# Divergence-Aware Warp Scheduling

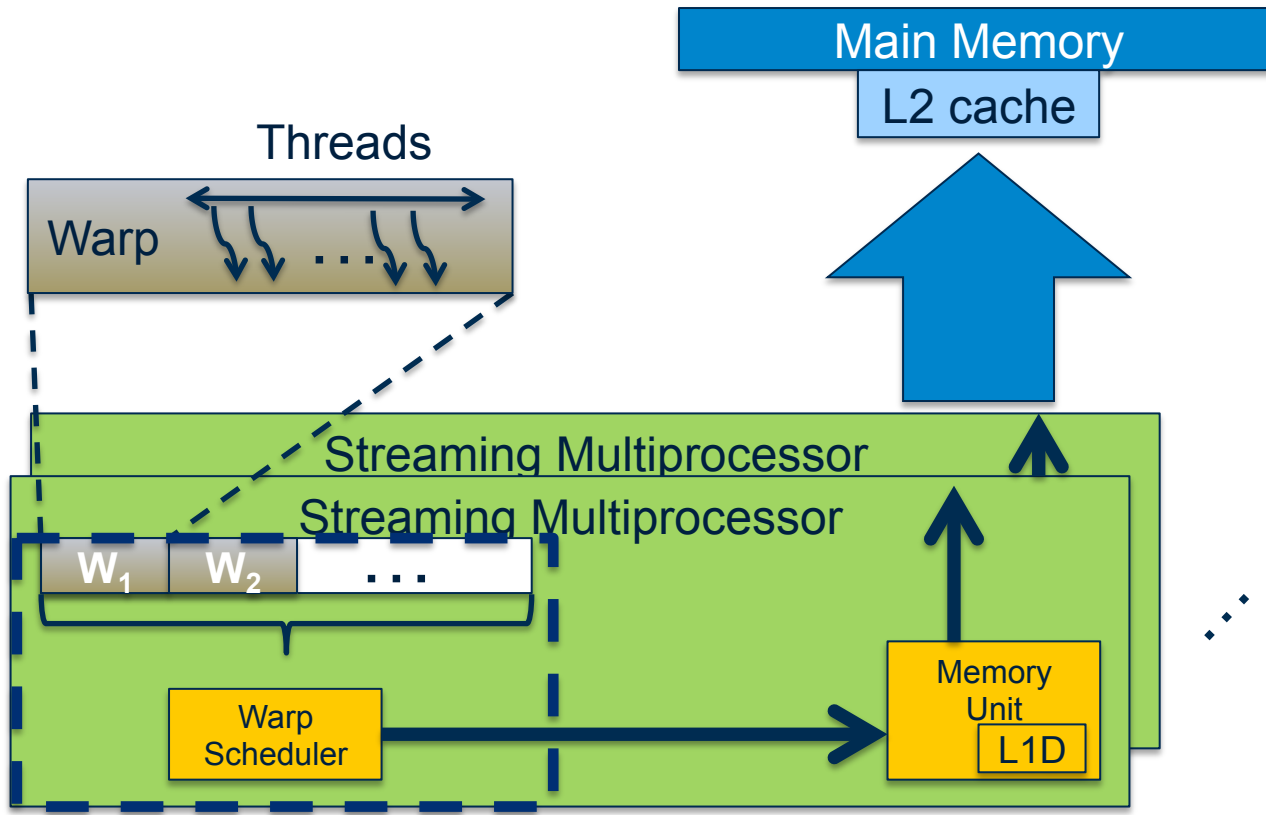**Timothy G. Rogers[1], Mike O'Connor[2], Tor M. Aamodt[1]**

**[1]The University of British Columbia**
**[2]NVIDIA Research**

# GPU

- 10000's concurrent threads
- Grouped into warps
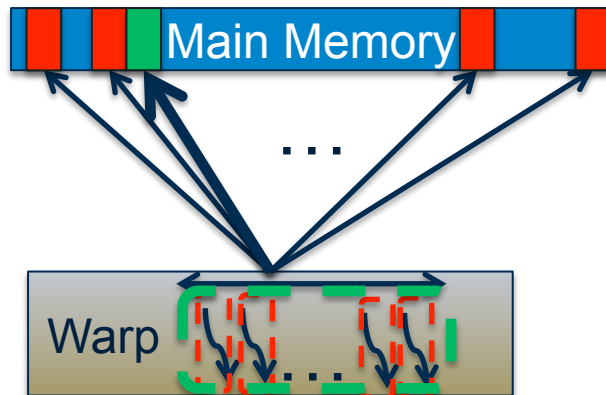- Scheduler picks warp to issue each cycle

# 2 Types of Divergence

## Branch Divergence

if(…) {
    …
}

Threads

Warp  | 1 | 0 | .. | 0 | 1 |

Effects functional unit utilization

Aware of branch divergence

## Memory Divergence

Main Memory

…

Warp

Can waste memory bandwidth

Aware of memory divergence

AND

Focus on improving performance

# Motivation

- **Improve performance of programs with memory divergence**
  - Parallel irregular applications
  - Economically important (server computing, big data)

- **Transfer locality management from SW to HW**

  - Software solutions:
    - Complicate programming
    - Not always performance portable
    - Not guaranteed to improve performance
    - Sometimes impossible

# Previous Work

- **Scheduling used to capture intra-thread locality (MICRO 2012)**

W₁  W₂  W  ...  W

Go  Go  Go  **Stop**  Go  **Stop**

Active Mask

Active Mask

Warp Scheduler

Memory Unit

L1D

**Lost Locality Detected**

| **Previous Work** | **Divergence-Aware Warp Scheduling** |
|---|---|
| **Reactive** • Detects interference then throttles | **Predict and be Proactive** |
| **Unaware of branch divergence** • All warps treated equally | **Adapt to branch divergence** |
| **Outperformed by profiled static throttling** | **Case Study: Divergent code 50% slowdown** | **Outperform static solution** | **Case Study: Divergent code <4% slowdown** |

## Divergence-Aware Warp Scheduling

### How to be proactive

- Identify where locality exists
- Limit the number of warps executing in high locality regions

### Adapt to branch divergence

- Create cache footprint prediction in high locality regions
- Account for number of active lanes to create **per-warp footprint prediction**.
- Change the prediction as branch divergence occurs.

# Where is the locality?

- **Examine every load instruction in program**



**Loop**

**Locality Concentrated in Loops**

# Locality In Loops Limit Study

## How much data should we keep around?

**Hits on data accessed in immediately previous trip**



**Average**

a place of mind

## DAWS Objectives

1. **Predict the amount of data accessed by each warp in a loop iteration.**

2. **Schedule warps in loops so that aggregate predicted footprint does not exceed L1D.**

# Observations that enable prediction

- Memory divergence in static instructions is predictable



**Main Memory**

**Divergence**

**Warp 1**

…
load
…

**Both Used To Create Cache Footprint Prediction**

- Data touched by divergent loads dependent on active mask



**Main Memory**

4 accesses

**Divergence**

**Warp**

| 1 | 1 | 1 | 1 |

**Main Memory**

2 accesses

**Divergence**

**Warp**

| 1 | 0 | 0 | 1 |

# Online characterization to create cache footprint prediction

## 1. Detect loops with locality

Some loops have locality          Some don't

**Limit multithreading here**

## 2. Classify loads in the loop

Loop with locality

```
while(…) {
    load 1   Diverged
    …
    load 2   Not Diverged
}
```

## 3. Compute footprint from active mask

Loop with locality

Warp 0

```
while(…) {
    load 1   Diverged        4 accesses
    …                            +
    load 2   Not Diverged    1 access
}
```

**Warp 0's Footprint = 5 cache lines**

# DAWS Operation Example

**Cache**

| | |
|---|---|
| A[0] | Hit |
| A[64] | Hit |
| A[96] | Hit |
| A[128] | Hit |

**Cache**

| | |
|---|---|
| A[0] | Hit x30 |
| A[64] | Hit x30 |
| A[96] | Hit x30 |
| A[128] | Hit x30 |

**Cache**

| |
|---|
| A[32] |
| A[160] |
| A[192] |
| A[224] |

## Example Compressed Sparse Row Kernel

```
int C[]={0,64,96,128,160,160,192,224,256};
void sum_row_csr(float* A, …)  {
    float sum = 0;
    int i =C[tid];
```
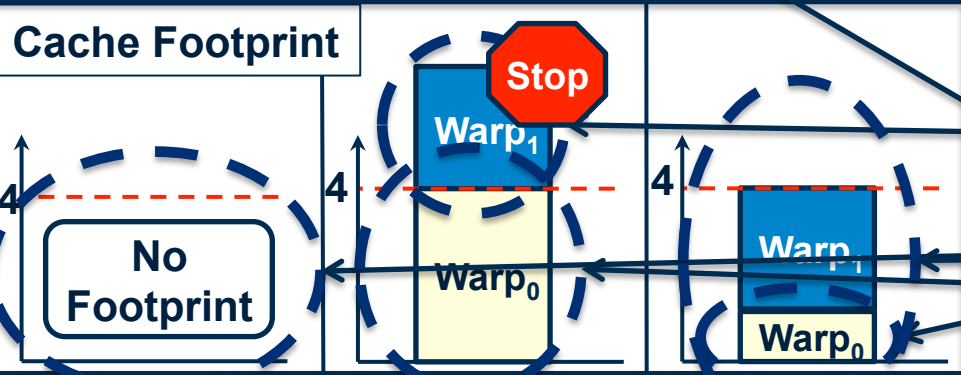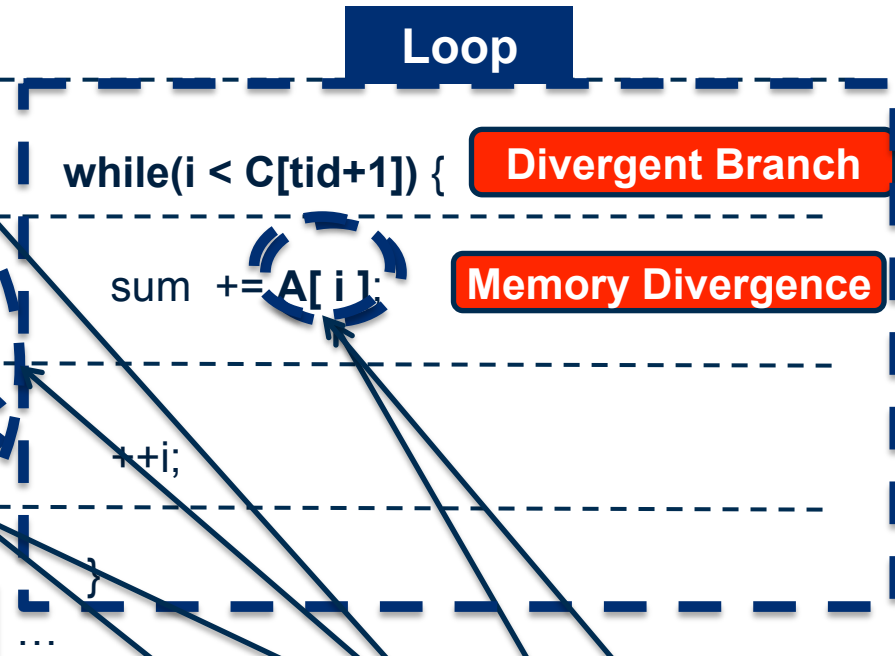
**Loop**

```
while(i < C[tid+1]) {        Divergent Branch

    sum  += A[ i ];          Memory Divergence

    ++i;

}

…
```

**1st Iter.**
**Warp₁**

Stop    Go

| 0 | 1 | 1 | 1 |

**Warp₀**    | 1 | 1 | 1 | 1 |    Go
**2nd Iter.**

**Warp₀**    | 1 | 0 | 0 | 0 |    Go
**33rd Iter.**

Go

Time₀        Time₁        Time₂

## Cache Footprint

4        4        Stop
                 **Warp₁**

4        **Warp₀**    **Warp₁**

No Footprint        **Warp₀**    **Warp₀**

Warp 0 has branch divergence        y warps        loop for        warps
Both warps capture        locality        er        = 4X1
Footprint decreased

# Methodology

## GPGPU-Sim (version 3.1.0)

- 30 Streaming Multiprocessors
  - 32 warp contexts (1024 threads total)
- 32k L1D per streaming multiprocessor
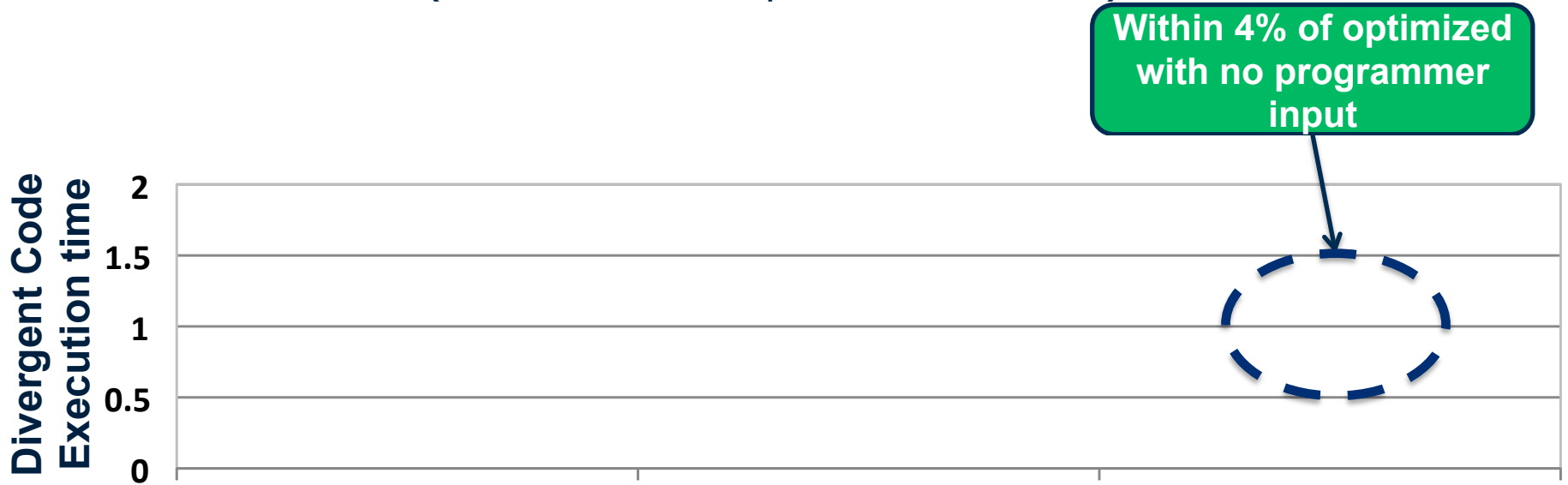- 1M L2 unified cache

## Compared Schedulers

- Cache-Conscious Wavefront Scheduling (CCWS)
- Profile based Best-SWL
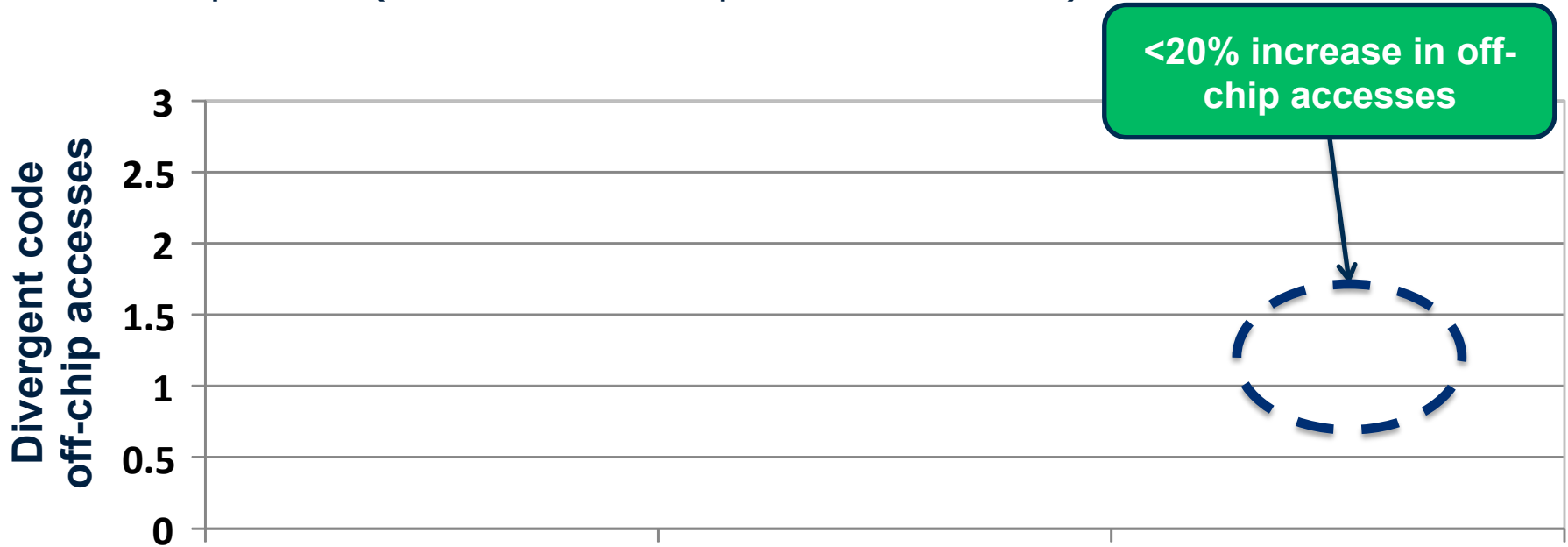- **Divergence-Aware Warp Scheduling (DAWS)**

More schedulers
in paper

# Sparse MM Case Study Results

- Performance (normalized to optimized version)

**Within 4% of optimized with no programmer input**

# Sparse MM Case Study Results

- Properties (normalized to optimized version)

**<20% increase in off-chip accesses**

**Divergent code issues 2.8x less instructions**

**Divergent version now has potential energy advantages**

*Y-axis: Divergent code off-chip accesses (0, 0.5, 1, 1.5, 2, 2.5, 3)*
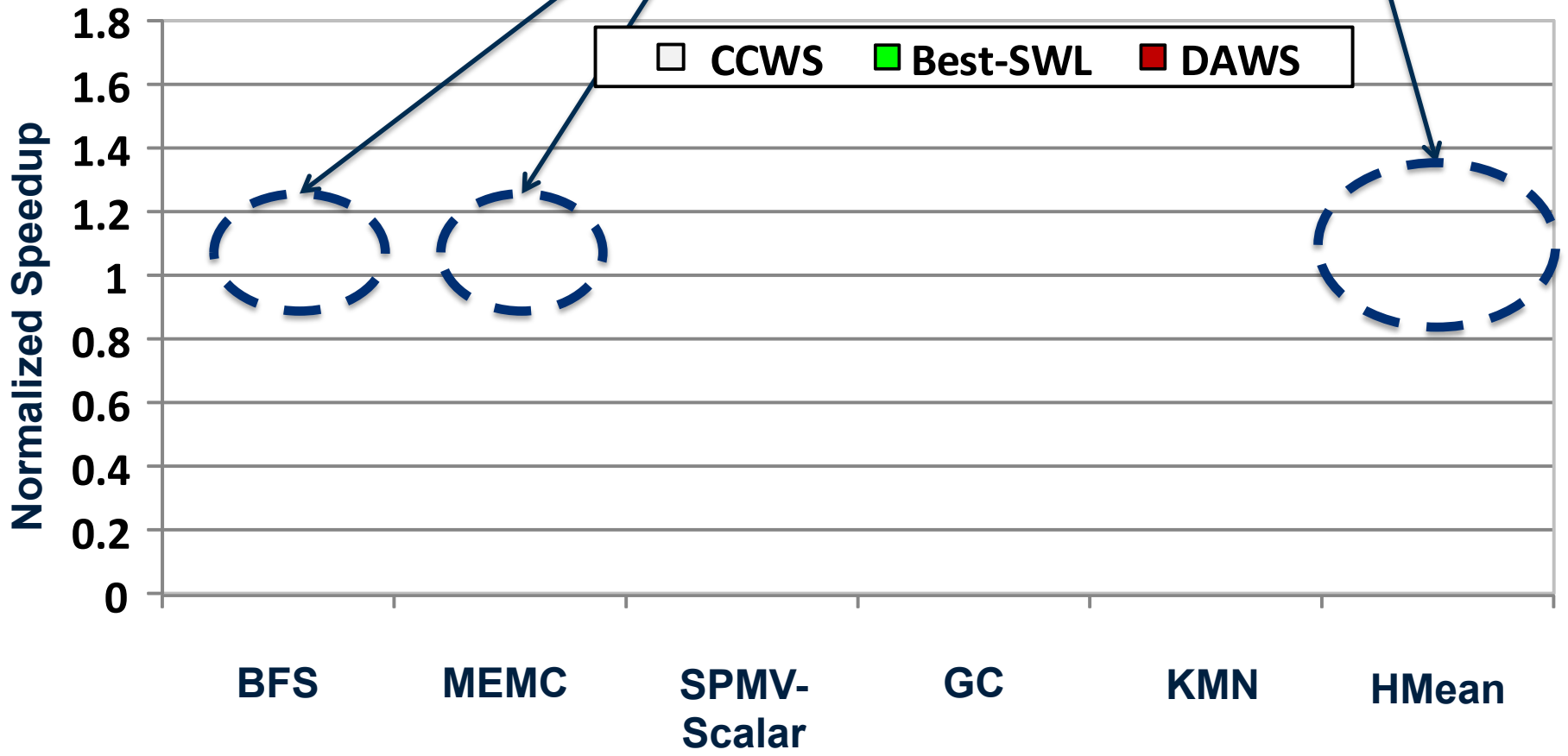
# Cache-Sensitive Applications

- Breadth First Search (BFS)
- Memcached-GPU (MEMC)
- Sparse Matrix-Vector Multiply (SPMV-Scalar)
- Garbage Collector (GC)
- K-Means Clustering (KMN)

Cache-Insensitive
Applications in paper

# Results

**Outperform Best-SWL in highly branch divergent**

**Overall 26% improvement over CCWS**

□ CCWS  ■ Best-SWL  ■ DAWS

Normalized Speedup: 1.8, 1.6, 1.4, 1.2, 1, 0.8, 0.6, 0.4, 0.2, 0

BFS  MEMC  SPMV-Scalar  GC  KMN  HMean

# Summary

**Questions?**

**Problem**

**Divergent loads in GPU programs.**
- **Software solutions complicate programming**

**Solution**

**DAWS**
- **Captures opportunities by accounting for divergence**

**Result**

Overall **26%** performance improvement over CCWS
Case Study: Divergent code performs **within 4%** code optimized to minimize divergence