# MIDDLE TENNESSEE

STATE UNIVERSITY

## INTRODUCTION

- Supercomputing facilitates and accelerates scientific discovery
- BIG DATA: from T to P, to E, to Z, to Y, and beyond...
- Simulation, Experimental, Observational
- Collaborative research and discovery
- ✤ A team of scientists with complementary expertise working jointly
- Massively distributed resources
  - Hardware Computing facilities, Storage system, i.e., Cloud, Special rendering engines, Display devices (tiled display, powerwall, etc.)
- Software Domain-specific analysis/processing tools/programs...
- Data Data canter in cloud
- Data-, computation-, and network-intensive scientific applications pose new challenges
- An increasing trend in many big data scientific applications to move computing workflow executions to a cloud environment
- A three-layer architecture for performance modeling and optimization of big data scientific workflows in distributed environments



### **PROBLEMS AND OBJECTIVES**

- Problem 1
  - Given a cloud environment modeled as a faulty heterogeneous computer network  $G_{cn} = (V_{cn}, E_{cn})$ , where each node or link is associated with an independent failure rate, and a bound *F* for the overall failure rate (OFR);
- Each computer node is DVFS enabled;
- A user submits a request of a streaming application modeled as a DAG-structured computing workflow  $G_w = (V_w, E_w)$ , and a bound **FR** on the frame rate (FR) or throughput.
- <u>Objective</u>: To find a workflow mapping schedule such that the total energy consumption (TEC) is minimized and overall performance satisfies the F and FR constraints.

$$F_{OFR} \leq F$$

$$\frac{1}{BT} \geq FR$$

• Problem 2

Given a computing workflow, a cloud model with different types of VMs, and a cost objective function.

$\sim$	processing nower (instructions/second)	VM Type	MIPS	Disk BW	Cost/hr
0	processing power (instructions/second)	T1	1000	10	\$0.10
0	storage access bandwidth (bytes/second	T2	2000	100	\$0.30
0	charge rate (\$/second)	T3	3000	10	\$0.30
		T4	4000	100	\$0.60

Objective: To assign workflow tasks to VM instances such that the total execution time and total cost are minimized.

## Resource and Performance Optimization of Big Data Scientific Workflows in **Distributed Network Environments**

Yi Gu, Dept. of Computer Science, Middle Tennessee State University









