# Networking and the Grid

Ahmed Abdelrahim

NeSC

PPARC e-Science Summer School

10th May 2005
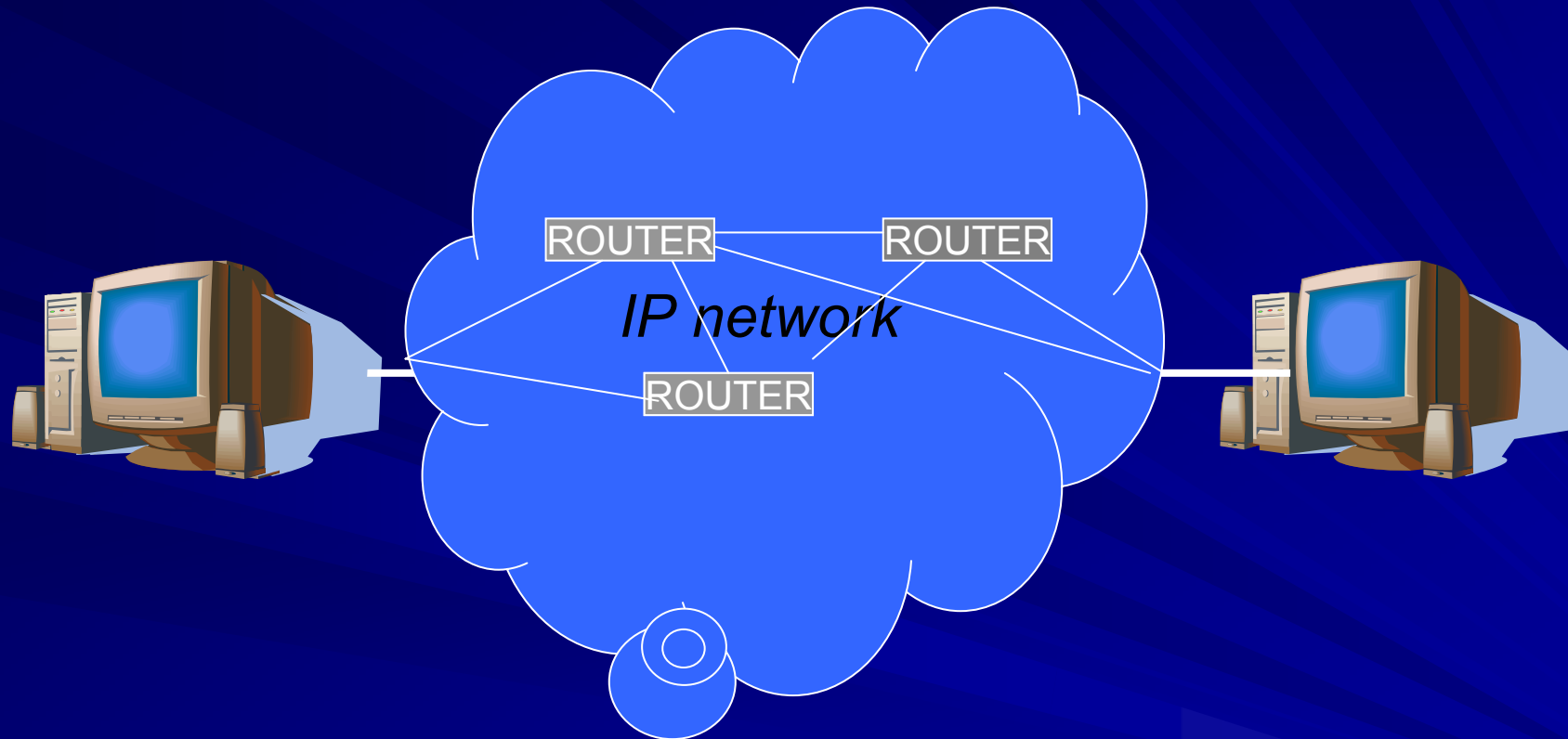
# Overview

- Project Area
- Networking
- High Performance Transfer and Problems
- EGEE JRA4 involvement
- Future Directions

# Project Area

- The Grid is interconnected through networks.
- The performance of a network affects the performance/operation of a Grid.
- How does it affect it ?
- Need to measure network performance metrics and correlate them with grid performance metrics.

# Networking



IP – Internet Protocol used to transfer packets from source to destination.

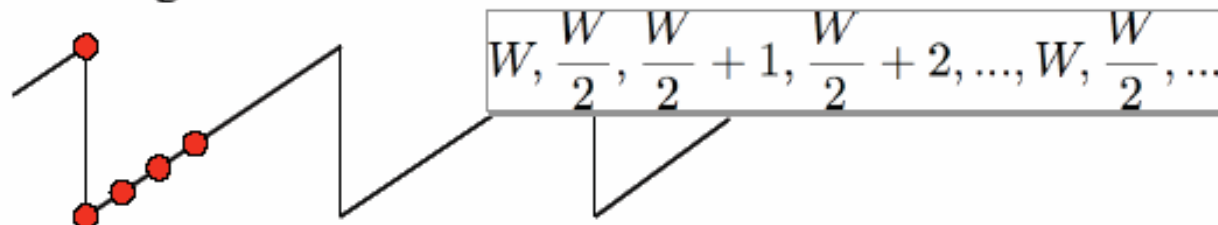Problems :Packets can get lost , delayed at router queues , arrive out of order!

# TCP

- TCP – reliable Transmission Control Protocol.
- Standardised in 1981
- Transfers majority of internet traffic – email, web ,ftp, ssh, telnet and GridFTP- all built upon TCP.
- Acknowledgement and sliding window based protocol
- Allows ordering of packets at receiver
- Retransmits lost packets.

# TCP Modelling: The "Steady State" Model

**The model:** Packet size $B$ bytes, round-trip time $R$ secs, no queue.

- A packet is dropped each time the window reaches W packets.
- TCP's congestion window:

$$W, \frac{W}{2}, \frac{W}{2}+1, \frac{W}{2}+2, ..., W, \frac{W}{2}, ...$$

- The maximum sending rate in packets per roundtrip time: $W$
- The maximum sending rate in bytes/sec: $W B / R$
- The average sending rate $T$:  $T = (3/4)W B / R$

- The packet drop rate $p$:  $p = \dfrac{1}{\frac{3}{8}W^2}$

- **The result:**  $T = \dfrac{\sqrt{6}B}{2R\sqrt{p}} = \dfrac{\sqrt{3/2}B}{R\sqrt{p}}$

Slide from NFNN  Workshop University College London July 2004
from presentation by Professor Mark Handley.

# TCP

Simple model  for TCP "Steady State Model"
 Average sending rate
- proportional to packet size
- Inversely proportional to packet loss
- Inversely proportional to Round Trip Time

# TCP

- TCP was designed for LANs (and not WANs)
- Get performance problems when using standard TCP in WANs because have high Round Trip Time (RTT).

# High Performance Transfers.

- Example transferring data from NeSC to CERN.
- Data is going to go across a number of networks. (in any single international path between two end points, at least five networks are involved:

Two university or campus networks

Two national backbones (the participating NRENs)

The European backbone (GÉANT or GÉANT2).

- Can have problems on any of these networks !

# What will affect my Transfer ?

- End host specification

Network interface card. (could be bottleneck)

Motherboard, CPU, memory, Disk Drive

- So need a high specification PC.

# What will affect my Transfer ?

- End host TCP tuning
- Need to calculate BDP (Bandwidth delay product) to determine adequate TCP buffer sizes at sender and receiver.
- EXAMPLE
- On a link where OC3 155Mbps and 100 BT Ethernet used throughout and RTT = 50 ms (obtained by ping)– require
- = 50 ms  x  (100Mbps / 8 bits)  = 625 KB

- (Note most OS default TCP buffers size is 24KB or 32 KB and Linux is only 8KB) (this is ok for LANs but not WANs)
- Using untuned buffers you often get less than 5% utilisation !

*(source:* TCP tuning guide for distributed application on wide area networks , Brian L. Tierney*)*

# What will affect my Transfer ?

- **Version of TCP used**

  - Can get newer versions of TCP which are suitable for High performance transfer.

  - Ongoing research – Fairness of these versions in terms of affect on other users.

# How can I monitor my end-host TCP parameters?

- Web100

Provides information on TCP parameters

# What else will affect my Transfer ?

- Performance of the networks my data traverses
- Delays
- Packet loss
- Jitter (IPDV)

# Network transfer problems
## (In order of most probable cause)

- End host (OS, architecture, disk, TCP)
- End host application itself
- Local network limitation (switches)
- Firewall
- And only then the WAN

- Understanding each area will help you to reduce the Wizard Gap (difference between transfer rates a network expert and normal user achieve).

# OK – how do I obtain information about the network?

■ Need Network Monitoring

- Active monitoring e.g., iperf, udpmon, PingER
  - (can put load on to the network though)
- Passive monitoring. Passive monitoring tools capture packets using either a standard NIC (network interface card) and libpcap library or using a specialized hardware monitoring adapters.
- (privacy issues)

# Network monitoring other networks?

- Each network has its own monitoring infrastructure so would be ideal to share networking information.

- Need standard interface to share information.

# Network performance Characteristics for Grids (METRICS) - GGF NMWG



Figure 3: Network characteristics that can be used to describe the behaviour of network entities.

GGF have designed an NMWG Schema

# EGEE JRA4 - NPM Aims

- **JRA4/NPM provides uniform access to network performance information from a heterogeneous set of monitoring frameworks**

# NPM Architecture

The current state of Network Monitoring Point deployment is shown below, with deployments at JRA4 partners (University of Edinburgh, CNRS, GARR and DANTE) and on the JRA1 testbed (at CERN and RAL).

Secondary producer – at CERN    Access requires a user certificate

RTT measured by PingER

# Demo of NPM prototype

```
Welcome to the Network Monitoring Prototype

Using proxy in: /tmp/x509up_u517
Please create your Network Measurement Request by selecting options
from the following menus

Select the route type:
1. Source/destination
2. Hop list
Choice?
```

- The demo client states what type of route it is, but the user need not be aware of the different network monitoring points contacted

selected route 1: GARR to Marseille

```
Select the source/destination route you want to examine:

1. GARR (it) to Marseille (fr) (endsite)

2. GARR (it) to MapCenter (fr) (endsite)

31. Nordunet Stockholm ttm08 (se) to GEANT Geneva ttm107 (ch)
(backbone)

32. SWITCH Zurich ttm85 (ch) to GEANT Geneva ttm107 (ch) (backbone)

Choice? 1
```

```
Select the characteristic you wish to measure/examine:

1. OWD - One Way Delay

2. IP (layer 3) available bandwidth

3. RTT

4. Packet Loss (round trip)

5. TCP achievable throughput

6. UDP achievable throughput

7. UDP jitter

Choice? 3
```

```
Select the query type for the characteristic:

1. last value three hours ago

2. all data (singletons) between 9am and 11am yesterday

3. average between 9am and 11am yesterday

4. min and max over yesterday

5. daily average over the last week

Choice? 3
```

- The client generates an appropriate NM-WG Request and outputs a summary to the screen
- It then sends the request to the appropriate NM-WG compliant network monitoring point (web service)

```
You selected to obtain RTT measurements (average between 9am and
11am yesterday) on Source/destination: GARR (it) to Marseille (fr)
(endsite)

Your selections resulted in a query for 1 result (mean over the
whole period) for path.delay.roundTrip for the route GARR (it)-
Marseille (fr) for measurements between Sun Apr 17 09:00:00 BST
2005 and Sun Apr 17 11:00:00 BST 2005
```

- Once the result is received, it is output to the user
- Note that the execution time is long – NPM team will be resolving this issue within the next development cycle

```
Request execution in progress ...

...complete (13574ms)


Displaying results for path.delay.roundTrip for route GARR (it)-
Marseille (fr)

Displaying results for mean:

   17 April 2005 09:58:48 BST 54.24155555555557ms
```

END OF NPM DEMO

# Future Project/Research Directions

- Work on Network diagnostics
- Correlation of Grid and Network performance metrics.
- Do practical tests e.g. submit jobs. Vs network performance metrics.
- Important to understand grid traffic requirements. (for planning, SLAs etc.)
- Continued involvement with NPM

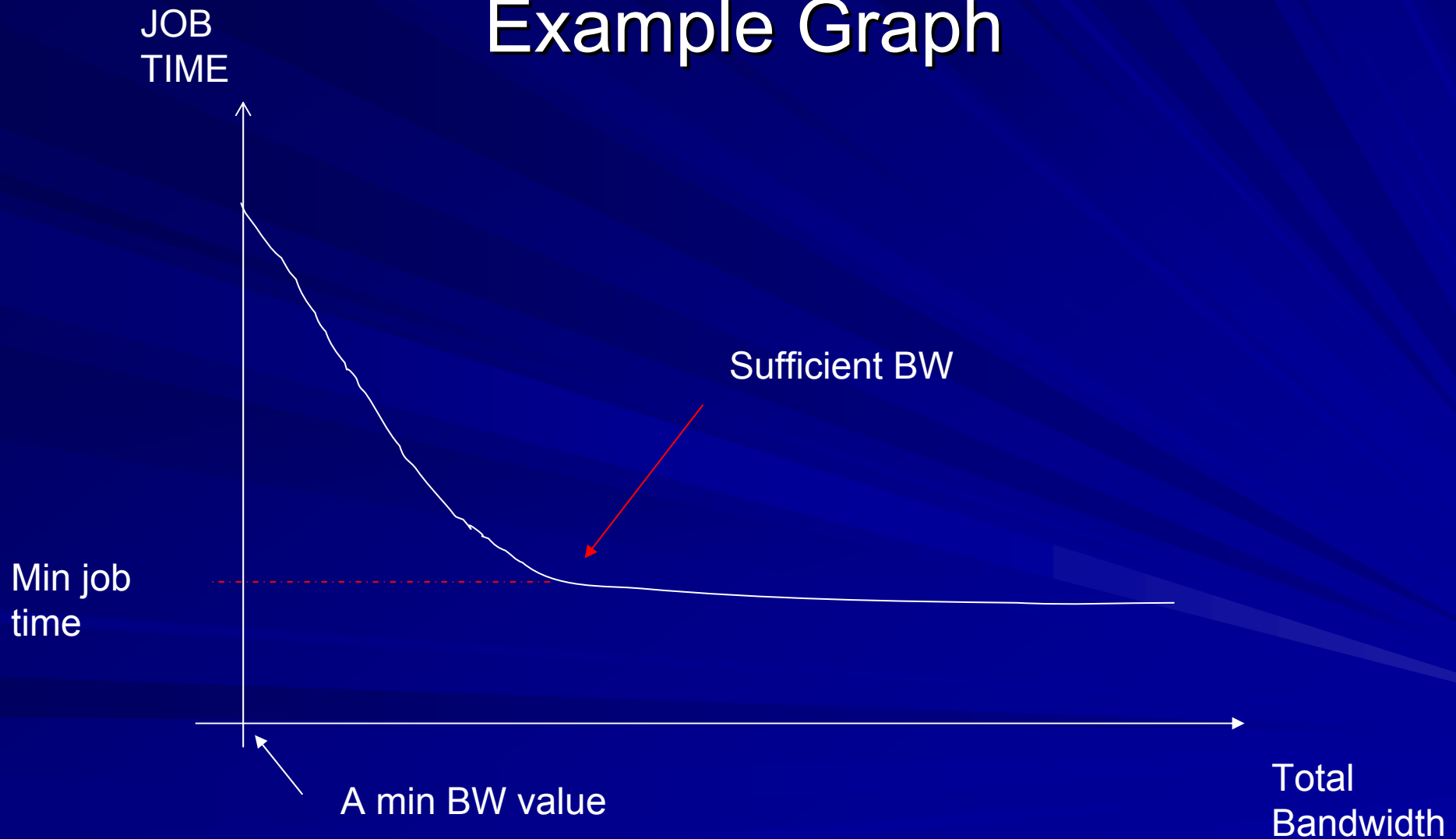# Future Project/Research Directions

- Use of networking information

Network cost functions

e.g. Replica Placement

WMS – where to submit jobs

# Contribution to NPM

1$^{st}$ Prototype
 – assisted with development

2$^{nd}$ prototype
–deployed monitoring node at NeSC.
 (Installed R-GMA server, client and WP7 tools)

Next Prototype (due September)
 – Will Work on Diagnostic tool functional specification.
- Maintain deployment node.

# Resources

• Networks for non-Networkers workshop UCL July 2004

-http://grid.ucl.ac.uk/NFNN_Programme.html


• TCP tuning guide for distributed application on wide area networks, Brian L. Tierney

-http://dsd.lbl.gov/TCP-tuning/tcp-wan-perf.pdf

THE END
and
Question Time!