

Next Generation Bioinformatics Tools



Fall 2012

Day 5 - Next Steps in Bioinformatics - Translational Bioinformatics, HPC, Security/Privacy, and the Cloud

Hesham H. Ali

UNO Bioinformatics Research Group
College of Information Science and Technology



What is next?

- Translational Bioinformatics
- Biomedical Informatics
- High Performance Computing
- Security/Privacy and Cloud Services
- Impact on Health Care



Biomedical Informatics?

- Bioinformatics
- Biomedical Imaging
- Health Informatics
- Medical/Clinical Informatics
- Public Health Informatics



Translational Informatics – Case Study in Aging Research



Aging Research

- Bioinformatics – The Genomic Aspect
- Wireless Technology – The Mobility Aspect
- Virtual Environments – Data Gathering
- Public Health Informatics – Data Integration/Analysis



Case Study in Aging

- With aging, certain behaviors decrease
 - Eating, drinking, activity levels
- Observed gene expression changes in the hypothalamus
 - Can we capture these expression changes?
 - Can we correlate these changes to behavioral decreases?
- Goal: Identify temporal biological relationships
 - Progression of disease
 - Effect of pharmaceuticals on systems of the body
 - Aging



Case Study in Aging

- 5 sets of temporal gene expression data

Strain	Gender	Tissue Type	Ages
BalbC	Male	Hypothalamus	Young, mid-age, aged
CBA	Male	Hypothalamus	Young, mid-age, aged
C57_J20	Male	Hypothalamus	Young, aged
BalbC	Female	Hypothalamus	Young, aged
BalbC	Female	Frontal cortex	Young, aged



Hubs

- **Hub:** a high-degree node in a network
- Node degree in filtered correlation networks follows power-law relationship
- Few nodes with high degree
- High degree nodes → highly essential

Albert et al 2005

Bergmann et al 2004
Carlson et al 2006



Hub Lethality

- Young Male BalbC Mouse
 - 12/20 hubs tested for *in vivo* knockout
 - 8/12 lethal phenotype pre-/peri-natally
 - 4/12 non-lethal but system-affecting
 - 0/12 no observed phenotype
- Aged Male BalbC Mouse
 - 11/20 hubs tested for *in vivo* knockout
 - 7/11 lethal phenotype pre-/peri-natally
 - 3/11 non-lethal but system-affecting
 - 1/11 no observed phenotype (Aldh3a1)



Hub Lethality

- Young Male BalbC Mouse
 - 12/20 hubs tested for *in vivo* knockout
 - 8/12 lethal phenotype pre-/peri-natally
 - 4/12 non-lethal but system-affected:
 - Hspa1a: cellular, growth/size, homeostasis
 - Dapk1: cellular, renal/urinary
 - Ffar2: Increased susceptibility to colitis, asthma, arthritis



Hub Lethality

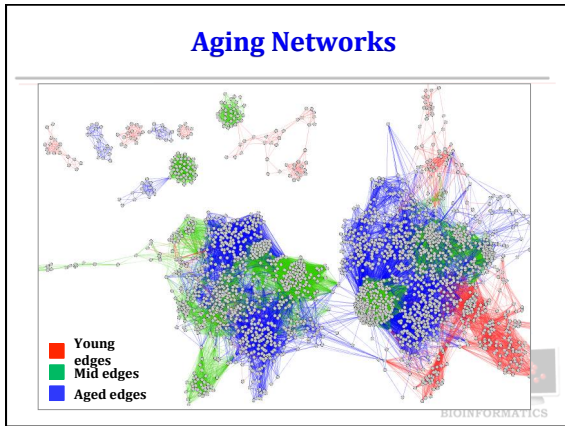
- Aged Male BalbC Mouse
 - 11/20 hubs tested for *in vivo* knockout
 - 7/11 lethal phenotype pre-/peri-natally
 - 3/11 non-lethal but system-affected:
 - Btn1a1: impaired lactation, impaired lipid accumulation in mammary gland
 - Bcl2l11: die later in life from auto-immune kidney disease
 - Rag2: arrested development of T and B cell maturation
 - 1/11 no observed phenotype (Aldh3a1)

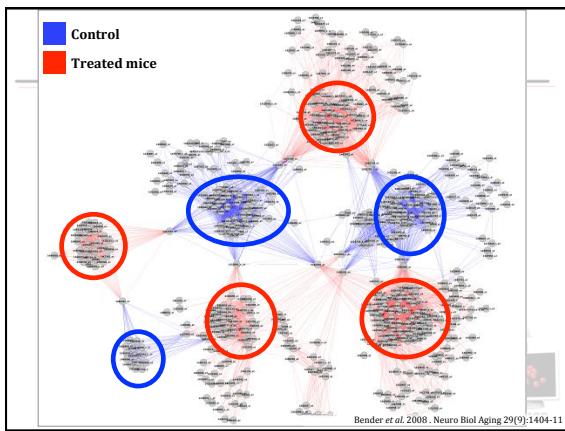


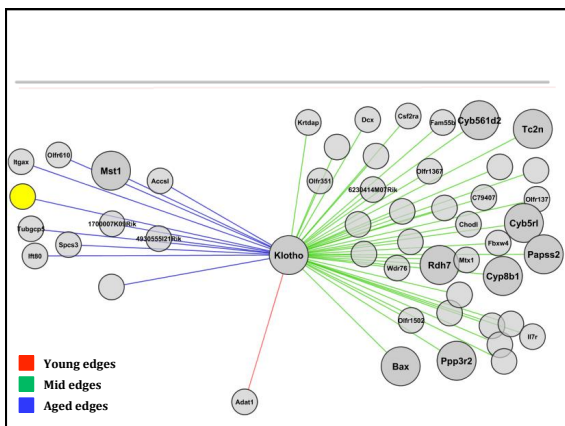
Hub Lethality

- Common hubs among top 250?
 - 7 genes found in common
 - 1/7 tested for pre-/peri-natal lethal knockout
 - Cdk2: cyclin dependent kinase 2
 - Young degree: 656
 - Aged degree: 1929









Centrality: Integrated Networks

- Networks representing multiple types/states
- Does centrality identify interesting nodes?
- Case study: aging



Structure Types

Elements (Nodes):

- Betweenness
- Closeness
- Degree
- BC: Highest betweenness + closeness
- CD: Highest degree + closeness
- BD: Highest betweenness + degree
- BCD: Betweenness + closeness + degree

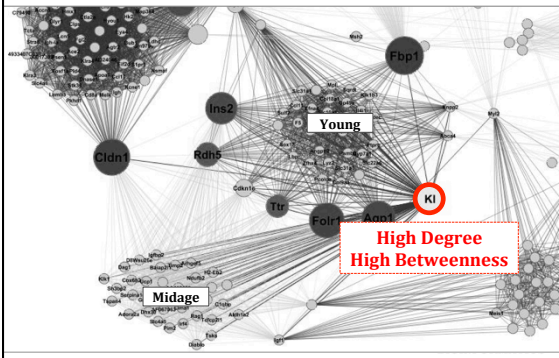
Our Focus

Subsystems (Relationships, groups) :


- Clusters, cliques
- Pathways
- Loops/cycles



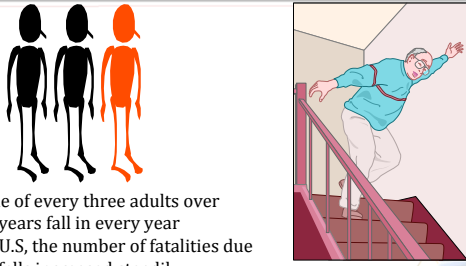
High BD Node: Klotho




Mobility Monitoring: Detecting and Predicting Falls



How serious is the problem of falls?




- One of every three adults over 65 years fall in every year
- In U.S, the number of fatalities due to falls increased steadily from 14,900 in the year 2000 to 17,700 in 2005.



Falls and Associated Problems

- Falls are the leading cause of accidental deaths in the United States among people over the age of 75
 - the number of fatalities due to falls increased steadily from 14,900 in the year 2000 to 17,700 in 2005.
- Nebraska's over age 65 population is 13.3% versus 12.4% for the national average.
 - Generally speaking, the more rural the area, the higher the percentage of older adults.
 - In Nebraska, approximately 78% of those hospitalized for fall related injuries were 65 years and older.

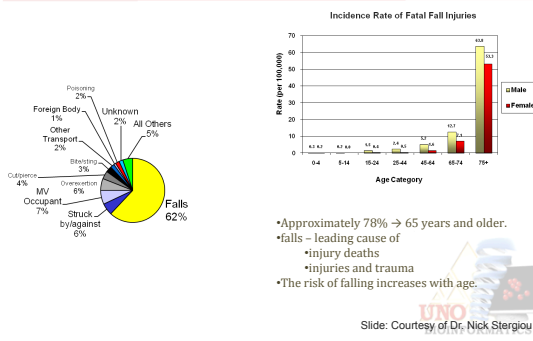


Older Adults in Nebraska

- Nebraska's over age 65 population is 13.3% versus 12.4% for the national average.
 - Generally speaking, the more rural the area, the higher the percentage of older adults.
- For people age 65 and older, falls are the leading cause of injury death.
 - Among older adults, falls are the leading cause of injury deaths and the most common cause of injuries and hospital admissions for trauma.



Significance of Falling



Traditional Gait Monitoring Methods

- Expensive
- Uncomfortable
- Limited mobility
- Complicated



http://www.vt.edu/spotlight/impact/2008-06-09_locomotion/2008-06-09_locomotion.html

Laboratory-based Gait Monitoring

- Expensive
- Uncomfortable
- Limited mobility
- Complicated



http://www.vt.edu/spotlight/impact/2008-06-09_locomotion/2008-06-09_locomotion.html

UNO BIOINFORMATICS

Wireless Sensor Based Mobility Monitoring

- Inexpensive
- Comfortable
- High mobility
- Simple



UNO BIOINFORMATICS

Wearable Fall Detection System - Philips Life Line System


There are several systems available on the market to detect falls



UNO BIOINFORMATICS


Goal of the Mobility Project

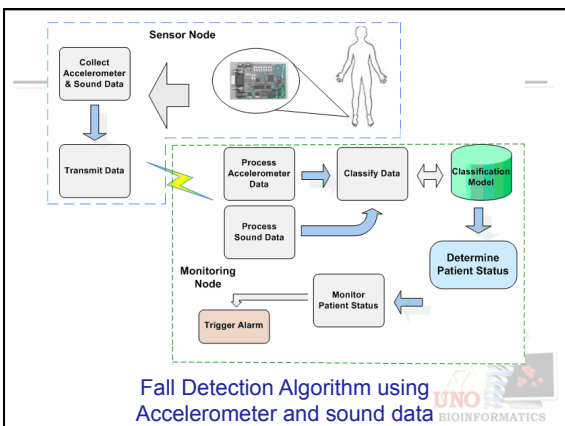
- Fall Prediction using Mobility Profiles
 - The system will identify anomalous movement and patterns that usually result in a fall or injury,
 - We would be able to take preemptive measures when such a pattern is detected, in order to reduce the occurrence of falls and prevent fall-related injuries.
 - We will develop an index that enables health care providers to determine how likely people are to fall
- Mobility Profile
 - Patient wearing a 3D-accelerometer will be monitored 24/7.
 - A complete mobility profile will be available for patients and care providers.



Four Project Phases


- Phase I Fall Detection
 - achieved over 95% of fall detection rate
- Phase II: Classification of ADLs
 - Running, Walking, Jumping, Stair Climbing, Standing, Sitting, and Lying.
- Phase III: Construction of Mobility Profiles
- Phase IV: Fall Prediction based on mobility profiles





Experiment

- Accelerometer
 - Impact detection
 - (unit: gravity)
- Gyroscope
 - Measure rate of rotation
 - (unit: degrees per second)




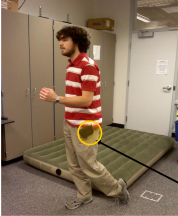


(shimmer-research.com)

UNO
BIOINFORMATICS

Phase I: Fall Detection

Accelerometer-based fall detection

- Measure acceleration in three orthogonal directions.

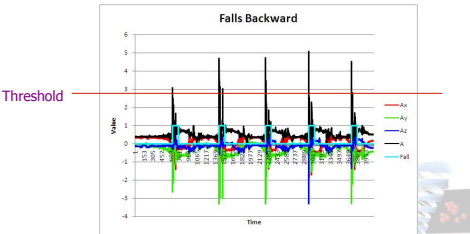
Shimmer Device

UNO
BIOINFORMATICS

Phase I: Fall Detection

Accelerometer-based fall detection


- Determine an acceleration threshold.
- Detect fall.



UNO
BIOINFORMATICS

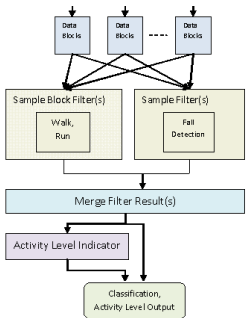
Phase II: Classification of ADLs


- Many Activities of Daily Living (ADLs) can be classified by analyzing the real-time acceleration data collected from sensors.
- Key Metrics
 - Inclination Angle
 - Standard deviation
 - Skewness
 - Signal Magnitude Area



Phase II: Classification of ADLs

An activity can be classified as walking or running based on the magnitude and frequency of a peak in Fourier transform.

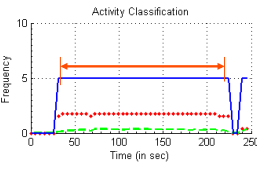




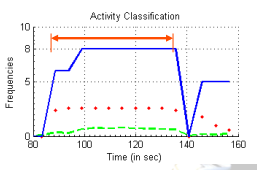
Phase II: Classification of ADLs

- Walking vs. Running based on the magnitude and frequency of a peak in Fourier transform.


Legend: — Activity Index, ●●●● Peak Frequency, —●●●● Peak Frequency Magnitude



walking



running



Phase II: Classification of ADLs

- Using the Fourier transform-based approach, differentiating some activities with similar periodic motions is not easy.
 - such as climbing stairs and walking.

The graph shows frequency components over time. There are two distinct periodic activity windows: one between 100-200 seconds and another between 300-400 seconds. Each window shows a primary peak at approximately 5 Hz and a secondary peak at approximately 1 Hz. The UNO BIOINFORMATICS logo is in the bottom right corner.

Inclination Angle

- Inclination angle is the measure between the x and y axis. It is assumed that if this value is around 180°, then the person would be standing.

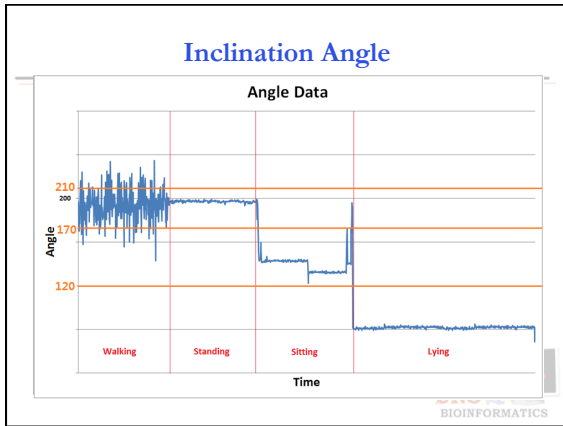
Standing Position diagram: A person is shown with a vertical z-axis. The inclination angle Φ is measured from the z-axis. Thresholds are marked at $\Phi = 0^\circ$ (High Lying Threshold), $\Phi = 90^\circ$ (Low Lying Threshold), $\Phi = 180^\circ$ (High Standing Threshold), and $\Phi = 270^\circ$ (Low Standing Threshold).

Lying Position diagram: A person is shown lying on their back with a vertical z-axis. The inclination angle Φ is measured from the z-axis. Thresholds are marked at $\Phi = 0^\circ$ (High Lying Threshold), $\Phi = 90^\circ$ (Low Lying Threshold), $\Phi = 180^\circ$ (High Standing Threshold), and $\Phi = 270^\circ$ (Low Standing Threshold).

- Can now differentiate standing, sitting, and lying.

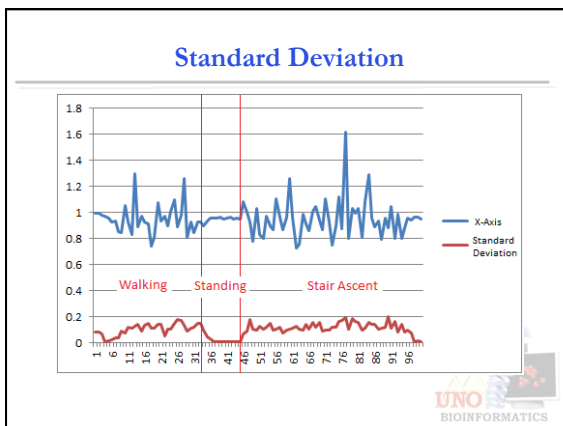
3D Acceleration Data vs. Activities

The graph shows 3D acceleration data (x, y, z axes) over time, divided into four activity segments: Walking, Standing, Sitting, and Lying. Gravity is on the y-axis (ranging from -1 to 2) and Time is on the x-axis. The x-axis (blue) shows high-frequency oscillations during walking and low-frequency oscillations during standing. The y-axis (red) shows a step change from ~1.0 to ~0.5 during sitting. The z-axis (green) shows a step change from ~0.0 to ~-0.5 during lying. The UNO BIOINFORMATICS logo is in the bottom right corner.



Standard Deviation

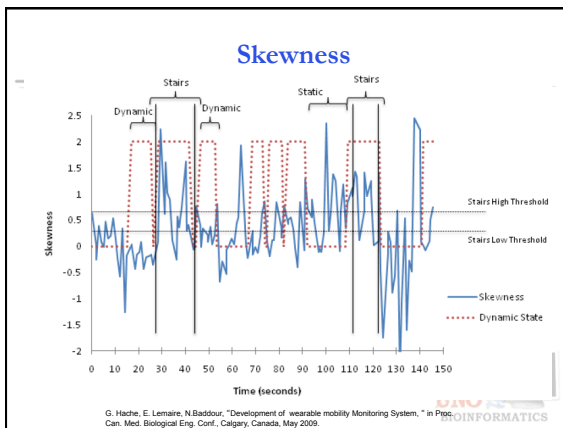
- Standard deviation of the x-axis acceleration helps in determining if the current mobility state is dynamic or static
 - Standard deviation measures the variability of data from the mean. Dynamic data will have measurably more variability than static data.
- When used with angle measurement, can now differentiate standing from walking/running.



Skewness

- Skewness helps in determining if the current mobility state is going up or down stairs.
 - Skewness measures the asymmetry of the distribution of x-axis acceleration values. It is assumed that going up/down stairs will produce data that has greater asymmetry than walking.
- When used with angle and standard deviation, can now differentiate walking/running from going up/down stairs.

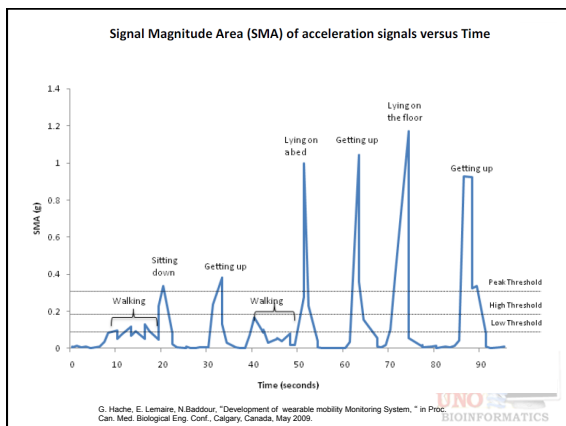




Signal Magnitude Area


- Measures amplitude and duration variation in the acceleration signal.
 - Assumed that amplitude and duration variation will be greater when the intensity of the activity changes. As such, SMA values will be greater when changing state as opposed to not changing state. Ex: getting out of bed vs walking
- When used with the prior four measurements, SMA will help differentiate transitions from other dynamic mobility states.






Phase III: Mobility Profiles

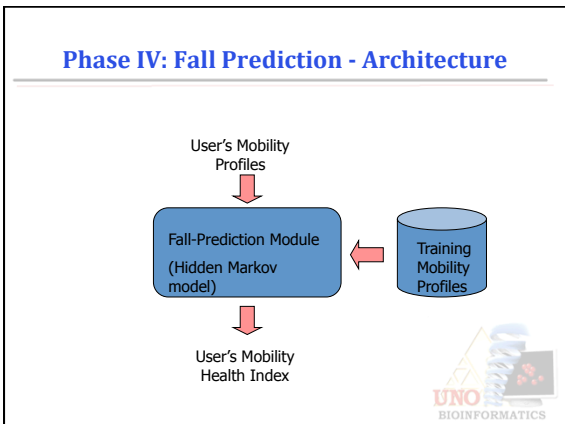
- Mobility Profile (between the given start-time and end-time)
 - Total # of steps
 - Average # of steps per second
 - Activity level with different precisions
 - # of rooms traveled (will be added)



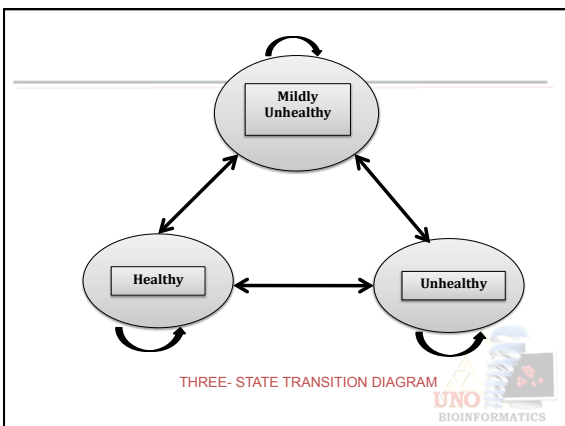
Phase IV: Fall Prediction

- Fall can injure the elderly in large scale.
- 10-15% falls cause some serious physical injury in older people.
- The early prediction of fall is an important step to alert and protect the subject to avoid injury.
- We employ Hidden Markov Models for detection and prediction of anomalous movement patterns among the human subjects.





- ### Prediction Model
- Develop a HMM model with
 - 3 states (Healthy, Unhealthy, Mildly Unhealthy)
 - 3 parameters (#steps moved, #rooms visited, # movement in arms)
 - Entire dataset was split into
 - Train data
 - Test data



Probability Calculations

- The model parameters for a HMM are generated from:
 - state transition probabilities
 - $a_{kl} = P(\pi_i = l \mid \pi_{i-1} = k)$, which is probability from state k to state l
 - emission probabilities
 - $e_l(b) = P(x_i = b \mid \pi_i = l)$, which is probability distribution over all the possible output symbols b for each state l .



Training

- For each observation the trained model parameters were calculated using maximum likelihood approach
- $a_{kl} = \frac{A_{kl}}{\sum_l A_{kl}}$
- $e_l(b) = \frac{E_l(b)}{\sum_b E_l(b)}$



Predicting Hidden State

- With the trained model parameters
 - predict hidden state path for a new sequence of observations using Forward-Backward (Posterior) algorithm.
 - predicts the hidden, probable state path



Assessment

- Performance evaluation
 - Accuracy
 - Sensitivity
 - Specificity

$$\text{specificity} = \frac{\text{number of True Negatives}}{\text{number of True Negatives} + \text{number of False Positives}}$$

$$\text{sensitivity} = \frac{\text{number of True Positives}}{\text{number of True Positives} + \text{number of False Negatives}}$$

$$\text{accuracy} = \frac{\text{number of true positives} + \text{number of true negatives}}{\text{numbers of true positives} + \text{false positives} + \text{false negatives} + \text{true negatives}}$$


UNO
BIOINFORMATICS

Translational Bioinformatics

- Correlation between genetic networks and mobility profiles
- Impact of mobility – or lack thereof – on biological networks
- How tied the mobility profiles with the biological network?
- Premature aging research

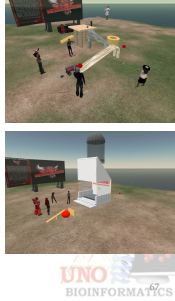


Visualizations and Virtual Environments



Avatars, Metaverses and Virtual Worlds

- Metaverses [Bainbridge, 2007] are 3-dimensional virtual worlds that are vivid, synthetic spaces where people can interact in rich and increasingly realistic ways
 - Metaverses provide vividness, interactivity, telepresence, immersion
 - Virtual Worlds such as World of Warcraft, SecondLife.com and There.com are instantiations of metaverses
 - Human actors interact as *avatars* with each other and with software agents
 - Uses the metaphor of the real world but without its physical limitations
 - Potential for simulating "face-to-face" interactions











High Performance Computing (HPC) and Bioinformatics



Explosion of HPC

- Increase in data centers due to internet explosion and now cloud computing.
- This past week: Grand opening of Scott Data Center
- Data Centers Models



HPC and Life Sciences

- HPC is fast becoming a key ally for the Life Sciences because of its ability to deal with large amounts of data.
- HPC typically have huge computing resources and are able to provide results quickly in relative terms.
- Life Sciences field of studies such as bio-informatics have large amounts of data that need to be processed.
- Most algorithms used in this field lend themselves to an HPC environment as they are highly parallelizable.
- Hence the natural alliance between HPC & Life Sciences.



HPC and Biological Networks

- Network creation: 2 weeks on PC
 - 10 hours in parallel, 50 nodes
 - 40,000 nodes = 800 million edges
- Network analysis: Best in parallel
 - Only 3% of entire genome forms complexes
- UNO Sapling Cluster
- Holland Computing Center: Firefly



Challenges Associated with Network Analysis

(1) Biological networks can be massive in size
 Supercomputing access may be limited
 Biological network knowledge may be limited.

(2) Noise within the network is likely
 Noise within the network cannot be ignored.

How to address these issues: Network filters

- Reduce network size
- Maintain biological signal
- Improve upon biological signal?



What about Energy?

- Energy generation, conservation & efficiency are and will continue to be key growth drivers.
- We constantly make decisions that influence energy utilization in our daily lives.
- Examples:
 - Use of laptops, We have power management schemes.
 - When should we recharge our cars (night time costs are lower).

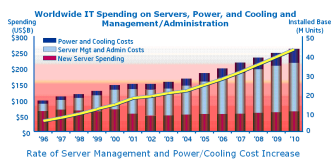


Energy as a key business driver

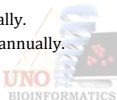
- Business are finding out that Energy costs have to be controlled to be profitable.
- Major business are locating their Data Centers where Energy Costs are lower
 - Google (Iowa & Oklahoma)
 - Yahoo (Nebraska)
 - Microsoft (Des Moines, Iowa)
- Even local communities are moving towards more sustainable models (Omaha's new Energy efficient Learning Community Center models).
- Clearly "**Energy**" is becoming a key business driver.



HPC and Energy Costs




- US Data centers consumed 5 MKW of energy.
- Equivalent to five 1000 MW power plants.
- Energy utility bills amounting to \$2.7 billion annually.
- World consumption estimated to cost \$7.2 billion annually.



Energy as a key design parameter (Lifetime versus Costs)

- The lifetime of mobile devices depends on energy levels they are used at (Power management, modes).
- Energy is changing the way devices are built.
- Energy is reflected as a key design parameter that the industries involved in building these devices will look at.
- Some of the technologies such as DVS have not completely been commercialized.

Hardware	Silicon Process Technology
	Chip Technology
	Power Management
Software	Operating System
	Applications



HPC and Energy Awareness

- As Energy becomes a key business driver:
 - Energy Aware Scheduling (EAS) has a key role to play in balancing the divergent goals of minimizing energy utilization and maximizing performance.
- There is a need for an EAS Model that addresses:
 - Small End Devices (Sensors, Cell-phones) to
 - High Performance Computing (HPC) environments and
 - Everything in between.
- The need to carefully develop a parallel model in a HPC environment, based on our understanding of the
 - Data, and
 - Application/Domain



Energy Awareness

- There is a major energy related problems associated with High Performance Computing
 - Should we be alarmed? Concerned?
 - Yes - Yes
- Life Sciences and HPC
 - Excitement - Opportunities - New line of Computing Research
 - Potential energy crisis
- The need to carefully develop a parallel model in a HPC environment, based on our understanding of the
 - Time to act is now
 - Application Domain



Computational abuse all over again

- First Abuse:
 - Blamed on the engineers: 70s – 80s
- Who is to blame now?
 - Bioscientists - High end domain experts
 - Personal clusters – enormous data levels – new devices like high throughput sequencers
- The price now is a lot higher



Is it justified?

- It has been estimated that over 60% of generated data by microarrays is unnecessary
- Why so much data? Because we can
- Problem: Energy consumption + data mining
- Examples:
 - BLAT Package: For some dataset, running time reduced from 17 hours to 3 minutes using 100 cluster nodes
 - Assembly of genomic short reads
 - Transcription Factor Binding Sites: Finding conserved motifs by brute force



Solutions?

- Some centralization is not too bad
- Better collaboration models – involvement of computational scientists early in the process
- Better understanding of current complex problems
- Better appreciation of computational facilities – unmask the clouds
- Better appreciation of biological (domain specific) problems
- Differentiate between home computing and scientific computing



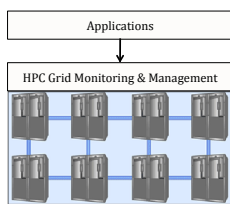
Technical Solutions

- Better application specific parallelization – custom solutions lead to better performance
- Smarter input data: better user level utilization plus integrated domain knowledge with computational tools
- Better scheduling model: multi-layer dynamic scheduling solution



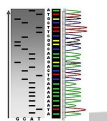
EAS in High Performance Computing

- The Energy Aware HPC layer is responsible for the scheduling of the domain specific applications on the grid with energy minimization as its main objective.
- Other parameters that will be adjusted include number of nodes, schedule length, etc.
- Application domain used will be bioinformatics.
- Computing domain used will be the Holland Computing Center.



HPC and Sequence Comparison

- The bioinformatics domain is rich in applications that require extracting useful information from very large and continuously growing sequence of databases.
- Automated techniques such as DNA sequencers, DNA microarrays & others are continually growing the dataset stored in large public databases such as GenBank.
- Sequence comparison remains the fundamental problem in bioinformatics.
- Most methods used for analyzing genetic/protein data have been found to be extremely computationally intensive (such as BLAST, BLAT, etc), providing motivation for the use of powerful computers or systems with high throughput characteristics.



A Proposed Solution

- To move from a simple speedup to the realm of energy awareness.
- Now energy awareness, places a new constraint on the scheduling system.
- The scheduling policy still has to be traditional performance focused and energy aware at the same time.
- The goal is to find the right harmony between these two, slightly divergent goals.
- The crucial question which follows is how one achieves the right balance between these two differing optimization criteria.



BLAT and HPC

- BLAT is a BLAST-like sequence comparison package.
- In this case, a set of files, each contains a set of motif queries to be searched for in each human chromosome.
- The human chromosome files used for these experiments were downloaded from the UCSC Genome bio-informatics website.
- The chromosomal sequences were assembled by the International Human Genome Project sequencing centers.
- We used MPI (GNU) to parallelize the runs on multiple nodes, which was a configurable parameter.
- Our experiments used sequences gathered from researchers at UNMC (University of Nebraska Medical Center) and parallelize the runs to study the performance characteristics under three different conditions. For our tests we used 24 query sequences.



Technical Solutions

- *Better application specific parallelization – custom solutions lead to better performance*
- Smarter input data: better user level utilization plus integrated domain knowledge with computational tools
- Better scheduling model: multi-layer dynamic scheduling solution



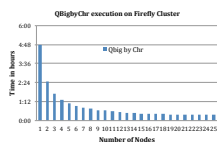
Exp1: All query sequences per chromosome



- When node is 1 we get a total execution time of 6:16 (hh:mm).
- When number of nodes = 25 we get a total execution time of 0:28,
 - which is a speedup of 13.
- Note however that when we vary nodes from 20 – 25, we do not see any additional gains,
 - this is because we have already used the inherent slack in the schedule and there are no additional gains to be made by increasing the number of processors.



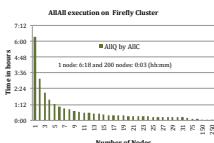
Exp2: Merged query sequences per chromosomes



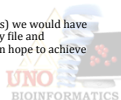
- When node is 1 we get a total execution time of 4:45 (hh:mm).
- When number of nodes = 25 we get a total execution time of 0:22,
 - which is a speedup of 12.
- Note however that when we vary nodes from 20 – 25, we do not see any additional gains,
 - this is because we have already used the inherent slack in the schedule and there are no additional gains to be made by increasing the number of processors.



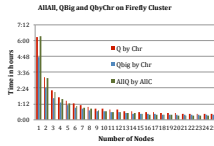
Exp3: All query files against all chromosome files



- When node is 1 we get a total execution time of 6:20 (hh:mm).
- When number of nodes = 25 we get a total execution time of 0:16,
 - which is a speedup of 22.
- With nodes = 150 we get execution time of 0:04,
 - which is a speedup of 86.
 - If we had 1176 processors (24 query files times 49 chromosome files) we would have seen this go down to the max execution for one combination of query file and chromosome file out of the 1176 combinations this is the best we can hope to achieve



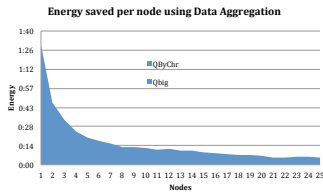
Comparison of the three approaches



- When node = 1, we see that the merged query approach is better than the other two approaches.
 - Note that this true when nodes 1 - 5.
 - After 5 nodes, the "All Query All Chromosome" approach gives us better results.
- With nodes between 25 - 30, we get twice the speedup with the "All Query All chromosome" approach.
- One can also note that the Merged Query approach always performs better than the Query by Chromosome approach.



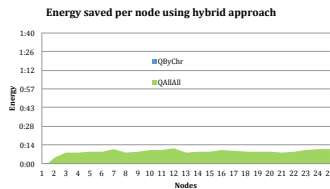
Energy saved using Data aggregation



- The chart shows energy saved using data aggregation (Qbig) approach.
- The baseline is the Database segmentation (QByChr) approach.



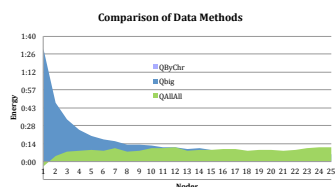
Energy saved using Hybrid approach



- Hybrid approach is combination of query segmentation and database segmentation.
- The chart shows energy saved using hybrid (QAllAll) approach.
- The baseline is the Database segmentation (QByChr) approach.



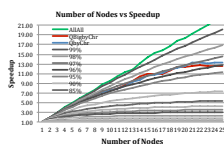
Comparison of Energy saved



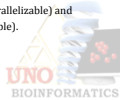
- The chart shows the data aggregation approach saves more energy when using about 10 nodes compared to the hybrid approach.
- When nodes used is 11 - 15 they both save about the amount of energy.
- However, the hybrid approach saves more energy when nodes used is > 15.



Number of Nodes versus Speedup



- Let us try and answer the question how parallelizable is the program?
- In-order to answer this question we try and plot the speedup for each experiment and place these by the curves based on Amdahl's Law.
- From the figure above we can conclude that
 - The QBigbyChr and QbyChr have a speedup of around 25 times (97% parallelizable) and
 - The AllAll approach has close to 100 times the speedup (99% parallelizable).




Summary of Results

- We found that the BLAT program is highly parallelizable and has a speedup of 99%.
- The experiments suggests that the merged query approach and the hybrid approach of all query segmentation and database segmentation consistently performs better than just the database segmentation approach.
- We also find that
 - When only about 5 nodes are available, it is better to use the merged query approach,
 - For number of nodes 6 - 10, we would be better off using the merged query approach, and

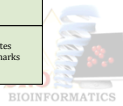


Energy Aware Scheduling (EAS) Model: A Dynamic Multi-Layer Approach




Stages in accomplishing Energy Efficiency

Hardware	Silicon Process Technology	<ul style="list-style-type: none"> • Second generation strained silicon • Improved interconnects
	Chip Technology	<ul style="list-style-type: none"> • Dynamic sleep transistors • Demand based switching • On-Die voltage regulation, Power Gating & Macro Fusion • Multi-core and clustered Micro-Architecture
	Power Management	<ul style="list-style-type: none"> • Voltage Regulation Technology • Improved Display Power Specs • Thermal design for advanced heat-sync technology
Software	Operating System	<ul style="list-style-type: none"> • Application code multi-threaded and multi-core ready • Power monitoring and analysis tools • Optimizing code for reducing CPU clock cycles • Energy Aware Scheduling of Application Tasks
	Applications	<ul style="list-style-type: none"> • Developing power conscious device drivers • Tuning OS for less interference with CPU's low-power states • Energy Aware Scheduling of application based on benchmarks



HPC Management

- High performance computing also requires sound management of energy at the node level
- Our intention is to use energy as one of the optimizing criteria
- The application domain to be used will be a bio-informatics domain which has many interesting problems that can benefit from Energy Aware Scheduling during computational analysis of sequences, alignment



Next Step: Cloud Computing – Lifting the Veil

The diagram illustrates the concept of cloud computing. On the left, four stylized human icons (a woman, a man, a person in a hard hat, and another man) have lines connecting them to a central cloud. Inside the cloud are logos for major technology companies: IBM, Google, Microsoft, Amazon, and Facebook. Below the cloud is the UNO Bioinformatics logo, which includes a stylized DNA helix and a computer monitor.

Technical Solutions

- Better application specific parallelization – custom solutions lead to better performance
- *Better scheduling model: multi-layer dynamic scheduling solution*
- Smarter input data: better user level utilization plus integrated domain knowledge with computational tools

The slide lists three technical solutions. The first is about application-specific parallelization. The second is a multi-layer dynamic scheduling solution. The third is about smarter input data, combining user-level utilization with domain knowledge. The UNO Bioinformatics logo is in the bottom right corner.

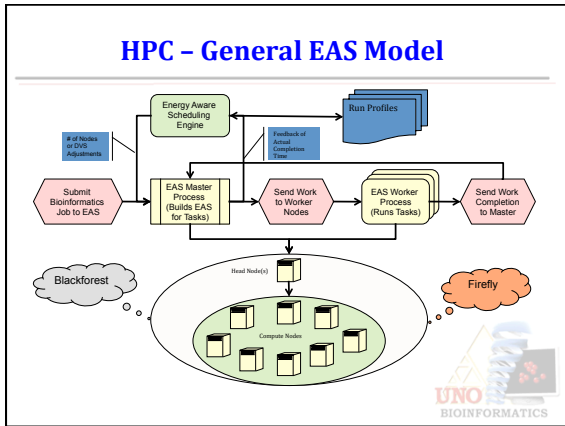
Energy Management – Research Vision

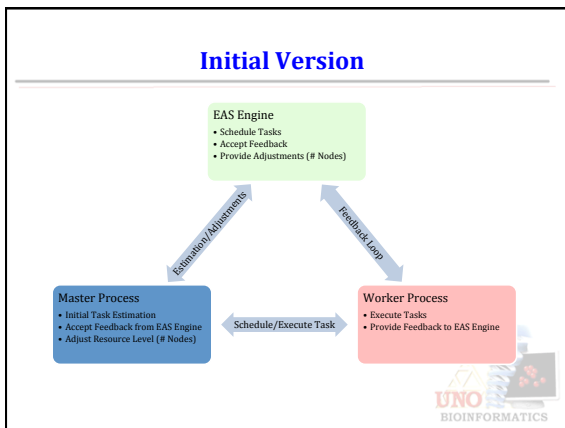
Energy Adjustment @ Cloud Level
• Using Energy Index for Datacenters

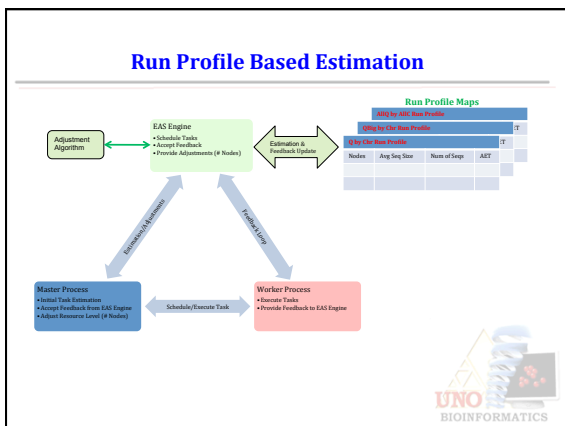
Energy Adjustment @ HPC Level
• Using Adjustments to # of Nodes

Energy Adjustment @ Node Level
• Using Dynamic Voltage Scaling (DVS)

The slide outlines a research vision for energy management across three levels: Cloud, HPC, and Node. Each level includes a specific strategy: Energy Index for Datacenters at the cloud level, adjusting the number of nodes at the HPC level, and Dynamic Voltage Scaling (DVS) at the node level. The UNO Bioinformatics logo is in the bottom right corner.



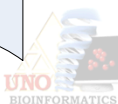




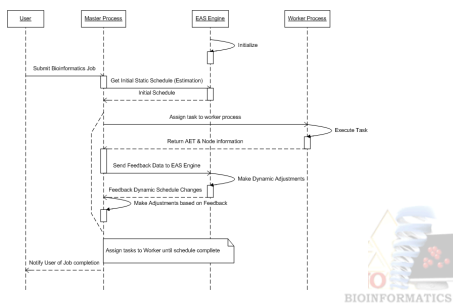
Algorithm Outline

```

Input: AET & Node Information for Task just completed
if (AET T_i < EET T_i)
  // Adjust down
  if (min time of task T_i < release time of task T_{i+1}) {
    Update the scaling level to absorb the slack
    Adjust Down();
  }
  Update Node Map (x_i // node, timing and sequence information
  if (original task periodicity shows a task T_i is available earlier) and
  (start time of T_{i+1} - T_{i+1} >= EET of T_{i+1}) or (T_{i+1} - T_{i+1} >= EET T_i + EET T_{i+1})
  Schedule task T_{i+1}
  else
    Wait
}
Else // Adjust up
  Update Node Map (x_i // node, timing and sequence information
  Adjust Up();
  Schedule task T_i, and remove it from Phase I schedule
}
Execute the task
  
```

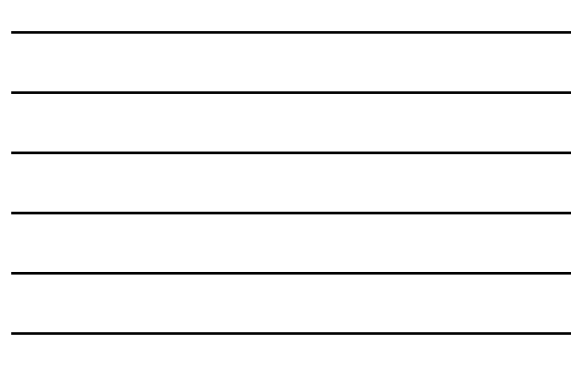


Algorithm - How it works?



Run Profile → Node Map → Adjustment

QUERY FILES	Job size (k/MB)	Job size (MB)	Job size (GB)	# of exp.	# Nodes	Min/Max	QoS	QoS Ctr
query/MC1_chr11.m	3311.1705	10891376	465	2993	2	0.20	4.45	6.10
query/MC1_chr12.m	2738.846	7823704	464	3127	2	0.00	2.06	3.11
query/MC1_chr13.m	1772560	5866979	464	3820	3	0.00	1.88	2.43
query/MC1_chr14.m	1461240	4664152	460	2965	4	1.04	1.19	1.43
query/MC1_chr15.m	1722390	5469172	486	3541	5	1.04	0.53	1.11
query/MC1_chr16.m	1273700	5020052	471	2760	7	0.01	0.61	1.00
query/MC1_chr17.m	1863880	6144152	460	4054	8	0.48	0.44	0.57
query/MC1_chr18.m	1470244	4598916	460	3220	9	0.43	0.31	0.52
query/MC1_chr19.m	1780560	5467068	459	3762	10	0.40	0.13	0.45
query/MC1_chr20.m	1486652	4920008	460	3210	12	0.33	0.33	0.45
query/MC1_chr21.m	2290620	7790137	461	4065	14	0.32	0.30	0.40
query/MC1_chr22.m	1881124	6092892	471	3927	14	0.29	0.27	0.38
query/MC1_chr23.m	181370	231056	475	1461	15	0.26	0.26	0.35
query/MC1_chr24.m	1380254	4304292	465	2759	16	0.24	0.26	0.35
query/MC1_chr25.m	1024192	3380116	461	2224	18	0.22	0.24	0.31
query/MC1_chr26.m	2309200	7628112	462	3602	19	0.21	0.18	0.26
query/MC1_chr27.m	2863504	9433736	463	4386	20	0.20	0.13	0.21
query/MC1_chr28.m	530800	1706700	467	1388	22	0.19	0.22	0.30
query/MC1_chr29.m	3584718	11338112	448	7297	24	0.17	0.22	0.30
query/MC1_chr30.m	1207135	4384115	451	2870	25	0.16	0.22	0.30
query/MC1_chr31.m	730752	2437704	460	1603	50	0.09	0.09	0.15
query/MC1_chr32.m	1236060	4104448	452	2710	70	0.06	0.06	0.10
query/MC1_chr33.m	1293990	4248620	459	2719	100	0.05	0.05	0.08
query/MC1_chr34.m	53608	170800	537	180	150	0.04	0.04	0.06
query/MC1_chr35.m	48025000	15192070	463	8012	200	0.03	0.03	0.04



Assessment of the Model

- We tested the algorithm using different groups of sequence files each with
 - Varying number of sequences.
 - Different deadline parameters.
 - Different Data clustering models.
 - Different run profiles.
- A key objective is to check whether the EAS model will be able to make adjustments whenever necessary to meet deadline.
- The other key objective is whether the adjustment will meet the deadline with minimum resources.



Scheduling – Energy & Deadline aware

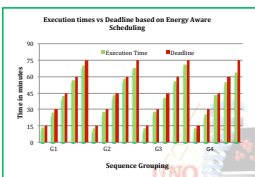
- Here we bring together our understanding of scheduling, High Performance Computing and our specific knowledge about BLAT in HPC.
- We develop a simple machine learning energy aware scheduling algorithm that takes into account
 - The run profile,
 - The number of sequences processed,
 - The number of nodes used for processing and
 - The time it took to execute.
- Now when new BLAT queries are submitted along with their desired deadline,
 - The algorithm uses information on the number of sequences to be processed,
 - To allocate the least number of nodes needed to meet that deadline,
 - Thus managing performance as well as energy to finish the tasks.



Scheduling – Energy & Deadline aware

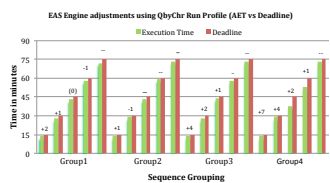
- We used 4 groups (G1, G2, G3, G4) of query files each group had (5, 10, 15, 20) files respectively with varying number of sequences.
- Each group of query sequence files was run against 5 different deadlines (15, 30, 45, 60, and 75 minutes).
- In each instance the actual execution time (AET) met the given deadline based on the minimum number of nodes assigned for each task group.
- Thus optimizing both performance and energy considerations.

Groups	Query Files	Total # of Sequences
G1	5	22566
G2	10	40530
G3	15	55946
G4	20	79222

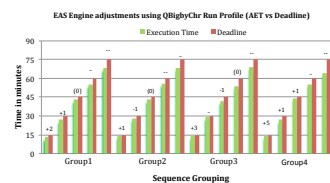


UNO BIOINFORMATICS

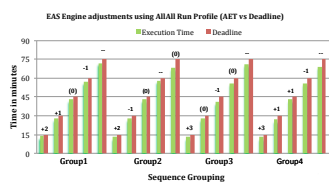
EAS Adjustments – QbyChr Profile



EAS Adjustments – QBigbyChr Profile




EAS Adjustments – AllAll Profile



HPC and Biological Networks

- Network creation: 2 weeks on PC
 - 10 hours in parallel, 50 nodes
 - 40,000 nodes = 800 million edges
- Network analysis: Best in parallel
 - Only 3% of entire genome forms complexes
- UNO Sapling Cluster
- Holland Computing Center: Firefly




Large-Scale Networks and Data Analysis

Many application domain rely on creating networks to model and analyze key relationships among data elements in the domain


Examples: Biological networks – social networks – inference networks - scheduling networks – Transportation networks

- ✓ **Modeling versus data mining**
 - Such networks are normally very large
 - They are susceptible to significant noise related problems
- ✓ **Sampling sub-networks (sub-graphs)**
 - Reduce network size
 - Reduce noise impact
 - Questions: does it preserve integrality of the original network

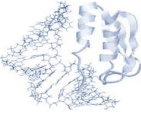


Motivation: Real World Networks

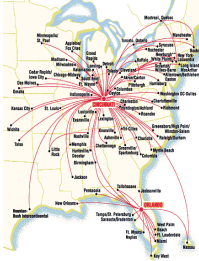
Social Network



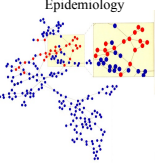
Bioinformatics




Air Transportation



Epidemiology







Motivation: Signal From Noise


facebook

Facebook helps you connect and share with the people in your life.




People You May Know





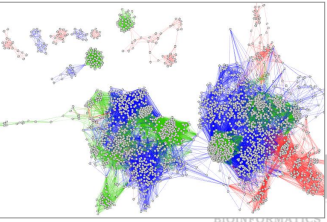
We need a technique that reveals true connections in the network by predicting missing and spurious interactions in a system.


We propose a method that separates wheat from chaff, the signal from the noise



Back to Correlation Networks


- Model for handling high-throughput biological data
- Network contains biologically relevant subgraphs:
 - Hubs
 - Clusters
 - Motifs
 - Bottlenecks





Size and Noise

- Network made from average gene expression experiment will have:
 - 40,000 nodes
 - 800 million edges
- Only 3% of genes in entire genome work together to form complexes
- Even with parallel computing resources, unfiltered networks are too noisy for biological discovery



Technical Solutions

- Better application specific parallelization – custom solutions lead to better performance
- Better scheduling model: multi-layer dynamic scheduling solution
- *Smarter input data: better user level utilization plus integrated domain knowledge with computational tools*



Network Filters

Design a network filter and obtain a sub-network of the original network such that:

- It maintains the important stuff – signal
- Remove unimportant stuff – noise
- Maintain network elements of biological relevance
- Uncover new ones



Network Filters

- Chordal graph sampling
 - Keep triangles in expression graphs
 - Remove large cycles, extra edges
 - Keep clusters, identify new clusters
- Spanning tree sampling
 - Keep high degree nodes (maybe?)
 - Remove up to 50% of edges
 - Enhance identification of lethal nodes
- Hybrid chordal-spanning tree method
 - Keep high degree nodes
 - Keep clusters
 - Remove 40-50% of edges
 - Proactively distort/enlarge network structures



HPC and Network Analysis

- Network sizes tend to be large
- Signal-to-noise ratio can be high
 - ID biologically relevant relationships?
 - Remove irrelevant nodes/edges?

Original network
ID noisy edges
Weight with literature
Enriched network

Young	Middle Aged	Aged
28,477 Nodes	29,371 Nodes	29,520 Nodes
77,807 Edges	382,628 Edges	126,354 Edges
1.9% Edge Density	8.8% Edge Density	2.9% Edge Density

• Original Nodes: 41,174
 • Original Edges: 847,628,551
 • Correlation: Only 1.00
 • P-value: $p < 0.005$

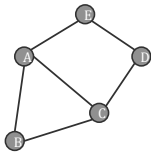
Chordal Graph Sampling

Goal: Develop a parallel network sampling technique that *filters noise*, while *preserving the important characteristics of the network*.

- ✓ **Maximal Chordal Subgraph**
 - Spanning subgraph of the network w
 - No cycles of length larger than three
- ✓ **Properties of Chordal Graph**
 - Preserves most cliques and highly connected regions of the network
 - Most NP hard problems can be solved in polynomial time
 - Complexity of finding maximal chordal subgraphs: $O(|E| \cdot \max_deg)$

Why chordal graphs?

- Chordal graphs are triangulated
 - We want to preserve K_3 subgraphs (triangle)
 - K_3 graphs/motifs are known to represent co-regulated genes
 - Use chordal graphs as a filter for finding co-regulated structures

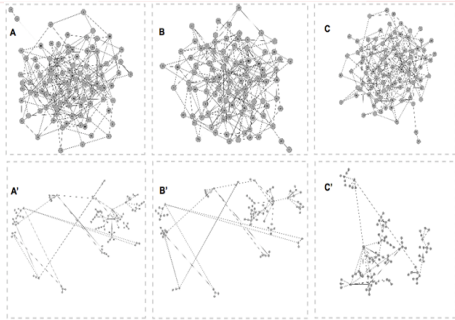


Subgraph formed by A,B,C is more likely to be biologically relevant.

If gene A and gene B are co-regulated, and if gene A and gene C are co-regulated, then genes B and C will be co-regulated.



Visual Representation



BIOINFORMATICS

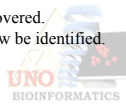
Proposed Approach

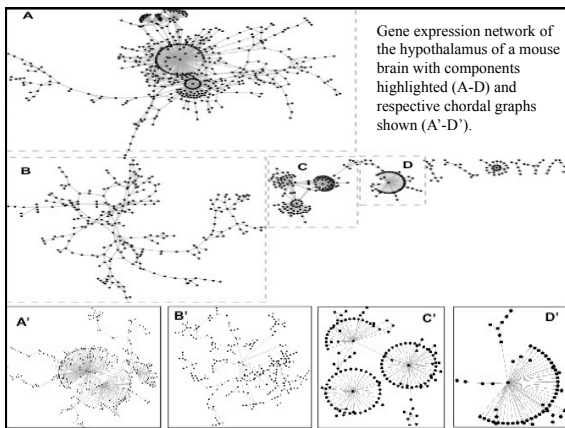
- ✓ Create networks from publicly available data
 - Aging mice gene expression data –
 - Young vs. middle-aged mice
- ✓ Test method on networks
- ✓ Assess results by examining biological relevance of network structures
 - Clusters enriched with function (Gene Ontology)
 - Do we maintain clusters and function in sampled graphs?
 - Do we find new functions in sampled graphs?

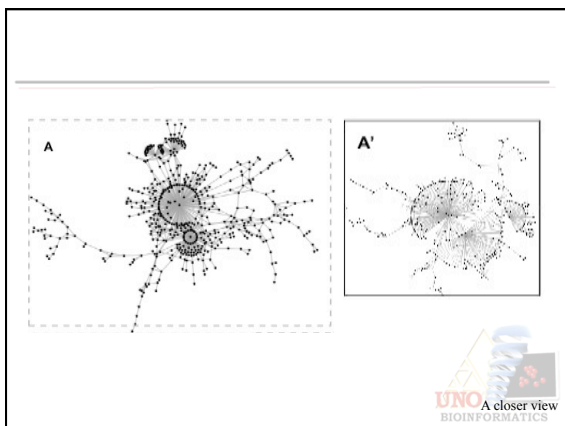


Hypothesis

- Hypothesis H_0 : Given a graph G representing a correlation network, maximal chordal subgraph G_1 will maintain most of the highly dense subgraphs of G while excluding edges representing noise-related relationships in the network.
 - H_{0a} - Key functional properties found in the clusters of unfiltered networks G are maintained in the sampled networks G_1
 - H_{0b} - New clusters with biological function are uncovered. Functional attributes previously lost in noise can now be identified.



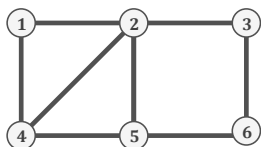




Algorithm

- Based on P. M. Dearing *et al.* "Maximal chordal subgraphs", Discrete Applied Mathematics 20(3), 1988.
- Method is based on growing the graph from a starting vertex and adding edges so long as they maintain the chordal characteristics.

1. Select V1
2. Select V2
3. Select V4
4. Select V5
5. Select V3
6. Select V6



Results: Combinatorial Properties

- Young mice, aged 2 months (GSE5078 via NCBI)
- 5,349 vertices and 7,277 edges

Combinatorial Properties	Original Network	Quasi Chordal Subgraph of young Mice(GSE5078) with				
		4,949	4,269	4,029	3,857	3,683
Number of edges	7277	4,949	4,269	4,029	3,857	3,683
Mean Clust. Coeff	.48	.39	.47	.40	.41	.41
High Degree vertices	146	73(50)	106(72)	96(65)	86(58)	95(65)
Core Numbers	46	26(56)	25(54)	39(84)	36(78)	26(56)



Functional Properties

Cluster Id	Original Network	QCS on 1 Partition	QCS on 2 Partitions	QCS on 4 Partitions	QCS on 8 Partitions	QCS on 16 Partitions
1	protein binding	protein binding	protein binding	protein binding	protein binding	protein binding
	catalytic activity	catalytic activity	catalytic activity	binding	enzyme regulator activity	binding
	binding	binding	binding	binding	enzyme regulator activity	binding
	binding	binding	binding	binding	enzyme regulator activity	binding
2	receptor binding	receptor binding	receptor binding	receptor binding	protein binding	protein binding
	protein binding	protein binding	protein binding	protein binding	protein binding	protein binding
	binding	binding	binding	binding	binding	binding
	binding	binding	binding	binding	transmembrane transport activity	transmembrane transport activity


- Most of the important functional units are preserved
- Sometimes reveals new clusters in the reduced networks
 - Ex: Enzyme regulator activity in QCS 4, 8
 - Ex: Transmembrane transport activity in QCS 8,16.

9/16/12



Maximal Chordal Subgraph


- Obtaining maximum chordal subgraph is NP-hard
- Algorithm for finding maximal chordal subgraph [Dearing et al. 1988]
 - Based on growing a subgraph with chordal edges
 - Depends on the order of the vertices
 - Unweighted Graph
- Complexity $O(E^2)$



Algorithm

Vertex 1 is selected
 Vertex 2 is selected
 Vertex 4 is selected
 Vertex 5 is selected
 Vertex 3 is selected
 Vertex 6 is selected

This is a traversal algorithm
Criteria: Seen neighbors of u is a subset of seen neighbors of v




Original Network

Chordal subgraph for each partition

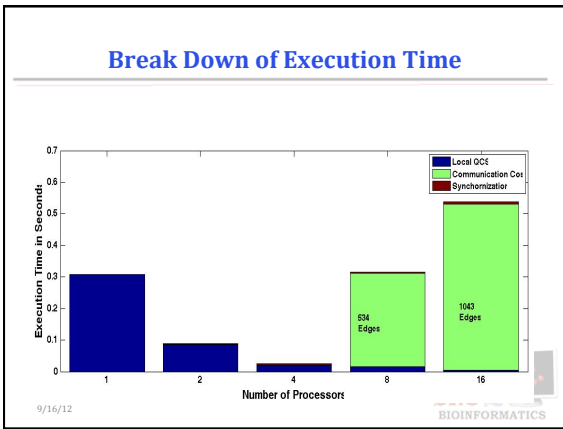
Adding the border edges

Nearly-chordal subgraph



Parallel Chordal Subgraphs (Coarse Grained)

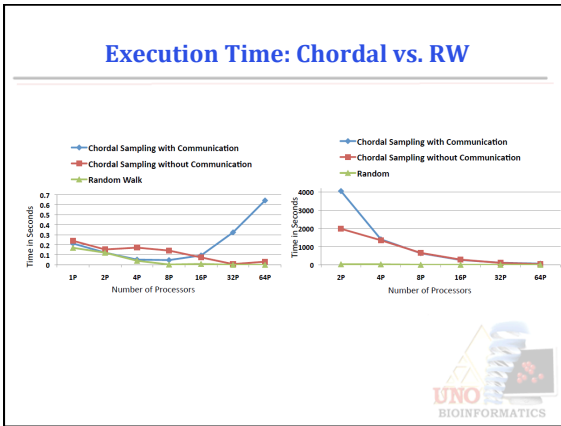
- Input: Original Graph, Output: Maximal Chordal Subgraph (MCS)
- Partition graph across processors
- In each processor P_i (In Parallel)
 - Find local MCS
- Handling Border edges (Communication)
 - Between 2 processors
 - If (na, nb) and (na, nc) are border edges across P_1 and P_2 , and (nb, nc) is a chordal edge in P_2
 - Between 3 processors
 - If (na, nb) is a border edge across P_1 and P_2
 - If (na, nc) are border edges across P_1 and P_3 , and
 - If (nb, nc) is a edge across P_2 and p_3
- Eliminating Cycles (In Parallel)
 - Run a cycle detection algorithm in each processor
 - Eliminate edge whose vertices have lowest degree

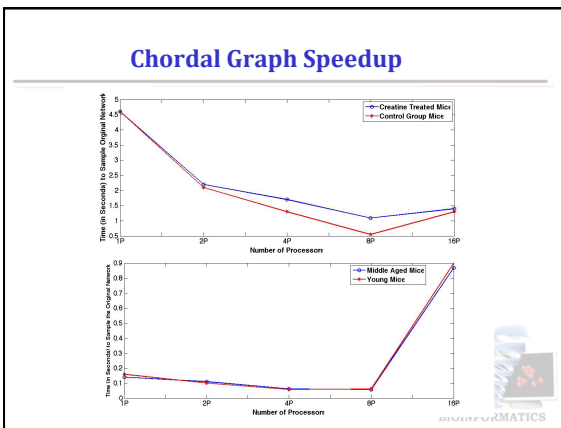


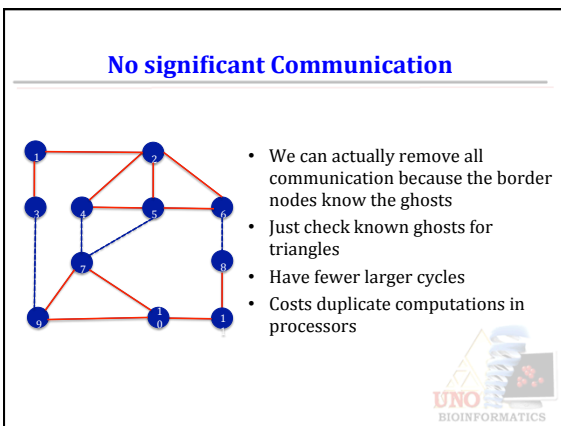
Execution Time

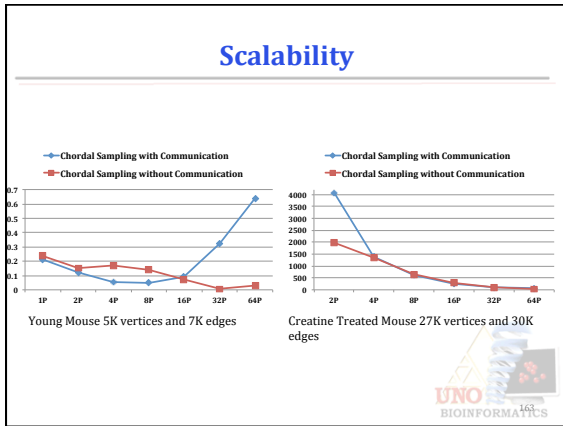
Number of Processors	Local MCS (s)	Communication Cost (s)	Synchronization (s)
1	0.30	0.00	0.00
2	0.08	0.02	0.00
4	0.02	0.02	0.00
8	0.01	0.29 (54 Edges)	0.00
16	0.01	0.52 (1043 Edges)	0.00

- Parallel algorithm loses efficiency as more processors are added
- Border edges increase with more processors
- We are currently investigating a master-worker approach that might reduce some of the communication costs







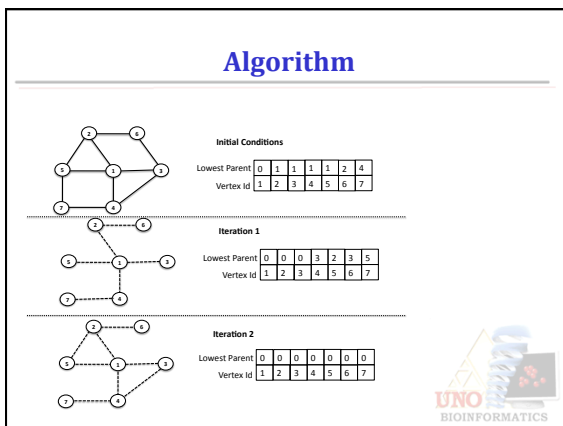


Parallel Chordal Subgraphs (Fine Grained)

- Every vertex associated with an id
- Vertex identifies **lowest parent**—smallest id number lower than itself
- Find chordal neighbors
- At each iteration check if own chordal neighbors are a subset of the chordal neighbors of lowest parent
 - If yes include LP as chordal neighbor
- Update LP to next highest vertex
- Terminate when no LPs remain

UNO BIOINFORMATICS

Influenced by Multithreaded Algorithms for Graph Coloring –Catalyurek, Feo, Gebremedhin, Pothan and Halappanavar (preprint)



Architecture

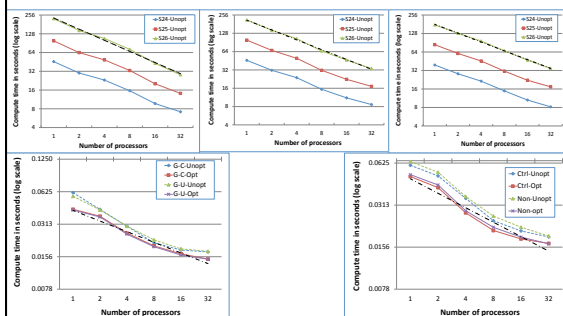
- Shared Memory (AMD Opteron)
 - 48 cores, (4 sockets with 12 processors)
 - 256 GB of globally addressable memory
 - 3 level Cache
- Multithreaded (Cray XMT)
 - 128 processors with 128 hardware threads
 - 8 GB per processor
 - No cache—uses hardware hashing for uniform memory access



Test Sets


Group	Vertices	Edges	Max Degree	Avg Degree
RMAT-ER (24)	16,777,216	134,217,654	42	16
RMAT-ER (25)	33,554,432	268,435,385	41	16
RMAT-ER(26)	67,108,864	536,870,837	48	16
RMAT-G(24)	16,777,216	134,181,095	1,278	416
RMAT-G (25)	33,554,432	268,385,483	1,489	442
RMAT-G(26)	67,108,864	536,803,101	1,800	469
RMAT-B(24)	16,777,216	133,658,229	38,143	8086
RMAT-B(25)	33,554,432	267,592,474	54,974	9,539
RMAT-B (26)	67,108,864	535,599,280	77,844	11,214
GSE5140 (CRT)	45,023	714,628	690	58
GSE5140(UNT)	45,020	1,311,396	657	32
GSE17072(CTL)	48,803	949,094	365	39
GSE17072(NON)	48,803	1,109,553	463	45

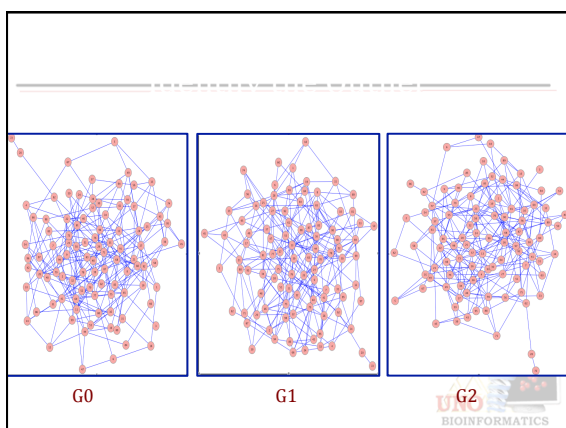
Results on AMD

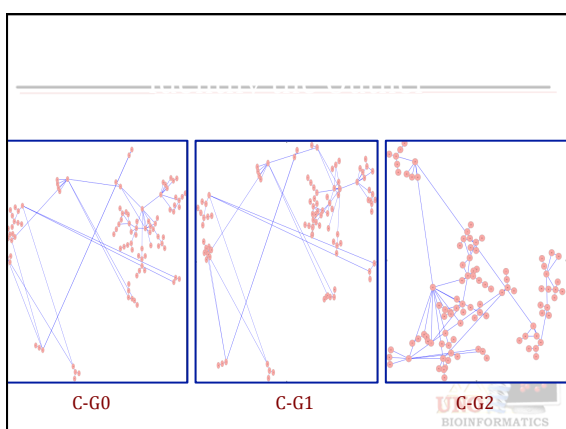


Speedup

Group	XMT (UnOpt)	XMT (Opt)	AMD (UnOpt)
RMAT-ER (24)	32.34	31.70	6.38
RMAT-ER (25)	40.29	29.57	6.96
RMAT-ER(26)	32.25	28.0	7.9
RMAT-G (24)	41.08	41.29	5.33
RMAT-G (25)	44.17	45.97	5.75
RMAT-G(26)	47.35	47.97	6.24
RMAT-B (24)	33.61	35.08	4.80
RMAT-B (25)	21.37	36.16	4.86
RMAT-B(26)	16.70	34.09	5.13
GSE5140 (CRT)	1.40	1.22	3.54
GSE5140(UNT)	1.43	1.14	2.85
GSE17072(CIL)	1.52	1.25	3.41
GSE17072(NON)	2.05	1.65	3.12

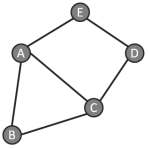







Chordal and other Filters

Chordal-based filters			Tree-based filters		
Filter	Name	Description	Filter	Name	Description
HD	High Degree	Traversal based on ascending order of vertices	ST	Spanning Tree	Tree determined by Prim's Algorithm
LD	Low Degree	Traversal based on descending order of vertices			



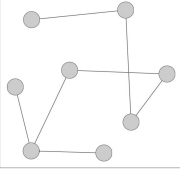
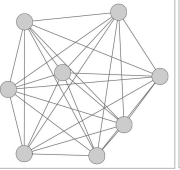
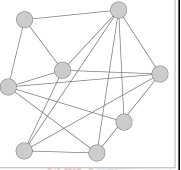
Chordal graphs

- Preserve K_3 subgraphs (triangle)
- Known to represent co-regulated genes
- Chordal graphs as a filter for finding co-regulated structures
- Maintain local dense connections (highly dense subgraphs)




Other Filters

Chordal-based filters			Tree-based filters		
Filter	Name	Description	Filter	Name	Description
HD	High Degree	Traversal based on ascending order of vertices	ST	Spanning Tree	Tree determined by Prim's Algorithm
LD	Low Degree	Traversal based on descending order of vertices			






Spanning Tree
Original
Chordal



Hypothesis

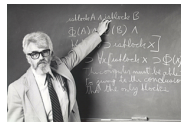
- H_0 : A tree-based network filter will identify essential high-degree nodes in some network G while reducing the number of edges comparably to a chordal-based filtered network
 - Propose a spanning tree network filter
 - Reduce the number of edges in the network
 - Nodes with high degree will maintain their high degree in the spanning tree
 - Unweighted edges
 - Edges chosen for the filter will connect to nodes are more likely to connect to central in the original network due to the "friendship paradox" described by [16] for modular networks
 - We use two controls, the original unfiltered network, and multiple chordal graph filters, to verify our findings.



Security/Privacy Issues - Working in the Cloud



Working in the Cloud



Computation may someday be organized as a public utility.

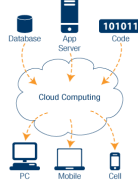
- John McCarthy, 1960



Working in the Cloud

• Cloud computing is Web-based processing and storage. Software and equipment are offered as a service over the Web.

- Data and applications can be accessed from any location
- Data and applications can easily be shared through a common platform
- Clouds need not be public; companies can introduce private cloud computing solutions



Cost Reduction & Convenience

- Flexible availability of resources
- Opportunity for developers to easily push their applications
- Targeted advertising
- Easy Software Upgrades for customers
 - Example: Webmail

UNO BIOINFORMATICS

The Future of Thin Clients

Browser (Thin Client) - Internet Connection - Cloud (Rackspace, Amazon, Salesforce)

UNO BIOINFORMATICS

Working in the Cloud

- Cloud Computing for Digital Nomads
 - Cloud computing allows road warriors, telecommuters, and freelancers to connect to their employers using affordable computing tools.
 - Digital Nomads enjoy location independency and the ability to set their own hours, but do not always have job security or benefits.

UNO BIOINFORMATICS



Cloud Security Incidents

- Cloud-based providers are an attractive target for hackers:
 - A hacker broke in to the personal Web services accounts of Twitter co-founder Evan Williams, his wife, and another Twitter employee, and used that access to steal a number of confidential company documents.
- Multiple levels of Cloud security are needed:
 - Encryption of all communications
 - Server-side security, including regular third-party audits
 - Client-side security, such as firewalls and anti-virus software
 - Client-side password security


TechCrunch Mar 7, 2009

In a privacy error that underscores some of the biggest problems surrounding cloud-based services, Google has sent a notice to a number of users of its Document and Spreadsheets products stating that it may have inadvertently shared some of their documents with contacts who were never granted access to them.



Cloud Security Incidents

- Dropbox's password nightmare highlights cloud risks
 - Cloud storage site - 25 million users
 - Glitch let visitors use any password to log in customer's accounts
- Blackberry outages
- Netflix accidentally revealed rental histories
- Insurer and Interviewers seeks patient's Facebook posts
- Sony PlayStation Network hacking



Cloud Security Challenges - Data


- Cloud-based providers need a data storage policy:
 - How is the data supplied to the service housed, protected, shared, manipulated, and disposed of?
 - Who owns the data? Who owns the meta-data?
 - Who has access to the provider's systems?
 - How can the data be used by parties other than the owner?
 - What is the provider's data retention policy?

Check out the Google Docs terms of service:

11. Content license from you

11.1 You retain copyright and any other rights you already hold in Content which you submit, post or display on or through, the Services. By submitting, posting or displaying the content you give Google a perpetual, irrevocable, worldwide, royalty-free, and non-exclusive license to reproduce, adapt, modify, translate, publish, publicly perform, publicly display and distribute any Content which you submit, post or display on or through, the Services. This license is for the sole purpose of enabling Google to display, distribute and promote the Services and may be needed for certain Services as defined in the Additional Terms of Use Services.


11.2 You agree that this license includes a right for Google to make such Content available to other companies, organizations or individuals with whom Google has relationships for the provision of syndicated services, and to use such Content in connection with the provision of those services.



Cloud Security Challenges - Availability

- Cloud-based computing requires an availability policy:
 - What is the provider doing to ensure that client data remains available, even in the event of a natural or human-induced disaster?
 - Provider back-ups should be frequent
 - Geographic redundancy reduces single-location disasters
 - Are client-side backups possible?
 - Use of open data standards allows data to be extractable or transferable to other storage locations

Flickr Accidentally Deletes a User's 4,000 Photos and Can't Get Them Back
THE NEW YORK OBSERVER February 1, 2011



Cloud Security Challenges - Availability



The popularity of Amazon's cheap, easily scalable hosting is showing its downside right now, with a number of popular websites and services throwing up errors or being down completely.

Foursquare, Quora, Reddit, Hooty and Hootsuite are among those affected by technical troubles on Amazon's servers. The company's status dashboard currently shows problems with the company's Elastic Compute Cloud and Relational Database Service operations, based in North Virginia, with connectivity issues confirmed.

We can confirm connectivity errors impacting EC2 instances and increased latencies impacting EB3 volumes in multiple availability zones in the US-EAST-1 region. Increased error rates are affecting EB3 CrossVolume API calls. We continue to work towards resolution.

FOROY TOOLBOX

f SHARE 345

Tweet 1,244

Like 60,000

Google+ 1,244

LinkedIn 1,244

StumbleUpon 1,244

Delicious 1,244

Diigo 1,244

Reddit 1,244

MyScoop 1,244

MyScoop 1,244



Cloud Security Challenges - Cancellation

- What if your service provider's business folds or is acquired?
 - Can data and business information be retrieved simply and easily, and in a format useful to the data owner?
 - Will cloud-based data be securely destroyed so others will not be able to access it later?



Cloud Security Challenges - Location

- Location of providers and data centers matter:
 - Local privacy requirements, libel laws, and obscenity regulations may be applied to cloud-based data.
 - If data is stored overseas, does that affect the client's ability to retrieve and produce information in litigation?
 - If data is stored overseas, is it now more accessible to domestic or international government investigators?
 - Does the service agreement limit the jurisdictions in which a client's data may reside?



Cloud Security Challenges – contd.

- Storage services provided by one vendor may be incompatible with another vendor's services should you decide to switch vendors
 - Example: Amazon's Simple Storage Service (S3) is incompatible with IBM's Blue Cloud, or Google, or Dell



Cloud Security Challenges – contd.

- Who controls the encryption/decryption keys?
 - Customer / cloud vendor

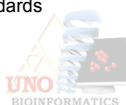
- Customers want:
 - SSL both ways across the Internet
 - Data encryption when data is at rest
 - Ideally customer must control the encryption/decryption keys



Cloud Security Challenges – contd.

- Data Integrity – assurance that data is identically maintained during any operation (such as transfer, storage, or retrieval)
 - Consistency and correctness

- Data must change only in response to authorized transactions
 - Unfortunately, there are no common standards



Cloud Security Challenges – contd.

- Proper fail-over technology
- Security at data-level so that data is secure wherever it goes
- Compliance standards do not envision compliance in a world of cloud computing
 - Issues of data privacy, segregation and security



Holy Grail for Cloud Security

- Homomorphic Encryption
 - An encryption system that allows computations to be performed on encrypted data
 - By enabling both multiplication and addition of encrypted data
 - Subsequent decryption yields the expected result in plaintext



Case Study – Porticor Cloud Security

- About: Porticor infuses trust into the cloud with secure, easy to use, and scalable solutions for data encryption and key management.
- Porticor's Virtual Private Data system combines state of the art encryption with patented key management to protect critical data in public, private and hybrid cloud environments. Within minutes, customers can encrypt their entire data layer and safely store the encryption keys. With breakthrough homomorphic split-key encryption technology, Porticor Virtual Private Data is the only system available that offers the convenience of cloud-based key management without sacrificing trust.



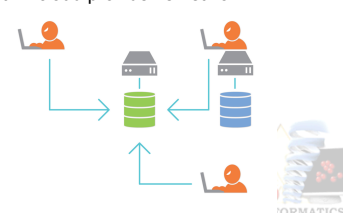
Porticor: Your Key to Cloud Data Security

- With groundbreaking technology and unparalleled manageability, Porticor is enabling businesses to rapidly protect their data in the cloud. Only Porticor offers homomorphic split-key encryption, to protect keys even when they are in use. At the same time, Porticor is breaking down barriers to entry, with installation and set-up in minutes, and pricing plans for organizations of all sizes.



3 Risks


1. Attacker breaches access control
2. Attacker Coming from another server in the same customer account
3. Attacker from within cloud provider's network



UNO BIOINFORMATICS


Risks, cont.

- Utilizes full disk encryption on virtual disk volumes and distributed storage (Amazon S3) in the cloud
- Each disk block written by the application (Database server) goes through a Porticor virtual appliance, where it is encrypted and sent to the disk volume. At no point is the plaintext data written to persistent storage
- All requests to read data from the disk similarly go through the Porticor appliance, which reads the encrypted data blocks, decrypts them and send the plaintext data back to the requesting application



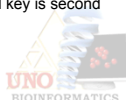
What about the keys?

- Ideally you don't want to store the keys in the same cloud as the encrypted data but you do need the keys to be in the cloud when needed
- Porticor Virtual Key Management (PVKM)
 - The keys are never exposed to anyone else. Not even to Porticor itself
 - For the first time, you can use encryption keys in the Cloud while maintaining strong security
 - You enjoy a simple automated system; while the protected disks are protected using strongest security: a key for each disk



The Banker Approach

- Porticor approach
 - Create a new Porticor protected project – generate a Master Key
 - Customer backs up the Master key as Porticor “never” sees it
 - Master key is kept in Porticor virtual appliance that lives within your own cloud account and never transferred to Porticor’s central database
- To Encrypt: encryption key is combination of Master key and a unique key created by PVKM
 - Thus Master key is one-half and PVKM generated key is second half



The Banker Approach

- Porticor virtual appliance generates the “second half” of the key for each new encryption
 - Master key, “first half”, remains common
- The second half is stored and protected in the PVKM system allowing access only to approved end user
 - By encrypting the second half by Master key and storing in the central PVKM



Does it Work?

- Did Porticor solve the problem?
 - Security dependent on Master key
 - Where is Master key stored?
 - In Porticor’s virtual appliance which lives within you own cloud account
 - Did we not just distrust the cloud interface against insider attacks and possible code vulnerabilities?
 - Where is the “second half” stored?
 - Central PVKM database outside your account
 - What is the security of that against
 - » Breaches – data deletion, modification, insider attacks
 - What is “second half” is lost?
 - » Master key (first half) becomes useless!



HPC – Best Practices

- Separation of special project central storage shares
- Use of front-end web-based servers to enable workflow execution within cluster. Having these servers hosted from a Virtual Private Cloud is the way to go!
- Identify opportunities for designing embarrassingly parallel applications – simpler to code.
- Combine MPI and OpenMP to make more efficient parallel programs.



Virtual Cloud Security

- Good old-fashioned basics:
 - Access Control (OS-level/Application-level)
 - Server/Machine firewall (OS-level)
 - Antivirus (OS-level)
 - Strong passwords (OS-level/Application-level)
 - VPN/IDS (Network-level)
- Private on-premise cloud
- Security Hardened virtual OS template
- Hypervisor security (cloud engine security)
- Encryption
 - Data In Transit (SSL/TLS, PKE)
 - Data At Rest (File system Encryption, Data Partitioning)
- Regulations:
 - HIPAA
 - HITTECH
 - FIPS 140-2
- Benefits: Scalability, quicker resource allocation, shared usage



Virtual Cloud Security Basics

- Strong passwords (OS-level/Application-level):
Best practices in creating, storing and periodically updating passwords.
- Access Control (OS-level/Application-level):
Control access to system or application using central Directory services like AD.
- Server/Machine firewall (OS-level):
Only services that need to be visible from outside should be opened up and by default, everything else should be inaccessible.
- Antivirus (OS-level):
An integral security part of any system.
- VPN/IDS (Network-level):
Network border-level access control using Virtual Private Network or Intrusion Detection System.



Private Cloud

- Core facilities need to acquire private infrastructure-level virtual cloud technology. Best vendor for such technology is VMware. The Bioinformatics Core facility at UNO uses *VMware vSphere Enterprise*.
- Public Clouds like Amazon EC2, RackSpace cannot be used in all cases due to various restrictions put forth by regulations (e.g. HIPAA data locality requirement). Such public clouds could only be used as a scalable platform for already anonymized data.
- Private Virtual Cloud is on-premise solution allowing all the benefits of virtualization technology both from an administrative and end-user perspective.

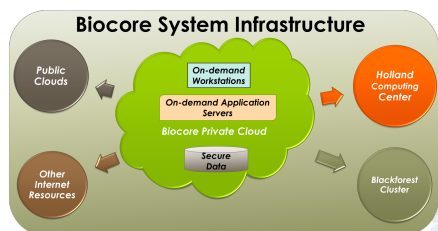


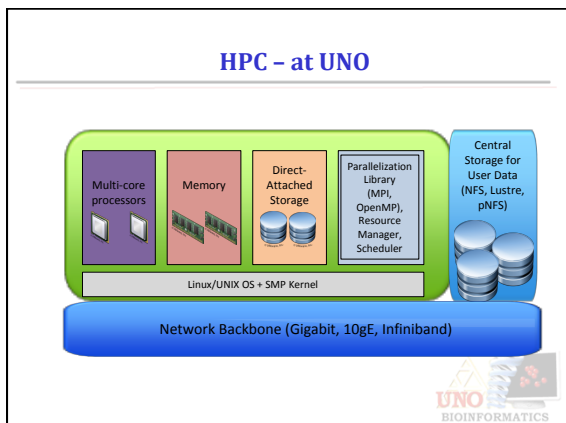
SecuRITY hardened OS template

- Virtual layer allows launching servers/machines from a template tested and hardened by the designated Security Office of the organization.
- Easier to isolate security issues on a virtual machine, resolve the issue on a clone of the virtual machine and deploy clone back to production environment.



Proposed Model





HPC - at UNO

Blackforest:

Consists of three hardware groups:

– Sapling:

- 17 x IBM xSeries 306m nodes with:
 - Intel Pentium D 3.00 GHz processor
 - 2 x 146GB SATA hard drives
 - and 4 GB of memory
- 3 TB of central NFS storage for user data
- Cisco Gigabit switch

HPC - at UNO

Blackforest:


– Morph:

- 16 x Dell PowerEdge 2970 nodes with
 - 2 x Quad-Core Opteron Processor 2378
 - 3 x 1TB SATA hard drives (2 TB raw)
 - 32 GB of memory
- 7 TB of central NFS storage
- Cisco Gigabit switch

HPC – at UNO

Blackforest:


- **Rapids:**
 - 8 x Dell PowerEdge R610 with
 - 2 x 6-core Intel Xeon X5660 2.8GHz Processor
 - 3 x 1TB SATA hard drives (2 TB raw)
 - and 128 GB of memory
 - 27 TB of central NFS storage
 - Cisco Gigabit switch



HPC – On-premise

Holland Computing Center:
Consists of two clusters:

- **Firedly:**
 - Was ranked 43rd (with 70 TeraFLOPS) in the world at 2010.
 - 871 x Dell SC1435 nodes with:
 - 2 x 2.8 GHz Dual Core AMD Opteron processors
 - 8 GB memory
 - 73 GB of direct storage
 - 280 x Dell SC1435 nodes with:
 - 2 x 2.2 GHz Quad Core AMD Opteron processors
 - 8 GB memory
 - 73 GB of direct storage
 - 150 TB of Central NFS storage
 - Infiniband connectivity




HPC – On-premise

Holland Computing Center:

- **Tusker:**
 - 106 Dell R815 Nodes with:
 - 4 x 16-core Opteron 6272 2.1GHz processors
 - 256/512GB of memory
 - 500 GB of direct storage
 - 350 TB of Terascale Lustre parallel storage
 - 40 Gb/s Infiniband connectivity

More information about HCC can be found at: <http://hcc.unl.edu>



Where do we go from here?

- Data Generation and Collection
- Data Access, Storage and Retrieval
- Data Integration
- Data Visualization
- Analysis and Data Mining
- Decision Support
- Validation and Discovery

A focus on Innovative Data Analysis

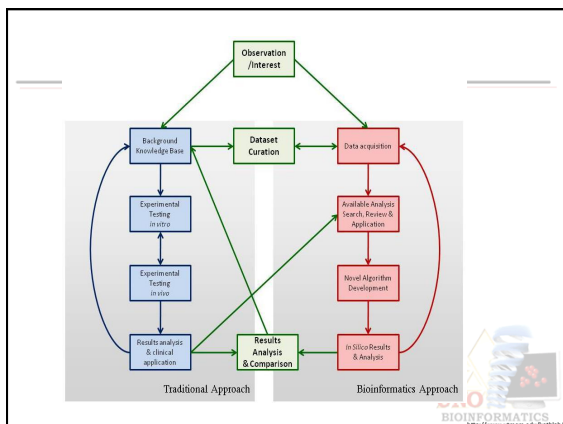


Data-Driven Decisions

- With high throughput data collection, Biology needs ways not only to store data but also to store knowledge (Smart data)
- Data: Things that are measured
- Information: Processed data
- Knowledge: Processed data plus meaningful relationships between measured entities
- Decision Support


Next Generation Tools: ICD Tools






Conclusions

- Next Generation Bioinformatics Tools need to be Intelligent, Collaborative, and Dynamic
- Biomedical scientists, Bioinformatics researchers and computer scientists need to work together to best utilize the combination of tools development and domain expertise
- HPC is critical to the success of the next phase of Biomedical research but again the integration needs to happen at a deeper level
- The outcome of collaboration has the potential of achieving explosive results with significant impact on human health and overall understanding of biological mysteries



Conclusions

- Many Scientific disciplines are now at crossroads
- The proper penetration of IT represent tremendous challenges and great opportunities
- Discovery are likely to take place at many places
- The importance of interdisciplinary approach to problem solving
- This may lead to scientific revolution





Acknowledgments



- UNO Bioinformatics Research Group
 - Kiran Bastola
 - Sanjukta Bhoomwick
 - Kate Dempsey
 - Jasjit Kaur
 - Ramez Mena
 - Sachin Pawaskar
 - Oliver Bonham-Carter
 - Ishwor Thapa
 - Dhawal Verma
 - Julia Warnke
- Former Members of the Group
 - Alexander Churbanov
 - Xutao Deng
 - Huiming Geng
 - Xiaolu Huang
 - Daniel Quest
- Biomedical Researchers
 - Steve Bonasera
 - Richard Hallworth
 - Steve Hinrichs
 - Howard Fox
 - Howard Gendelman
- Funding Sources
 - NIH INBRE
 - NIH NIA
 - NSF EPSCoR
 - NSF STEP
 - Nebraska Research Initiative