

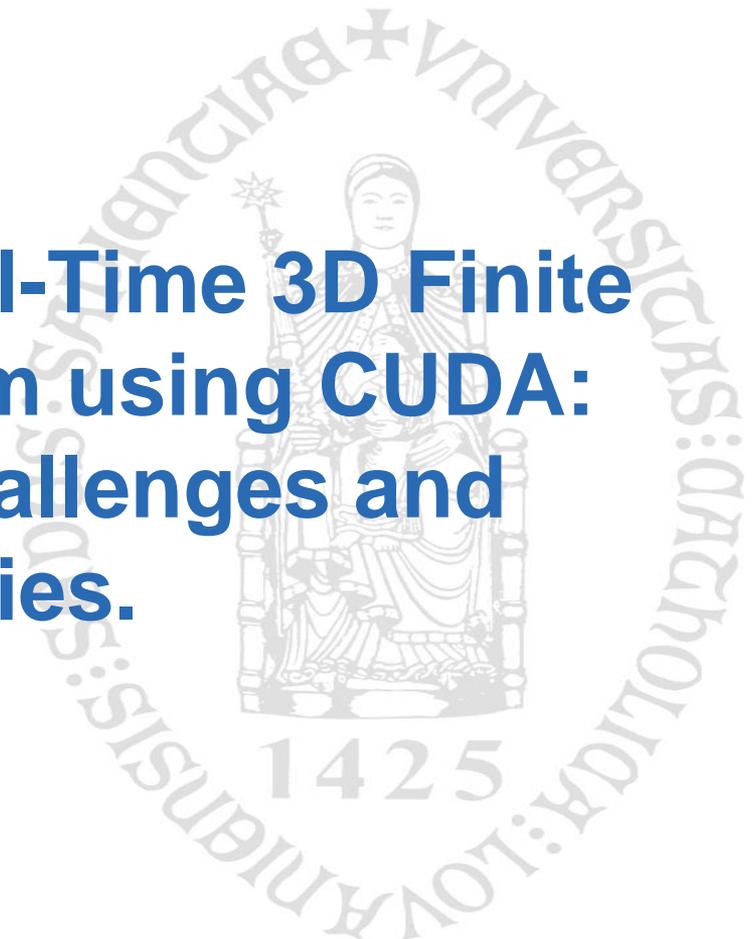


# Parallelizing a Real-Time 3D Finite Element Algorithm using CUDA: Limitations, Challenges and Opportunities.

Vukasin Strbac

Biomechanics section

KU Leuven



# The Finite Element Method and the GPU

---

## Implicit

Basically two phases:

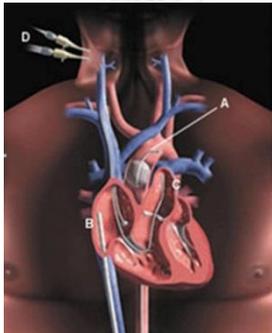
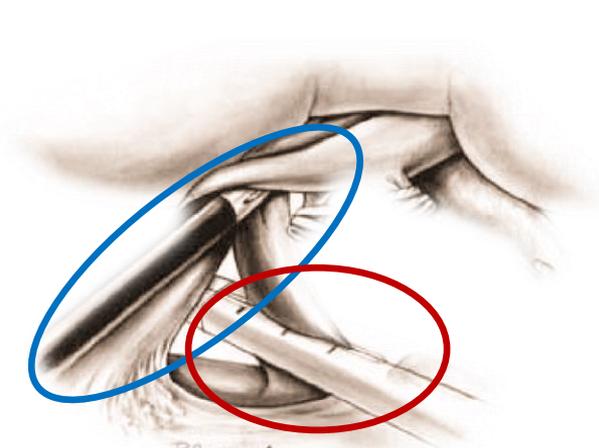
1. Assembly: computation of elemental stiffness matrices and scatter into large global stiffness matrix (a bit more interesting cfr. Luitjens GTC2012)
2. Solve the system to obtain displacements (well treated already)

## Explicit

Loop:

1. Compute forces directly from a material model on a per-element basis (similar to element stiffness matrix computation)
2. Update displacements explicitly (e.g. central differences)
3. Impose boundary conditions, incrementally

# Real-time application of Finite Elements



Fabrizio1999

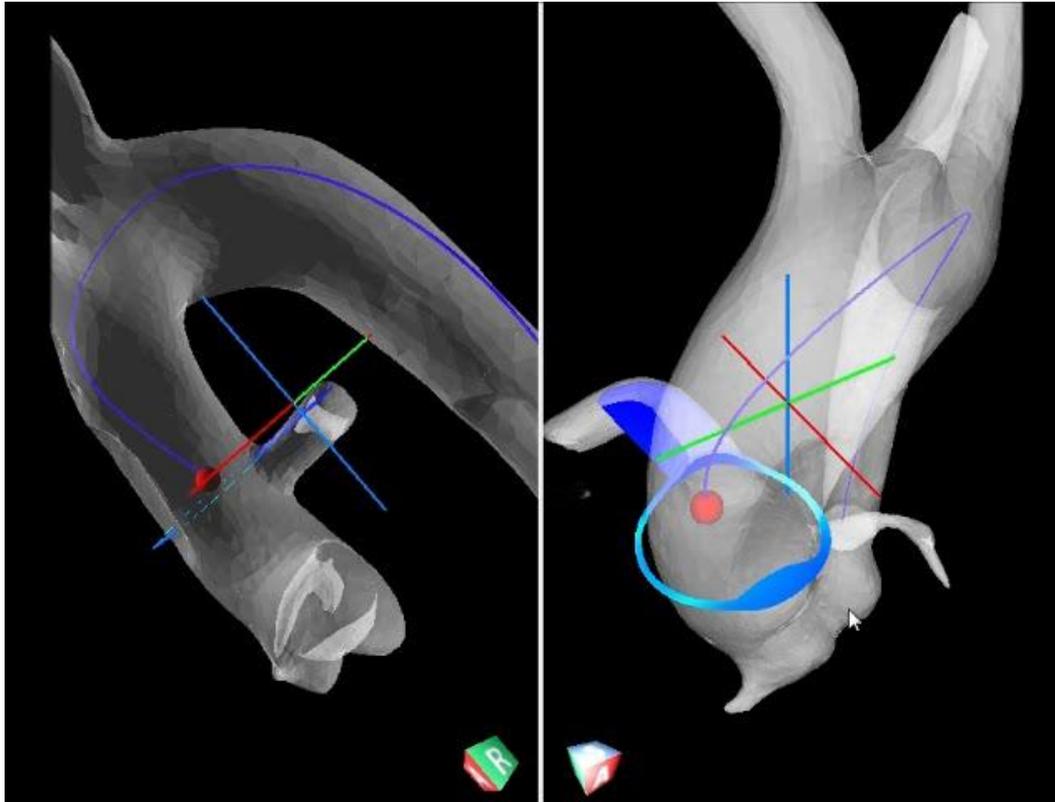
- Pulling aside the renal vein
  - **Clamping the renal artery**
  - Inflation of an endoclamp balloon
- 
- Active constraints:
    - Design a new, safer clamp?
    - Monitor inflation pressure in the balloon?
    - ...

## Real-time FE in the surgical theater

---

- 'Limited' number of DOFs is fine
  - Driven by pathological conditions of tissue which can emerge even on coarser meshes.
- Nonlinear, both in geometry and materials
- Quick response
  - Explicit solution gives more opportunity to check values
- Why not use a CPU cluster?
  - Benefit from a tightly-coupled system (especially explicit computation)

## Real-time TLED, aortic expansion example



- Simulation at 17 fps  
~70ms per solution, 1000 timesteps each
- Loading is 0-500mmHg sine wave
- Mimicking endoclamp balloon expansion
- Neo-Hookean material with  $E=3000\text{Pa}$ ,  $\nu=0,49$

# Total Lagrangian Explicit Dynamic (TLED)

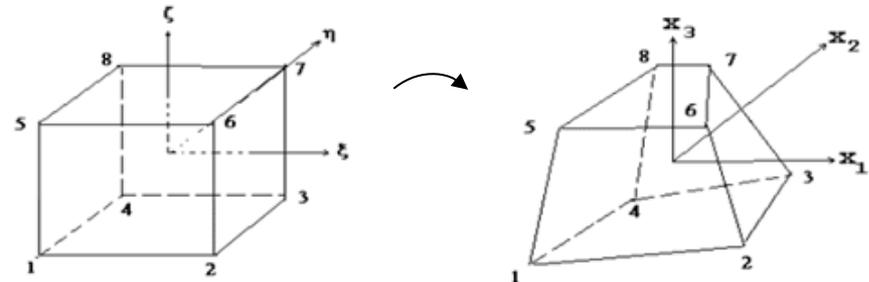
```

Compute time-step;
Compute shape function derivatives;
for each time-step{
    computeForces () ;
    updateDisplacements () ;
    enforceBoundaryConditions () ;
}
    
```

CPU

GPU

- 3D low order elements:
  - 8-node hexahedron
  - Trilinear shape functions
- Low order integration: single integration point



Avezedo 1999

## TLED forces computation

---

```

computeForces(u,  $\nabla N$ , E1, V, matP){
    load(E1); //element info
    load(u, (E1)); //displacements of that element
    load( $\nabla N$ ); //shape function derivatives
    X=f(u,  $\nabla N$ ); //deformation gradient
    C=f(X); //cauchy-green deformation tensor
    Cinv = f(C); //its inverse
    S=f(Cinv, matP); //2PK stress
    load(V); //volumes
    F=f(S, V, X,  $\nabla N$ ); //final forces
    store(F, (E1)); //store
}
  
```

# TLED memory considerations

```

computeForces(u, ∇N, E1, V, matP) {
    load(E1);                coalesced read +8.
    load(u, (E1));          uncoalesced read+24
    load(∇N);               coalesced read +24
    X=f(u, ∇N);             +9
    C=f(X);                 +9
    Cinv = f(C);            +9
    S=f(Cinv, matP);       +6
    load(V);                coalesced +1
    F=f(S, V, X, ∇N);       +24
    store(F, (E1));        Uncoalesced store
}
  
```

+intermediate  
 = ~128+ registers  
 for the simplest case

# Memory analysis on C2075

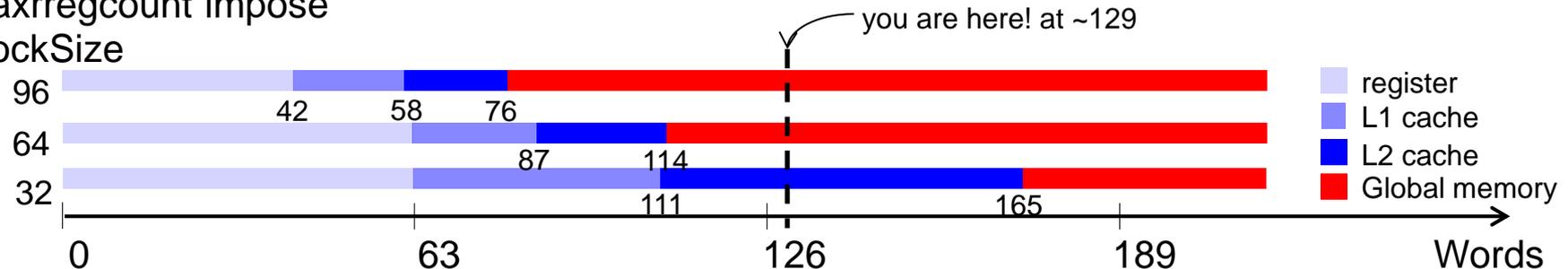
\*Single Precision

Memtype	32	64	96
Register words	63	63	42
L1(48kB per SM)	192B	96B	64B
L2 (768kB total)	219B	109B	19B
L1 words	48	87	58
L2 words	54	114	76

- Registers->L1->L2->global
- Registers spilling everywhere!
- Currently blockSize = 32
- blockSize of 64 would spill all the way to global memory.

Keeping local memory away from gmem very beneficial

--maxrregcount impose  
blockSize



## Conclusion for c2075:

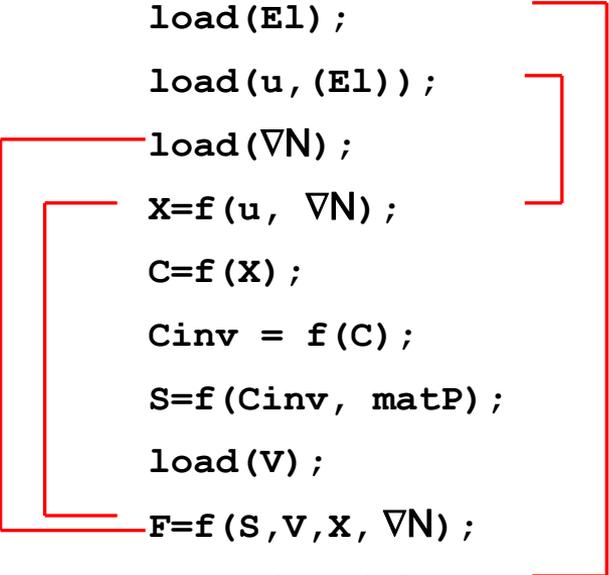
---

- Currently lots of execution dependency stalls due to “few” warps in flight (IPC:  $\sim 0.73$ , C2075)
- Increasing number of resident warps  $\rightarrow$  better latency hiding
  - But have to get data from global memory!
- **Does not pay off.**

## TLED variables lifetime

```

computeForces(u, ∇N, E1, V, matP) {
    load(E1);
    load(u, (E1));
    load(∇N);
    X=f(u, ∇N);
    C=f(X);
    Cinv = f(C);
    S=f(Cinv, matP);
    load(V);
    F=f(S, V, X, ∇N);
    store(F, (E1));
}
  
```



Lifetime of some variables not friendly!

# Pseudocode, computing forces

---

```

computeForces(u,  $\nabla N$ , E1, V, matP) {
    load(E1);
    load(u, (E1));
    load( $\nabla N$ );
    X=f(u,  $\nabla N$ );
    C=f(X);
    Cinv = f(C);
    S=f(Cinv, matP);
    -----
    load(V);
    F=f(S, V, X,  $\nabla N$ );
    store(F, (E1));
}
  
```

Lifetime of variables not friendly!

Can we split the kernel?

Store here

Load here

# Pseudocode, computing forces, splitting

```

computeFstore(u, ∇N, E1, v, matP) {
    computeFread(S, X, ∇N, v, E1) {
        load(E1);
        load(u, (E1));
        load(∇N);
        X=f(u, ∇N);
        store(X);
        C=f(X);
        Cinv = f(C);
        S=f(Cinv, matP);
        store(S);
    }
}
    
```

Single, large kernel still faster...

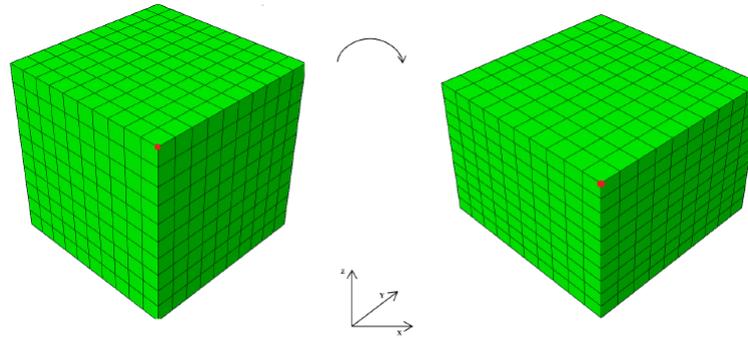
- Register req. go down
- Occupancy goes up
- **But added +47/15 fetches/stores**
- Doesn't pay off
  - + maxrregcount = 41
- Register req. forced down
- **Register spills**
- **Occupancy goes further up but doesn't hide latency**
- Doesn't pay off

## Challenges and limitations

---

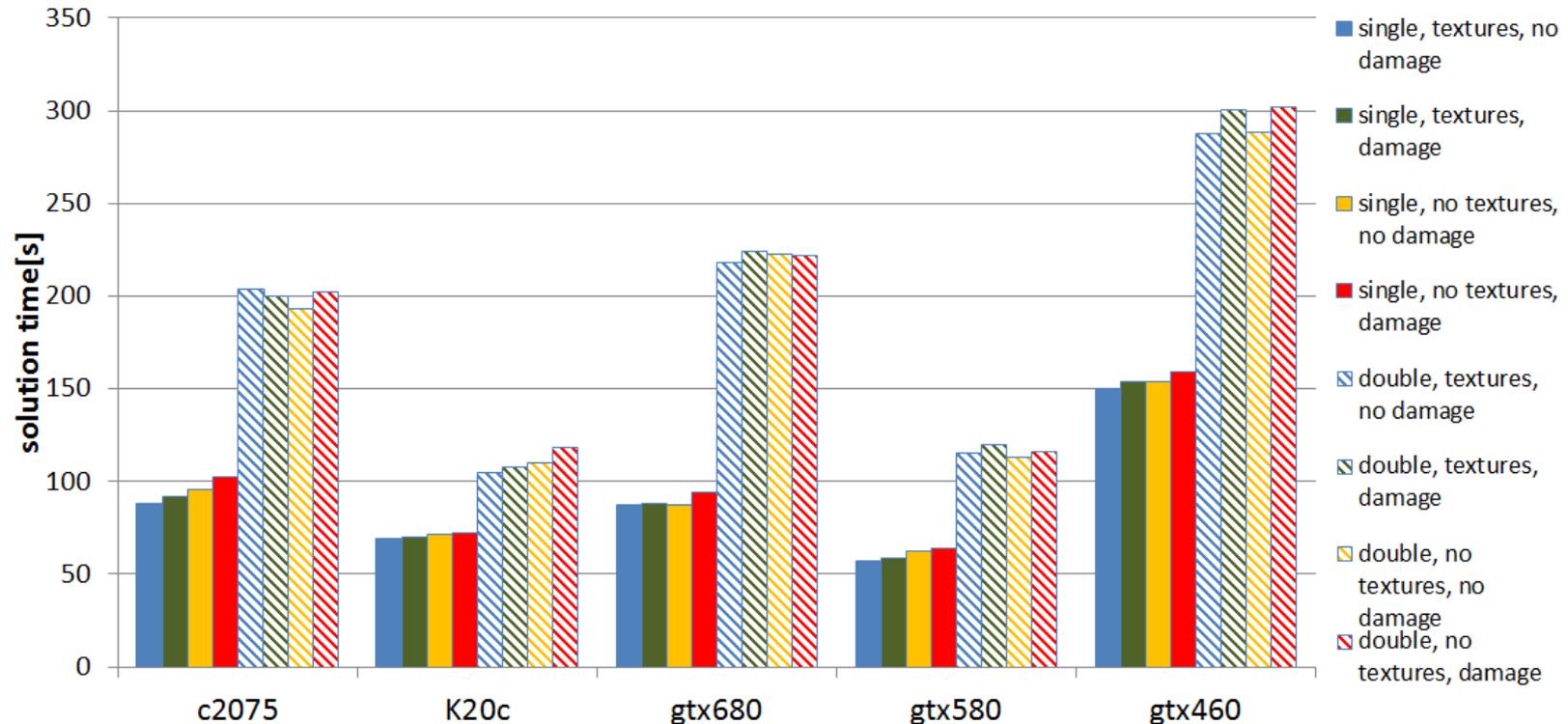
- In explicit FE granularity is intuitive and parallelism trivial but limited
  - Element kernels
  - Nodal kernels
- Therefore lots of memory per element
  - Especially for Lagrangian methods
  - Even more so on higher order elements and elasticity problems.
- High register req. -> occupancy issues
  - Low occupancy -> difficult latency hiding
- High register req.-> register spilling
  - Additional latency, can be hidden only by ILP
- Splitting the kernel unlikely to improve performance

# Testing and Validation



- Simple compression problem on truth cube
- Boundary conditions:
  - Top node set goes down Z, rest is free
  - Bottom node constrained in Z, rest is free
- Models range from 5x5x5 to 45x45x45
- Comparing against Abaqus
  - Accuracy and speed
- Model: Neo-Hookean solid
- GPUs used in testing:
  - K20c
  - C2075
  - GTX680
  - GTX580
  - GTX460
- **Testing modes**
  - Single/Double precision
  - Using/not using texture memory

# Testing and validation: 90k elements, 50k timesteps

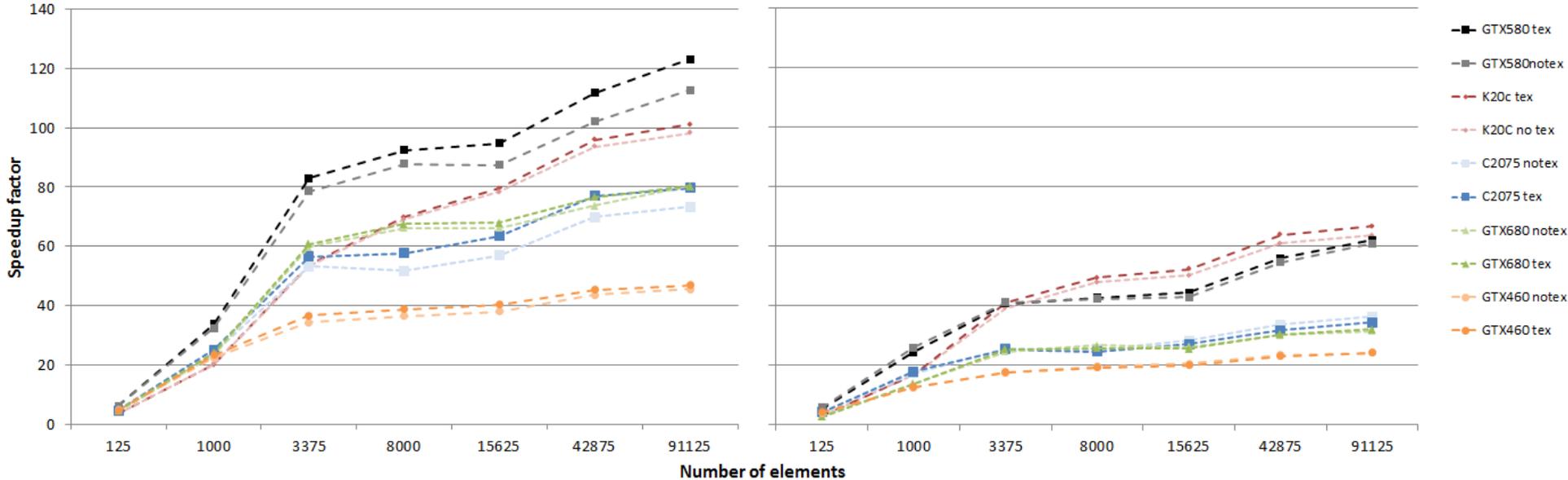


Single card solution

# Testing and validation: per time-step view

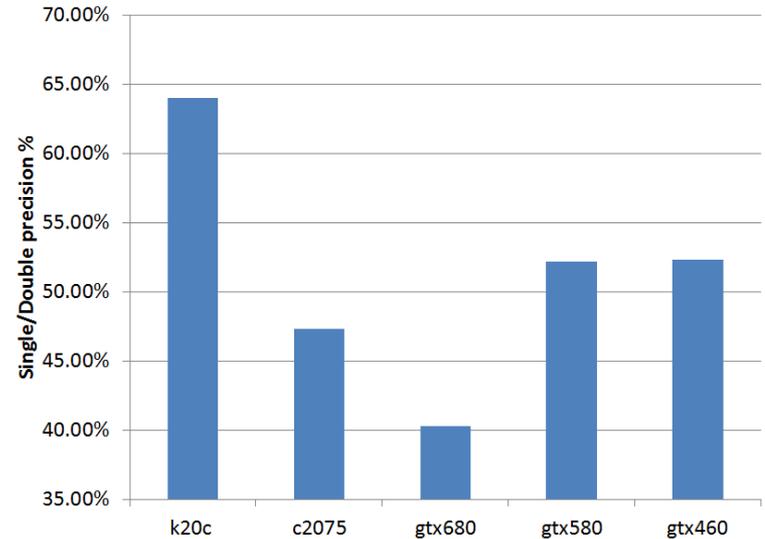
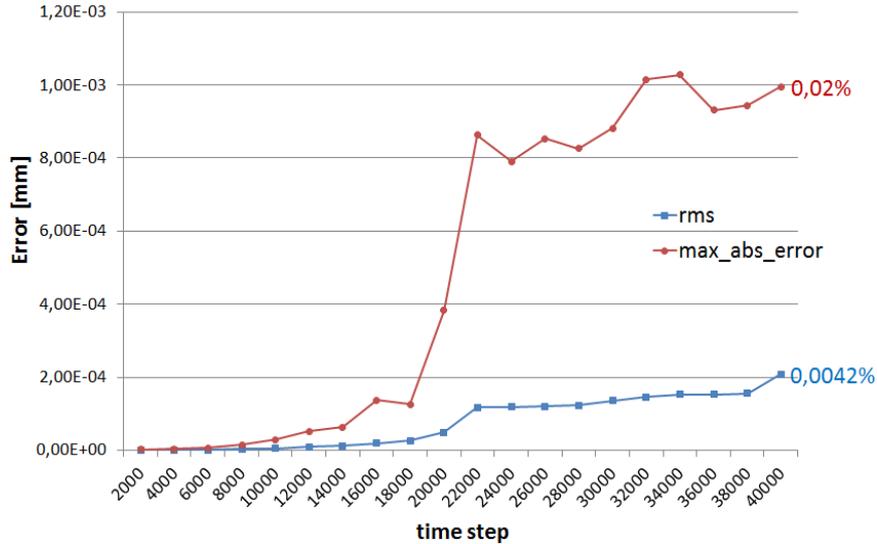
General speedup vs Abaqus, single

General speedup vs Abaqus, double



But can we use the single precision solution?

# Single precision viability



- Round-off error increases with
  - number of time-steps
  - element density

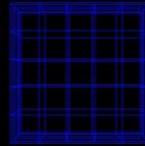
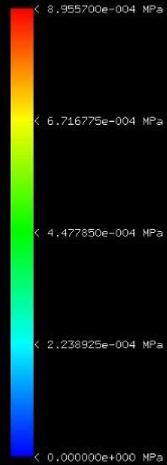
## Final comments

---

- Generally:
  - Decomposition at the element level is necessary
  - Problem driven by memory, bound by latency
- Practically:
  - Prefer more, lower order elements
  - Lower occupancy may provide better performance
  - Splitting a large force computation kernel unlikely to increase performance
  - Use atomics on Kepler, no need for mini-reduction of force contributions

Use left/right mouse button to rotate/zoom  
Use numbers 1-7 to switch/reset models  
Use 'c' to switch stresses  
Use 'l' to switch autonomous compression/tension  
Use 'k' to switch to free mode  
Use '+' / '-' to compress/extend in free mode  
Use 'p' to lock/unlock the top and bottom nodes  
Use 'f' to switch between wireframe/fill rendering  
Use 'n' for colormap  
Use 'b' for colormap values  
Use 's' for status  
Use 'h' for help. Turn off TEXT for bet performance

Model: 125  
Stress: vonMises  
Load: compression  
Load: Line



---

Thank you.

Questions?

## computeForces()

- Create 2 additional versions of the kernel, memory only and math only.
- Tricking the compiler can be tricky: mem-only version must contain “some” compute, and math-only must waste cycles on generating data instead of fetching.
- Poor overlap between mem and math.

CONCLUSION:

Memory and latency bound.

### computeForces()

