Assembling a Cassava Transcriptome using Galaxy on a High Performance Computing Cluster



Aobakwe Matshidiso

Supervisor: Prof Chrissie Rey Co-Supervisor: Prof Scott Hazelhurst

Next Generation Sequencing

The Good



Next Generation Sequencing The Bad The Challenges A plethora of data processing tools Millions of short reads **Steep Learning Curves** Large memory requirements Gigabytes of disk space



- Investigate the efficiency of transcriptome assembly and alignment tools
- Use the above-mentioned assembly and alignment tools to establish a cassava transcriptome



- RNA extracted from leaf tissue at three time points: 12, 32, 67 days post inoculation with South African Cassava Mosaic Virus (SACMV)
- From CMD-resistant TME3 and CMD-Susceptible T200 cassava cultivars
- Using SOLiD 4 System from Applied Biosciences

The High Performance Computing Platform

- Ubuntu "Lucid" Virtual Machine
- Eight 3 Ghz Processing Cores, 72 GB RAM
- accessible by SSH over high-speed internet
- On top of the ZA-Wits-Core Cluster: 100 processor cores
- MPI Parallelization

Tools

- DE NOVO ASSEMBLERS De Bruijn Graphs: ABySS, Velvet Overlap Layout Consensus: Shore, Shorty Greedy: SSAKE, SHARCGS
- REFERENCE-BASED ALIGNMENT BWT: Bowtie, BWA Reads Hashing: SHRiMP, PerM, SOCS Genome Hashing: PASS, MOSAIK

225+193= GLASBERGEN

Copyright 2005 by Randy Glasbergen. www.glasbergen.com

"You have to solve this problem by yourself. You can't call tech support."

Workflow Pipeline: Galaxy

- Workflow planning
- Centralize the Investigation
- Easy to install, Easy to configure [???]
- Large Datasets
- Track Workflow Histories

Workflow Pipeline: Galaxy



Quality Control

Dataset: EA – Resistant, 12 dpi, SACMV negative

	Raw Forward	Raw Reverse				
# of Reads	58, 133, 361	58, 133, 361				
% GC	48%	50%				
Sequence Length	50	35				

Overrepresented Sequences

Sequence	Count	Percentage
CAAAACGACTCTCGGCAACGGATAT CTCGGCTCTCGCATCGATGAAGAA	884 830	2%
AACGACTCTCGGCAACGGATATCTC GGCTCTCGCATCGATGAAGAACGC	339 864	1%

Quality Control

EA RAW Forward Reads

Quality Scores

read position

Alignment Statistics

Bowtie Alignment: EA Raw with Default Settings

No. CPU Cores	1	2	4	6	8	
User time	31, 261s	31, 179s	31, 869s	33, 322s	33, 530s	
System Time	54.33 s	48s	52.48s	53.21s	62.66s	
% of CPU	99	199%	399%	598%	791%	
Wall Clock Time	8:43:58	4:21:20	2:13:16	1:32:56	1:10:41	
Reads Processed	58, 133, 361	58, 133, 361	58, 133, 361	58, 133, 361	58, 133, 361	
Reads with Alignment	193	193	193	193	193	
Reads with Failed Alignments	58, 133, 168	58, 133, 168	58, 133, 168	58, 133, 168	58, 133, 168	

DBG Assemblers

	Ν	Min	Med	Mean	Мах	N50	N50 length
Velvet19	743, 122	37	41	44	256	310, 736	42
Velvet25	41, 927	49	45	69	436	14, 497	74
ABySS 19	19, 763, 844	19	19	20	382	8, 928, 450	19
ABySS 25	6, 333, 943	25	26	28	378	2, 791, 286	26

Alignment Tools

Time Lapsed per processor cores

De Bruijn Assemblers

Aligners

% of CPU Used per Processor Core

Mapping Quality – Bowtie, EA Enhanced, Default Settings

Mapping Quality – Bowtie EA Enhanced, Strict/ Enhanced Settings

MAPQ >= 30 (51.5%, 44485825)

MAPQ < 3 (0.0%, 0)

Unmapped (48.5%, 41845688)

Mapping Quality – Bowtie Enhanced Data, Default Settings, Normal Reference

Bowtie Enhanced Data, Enhanced Settings, Normal Reference

Mapping Quality – BWA Enhanced Data, Default Settings, Softmasked Reference

Bowtie Enhanced Data, Enhanced Settings, Normal Reference

The Transcriptome

184,763 to 184,820 (58

182,663 to 207,662 (25 Kb)

K F F W	G	F	e e	÷	÷	F	G	G	a	G	G	G	F I	q
	G G I I I							a a a	42.42.41	a a a a				a a A
184,763 U184,763	4				1.7		-	+		T.		184	,820 L	J184,82
									a <td></td> <td>1 Q Q 2 Q Q 2 Q Q 3 Q Q 4 Q Q</td> <td></td> <td></td> <td></td>		1 Q Q 2 Q Q 2 Q Q 3 Q Q 4 Q Q			

Visualization of reads mapped to a scaffold in Tablet 2.6

Quantile Normalization

Gene Counts Across Datasets

Normal State versus Disease Challenged (Resistant vs Susceptible, 67 dpi)

Naïve Expression Analysis

Distribution of Log2 Ratios of Gene Expression Counts

Normal State versus Disease Challenged (Resistant vs Susceptible, 67 dpi)

In Progress

- Gaining new knowledge relating to host-pathogen interactions, in the cassava-SACMV relationship
- Proceeding to build a cassava diseasechallenged transcriptome
- Annotation of the cassava genome in progress

Conclusion

- Quality Control is Important, Pre-processing is critical
- Galaxy can enhance and facilitate bioinformatics analyses
- It's important to optimise your software to obtain quality results
- A High Performance system will greatly enhance your studies

Funding

