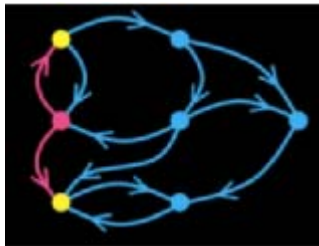# GBTL-CUDA: Graph Algorithms and Primitives for GPUs

Peter Zhang, Samantha Misurda, Marcin Zalewski, Scott McMillan, Andrew Lumsdaine

# What is this talk about?

- ## Graph BLAS

  – an emerging paradigm for graph computation

  – programs new graph algorithms in a highly abstract *language of linear algebra*.

  – executes in a wide variety of programming environments

- ## Our implementation of Graph BLAS

  – Graph BLAS Template Library (GBTL)

  – High-level C++ frontend

  – Switchable backends: CUDA and sequential

  – Released at: https://github.com/cmu-sei/gbtl
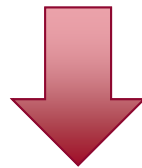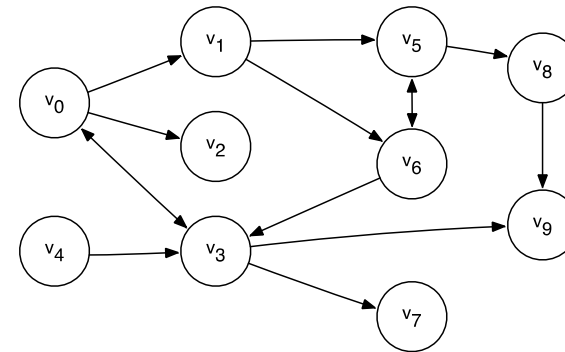
Graph BLAS Forum
http://www.graphblas.org

Software Engineering Institute
Carnegie Mellon University

CREST
Indiana University

# Graph algorithms meet hardware: love lost in translation?

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies
CREST

3

INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington

# Current Graph Algorithms

High Level Algorithm
(BFS, MIS, MST, SSSP)

Implementation Concerns

**Sequential**
- for loops
- contiguous memory
- …

**Data Parallel**
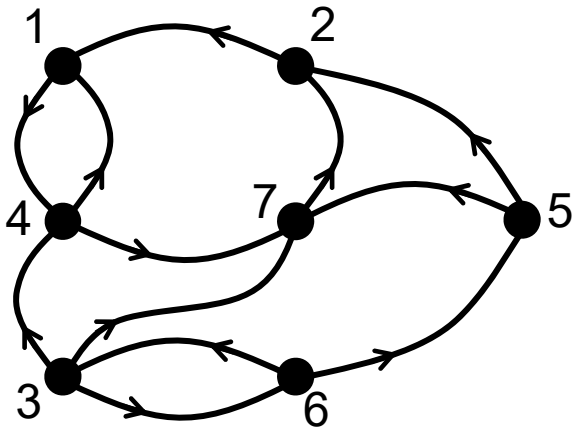- Parallel kernels
- PCI Memory transfer
- …

**Distributed**
- Message passing
- Workload balancing
- …

CPU

CUDA

MPI

# It's the representation…



- Adjacency matrix representation of graphs in graph theory (Kőnig 1931)
- Matrices?!

# BLAS for Linear Algebra

| Application |
|:---:|

| Unified BLAS Interface |
|:---:|

| X86-specific BLAS Optimization | ARM-specific BLAS Optimization | PPC-specific BLAS Optimization |
|:---:|:---:|:---:|
| X86 Architecture | ARM Architecture | PPC Architecture |

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies
6
INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington
CREST

# Graph BLAS

Same Application Cross Platform

**Application**

Same Interface Cross Platform

**Unified Graph BLAS Interface**

Architecture Specific Concerns

- OpenMP
- Pthreads
- Shared memory

- Asynchronous memory transfer
- Texture memory
- Dynamic Parallelism

- Memory coalescing
- Barrier control

**CPU**

**CUDA**

**MPI**

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies
7
INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington
CREST

# Graph BLAS

- A community effort to define a set of primitives used to describe graph algorithms using sparse linear algebra
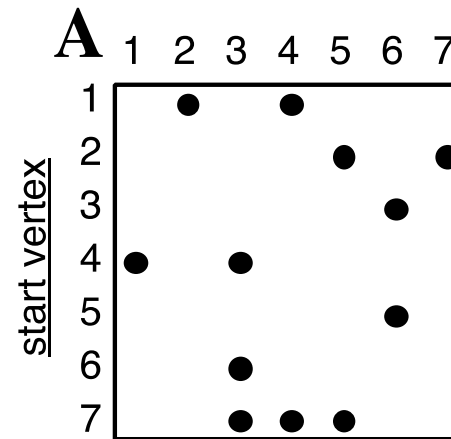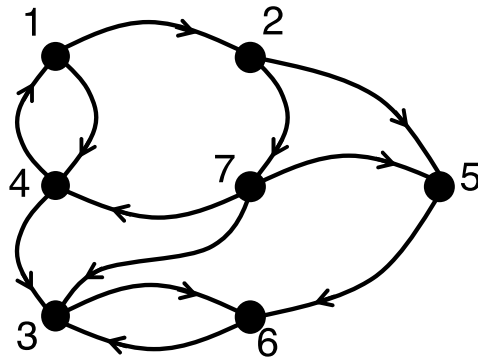
- Rich data structures, algebraic abstraction
  - Sparse adjacency matrices represent graphs
  - Semirings to define specific behavior

# Graph BLAS: Semiring

- The standard arithmetic semiring:

Multiplication

Domain D ——————— $< D, +, \times, 0, 1 >$ ——————— Multiplicative Identity

Addition

Additive Identity

- The "*minPlus*" semiring, for BFS with parents:

"Multiplication"

Domain D ——————— $< D, min, +, \infty, 0 >$ ——————— Multiplicative Identity

"Addition"

Additive Identity

# Graph BLAS: BFS example

- Breadth-first Search (BFS): can be represented by matrix-vector multiplications in linear algebra

- Wavefront: vector

- One multiplication operation results in the next wave front

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies
10
INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington
CREST

# Graph BLAS: BFS traversal

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies
11
INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington
CREST

# Graph BLAS Template Library (GBTL)

- ## A C++ implementation of Graph BLAS
  - Allows for generic programming and metaprogramming

- ## A frontend-backend design
  - Uniform frontend for algorithm abstraction
    - generic semantic checks
    - simplifies the templates
  - Hardware-specific backend optimized for different architecture

- ## A separation of concerns: render unto hardware experts hardware-specific optimizations

| Graph Analytic Applications |
| Graph Algorithms |
| GraphBLAS API (Separation of Concerns) |
| Graph Primitives (tuned for hardware) |
| Hardware Architecture |

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies
CREST

12

INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington

# GBTL: Frontend and Backend

- Graph BLAS API: boundary between the algorithms and hardware-specific implementations

- Frontend forwards calls to backend namespace via C++ templates

  - performs generic semantic checks and implementation independent operations
  - simplifies templates passed in by user for meta programming

| Graph Analytic Applications |
| Graph Algorithms |
| GraphBLAS API (Separation of Concerns) |
| Graph Primitives (tuned for hardware) |
| Hardware Architecture |

graphblas namespace

Forward

graphblas::backend namespace

# GBTL: Algorithm Example

```
1  // wavefront initialized with root vertex = 0
2  bfs(graph, wavefront) {
3      vector visited = wavefront;
4
5      while(!wavefront.empty()) {
6          // increment level in next wavefront
7          wavefront = vXm(wavefront,
8                          graph,
9                          Add1NotZero);
10
11         // filter out already visited vertices
12         // if the vertices have values less than
13         // current level in visted vector
14         wavefront = eWiseMult(wavefront,
15                              visited,
16                              Mult);
17
18         //update visited vector by filtered wavefront
19         visited = eWiseAdd(wavefront,
20                              visited,
21                              throwException);
22     }
23     return visited;
24 }
```

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies
CREST

14

INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington

# GBTL: Frontend Matrix

- Frontend `Matrix` class: an opaque data structure, uniform across backends

  Frontend Matrix Object Construction:

  ```
  Matrix <double, DenseMatrixTag,
  DirectedMatrixTag> matrix(...);
  ```

- User can provide hints to frontend Matrix at construction time through parameter pack, backend can make decisions based on hints

- Backend `Matrix` classes: specialized for hardware and implementation

# GBTL: Frontend Matrix Class

```
1   // TagsT template parameters provide hints
2   template <typename ScalarT, typename... TagsT>
3   class Matrix :
4       public backend::Matrix<ScalarT, TagsT...>
5   {
6   public:
7       typedef ScalarT ScalarType;
8
9       // Empty construction; fixed dimensions
10      Matrix(IndexType num_rows, IndexType num_cols);
11
12      // Other frontend matrix interface...
13
14  private:
15      // immutable dimensions:
16      IndexType const m_num_rows, m_num_cols;
17
18      // opaque backend implementation
19      backend::Matrix<
20              ScalarType,
21              detail::SparsenessCategoryTagT,
22              TagsT...>
23          m_matrix;
24  };
```
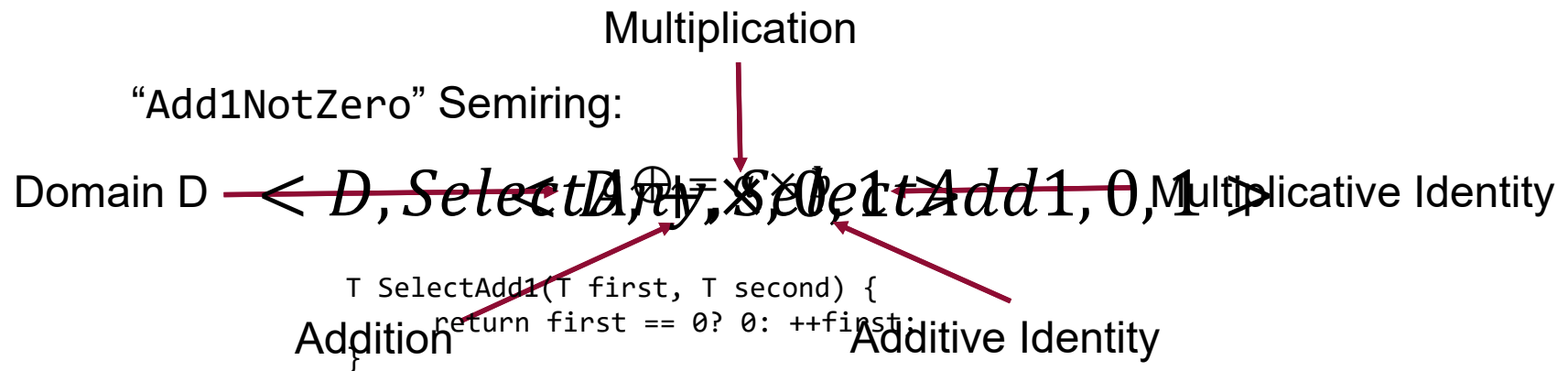
# GBTL: Algorithm Example

```
1  // wavefront initialized with root vertex = 0
2  bfs(graph, wavefront) {
3      vector visited = wavefront;
4
5      while(!wavefront.empty()) {
6          // increment level in next wavefront
7          wavefront = vXm(wavefront,
8                          graph,
9                          Add1NotZero);
10
11         // filter out already visited vertices
12         // if the vertices have values less than
13         // current level in visted vector
14         wavefront = eWiseMult(wavefront,
15                               visited,
16                               Mult);
17
18         //update visited vector by filtered wavefront
19         visited = eWiseAdd(wavefront,
20                            visited,
21                            throwException);
22     }
23     return visited;
24 }
```

# GBTL: vXm, semiring overloading

```
template<typename AVectorT,
         typename BMatrixT,
         typename CVectorT,
         typename SemiringT = graphblas::ArithmeticSemiring<T>,
         typename AccumT = graphblas::math::Assign<T> >         ⊕
inline void vxm(AVectorT const &a,
                BMatrixT const &b,
                CVectorT &c,
                SemiringT s = SemiringT(),
                AccumT accum = AccumT())
{
    vector_multiply_dimension_check(a, b.get_shape());
    backend::vxm(a.m_vec, b.m_mat, c.m_vec, s, accum);
}
```

Multiplication

"Add1NotZero" Semiring:

Domain D — $< D, SelectArg, SelectAdd1, 0, 1 >$ Multiplicative Identity

```
T SelectAdd1(T first, T second) {
    return first == 0? 0: ++first;
}
```

Addition                                    Additive Identity

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies
18
INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington
CREST

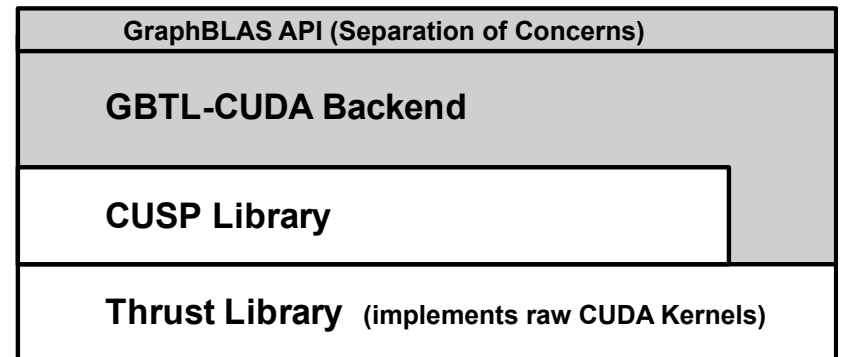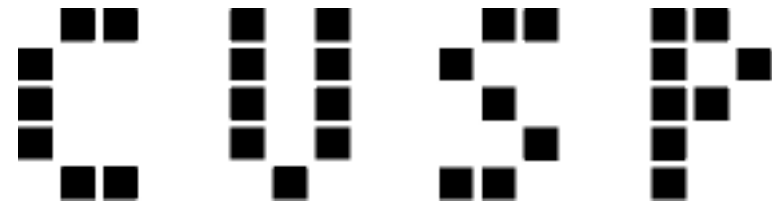# GBTL: Algorithm Example

```
1  // wavefront initialized with root vertex = 0
2  bfs(graph, wavefront) {
3      vector visited                    wavefront
4
5    while(!wavefro
6         // increme
7         wavefront
8                                      g
9                              A      tZero);
10
11        // filter out already visited vertices
12        // if the vertices have values less than
13        // curren
14        wavefront
15
16
17
18        //update visited vector by filtered wavefront
19        visited = eWiseAdd(wavefront,
20                           visited,
21                           throwException);
22    }
23    return visited;
24 }
```

wavefront =
wavefront $\otimes$ visited

visited =
wavefront $\oplus$ visited

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies

19

INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington

CREST
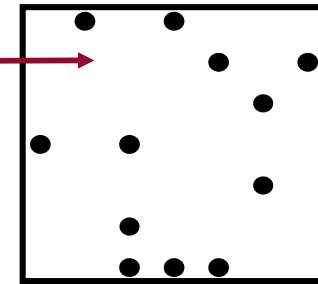
# GBTL: GPU Backend

- Implemented using CUSP: parallel algorithms and data structures for sparse linear algebra
- CUSP: built on top of the Thrust, a C++ library with GPU programming primitives
- Generalization meets performance



| GraphBLAS API (Separation of Concerns) | | |
|---|---|---|
| GBTL-CUDA Backend | | |
| CUSP Library | | |
| Thrust Library  (implements raw CUDA Kernels) | | |

# GPU Backend Data Type: Sparse Matrix

- We use sparse matrices to improve storage efficiency

- Sparse matrices have unstored elements called *structural zeros*
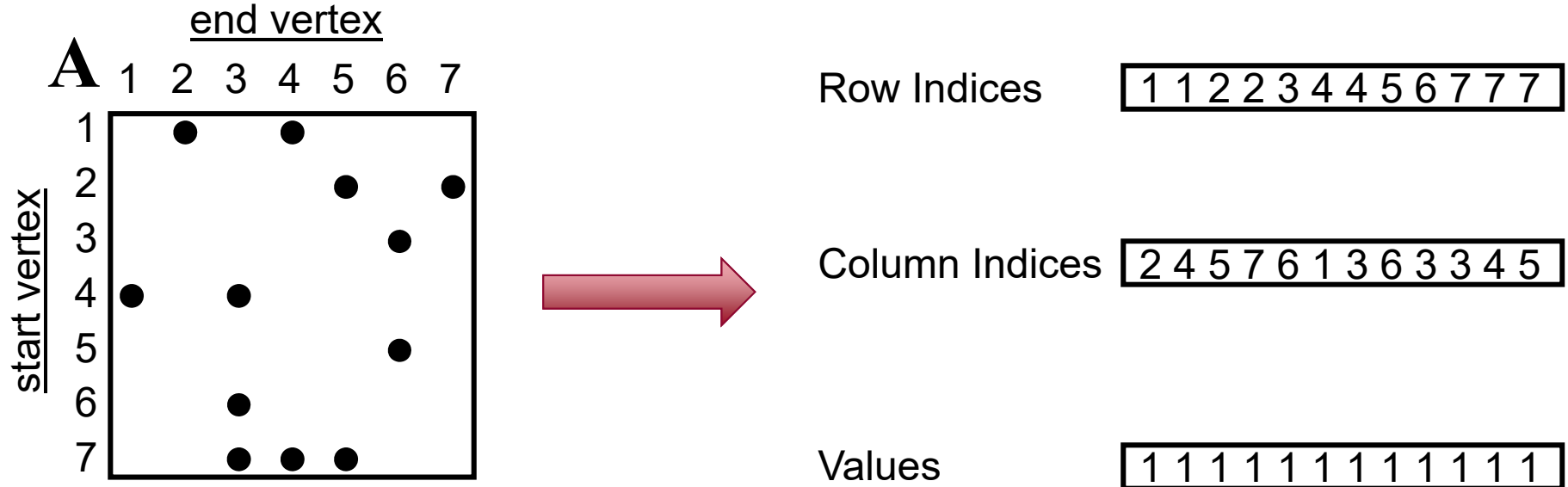
Blank space in matrix

- Different sparse matrix formats: Compressed Sparse Row (CSR), Coordinate (COO), List Of Lists (LIL)

- Backend makes decision on format-of-choice, based on hardware layout

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies
21
INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington
CREST

# Sparse Matrix Example: COO
## A tale of three vectors

- ## COO matrices: enables easy stream processing
  - Regularity in data layout



**end vertex**

A

|       | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|-------|---|---|---|---|---|---|---|

start vertex

Row Indices: 1 1 2 2 3 4 4 5 6 7 7 7

Column Indices: 2 4 5 7 6 1 3 6 3 3 4 5

Values: 1 1 1 1 1 1 1 1 1 1 1 1

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies

22

INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington

CREST

# GPU Backend: mXm Scaling



mXm scaling in Millions of Floating-point Operations Per Second (MFLOPS)

- Erdős–Rényi graphs
- Average of 16 runs

# GPU Backend: BFS Performance

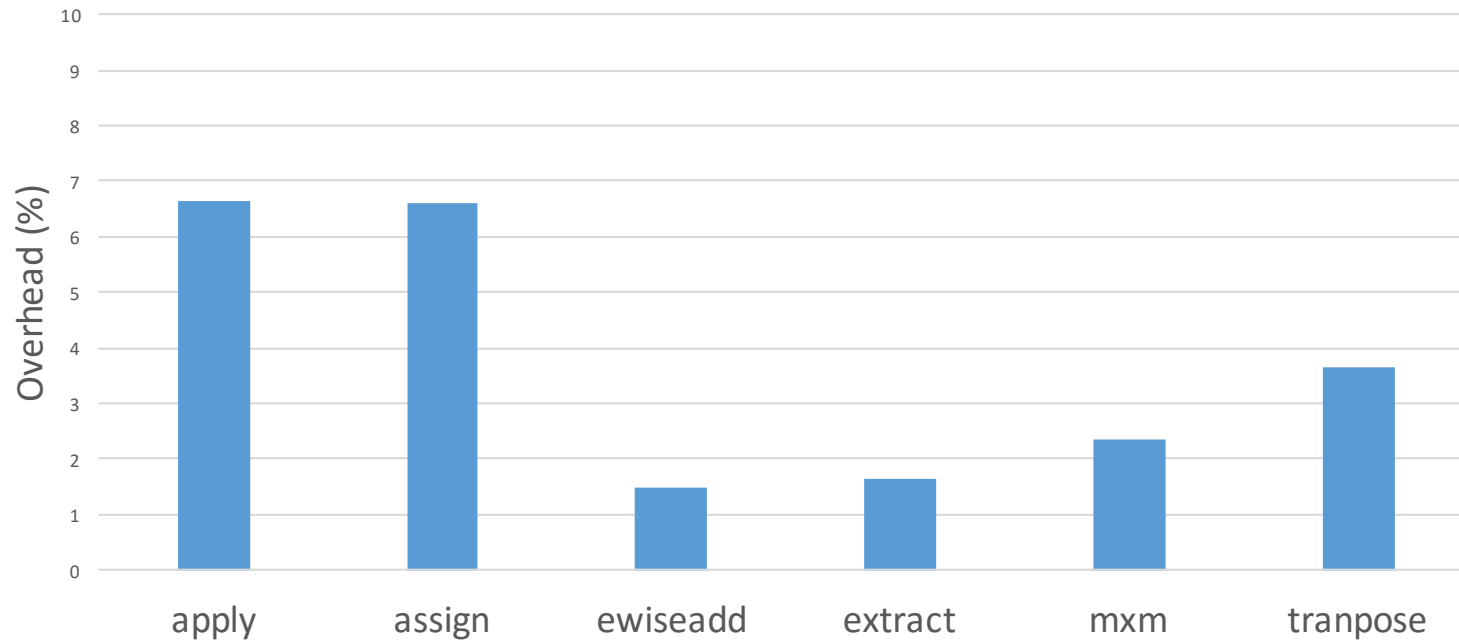- We test runtime of our BFS algorithm on several real world graphs

| Graph Name | # Vertices | # Edges | Runtime(ms) | MTEPS(*) |
|---|---|---|---|---|
| Journals | 124 | 12,068 | 5.76 | 2.1 |
| G43 | 1,000 | 19,980 | 14.61 | 1.4 |
| ship_003 | 121,728 | 3,777,036 | 558.95 | 6.8 |
| belgium_osm | 1,441,295 | 3,099,940 | 10502.4 | 0.3 |
| roadNet-CA | 1,971,281 | 5,533,214 | 4726.21 | 1.2 |
| delaunay_n24 | 16,777,216 | 100,663,202 | 65507.7 | 1.5 |

*MTEPS = Millions of Traversed Edges Per Second

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies
CREST
24
INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington

# GBTL: API Overhead

GraphBLAS Template Library Overhead



Environment:
- NVIDIA GPU
- Overhead of API call compared against direct CUSP call

Methodology:
- Average of 16 runs on Erdős–Rényi graphs generated using the same dimension and sparsity

# Recap and Future Plans

- ## Graph BLAS:
  - Graph algorithms on sparse linear algebra primitives

- ## Graph BLAS Template Library (GBTL):
  - Extensibility meets performance
  - Abstraction layer: translator with low overhead penalty
  - Proof-of-concept: it works well!

- ## Future Plans
  - Multi-GPU backend
  - Distributed CPU/GPU backend
  - Participate in community discussion on future specifications

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies

26

CREST

INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington

# Thank You

# Questions?

INDIANA UNIVERSITY
Center for Research in Extreme Scale Technologies

27

INDIANA UNIVERSITY
SCHOOL OF INFORMATICS AND COMPUTING
Bloomington

CREST