# Unveiling The Hidden Life Of A Very Visible Phytoplankton:

## Use of normalized Sanger sequenced libraries and non-normalized 454 libraries to compare the transcriptomes of haploid and diploid *Emiliania huxleyi*.

**Peter von Dassow[1], Hiroyuki Ogata[2], Ian Probert[1], Stéphane Audic[2], Jean-Michel Claverie[2], Patrick Wincker[3], Corinne Da Silva[3] and Colomban de Vargas[1]**

[1]Evolution du Plancton et PaleOceans, Station Biologique de Roscoff, CNRS UPMC UMR7144, 29682 Roscoff, France
[2]Structural and Genomic Information, CNRS UPR2589, 13288 Marseille, France
[3]Genoscope Centre National de Sequençage, 91057 Evry, France

EPPO - Evolution du Plancton et Paleo-Oceans
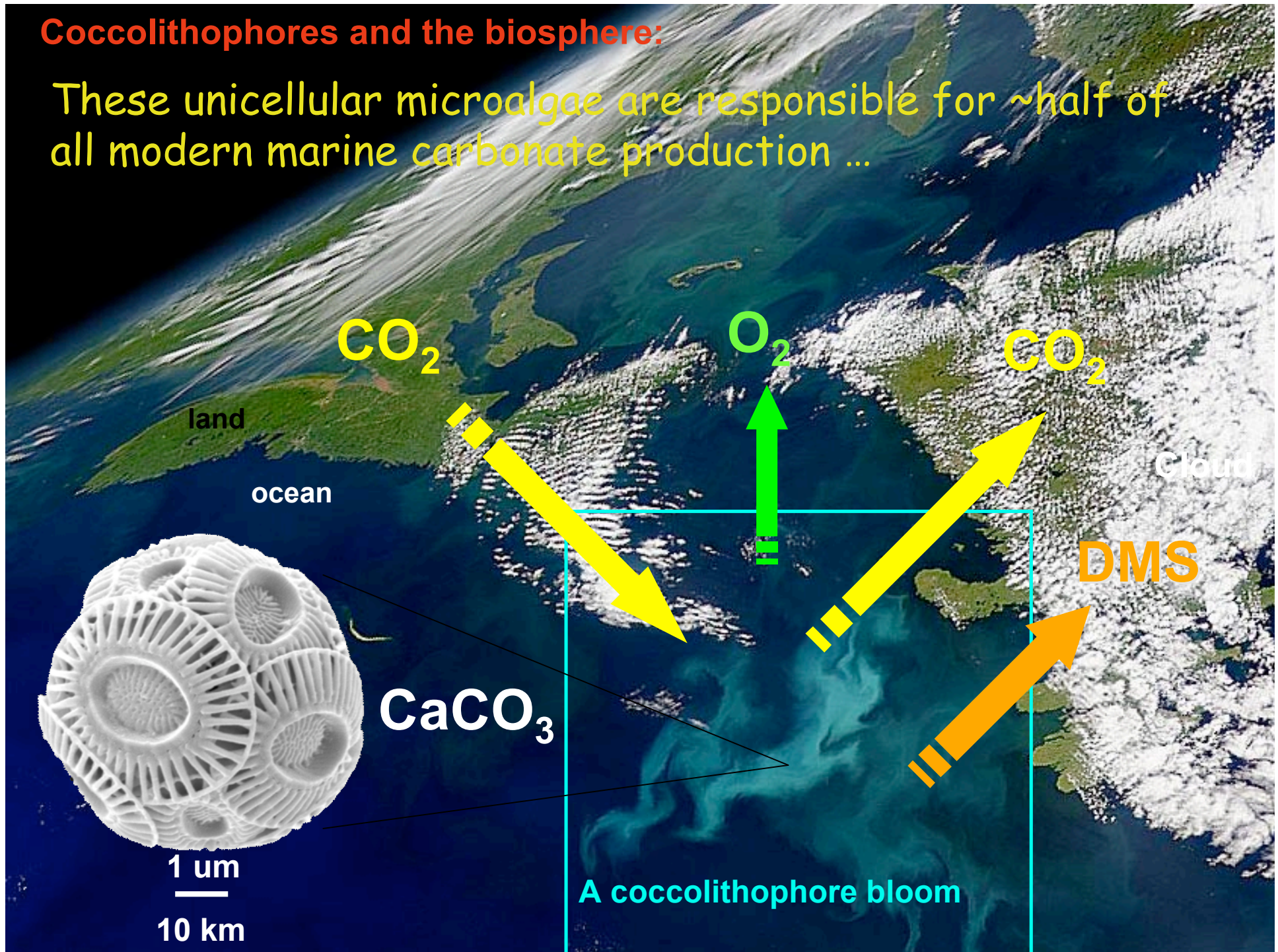
Sequencage haut-debit:
2008-2012

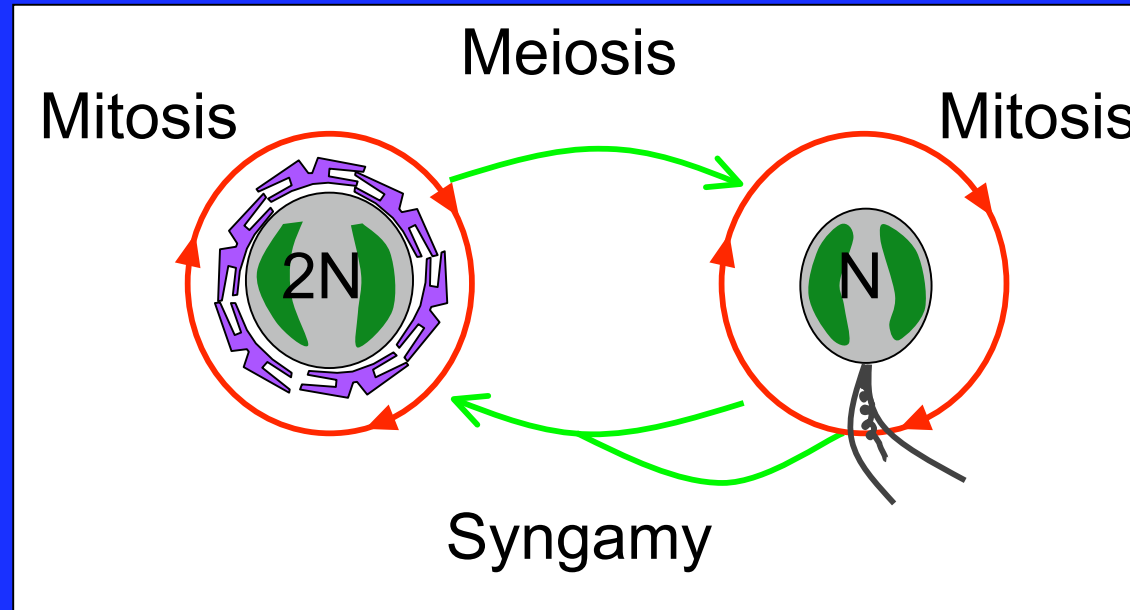The grand diversity of eukaryotic life:
Much more than animals, fungi, and plants!

**Coccolithophores and the biosphere:**

These unicellular microalgae are responsible for ~half of all modern marine carbonate production ...

$CO_2$

$O_2$

$CO_2$

land

Cloud

ocean

DMS

$CaCO_3$

1 um

A coccolithophore bloom

10 km

# *Emiliania huxleyi* life cycle



Diploid (2N):

Non-motile

Calcified

**Forms massive blooms**

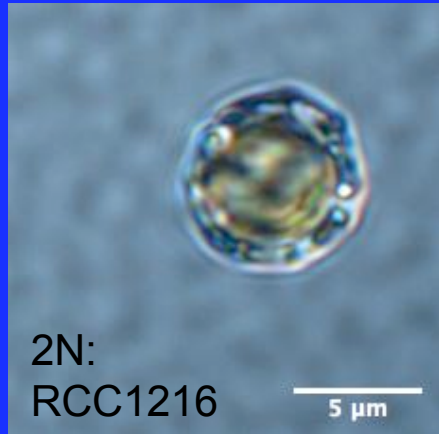**Not photoinhibited**

**Killed by *EhV*s**

Haploid (1N):

Motile

Non-calcified

**Doesn't bloom?**

**Photoinhibited**

**Resistant to *EhV*s**

Frada et al., 2008

# 1N (motile) and 2N (calcified) strains with same parent



2N:
RCC1216
5 µm



1N:
RCC1217
5 µm

Generate Sanger and 454 EST libraries from both.

Compare to JGI's genome assembly of a different *E. huxleyi* strain, CCMP 1516

# Two types of library combine sequence coverage and depth

Sanger-sequenced libraries

1. Oligo-dT primed cDNA

2. Normalized

3. ≈19000 longer reads

454-sequenced libraries

1. random primed cDNA

2. *NOT* normalized

3. ≈255000 shorter reads
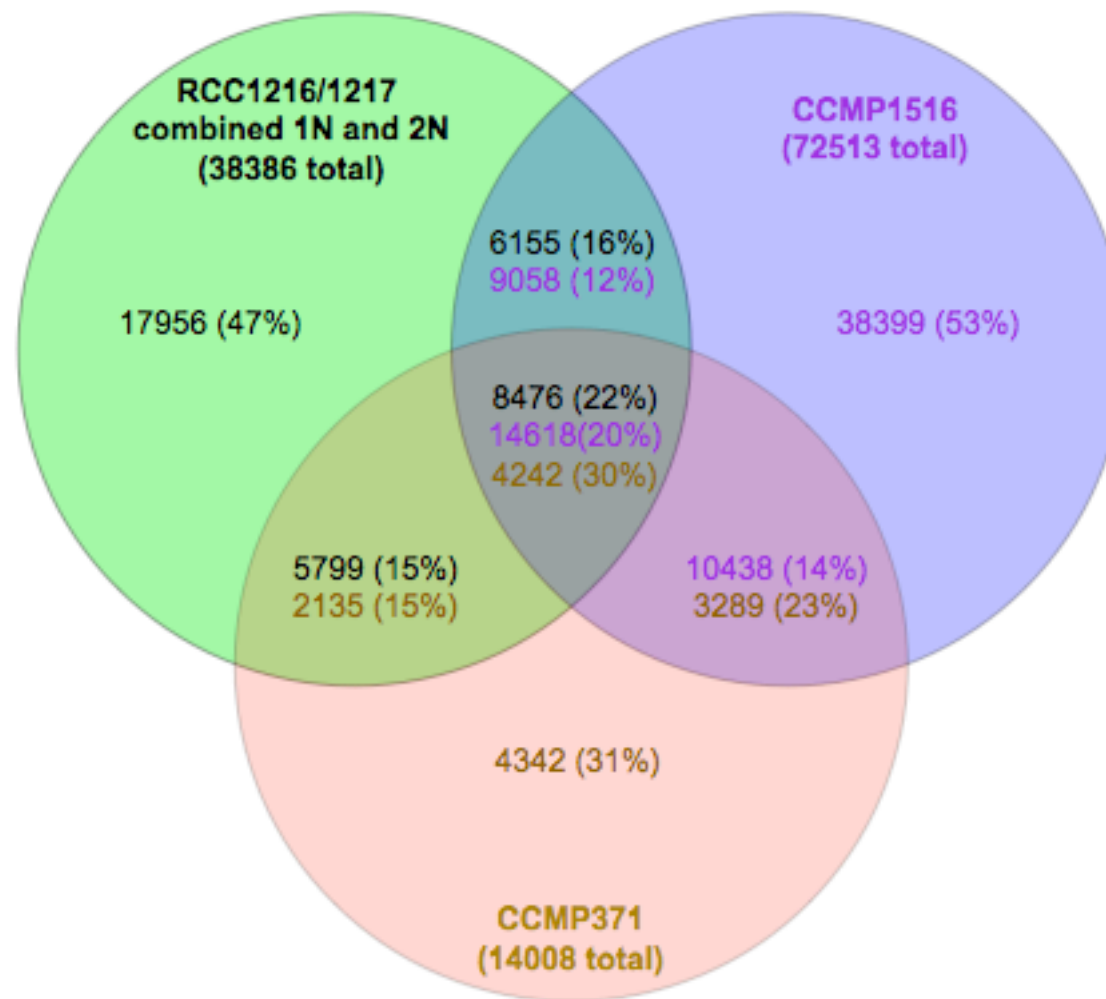
Gene A
*454 reads*
*Sanger EST reads*
*Real transcript A*
5'————————————3'
**Cluster A**

Gene B
*454 reads*
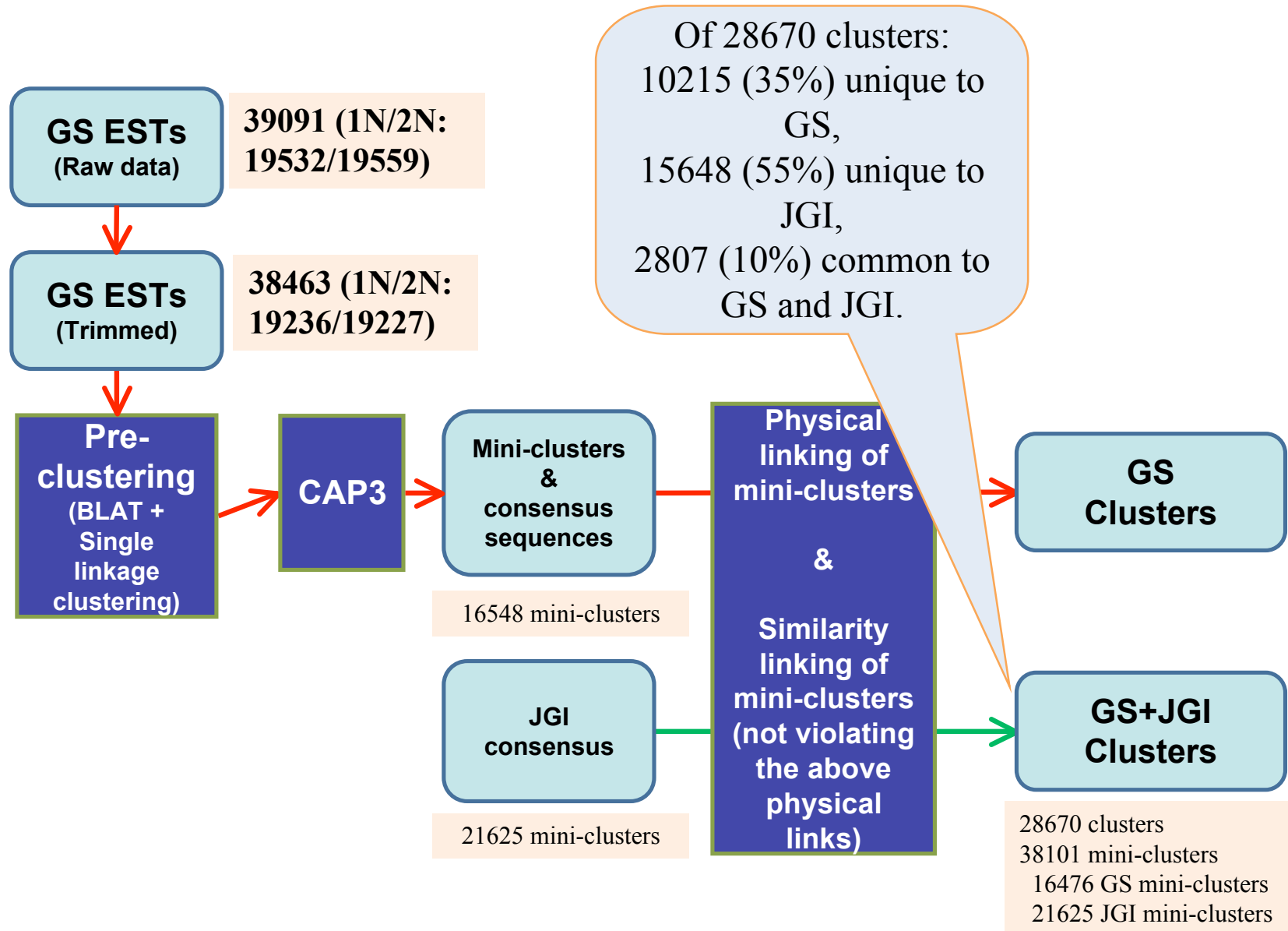*Sanger EST reads*
*Real transcript B*
5'————————————3'
**Cluster B**

# Sanger ESTs reveal a large amount of new transcriptomic information for *Emiliania huxleyi*

# Sanger read mapping/clustering statistics

**GS ESTs** (Raw data)

**39091 (1N/2N: 19532/19559)**

**GS ESTs** (Trimmed)

**38463 (1N/2N: 19236/19227)**

**Pre-clustering (BLAT + Single linkage clustering)**

**CAP3**

**Mini-clusters & consensus sequences**

16548 mini-clusters

**JGI consensus**

21625 mini-clusters

**Physical linking of mini-clusters**

**&**

**Similarity linking of mini-clusters (not violating the above physical links)**

Of 28670 clusters:
10215 (35%) unique to GS,
15648 (55%) unique to JGI,
2807 (10%) common to GS and JGI.

**GS Clusters**

**GS+JGI Clusters**
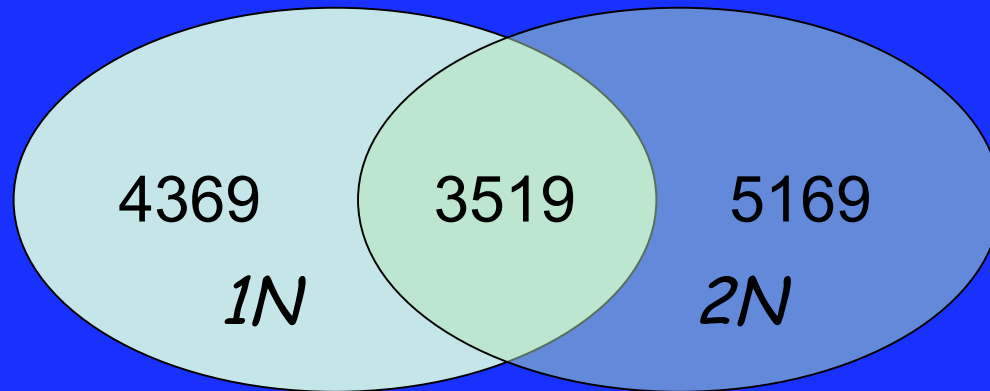
28670 clusters
38101 mini-clusters
  16476 GS mini-clusters
  21625 JGI mini-clusters

# Comparing alternate phases of the life cycle greatly increases transcripts detected

39000 Sanger ESTs:

13057 clusters total

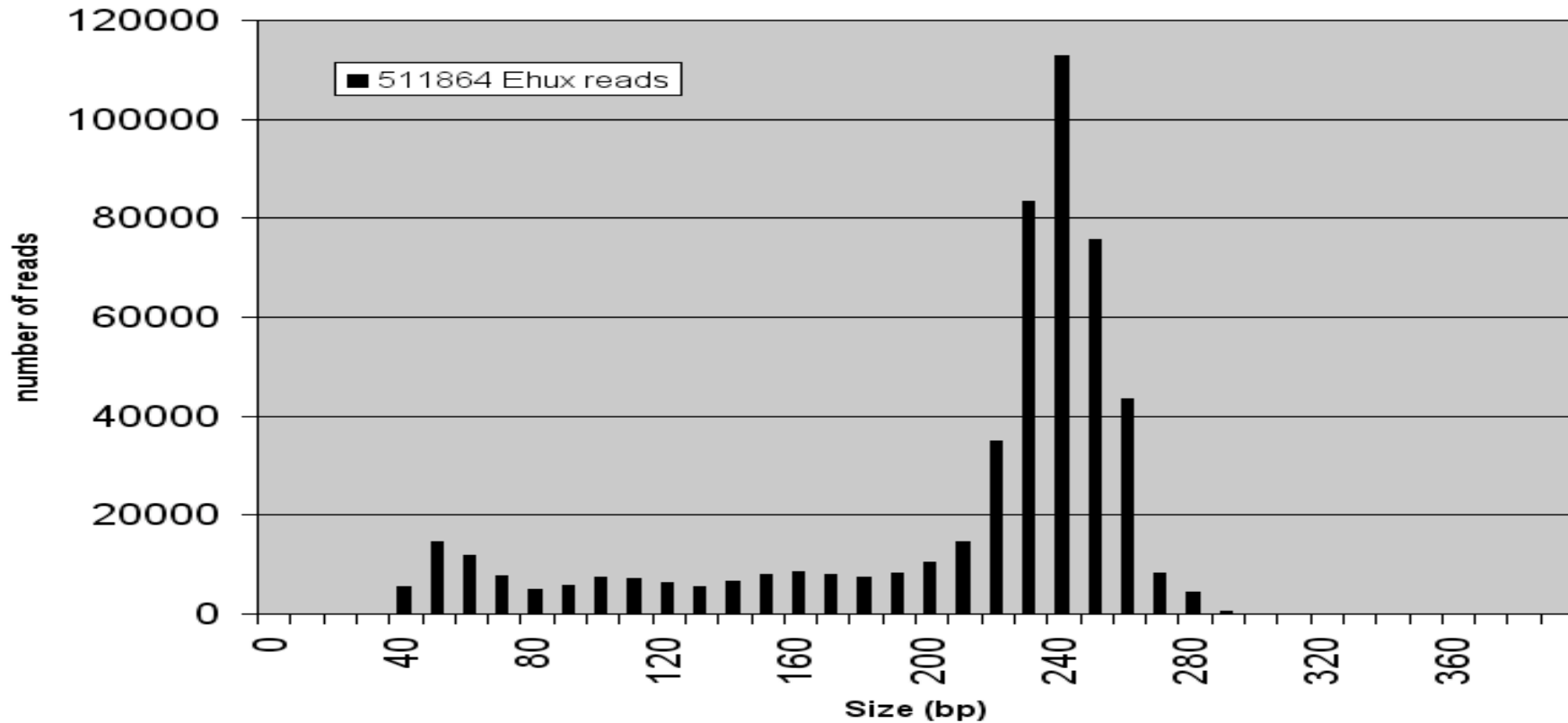4369 | 3519 | 5169

*1N* | | *2N*

| | 1N | 2N |
|---|---|---|
| # clusters | 7888 | 8688 |
| ML estimate of transcriptome richness | 10039 | 11988 |
| Chao1 estimate of transcriptome richness | 12840±214 (12438,13278) | 15931±289 (15385,16522) |
| Coverage | 61.4-78.6% | 54.5-72.5% |

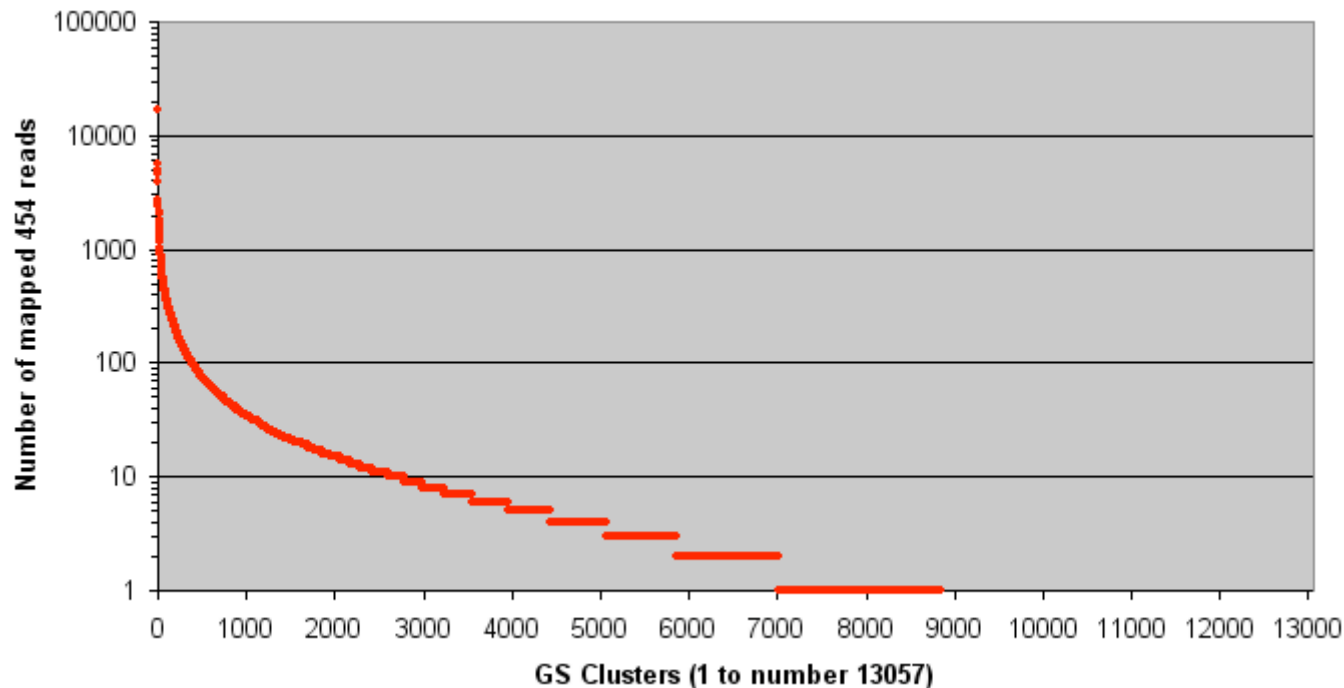*Chao Jaccard-type similarity: ≤50% of expressed genes shared!*

# Ehux 454 read data

- 1N (file: FGFGJ1101)
  - 256484 reads
  - Average size: 216.30 (S.D.: 59.09)
  - G+C: 62.76%
- 2N (file: FG5FMAE01)
  - 255380 reads
  - Average size: 209.84 (S.D.: 62.23)
  - G+C: 62.55%

# Mapping of reads on Genoscope ESTs

- (Simple criteria) BLAT, AL>=90nt, Identity>=95%, Coverage of read by alignment>=70%
- Of 511864 reads, 262023 (51%) were mapped on the Ehux mini-clusters derived from Genoscope EST data sets
- Corresponding ESTs
  - 10611 mini-clusters out of 16471 total mini-clusters (64%)
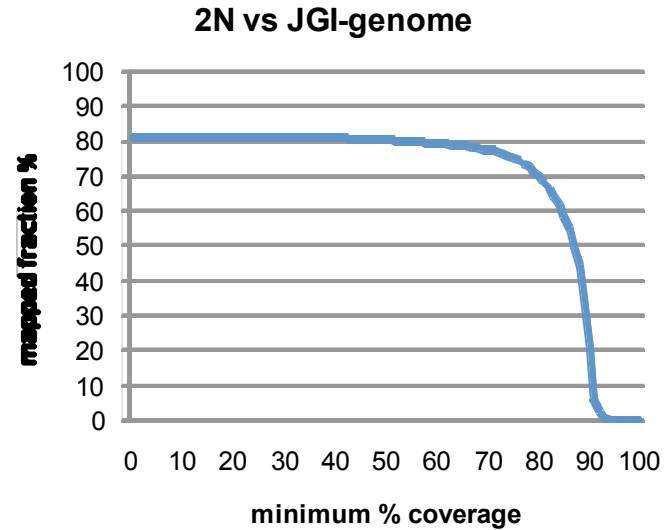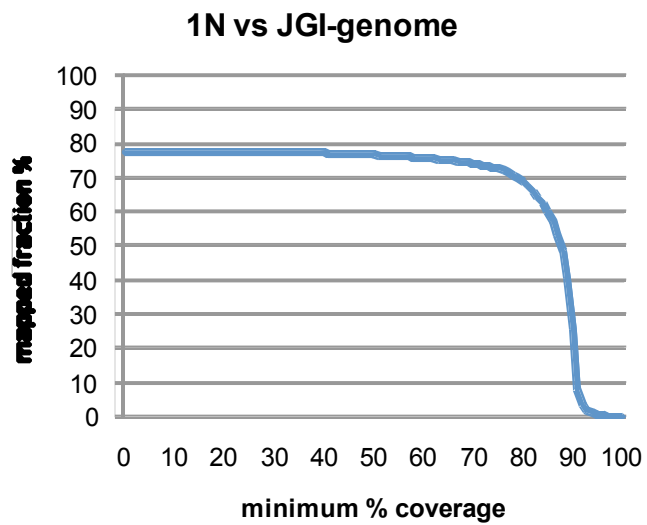  - 8844 clusters out of 13057 total clusters (68%)

Ehux 454 reads mapping statistics
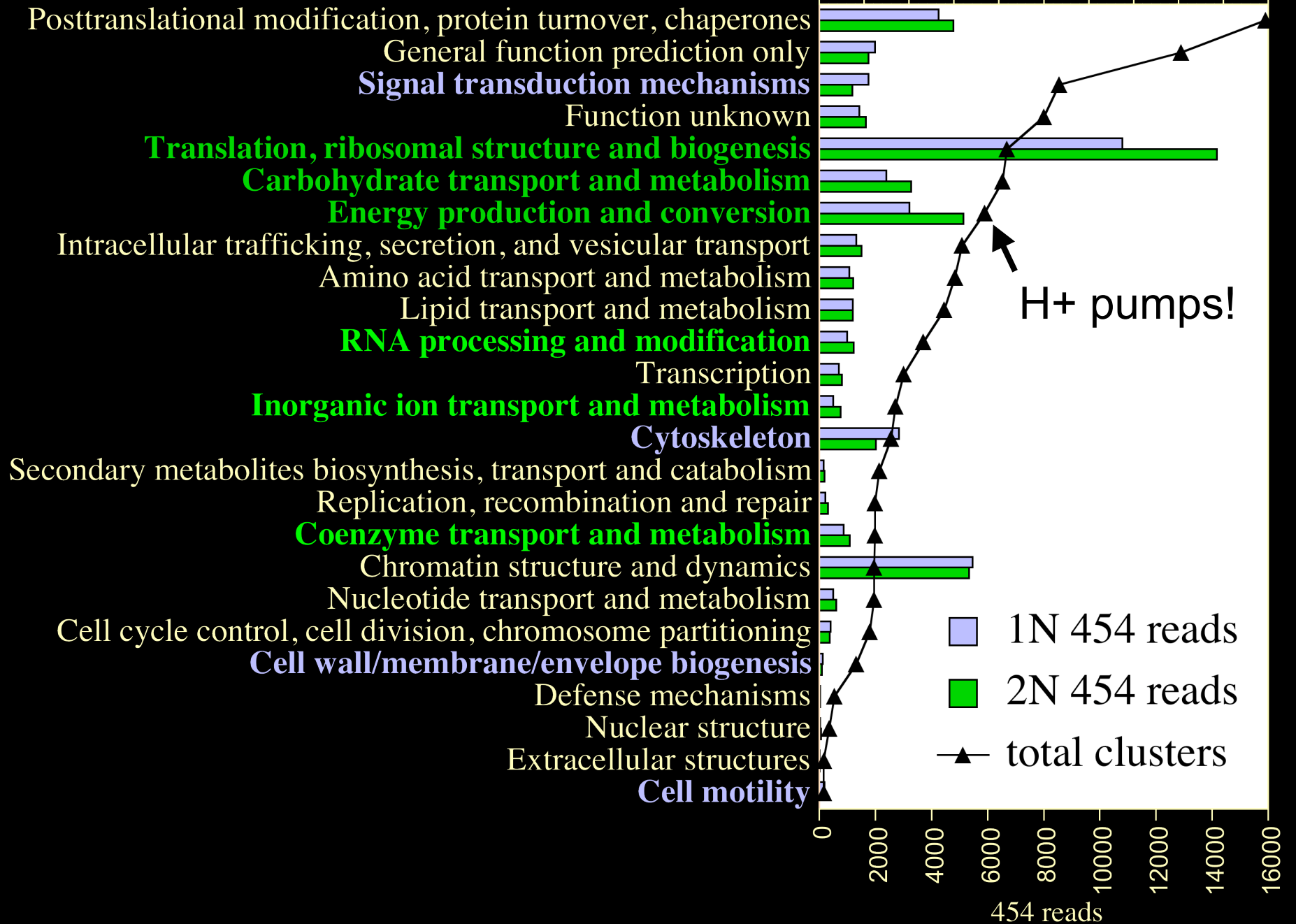BLAT + alignment length≥90nt
by minimum %-coverage of read

1N                                    2N

Against JGI genomic sequences

1N vs JGI-genome

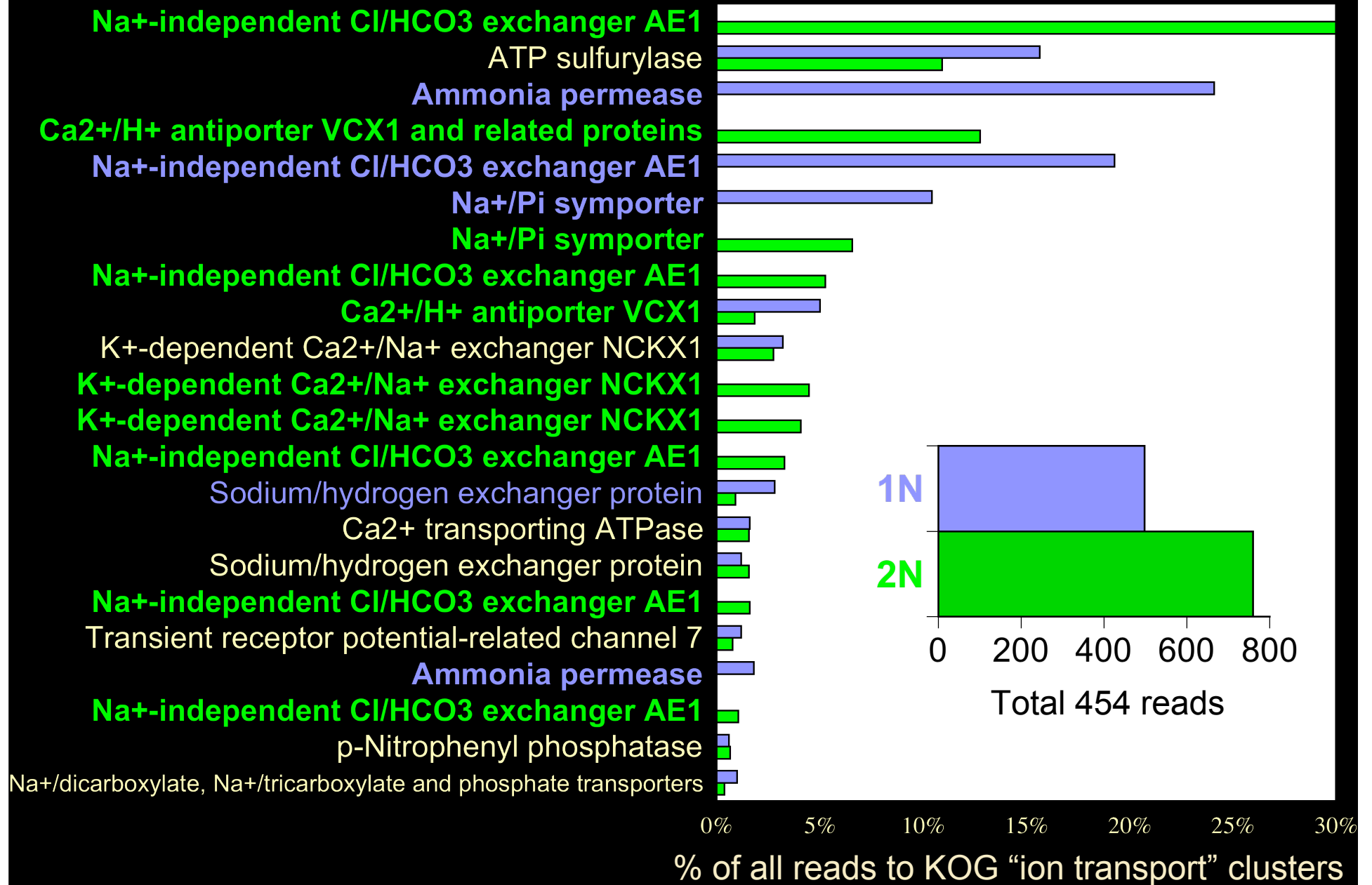2N vs JGI-genome
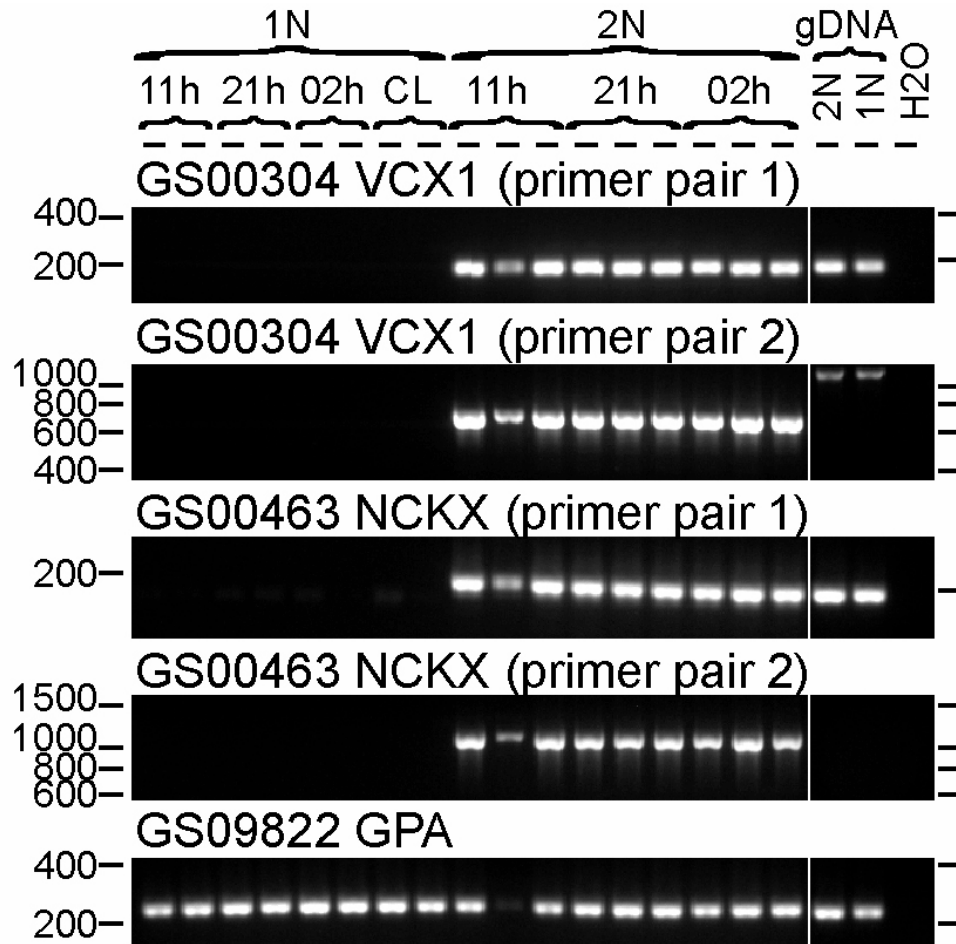
Calcifying 2N cells pump ions more

# RT-PCR tests of potential calcification genes



Confirm highly 2N-specific expression of a VCX1 and NCKX gene

Surprise!
GPA is not 2N-specific!
9 454 reads from 1N
Only 1 454 read from 2N

# Motile 1N cells: Flagella and sensor systems

154 flagellar-related clusters

Expressed <u>only</u> in 1N cells:

86 flagellar or basal body structural elements
with no known cytoplasmic role

13 flagellar dynein heavy chains

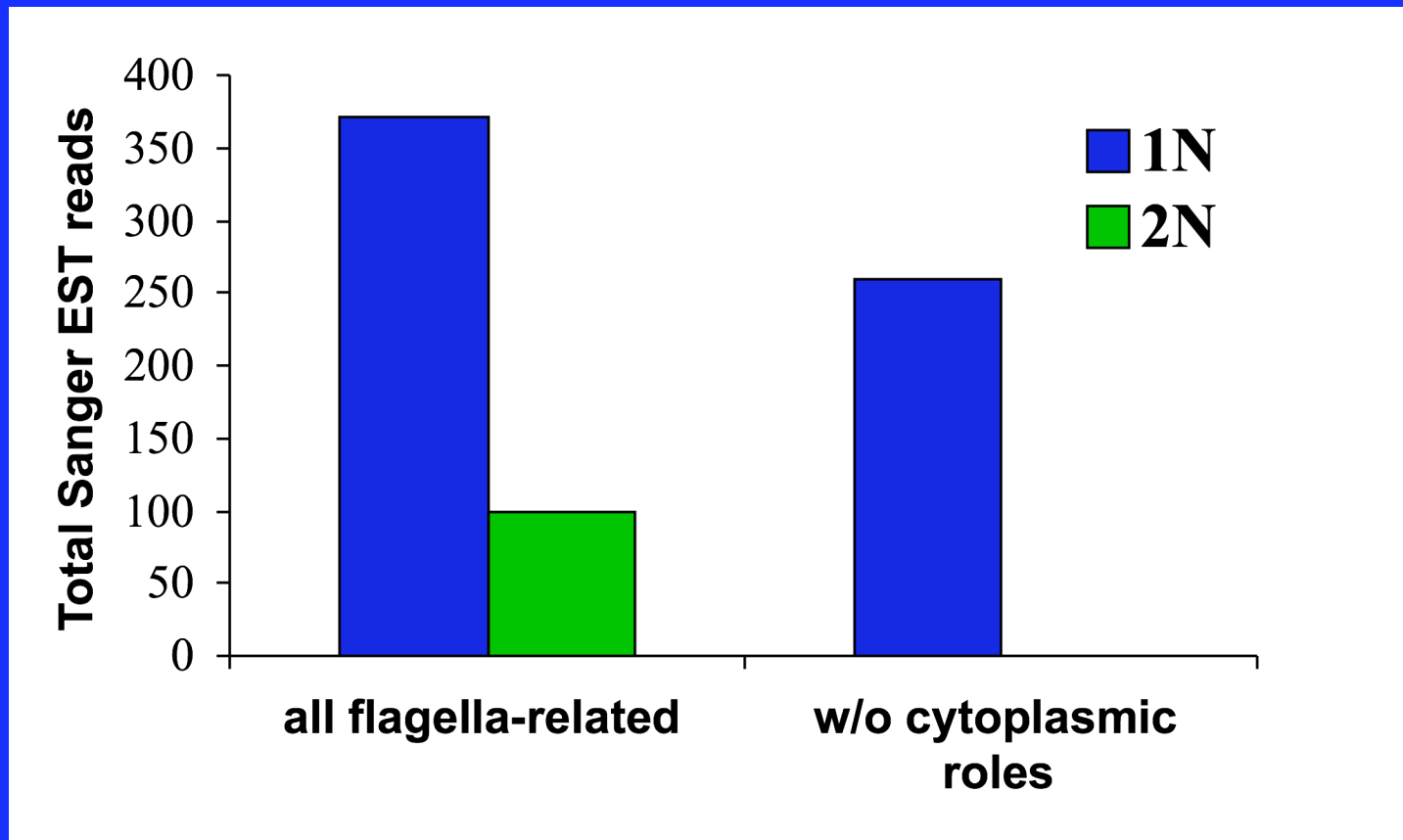1 cytoplasmic dynein heavy chains



Sensors:
Two 1N-specific phototropin-like LOV2 proteins
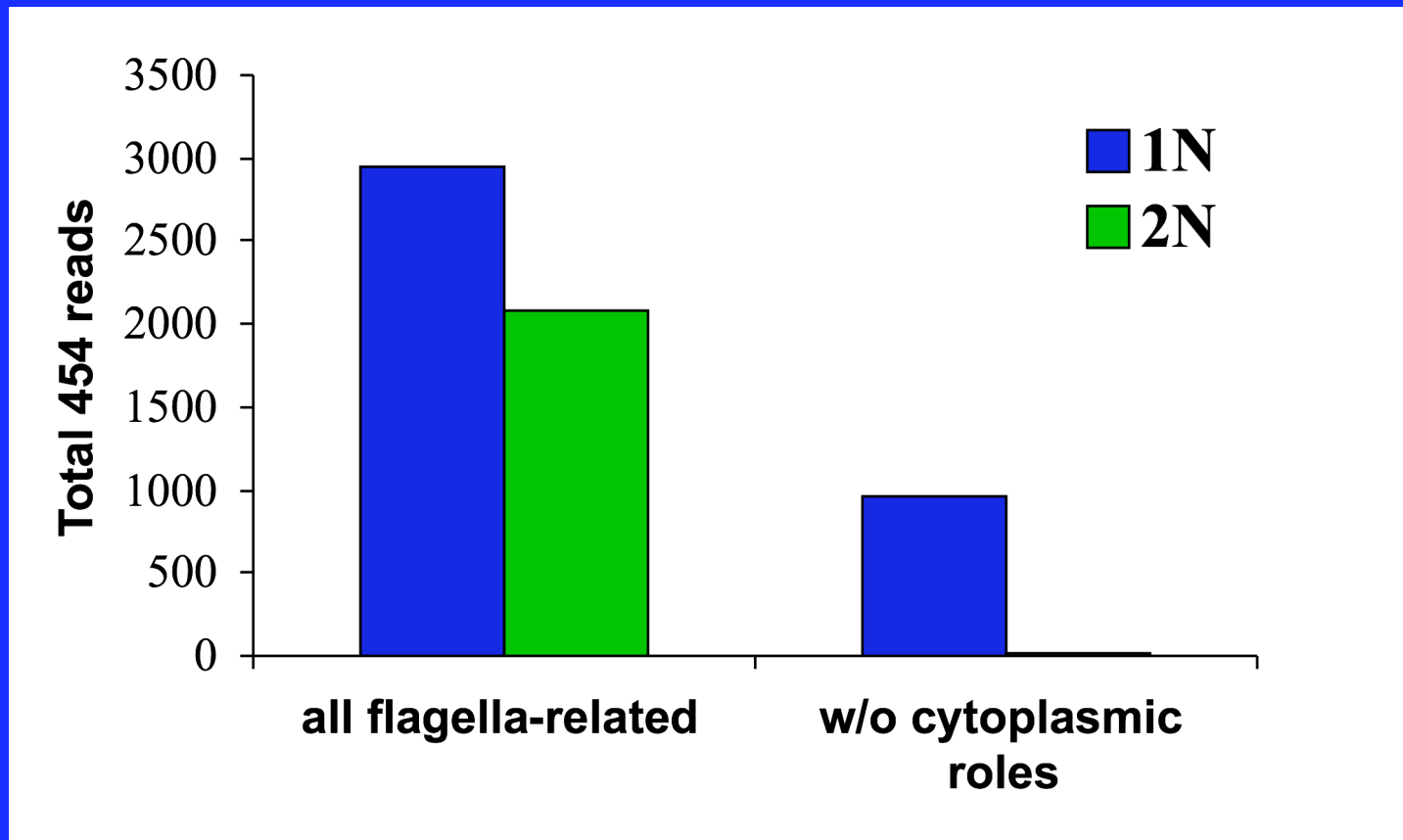1N-specific cGMP protein kinase homolog

# RCC1216/1217 has a lot of genes not found in the JGI whole genome sequence of CCMP 1516!

| | Combined (1N+2N) | Highly 1N-specific (>10x difference $p < 10^{-5}$) | Highly 2N-specific (>10x difference $p < 10^{-5}$) |
|---|---|---|---|
| Sanger clusters | 22.6% | 33.3% | 23.9% |
| 454 reads | 14.9% | 39.0% | 20.6% |

**Flagellar genes lost!**

| | RCC1216/1217 | CCMP1516 |
|---|---|---|
| Calcification | Well calcified | Poorly/non-calcified |
| Life cycle | Forms motile 1N cells | Does NOT form motile 1N cells |
| G1 DNA content | ≈4x Isochrysis | ≈2.9x Isochrysis |

# *Genomic variation between Ehux strains*

22.6% of Sanger-sequenced clusters do not map to JGI genome assembly!

47% of flagellar-associated genes do not map

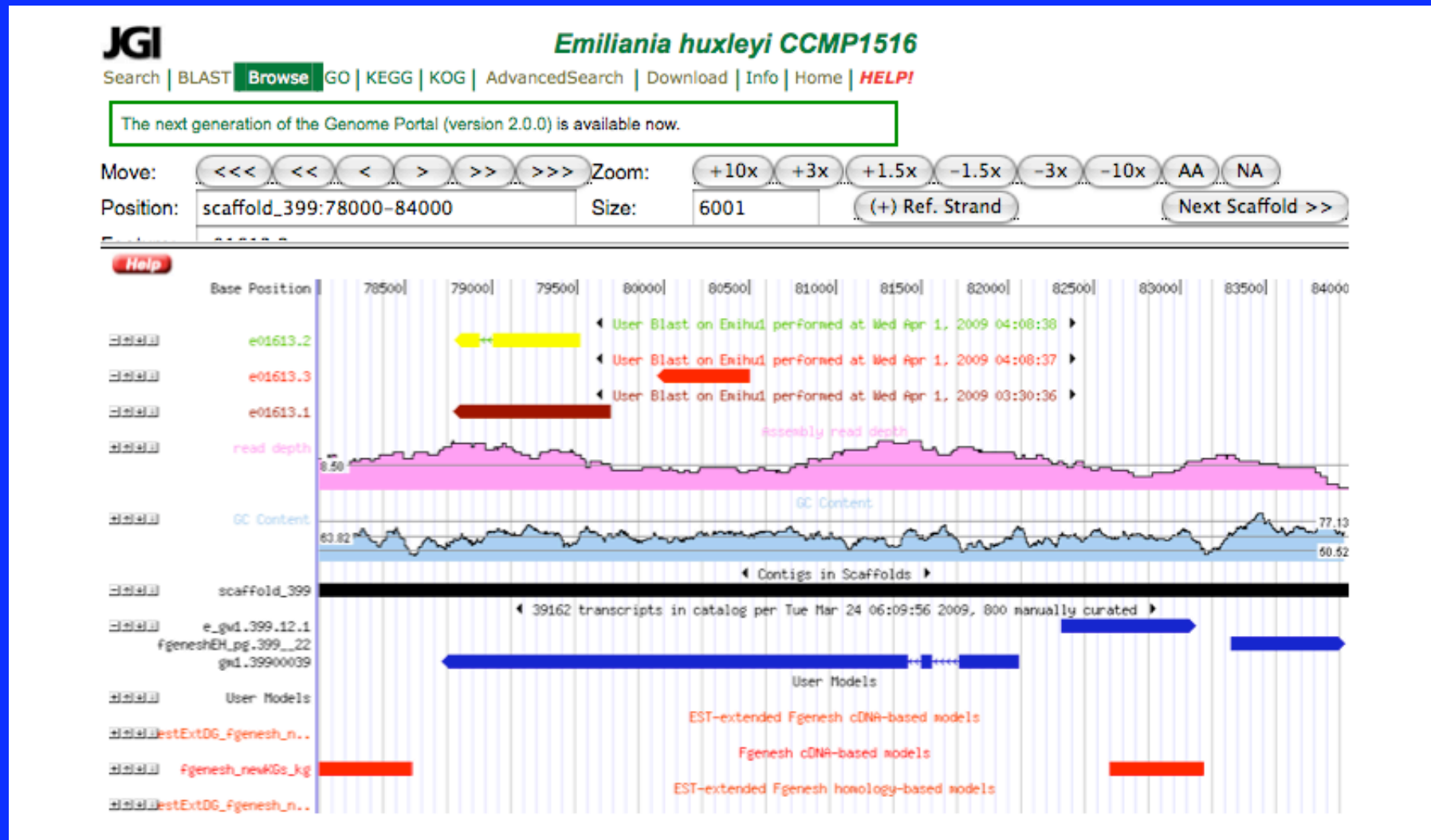5 out of 13 distinct dynein heavy chain genes do not map

Only three loci in JGI/CCMP 1516 assembly have sufficient space to encode the ≈4000 aa dynein heavy chain genes

Several regions of fragmented homology
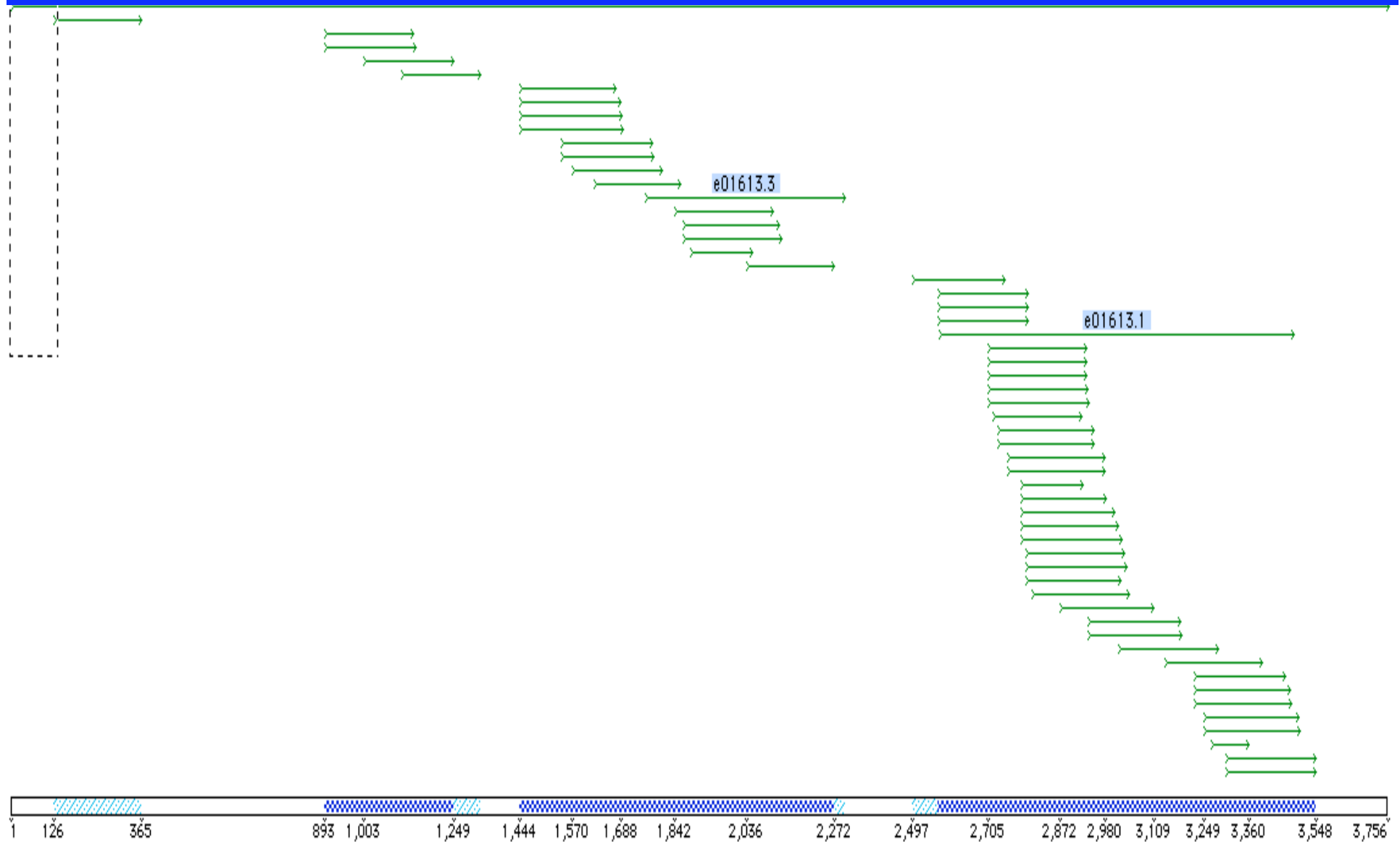
CCMP 1516 never observed to form flagellated cells

PCR suggests other strains may lose flagellar dyneins too!

# The CCMP1516 genome assembly does not have room to encode 4000 amino acid dynein heavy chains?



# Can we use 454 reads to extend transcript models when part of the gene is missing from the assembly??

454 read distribution across dynein heavy chain genes
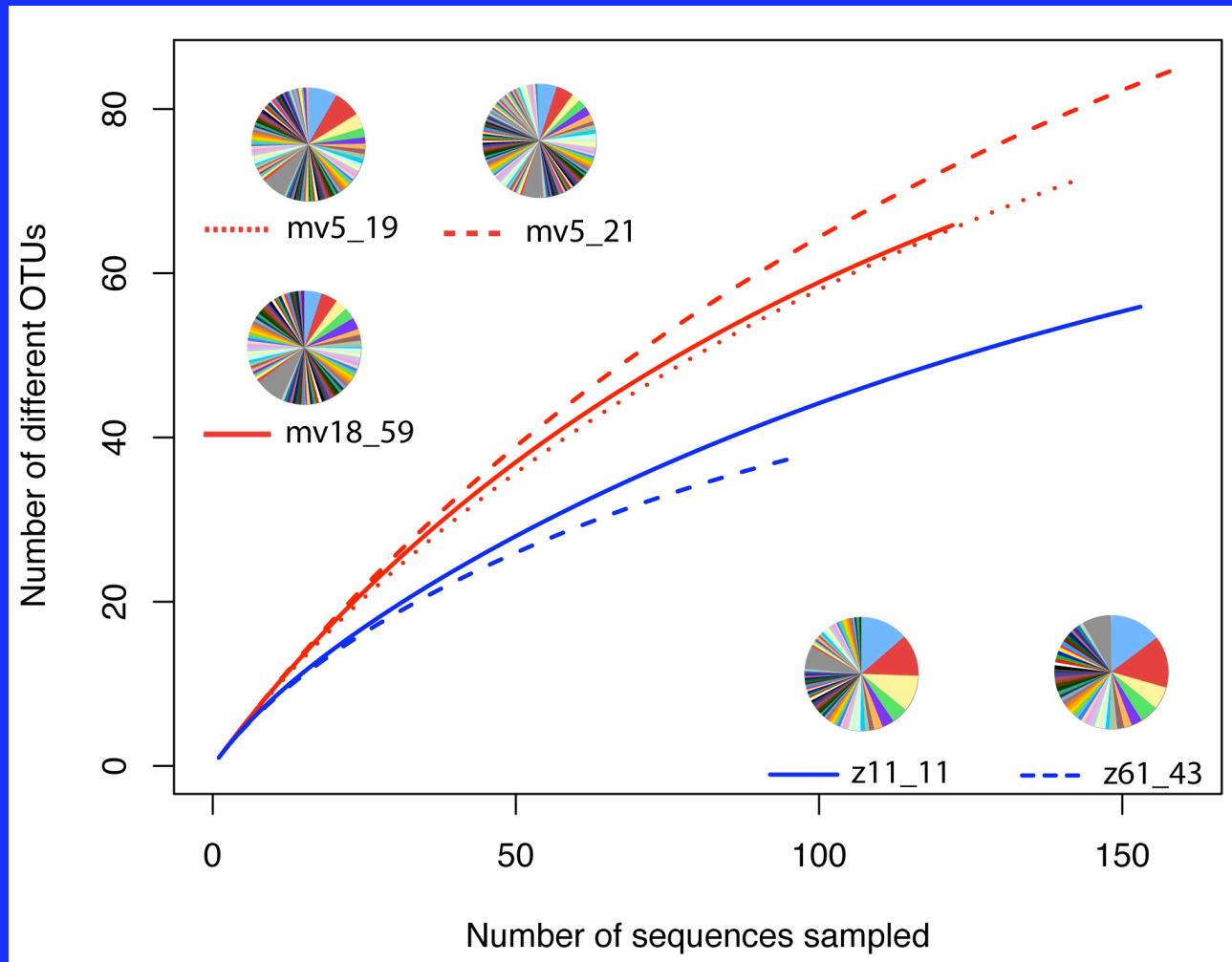
Summary of results so far:

1. Combination of 454 and Sanger technology <u>and</u> proper biology allows deep comparison of transcriptomes

2. We discovered the JGI genome assembly is selectively missing haploid genes!

3. We have to rely much less than planned on existing genome assembly

4. Hybrid 454-Sanger transcript models might help to retrieve missing information
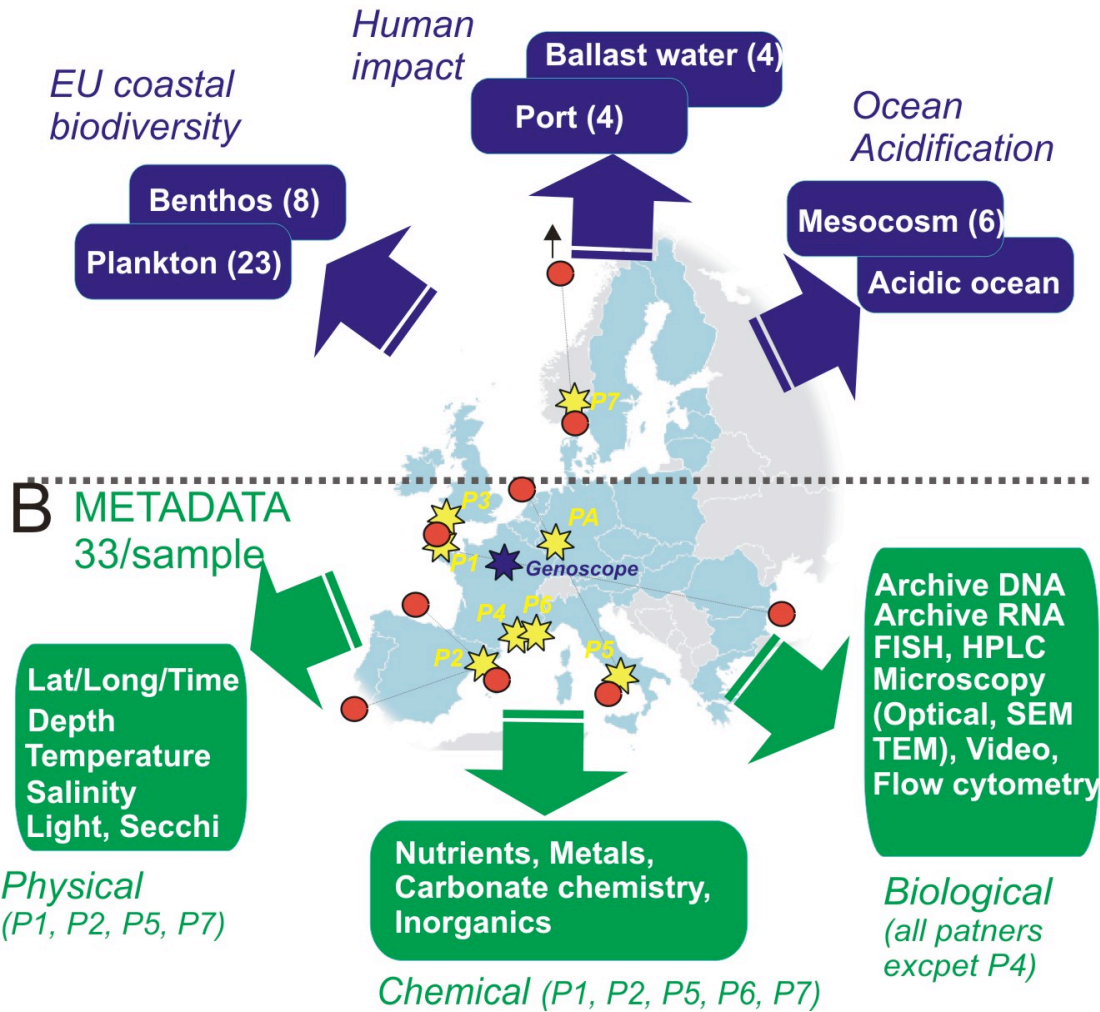
What we still need to do and to figure out:

1. Microarray verification (new arrays ready)

2. Finish mapping 454 reads to JGI genome

3. Create hybrid 454-Sanger clusters with and without using JGI's genome assembly

4. Global statistical description of transcriptomes and their differences based on 454:

   i. Chao1, ML, and Shannon-diversity estimates of transcriptome complexity

   ii. Chao Jaccard-type estimates of transcriptome differences

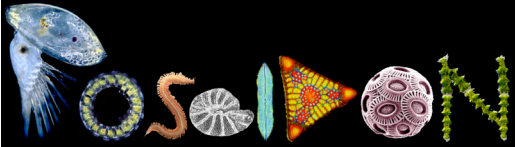   iii. Account for gene or cluster length in statistical analyses of 454 data

POSEIDON

TARA
OCEANS

OCEAN GLACIAL ARCTIQUE

PETRO-PAVLOVSK

LORIENT

VLADIVOSTOK

TOKYO

ANCHORAGE

SEATTLE

TUNIS

ALEXANDRIE

TAIPEI

OCEAN
PACIFIQUE
NORD

OCEAN
ATLANTIQUE
NORD

ABU DHABI

BOMBAY

HONG KONG

MANILLE

DJIBOUTI

MALE

JAKARTA

PORT MORESBY

ILES GALAPAGOS

ILES MARQUISES

ILE ST-HELENE

ARCHIPEL
DES CHAGOS

MADAGASCAR

RIO DE JANEIRO

SYDNEY

AUCKLAND

VALPARAISO

BUENOS AIR

LE CAP

OCEAN
PACIFIQUE
SUD

PUERTO MONTT

OCEAN
ATLANTIQUE
SUD

OCEAN
INDIEN

PORT STANLEY

PUERTO WILLIAMS

TARA OCEANS

Fin!

Merci!