## Harnessing GPU Computing in System Level Software

Weibin Sun

System functionality can be expensive!

System functionality can be expensive!

- <u>Crypto</u>: AES, SHAI, MD5, RSA
- Lookup: IP routing, DNS, key-value store, DB
- Erasure coding: RAID, distributed storage,

wireless

- **Pattern matching**: NIDS, virus, searching
- Compression: storage, network

▶ ...

System functionality can be expensive!

System functionality can be expensive!

Modern digital workloads are big!

















GPUs are the de-facto standard parallel processors!





GPUs are the de-facto standard parallel processors!



## Thesis work

- Generic principles
- Concretely:
  - Two generic frameworks for representative system tasks: storage and network
  - Example applications on top of both frameworks
- Literature survey \*

## GPU essentials







32-core warp





#### Related work: pioneers

- Gnort [RAID '2008]
- PacketShader [SIGCOMM '2010]
- SSLShader [NSDI '2011]
- EigenCFA [POPL '2011]
- Gibraltar [ICPP '2010]
- CrystalGPU [HPDC '2010]

#### Related work: pioneers

- Gnort [RAID '2008]
- PacketShader [SIGCOMM '2010]
- SSLShader [NSDI '2011]
- EigenCFA [POPL '2011]
- Gibraltar [ICPP '2010]
- CrystalGPU [HPDC '2010]

Specialized and Ad-hoc solutions to one particular application!

#### GPU computing frameworks/models

- Hydra [ASPLOS'08]
- CUDA-Lite [LCPC'08]
- hiCUDA [GPGPU'09]
- ASDSM [ASPLOS'10]
- Sponge [ASPLOS'11]
- CGCM [PLDI'II]
- PTask [SOSP'11]

٠

#### GPU computing frameworks/models

- Hydra [ASPLOS'08]
- CUDA-Lite [LCPC'08]
- hiCUDA [GPGPU'09]
- ASDSM [ASPLOS'10]
- Sponge [ASPLOS'II]
- · CGCM [PLDI'II]
- PTask [SOSP'11]



- Simple workflow
- Long computation, PCIe copy and synchronization not matter
- No batched data divergence

#### GPU computing frameworks/models



## Rest of this talk

- Generic principles
- GPUstore : for storage systems
- Snap : for packet processing

## System level GPU computing: generic principles
Latency-oriented system code V.S. throughput-oriented GPU

CPU-GPU synchronization hurts asynchronous systems

Wasted PCIe data transfer

Double buffering for GPU DMA

Latency-oriented system code V.S. throughput-oriented GPU

CPU-GPU synchronization hurts asynchronous systems

Wasted PCIe data transfer

Double buffering for GPU DMA

Batched processing: wait for enough workload

Batched processing: wait for enough workload

how much is enough?

Batched processing: wait for enough workload

- how much is enough?
- increased latency

Latency-oriented system code V.S. throughput-oriented GPU

CPU-GPU synchronization hurts asynchronous systems

Wasted PCIe data transfer

Double buffering for GPU DMA









0.





















Do you need the entire packet to do routing lookup?

Do you need the entire packet to do routing lookup?

What you need:

Do you need the entire packet to do routing lookup?

What you need:



Do you need the entire packet to do routing lookup?

What you need:



What you may transfer:

Do you need the entire packet to do routing lookup?

What you need:



What you may transfer:



# Compact your workload data

For ``use-partial" data usage pattern









DMA only on locked memory pages

- DMA only on locked memory pages
- GPU drivers only recognize their own locked memory pool



- DMA only on locked memory pages
- GPU drivers only recognize their own locked memory pool



- DMA only on locked memory pages
- GPU drivers only recognize their own locked memory pool



DMA only on locked mer ory pages Make GPU driver and system 'V code use the same locked memory for zero-copy DMA! Externa **GPU** mem mem DMA locked

20

mem


#### Problems



## Generic principles



Asynchronous non-blocking GPU programming

Reduce PCIe data transfer with compacted workload

#### GPUstore: Harnessing GPU Computing in OS Storage Systems

[Weibin Sun, Robert Ricci, Matthew L. Curry @ SYSTOR'12]

• A storage framework for Linux kernel to use CUDA GPUs in filesystems, block drivers, etc.

• A storage framework for Linux kernel to use CUDA GPUs in filesystems, block drivers, etc.



- A storage framework for Linux kernel to use CUDA GPUs in filesystems, block drivers, etc.
- Minimally invasive GPU integration



- A storage framework for Linux kernel to use CUDA GPUs in filesystems, block drivers, etc.
- Minimally invasive GPU integration
  - Small changes



- A storage framework for Linux kernel to use CUDA GPUs in filesystems, block drivers, etc.
- Minimally invasive GPU integration
  - Small changes
  - Preserve interface and semantics



- A storage framework for Linux kernel to use CUDA GPUs in filesystems, block drivers, etc.
- Minimally invasive GPU integration
  - Small changes
  - Preserve interface and semantics
  - Keep efficient











#### GPUstore request workflow



#### GPUstore request workflow



Batched processing

Asynchronous non-blocking GPU programming

Reduce PCIe data transfer with compacted workload

Batched processing

Asynchronous non-blocking GPU programming







### Request scheduling

Merge small requests

Split large requests

## Request scheduling



## Request scheduling



Split large requests

How to define "small" and "large"?

• Up to specific services.

Batched processing

Asynchronous non-blocking GPU programming

Reduce PCIe data transfer with compacted workload

Reduce PCIe data transfer with compacted workload

All use the same locked memory for zero-copy DMA

User space CUDA library dependency

All use the same locked memory for zero-copy DMA

User space CUDA library dependency

In-kernel CUDA's locked memory allocator

#### In-kernel CUDA locked memory allocator







### Why infeasible?



### Why infeasible?

→ Who allocates the memory?



## Why infeasible?

- ➡ Who allocates the memory?
  - "pass-by-ref" interface


# Why infeasible?

- Who allocates the memory?
  - "pass-by-ref" interface
- Can NOT modify highly-depended cache allocators:



# Why infeasible?

- Who allocates the memory?
  - "pass-by-ref" interface
- Can NOT modify highly-depended cache allocators:
  - page cache, buffer cache
  - object cache
  - packet pool
  - ...







GPU mem







All use the same locked memory for zero-copy DMA

All use the same locked memory for zero-copy DMA

In-kernel CUDA locked memory allocator

All use the same locked memory for zero-copy DMA

In-kernel CUDA locked memory allocator

Remap external memory into GPU driver's locked memory area

### Implementation and evaluation

### Implementation and evaluation

- Case studies of major storage residents:
  - <u>dm-crypt</u>: disk encryption layer
  - <u>eCryptfs</u>: encrypted filesystem
  - <u>md</u>: software RAID6



- <u>eCryptfs</u>: encrypted filesystem
- <u>md</u>: software RAID6
- Integration cost:

| Subsystem | Total LOC | Modified LOC | Percent |
|-----------|-----------|--------------|---------|
| dm-crypt  | $1,\!800$ | 50           | 3%      |
| eCryptfs  | 11,000    | 200          | 2%      |
| md        | 6,000     | 20           | 0.3%    |









### Faster than SSD: dm-crypt



Throughput (MB/s)

### Faster than SSD: dm-crypt



### Working with existing optimization

eCryptfs on RAM disks

#### Linux max 128KB read-ahead effects on read



# GPUstore Summary

- Enables efficient GPU-accelerated storage in Linux kernel with small changes
- <u>https://github.com/wbsun/kgpu</u>
- Details not presented: user-space helper, non-blocking K-U comm, more exp result

### Snap: Fast and Flexible Packet Processing With GPUs and Click

[Weibin Sun, Robert Ricci @ ANCS'13]

### Click background [Kohler, et.al. TOCS'00]



- Elements
- Ports
  - push/pull

## Why using GPUs in Click?





### Why using GPUs in Click?

### Why using GPUs in Click?





### Snap: the idea

- Moving parts of Click pipeline onto GPUs
  - CPU for sequential
  - GPU for parallel
- Keep Click's modular style



Batched processing

Asynchronous non-blocking GPU programming

Reduce PCIe data transfer with compacted workload

All use the same locked memory for zero-copy DMA

Batched processing






• Packets in a batch go to different paths



• Packets in a batch go to different paths



• Packets in a batch go to different paths



• Packets in a batch go to different paths



#### Two choices:

Split the batch into two copies

• Packets in a batch go to different paths



#### <u>Two choices:</u>

- Split the batch into two copies
  - needs PCIe copy, sync!

• Packets in a batch go to different paths



#### <u>Two choices:</u>

- I. Split the batch into two copies
  - needs PCIe copy, sync!
- 2. Keep the batch, skip some packets on GPUs









Batched processing

Asynchronous non-blocking GPU programming

Reduce PCIe data transfer with compacted workload

Asynchronous non-blocking GPU programming

Reduce PCIe data transfer with compacted workload

Asynchronous non-blocking GPU programming

Reduce PCIe data transfer with compacted workload

Asynchronous non-blocking GPU programming

Each batch binds a CUDA stream

Reduce PCIe data transfer with compacted workload

Asynchronous non-blocking GPU programming

Each batch binds a CUDA stream

BElement only does async GPU ops

Reduce PCIe data transfer with compacted workload

Asynchronous non-blocking GPU programming

- Each batch binds a CUDA stream
- BElement only does async GPU ops
- Use GPUCompletionQueue to check stream status

Reduce PCIe data transfer with compacted workload

Reduce PCIe data transfer with compacted workload



Packet Classification's regions-of-interest (ROI):

Packet

















Batched processing

Asynchronous non-blocking GPU programming

Reduce PCIe data transfer with compacted workload

All use the same locked memory for zero-copy DMA

ROIs stored in CUDA locked memory

#### Evaluation

- Basic forwarding I/O improvement
- Example applications:
  - Classification
  - IP routing
  - Pattern matching for IDS
- Latency, re-order
- Flexibility and modularity

#### Classifier+routing+pattern matching



#### Classifier+routing+pattern matching



#### Evaluation



### Evaluation

Snap-CD

- GPU reached 40Gb/s line rate at 128B
- CPU just I/3 or I/4

Sn Sn

 Latency tolerable in LAN for non-latencysensitive app, negligible in WAN

 E
 0
 64
 256
 512

 0
 64
 packet Size (Bytes)

 (a)
 SDN Forwarder

 (a)
 SDN Forwarder
### Preserving Click's flexibility: full IP router



### Preserving Click's flexibility: full IP router



### Preserving Click's flexibility: full IP router



# Snap Summary

- A generic parallel packet processing framework
  - flexibility of Click
  - fast parallel power from GPUs
- <u>https://github.com/wbsun/snap</u>
- Details not presented: network I/O, async scheduling

### Thesis statement

The throughput of system software with parallelizable, computationally expensive tasks can be improved by using GPUs and frameworks with memory-efficient and throughput-oriented designs.

### Conclusion

## Conclusion

#### **Generic principles:**

Batched processing

Asynchronous non-blocking GPU programming

Reduce PCIe data transfer with compacted workload

All use the same locked memory for zero-copy DMA

## Conclusion

#### **Generic principles:**

Batched processing

Asynchronous non-blocking GPU programming

Reduce PCIe data transfer with compacted workload

All use the same locked memory for zero-copy DMA

#### **Concrete frameworks:**

• GPUstore with high throughput storage applications

Snap with high throughput network packet processing

# Thanks! Q&A

# Backup slides

### ROI for memory access coalescing



# How ROI slicing works?

- To Click insiders:
  - Batcher accepts ROI requests from BElements
  - Batcher merges requested ROIs into result ROIs
  - Each BElement asks for its ROIs' offsets
  - GPU kernels invoked by BElements use variables for offsets

## How predicated execution works?

- To Click insiders:
  - Manually assigned: where and what!

## Path-encoded predicate

