**Distributed Systems Laboratory**

**Technion**

**Computational Biology Laboratory**

# Superlink-online

A Distributed System For Genetic Linkage Analysis
using EGEE and BOINC

*Mark Silberstein, Artyom Sharov, Assaf Schuster, Dan Geiger*
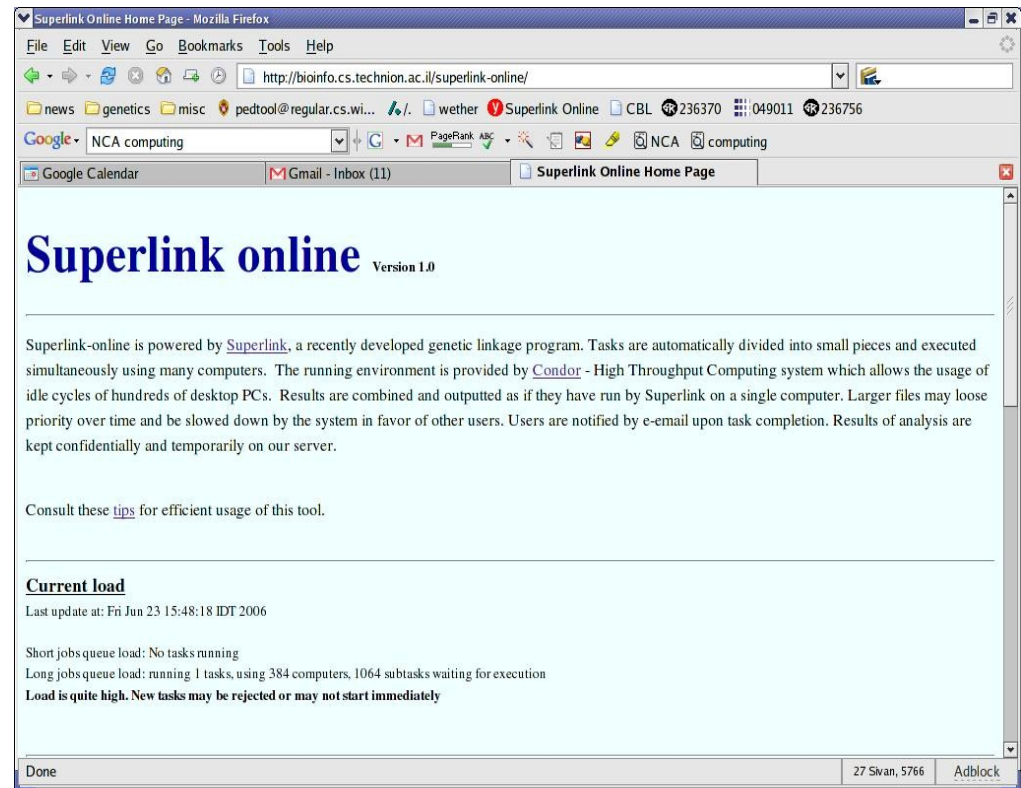
*CS  Department, Technion, Haifa*

**www.eu-egee.org**

Information Society
and Media

# Genetic Linkage Analysis

- **Purpose**: to obtain crude chromosomal location of gene(s) associated with a phenotype of interest

  – examples: Cystic fibrosis (found), diabetes, Alzheimer, blood pressure

- We focus on parametric linkage analysis on pedigrees

Artyom Sharov, Mark Silberstein

**Enabling Grids for E-sciencE**

- **Problem:**
  - Many analyses are infeasible due to the high computational demands

- **Reason:**
  - Exponential nature due to inference in Bayesian networks

- **Solution:**
  - Split the task into substantially smaller subtasks
  - Execute subtasks on multiple CPUs

**http://bioinfo.cs.technion.ac.il/superlink-online**

- **User submits his/her data for analysis**

  – No specification of running time

- **Secured user web interface**

  – Monitoring of partial results

  – Cancellation

- **E-mail notifications**

Enabling Grids for E-sciencE

- ~14,000 tasks, > 250 CPU years utilized since 2006

- Over 30 citations in leading genetics journals

- Over 200 users from universities and research centers in US, France, Germany, UK, Italy, Austria, Spain, Taiwan, Australia, and others

Artyom Sharov, Mark Silberstein

# Success stories

- ~14,000 tasks, > 250 CPU years utilized since 2006

- Over 30 citations in leading genetics journals

- Over 200 users from universities and research centers in US, France, Germany, UK, Italy, Austria, Spain, Taiwan, Australia, and others

Using over 3,500 computers in Israel and US
But need much more to allow comprehensive analysis of more complex data

- Few KB input/output

- **High RAM footprint**

- **Single job running time**

  – Seconds to hours (cannot be estimated exactly)

- **Single task**

  – up to **1M** jobs

- Few KB input/output

- **High RAM** footprint

- **Single job running time**

  – Seconds to hours (cannot be estimated exactly)

- **Single task**

  – up to **1M** jobs

Challenges:
Automatic, reliable, efficient execution
Acquiring many resources without VO coordination

Artyom Sharov, Mark Silberstein

- **A typical run:**

  - 0.5M jobs from several seconds to 1 hour (~20 minutes on average)

  - accomplished within 10 days

  - up to 2000 (1300 on average) concurrently executing clients in BIOMED VO

  - **600 PFLOPs** consumed

  - **fully automated**

  - **no prior coordination** with EGEE admins

  - 19 CPU years within 10 days utilized

    - Note: WISDOM project's recent **prioritized** data challenge utilized 144.21 years within 37 days

**eGee**

Agents wrapped as batch jobs

Agent submission and scheduling machine

Applciation jobs

Job dispatching server

Cluster

Artyom Sharov, Mark Silberstein

Enabling Grids for E-sciencE

Agents wrapped as batch jobs

Agent submission and scheduling machine

On hold

Applciation jobs

Job dispatching server

Cluster

Artyom Sharov, Mark Silberstein

Agents wrapped as batch jobs

Applciation jobs

Agent submission and scheduling machine

On hold

Job dispatching server

Cluster

Agents wrapped as batch jobs

Agent submission and scheduling machine

Applciation jobs

Job dispatching server

Cluster

BOINC clients wrapped as EGEE jobs

BOINC client submission and scheduling machine

Applciation jobs

Job dispatching server by BOINC

- Berkeley Open Infrastructure for Network Computing

- Out-of-the-box solution with 9 years of reputation

    – Backend of SETI@HOME

- Scalability: up to 2M hosts, billions of jobs in the queue

- Advanced scheduling

- Fault-tolerant

    – built for opportunistic environments

- Firewall-friendly

    – Clients pull jobs via HTTP

- Built in mechanisms to verify integrity and validity of results

- Built in accounting and statistics

- Submitted as an ordinary EGEE job

- Runs as long as there are jobs on the server, self-terminates if idle

- Restricted to 1 core to comply with batch system allocation (policy can be adjusted)

- Performance benefits:

  - run jobs back-to-back, caching the executable and constant data

  - no batch system scheduling overhead

  - beneficial even for seconds-long jobs

Enabling Grids for E-sciencE

- Injects clients into the batch system

- Keeps track of the running clients

  – finalizes output of finished clients

  – kills long-waiting clients

  – maintains *virtual cluster:* the required number of running clients

- Avoids Resource Broker overload

  – Use of multiple resource brokers

# Interface to administrators

- Centralized database for client policy

- Accounting and statistics

- Jobs results

**Generalized solution**
multiple resource pools
Enabling Grids for E-sciencE

- http://cbl-boinc-server2.cs.technion.ac.il/superlinkattechnion

Artyom Sharov, Mark Silberstein

Enabling Grids for E-sciencE

- Superlink-online genetic analysis portal
  - http://bioinfo.cs.technion.ac.il/superlink-online
- Superlink@Technion Community computing backend
  - http://cbl-boinc-server2.cs.technion.ac.il/superlinkattechnion
- Superlink@clusters EGEE and other clusters integration backend
  - http://cbl-boinc-server1.cs.technion.ac.il/superlinkatclusters
- Contact us:
  - Mark Silberstein - marks@cs.technion.ac.il
  - Artyom Sharov - sharov@cs.technion.ac.il

Artyom Sharov, Mark Silberstein

**eGee**

?

Artyom Sharov, Mark Silberstein

**egee**

- Execution hierarchy:
  - waterfall principle
  - more unreliable resources downhill

2-4 CPUs
**Dedicated server**

1,000 CPUs
**Clusters of Workstations**

**Computational Grids**        100,000 CPUs

**Community Grids**            1,000,000 CPUs

**eGee**

Enabling Grids for E-sciencE

~15,000 CPUs

**Computational Grids (EGEE-II BIOMED) via** BOINC