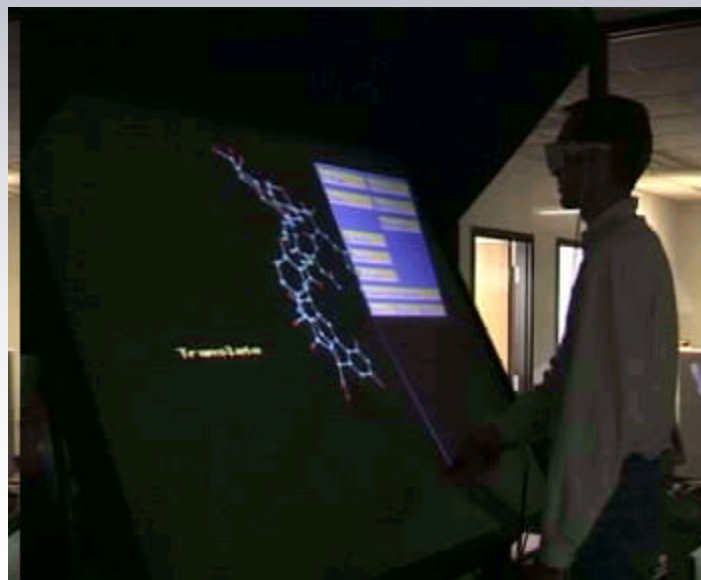# Molecular Structure Determination on a Computational & Data Grid

## Mark Green & Russ Miller

**Center for Computational Research, SUNY-Buffalo**

**Computer Science & Engineering, SUNY-Buffalo**

**Hauptman-Woodward Medical Inst**

**University at Buffalo**

*The State University of New York*

# Biomedical Advances

- **PSA Test (screen for Prostate Cancer)**
- **Avonex: Interferon Treatment for Multiple Sclerosis**
- **Artificial Blood**
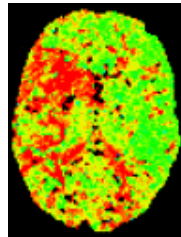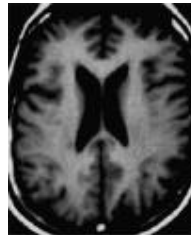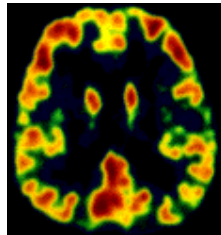- **Nicorette Gum**
- **Fetal Viability Test**
- **Implantable Pacemaker**
- **Edible Vaccine for Hepatitis C**
- **Timed-Release Insulin Therapy**
- **Anti-Arrythmia Therapy**
  - **Tarantula venom**

- **Direct Methods Structure Determination**
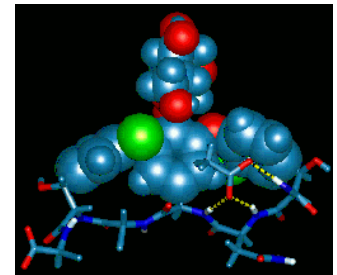  - **Listed on "Top Ten Algorithms of the 20th Century"**
  - **Vancomycin**
  - **Gramacidin A**
- **High Throughput Crystallization Method: Patented**
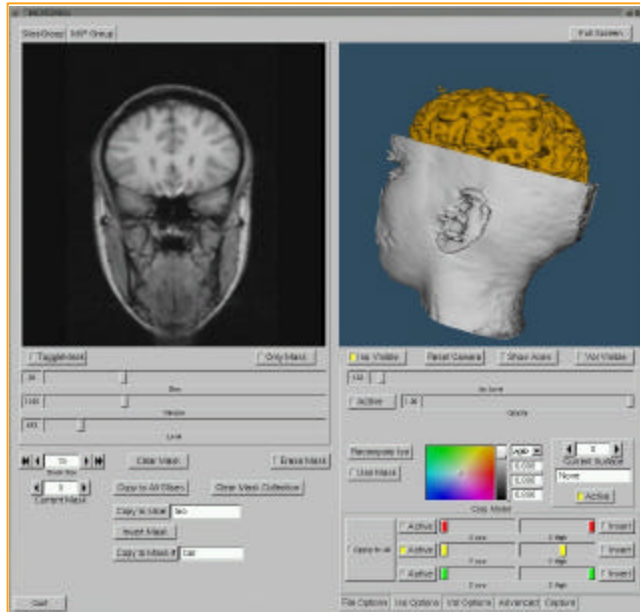- **NIH National Genomics Center: Northeast Consortium**
- **Howard Hughes Medical Institute: Center for Genomics & Proteomics**

# Bioinformatics in Buffalo
## A $360M Initiative

- **New York State: $121M**
- **Federal Appropriations: $13M**
- **Corporate: $146**
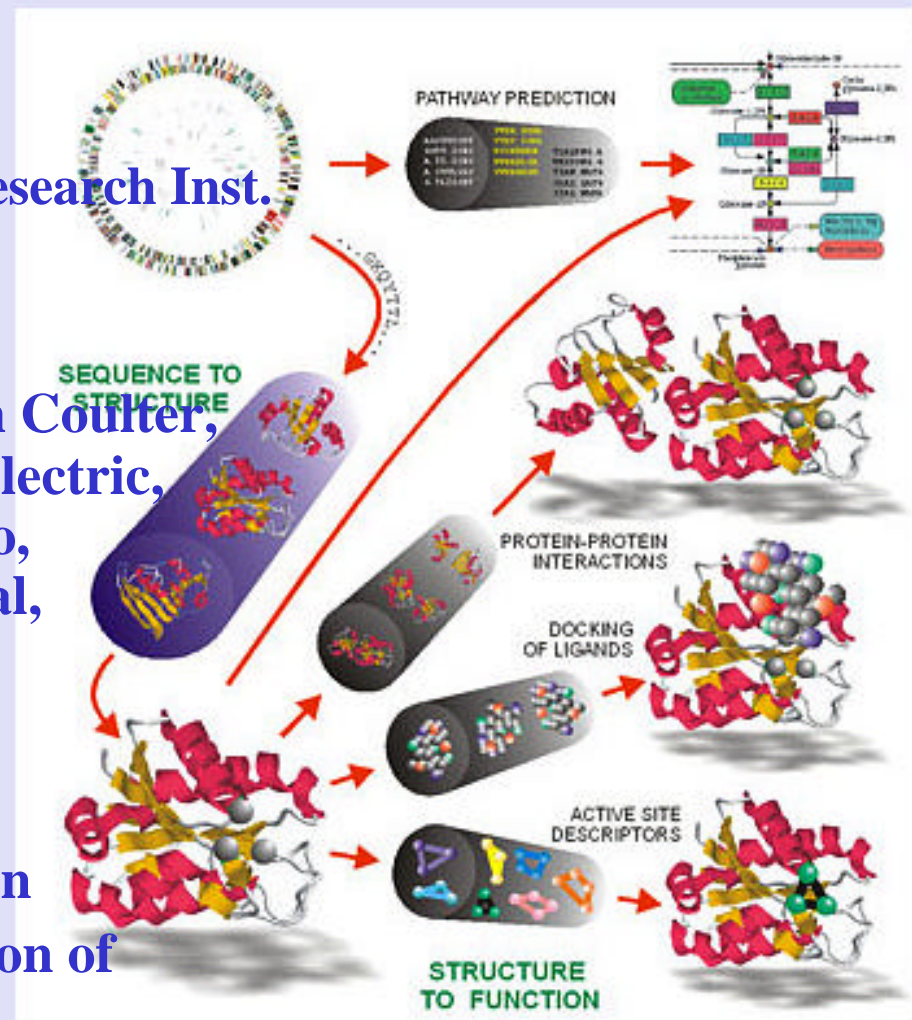- **Foundation: $15M**
- **Grants & Contracts: $64M**

# Bioinformatics Partners

- **Lead Institutions**
  - ❑ **University at Buffalo (UB)**
  - ❑ **Hauptman-Woodward Medical Research Inst.**
  - ❑ **Roswell Park Cancer Institute**
- **Corporate Partners**
  - ❑ **Amersham Pharmacia, Beckman Coulter, Bristol Myers Squibb, General Electric, Human Genome Sciences, Immco, Invitrogen, Pfizer Pharmaceutical, Wyeth Lederle, Zeptometrix**
  - ❑ **Dell, HP, SGI, Stryker, Sun**
  - ❑ **AT&T, Sloan Foundation**
  - ❑ **InforMax, Q-Chem, 3M, Veridian**
  - ❑ **BioPharma Ireland, Confederation of Indian Industries**

# Center for Computational Research 1999-2004 Snapshot

- **High-Performance Computing and High-End Visualization**
  - ❑ **110 Research Groups in 27 Depts**
  - ❑ **13 Local Companies**
  - ❑ **10 Local Institutions**
- **External Funding**
  - ❑ **$111M External Funding**
    - ❍ **$13.5M as lead**
    - ❍ **$97.5M in support**
  - ❑ **$41.8M Vendor Donations**
  - ❑ **$360M Bioinformatics Initiative**
- **Deliverables**
  - ❑ **350+ Publications**
  - ❑ **Software, Media, Algorithms, Consulting, Training, CPU Cycles…**

# Major CCR Resources (12TF & 80TB)

- **Dell Linux Cluster: #22 ® #25 ® #38**
  - ❑ **600 P4 Processors (2.4 GHz)**
  - ❑ **600 GB RAM; 40 TB Disk; Myrinet**
- **Dell Linux Cluster: #187 ® #368 ® off**
  - ❑ **4036 Processors (PIII 1.2 GHz)**
  - ❑ **2TB RAM; 160TB Disk; 16TB SN**
  - ❑ **Restricted Use (Skolnick)**

- **IBM BladeCenter Cluster**
  - ❑ **532 P4 Processors (2.8 GHz)**
  - ❑ **5TB SAN**
- **Apex Bioinformatics System**
  - ❑ **Sun V880 (3), 6800, 280R (2), PIIIs**
  - ❑ **Sun 3960: 7 TB Disk Storage**
- **HP/Compaq SAN**
  - ❑ **75 TB Disk; 190 TB Tape SGI Origin3800**
  - ❑ **64 Alpha Processors (400 MHz)**
  - ❑ **32 GB RAM; 400 GB Disk**
- **IBM RS/6000 SP: 78 Processor**
- **Sun Cluster: 80 Processors**
- **SGI Intel Linux Cluster**
  - ❑ **150 PIII Processors (1 GHz)**
  - ❑ **Myrinet**

# CCR Visualization Resources

- **Fakespace ImmersaDesk R2**
  - ❑ **Portable 3D Device**
- **Tiled-Display Wall**
  - ❑ **20 NEC projectors: 15.7M pixels**
  - ❑ **Screen is 11'´7'**
  - ❑ **Dell PCs with Myrinet2000**
- **Access Grid Node**
  - ❑ **Group-to-Group Communication**
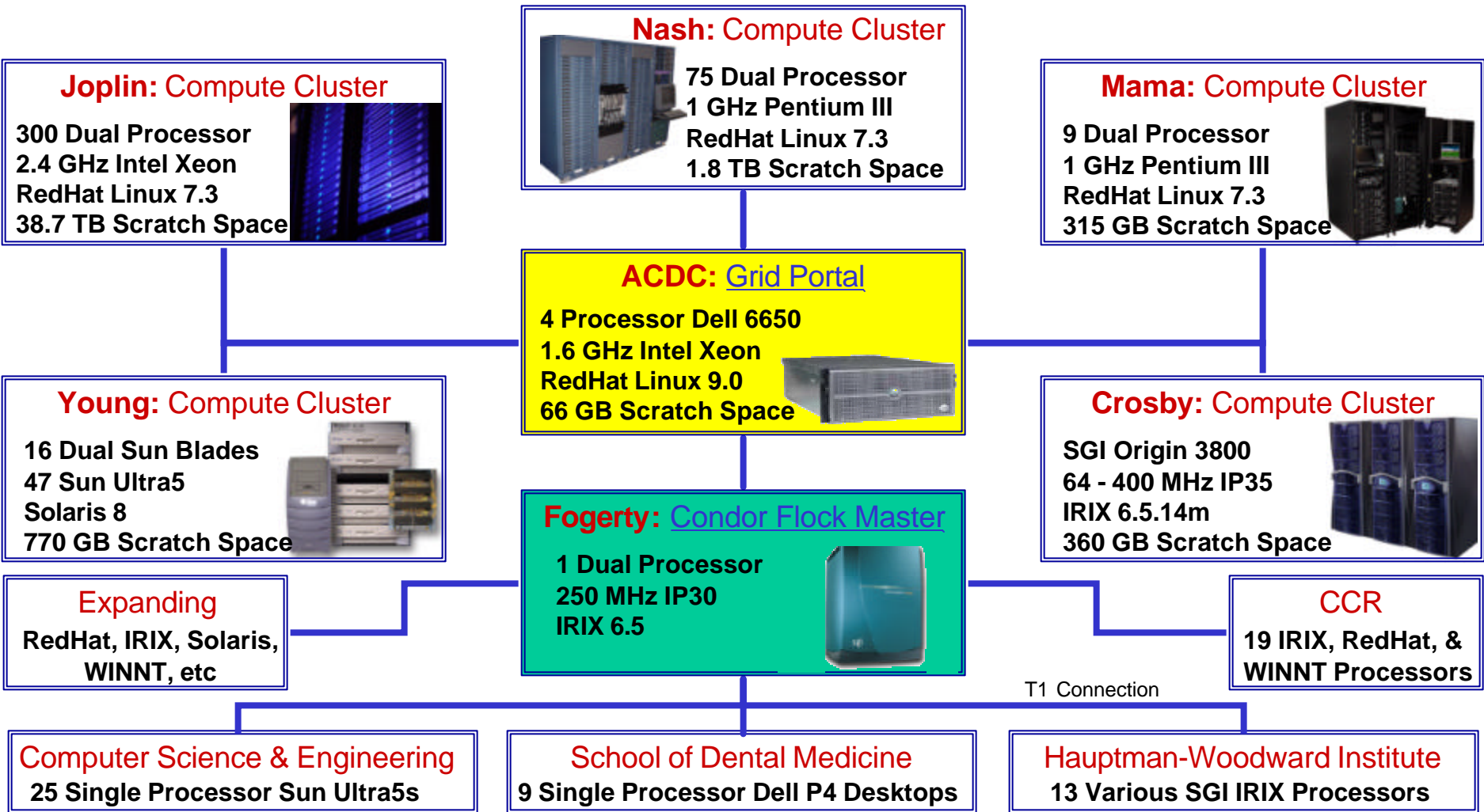  - ❑ **Commodity components**
- **SGI Reality Center 3300W**
  - ❑ **Dual Barco's on 8'´4' screen**
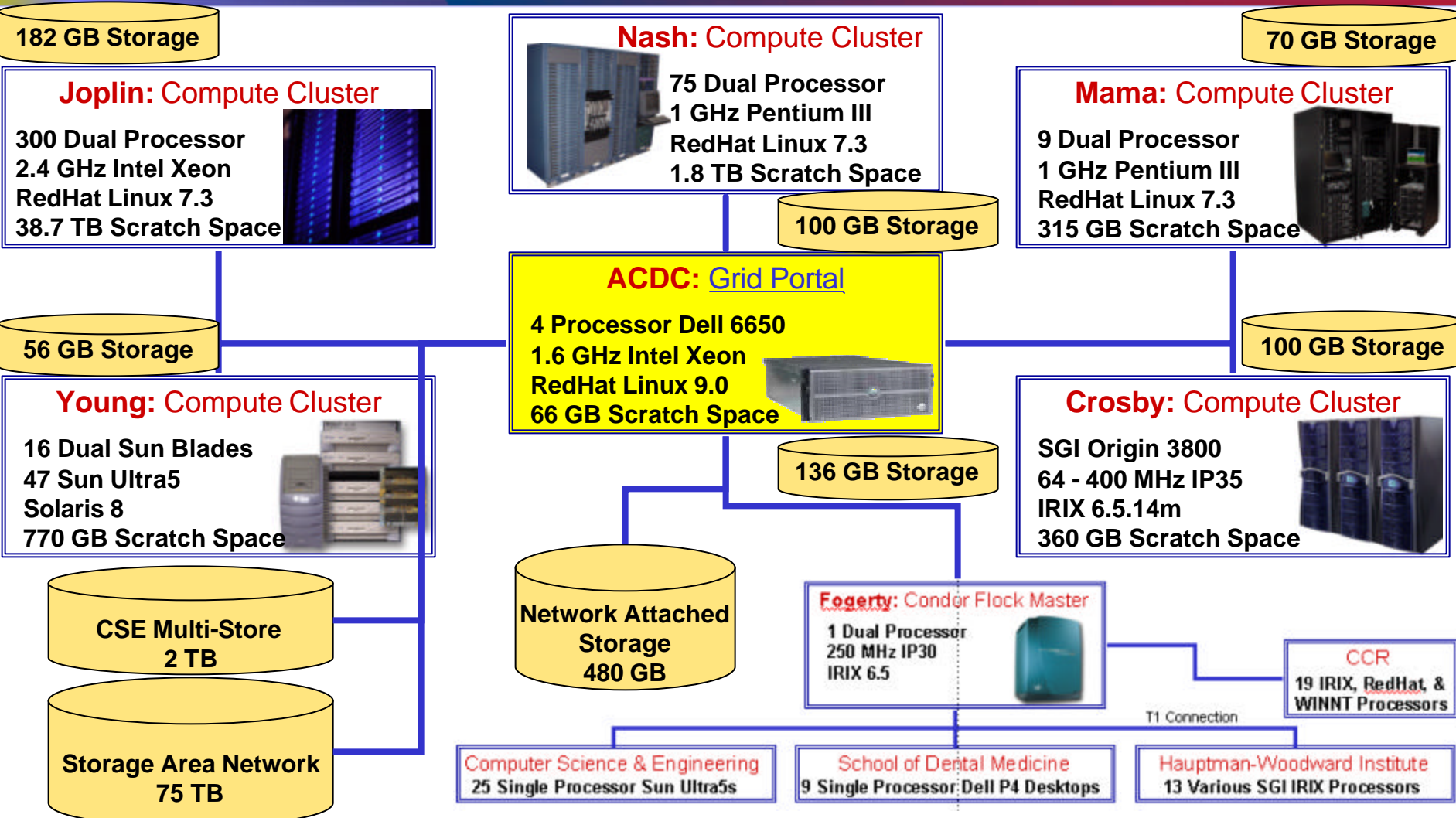
# Network Connections

CCR CENTER FOR COMPUTATIONAL RESEARCH
www.ccr.buffalo.edu

1000 Mbps — cse@buffalo

100 Mbps

FDDI

UB

1.54 Mbps (T1) - RPCI

ROSWELL PARK CANCER INSTITUTE

44.7 Mbps (T3) - BCOEB

1.54 Mbps (T1) - HWI

BCOEB

HW

100 Mbps

Medical/Dental

OC-3 - I1

155 Mbps (OC-3) I2

NYSERNet
350 Main St

Abilene

622 Mbps (OC-12)

Commercial

# Advanced CCR Data Center (ACDC) Computational Grid Overview

**Nash:** Compute Cluster

75 Dual Processor
1 GHz Pentium III
RedHat Linux 7.3
1.8 TB Scratch Space

**Joplin:** Compute Cluster

300 Dual Processor
2.4 GHz Intel Xeon
RedHat Linux 7.3
38.7 TB Scratch Space

**Mama:** Compute Cluster

9 Dual Processor
1 GHz Pentium III
RedHat Linux 7.3
315 GB Scratch Space

**ACDC:** Grid Portal

4 Processor Dell 6650
1.6 GHz Intel Xeon
RedHat Linux 9.0
66 GB Scratch Space

**Young:** Compute Cluster

16 Dual Sun Blades
47 Sun Ultra5
Solaris 8
770 GB Scratch Space

**Crosby:** Compute Cluster

SGI Origin 3800
64 - 400 MHz IP35
IRIX 6.5.14m
360 GB Scratch Space

**Fogerty:** Condor Flock Master

1 Dual Processor
250 MHz IP30
IRIX 6.5

Expanding
**RedHat, IRIX, Solaris, WINNT, etc**

CCR
**19 IRIX, RedHat, & WINNT Processors**

T1 Connection

Computer Science & Engineering
**25 Single Processor Sun Ultra5s**

School of Dental Medicine
**9 Single Processor Dell P4 Desktops**

Hauptman-Woodward Institute
**13 Various SGI IRIX Processors**

Note: Network connections are 100 Mbps unless otherwise noted.

**University at Buffalo** *The State University of New York* **Center for Computational Research**   **CCR**

# ACDC Data Grid Overview

**182 GB Storage**

## Joplin: Compute Cluster

300 Dual Processor
2.4 GHz Intel Xeon
RedHat Linux 7.3
38.7 TB Scratch Space

## Nash: Compute Cluster

75 Dual Processor
1 GHz Pentium III
RedHat Linux 7.3
1.8 TB Scratch Space

**70 GB Storage**

## Mama: Compute Cluster

9 Dual Processor
1 GHz Pentium III
RedHat Linux 7.3
315 GB Scratch Space

**100 GB Storage**

## ACDC: Grid Portal

4 Processor Dell 6650
1.6 GHz Intel Xeon
RedHat Linux 9.0
66 GB Scratch Space

**56 GB Storage**

## Young: Compute Cluster

16 Dual Sun Blades
47 Sun Ultra5
Solaris 8
770 GB Scratch Space

**100 GB Storage**

## Crosby: Compute Cluster

SGI Origin 3800
64 - 400 MHz IP35
IRIX 6.5.14m
360 GB Scratch Space

**136 GB Storage**

**CSE Multi-Store 2 TB**

**Storage Area Network 75 TB**

**Network Attached Storage 480 GB**

**Fogerty: Condor Flock Master**

1 Dual Processor
250 MHz IP30
IRIX 6.5

**CCR**
19 IRIX, RedHat, & WINNT Processors

**Computer Science & Engineering**
25 Single Processor Sun Ultra5s

**School of Dental Medicine**
9 Single Processor Dell P4 Desktops

**Hauptman-Woodward Institute**
13 Various SGI IRIX Processors

T1 Connection

Note: Network connections are 100 Mbps unless otherwise noted.

University at Buffalo  *The State University of New York*  Center for Computational Research  **CCR**

# WNY Grid Highlights

- **Heterogeneous Computational & Data Grid**
- **Currently in Beta with *Shake-and-Bake***
- **WNY Release in 2H04**
- **Bottom-Up General Purpose Implementation**
  - ❑ **Ease-of-Use User Tools**
  - ❑ **Administrative Tools**
- **Back-End Intelligence**
  - ❑ **Backfill Operations**
  - ❑ **Prediction and Analysis of Resources to Run Jobs (Compute Nodes + Requisite Data)**

# X-Ray Crystallography

- **Objective: Provide a 3-D mapping of the atoms in a crystal.**

- **Procedure:**

   1. **Isolate a single crystal.**

   2. **Perform the X-Ray diffraction experiment.**

   3. **Determine molecular structure that agrees with diffration data.**



Source of X-rays    Crystal

# X-Ray Data & Corresponding Molecular Structure

**Underlying atomic arrangement is related to the reflections by a 3-D Fourier transform.**

**Reciprocal Space**

FFT

FFT⁻¹

**Real Space**

**X-Ray Data**

**Molecular Structure**

- **Phases lost during the crystallographic experiment.**
- ***Phase Problem*: Determine phases of the reflections.**

# *Shake-and-Bake* Method: Dual-Space Refinement

**Shake-and-Bake**  **DeTitta, Hauptman, Miller, Weeks**

**Trial Structures**

**Structure Factors**

**Trial**

**Phases**

Tangent Formula

**FFT**

**Phase Refinement**

**Density Modification (Peak Picking) (LDE)**

Parameter Shift

**FFT$^{-1}$**

**?**

**Solutions**

*Reciprocal Space* **"Shake"**          *Real Space* **"Bake"**

# Ph8755: *SnB* Histogram



**Atoms: 74**        **Phases:  740**
**Space Group: P1**   **Triples:  7,400**

**Trials: 100**

**Cycles: 40**

**Rmin range: 0.243 - 0.429**

# Grid-Based *SnB* Objectives

- **Install Grid-Enabled Version of *SnB***
- **Job Submission and Monitoring over Internet**
- ***SnB* Output Stored in Database**
- ***SnB* Output Mined through Internet-Based Integrated Querying Tool**

- **Serve as Template for Chem-Grid & Bio-Grid**
- **Experience with Globus and Related Tools**

# Grid Services and Applications

**Applications**

| Shake-and-Bake | Apache | MySQL | Oracle |

**ACDC-Grid Computational Resources**

**High-level Services and Tools**

Globus Toolkit

NWS

| MPI | MPI-IO | C, C++, Fortran, PHP | globusrun |

**Core Services**

| Metacomputing Directory Service | Globus Security Interface | GRAM | GASS |

**ACDC-Grid Data Resources**

**Local Services**

| Condor | Stork | MPI | RedHat Linux | WINNT |
| LSF | PBS | Maui Scheduler | TCP | UDP | Irix | Solaris |

Adapted from Ian Foster and Carl Kesselman

# ACDC-Grid Portal Login

# Data Grid Capabilities

# Grid Portal Job Status

- **Grid-enabled jobs can be monitored using the Grid Portal web interface dynamically.**
  - ❑ **Charts are based on:**
    - ⭕ **total CPU hours, or**
    - ⭕ **total jobs, or**
    - ⭕ **total runtime.**
  - ❑ **Usage data for:**
    - ⭕ **running jobs, or**
    - ⭕ **queued jobs.**
  - ❑ **Individual or all resources.**
  - ❑ **Grouped by:**
    - ⭕ **group, or**
    - ⭕ **user, or**
    - ⭕ **queue.**

# Grid Portal Job Status

# ACDC-Grid Portal User Management

# ACDC-Grid Portal
# Resource Management



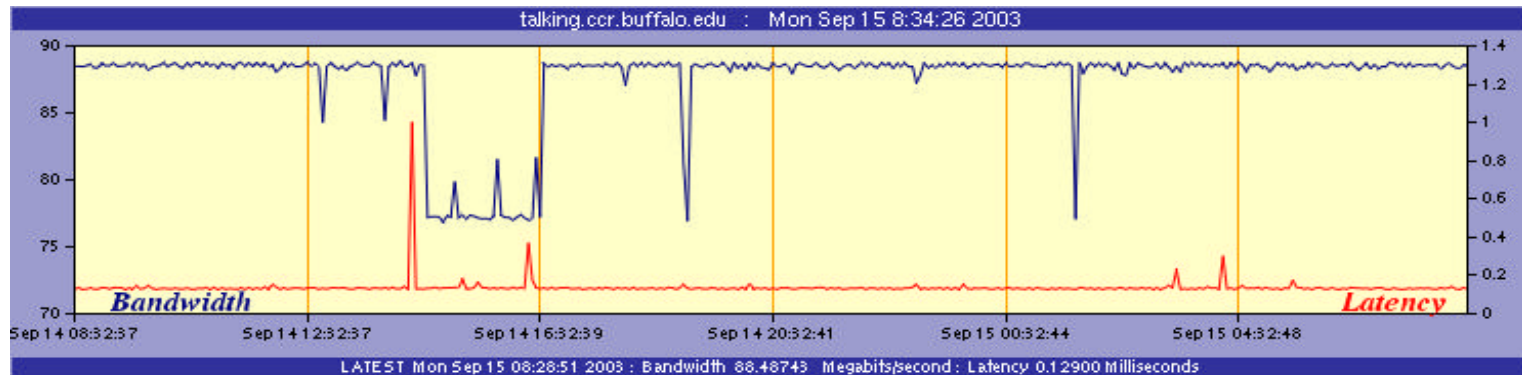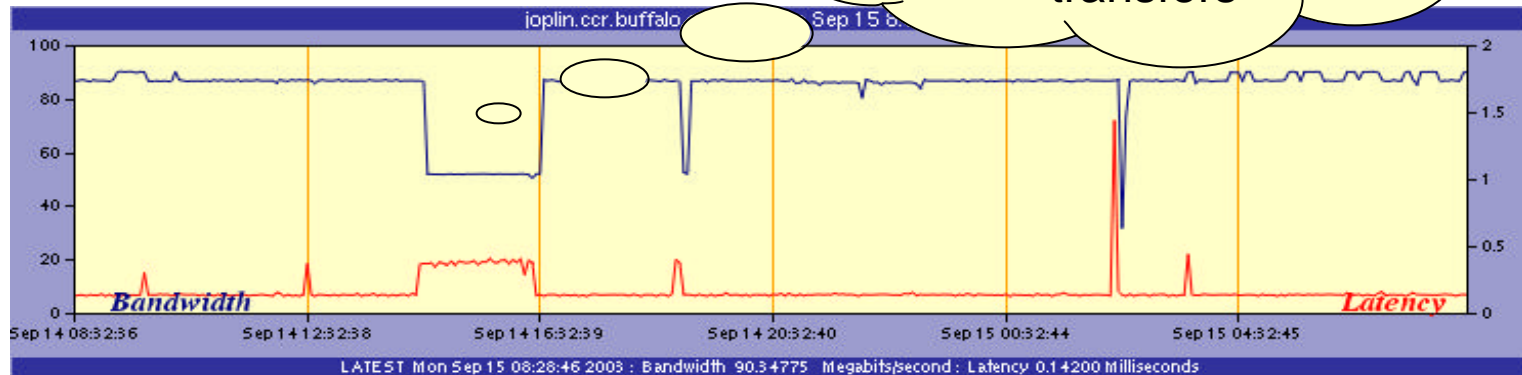■ **Administrator grants a user access to ACDC-Grid**
  ❑ **resources,**
  ❑ **software, and**
  ❑ **web pages.**

# Data Grid Resource Info

# Data Grid Resource Info

Both platforms have reduced bandwidth available for additional transfers

# Data Grid File Age



**File age, access time, and resource id denote:**

- ❑ **the amount of time since a file was accessed,**
- ❑ **when the file was accessed, and**
- ❑ **where the file currently resides respectively.**

# ACDC-Grid Development/Maintenance

- **Development Requirements**
  - ❑ **7 – Person months for Grid Services Coordinator**
    - ○ Including Grid and Database conceptual design and implementation
  - ❑ **5 – Person months for Grid Services Programmer**
    - ○ Web portal programming
  - ❑ **5 – Person months for System Administrator**
    - ○ Globus, NWS, MDS, etc. installations
  - ❑ **3 – Person months for Database Administrator**
    - ○ Grid Portal Database implementation

- **Minimum Maintenance Requirements**
  - ❑ **1 – Grid Services Coordinator**
    - ○ 100% level of effort
  - ❑ **1 – Grid Services Programmer**
    - ○ 100% level of effort
  - ❑ **1 – System Administrator**
    - ○ 50% level of effort
  - ❑ **1 – Database Administrator**
    - ○ 10% level of effort

# Acknowledgments

- **Steve Gallo**
- **Jason Rappleye**
- **Jeff Tilson**
- **Martins Innus**
- **Cynthia Cornelius**

- **George DeTitta**
- **Herb Hauptman**
- **Charles Weeks**
- **Steve Potter**

- **National Science Foundation, National Institutes of Health, Oishei Foundation, Wendt Foundation, Sloan Foundation, Verizon, NYS**

- **Gov Pataki, Congressman Reynolds, Senator Clinton, Senator Schumer, Congressman Quinn**

www.ccr.buffalo.edu