

Qian Zhang

LSU Department of Computer Science  
Fall 2010 Final Exam  
CSC7700 Scientific Computing  
December 6th 2010, 5.30pm to 7.30pm

**General Instructions**

- This is a closed book exam.
- No calculators or electronic devices.
- Part I of the exam covers all the five course modules and is designed to take 80 minutes to complete. Part II of the exam is for the Networks and Data module and is designed to take 40 minutes to complete.
- Part I is worth 20% of the final grade. Each module includes 5 questions. All questions have equal weight. Answer all questions.
- Part II is worth 10% of your final grade. Answer only four out of five questions. If you answer all five, only the lowest graded four will be taken into consideration. Questions have two parts, you need to answer both parts of the four questions you select.

## Part I

## Module A: Basic Skills

1. Provide two reasons why the same text file can look different when viewed on different systems or within different tools.

Different editors have different rules of display.

2. In the context of numerical simulations, explain what is meant by discretization and why it is used.

Discretization is to discretize the continuous range into discrete points/parts. For numerical simulations, they are mostly PDE and to solve PDE, we should discretize to approximate continuous part.

3. Briefly describe what a pseudo random-number generator is, and name three disadvantages over real random-number generators. Name two reasons why pseudo random-number generators are often used despite these disadvantages?

PRNG is an algorithm to use a sequence of numbers to approximate the properties of random numbers.

Disadvantages: ① not necessarily random ② necessarily periodic and shorter than expected ③ not robust

Because ① some numerical problems is not grid-based.

② IC & BC are not obtainable even can't generate a evolution eqn  $PDE$ .

4. Name one advantage and two potential disadvantages of the Newton-Raphson method over the bisection method for root-finding.

advantage: reliable, accurate

Disadvantage: too slow, sometimes may not find the root.

5. Explain the difference between centralized and distributed version control systems, including one advantage and one disadvantage for each. Name one software implementation example for each kind of system.

version control is the management of changes on document, programs in computer. update - export - commit

Centralized version control

take a longer time

stable

distributed version control

get update notification in time

may cause conflicts

## Module B: Networks and Data

1. List two TCP parameters used in iperf and briefly describe their influence on the performance of TCP.

-P : number of streams paralleled, parallel  $\uparrow \rightarrow$  bandwidth  $\uparrow$   
but reach a maximum when  $p=8$

-W : window size, window size larger  $\rightarrow$  bandwidth  $\uparrow$

2. Briefly describe what the server-side data processing plug-in included in the standard GridFTP installation does and what it can be used for (hint - you used it in your homework)

GridFTP is a file transfer protocol on Grid forum.

Gridftp is parallel, striping, restartable transfer, partial file transfer, support a third-party transfer and support other protocols besides TCP.

Globus-url-copying  
transfer large file.

3. List two benefits that middleware provides to developers of distributed applications.

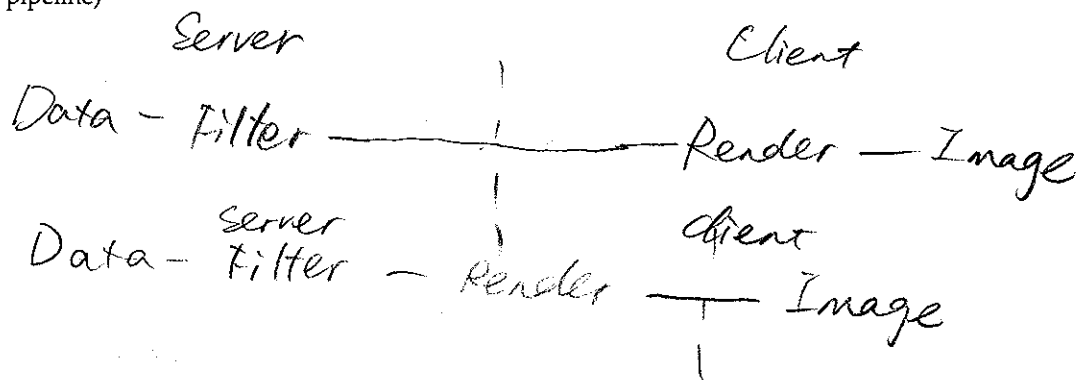
Middleware is a software between OS and application.

4. Briefly outline two methods for accessing remote data in a distributed application.

OpenSSH: most popular; work on all platforms; feature-rich

CopSSH: only for 'microsoft'!

5. Briefly outline two methods of doing remote visualization (based on distribution of the visualization pipeline)



## Module C: Simulations and Application Frameworks

1. What determines the accuracy of a simulation? List two ways in which accuracy can be improved.

Resolution

Decrease discretized unit ( $dx$ )

(Crank-Nicholson)

Choose higher-order algorithms: eg from 2nd order to

4th order accuracy  
(RK4)

2. What is MPI, and what is it used for? Assume there are two processes, and process A needs to access an array element stored on process B. Schematically, how does this work?

MPI is an API, Message passing interface, work on all HPC systems  
MPI: assign the copy of the program to different processors, each processor has a unique name, each processor work on the program independently, only communicate when exchanging message, MPI hides low-level system-computer message. MPI is reliable and ordered.  
processor A send a message to B, B detect and receive the message.  
Then B send array to A, A received. MPI Send, MPI Rec

3. What is a software framework? Name one software framework, and provide three characteristic elements of a software framework.

framework provides glue between components

Cactus is a framework (flesh), the components are thorns.

4. What are CCL files in Cactus? List which CCL files exist, and what they define.

CCL is configuration files

activated/synchronized

schedule.ccl : When flesh schedule which functions; When which variable is freed/

param.ccl : use of variables of other thorn.

.ccl : implementation name; global functions; relation between thorns;  
thorn variables

5. Name and briefly describe five tools that support code development in large, distributed, international collaborations.

Reliable: Code is right

expensible: Others can add or change the code

scalable =

performance: can achieve some function

maintainable:



Qia

## Module D. Scientific Visualization

1. Define and describe a "Visualization Pipeline".

Data  $\rightarrow$  Filter  $\rightarrow$  Geometry  $\rightarrow$  Rendering  $\rightarrow$  Image

2. What is the difference between the "push model" and the "pull model"?

PUSH Model

data available as soon as possible  
Traverse viz by loading  
Filter has information <sup>at early stage</sup>  
about data  
Load data even not used

PULL Model

Data available as late as possible  
Traverse viz at rendering time  
Filter has no information about data  
until output request  
Load only when used

3. Describe the three atomic elements ("building blocks") in a visualization network.

Vish + Fiber bundle = fish

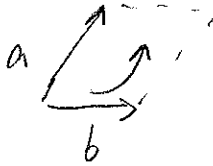
Ocean is kernel

Structure: Ocean, fish, quish, module, tutorial

4. Define and describe the purpose of a bi-vector.

$a \wedge b = \text{vector.}$

outer product



5. Which are the three property objects ("communication types") in the "F5" fiber bundle data model that are visible to the end user?

Group

Dataset

Data element

## Module E: Distributed Scientific Computing

1. We discussed five applications – Montage, Nektar, Climateprediction.net, SCOOP and Ensemble-based/Replica-Exchange simulations. For any THREE of these (you choose which three), answer any ONE of the following: Why they were distributed? How they were distributed? The Challenges &/or success in distributing them?

	Why distributed	How distributed	Challenge
Montage	scale processing > local limits	DAG created & ② DAGMAN executed	map to right resources
Ensemble-based/Replica Exchange	① many uncoupled units ② many resources	① many existing implementations ② SAGA Based "Pilot-job" to use on distributed TG.	① get SAGA to work on all machines ② get the best resources ③ coordinate across resources
SCOOP	① distributed / naturally ② require simulation faster than real world	customized workflow	① Not robust ② co-scheduling ③ coordinate resources across

2. Estimate to within an order of magnitude the number of jobs that are executed in the Worldwide LHC Computing Grid (WLCG) per day. Estimate to within an order of magnitude the number of bytes of data generated (overall) by the WLCG. Estimate the cost of the LHC Experiment. Therefore what is the cost of generating a byte of data from the LHC experiment?

100M jobs/day

10 TBytes

\$10M

$$\frac{\$10^7}{10^{12}} = \$10^{-5}$$

3. Using your estimate (whatever it was) of number of jobs (on the WLCG) from the previous answer, given that there are approximately 250,000 cores as part of the WLCG, and that it has a typical utilization factor of 50%, estimate the average time each job takes. (assume: each job is a single-core job).

$$\frac{2.5 \times 10^5}{10^8} = \frac{2.5}{10^3} = \frac{1}{400}$$

$$\frac{1}{400} \times \frac{1}{2} = \frac{1}{800} \text{ day} = \frac{3600 \times 24}{800} = 108 \text{ s}$$

4. List two factors – technological or non-technological, driving Cloud Computing. Provide a “real production” example of a Cloud offering. Is the Cloud offering an example of IaaS, PaaS or SaaS?

The space of DA is large but the effective number is small

Developing DA is difficult.

Embracing DA

Yes

5. Provide one difference between predominantly HTC and HPC Grids. Provide a “real production” example of a HPC and HTC Grid.

number of jobs

HPC = OCG

HTC = EGI

## Part II

## Networks and Data

### Question 1

- A) How are layers used in network implementations?

- ① Physical layer
- ② Data link layer
- ③ Network
- ④ Transport
- ⑤ Session
- ⑥ Presentation
- ⑦ Application

- B) What are the major differences between TCP and UDP?

TCP	UDP
ordered	unordered
reliable	unreliable
congestion control	No congestion control
connection-based	connectionless
byte-based	packet-based

## Question 2

- A) What data transmission protocol would you use for bulk data transmission and why? What protocol would you use for video or audio conference and why?

Gtftp TP: ① Parallel ② striping ③ Restartable transfer ④ Partial file transfer ⑤ Third-Party transfer ⑥ support other protocols besides TCP.

UTC: too much delay makes noninteractive  
compression increase delay but decrease bandwidth

- B) Describe circuit network services and their advantage.

### Question 3

- A) Describe what a naming service is (in middleware implementations) and what is it used for.
- B) In your own words, describe the "end-to-end" argument.



Qia

#### Question 4

- A) List the usual sequence of operations for accessing data in a distributed file system.

Distribute file system: Home directory, data directory  
and scratch directory

- B) Briefly describe the two possible (and sometimes conflicting) optimization goals of a scheduling system.

### Question 5

- A) Describe use case scenarios where remote visualization is useful or needed.

- ① local visualization is not powerful enough
- ② Data copy is impossible at local.

- B) Describe some of the possible benefits of distributed visualization.

1. Don't need to copy data