

Data Staging: Moving large amounts of data around, and moving it close to compute resources

Digital Preservation Advanced Practitioner Course Glasgow, July 19th 2013

c.cacciari@cineca.it







Outline

- Definition
- Starting point
- Moving data around
- Size/Performances
- The user
- Data selection
- Transfer options
- Numbers
- Registered data



Data movement: staging or replication?

[] 11010010







Staging

Safe Replication to enable communities easily create replicas of their scientific datasets in multiple data centres for improving data curation and accessibility data curation and accessibility.
Data Staging to facilitate unities to stage stored data



Data Staging to facilitate communities to stage stored data onto external computational facilities, such as HPC resources





Moving large amounts of data around

111010010**1**



All data are gray in the data staging

1101001010

The data types are irrelevant

Except when they can affect the size or the number of the files



http://www.flickr.com/photos/of

When I should care about data types?

110100101

When you can compress them



Data transfer efficiency is crucial in data staging





Time to copy 1TB

- 10 Mb/s network: 300 hrs (12.5 days)
- 100 Mb/s network: 30 hrs
- 1 Gb/s network: 3 hrs (are your disks fast enough?)
- **10 Gb/s network:** 20 minutes (need *really fast disks and* file system)
- Compare these speeds to:
 - USB 2.0 portable disk
 - 60 MB/sec (480 Mbps) peak
 - 20 MB/sec (160 Mbps) reported on line
 - 15-40 hours to load 1 Terabyte



Data Throughput – Transfer Times

111010010**1**

Bandwidth Requrements to move Y Bytes of data in Time X

Bits per Second Requirements											
10PB	25,020.0 Gbps	3,127.5 Gbps	1,042.5 Gbps	148.9 Gbps	34.7 Gbps						
1PB	2,502.0 Gbps	312.7 Gbps	104.2 Gbps	14.9 Gbps	3.5 Gbps						
100TB	244.3 Gbps	30.5 Gbps	10.2 Gbps	1.5 Gbps	339.4 Mbps						
10TB	24.4 Gbps	3.1 Gbps	1.0 Gbps	145.4 Mbps	33.9 Mbps						
1TB	2.4 Gbps	305.4 Mbps	101.8 Mbps	14.5 Mbps	3.4 Mbps						
100GB	238.6 Mbps	29.8 Mbps	9.9 Mbps	1.4 Mbps	331.4 Kbps						
10GB	23.9 Mbps	3.0 Mbps	994.2 Kbps	142.0 Kbps	33.1 Kbps						
1GB	2.4 Mbps	298.3 Kbps	99.4 Kbps	14.2 Kbps	3.3 Kbps						
100MB	233.0 Kbps	29.1 Kbps	9.7 Kbps	1.4 Kbps	0.3 Kbps						
	1H	8H	24H	7Days	30Days						

This table available at http://fasterdata.es.net



Data staging is Just in Time

111010010**1**

You do not plan it

You can pre-stage... sometimes

There are techniques to improve the efficiency





Pipelining (of commands): speeds up lots of tiny files by stuffing multiple commands into each login session back-to-back without waiting for the first command's response





Parallelizing: on wide-area links, using multiple TCP streams in parallel (even between the same source and destination) can improve aggregate bandwidth over using a single TCP stream







Striping: data may be striped or interleaved across multiple servers

11010010**1**









Parallelism and TCP tuning are the keys

- It is much easier to achieve a given performance level with four parallel connections than with one
- A good TCP tuning can improve drastically performances



Latency interaction is critical

- Wide area data transfers have much higher latency than LAN transfers
- Many tools and protocols assume a LAN
- Examples: SCP/SFTP, HPSS mover protocol



Use the right tool

SCP/SFTP: 10 Mb/s

- standard Unix file copy tools
- fixed 1 MB TCP window in OpenSSH
 - only 64 KB in OpenSSH versions < 4.7</p>

FTP: 400-500 Mb/s

- assumes TCP buffer autotuning
- Parallel stream FTP: 800-900 Mbps



But efficiency is not all

- Easiness of use
- High reliability
- Third-party transfer
- Possibility to control/limit the transfer throughput to avoid engulfing the network

Which priority? Different scenarios, different needs

11010010



... and moving it close to compute resources









Who will move the data?

- Any user, part of a community or citizen scientist.
- In most part of the cases, a domain expert, such as a biologist, a linguist, a seismologist, not an expert of data transfer software.





As Ethernet speeds have increased, there has been a widening gap between ability of the novice to fully use bandwidth capabilities, compared to that of network capabilites achieved by users who have had their network tuned by an expert ('Wizard").



[Diagram c/o Matt Mathis, http://www.psc.edu/%7Emathis/papers/] http://www.minifigure.org/wp-content/uploads/2011/05/wizard_s.jpg



Data selection





Data Staging Overview

How the data sets are selected by the domain experts?

111010010**1**

Through their **local tools**, where "local" means on the community side

Through **remote services**, generic or specific per community







Move the data







EUDAT





Globus OnLine Service

11010010

globusonline.org https://www.globusonline.org/xfer/StartTransference	🖙 🖝 🖉 🚼 🕶 Google					
attiva 🔻 👤 Cookie 👻 🎢 CSS 👻 🗮 Moduli 👻 💼 Immagini 👻 🕦 Informaz	zioni 👻 🚨 Va	arie 👻 🥜 Contorna 🤋	🔻 🖪 Ridimensiona 🔻	🔆 Strumenti 🔻	🖹 Visualizza sorgente 🔻	A Opzie
🕒 globus online 🛛 🕫	o To: Start Tra	nsfer 💌		mcarpene	Sign Out	
Transfer Files - source overwrites files on destination				View	Transfer Activity	
Endpoint mcarpene#pdl 🔍 Go		End	point mcarpene#PLX		Go	
Path /~/ Go			Path /~///		Go	
select all none 🌪 up one folder 💍 refresh list 📑	¢-	select all none	L up one folder 🔥 ref	resh list	* •	
Documenti	Folder	asdata			Folder	
GSI-SSHTerm_IGE_for_PRACE_DGRID_LRZ-v1.3.2	Folder	🚞 cineca			Folder	
Immagini	Folder	prod			Folder	
Modelli	Folder	user_test			Folder	
Musica	Folder	auseragip			Folder	
Scaricati	Folder	userbmwor			Folder	
Scrivania	Folder	usercorsi			Folder	
Ubuntu One	Folder	auserdeisa			Folder	
Video	Folder	userdompe			Folder	
VirtualBox VMs	Folder	userexternal			Folder	
globusconnect-1.4	Folder	userfercfd			Folder	
rpmbuild	Folder	userforfait			Folder	
workspace	Folder	userfranc			Folder	
GSI-SSHTerm_IGE_for_PRACE_DGRID_LRZ-v1.3.2.zip	7.28MB	usergrant			Folder	
examples.desktop	179b	userhpe			Folder	
getskype-linux-beta-ubuntu-64	22.5MB	userhyper			Folder	
globusconnect-latest.tgz	7.91MB	userinaf			Folder	
skype_2.2.0.35-0natty1_amd64.deb	22.49MB	userincm			Folder	
		userinternal			Folder	
		userjrc			Folder 🥃	
			Get Globus Connor	t		
Laber mis Transfer			Turn your computer in	to an ondepaint. The ea	reject and	
This will be displayed in your transfer activi	ity.		most convenient way t machine.	o send and receive file:	s on your	

http://www.globusonline.org













Local staging targets







110100101

• Sample Results:

– RTT = 53 ms, network capacity = 10Gb/s.

- Tool Throughput
 - scp: 140 Mb/s
 - HPN patched scp: 1.2 Gb/s
 - FTP: 1.4 Gb/s
 - GridFTP, 4 streams 5.4 Gb/s
 - GridFTP, 8 streams 6.6 Gb/s





iRODS

- iRODS is a data management system, which integrates a **rule engine**
- The rules can be triggered automatically based on specific events (for example a data set is moved to a particular directory)
- Invoked from remote via iRODS command line client or integrated into applications based on iRODS java (Jargon) and python (PyRods) API





Metadata!





"Managing Large Datasets with iRODS — a Performance Analysis" - Proceedings of the International Multiconference on Computer Science and Information Technology pp. 647–654



Transfer Time vs File Size (Tar-gzipped Collection)





https://openwiki.uninett.no/norstore:irods:workshop2012

Data Staging and registered data







Conclusions

- The way to move data has to be enough **flexible** to accomodate different transfer protocols, different access mechanisms.
- Flexibility means also that the transfer tools can be used as they are, with default parameters, for average performances, but also **fine tuned** by experts **for faster transfers**.
- No solution fits all, so different services are provided



Contact Us





Project Coordinator: Kimmo Koski kimmo.koski@csc.fi

Scientific Coordinator: Peter Wittenburg peter.wittenburg@mpi.nl

Project Manager: Damien Lecarpentier damien.lecarpentier@csc.fi

Dissemination Manager: Nagham Salman nagham.salman@bsc.es

> Industry Task Force: David Manset dmanset@maatg.fr



37



Thank you!

