# Advanced RDMA-based Admission Control for Modern Data-Centers

Ping Lai    Sundeep Narravula    Karthikeyan Vaidyanathan
Dhabaleswar. K. Panda
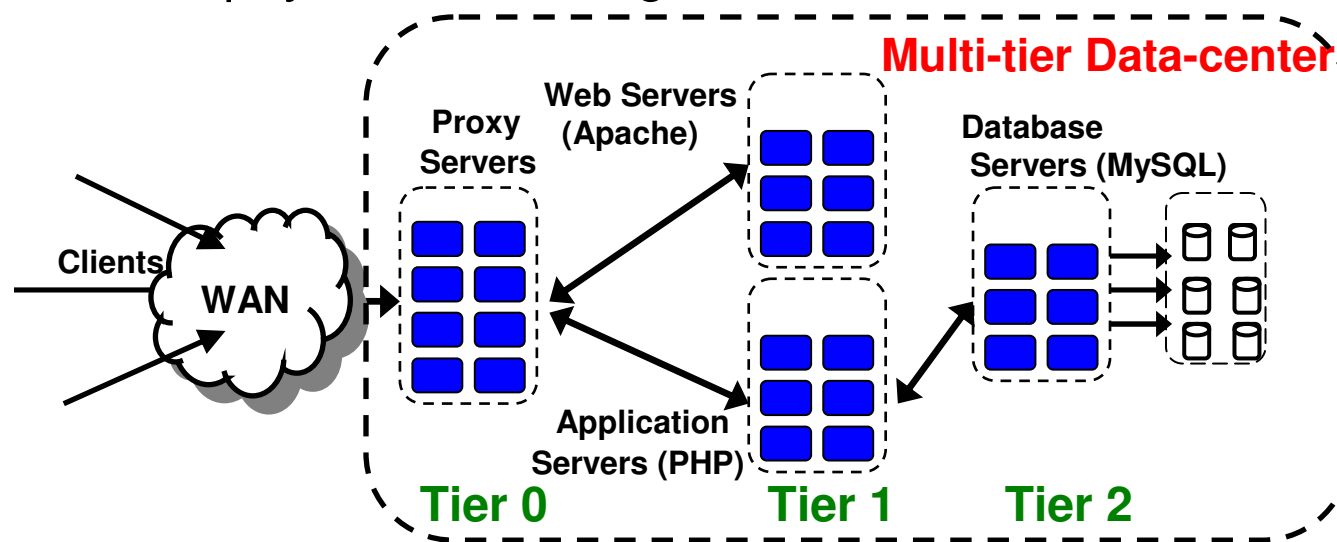
Computer Science & Engineering Department
Ohio State University

# Outline

- **Introduction & Motivation**
- Proposed Design
- Experimental Results
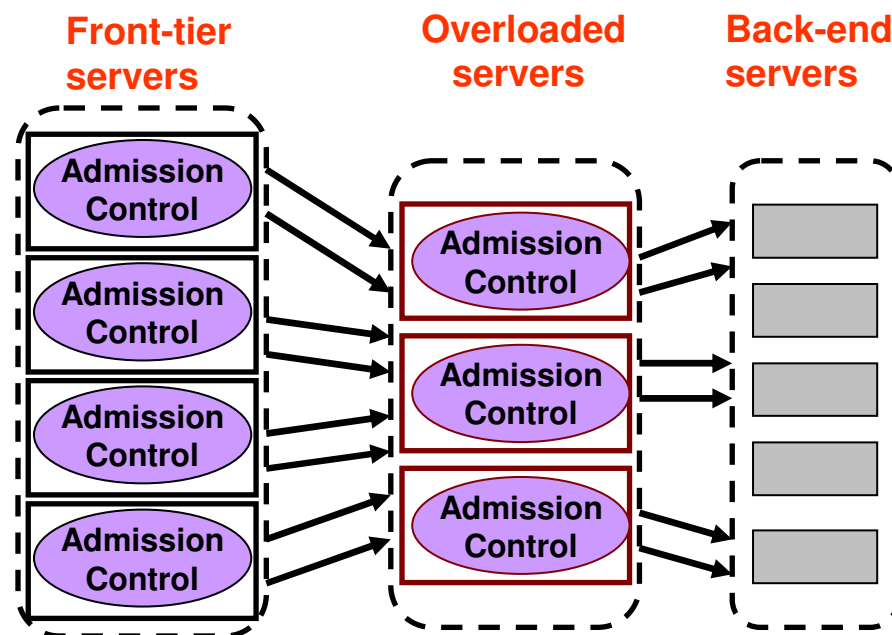- Conclusions & Future Work

# Introduction

- ## Internet grows
  - Number of users, various type of services, huge amount of data
  - Typical apps: e-commerce, bio-informatics, online banking etc.

- ## Web-based multi-tier data-centers
  - Huge bursts of requests →server overloaded
  - Clients pay for service →guaranteed QoS

} Efficient admission control needed!

# General Admission Control

- ## What is admission control?
  - determine whether to accept/drop the incoming requests while guaranteeing the performance (or QoS requirements) of some already existing connections in the overloaded situation

- ## Typical approaches
  - Internal approach: on the overloaded servers
  - External approach: on the front-tier nodes. Main advantages are:
    - Make global decisions
    - More transparent to the overloaded servers
    - Easily applicable to any tier

**Front-tier servers**  **Overloaded servers**  **Back-end servers**



4

# Motivation

- ## External approach
  - Front-tier proxy servers need to get load information from back-end servers

- ## Problems with the existing designs
  - Use TCP/IP – coarse-grained and high overhead; responsiveness depends on load
  - Workload is divergent and unpredictable – require fine-grained and low overhead

OHIO
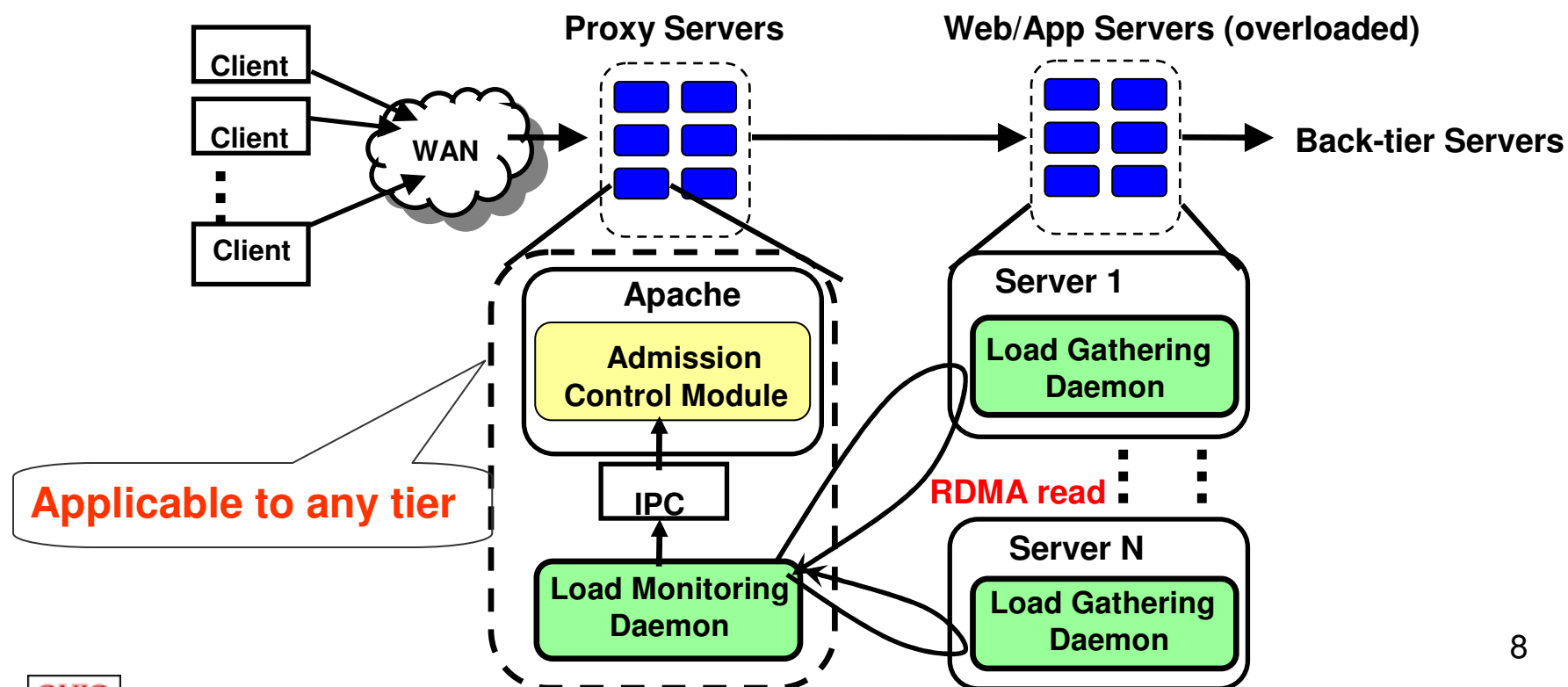STATE

# Opportunity & Objective

- ## Opportunity: modern high-speed interconnects
  - iWARP/10-Gigabit Ethernet, InfiniBand, Quadrics etc.
  - High performance: low latency & high bandiwidth
  - Novel features: atomic operation, protocol offloading, RDMA operations etc.
  - RDMA: low latency & no communication overhead on the remote node

- ## **Objective**
  - Leverage the advanced features (RDMA operation) to design more efficient, lower overhead and better QoS guaranteed admission control

OHIO
STATE

# Outline

- Introduction & Motivation
- **Proposed Design**
- Experimental Results
- Conclusions & Future Work

OHIO
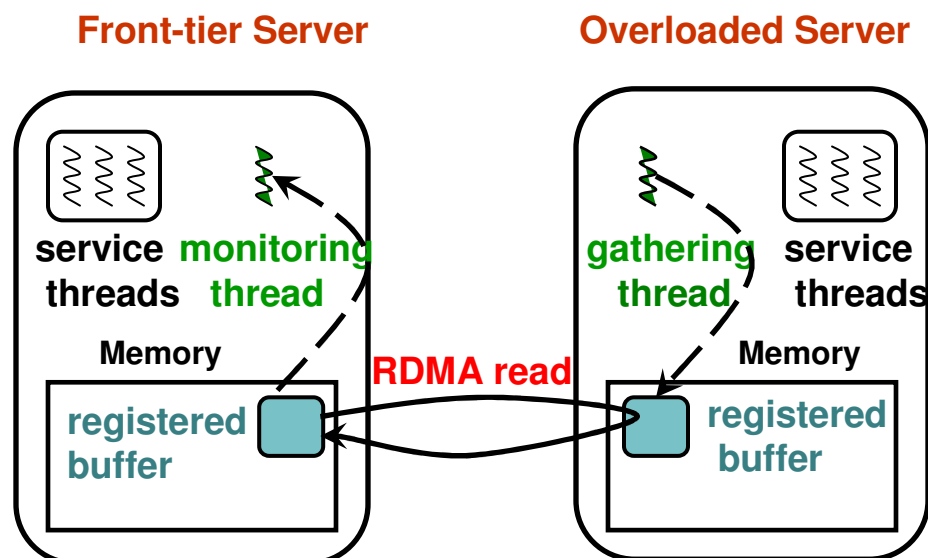STATE

# System Architecture

- Load gathering daemon running on overloaded web servers
- Load monitoring daemon running on front-tier proxy servers
- Admission control module running on front-tier proxy servers

**Client**

**Client**

**Client**

**WAN**

**Proxy Servers**

**Web/App Servers (overloaded)**

**Back-tier Servers**

**Apache**

**Admission Control Module**

**IPC**

**Load Monitoring Daemon**

**Applicable to any tier**

**Server 1**

**Load Gathering Daemon**

**RDMA read**

**Server N**

**Load Gathering Daemon**

8

# Load Gathering and Monitoring Daemon

- ## Load gathering daemon
  - Running on each of the overloaded servers in background – low overhead
  - Gather instantaneous load information

- ## Load monitoring daemon
  - Running on each of the front-tier proxy servers
  - Retrieve load information from all the load gathering daemons

- ## Communication is important!
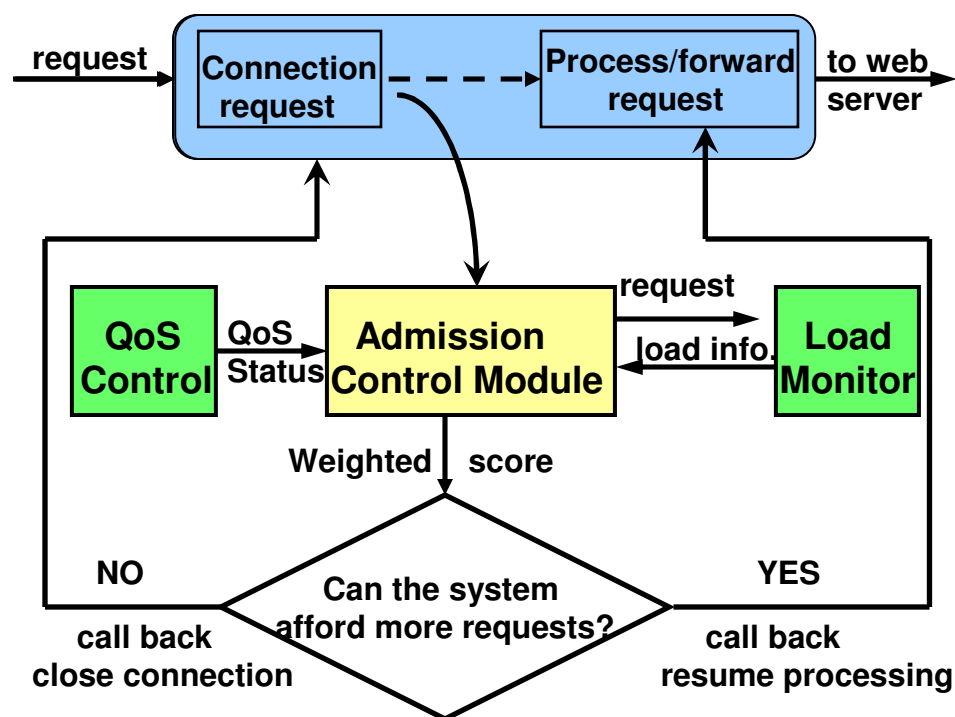  - TCP/IP is not good, so?

OHIO
STATE

# Gathering and Monitoring Daemon Cont.



**Front-tier Server**          **Overloaded Server**

service  **monitoring**          **gathering**  service
threads  **thread**          **thread**  threads

Memory          **RDMA read**          Memory

**registered buffer**          **registered buffer**

- ## Use RDMA read

    – Monitoring daemon issues RDMA read to gathering daemon
        - Buffer must be registered and pinned down before the operation
        - Monitoring daemon has to know the memory address of the remote buffer
    – Retrieve load information at high granularity under overload – better decisions
    – No CPU involvement on the loaded servers – low overhead

10

# Admission Control Module

- Use *shared memory* to communicate with load monitoring daemon
- Attach to Apache: dynamically loadable; trap into Apache request processing

- New processing procedure
  - Apache main thread call the admission control module after TCP connection is established
  - Admission control module uses weighted score to make decisions
  - If all of the back-end servers are overloaded, call back to Apache thread to close the new connections; otherwise, call back to resume the processing



11

# Outline

- Introduction & Motivation
- Proposed Design
- Experimental Results
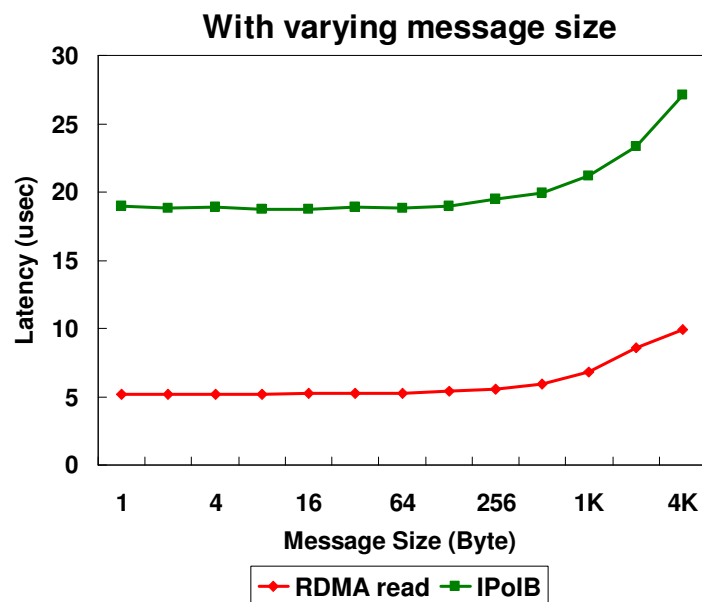- Conclusions & Future Work

OHIO
STATE

# Experimental Platforms

- 32 Compute nodes
  - Dual Intel 64-bit Xeon 3.6 GHz CPU, 2 GB memory
  - Mellanox MT25208 InfiniBand HCA, OFED 1.2 driver
  - Linux 2.6

- Two-tier data-center including proxy servers and web servers; web servers are potentially overloaded
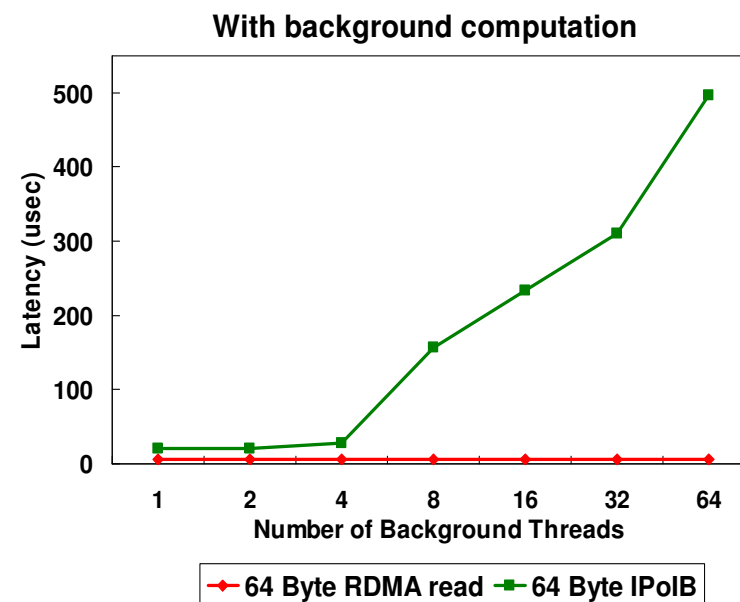
- Apache 2.2.4 for proxy servers and web servers

13

OHIO
STATE

# Experiment Results Outline

- Micro-benchmarks: basic IBA performance

- Data-center level evaluation
  - Single file trace
    - Average response time and aggregate TPS
    - Instant performance analysis
    - QoS analysis
  - Worldcup trace and Zipf trace
    - Worldcup trace: real data from world cup 1998
    - Zipf trace: workloads follow Zipf-like distribution (probability of i'th most popular file $\propto 1/i^{\alpha}$)

OHIO
STATE

# Performance of RDMA read and IPoIB (TCP/IP over IBA)

**With varying message size**

**With background computation**

- 1 Byte message
  - RDMA read: 5.2 us
  - IPoIB: 18.9 us
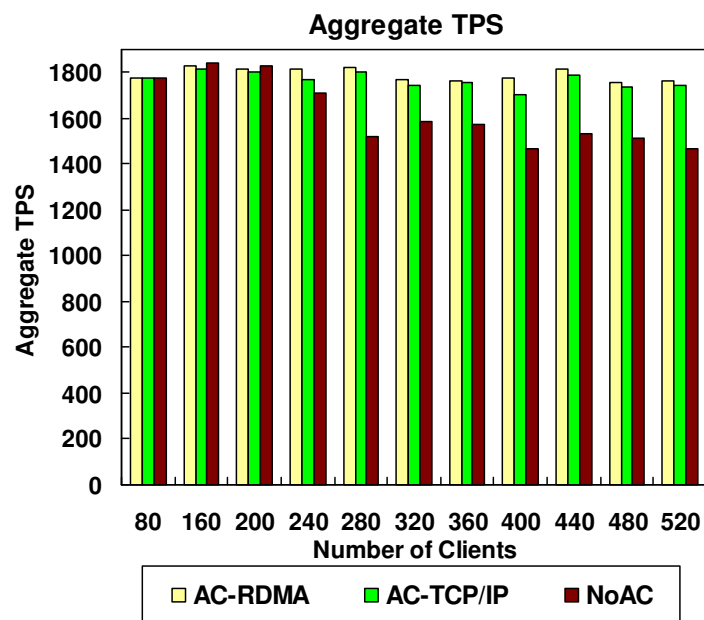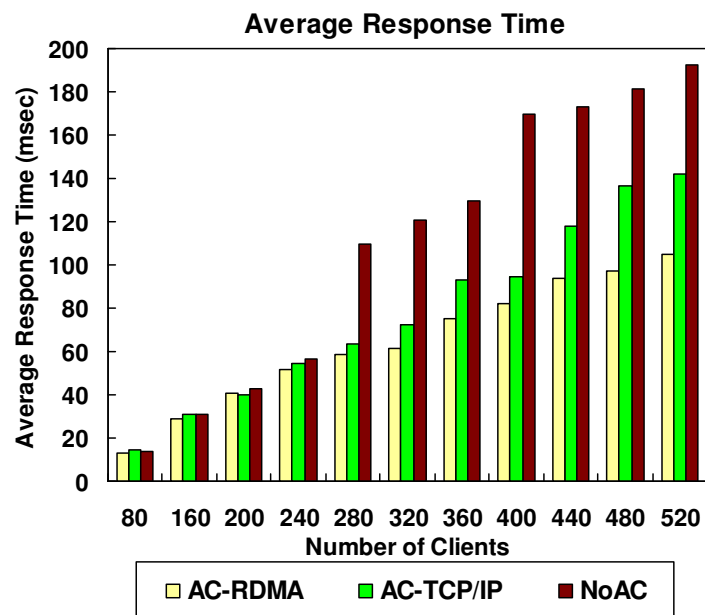- Improvement using RDMA increases when message size increases

- IPoIB significantly degrades
- RDMA read keeps constant latency

Performance of IPoIB depends on load, while RDMA NOT!

15

# Data-Center level Evaluation

- Configuration
  - 4 nodes used as proxy servers
  - 1 node used as web server
  - Remaining nodes are clients

- Load information updated every 1 ms

- Measured average client-perceived response time (for successful request) and aggregate system TPS

- Comparing performance of three systems
  - **AC-RDMA**: system with admission control based on RDMA read (the proposed approach)
  - **AC-TCP/IP**: system with admission control based on TCP/IP
  - **NoAC**: system without any admission control

OHIO
STATE

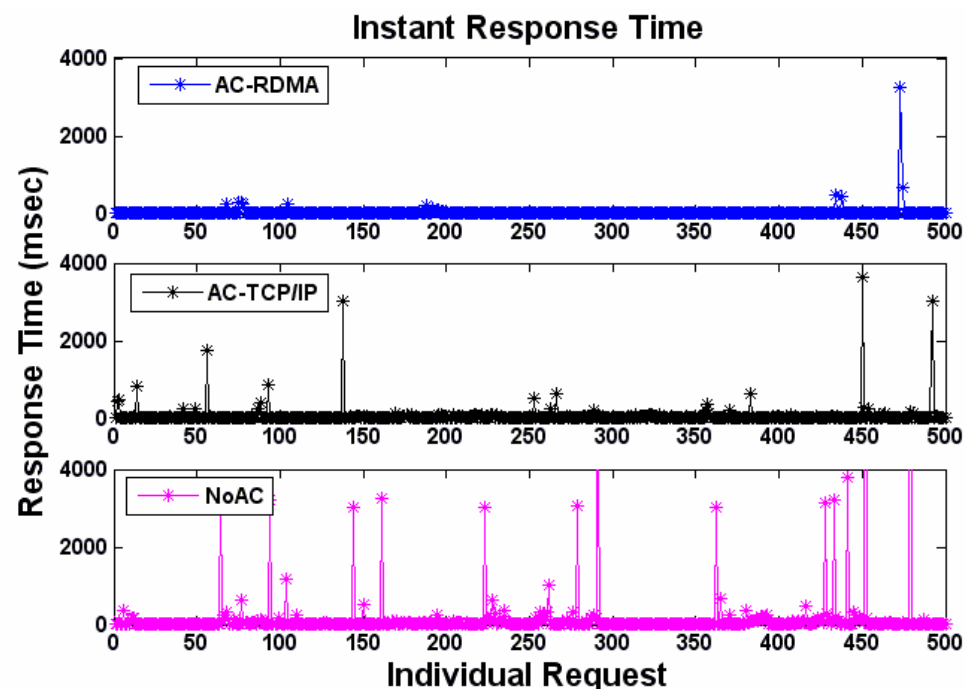# Performance with Single File Trace (16 KB)



- With 520 clients
  - NoAC: 192.31ms
  - AC-TCP/IP: 142.29ms -26% improvement
  - AC-RDMA: 105.03ms - 26% improvement over AC-TCP/IP (45% improvement over NoAC)

- AC-RDMA and AC-TCP/IP are comparable
- System with admission control has higher TPS than the original system
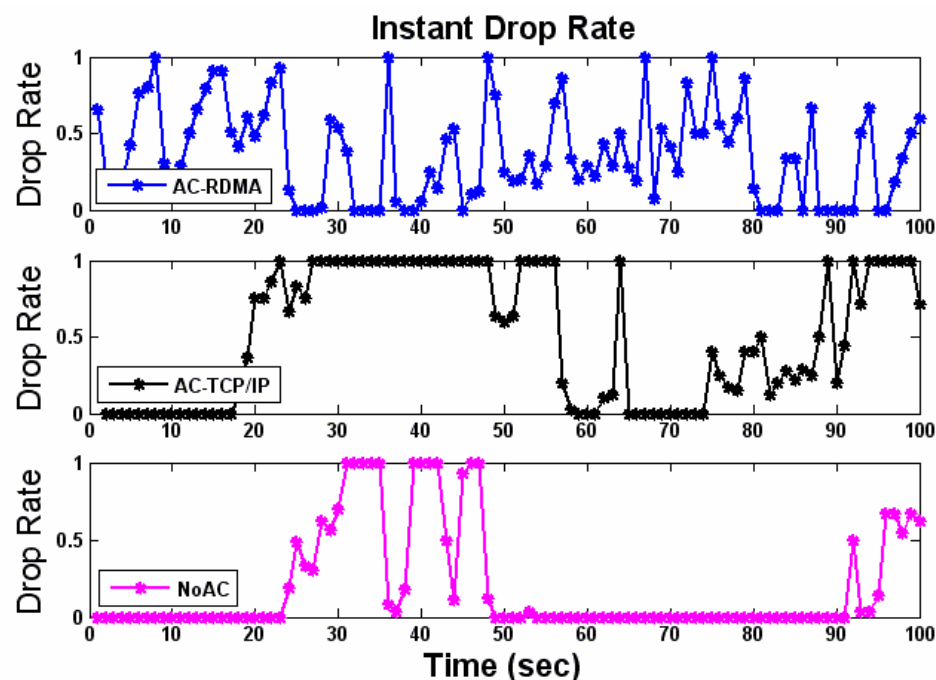
17

# Instant Performance

- ## Workload: 400 clients

- ## Instant response time

  - NoAC: many requests served with very long time

  - AC-RDMA: almost no such requests

  - AC-TCP/IP: some requests with long response time



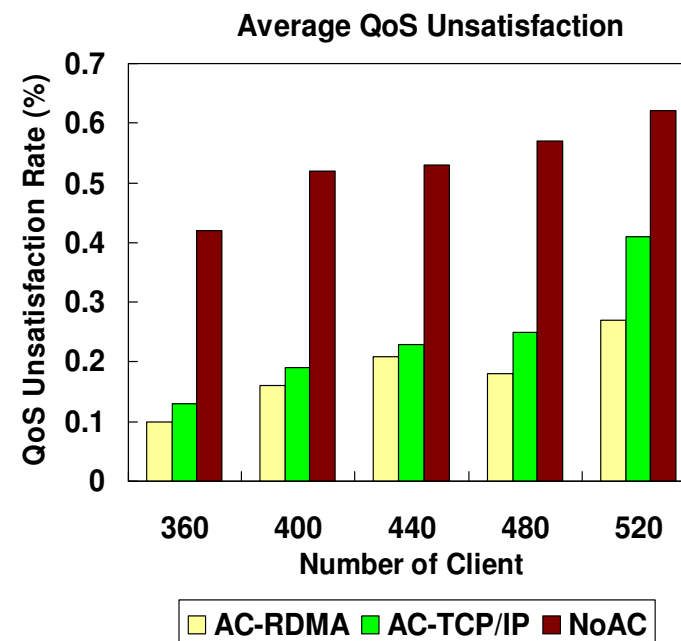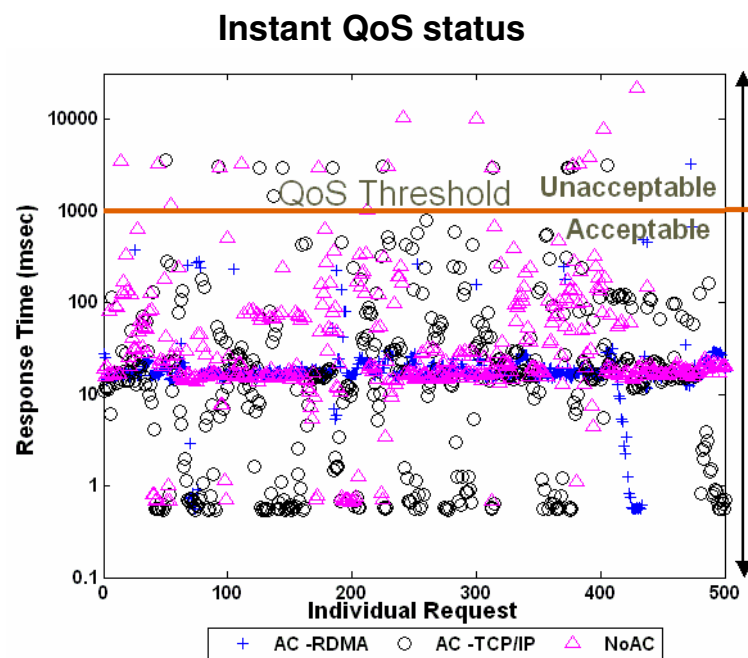Instant performance is consistent with the trend of average response time

# Instant Performance Cont.

- Instant drop rate
  - AC-RDMA: closely reflects the instantaneous changing load on web servers
  - AC-TCP/IP: longer streak of continuous drops or acceptance
  - NoAC: a lot of acceptance; some drops because of timeout



Instant Drop Rate

AC-RDMA gets the load information timely, while AC-TCP/IP sometimes reads the stale information due to the slow response from overloaded servers in TCP/IP communication

19

OHIO STATE

# QoS Analysis

**Instant QoS status**



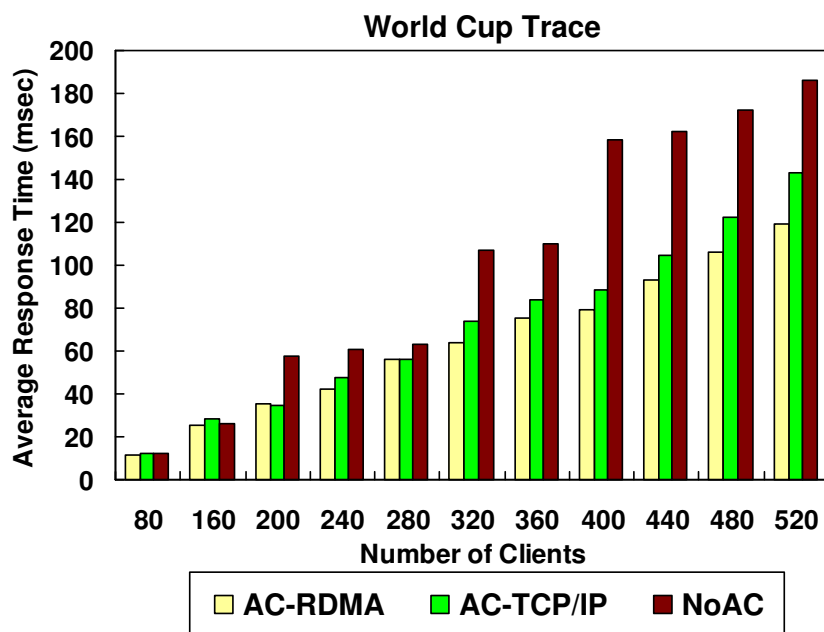**Average QoS Unsatisfaction**



- Instant QoS status
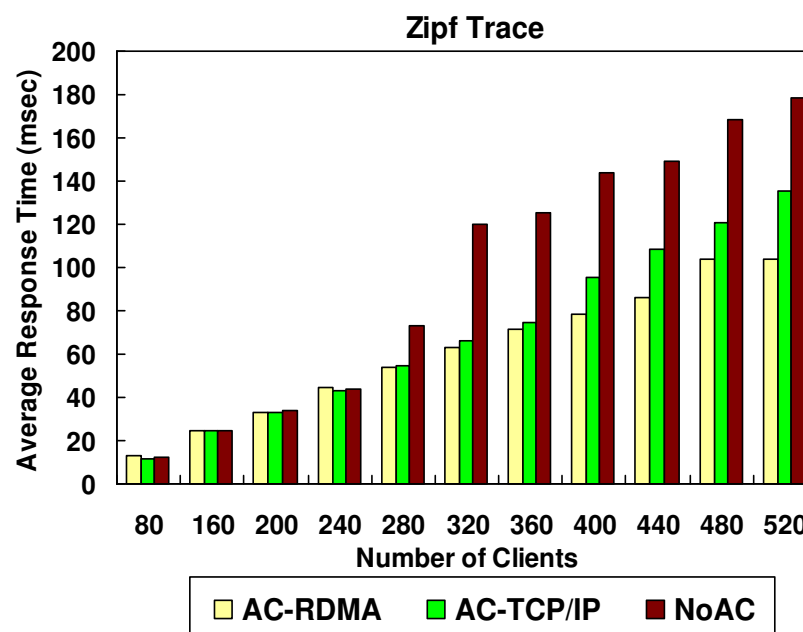  - AC-RDMA has much better capability of satisfying the QoS requirement

- Average QoS unsatisfaction
  - AC-RDMA is the best

With the same requirement of QoS (e.g., response time), AC-RDMA can serve more clients than the other two systems

20

# Performance with Worldcup and Zipf Trace

- ## AC-RDMA is better
  - Compared to AC-TCP/IP: 17%
  - Compared to NoAC: 36%

- ## AC-RDMA is better
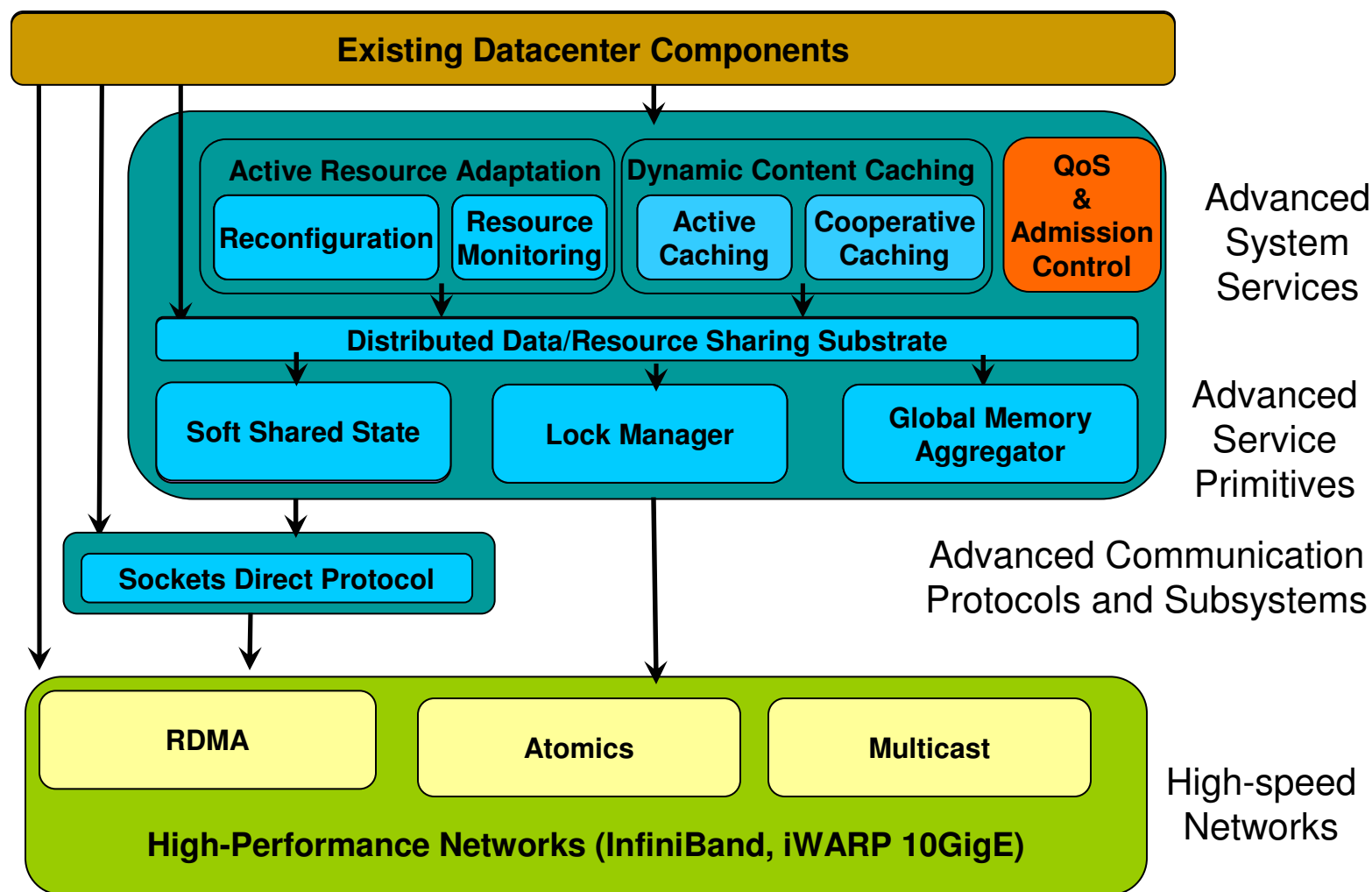  - Compared to AC-TCP/IP: 23%
  - Compared to NoAC: 42%

# Outline

- Introduction & Motivation

- Proposed Design

- Experimental Results

- **Conclusions & Future Work**
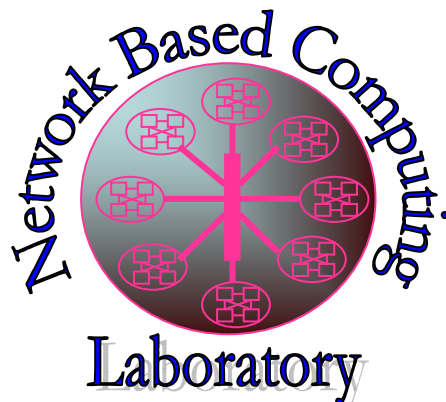
# Conclusions & Future Work

- Leveraged RDMA read in designing efficient admission control mechanism used in multi-tier data-centers

- Implemented the design in a two-tier data-center over InfiniBand

- Evaluated with single file, worldcup trace and Zipf trace
  - AC-RDMA outperforms AC-TCP/IP up to 28%, outperforms NoAC up to 51%
  - AC-RDMA can provide better QoS satisfaction
  - Main reasons
    - Update load information timely
    - No extra overhead on the already overloaded servers

- Future work: study the scalability performance, incorporate other earlier work for integrated resource management service etc.

# Overall Datacenter Framework

**Existing Datacenter Components**

**Active Resource Adaptation**
- Reconfiguration
- Resource Monitoring

**Dynamic Content Caching**
- Active Caching
- Cooperative Caching

**QoS & Admission Control**

Advanced System Services

**Distributed Data/Resource Sharing Substrate**

- Soft Shared State
- Lock Manager
- Global Memory Aggregator

Advanced Service Primitives

**Sockets Direct Protocol**

Advanced Communication Protocols and Subsystems

- RDMA
- Atomics
- Multicast

**High-Performance Networks (InfiniBand, iWARP 10GigE)**

High-speed Networks

24

OHIO STATE

# Thank you

{laipi, narravul, vaidyana, panda}@cse.ohio-state.edu

## NBC-LAB

**Network-Based Computing Laboratory**
http://nowlab.cse.ohio-state.edu/
Data-Center Web Page
http://nowlab.cse.ohio-state.edu/projects/data-centers/index.html

OHIO STATE