# Alternative GPU friendly assignment algorithms

**Paul Richmond and Peter Heywood**
Department of Computer Science
The University of Sheffield

# Graphics Processing Units (GPUs)

# Context: GPU Performance

## Serial Computing

~40
GigaFLOPS

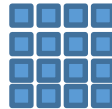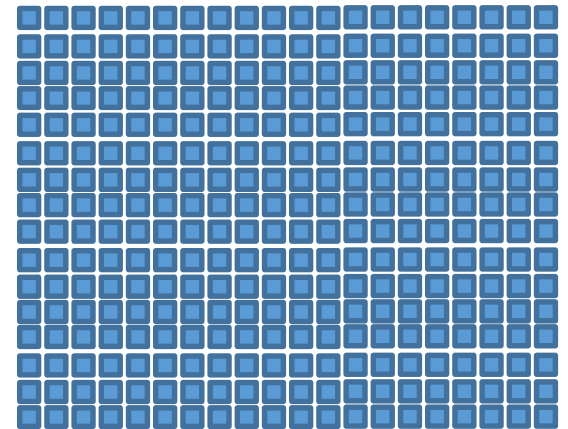1 core

## Parallel Computing
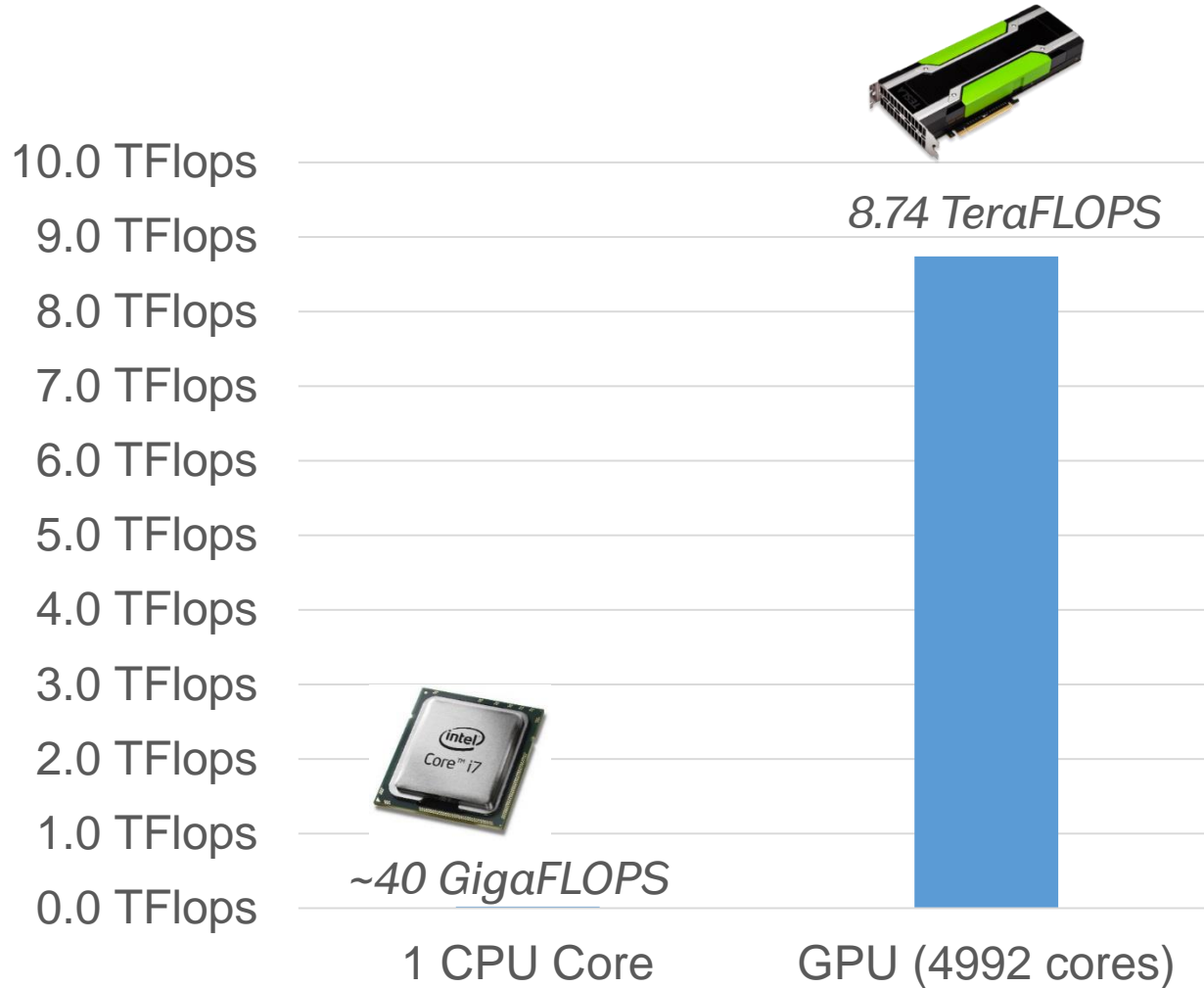
384
GigaFLOPS

16
cores

## Accelerated Computing

8.74
TeraFLOPS

**4992 cores**

10.0 TFlops
9.0 TFlops
8.0 TFlops
7.0 TFlops
6.0 TFlops
5.0 TFlops
4.0 TFlops
3.0 TFlops
2.0 TFlops
1.0 TFlops
0.0 TFlops

*8.74 TeraFLOPS*

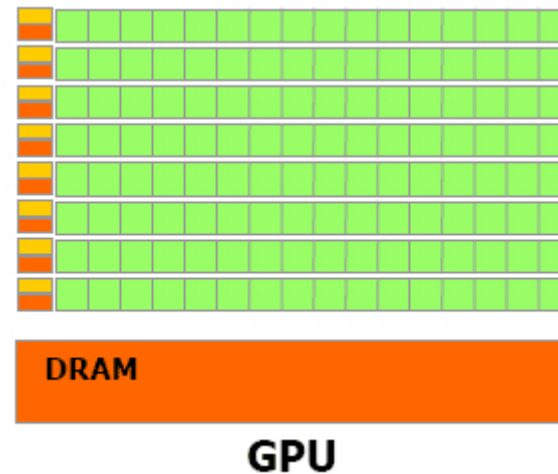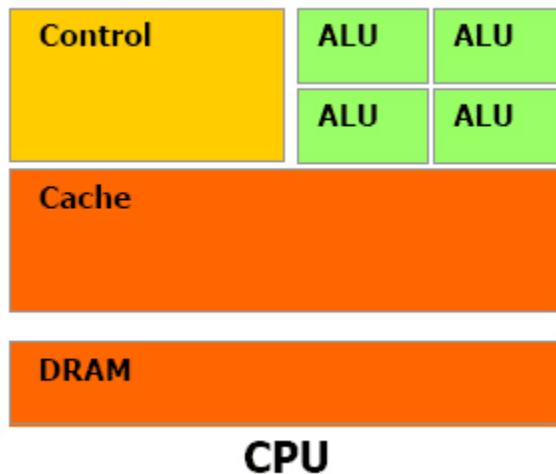*~40 GigaFLOPS*

1 CPU Core          GPU (4992 cores)
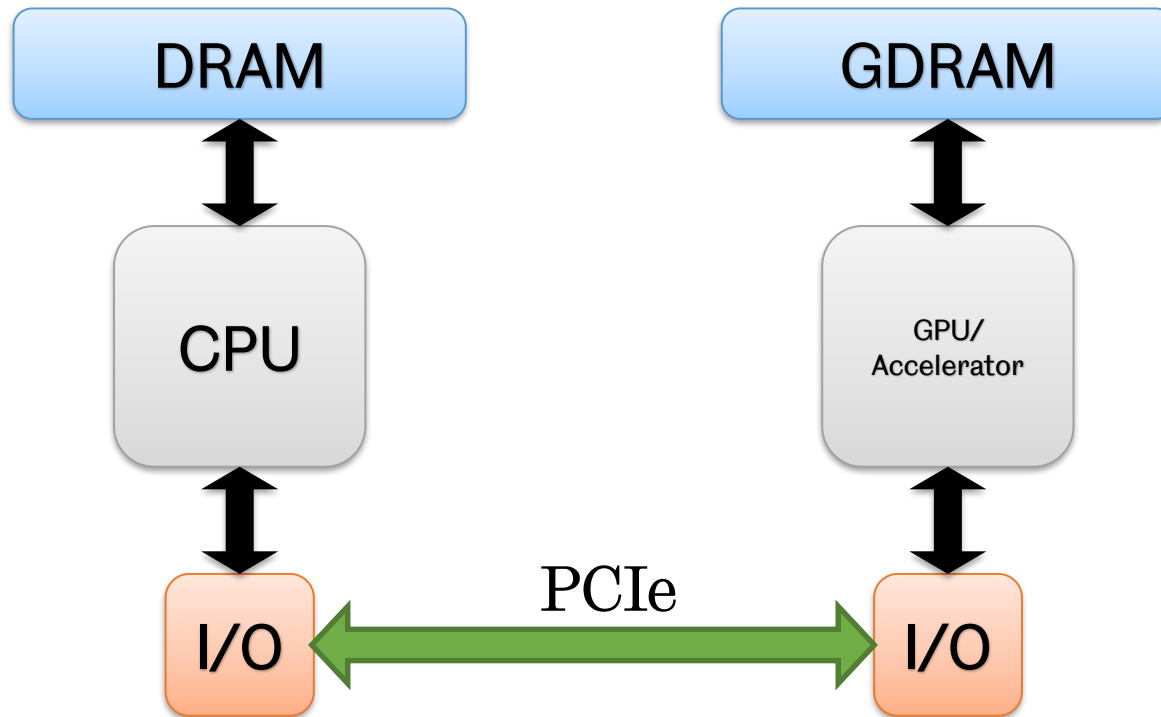
6 hours *CPU* time
vs.
**1 minute *GPU* time**

# Accelerators

- Much of the functionality of CPUs is unused for HPC
  - Complex Pipelines, Branch prediction, out of order execution, etc.

- Ideally for HPC we want: **Simple**, **Low Power** and **Highly Parallel** cores

# An accelerated system



- Co-processor not a CPU replacement
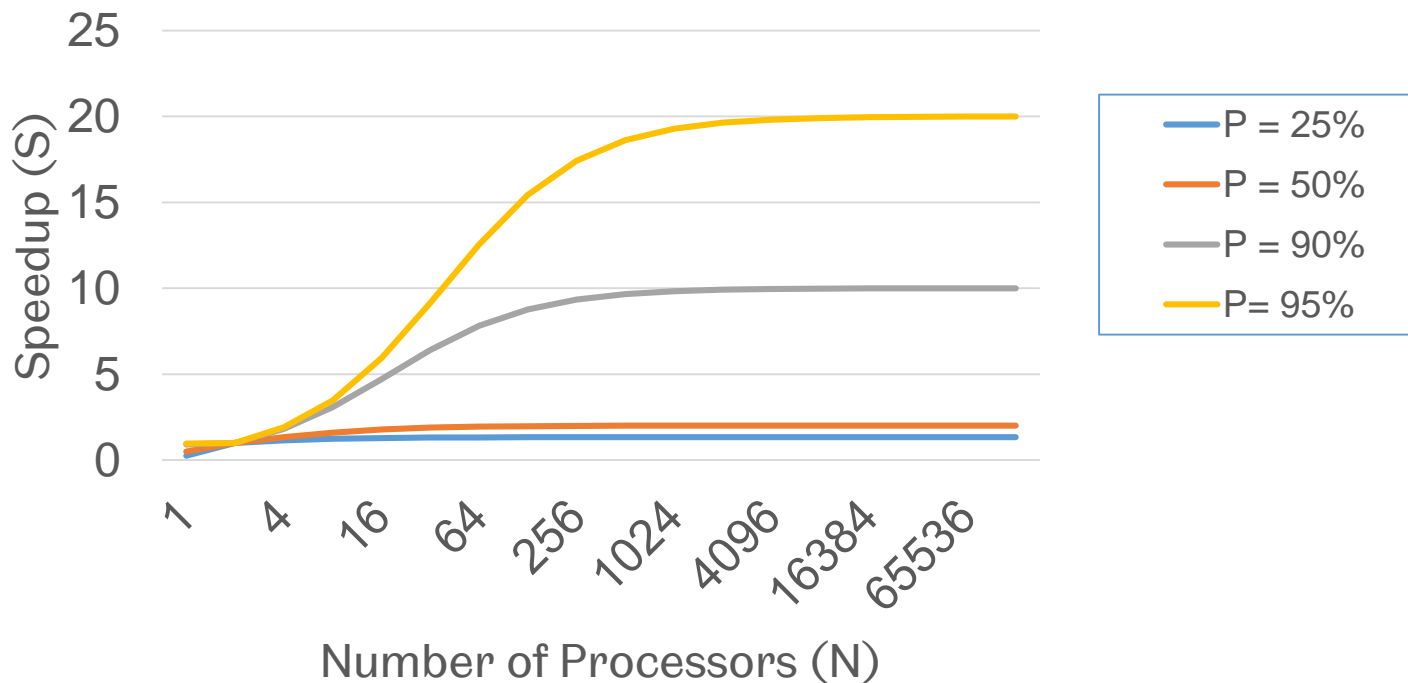
# Thinking Parallel

- Hardware considerations
  - High Memory Latency (PCI-e)
  - Huge Numbers of processing cores
- Algorithmic changes required
  - High level of parallelism required
  - Data parallel != task parallel

"If your problem is not parallel then think again"

# Amdahl's Law

$$Speedup\ (S) = \frac{1}{\frac{P}{N} - (1-P)}$$



- Speedup of a program is limited by the proportion than can be parallelised

- Addition of processing cores gives diminishing returns
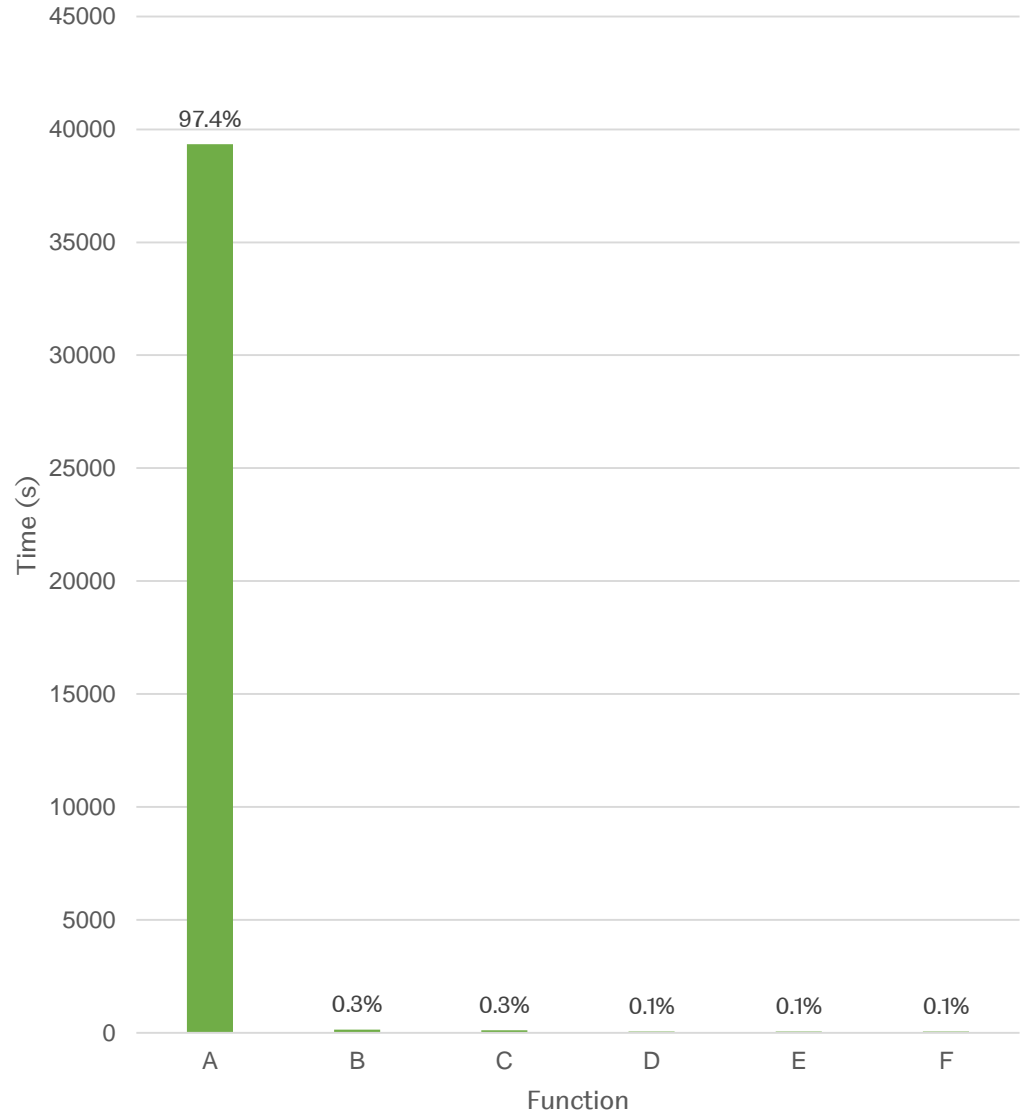
# SATALL Optimisation

## Profile the Application

- 11 hour runtime

- Function A
  - 97.4% runtime
  - 2000 calls

- Hardware
  - Intel Core i7-4770k 3.50GHz
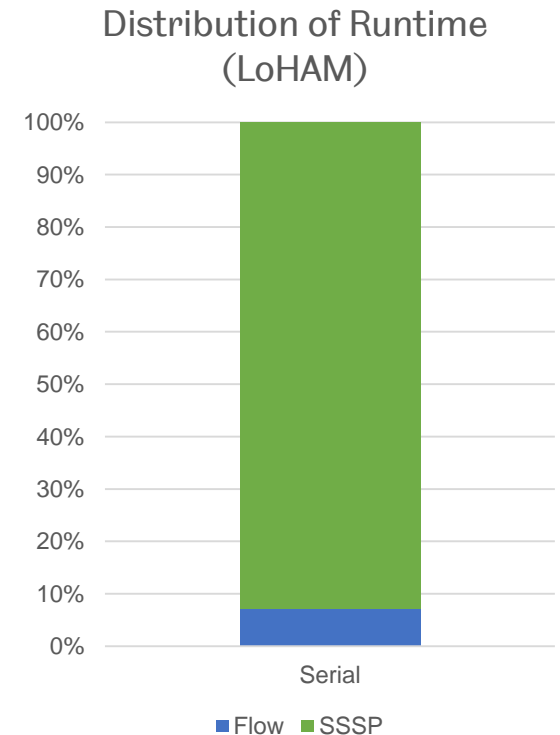  - 16GB DDR3
  - Nvidia GeForce Titan X

**Time per function for Largest Network (LoHAM)**



Bar chart: Time (s) vs Function. A: 97.4%, B: 0.3%, C: 0.3%, D: 0.1%, E: 0.1%, F: 0.1%

# Function A

- Input
  - Network (directed weighted graph)
  - Origin-Destination Matrix

- Output
  - Traffic flow per edge

- 2 Distinct Steps

  1. **Single Source Shortest Path (SSSP)**

     *All-or-Nothing Path*

     For each origin in the O-D matrix

  2. **Flow Accumulation**

     Apply the OD value for each trip to each link on the route

### Distribution of Runtime (LoHAM)

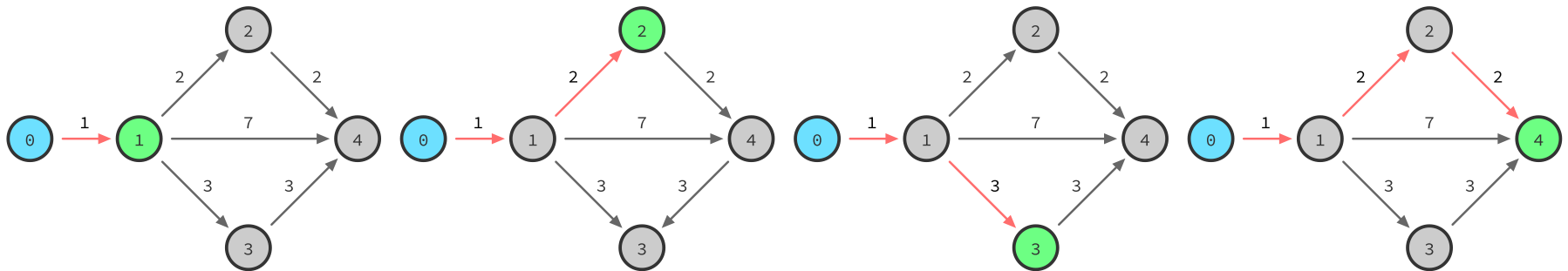| | Serial |
|---|---|
| 100% | |
| 90% | |
| 80% | |
| 70% | |
| 60% | |
| 50% | |
| 40% | |
| 30% | |
| 20% | |
| 10% | |
| 0% | |

■ Flow  ■ SSSP

# Single Source Shortest Path

For a single **Origin Vertex** (Centroid)

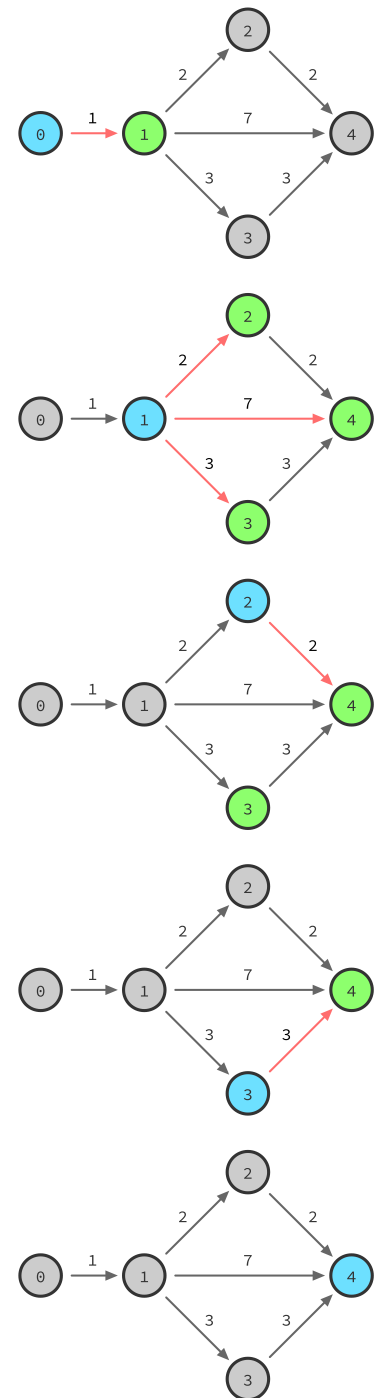Find the route to each **Destination Vertex**

With the **Lowest Cumulative Weight** (Cost)
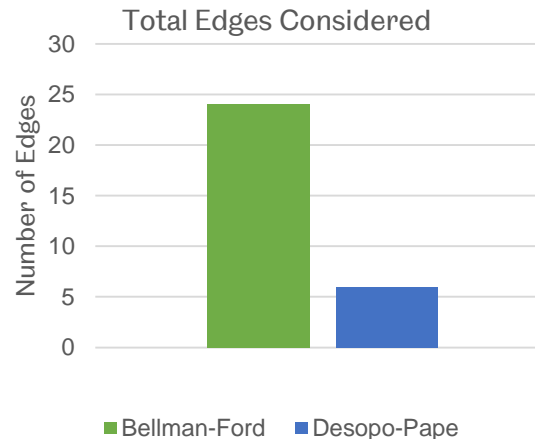
# Serial SSSP Algorithm

- D'Esopo-Pape (1974)
  - Maintains a **priority queue** of vertices to explore
    Highly Serial
  - Not a **Data-Parallel** Algorithm

- We must change algorithm to match the hardware

Pape, U. "Implementation and efficiency of Moore-algorithms for the shortest route problem." Mathematical Programming 7.1 (1974): 212-222.
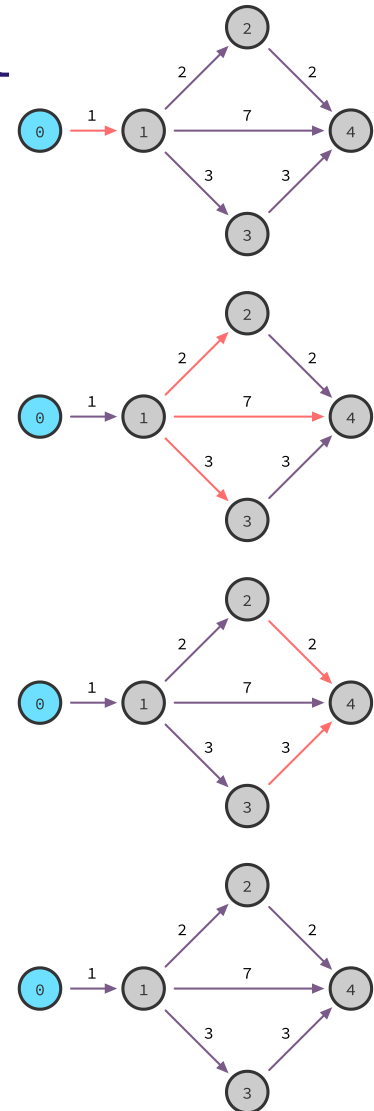
# Parallel SSSP Algorithm

- Bellman-Ford Algorithm (1956)
  - Poor serial performance & time complexity
  - Performs significantly more work
  - **Highly Parallel**
  - Suitable for GPU acceleration

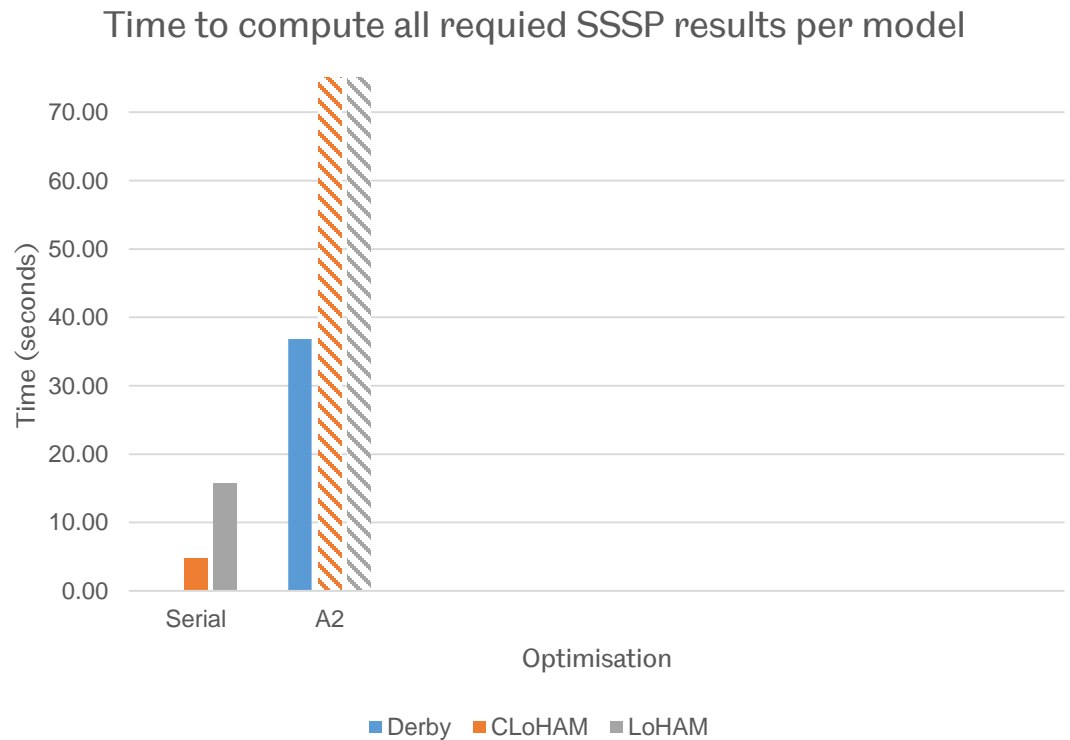Bellman, Richard. On a routing problem. 1956.

**Total Edges Considered**



Number of Edges

■ Bellman-Ford  ■ Desopo-Pape

# Implementation

- **A2** - Naïve Bellman-Ford using Cuda

- Up to 369x slower

- Striped bars continue off the scale

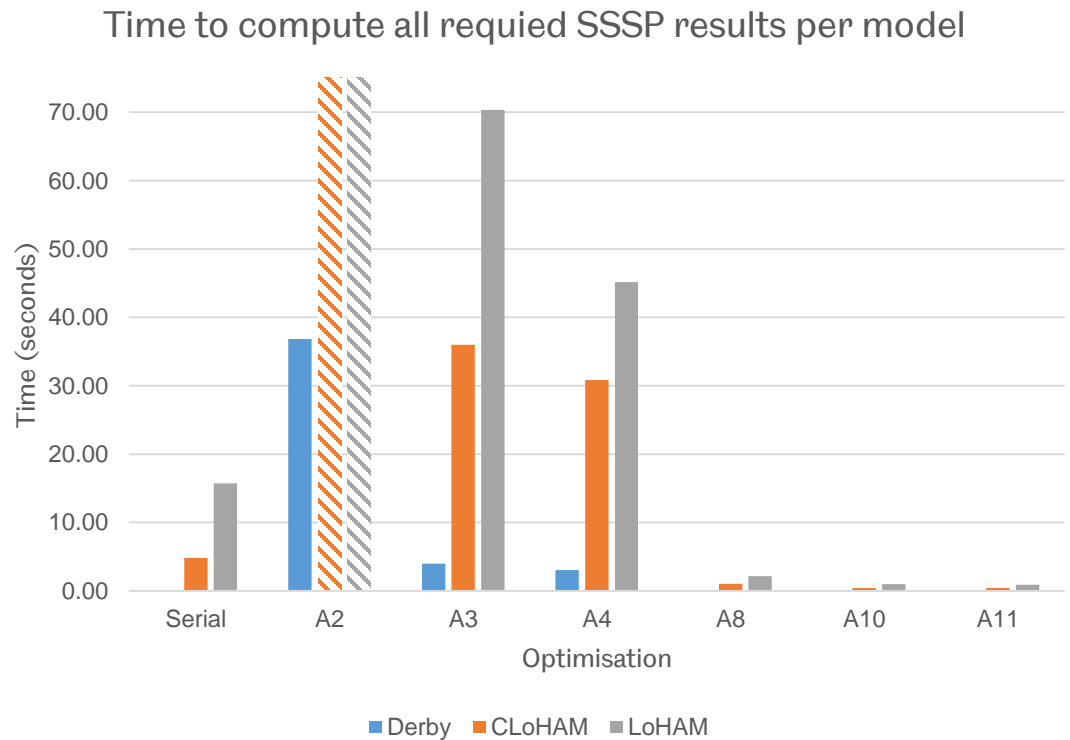| | |
|---|---|
| **Derby** | 36.5s |
| **CLoHAM** | 1777.2s |
| **LoHAM** | 5712.6s |

Time to compute all requied SSSP results per model

Legend: Derby, CLoHAM, LoHAM

X-axis: Optimisation (Serial, A2)
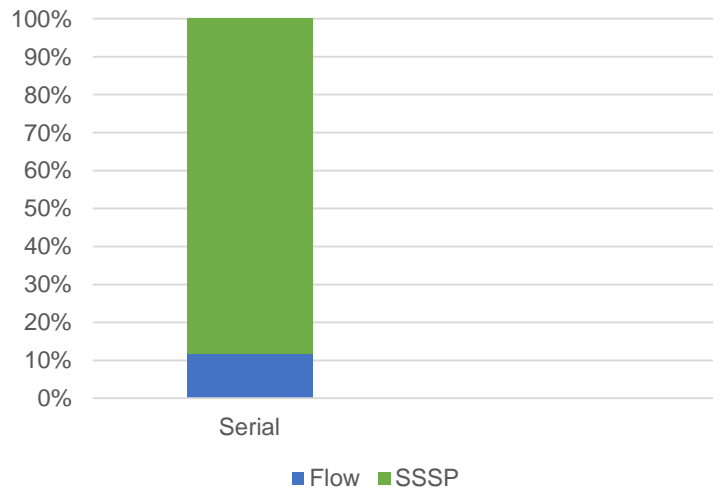Y-axis: Time (seconds)

# Implementation

- Followed iterative cycle of performance optimisations
- A3 – Early Termination
- A4 – Node Frontier
- **A8 – Multiple origins Concurrently**
  - SSSP for each Origin in the OD matrix
- A10 – Improved load Balancing
  - Cooperative Thread Array
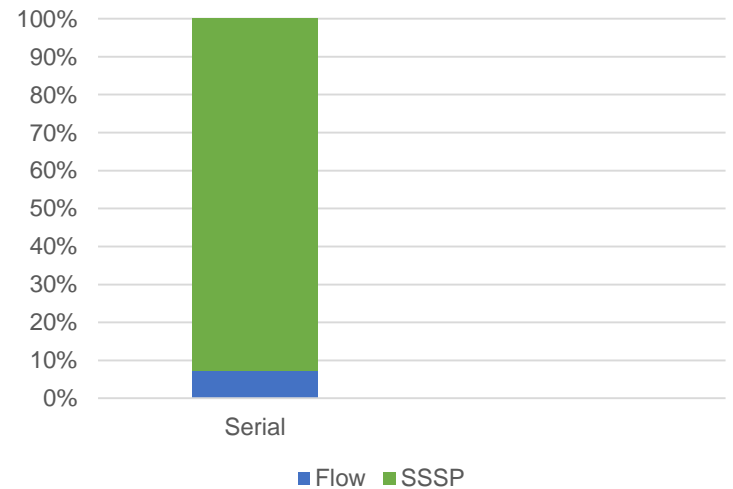- A11 – Improved array access

Time to compute all requied SSSP results per model



Legend: Derby, CLoHAM, LoHAM

# Limiting Factor (Function A)


Distribution of Runtime (CLoHAM)
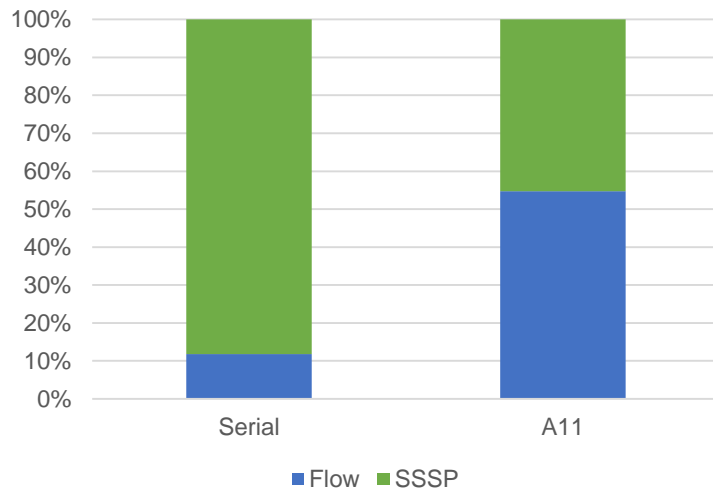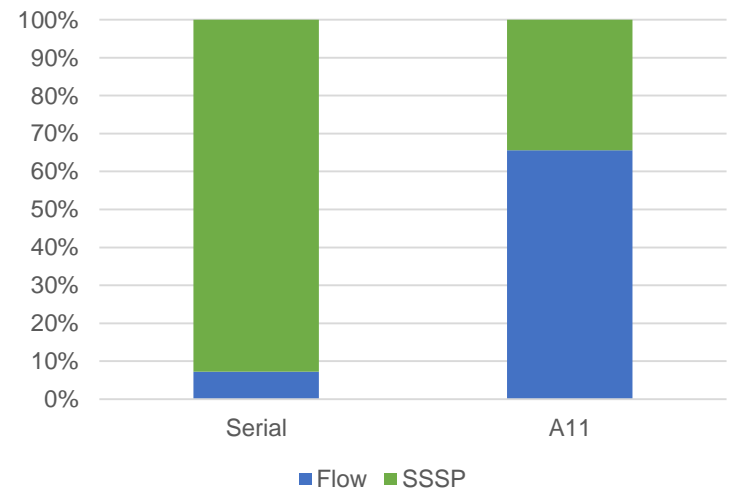

Distribution of Runtime (LoHAM)

# Limiting Factor (Function A)

- Limiting Factor has now changed
- Need to parallelise Flow Accumulation

Distribution of Runtime (CLoHAM)

Distribution of Runtime (LoHAM)
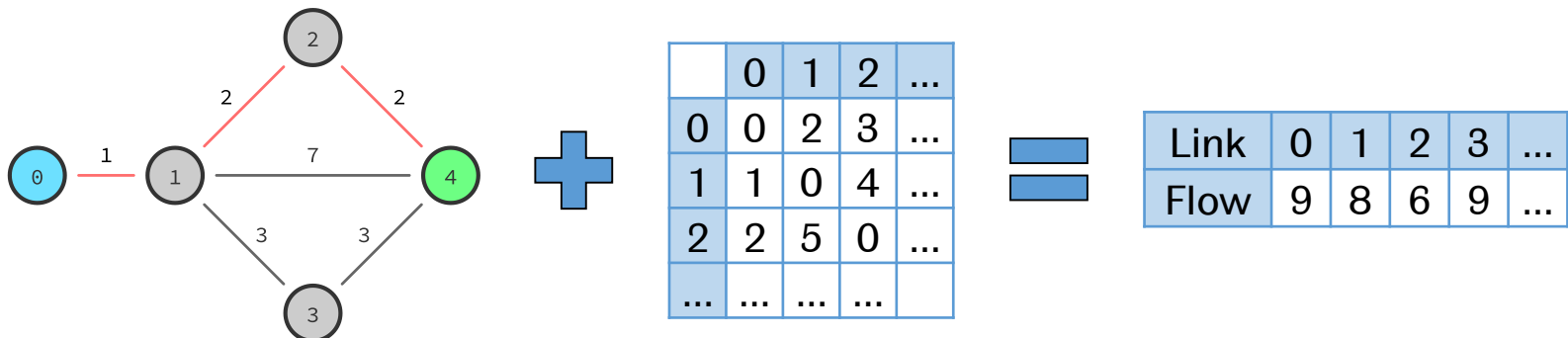
# Flow Accumulation

- Shortest Path + OD = Flow-per-link

    For each origin-destination pair

    > Trace the route from the destination to the origin increasing the flow value for each link visited
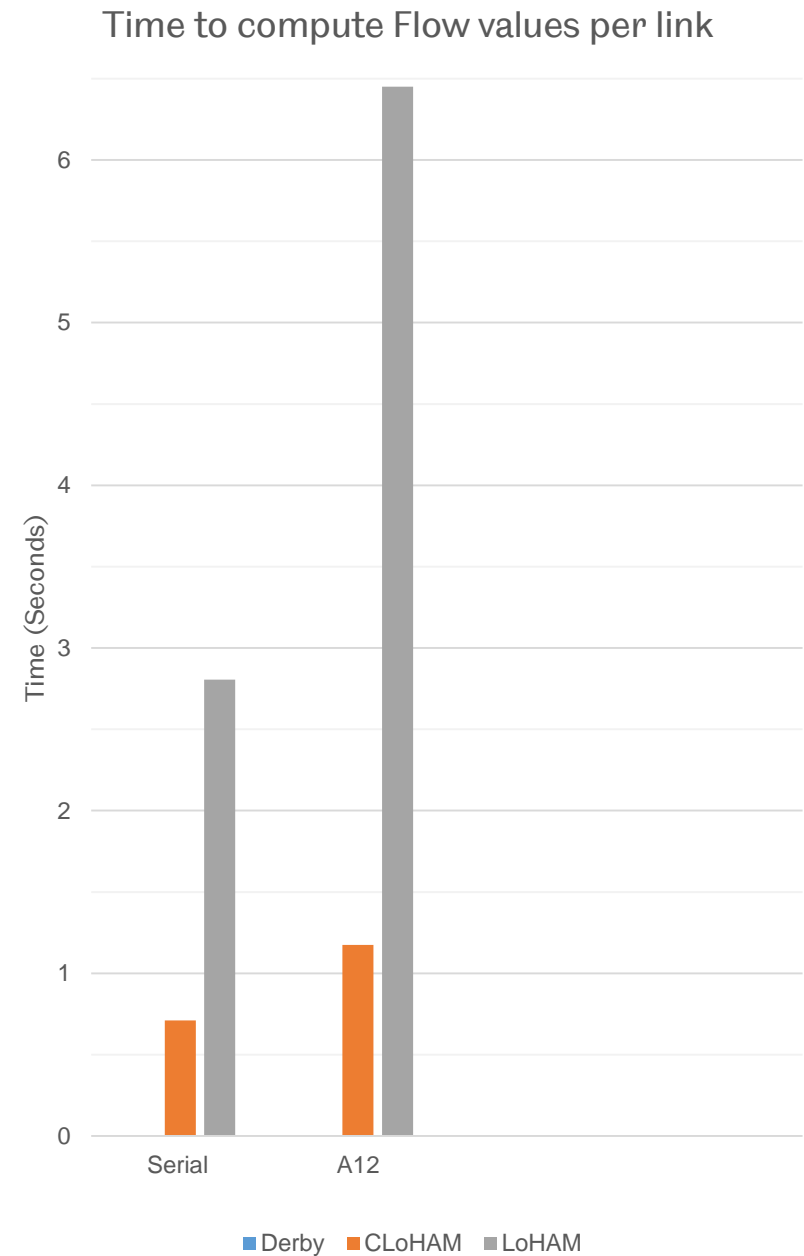
- Parallel problem

    - **But** requires synchronised access to shared data structure for all trips (atomic operations)
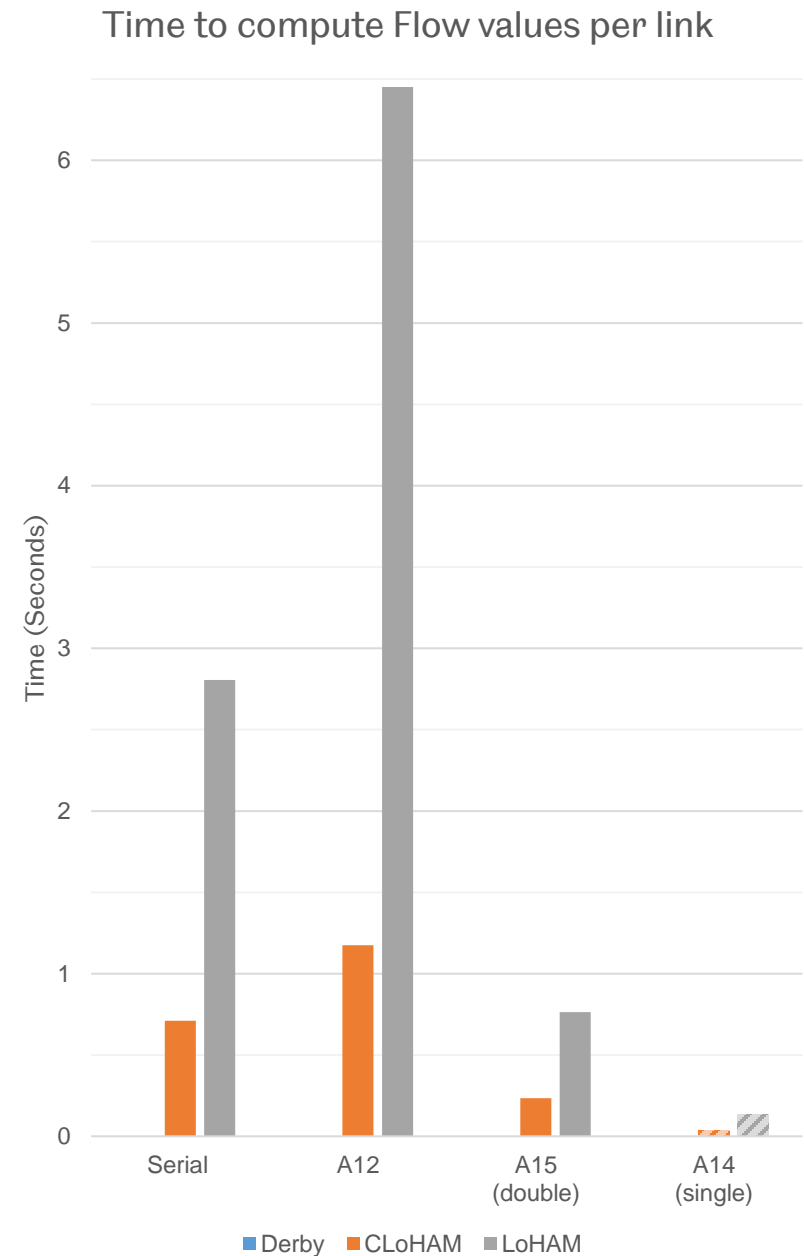
# Flow Accumulation

- **Problem:**
  - **A12** - lots of atomic operations serialise the execution

## Time to compute Flow values per link



Legend: Derby, CLoHAM, LoHAM

Y-axis: Time (Seconds)

X-axis categories: Serial, A12

# Flow Accumulation

- **Problem:**
  - **A12** - lots of atomic operations serialise the execution
- **Solutions:**
  - **A15** - Reduce number of atomic operations
    - Solve in batches using parallel reduction
  - **A14** - Use fast hardware-supported single precision atomics
    - Minimise loss of precision using multiple 32-bit summations
    - 0.000022% total error

## Time to compute Flow values per link

# Integrated Results

# Assignment Speedup relative to Serial

Serial
- LoHAM – 12h 12m
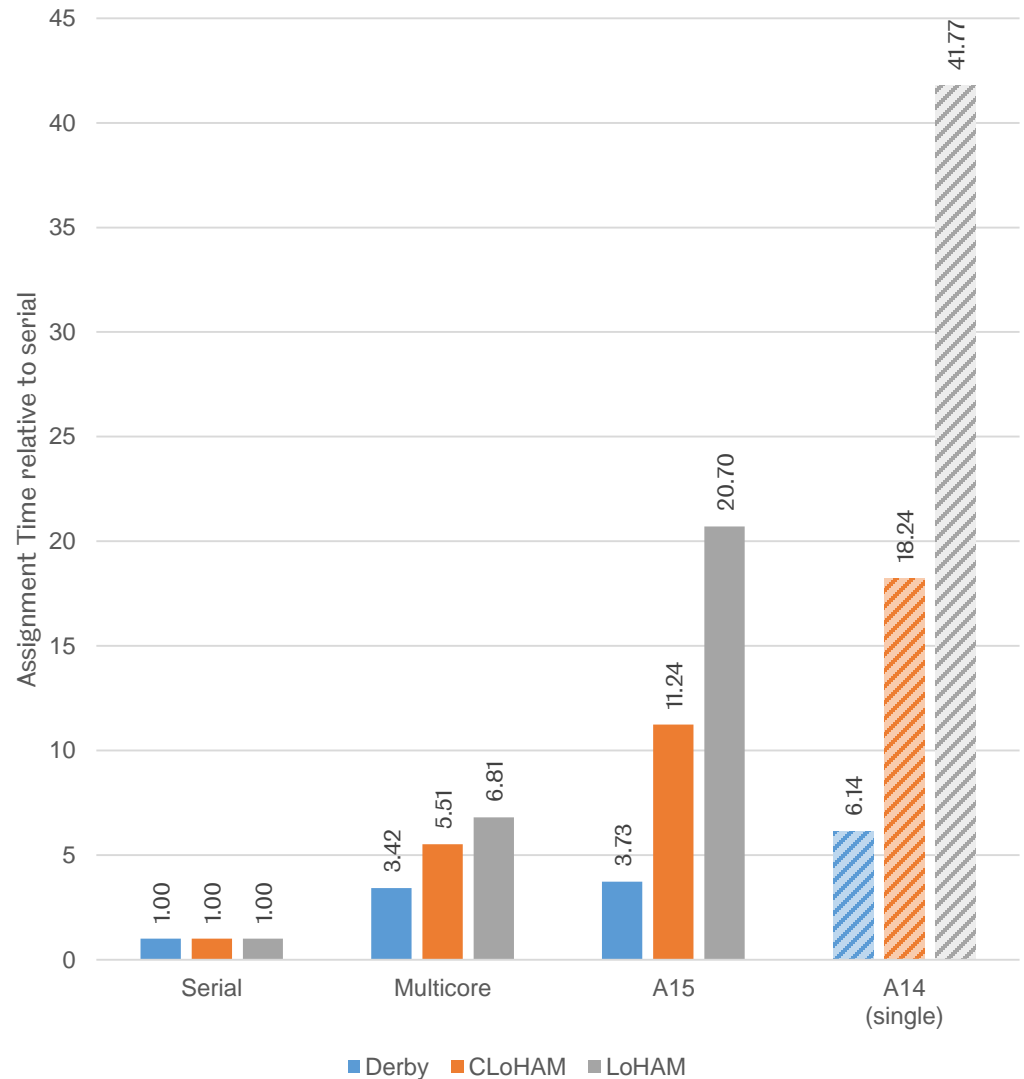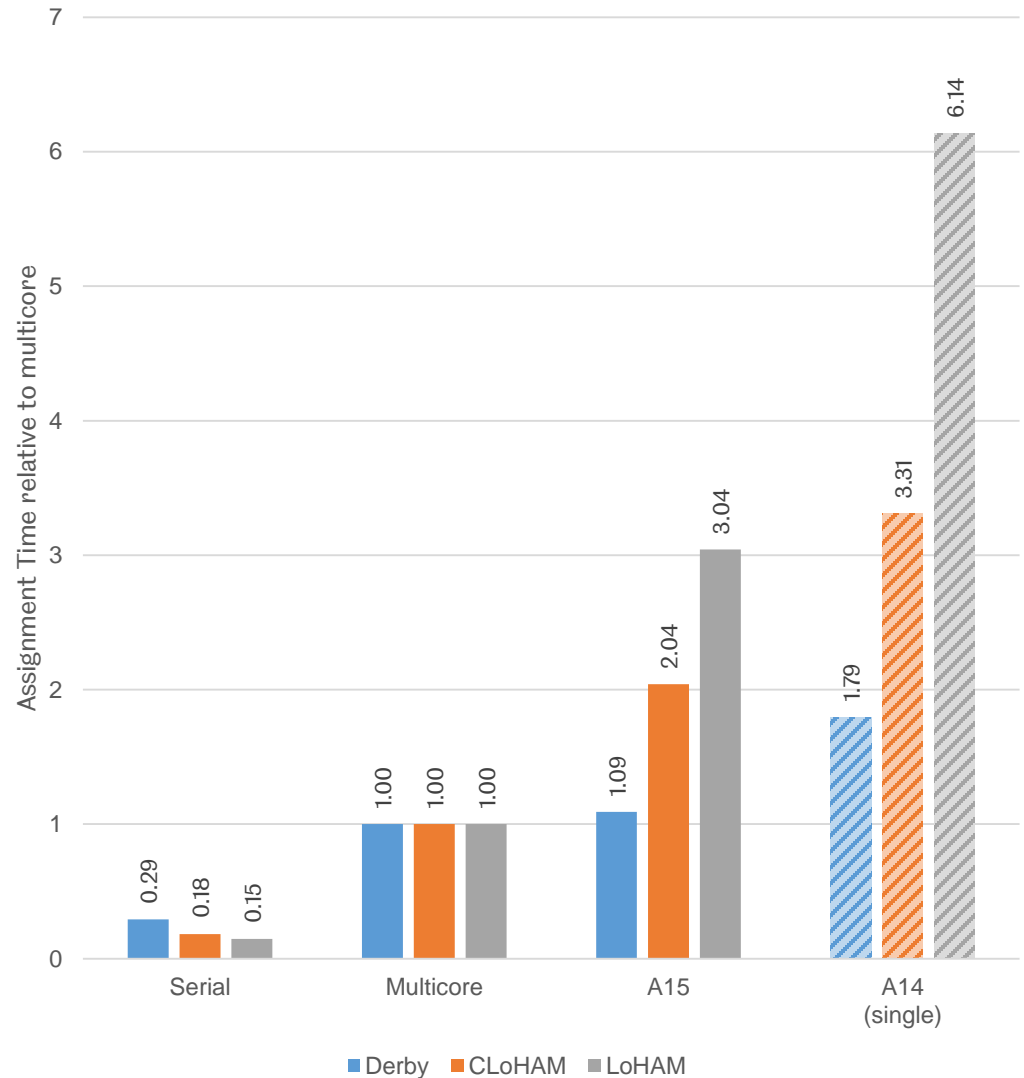
Double precision
- LoHAM – 35m 22s

Single precision
- Reduced loss of precision
- LoHAM – 17m 32s

Hardware:
- Intel Core i7-4770K
- 16GB DDR3
- Nvidia GeForce Titan X
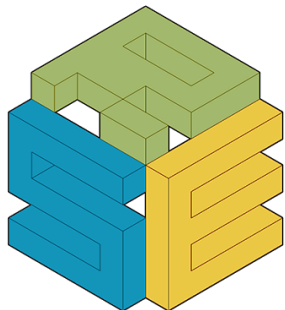
## Relative Assignment Runtime Performance vs Serial



| | Derby | CLoHAM | LoHAM |
|---|---|---|---|
| Serial | 1.00 | 1.00 | 1.00 |
| Multicore | 3.42 | 5.51 | 6.81 |
| A15 | 3.73 | 11.24 | 20.70 |
| A14 (single) | 6.14 | 18.24 | 41.77 |

Assignment Time relative to serial

# Assignment Speedup relative to Multicore

**Multicore**

- LoHAM – 1h 47m

**Double precision**

- LoHAM – 35m 22s

**Single precision**

- Reduced loss of precision
- LoHAM – 17m 32s

**Hardware:**

- Intel Core i7-4770K
- 16GB DDR3
- Nvidia GeForce Titan X

## Relative Assignment Runtime Performance vs Multicore

Assignment Time relative to multicore

| Category | Derby | CLoHAM | LoHAM |
|---|---|---|---|
| Serial | 0.29 | 0.18 | 0.15 |
| Multicore | 1.00 | 1.00 | 1.00 |
| A15 | 1.09 | 2.04 | 3.04 |
| A14 (single) | 1.79 | 3.31 | 6.14 |

Legend: ■ Derby ■ CLoHAM ■ LoHAM

# GPU Computing at UoS

# Expertise at Sheffield

- Specialists in GPU Computing and performance optimisation

- Complex Systems Simulations via FLAME and FLAME GPU

- Visual Simulation, Computer Graphics and Virtual Reality
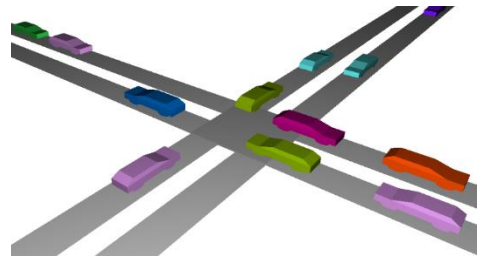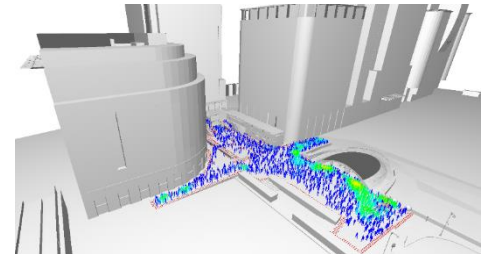
- Training and Education for GPU Computing

# Thank You

# Largest Model (LoHAM) results

**Paul Richmond**

p.richmond@sheffield.ac.uk

paulrichmond.shef.ac.uk

**Peter Heywood**

p.heywood@sheffield.ac.uk

ptheywood.uk

| | Runtime | Speedup Serial | Speedup Multicore |
|---|---|---|---|
| **Serial** | 12:12:24 | 1.00 | 0.15 |
| **Multicore** | 01:47:36 | 6.81 | 1.00 |
| **A15 (double precision)** | 00:35:22 | 20.70 | 3.04 |
| **A14 (single precision)** | 00:17:32 | **41.77** | **6.14** |

# Backup Slides

# Benchmark Models

- 3 Benchmark networks
  - Range of sizes
  - Small to V. Large
  - Up to 12 hour runtime

| Model | Vertices (Nodes) | Edges (Links) | O-D trips |
|---|---|---|---|
| Derby | 2700 | 25385 | $547^2$ |
| CLoHAM | 15179 | 132600 | $2548^2$ |
| LoHAM | 18427 | 192711 | $5194^2$ |



Benchmark model performance

# Edges considered per algorithm
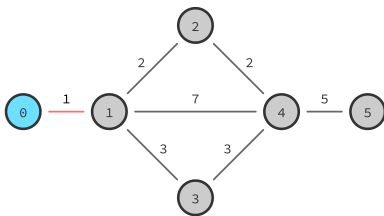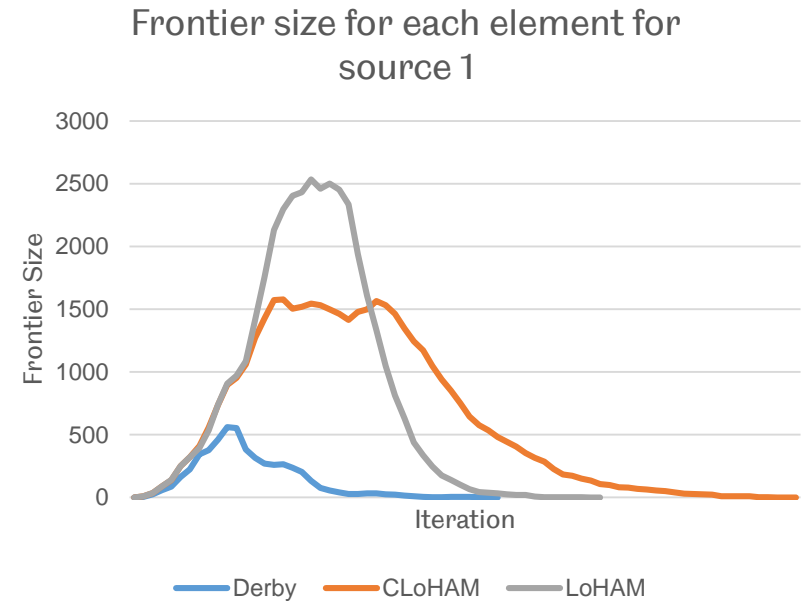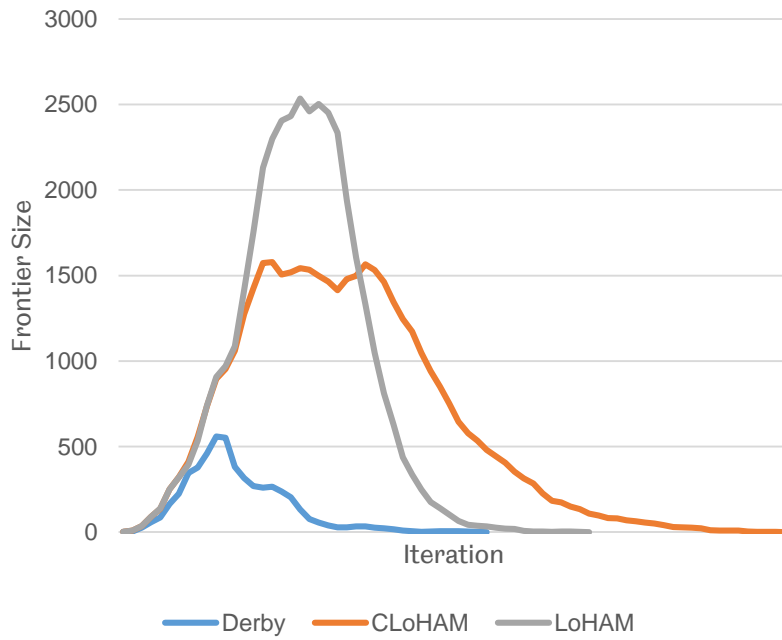
# Vertex Frontier (A4)

- Only Vertices which were updated in the previous iteration can result in an update

- Much fewer threads launched per iteration
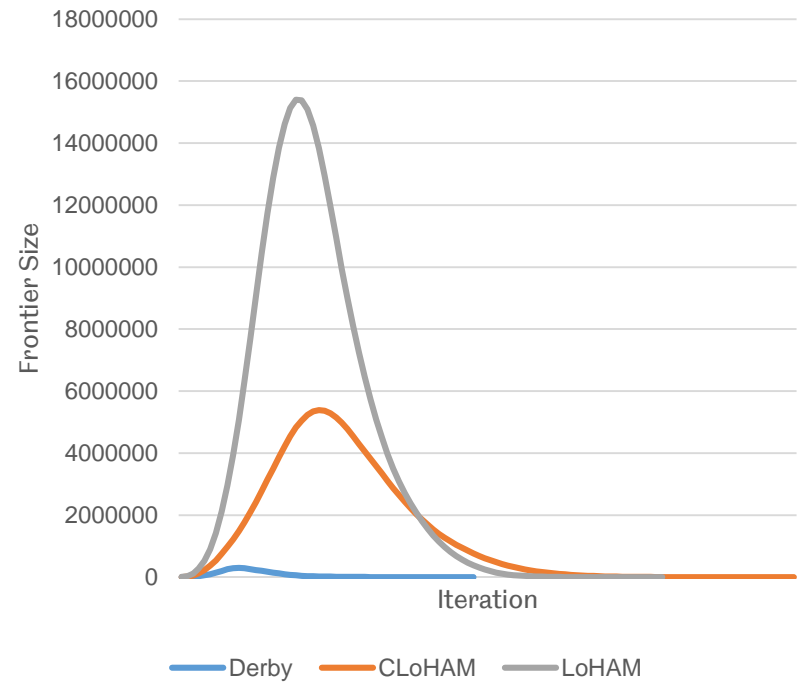
  - Up to 2500 instead of 18427 per iteration

Frontier size for each element for source 1

# Multiple Concurrent Origins (A8)

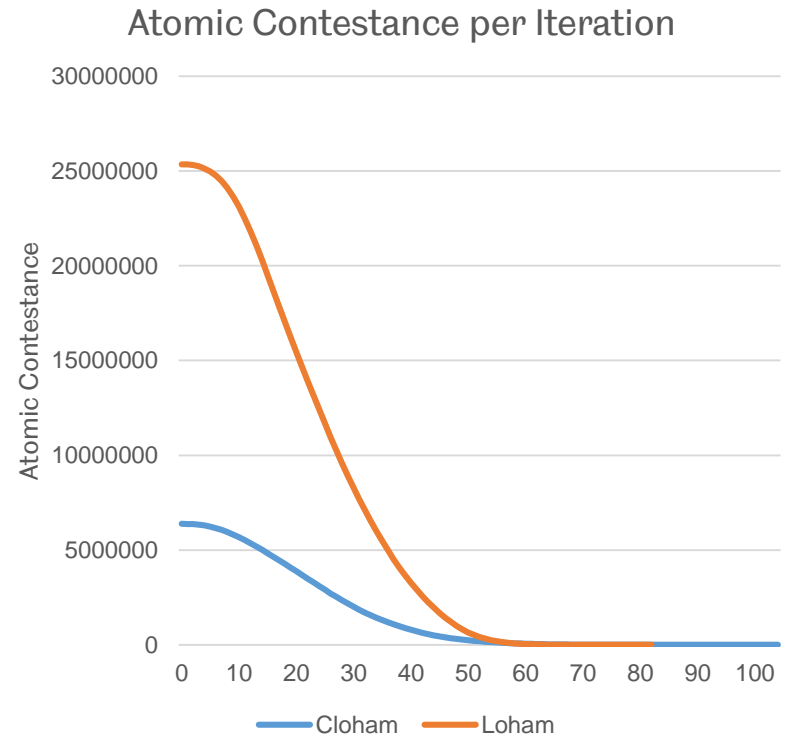Frontier size for each element for source 1



Frontier Size for all concurrent sources

# Atomic Contention

- Atomic operations are guaranteed to occur

- Atomic Contention
  - multiple threads atomically modify same address
  - Serialised!

- `atomicAdd(double)` not implemented in hardware
  - Not yet

- Solutions
  1. Algorithmic change to minimise atomic contention
  2. Single precision

### Atomic Contestance per Iteration

# Raw Performance

Hardware:

- Intel Core i7-4770K
- 16GB DDR3
- Nvidia GeForce Titan X

## Assignment runtime per algorithm

Time (seconds)

Serial | Multicore | A8 | A10 | A11 | A15 | A13 (single) | A14 (single)

■ Derby  ■ CLoHAM  ■ LoHAM