

Infrastructure, Data Cleansing and Mining for Scientific Simulations

Committee Members:

Dr. Bowyer
Dr. Flynn
Dr. Madey
Dr. Uhran

Yingping Huang

Agenda

- ◆ **Overview**

- ◆ Background

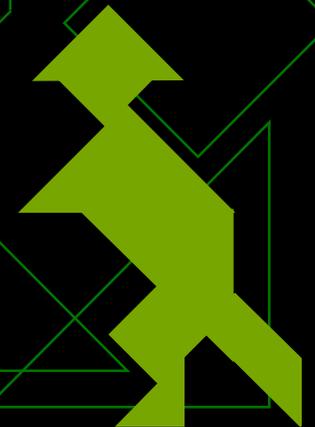
- ◆ Multi-tier infrastructure

- ◆ Data cleansing algorithms

- ◆ Data mining applications

- ◆ Summarize

- ◆ Timeframe

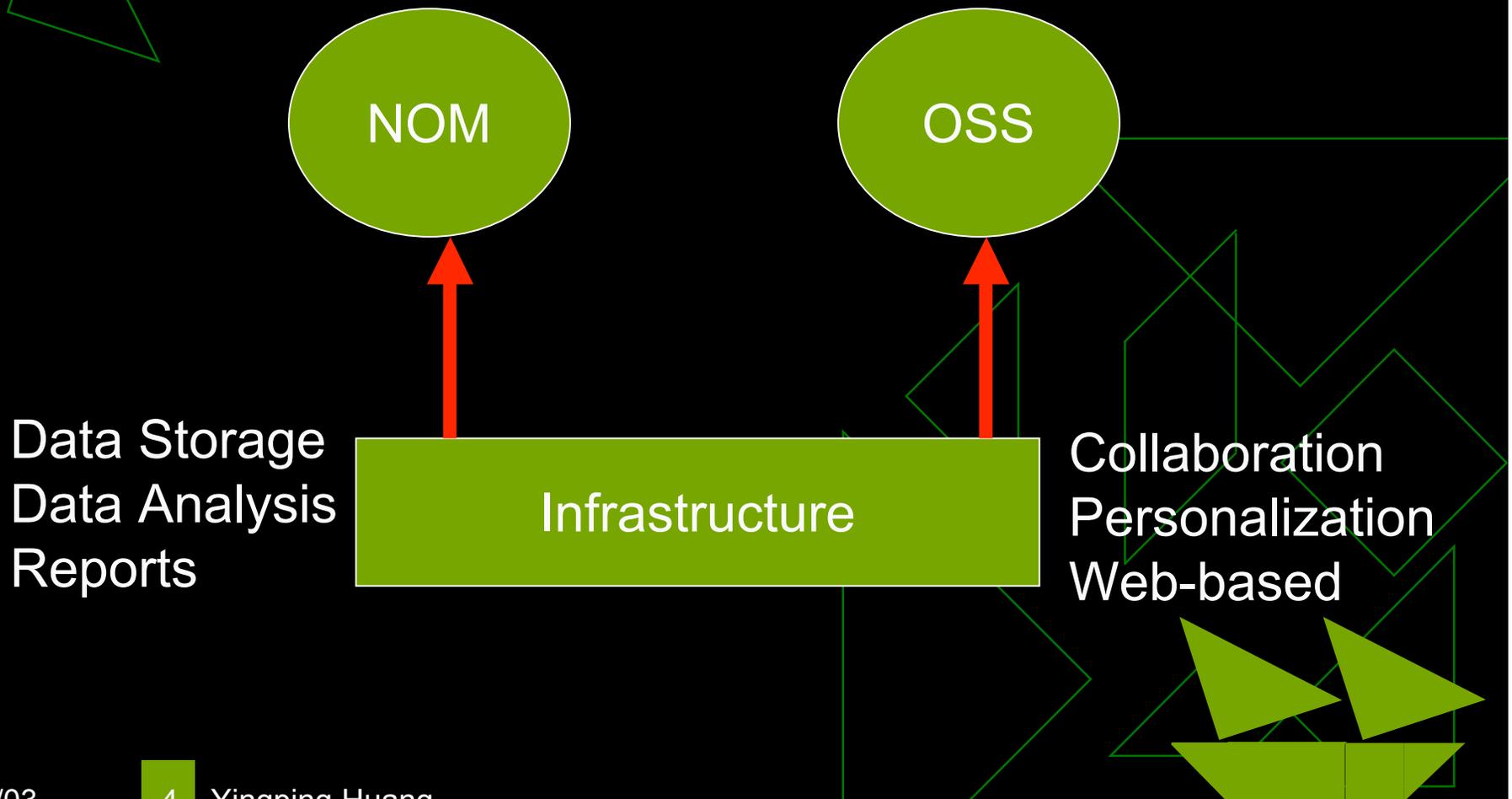


Overview

- ◆ Multi-tier infrastructure powers scientific simulations.
- ◆ Data cleansing algorithms result in better data quality.
- ◆ Data mining applications discover hidden knowledge in environmental and social science.

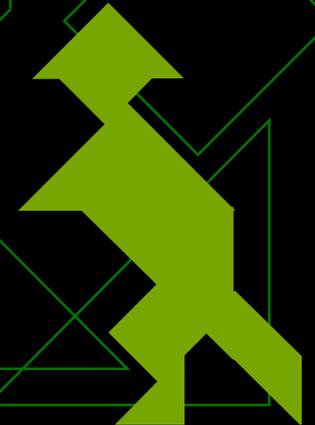
Motivation

Simulation Anytime and Anywhere



Agenda

- ◆ Overview
- ◆ **Background**
- ◆ Multi-tier infrastructure
- ◆ Data cleansing algorithms
- ◆ Data mining applications
- ◆ Summarize
- ◆ Timeframe



Background

◆ Projects under way

■ NOM

- ◆ Research on natural organic matter (NOM)
- ◆ Study evolution of NOM over time
- ◆ Joint work of scientists across disciplines including chemists, biochemists, environmental scientists

■ OSS

- ◆ Research on the open source software (OSS) development phenomenon
- ◆ Study the behavior of OSS developers and their motivations
- ◆ Joint work with social scientists

Simulation Models

- ◆ Standalone or traditional client-server
 - Software needs to be installed on clients
 - Incompatibility makes installation difficult
- ◆ Web-based using applets
 - Security – file permission, firewall
 - Inconvenience – plug-ins download
 - Network traffic – download before executing
 - Incompatibility – Swarm
- ◆ What should be done?
 - Web-based server-side simulation models
 - Centralized simulation management
 - Collaboration and personalization

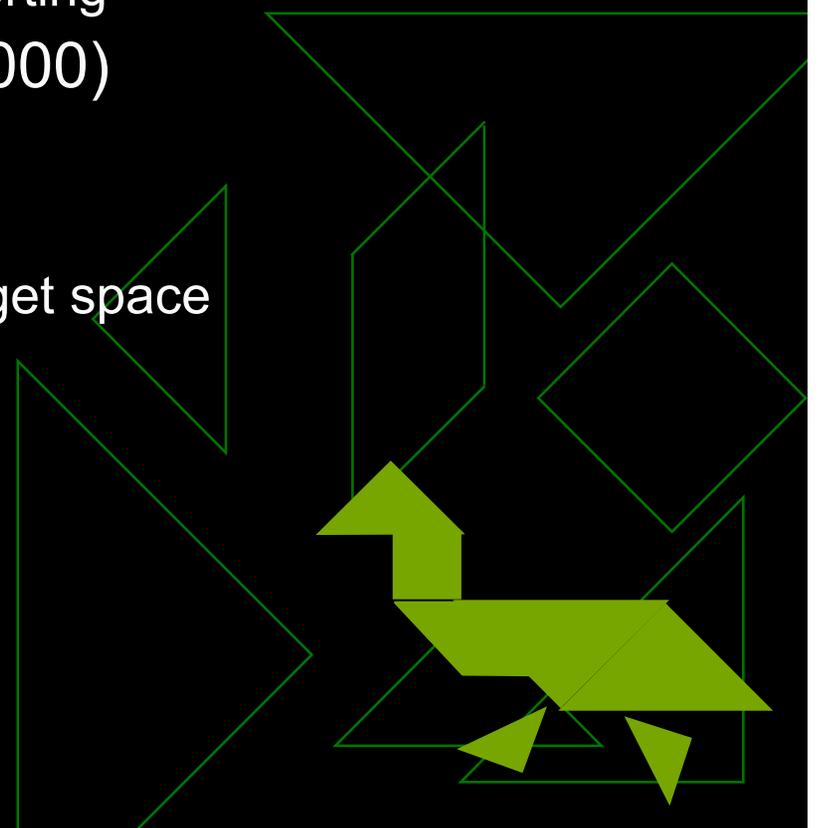
Data Cleansing

◆ Known approaches

- Sorted neighborhood (Stolfo 1995/1998)
 - ◆ Domain dependent keys for sorting
- Record matching (Monge, 2000)
 - ◆ Edit distance only
- String mapping (Li, 2003)
 - ◆ Potential high dimensional target space

◆ Our approaches

- Sample database
- Lipschitz mapping



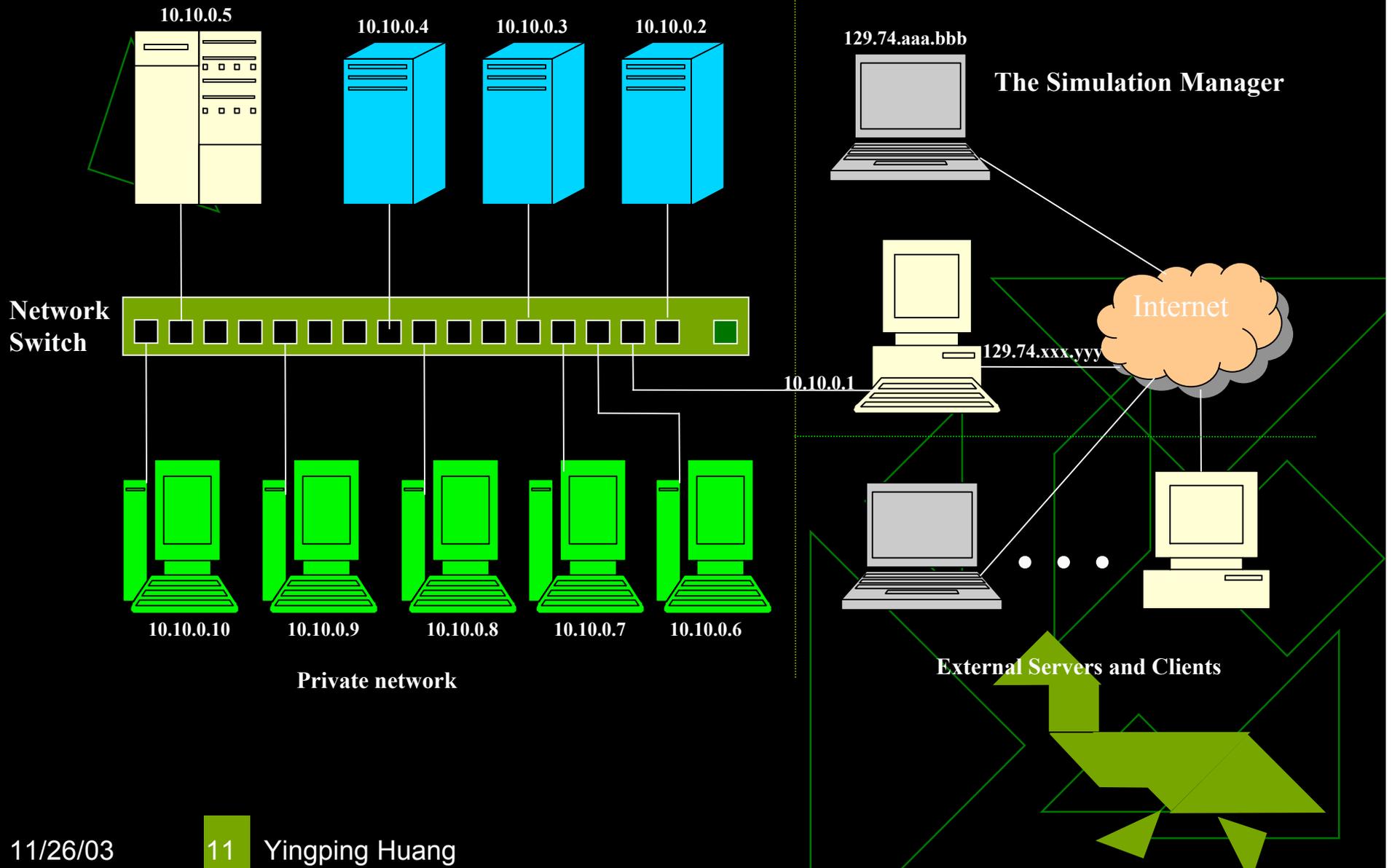
Data Mining

- ◆ Data mining in astronomy
 - SKYCAT: star/galaxy classification (Fayyad, 1996)
 - JARTool: detect volcanoes on Venus (Burl, 1998)
 - Sapphire: find galaxies (Kamath, 2001)
- ◆ Data mining in biology
 - Bioinformatics
 - SARS diagnosis (ehealth.org)
- ◆ What should be done?
 - Data mining for social science (OSS)
 - Data mining for environmental science (NOM)
 - Add intelligence to simulation models by applying data mining results

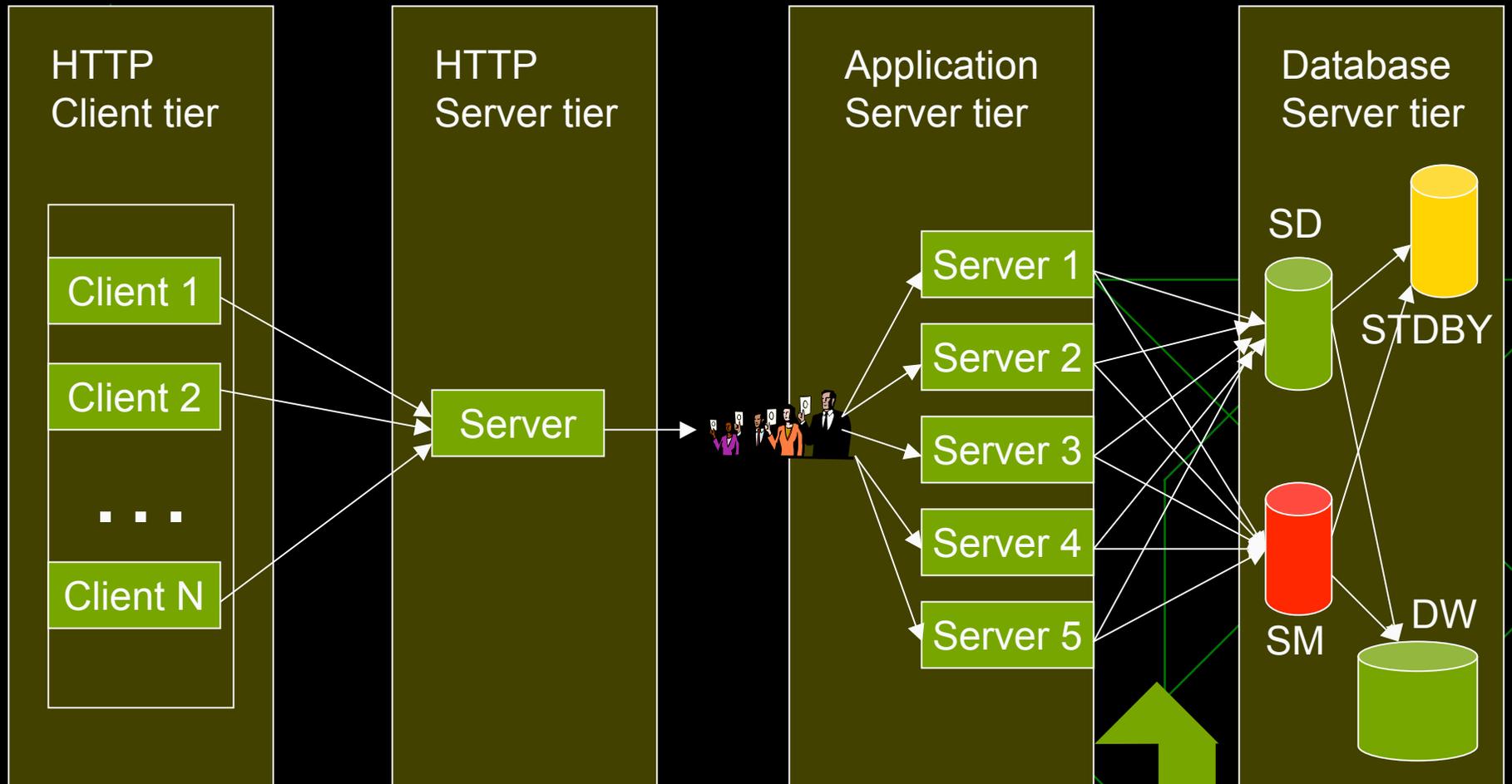
Agenda

- ◆ Overview
- ◆ Background
- ◆ **Multi-tier infrastructure**
- ◆ Data cleansing algorithms
- ◆ Data mining applications
- ◆ Summarize
- ◆ Timeframe

Physical Layout



Multi-tier Architecture



Two Features

- ◆ Load-balancing
 - Scalability achieved
 - Implementation using JMS, AQ & EJB
 - Implementation using Shell scripts & PL/SQL
- ◆ Simulation-resuming
 - Reliability achieved
 - Checkpoint
 - Implementation using JTA/JTS

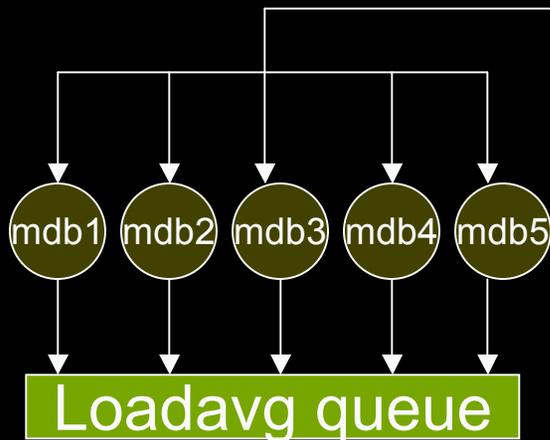
Load-balancing Using JMS & AQ



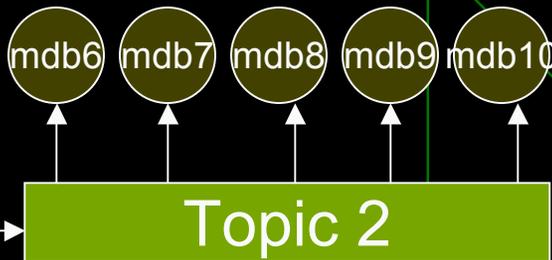
job_id	resumed	checkpoint	status
100	0	0	

Job queue

Topic 1



Invoke simulation & update job queue



Judge bean

Shell Scripts & PL/SQL

- ◆ Dispatcher (HTTP server)
 - Dispatch simulations
 - Send KEEPALIVE messages to running simulations
- ◆ Intelligent agent (application server)
 - Upload load averages
 - Check simulations
 - Send ACK to KEEPALIVE messages

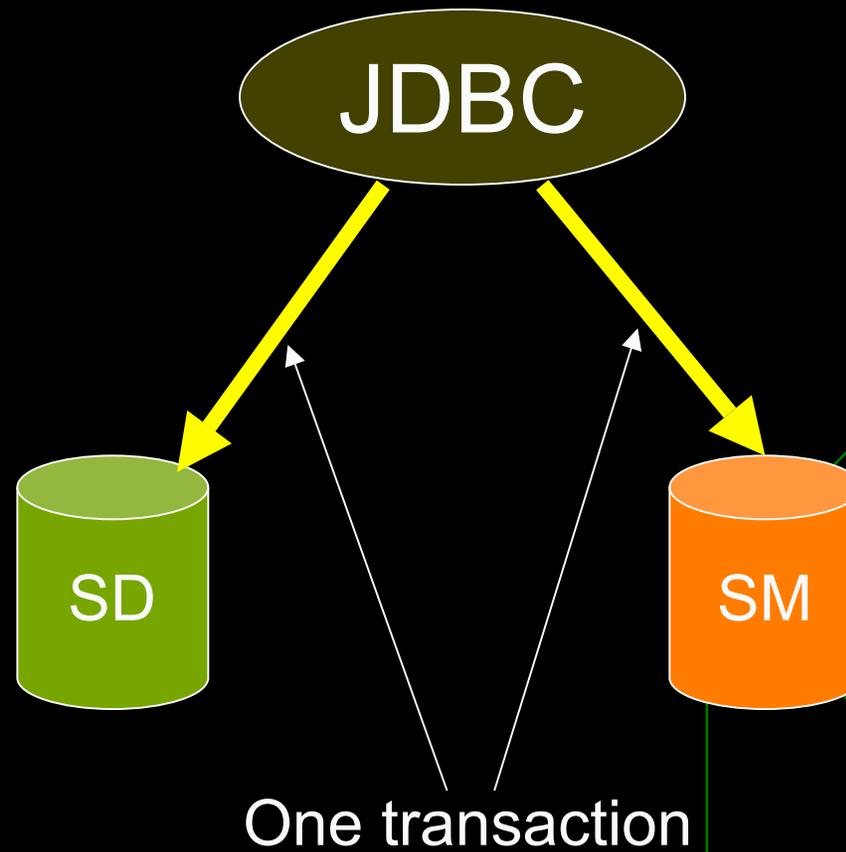
Load-balancing Algorithm

- ◆ Instance learning approach
 - Based on completion time prediction
- ◆ Two step completion time prediction
 - Completion time estimation
 - ◆ Load average
 - ◆ Data amount
 - Completion time prediction
 - ◆ Nearest neighborhood

Completion Time Estimation

- ◆ Completion time estimation formula

Checkpoint



JTA/JTS

Checkpoint Issues

- ◆ Checkpoint data
 - All data for restarting the simulation
 - Size depends on number of agents
- ◆ Checkpoint frequency
 - Checkpoint-interval
 - ◆ # of MB data
 - Checkpoint-timeout
 - ◆ # of minutes

Simulation-resuming

- ◆ To restart a terminated simulation
 - A new simulation with same job_id inserted into the job queue
 - A terminated simulation has smaller job_id than new simulations, higher priority
- ◆ In case of application server failure
 - All simulations' job_ids inserted into the job queue
 - All simulations will be running on other application servers

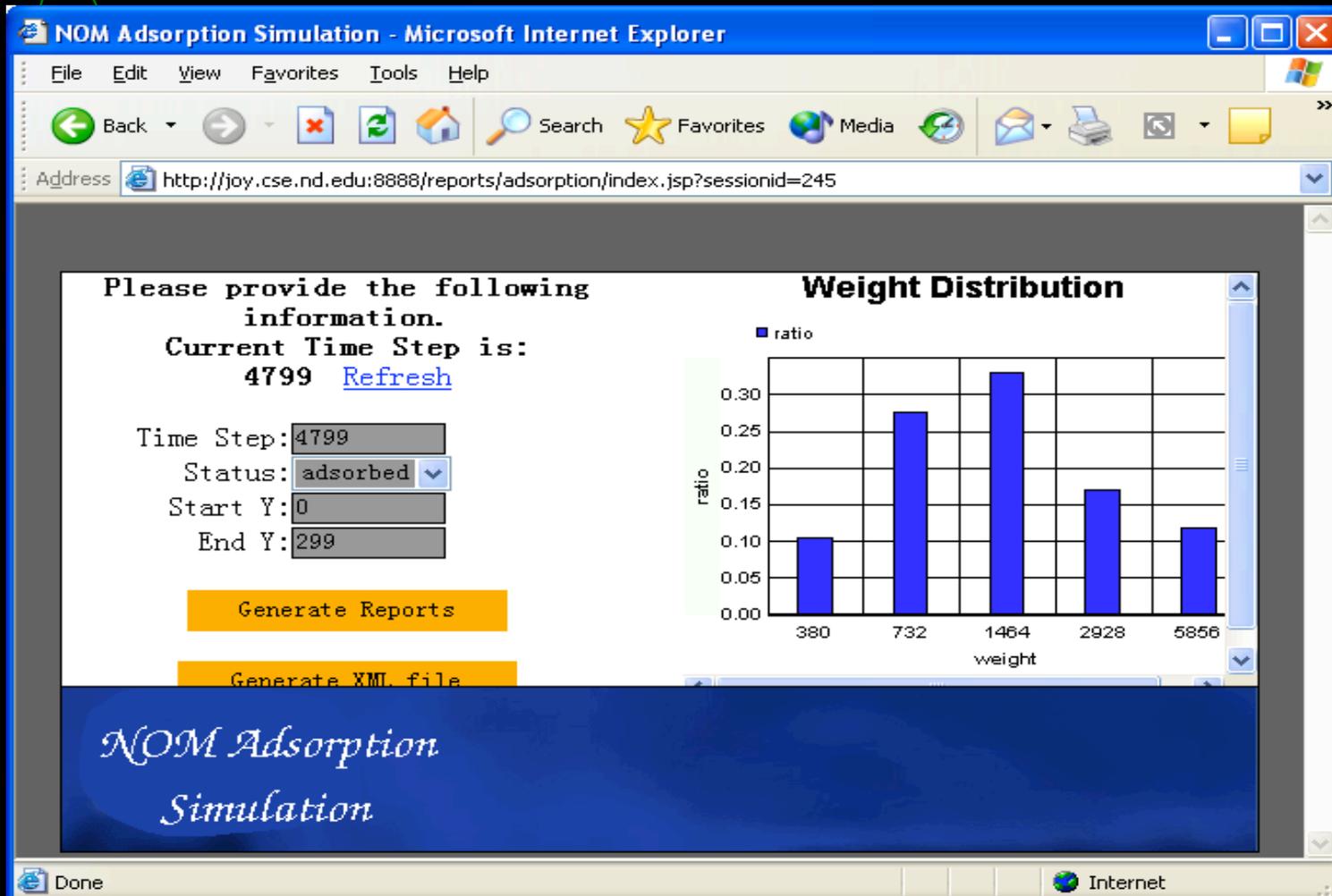
Collaboration Suite

The screenshot shows a Microsoft Internet Explorer browser window with the address bar set to <http://tobit.cse.nd.edu>. The page content includes:

- NOM** logo with a thought bubble icon.
- NOM DISCUSSION BOARD** and **NOM CHAT ROOM** buttons.
- NOM BBS** section with a globe icon and text: "NOM researchers can discuss modeling, programming, reporting and so on."
- NOM SIMULATOR** button.
- NOM Simulator** section with a gear icon and text: "NOM researchers can start simulations and view information such as reports for their simulations."
- Registration form** for the chat room with fields for "Your Name:", "Password:", "Member" (radio button), "Not Member" (radio button), "Register" link, and "Visitor:" field, followed by a "go!" button.

The browser interface includes a menu bar (File, Edit, View, Favorites, Tools, Help), a toolbar with navigation buttons (Back, Forward, Stop, Refresh, Home, Search, Favorites, Media), and a status bar at the bottom showing "Done" and "Internet".

Graphical Reports



XML Reports

Please provide the following information.
Current Time Step is: 4799 [Refresh](#)

Time Step:
Status:
Start Y:
End Y:

Weight Distribution

weight	count	total	ratio
380	8	76	0.1053
732	21	76	0.2763
1464	25	76	0.3289
2928	13	76	0.1711
5856	9	76	0.1184

NOM Adsorption Simulation

Agenda

- ◆ Overview
- ◆ Background
- ◆ Multi-tier information system
- ◆ **Data mining applications**
- ◆ Summarize
- ◆ Timeframe

Methodology

- ◆ Traditional approach
 - Form hypotheses
 - Verify hypotheses by finding patterns in data
- ◆ Data mining approach
 - Find patterns in data
 - Form hypotheses
 - Design simulation models
 - Verify hypotheses
- ◆ U. Fayyad, J. Gray at Microsoft Research

Technology & Software

◆ Data mining technology

- Clustering
 - ◆ K-means
 - ◆ Orthogonal cluster
- Classification
 - ◆ Decision tree
 - ◆ Naïve Bayes
- Association rules
 - ◆ Apriori

◆ Data mining software

- Oracle Data Mining Suite
- DM4J
- JDeveloper

OSS

- ◆ Study behavior of open source software (OSS) developers
 - Agent-based
 - Stochastic
- ◆ Data mining involving
 - Clustering
 - Classification
 - ◆ Churn prediction
 - ◆ Acquisition prediction
 - Association rules

OSS Data Warehousing

- ◆ Data from sourceforge.com
 - Developers
 - Projects
- ◆ Data warehousing
 - Table partitioning
 - Aggregation
 - Star schema
 - Analysis SQL
 - ETL tools → Warehouse Builder

NOM

- ◆ Study behavior of natural organic matter (NOM)
 - Agent-based
 - Stochastic
- ◆ Data mining involving
 - Clustering
 - ◆ Micelle formation
 - Classification
 - ◆ Transportation prediction
 - ◆ Adsorption prediction
 - Association rules

Agenda

- ◆ Overview
- ◆ Background
- ◆ Multi-tier information system
- ◆ Data mining applications
- ◆ **Summarize**
- ◆ Timeframe

Summarize

- ◆ Multi-tier information system integrates
 - Application servers & reports server
 - Database servers
 - Data warehousing & data mining
 - Swarm
- ◆ Collaboration suite
- ◆ Data mining guided model-design

Insights & Impacts

- ◆ Server-side simulation models
 - Centralized simulation management
 - Centralized data repository
- ◆ Collaboration suite
 - Simulation sharing
 - Knowledge sharing
- ◆ Data mining applications
 - Find patterns in data
 - Model deployment for simulation-design

Agenda

- ◆ Overview
- ◆ Background
- ◆ Multi-tier information system
- ◆ Data mining applications
- ◆ Summarize
- ◆ **Timeframe**

Timeframe

May 2003 ~ May 2004



Implement infrastructure

Data collection & statistical analysis

Data mining model design

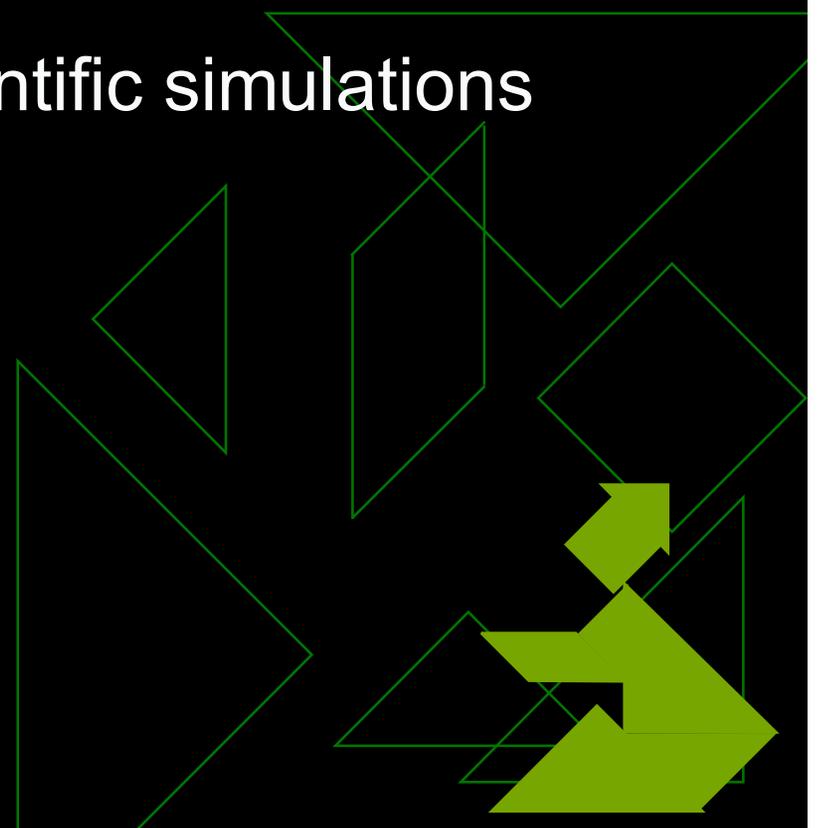
Data mining model evaluation

Deployment

Writing up

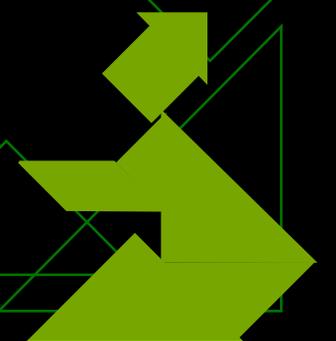
Expected Publications

- ◆ Information system design for scientific simulations
 - By August 2003
- ◆ Data warehousing for scientific simulations
 - By November 2003
- ◆ Data mining for OSS
 - By February 2004
- ◆ Data mining for NOM
 - By March 2004



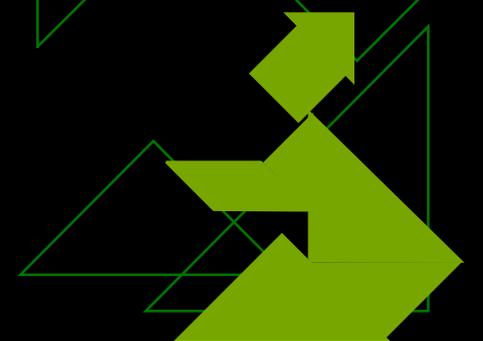
Demo

Demonstration



Finally

Thank you!



Features

- ◆ Multi-tier information system
 - HTTP client tier → HTTP server tier → Application server tier → EIS tier
- ◆ Scalability at the application server tier
 - Load-balancing
- ◆ Reliability at the application server tier
 - Simulation-resuming
- ◆ Reliability at the database tier
 - Standby databases

Features (cont.)

- ◆ Data mining models
 - Stored in database
 - Stored Java procedures
 - PL/SQL procedure call using JDBC
- ◆ Simulation models
 - Agent-based
 - Stochastic
 - Data mining guided