# Fast triangular solve on GPUs

Jonathan Hogg

STFC Rutherford Appleton Laboratory

# ASEArch flagship grant

Aims:

- ▶ Deliver a sparse linear solver on GPUs
- ▶ Deliver an interior point solver for linear/quadratic programs on GPUs
- ▶ Do so in such a way that they can be easily ported to other architectures

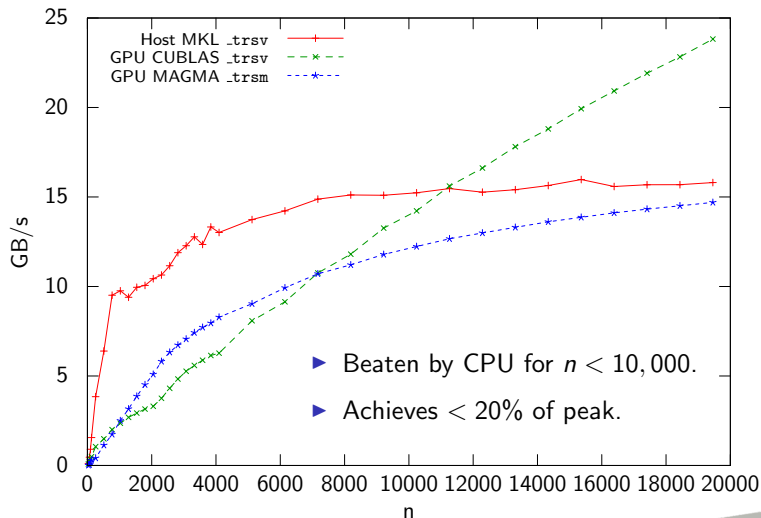Science & Technology
Facilities Council

# ASEArch flagship grant

Aims:

- ▶ Deliver a sparse linear solver on GPUs
- ▶ Deliver an interior point solver for linear/quadratic programs on GPUs
- ▶ Do so in such a way that they can be easily ported to other architectures

Relation of this talk:

- ▶ Learning project
- ▶ Base kernel we need to perform well — current CUBLAS implementation is poor.

Science & Technology
Facilities Council

# Current libraries



- Beaten by CPU for $n < 10,000$.
- Achieves $< 20\%$ of peak.

# Basic (in-place) Algorithm

**Input:** Lower-triangular $n \times n$ matrix $L$, right-hand-side vector $x$.
**for** $i = 1, n$ **do**
$$\boxed{x(i+1:n)} = \boxed{x(i+1:n)} - \boxed{L(i+1:n,i)} * \boxed{x(i)}$$
**end for**
**Output:** solution vector $x$.

$$\begin{pmatrix} 1 & & & \\ l_{21} & 1 & & \\ l_{31} & l_{32} & 1 & \\ l_{41} & l_{42} & l_{43} & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix}$$

Science & Technology
Facilities Council

# Small matrices are latency bound

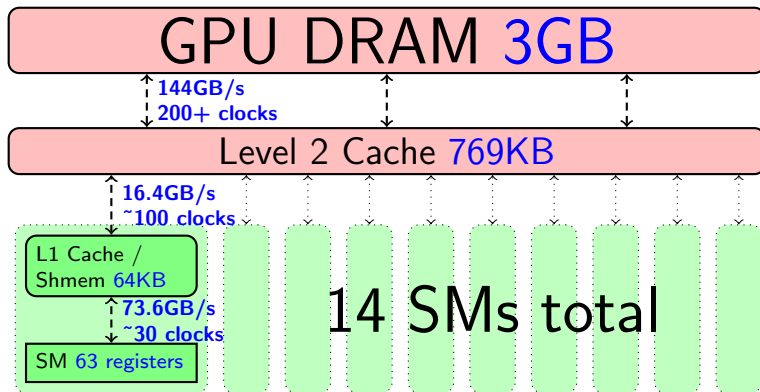**1 fmad per entry in $L$ $\Rightarrow$ memory-bound.**

- ▶ C2050 can deliver approx 9 doubles/sec from main memory
- ▶ Global memory latency 200 cycles (optimistic?)
- ▶ $n = 32 \Rightarrow 195$ cycles per column waiting for data
- ▶ Require $n > 1800$ to fully hide latency
- ▶ Cache doesn't help — no hardware prefetch.

**What can we do?**

Science & Technology
Facilities Council

# Small matrices are latency bound

**1 fmad per entry in $L \Rightarrow$ memory-bound.**

- C2050 can deliver approx 9 doubles/sec from main memory
- Global memory latency 200 cycles (optimistic?)
- $n = 32 \Rightarrow 195$ cycles per column waiting for data
- Require $n > 1800$ to fully hide latency
- Cache doesn't help — no hardware prefetch.

**What can we do?**

Bring data closer to core, reducing latency

- Shared memory; or
- Registers

# C2050 Memory layout

# Shared memory $n > 32$

Quickly run out of shared memory if we try and hold entire matrix!
Instead:

- Cache only $32 \times 32$ tiles down diagonal
- Cache next col while solve performed on diagonal

$$\begin{pmatrix} L_{11} & & & \\ L_{21} & L_{22} & & \\ L_{31} & L_{32} & L_{33} & \\ L_{41} & L_{42} & L_{43} & L_{44} \end{pmatrix}$$

Science & Technology
Facilities Council

# Shared memory $n > 32$

Quickly run out of shared memory if we try and hold entire matrix!

Instead:

- ▶ Cache only $32 \times 32$ tiles down diagonal
- ▶ Cache next col while solve performed on diagonal

$$\begin{pmatrix} L_{11} & & & \\ L_{21} & L_{22} & & \\ L_{31} & L_{32} & L_{33} & \\ L_{41} & L_{42} & L_{43} & L_{44} \end{pmatrix}$$

Execution trace ($128 \times 128$):

# Small matrix results

| $n =$ | 32 | 64 | 96 | 128 |
|---|---|---|---|---|
| Shared-memory | 7 | 13 | 19 | 25 |
| Registers | 17 | 37 | 68 | 149* |
| CUBLAS `dtrsv()` | 31 | 58 | 85 | 113 |

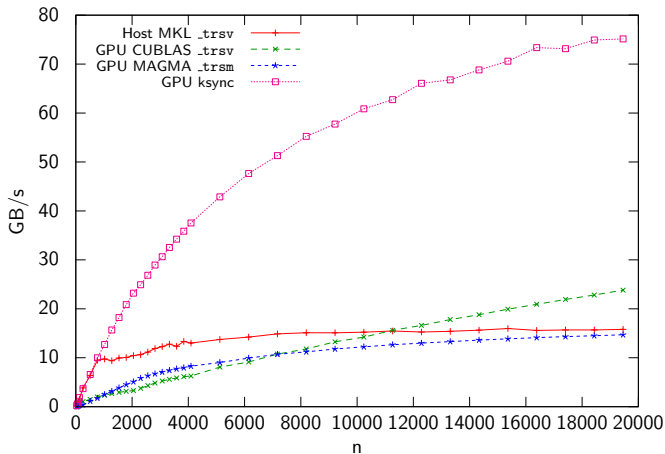* indicates register spill occurred

Science & Technology
Facilities Council

# Larger matrices

**So far using a single SM.**

- ▶ Quickly L1⟷L2 bandwidth becomes bounding (only 16.4GB/s vs 144GB/s global)
- ▶ Need to use multiple SMs!

**Why not use small matrix kernel then efficient matrix-vector?**

- ▶ Driver handles synchronization (different kernels)
- ▶ Matrix-vector achieves high bandwidth

# Kernel-synchronized results

# We can do better!

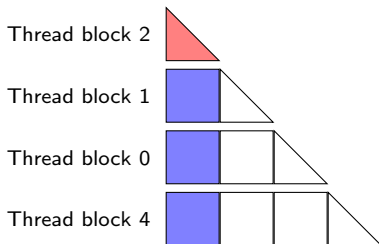| $n =$ | 512 | 1024 | 4096 |
|---|---|---|---|
| `blkSolve()` ($\mu$s) | 108.3 | 217.3 | 904.7 |
| `dgemv()` ($\mu$s) | 37.8 | 95.1 | 842.0 |
| Execution time ($\mu$s) | 171.0 | 370.8 | 2006.5 |
| Launch overhead | **17.0%** | **18.7%** | **14.9%** |
| Work in `blkSolve()` | **18%** | **9%** | **2%** |

▶ Substantial overheads from using kernel launches for synchronization

▶ Amount of time in blkSolve() — Amdahl strikes again!

Science & Technology
Facilities Council

# Global-memory synchronized

### Aim: Single kernel-launch

- ▶ Use global memory for synchronization — costs l2 cache miss + `__threadfence()`.
  (Much cheaper than using kernel launches)
- ▶ Fine grained synchronization...
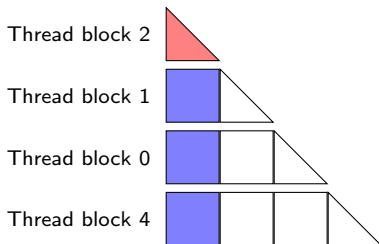- ▶ ...hence matrix-vector product runs concurrently with solve.

Science & Technology
Facilities Council

# Thread block $\Rightarrow$ block row



**CAUTION**
Thread blocks are not
scheduled in order!

# Thread block ⇒ block row



**CAUTION**
Thread blocks are not
scheduled in order!

Dynamically pick row to
avoid deadlock

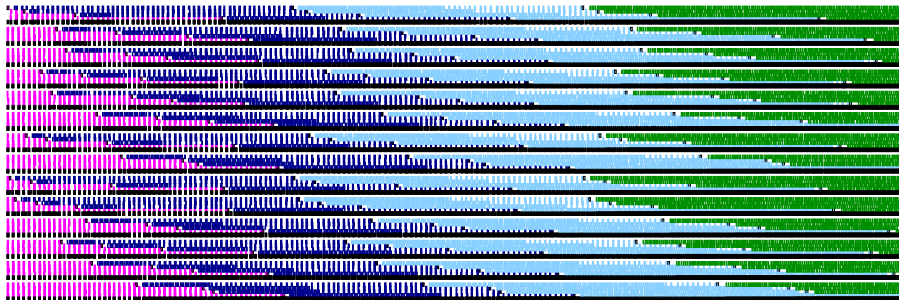# Thread block ⇒ block row



**CAUTION**
Thread blocks are not
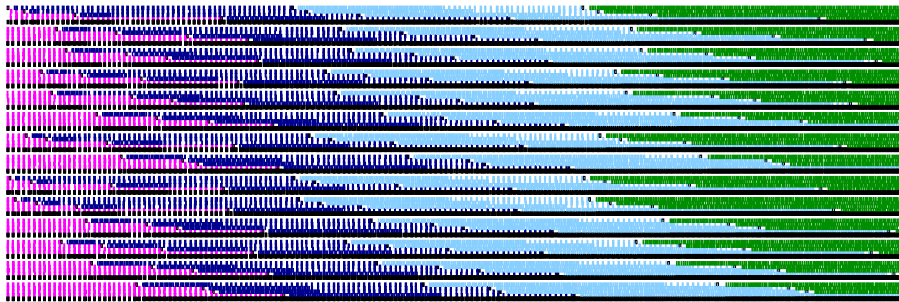scheduled in order!

Dynamically pick row to
avoid deadlock

Only need two scalars for synchronization:

- ▶ Row for next thread block
- ▶ Latest column for which solution is available
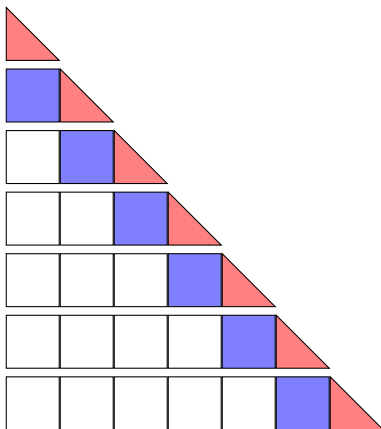
# Execution trace

# Execution trace



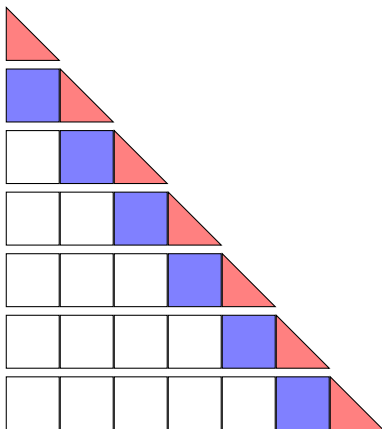Mode 1 Not waiting on data, constant computation.
Mode 2 Stops and starts as each column completes.

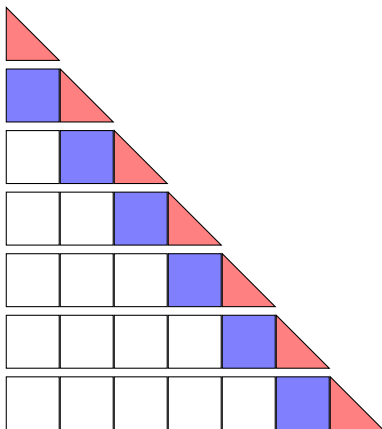# Critical path



Critical path is coloured;
Executes serially

# Critical path



Critical path is coloured;
Executes serially

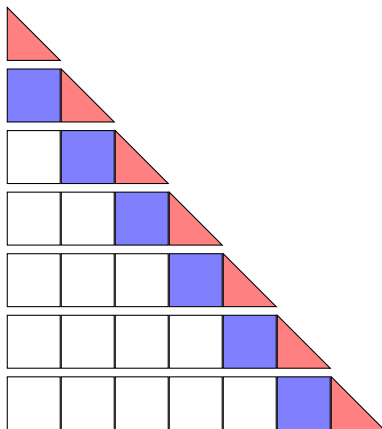Use tricks from before:
**pre-cache values**

# Critical path



Critical path is coloured;
Executes serially

Use tricks from before:
**pre-cache values**
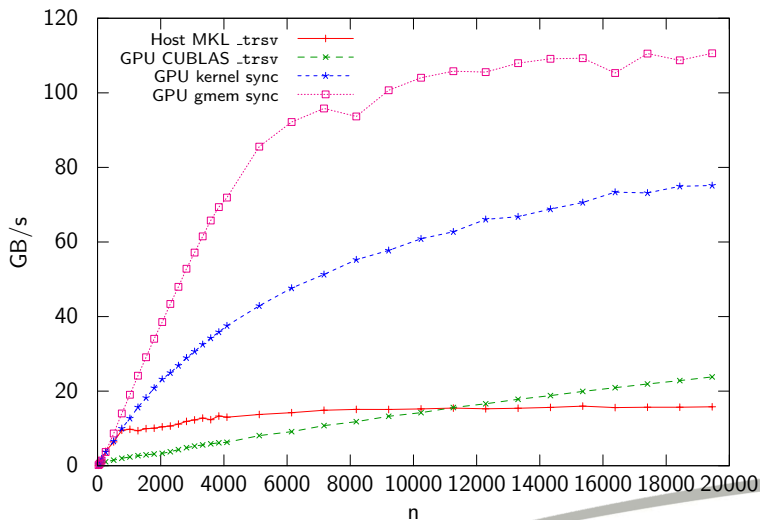
**BUT:**
Maintain high occupancy!

# Critical path



48k shmem $\Rightarrow$ At most 5 $32 \times 32$ tiles
Want 4 thread blocks/SM!

▶ Use shared memory for  diagonal  tiles.

▶ Use registers for  subdiagonal  tiles.

Science & Technology
Facilities Council

# Global-memory synchronization results

Science & Technology
Facilities Council

# Better yet!

## Memory-bound $\Rightarrow$ spare flops

Can we do redundant computation to speed the critical path?

# Better yet!

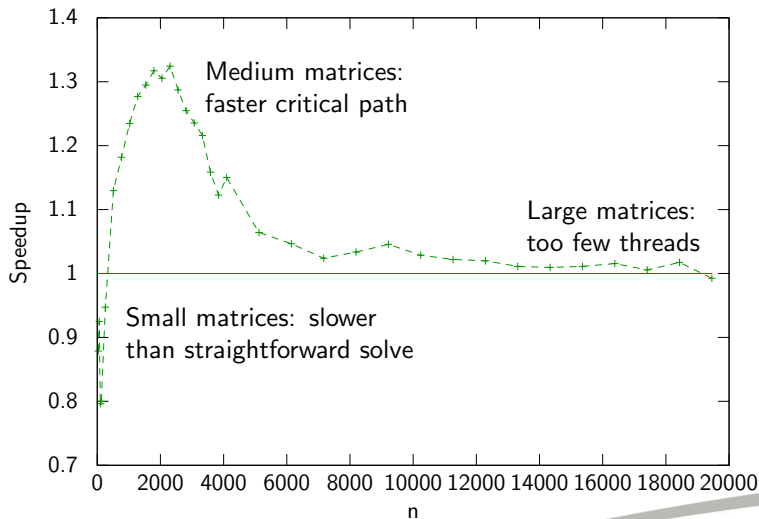## Memory-bound $\Rightarrow$ spare flops

Can we do redundant computation to speed the critical path?

### YES

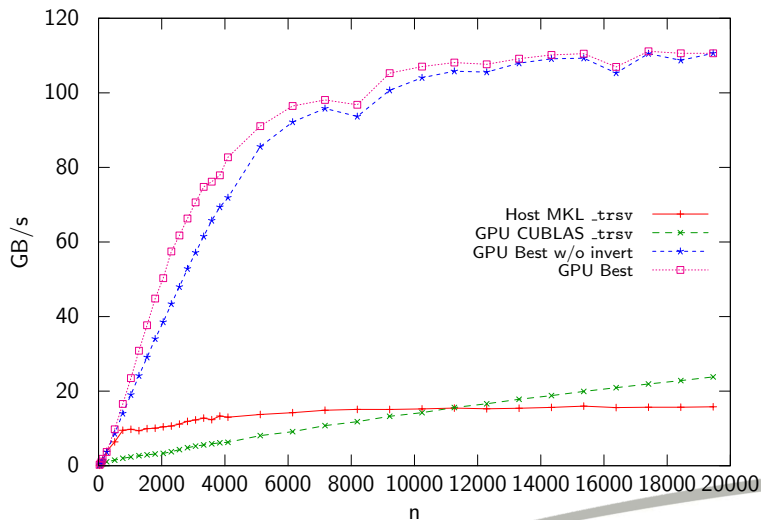Explicit inversion of diagonal blocks (stable see Higham 1995)

- Diagonal solve $\rightarrow$ Matrix-vector multiply
- Same number of memory accesses, *less communication*!
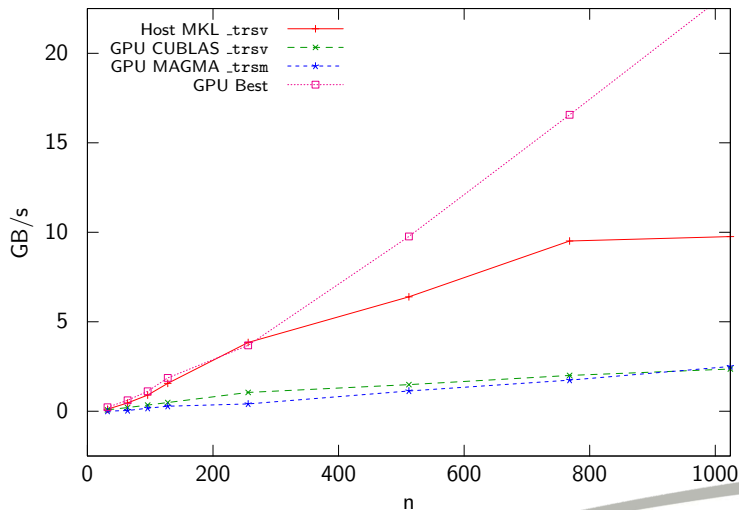
Science & Technology
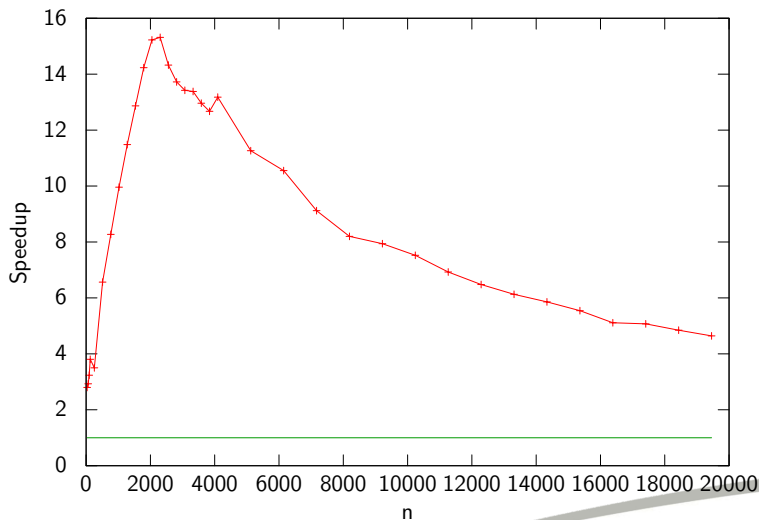Facilities Council

# Speedup over previous version

# Overall best performance

# Overall best performance (zoomed)

# Speedup vs CUBLAS

# Conclusions

We've beaten CUBLAS soundly.
Achieved 75% of peak bandwidth.

- ▶ Can we do even better somehow?
- ▶ Could use tasks — but register pressure!

Next step is the sparse case

Code will be available under BSD licence

# Questions?

# Explicit inversion

$$
\begin{pmatrix} L_{11} & \\ L_{21} & L_{22} \end{pmatrix}
\begin{pmatrix} X_{11} & \\ X_{21} & X_{22} \end{pmatrix}
=
\begin{pmatrix} L_{11}X_{11} & \\ L_{21}X_{11} + L_{22}X_{21} & L_{22}X_{22} \end{pmatrix}
$$

Equate to identity.

$$
\begin{aligned}
X_{11} &= L_{11}^{-1} && \text{by recursion} \\
X_{22} &= L_{22}^{-1} && \text{by recursion} \\
L_{22}X_{21} &= -L_{21}X_{11} && \text{solve is stable - Higham 1995}
\end{aligned}
$$

Doesn't require right-hand-side — can be done before needed

BUT: takes considerably longer than a solve: useless for small $n$.

Science & Technology
Facilities Council

# Small matrix — Registers

- Block on use, not on load.
- Allow Instruction Level Parallelism (ILP).
- See Volkov's *Better Performance at Lower Occupancy*.

Each thread only has 63 registers!
… typically need half of these for normal operation.

# Small matrix — Registers

- Block on use, not on load.
- Allow Instruction Level Parallelism (ILP).
- See Volkov's *Better Performance at Lower Occupancy*.

Each thread only has 63 registers!
… typically need half of these for normal operation.

However, doesn't help:

- To use more than 1 thread, need to communicate via shared memory
  (so no latency gain).
- Adds complications to code $\Rightarrow$ extra overheads.
- Quite quickly leads to register spill $\Rightarrow$ slowdown.

Science & Technology
Facilities Council

# Small matrix — Shared Memory

A $32 \times 32$ matrix of doubles requires 8KiB $\Rightarrow$ lots of room.
Simple code ($\text{blkSize} = 32$):

```
template <int blkSize>
void __device__ dblkSolve(const double *a, int lda,
        double &val) {

    volatile double __shared__ xs;

#pragma unroll 16
    for(int i=0; i<blkSize; i++) {
        if(threadIdx.x==i) xs = val;
        if(threadIdx.x>=i+1)
            val -= a[i*lda+threadIdx.x] * xs;
    }
}
```

Just precache a in shared memory!

Science & Technology
Facilities Council