# Condor Team 2009

## Established 1985

# Welcome to Condor Week #11
## (year #26 for the project)

CONDOR
high throughput computing

www.cs.wisc.edu/Condor

THE UNIVERSITY of WISCONSIN
MADISON

# Main Threads of Activities

> **Distributed Computing Research** – develop and evaluate new concepts, frameworks and technologies that are based on distributed computing principals

> **Software Development & Support** - keep Condor software "**flight worthy**" and support our users

> **The UW Center for High Throughput Computing (CHTC)** – provide HTC capabilities to the UW campus via a distributed computing and storage infrastructure

> **The Open Science Grid (OSG)** – a national, distributed computing partnership for data-intensive research

> **The NSF Middleware Initiative (NMI)** – develop, build and operate a national Build and Test facility

MultiCore

Cyber Infrastructure

SaaS

WorkFlows

Workforce

HDFS

HTPC

IaaS

Grids

eResearch

HPC

Cyber Security

100Gb

SLA

Virtual Machines

QoS

Web Services

eScience

Open Source

HTC

Map-Reduce

Clouds

Master-Worker

High Availability

Green Computing

GPUs

# Center for Enabling Distributed Petascale Science
## A SciDAC Center for Enabling Technology
(http://www.cedps.net/images/8/81/CEDPS_ProjectPlan_Oct2009.pdf)

We provide a project plan for the Center for Enabling Distributed Petascale Science (CEDPS), covering the second phase of the project, from October 2009 to June 2011. Our goal in this second phase is to *develop, evaluate, and apply a fully integrated approach to data movement and caching problems that provides an end--to--end solution for extreme-scale data-intensive applications*. This solution uses highly available services, termed **Globus.org**, to provide the needed services; integrates these services into resource managers, such as Condor and the FermiLab glide-in workload management service; and incorporates a state of-the-art logging infrastructure to enhance performance and reliability.

# Claims for "benefits" provided by Distributed Processing Systems

P.H. Enslow, *"What is a Distributed Data Processing System?"* Computer, January 1978

- High Availability and Reliability
- High System Performance
- Ease of Modular and Incremental Growth
- Automatic Load and Resource Sharing
- Good Response to Temporary Overloads
- Easy Expansion in Capacity and/or Function

# Definitional Criteria for a Distributed Processing System

P.H. Enslow and T. G. Saponas ""*Distributed and Decentralized Control in Fully Distributed Processing Systems*" Technical Report, 1981

- Multiplicity of resources
- Component interconnection
- Unity of control
- System transparency
- Component autonomy

TIME *A thinker's guide to the most important trends of the new decade*

# In Defense of Failure

### By Megan McArdle Thursday, Mar. 11, 2010

"The goal shouldn't be to eliminate failure; it should be to build a system resilient enough to withstand it"

"The real secret of our success is that we learn from the past, and then we forget it. Unfortunately, we're dangerously close to forgetting the most important lessons of our own history: how to fail gracefully and how to get back on our feet with equal grace."

CONDOR
high throughput computing

THE UNIVERSITY of WISCONSIN
MADISON

# From Forbes Magazine …

## Open Source Energy Savings

*Dan Woods*, 03.02.10, 06:00 AM EST **Software for spreading work over huge collections of computers can be used to cut power costs.**

**Condor** supports all the operating systems a typical company or research institution would have and is **rock solid** in terms of stability and functions for its intended purpose, which is carving up work and sending it out to any number of computers for processing.

www.cs.wisc.edu/~miron

# What Did We Learn From Serving a Quarter of a Million Batch Jobs on a Cluster of Privately Owned Workstations

# 1992

*Miron Livny*

Computer Sciences Department
University of Wisconsin — Madison
Madison, Wisconsin
{miron@cs.wisc.edu}

---

# User Prospective

- Maximize the capacity of resources accessible via a single interface
- Minimize overhead of accessing remote capacity
- Preserve local computation environment

# Global
# Scientific Computing
## via a
# Flock of Condors

**CERN 92**

*Miron Livny*

Computer Sciences Department
University of Wisconsin — Madison
Madison, Wisconsin
{miron@cs.wisc.edu}

---

## MISSION

Give scientists effective and efficient access to large amounts of cheap (if possible free) CPU cycles and main memory storage

---

## THE CHALLENGE

How to turn existing privitely owned clusters of *workstations*, *farms*, *multiprocessors*, and *supercomputers* into an efficient and effective Global Computing Environment?

In other words, how to minimize wait while idle?

---

## APPROACH

Use wide-area networks to transfer batch jobs between Condor systems

- Boundaries of each Condor system will be determined by physical or administrative considerations

---

## TWO EFFORTS

☐ *UW CAMPUS*

Condor systems at Engineering, Statistics, and Computer Sciences

☐ *INTERNATIONAL*

We have started a collaboration between CERN-SMC-NIKHEF-Univ. of Amsterdam, and University of Wisconsin-Madison

# Open Science Grid (OSG)

# Native Package in Condor

> Beginning in Condor 7.4.x, we have new native packages:
> - Improved RPM with a yum repository
>   - Condor "just works" out of the box
>   - Condor is now installed in standard locations (/usr/bin, /var/log, ...)
>   - Creates "condor" user
>   - Installs Condor init script
> - New Debian package with deb repository
>   - Matches features in RPM (above)

> Condor is included the Fedora Distribution
> Condor is the "G" in REHL MRG
> Condor is a Ubuntu project

CONDOR
high throughput computing

THE UNIVERSITY of
WISCONSIN
MADISON

## HPCwire: Microsoft Injects More Goodies into Windows HPC (Published 09 April 10 02:27 PM)

<...>

For example, Microsoft has added the ability to aggregate Windows 7 workstations into an HPC cluster. Essentially, each workstation is monitored and managed as an ad-hoc compute node within a traditional cluster. Capabilities like time-of-day scheduling and shutting down of preemptive jobs are included so that the machine can be made available to a real live person when required. The common scenario is one where an organization has a small cluster made up of say dozens of servers, along with perhaps hundreds of Windows 7 PCs sitting idle at night.

**114M Hours**

# The **Qiang Cui (QC) Group**

"We develop and apply theoretical tools (electronic structure, nuclear dynamics and statistical mechanics) to biophysical problems (enzyme catalysis, bioenergy transduction and biomaterials etc.)."

http://kandinsky.chem.wisc.edu/~qiang/

**Started to use the High Through Parallel Computing (HTPC) provided by CHTC in 05/09 and expended to OSG in 11/09**

# QC Group Comsumption

# Impact of these cycles

1. Iron-Catalyzed Oxidation Intermediates Captured in A DNA Repair Monooxygenase, C. Yi, G. Jia, G. Hou, Q. Dai, G. Zheng, X. Jian, C. G. Yang, Q. Cui, and C. He, **Science**, Submitted

2. Disruption and formation of surface salt bridges are coupled to DNA binding in integration host factor (IHF): a computational analysis, L. Ma, M. T. Record, Jr., N. Sundlass, R. T. Raines and Q. Cui, **J. Mol. Biol,** Submitted

3. An implicit solvent model for SCC-DFTB with Charge-Dependent Radii, G. Hou, X. Zhu and Q. Cui, **J. Chem. Theo. Comp.**, Submitted

4. Sequence-dependent interaction of $\beta$-peptides with membranes, J. Mondal, X. Zhu, Q. Cui and A. Yethiraj, **J. Am. Chem. Soc**., Submitted

5. A new coarse-grained model for water: The importance of electrostatic interactions, Z. Wu, Q. Cui and A. Yethiraj, **J. Phys. Chem**. Submitted

6. How does bone sialoprotein promote the nucleation of hydroxyapatite? A molecular dynamics study using model peptides of different conformations, Y. Yang, Q. Cui, and N. Sahai, {\it Langmuir}, Submitted

7. Preferential interactions between small solutes and the protein backbone: A computational analysis, L. Ma, L. Pegram, M. T. Record, Jr., Q. Cui, **Biochem**., 49, 1954-1962 (2010)

8. Establishing effective simulation protocols for $\beta$- and $\alpha/\beta$-peptides. III. Molecular Mechanical (MM) model for a non-cyclic $\beta$-residue, X. Zhu, P. K\"onig, M. Hoffman, A. Yethiraj and Q. Cui, **J. Comp. Chem.**, In press (DOI: 10.1002/jcc.21493)

9. Curvature Generation and Pressure Profile in Membrane with lysolipids: Insights from coarse-grained

   simulations, J. Yoo and Q. Cui, **Biophys. J**. 97, 2267-2276 (2009)

# MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat

jeff@google.com, sanjay@google.com

*Google, Inc.*

## Abstract

MapReduce is a programming model
ated implementation for processing and
data sets. Users specify a *map* function
key/value pair to generate a set of interme
pairs, and a *reduce* function that merges
values associated with the same interme
real world tasks are expressible in this n
in the paper.

Programs written in this functional sty
cally parallelized and executed on a large
modity machines. The run-time system t
details of partitioning the input data, sch
gram's execution across a set of machine
chine failures, and managing the require
communication. This allows programmers without any
experience with parallel and distributed systems to eas-
ily utilize the resources of a large distributed system.

The MapReduce implementation relies on an in-house cluster management system that is responsible for distributing and running user tasks on a large collection of shared machines. Though not the focus of this paper, the cluster management system is similar in spirit to other systems such as Condor [16].

and many other functional languages. We realized that most of our computations involved applying a *map* operation to each logical "record" in our input in order to compute a set of intermediate key/value pairs, and then apply

BAD-FS [5] has a very different programming model from MapReduce, and unlike MapReduce, is targeted to the execution of jobs across a wide-area network. However, there are two fundamental similarities. (1) Both systems use redundant execution to recover from data loss caused by failures. (2) Both use locality-aware scheduling to reduce the amount of data sent across congested network links.

measurements of our implementation for a variety of tasks. Section 6 explores the use of MapReduce within Google including our experiences in using it as the basis

Pete,

  Here are some numbers you ask about for LIGO's use of DAGs to manage large data analysis tasks broken down by the largest number of jobs managed in different categories:

1) DAG Instance--one condor_dagman process: **196,862.**
2) DAG Workflow--launched from a single condor_submit_dag but may include multiple automatic sub- or spliced DAGs: **1,120,659.**
3) DAG Analysis--multiple instances of condor_submit_dag to analyze a common dataset with results combined into a single coherent scientific result: **6,200,000.**
4) DAG Total--sum over all instances of condor dagman run: **O (100M).**

P.S. These are lower bounds as I did not perform an exhaustive survey/search, but they are probably close.

Thanks.

# Thank you for building such



## a wonderful community