

Improving HPC resource utilization in the genome assembly of a biofuel producing green algae

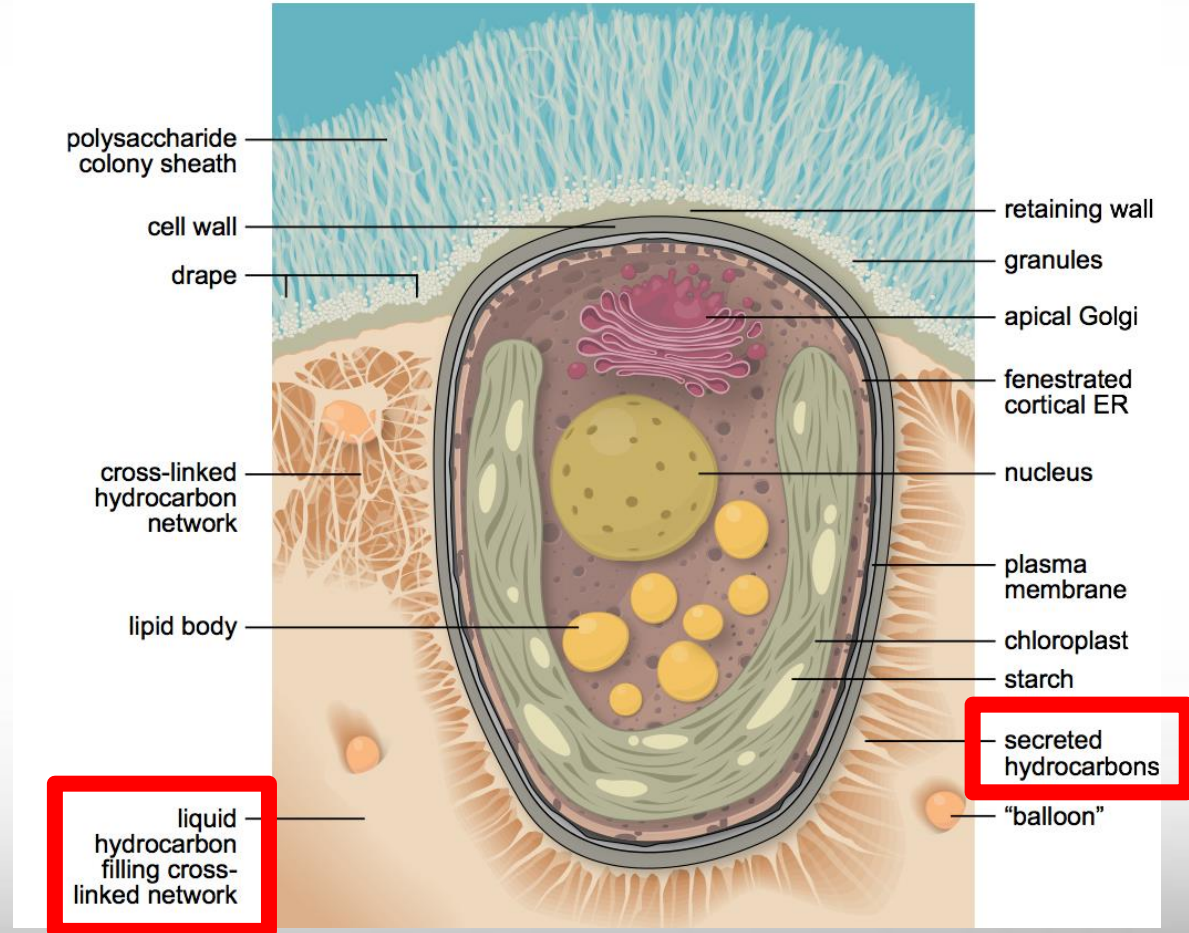
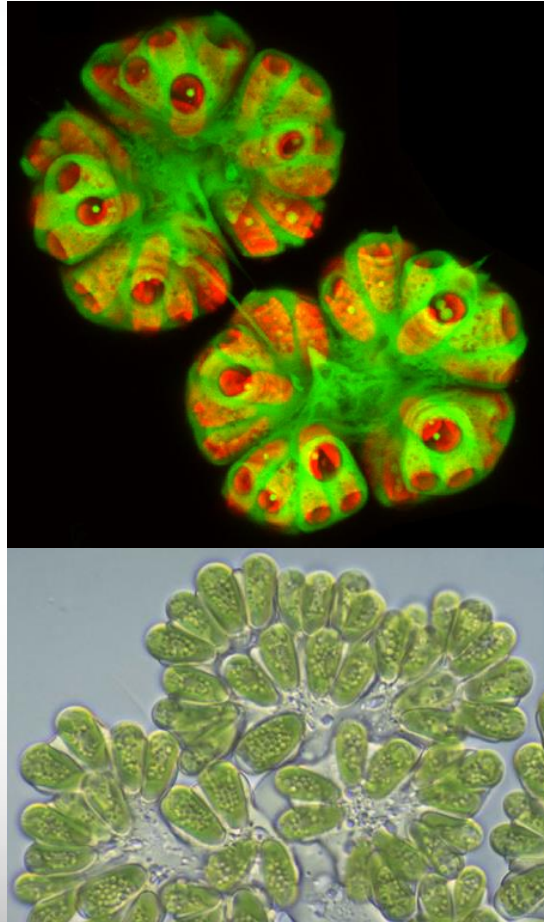
Project by Dan Browne, PhD Candidate, Devarenne Lab, Biochemistry & Biophysics
Presented by Michael Dickens, High Performance Research Computing
Texas A&M University



Improving HPC resource utilization in the genome assembly of a biofuel producing green algae

Dan Browne, PhD Candidate, Devarenne Lab, Biochemistry & Biophysics Department, Texas A&M University

Basic model of *Botryococcus braunii* cell biology



Improving HPC resource utilization in the genome assembly of a biofuel producing green algae

Dan Browne, PhD Candidate, Devarenne Lab, Biochemistry & Biophysics Department, Texas A&M University

Why sequence the *B. braunii* genome?



- *B. braunii* is a potential source of renewable fuels and chemicals
- *B. braunii* is found worldwide, most notably in oil and coal shale deposits
- *B. braunii* has a very high oil content, ~40% of dry weight
- *B. braunii* oils can be processed with conventional petroleum technology

Main project organizers:



Andy Koppisch
Northern Arizona University



Joe Chappell
University of Kentucky



Tim Devarenne
Texas A&M University



Shigeru Okada
Tokyo University



Improving HPC resource utilization in the genome assembly of a biofuel producing green algae

Dan Browne, PhD Candidate, Devarenne Lab, Biochemistry & Biophysics Department, Texas A&M University

B. braunii whole-genome sequencing with Illumina

| Library Name | Library Type | Insert Size | Total Sequence Reads | Read Length | Genome Size | Coverage |
|--------------|--------------|-------------|----------------------|-------------|-------------|----------|
| SXPX | Paired End | 800 bp | 499,073,402 | 250 bp | 166 Mb | ~750x |

Genome sequence can be used to identify genes involved in hydrocarbon production



Genomic DNA

← AATAATGTCAATTTGGTAGATATCAGAGAGTTTTATGTTGACAAAGATGG

| | | |
|--------------|--------------|-------------|
| AATAATGTCA | GATATCAGAGA | ATGTTGACAAA |
| AATAATGTCAA | GATATCAGAGAG | ATGTTGACAAA |
| ATAATGTCAAT | TATCAGAGAGT | GTTGACAAAG |
| ATAATGTCAAT | TATCAGAGAGT | GTTGACAAAG |
| TAATGTCAATT | TCAGAGAGT | TTGACAAAGAT |
| TGCAATTTGG | CAGAGAGT | TTGACAAAGAT |
| TGCAATTTGGT | CAGAGAGT | TGACAAAGATG |
| AATTTGGTAGAT | GAGAGT | GACAAAGATGG |
| TTGGTAGATAT | | CAAAGATGG |
| TGGTAGATATC | | AAAGATGG |

Computational Assembly

→ AATAATGTCAATTTGGTAGATATCAGAGAGTNNNATGTTGACAAAGATGG

Reconstructed DNA Sequence



Improving HPC resource utilization in the genome assembly of a biofuel producing green algae

Dan Browne, PhD Candidate, Devarenne Lab, Biochemistry & Biophysics Department, Texas A&M University

B. braunii whole-genome sequencing with Illumina

| Library Name | Library Type | Insert Size | Total Sequence Reads | Read Length | Genome Size | Coverage |
|--------------|--------------|-------------|----------------------|-------------|-------------|----------|
| SXPX | Paired End | 800 bp | 499,073,402 | 250 bp | 166 Mb | ~750x |

Genome sequence can be used to identify genes involved in hydrocarbon production



Genomic DNA

AATAATGTCAATTTGGTAGATATCAGAGAGTTTTATGTTGACAAAGATGG

```

AATAATGTCA      GATATCAGAGA      ATGTTGACAAA
AATAATGTCAA      GATATCAGAGAG     ATGTTGACAAA
ATAATGTCAAT      TATCAGAGAGT      GTTGACAAAG
ATAATGTCAAT      TATCAGAGAGT      GTTGACAAAG
TAATGTCAATT      TCAGAGAGT        TTGACAAAGAT
TGCAATTTGG       CAGAGAGT         TTGACAAAGAT
TGCAATTTGGT      CAGAGAGT         TGACAAAGATG
                  AATTTGGTAGAT     GAGAGT           GACAAAGATGG
                  TTGGTAGATAT      CAAAGATGG
                  TGGTAGATATC      AAAGATGG
    
```

Computational
Assembly

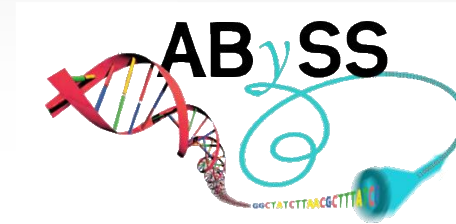
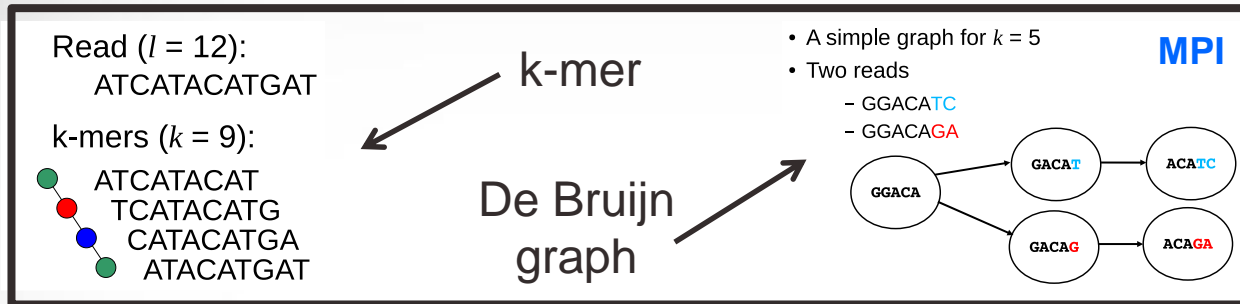
AATAATGTCAATTTGGTAGATATCAGAGAGTNNNATGTTGACAAAGATGG

Reconstructed DNA Sequence

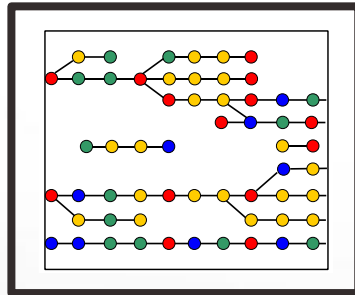


Workflow of Assembly By Short Sequences (ABySS): A parallel *de novo* genome assembler with MPI support

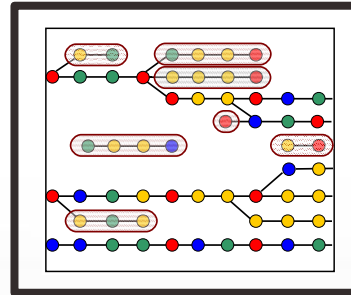
(1) ABYSS-P



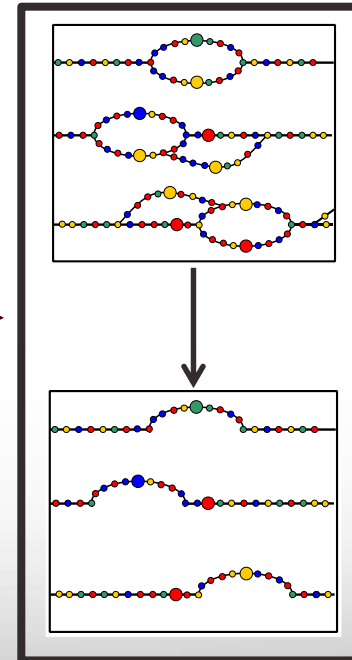
(2) AdjList



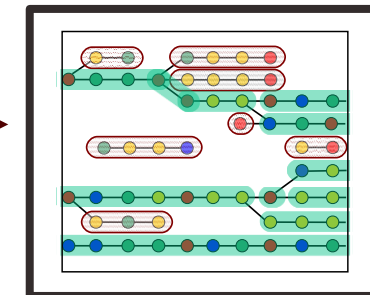
(3) Prune tips



(4) Pop bubbles



(5) Generate contigs



map 1

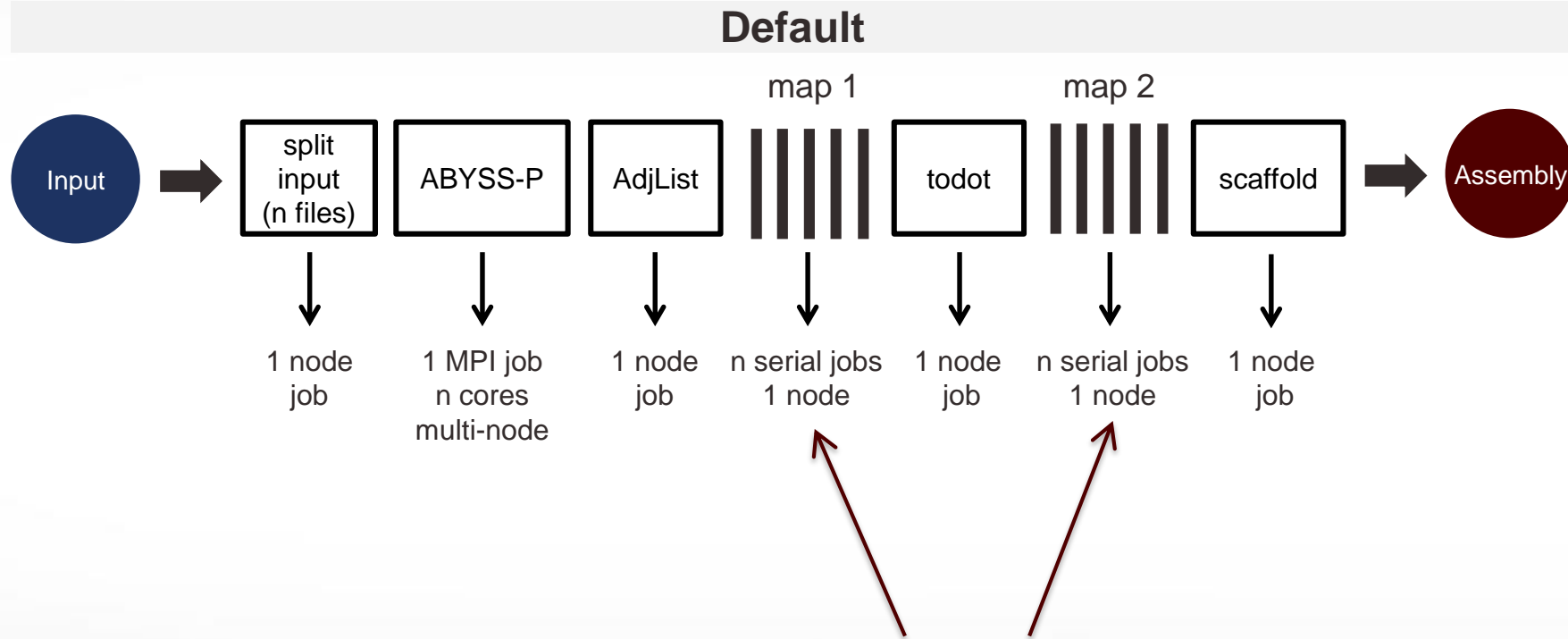
map 2

scaffold

Improving HPC resource utilization in the genome assembly of a biofuel producing green algae

Dan Browne, PhD Candidate, Devarenne Lab, Biochemistry & Biophysics Department, Texas A&M University

Default and modified ABySS execution pipelines

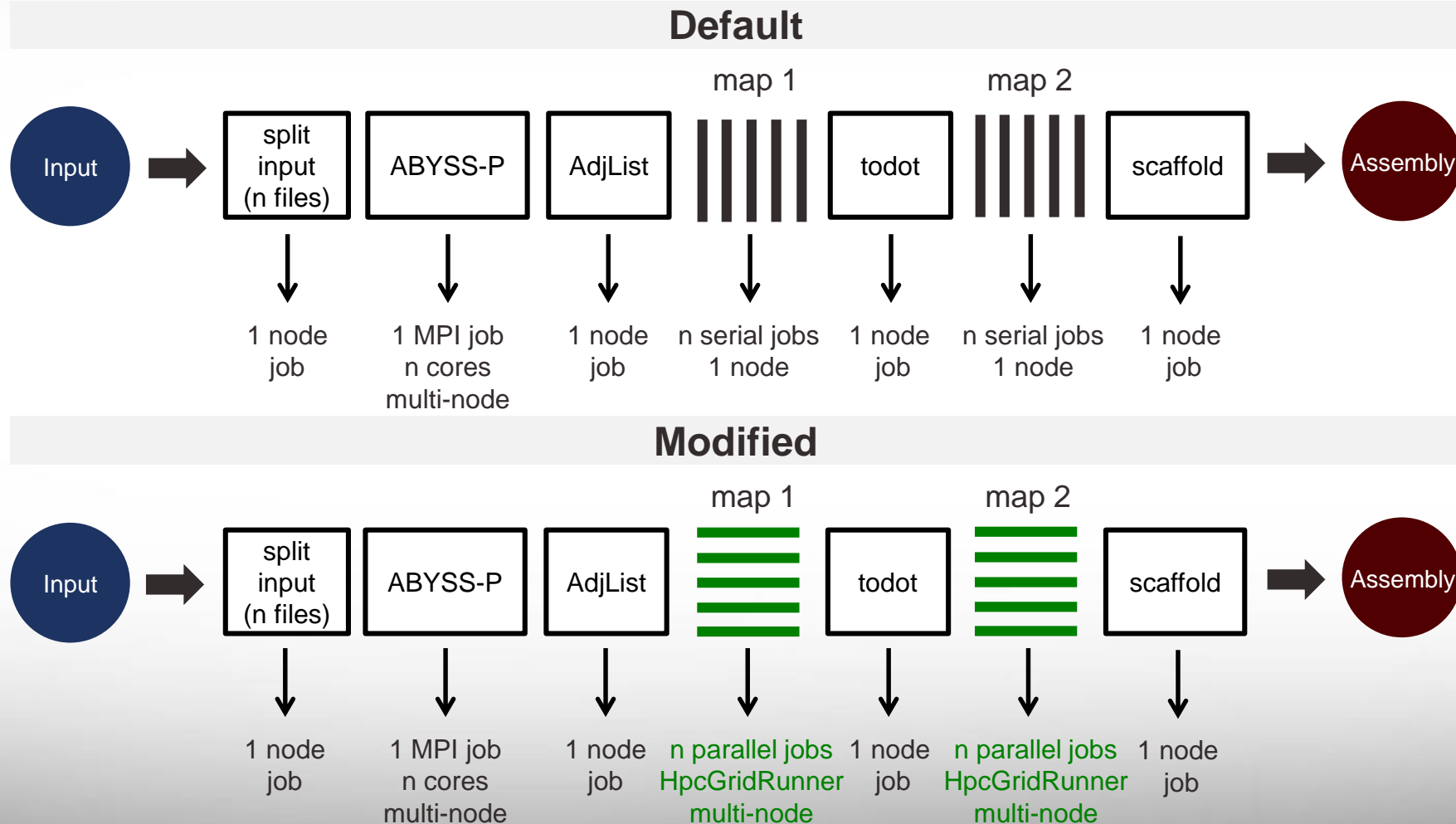


All commands of each mapping step run in serial and limited to one compute node

Improving HPC resource utilization in the genome assembly of a biofuel producing green algae

Dan Browne, PhD Candidate, Devarenne Lab, Biochemistry & Biophysics Department, Texas A&M University

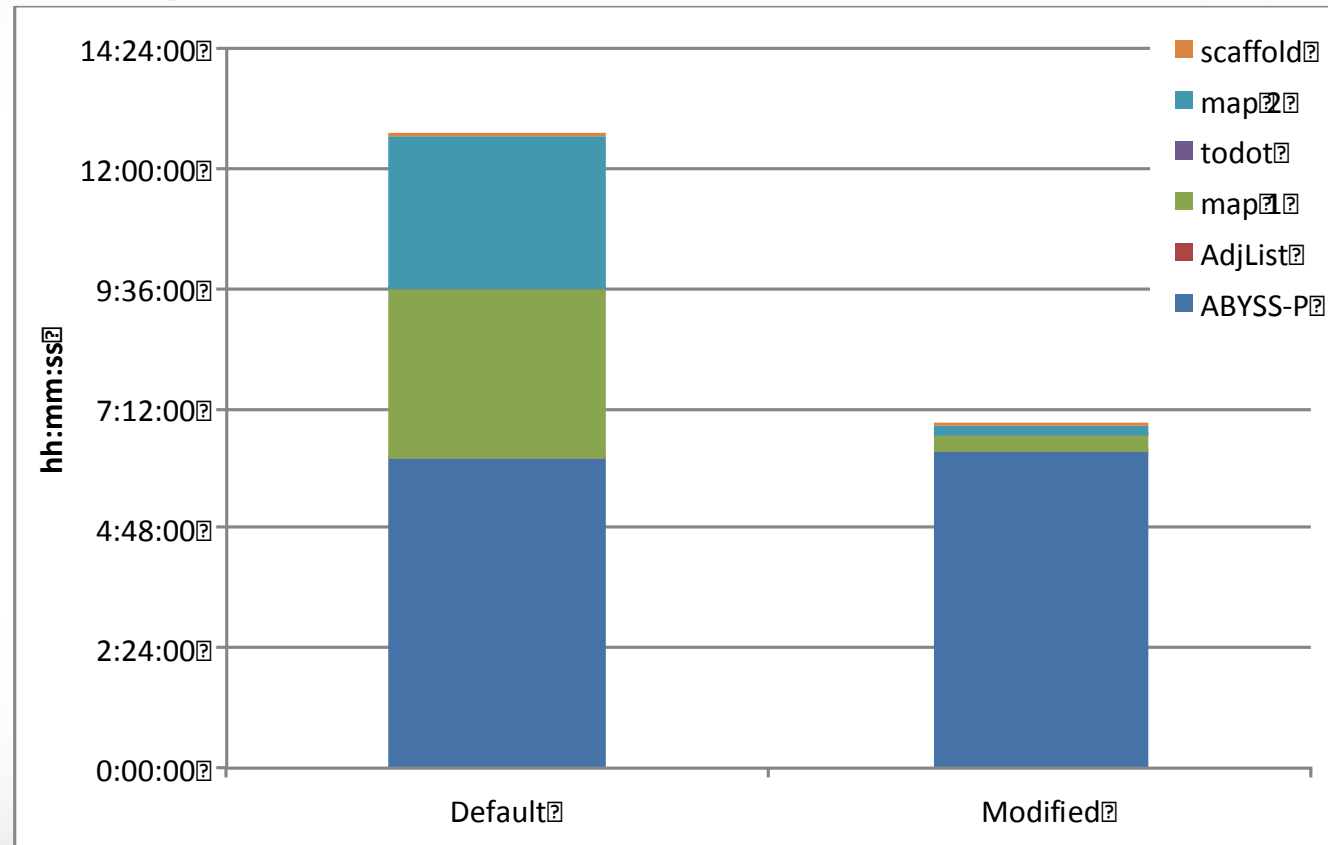
Default and modified ABySS execution pipelines



Improving HPC resource utilization in the genome assembly of a biofuel producing green algae

Dan Browne, PhD Candidate, Devarenne Lab, Biochemistry & Biophysics Department, Texas A&M University

Assembly times of default and modified pipelines



- HPC resource utilization: 50 cores (5 cores/node * 10 nodes)
- Assembly time reduced by 46% using modified ABySS pipeline.
- Modified pipeline eliminated 45 cores being idle for almost 6 hours.

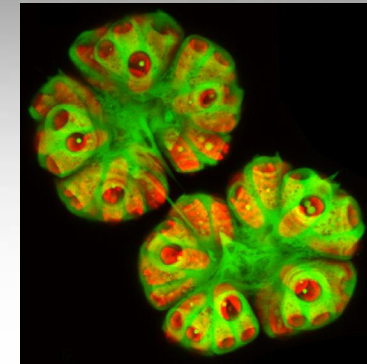




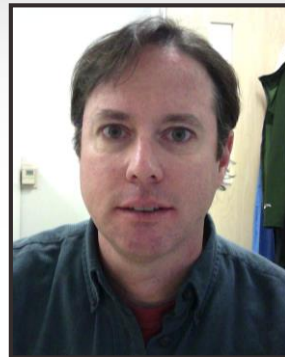
Department of
Biochemistry &
Biophysics

Devarenne Lab 2015

<http://devarennelab.tamu.edu>



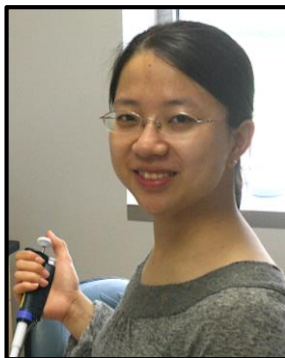
*Botryococcus
braunii*



Tim Devarenne, PhD
Associate Professor



Hem Thapa
Grad Student



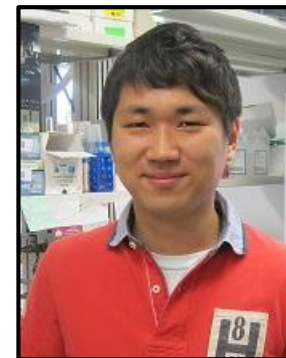
Dongyin Su
Grad Student



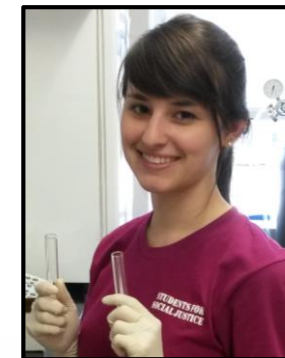
Mehmet Tatli
Grad Student



Dan Browne
Grad Student



Incheol Yeo
Grad Student



Victoria Yell
Undergrad Student



U.S. DEPARTMENT OF
ENERGY

USDA



NIFA

