# Volunteer Computing and Protein-Ligand Docking

Michela Taufer
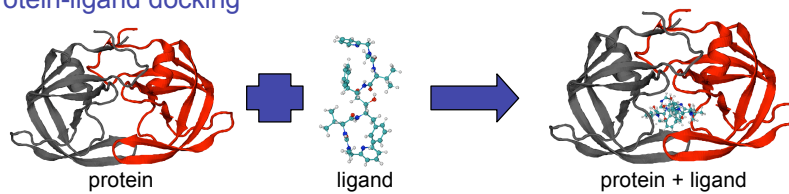
**GCL**

Global Computing Lab
University of Delaware

---

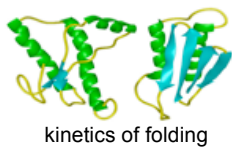# Simulating Biological Systems

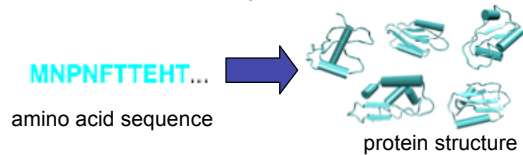**Protein-ligand docking**



protein          ligand          protein + ligand

**Protein folding**                **Protein structure prediction**



MNPNFTTEHT...

kinetics of folding          amino acid sequence          protein structure

- Search in large conformational spaces based on e.g., Monte Carlo and/or Molecular Dynamics simulations
- Different models of the biological systems require different amounts of resources/time and provide different levels of accuracy
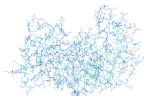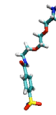
# Protein-ligand Docking

- Procedure to identify candidates for drug development by screening large protein-ligand databases
- Protein-ligand docking as a set of random attempts using Molecular Dynamics (MD) simulations and CHARMM force field (*M.Taufer et al., Concurrency and Computation'05*)
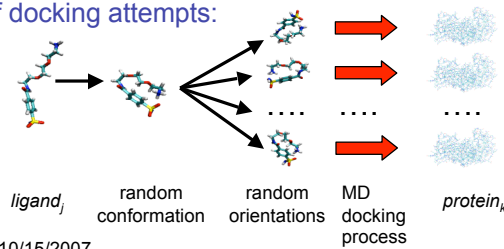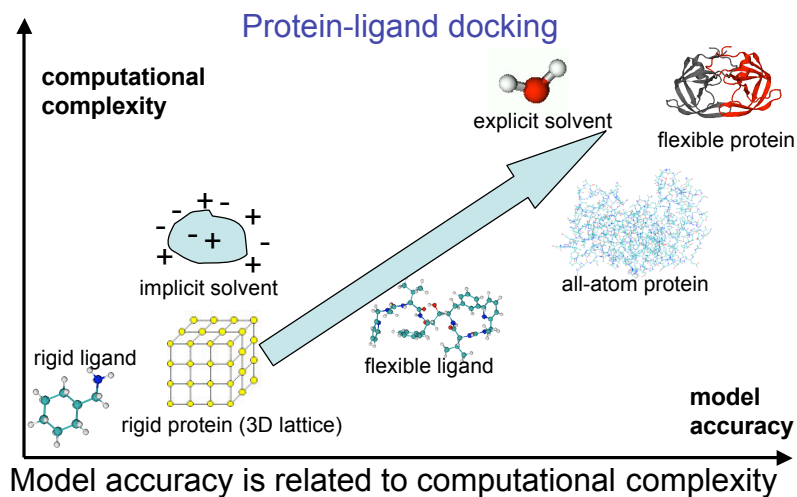
Given a protein

Given a ligand

Example of docking attempts:

….       ….       ….

*ligand$_j$*    random conformation    random orientations    MD docking process    *protein$_k$*

Michela Taufer - 10/15/2007

3

---

# Complexity vs Accuracy

Protein-ligand docking

**computational complexity**

explicit solvent

flexible protein

-   +   -
-   +  
-   +   -
-   +   -

implicit solvent

all-atom protein

rigid ligand

flexible ligand

rigid protein (3D lattice)

**model accuracy**

Model accuracy is related to computational complexity

Michela Taufer - 10/15/2007

4

# Using Supercomputers for many Months



*Waiting for resources*        *Resource contention*        *Resource utilization*        *time*

*We need much more compute resources than are available at current supercomputer centers*

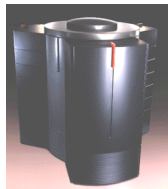Michela Taufer - 10/15/2007                                                                    5

---

# Searching for New Compute Resources: Volunteer Computing

*Late 80s*                *Middle 90s*                *2000*                *time*



**Cray C90**            **Cluster of PCs**            **Volunteer computing**
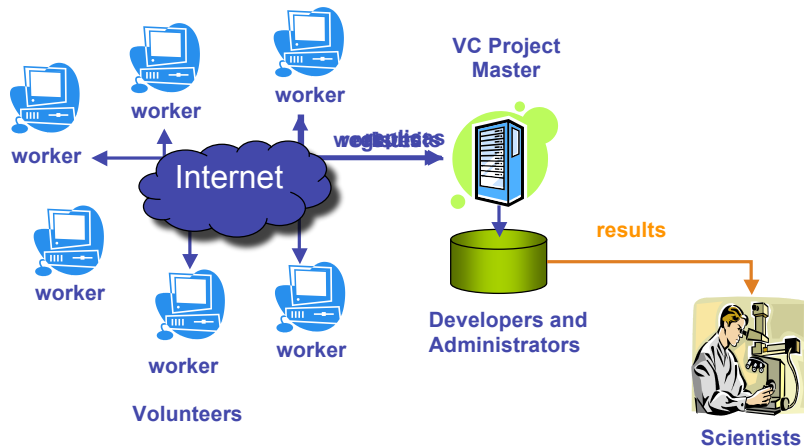
Michela Taufer - 10/15/2007                                                                    6
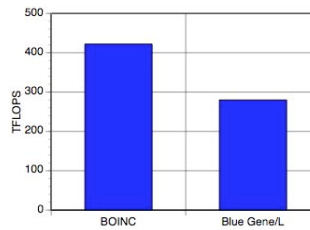
# Volunteer Computing (VC)



BOINC (Berkeley Open Infrastructure for Network Computing) is a volunteer computing middleware that manages volunteer resources

# Strengths and Challenges

- BOINC provides a powerful environment for large-scale simulations:
  - Collectively BOINC projects provide higher sustained compute power than Blue Gene/L (LLNL)
  - Loosely coupled resources
- Resources are anonymous and untrusted
  - Invalid results due to over-clockers and malicious attackers
  - Replication of computation is used to validate results
- Resources are heterogeneous and volatile
  - Divergence of results due to e.g., different architectures and OSs
  - Timed-out results due to e.g., resource sharing among projects and worker disconnections
- Incentives are needed to attract and retain volunteers
  - Credits are assigned for completed computation

# VC Projects

- Predictor@home (P@h)
  - TSRI; protein structure prediction
- Climateprediction.net
  - Oxford,UK; global climate change study
- Rosetta@home
  - U. Washington; protein study
- SETI@home
  - U.C. Berkeley; SETI
- World Community Grid
  - IBM; several life science applications

...and about 30 others

# Predictor@home

# Climateprediction.net



Michela Taufer - 10/15/2007

11

# SETI@home



Michela Taufer - 10/15/2007

12

6

# Docking@Home: DAPLDS Portal

# DAPLDS Overview



ligands + proteins     compute simulations     best docked ligands

- Search for algorithms that provide accurate results
- Search for an computers that can execute these algorithms

# Molecular Characterization

- Size of protein and ligand
- Number of ligand rotable bonds

Ligand in
1tng complex

Ligand in
1hvi complex

- Partial charges and metal atoms in the ligand
- Empirical knowledge of protein flexibility

Protein in
1tng complex

Protein in
1hvi complex

5

# Algorithmic Adaptation

- Implement multi-scale docking models:

    $model_i= f(protein-ligand\ representation,\ solvent\ treatment,\ sampling\ strategy)$

- Cluster protein-ligand complexes in classes based on characteristics:

    $class_l= \{complex_h\}\ with\ h= 1,\ ...,\ N\ and\ N>>1$

- Define adaptive techniques based on simple heuristics and machine learning techniques to match models to classes dynamically:

    $model_{0\ |DA > p} -> \{class_a,\ ...\}$

    ...

    $model_{i-1\ |DA > p} -> \{class_a,\ class_b,\ ...\}$

    $model_{i\ |DA > p} -> \{class_b,\ class_d,\ ...\}$

- Matching based on quantitative values, e.g., free energy of binding and RMSDs

volunteer's computers

solvent treatment

$model_{i+1}$

$model_i$

initial model

sampling strategy

protein-ligand representation

16

8

# Model



| Model protein surface | Model ligand surface |

Alter ligand configuration and orientation

Dock ligand into active protein site

**MD simulation**

Energy minimization

Calculate score

Evaluate candidate solutions

---

# Task Parallelism in VC

- VC projects search for computational solutions that are close to results in Nature
- Searches are partitioned into large numbers of work-units (*WU*s) and implemented via large-scale simulations:

$$simulation \equiv \left\langle WU_j \middle| j = 1,...,N \right\rangle$$

- *WU$_j$* are characterized by a computational model and initial conditions
- Scientists can quantify quality of *WU* results
  - Quantitative measurement of binding affinity: e.g., free energy
- *WU*s are replicated to assure validity of final results
  - *WU* replicas assigned to different volunteer computers
  - Sanity check is performed on the replicas
  - New replicas are automatically generated if errors or time-outs occur

# Task Parallelism



Application

Middleware - BOINC

Volunteer' computers

---

# Taxonomy of VC Work-units

$$WU_{generated} \geq WU_{distributed} \geq WU_{completed} \geq WU_{validated}$$

- $WU_{generated}$ → number of $WU$s generated by the VC master

- $WU_{distributed}$ → number of $WU$s distributed to VC workers so far (incl. in progress, completed, or failed $WU$s)

- $WU_{completed}$ → number of $WU$s completed (and, therefore, returned to the VC master)

- $WU_{validated}$ → the number of $WU$s with reliable results (that can be trusted by the scientists)

# Availability

- An available worker is a productive worker
  - It returns a large number of results over a certain interval of time $\rightarrow$ high throughput
- Returning results may be delayed or timed-out by:
  - Resource sharing among projects
  - Worker disconnection from the network or project
  - Sole use of the worker by the owner
- Availability of $worker_i$:

$$availability_i = WU_{completed\ i} / WU_{distributed\ i}$$

# Reliability

- An available worker is not necessarily a reliable worker
- Results may be affected by:
  - Hardware malfunctions
  - Incorrect software modifications
  - Malicious attacks
- Redundancy computing is used to validate results
  - Several replicas of a *WU* are generated and distributed
  - Sanity check on multiple replicas of the same *WU* is performed
- Reliability of *worker$_i$*:

$$reliability_i = WU_{validated\ i} / WU_{completed\ i}$$

# Characterization of Workers

- Characterize the volunteer's workers:
  - Is the volunteer's worker extensively working for the project?
  - Is the result trustworthy for the scientists?
- Availability and reliability of *worker$_i$*:

$$\text{worker}_i \Rightarrow \left\langle \text{availability}_i \quad \text{reliability}_i \right\rangle$$

- When workers apply for new *WU* replicas, the VC master updates:
  - Worker availability and reliability
  - Availability and reliability of whole worker population
- Availability and reliability thresholds for the entire worker population:

$$\text{threshold}_{\text{availablity}} = f_{i=0,\ldots,\text{num\_workers}} \text{availability}_i$$

$$\text{threshold}_{\text{reliability}} = f_{i=0,\ldots,\text{num\_workers}} \text{reliability}_i$$

Michela Taufer - 10/15/2007

23

# Classification of Workers



**Availability**

Threshold $_{availability}$

100%

| High availability Low reliability **HA/LR** | High availability High reliability **HA/HR** |
| Low availability Low reliability **LA/LR** | Low availability High reliability **LA/HR** |

100%

**Reliability**

Threshold $_{reliability}$

Michela Taufer - 10/15/2007

24

# Scheduling Policies

- First-come first-serve:
  - Replicas are assigned to any worker - similar to having availability and reliability thresholds both equal to 0
- Fixed thresholds:
  - Replicas are assigned to workers that have availability and reliability above fixed thresholds
- Variable thresholds:
  - The assignment of replicas reflects runtime changes in the volunteer community
- Compound of simple heuristics driven by genetic algorithms
  - Population of chromosomes → conditional rules for work-unit assignments
  - Mutation and crossover on parameters and metrics that characterize workers in chromosomes

- Homogenous redundancy
  - Send replicas to numerically equivalent workers, demand identical answers (*M. Taufer et al., HWC'05*)

---

# Evaluation

- Challenges in testing new policies on VC projects:
  - Designers cannot predict the effect of their decisions
  - Time to measure and compare performance would be too long
  - Problems due to testing might upset volunteers
  - Every experiment is unique and unrepeatable
- Testing on a reduced-size system is not indicative of real VC projects
- Need for a simulation environment that allows project designers to tailor the simulation to the volunteer computing community
  - Simulator of BOINC Applications or SimBA

# SimBA

- SimBA (Simulator of BOINC Applications) is a discrete event simulator written in Python, that models the behavior of a BOINC master (*T. Estrada et al., eScience'06*)

- Entities are:
  - work-unit generator
  - worker generator
  - worker
  - work-unit
  - work-unit replica
  - sampler (monitor)

- Events are:
  - generate work-unit
  - generate worker
  - request instances
  - generate instances
  - determine instance output
  - check if simulation is over

Michela Taufer - 10/15/2007

27

# BOINC vs. SimBA



28

14

# SimBA User Interface



• Text-based output that gives a detailed overview of each step the simulator is performing

• Real-time graphs that show the current *WU*s generated, distributed, returned, and validated

Michela Taufer - 10/15/2007                                                                 29

---

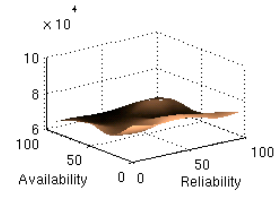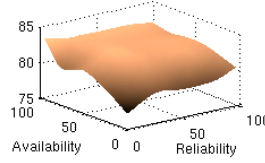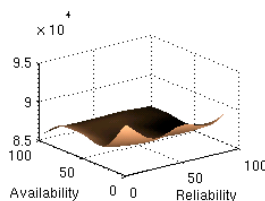# Samples of SimBA Outputs

**P@h**
**Monte Carlo application**

**World Community Grid - IBM**
**F@h                              Defeat Cancer**



Michela Taufer - 10/15/2007                                                                 30

15

# A Case Study: Description

- Traces from an existing VC project: P@h
- Minimum number of replicas per *WU*: 3
- Number of workers in P@h traces: 14,000
- Duration of VC project: 12 days
- Scheduling policies:
    - First-Come First-Serve
    - Fixed thresholds - 75% availability and 75% reliability
    - Variable thresholds - availability and reliability change at runtime based on number of pending *WU*s

Our Question: What scheduling policy does provide us with a better throughput and resource utilization over 12 days of VC project?

# A Case Study: Results

P@h performance of different scheduling polices based on SimBA results

| P@h | FCFS | Fixed Thresholds (75% - 75%) | | Variable Thresholds | |
|---|---|---|---|---|---|
| Generated WU | 78,658 | 78,886 | +0.2% | 78621 | ~0% |
| Generated WU replicas | 284,140 | 252,258 | -11.2% | 253,611 | -10.7% |
| Error replicas | 38,491 | 13,802 | -64.1% | 14,424 | -62.5% |
| Valid WU (throughput) | 70,948 | 71,201 | +0.4% | 72,929 | +2.8% |
| Avg. replicas per WU | 3.6 | 3.2 | -11.1% | 3.2 | -11.1% |

- Dynamically adaptable thresholds keep track of runtime changes in the VC community (*M. Taufer et al., PCGrid'07*)
    - Average replicas per *WU* drops from 3.6 to 3.2 for thresholds based policies
    - Previously ineligible volunteer computers are regained for the project

# Conclusions

- Volunteer computing is a powerful paradigm for applications that use task parallelism
- Interesting research topics in volunteer computing include:
    - Improve accuracy of models for the simulations of protein-ligand interactions
    - Study of effective scheduling polices for the selection of volunteer's computers
- Several applications in science can benefit from effectives volunteer computer systems
    - Geology
    - Biology
    - Climate studies
- Computer scientists working in this field are exposed to interdisciplinary research
- If you want to know more, contact me at taufer@udel.edu or come by my office at Smith Hall 406

Michela Taufer - 10/15/2007

33