# A Task-Based Parallel Programming Framework with Modularity, Scalability and Adaptability Features

Afshin Zafari, Martin Tillenius, Elisabeth Larsson
February 2014

Division of Scientific Computing, IT Department
Uppsala University, Sweden

Informationsteknologi

UPPSALA UNIVERSITET

# Outline

Informationsteknologi

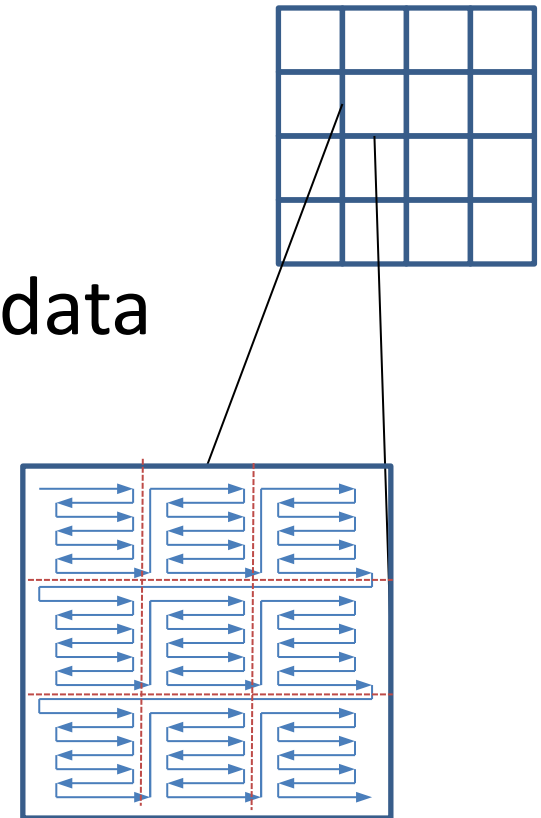# Task based parallel programming

- Program = <Operations , Operands>

- Algorithm := <Tasks , Data>

- Tasks := <Operations, In/Out Data>

- {Tasks} , {Data} → Scheduler → Run tasks in parallel

- Kernels: Actual computations

- **SuperGlue** and **DuctTeip** frameworks
  - (www.it.uu.se/research/scicomp/software/superglue)
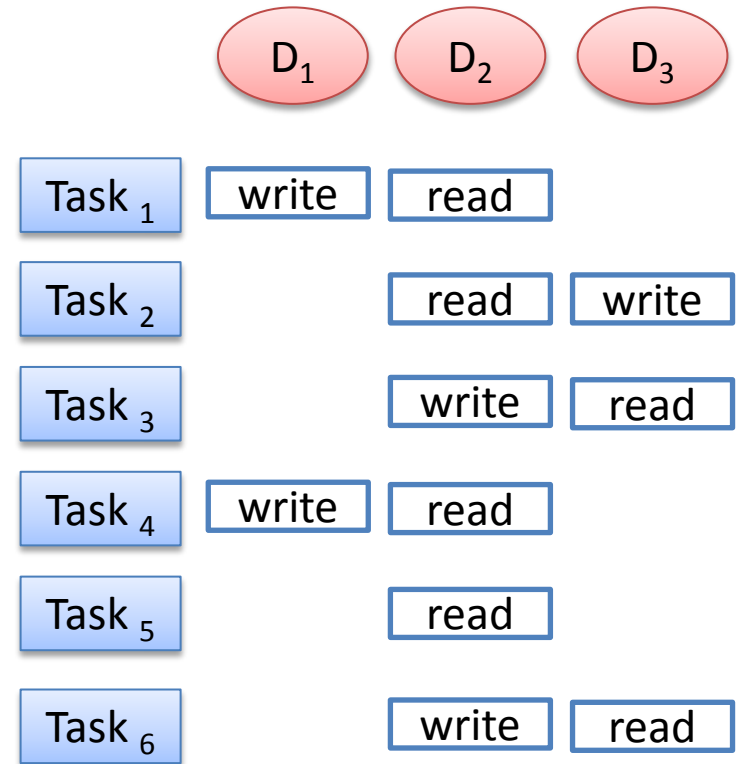
# Data in DuctTeip framework

- Processors are aligned in a virtual grid
- Data are partitioned in large/small scales
- Large data → communication
- Small data → computations
- Separate tasks for large/small data
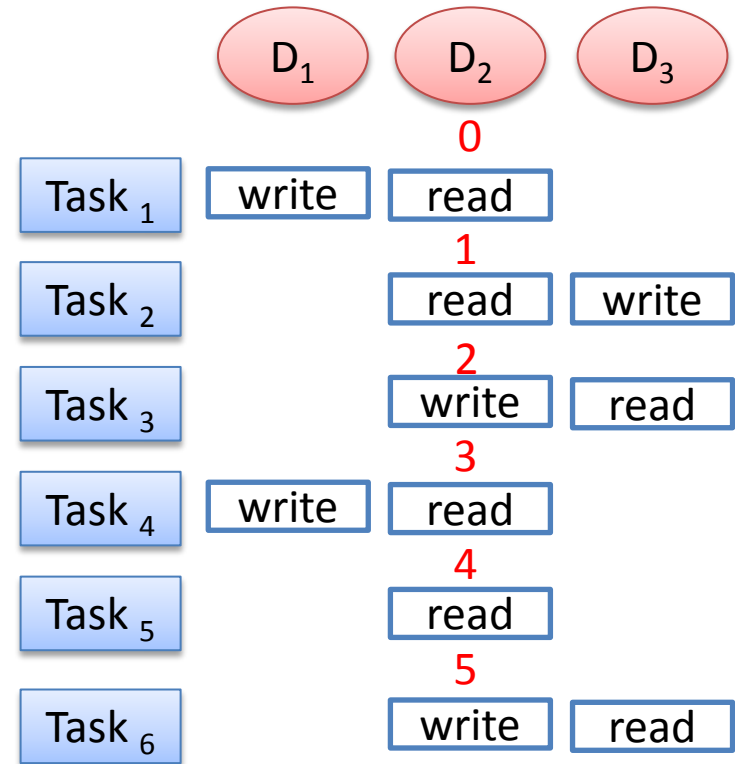- Efficient storage

# Data Versions

- Task-data dependency
- Data has *versions*
- Versions incremented after any access
- When versions of all read/write data are ready, task can run
- All ready tasks can run in parallel

$D_1$ $D_2$ $D_3$

| Task $1$ | write | read | |
| Task $2$ | | read | write |
| Task $3$ | | write | read |
| Task $4$ | write | read | |
| Task $5$ | | read | |
| Task $6$ | | write | read |

Informationsteknologi
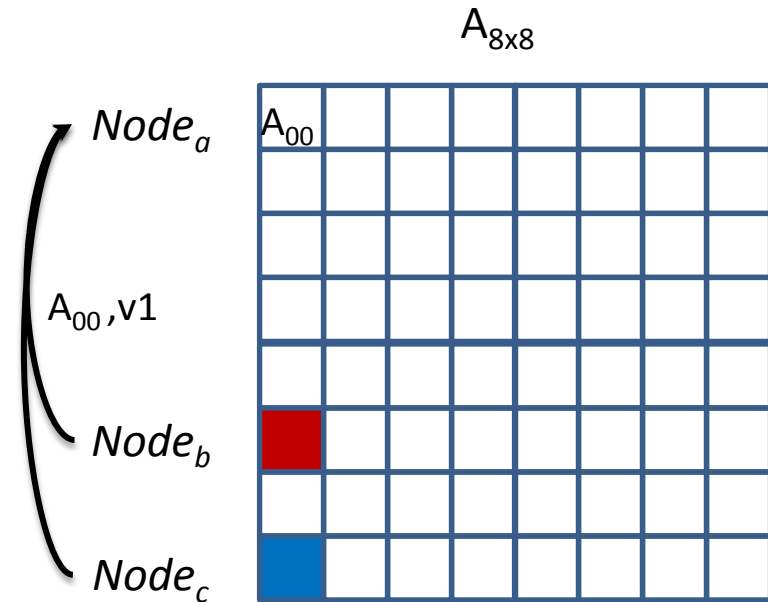
# Data Versions

- Task-data dependency

- Data has *versions*

- Versions incremented after any access

- When versions of all read/write data are ready, task can run

- All ready tasks can run in parallel

| | $D_1$ | $D_2$ | $D_3$ |
|---|---|---|---|
| | | 0 | |
| Task $_1$ | write | read | |
| | | 1 | |
| Task $_2$ | | read | write |
| | | 2 | |
| Task $_3$ | | write | read |
| | | 3 | |
| Task $_4$ | write | read | |
| | | 4 | |
| Task $_5$ | | read | |
| | | 5 | |
| Task $_6$ | | write | read |

6

# Distributed Environments

- Request for remote data → *listener*

- Data owner sends requested version of data, when it's ready

$A_{8x8}$

$Node_a$

$A_{00}$

$A_{00}, v1$

$Node_b$

$Node_c$

- Versions upgraded after listeners replied
- Duplicate listeners are replied once
- Requesters can handle many data and versions  ($D_1v_1$  $D_2v_1$  $D_3v_1$   $D_1v_2$  ...)
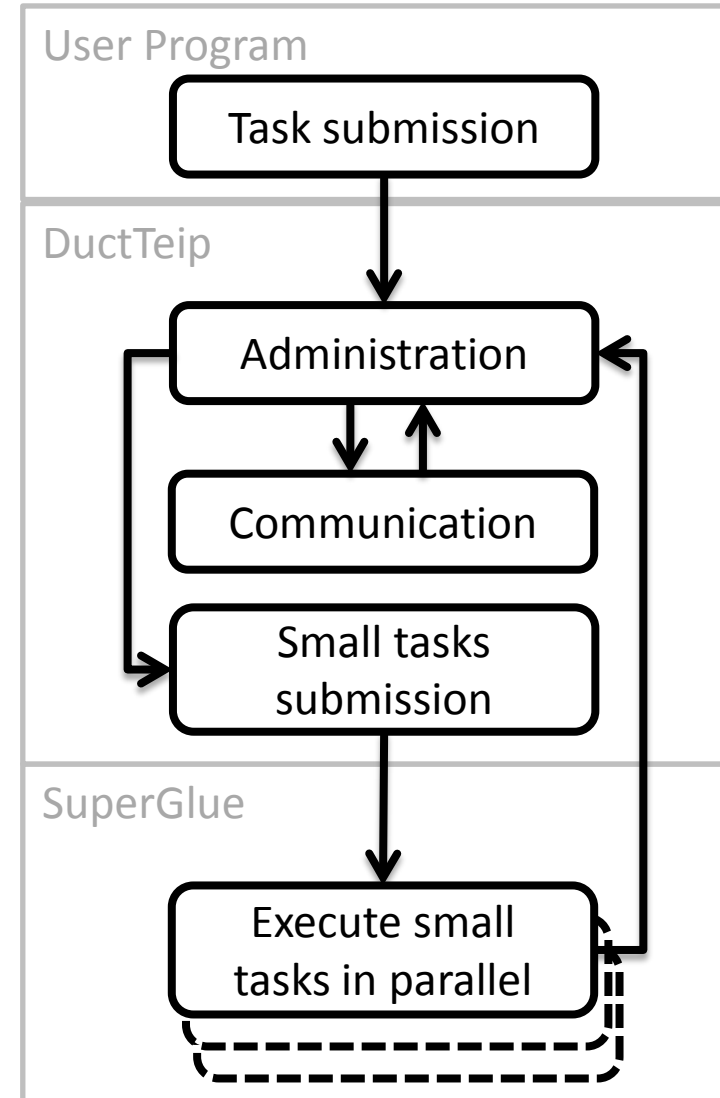
7

# How to use DuctTeip framework

- Configurations
  - Process grid (1D,2D,3D,...)
  - Two-level data partitioning: row/col/block cyclic
  - Row/col major ordering of data (e.g. for BLAS)
  - Who reads and who runs tasks: all/some/one
- User Program
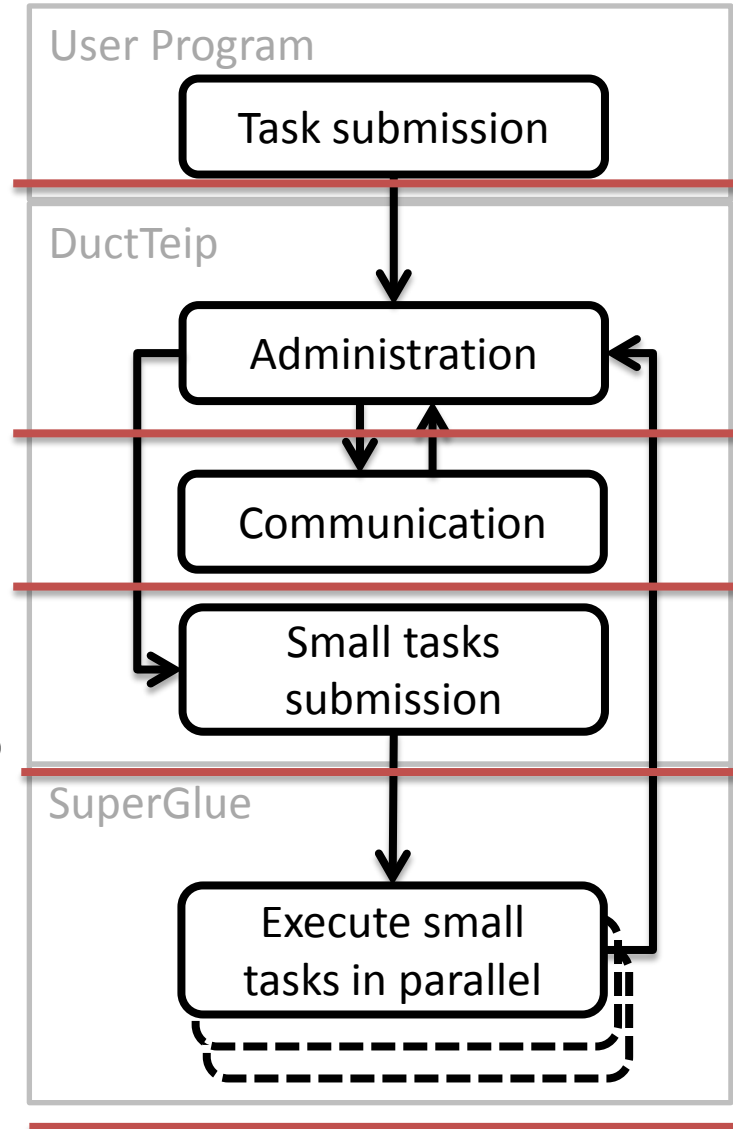  - Taskifies Algorithms
  - Implements kernels

Informationsteknologi

UPPSALA
UNIVERSITET

# How DuctTeip works

- Administration
  - Tracking versions
  - Handling tasks, listeners
- Communication
  -  tasks, listeners, data
- Execution
  - Submitting smaller tasks to **SuperGlue** framework

User Program
| Task submission |

DuctTeip
| Administration |
| Communication |
| Small tasks submission |

SuperGlue
| Execute small tasks in parallel |

# How DuctTeip works

- Administration
  - Tracking versions
  - Handling tasks, listeners
- Communication
  -  tasks, listeners, data
- Execution
  - Submitting smaller tasks to **SuperGlue** framework



Informationsteknologi

# Experiments

## Software

- Cholesky algorithm

1. ScaLAPACK
   - pgi 2013 + acml
   - openmpi 1.6.5
   - scalapack 2.0.2
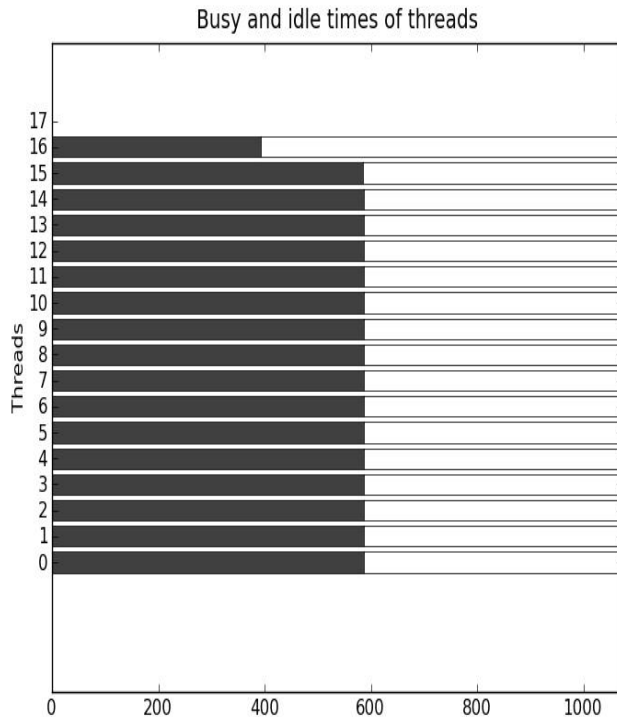
2. DuctTeip
   - Intel 13.1 + acml
   - openmpi 1.6

## Hardware

- UPPMAX Cluster
  - 166 Nodes
  - 2 Sockets/Node
  - 8 Cores/Socket
  - AMD 6220,3.0GHz
  - 32 GB RAM/Node
  - QDR Infiniband

Informationsteknologi

# Results – Execution Time

- Matrix Size : $142080^2$ , Process Grid:5x2
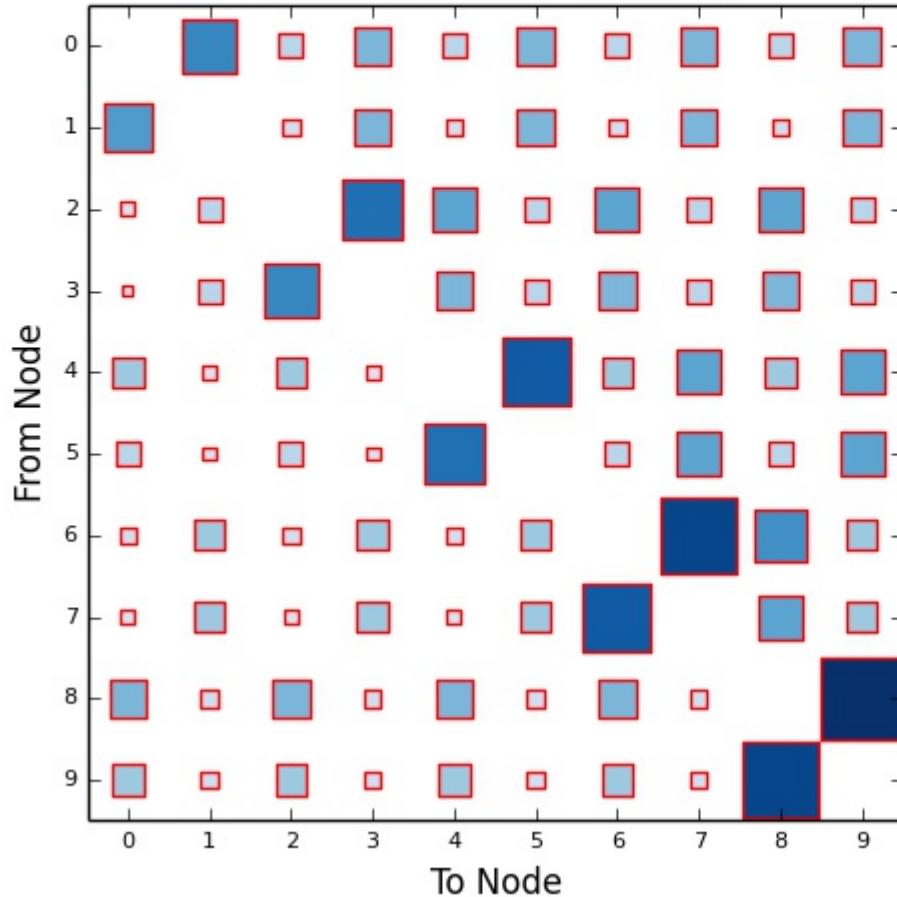- 1540 Large Tasks, 43,680,640 GEMM



DuctTeip

SuperGlue

- 110 Large Tasks
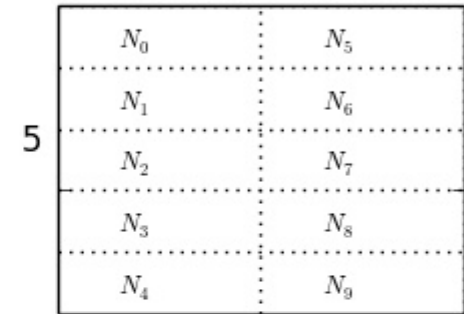- Overhead: 10%

- 202 Large Tasks
- Overhead: 3%

# Results – Communication



Informationsteknologi

# Results – Communication



Informationsteknologi

# Results – Strong Scaling



Strong scaling (10 nodes, 160 cores)

# Results – Weak Scaling



Weak scaling (starting matrix size $=12160^2$ )

Legend:
- GEMM Peak (Cores scale-up)
- GEMM Peak (Nodes scale-up)
- DuctTeip
- ScaLAPACK

Informationsteknologi

16

# **Conclusion**

- A Framework with:
  - Low Overhead
  - Scalability
  - Flexibility
    - Hybrid Parallel (Shared/Distributed Memory)
    - Hierarchical (two levels)
  - Modularity
    - Decoupled processes
  - Adaptability
    - Specific task, data objects