# Computational Teams Cybertools / Cyberinfrastructure (CTCI)
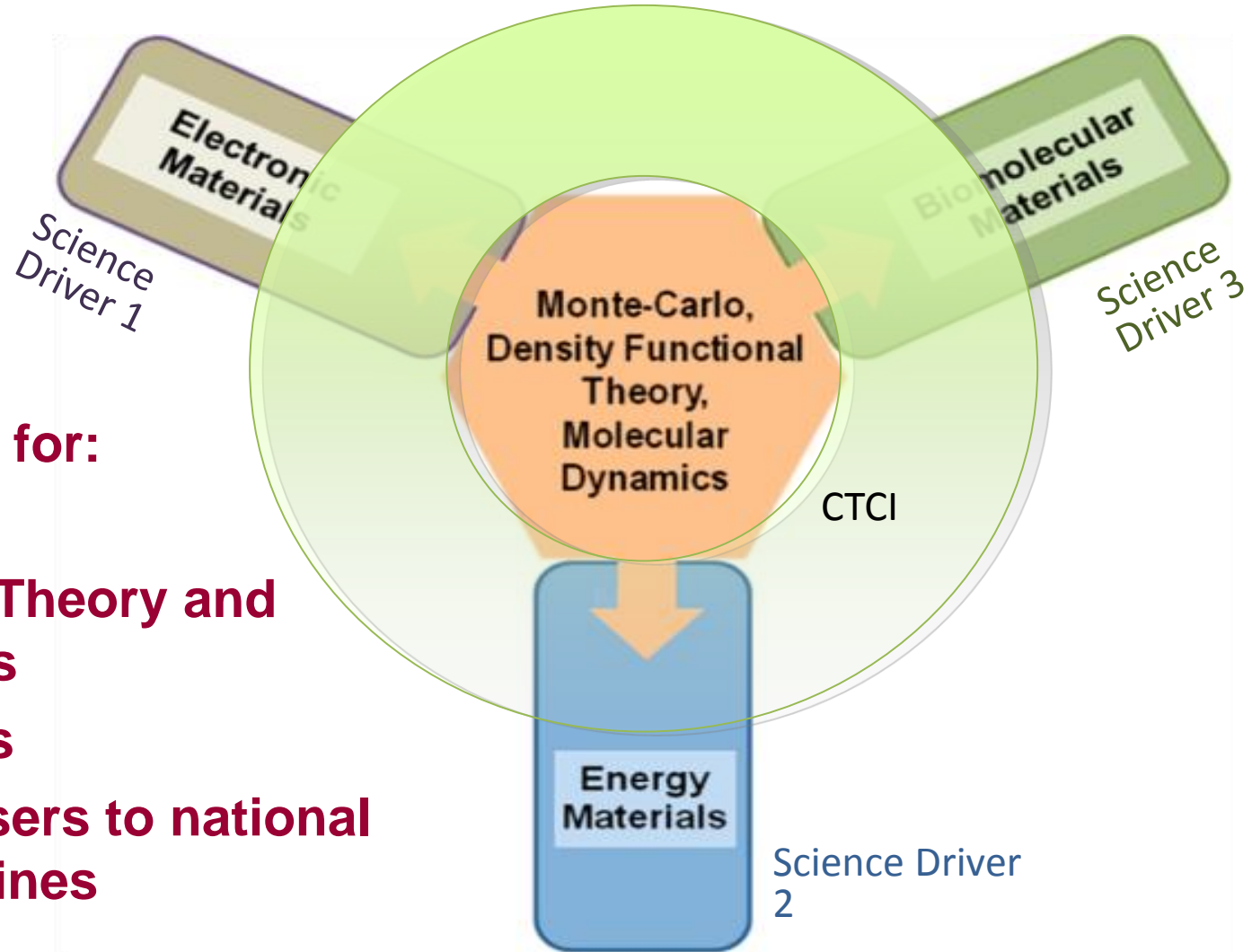
## "The glue"

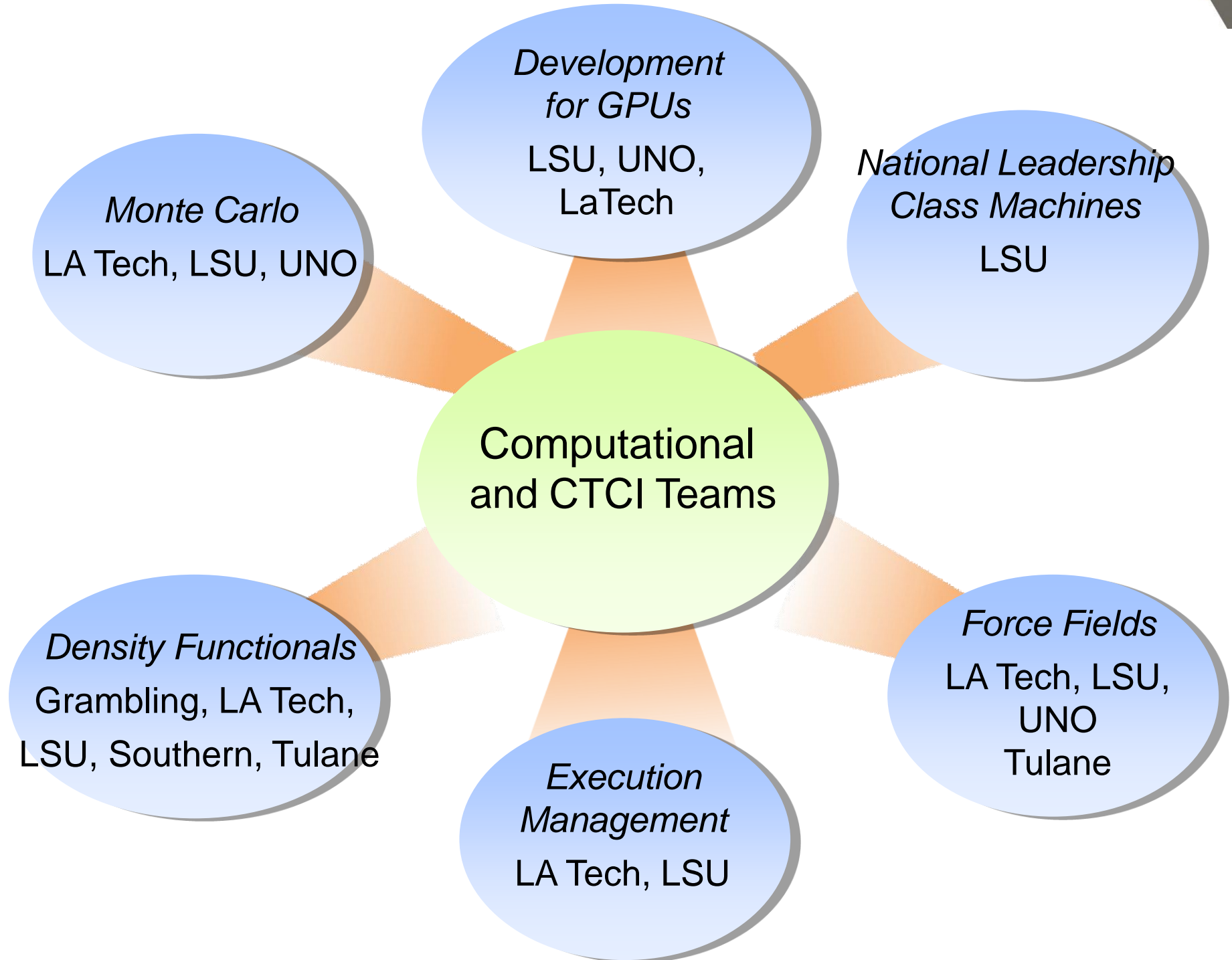LA-SiGMA Symposium, Baton Rouge: July 23, 2012

# The Goal

## Computational and CTCI Teams: Build transformational common toolkits

- **Build common toolkits for:**
  - **Monte Carlo**
  - **Density Functional Theory and Force Field Methods**
  - **Molecular Dynamics**
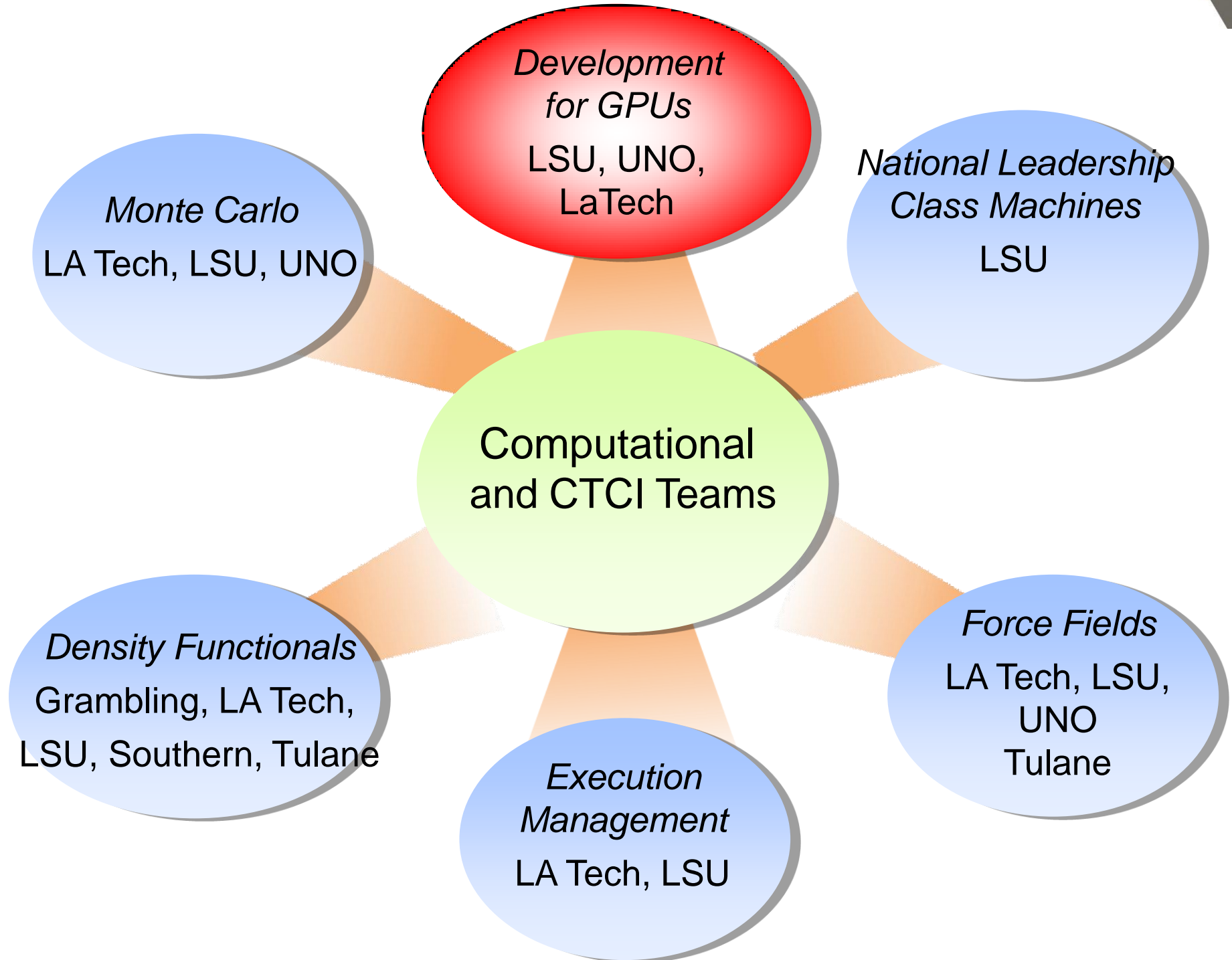- **Graduate LA-SiGMA users to national leadership class machines**

Electronic Materials

Science Driver 1

Biomolecular Materials

Science Driver 3

Monte-Carlo, Density Functional Theory, Molecular Dynamics

CTCI

Energy Materials

Science Driver 2

# Computational/CTCI Research Themes



**Development for GPUs**
LSU, UNO, LaTech

**National Leadership Class Machines**
LSU

**Monte Carlo**
LA Tech, LSU, UNO

**Computational and CTCI Teams**

**Force Fields**
LA Tech, LSU, UNO Tulane

**Density Functionals**
Grambling, LA Tech, LSU, Southern, Tulane

**Execution Management**
LA Tech, LSU

# Computational/CTCI Research Themes



Development for GPUs
LSU, UNO, LaTech

National Leadership Class Machines
LSU

Monte Carlo
LA Tech, LSU, UNO

Computational and CTCI Teams

Density Functionals
Grambling, LA Tech, LSU, Southern, Tulane

Force Fields
LA Tech, LSU, UNO Tulane

Execution Management
LA Tech, LSU

# GPU: Program Optimization

- Programming GPUs (for example, the NVIDIA GPUs using CUDA or OpenCL) is still tedious:
  - Performance of GPU highly sensitive to the formulation of the kernel; needs significant experimentation
  - Programmers may like this low level of control (suitable for library development; compilers and tools are not highly helpful here)
- **Strategy at LSU (Yun, Ramanujam):**
  - Understand the impact of and interactions among program optimizations for HF-QMC, PT and VMC
  - Develop and use effective transformation and optimization strategies
  - Code partitioning between CPU and GPU
- **Strategy at LaTech (Leungsuksun):**
  - Parallel Programming Tool Development based on Single Assignment C (SaC) toolset that enables parallel application developers expressing their problems in a high-level language

# GPU Programming Team

- GPU Programming Team of roughly 25 faculty, students, and postdocs from LSU, LA Tech & Lousiana School for Math, Sciences, and the Arts (HS).
- Housed in the _Collaboratorium_ at LSU and containing 12 GPU-enabled desktop computers.
- LSU "Condo" – LA-SiGMA to purchase GPU-enabled nodes on CCT/LSU's Tezpur upgrade.
- Another GPU cluster (_Shelob_) from NSF CRI funding (significant LA-SiGMA role)





Nvidia GTX 460, http://www.nvidia.com/object/product-geforce-gtx-460-us.html

# Parallel Tempering (See poster)

- The main goal here is to develop an efficient Parallel Tempering Monte Carlo code on GPU, with which we can study systems with complex energy landscapes.

- Developed a full-featured Ising model simulation program for CUDA GPUs.

- Studying phase transition of spin glass in a finite magnetization field.

- Results:
  - The FPGA design (custom hardware) of Montovani et al. delivers the best time of 16 picoseconds per spin flip proposal (PS/spin).
  - We achieve 39 picoseconds per spin flip proposal (PS/spin) on a single commodity GPU card, which is 3X better than other 2D GPU implementations.
  - Our GPU version is about 600 times faster than our prototype CPU implementation.

# Hirsch-Fye Quantum Monte Carlo-1 (See poster)

- This project simulates the interaction of conducting electrons in a metal.

- Using the Hirsch-Fye method mapped to a problem of electrons scattering of an Ising field in space and (imaginary) time.

- The configurations of the Ising field are sampled using Monte Carlo techniques.

# Hirsch-Fye Quantum Monte Carlo-2 (See poster)

- A key optimization for GPUs is replacing a single outer product with a panel-panel matrix multiplication by employing the technique of delayed updating. This increases the computation to memory access ratio.

$$G_{ij}^{(M)} = G_{ij}^{(1)} + \sum_{N=1}^{M} A_i^{(N)} B_j^{(N)} \rightarrow$$

# Seeding the Experimental Inverse Problem (see Poster)



From X-ray and neutron scattering data to images of flame retardants in polymer blends

Using results of actual observations to infer the values of the parameters characterizing the system under investigation

# Computational/CTCI Research Themes

# Graduating to National Leadership Class (NLC) Machines

- Explore parallelism and scalability

- Get experience with code development on smaller clusters such as LONI/TeraGrid-XSEDE

- Demonstrate how your codes will scale to the NLC machines

- Apply for compute time on NLC machines

# Current XSEDE Allocations 2012

- Bishop @LaTech:    7.8MSU
- Mobley @UNO:       1.2MSU
- Morena @LSU:       7.5MSU
- Sun @Tulane :      200,000 SU

# The Little Fe Project:
# Recruiting the next users.

SuperComputing 2011: Little Fe Build Out Session.

    Groups Selected by Application/proposal process

    12 Little Fe's Awarded: 3 in Louisiana

Louisiana Tech

Louisiana School Science Math and the Arts (High School)

LSU

# Computers LA-SiGMA can access

- **A 8-node Dell PowerEdge R70 GPU cluster:**
  - **2 Sandy Bridge 2.5GHz CPUs, 64GB mem**
  - **2 NVIDIA Tesla M2090 6GB GPU**
- **The new SuperMike:**
  - **Minimum of 146 CPU Tflops in 382 CPU nodes:**
    - **2x Sandy Bridge 8-core 2.6GHz, 32GB mem**
  - **66 GPU Tflops from 50 GPU M2090 nodes**
    - **2x Sandy Bridge 8-core 2.6GHz, 64GB mem**
    - **2x NVIDIA M2090 GPUs**



| Application | Model | 2x CPUs (16 cores) | 2x CPUs + 1x GPU | 2x CPUs + 2x GPUs | Speedup with 2 GPUs |
|---|---|---|---|---|---|
| NAMD (days/ns) | ap0a1 | 1.78 | 0.98 | 0.50 | 3.52 |
| | f1atpase | 5.28 | 1.26 | 0.92 | 5.72 |
| | stmv | 19.58 | 4.95 | 2.94 | 6.66 |
| AMBER (ns/day) | JAC_NVE | 10.72 | 33.29 | 46.18 | 4.31 |
| | JAC_NPT | 9.42 | 28.79 | 38.79 | 4.12 |
| | FactorIX_NVE | 2.50 | 9.39 | 13.07 | 5.23 |
| | FactorIX_NPT | 2.28 | 8.30 | 11.41 | 5.00 |
| | Cellulose_NVE | 0.57 | 2.04 | 2.90 | 5.09 |
| | Cellulose_NPT | 0.55 | 1.99 | 2.79 | 5.07 |
| LAMMPS (loop time sec) | EAM | 600.94 | | 129.36 | 4.65 |
| | AU | 458.20 | | 113.31 | 4.04 |

- **NSF CRI (H. Liu, PI) purchase an NVIDIA Kepler (K20) system**
  - **At least 24 GPU nodes with:**
    - **2 Intel Sandy Bridge-EP processors, 64GB memory,**
    - **At least 2 Kepler K20 GPUs.**

R

# Computational/CTCI Research Themes

# Density Functionals and Force Fields (rsv)

Perdew group has developed a "work-horse semilocal functional" [*Phys. Rev. Lett.* 103, 026403 (2009)] for large electronic systems that yields accurate lattice constants, surface energies, and atomization energies.

• This functional has been incorporated into VASP, a massively parallel DFT code.
• We are constructing force fields combining *ab initio* calculations of small clusters with different DFT functionals and bulk simulations/calculations.



$Cu_4O_4H_4$ cluster
*LSU (Hall, Dellinger), La Tech (Wick, Ramachandran)*

# Computational/CTCI Research Themes

# Ensemble Based MD Simulation Challenges

**Nucleosome Simulations on XSEDE**

336 Systems * 20ns = 6,720 tasks

64-128cpu/task * 8hrs ~ 3.5 MSU

5Gb/task * 6720 ~ 34 TB Data

NAMD with BigJobs



16 nucleosomes * 21 positions = 336 systems
160,000 atoms per system

**DNA: Simulations on LONI**

4 Systems * 1000ns = 4,000 tasks

128 cpu/task * 2.5hrs ~ 1.3 MSU

1Gb/task * 4000 ~ 4TB Data

Amber with ManyJobs

ABC: International Collaboration



4 sequences (18bp) = 4 systems
47,000 atoms per system

# Scale Across
## (Bishop & Jha)

- efficiently distribute computations across computers:
  XSEDE
  LONI
  Local Clusters
  Whatever

Pilot Job Concept
Two Implementations:
  BigJobs:  SAGA
  ManyJobs: Python

# BigJobs

- SAGA: Slide in API for Grid Applications

- LONI, XSEDE and
national grade infrastructure

- Recently restructured
  and deployed on XSEDE
Advert service on XSEDE VM Data
  Quarry.

Updated documentation and examples
https://github.com/sagaproject/BigJob/wiki



**Nucleosome Ensemble**
63 Simulations * 192 Core
12,096 CPU on Kraken
Min 4hr Run time
1 ns of 160,000 atom system

# ManyJobs

- Python Based with ssh (gsi-ssh)
  - "no prerequisites"
  - Easy deployment

- Clusters, LONI, XSEDE

**DNA Ensembles**
4 Simulations * 128 Core
5 LONI Machines
Min 2hr run time for 1ns
3,600 ns of 40,000 atom systems

# ManyJobs



The Ascona B-DNA Consortium Simulation Effort
Coordinating > 20 Int.'l research groups and > 100,000 Simulations
 Bishop: 3 month run period ~ 3500 simulations ~ 900,000SU on LONI

# Progress: XSEDE 12

**Running Many Molecular Dynamics Simulations on Many Supercomputers**

Rajib Mukherjee, Abhinav Thota, Hideki Fujioka, Thomas C. Bishop and Shantenu Jha

**The Anatomy of a Successful ECSS Project: Lessons of Supporting High Throughput High Performance Ensembles on XSEDE**

Melissa Romanus, Pradeep Mantha, Yaakoub El Khamra, Andre Merzky,Shantenu Jha, Matt McKenzie and Thomas C Bishop

**XSEDE Campus Bridge Early Adopter Program:**

Global Federated File System (GFFS) Pilot Project:

Goal to incorporate GFFS technology into High Performance High Throughput Simulation Workflow.

C.Stewart, R.Knepper, T.Miller, A. Grimshaw, T.C. Bishop, S. Jha.

# Posters

Rocky Brown, REU from Radford

Victoria Bamburg, REU from LSMA

# Computational Tools for Multi-scale Simulations (Dua, LaTech)

**Goal:** To develop techniques, algorithms, and strategies for extracting information and knowledge from data generated by Science Drivers and create Computational Tools related efforts.



Fig. 1 Avg. pairwise Euclidean distance v/s dimensions

**Efforts over the past year:**

- Data adaptive rule based approach to supervised learning.
- A grid based agglomerative approach to unsupervised learning.

**Future Efforts:**

- Distributed data mining frameworks.
- Proposed system architecture.
- Integration of variants of the above approaches to proposed architecture

# Data adaptive rule based supervised learning

**Goal:** To develop a data-adaptive partitioning schema of feature space for rule-based classification.

**Objectives :**

To develop a data adaptive partitioning scheme

To develop a method for rule extraction

To exploit the extracted rules for supervised learning / classification



Fig. 2. Histogram plotted for one feature of a dataset.

**Significance & Applications:**

- Data adaptive partitioning ensures reduction in the number of rules
- Ensure the choice of rules that are both high in sensitivity and specificity
- Modular development of algorithm for easy distribution

# Initial Results

Overall accuracy of proposed data adaptive partitioning classification results have compared with rule-based classifiers and non-rule based classifiers.

| Classifiers | Overall Accuracy (%) |
|---|---|
| **Rule Based Classifiers** | |
| Conjuctive Rule | 66 |
| Decision table | 77 |
| DTNB | **82** |
| JRIP | 66 |
| NNGE | 75 |
| One R | 62 |
| PART | 77 |
| Ridor | **82** |
| ZeroR | 33 |

| Classifiers | Overall Accuracy (%) |
|---|---|
| **Non-rule based classifiers** | |
| Naïve Bayes | 84 |
| Logistic | 73 |
| Multi Layer Perceptron | 59 |
| RBF Network | **86** |
| Simple Logistic | 77 |
| SMO | **86** |
| Random Forest | 80 |

| Classifiers | Overall Accuracy (%) |
|---|---|
| **Proposed data adaptive partitioning** | |
| Slope based partitioning | 82.2 |
| Non-slope based partitioning | 86 |

# Grid-based agglomerative clustering algorithm

**Goal:** The goal of this research is to develop a data mining algorithm for clustering multidimensional datasets.

**Objectives:**

- To develop an algorithm for multi-level data adaptive grid generation.
- To develop a data preprocessing algorithm for sparseness reduction.
- To develop a grid based agglomerative hierarchical clustering algorithm.



Fig. 7. (a) A uniform grid for 2D data, (b) A non-uniform grid for 2D data



Fig. 8. A two dimensional grid with grid cell numbering

**Significance & Applications:**

- Clustering algorithms augmented with a data preprocessing through sparseness reduction are more accurate and produce better clustering results.
- Our developed algorithm is generic, and is easily adaptable for other scientific applications.

# Initial Results - scalability analysis

## Grid generation



Fig.11. Execution time v/s dataset size

## Clustering algorithm



Fig.12: Execution time v/s dimensions

# Proposed extension

Based on the MapReduce programming paradigm.

Apache Hadoop
- Hadoop - Distributed File System
- Hadoop - MapReduce.
  - Map function.
  - Reduce function.

# Data Mining using MapReduce

Requirements:

**Scalability:** We mean that the system can easily be altered to accommodate changes in the number of users, resources and computing entities.

**Reliability:** Difficult to achieve as it is closely related to the complexity of the interactions between simultaneously running components.

**Availability:** The system can restore operations, permitting it to resume providing services even when some components have failed.

**Evolution:** Keeping up with changes to the system with newer computational features and newer requirements.

# Funding and Outreach

- INCITE proposal (in collaboration with Pacific Northwest National Laboratories) for compute cycles on Jaguar and Titan
- XSEDE Allocation and GFFS incorporation into BigJobs
- ManyJobs with LONI-CS: Hideki Fujioka
- NSF CRI proposal for GPU cluster (*Shelob*) funded (includes several LA-SiGMA faculty members)
- NSF proposal for ScaleMS Bishop and Jha
- Indo-US Center (IUSSTF)
- SCiDAC and other DOE proposals
- Outreach:
  - Summer REU and RET programs
  - Beowulf <u>Boot Camp</u> for High School Students and Teachers
  - Little Fe
  - FEScUE with Colorado State University (Bishop)
  - GPU and Execution Management regular video meetings
  - Conference tutorials on GPUs:
    - International Symp. on Code Gen. & Opt., April 2012
    - Intl. Conf. Parallel Arch. & Comp. Tech., Oct. 2011

# Summary

- CTCI: Leveraging the exponential increase in computer power
  - Recruiting new and graduating existing users to national leadership class machines
  - Preparing users for next-generation computers
  - Developing common computational toolkits
  - Expanding collaborations within LA-SiGMA and developing partnerships with national labs

"The glue" that binds the SDs

CTCI

# Thank You