# DjiNN and Tonic: DNN as a Service and Its Implications for Future Warehouse Scale Computers

**Johann Hauswald**, Yiping Kang, Michael A. Laurenzano, Quan Chen, Cheng Li, Trevor Mudge, Ronald G. Dreslinski, Jason Mars, Lingjia Tang

University of Michigan — Ann Arbor, MI

DjiNN
& Tonic

**M**

**COMPUTER SCIENCE
& ENGINEERING**

**UNIVERSITY OF MICHIGAN**

ClarityLab

# Intelligent Personal Assistant (IPA) Queries

# Intelligent Personal Assistant (IPA) Queries



"Set my alarm
for 6am."

# Intelligent Personal Assistant (IPA) Queries
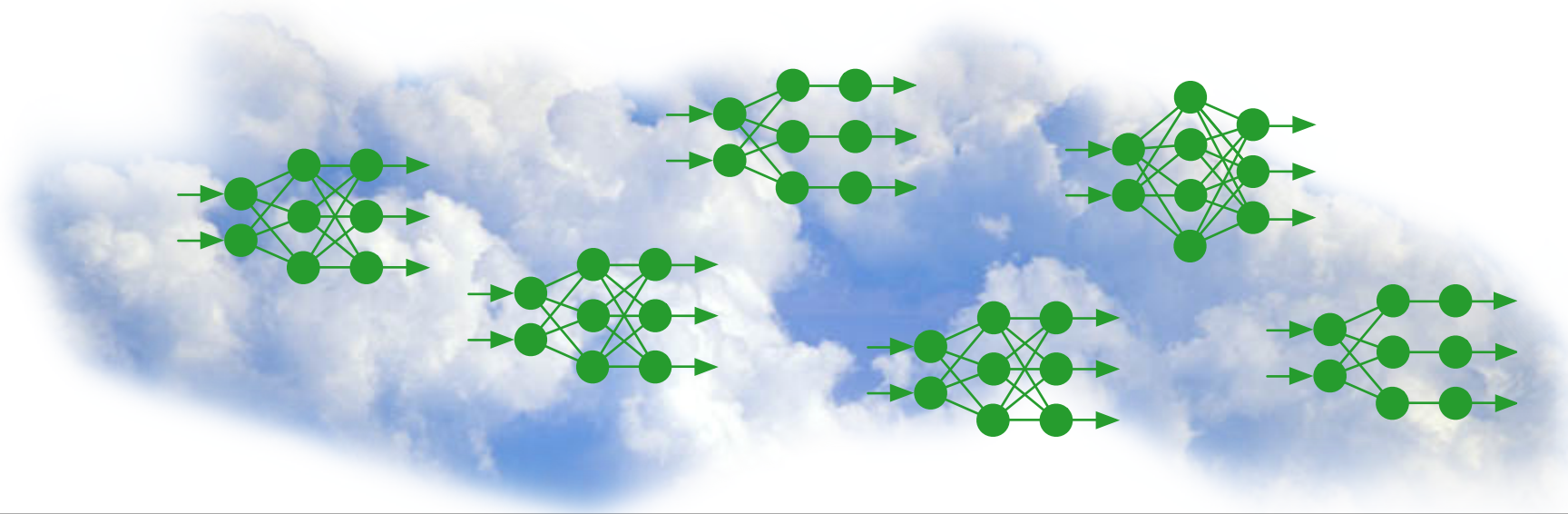
# Intelligent Personal Assistant (IPA) Queries

"What's the speed of light?"

DjiNN & Tonic

ClarityLab

# Intelligent Personal Assistant (IPA) Queries

# Intelligent Personal Assistant (IPA) Queries



"Who is this?"

# Intelligent Personal Assistant (IPA) Queries

# Intelligent Personal Assistant (IPA) Queries

# Deep Neural Networks (DNNs)



Inference

**Feature Maps**

Input   Convolutional layer   Pooling layer

Fully Connected layer

0.1 "Spiderman"

**0.9 "Superman"**

0.5 "Batman"

Network Architecture

# Deep Neural Networks (DNNs)



speech features

word vectors

"who" $\longrightarrow$ $\begin{bmatrix} w_0 \, w_1 \dots w_k \end{bmatrix}$

"is" $\longrightarrow$ $\begin{bmatrix} w_0 \, w_1 \dots w_k \end{bmatrix}$

"this" $\longrightarrow$ $\begin{bmatrix} w_0 \, w_1 \dots w_k \end{bmatrix}$

Inference

"Who", "is", "this"

"Who" (PRONOUN)
"is" (VERB)
"this" (PRONOUN)

Fully Connected layer

Network Architecture

# DNN as a Service



Image Classification

Digit Recognition

Facial Recognition

Speech Recognition

Natural Language Processing

# DNN as a Service

Image Classification

Digit Recognition

Facial Recognition

Speech Recognition

Natural Language Processing

**Unified, highly optimized appliance for DNN**

DjiNN & Tonic
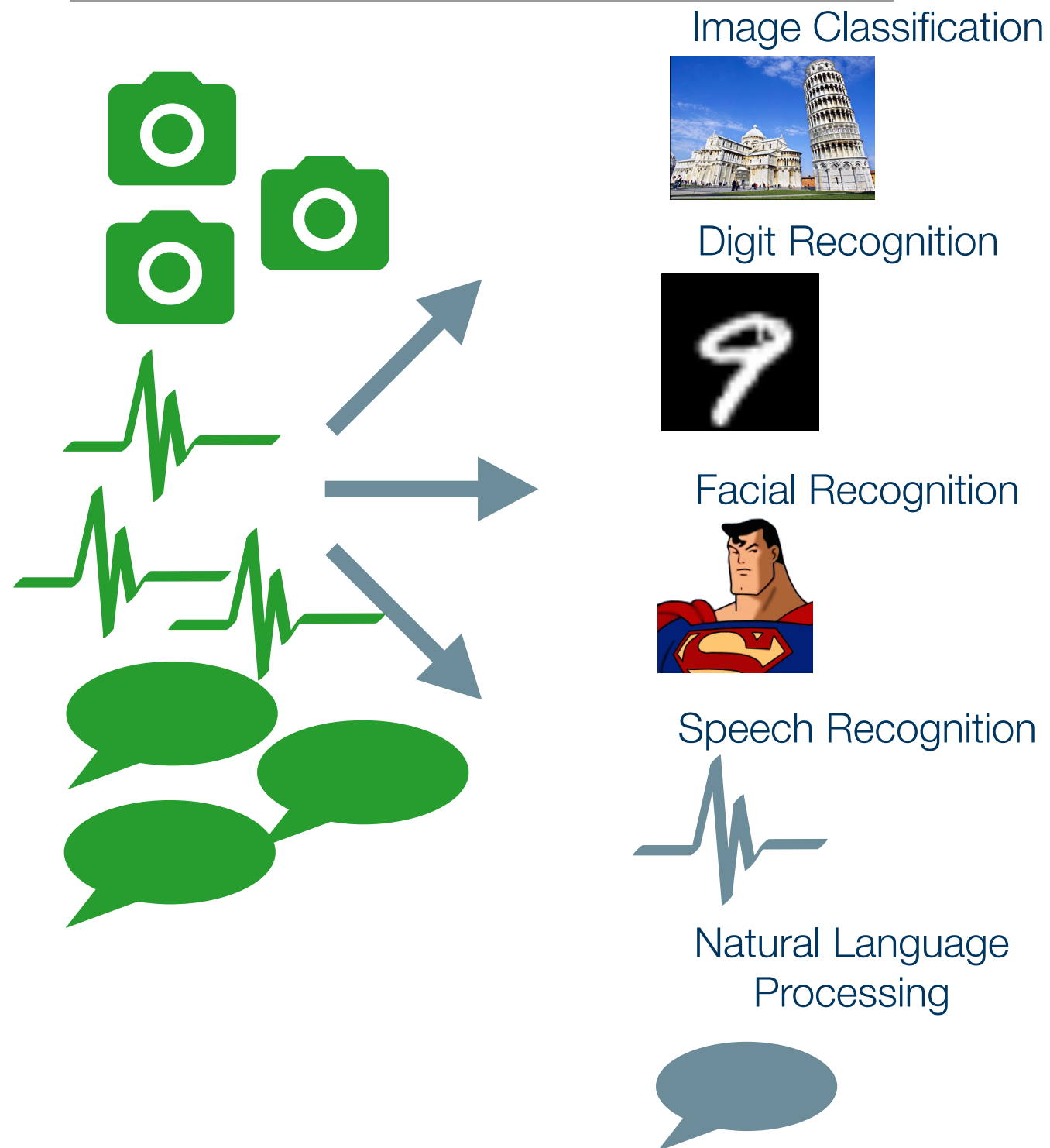
ROBERT MCMILLAN     07.16.14   6:30 AM
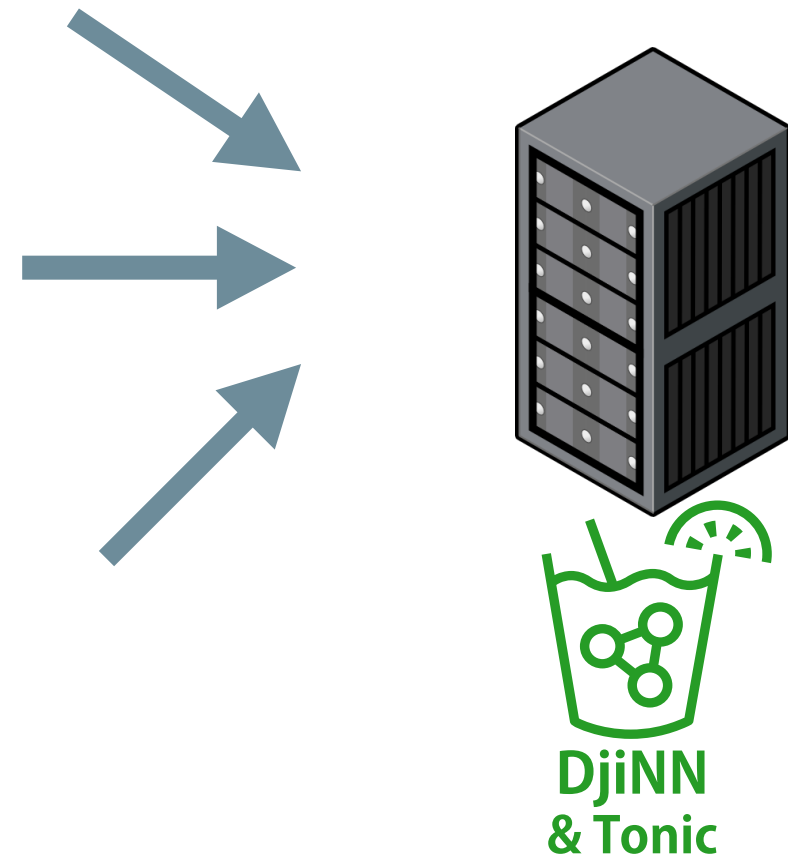
# INSIDE THE ARTIFICIAL BRAIN THAT'S REMAKING THE GOOGLE EMPIRE

*"We now have probably 30 or 40 different teams at Google using our infrastructure."*
-Jeff Dean

# DNN as a Service

Image Classification

Digit Recognition

Facial Recognition

Speech Recognition

Natural Language Processing

**Unified, highly optimized appliance for DNN**

DjiNN & Tonic

# Challenge

Design a high throughput Warehouse Scale Computer (WSC)
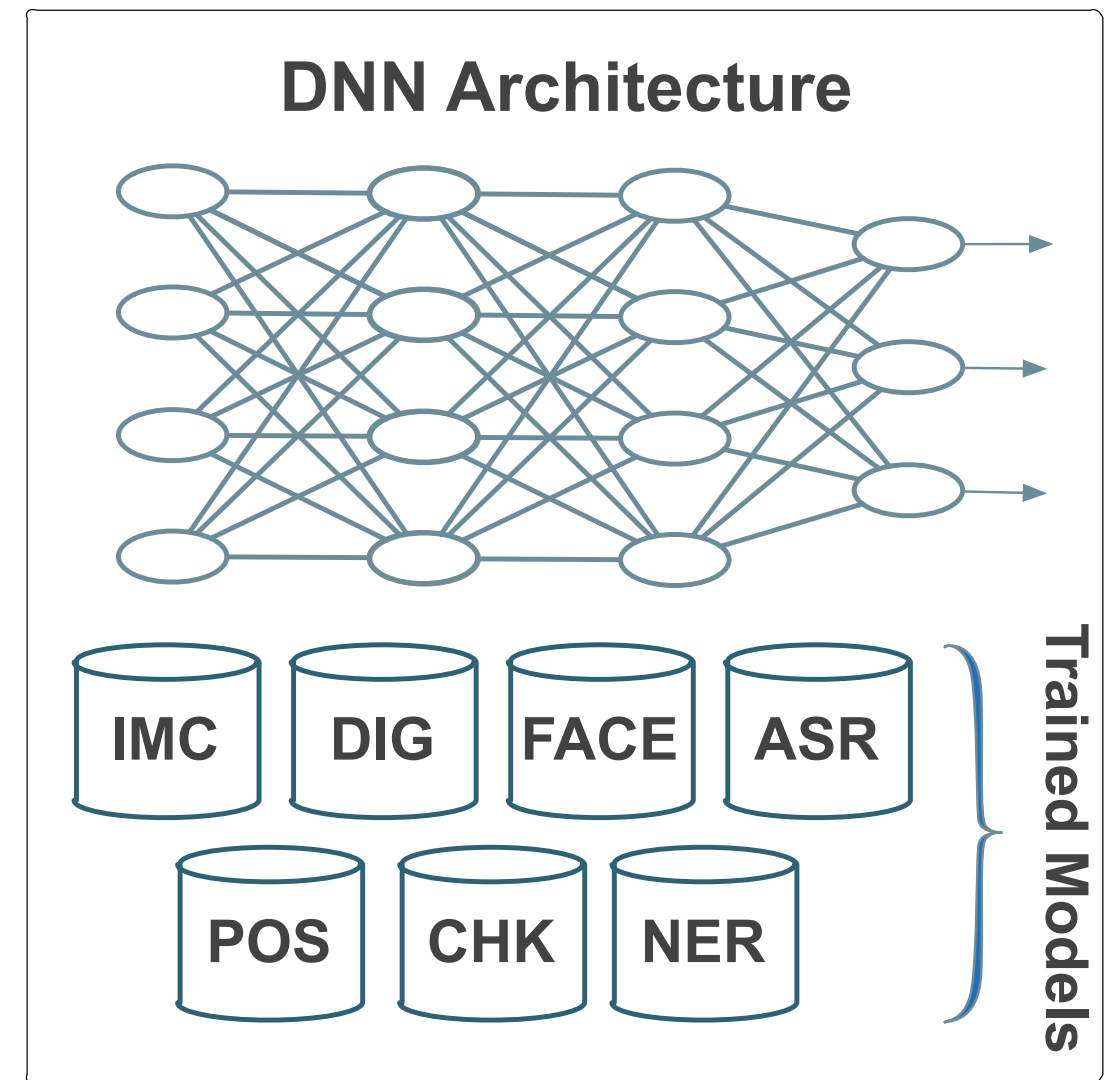for DNN as a Service

# Outline

- DjiNN and Tonic: DNN as a Service

- Identifying Bottlenecks for DNN as a Service

- Designing a High Throughput System

- Future Warehouse Scale Computer (WSC) Designs

- Conclusion

# DjiNN and Tonic: DNN as a Service

# DjiNN Design Goals

- Single web-service for DNN

- Diverse applications
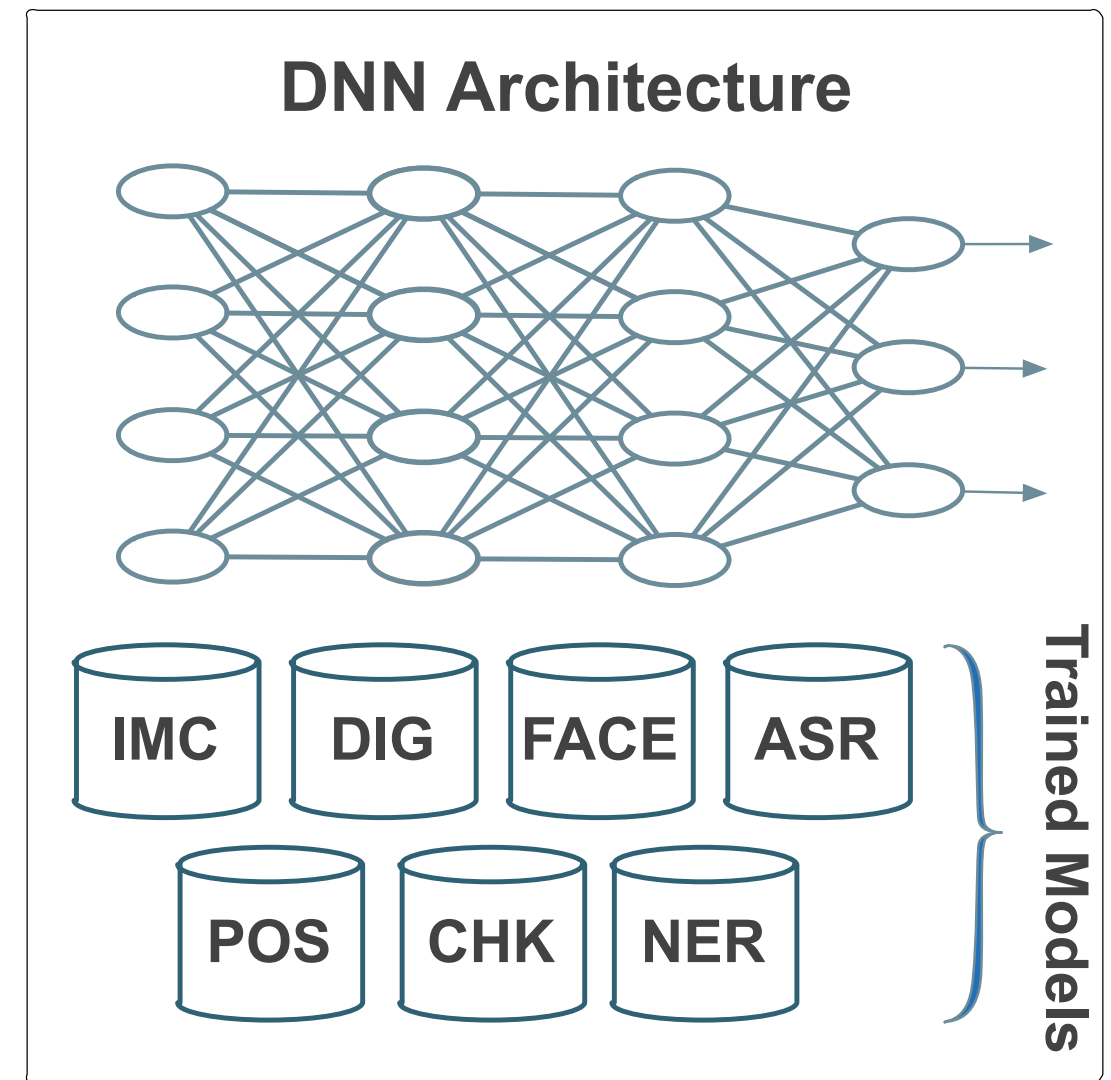
- Low overhead request processing

### DjiNN DNN Service

**DNN Architecture**



IMC | DIG | FACE | ASR

POS | CHK | NER

Trained Models

# DjiNN Implementation

- Decoupled architecture

- Arbitrary network architecture support

- Memory resident models for thread pool

- More details in paper

# Tonic Suite

- End-to-end applications that make requests to the DjiNN Service

- Span image, speech, and natural language processing

- State-of-the-art neural network architectures

**Image Task**

IMC    DIG    FACE



**Speech Recognition (ASR) Task**

    *"It's business, Superman"*

**Natural Language Processing Task**

POS    *"business" (noun)*
       *"Superman" (P. noun)*

CHK    *"It's" (VP, B-NP)*
       *"business" (NP, I-NP)*

NER    *"Superman" (PERSON)*

DjiNN & Tonic

ClarityLab

Release includes: inputs, pre-trained models, and modified Caffe

# djinn.clarity-lab.org

# Identifying Bottlenecks for DNN as a Service

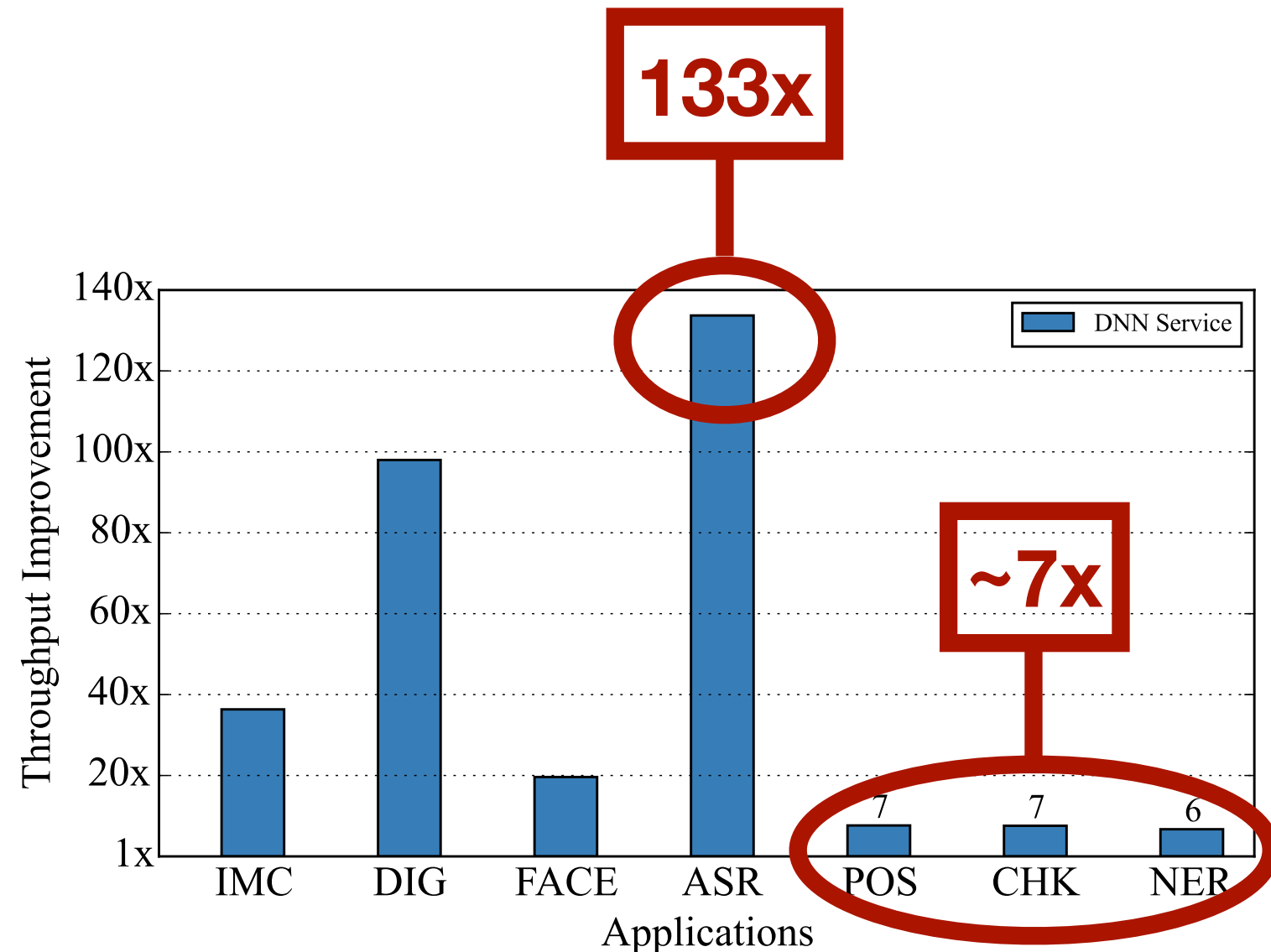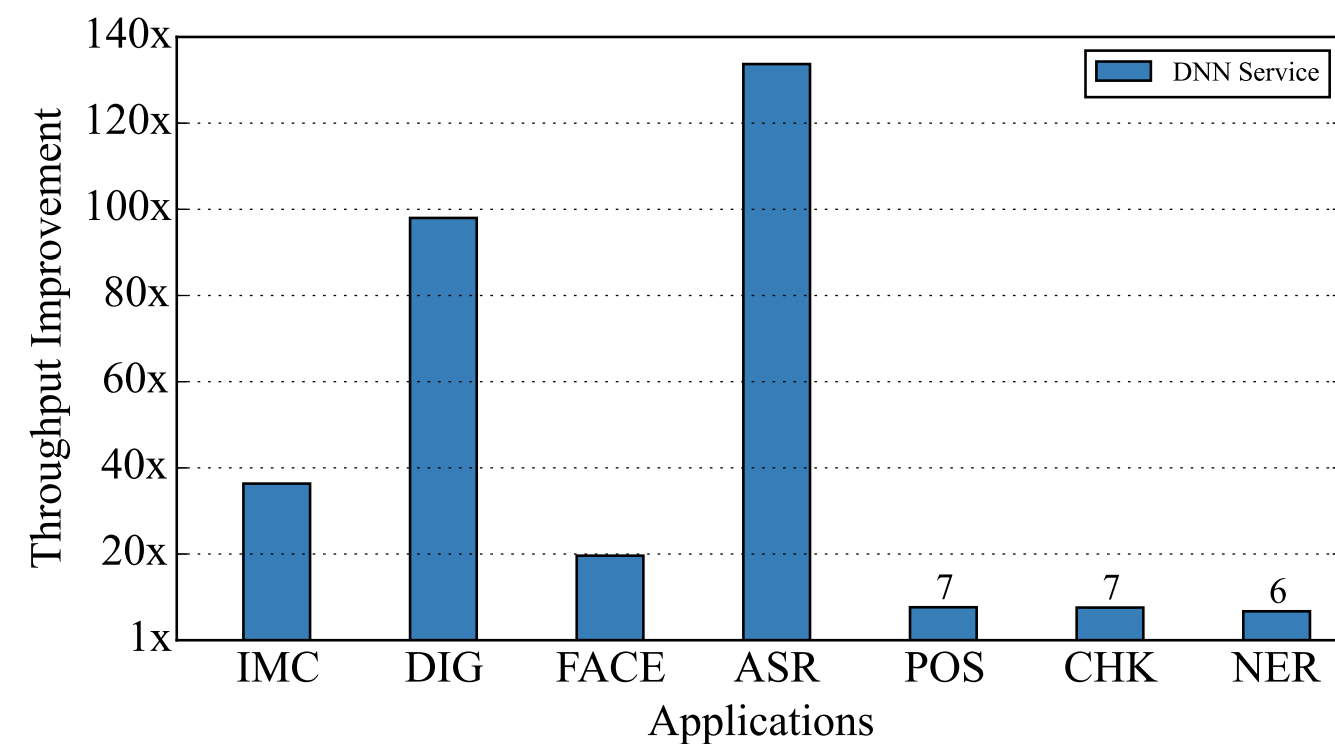# Identifying Bottlenecks for DjiNN and Tonic

- Software: Caffe (modified)

- CPU: Intel Xeon E5-2620 2.10GHz

  - ATLAS (vectorized)



(a) IMC     (b) DIG     (c) FACE     (d) ASR

(e) POS     (e) NER     (f) CHK

Legend: DNN, Other

**DNN: More than 80% of cycles**

# Identifying Bottlenecks for DjiNN and Tonic

- Software: Caffe (modified)

- CPU: Intel Xeon E5-2620 2.10GHz

  - ATLAS (vectorized)

- GPU: NVIDIA Tesla K40M
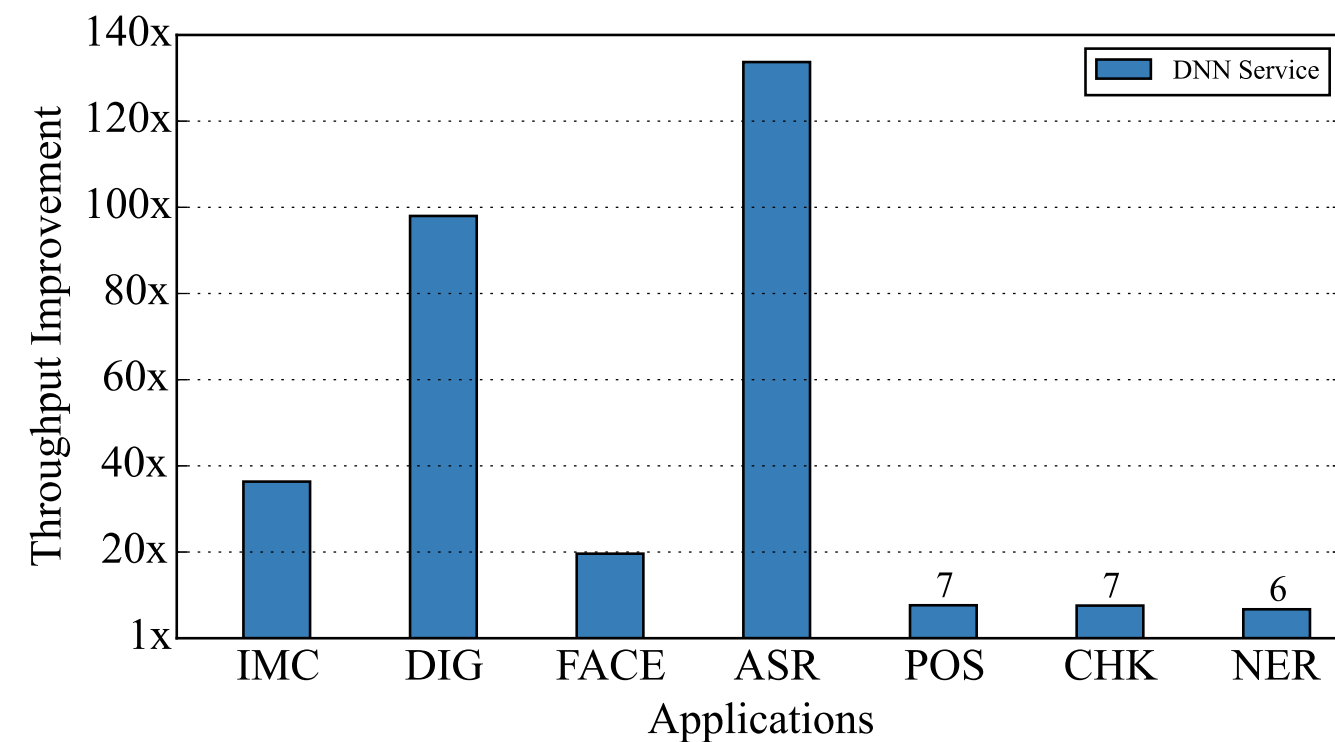
  - cuDNN v1 and Caffe

# Identifying Bottlenecks for DjiNN and Tonic

- Software: Caffe (modified)

- CPU: Intel Xeon E5-2620 2.10GHz

  - ATLAS (vectorized)

- GPU: NVIDIA Tesla K40M

  - cuDNN v1 and Caffe

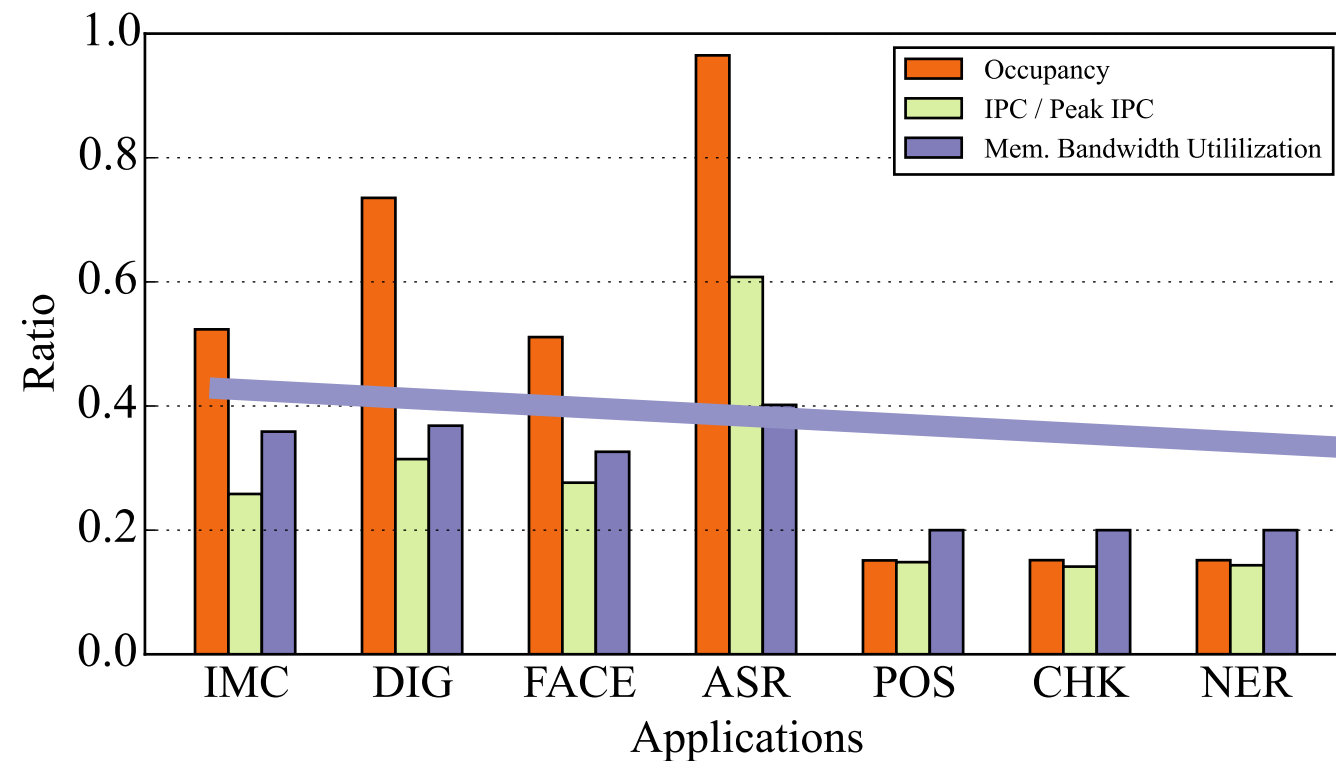# Identifying Bottlenecks for DjiNN and Tonic



Throughput Improvement

GPU Profiling

# Identifying Bottlenecks for DjiNN and Tonic

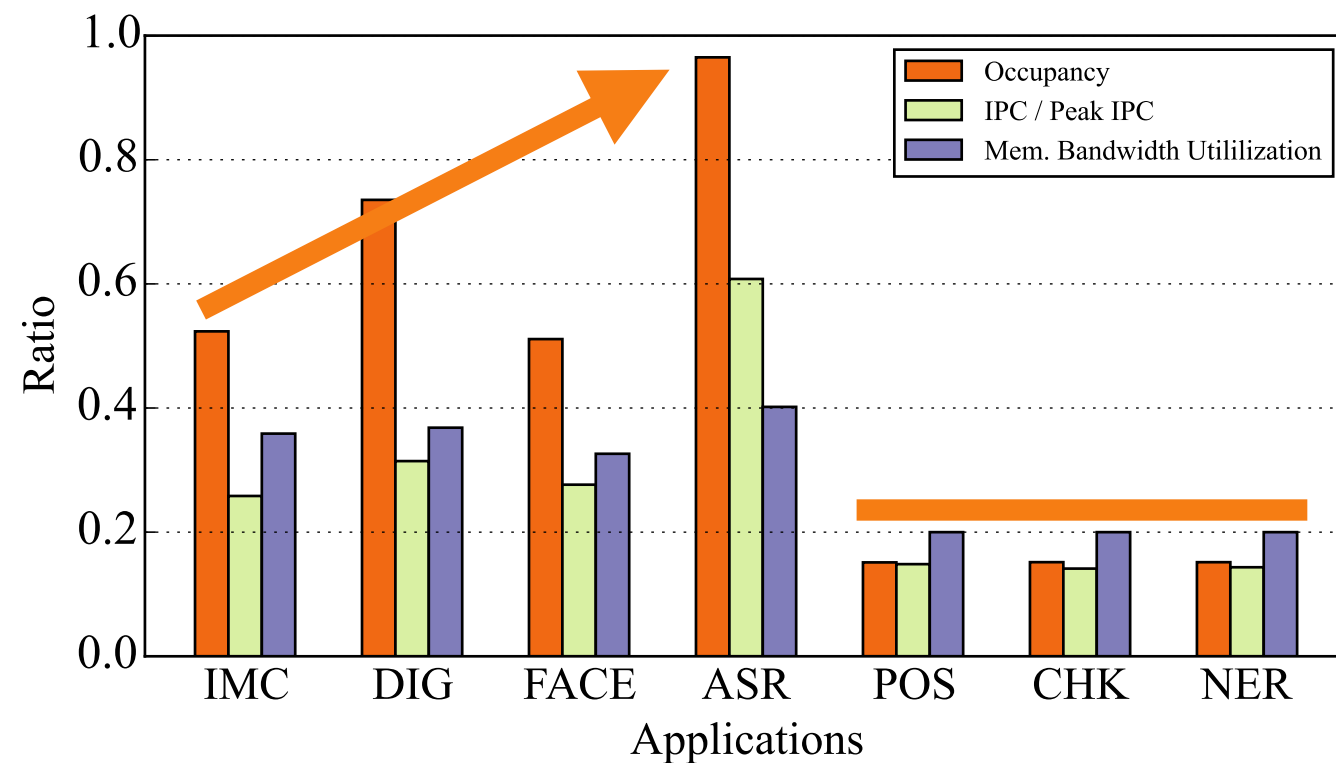**Low memory bandwidth utilization**



Throughput Improvement



GPU Profiling

# Identifying Bottlenecks
# for DjiNN and Tonic



Throughput Improvement
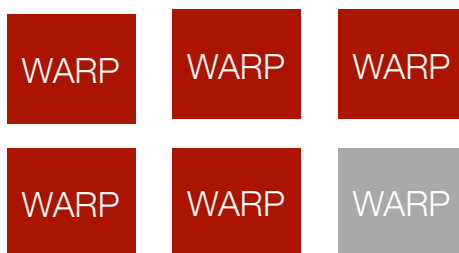
GPU Profiling

**GPU is not fully utilized**

# Designing a High Throughput System

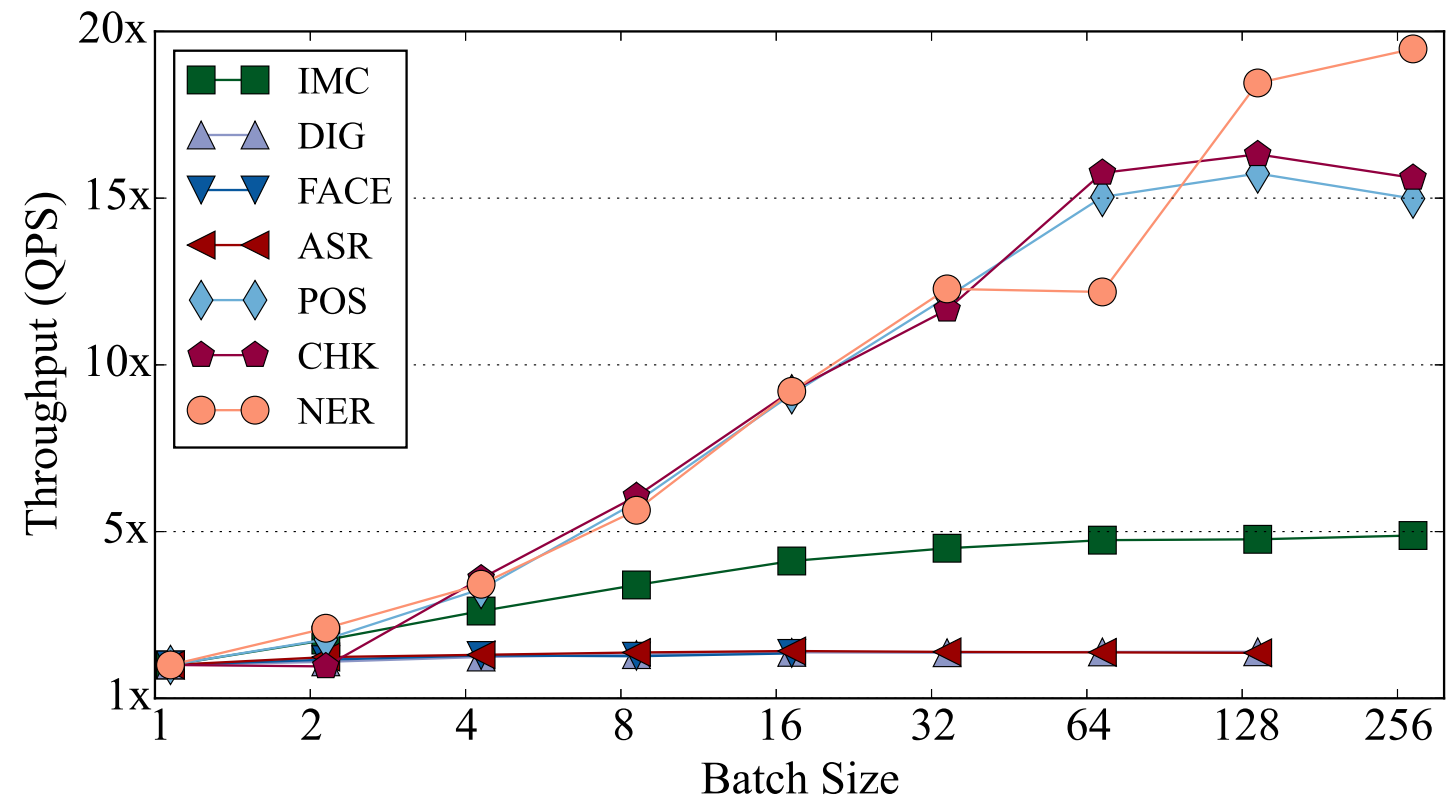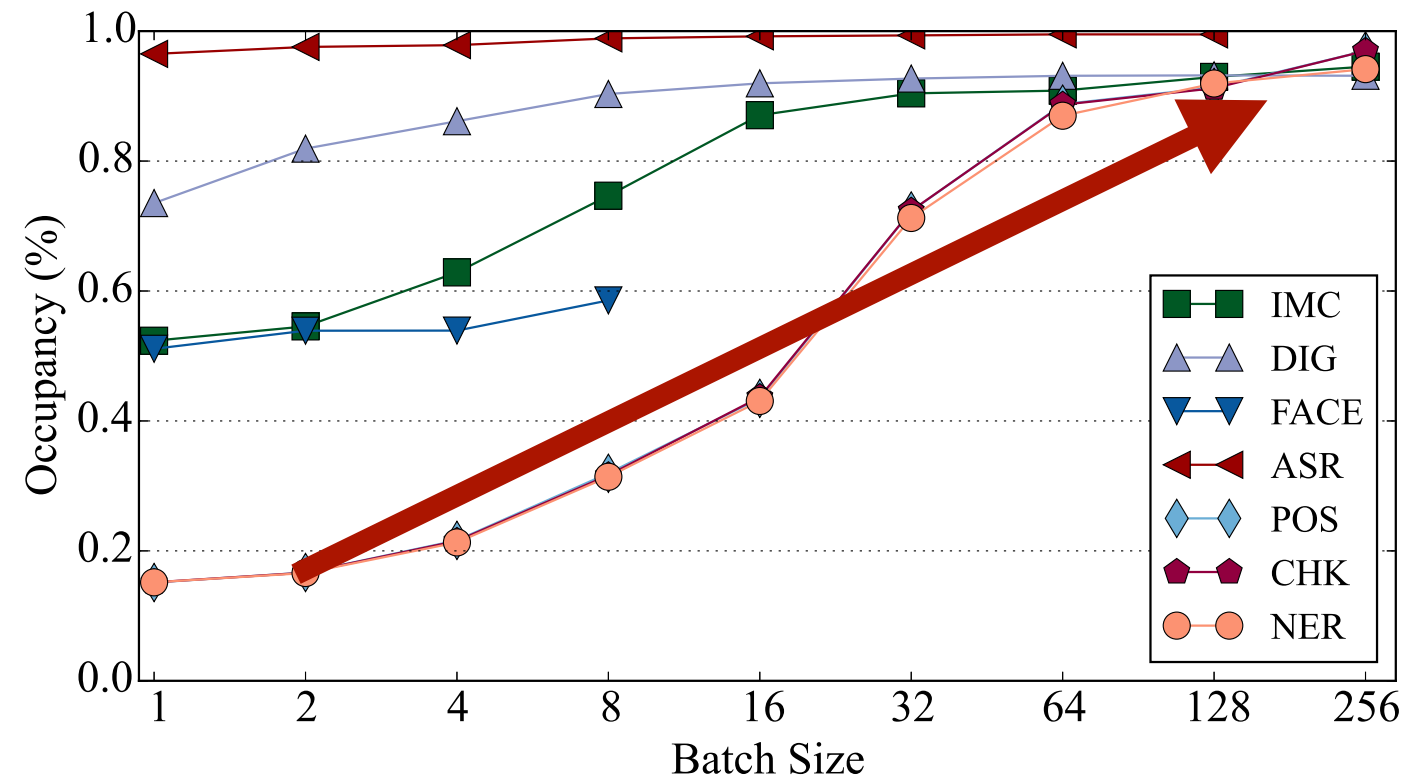# Designing a High Throughput System — Batching

Batch Size:

$$\begin{matrix} X_0 & X_{01} & X_{02} & X_{03} & X_{04} & X_{05} \\ X_1 & X_{11} & X_{12} & X_{13} & X_{14} & X_{15} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ X_d & X_d & X_d & X_d & X_d & X_d \end{matrix}$$

Streaming Multiprocessor

| WARP | WARP | WARP |
| WARP | WARP | WARP |

✔ High Occupancy

# Designing a High Throughput System — Batching

Batch Size:

$$
\begin{matrix}
X_0 & X_{01} & X_{02} & X_{03} & X_{04} & X_{05} \\
X_1 & X_{11} & X_{12} & X_{13} & X_{14} & X_{15} \\
\cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
\cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
\cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
X_d & X_d & X_d & X_d & X_d & X_d
\end{matrix}
$$

Streaming Multiprocessor

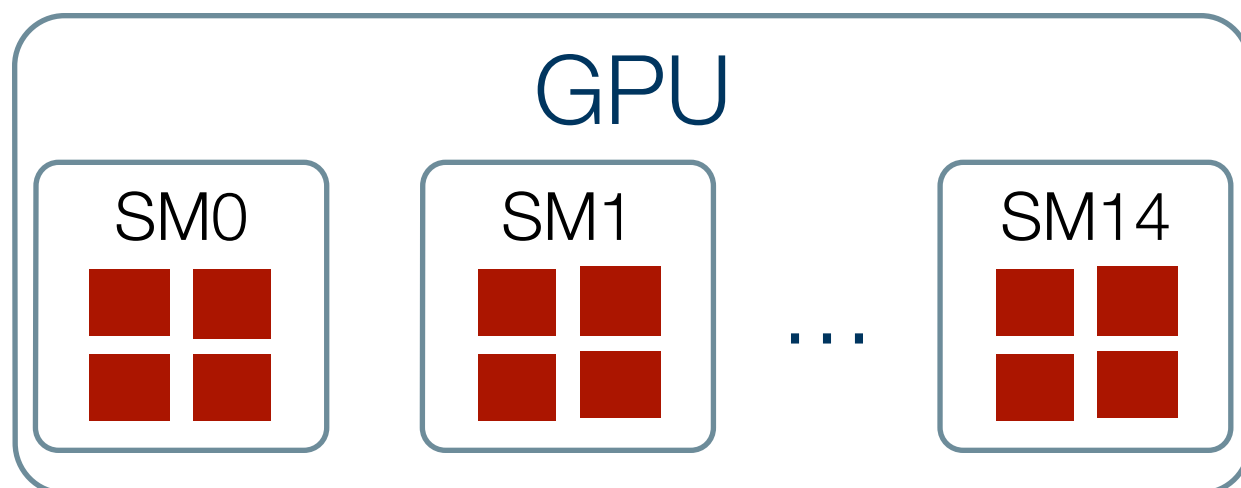| WARP | WARP | WARP |
| WARP | WARP | WARP |

✔ High Occupancy



~15x

~5x

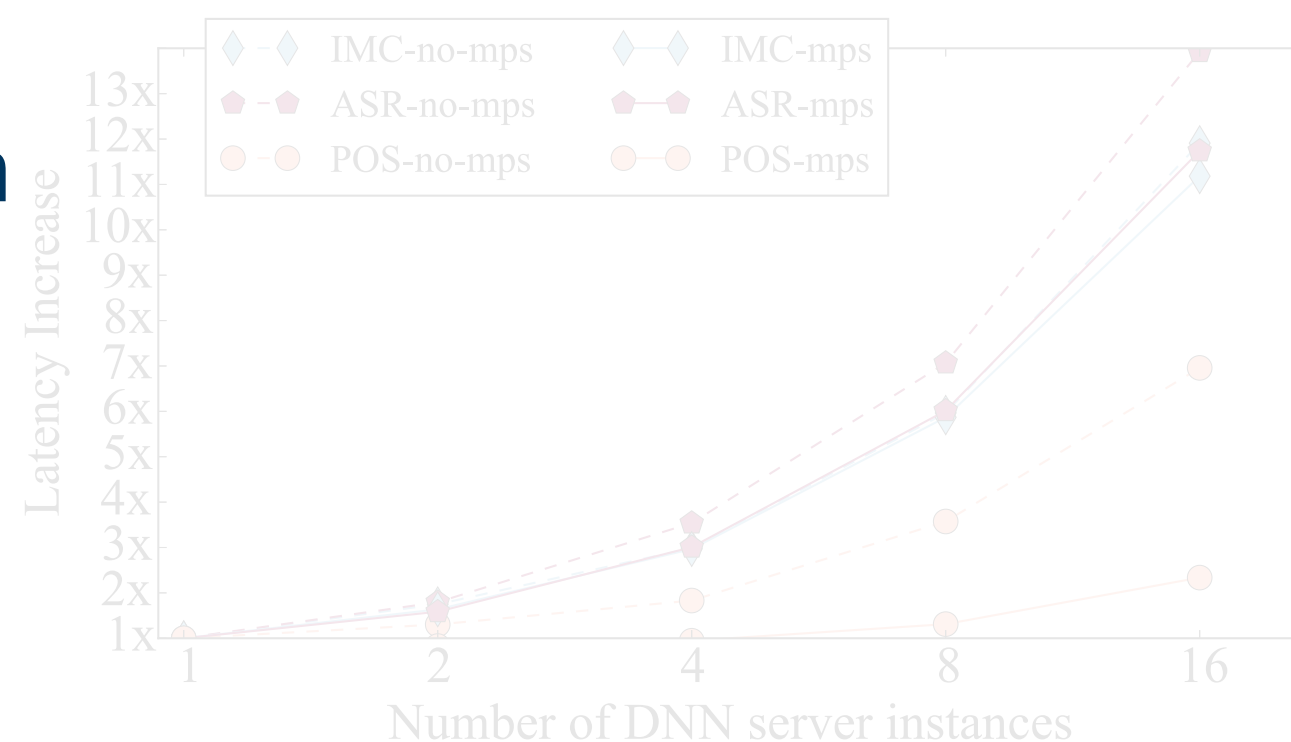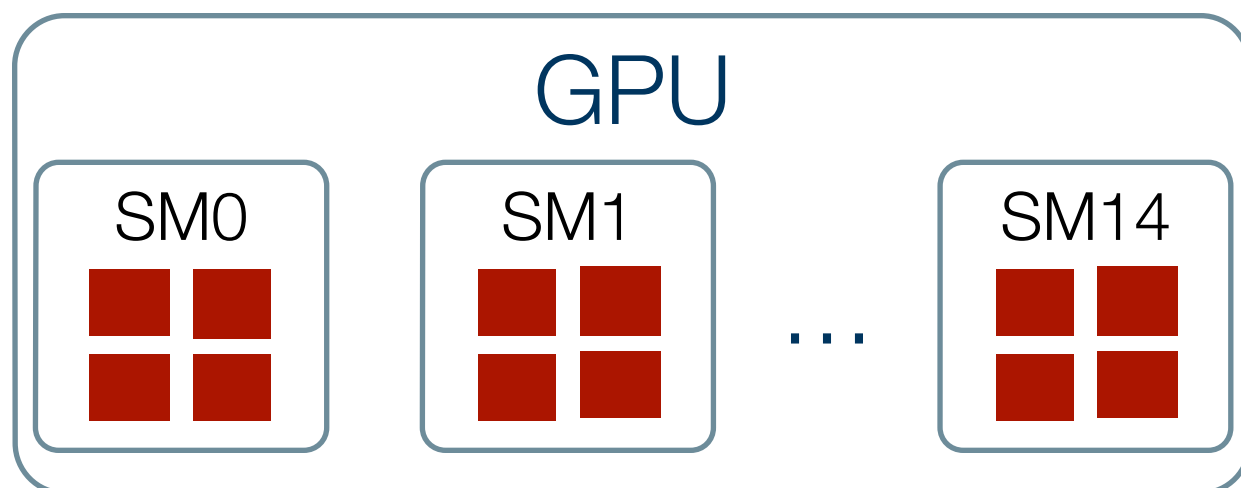# Designing a High Throughput System — Concurrent Execution

- Launch concurrent DNN services on the GPU
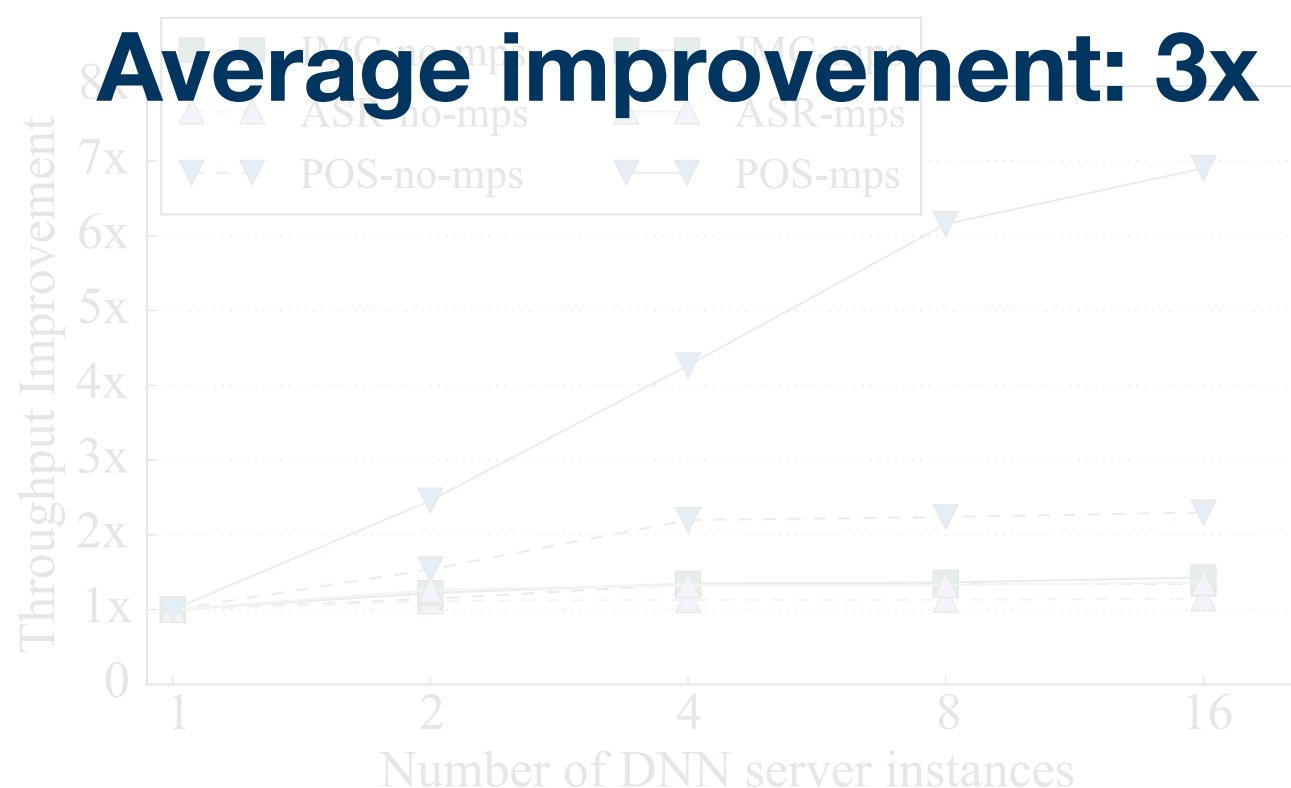
- Leverage NVIDIA Multi-Process Service (MPS) [1]



GPU

SM0  SM1  ...  SM14

[1] "Multi-Process Service" https://docs.nvidia.com/deploy/pdf/CUDA_ Multi_Process_Service_Overview.pdf

# Designing a High Throughput System — Concurrent Execution

- Launch concurrent DNN services on the GPU

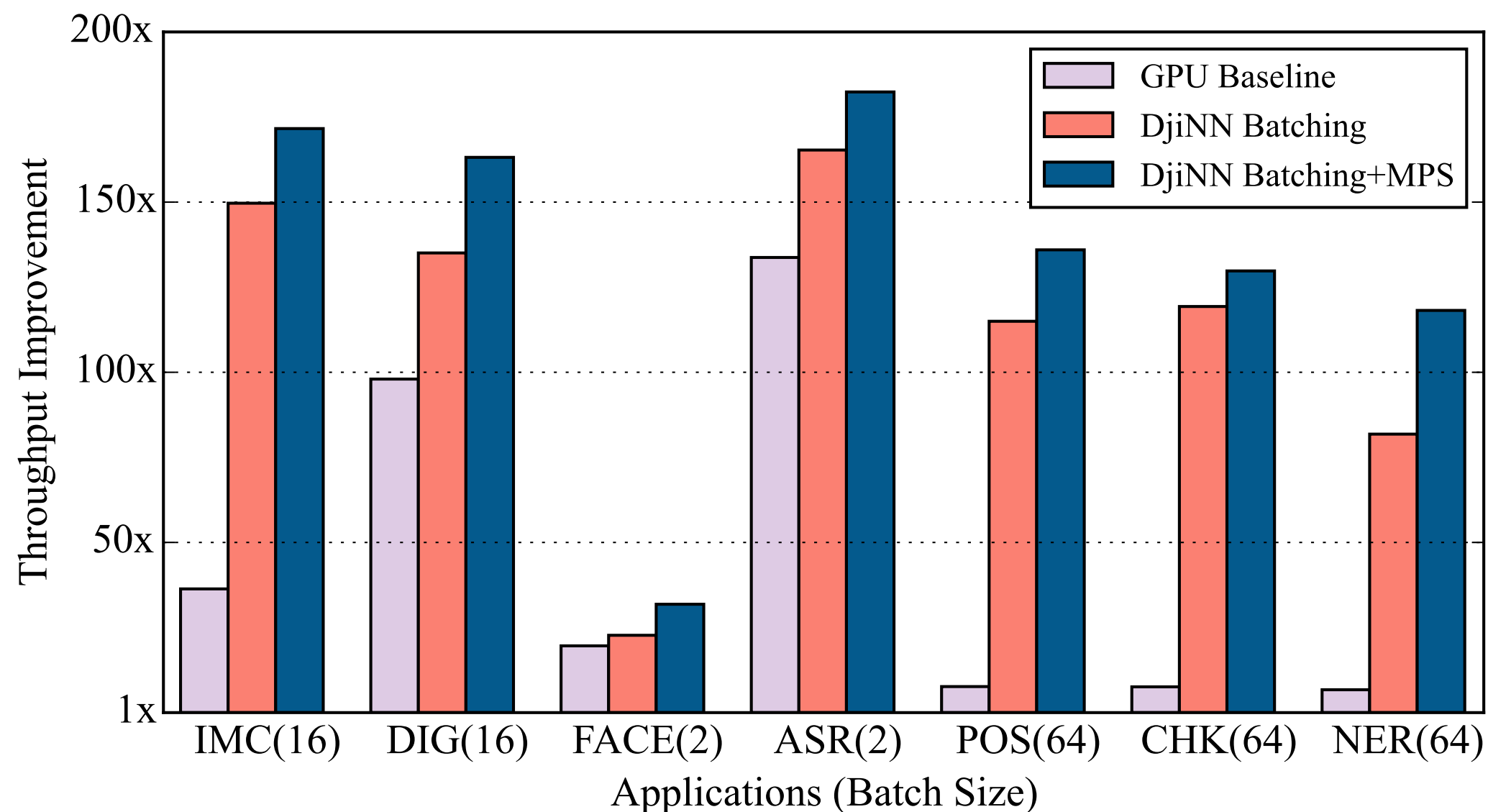- Leverage NVIDIA Multi-Process Service (MPS) [1]



**4 DNN concurrent instances Average improvement: 3x**

## GPU

| SM0 | | SM1 | ... | SM14 |



[1] "Multi-Process Service" https://docs.nvidia.com/deploy/pdf/CUDA_ Multi_Process_Service_Overview.pdf

# Designing a High Throughput System



**Average throughput improvement: 133x**
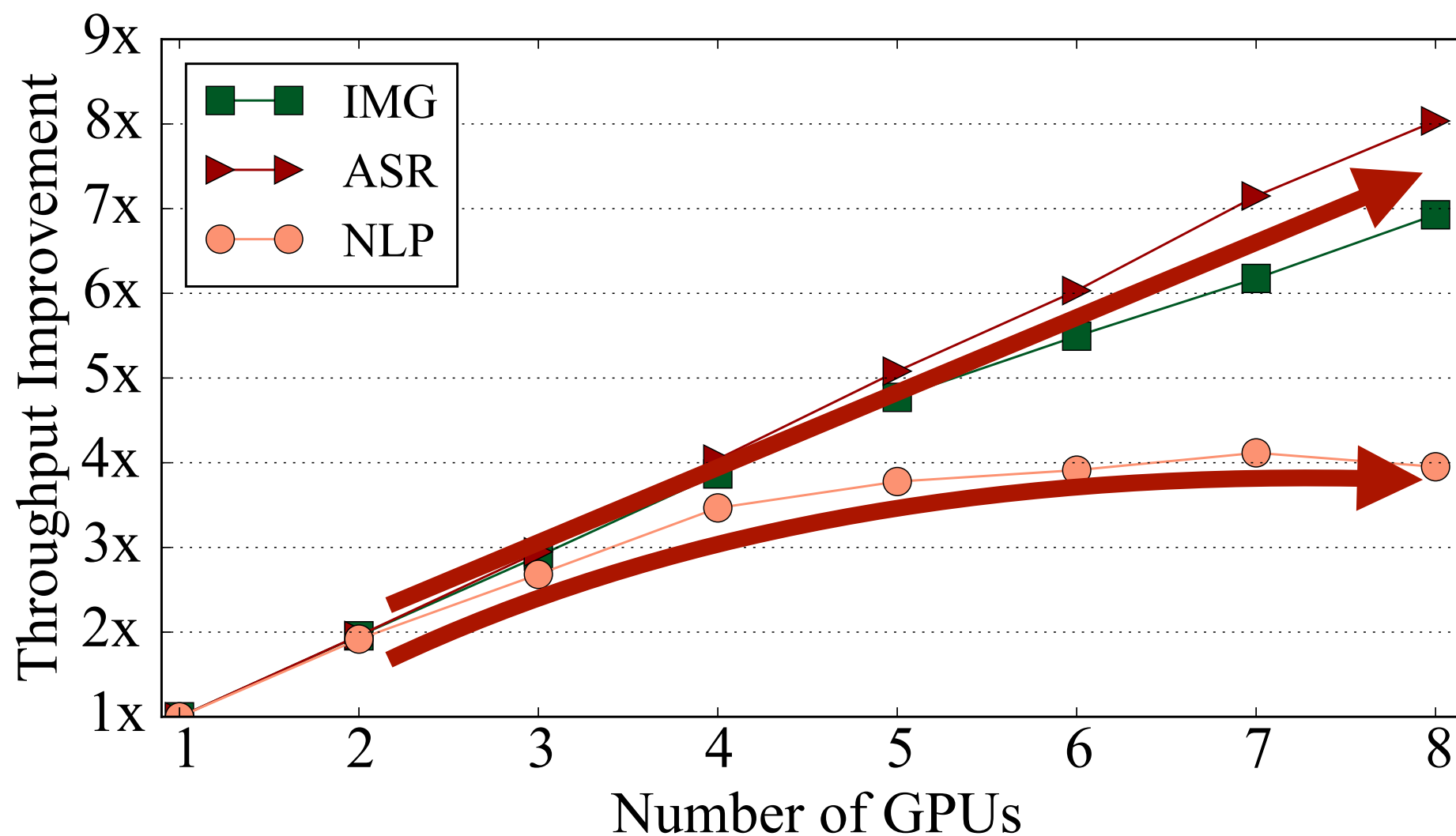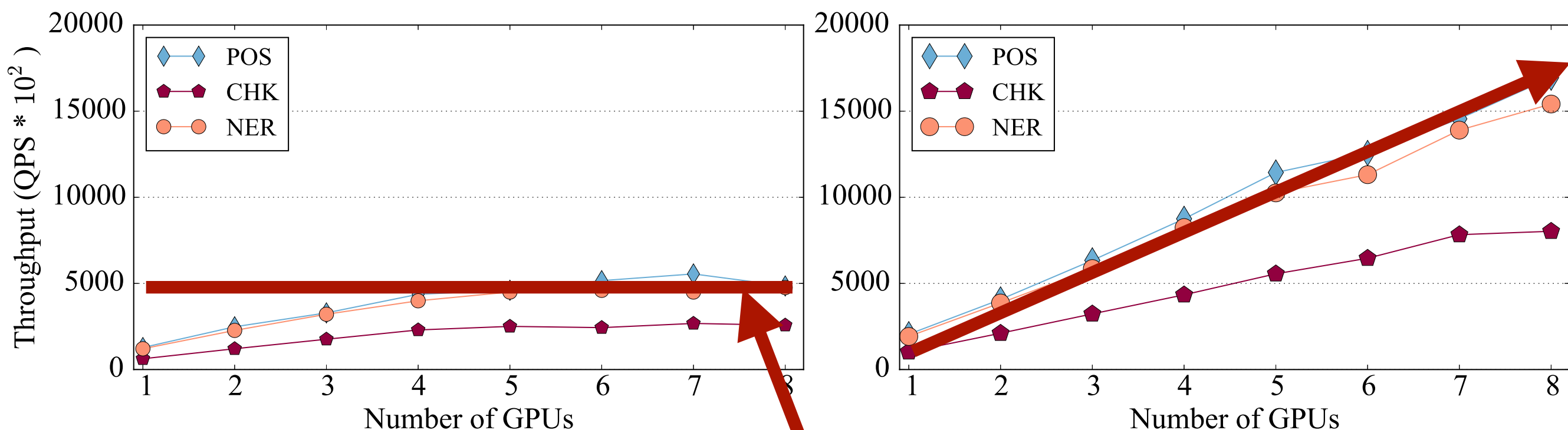
1 GPU

# GPU Scaling



**Average throughput improvement using optimizations: 771x**

# Bandwidth Requirements for Peak Throughput

Experimental setup: eliminate any data transfer to GPU
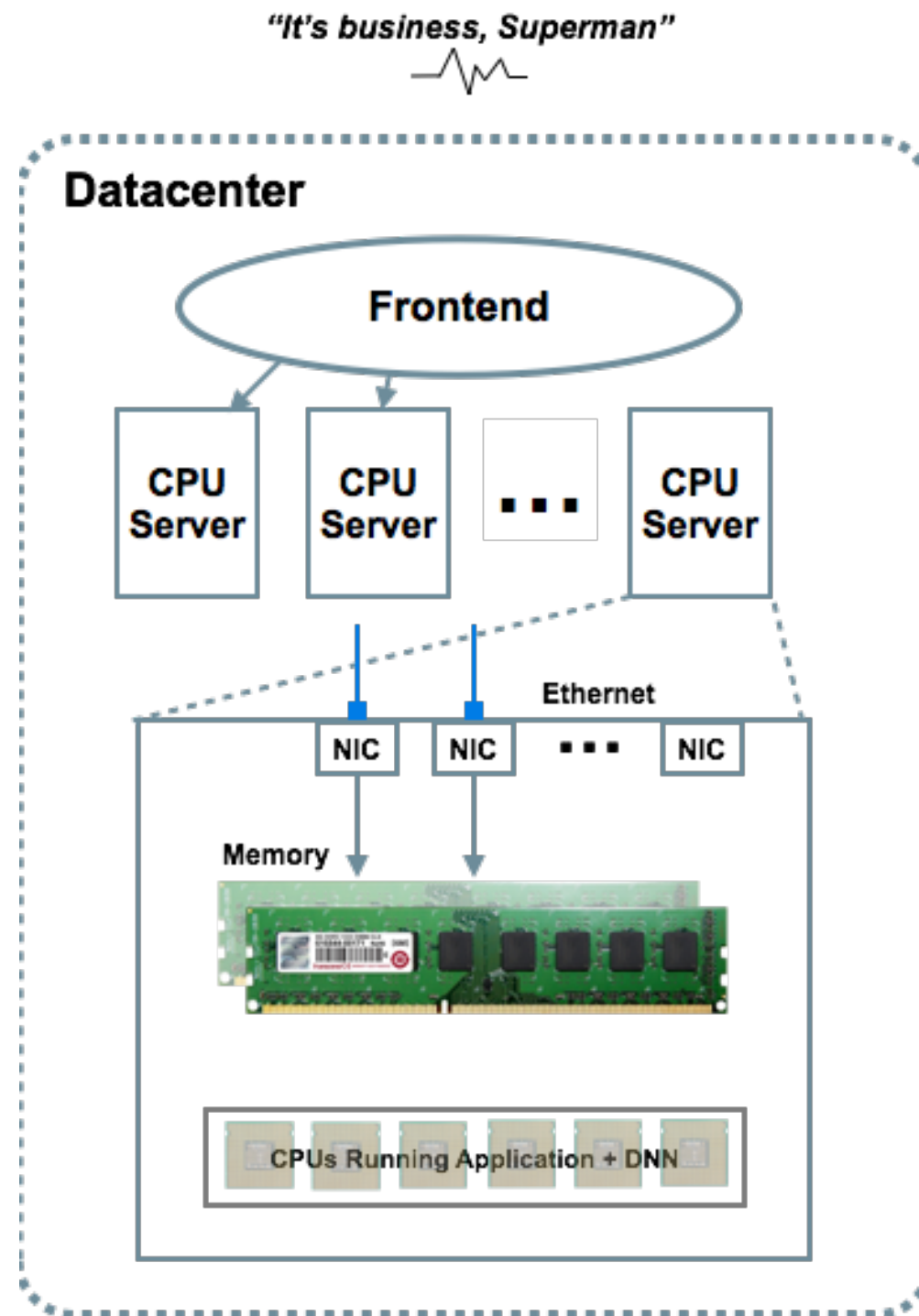


**Bandwidth to GPUs saturated**

**NLP requires more bandwidth to GPUs**
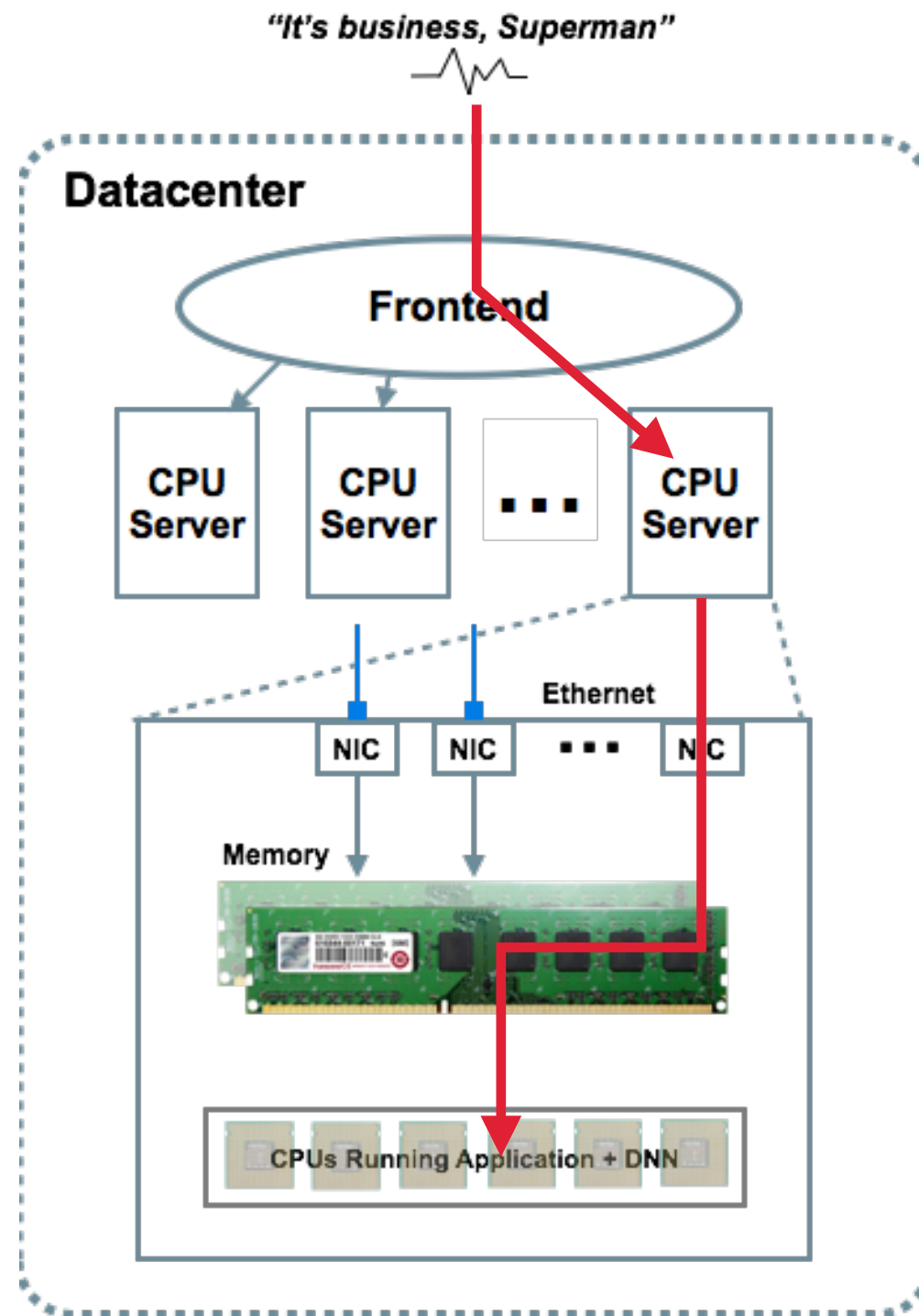
# Single Server Key Insights

- DNNs do not benefit equally from optimizations

  - High communication Natural Language Processing (NLP) tasks require far more bandwidth

- Optimizing compute platforms depends on DNN's computation and communication characteristics

# Future Warehouse Scale Computer (WSC) Designs
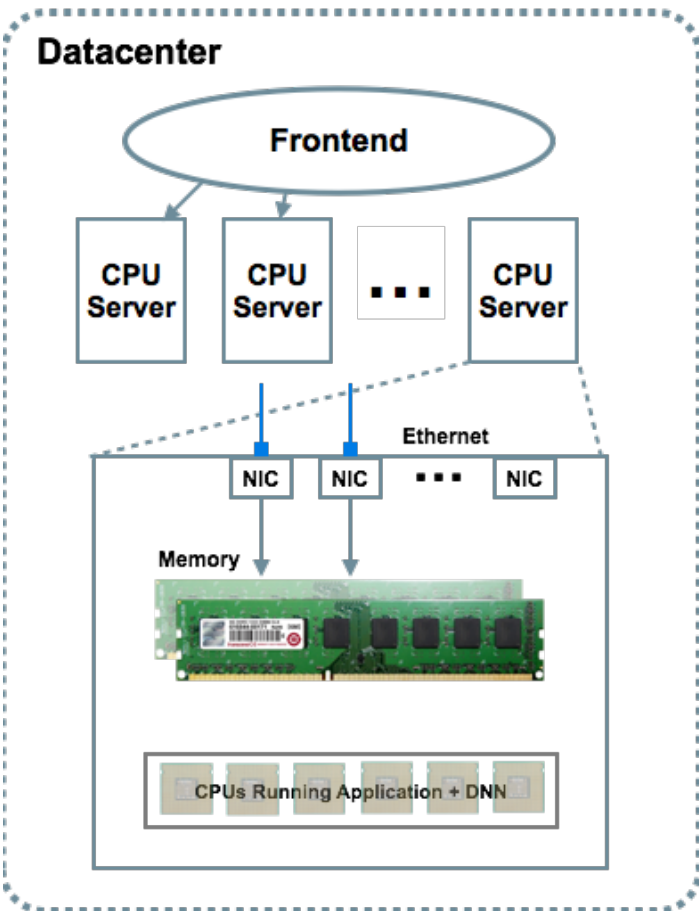
# Future WSC Designs — Server Designs

# Future WSC Designs — Server Designs
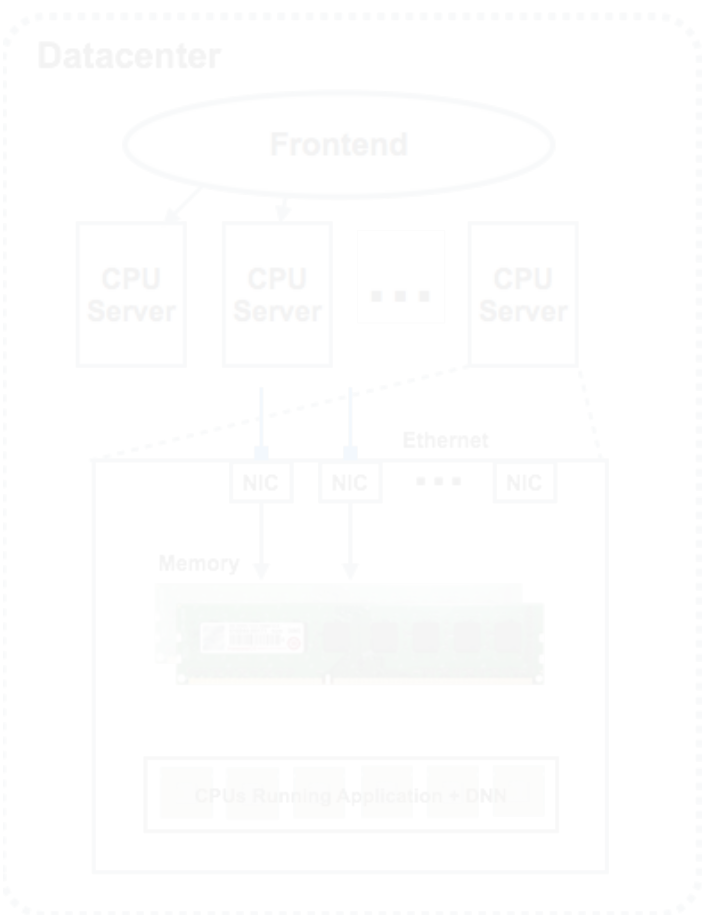
# Future WSC Designs — Server Designs



"It's business, Superman"

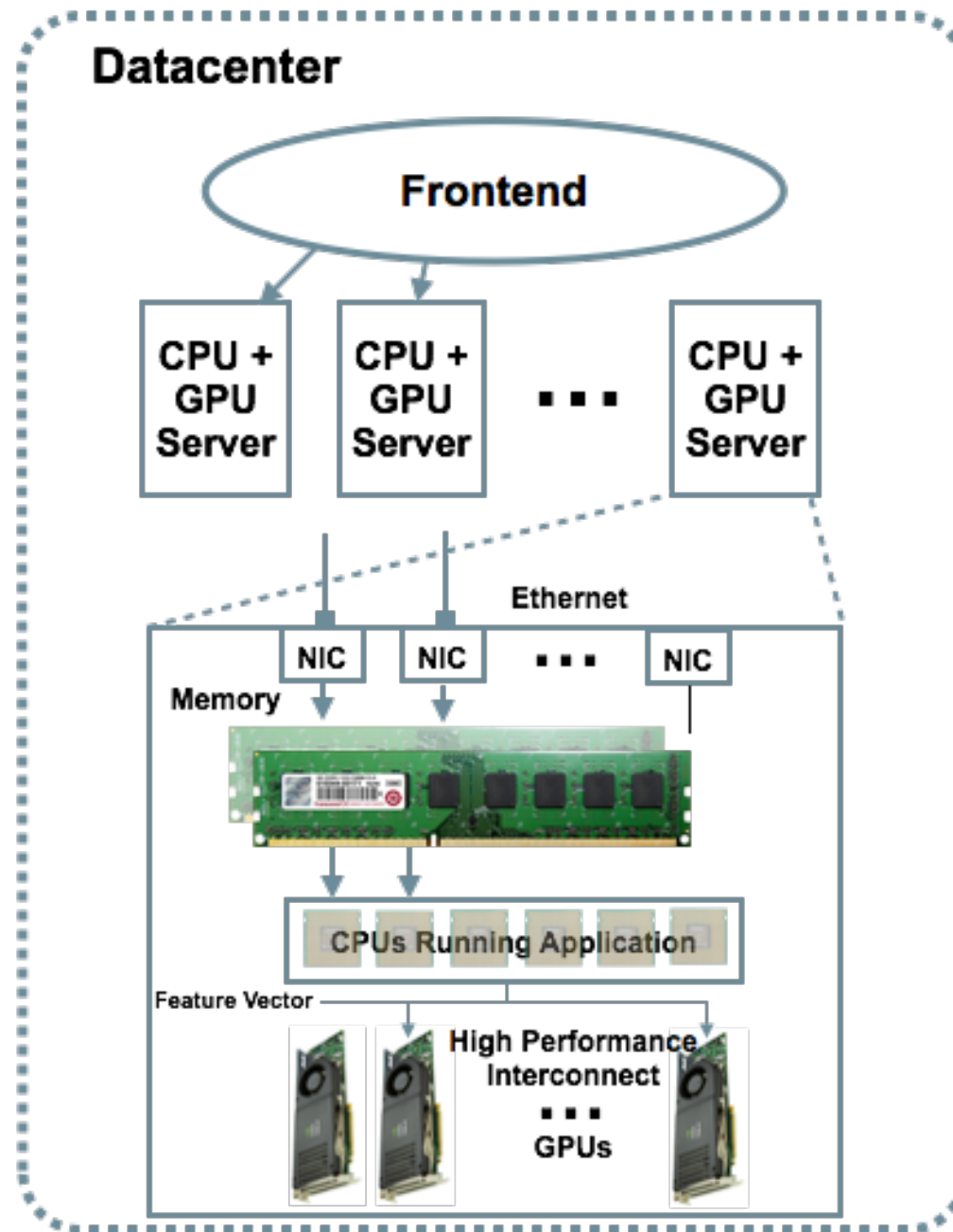CPU-only:
+no extra HW
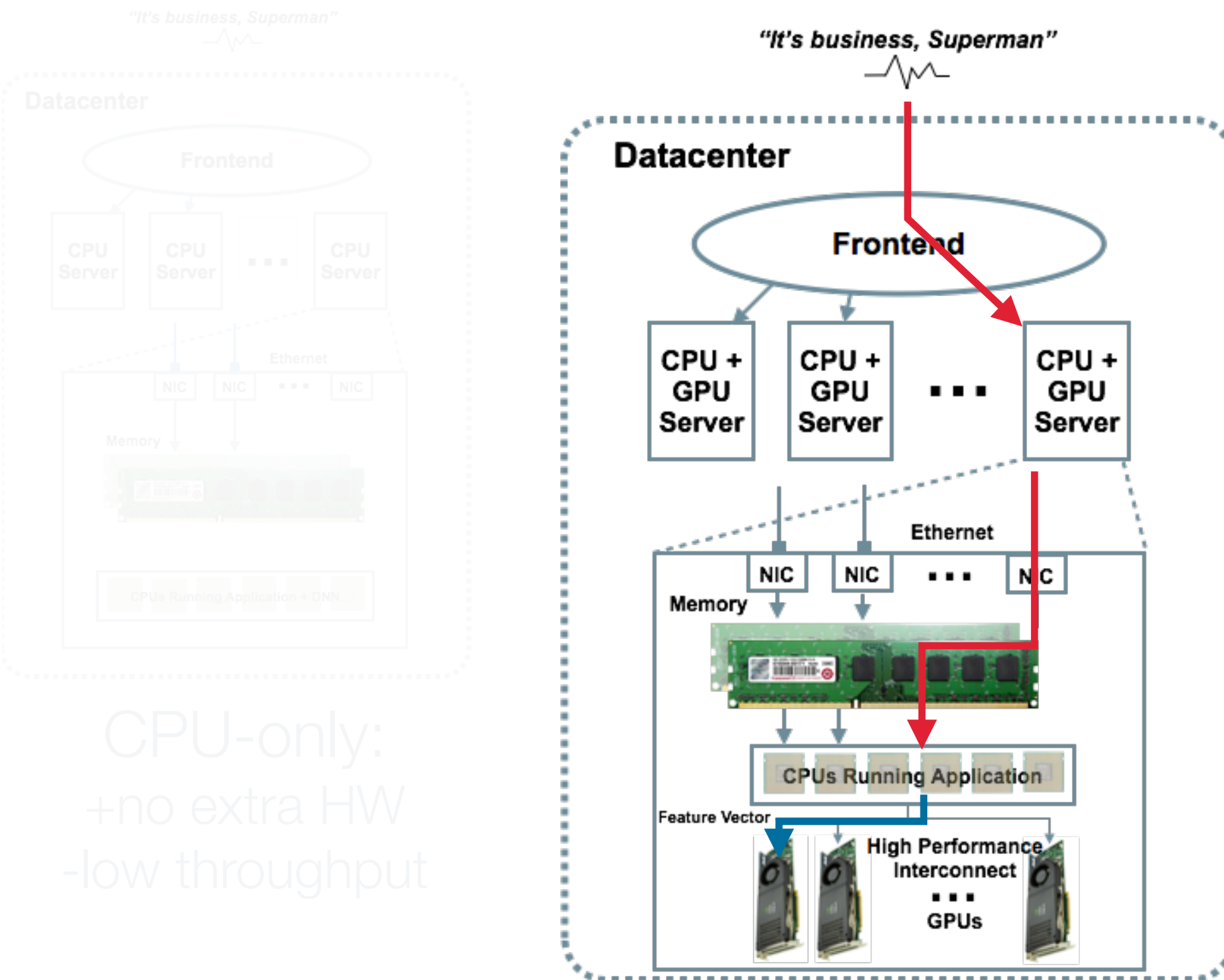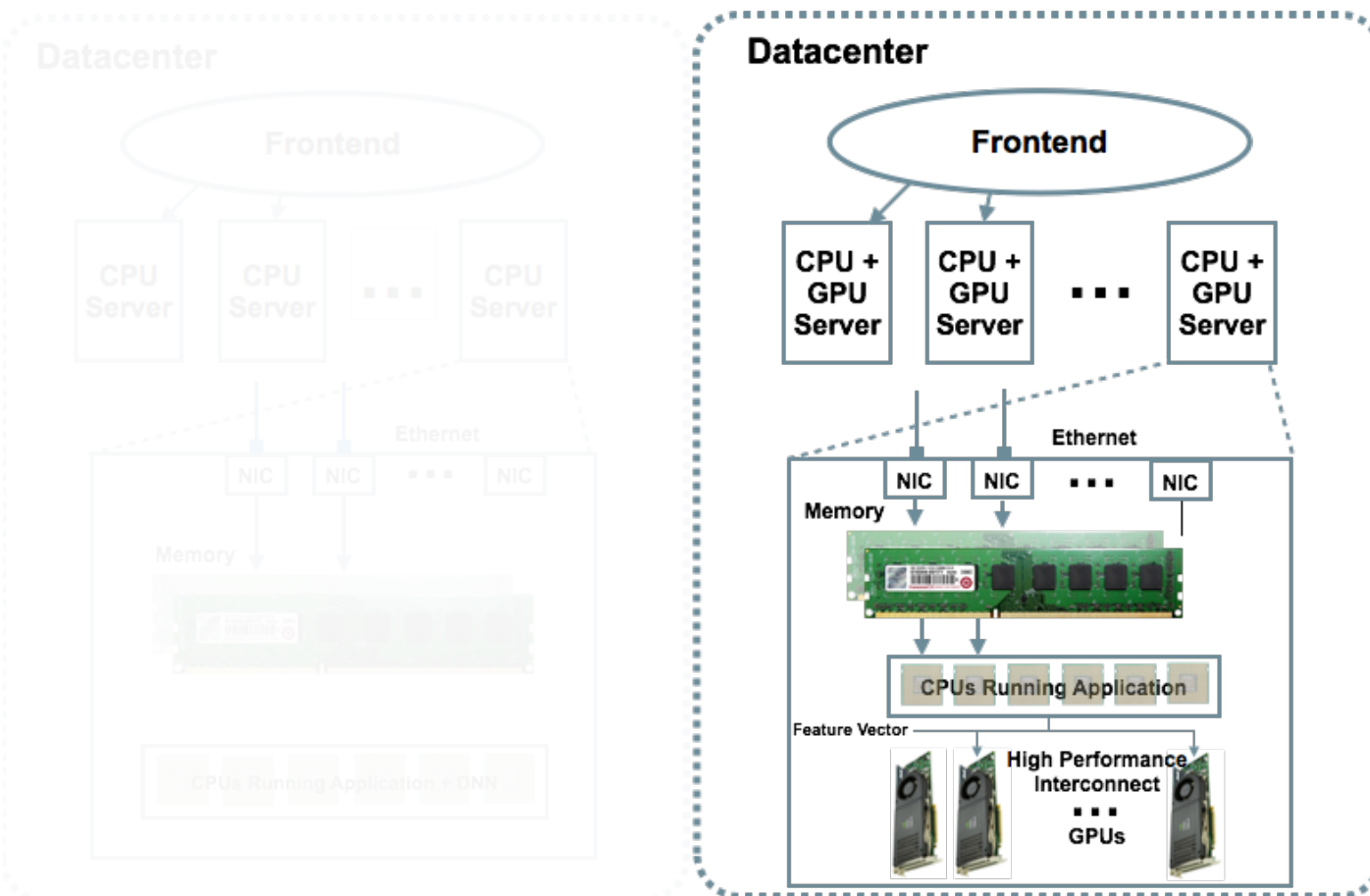-low throughput

# Future WSC Designs — Server Designs

# Future WSC Designs — Server Designs

# Future WSC Designs — Server Designs



**Integrated GPU:**
+homogeneous
-over-provision GPU

# Future WSC Designs — Server Designs

# Future WSC Designs — Server Designs

# Future WSC Designs — Server Designs



CPU-only:
+no extra HW
-low throughput

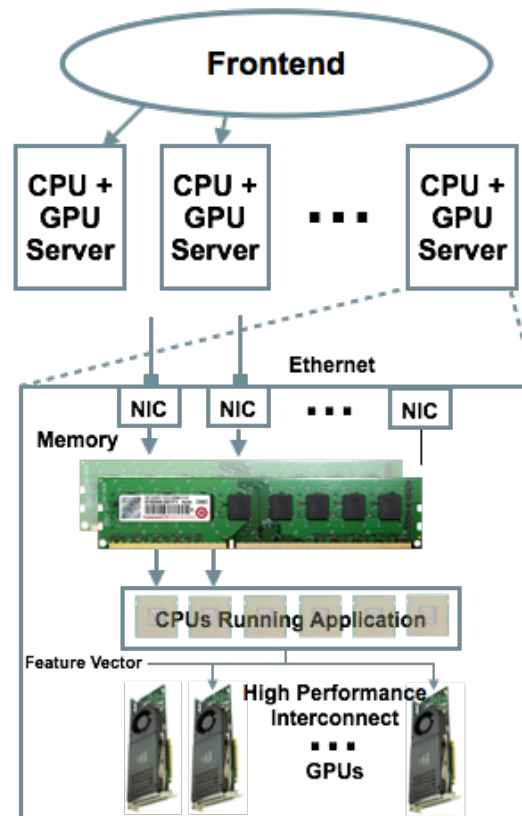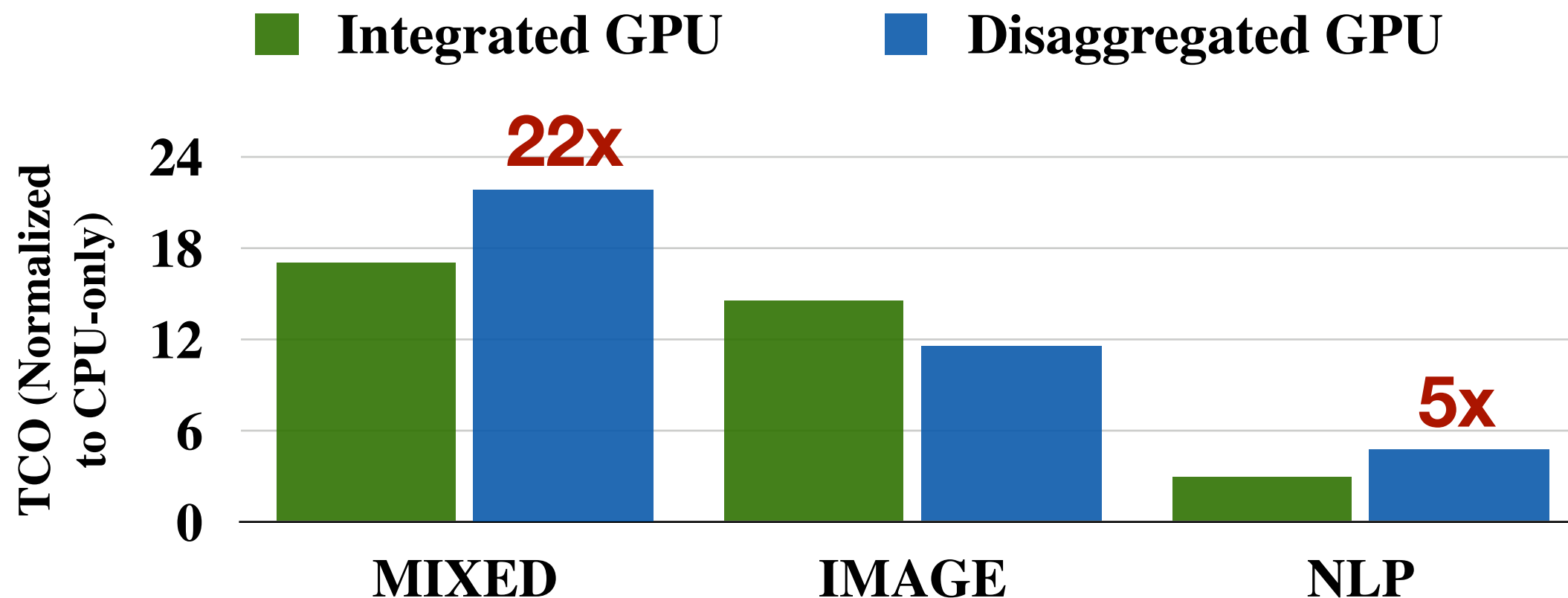Integrated GPU:
+homogeneous
-over-provision GPU

Disaggregated GPU:
+decouple CPU/GPU
-data transfer

# Future WSC Designs — Total Cost of Ownership (TCO)

- Expand Barroso [1] model with GPU and networking costs

- Addressing the bandwidth bottleneck further improves NLP TCO (more details in paper)



Legend: ■ Integrated GPU (green)  ■ Disaggregated GPU (blue)

Bar chart with y-axis "TCO (Normalized to CPU-only)" ranging 0 to 24, x-axis categories MIXED, IMAGE, NLP. Annotations: **22x** above MIXED disaggregated bar, **5x** above NLP disaggregated bar.

[1] Barroso, Luiz André, et. al. "The datacenter as a computer: An introduction to the design of warehouse-scale machines."

# WSC Scale Key Insights

- TCO improvement dependent on the DNN workload composition

  - Bandwidth constrained workloads underutilize GPUs in Integrated design

  - Sufficient bandwidth is critical to fully utilize resources

DjiNN & Tonic

ClarityLab

# Conclusion

- Unified, **state-of-the-art**, highly **optimized** DNN as a Service

- DNNs have **different** compute and communication characteristics

  - **Do not benefit equally** from optimizations

  - Characteristics **impact** system designs

DjiNN
& Tonic

ClarityLab

# Thank you



**DjiNN & Tonic**

# djinn.clarity-lab.org