# Numerical optimisation and parallelisation

## How to speed up your code



**PSG group meeting**

**Tobias Wittwer**

**November 17, 2008**

**Delft Institute of Earth Observation and Space Systems**

**T U Delft**

Technische Universiteit Delft

# Motivation

Why numerical issues are important

- Our research is based on the results of computations
- Faster programs allow you to
    - Generate results quicker
    - Test different ideas in less time
    - Make more efficient use of your time
- Faster programs can be achieved by
    - Optimisation (increasing the speed of the code)
    - Parallelisation (using more computing power)

**T̃U**Delft

# Overview

Contents of the presentation

- Timing and profiling
- Some optimisation issues
- Matrix-matrix multiplication
- Choosing the right BLAS/LAPACK library
- Why packed storage is bad (for performance)
- Shared memory parallelisation
- Distributed memory parallelisation
- Computing resources

**TU**Delft

# Timing and profiling

Identifying computationally intensive program parts

- Only programs parts that require a lot of runtime need to be optimised

- `omp_get_wtime()` calls can be placed to identify slow program parts

- Profiling yields number of routine calls and time spent for executing routines, shows memory access problems

- Profiling usually requires unoptimised compilation and may thus lead to inaccurate numbers
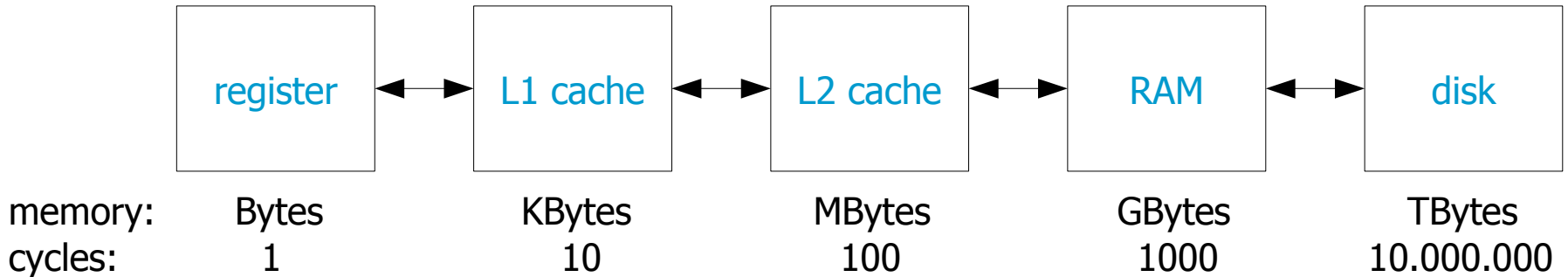
TUDelft

# Numerical optimisation

Some basic guidelines

- Compile with optimisations turned on (`-O3 -xP` on Intel systems, `-O3 -xW` on Cleopatra)
- Avoid *if* in loops
- Compute constant expressions only once and store them
- Expensive operators are:
  - trigonometric functions
  - square roots
  - exponentiation
  - division

**T**UDelft

# Numerical optimisation

## Minimise memory access

| register | | L1 cache | | L2 cache | | RAM | | disk |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | ↔ | | ↔ | | ↔ | | ↔ | |

|  | | | | |
|---|---|---|---|---|
| memory: | Bytes | KBytes | MBytes | GBytes | TBytes |
| cycles: | 1 | 10 | 100 | 1000 | 10.000.000 |

- Memory access is a bottleneck

- Store even results of small computations in scalars when used repeatedly in a loop

- Computing something can be faster than retrieving it from memory

TUDelft

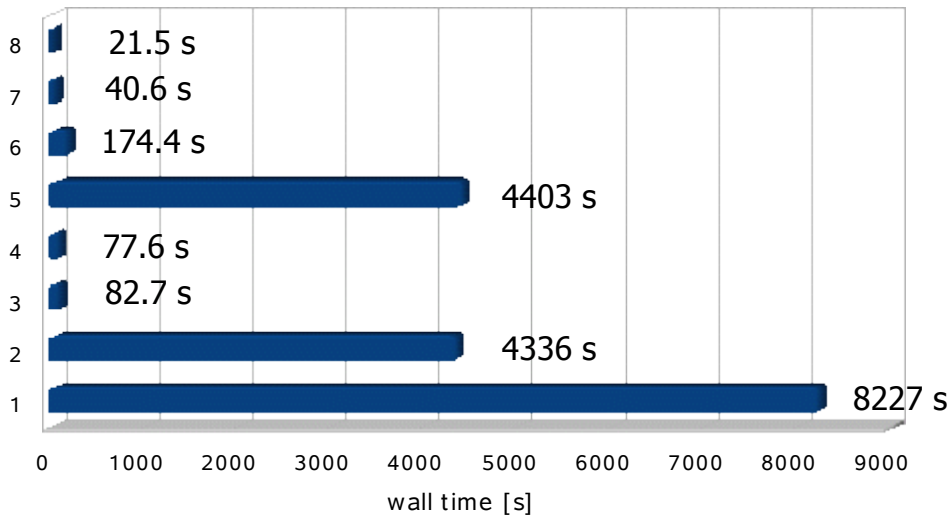# Numerical optimisation

## Series computations

$$\beta(x) = \sum_l \frac{2l+1}{4\pi} \left(\frac{R}{r}\right)^l \alpha(x, l)$$

- Constant expressions should be computed only once
- Exponentiation is expensive, can be computed recursively:

$$\left(\frac{R}{r}\right)^l = \left(\frac{R}{r}\right)^{l-1} \left(\frac{R}{r}\right)$$
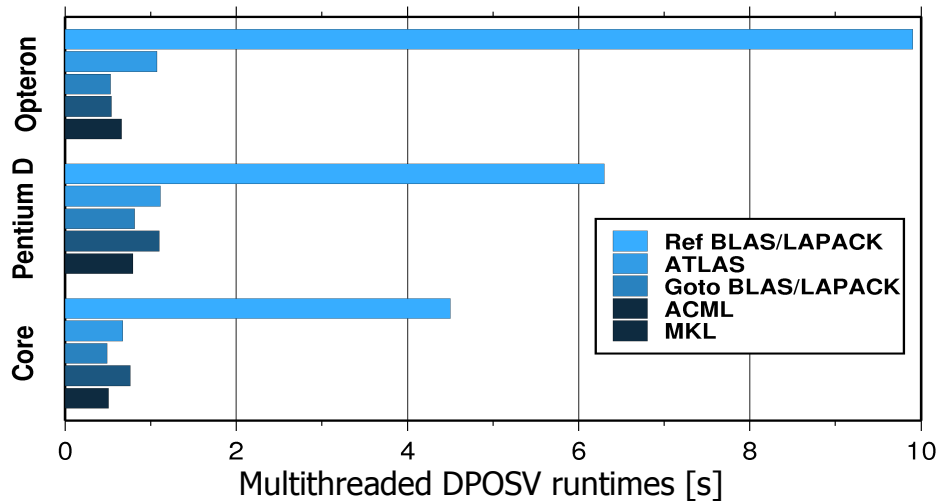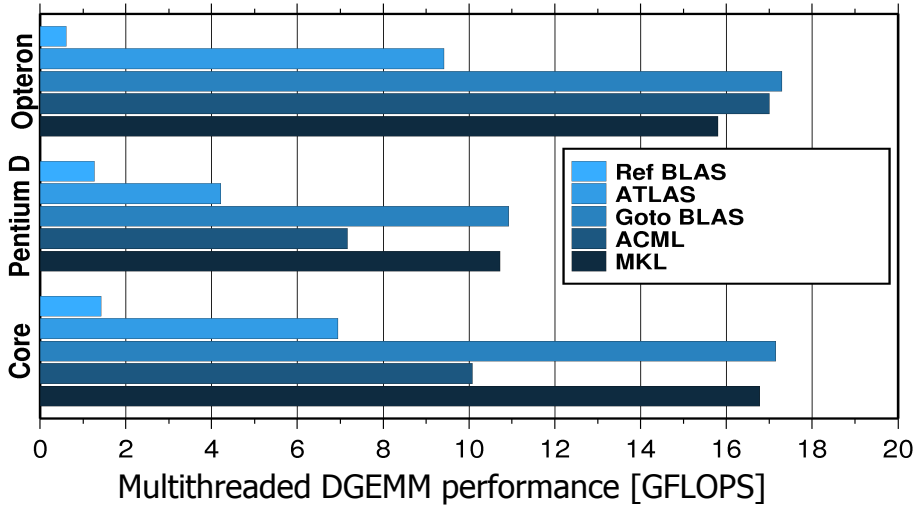
**T**UDelft

# Matrix-matrix multiplication
## Runtime comparison

21.5 s
40.6 s
174.4 s
4403 s
77.6 s
82.7 s
4336 s
8227 s

0  1000  2000  3000  4000  5000  6000  7000  8000  9000
wall time [s]

8 – DGEMM MKL 2 threads
7 – DGEMM MKL 1 thread
6 – DGEMM reference BLAS
5 – matmul()
4 – triple loop -O3 -xP -parallel
3 – triple loop -O3 -xP
2 – triple loop -O2
1 – triple loop -O0

- dual-core Pentium D, Intel FC 10.1
- vectorisation is fast (`-O3 -xP`)
- `matmul()` is very slow!
- reference BLAS is slow
- MKL is fast and parallelised
- BLAS can use `DSYRK` for $N=A'A$ (50% faster)

TUDelft

# Choosing the right BLAS/LAPACK lib



- Reference BLAS is slow!

- Goto BLAS or vendor-specific libraries should be used – optimised and parallelised

- LAPACK performance is governed by BLAS performance

TUDelft

# Why packed storage is bad

Packed storage LAPACK routines are slow

- Cholesky factorisation of 10201x10201 matrix
- Cleopatra, 4 threads
- Goto BLAS + LAPACK
- Unpacked: DPOTRF, 26.7 s
- Packed: DPPTRF, 631 s
- Packed storage offers less than 50% memory benefit, but is more than 20 times slower
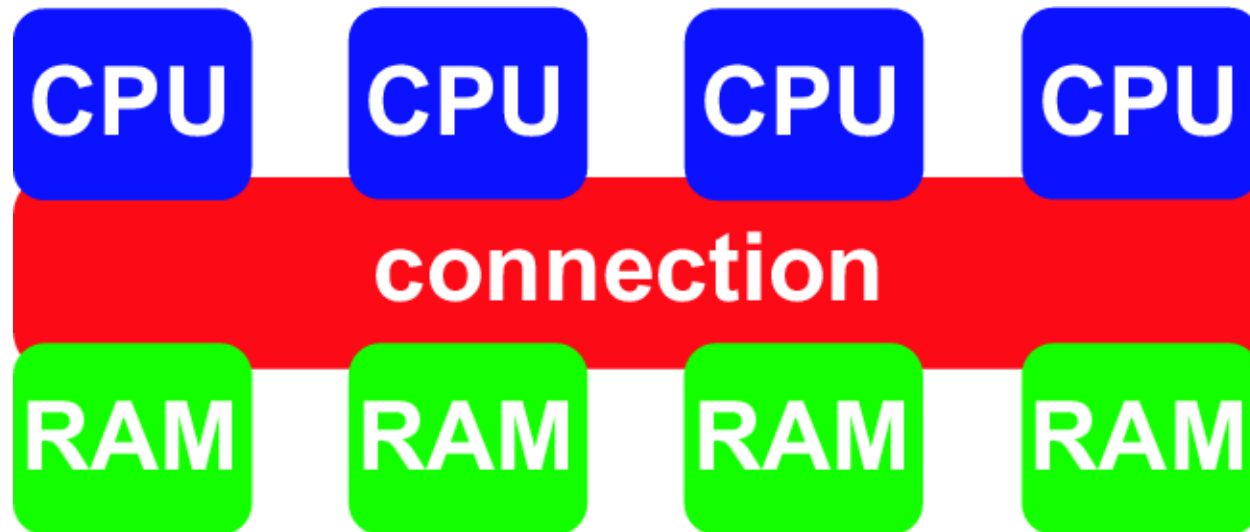- Avoid packed storage when runtime is important!

**T**U Delft

# Parallelisation

Why parallelisation is necessary

- Gigahertz-race has stopped
- Further developments aim at increasing the number of "cores" per CPU
- Parallel programming is required to make use of these parallel computer architectures
- Good parallel programming almost eliminates memory and computational restrictions
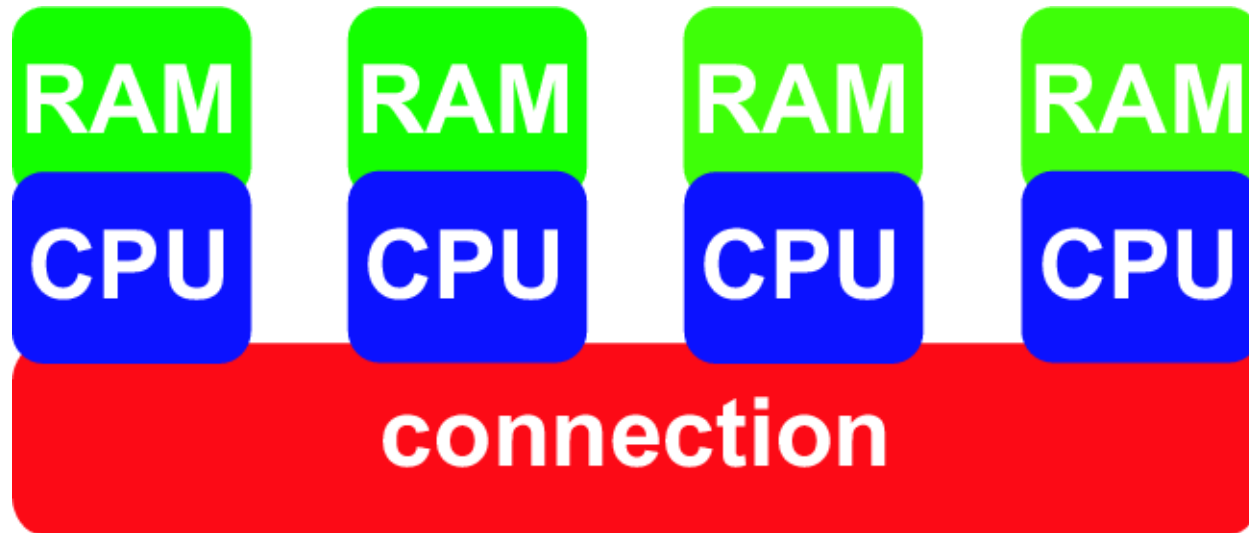
TUDelft

# Parallel computer architectures

Shared memory



- Easy to program, fast
- Expensive, limited maximum size
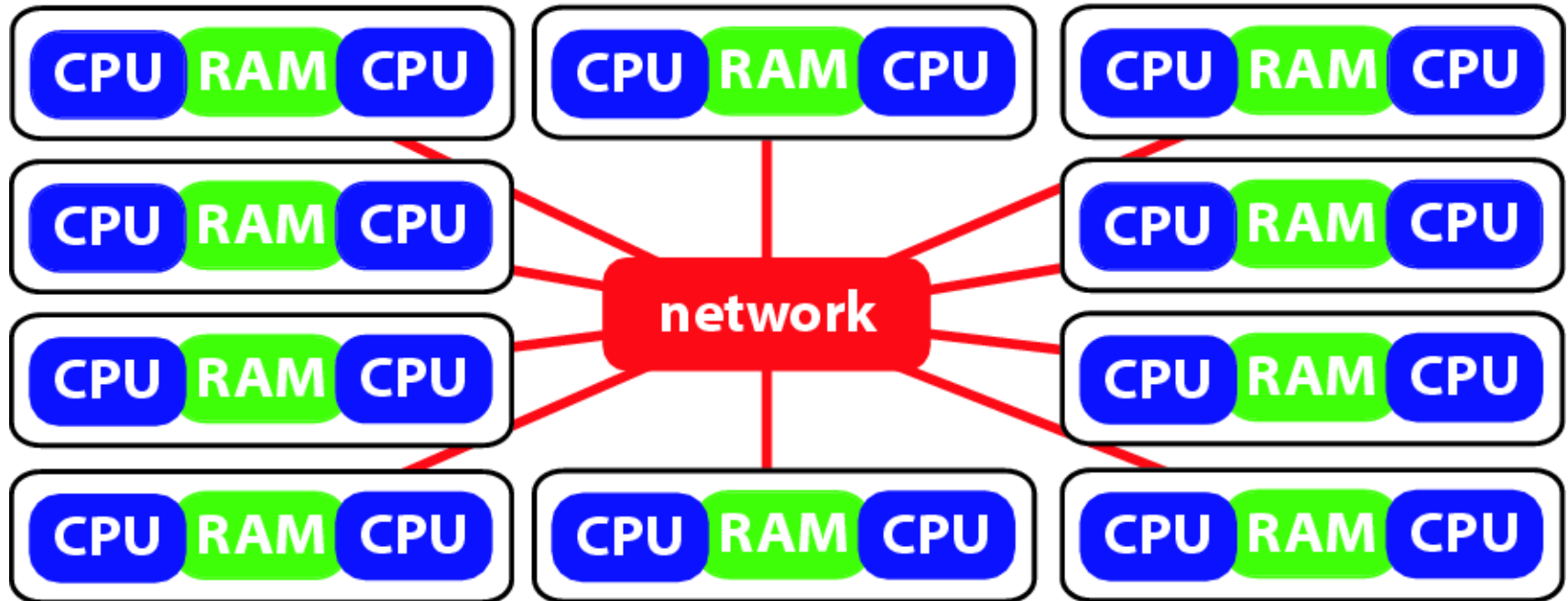
TUDelft

# Parallel computer architectures

Distributed memory



- Almost unlimited maximum size
- More difficult to program, slow if lots of communication is required

**T**UDelft

# Parallel computer architectures

Clusters



- Hybrid architecture, can be built from cheap components
- Similar benefits and limitations as distributed mem systems

**T**U Delft

# Shared memory parallelisation

## OpenMP

- OpenMP is used to parallelise loops
- Needs to be supported by the compiler
- Very easy to use

```
!$omp parallel do
do i=1,n
  A(i) = f(i)
end do
!$omp end parallel do
```

- Load balancing is done by operating system
- Memory-intensive programs might not benefit

TUDelft

# Distributed memory parallelisation

MPI

- MPI is used for explicitly exchanging data between processes

```
do i=myrank+1,n,nprocs
  A(i) = f(i)
end do
call mpi_allreduce(A,...)
```

- No automatic load balancing
- Performance is dependent on communications network
- Communication should be kept to a minimum

**T**U Delft

# Distributed memory parallelisation

## ScaLAPACK

- ACML, MKL, Goto BLAS & LAPACK routines are parallelised for shared memory only

- ScaLAPACK contains efficient distributed-memory routines of BLAS and LAPACK functionality

- Matrices are distributed among processes, full system memory can be used

- Uses MPI for communication:

  - Performance is dependent on network

  - Communication should be kept to a minimum

- Relatively complicated to use

**TU**Delft

# Computing resources

## Cleopatra

- 33 nodes

- 4 Opteron 280 nodes, 2.4 GHz

- 32 nodes with 8 GB RAM, 1 node with 16 GB RAM

- 272 GB total RAM

- Infiniband network

- 633 GFLOPS peak performance

- exclusively available to our group

**TU**Delft

# Computing resources

Huygens 2

- 104 nodes
- 32 Power 6 cores, 4.7 GHz
- 83 nodes with 128 GB RAM, 18 nodes with 256 GB RAM
- 15.6 TB total RAM
- Infiniband network
- 60 TFLOPS peak performance
- limited CPU time budget!

TUDelft

# Conclusions

Things to remember

- Compile with `-O3 -Xp` / `-O3 -Xw`
- Memory access is slow
- Small optimisations can speed up your code significantly
- Use BLAS and LAPACK routines for linear algebra
- Use ACML, MKL or Goto implementations
- Don't use packed storage for runtime-critical code
- Don't use Matlab for runtime-critical code
- Parallelise when necessary

**TU**Delft