

Charles A. Goodman and Charles F. Delwiche
Cell Biology and Molecular Genetics, University of Maryland

Introduction

Marine and freshwater environments are vital resources to all life on earth and are among the biomes most severely affected by human impact, and the ability to responsibly manage these ecosystems is impeded by a need for deeper understanding of their diverse networks of interactions. Green algae are ecologically important as basal members of food webs¹ and in terms of net primary productivity^{2,3}. Clade-specific trends in invasion and competition have previously been demonstrated^{4,5} (Figure 1), though community dynamics among primary producers are poorly understood, and historic paradigms don't adequately explain trends in competition^{5,6}. As part of a larger collaborative effort to determine the potential for phylogenetic relatedness to predict outcome of competition among planktonic green algae, RNAseq data were collected to track gene expression through a series of inter- and intra-specific competition growth curves. Here, we present eight de novo transcriptome assemblies and corresponding timecourse differential gene expression analyses from species broadly distributed across the green-algal tree. Our results show distinct phases of gene expression across the monoculture growth curves, reflecting the effects of intraspecific competition. In addition, they indicate substantial differences among species in patterns of gene expression within commonly shared orthogroups. These observations provide baseline patterns of expression for a planned subsequent analysis of biculture competition data.

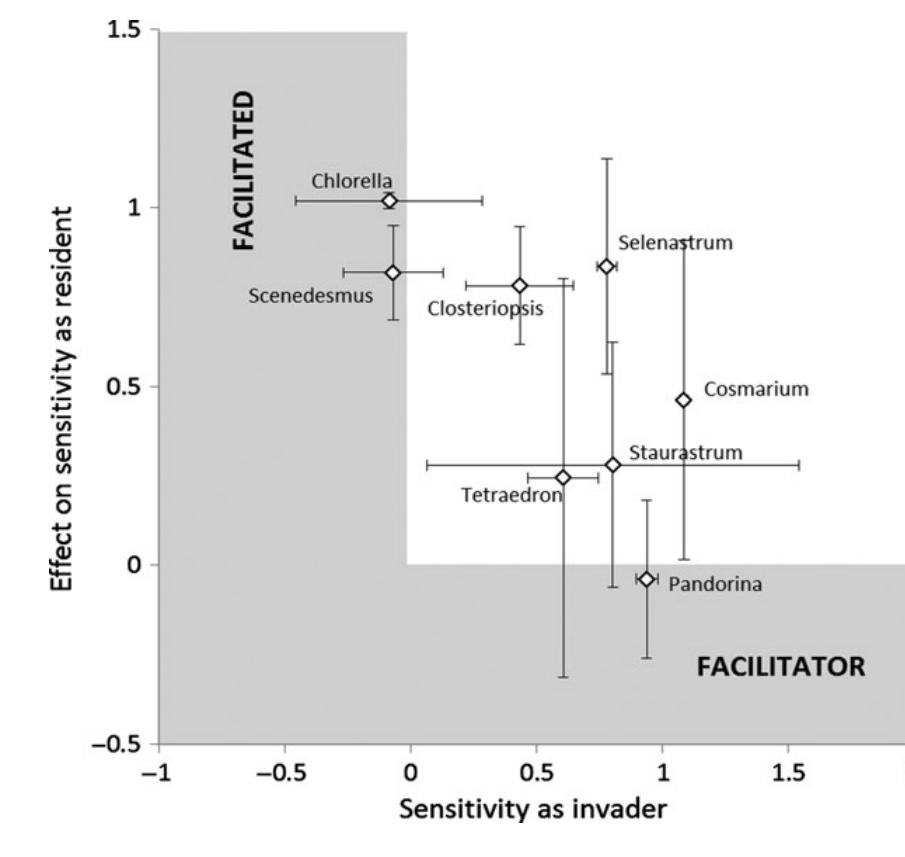


Figure 1. Sensitivity to Invasion - from Venail et al. 2014

A related study focused on the effects of invasion on growth. The X-axis shows effect on growth rate when invading an established steady-state culture: increasing sensitivity indicates a decreasing growth rate as compared to monoculture. Values approaching 0 indicate an aggressive invader, values approaching 1 indicate decreasing ability to invade.

$$Si = (g[\text{mono}_i] - g[\text{inv}_i])/g[\text{mono}_i]$$

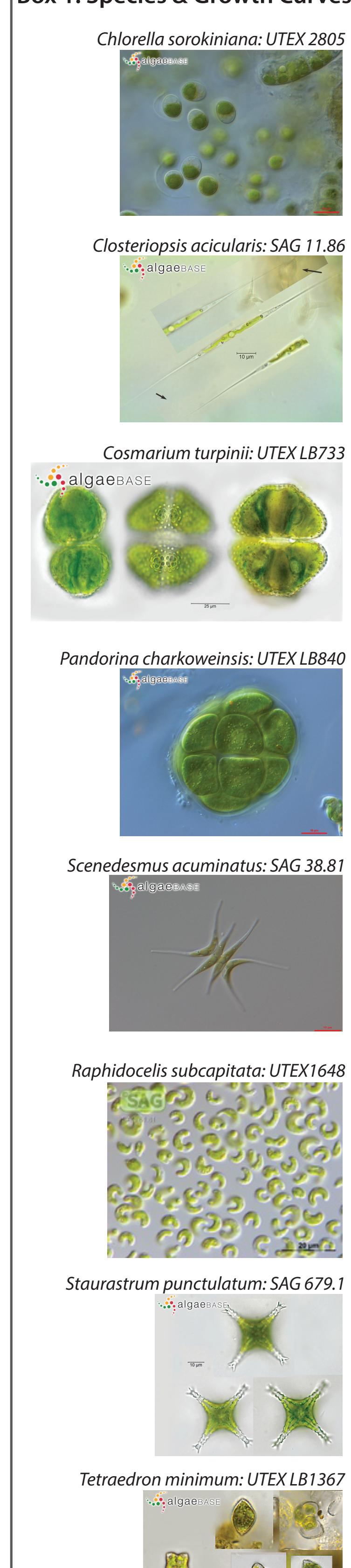
The Y-axis is an aggregate measure of how much each species affected invader's growth rate. A high effect suggests that invader's rates were not affected, a low effect indicates that invader's rates were affected.

$$Ei = (S_i[\text{mono}_i] + S_i[\text{inv}_i] + \dots + S_i[\text{mono}_j]) / n$$

Detecting Orthology

While looking at reciprocal best blast hits is a known strategy in establishing orthology, evidence suggests it isn't universally reliable in predicting nearest-neighbor relationships^{7,8}. I therefore used OrthoFinder to identify and bin orthologous sequence across our eight species – this has the added advantage of considering same-species orthology and gene isoforms. OrthoFinder's basic output bins genes by "orthogroups", which does not automatically distinguish ortholog from paralog, but does provide a factor for doing so with additional phylogenetic analysis. I've identified 3156 orthogroups which are common to all 8 species, with between 38.4-50.1% of genes belonging to these groups, and also identified a number of species-unique orthogroups. Groupings of shared OrthoGroups are shown in Figure 4.

Box 1. Species & Growth Curves



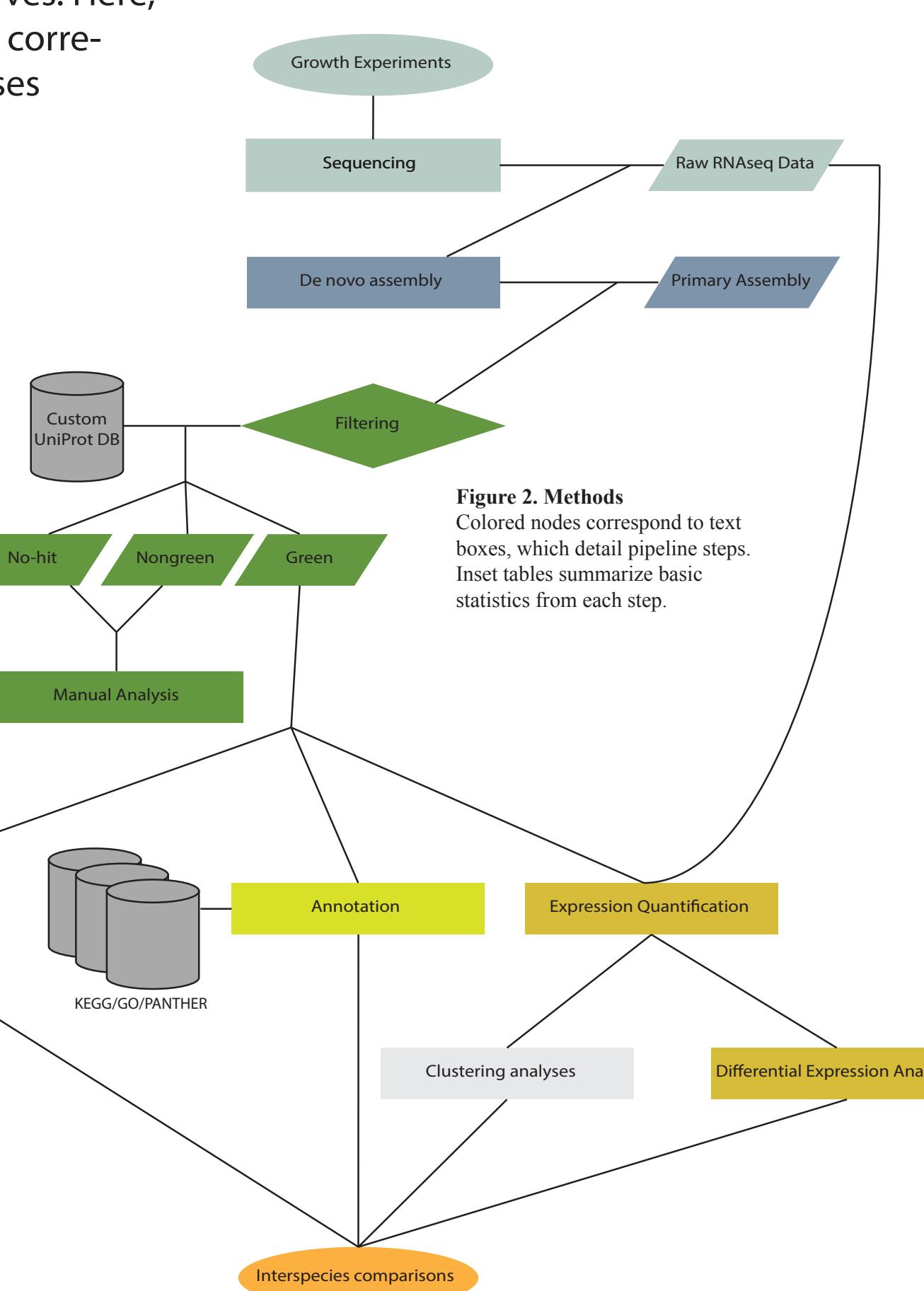
cagood@umd.edu, delwiche@umd.edu

Link: C. Goodman GitHub repositories

Currently, data and pipeline steps are privately held, pending publication.

Growth

This study is a novel analysis of RNAseq data coming from competition growth assays done in collaboration with Narwani et al. (2017). Eight species of unicellular, planktonic, freshwater green algae (table 1) were selected for their co-occurrence in nature, compatible culture requirements, and visually distinct morphologies. Species were grown as monocultures and in their 28 possible pairwise combinations, in triplicate (Figure 3a). 100-ml bottles containing enriched COMBO growth medium were inoculated at 200 cells ml⁻¹, and cultures were grown at 20°C under a 16h:8h light-dark cycle for 46 days. 10% of each culture was removed every other day, used for density monitoring (via chlorophyll fluorescence), cell/colony-counts (via flow-cam), and preservation for RNA extraction. Growth curves were calculated based on cell-count data. RNA was extracted and Illumina HT sequencing was performed. For monocultures, RNA was sequenced from four time-points along the 46-day growth curve, chosen to reflect early (T1), pre-infection log phase (T2), post-infection log phase (T3), and carrying-capacity (T4) growth densities (Figure 3b). The scope of this analysis includes only monoculture data; bicultures will be assessed in future efforts. Species images and growth curves are shown in Box 1.



De novo assembly

Reference genomes remain unavailable for these taxa, and so we rely on de novo assembly. To generate reference transcriptomes, paired-end sequence data from all time-points were concatenated into a single pair of fastq files per species. FastQC was used to estimate sequence quality; reads with quality threshold and artificial adapter sequences were trimmed from the set using bbduk⁸. RNAspades⁹ was used for de novo assembly, and transcriptome completeness was assessed using BUSCO¹⁰ Chlorophyte and Eukaryote databases.

"Green" Filtering

We obtained transcript annotations with an integrated filtering step, using BLAST^{11,12} to compare each transcript against a protein database containing both "green" plant and algal sequence, and "nongreen" sequence from representative taxa across the rest of the tree. Protein data were drawn from UniProtKB's reviewed Swiss-Prot and unreviewed TrEMBL databases^{13,14}. Queries that preferentially matched non-green subject sequences were separated from the green dataset for additional processing; particular attention will be paid to non-green data that indicates culture contamination and/or fellow-travelers. Queries that preferentially matched green subjects were carried forward in the presented analysis.

Quantification and Differential Expression Analysis

Culture densities were previously measured as described in Narwani et al. 2017. Preliminary growth curves were estimated using the original data, via GrowthCurver²³, an R package which fits cell density data to a standard logistic curve. Improved growth models are currently under consideration, including the use of a discrete (rather than continuous) curve, as well as accounting for periodic harvesting²⁴. To quantify gene expression, I employed Kallisto²⁵, chosen for its use of a hash-based quantification approach which has benchmarked well against other common methods, running orders of magnitude more quickly than comparable software packages^{26,27} and integrating bootstrapped estimates of expression variance. I obtained counts for each of the monocultures by counting trimmed reads against their respective reference transcriptomes, running Kallisto with 100 bootstrapped replicates.

I used sleuth^{24,27,28} to identify significant differential expression of transcripts between all consecutive pairs of time points in each dataset. Sleuth is unique among DE analytical packages in that it combines biological and "inferential" variance in its estimate, basing the latter on the bootstrapped count values produced with Kallisto. I separately compared T1/T2, T2/T3, and T3/T4 – where the early time point was taken as the "control" value, and the late time point was the "test" value. Sleuth has been shown to overestimate false-discovery rate²⁹, though a conservative estimate is desirable as fewer genes are reported as "significant"; but those reported tend to be enriched for true differential expression. I define "significant DE" as FDR < 0.1 and |Z| > 1, and "constitutive expression" as transcripts for which |Z| >= 0.2 across all growth-time curve points. Summary figures are shown in Box 2, Fig. 3c, Box 3, and via QRC 1.

Table 1. Experimental IDs

Species	ID
<i>Chlorella sorokiniana</i>	DC10
<i>Closteriopsis acicularis</i>	DC20
<i>Cosmarium turpinii</i>	DC30
<i>Pandorina charkweinensis</i>	DC40
<i>Scenedesmus acuminatus</i>	DC50
<i>Selenastrum capricornutum</i>	DC60
<i>Staurastrum punctulatum</i>	DC70
<i>Tetraedron minimum</i>	DC80

Table 2. Reference Transcriptome Assembly

ID	DC10	DC20	DC30	DC40	DC50	DC60	DC70	DC80
Reads (90bp)	2,538 ± 0.7	2,686 ± 0.7	2,856 ± 0.7	1,888 ± 0.7	2,926 ± 0.7	2,756 ± 0.7	2,426 ± 0.7	2,508 ± 0.7
Unique Contigs	80,151	124,149	219,575	110,843	118,654	153,536	221,016	243,395
Total Length (bp)	6,938 ± 0.7	8,476 ± 0.7	12,666 ± 0.8	9,686 ± 0.7	9,516 ± 0.7	7,076 ± 0.7	14,336 ± 0.8	13,561 ± 0.8
Coverage	32.8X	28.4X	20.6X	17.5X	27.6X	34.9X	15.2X	17.1X
N50 (bp)	1930	1842	1235	2309	2554	825	1090	1887
Missing BUSCOs [†]	8.63%	8.12%	28.18%	8.76%	10.65%	11.35%	43.73%	8.63%

[†] BUSCO database: chlorophyta_odb10, n= 2168

* Chlorophyte

Table 3. "Green" Filtered Transcriptomes

ID	DC10	DC20	DC30	DC40	DC50	DC60	DC70	DC80
Total Hits	34,657	32,769	62,993	37,852	32,120	56,073	61,529	45,492
Green Hits	27,362	25,820	40,883	33,141	25,297	48,508	54,073	20,273
Green Tot. Length (bp)	4,906 ± 0.7	5,016 ± 0.7	5,226 ± 0.7	6,746 ± 0.7	5,576 ± 0.7	3,936 ± 0.7	6,156 ± 0.7	3,656 ± 0.7
% of Reference Trans.	70.73%	59.13%	41.30%	69.59%	58.54%	55.53%	42.97%	26.97%
N50 (bp)	2,471	2,616	1,872	3,169	3,539	1,373	1,514	2,502
Green Missing BUSCOs [†]	9.36%	8.21%	29.38%	8.81%	10.70%	11.35%	49.31%	9.64%

[†] BUSCO database chlorophyta_odb10, n= 2168

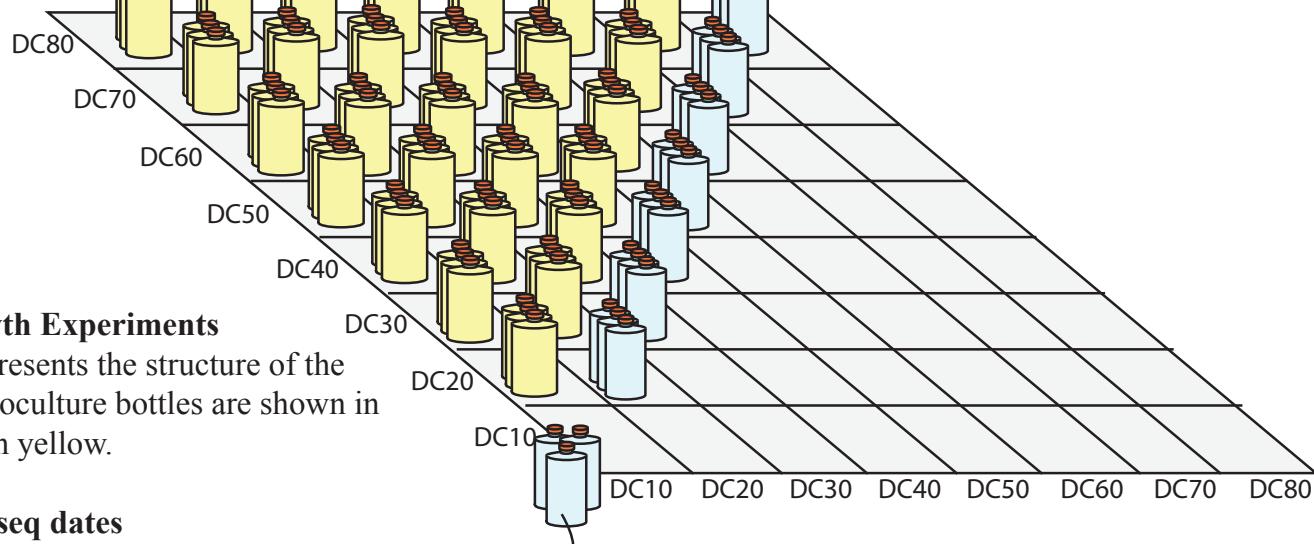
* Chlorophyte

Figure 3a. Growth Experiments

This diagram represents the structure of the experiment. Monoculture bottles are shown in blue, bicultures in yellow.

Figure 3b. RNAseq data

The red growth curve was calculated prior to choosing samples for sequencing. Sequenced RNA sample dates are shown in dark blue, and were chosen based on position in the curve.



Chlorella sorokiniana

Day

Cell Count

2e+06
1e+06
0e+00

0e+00
1e+06
2e+06

3e+06
4e+06

5e+06
6e+06

7e+06
8e+06

9e+06
1e+07

1e+07
2e+07
3e+07

4e+07
5e+07
6e+07

7e+07
8e+07
9e+07

1e+08
2e+08
3e+08

4e+08
5e+08
6e+08

7e+08
8e+08
9e+08

1e+09
2e+09
3e+09

4e+09
5e+09
6e+09

7e+09
8e+09
9e+09

1e+10
2e+10
3e+10

4e+10
5e+10
6e+10

7e+10
8e+10
9e+10

1e+11
2e+11
3e+11

4e+11
5e+11
6e+11

7e+11
8e+11
9e+11

1e+12
2e+12
3e+12

4e+12
5e+12
6e+12

7e+12
8e+12
9e+12

1e+13
2e+13
3e+13