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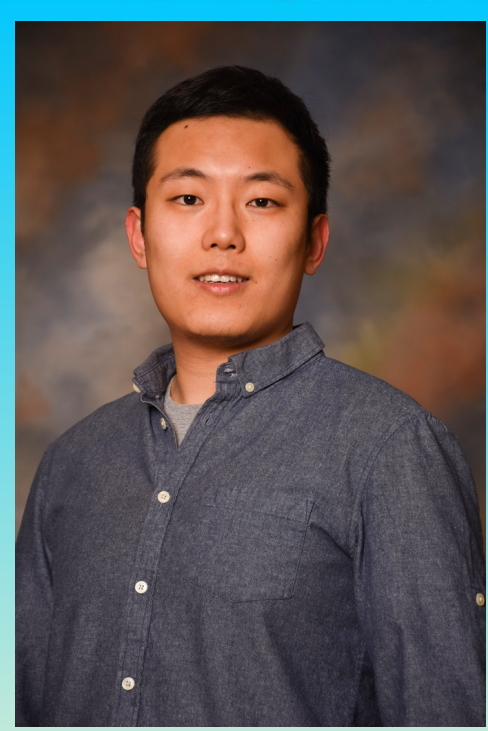## Performance Predictability from the Bottom Up



**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

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



**Requests arrive at different rates and regularity**
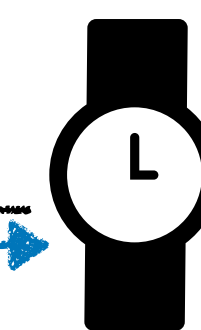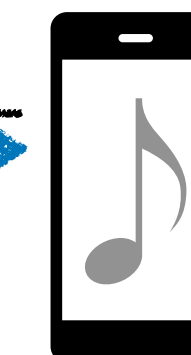


**Each request has an inherent deadline**

## 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

**Inference latency**

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

Clipper
100ms SLO

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

CDF

**200 ms**

Latency (ms): 0 100 500

**Tail latency >> SLO**

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

Clockwork
100ms SLO

Percentile: 99.9999, 99.999, 99.99, 99.9, 99, 90, 0
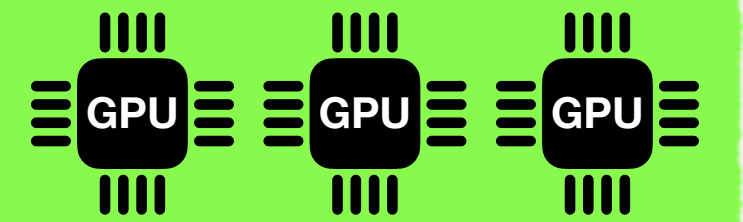
**95 ms**

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**

**Throughput gains only 25%**

**Single-thread latency is extremely predictable**

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

# 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)

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

# Designing a Predictable Worker (2/2)

**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



**Users**

**Centralized Controller**

**Worker processes**

Smarter load balancing & scheduling decisions

LOADs

INFERs

RAM

GPU Memory
**PageCache**

GPU Exec

**GPU Worker Node W₁**

**Global State Manager**

Latency Profiles

Pending Tasks

Memory State

# Evaluation

# Questions

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

Can Clockwork serve thousands of model instances?

How low can Clockw**This talk**s 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**

# 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.

# Are Clockwork Workers Predictable?

Clockwork relies on predicting the model inference latency for scheduling

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

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

## INFER

Experiment duration = 6 hours,
Offered load ~ 10,000 r/s

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

— Overpredict
— Underpredict

Percentile (y-axis): 0, 90, 99, 99.9, 99.99, 99.999, 99.9999, 99.99999, 99.999999, 99.9999999

Error (ms) (x-axis): 0, 5, 10, 15

Clockwork consistently overpredicts more than its underpredicts

Errors are significant only in extremely rare cases

# Does Consolidating Choice Help?

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

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

**Goodput = SLO compliant throughput**

**The workload is successfully scheduled by Clockwork**

- Goodput $\approx$ offered load

**Latency of all completed requests**

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

- All others completed within SLO

**Batching prioritized, absorbs spikes**

**Many cold starts**

**Cold requests = 1.3% of all requests**

# Does Clockwork Controller Scale?



**40 emulated workers**

Offered Load
Goodput
**Peak**

Req/s: 100000, 75000, 50000, 25000, 0

Time (mins): 0 1 2 3 4 5 6 7 8

**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**

Peak Goodput (r/s) vs Number of Workers
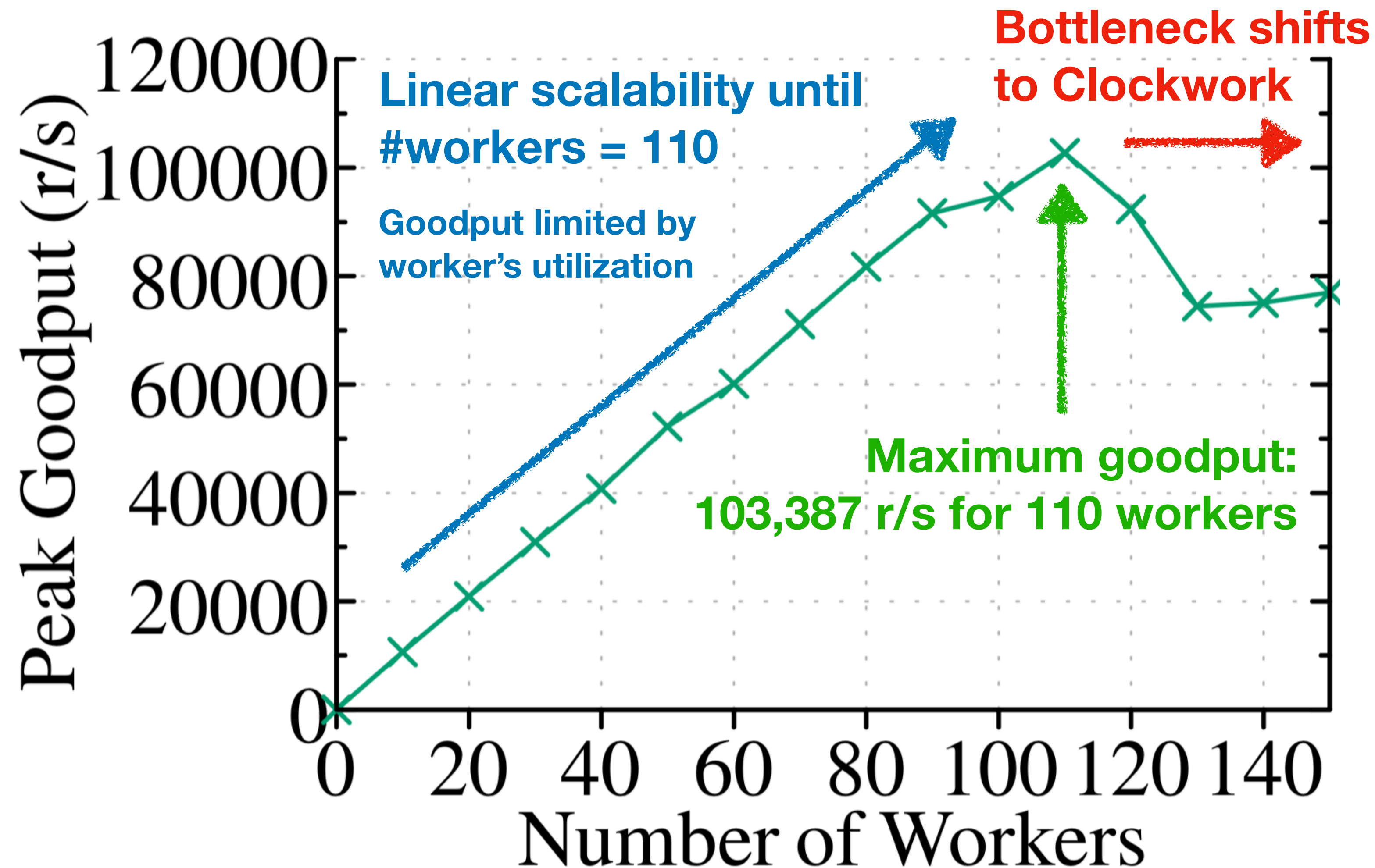
**Methodology**

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

# 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**