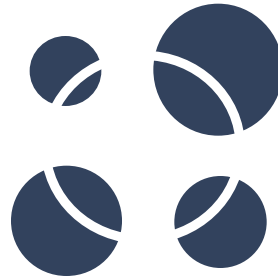
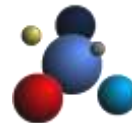


IberGrid 2007
Santiago de Compostela, Spain
May 15, 2007



Advanced Strategies for Efficient Workflow Management with GridWay

José Luis Vázquez-Poletti
Eduardo Huedo Cuesta
Rubén Santiago Montero
Ignacio Martín Llorente
<http://asds.dacya.ucm.es>



Distributed Systems Architecture Group
Universidad Complutense de Madrid



Contents

1. Protein Clustering Application (CD-HIT)
2. Parallel Execution of CD-HIT (cd-hit-para)
3. Porting cd-hit-para to the Grid
4. Results
5. Current Work



1. Protein Clustering Application (CD-HIT)

CD-HIT

- “*Cluster Database at High Identity with Tolerance*”
- Protein (and also DNA) clustering
 - Compares protein DB entries
 - Eliminates redundancies
- Example: Used in UniProt for generating UniRef data sets
- Our case: Widely used in the Spanish National Oncology Research Center (CNIO)
 - Input DB: 504,876 proteins / 435MB
- Infeasible to be executed on single machine
 - **Memory requirements**
 - Total execution time

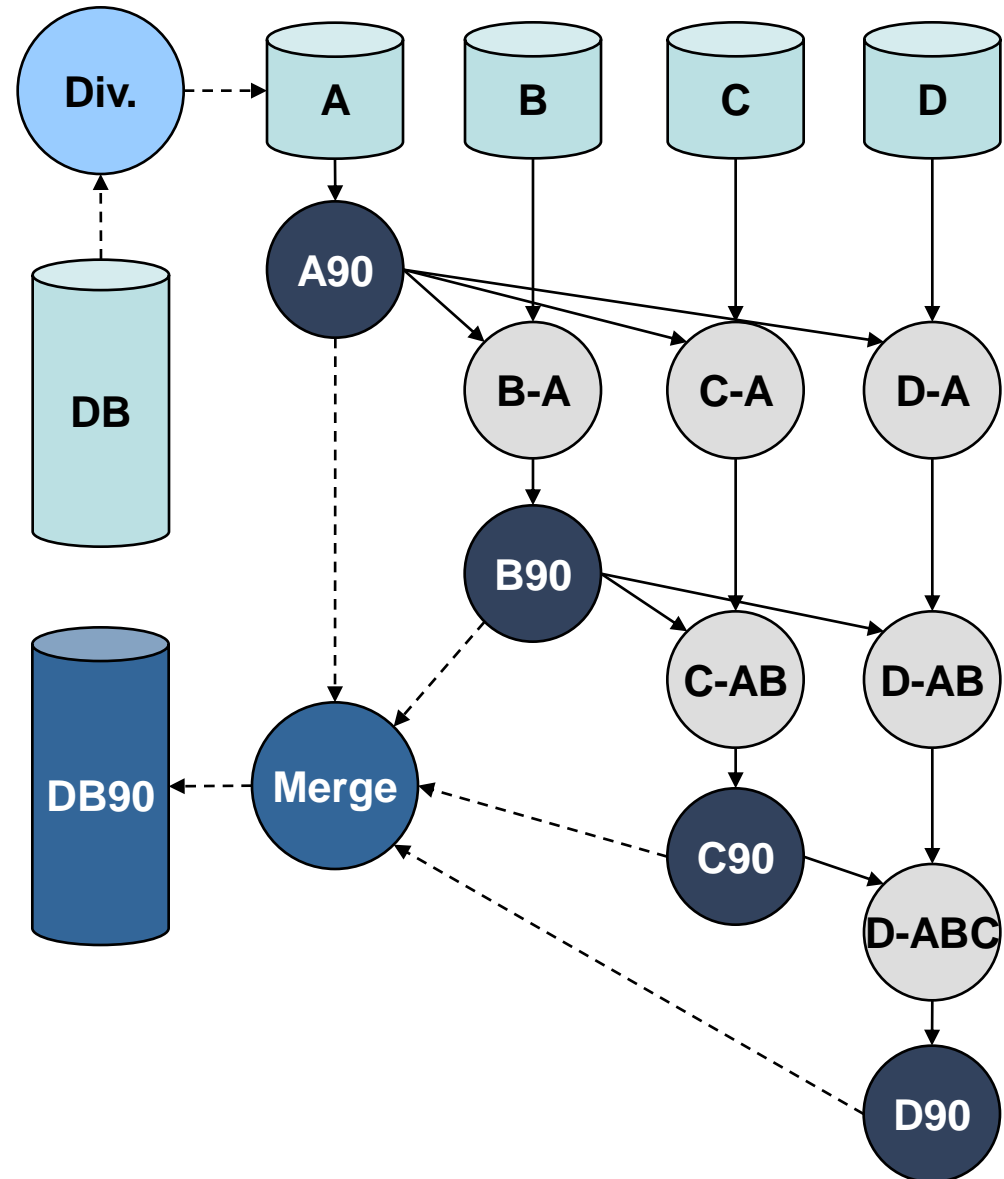




2. Parallel Execution of CD-HIT (cd-hit-para)

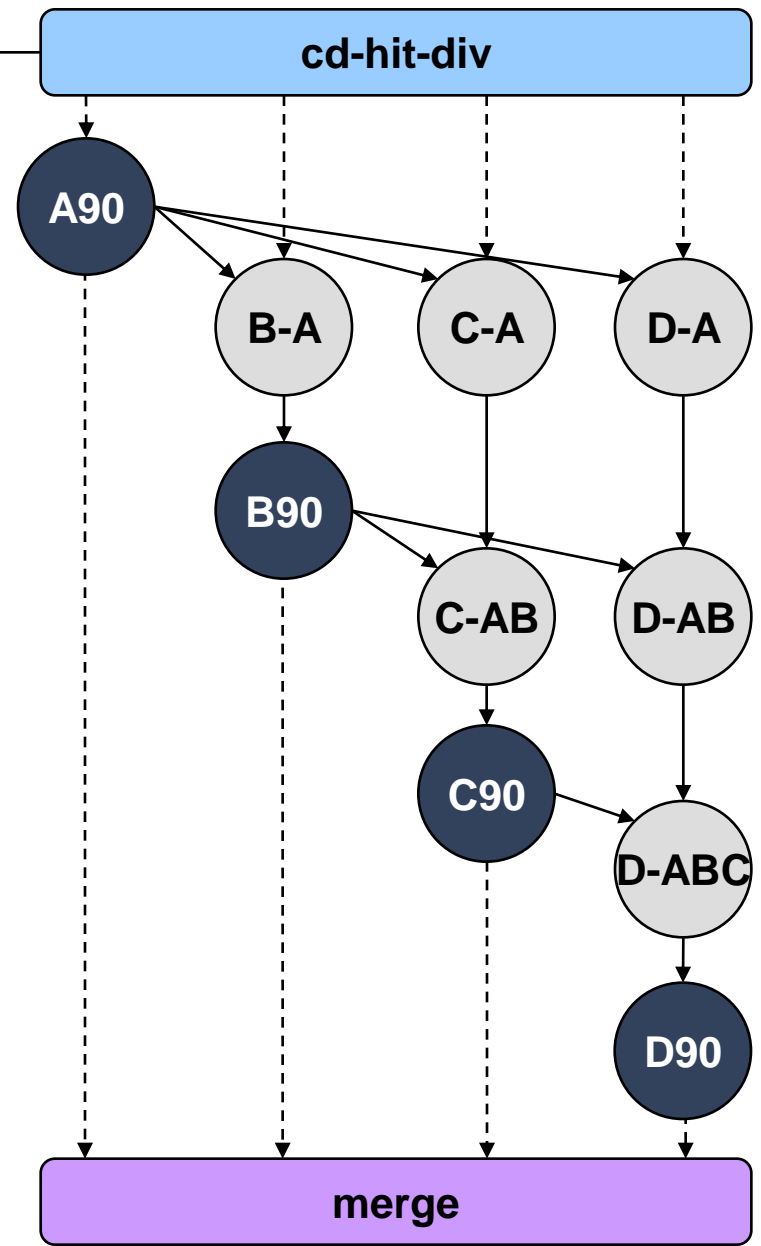
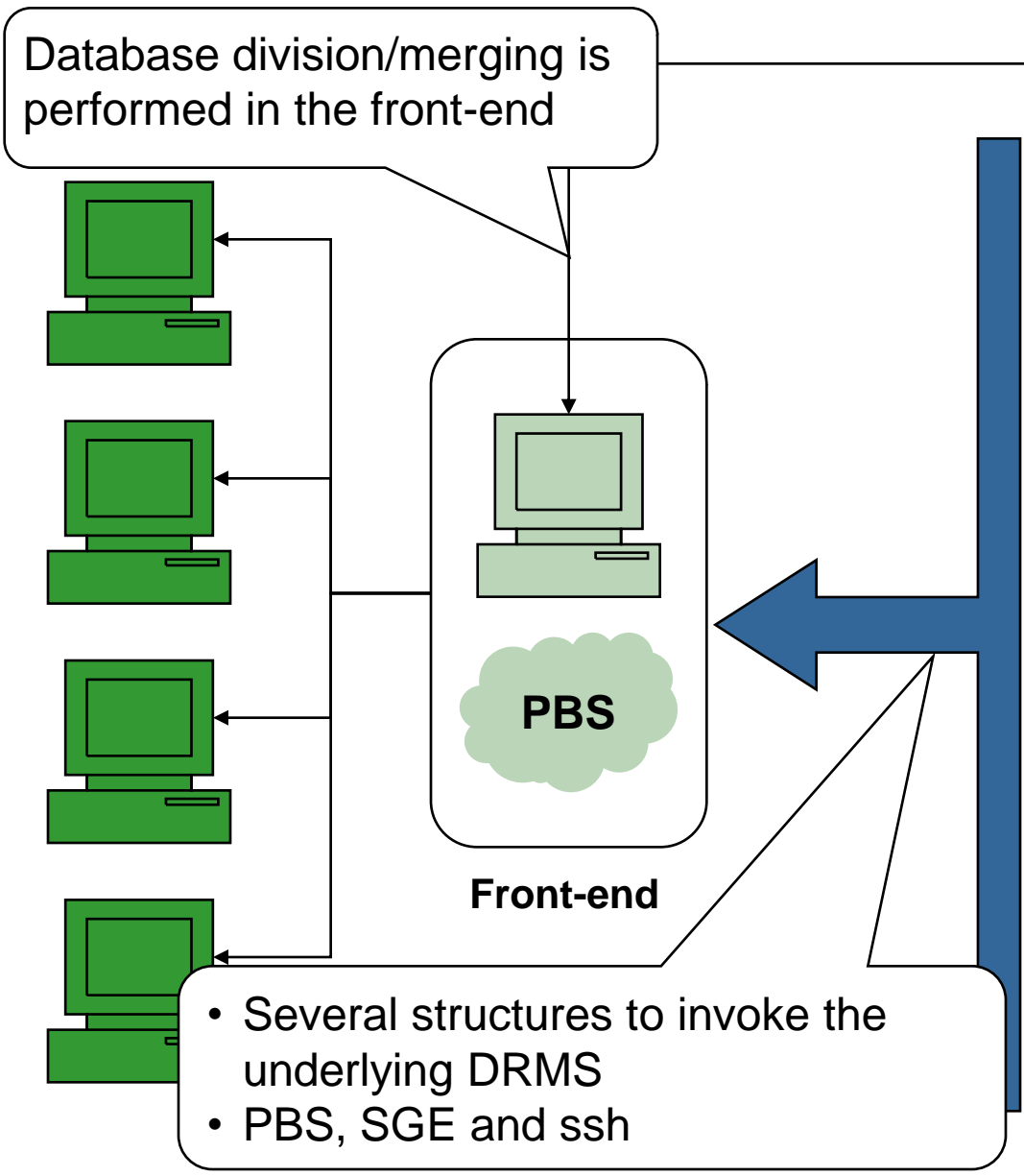
cd-hit-para

- Execute cd-hit in **parallel mode**
- **Idea:** divide the input database to compare each division in parallel
 - Divide the input db
 - Repeat
 - Cluster the first division (cd-hit)
 - Compare others against this one (cd-hit-2d)
 - Merge results
- Speed-up the process and deal with **larger databases**
- **Computational characteristics**
 - Variable degree of parallelism
 - Grain must be adjusted



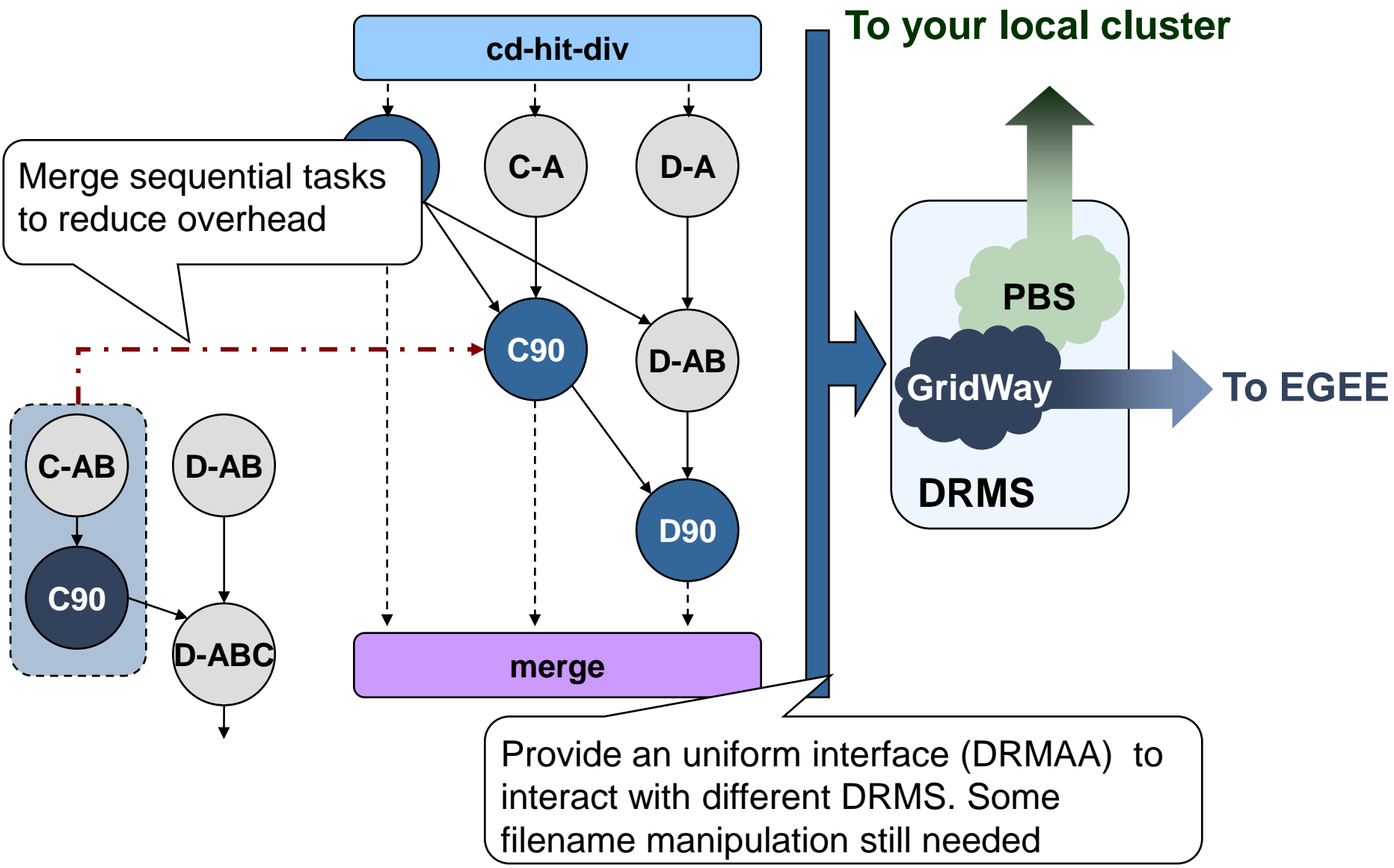


2. Parallel Execution of CD-HIT (cd-hit-para)





3. Porting cd-hit-para to the Grid





Grid Infrastructure

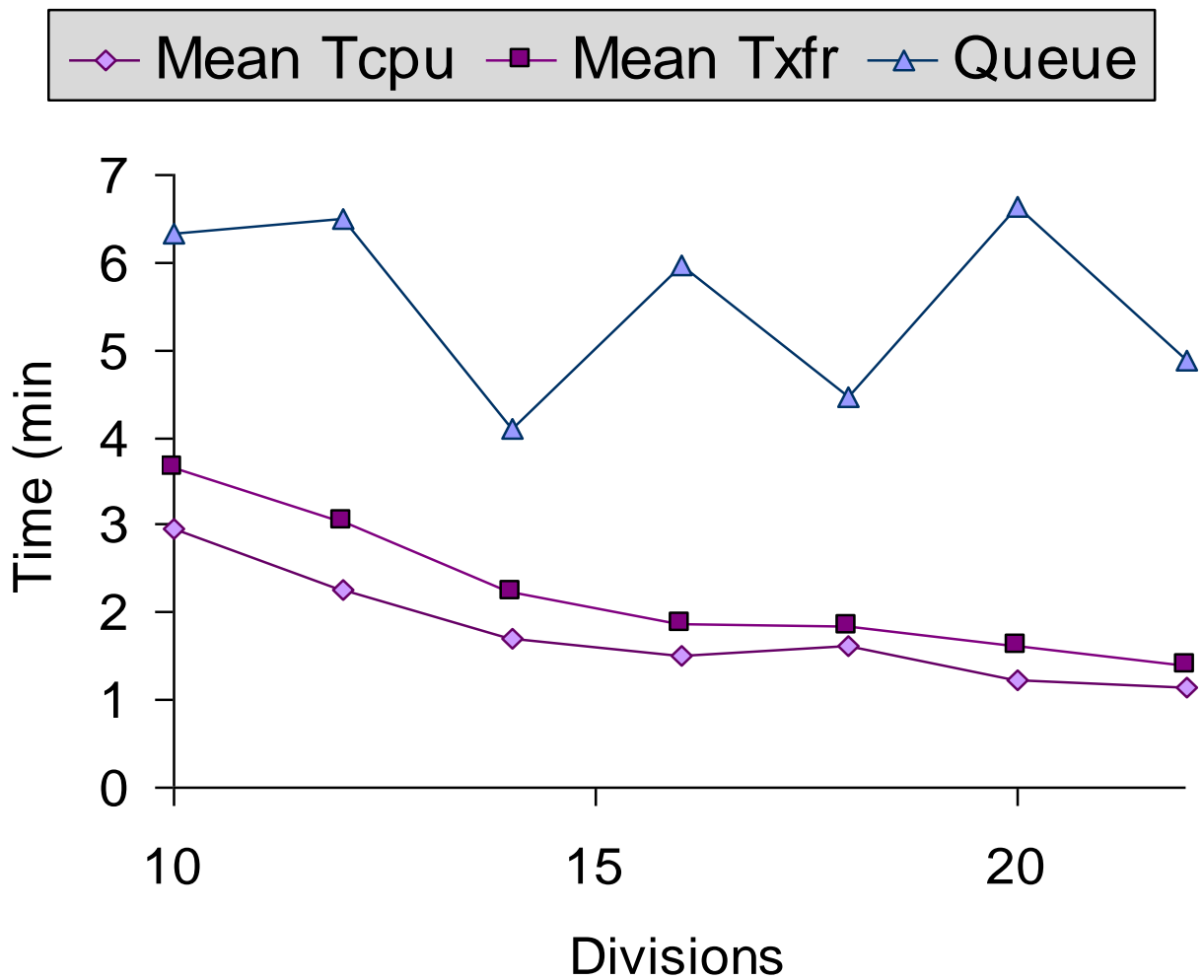
BIOMED sites

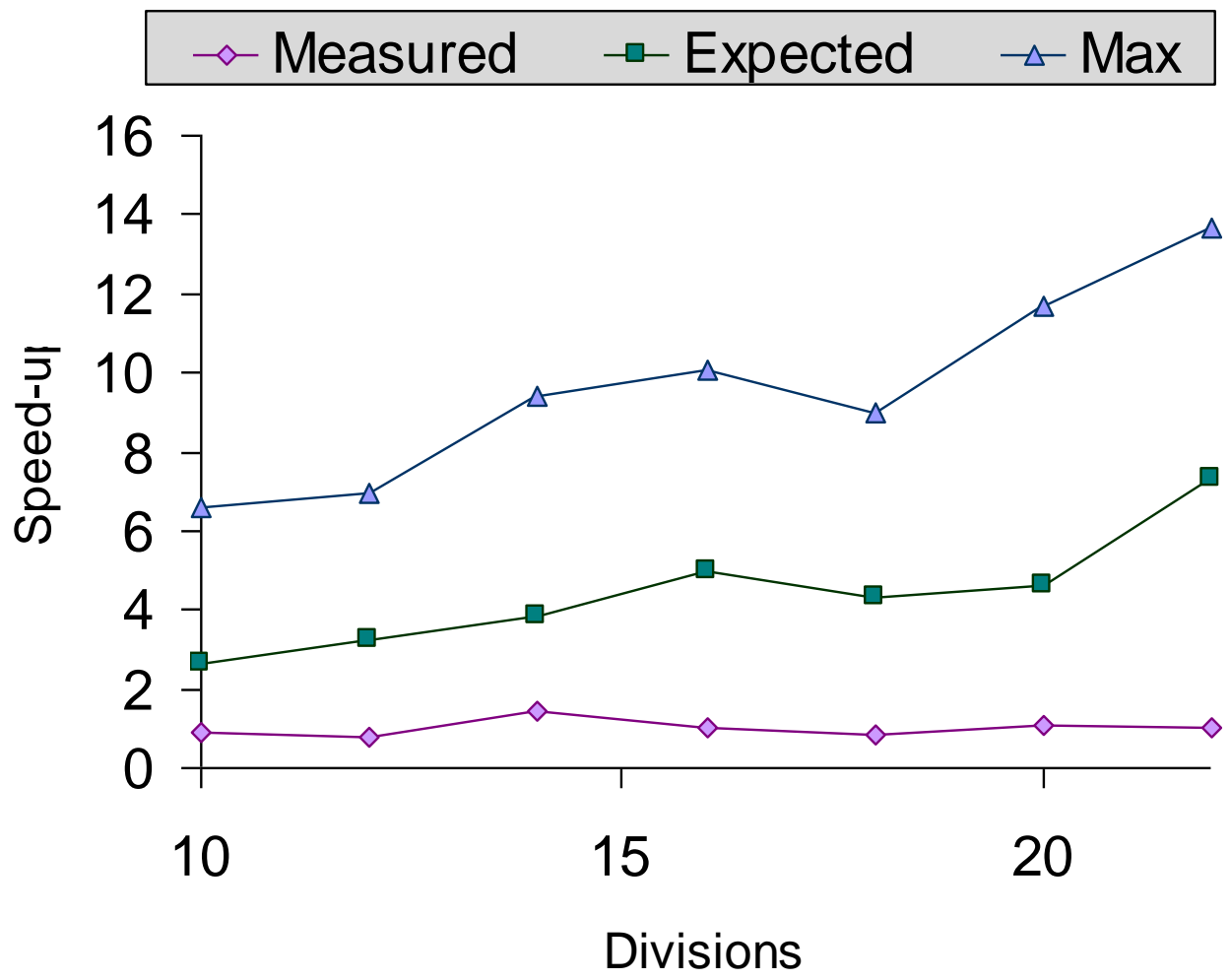


Site	Processor	Nodes	Speed
BIFI	Intel P4	56	3.2GHz
CESGA	Intel PIII	16	500MHz
CGG	Intel PIII	58	1.2GHz
CIEMAT	Intel Xeon	226	3.2GHz
GRIF	Intel P4	14	2.8GHz
JINR	Intel PD	30	2.8GHz
L.-HEP	Intel P4	374	3GHz
PNPI	Intel P4	60	3GHz
RAL	Intel P4	62	2.8GHz
RALPP	Intel PIII	1064	1GHz
ScotGRID	Intel Xeon	6	2.8GHz
SINP	Intel Xeon	94	2.8GHz



DB div.	Mean size	Tasks
10	44MB	45
12	36.5MB	66
14	31.5MB	91
16	27.5MB	120
18	24.5MB	153
20	22MB	190
22	20MB	231





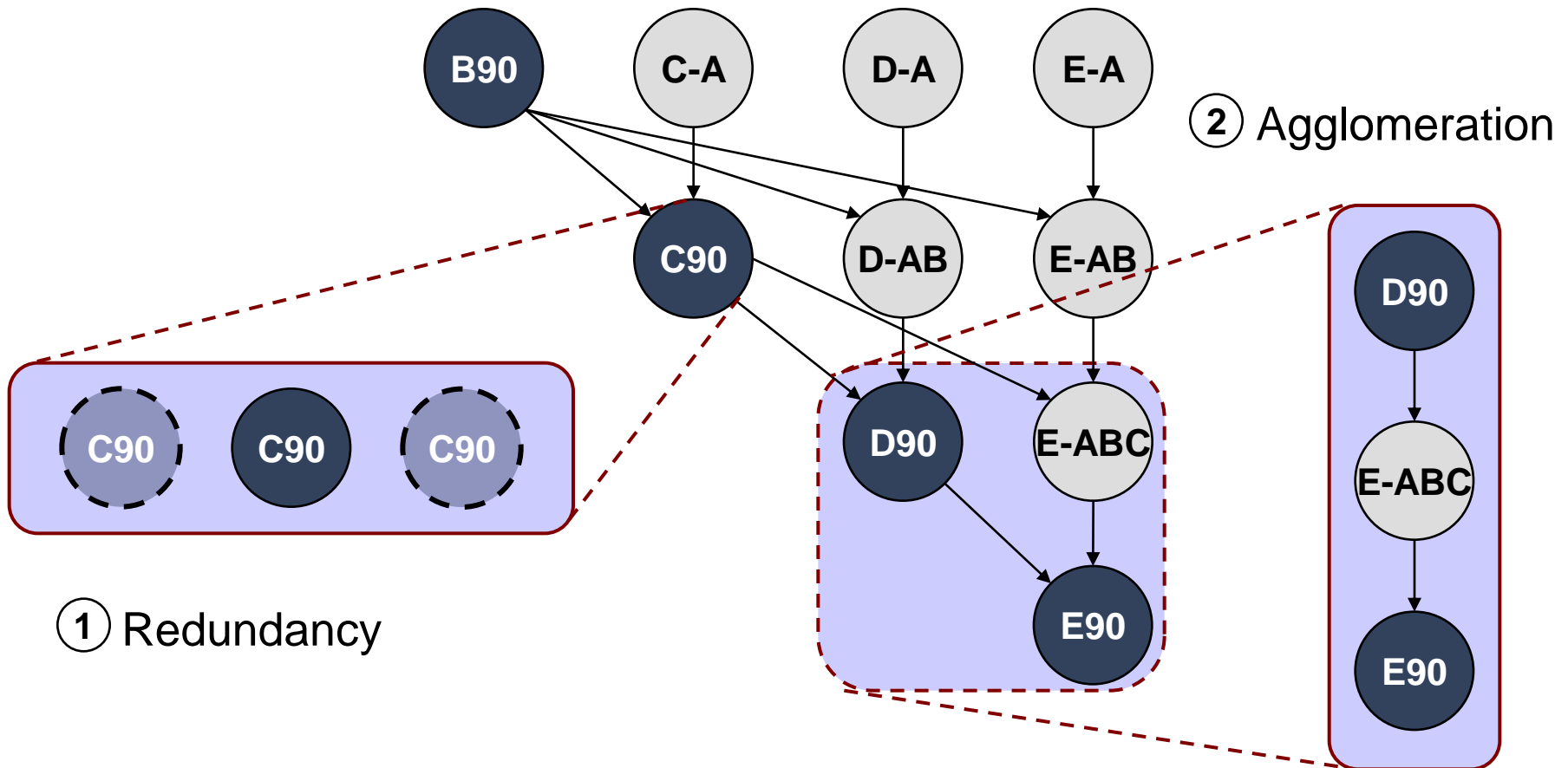
Expected: without considering either queue wait times or job failures
Max: taking into account just the *critical path*



5. Current Work

cd-hit-grid

- Larger database 4,186,284 proteins (1.7 GB) from CNIO
- Improve the efficiency of the workflow execution with redundancy and agglomeration





Information and download at <http://www.GridWay.org>

The screenshot shows a Mozilla Firefox browser window displaying the GridWay Metascheduler website. The browser's address bar shows the URL <http://www.gridway.org/>. The website has a dark blue sidebar on the left with the GridWay logo and a 'Contents' menu. The main content area features the title 'GridWay Metascheduler' and 'Metascheduling Technologies for the Grid'. It includes navigation links for 'GridWay5.2 Flier', 'FAQ', 'Sitemap', and 'Contact us', along with the 'the globus alliance' logo. The page contains a 'WELCOME TO GRIDWAY' section with a paragraph describing the scheduler's capabilities, a 'WHY GRIDWAY?' section with a bulleted list of benefits, and a concluding paragraph about its use in grid infrastructures.

GridWay Metascheduler
Metascheduling Technologies for the Grid

GridWay5.2 Flier FAQ Sitemap Contact us

WELCOME TO GRIDWAY

The GridWay Metascheduler enables large-scale, reliable and efficient sharing of computing resources (clusters, computing farms, servers, supercomputers...), managed by different LRM (Local Resource Management) systems, such as PBS, SGE, LSF, Condor..., within a single organization (enterprise grid) or scattered across several administrative domains (partner or supply-chain grid). GridWay is a Globus project, adhering to Globus philosophy and guidelines for collaborative development and so welcoming code and support contributions from individuals and corporations around the world.

WHY GRIDWAY?

There exist a number of commercial and open source workload management and scheduling systems available today, each one suitable for different underlying computer infrastructures and execution profiles. GridWay stands out from other metascheduling systems because it has been specifically designed to work on top of Globus services, offering the highest functionality, quality of service and reliability on this kind of infrastructures, namely:

- ◆ **For project and infrastructure directors.** GridWay is an open-source community project, adhering to Globus philosophy and guidelines for collaborative development.
- ◆ **For system integrators.** GridWay is highly modular, allowing adaptation to different grid infrastructures, and supports several OGF standards.
- ◆ **For system managers.** GridWay gives a scheduling framework similar to that found on local LRM systems, supporting resource accounting and the definition of state-of-the-art scheduling policies.
- ◆ **For application developers.** GridWay implements the OGF standard DRMAA API (C and JAVA bindings), assuring compatibility of applications with LRM systems that implement the standard, such as SGE, Condor, Torque,...
- ◆ **For end users.** GridWay provides a LRM-like CLI for submitting, monitoring, synchronizing and controlling jobs, that could be described using the OGF standard JSDL.

With GridWay, a Grid infrastructure can be exploited and managed in the same way as a local computing cluster. We invite you to check its Metascheduling Functionality Checklist and its benefits for end users and system administrators.

Terminado



**Thank you
for your attention!**

