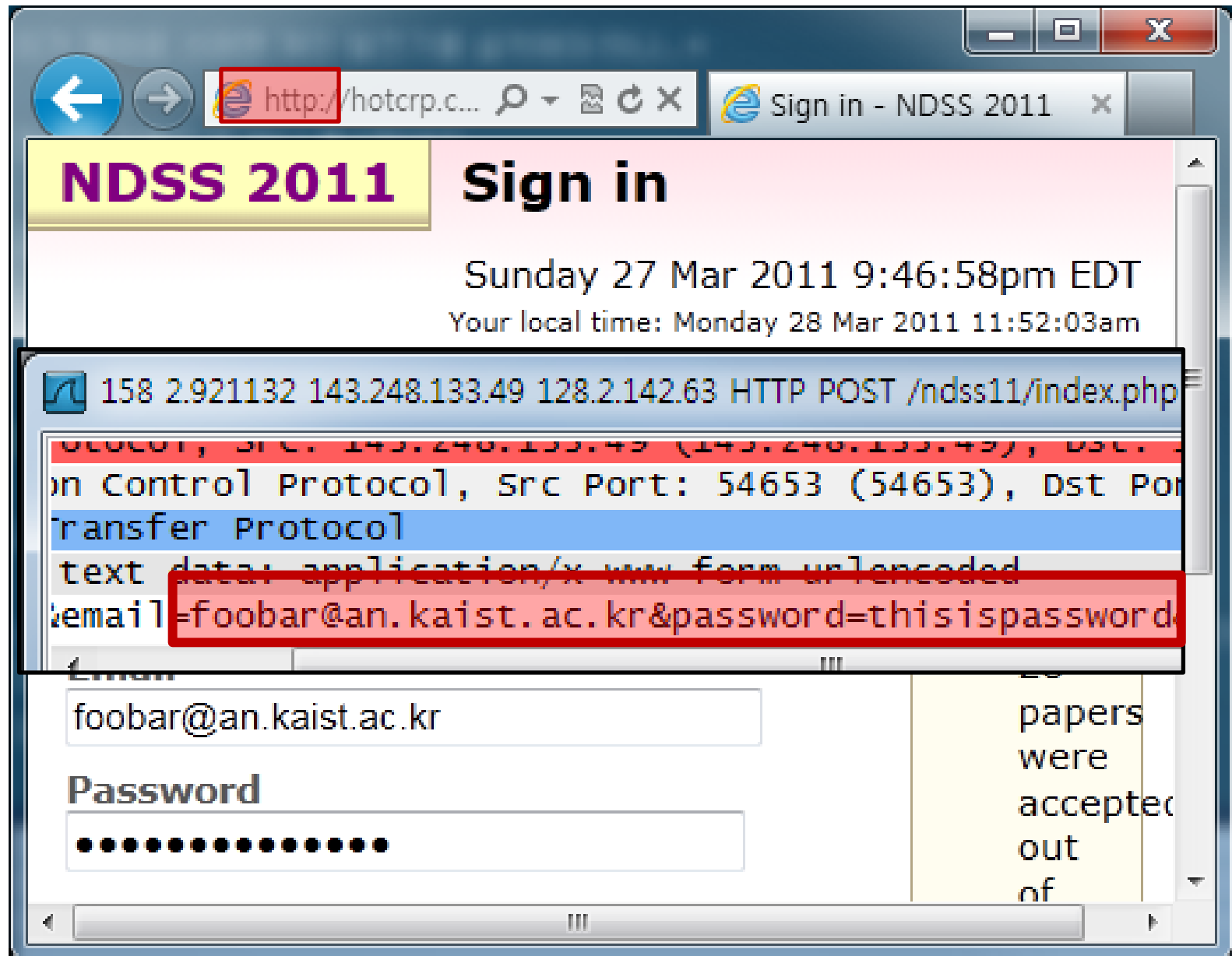# SSLShader: Cheap SSL Acceleration with Commodity Processors

**Keon Jang**[+], Sangjin Han[+], Seungyeop Han[*],
Sue Moon[+], and KyoungSoo Park[+]

KAIST[+] and University of Washington[*]
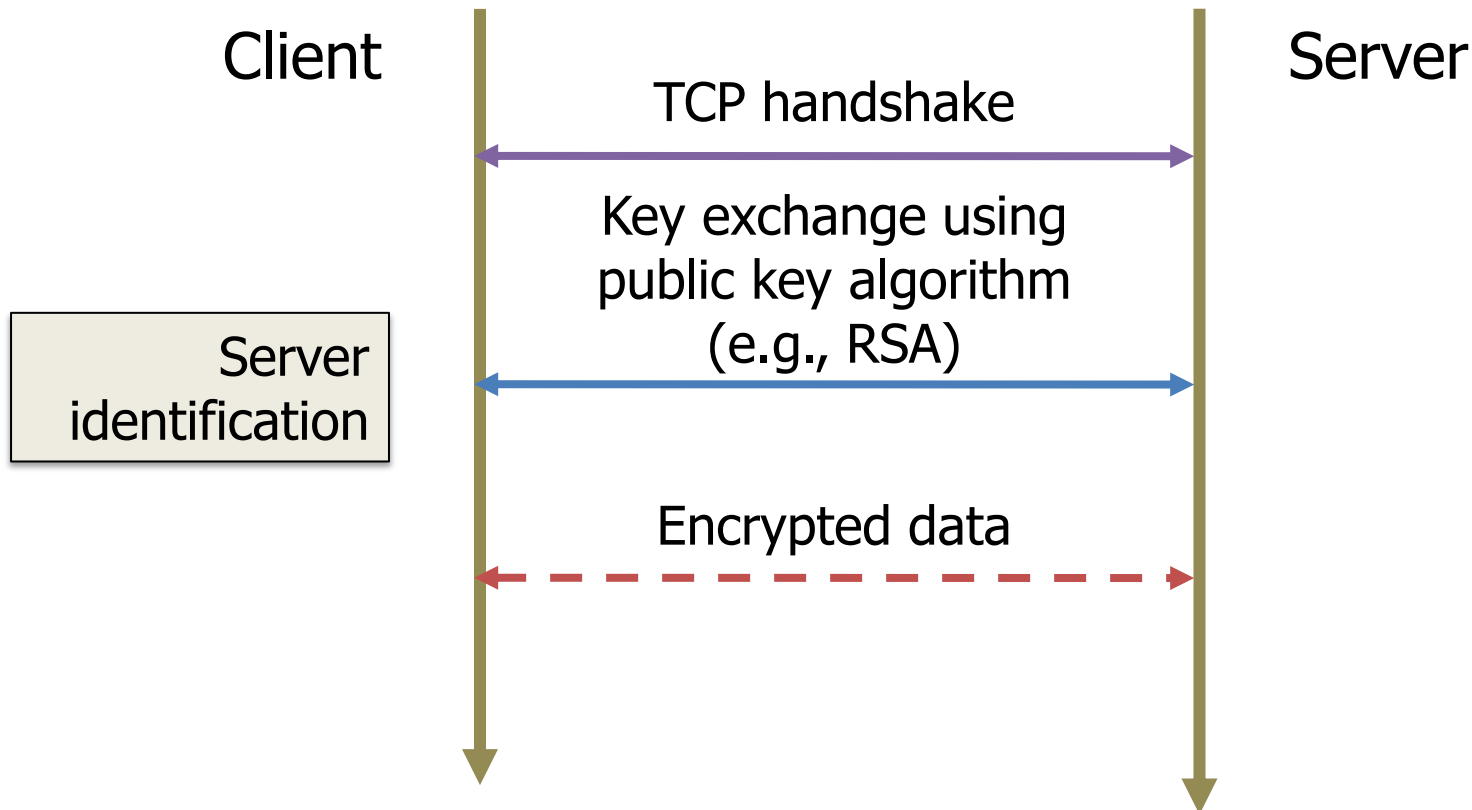
# Security of Paper Submission Websites

# Security Threats in the Internet

- Public WiFi without encryption
  - Easy target that requires almost no effort

- Deep packet inspection by governments
  - Used for censorship
  - In the name of national security

- NebuAd's targeted advertisement
  - Modify user's Web traffic in the middle

# Secure Sockets Layer (SSL)

- A de-facto standard for secure communication
  - Authentication, Confidentiality, Content integrity



Client                                    Server

TCP handshake

Key exchange using
public key algorithm
(e.g., RSA)

Server
identification

Encrypted data

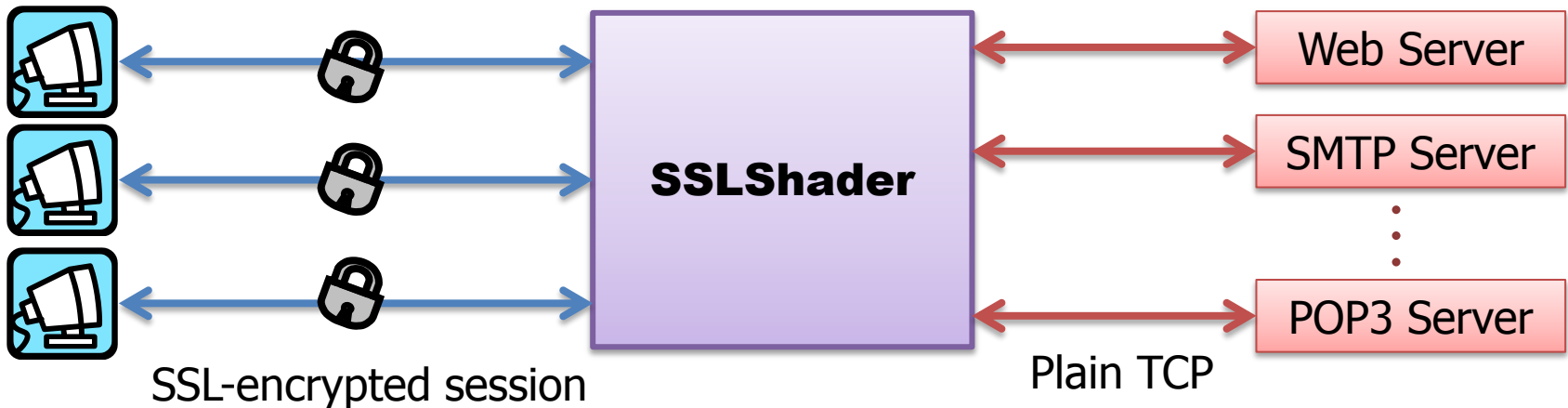# SSL Deployment Status

- Most of Web-sites are not SSL-protected
  - Less than <span style="color:red">0.5%</span>
    - [NETCRAFT Survey Jan '09]

- Why is SSL not ubiquitous?
  - Small sites: lack of recognition, manageability, etc.
  - <span style="color:red">Large sites: cost</span>
    - SSL requires lots of computation power

# SSL Computation Overhead

- Performance overhead (HTTPS vs. HTTP)
  - Connection setup           22x

  - Data transfer               50x

- Good privacy is expensive
  - More servers
  - H/W SSL accelerators

- Our suggestion:
  - Offload SSL computation to GPU

# SSLShader

- SSL-accelerator leveraging GPU
  - High-performance
  - Cost-effective

- SSL reverse proxy
  - No modification on existing servers



SSL-encrypted session

SSLShader

Plain TCP

Web Server
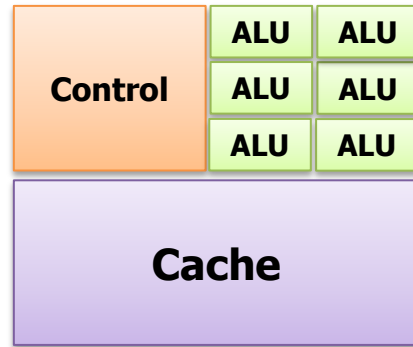
SMTP Server

POP3 Server

# Our Contributions

- GPU cryptography optimization
  - The fastest RSA on GPU
  - Superior to high-end hardware accelerators
  - Low latency

- SSLShader
  - Complete system exploiting GPU for SSL processing
    - Batch processing
    - Pipelining
    - Opportunistic offloading
    - Scaling with multiple cores and NUMA nodes

# CRYPTOGRAPHIC PROCESSING WITH GPU

# How GPU Differs From CPU?

| Control | ALU | ALU |
|---------|-----|-----|
|         | ALU | ALU |
|         | ALU | ALU |

**Cache**

**ALU**

Intel Xeon 5650 CPU:
6 cores

NVIDIA GTX580 GPU:
512 cores

Instructions / sec

$62 \times 10^9$  <  $870 \times 10^9$

# Single Instruction Multiple Threads (SIMT)

Example code: vector addition (C = A + B)

CPU code

GPU code

```
void VecAdd(
int *A, int *B, int *C, int N)
{
    //iterate over N elements
    for(int i = 0; i < N; i++)
        C[i] = A[i] + B[i]

}


VecAdd(A, B, C, N);
```
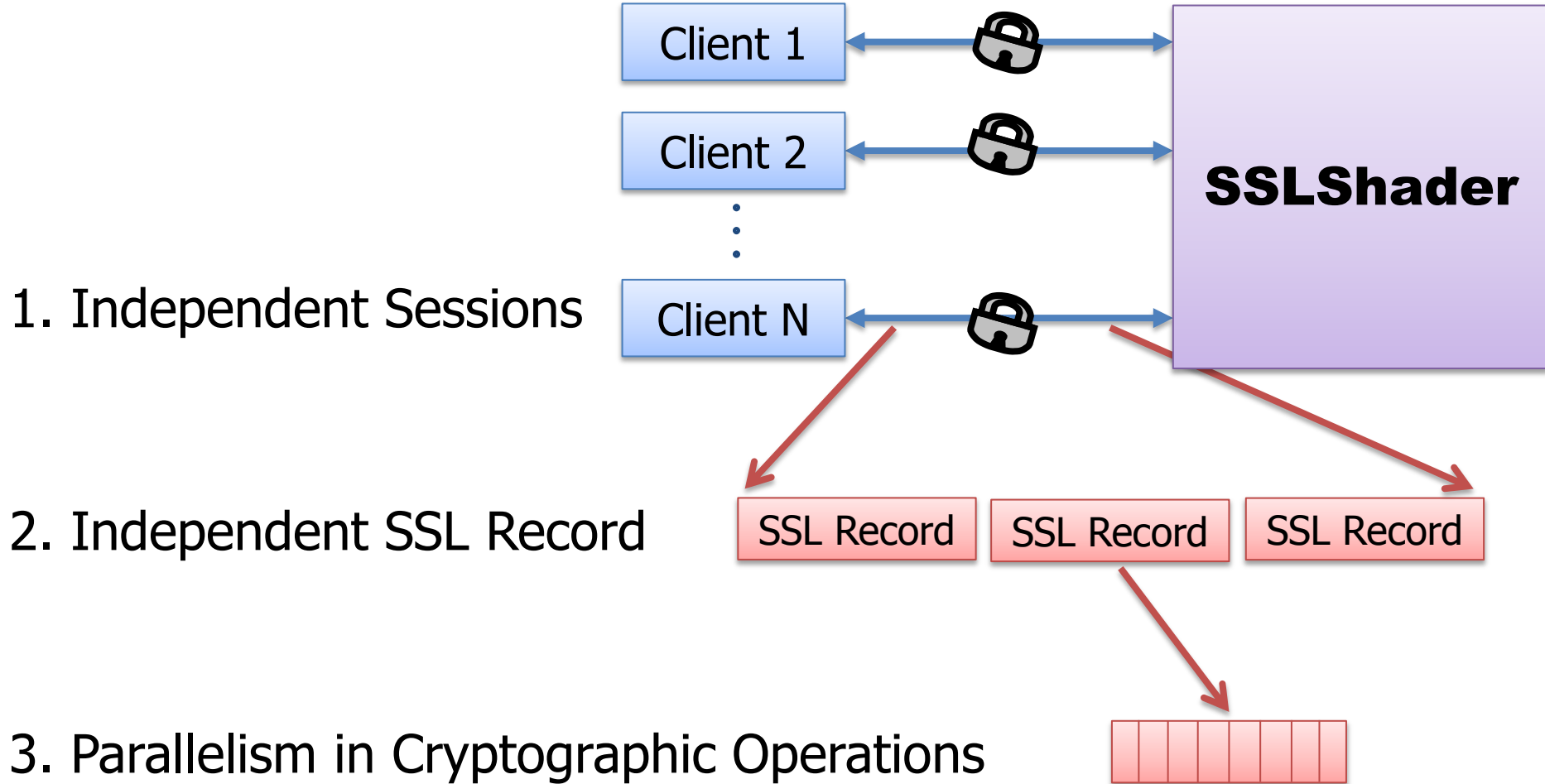
```
__global__ void VecAdd(
int *A, int *B, int *C)
{
    int i = threadIdx.x;
    C[i] = A[i] + B[i]
}

//Launch N threads
VecAdd<<<1, N>>>(A, B, C);
```
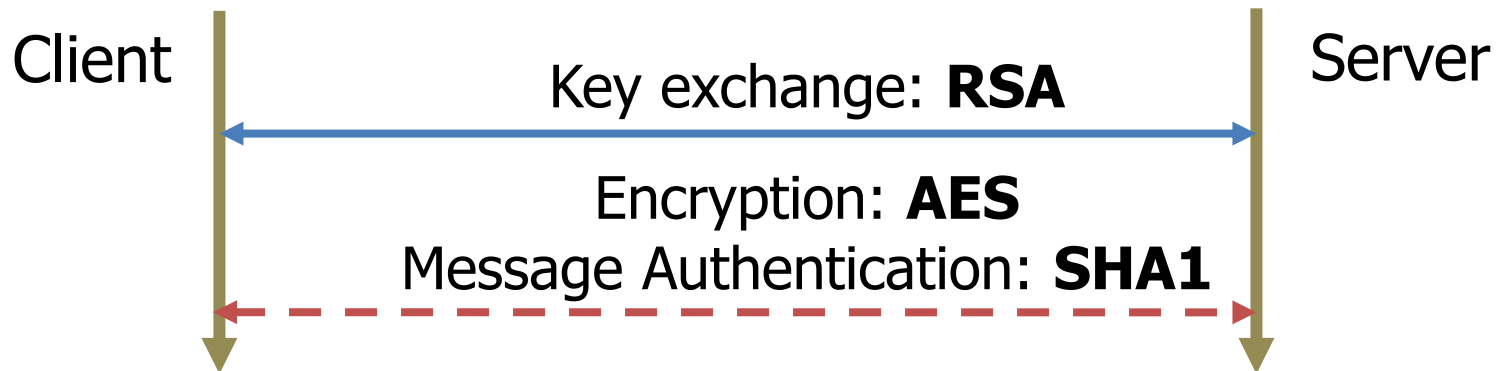
# Parallelism in SSL Processing

Client 1

Client 2

⋮

Client N

**SSLShader**

1. Independent Sessions

2. Independent SSL Record

SSL Record    SSL Record    SSL Record

3. Parallelism in Cryptographic Operations

# Our GPU Implementation

- Choices of cipher-suite

Client              Server

Key exchange: **RSA**

Encryption: **AES**
Message Authentication: **SHA1**

- Optimization of GPU algorithms
  - Exploiting massive parallel processing
    - Parallelization of algorithms
    - Batch processing
  - Data copy overhead is significant
    - Concurrent copy and execution

# Basic RSA Operations

- $M$: plain-text, $C$: cipher-text
- $(e, n)$: public key, $(d, n)$: private key

- Encryption:
  →Client

Small number: 3, 17, 65537

$$C = M^e \bmod n$$

1024/2048 bits integer (300 ~ 600 digits)

- Decryption:
  →Server

$$M = C^d \bmod n$$

**Exponentiation → many multiplications**

# Breakdown of Large Integer Multiplication

Schoolbook
multiplication

Accumulation is difficult to parallelize due to

**"overlapping digits"**

**"carry propagation"**

```
        649
  x     627
----------
         63
        280
       4200
       1800
       800
      12000
      5400
     32000
 +  360000
----------
     406923
```
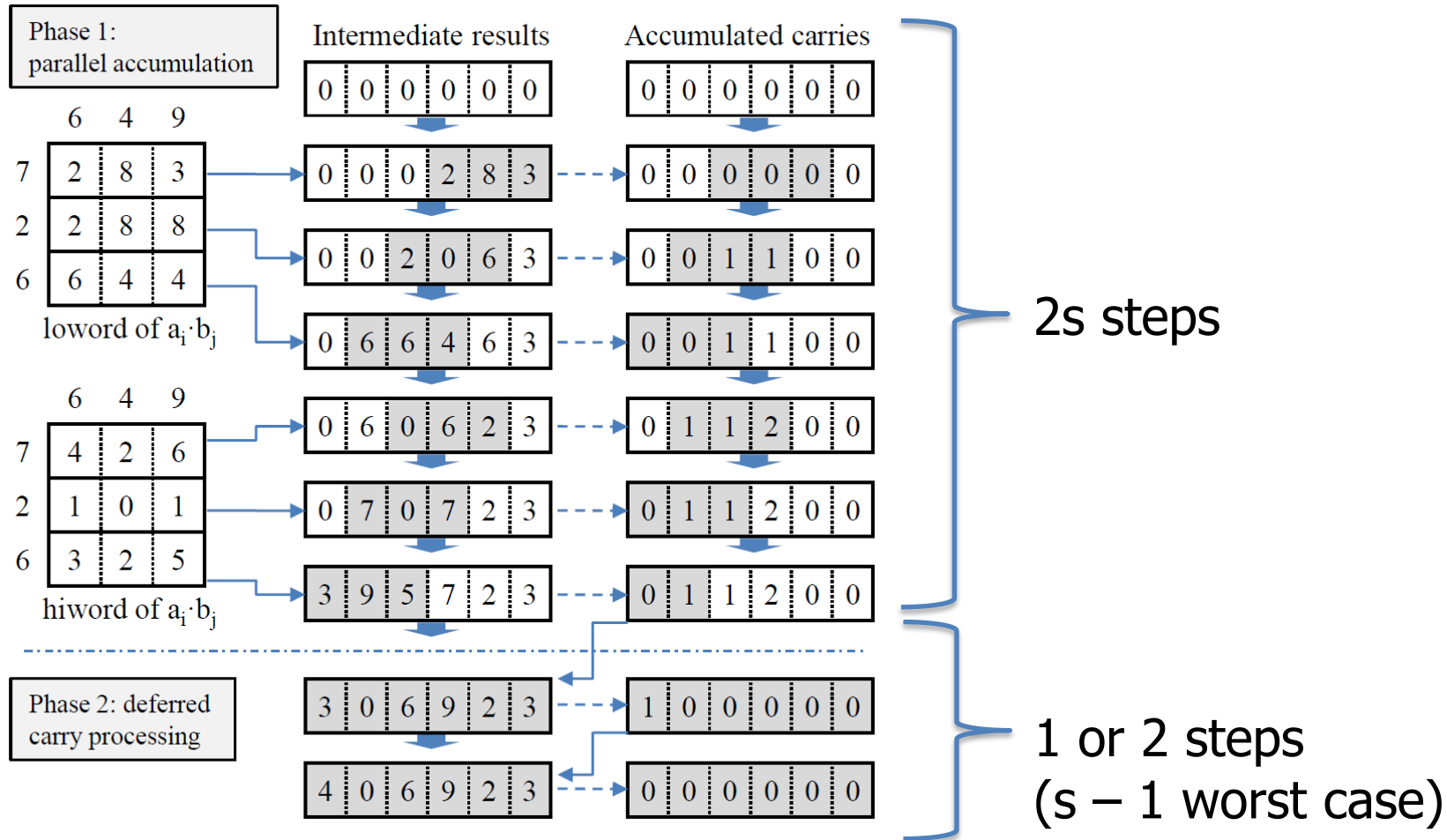
3 x 3 = 9 multiplications
9 addition of 6-digits integers

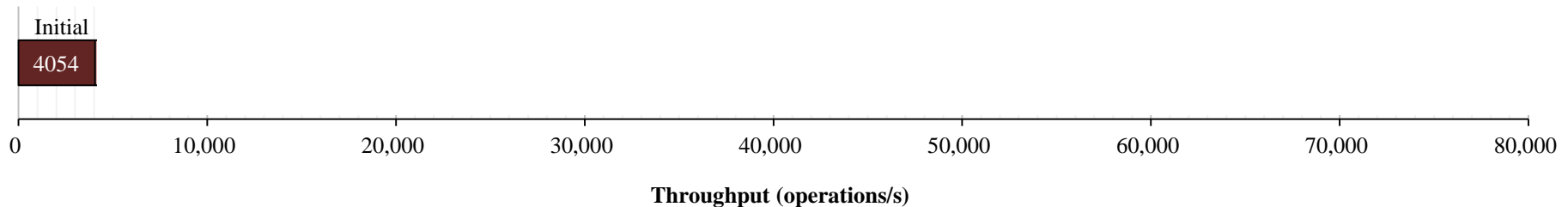# *O(s)* **Parallel Multiplications**

s = # of words in a large integer
(E.g., 1024-bits = 16 x 64 bits word)

Example of
$649 \times 627 = 406{,}923$



2s steps

1 or 2 steps
(s − 1 worst case)

# More Optimizations on RSA

Initial
4054

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 0 | 10,000 | 20,000 | 30,000 | 40,000 | 50,000 | 60,000 | 70,000 | 80,000 |

**Throughput (operations/s)**

- Common optimizations for RSA
  - Chinese Remainder Theorem (CRT)
  - Montgomery Multiplication

## Read our paper for details ☺

  - Faster calculation of M×n
  - Interleaving of T + M×n
  - Mixed-Radix Conversion Offloading
- GPU specific optimizations
  - Warp Utilization
  - Loop Unrolling
  - Elimination of Divergence
  - Avoiding Bank Conflicts
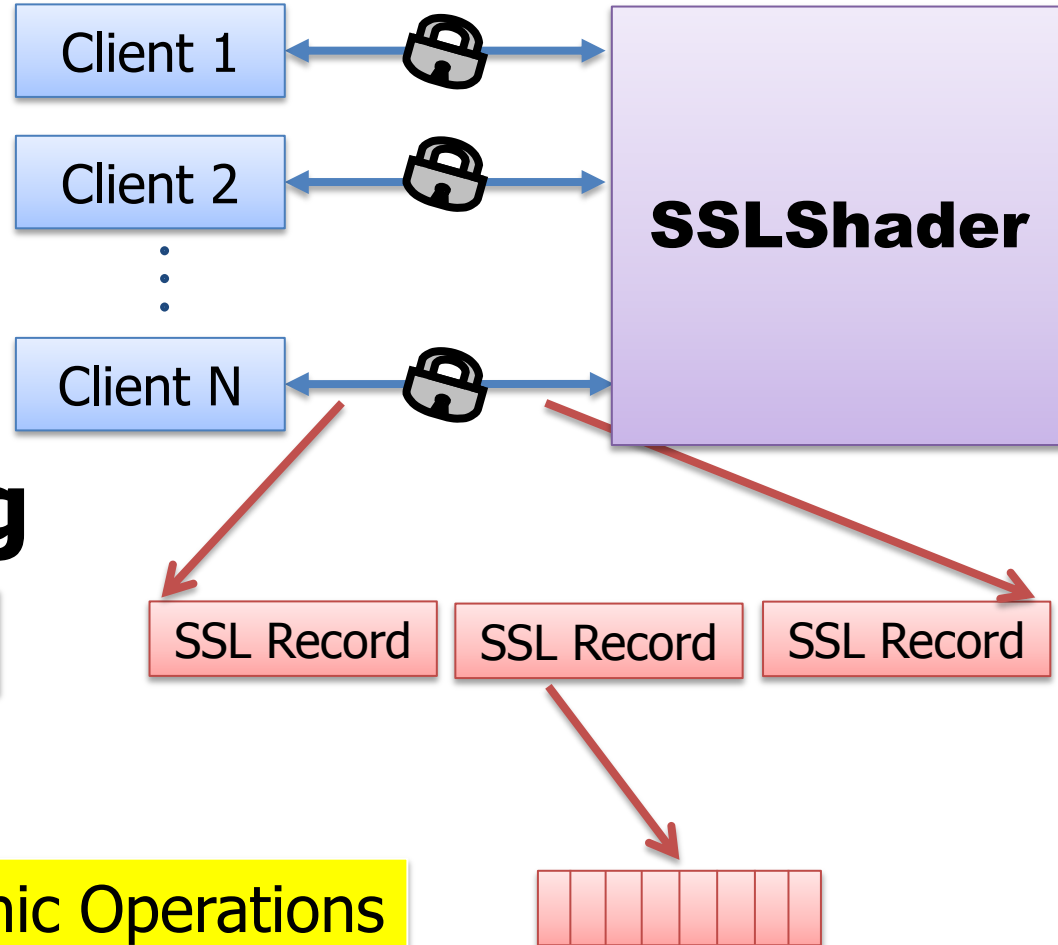  - Instruction-Level Optimization

# Parallelism in SSL Processing

Client 1

Client 2

⋮

Client N

**SSLShader**

1. Independent Sessions

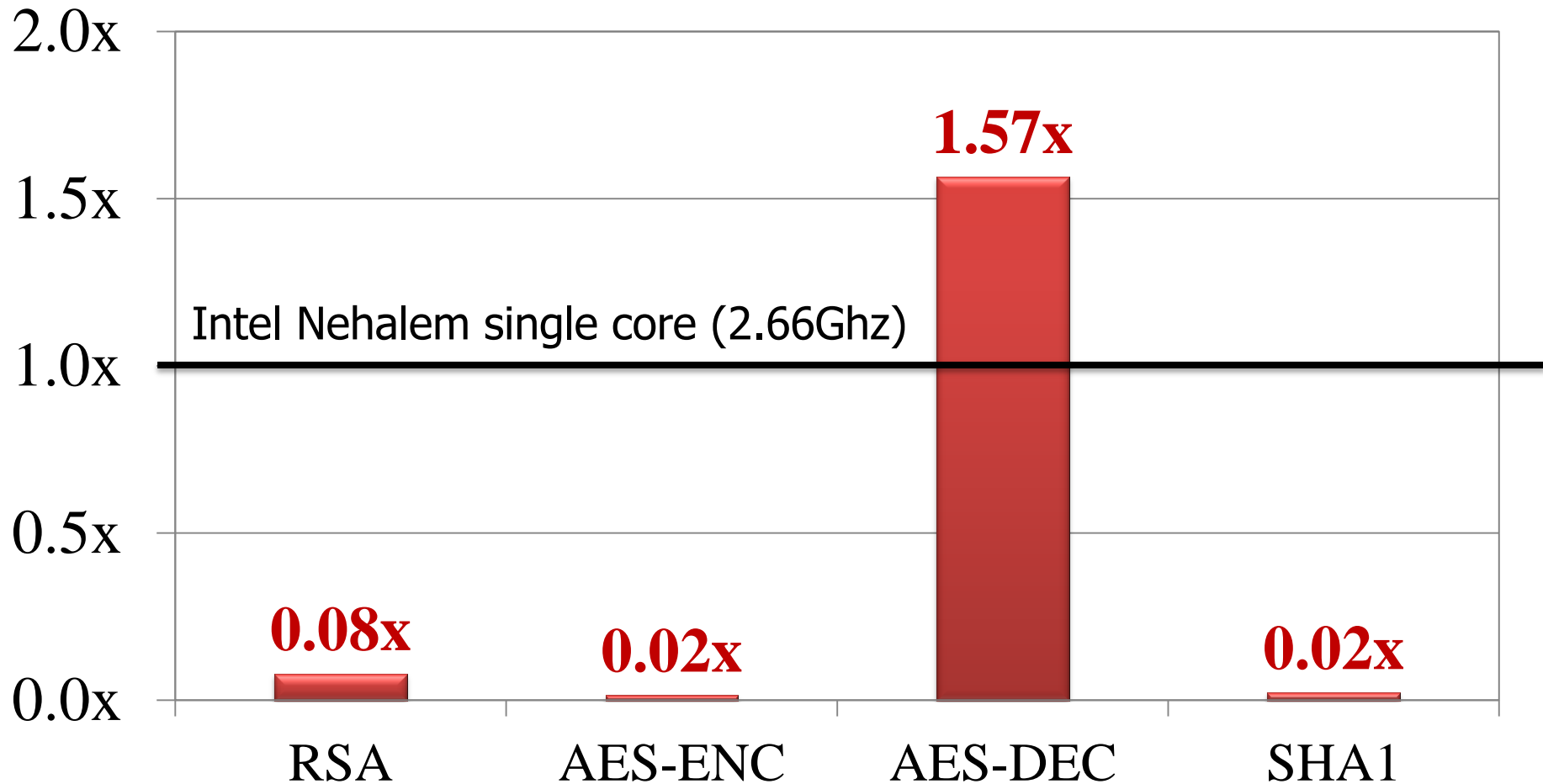# Batch Processing

2. Independent SSL Record

SSL Record  SSL Record  SSL Record

3. Parallelism in Cryptographic Operations
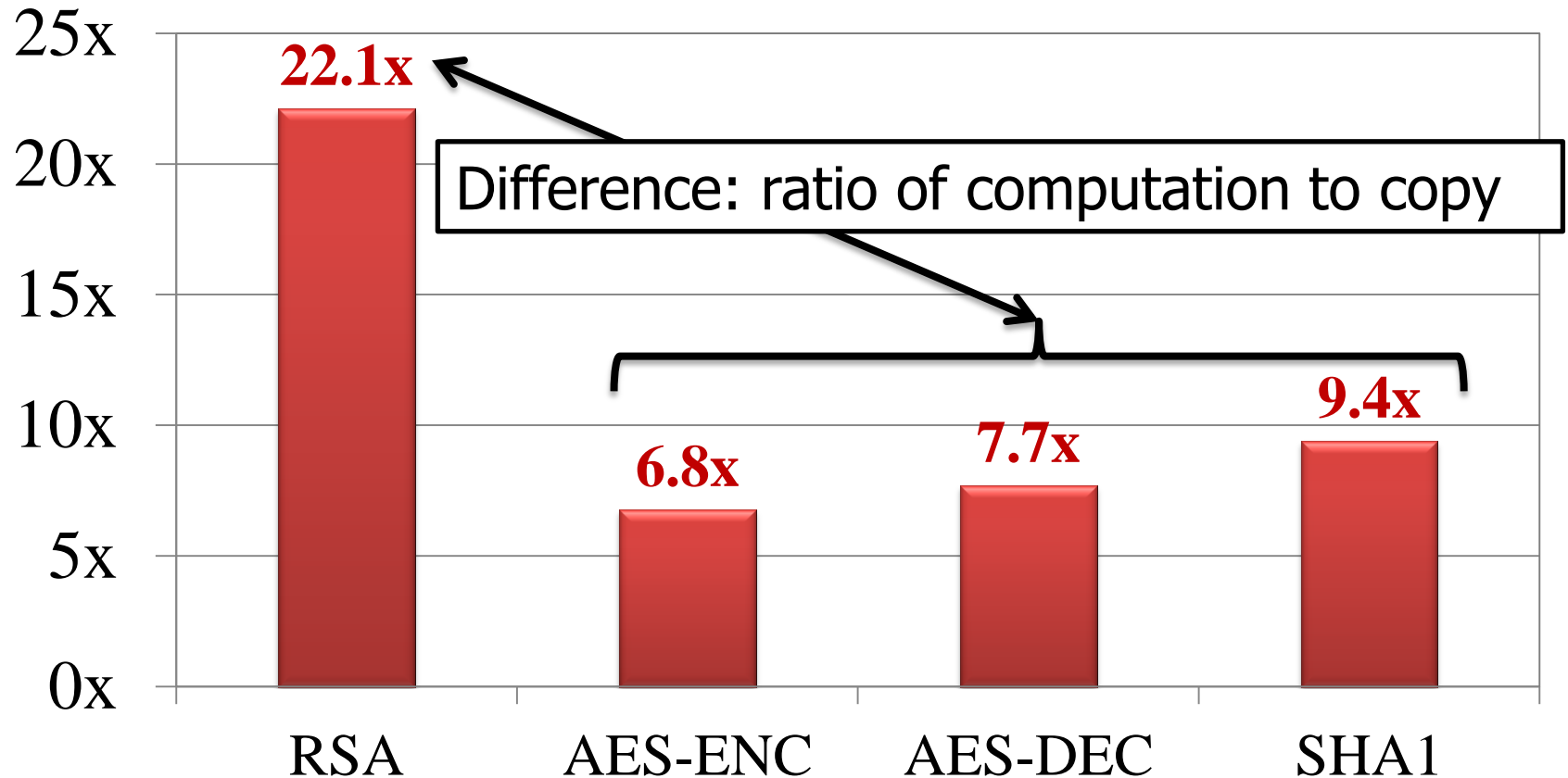
# GTX580 Throughput w/o Batching

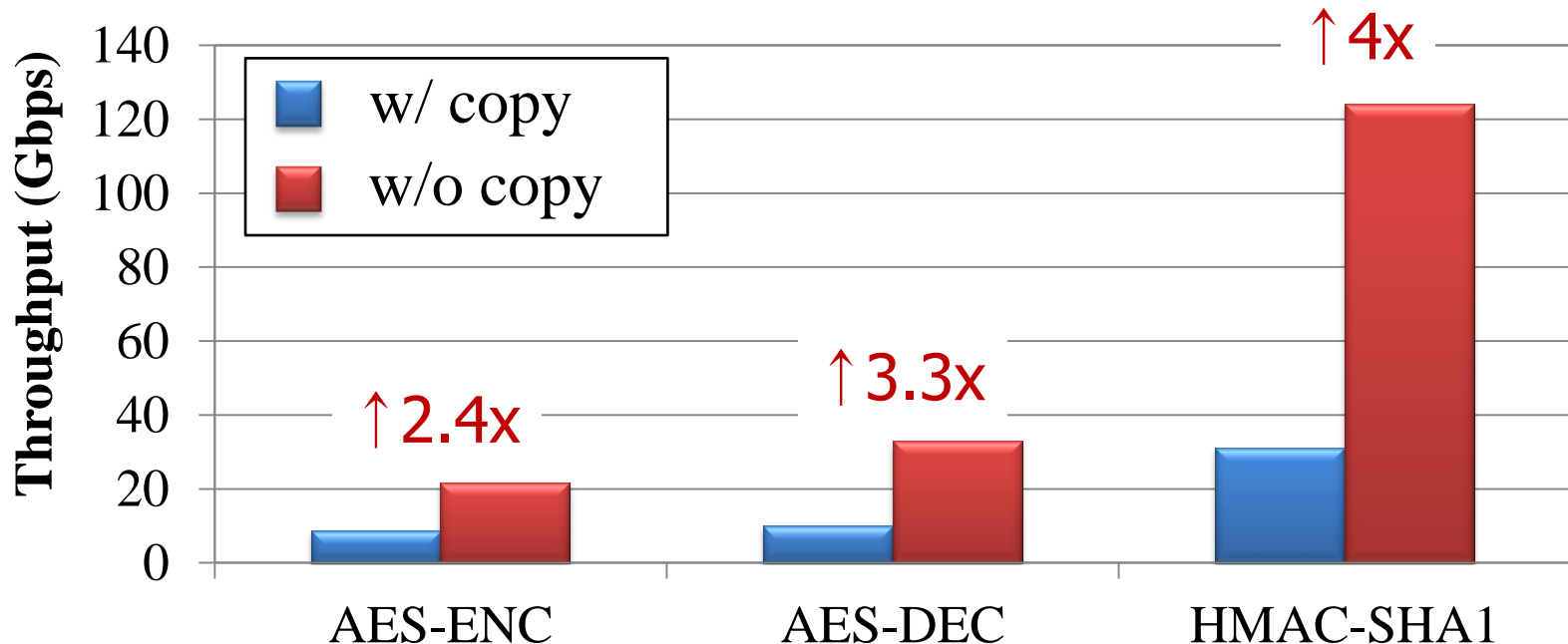**Throughput relative to a "single CPU core"**

# GTX580 Throughput w/ Batching

Batch size: **32~4096** depending on the algorithm

**Throughput relative to a "single CPU core"**



Difference: ratio of computation to copy

- RSA: 22.1x
- AES-ENC: 6.8x
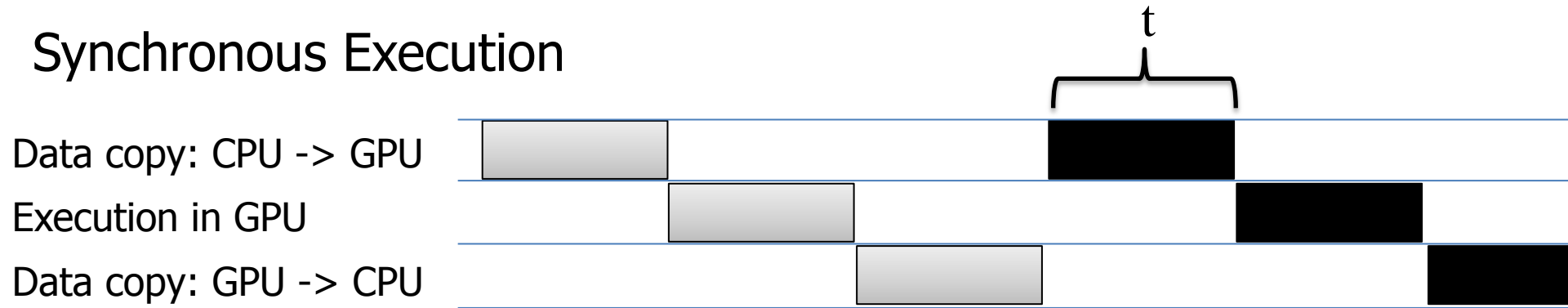- AES-DEC: 7.7x
- SHA1: 9.4x

# Copy Overhead in GPU Cryptography

- GPU processing works by
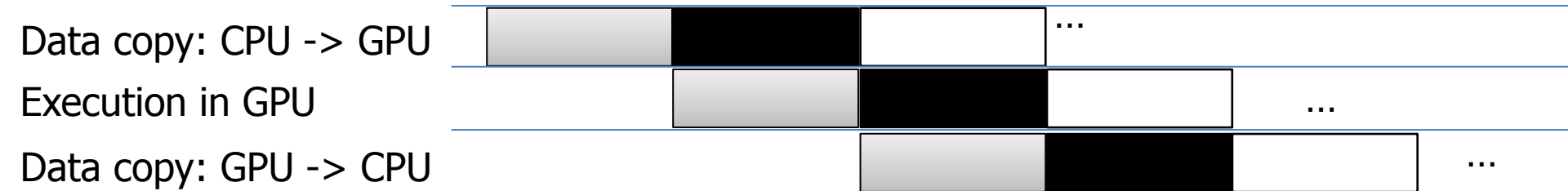  - Data copy: CPU → GPU
  - Execution in GPU
  - Data copy: GPU -> CPU

# Hiding Copy Overhead

## Synchronous Execution

t

Data copy: CPU -> GPU

Execution in GPU

Data copy: GPU -> CPU

Processing time : $3t$

## Pipelining

Data copy: CPU -> GPU    ...

Execution in GPU    ...
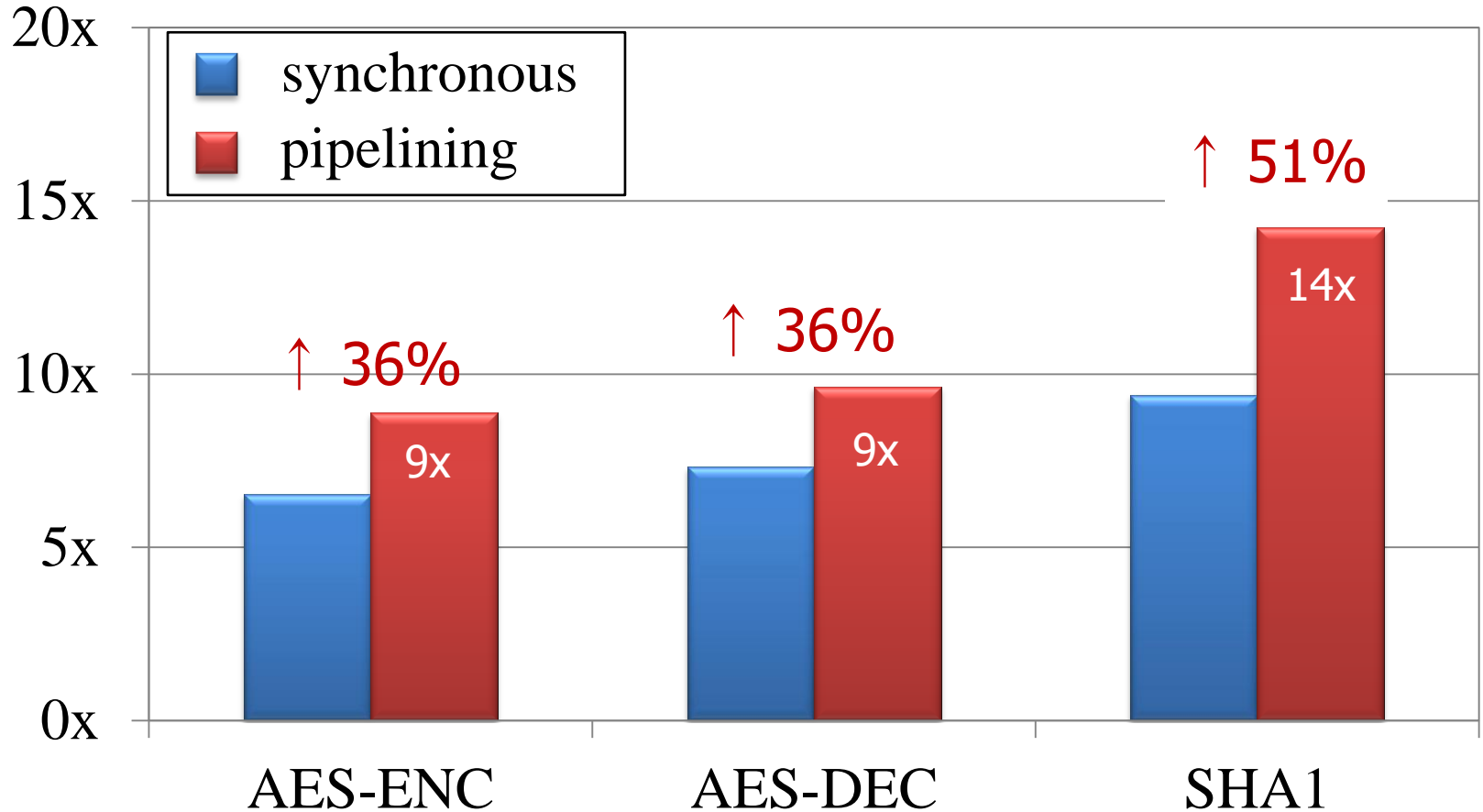
Data copy: GPU -> CPU    ...

Amortized processing time : $t$

# GTX580 Performance w/ Pipelining



**Throughput relative to a single core**

# Summary of GPU Cryptography

- Performance gain from GTX580

  - GPU performs as fast as **9 ~ 28** CPU cores
  - Superior to high-end hardware accelerators

| | RSA-1024 (ops/sec) | AES-ENC (Gbps) | AES-DEC (Gbps) | SHA1 (Gbps) |
|---|---|---|---|---|
| GTX580 | 91.9K | 11.5 | 12.5 | 47.1 |
| CPU core | 3.3K | 1.3 | 1.3 | 3.3 |

- Lessons

  - Batch processing is essential to fully utilize a GPU
  - AES and SHA1 are bottlenecked by data copy
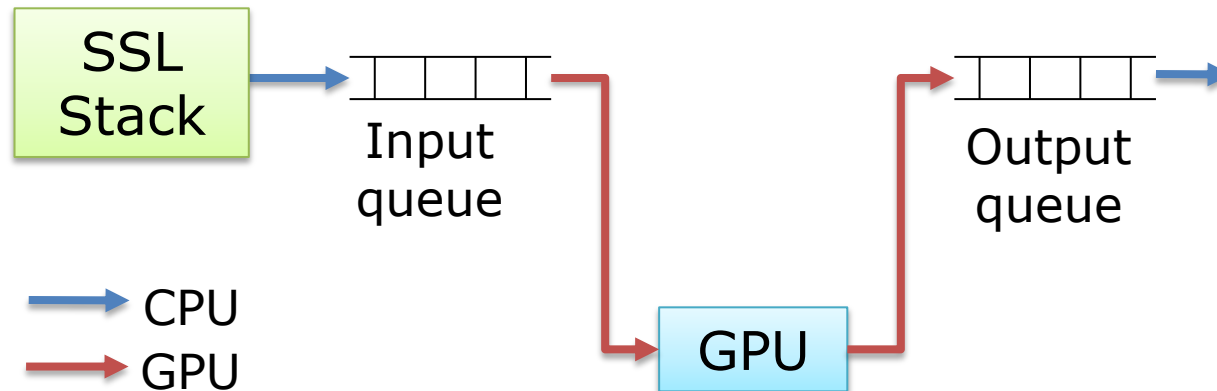    - PCIe 3.0
    - Integrated GPU and CPU

# BUILDING SSL-PROXY THAT LEVERAGES GPU

# SSLShader Design Goals

- Use existing application without modification
  - SSL reverse proxy

- Effectively leverage GPU
  - Batching cryptographic operations
  - Load balancing between CPU and GPU

- Scale performance with architecture evolution
  - Multi-core CPUs
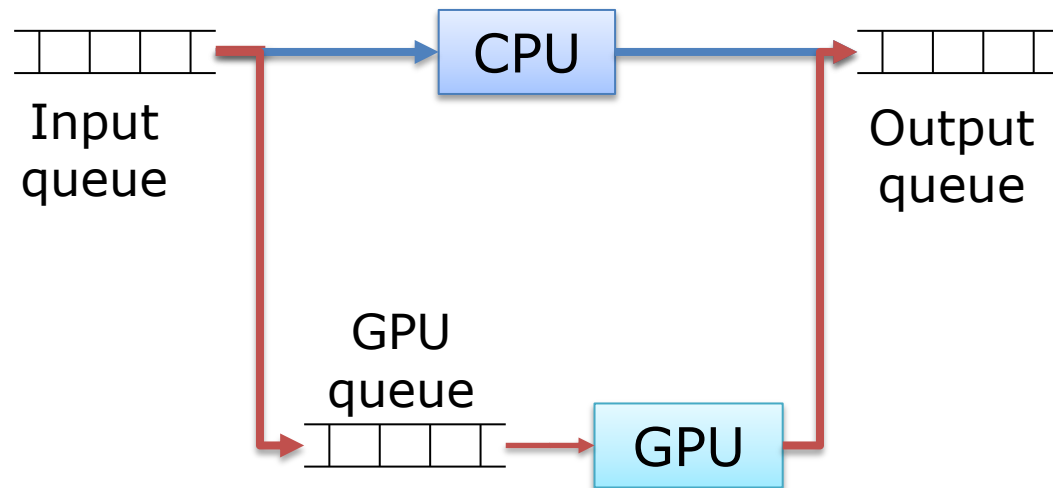  - Multiple NUMA nodes

# Batching Crypto Operations

- Network workloads vary over time
  - Waiting for fixed batch size doesn't work



- Batch size is dynamically adjusted to queue length
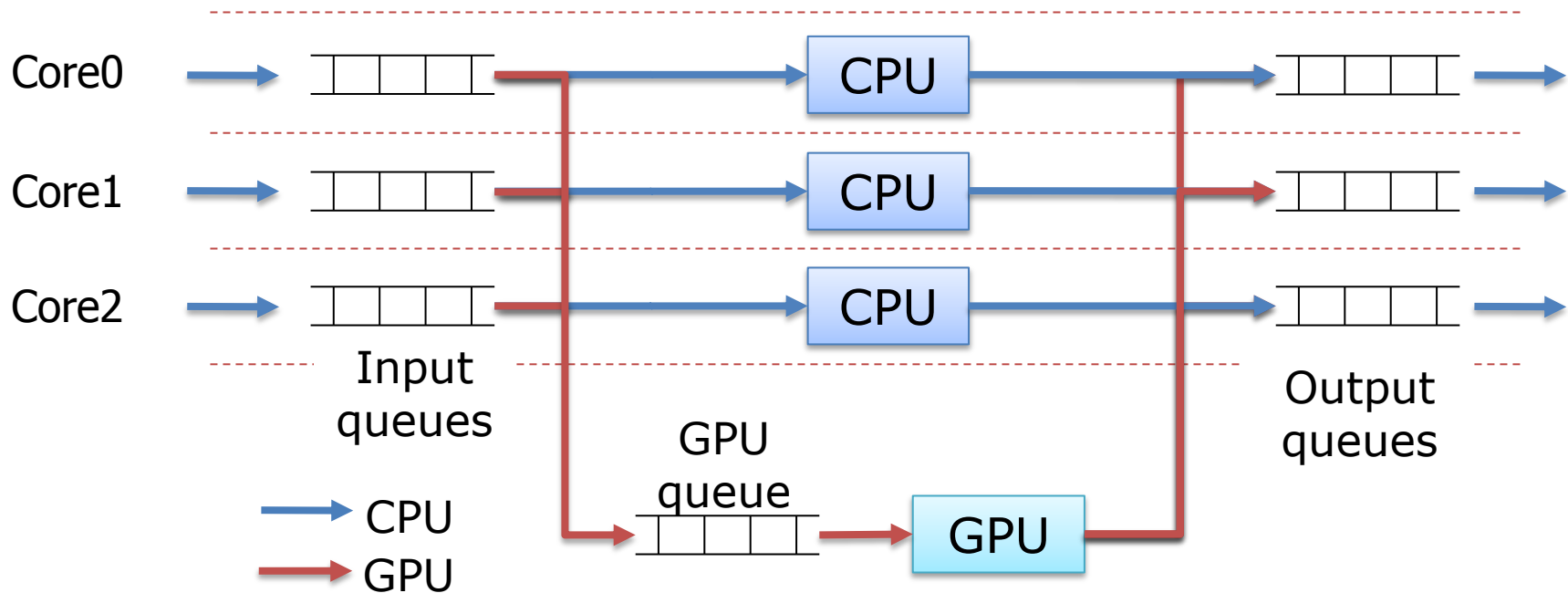
# Balancing Load Between CPU and GPU

- ## For small batch, CPU is faster than GPU
  - ### Opportunistic offloading



Input queue

CPU

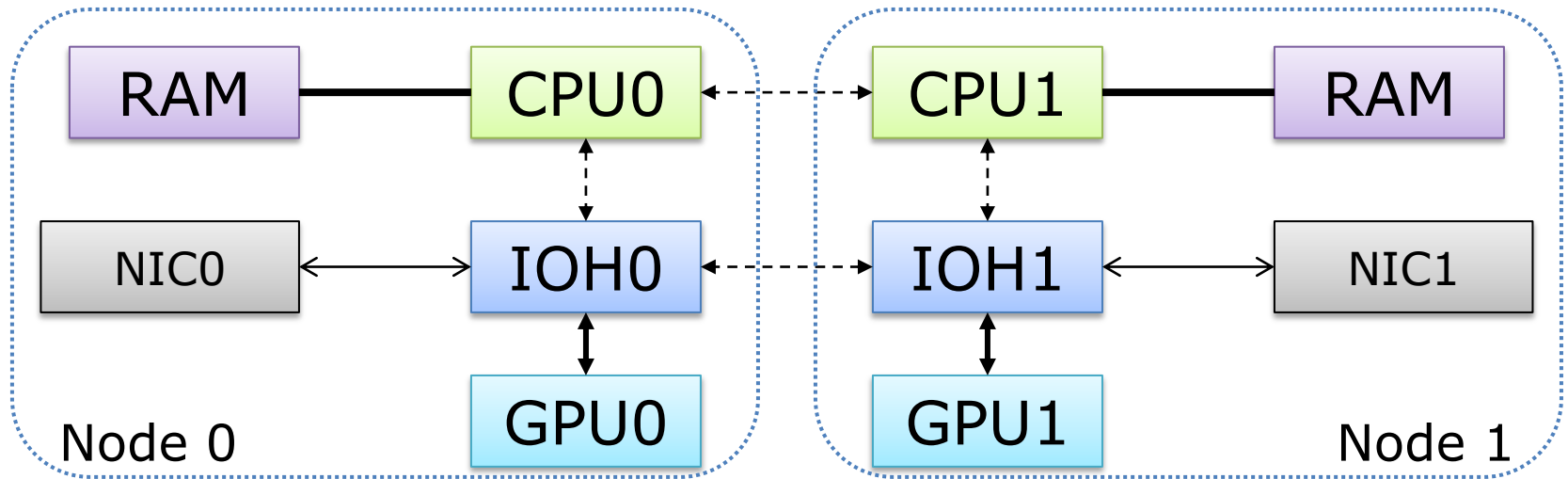Output queue

GPU queue

GPU

→ CPU processing

→ GPU processing
when input queue length > threshold

# Scaling with Multiple Cores



- Per-core worker threads
  - Network I/O, cryptographic operation
- Sharing a GPU with multiple cores
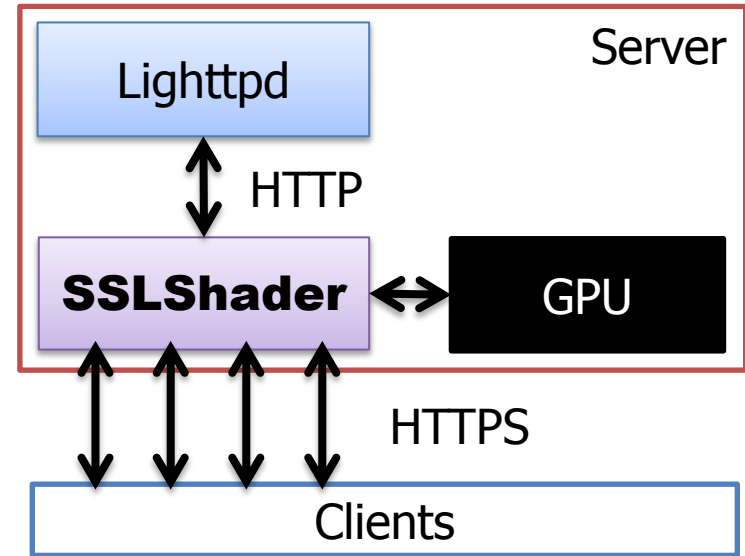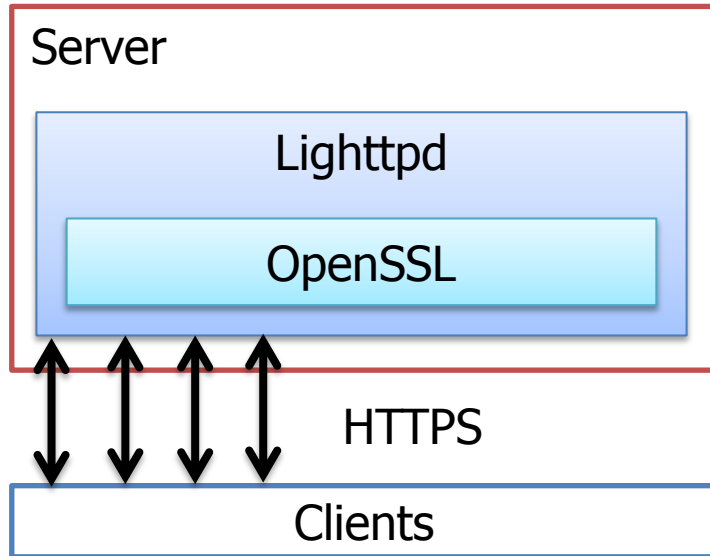  - More parallelism with larger batch size

# Scaling with NUMA systems

```
┌─────────────────────────────────┐   ┌─────────────────────────────────┐
│  ┌─────────┐      ┌─────────┐    │   │   ┌─────────┐      ┌─────────┐   │
│  │   RAM   │──────│  CPU0   │◄ ─ ─│─ ─►│   CPU1   │──────│   RAM   │   │
│  └─────────┘      └─────────┘    │   │   └─────────┘      └─────────┘   │
│                        ▲         │   │        ▲                         │
│  ┌─────────┐      ┌─────────┐    │   │   ┌─────────┐      ┌─────────┐   │
│  │  NIC0   │◄────►│  IOH0   │◄ ─ ─│─ ─►│   IOH1   │◄────►│  NIC1   │   │
│  └─────────┘      └─────────┘    │   │   └─────────┘      └─────────┘   │
│                        ▲         │   │        ▲                         │
│                   ┌─────────┐    │   │   ┌─────────┐                    │
│  Node 0           │  GPU0   │    │   │   │  GPU1   │           Node 1   │
│                   └─────────┘    │   │   └─────────┘                    │
└─────────────────────────────────┘   └─────────────────────────────────┘
```

- A process = worker threads + a GPU thread
  - Separate process per NUMA node
  - Minimizes data sharing across NUMA nodes

# Evaluation

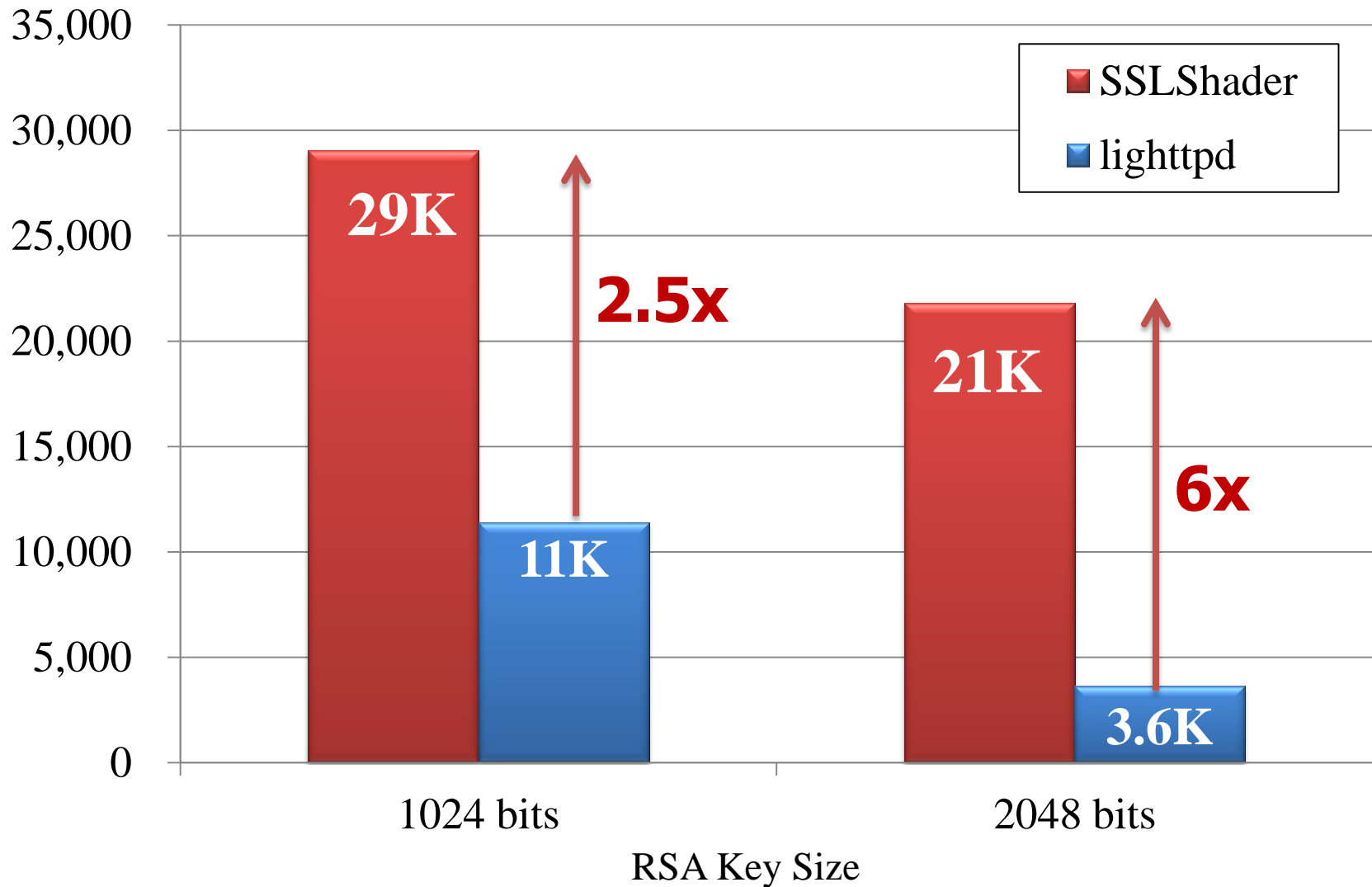- Experimental configurations



Server Specification

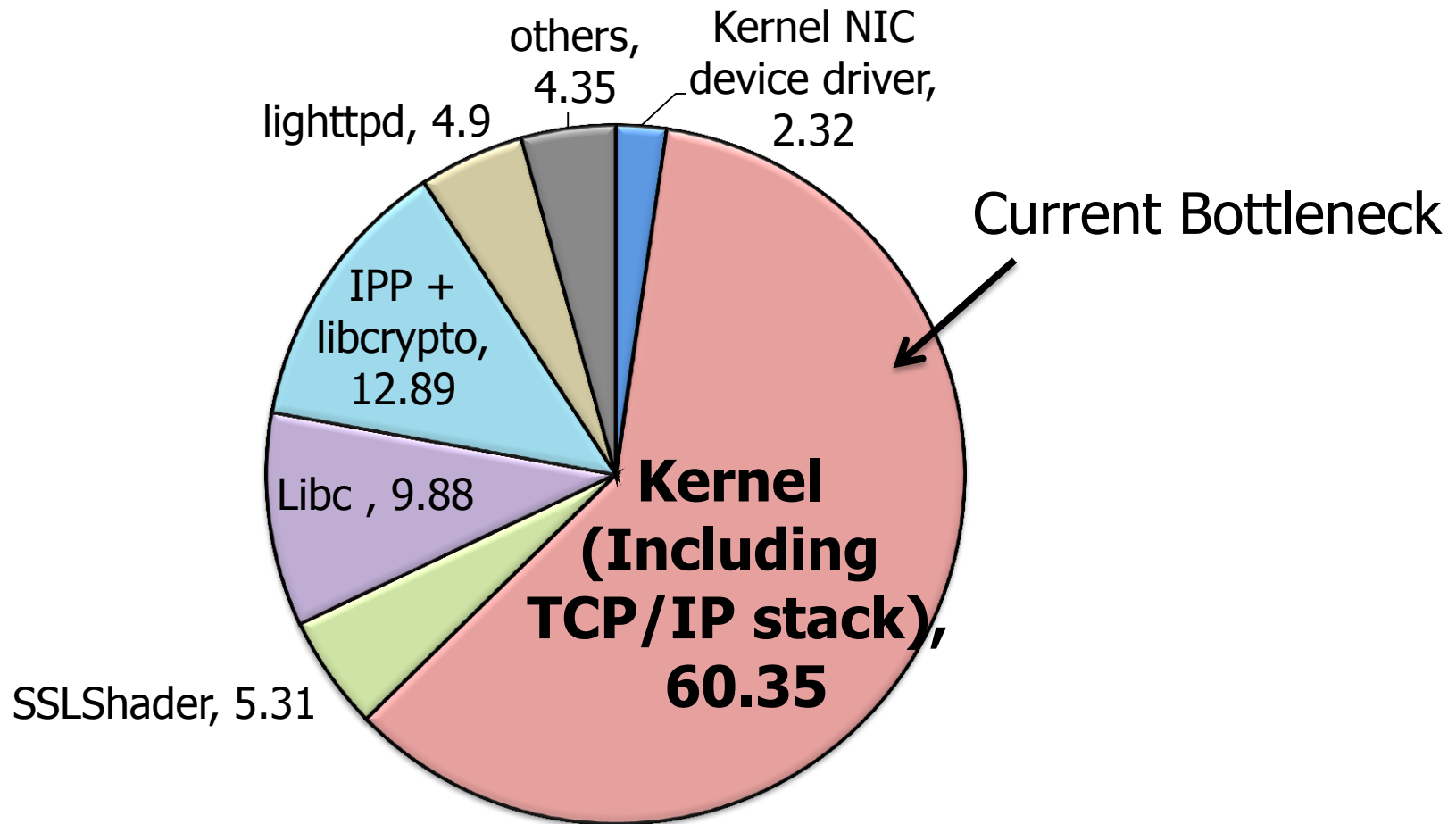| | Model | Spec | Qty |
|---|---|---|---|
| CPU | Intel X5650 | 2.66Ghz x 6 croes | 2 |
| GPU | NVIDIA GTX580 | 1.5Ghz x 512 cores | 2 |
| NIC | Intel X520-DA2 | 10GbE x 2 | 2 |

# Evaluation Metrics

- HTTPS connection handling performance
  - Use small content size
  - Stress on RSA computation

- Latency distribution at different loads
  - Test opportunistic offloading

- Data transfer rate at various content size
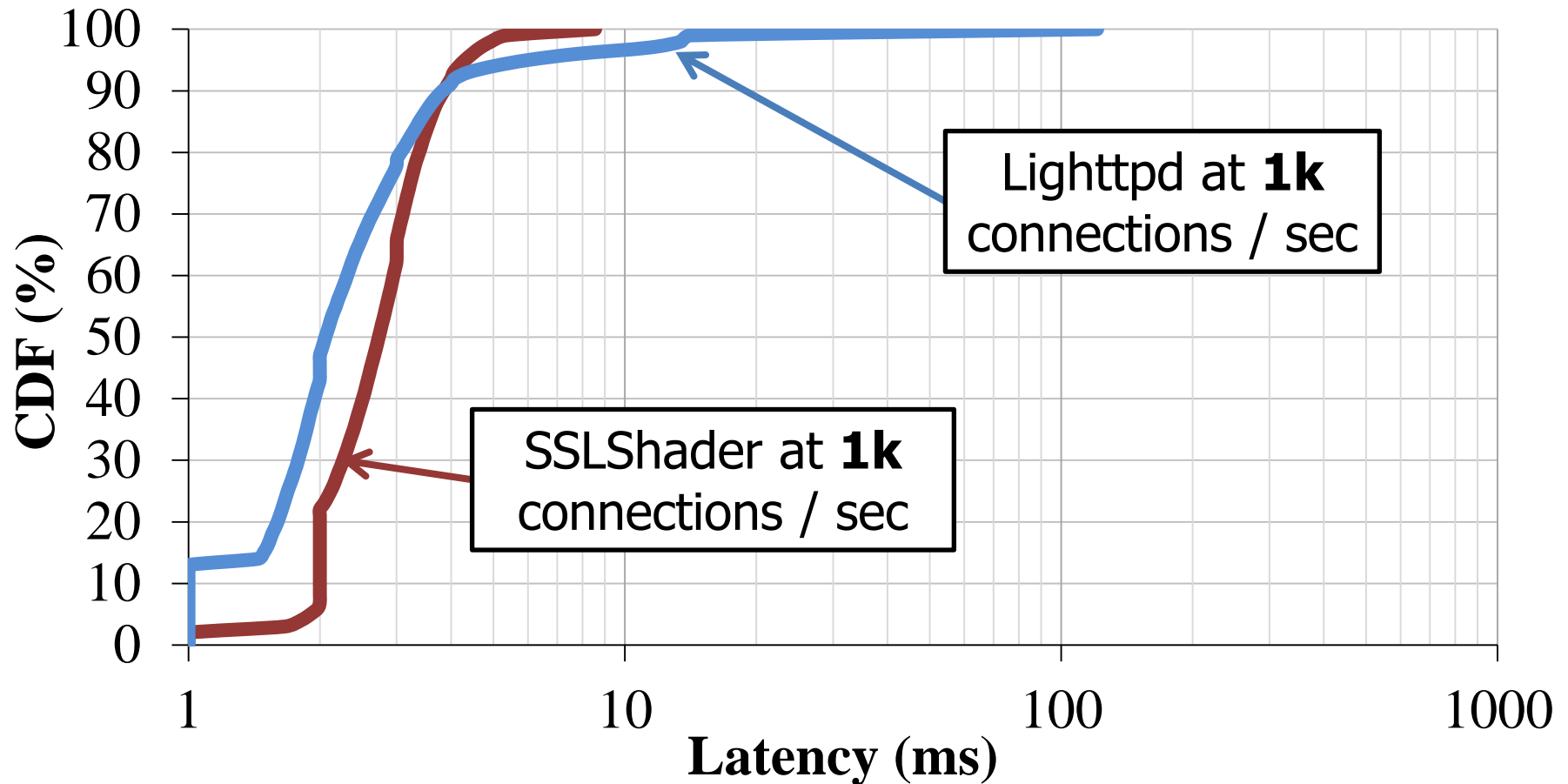
# HTTPS Connection Rate

Connections / sec

RSA Key Size

- SSLShader
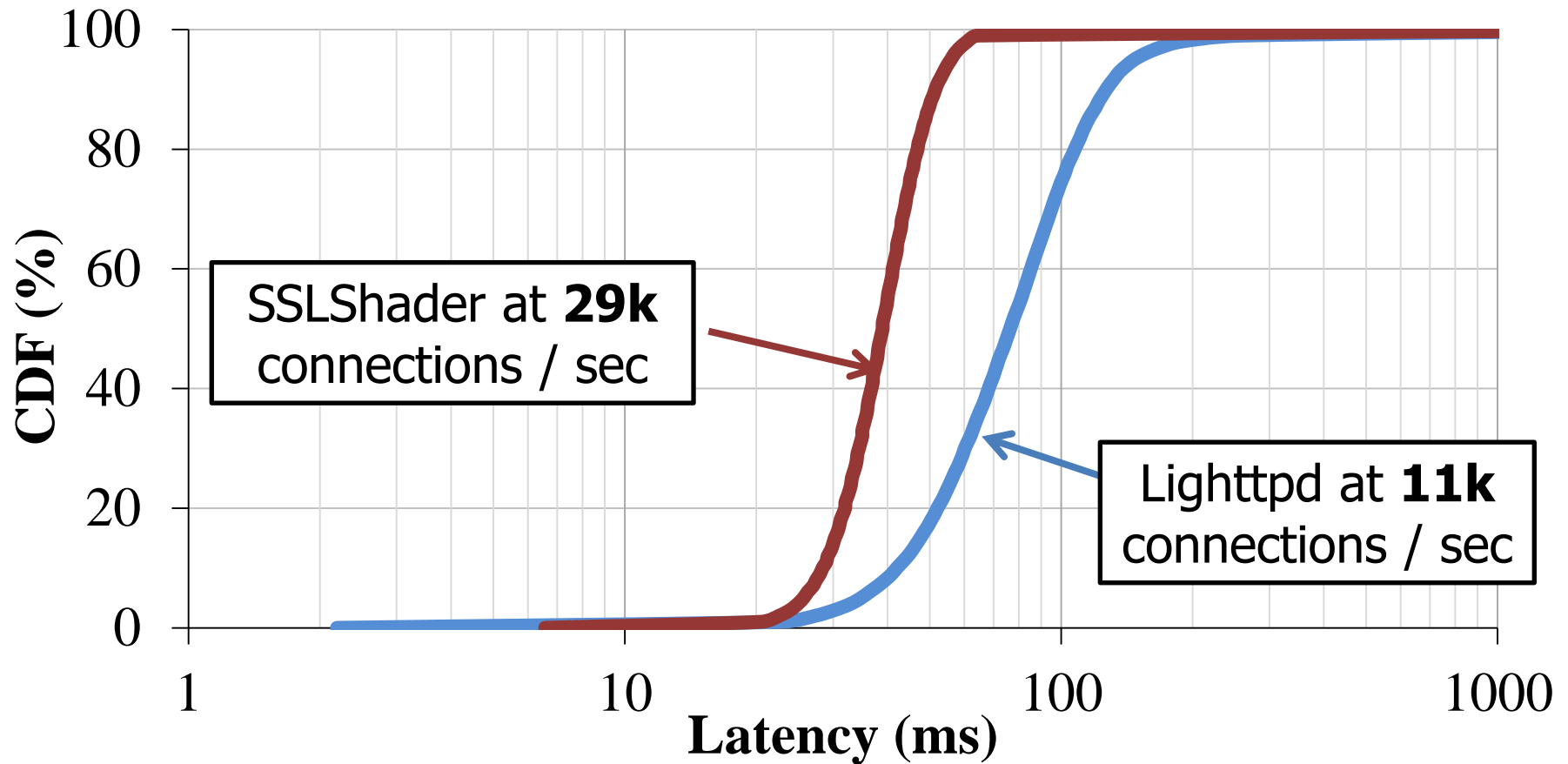- lighttpd

1024 bits: 29K, 11K — 2.5x
2048 bits: 21K, 3.6K — 6x

# CPU Usage Breakdown (RSA 1024)



others, 4.35

Kernel NIC device driver, 2.32

lighttpd, 4.9

Current Bottleneck

IPP + libcrypto, 12.89

Libc , 9.88

Kernel (Including TCP/IP stack), 60.35

SSLShader, 5.31

# Latency at Light Load



Lighttpd at **1k** connections / sec

SSLShader at **1k** connections / sec

**Similar latency at light load**

# Latency at Heavy Load



SSLShader at **29k** connections / sec

Lighttpd at **11k** connections / sec

CDF (%) — Latency (ms)

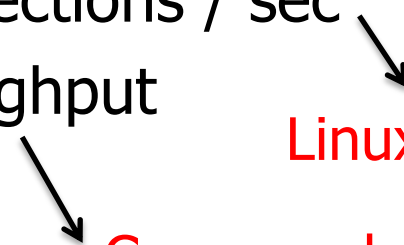**Lower latency and higher throughput at heavy load**

# Data Transfer Performance



Typical web content size is under 100KB

# CONCLUSIONS

# Summary

- Cryptographic algorithms in GPU
  - Fast RSA, AES, and SHA1
  - Superior to high-end hardware accelerators

- **SSLShader**
  - Transparent integration
  - Effective utilization of GPU for SSL processing
    - Up to 6x connections / sec
    - 13 Gbps throughput

Linux network stack performance

Copy overhead

For more details

https://shader.kaist.edu/sslshader

# QUESTIONS?

# THANK YOU!