



# The challenge of the solve phase of a multicore solver

**Jonathan Hogg**

Jennifer Scott

Rutherford Appleton Laboratory

PMAA 2010, Basel, 2 July 2010

# Sparse Direct Solvers

Solve  $A\mathbf{x} = \mathbf{b}$  using a Cholesky factorization:

$$PAP^T = LL^T$$

Where  $A$  is...

- ▶ Large
- ▶ Sparse
- ▶ Positive Definite

Variants for more general matrices —  $LDL^T$ ,  $LU$ .

# Sparse Direct Solvers

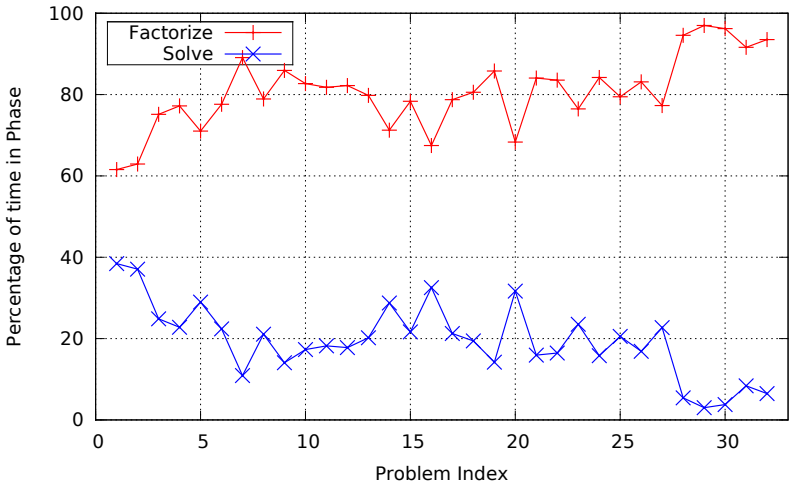
Four phases:

- Ordering** Reduce fill-in (AMD, MeTiS, SCOTCH, ...)
- Analyse** Construct data structures, plan computation
- Factorize** Numerical computation
- Solve** Triangular solves with  $L$  and  $L^T$

**Traditionally:** Factorize phase by far the most time consuming:

1 Factorize = 50-100 Solves

# Multicore. Oh dear.



Factorize parallel. Solve serial.

## So? Just parallelize the solve!

Lots of ways of doing this:

- ▶ Threaded BLAS/OpenMP do loops

## So? Just parallelize the solve!

Lots of ways of doing this:

- ▶ Threaded BLAS/OpenMP do loops
- ▶ Assembly tree parallelism

## So? Just parallelize the solve!

Lots of ways of doing this:

- ▶ Threaded BLAS/OpenMP do loops
- ▶ Assembly tree parallelism
- ▶ DAG-based methods

## So? Just parallelize the solve!

Lots of ways of doing this:

- ▶ Threaded BLAS/OpenMP do loops
- ▶ Assembly tree parallelism
- ▶ DAG-based methods



## So? Just parallelize the solve!

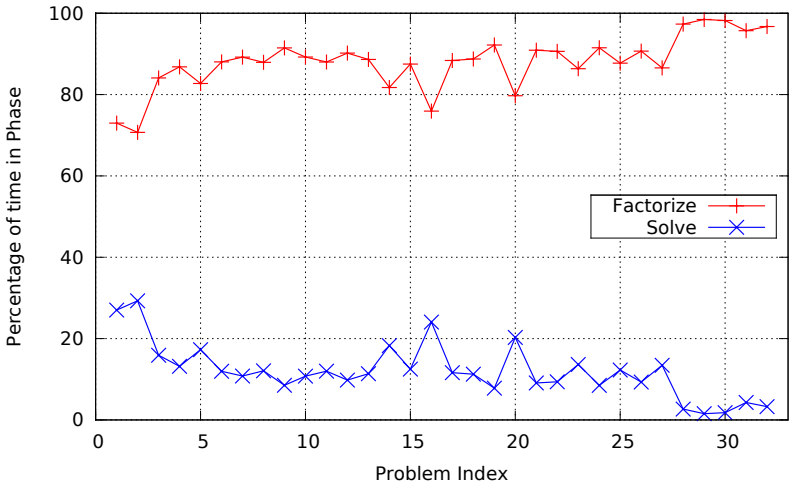
Lots of ways of doing this:

- ▶ Threaded BLAS/OpenMP do loops
- ▶ Assembly tree parallelism
- ▶ DAG-based methods

Regardless of how we did this...

...only a 2–3 times speedup

# Not much better...



Factorize and solve parallel.

# Why?

Typically:

	Flops	Cache Misses
Factor	> 98%	60–80%
Solve	< 2%	20–40%

Speedup on a single quad-core: < 1.30

Speedup on two quad-cores: < 3.00

**Solve is memory bound.**

Multicore = more cores, same memory bandwidth.

# Is this a problem?

## Is this a problem?

Yes.

Lots of use cases for multiple sequential solves following a factorize.

## Is this a problem?

Yes.

Lots of use cases for multiple sequential solves following a factorize.

**Iterative Refinement** classical case. Use multiple solves to reduce the error.

## Is this a problem?

Yes.

Lots of use cases for multiple sequential solves following a factorize.

**Iterative Refinement** classical case. Use multiple solves to reduce the error.

**Mixed Precision** precondition high precision iterative methods with low precision factorization.

## Is this a problem?

Yes.

Lots of use cases for multiple sequential solves following a factorize.

**Iterative Refinement** classical case. Use multiple solves to reduce the error.

**Mixed Precision** precondition high precision iterative methods with low precision factorization.

**Preconditioning** in general often uses incomplete factorizations.



## Is this a problem?

Yes.

Lots of use cases for multiple sequential solves following a factorize.

**Iterative Refinement** classical case. Use multiple solves to reduce the error.

**Mixed Precision** precondition high precision iterative methods with low precision factorization.

**Preconditioning** in general often uses incomplete factorizations.

**Interior-point Methods** multiple correctors and often aggressive iterative refinement.

## Can we fix it?

Need to exploit the cache.

*or*

Need to trade off **more computation** for **less memory bandwidth**.

## Can we fix it?

Need to exploit the cache.

*or*

Need to trade off **more computation** for **less memory bandwidth**.

**Combined factor-solve?**

## Can we fix it?

Need to exploit the cache.

*or*

Need to trade off **more computation** for **less memory bandwidth**.

Combined factor-solve?

Compression?

## Combined factor-solve

Forward substitution can be done at same time as factorization.  
(data is already in cache!)

Well established in out-of-core solvers.

Halves time for first solve

**But...** doesn't help on subsequent solves.

# Compression

Two *theoretical* limits:

**Compression ratio** How much bandwidth do we save?

**Decompression rate** How many extra cycles does this take?

Will **never go faster** by more than the compression ratio.

## Compression

Two *theoretical* limits:

**Compression ratio** How much bandwidth do we save?

**Decompression rate** How many extra cycles does this take?

Will **never go faster** by more than the compression ratio.

Exploit two levels of compression:

**Algorithm specific** We have deliberately introduced extra zeros to make the factorize go faster, does removing them help?

**Generic** Use a generic compression algorithm such as gzip, bzip2, LZO...

## Compression ratios

	LZO		LZO		LZO	
	nemin = 32 size	nemin = 32 size <b>CR</b>	nemin = 1 size	nemin = 1 <b>CR</b>	nemin = 1 size	nemin = 1 <b>CR</b>
1.	579	333 <b>1.74</b>	283	<b>2.04</b>	269	<b>2.14</b>
2.	981	567 <b>1.73</b>	489	<b>2.01</b>	463	<b>2.12</b>
3.	315	283 <b>1.11</b>	292	<b>1.08</b>	269	<b>1.17</b>
4.	407	328 <b>1.24</b>	342	<b>1.19</b>	303	<b>1.34</b>
5.	10829	10310 <b>1.05</b>	10304	<b>1.05</b>	10009	<b>1.08</b>

(size is storage for  $L$  in Megabytes)

1. CEMW/tmt\_sym
2. Schmid/thermal2
3. GHS\_psdef/crankseg\_1
4. DNVS/shipsec1
5. GHS\_psdef/audikw\_1



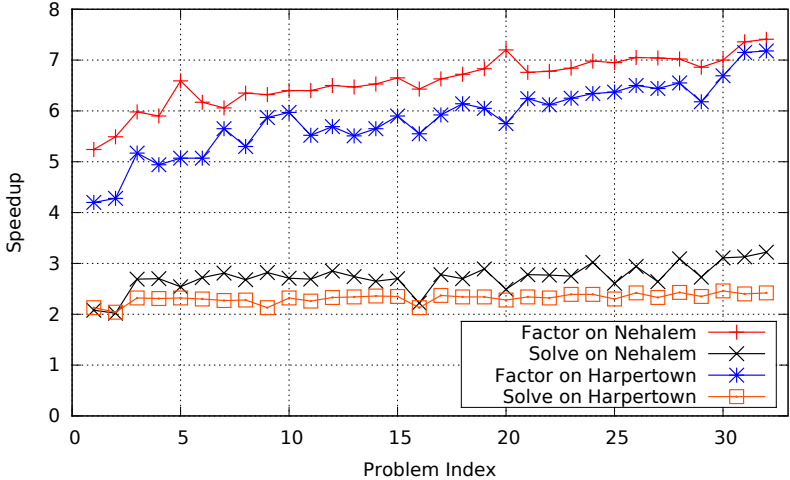
## In practice...

Not even getting close to theoretical maxima.

	4 cores			8 cores			CR
	without	with	ratio	without	with	ratio	
1.	0.44	0.38	<b>1.16</b>	0.30	0.27	<b>1.11</b>	<b>1.74</b>
2.	0.76	0.67	<b>1.13</b>	0.53	0.48	<b>1.10</b>	<b>1.73</b>
3.	0.19	0.19	<b>1.00</b>	0.11	0.11	<b>1.00</b>	<b>1.11</b>
4.	0.25	0.24	<b>1.04</b>	0.15	0.14	<b>1.07</b>	<b>1.24</b>
5.	6.90	6.92	<b>1.00</b>	3.61	3.59	<b>1.01</b>	<b>1.05</b>

(times in seconds)

# Are architectures improving?





## Conclusions

- ▶ Multicore is not SMP
- ▶ Memory-bound operations are a problem
- ▶ Solve phase of direct methods now significant
- ▶ No easy fix

Any questions?  
Any solutions?

# Time distribution, all phases.

