

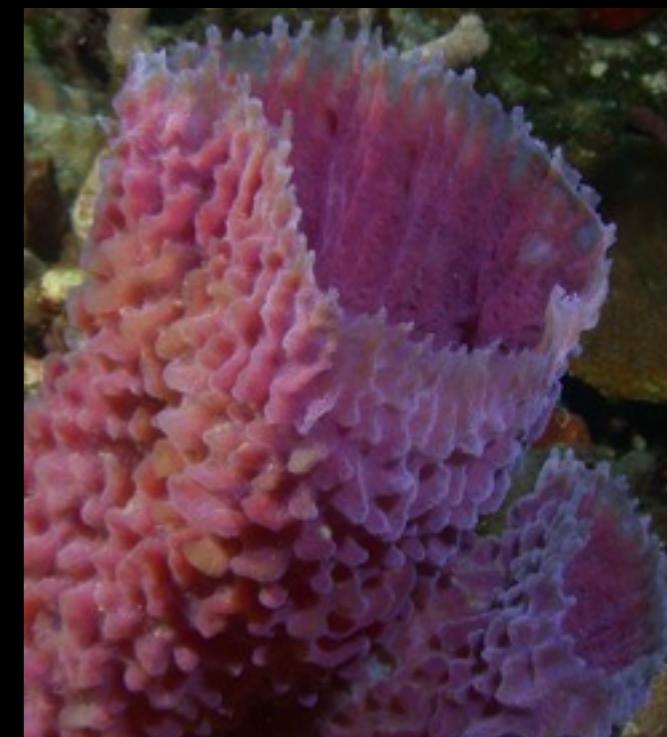
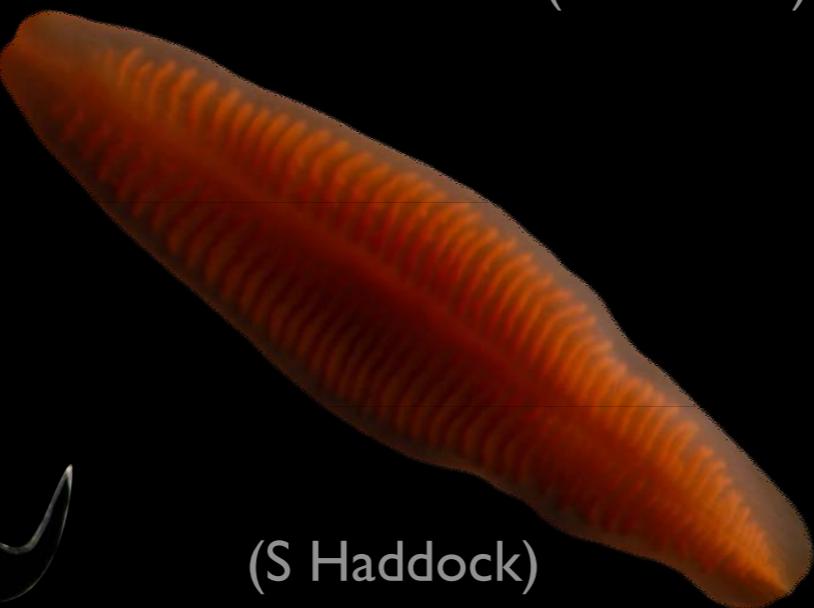
Phylogenomics

Casey Dunn
Assistant Professor
Ecology and Evolutionary Biology



BROWN

Casey Dunn



Casey Dunn



Stefan Siebert
(postdoc)



Stephen Smith
(postdoc)



Rebecca Helm
(grad student)



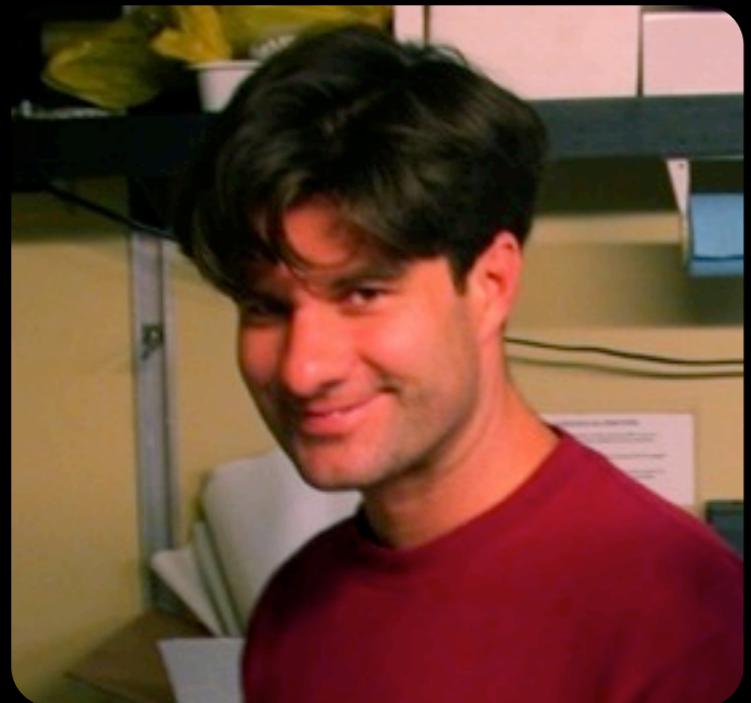
Sophia Tintori
(tech)



Freya Goetz
(tech)

Casey Dunn

Collaborators



Steve Haddock

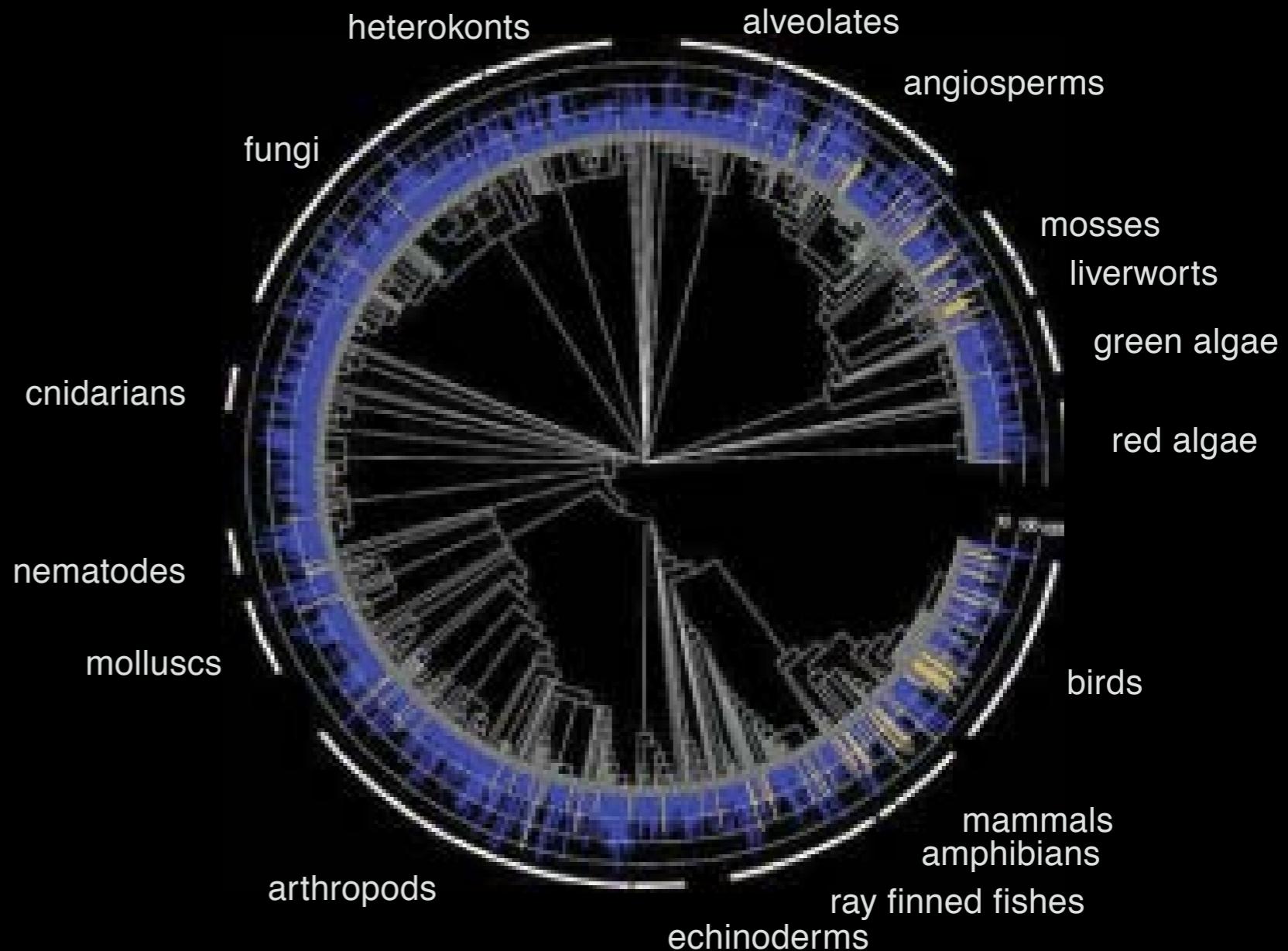
Gonzalo Giribet, A.
Stamatakis, Protostome
AToL team, Cnidarian
AToL team, Mark
Robinson



What does “phylogenomics” mean?

1. The study of genome evolution in a phylogenetic context
2. The inference of species phylogenies with genome data
3. The inference of species phylogenies with data from lots of genes

"...stronger sampling effort aimed at **genomic depth**, in addition to **taxonomic breadth**, will be required to build high-resolution phylogenetic trees at [a broad] scale."



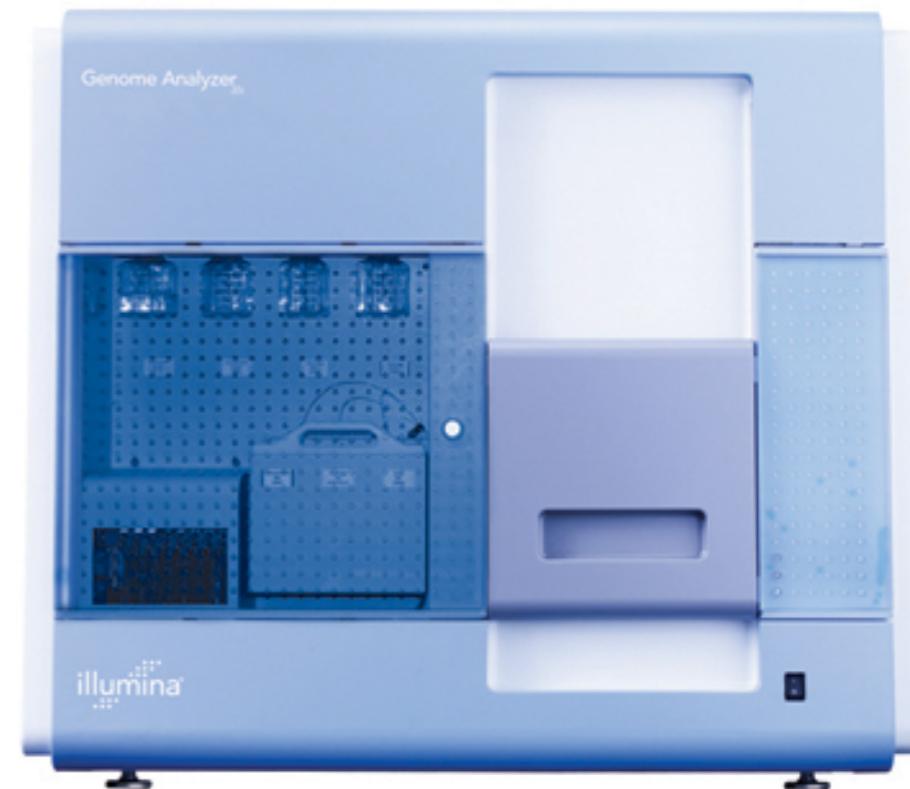
Sanderson, 2008
doi:10.1126/science.1154449

2.6 million sequences
1127 taxa

Why collect data from lots of genes?

- Many hard problems will require lots of data
- Lots of data makes some aspects of inference easier
- These data are useful for things besides building trees
- It can be much cheaper to collect a lot of data than a little bit of data

DNA sequencing



Helicos

Roche

Illumina

Current Illumina costs:

\$2,095 for one lane (HiSeq)

Paired-end 100bp

~150 million clusters

30 gigabases of data

Current Illumina costs:

~\$120 per sample to prepare a library

Current Illumina costs:

Samples per lane	Cost per sample	Clusters per sample (millions)	Gigabases per sample
1	\$2,215	150	30
4	\$644	37.5	7.5
8	\$382	18.75	3.75
12	\$295	12.5	2.5

Current Illumina costs:

~\$100 a gigabase

\$0.000001 per base

Will cheap sequence data
allow us to answer all our
questions?

Of course not.

Should we approach
problems with more data or
improved analysis methods?

This is a false dichotomy.

We need both!

Are other types of data now
obsolete?

No!

These data open entirely
new opportunities for
integrating genomic,
morphological, and
functional perspectives

Marrus claudanielis

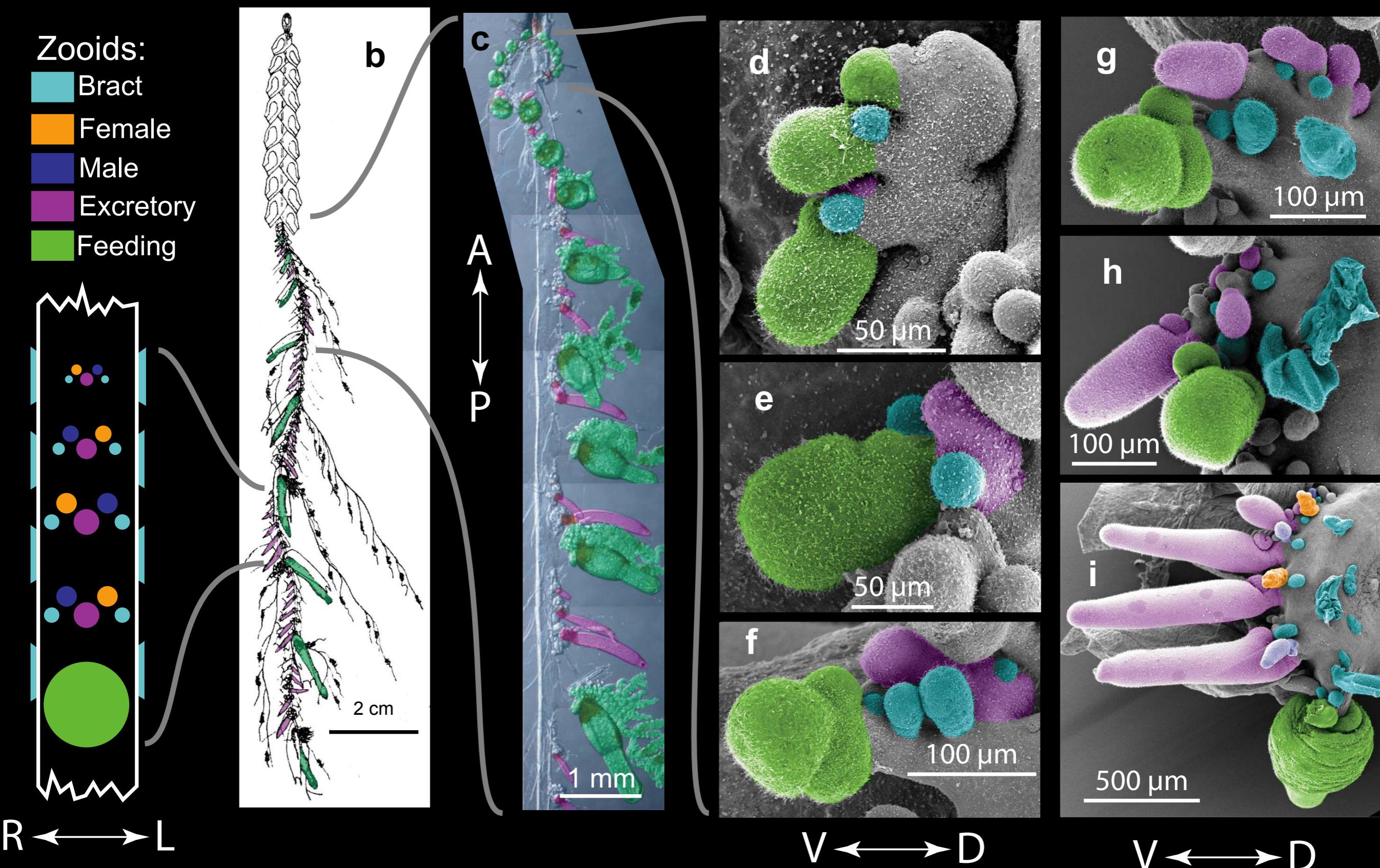


Dunn, Pugh, and Haddock (2005)
Bull. Mar. Sci. 76:699-714

1cm
(MBARI)

Casey Dunn

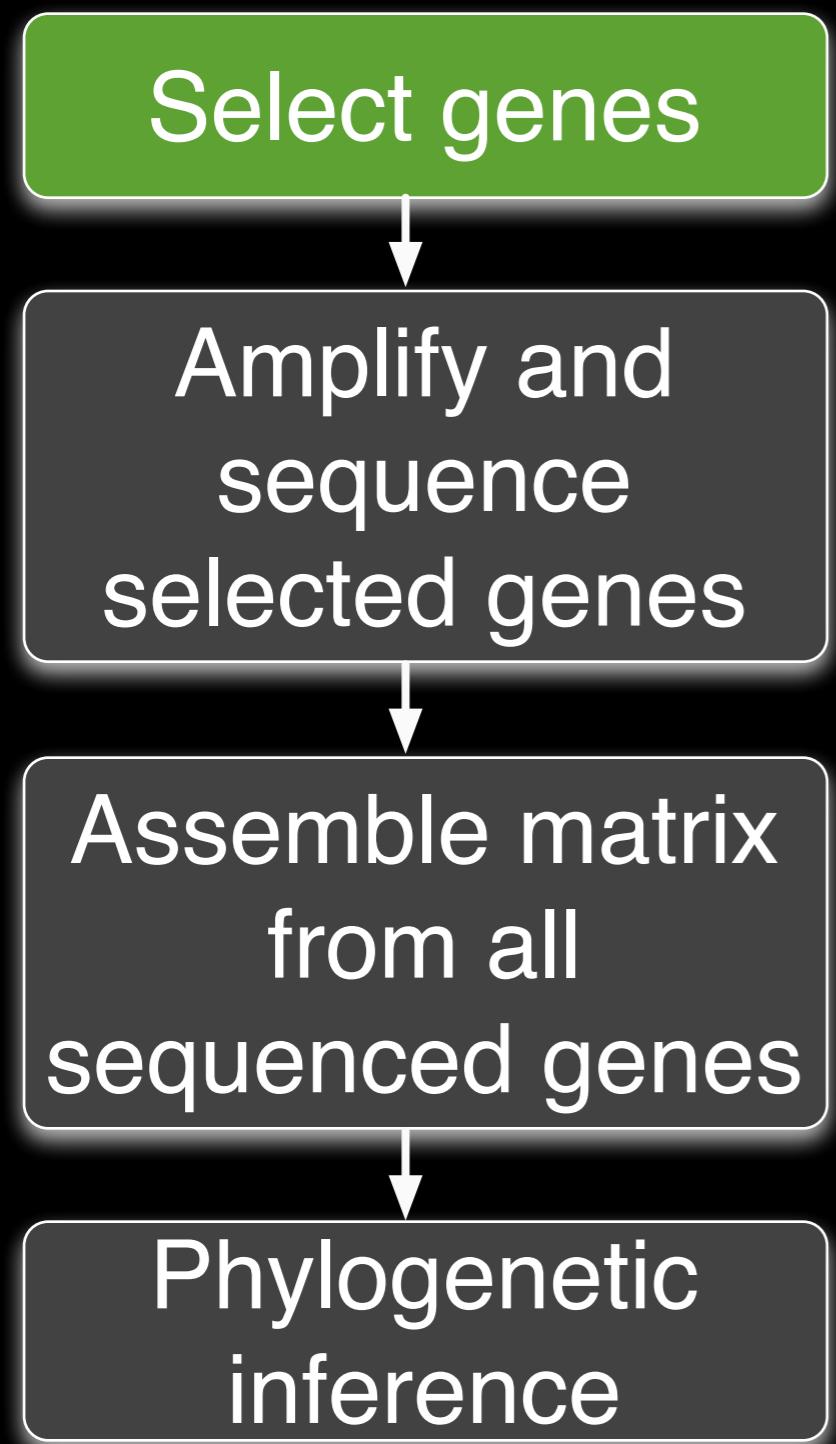
“Physonects”: *Nanomia bijuga*



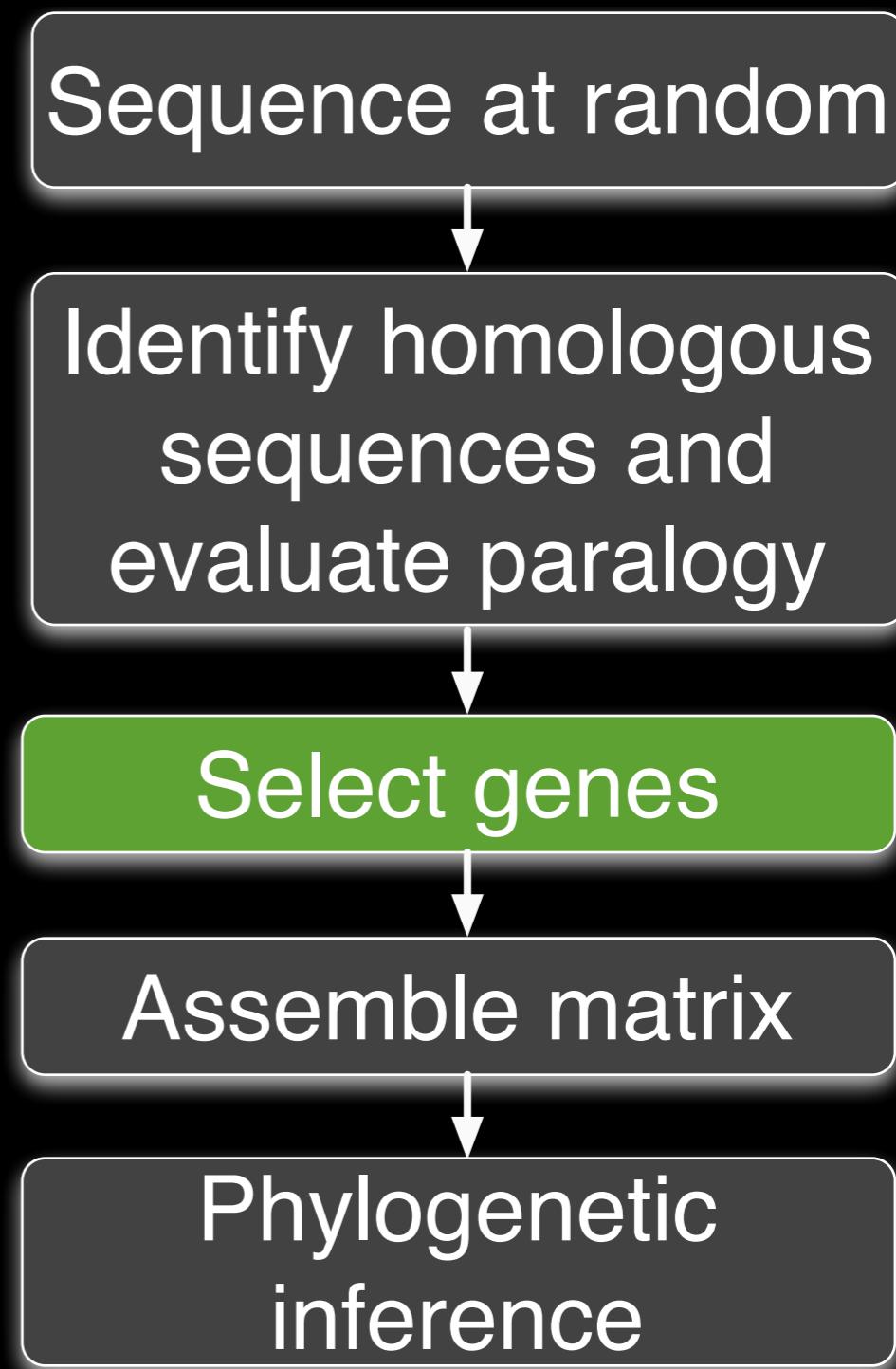
“More Isn’t
Just More—
More Is
Different”

Wired, June 23, 2008

Gene selection as part of project design (Directed PCR)



Gene selection as part of data analysis (ESTs, shotgun genomes)



Getting data...

DNA sequencing



Helicos

Roche

Illumina

Read (sequence data)



Read (sequence data)



Fragment of DNA

DNA Fragments can be:

Amplified/ enriched gene regions

Genomic DNA

cDNA (Transcriptomes)

Genome

Transcriptome

Start with DNA

Get all genes,
regulatory regions, etc

Genomes can be
really big

Can be hard to
identify genes

Start with mRNA

Get a snapshot of
active genes

Almost all data is from
coding regions

Handling RNA is tricky

Overview of sequencing:

Get DNA or RNA

Make a library (chunks of DNA with adapters)

Prepare library for sequencing

Sequence

Process data into raw reads



Casey Dunn

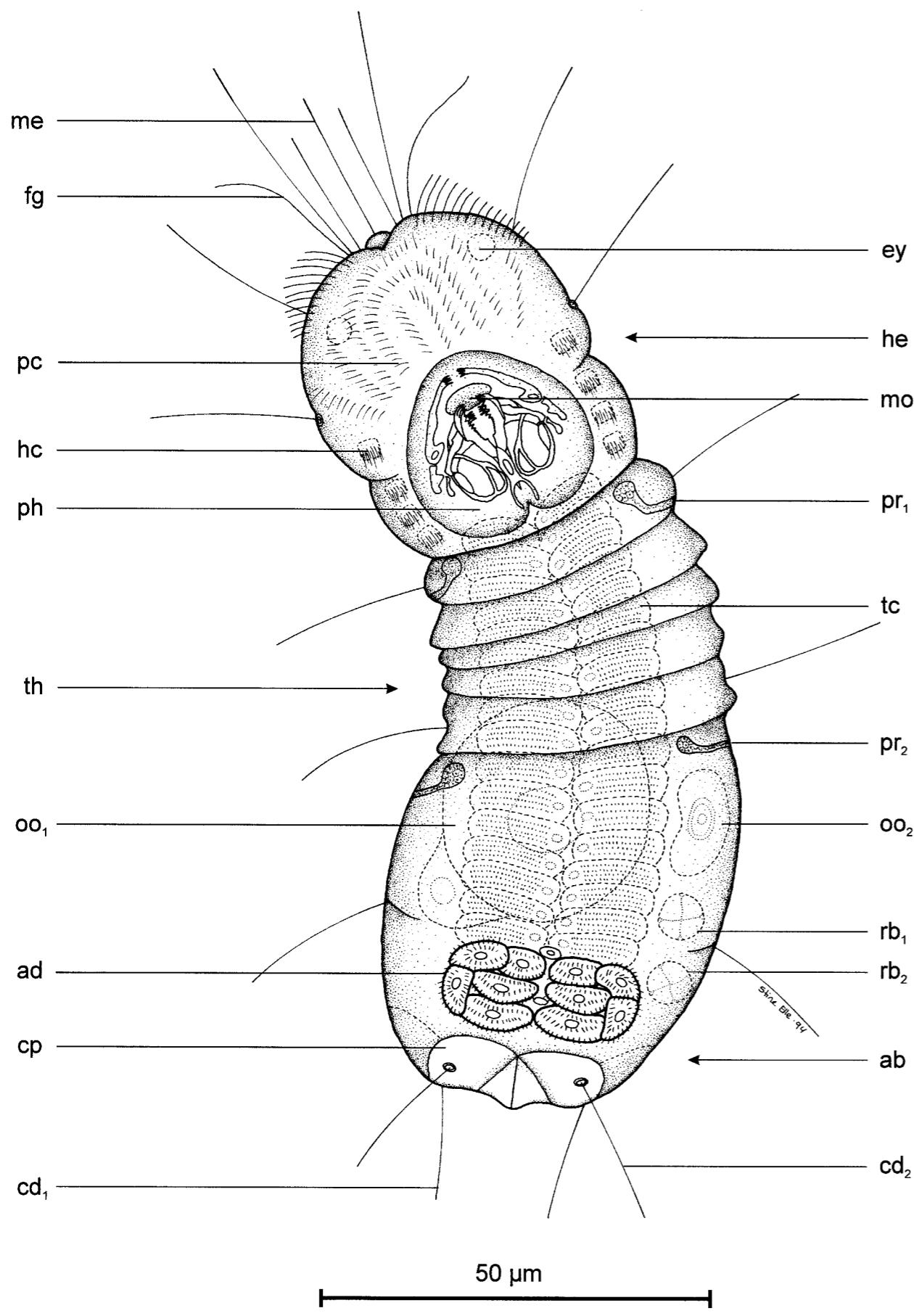
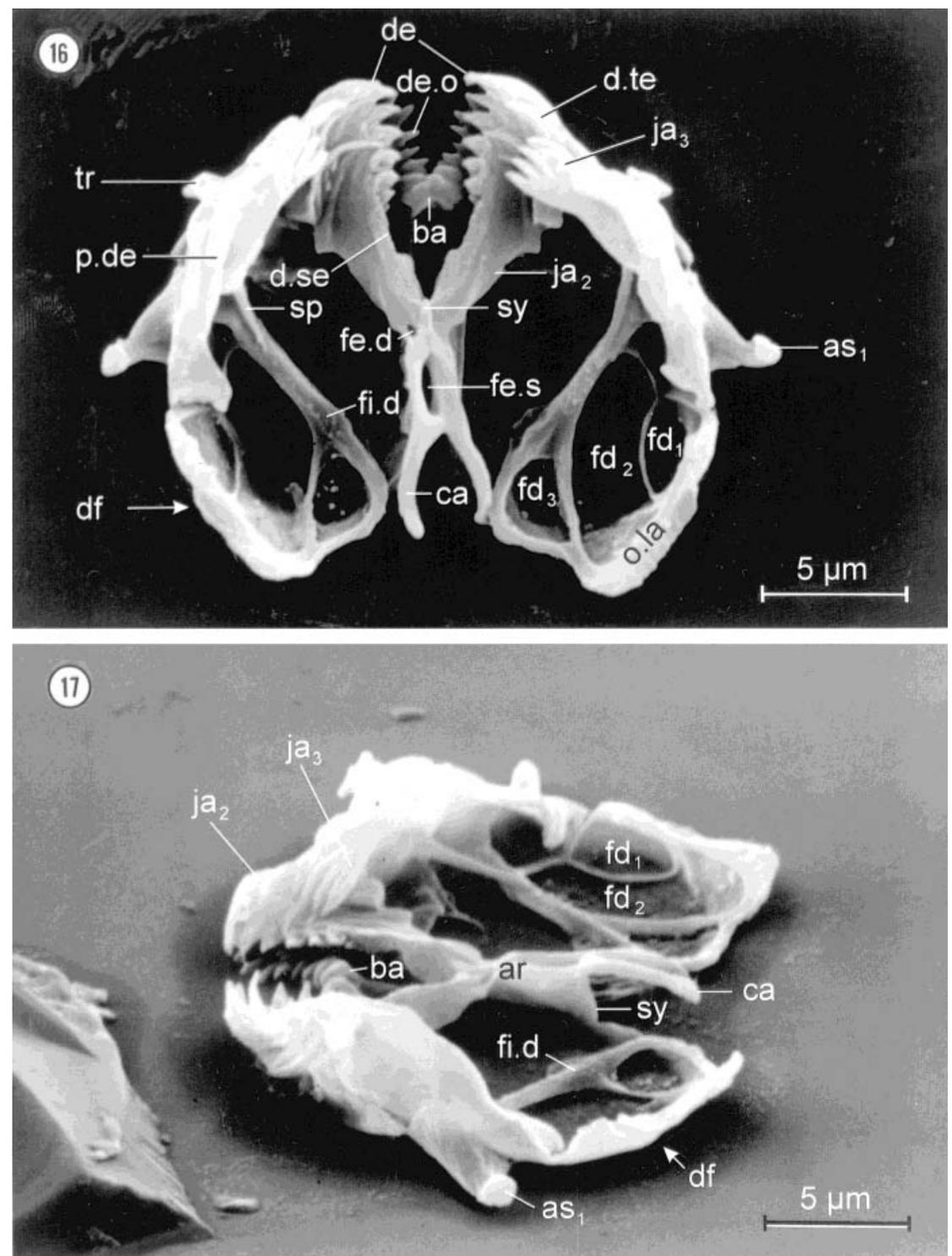


Figure 2



(Kristensen and Funch, 2000)

Casey Dunn



Casey Dunn



Casey Dunn



Casey Dunn

Some options for preservation

Freeze tissue (-80C or colder)

RNALater (Ambion), kept cold

Extract RNA in the field

Homogenize in Trizol, keep cold



Casey Dunn

mRNA isolation - Lots of tissue

Isolate Total RNA with Trizol

Digest DNA

Isolate mRNA (eg NEB S1550S,
Dynabeads mRNA Kit)

mRNA isolation - Small amount of tissue

mRNA straight from tissue (eg NEB S1550S, Dynabeads mRNA DIRECT Kit)

mRNA isolation - Minute amount of tissue

Linear cDNA amplification (eg NuGen Ovation)

RNA quality is (almost) Everything!

Avoid contamination

Reduced sample size requirements
have improved this

RNA quality is (almost) Everything!

Quantity matters - be cautious
working at the bottom range of
sample requirements

RNA quality is (almost) Everything!

Amount of ribosomal RNA matters

There are tradeoffs between rRNA fraction and yield. If material is limiting, purify less and sequence more

RNA quality is (almost) Everything!

If you have enough RNA, look at its size distribution with a BioAnalyzer or similar tool.

What do you do once you have
high quality DNA or RNA?

Break it into little pieces!

Read



DNA

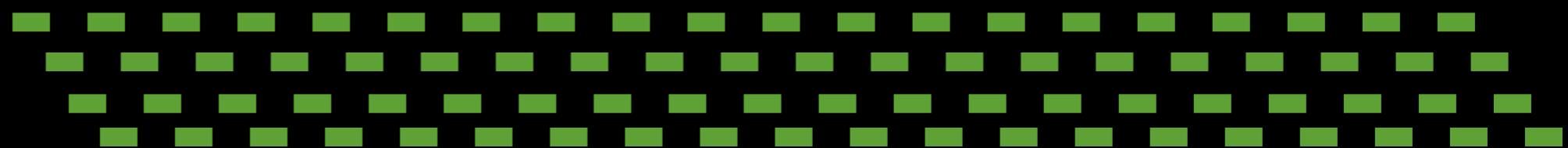
Read

Starting
material

Fragment ↓



Prepare library,
sequence ↓



Library preparation options

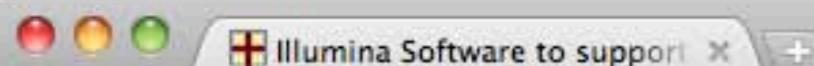
Get a library preparation kit
from the sequencer vendor

Get a third party library
preparation kit

Make the library from scratch

These days, in my lab we use:

- TruSeq RNA kit (Illumina)
- NEBNext (NEB)



Division of Biology and Medicine

Center for Genomics and Proteomics Genomics Core Facility

[About](#)[Contact](#)[Illumina Sequencing](#)[Equipment](#)[Equipment Sign Up](#)[Services](#)[Acknowledgments](#)[DNA Sanger Sequencing](#)[Current Prices](#)

» Biomed Research » Biomed Core Facilities

Overview of Next-Gen Sequencing[HiSeq2000/GAIIX](#)[Listserv/ Discussion Group](#)[Timeline for Implementation](#)[Price Structure](#)[Sample requirements](#)[Covaris 220 Recipes and Guides](#)[Sample Submission Form](#)[Bioinformatics](#)[Illumina Software to support the GAIIX system](#)[Seminars](#)

Welcome!

Thank you for your interest in Next Generation Sequencing at Brown University. The following information is intended to provide users with an introduction on how to get started with plans for a sequencing project.

The Genomics Core Facility, located at 70 Ship Street in the Laboratories for Molecular Medicine, is performing high throughput sequencing using Illumina instruments.

The Genomics Core Facility introduced Next Generation Sequencing Service (NGS) in June 2010 with a GAIIX sequencer. In April 2011 we expanded our successful service and we now offer sequencing on a HiSeq2000 instrument in addition to the GAIIX. For more information on our instruments go to the link "GAIIX/HiSeq2000".

Read about our implementation of the NGS services on the "Timeline for Implementation" link on the left. You can also view our GAIIX install progress on the "Track our Progress" link on the right.

[Track Our Progress](#)

Prices at:

www.brown.edu/Research/CGP/core/illumina/price

Casey Dunn

Data are usually delivered
in fastq or qseq format.

fastq example:

```
@HWI-ST625:51:C02UNACXX:7:1101:1179:1962 1:N:0:TTAGGC
CTAGNTGTTGAAGAGAAGGTTCAAGAACCAAAAGAAAGCTCACAAACACATATGGT
+
=AAA#DFDDDHHFDGHEHIAFHIIIIIGICDGAGDHGGIHG@A@BFIFIHIIIGC@@8

@HWI-ST625:51:C02UNACXX:7:1101:1242:1983 1:N:0:TTAGGC
ATAATTCAATGACTGGAGTAGTGAAAATGAACATAGATATGAGAATAACCGTAGA
+
ACCCFFFFFGHHHHJJJIJEHIFHIJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ
```

These
instruments
generate a lot of
data.

The data files from lane of
an Illumina HiSeq are
more than 70 GB.

A HiSeq has 16 lanes.



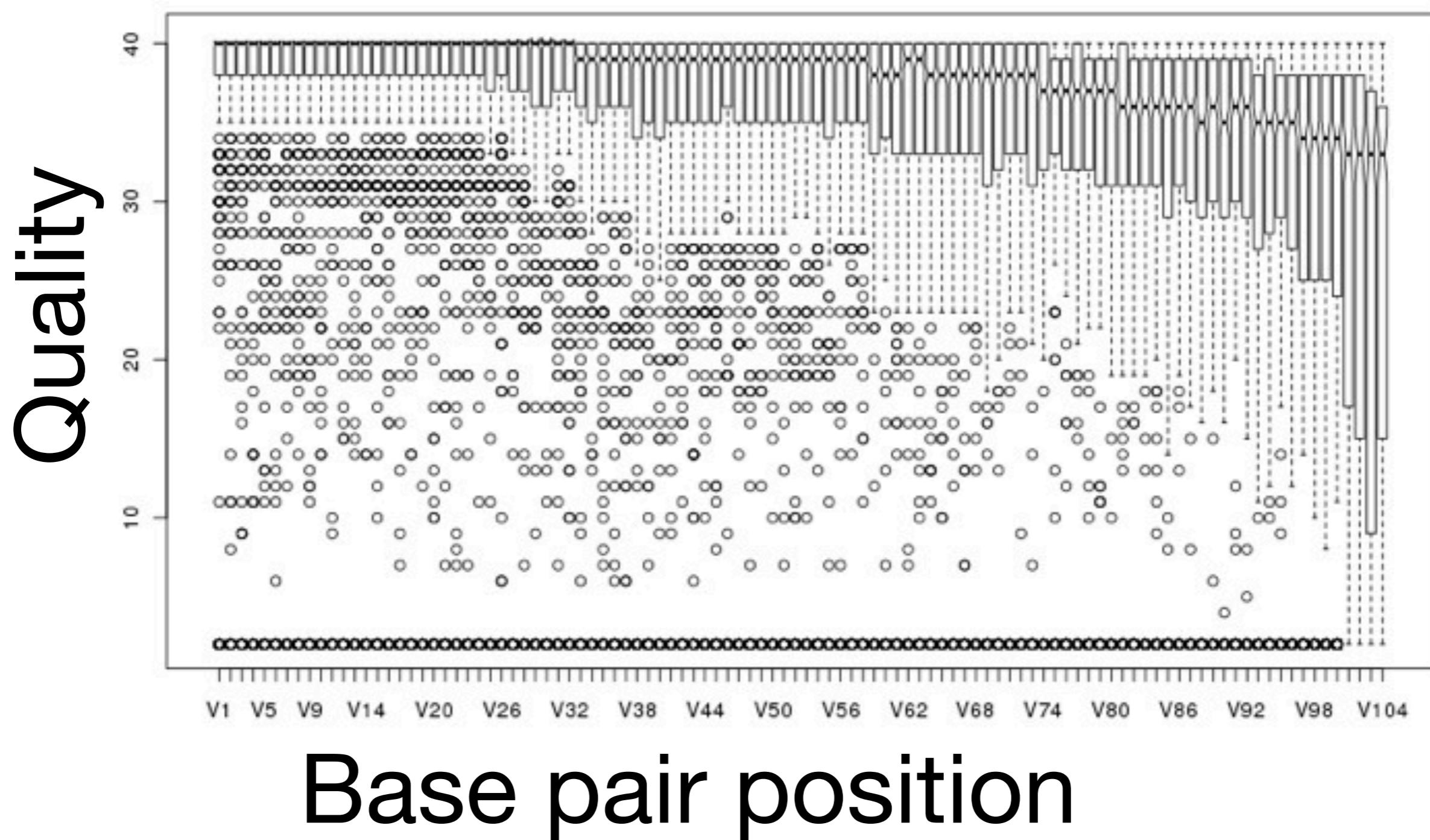
(Sam Fulcomer)

Casey Dunn

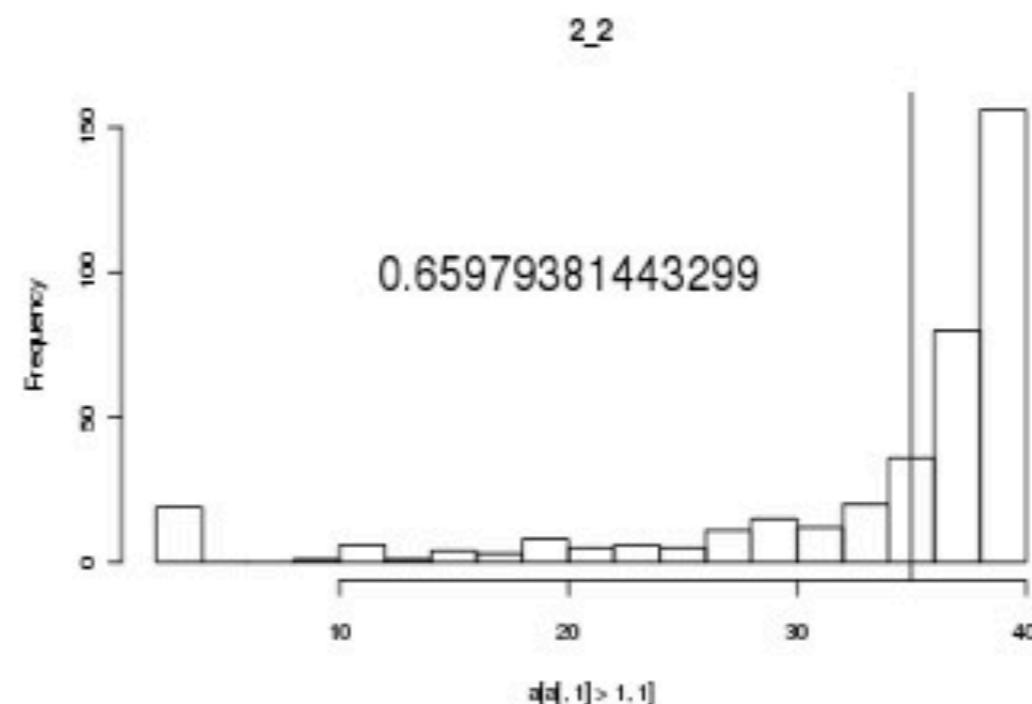
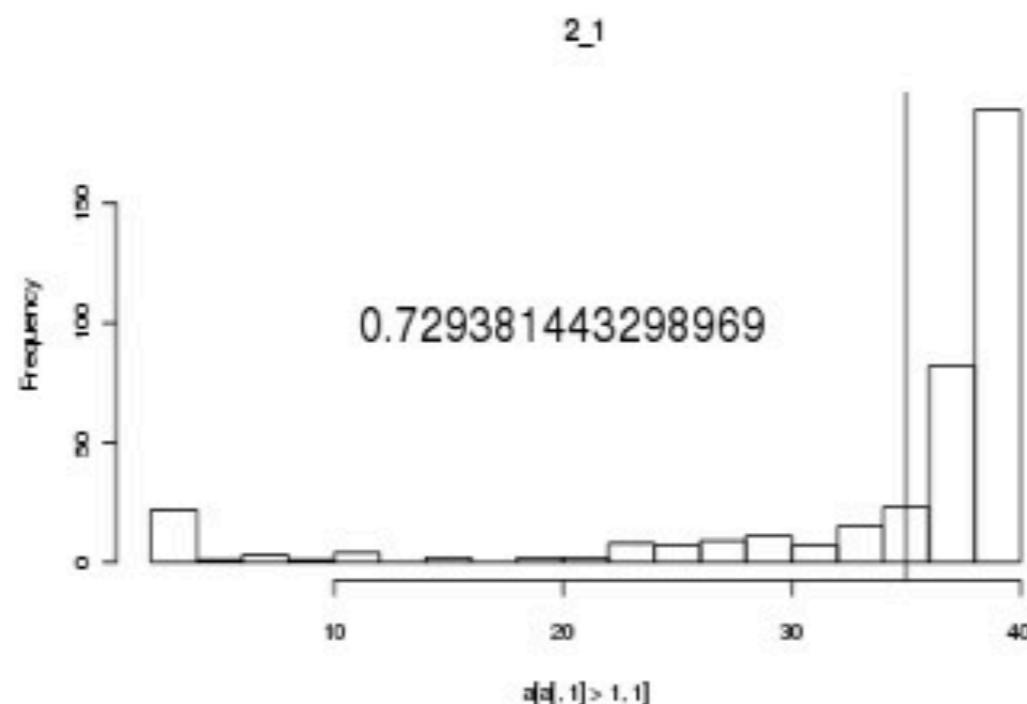
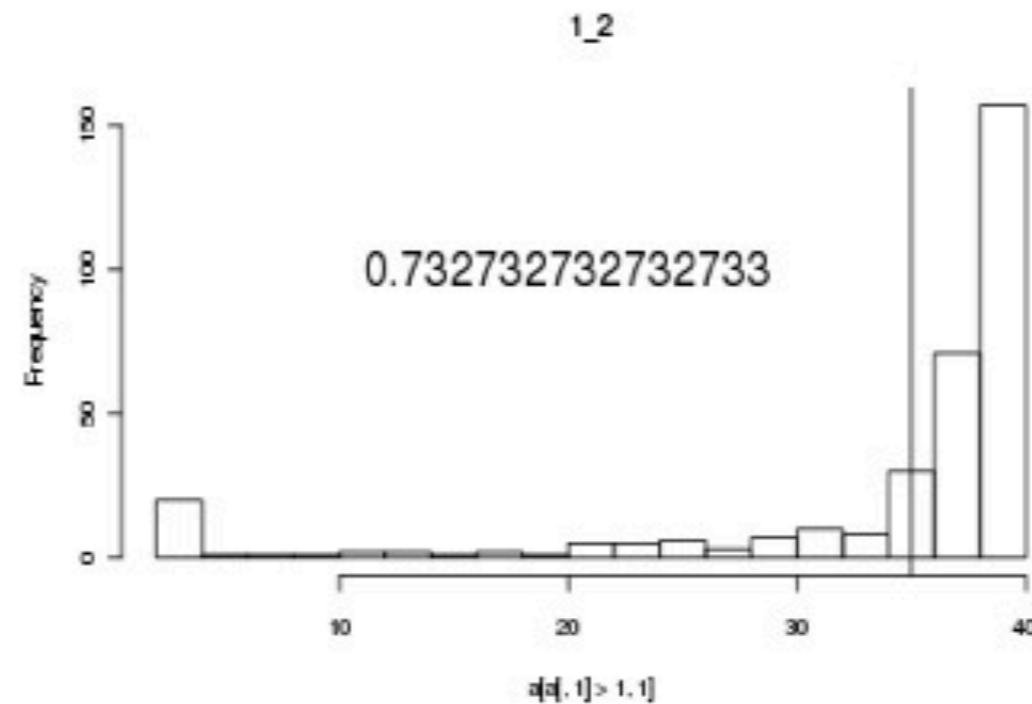
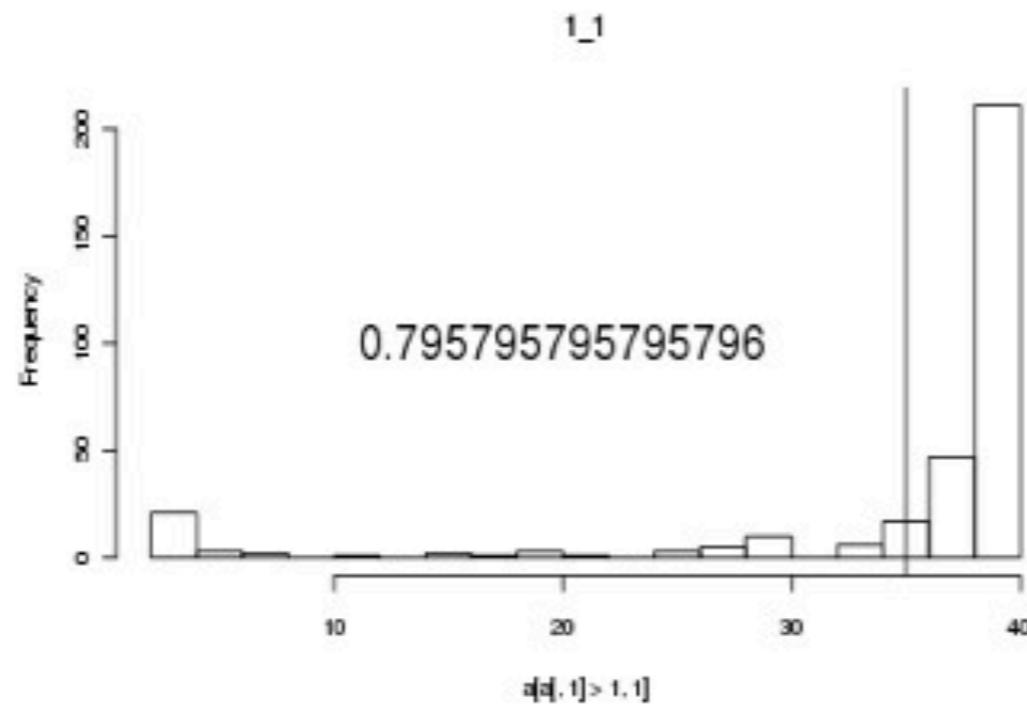
The first thing you'll want to do is to examine the quality profiles of your data.

There are many tools that do this, we use python and R.

Plot the quality of each base across reads:



Make a histogram of the mean quality of each read:



Use these plots to decide if you need to trim off the end of reads or discard low quality reads.

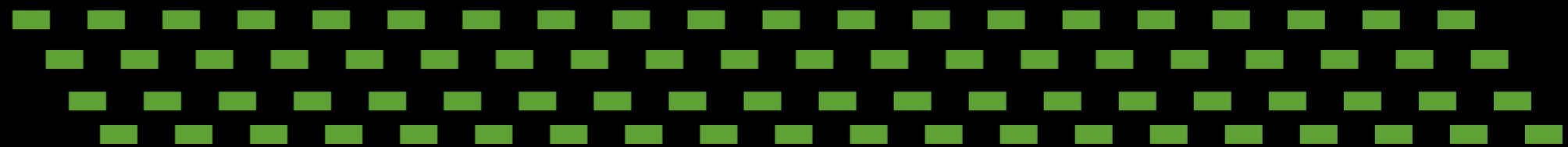
Assembly

Starting
material

Fragment ↓



Prepare library,
sequence ↓



Assembly ↓

Final
product



Overlap assemblers that work fine
on large sanger datasets don't
scale to these very large data sets

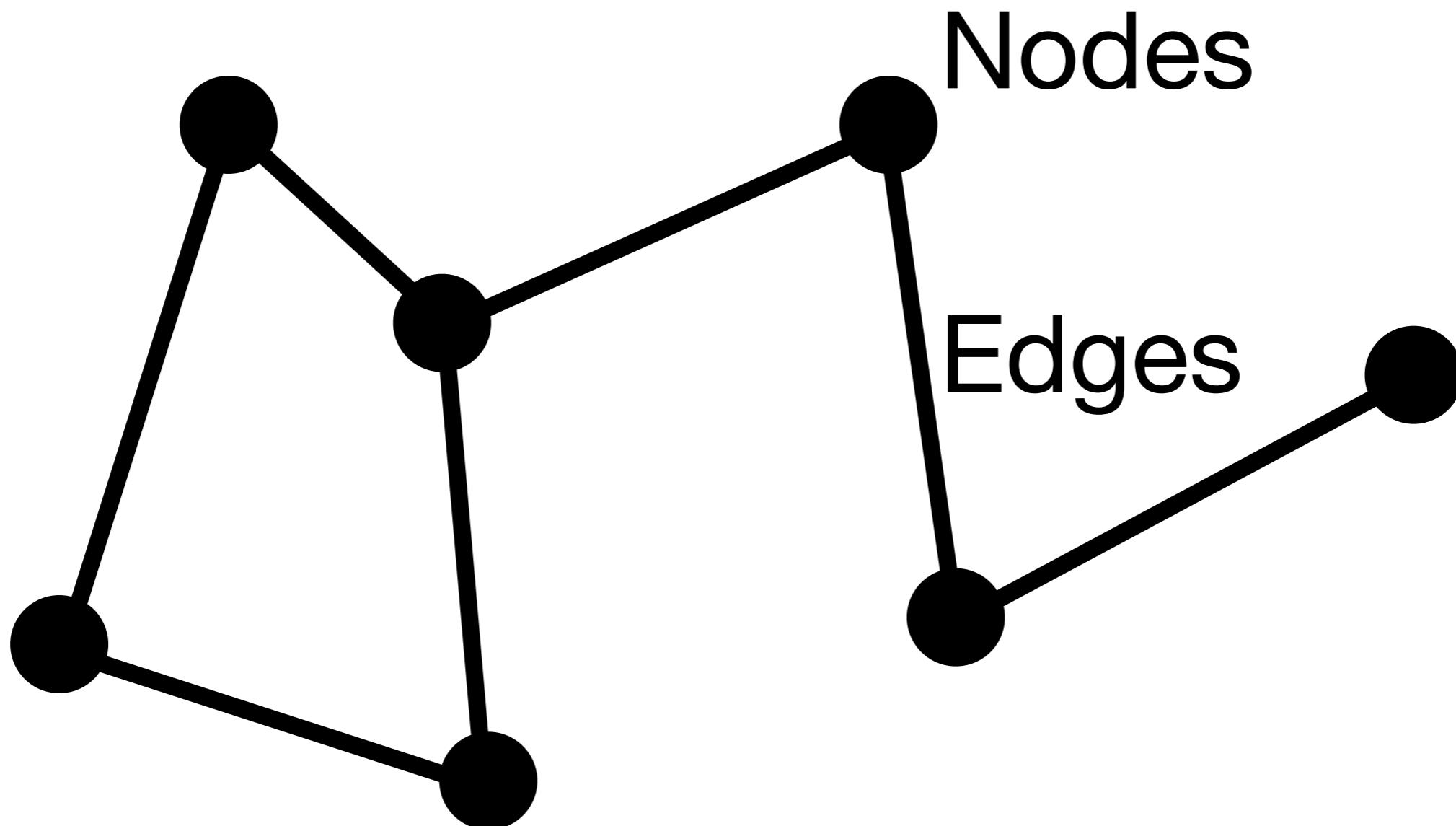
The number of pairwise
comparisons that are needed to
detect overlap become intractable

A new generation of de Bruijn graph assemblers have been developed to meet these challenges

Better defined memory footprint

Simpler comparisons between sequences

What is a graph?



The first step in de Bruijn graph assembly is breaking each read down into all sequences of k length

actgtcat →

actg
ctgt
tgtc
gtca
tcat

There are 4^k possible k-mers

In practice, k is often in the 25-70 range

The k-mers are loaded into a hash table:

actg	1
ctgt	1
tgtc	1
gtca	1
tcat	1

A de Bruijn graph is constructed from the has table

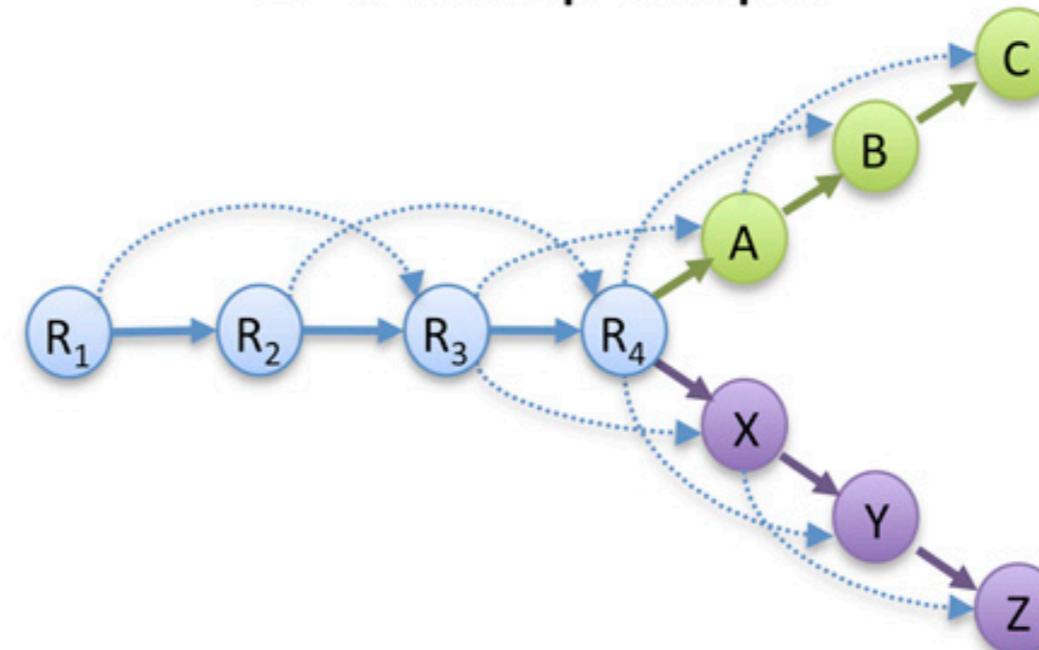
Each node corresponds to a k-mer sequence from the hash table

An edge unites each node that extends another node by one base pair

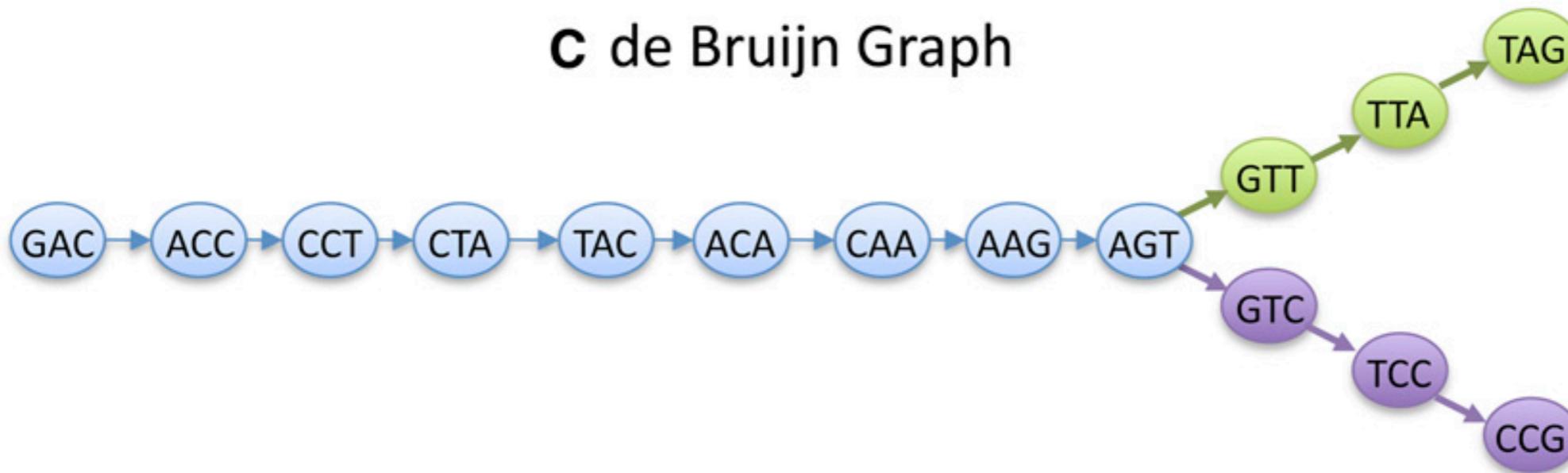
A Read Layout

$R_1:$	GACCTACA
$R_2:$	ACCTACAA
$R_3:$	CCTACAAG
$R_4:$	CTACAAGT
A:	TACAAGTT
B:	ACAAGTTA
C:	CAAGTTAG
X:	TACAAGTC
Y:	ACAAGTCC
Z:	CAAGTCCG

B Overlap Graph



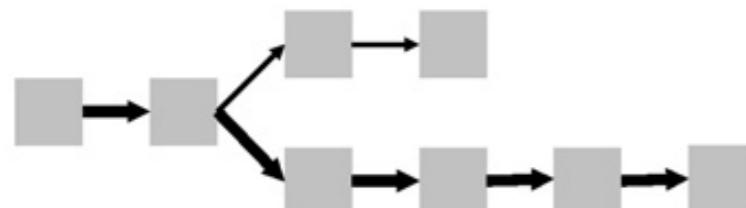
C de Bruijn Graph



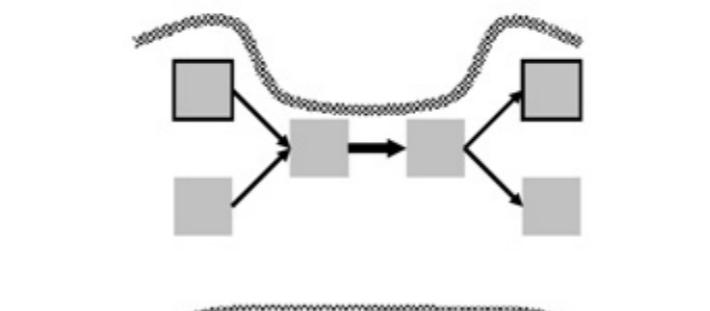
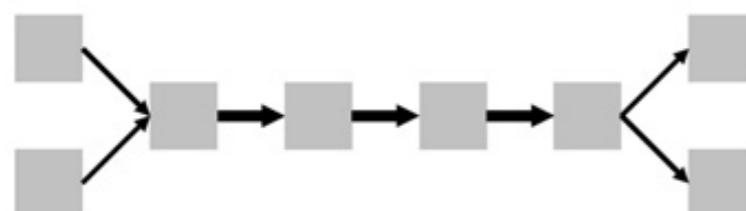
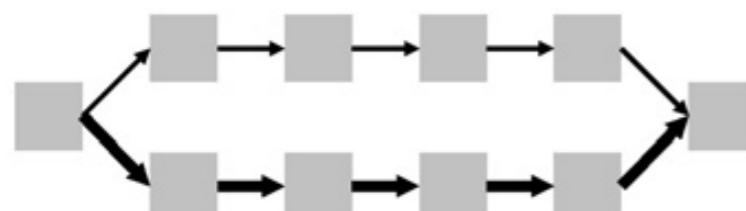
Paths through the de Bruijn graph are assembled sequences

These paths can be very complicated due to sequencing error, snp's, splicing variants, repeats, etc

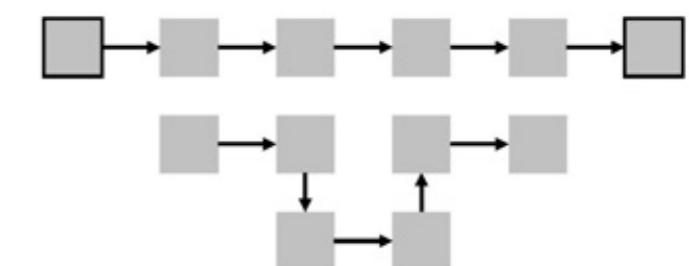
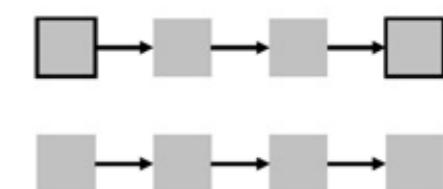
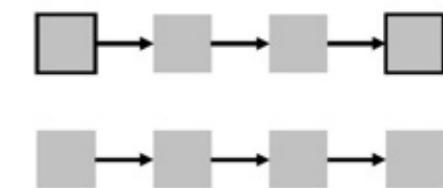
The graphs require considerable post-processing to simplify them (pop bubbles, trim dead ends, etc)



(before)



(after)



Assembly takes a lot of RAM!

One lane of Illumina HiSeq data can require hundreds of gigabytes of RAM to assemble

This is one of the largest challenges for using next-generation sequencing data to build trees

Eliminating low-quality data can greatly reduce RAM requirements

Genome

Transcriptome

Start with DNA

Get all genes,
regulatory regions, etc

Genomes can be
really big

Can be hard to
identify genes

Start with mRNA

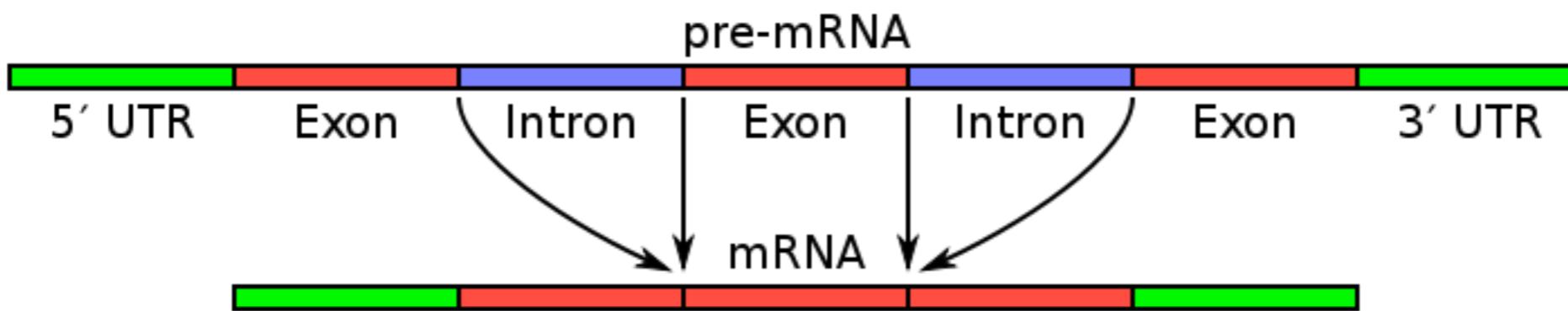
Get a snapshot of
active genes

Almost all data is from
coding regions

Handling RNA is tricky

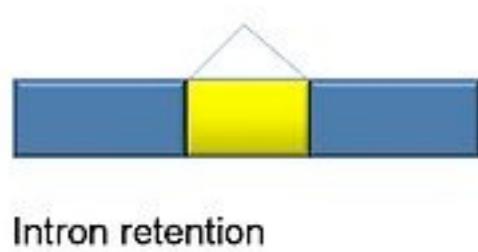
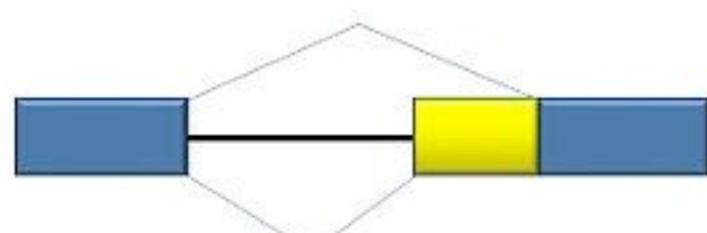
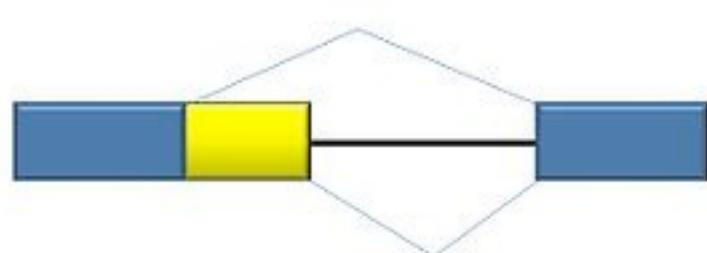
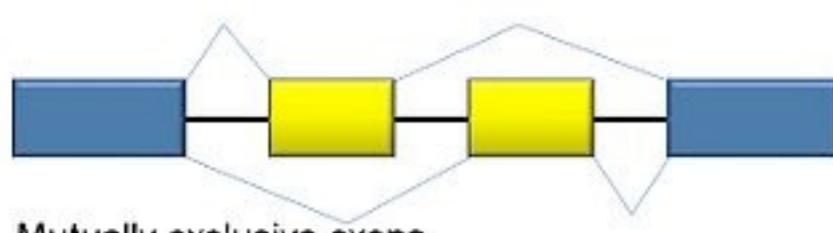
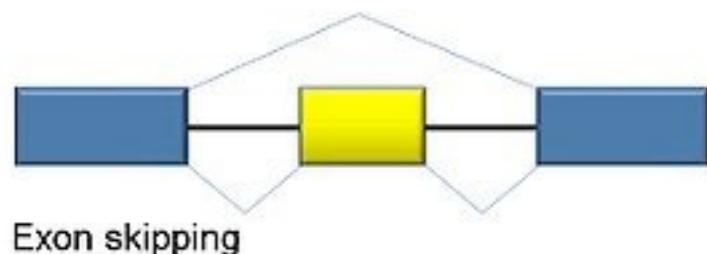
Transcript splicing

mRNA's are spliced before leaving the nucleus



en.wikipedia.org/wiki/File:Pre-mRNA_to_mRNA.svg

Transcript splicing



With deep sequencing,
many splice variants
are sequenced for
each gene

Assembly results...

Genome

...aagtcagtggagatgcaccatgagacaccttggagaagaagctgtccctggagacaatgtgggt...

Transcript

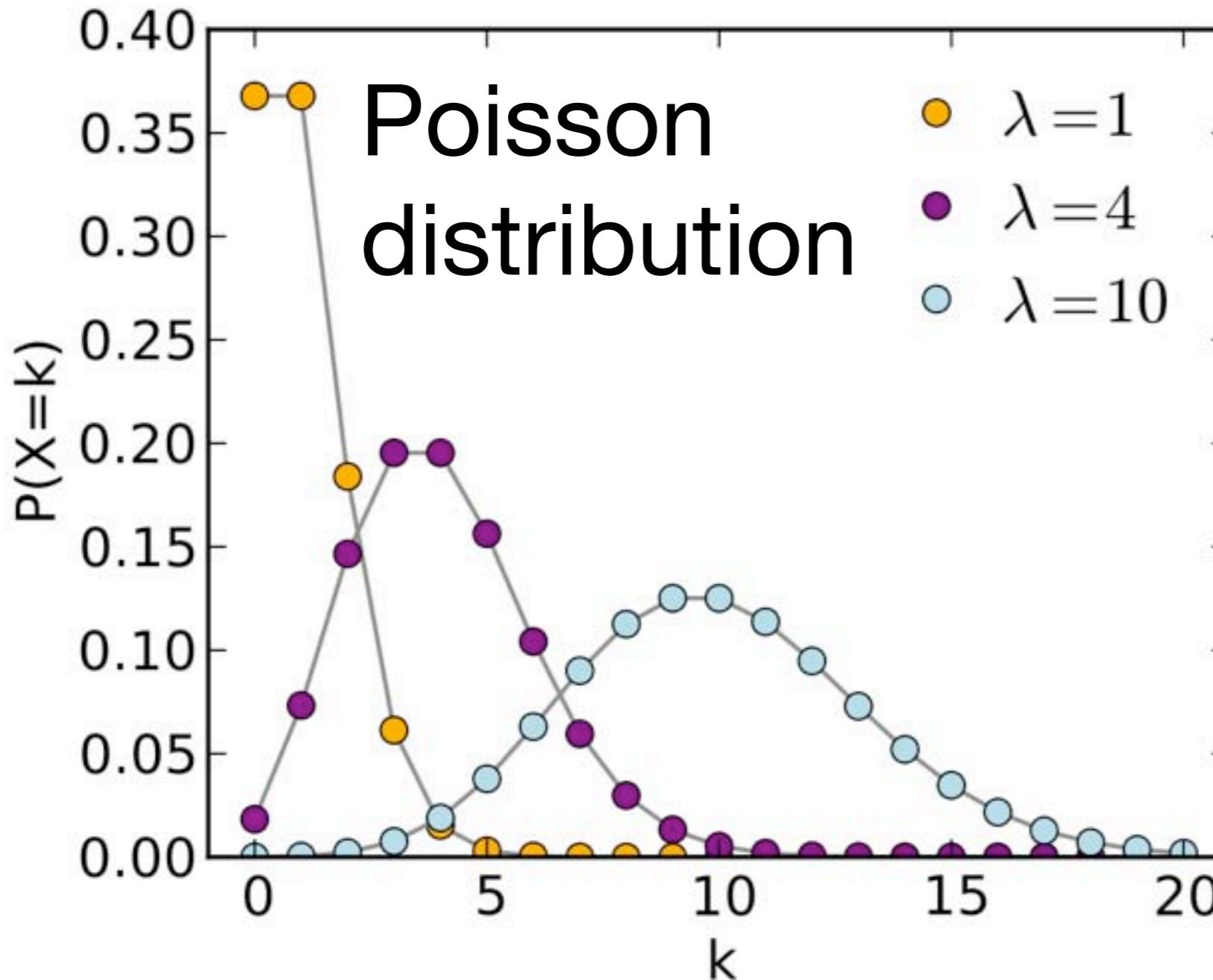
...aagtcagta ggagatgcaccatgag
ccttggagaagaag ctgtccctgg gtcct
agacaatgtgggt...

Splice variants

- Different splice variants for a given gene can vary widely in abundance
- Deep sequencing captures some “intermediate splice variants”, molecules in the process of being spliced
- Sequencing and assembly errors can be misinterpreted as splice variants
- Data may be insufficient to predict splice variants

It gets worse...

Genomes are uniform depth



[en.wikipedia.org/wiki/
File:Poisson_pmf.svg](https://en.wikipedia.org/wiki/File:Poisson_pmf.svg)

Assemblers can make assumptions about uniform distribution of sequencing effort

Expression differences mean:

- Can't assume that the expected frequency of sequences is uniform across or even within genes
- Low copy number doesn't necessarily indicate an error
- High copy number doesn't necessarily indicate a repeat
- Sequencing error is hard to accomodate in transcriptomes

When assembling transcriptomes, it is essential to use an assembler that can explicitly accommodate splice variants and expression differences!!!!

Transcriptome assemblers
include:

Newbler (Roche)

Oases (www.ebi.ac.uk/~zerbino/oases)

Trinity (trinityrnaseq.sourceforge.net)

TransAbyss (www.bcgsc.ca/platform/bioinfo/software/trans-abyss)

Biological
concept

Newbler
term

Oases
term

Gene

isogroup

locus

Splice
variant

isotig

transcript

exon

contig

contig

Post-assembly annotation:

- Selection of exemplar transcripts for each gene
- Blastx to a taxon restricted subset of the NCBI nr database
- Translation with prot4est

Identifying homologs

Phylogenetic tools build trees
from homologous characters

Most phylogenetic tools
assume character homology,
they can't evaluate it

We need to make a first pass
with phenetic tools

Throw all sequences for all taxa in a study into a hat

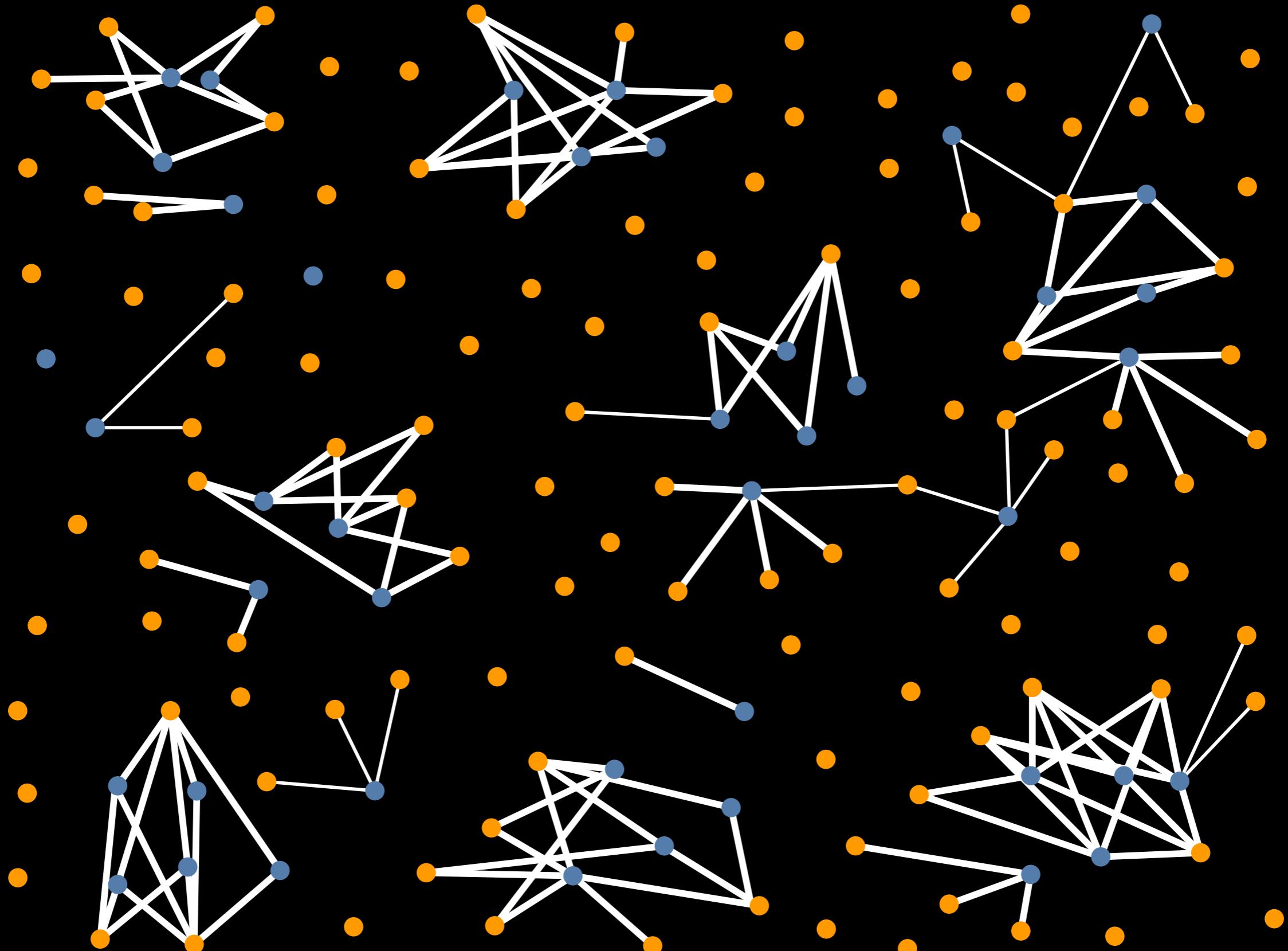
Make all pairwise sequence comparisons (eg blast)

Construct a graph where nodes are sequences and edges indicate similarity



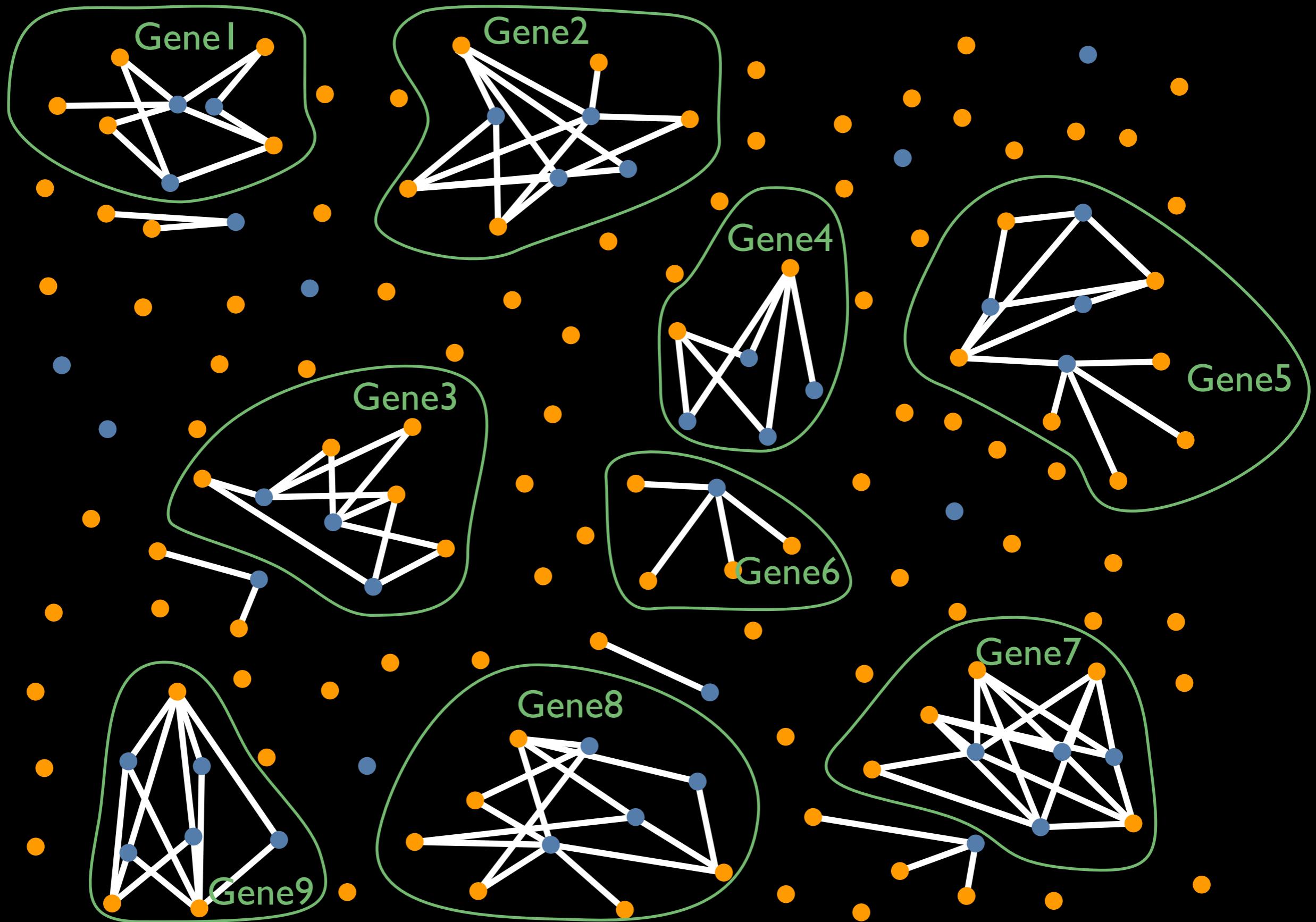
Nodes are sequences, thickness of edges indicate similarity

Casey Dunn



Nodes are sequences, thickness of edges indicate similarity

Casey Dunn



Nodes are sequences, thickness of edges indicate similarity

Casey Dunn

We use:

Blastp on a cuda array, assign
-log₁₀(e-value) for edge weight

Throw away edges < 20

mcl for clustering (inflation ~2.1)

Apply taxon sampling criteria

Sequence analysis

Advance Access publication November 29, 2010

Improving the quality of protein similarity network clustering algorithms using the network edge weight distribution

Leonard Apeltsin¹, John H. Morris¹, Patricia C. Babbitt^{1,2} and Thomas E. Ferrin^{1,2,*}

¹Department of Pharmaceutical Chemistry and ²Department of Bioengineering and Therapeutic Sciences,
University of California, San Francisco, CA, USA

Associate Editor: Burkhard Rost

Identifying orthologs

“The paralogy problem”

But paralogs aren’t inherently
a problem

The problem is miscribing
paralogs as orthologs

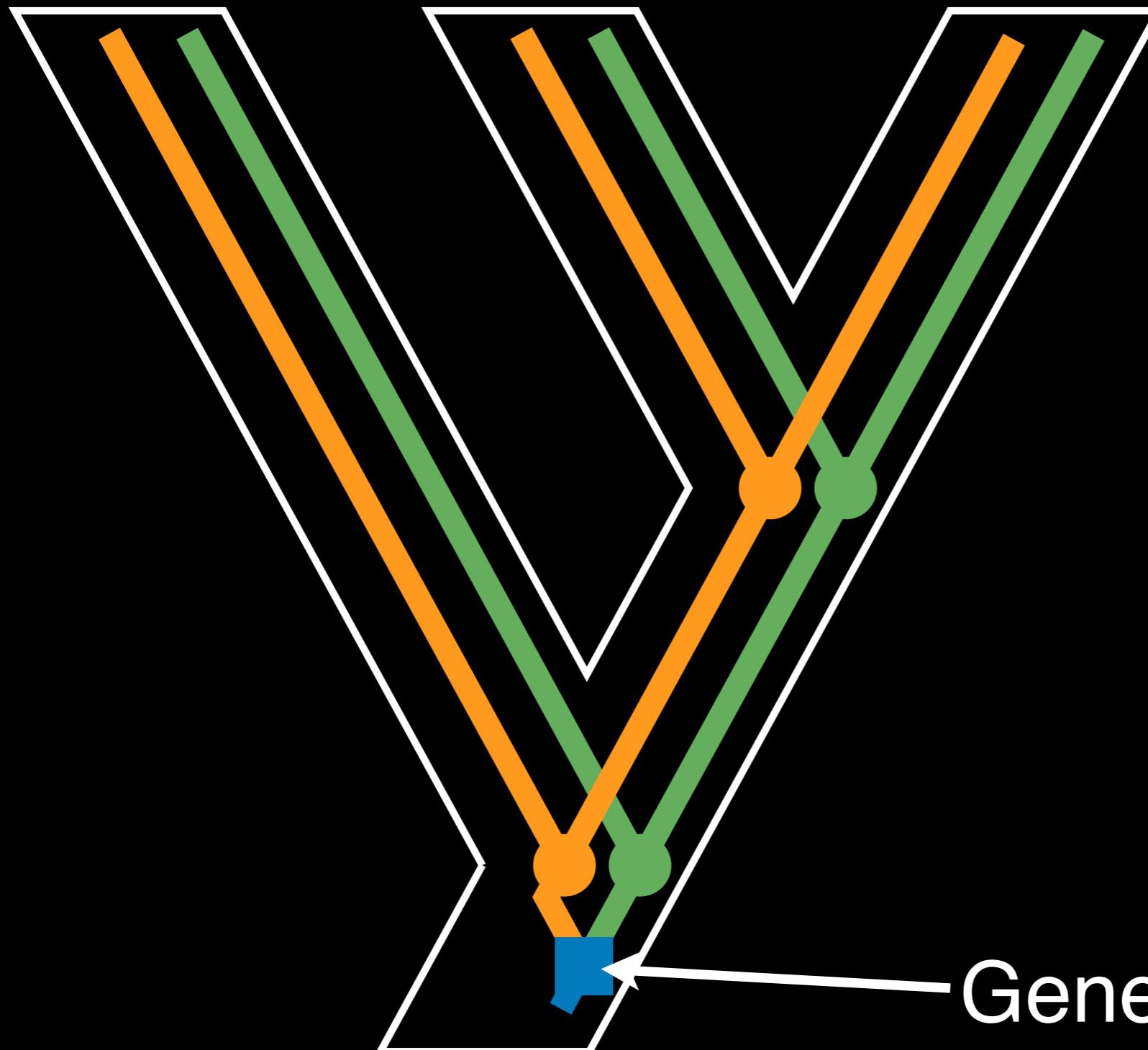
Species A



Species B



Species C



Gene divergence
due to duplication

Casey Dunn

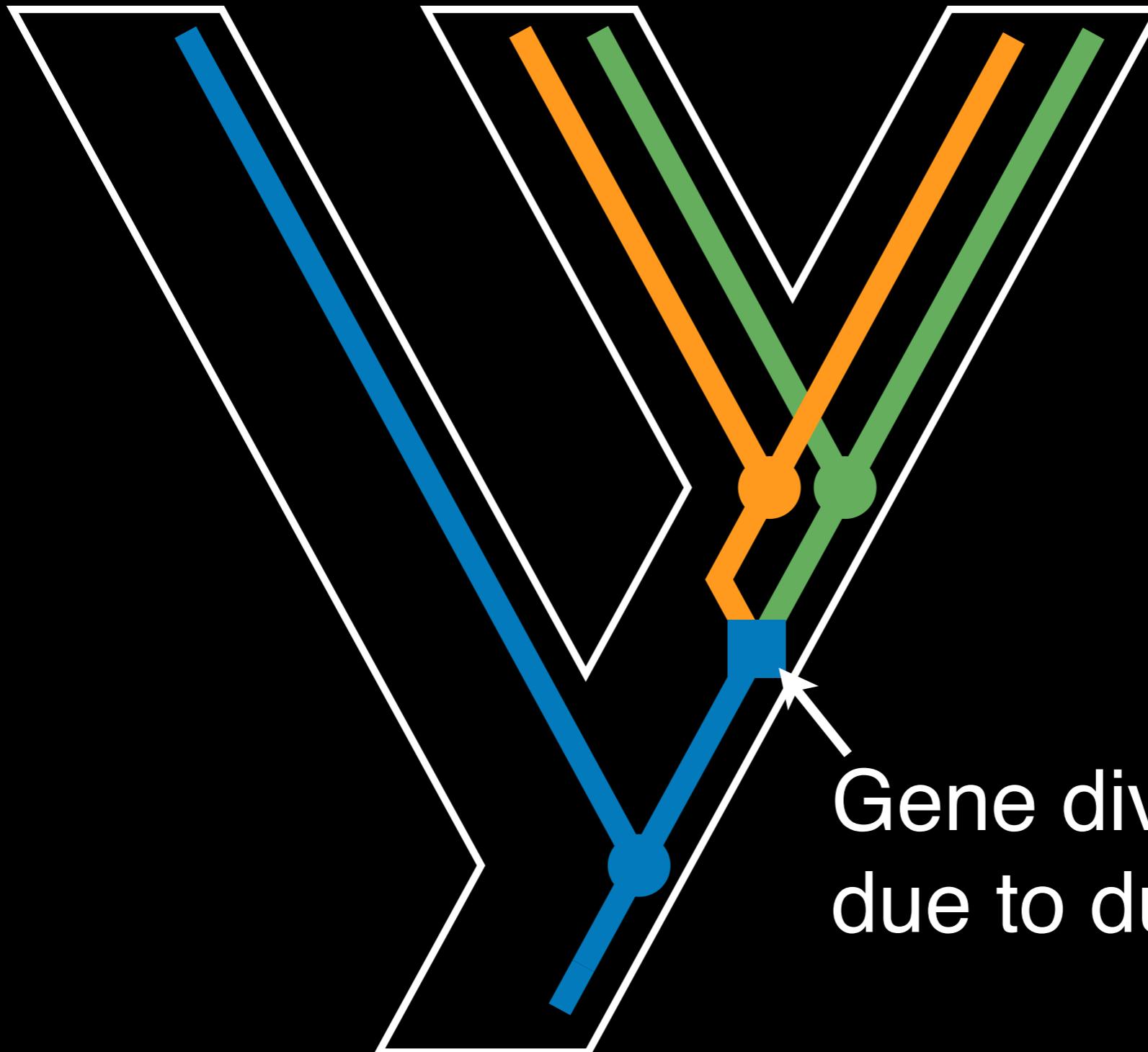
Species A



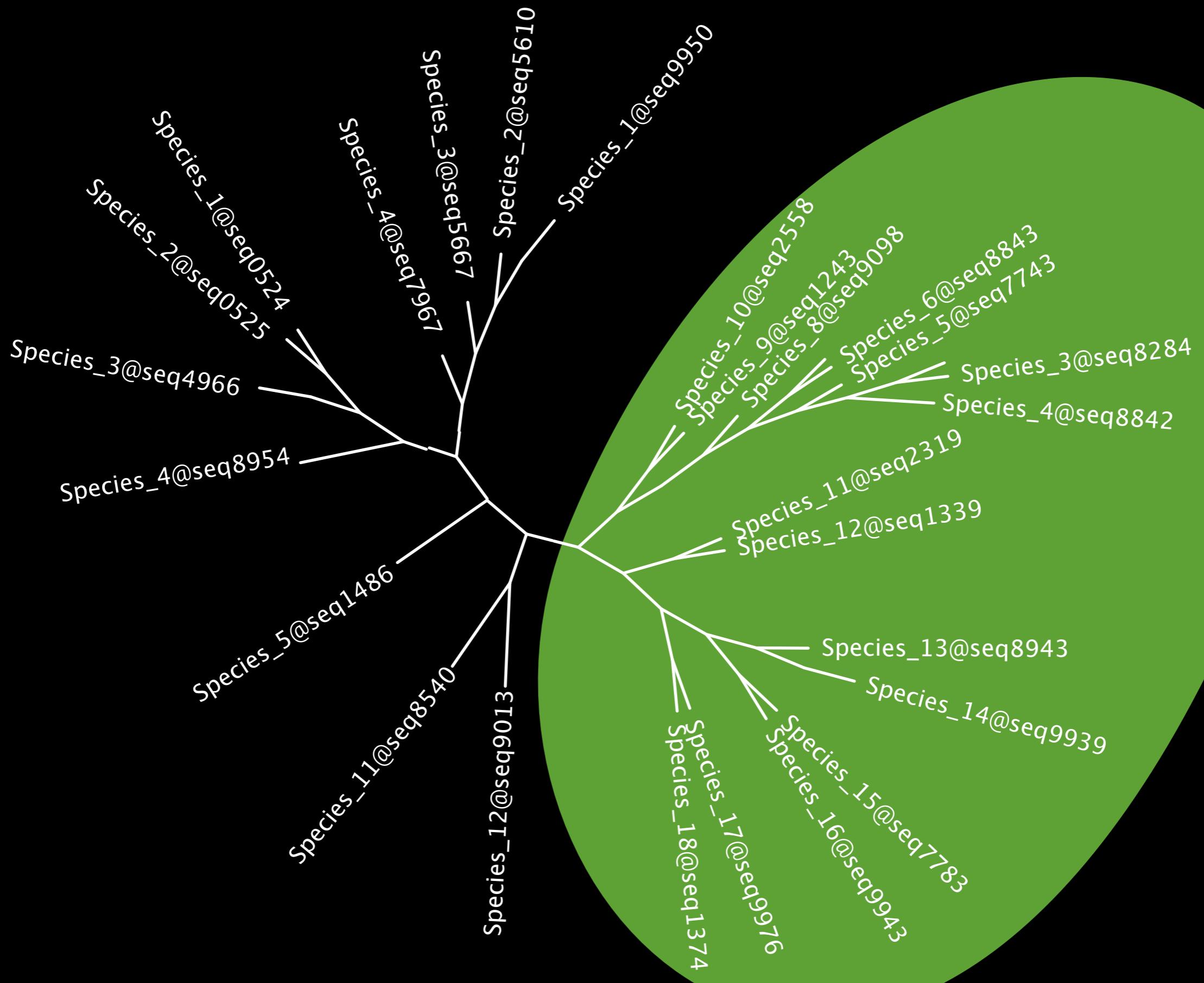
Species B

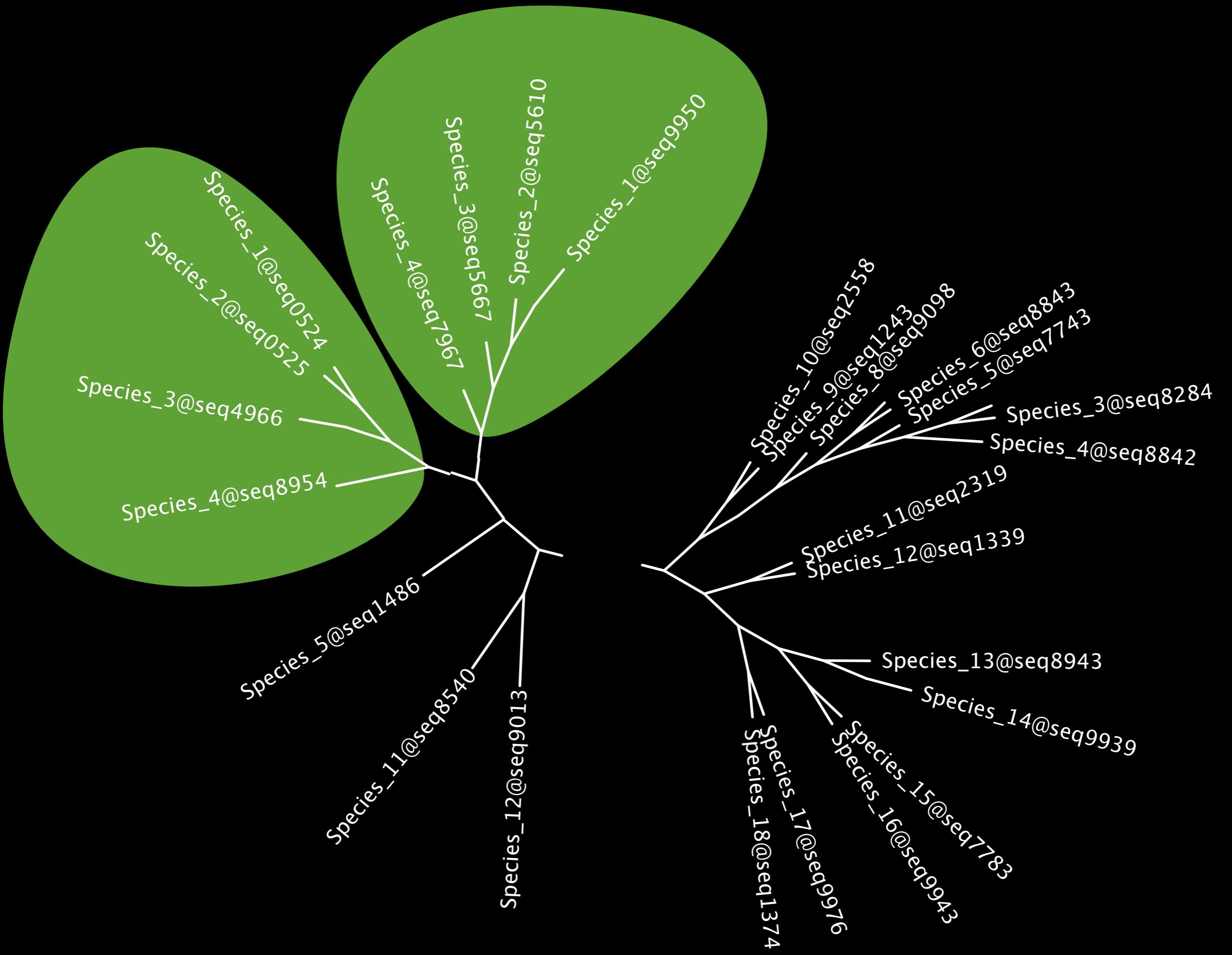


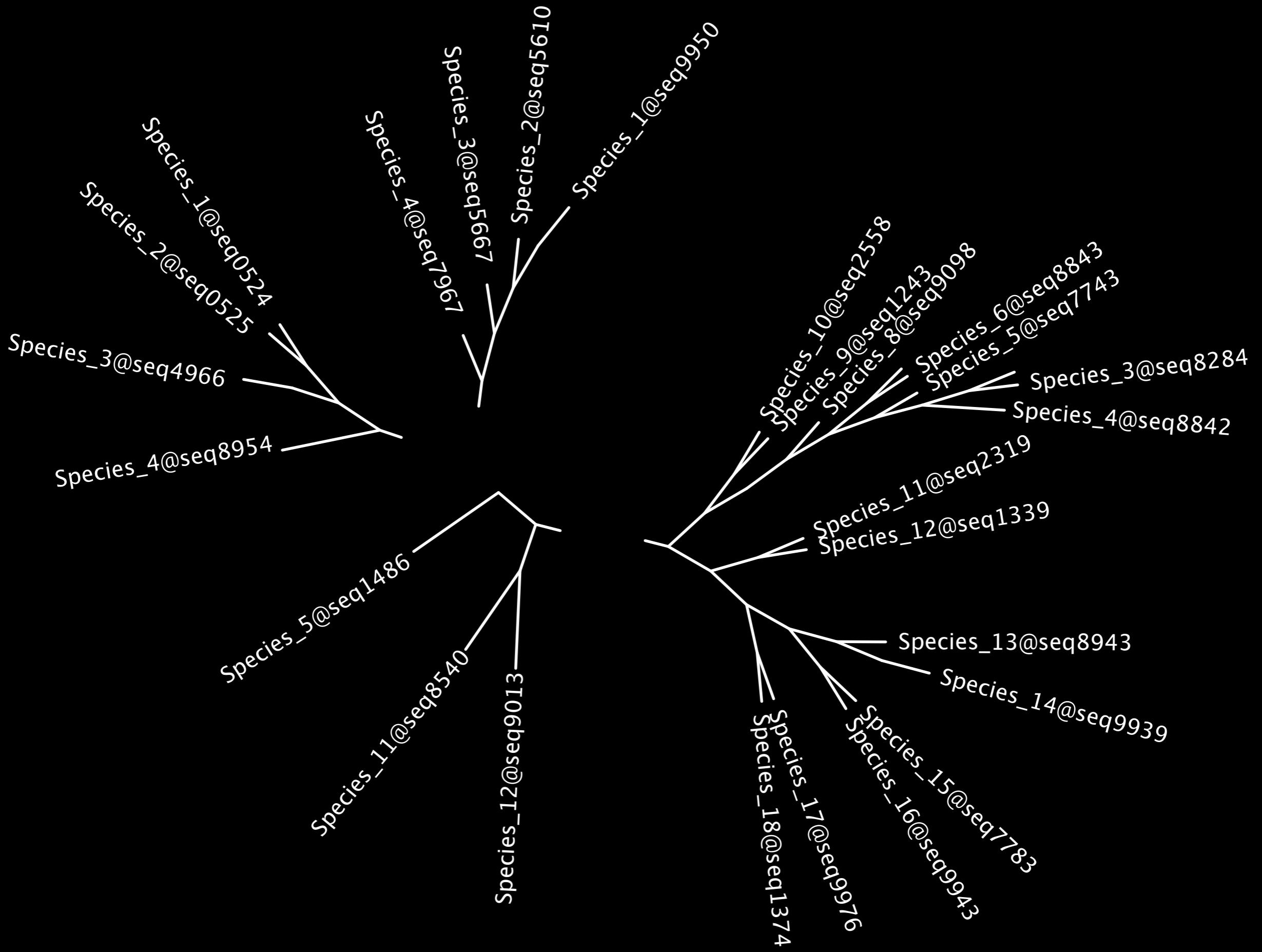
Species C



Gene divergence
due to duplication







Isolation of
Homologs

Isolation of
Orthologs

Evaluation of
Orthology

Phenetic

Phylogenetic

Phenetic

Phylogenetic

Once we have subtrees of orthologs...

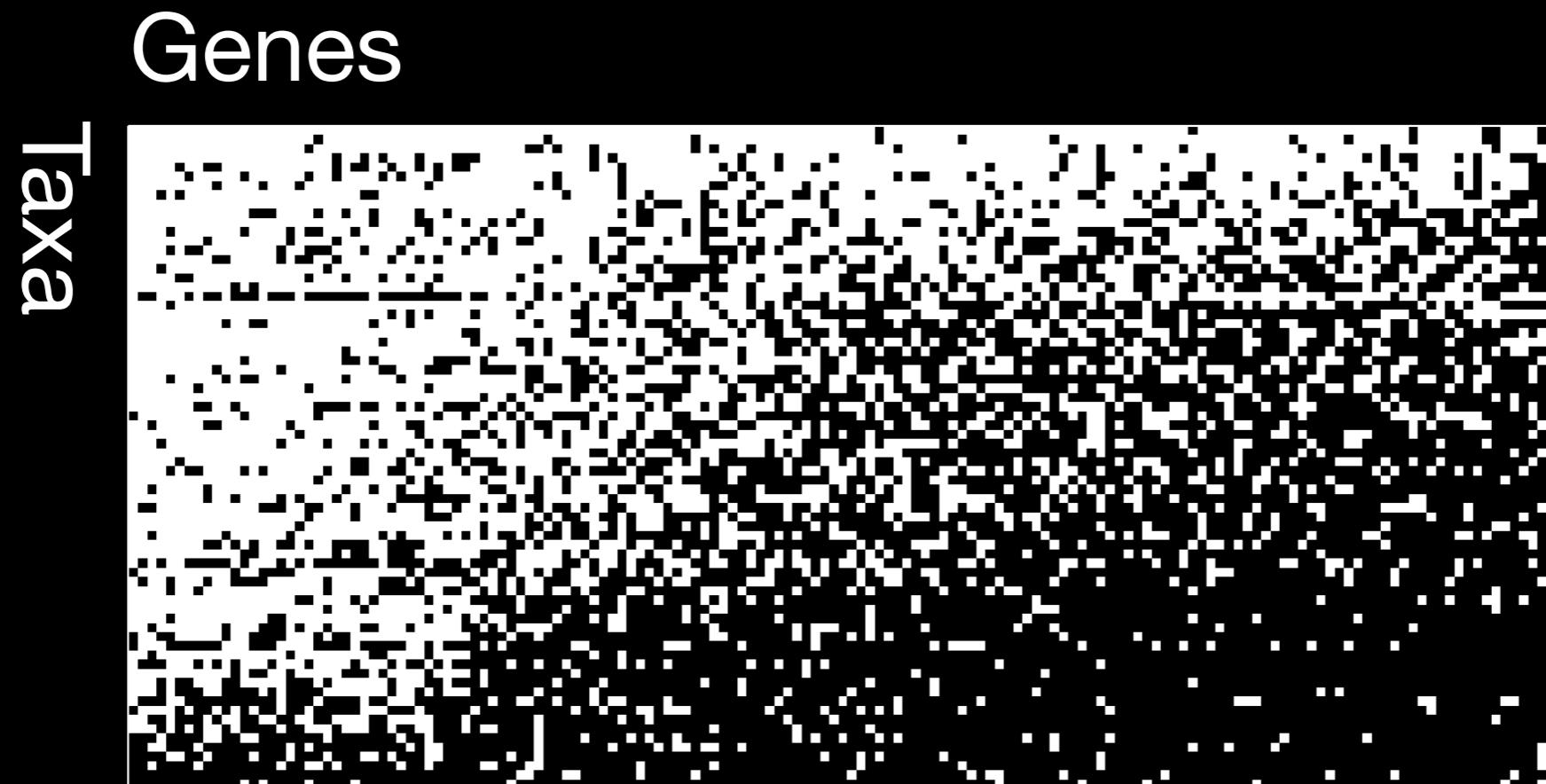
Align each ortholog

Concatenate!

There are many exciting alternatives to concatenation

As these become more computationally efficient, robust to missing data, etc they will be exciting to apply to these datasets

77 taxa, 150 Genes, >20k aa



White cells indicates sampled gene
50.9% gene sampling

Dunn *et al.*, 2008
doi:10.1038/nature06614

Casey Dunn

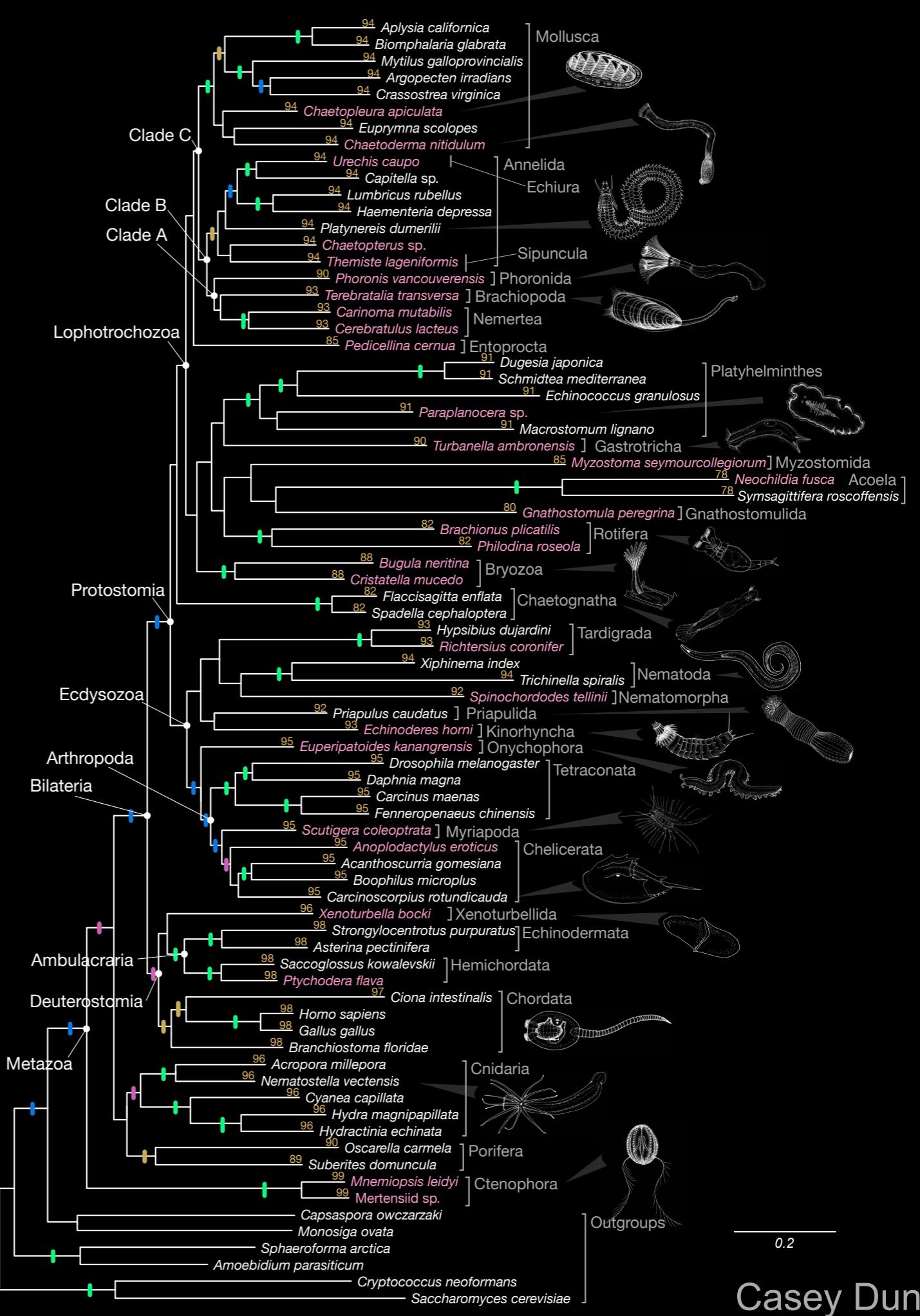
Dunn et al., 2008
doi:10.1038/
nature06614

150 Genes, >20k aa

Bootstrap support



raxML
1,000 BS replicates
WAG+Γ



0.2

Casey Dunn

Extracting more information from analyses

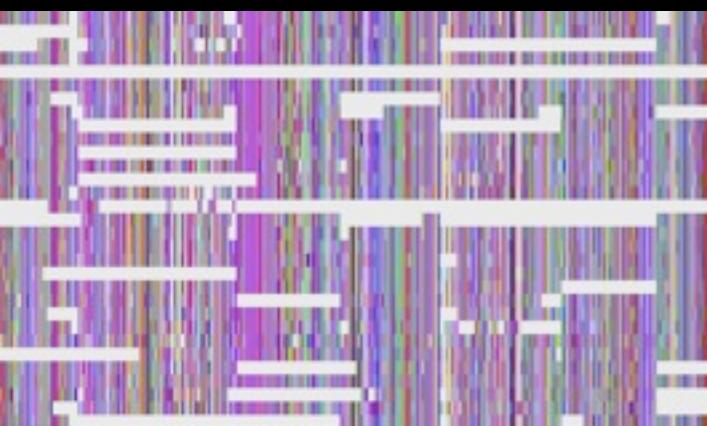
Mess with
your dataset
and analysis
settings

Rerun
your
analyses

But this doesn't work well for large analyses

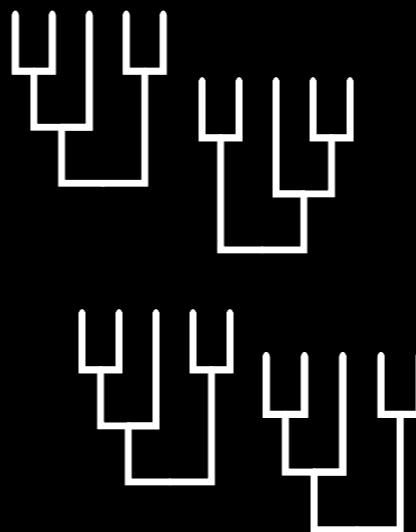
Information attrition

Matrix



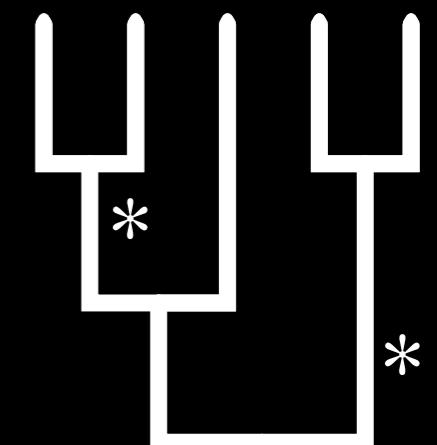
Inference

Tree set



Summary

“Final Product”



We throw away:

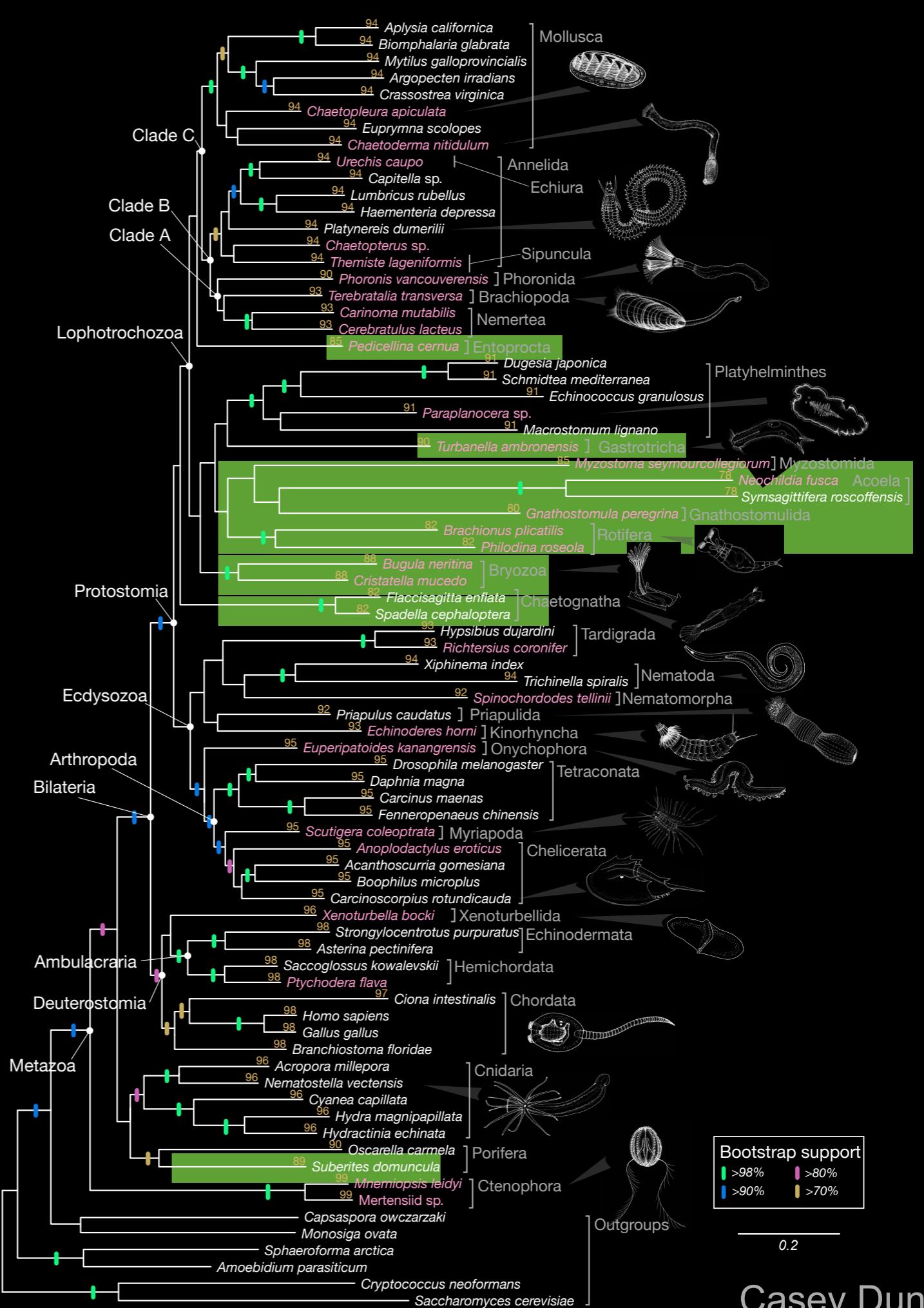
Explicit information on distribution of support and conflict

Non-consensus topological variation, information on variation in taxon placement

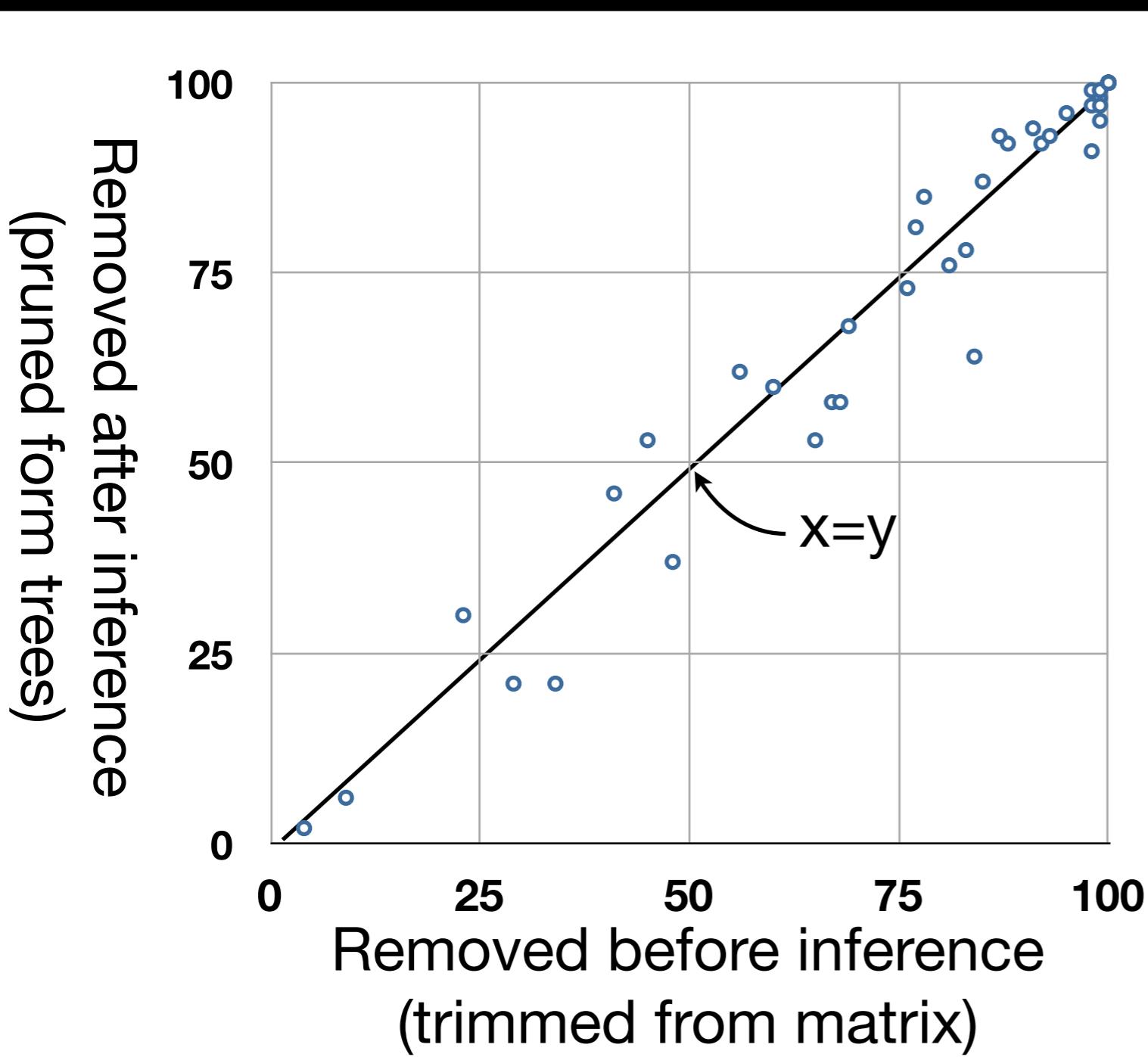
Unstable taxa can obscure support for relationships between stable taxa.

Leaf stability indeces
(Thorley & Wilkinson, 1999)
quantify the stability of each taxon.

Leaf Stability < 90%



Split frequencies following removal of unstable taxa

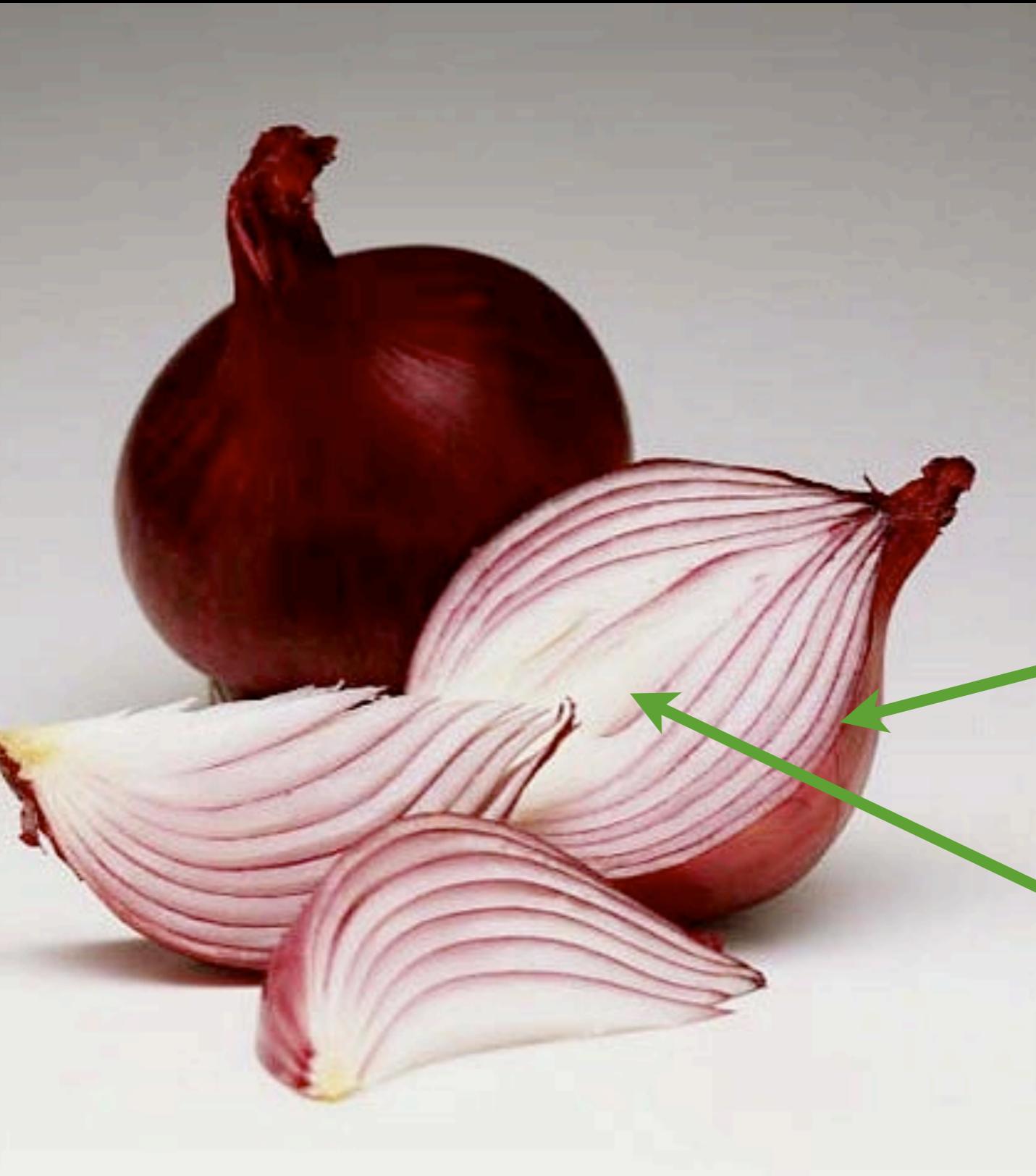


Unstable taxa have little or no influence on relationships between stable taxa

1,000 BS replicates

Casey Dunn

Treesets as onions



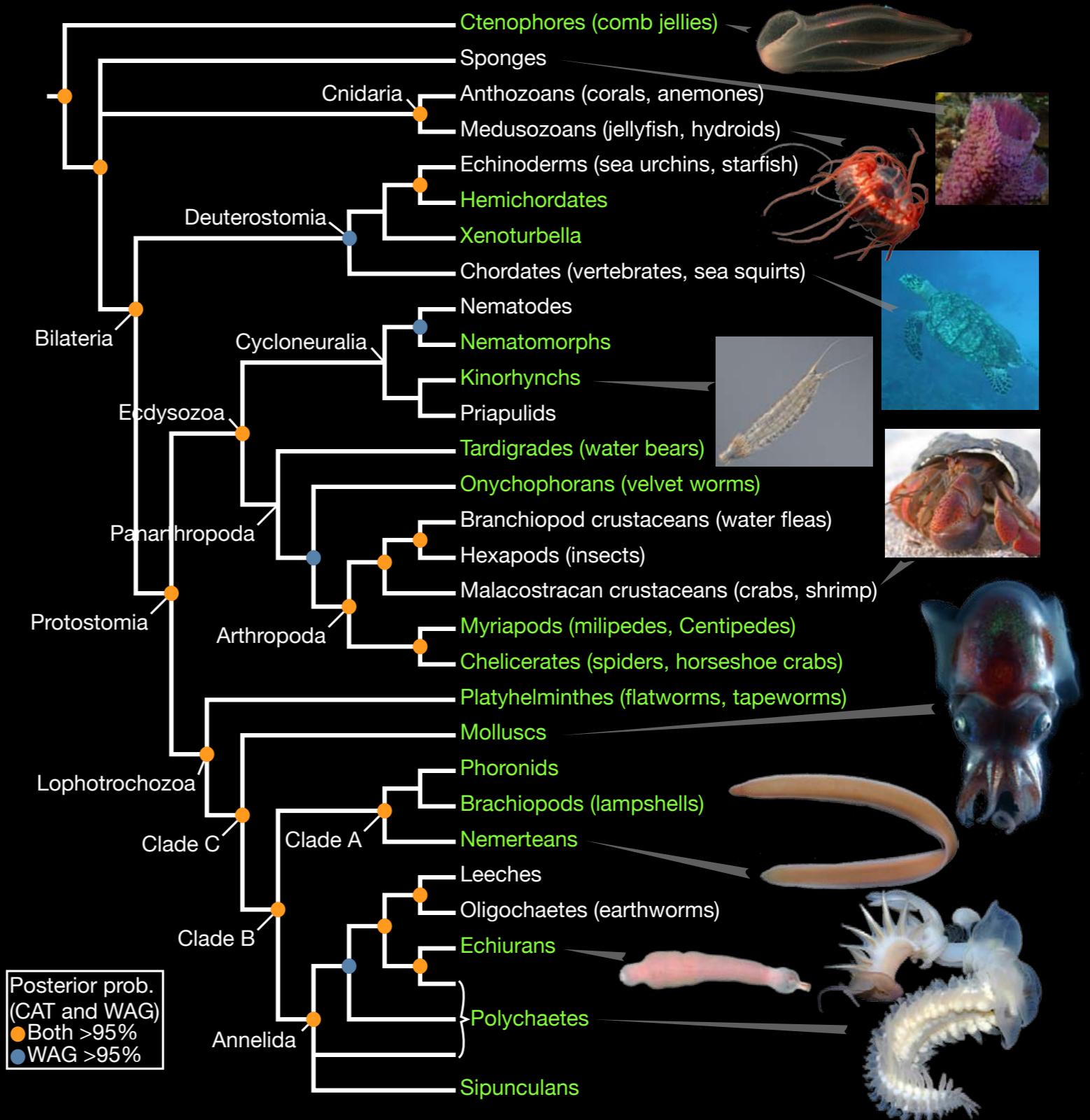
Least stable taxa

Most stable taxa

Investigation of stable taxa

Dunn *et al.*, 2008
doi:10.1038/
nature06614

150 Genes
>20k aa



Posterior prob.
(CAT and WAG)
● Both >95%
● WAG >95%

Casey Dunn

Tools for identifying and visualizing relationships between stable taxa



phyutility

code.google.com/p/phyutility/

Smith & Dunn (2008), Bioinformatics

doi:10.1093/bioinformatics/btm619

- Calculate leaf stabilities
- Prune taxa from treesets
- Explore the positions of unstable taxa

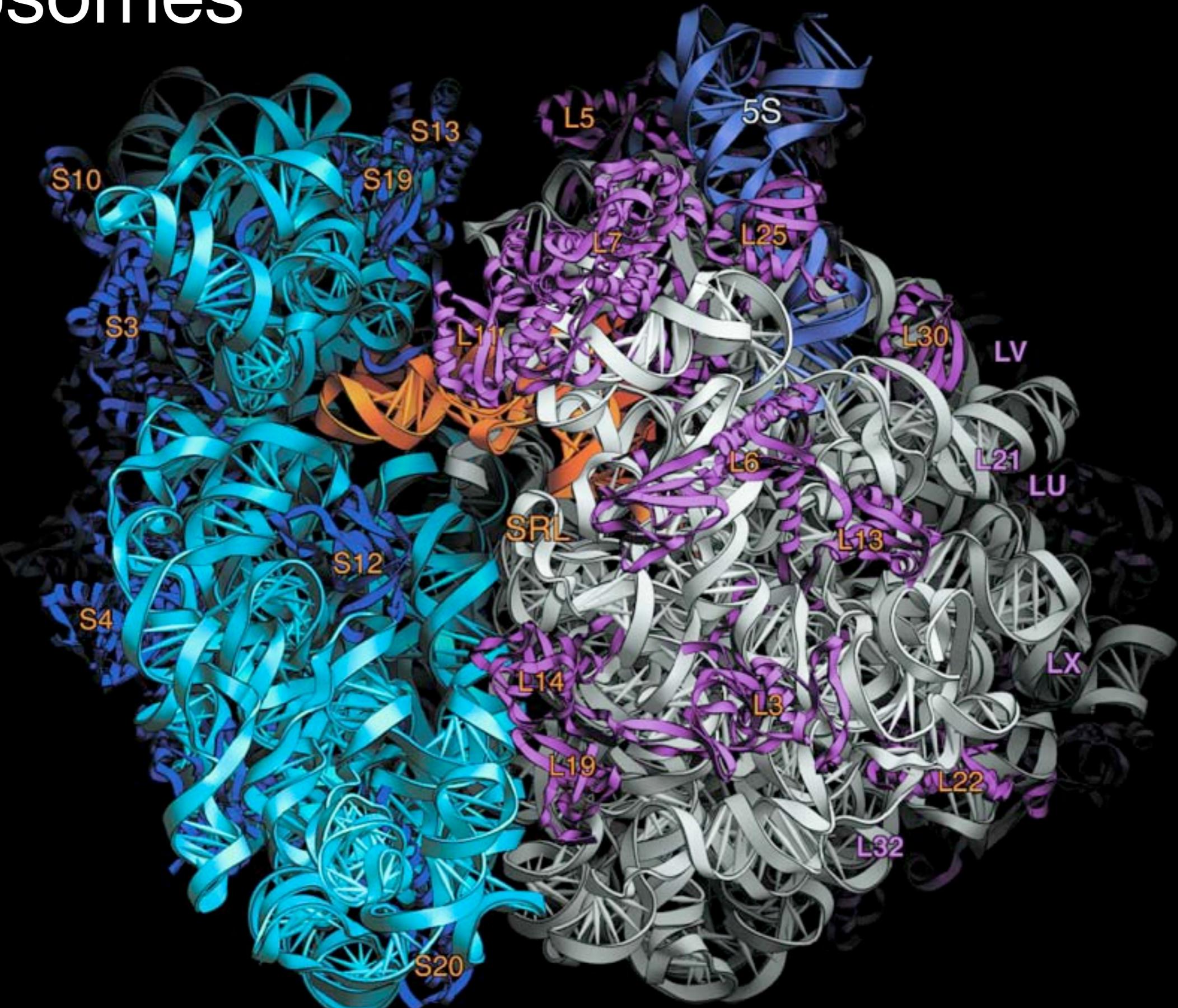
Casey Dunn

Beyond species trees

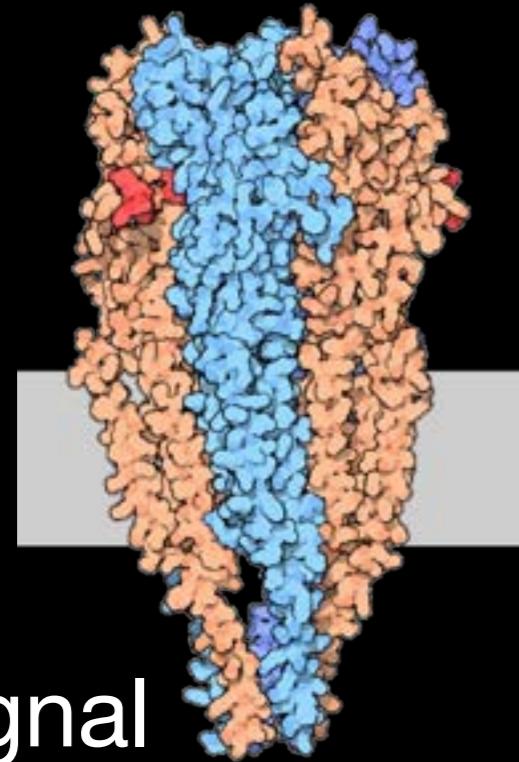
Phylogenies now generate:

- Species trees
- Extensive gene sequence data
- Well sampled gene trees

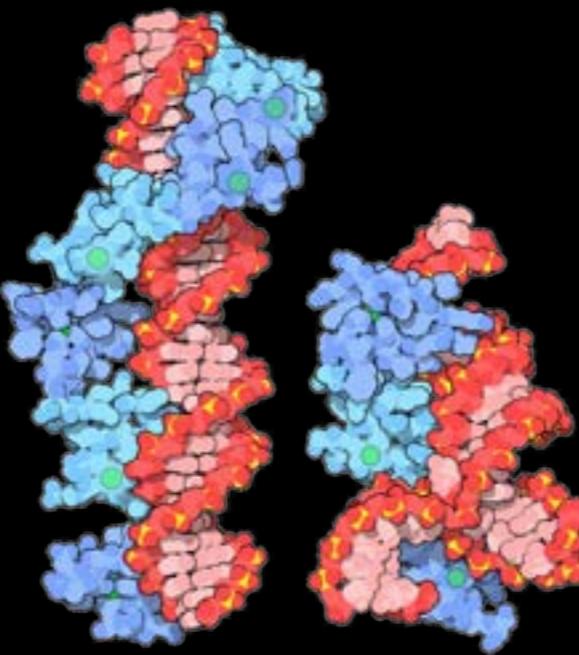
Ribosomes



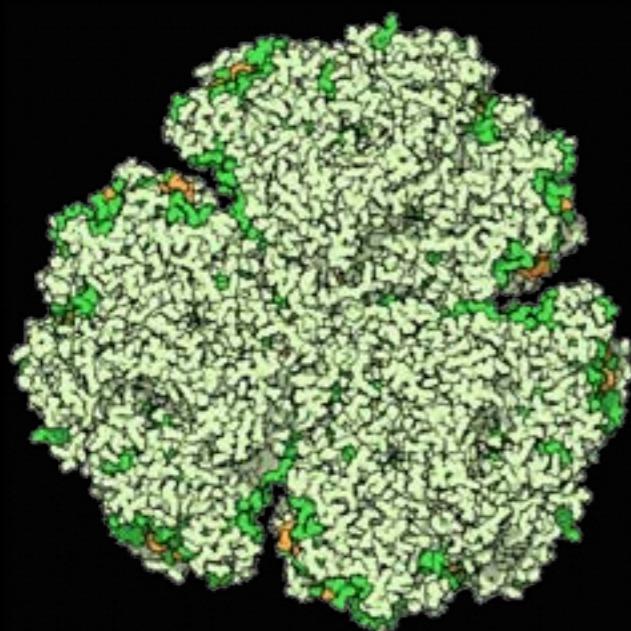
Sequences relevant to focal phenotypes



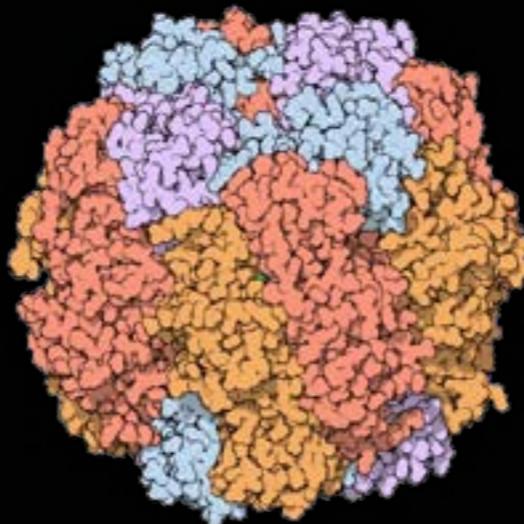
signal
transduction



morphogenesis



photosynthesis



carbon
sequestration

But we don't know which
genes are relevant to
which phenotypes

A small glimpse of a much greater schism

Genomes



```
>FZTB7Y04I0Z6U rank=0418088 x=3584.5 y=3492.0 length=457
CAAGGTCTGAACCAACAGTGGATACAATTCAAATGACCGGAATGAAAGAAATCCATT
CACTGGTCATTTGACTTCCACGTGTTGGATCGATTTCTCTTTAACCTCTCCT
GAAGCTTGTGTTGCTCTGGCATGTTGGCATGTCAGTACACAACATATAACC
TCGGTTCTGTTGAACCTCACGTCGTTGCCTCAAGATGCTCGAAATACCGCCCG
TCGAAGAAACCTGGTAAGCCTCCAGCAGAAGGCCATTAGTGTATTCCACCGGTAT
CTACTTGACGTATCTGTAGAAAAGGGAAAGCAGAGGGTGCCTGCTCAACTGAAACGC
ATTGACCACTACTCTTGAAGAAGACAAACTGCCAGCACTGGGAGACTGTCC
CTTCGGAACTCTCGCCGAGTTTGGAACACCTTGTTC
>FZTB7Y04I05F0 rank=0418094 x=3472.5 y=2494.5 length=288
AATGAAATATGCTGAGCAGTCAAGTTCTACTCACGAAAGAAACACATTGTAGATGG
TTTCATCGAACCAACAAATGAAGAGGGTTGGTTGATCCTTAGAAGAATTGGTGA
ACAGTTGAATAAGGGTGTGAAGAAAAGCTGAATCTGTAGAAAAGTGAAGAAGAGAAATTG
GCTGGATGGTGTGAAACACTTATCATTTGGTAAGAACACAAAAGGTATTCTGAATTTT
GGCTCACTGCAATGAAGAACGTTGAAATACTTGAAGATATGATTCAAGG
>FZTB7Y04I07J9 rank=0418096 x=3473.0 y=1143.0 length=421
AGGCCGGGCCCTTCGATTAAGATATCTAAAGAGTTGGTCTCCACGGAGCTAAGGCT
AACAAATCTAGTAAATCTGCATTGGTGAACCTTCTCTTAAATGCTGACACA
TCTGTATCCGAACTTGCCTGTATACAGTCCCTTATTTTATACGATGATATCGAT
ATCACAGGAGAAAAAAATGGCTAAATCATGGCTGCTGCCAACGTAAACGTAGAAACCTT
CTGGCCTGGACTCTCGCCAAGGCTCTCCAAAGGACGTAACATCGGTGACCTTATCTGCAA
TGAGGATCCTCGCAGCCGCTGCTCCAGGCCGCTGCTGCTGCTGGTGTGCTCCAGC
TGCTGCTGAAGAGAAGAAGAAAAGAGTCAGTTCAAGATGAGGATCAGATGATGATA
```

Evolutionary functional genomics

Morphology, function, ecology, development



Casey Dunn

Measuring expression



(S Haddock)



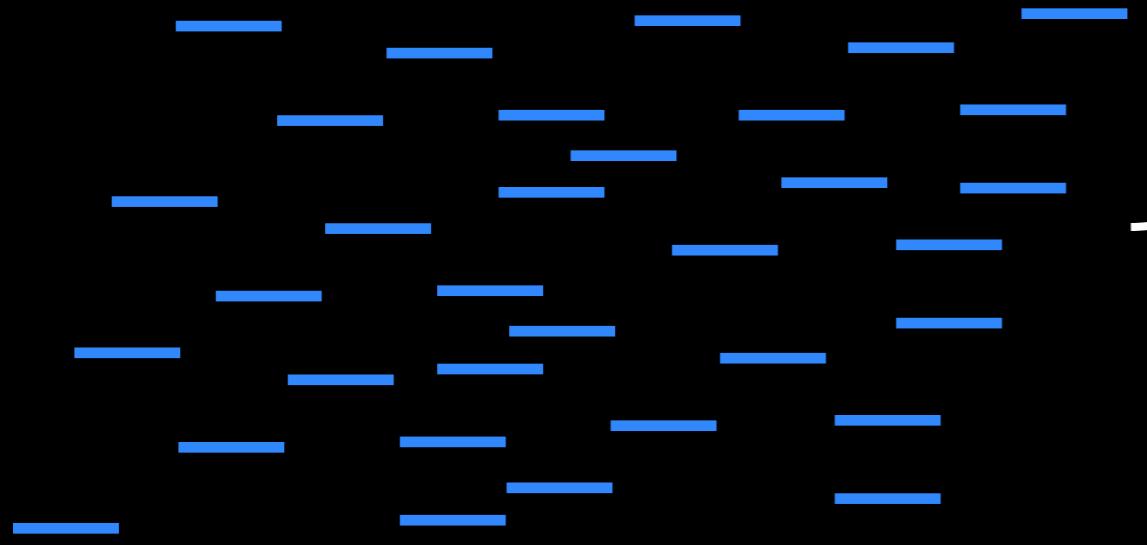
Casey Dunn

Which genes are
differentially expressed
between bodies in a
siphonophore colony?

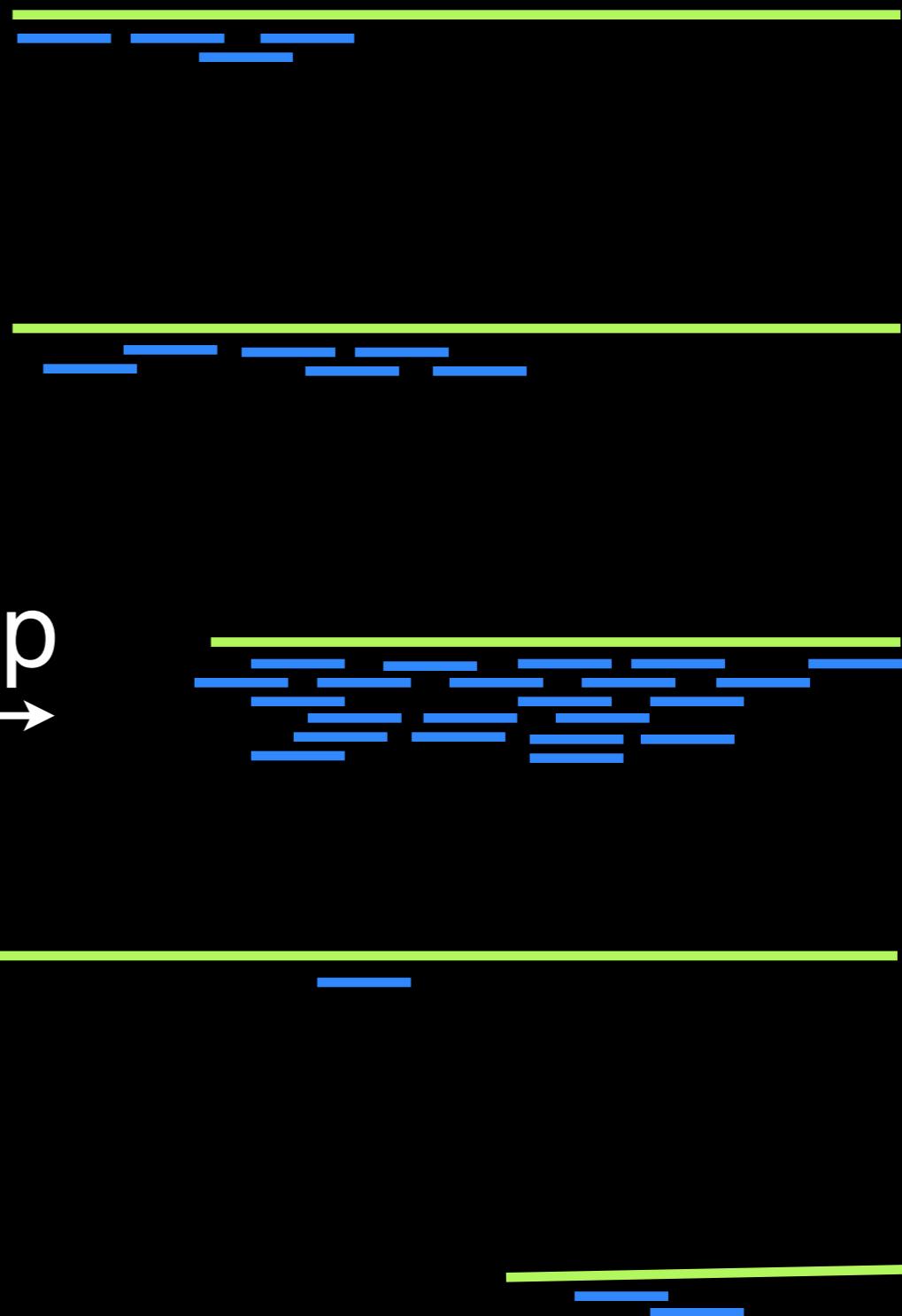
Reference



Reads



Map



Gene Count

Gene001 4

Gene002 6

→ Gene003 22

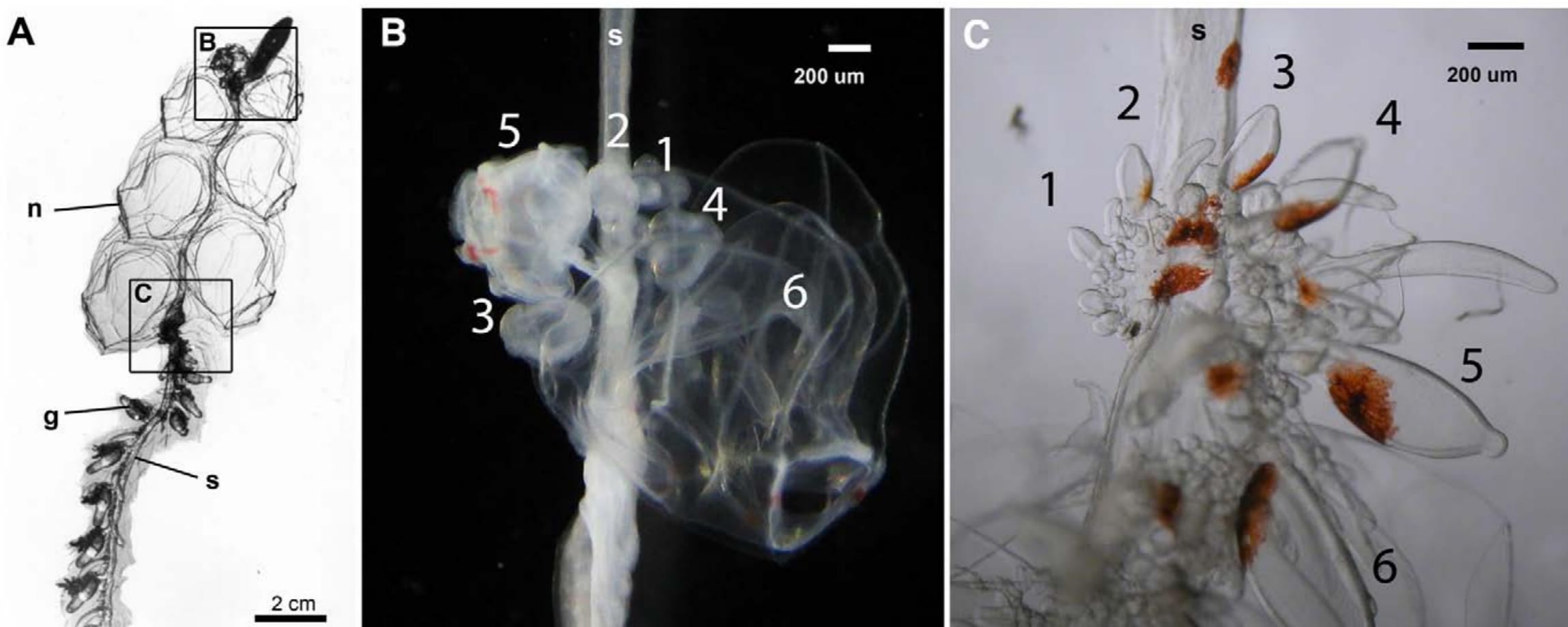
Gene004 1

Gene005 2

Differential Gene Expression in the Siphonophore *Nanomia bijuga* (Cnidaria) Assessed with Multiple Next-Generation Sequencing Workflows

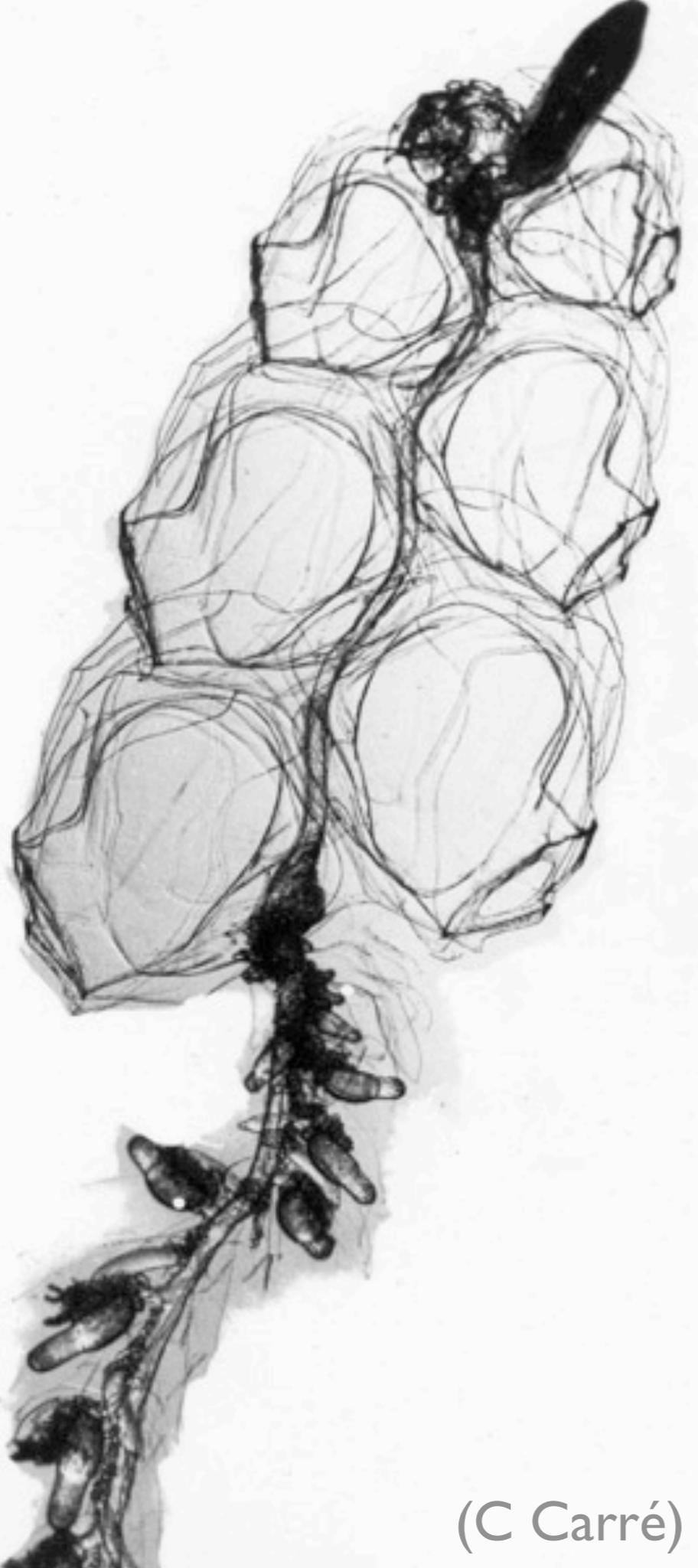
Stefan Siebert^{1*}, Mark D. Robinson^{2,3}, Sophia C. Tintori¹, Freya Goetz¹, Rebecca R. Helm¹, Stephen A. Smith^{1,4}, Nathan Shaner⁵, Steven H. D. Haddock⁵, Casey W. Dunn^{1*}

1 Department of Ecology and Evolutionary Biology, Brown University, Providence, Rhode Island, United States of America, **2** Epigenetics Laboratory, Cancer Research Program, Garvan Institute of Medical Research, Sydney, New South Wales, Australia, **3** Bioinformatics Division, Walter and Eliza Hall Institute of Medical Research, Parkville, Victoria, Australia, **4** Heidelberg Institute for Theoretical Studies, Heidelberg, Germany, **5** Monterey Bay Aquarium Research Institute, Moss Landing, California, United States of America



(dx.doi.org/10.1371/journal.pone.0022953)

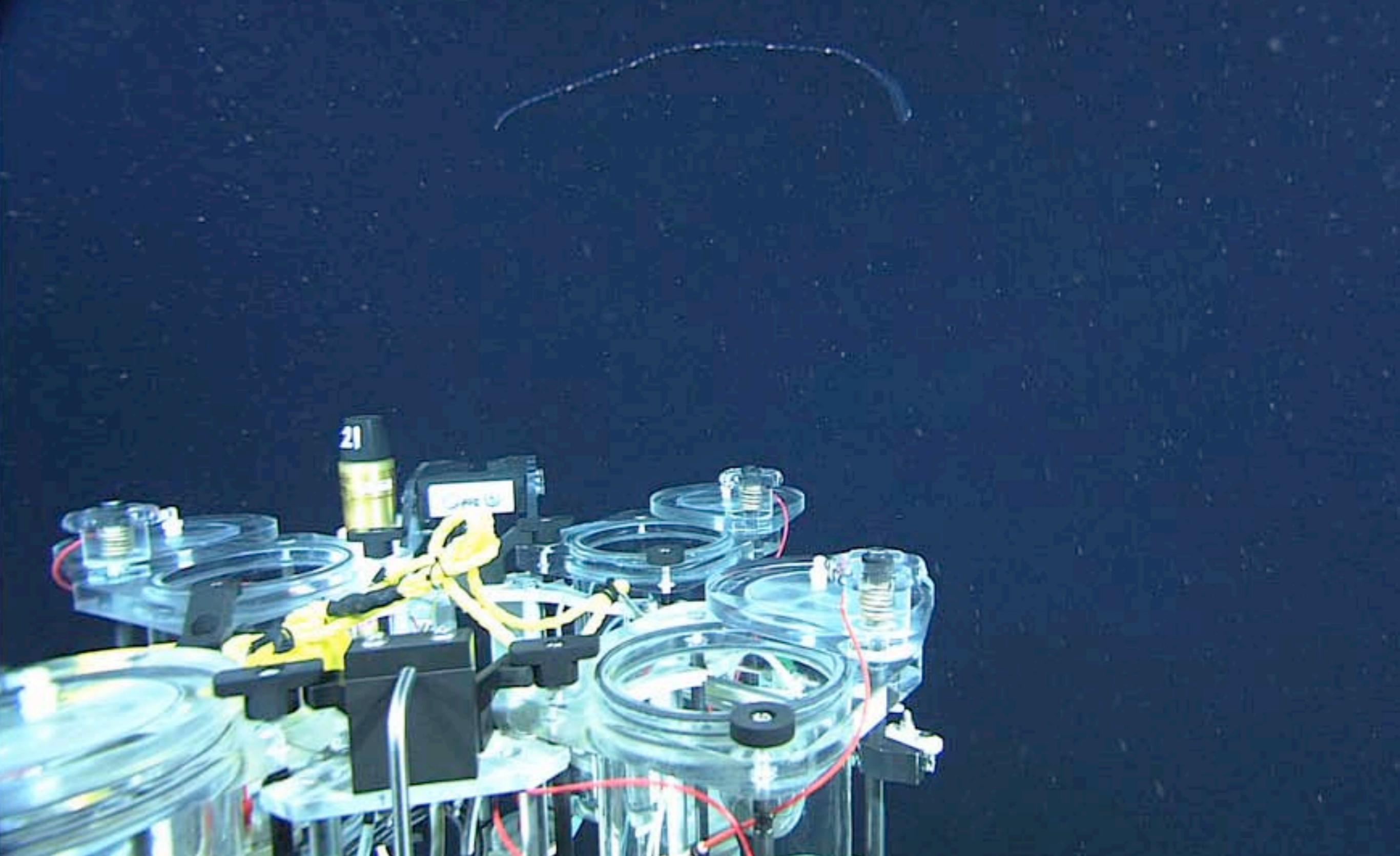
Casey Dunn



(C Carré)

Nanomia bijuga

Casey Dunn



(MBARI)

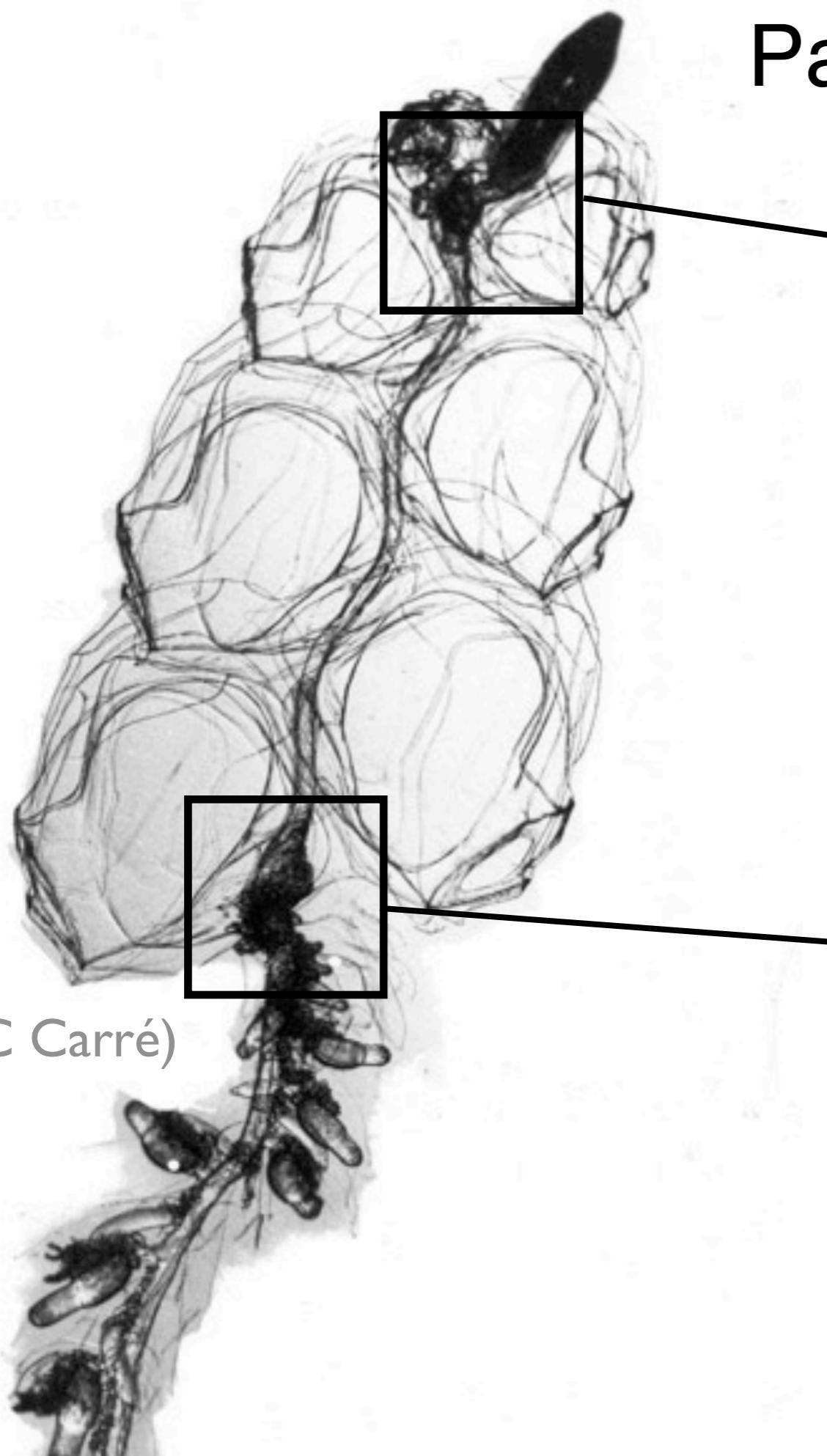
Casey Dunn

Nanomia 454 sequencing

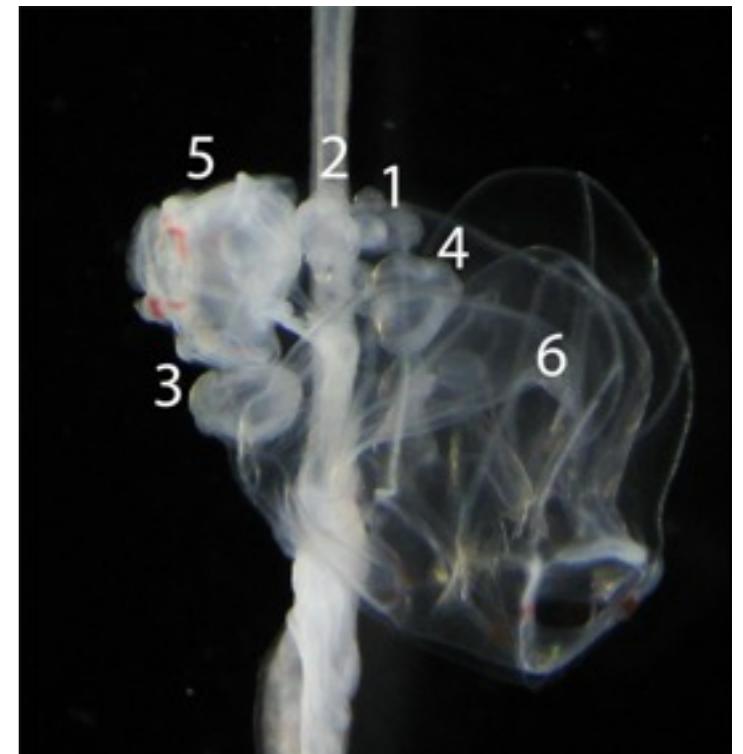
589k reads sequenced
(454 Titanium)

19,925 “genes” in reference
(Newbler, cap3)

Paired samples, 3 specimens



Swimming



Feeding



Casey Dunn

Replicated design

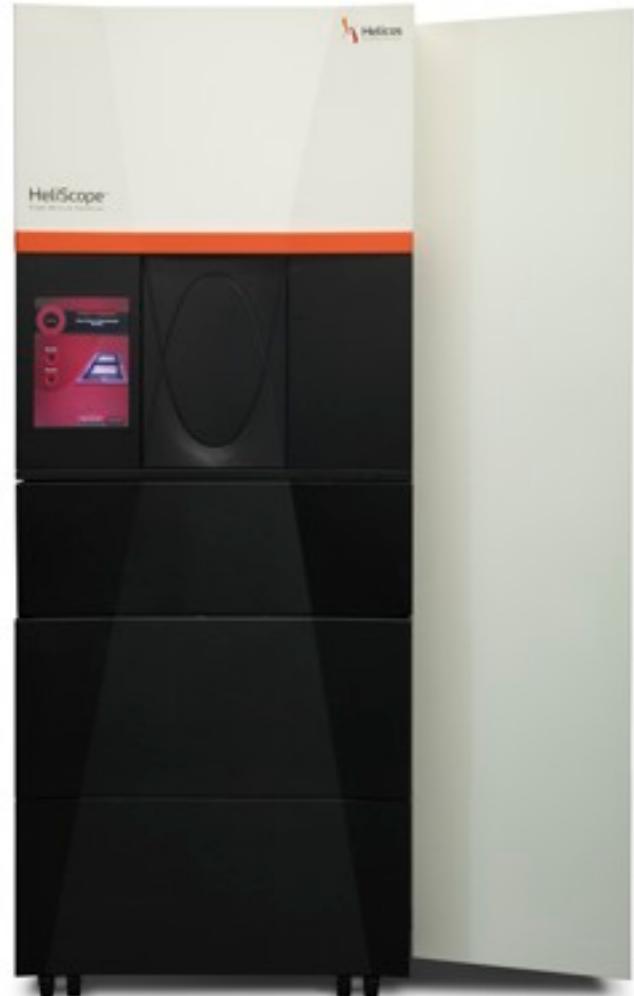
	Tissue A	Tissue B
Specimen 1	X Reads	X Reads
Specimen 2	X Reads	X Reads
Specimen 3	X Reads	X Reads



Casey Dunn



Casey Dunn



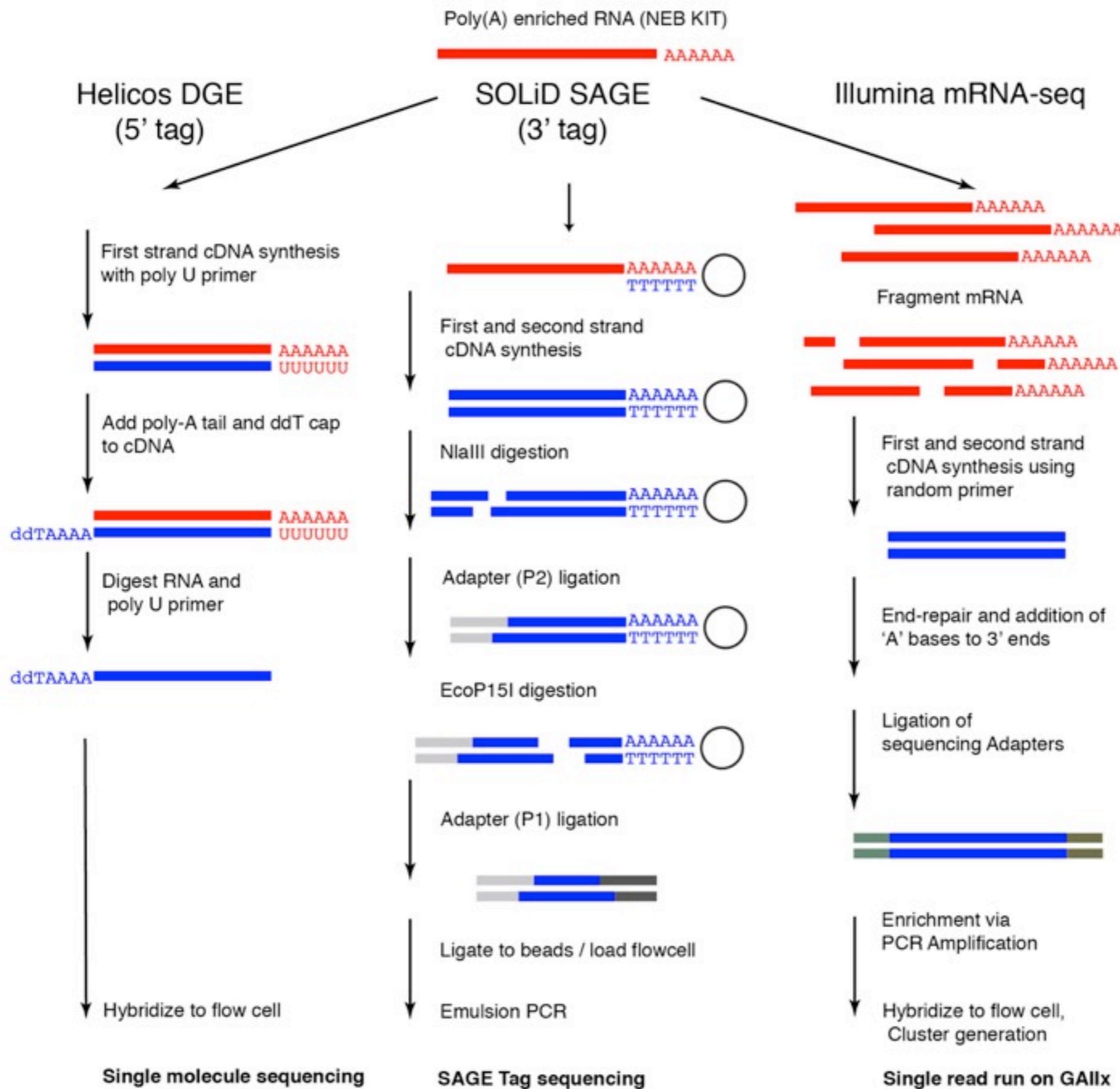
Helicos



SOLiD



Illumina

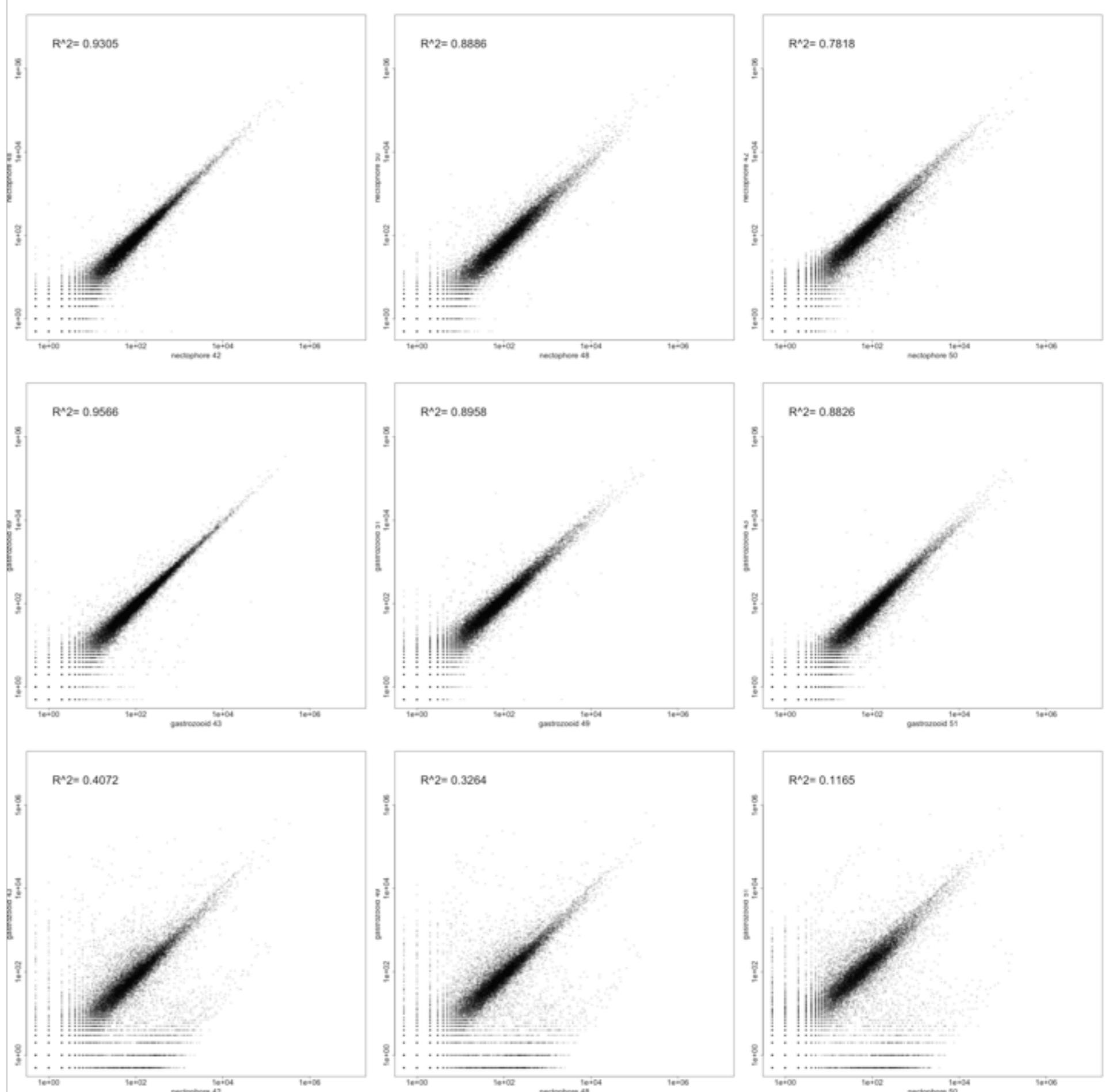


Illumina

Swimming

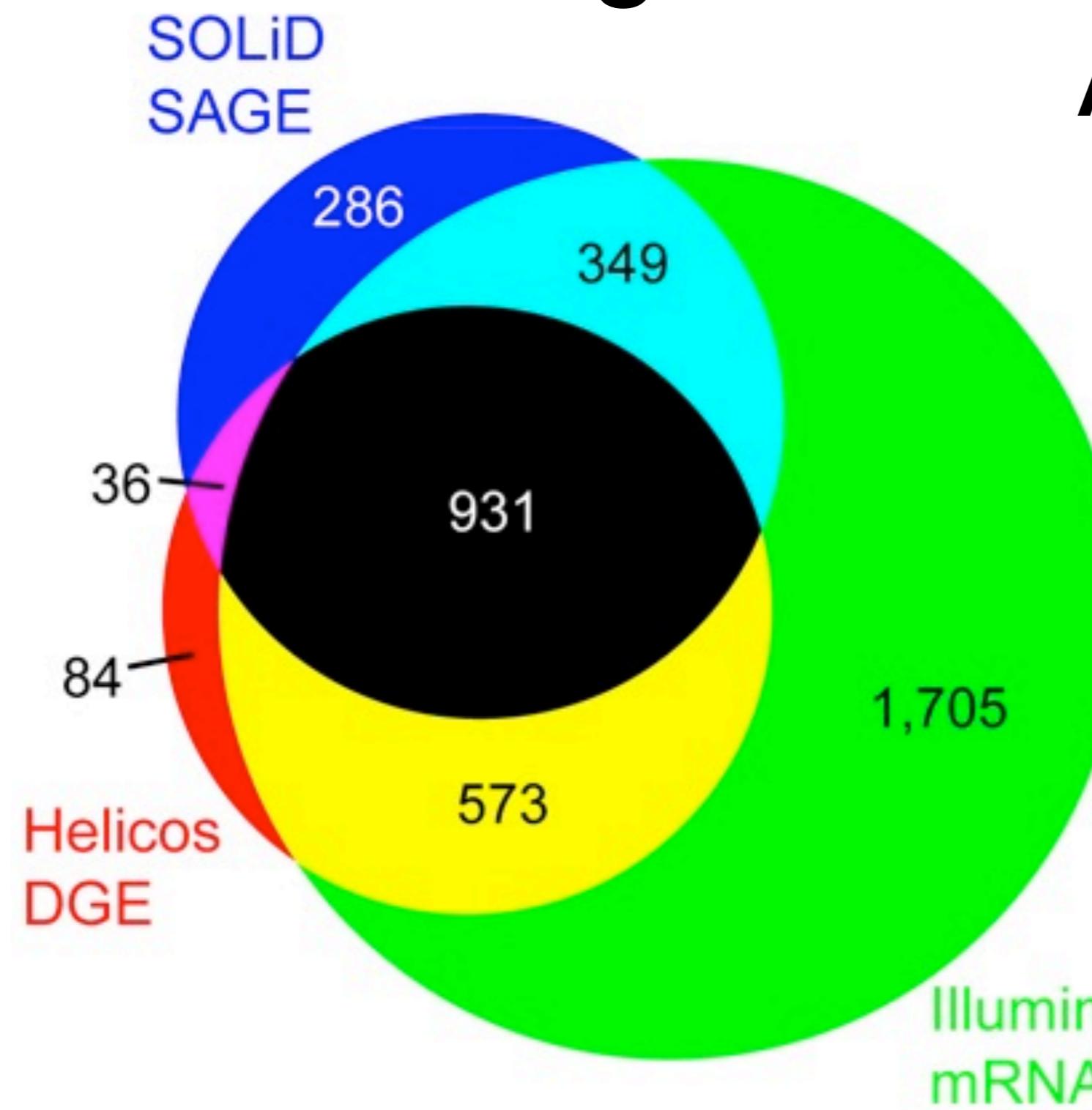
Feeding

Swimming v. Feeding



Genes with significant DE

All genes

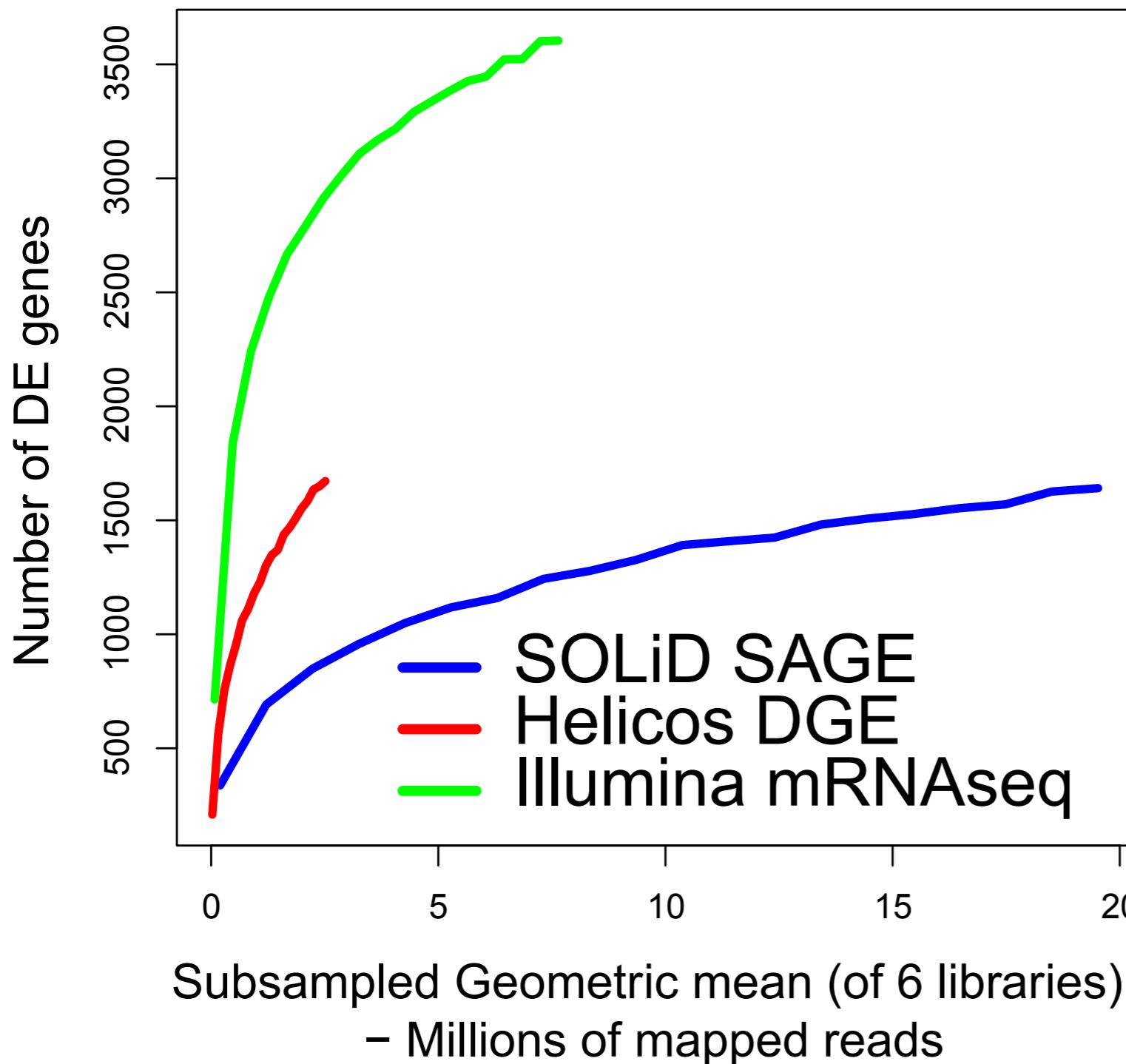


EdgeR, Bonferroni corrected $p < 0.05$

(dx.doi.org/10.1371/journal.pone.0022953)

Casey Dunn

Are these differences due to differences in read numbers across workflows?



No.

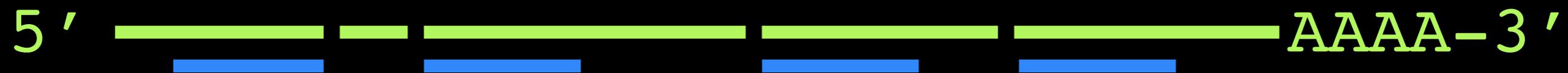
Helicos (Tag-based)



SOLiD (Tag-based)



Illumina (RNAseq)



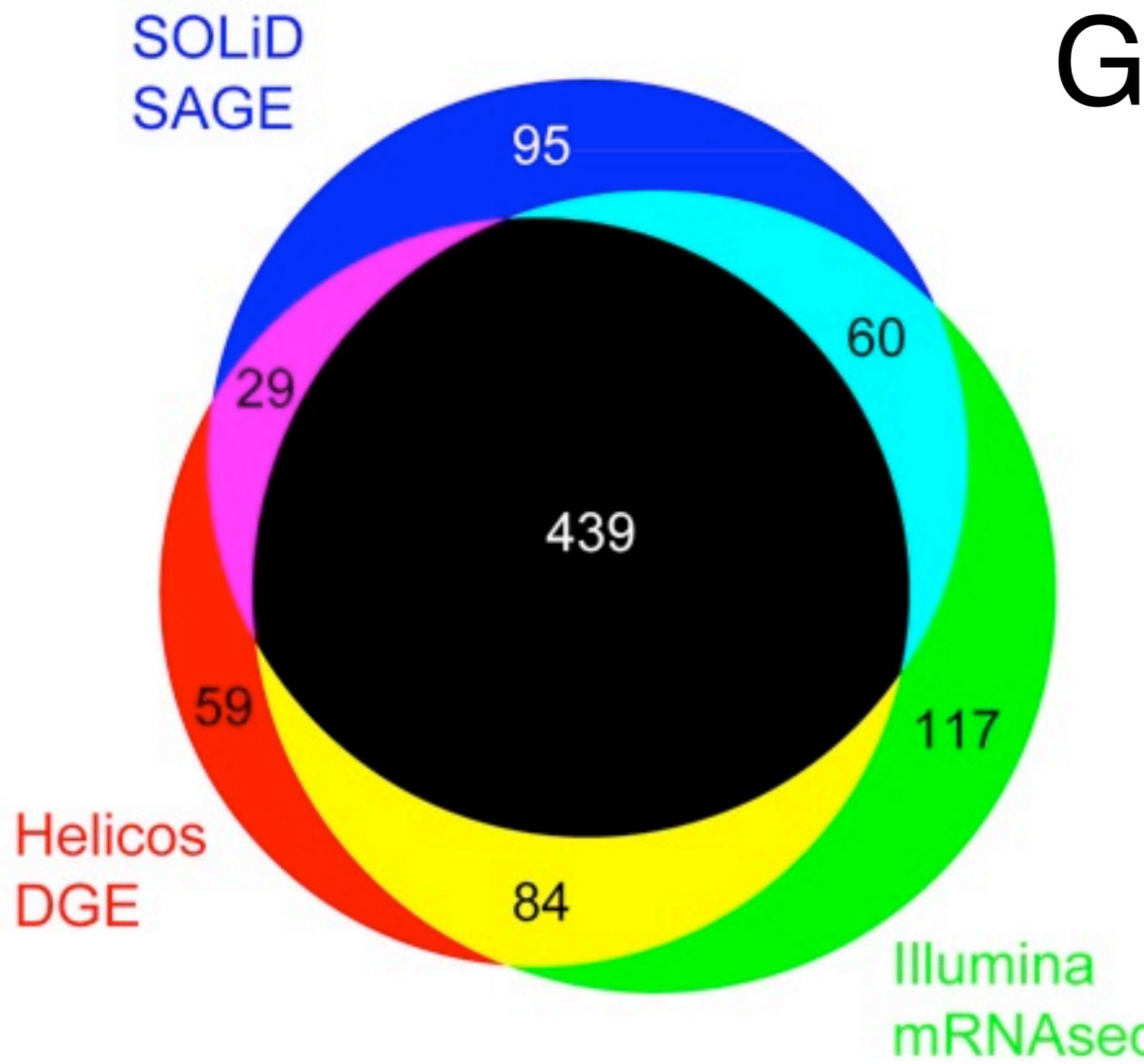
Original 454-derived reference:

19,925 genes

Reference sequences that
unambiguously have tag sites:

4,255 genes

Genes with significant DE



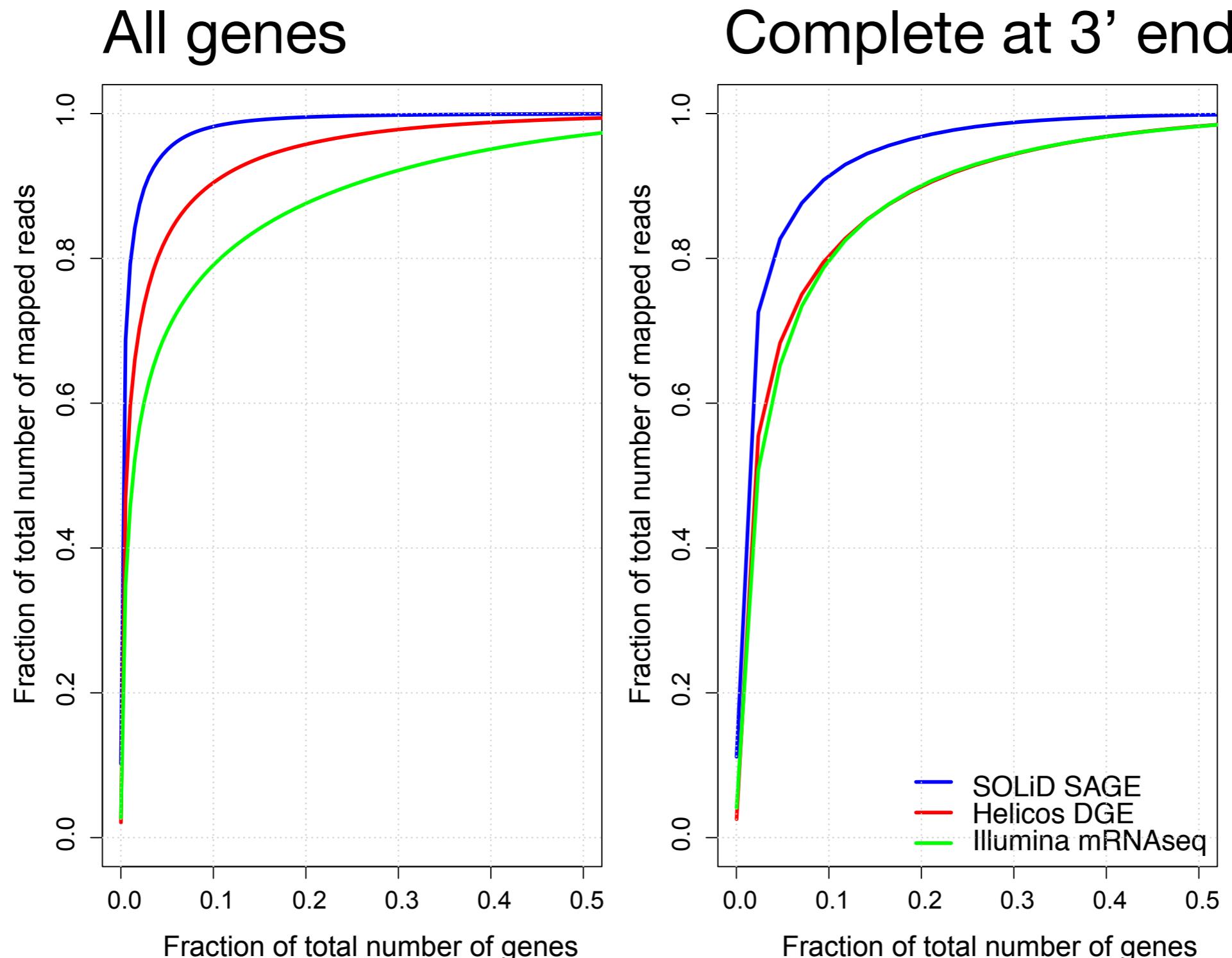
Genes with
complete
3' end

EdgeR, Bonferroni corrected $p < 0.05$

(dx.doi.org/10.1371/journal.pone.0022953)

Casey Dunn

Distribution of sequencing effort



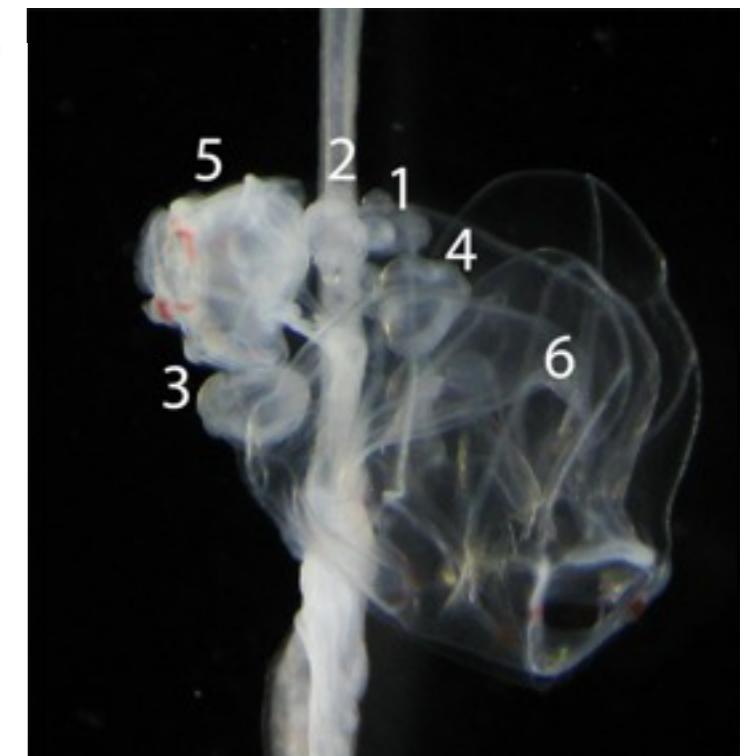
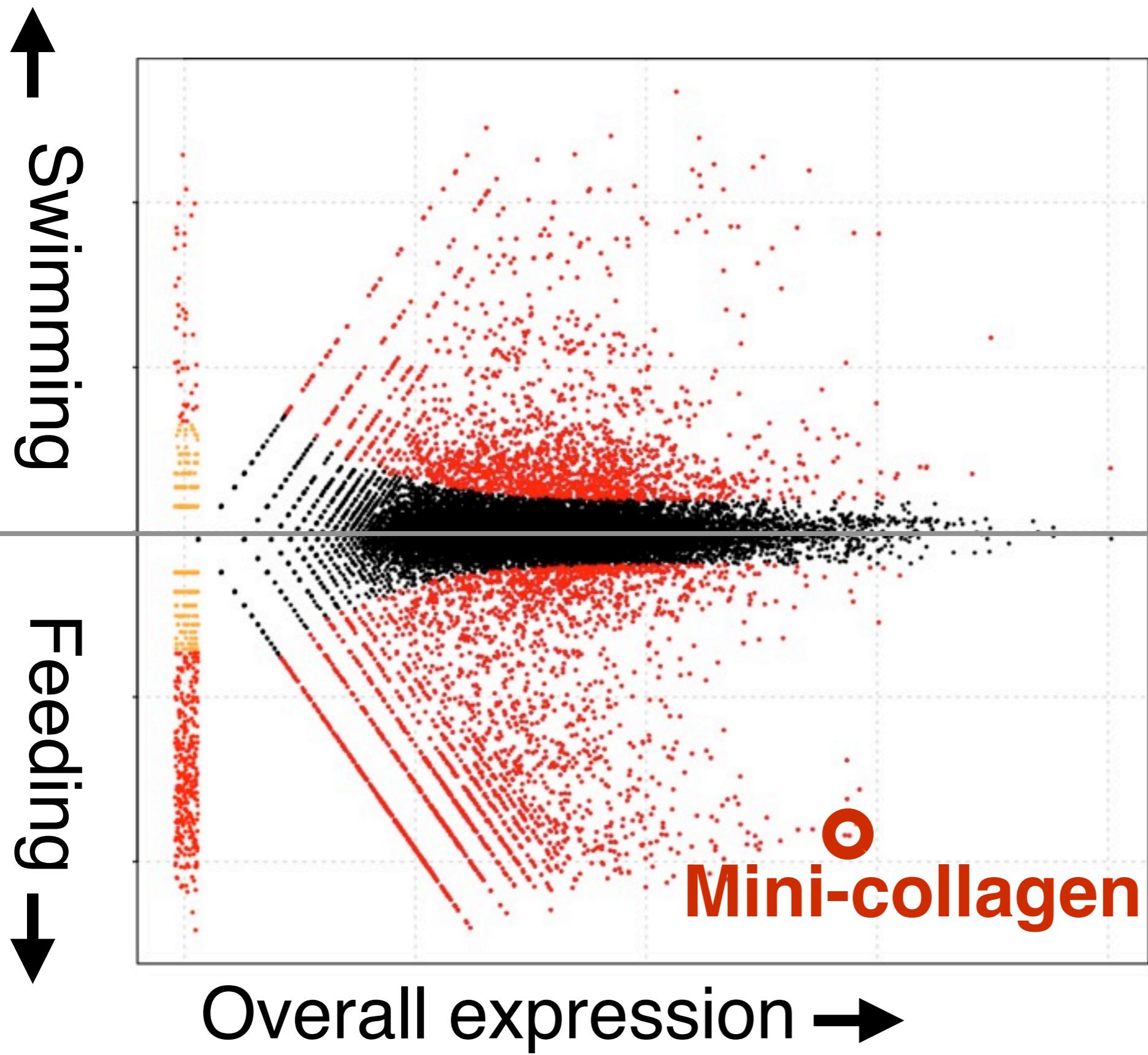
(dx.doi.org/10.1371/journal.pone.0022953)

Casey Dunn

Where to next?

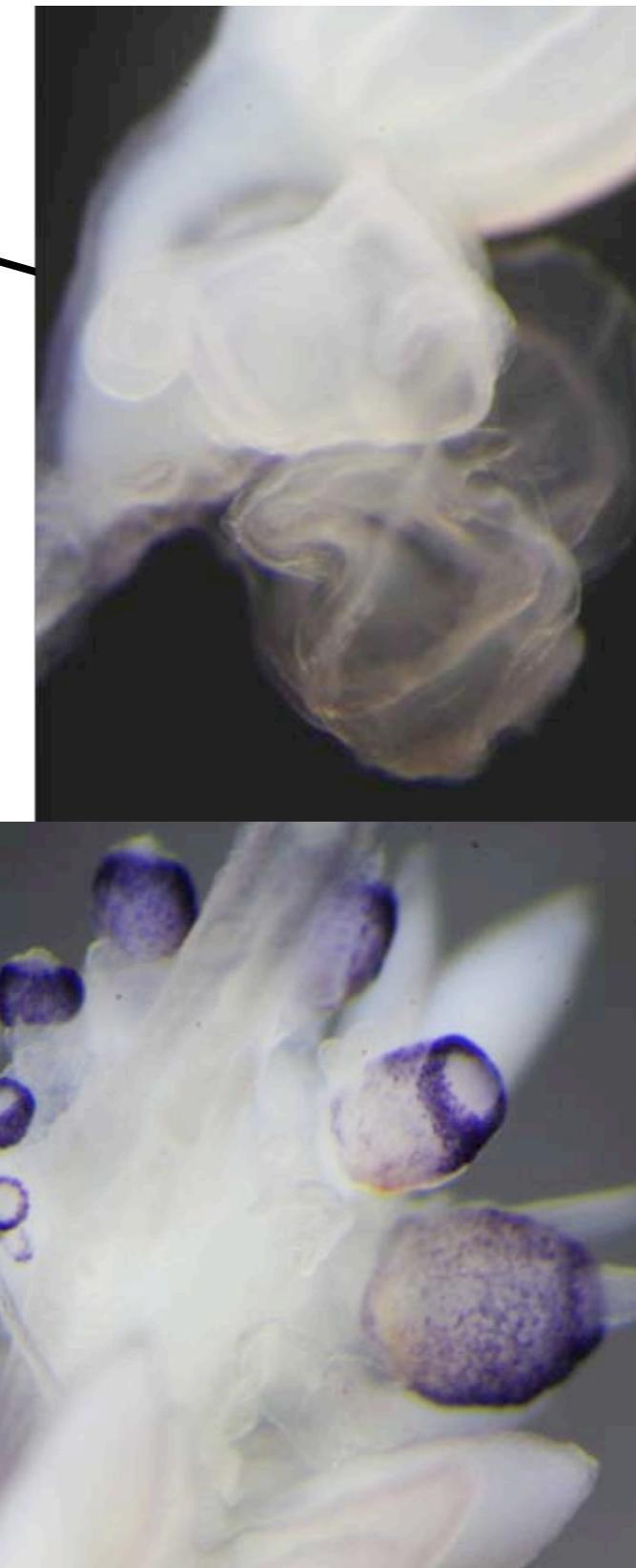
Characterization of genes with
significant differential
expression

Red genes have significant differential expression





(C Carré)



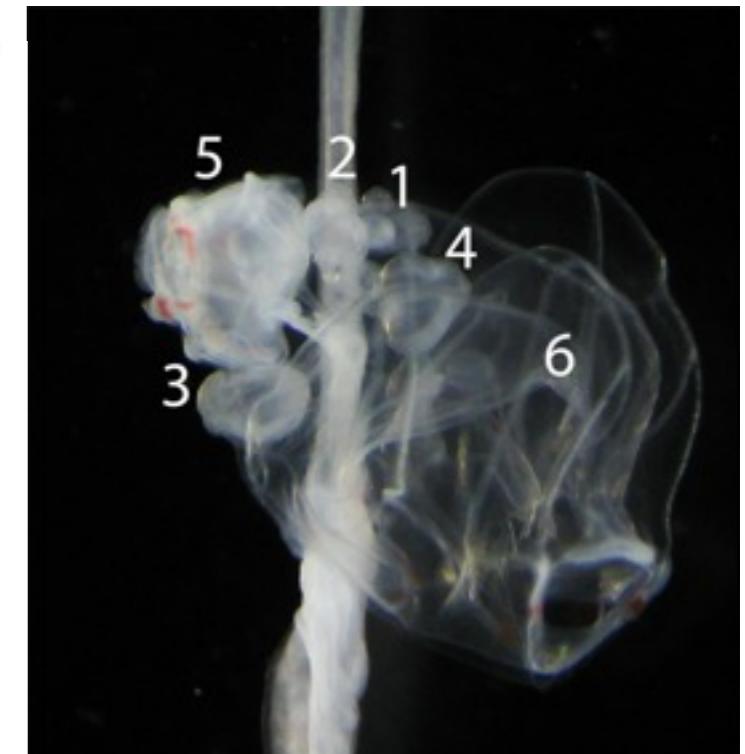
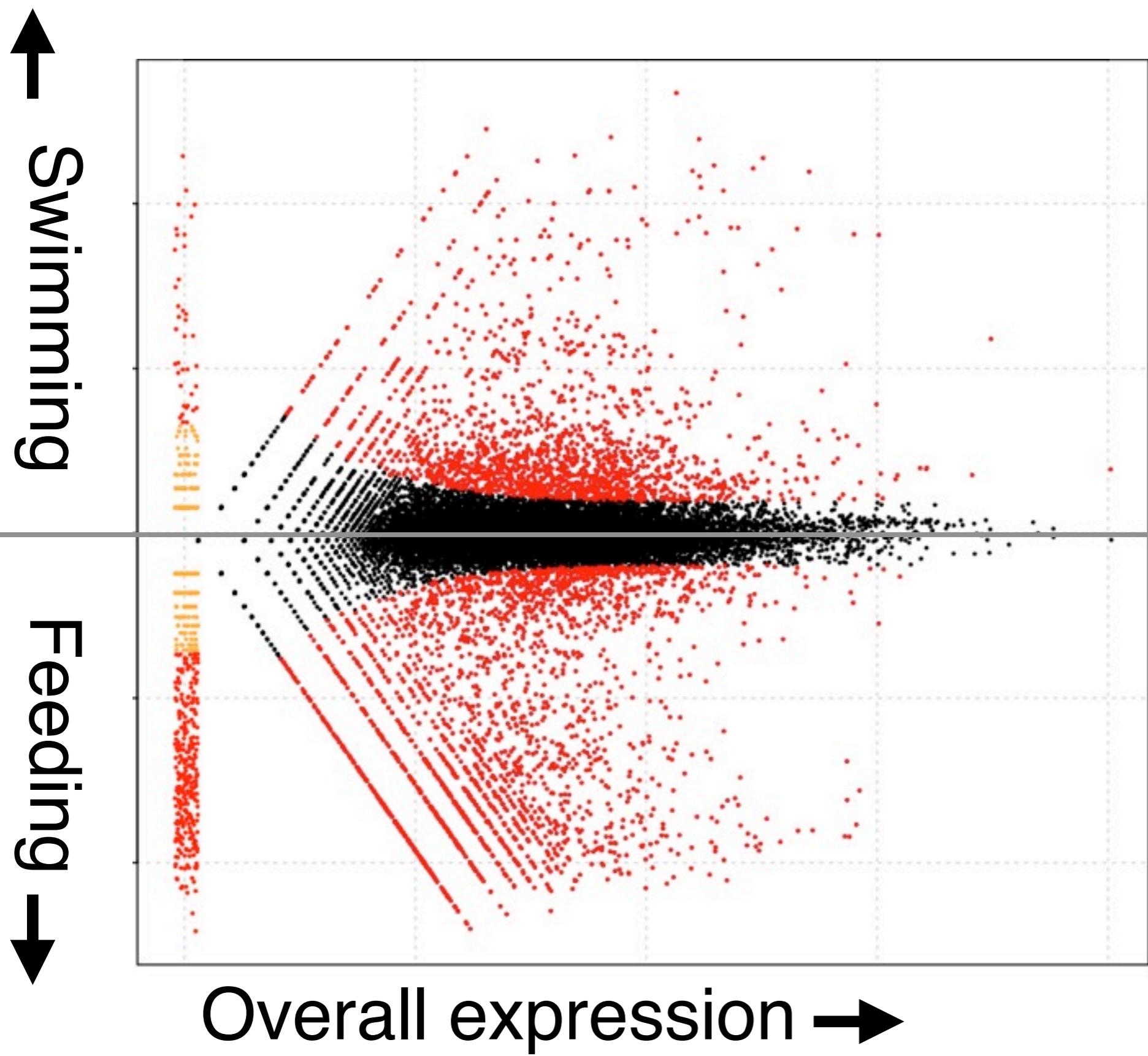
swimming
bodies

feeding
bodies

Cells expressing mini-collagen are **blue**

Casey Dunn

Red genes have significant differential expression



Uh oh.

“Data deluge”

“Firehose of data”

“I’m drowning in data.”

“Data overload”

The problem isn't too
much data.

We need more data that
tell us about our data

What other data do we need?

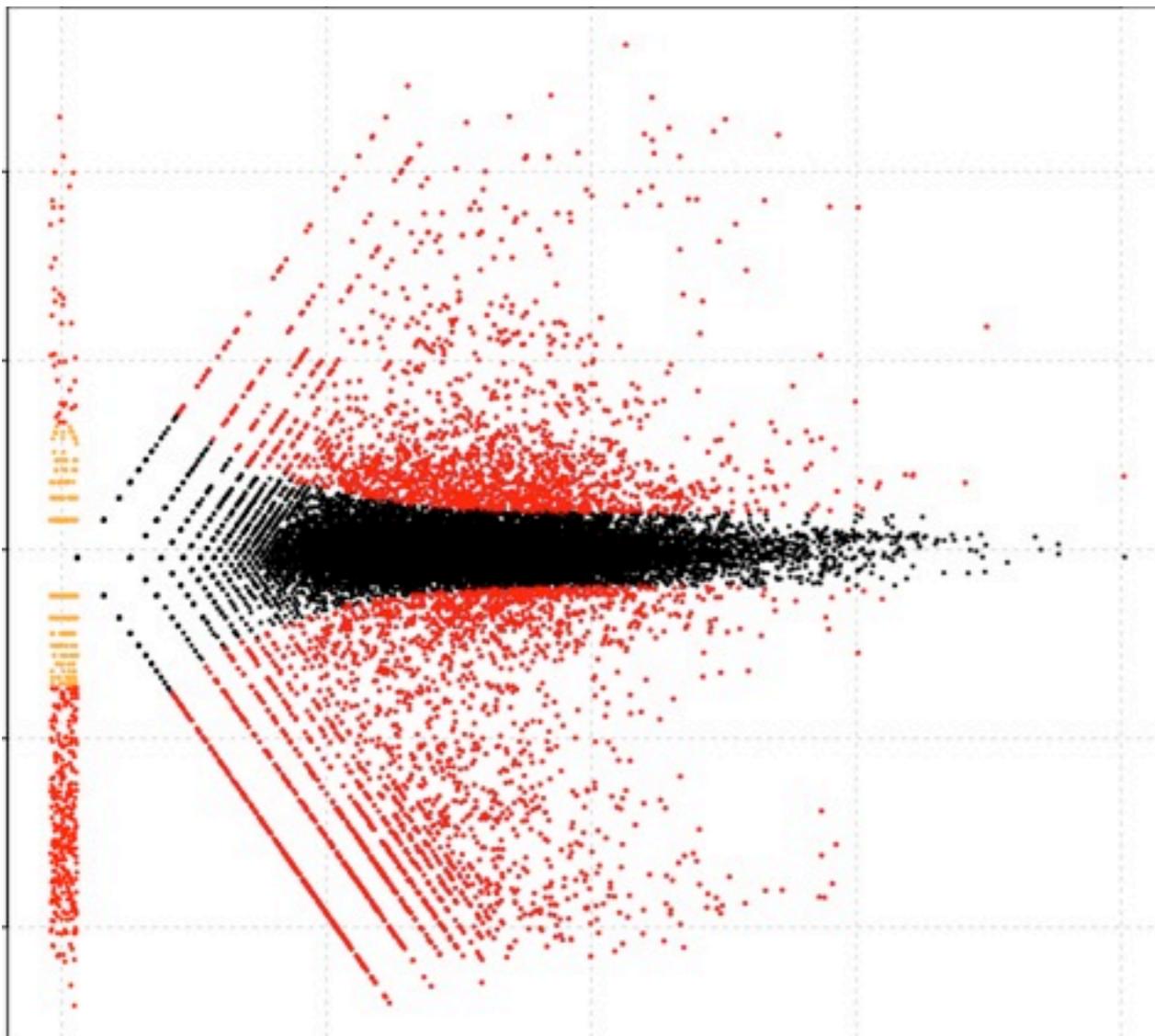
Comparative data - we need to be looking at a lot more than one species at a time.

Current approach:
Which genes have expression correlated with my phenotype of interest?

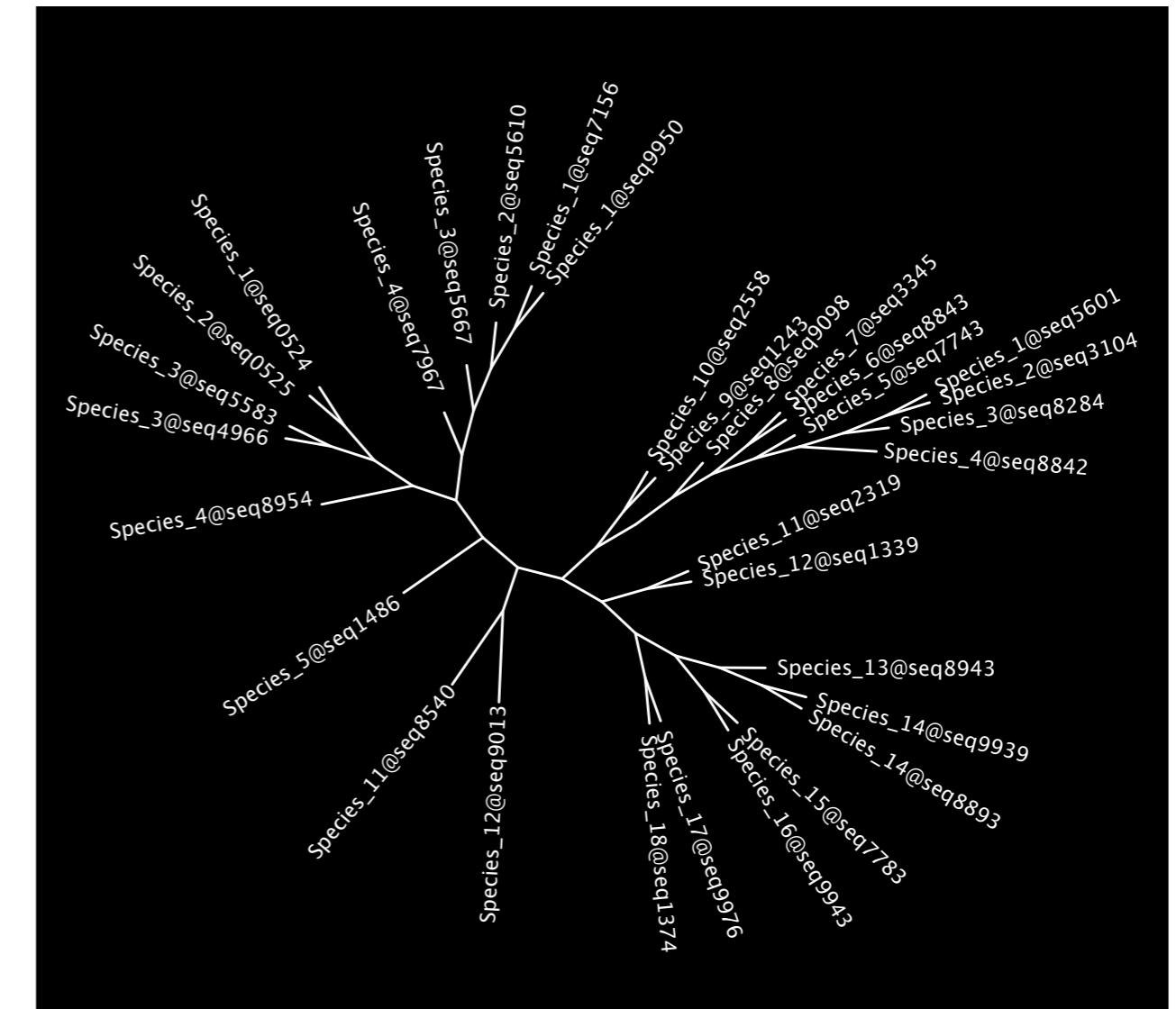
New approach:
Which genes have evolutionary changes in expression that are coincident with changes in my phenotype of interest?

Analyze expression data on phylogenies

Expression data



Gene trees



The state of comparative biology

New tools are going to transform
comparative biology

But the biggest impact will be how
they enable comparative biology
to transform the rest of biology

Mechanisms and diversity

20th Century

Experimental
work in model
organism



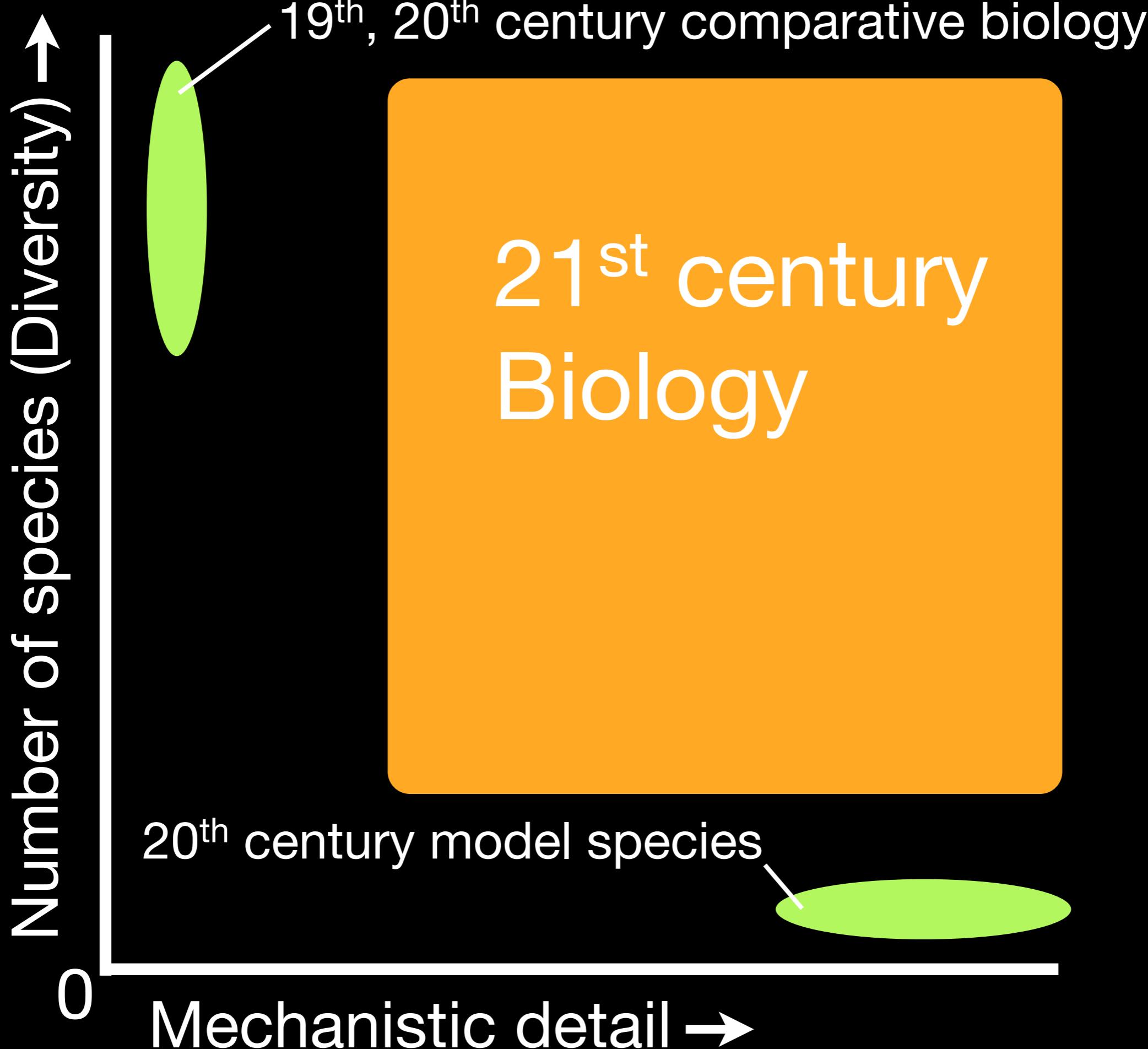
Comparative
work in
nonmodel
organisms

21st Century

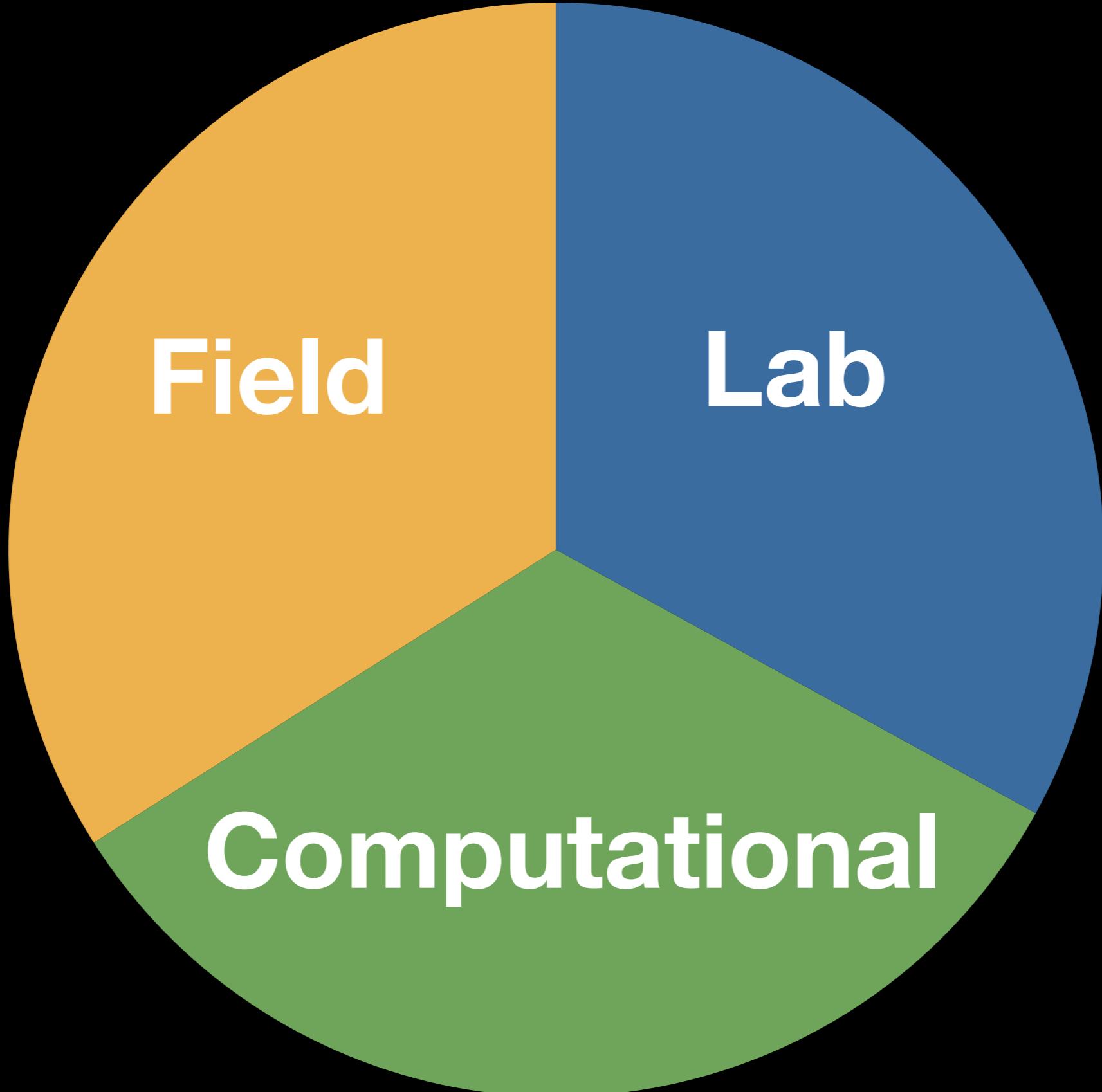
Comparative
work in
nonmodel
organisms



Experimental
work in model
organism



Computation





practical computing for biologists

Steven H. D. Haddock

The Monterey Bay Aquarium Research Institute,
and University of California, Santa Cruz

Casey W. Dunn

Department of Ecology and Evolutionary Biology,
Brown University



Sinauer
Associates, Inc.

Casey Dunn

goals

To show you how to use general tools to address the day-to-day computational challenges faced by biologists.

goals

We focus on the entire computer as a **general analysis environment**, rather than focus on one particular type of analysis or analysis tool.

The material we present can be thought of as **glue** to hold together and **automate** your existing analysis tools, and as a **general purpose workbench** for creating new analysis tools.

goals

Will we show you how to **convert file formats** so that you can use the output of one program as input to another?

Yes.

Will we show you how to use command-line tools to **automate a series of analyses**?

Yes.

Will we show you how to **use a remote cluster** to run programs?

Yes.

goals

Will we show you how to optimize an algorithm to speed it up ten fold?

No.

Will we explain maximum likelihood?

No.

Will we walk you through microarray data analysis?

No.

contents

PART I: Text Files 7

PART II: The Shell 45

PART III: Programming 103

PART IV: Combining Methods 243

PART V: Graphics 321

PART VI: Advanced Topics 381

Appendices 449

sample content

```
host:~ lucy$ cd ~/pcfbsandbox
host:sandbox lucy$ ls
host:sandbox lucy$ cp ../examples/*.txt ./
host:sandbox lucy$ ls
reflist.txt ← You will see others too...
```

```
#!/usr/bin/env python

DNASeq = 'ATGAAC'
print 'Sequence:', DNASeq

SeqLength = float(len(DNASeq))

print 'Sequence Length:', SeqLength

NumberA = DNASeq.count('A')
NumberC = DNASeq.count('C')
NumberG = DNASeq.count('G')
NumberT = DNASeq.count('T')

print 'A:', NumberA/SeqLength
print 'C:', NumberC/SeqLength
print 'G:', NumberG/SeqLength
print 'T:', NumberT/SeqLength
```

sample content

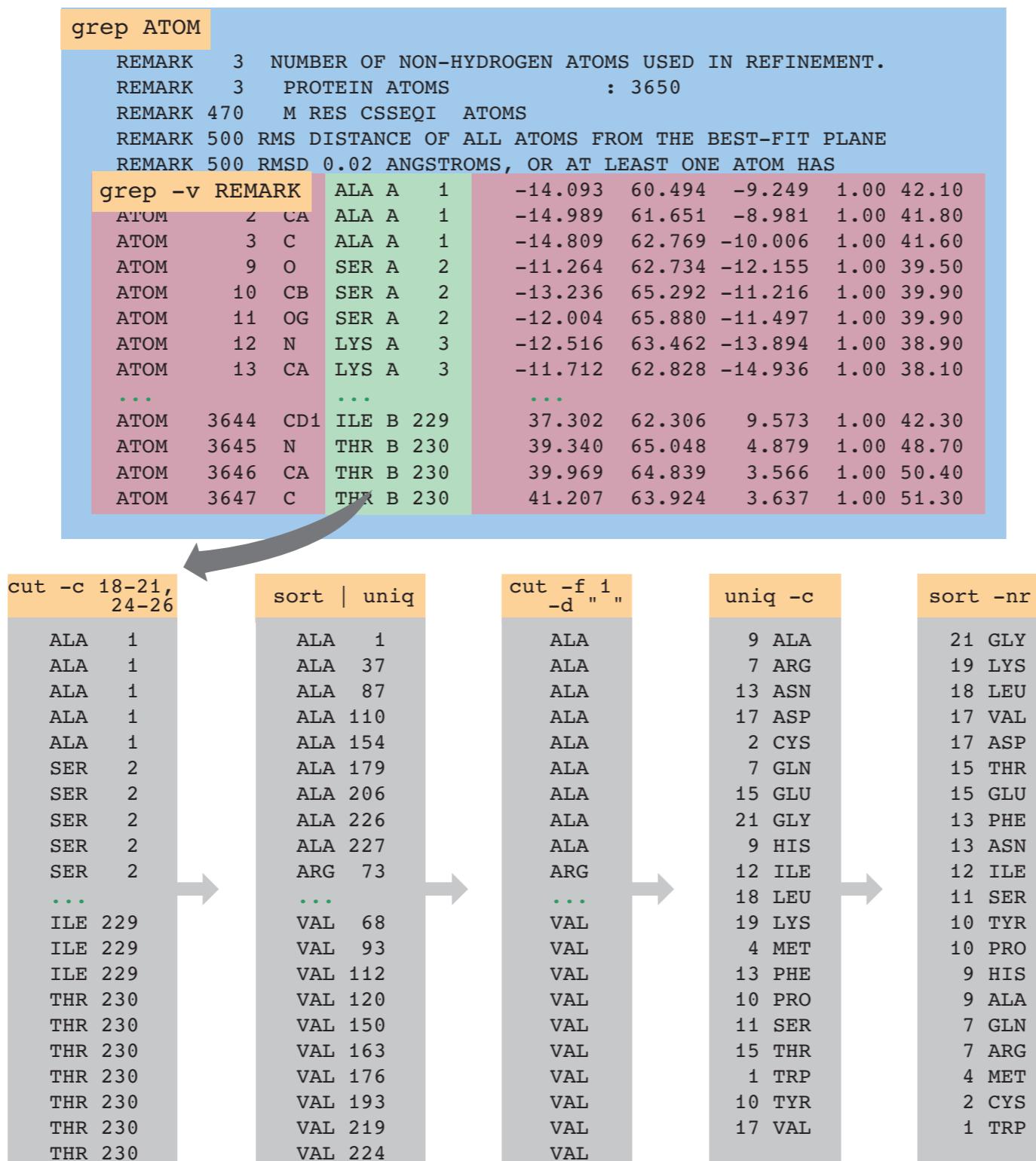
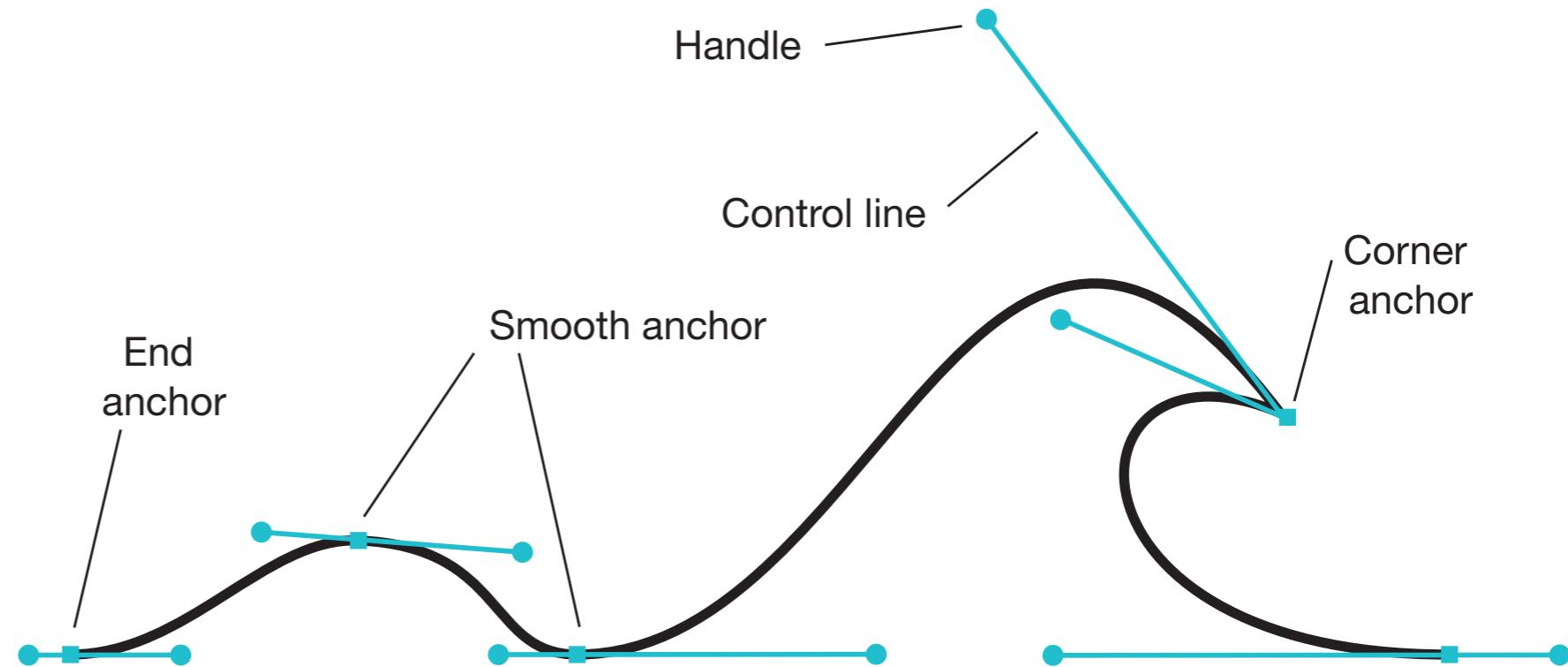


FIGURE 16.1 The successive extractions and modifications made by each command in the example pipeline. Orange boxes show bash commands and other boxes show the output once those commands have been added to the pipeline.

sample content

FIGURE 18.2 Bézier curve showing anchor points, handles, and control lines





Casey Dunn