

Interactively exploring distributed computational models of biology

James Watson & Janet Wiles

ARC Centre for Complex Systems

ARC Centre of Excellence in Bioinformatics



Copyright

Permission is granted for this material, presented at the 8th Asia-Pacific Complex Systems Conference (Complex'07), 2-5 July 2007, Surfers Paradise Marriott Resort, Queensland, to be available on the Complex'07 website to be shared for non-commercial, educational purposes, provided that this copyright statement appears on the reproduced material, and notice is given that the copying is by permission of the author(s). To disseminate otherwise or to republish requires written permission from the author(s).

ARC Centre for Complex Systems

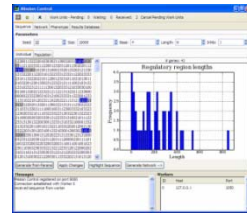
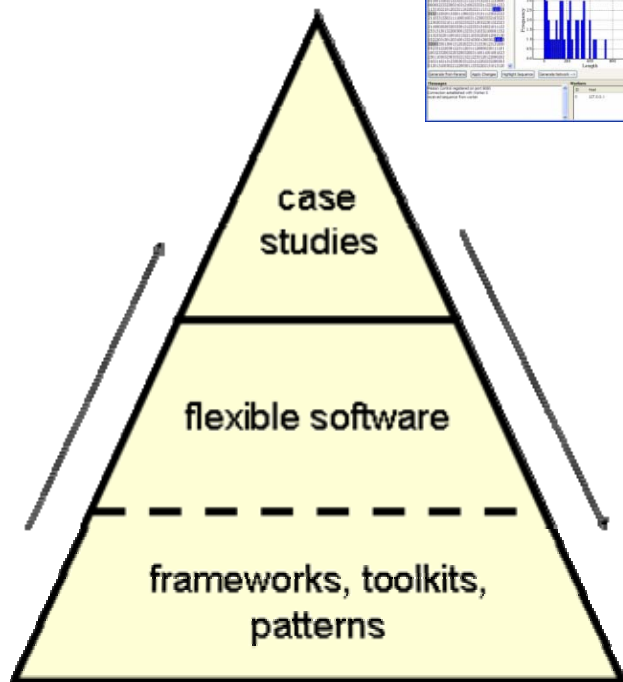
School of ITEE | The University of Queensland | ST LUCIA QLD 4069 | AUSTRALIA

T: +61 7 3365 1003 | F: +61 7 3365 1533 | E: outreach@accs.edu.au

www.complex07.org

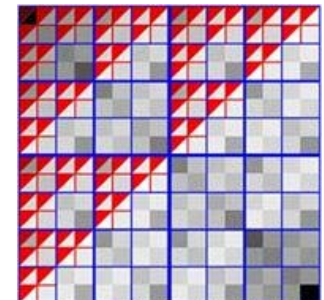
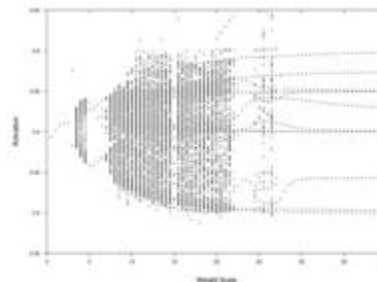
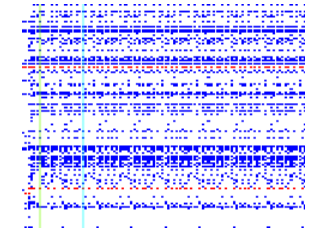
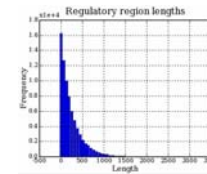
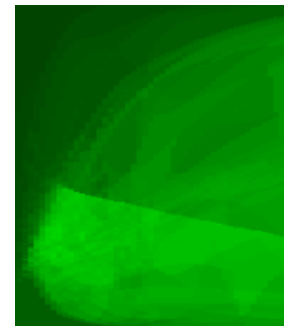
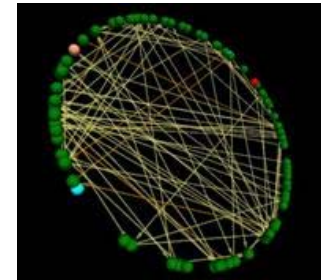
Biological modeling techniques

Small team sizes
 Transient models
 Reused components
 Reinvention
 Team members move on



```

3312200333002132110100023013032223201213
101203301331032332332000111322231010113
3123112203332123033300122322212113201113
3320031100203232033220023033103202001122
3032322101031013200120001200230121322121
1323323001010132001100211120032223231012
1231222010110031011220020131200332031322
2231132131101311210031210111033331321031
3211222002032121313010132010320113103120
0101300023231332012122012312012222010202
1012303203213112023121002003231110020111
    
```



Interactive modeling

- The ACCS GRN group has found dynamic visualization and interaction particularly useful in biological modeling preferable to batched simulations
 - Cross-disciplinary makes model shared language
 - Projects (e.g., Neurosphere Lab) have been successful when collaborator interacts / critiques through visualizations / GUI's
 - User exploration facilitates insight into dynamics
 - Quickly hone in on points of interest by dynamically prioritizing (steering) computations to those regions

Simulations take time

Speed is necessary for interactive modeling

Options

- Improve software
- Faster hardware
 - becoming multi-core
- More hardware

Why distributed?

- Supercomputing is not ideal for interactivity
- The Grid
 - Many machines currently available
- Idle CPU's in many research institutes
- This project: explore feasibility of interactive simulations using distributed computing

Available hardware

- ~400 Desktop machines in student labs
- Fast local network

GX150SD - P3 1GHz, 256MB RAM
GX270SD - P4 2.8GHz, 512MB RAM
GX280SD - P4 3GHZ HT, 512MB RAM
GX280DT - P4 3GHZ HT, 512MB RAM
GX520DT - P4 3GHZ HT, 1024MB RAM
OP745DT - Core2 Duo 2.13GHz, 2048MB RAM

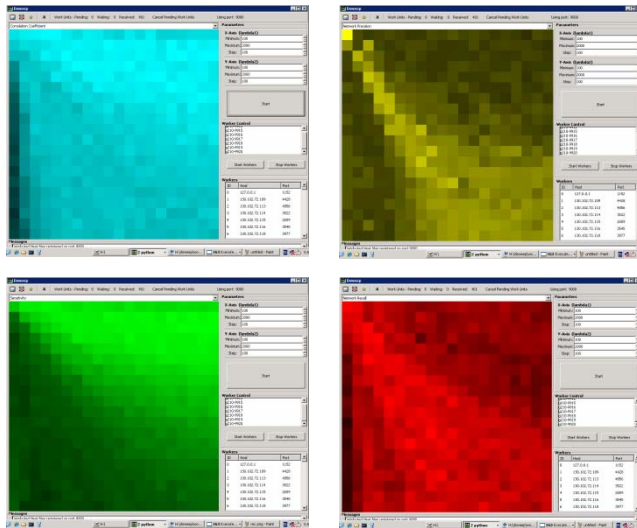
"All hardware types have either 100mbit or 1000mbit network cards, all of them are connected at 100mbit speed with gigabit backbones."

Lab Characteristics

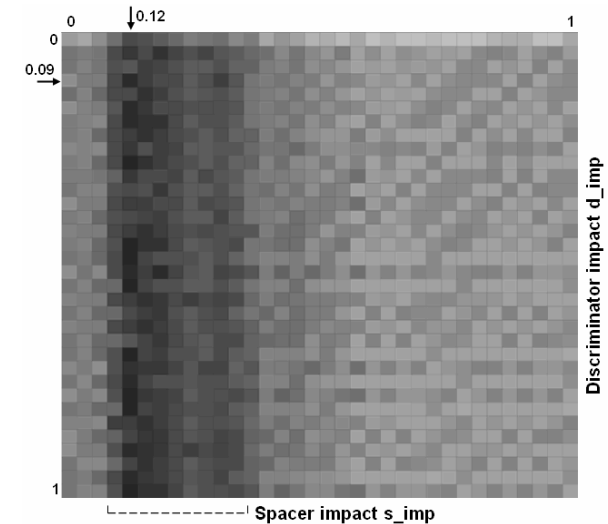
- Constantly changing usage
- Re-imaged on boot
- Have a 'failure' rate of ~5 machines an hour

Distributed speed-up

- Search for regularisation parameter values in multi-purpose machine learning
 - Up to 20 minutes per run



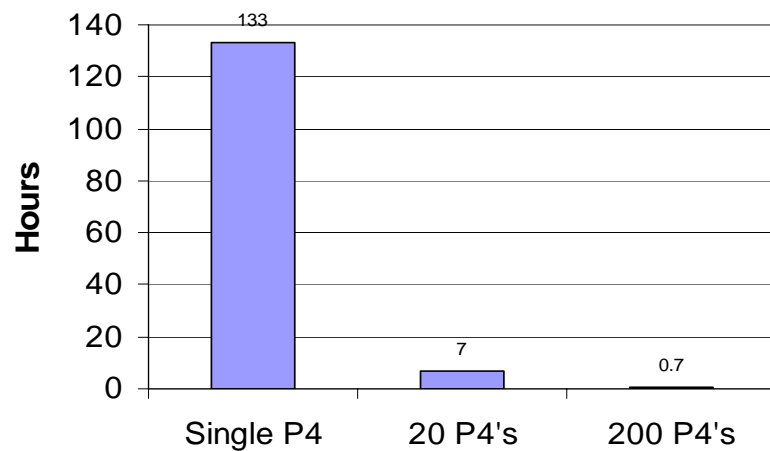
- Promoter model of *Bacillus subtilis* for transcription start site prediction
 - ~3 minutes per run
 - We performed 1000's of runs for a recent submission



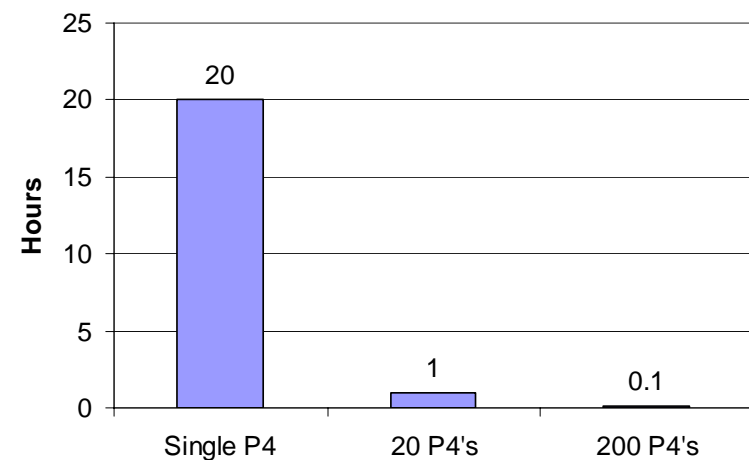
Distributed speed-up

- Search for regularisation parameter values in multi-purpose machine learning
 - Up to 20 minutes per run
- Promoter model of *Bacillus subtilis* for transcription start site prediction
 - ~3 minutes per run
 - We performed 1000's of runs for a recent submission

Approximate running time for
20x20 20 minute simulation



Approximate running time for
20x20 3 minute simulation



Dsweep [Min] [Max] [Close]

Work Units - Pending: 0 Waiting: 0 Received: 1 Cancel Pending Work Units Using port: 9000

error value

Parameters

X-Axis (x1)
Minimum: 12
Maximum: 25
Step: 1

Y-Axis (x2)
Minimum: 2
Maximum: 15
Step: 1

Start

Worker Control

Start Workers Stop Workers

Workers

ID	Host	Port
0	127.0.0.1	1804

Messages

Distributed Heat Map registered on port 9000
Connection established with Worker 0 on 127.0.0.1:1804

Beyond parameter sweeps

- Many biological models are “embarrassingly parallel”
 - Individuals / populations / etc. can be easily divided into distinct units of work
 - Generate regulatory network from genetic sequence
 - Work units often small, but many

Evolving the Artificial Genome 2.0

Model Help

Explore Mapping

Evolution

Genome

Seed: 10 Size: 10000

Base: 4 Length: 6

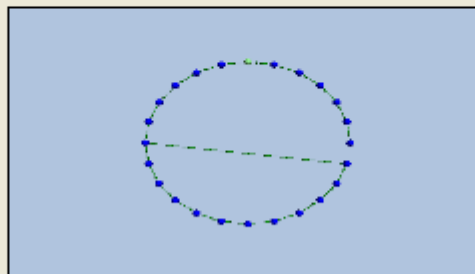
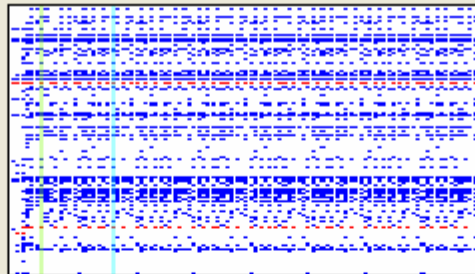
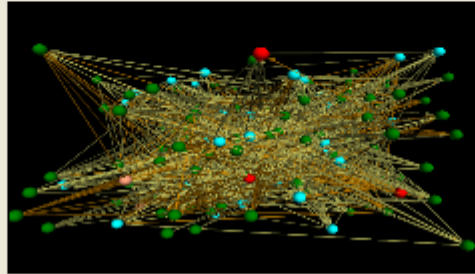
Inhibition: 1 Time Steps: 100

Generate

```
0122011323332033030313332330203120001130
3021010003221133120031110120321000300031
3000211333332200003222113320120133030322
1330320220031021331203133113013110132211
3312200333002132110100023013032223201213
1012033013310323323320000111322231010113
3123112203332123033300122322212113201113
3320031100203232033220023033103202001122
3032322101031013200120001200230121322121
1323323001010132001100211120032223231012
1231222010110031011220020131200332031322
2231132131101311210031210111033331321031
3211222002032121313010132010320113103120
0101300023231332012122012312012222010202
1012303203213112023121002003231110020111
3110102113000003202310301203130112122031
1303301000132212123303012202010011102311
1031203303100311300333323121003102200311
3112003102022001001022201332112112120013
2002311030331030111000320301221012121300
2311230300103123022011120132033110112230
3002020233220321130331221032312111100021
2120211120332101222031321102112203222103
2202103001233120133122300123022130111112
2011222000120020011000000000000000000000
```

Network

nodes: 200 # links: 2336 K: 11.680



Phenotype

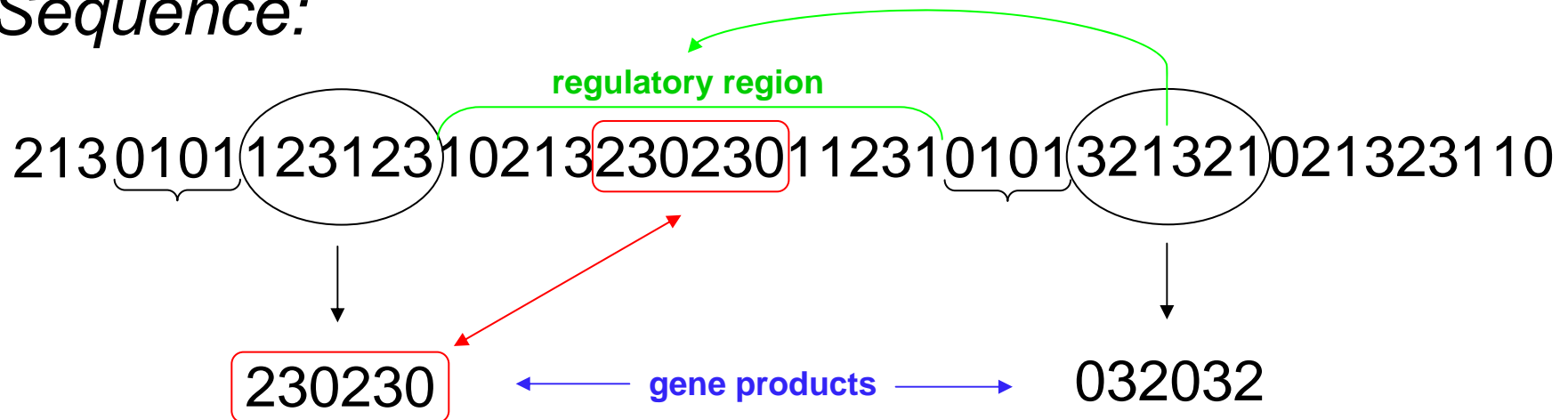
- Light Interception: 71.57 %
- Reproductive Success: 76.93 %
- Surface Area: 100.00 %

Total Relative Fitness: 3.80



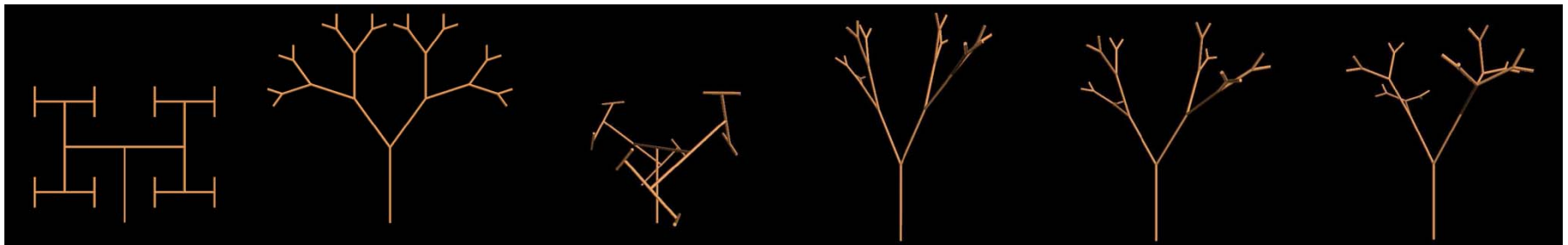
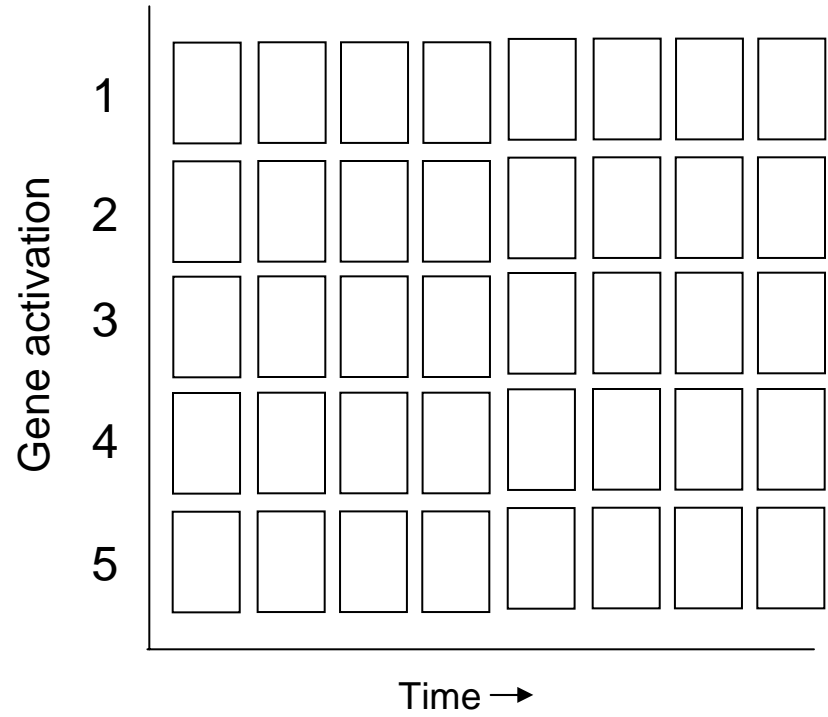
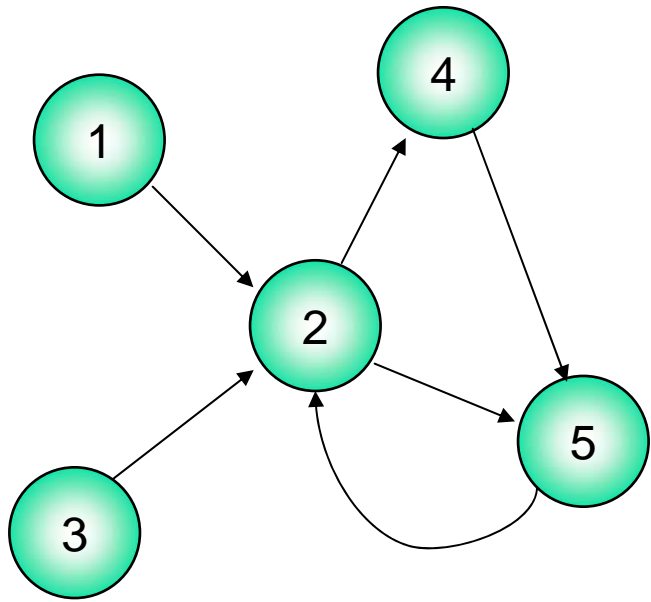
The Artificial Genome

Sequence:



Defines network:





Evolving the Artificial Genome 2.0

Model Help

Explore Mapping

Evolution

Genome

Seed: Size: Base:
 Length: Inhib.: Steps:

Evolve

Population Size: # Generations:

Selection Algorithm:

- Roulette Wheel
- Ranked Roulette Wheel
- Tournament K:
- None

Mutation:

- Single Point
- Duplication
- Transposition
- Deletion
- Inversion
- Dupl. + Deletion

Times per event:

- Use Crossover
- Use Elitism

Consider the performance of: L R S

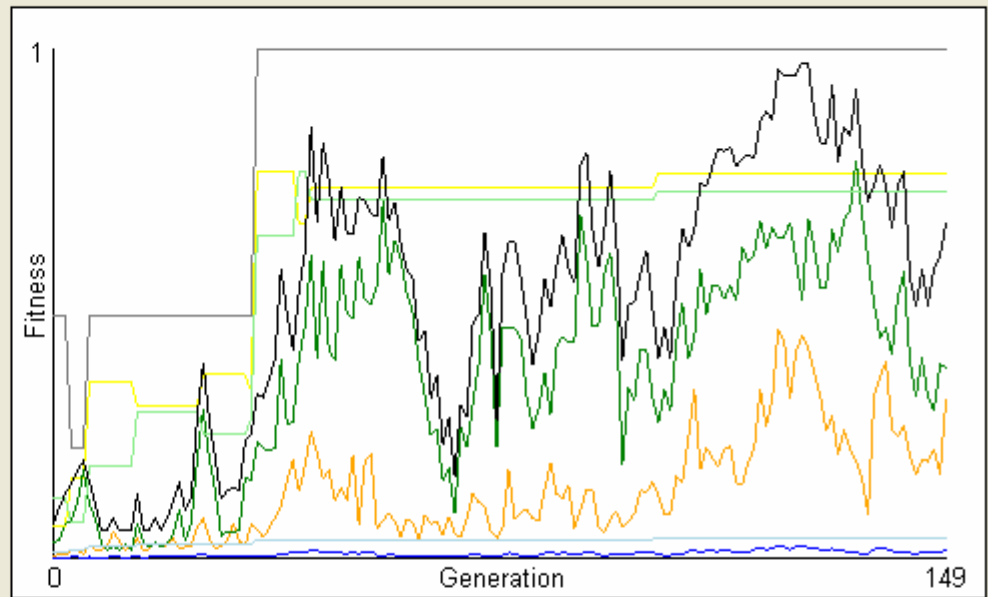
Display

	Max.	Avg.		Max.	Avg.
Total Fit.	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Rep. Suc.	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Light Inter.	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Surf. Area	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

Start

Stop

Results



Best Phenotypes

L: 15.79 %	R: 6.95 %	S: 21.76 %	Total: 1.71
L: 34.52 %	R: 18.05 %	S: 47.57 %	Total: 2.36
L: 29.93 %	R: 28.68 %	S: 47.57 %	Total: 2.62
L: 36.07 %	R: 24.49 %	S: 47.57 %	Total: 2.65
L: 32.87 %	R: 27.14 %	S: 47.57 %	Total: 2.66
L: 76.06 %	R: 63.20 %	S: 100.00 %	Total: 3.64
L: 65.74 %	R: 76.11 %	S: 100.00 %	Total: 3.68
L: 72.91 %	R: 70.42 %	S: 100.00 %	Total: 3.72
L: 75.59 %	R: 72.21 %	S: 100.00 %	Total: 3.79
L: 71.57 %	R: 76.93 %	S: 100.00 %	Total: 3.80

Number of best phenotypes: 11

Last found at gen.: 149

Save All...

Distributing the model

Change key method calls to messages

 Create_Seq() → request message 'Create_Seq'

Distribute workers that understand these messages

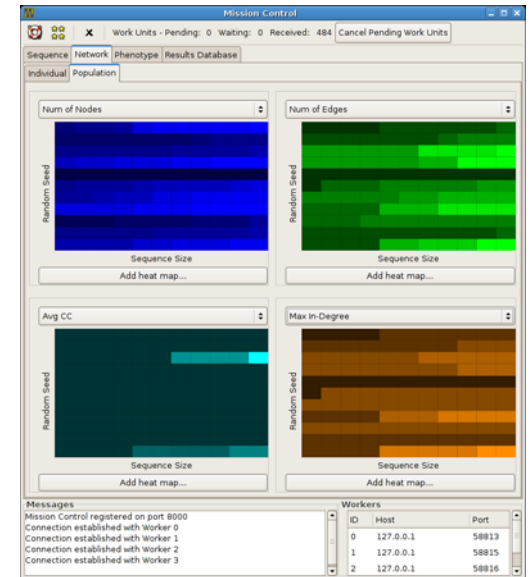
Maintain list of requests and their parameters

Send requests to idle workers

Update GUI as results are returned

Prototype

- Minimal image impact
 - Workers are stand-alone binaries
 - Bundle simulation code / data
 - Linux & Windows
- Work requests made dynamically according to GUI interaction
- Worker machines can be dynamically added / removed
- No administrative privileges required
 - no services, frameworks, etc.
 - but need network access to machines
- Practical, but no replacement for existing distributed tools
 - Manual fault recovery
 - No sophisticated scheduling (for this, see Condor / Nimrod)



Mission Control [Work Units - Pending: 0 Waiting: 0 Received: 0 Cancel Pending Work Units]

Sequence | Network | Phenotype | Results Database

Parameters

Seed: 12 Size: 10000 Base: 4 Length: 6 Inhib: 1

Individual | Population

genes:

Regulatory region lengths

Frequency

Length

Generate from Params | Apply Changes | Highlight Sequence | Generate Network -->

Messages

Mission Control registered on port 8000
 Connection established with Worker 0
 Connection established with Worker 1
 Connection established with Worker 2
 Connection established with Worker 3
 Connection established with Worker 4
 Connection established with Worker 5
 Connection established with Worker 6
 Connection established with Worker 7
 Connection established with Worker 8
 Connection established with Worker 9
 Connection established with Worker 10
 Connection established with Worker 11

Workers

ID	Host	Port
0	130.102.74.134	4609
1	130.102.73.193	4902
2	130.102.73.201	4751
3	130.102.73.197	4706
4	130.102.73.182	4081
5	130.102.73.196	4555
6	130.102.73.203	4703

Mission Control

Work Units - Pending: 0 Waiting: 0 Received: 0 Cancel Pending Work Units

Sequence Network Phenotype Results Database

Parameters

Seed: 12 Size: 10000 Base: 4 Length: 6 Inhib: 1

Individual Population

genes:

Regulatory region lengths

Frequency

Length

Generate from Params Apply Changes Highlight Sequence Generate Network -->

Messages

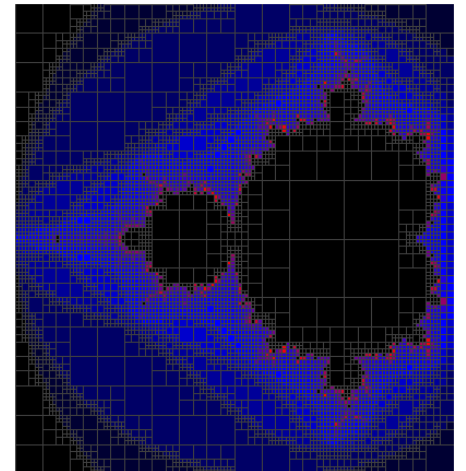
Mission Control registered on port 8000
Connection established with Worker 0

Workers

ID	Host	Port
0	127.0.0.1	4574

Work in progress

- Integration with Condor
 - Resilient to failure
 - Improved CPU monitoring
- Automated searches of parameter space
 - Nimrod/O integration



Distributed computing is worthwhile!

- Interactivity can be incorporated into distributed models of biology
- Not all models easily distributed
- It's more effort
 - Distribute design from the start
 - Machine failures, shared usage, etc.
 - Getting easier
- But can greatly expand range of feasible runs, collaborative value

Acknowledgements

- ARC Centre for Complex Systems
- ARC Centre in Bioinformatics
- Queensland Cyber Infrastructure Foundation
- Jon Kloske and Rik Taylor (ITIG)

Further Information

James Watson

jwatson@itee.uq.edu.au

<http://www.itee.uq.edu.au/~jwatson/>