements



### Requests arrive at different rates and regularity

# Inference Serving at the Cloud Scale is Difficult



**1000s of trained models of different types and resource requirements**

**Requests arrive at different rates and regularity**

# Inference Serving at the Cloud Scale is Difficult

**1000s of trained models of different types and resource requirements**



**Requests arrive at different rates and regularity**



**Each request has an inherent deadline**

**Latency SLOs**
**(e.g., 100ms)**

# Inference Serving at the Cloud Scale is Difficult

**1000s of trained models of different types and resource requirements**



**Requests arrive at different rates and regularity**



**Each request has an inherent deadline**

### Latency SLOs
**(e.g., 100ms)**

**HW accelerators are necessary!**

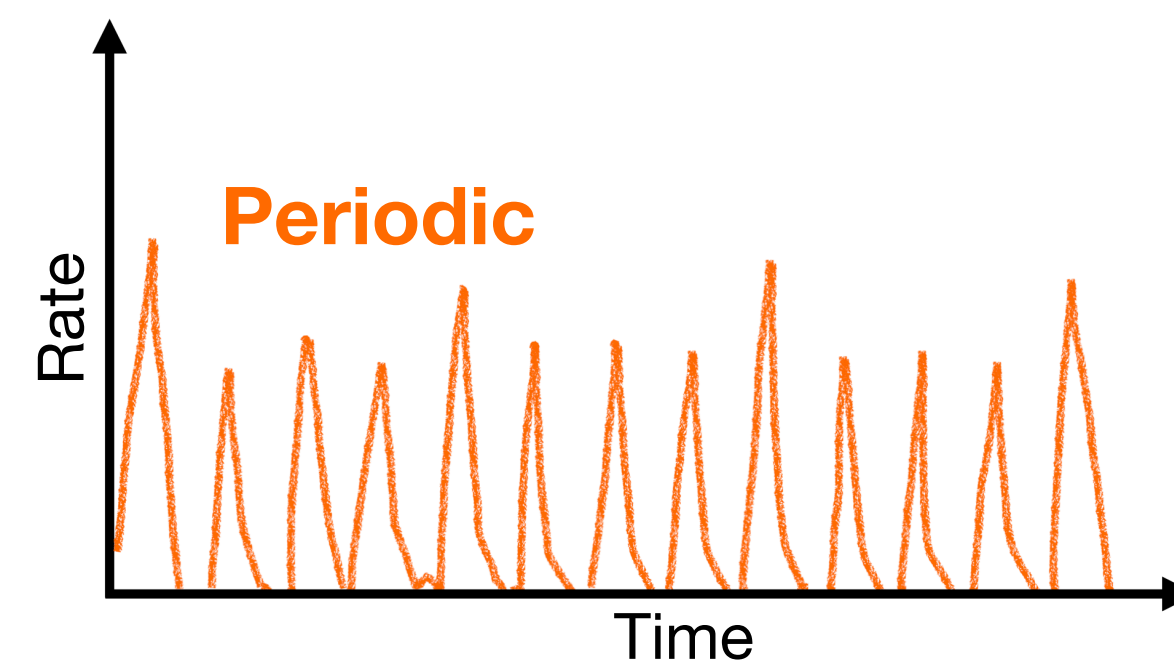| ResNet-50 | Latency | Throughput |
|-----------|---------|------------|
| CPU | 175 ms | 6 req/s |
| GPU | 2.8 ms | 350 req/s |

# Inference Serving at the Cloud Scale is Difficult

**1000s of trained models of different types and resource requirements**



**Requests arrive at different rates and regularity**



**Each request has an inherent deadline**

**Latency SLOs**
**(e.g., 100ms)**

**HW accelerators are necessary!**

| ResNet-50 | Latency | Throughput | Cost |
|-----------|---------|------------|------|
| CPU | 175 ms | 6 req/s | $ |
| GPU | 2.8 ms | 350 req/s | $$$ |

# Inference Serving at the Cloud Scale is Difficult

**1000s of trained models of different types and resource requirements**



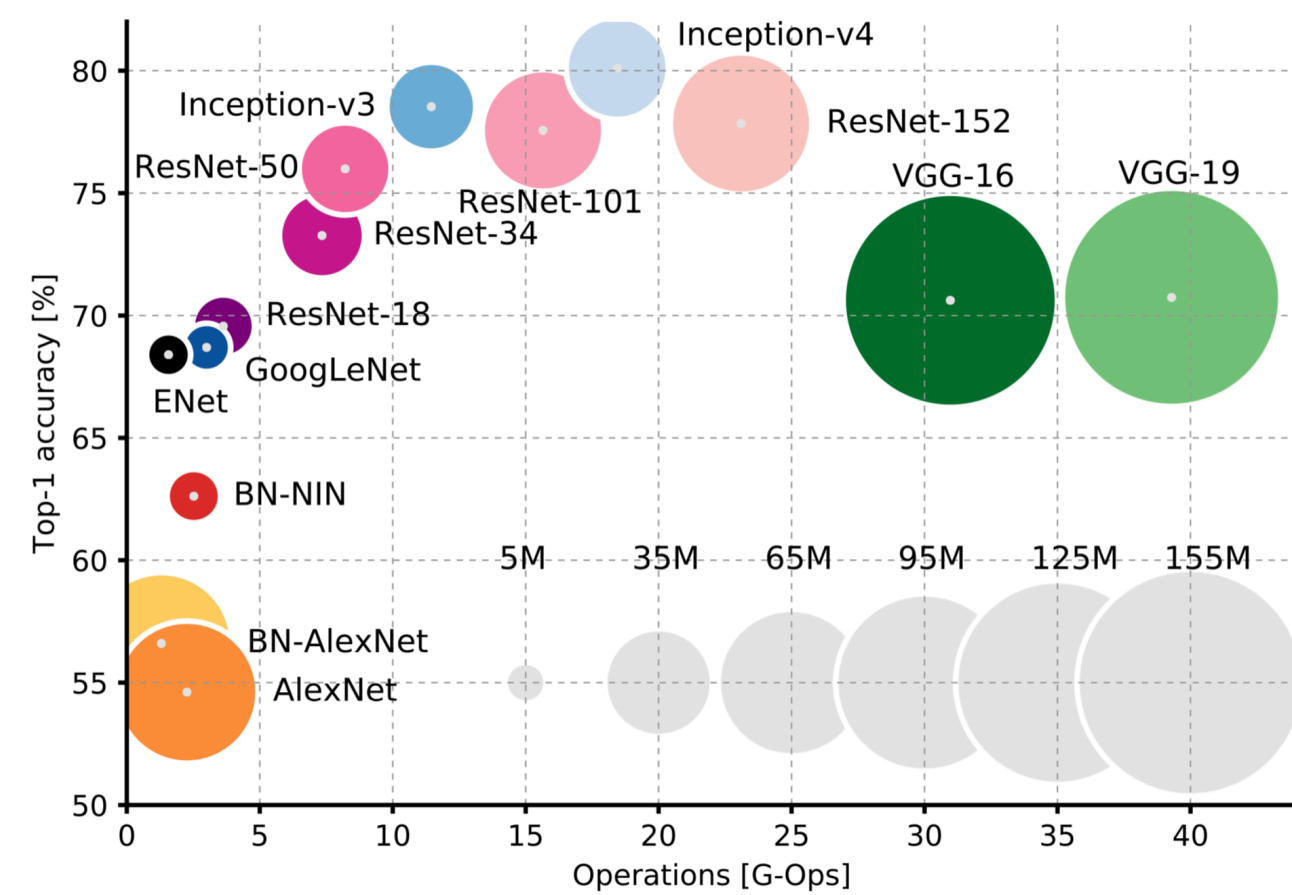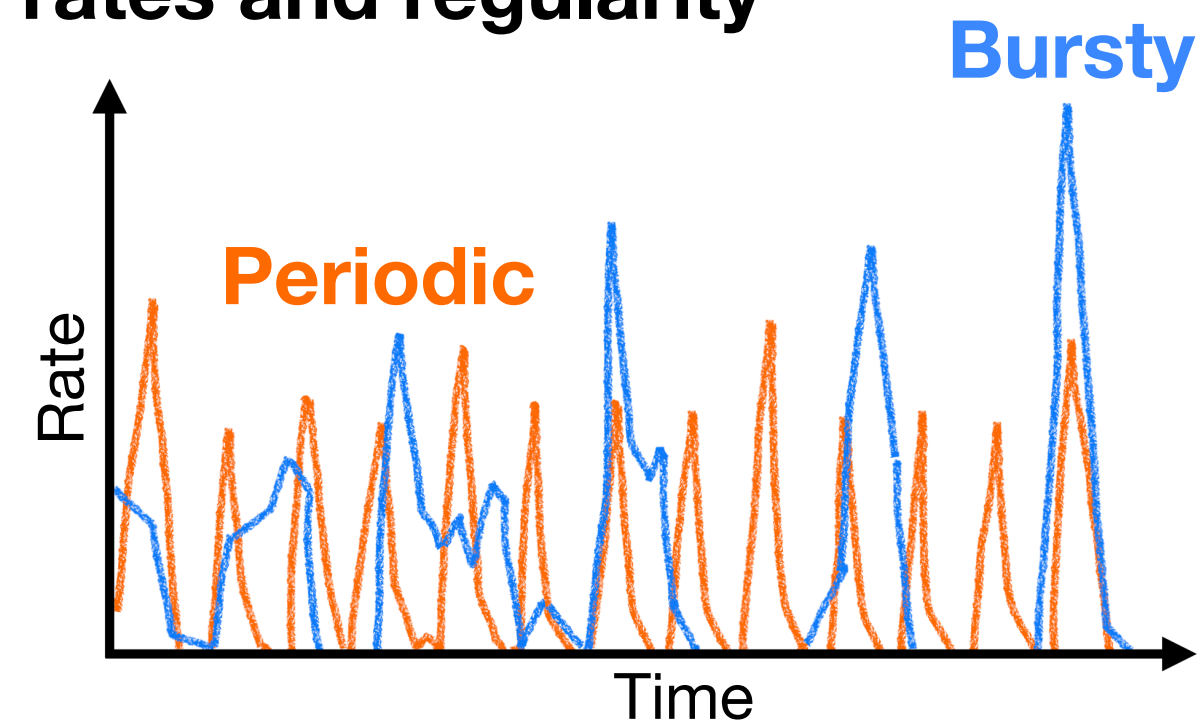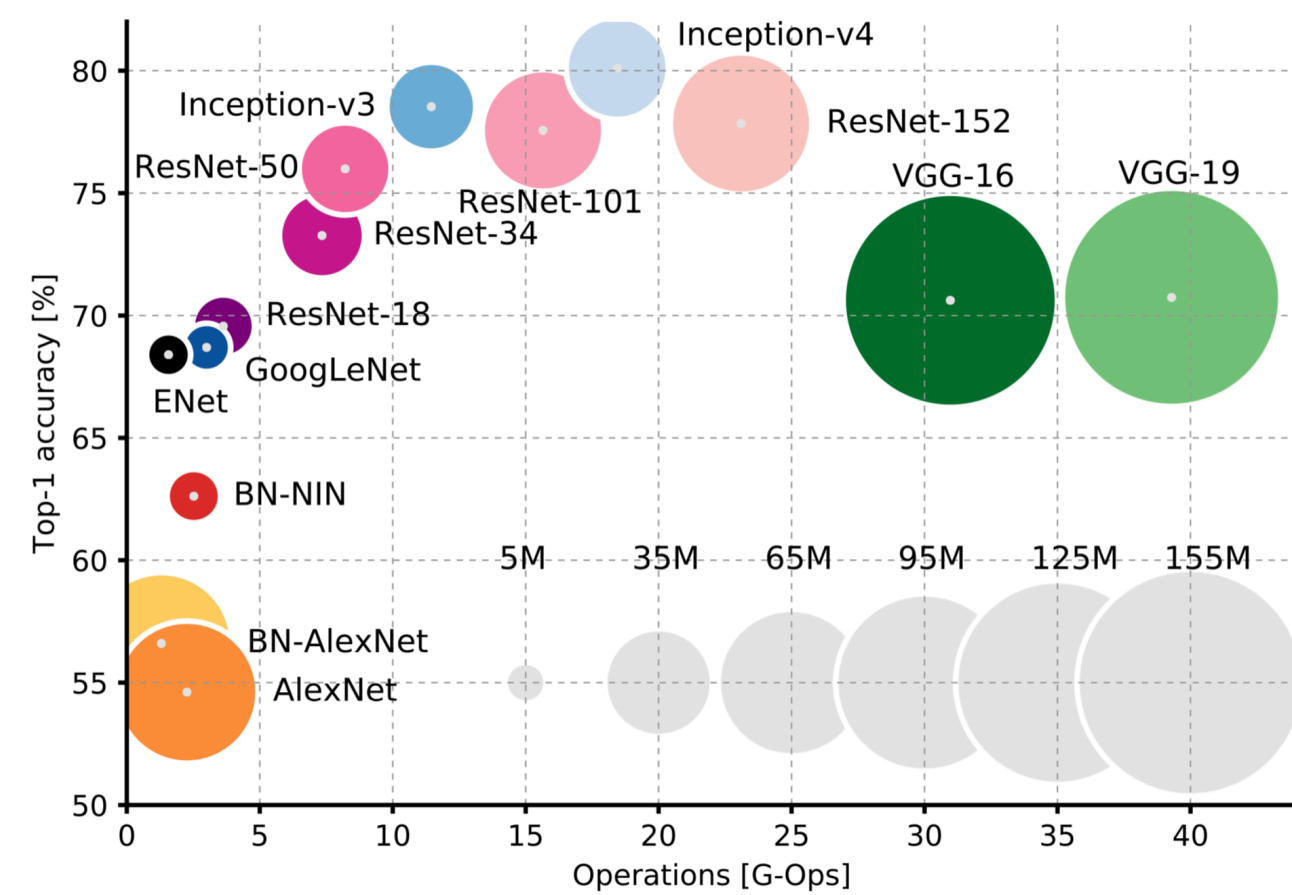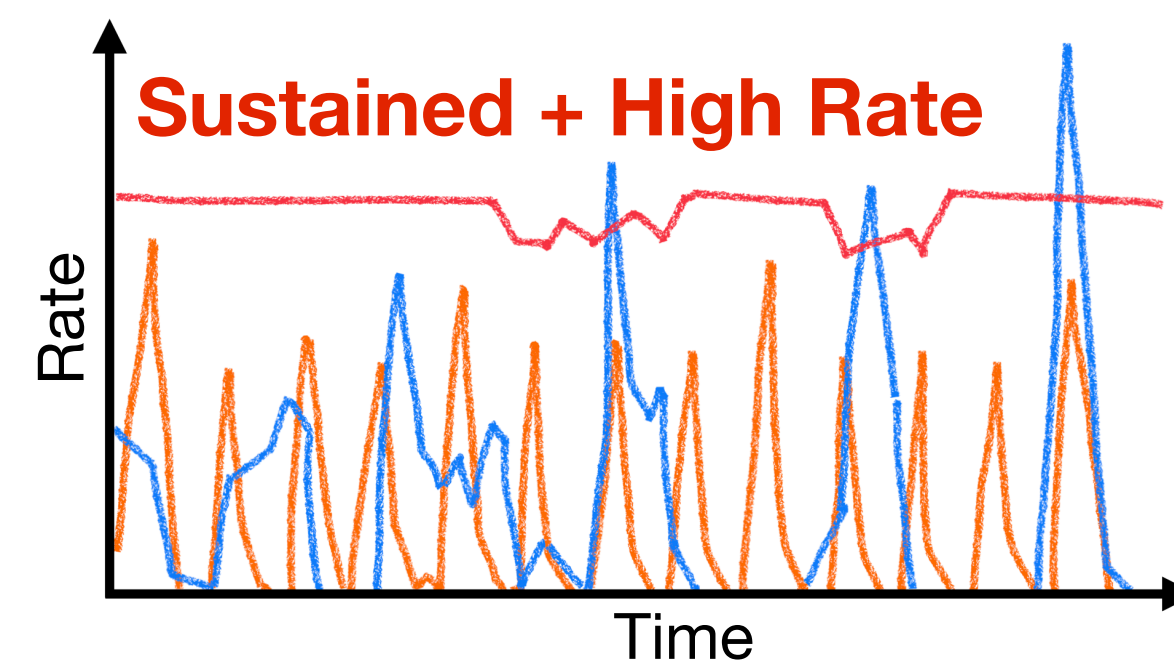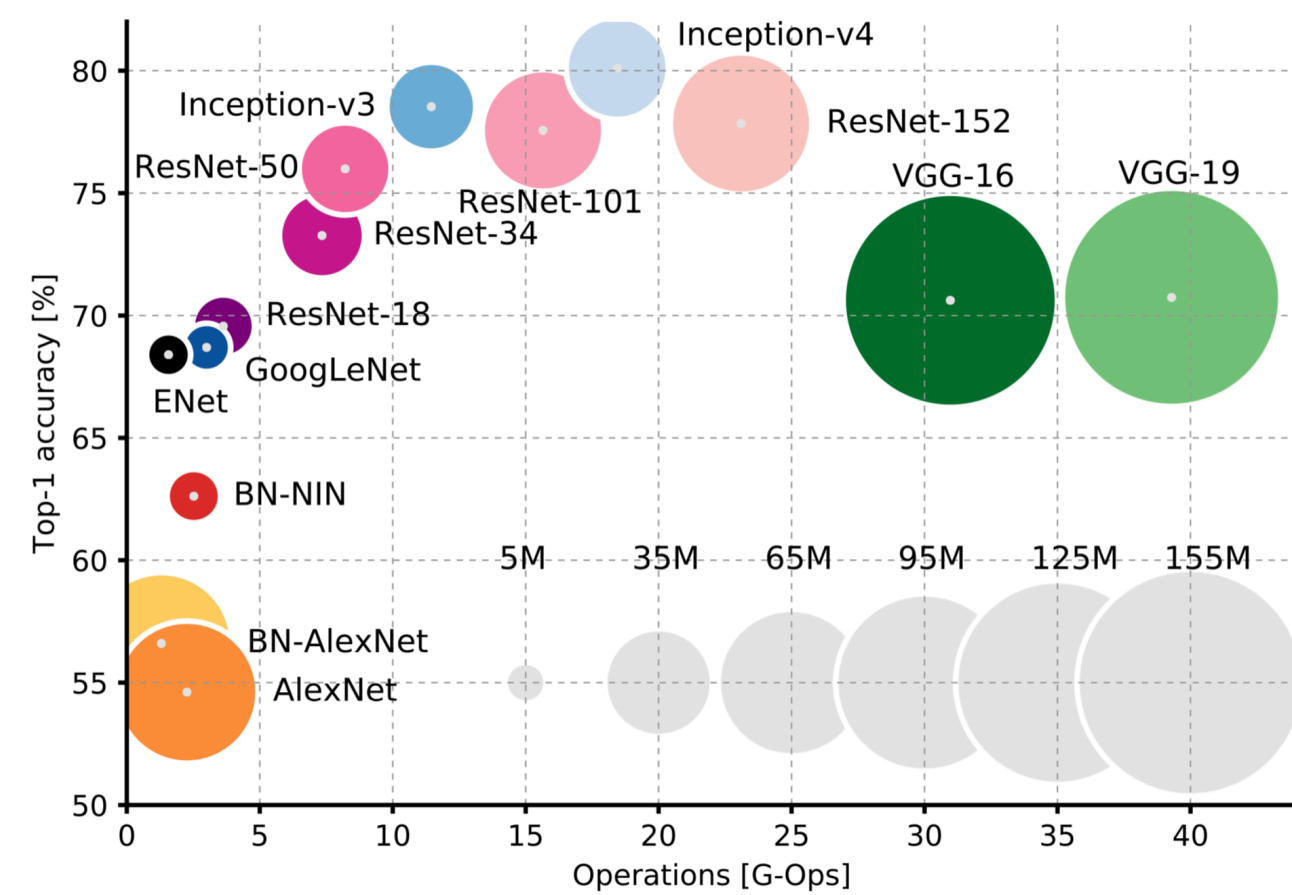**Requests arrive at different rates and regularity**



**Each request has an inherent deadline**

## Latency SLOs
**(e.g., 100ms)**

**HW accelerators are necessary!**



| ResNet-50 | Latency | Throughput | Cost |
|-----------|---------|------------|------|
| CPU | 175 ms | 6 req/s | $ |
| GPU | 2.8 ms | 350 req/s | $$$ |

## Problem

**How can cloud providers efficiently share resources while meeting SLOs?**

# Existing Systems Incur Very High Tail Latency

# Existing Systems Incur Very High Tail Latency

**Inference latency**

- **15 trained ResNet50**

- **Single GPU worker**

- **16 concurrent requests per model**

Clipper
100ms SLO

CDF

Percentile

99.9999
99.999
99.99
99.9
99
90
0

0  100            500

Latency (ms)

# Existing Systems Incur Very High Tail Latency

**Inference latency**

- **15 trained ResNet50**
- **Single GPU worker**
- **16 concurrent requests per model**

**Tail latency >> SLO**

Clipper
100ms SLO

# Existing Systems Incur Very High Tail Latency

**Inference latency**

- **15 trained ResNet50**

- **Single GPU worker**

- **16 concurrent requests per model**

**Tail latency >> SLO**

Clipper
100ms SLO



**200 ms**

**Concurrent DNN inference over GPU**

# Existing Systems Incur Very High Tail Latency

**Inference latency**

- **15 trained ResNet50**
- **Single GPU worker**
- **16 concurrent requests per model**

**Tail latency >> SLO**



Clipper
100ms SLO

Percentile: 99.9999, 99.999, 99.99, 99.9, 99, 90, 0

**200 ms**

Latency (ms): 0   100   500



Latency (ms): 100, 10, 1

**100x**

Concurrency: 1 2 4 8 16

**Concurrent DNN inference over GPU**

**High variance in latency**

**Throughput gains only 25%**

5

# Existing Systems Incur Very High Tail Latency

**Inference latency**

- **15 trained ResNet50**
- **Single GPU worker**
- **16 concurrent requests per model**

**Clipper 100ms SLO**

**200 ms**

Percentile vs Latency (ms) — 0, 90, 99, 99.9, 99.99, 99.999, 99.9999; Latency axis: 0, 100, 500

**Tail latency >> SLO**

**100x**

Latency (ms): 1, 10, 100 vs Concurrency: 1 2 4 8 16

**Concurrent DNN inference over GPU**

**High variance in latency**

**Single-thread latency is extremely predictable**

**Throughput gains only 25%**

# Existing Systems Incur Very High Tail Latency

**Inference latency**

- **15 trained ResNet50**
- **Single GPU worker**
- **16 concurrent requests per model**

Clipper
100ms SLO

**200 ms**

Percentile: 99.9999, 99.999, 99.99, 99.9, 99, 90, 0

Latency (ms): 0 100 500

**Tail latency >> SLO**

Preserves DNN predictability at every stage of model serving

**Clockwork adopts a contrasting approach!**

**100x**

Latency (ms): 100, 10, 1

Concurrency: 1 2 4 8 16

**Single-thread latency is extremely predictable**

**Concurrent DNN inference over GPU**

**High variance in latency**
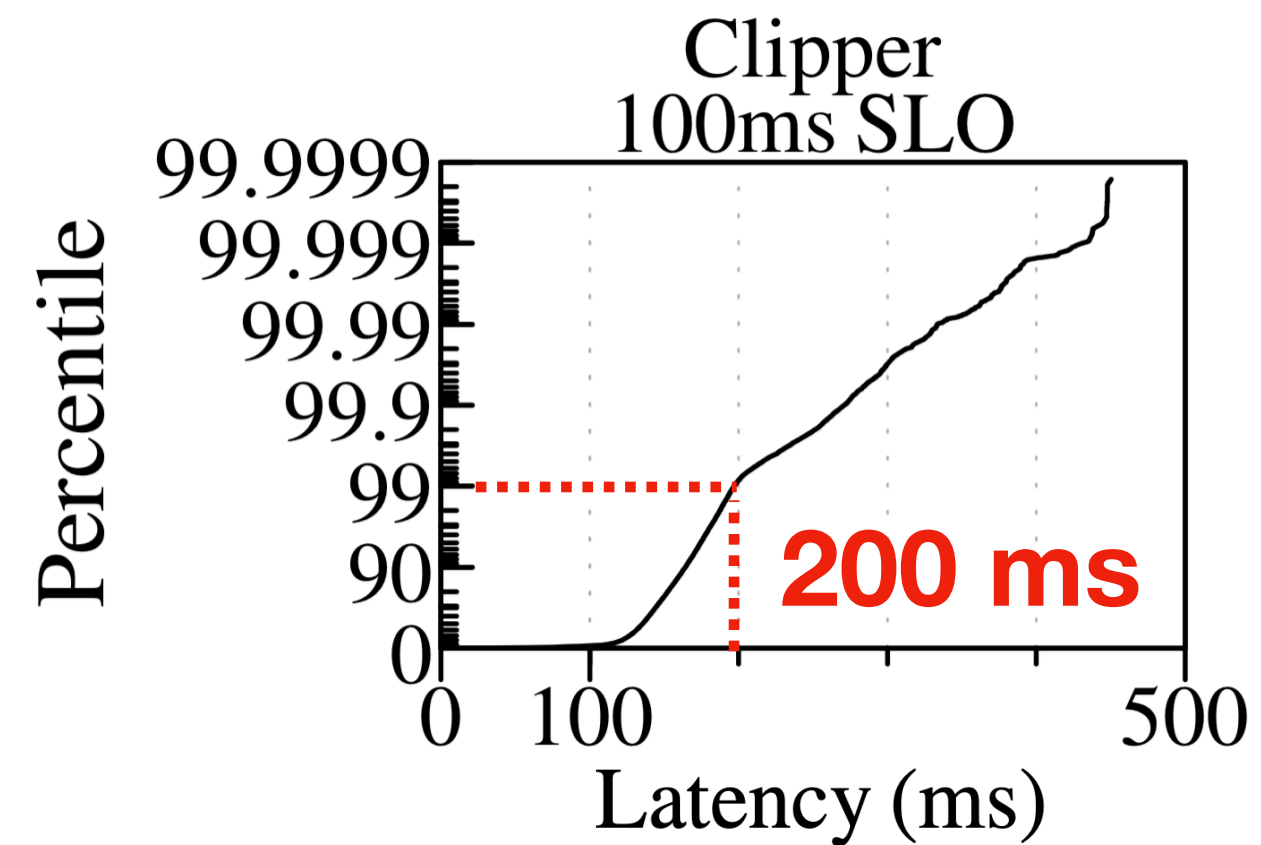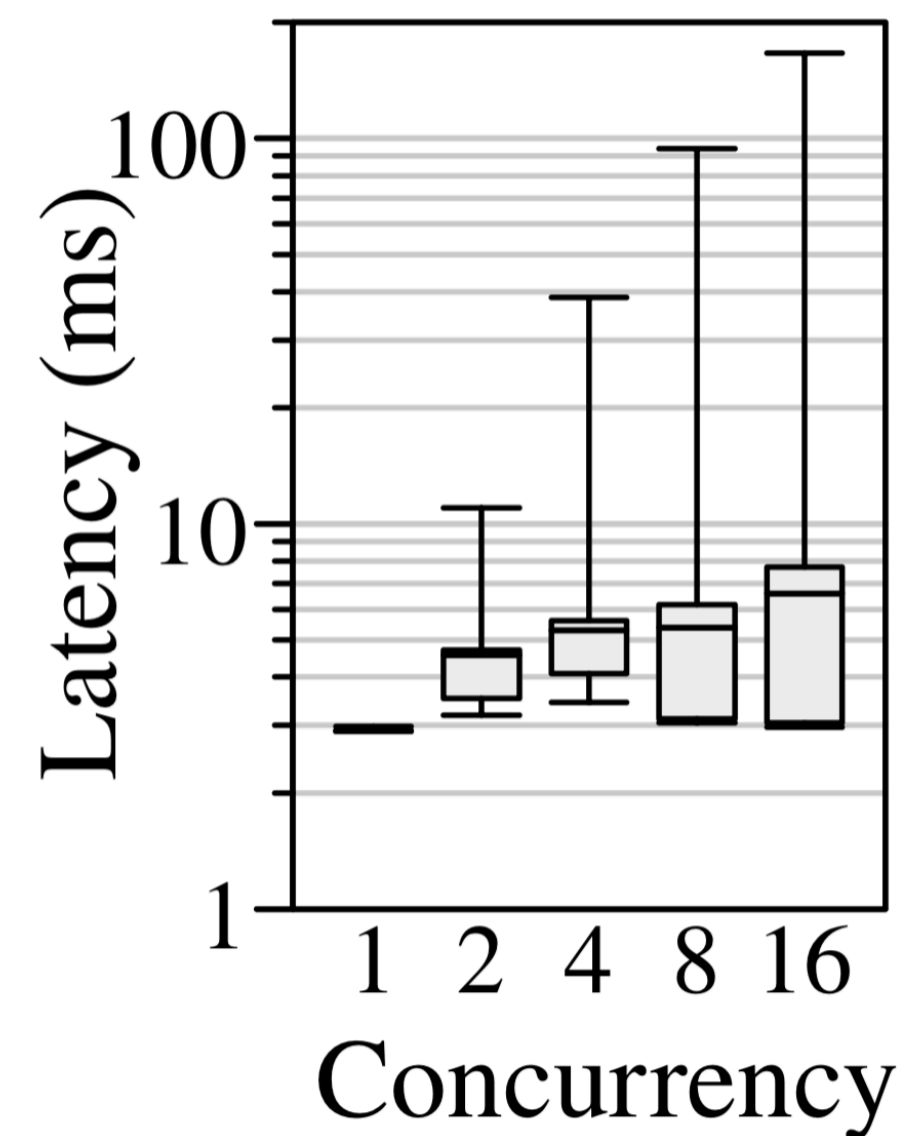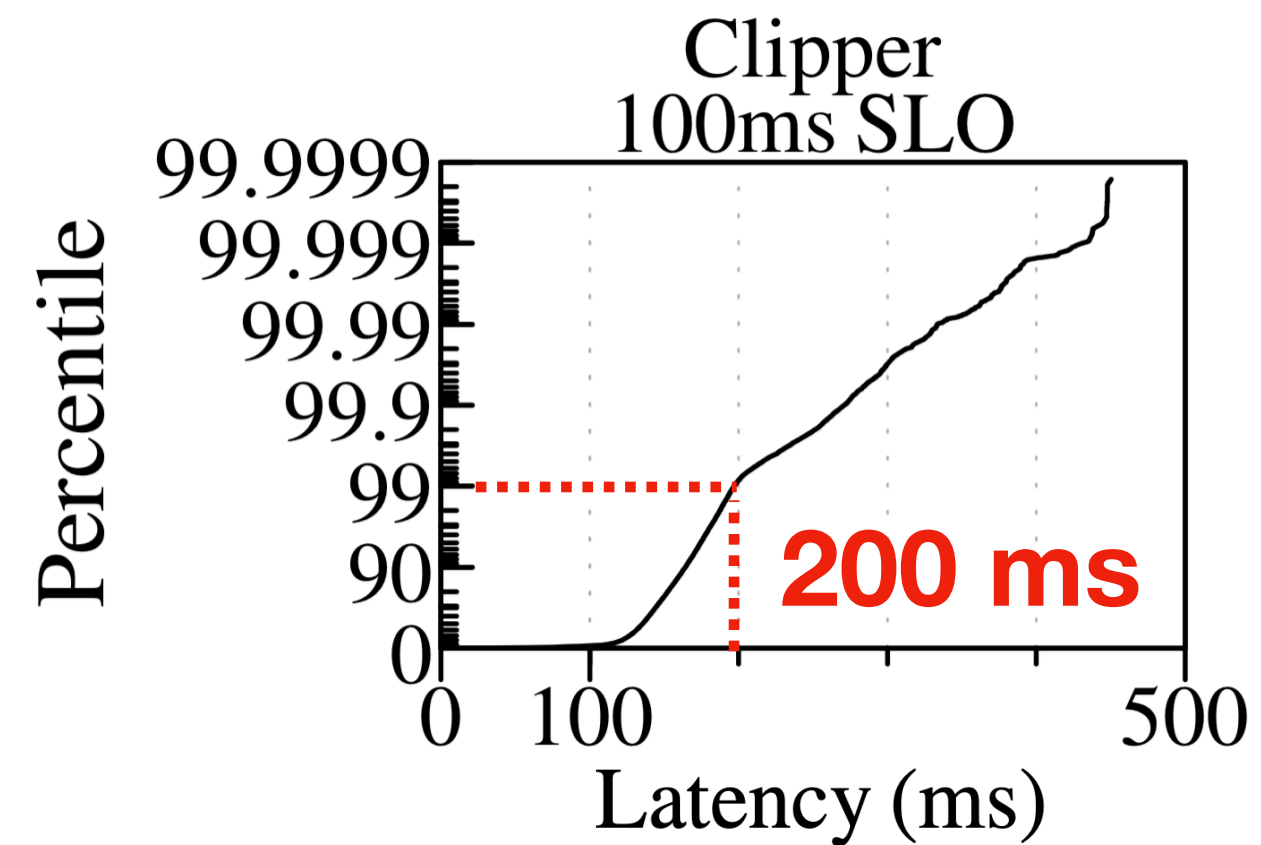
**Throughput gains only 25%**

# Existing Systems Incur Very High Tail Latency

**Inference latency**

- **15 trained ResNet50**
- **Single GPU worker**
- **16 concurrent requests per model**

Clipper 100ms SLO

**200 ms**

Percentile: 99.9999, 99.999, 99.99, 99.9, 99, 90, 0

Latency (ms): 0, 100, 500

**Tail latency >> SLO**

**Preserves DNN predictability at every stage of model serving**

Clockwork 100ms SLO

**95 ms**

Percentile: 99.9999, 99.999, 99.99, 99.9, 99, 90, 0

Latency (ms): 94.5, 95, 95.5

**Tail latency within SLO**

**Clockwork adopts a contrasting approach!**

**100x**

Latency (ms): 100, 10, 1

Concurrency: 1 2 4 8 16

**Concurrent DNN inference over GPU**

**High variance in latency**

**Single-thread latency is extremely predictable**

**Throughput gains only 25%**

# How does Clockwork Achieve End-to-End Predictability?

# Design Principles

Goal: 1000s of models, many users, limited resources

# Design Principles

Goal: 1000s of models, many users, limited resources

Maximize sharing

# Design Principles

Goal: 1000s of models, many users, limited resources

1. Predictable worker with no choices

Maximize sharing

# Design Principles

Goal: 1000s of models, many users, limited resources

Maximize sharing

1. Predictable worker with no choices

2. Consolidating choices at a central controller

# Design Principles

Goal: 1000s of models, many users, limited resources

Maximize sharing

1. Predictable worker with no choices

2. Consolidating choices at a central controller

3. Deadline-aware scheduling for SLO compliance

# Designing a Predictable Worker (1/2)



Worker Node

# Designing a Predictable Worker (1/2)

**Users upload pre-trained models in advance:** ● ▲ ■ ⬟ ★ ◆ **...**



**Worker Node**

# Designing a Predictable Worker (1/2)



Users upload pre-trained models in advance: ●△■⬠★◆ ...

Inference request for ★

Cold

Allocate memory for ★ ...

RAM

4 TB

GPU Memory

GPU Exec

GPU

32 GB

Worker Node

**Users upload pre-trained models in advance:** ●▲■⬟★◆ **...**

**Inference request for** ★

**Cold**

**Allocate memory for** ★ **...**

**Execute inference**

**RAM**

**GPU Memory**

**GPU Exec**

**GPU**

**4 TB**

**32 GB**

**Worker Node**

**Users upload pre-trained models in advance:** ●△■⬟★◆ ...

**Inference request for** ★

**Cold**

**Inference request for** ★
**(execute, since already in GPU memory)**

**Warm**

**RAM**

**4 TB**

**Allocate memory for** ★ ...

**GPU Memory**

**32 GB**

**Execute inference**

**GPU Exec**

**GPU**

**Worker Node**

**Users upload pre-trained models in advance:** ●△■●★◆ **...**

**Inference request for** ★

**Cold**

**Allocate memory for** ★ **...**

**Execute inference**

**Inference request for** ★ **(execute, since already in GPU memory)**

**Warm**

**RAM**

**GPU Memory**

**GPU Exec**

**GPU**

**Worker Node**

**4 TB**

**32 GB**

**Users upload pre-trained models in advance:** ● △ ■ ⬠ ★ ◆ **...**

**Inference request for ★**

**Cold**

**Allocate memory for ★ ...**

**Queues**

**Execute inference**

**Inference request for ★ (execute, since already in GPU memory)**

**Warm**

**RAM**

**4 TB**

**GPU Memory**

**32 GB**

**GPU Exec**

**GPU**

**Worker Node**

# Designing a Predictable Worker (1/2)



**Users upload pre-trained models in advance:** ● ▲ ■ ⬠ ★ ◆ ...

**Inference request for** ★

**Cold**

**Queues**

**Allocate memory for** ★ ...

**Inference request for** ★
**(execute, since already in GPU memory)**

**Warm**

**RAM**

**4 TB**

**GPU Memory**

**Execute inference**

**GPU Exec**

**GPU**

**32 GB**

**Worker Node**

**Concurrent inferences**

**+** **Proprietary & undocumented policies**

**➜** **Unpredictable response times**

100x

Latency (ms)

Concurrency

# Designing a Predictable Worker (1/2)



Users upload pre-trained models in advance: ●△■⬠★◆ ...

**Queues**

Inference request for ★

**Cold**

Allocate memory for ★ ...

Execute inference

Inference request for ★ (execute, since already in GPU memory)

**Warm**

**RAM**

**4 TB**

**GPU Memory**

**32 GB**

**GPU Exec**

**GPU**

**Worker Node**

**Managed memory can be unpredictable**

- GPU memory (cache) hits & misses

**ResNet-50 — Hit: 2.3 ms | Miss: 10.6 ms**

**Concurrent inferences**

**+** Proprietary & undocumented policies

**➡** Unpredictable response times

100x

Latency (ms)

100

10

1

1  2  4  8  16
Concurrency

8

# Designing a Predictable Worker (2/2)

**Predictable Clockwork worker process**



RAM

GPU Memory

GPU Exec

GPU

**Worker Node**

**Predictable Clockwork worker process**

RAM

GPU Memory

GPU Exec

GPU

**Worker Node**

**Concurrent inferences**

+ Proprietary & undocumented policies

→ Unpredictable response times

**Solution**

Execute inference one at a time

# Designing a Predictable Worker (2/2)

**Predictable Clockwork worker process**

RAM

GPU Memory

**PageCache**

GPU

GPU Exec

**Worker Node**

**Managed memory can be unpredictable**

**Solution**

Preallocate GPU memory & manage it explicitly using LOAD/UNLOAD actions

**Concurrent inferences**

**+** Proprietary & undocumented policies

**➡** Unpredictable response times

**Solution**

Execute inference one at a time

9

# Designing a Predictable Worker (2/2)



**Predictable Clockwork worker process**

Earliest Deadline First

PCI

Time

GPU

Time

RAM

GPU Memory

**PageCache**

GPU

GPU Exec

GPU

**Worker Node**

**Managed memory can be unpredictable**

Solution

Preallocate GPU memory & manage it explicitly using LOAD/UNLOAD actions

**Concurrent inferences**

+ Proprietary & undocumented policies

→ Unpredictable response times

Solution

Execute inference one at a time

**Choices outsourced via action APIs**

**Predictable Clockwork worker process**

Earliest Deadline First

LOAD/UNLOAD (◆, Deadline)

INFER (★, I/P, Deadline)

PCI
Time

GPU
Time

**RAM**

**GPU Memory**

**PageCache**

**GPU Exec**

**GPU**

**Worker Node**

**Managed memory can be unpredictable**

**Solution**

**Preallocate GPU memory & manage it explicitly using LOAD/UNLOAD actions**

**Concurrent inferences**

**+** Proprietary & undocumented policies

**→** Unpredictable response times

**Solution**

Execute inference one at a time

9

# Consolidating Choices

# Consolidating Choices

# Consolidating Choices



**Users**

**Centralized Controller**

**Worker processes**

RAM

GPU Memory
**PageCache**

GPU Exec

**GPU Worker Node W₁**

LOADs

INFERs

**Smarter load balancing & scheduling decisions**

**Global State Manager**

**Latency Profiles**

**Pending Tasks**

**Memory State**

# SLO-aware Scheduling

# SLO-aware Scheduling

# SLO-aware Scheduling

# SLO-aware Scheduling

# SLO-aware Scheduling

# SLO-aware Scheduling

# SLO-aware Scheduling

# SLO-aware Scheduling

**Users**

**Centralized Controller**

**Worker processes**

LOADs

INFERs

RAM

GPU Memory

GPU Exec

**GPU Worker Node W$_1$**

GPU

RAM
GPU Memory
PageCache
GPU Exec
GPU

RAM
GPU Memory
PageCache
GPU Exec
GPU

**Many benefits**

- **Prevent wasteful work**

- **Manage LOAD → INFER dependencies**

- **Choosing the best batching strategy**

# Evaluation

# Questions

# Questions

**How does Clockwork compare to prior model serving systems Clipper and INFaaS?**

# Questions

**How does Clockwork compare to prior model serving systems Clipper and INFaaS?**

**Can Clockwork serve thousands of model instances?**

# Questions

**How does Clockwork compare to prior model serving systems Clipper and INFaaS?**

**Can Clockwork serve thousands of model instances?**

**How low can Clockwork go in terms of the latency SLOs it can satisfy?**

# Questions

**How does Clockwork compare to prior model serving systems Clipper and INFaaS?**

**Can Clockwork serve thousands of model instances?**

**How low can Clockwork go in terms of the latency SLOs it can satisfy?**

**Can Clockwork isolate the performance of latency-sensitive clients from batch requests without latency SLOs?**

# Questions

How does Clockwork compare to prior model serving systems Clipper and INFaaS?

Can Clockwork serve thousands of model instances?

How low can Clockwork go in terms of the latency SLOs it can satisfy?

Can Clockwork isolate the performance of latency-sensitive clients
from batch requests without latency SLOs?

# Questions

**Simple workloads in controlled settings**

How does Clockwork compare to prior model serving systems Clipper and INFaaS?

Can Clockwork serve thousands of model instances?

How low can Clockwork go in terms of the latency SLOs it can satisfy?

Can Clockwork isolate the performance of latency-sensitive clients from batch requests without latency SLOs?

Are Clockwork workers predictable?

Does consolidating choice help achieve end-to-end predictability?

Can Clockwork controller Scale?

**Workloads from production traces**

# Questions

How does Clockwork compare to prior model serving systems Clipper and INFaaS?

Can Clockwork serve thousands of model instances?

How low can Clockwork                     of the latency SLOs it can satisfy?

**This talk**

Can Clockwork isolate the performance of latency-sensitive clients from batch requests without latency SLOs?

**Are Clockwork workers predictable?**

**Does consolidating choice help achieve end-to-end predictability?**

**Can Clockwork controller Scale?**

**Workloads from production traces**

# Experiment Setup

**12 Workers: NVIDIA Tesla v100 GPU | 32 GB GPU Memory** **+** **1 Controller** **+** **1 Client**

# Experiment Setup

| 12 Workers: NVIDIA Tesla v100 GPU | 32 GB GPU Memory | **+** | 1 Controller | **+** | 1 Client |

Microsoft's Azure Functions ➡️

Shahrad et al. *"Serverless in the Wild: Characterizing and Optimizing the Serverless Workload at a Large Cloud Provider."* USENIX ATC 2020

46,000 functions, 2 weeks
- Heavy sustained workloads
- Low utilization cold workloads
- Workloads with periodic spikes
- Bursty workloads



# Workload

# Experiment Setup

| 12 Workers: NVIDIA Tesla v100 GPU | 32 GB GPU Memory | **+** | 1 Controller | **+** | 1 Client |

Microsoft's Azure Functions ➡️

Shahrad et al. *"Serverless in the Wild: Characterizing and Optimizing the Serverless Workload at a Large Cloud Provider."* USENIX ATC 2020

46,000 functions, 2 weeks
- Heavy sustained workloads
- Low utilization cold workloads
- Workloads with periodic spikes
- Bursty workloads

# Workload

4026 model instances
- Saturates 768 GB RAM
- 61 different model architectures
- ResNet, DenseNet, Inception, etc.



14

# Are Clockwork Workers Predictable?

# Are Clockwork Workers Predictable?

Clockwork relies on predicting the model
inference latency for scheduling

**Overpredictions** ⟶ **Idle resources**

**Underpredictions** ⟶ **SLO violations**

# Are Clockwork Workers Predictable?

Clockwork relies on predicting the model inference latency for scheduling

**Overpredictions** ⟶ **Idle resources**

**Underpredictions** ⟶ **SLO violations**

# Are Clockwork Workers Predictable?

Clockwork relies on predicting the model inference latency for scheduling

**Overpredictions** → **Idle resources**

**Underpredictions** → **SLO violations**



Clockwork consistently overpredicts more than its underpredicts

# Are Clockwork Workers Predictable?

Clockwork relies on predicting the model inference latency for scheduling

**Overpredictions** ⟶ **Idle resources**

**Underpredictions** ⟶ **SLO violations**

INFER

Clockwork consistently overpredicts more than its underpredicts

Errors are significant only in extremely rare cases

**Underprediction error = 55us**
**Overprediction error = 144us**

Overpredict
Underpredict

# Does Consolidating Choice Help?

Offered load ~10,000 r/s, periodic spikes ~12,000 r/s

Latency SLO = 100 ms deadline for each request

# Does Consolidating Choice Help?

**Goodput =**
**SLO compliant**
**throughput**

**Offered load ~10,000 r/s, periodic spikes ~12,000 r/s**

**Latency SLO = 100 ms deadline for each request**

# Does Consolidating Choice Help?

**Goodput =
SLO compliant
throughput**

**Latency of all
completed
requests**



Offered load ~10,000 r/s, periodic spikes ~12,000 r/s

Latency SLO = 100 ms deadline for each request

# Does Consolidating Choice Help?

**Goodput = SLO compliant throughput**

**Latency of all completed requests**

Time (Minutes)



**Offered load ~10,000 r/s, periodic spikes ~12,000 r/s**

**Latency SLO = 100 ms deadline for each request**

**The workload is successfully scheduled by Clockwork**

- Goodput ≈ offered load
- Out of 208 million requests, only 58 failed due to mispredictions
- All others completed within SLO

# Does Consolidating Choice Help?

**Goodput =**
**SLO compliant**
**throughput**

**Latency of all**
**completed**
**requests**
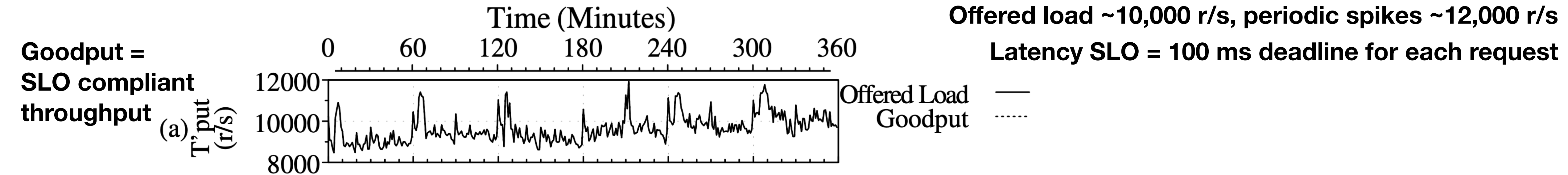


**Offered load ~10,000 r/s, periodic spikes ~12,000 r/s**

**Latency SLO = 100 ms deadline for each request**

**The workload is successfully**
**scheduled by Clockwork**

- Goodput ≈ offered load

- Out of 208 million requests, only
  58 failed due to mispredictions

- All others completed within SLO

# Does Consolidating Choice Help?

Goodput =
SLO compliant
throughput

Latency of all
completed
requests

**Offered load ~10,000 r/s, periodic spikes ~12,000 r/s**
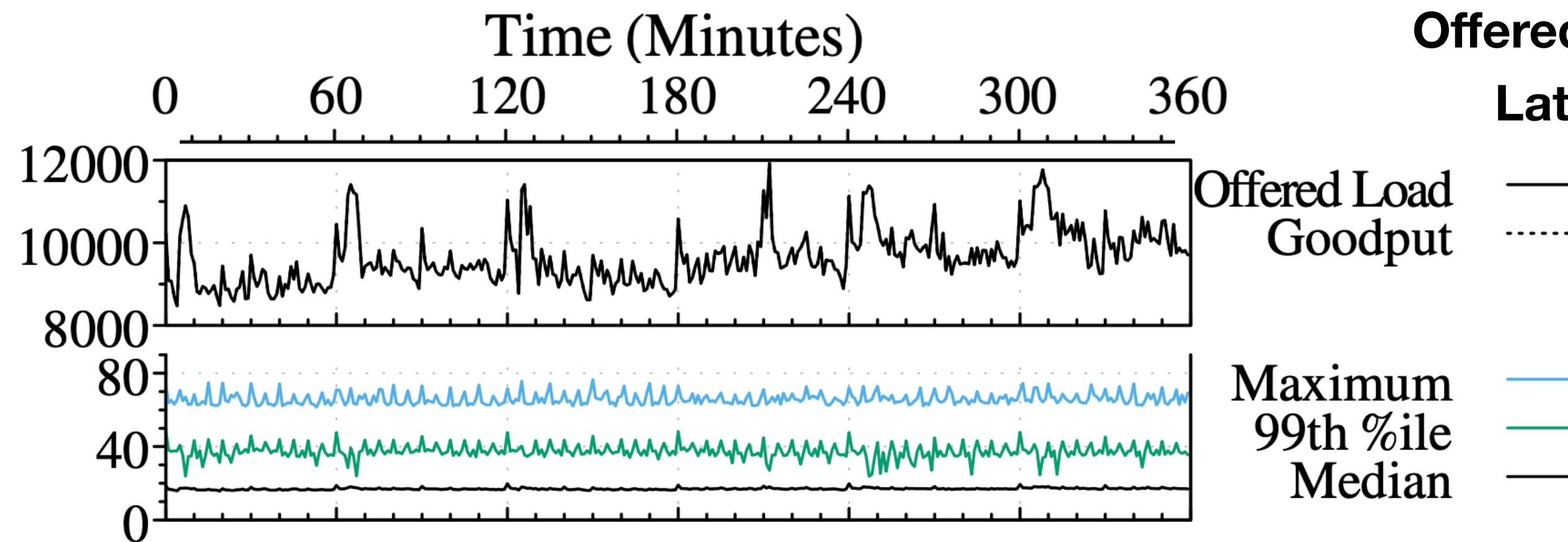
**Latency SLO = 100 ms deadline for each request**



**The workload is successfully scheduled by Clockwork**

- Goodput $\approx$ offered load
- Out of 208 million requests, only 58 failed due to mispredictions
- All others completed within SLO

# Does Consolidating Choice Help?

Goodput =
SLO compliant
throughput

Latency of all
completed
requests



Offered load ~10,000 r/s, periodic spikes ~12,000 r/s

Latency SLO = 100 ms deadline for each request
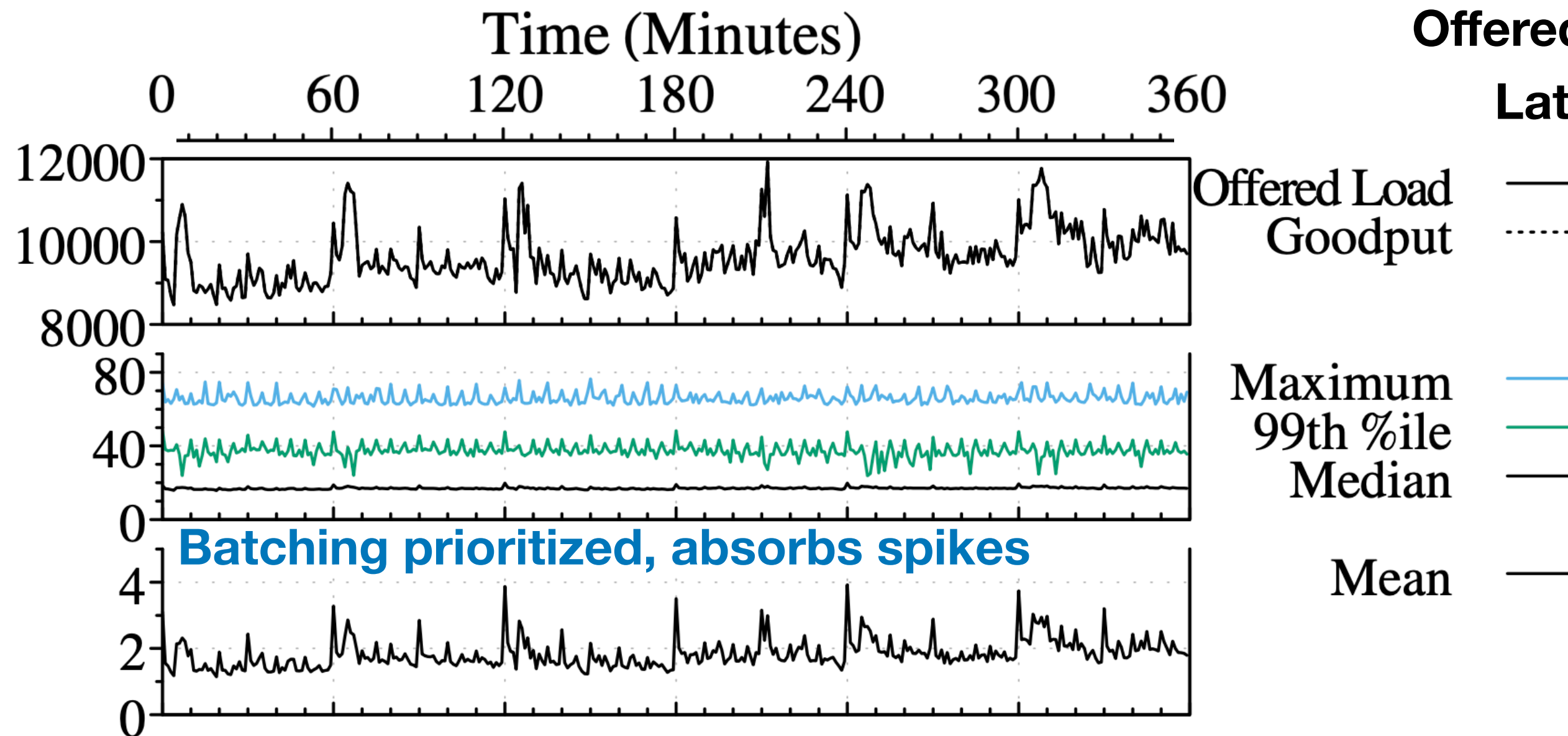
**The workload is successfully scheduled by Clockwork**

- Goodput ≈ offered load
- Out of 208 million requests, only 58 failed due to mispredictions
- All others completed within SLO

# Does Consolidating Choice Help?

**Goodput = SLO compliant throughput**

**Latency of all completed requests**

Time (Minutes)



**Offered load ~10,000 r/s, periodic spikes ~12,000 r/s**

**Latency SLO = 100 ms deadline for each request**

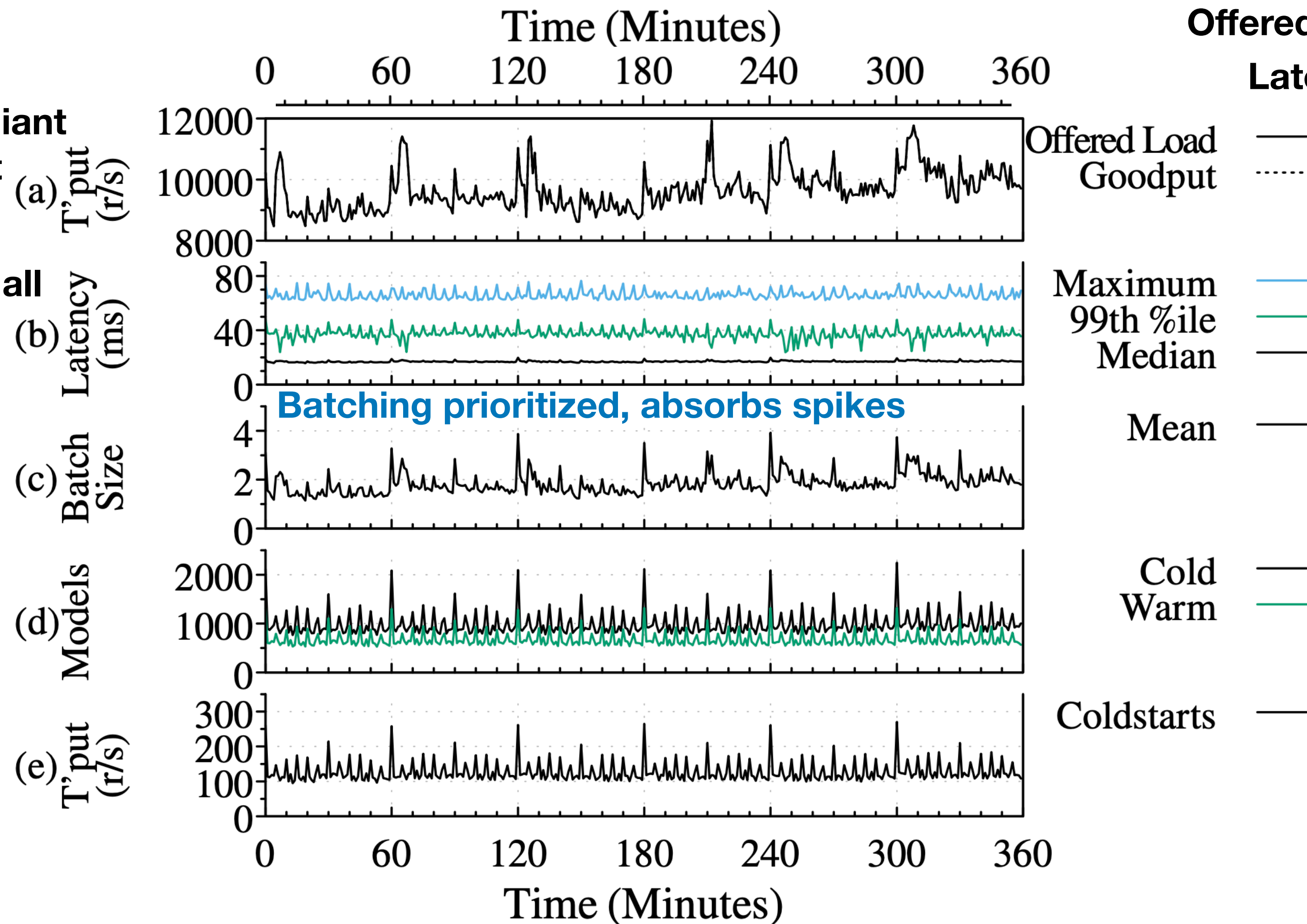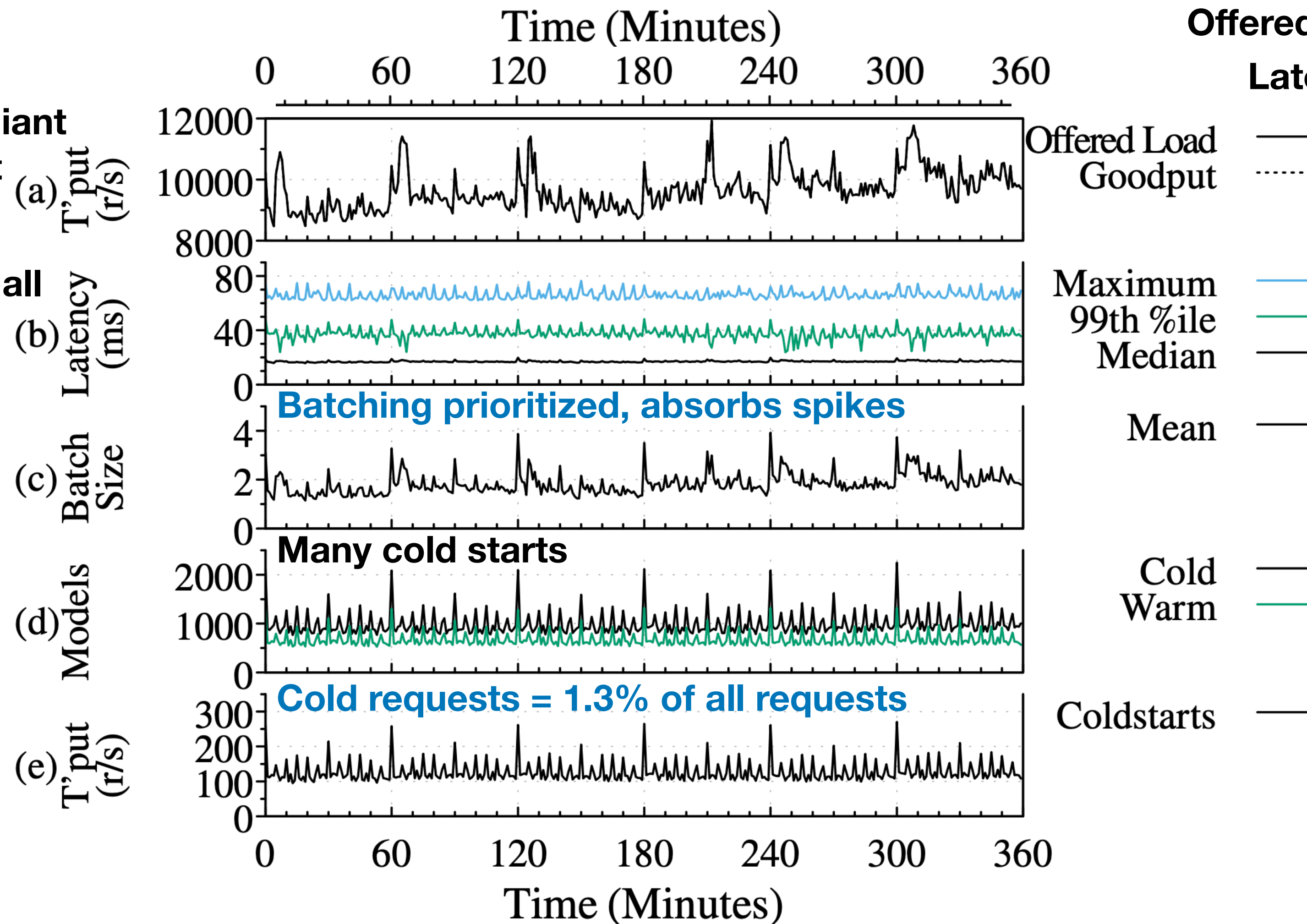**The workload is successfully scheduled by Clockwork**

- Goodput ≈ offered load

- Out of 208 million requests, only 58 failed due to mispredictions

- All others completed within SLO

# Does Clockwork Controller Scale?

# Does Clockwork Controller Scale?

**Methodology**

- **Replace GPU workers with emulated workers**

- **From the controller's vantage point, nothing changes**

- **Measure the peak goodput as we vary #workers**

# Does Clockwork Controller Scale?



**40 emulated workers**

**Methodology**

- Replace GPU workers with emulated workers
- From the controller's vantage point, nothing changes
- Measure the peak goodput as we vary #workers

# Does Clockwork Controller Scale?



**40 emulated workers**

Req/s vs Time (mins), with legend showing Offered Load and Goodput, and a "Peak" annotation.

**Methodology**

- Replace GPU workers with emulated workers

- From the controller's vantage point, nothing changes

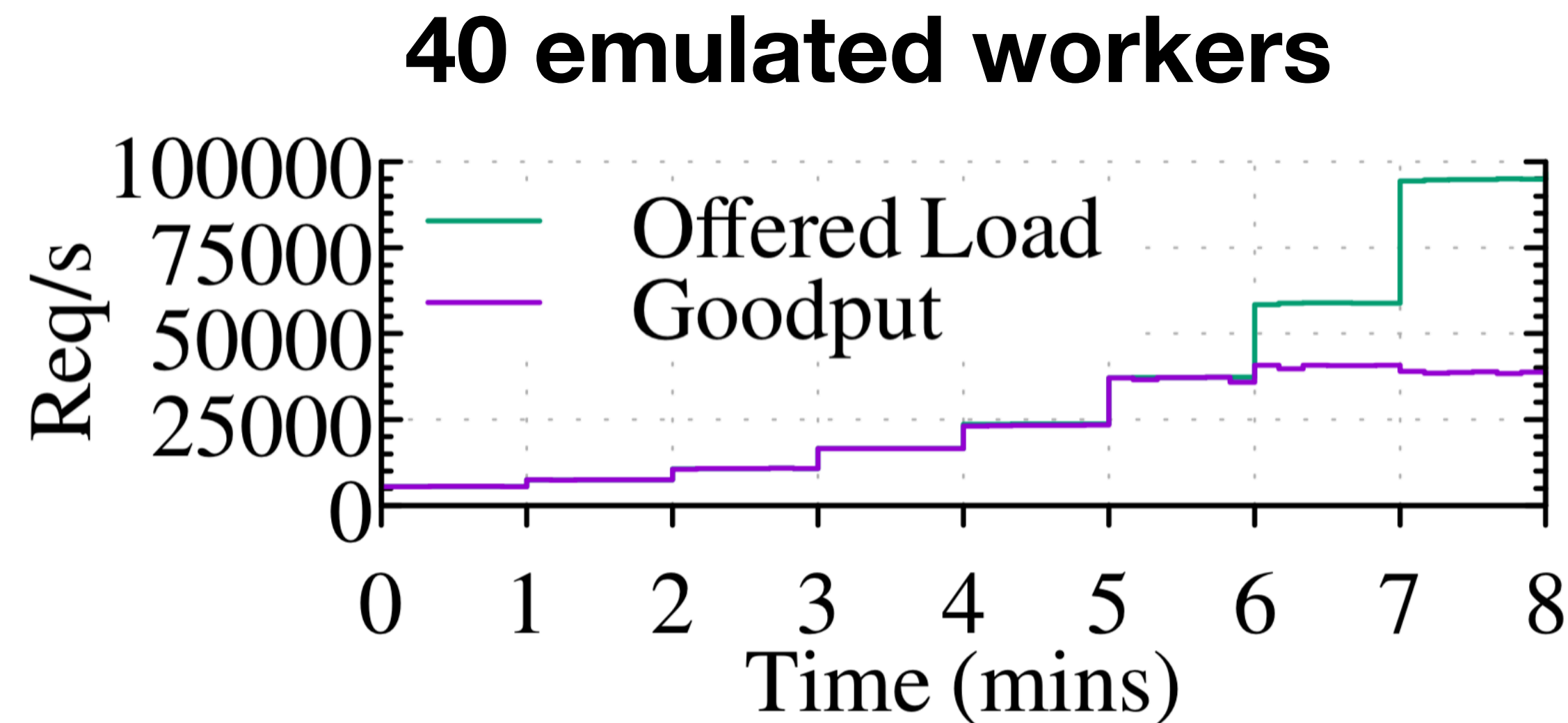- Measure the peak goodput as we vary #workers

# Does Clockwork Controller Scale?

**Methodology**

- **Replace GPU workers with emulated workers**

- **From the controller's vantage point, nothing changes**

- **Measure the peak goodput as we vary #workers**

# Does Clockwork Controller Scale?



**Linear scalability until #workers = 110**

Goodput limited by worker's utilization

**Methodology**

- Replace GPU workers with emulated workers
- From the controller's vantage point, nothing changes
- Measure the peak goodput as we vary #workers
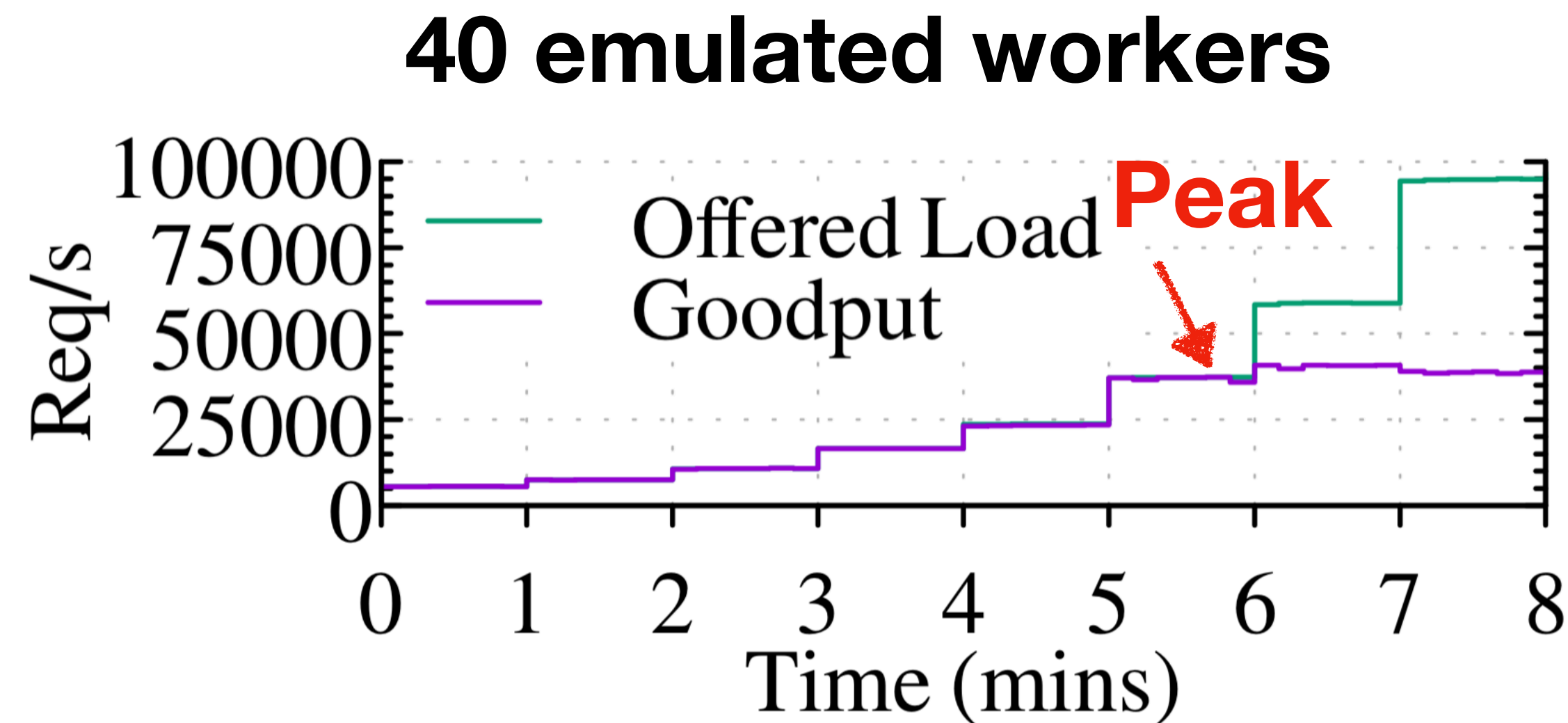
# Does Clockwork Controller Scale?



**Bottleneck shifts to Clockwork**

**Linear scalability until #workers = 110**

Goodput limited by worker's utilization

**Maximum goodput: 103,387 r/s for 110 workers**

*Y-axis: Peak Goodput (r/s)* — 0, 20000, 40000, 60000, 80000, 100000, 120000

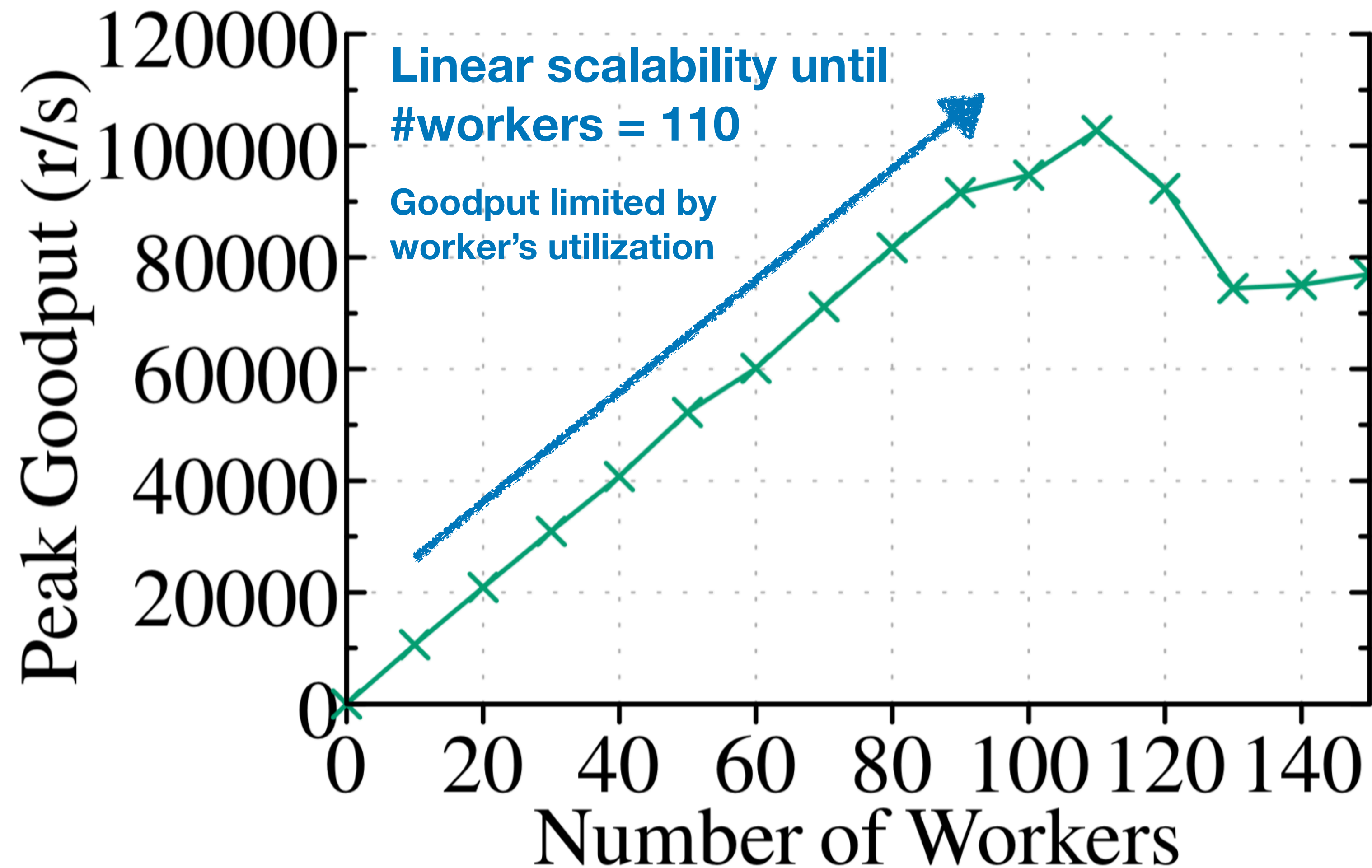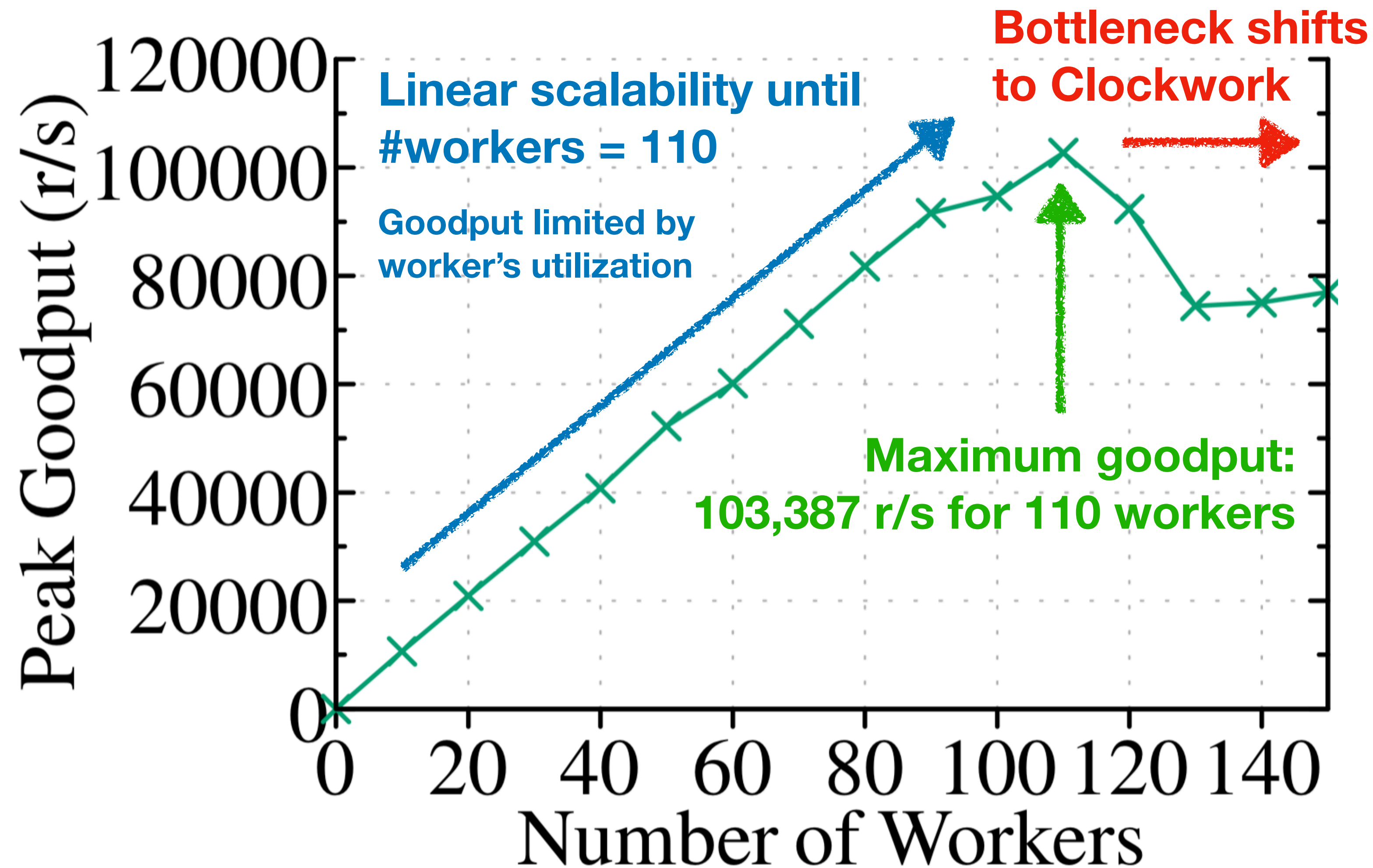*X-axis: Number of Workers* — 0, 20, 40, 60, 80, 100, 120, 140

**Methodology**

- Replace GPU workers with emulated workers
- From the controller's vantage point, nothing changes
- Measure the peak goodput as we vary #workers

# Summary

# Summary

**Key idea: DNN executions on GPUs exhibit negligible latency variability**

- Intuitive – DNN inferences involve no conditional branches – and demonstrable in practice

# Summary

**Key idea: DNN executions on GPUs exhibit negligible latency variability**

- Intuitive – DNN inferences involve no conditional branches – and demonstrable in practice

**Clockwork: From DNN predictability to an E2E predictable DNN serving platform**

- Recursively ensures that all internal architecture components have predictable performance
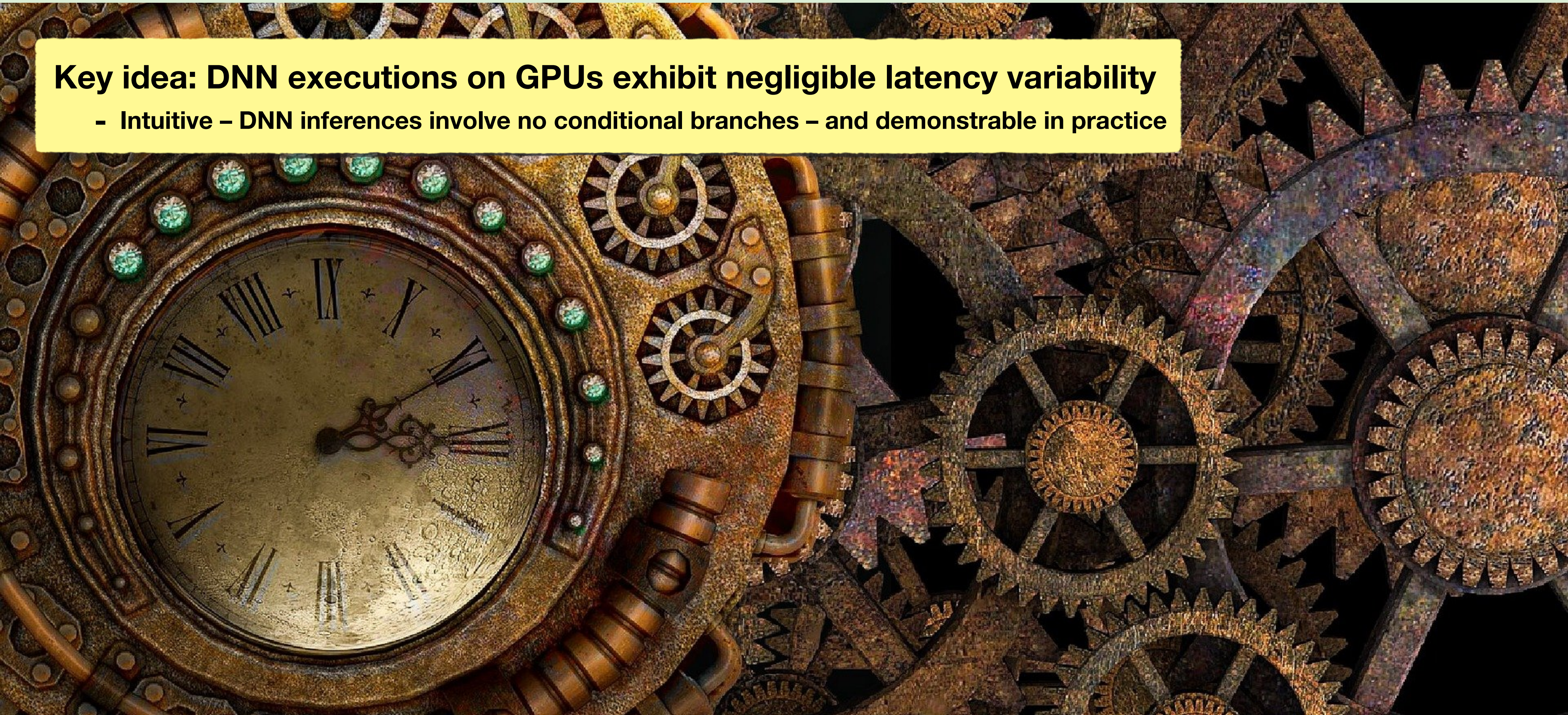- Concentrating all choices in a centralized controller

# Summary

**Key idea: DNN executions on GPUs exhibit negligible latency variability**

- Intuitive – DNN inferences involve no conditional branches – and demonstrable in practice

**Clockwork: From DNN predictability to an E2E predictable DNN serving platform**

- Recursively ensures that all internal architecture components have predictable performance
- Concentrating all choices in a centralized controller

**Outperforms state-of-the-art DNN serving platforms**

- Efficiently fulfills aggressive tail-latency SLOs
- Supports 1000s of DNN models with varying workload characteristics concurrently on each GPU

# Summary

**Key idea: DNN executions on GPUs exhibit negligible latency variability**
- Intuitive – DNN inferences involve no conditional branches – and demonstrable in practice

**Clockwork: From DNN predictability to an E2E predictable DNN serving platform**
- Recursively ensures that all internal architecture components have predictable performance
- Concentrating all choices in a centralized controller

**Outperforms state-of-the-art DNN serving platforms**
- Efficiently fulfills aggressive tail-latency SLOs
- Supports 1000s of DNN models with varying workload characteristics concurrently on each GPU

**https://gitlab.mpi-sws.org/cld/ml/clockwork**

ARTIFACT EVALUATED
usenix ASSOCIATION
AVAILABLE

ARTIFACT EVALUATED
usenix ASSOCIATION
FUNCTIONAL

ARTIFACT EVALUATED
usenix ASSOCIATION
REPRODUCED