EUDAT Collaborative Data Infrastructure

PRACE

# Data movement

## General concepts and specific tools

*Claudio Cacciari | SURFsara*

**How can I move my data?**

Where do you want to move your data from?

And where to?

And which amount of data?

# what are you using to move data?

# Poll: first question



http://etc.ch/riGa

# basic concepts

# Bandwidth, throughput and speed

https://www.youtube.com/embed/TVpg7StOxgg

# Advanced concepts

- https://www.youtube.com/embed/TVpg7StOxgg

- https://www.youtube.com/embed/aD_yi5VjF78

- https://www.youtube.com/embed/6IP0ow8Voe0

# Your experience

## How much bandwidth do you have?

# Poll: second question

http://etc.ch/riGa

# Network requirements and expectations

http://fasterdata.es.net/home/requirements-and-expectations

# Real world

https://viz.measurementlab.net

# A nice introduction

https://princetonuniversity.github.io/PUbootcamp/sessions/data-transfer-basics/ PUBootCamp_20181031_DataTransfer.pdf

## Transfer tools: scp vs. GridFTP

Sample Results: disk-to-disk testing from Berkeley, CA to Argonne, IL (near Chicago). RTT = 53 ms, network capacity = 10Gbps, RAID = 4 disks, RAID Level-0. **Note that to get more than 1 Gbps (125 MB/s) disk to disk requires RAID.**

| Tool | Throughput | Downloading 500 GB data |
|---|---|---|
| scp | 140 Mbps (17.5 MB/s) | **8 hours** |
| HPN patched scp, 1 disk | 760 Mbps (95 MB/s) | |
| HPN patched scp, RAID disk | 1.2 Gbps (150 MB/s) | |
| GridFTP, 1 stream, 1 disk | 760 Mbps (95 MB/s) | **1.5 hours** |
| GridFTP, 1 stream, RAID disk | 1.4 Gbps (175 MB/s) | |
| GridFTP, 4 streams, RAID disk | 5.4 Gbps (675 MB/s) | |
| GridFTP, 8 streams, RAID disk | 6.6 Gbps (825 MB/s) | **10 minutes** |

(ref: http://fasterdata.es.net/data-transfer-tools/)

OFFICE of INFORMATION TECHNOLOGY

PICSciE 24

## How to select a transfer tool

| • Transfer takes less than 10 mins<br>• Not that frequent | Other than that, your transfer job is a noticeable chunk in your workflow |
|---|---|
| **SCP (WinSCP)**<br>**FTP (FileZilla)**<br>**rsync** | **Globus (GridFTP)**<br>**BBCP**<br>**LFTP**<br>**Aria2c**<br>**FDT**<br>**…** |

OFFICE of
INFORMATION TECHNOLOGY

PICSciE

28

## (extra) Transfer settings: Encryption

| Tool | Encrypted Control | Encrypted Data |
|---|---|---|
| FTP<br>HTTP (even password-based access) | | |
| BBCP<br>BBFTP<br>Globus/GridFTP | ✔ | |
| SCP<br>SFTP<br>rsync over SSH<br>Globus/GridFTP with encryption-on<br>HTTPS | ✔ | ✔ |

Data encryption provides best security,
but negatively impacts transfer speed (10-50% slower)

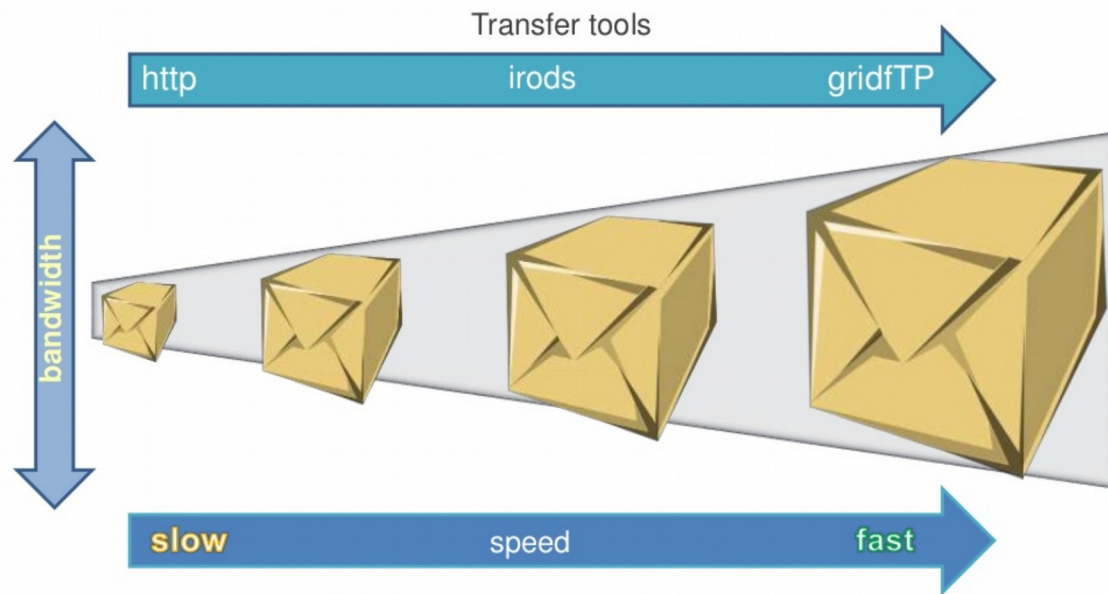OFFICE of
INFORMATION TECHNOLOGY

PICSciE  **29**

## Summary and Best Practices

- **Data transfer speed is affected by: Endpoints, network, and transfer tool**

- **Know the limitation of your endpoints**

- **Used wired instead of wireless for large transfers**

- **Seek better transfer tools if transfer takes > 10 minutes and happens frequently**
  - e.g., Globus, BBCP

- **Ask for help**
  - Your department IT staff
  - About using RC resources: cses@princeton.edu

OFFICE of
INFORMATION TECHNOLOGY

PICSciE

30

# The right tools



High Performance Transfers

Transfer tools: http, irods, gridfTP

bandwidth

slow — speed — fast

EUDAT Summer School, 3-7 July 2017, Crete

44

Moving large amounts of data around

Data sets can be large both in terms of
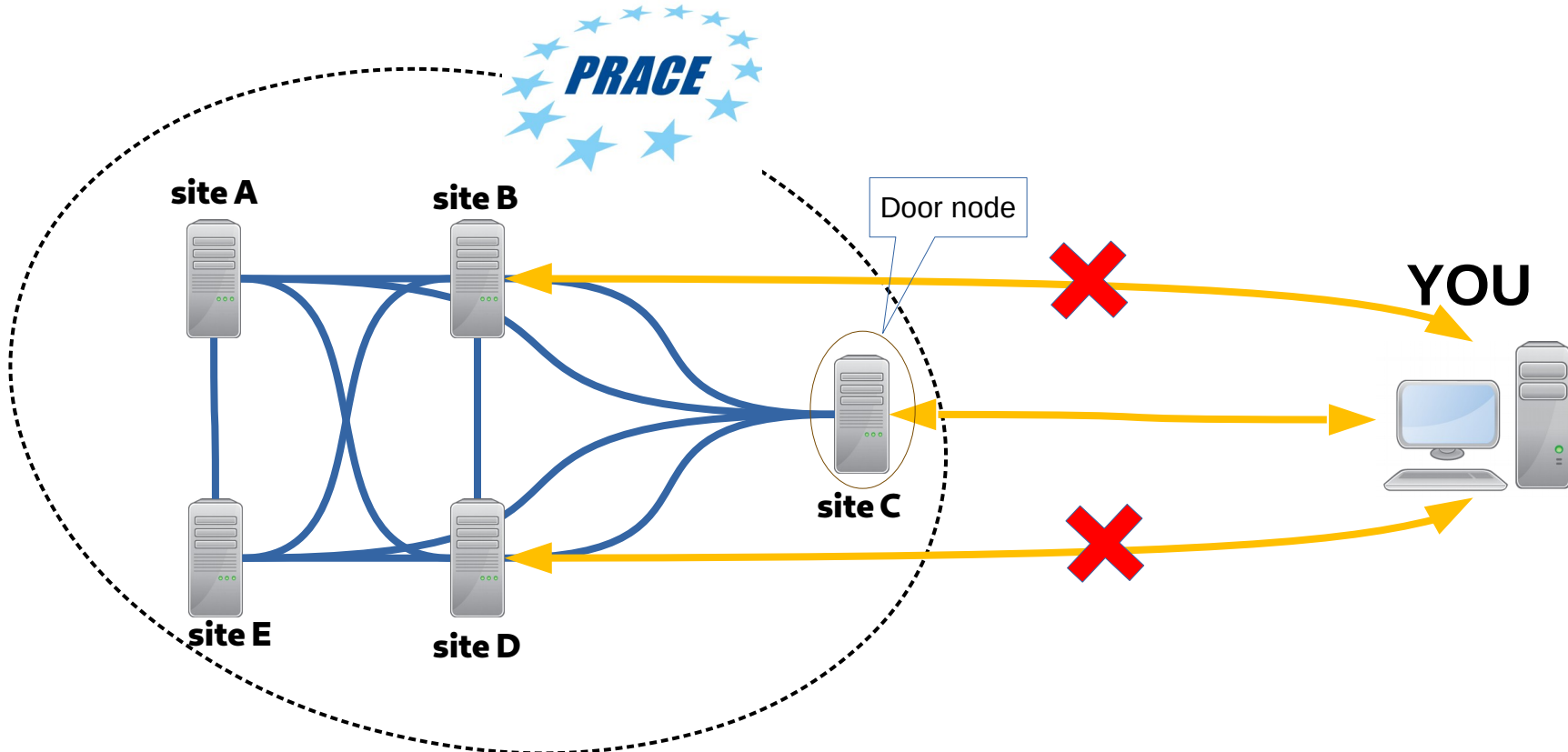
Numbers of objects

Single object size

# PRACE: which tools are available?
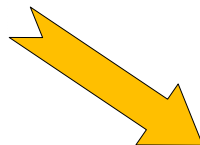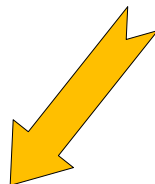
http://www.prace-ri.eu/data-transfer-with-gridftp

## GridFTP: The Protocol

- **Existing standards**
  - RFC 959: File Transfer Protocol
  - RFC 2228: FTP Security Extensions
  - RFC 2389: Feature Negotiation for the File Transfer Protocol
  - Draft: FTP Extensions
  - GridFTP: Protocol Extensions to FTP for the Grid
    - Grid Forum Recommendation
    - GFD.20
    - http://www.ggf.org/documents/GWD-R/GFD-R.020.pdf

# GridFTP software

Free and open source

https://gridcf.org/

**Grid Community Forum**

Community-based support for core software packages in grid computing

## Overview

The Grid Community Forum (GridCF) is a global community that provides support for core grid software.

Specifically, the GridCF is attempting to support a software stack christened the Grid Community Toolkit (GCT). The GCT is an open-source fork of the venerable Globus Toolkit created by the Globus Alliance. The GCT is *derived* from the Globus Toolkit, but is not the Globus Toolkit. Further, the GridCF is not a part of the Globus Alliance.

The GridCF is a nascent organization: we are looking for energetic contributors across a broad range of technical skills. Check out our governance doc and join us on GitHub!

https://www.globus.org/

Free and open source

On the GridFTP Door Node server

- Globus toolkit has been installed,
- connections to the PRACE network and thus to the GridFTP servers at every PRACE site
- the machine can be accessed from the public internet,

This requires access to PRACE systems, that is the user should:

- have a valid PRACE account
- obtain a X.509 certificate (please refer to PRACE Certificate FAQ for more details)
- access the PRACE infrastructure, as explained in the Interactive Access to HPC resources section of the PRACE User Documentation

# Poll: third question



http://etc.ch/riGa

- DEMO/ hands on:

    − Globus connect (personal and server)

    − Grid Community Toolkit (globus-url-copy)

# Do you want to know more?

- Globus:
  - https://www.slideshare.net/globusonline/introduction-to-globus-for-system-administrators-globusworld-tour-umich-159124701
  - https://www.globus.org/data-transfer
  - https://docs.globus.org/globus-connect-server-installation-guide

# Grid Community Toolkit

- https://www.mcs.anl.gov/~mlink/tutorials

- "Configuring and Deploying GridFTP for Managing Data Movement in Grid/HPC Environments," 21st ACM/IEEE annual SuperComputing Conference (SC 2008) Tutorial, Austin, TX, November 2008. (Slides)
- "Distributed Data Management in Grid Environments," Midwest Grid School 2008 Data Management Tutorial, Chicago, IL, September 2008. (Slides)
- "Managing Data Movement with Globus GridFTP," Open Source Grid and Cluster Conference 2008 GridFTP Tutorial, Oakland, CA, May 2008. (Slides)
- "Configuring and Deploying GridFTP for Managing Data Movement in Grid/HPC Environments," 20th ACM/IEEE annual SuperComputing Conference (SC 2007) Tutorial, Reno, NV, November 2007. (Slides)
- "Optimizing Data Transport: A Tutorial on Deploying GridFTP - From Simple to Advanced Feature Configurations," The 8th LCI (Linux Clusters Institute) Conference on High-Performance Clustered Computing, South Lake Tahoe, CA, May 2007. (Slides)
- "GT4 GridFTP for Administrators: The New GridFTP Server," National eScience Centre (NeSC), Edinburgh, Scotland, January 2005. (Slides)
- "GT4 GridFTP for Users: The New GridFTP Server," National eScience Centre (NeSC) Edinburgh, Scotland, January 2005. (Slides)
- "GT4 GridFTP for Developers: The New GridFTP Server," National eScience Centre (NeSC), Edinburgh, Scotland, January 2005. (Slides)
- "A Tutorial Introduction to High Performance Data Transport," 16th ACM/IEEE annual SuperComputing Conference (SC 2003) Phoenix, AZ, November 2003. (Slides)
- "Data Management using GridFTP," The 7th Global Grid Forum Meeting, Tokyo, Japan, March, 2003. (Slides)

- https://gridcf.org/gct-docs/6.0/gridftp/user/index.html
- https://gridcf.org/gct-docs/6.0/gridftp/admin/index.html

# Globus Toolkit

- https://github.com/globus/globus-toolkit

**Can you describe what this announcement means for each Globus Toolkit component and speak to the recommended migration path for each?**

- GridFTP server: The Globus Connect software distributed with the Globus cloud service provides all of the functionality of Globus Toolkit GridFTP, and a growing set of additional capabilities as well. We recommend transitioning to Globus Connect and the Globus cloud service for data transfer.
- globus-url-copy (GridFTP client): Transition to using the Globus CLI or Python SDK to transfer data via the Globus cloud transfer service.
- GSI: Globus Auth provides a more modern, secure, web-friendly, OAuth2-based security approach than the 1990s-era X.509 of GSI. See https://docs.globus.org/api/auth/ for more information.
- GSI-OpenSSH: We will release Globus Auth-based authentication support for SSH later in 2017. Unlike the current GSI-OpenSSH, this will not require replacing your SSH server, but instead is implemented as a PAM module for use with your existing SSH server. Please contact us for more details if you are eager to use this feature.
- MyProxy: MyProxy provides X.509 certificate management that is used with Globus Connect Server version 4. Globus Connect Server version 5 will use Globus Auth for security and MyProxy will not be needed. We recommend transitioning to Globus Auth when Globus Connect Server version 5 is released..
- GRAM: The GRAM job submission tool is no longer widely used. We recommend performing remote execution via SSH with Globus Auth, once that is available.

**EUDAT** Collaborative Data Infrastructure

**PRACE**



## Whats is B2STAGE?

- **B2STAGE** is a reliable, efficient and easy-to-use **service to transfer research data sets between EUDAT storage resources and high-performance computing (HPC) workspaces.**

- The service allows users to:
  - **transfer large data collections** from EUDAT storage facilities to external HPC facilities for processing
  - **ingest computation results** onto the EUDAT infrastructure
  - **access** stored data sets through **associated PIDs**
  - in conjunction with **B2SAFE, replicate community data sets**, ingesting them onto EUDAT storage resources for long-term preservation
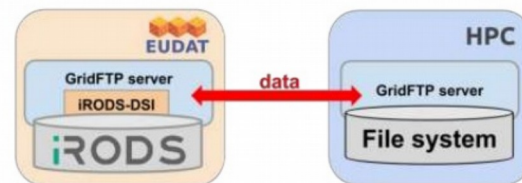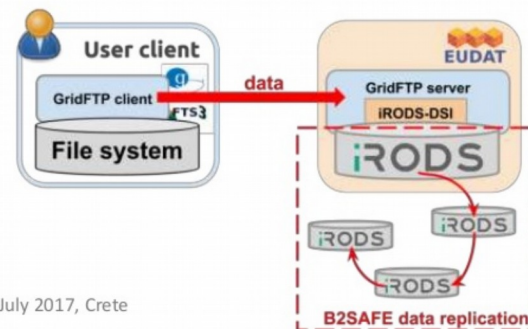
EUDAT Summer School, 3-7 July 2017, Crete

Who can use B2STAGE & Why?

- **Researchers**: can transfer large data collections from EUDAT storage resources to HPC facilities **for processing**.

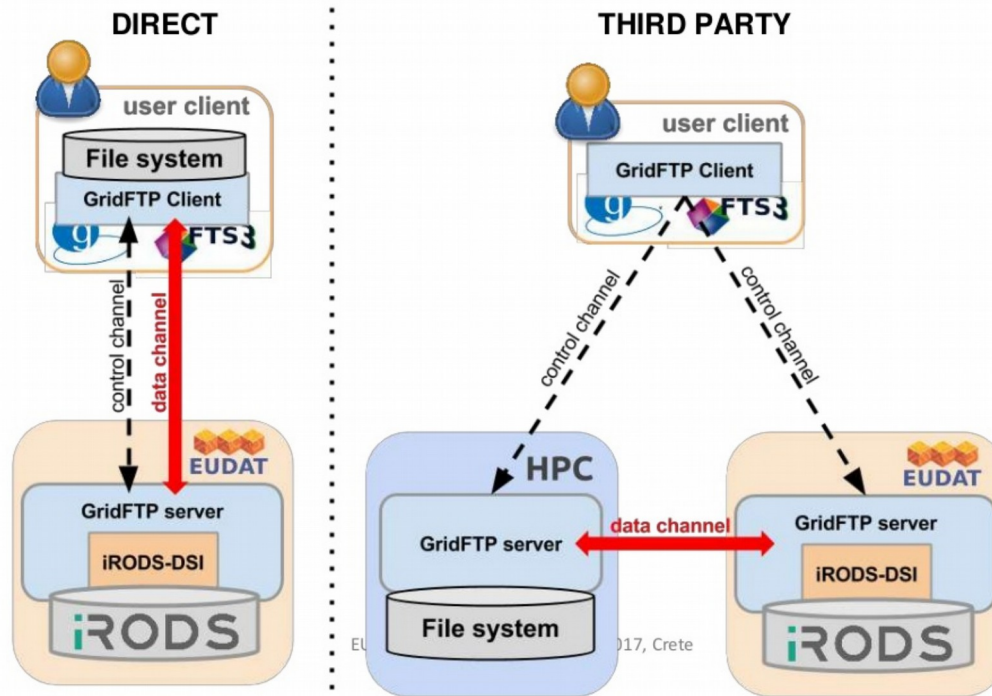- **Community Data Managers**: can ingest data sets onto EUDAT storage resources **for long term preservation** (in combination with the **B2SAFE**).

EUDAT Summer School, 3-7 July 2017, Crete

# Beyond GridFTP

- **FTS**: https://fts.web.cern.ch/#blog

- **MDTM**: https://mdtm.fnal.gov

- **LFTP**: https://lftp.tech

- **jGlobus with AdaptiveGridFTPClient**: https://github.com/earslan58/JGlobus

# Questions?