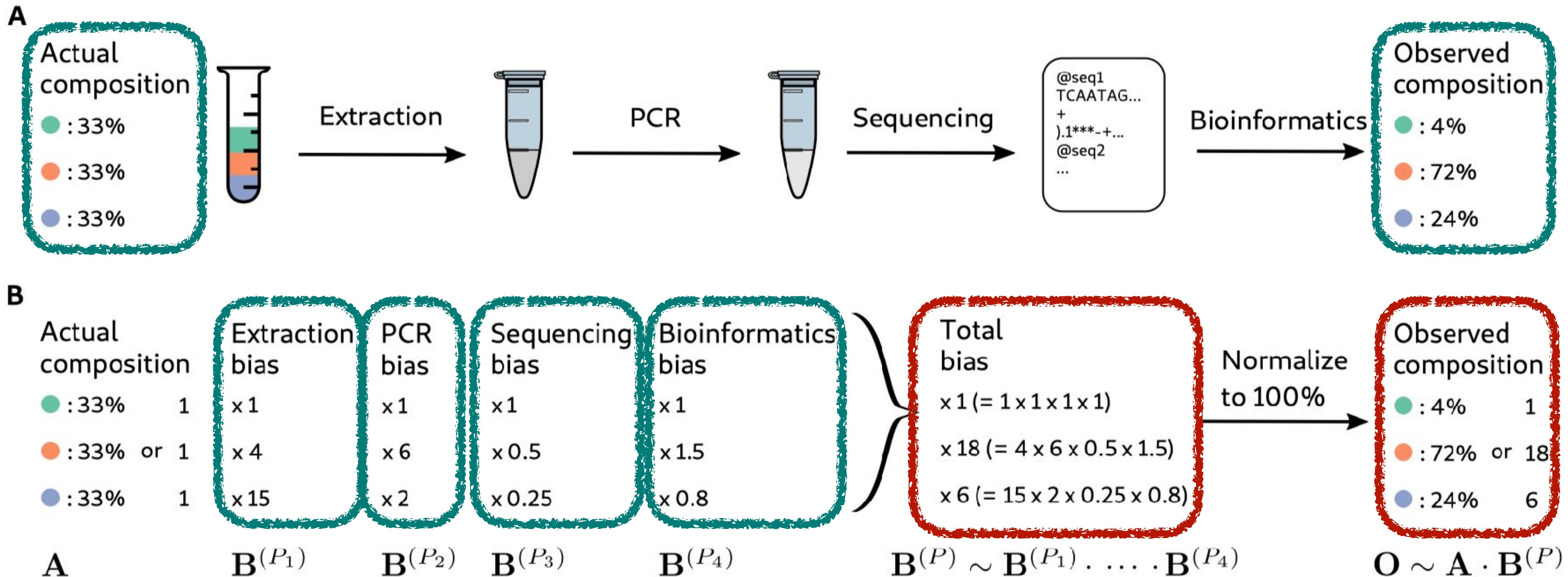


Estimating abundance from DNA-based data

Douglas Yu

Across-species quantification is difficult



eLife 8:e46923.

Consistent and correctable bias in metagenomic sequencing experiments

Michael R McLaren¹, Amy D Willis², Benjamin J Callahan^{1,3*}

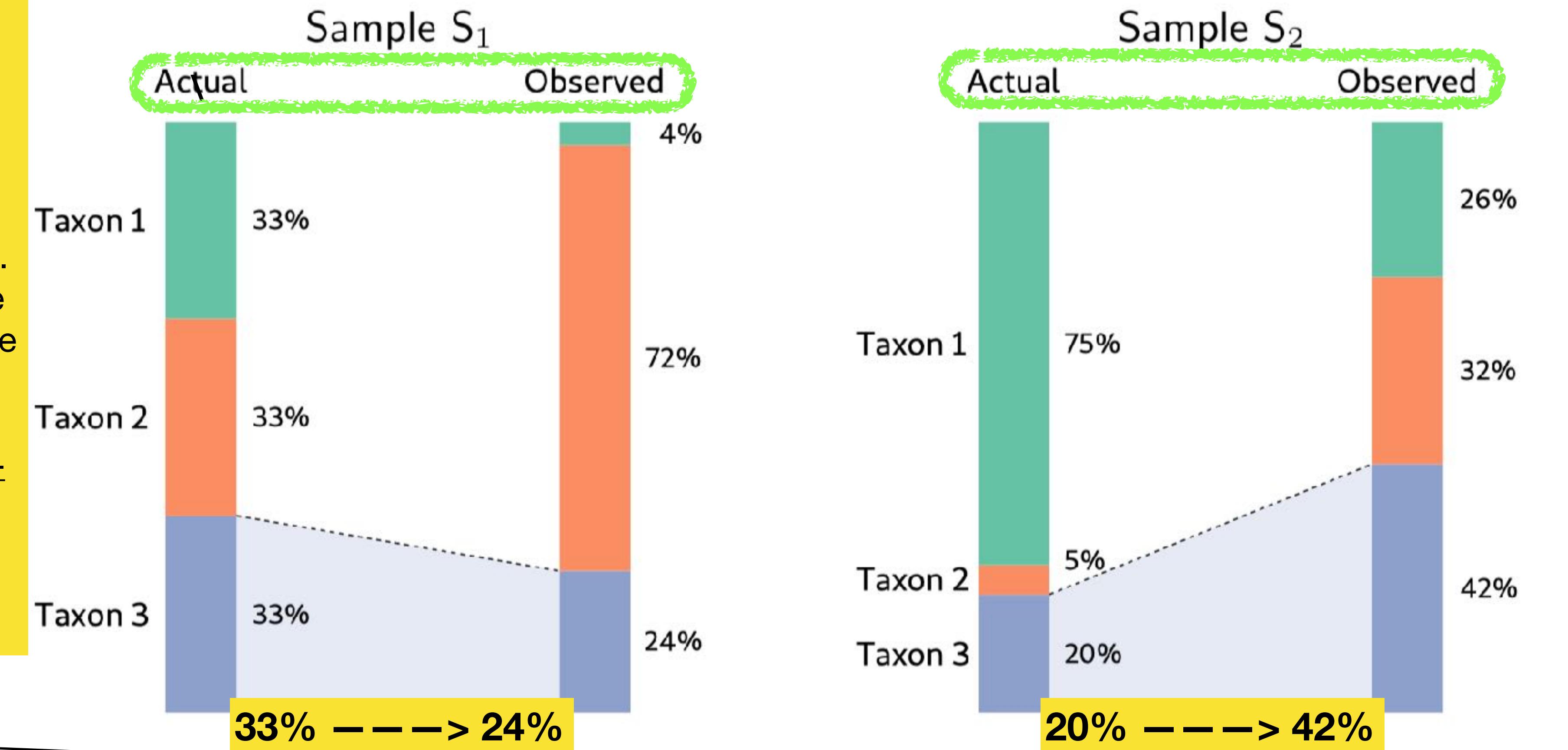
Across-species quantification is really difficult!

Even with the same underlying biases, it is possible for a given species to become more or less abundant in a sample after metabarcoding.

Imagine this is a diet study. You only see the Observed proportions. Your conclusions about the relative abundances of prey items would be incorrect.

Diet studies try to achieve “across-species” quantification:

“Taxon 3 has higher/lower relative abundance than Taxon 2”.



$$\begin{aligned}
 \text{Total bias} &= x_1 (= 1 \times 1 \times 1 \times 1) \\
 &\quad \times 18 (= 4 \times 6 \times 0.5 \times 1.5) \\
 &\quad \times 6 (= 15 \times 2 \times 0.25 \times 0.8) \\
 B^{(P)} &\sim B^{(P_1)} \dots B^{(P_4)}
 \end{aligned}$$

Ratios to Taxon 1

$$\begin{aligned}
 \text{Actual} &= \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \\
 \text{B} &= \begin{pmatrix} 1 \\ 18 \\ 6 \end{pmatrix} \\
 \text{A}(S_1) &= \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 18 \\ 6 \end{pmatrix} = \begin{pmatrix} 1 \\ 18 \\ 6 \end{pmatrix} \\
 \text{O}(S_1) &\sim \begin{pmatrix} 0.04 \\ 0.72 \\ 0.24 \end{pmatrix}
 \end{aligned}$$

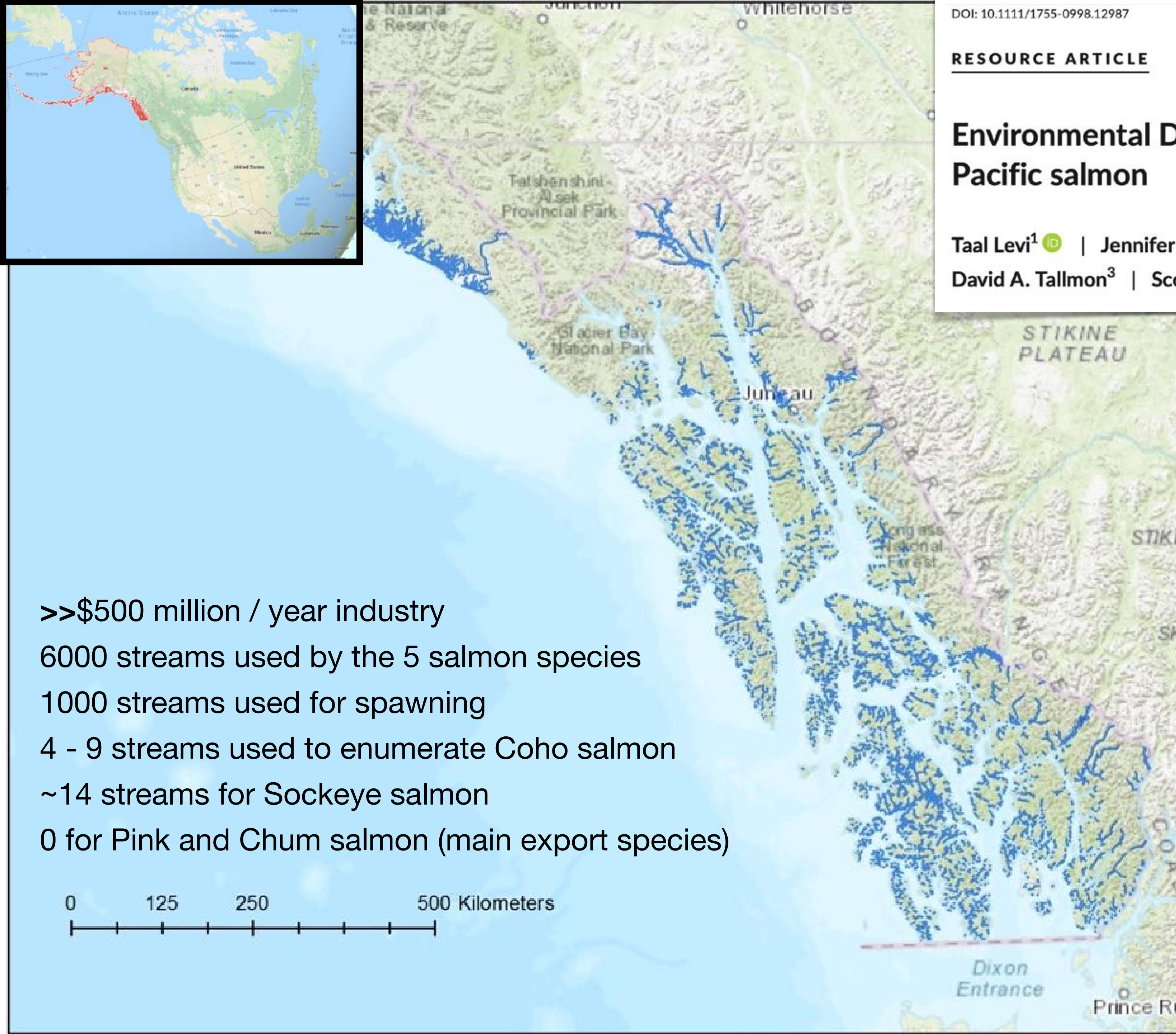
$$\begin{aligned}
 \text{Actual} &= \begin{pmatrix} 1 \\ 1/15 \\ 4/15 \end{pmatrix} \\
 \text{B} &= \begin{pmatrix} 1 \\ 18 \\ 6 \end{pmatrix} \\
 \text{A}(S_2) &= \begin{pmatrix} 1 \\ 1/15 \\ 4/15 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 18 \\ 6 \end{pmatrix} = \begin{pmatrix} 18/15 \\ 24/15 \end{pmatrix} \\
 \text{O}(S_2) &\sim \begin{pmatrix} 0.26 \\ 0.32 \\ 0.42 \end{pmatrix}
 \end{aligned}$$

Methods to extract abundance information from DNA data

- Single-species quantitative PCR (qPCR)
- Multiplexed specimen barcoding (mBRAVE)
- Mitogenomics and DNA spike-in (SPIKEPIPE)
- Metabarcoding and DNA spike-in (qSeq)
- Reverse metagenomics (RevMet)

Methods to extract abundance information from DNA data

- Single-species quantitative PCR (qPCR)
- Multiplexed specimen barcoding (mBRAVE)
- Mitogenomics and DNA spike-in (SPIKEPIPE)
- Metabarcoding and DNA spike-in (qSeq)
- Reverse metagenomics (RevMet)



Environmental DNA for the enumeration and management of Pacific salmon

Taal Levi¹ | Jennifer M. Allen¹ | Donovan Bell² | John Joyce² | Joshua R. Russell² | David A. Tallmon³ | Scott C. Vulstek² | Chunyan Yang⁴ | Douglas W. Yu^{4,5,6}

Alaskan salmon fishery



The ‘escapement’
is the breeding
population



Is the fishery allowing enough salmon to escape?

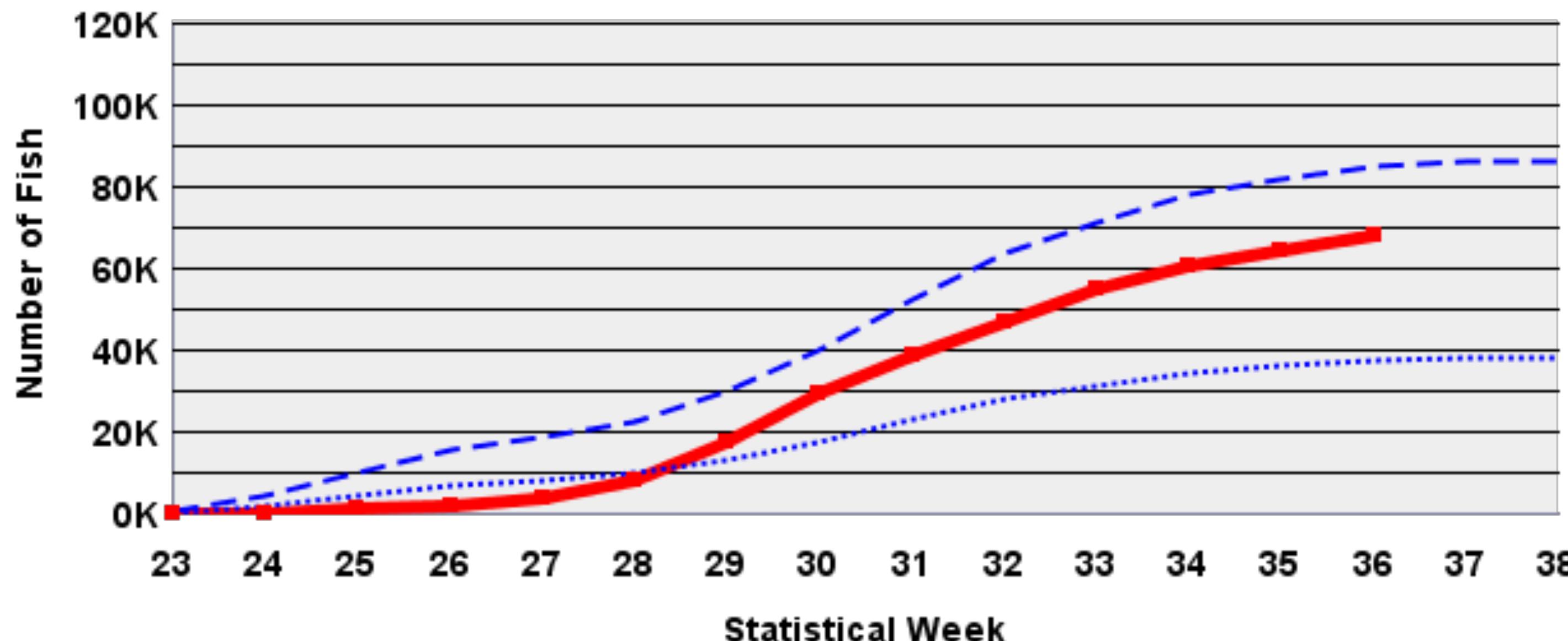
**Commercial Fishing**[Commercial Fishing Home](#)[News Releases &
Announcements](#)[Information By Area](#)

- Arctic-Yukon-Kuskokwim
 - Kuskokwim
 - Arctic Area
 - Norton Sound & Kotzebue
 - Yukon
- Central Region
 - Bristol Bay
 - Copper River
 - Cook Inlet
 - Lower Cook Inlet
 - Upper Cook Inlet
 - Prince William Sound
- Southeast Region
- Westward Region
 - Alaska Peninsula
 - Bering Sea & Aleutian Islands
 - Chignik
 - Kodiak Island

[ADF&G Home](#) » [Fishing](#) » [Commercial](#) » [Information By Area](#) » [Southeast](#)

Commercial Salmon Fisheries

Chilkoot Lake Weir - Sockeye Counts

Cumulative Counts Compared to Escapement Goals

ALASKA DEPARTMENT OF FISH AND GAME

DIVISION OF COMMERCIAL FISHERIES

NEWS RELEASE



*Sam Cotten, Commissioner
Jeff Regnart, Director*



Contact:

Gordie Woods

Phone: (907) 784-3255

Fax: (907) 784-3254

Yakutat Area Office

P.O. Box 49

Yakutat, AK, 99689

Date: September 4, 2015

Time: 8:30 a.m.

YAKUTAT COMMERCIAL SET GILLNET OPENING ANNOUNCEMENT

Dangerous River: will be open from 12:01 p.m., Sunday, September 6 through 12:00 noon, Wednesday, September 9.

Akwe River: will be open from 12:01 p.m., Sunday, September 6 through 12:00 noon, Wednesday, September 9.

Lost River: will remain closed until further notice.

Situk-Ahrnklin Inlet: will be open from 12:01 p.m., Sunday, September 6 through 12:00 noon, Wednesday, September 9 with the following restrictions:

North Bank of Situk-Ahrnklin Inlet: commercial set gillnet fishing will be prohibited along the north bank of the Situk-Ahrnklin Inlet between two ADF&G regulation markers located 500 yards above and 500 yards below the confluence of the Situk-Ahrnklin Inlet and the Lost River;



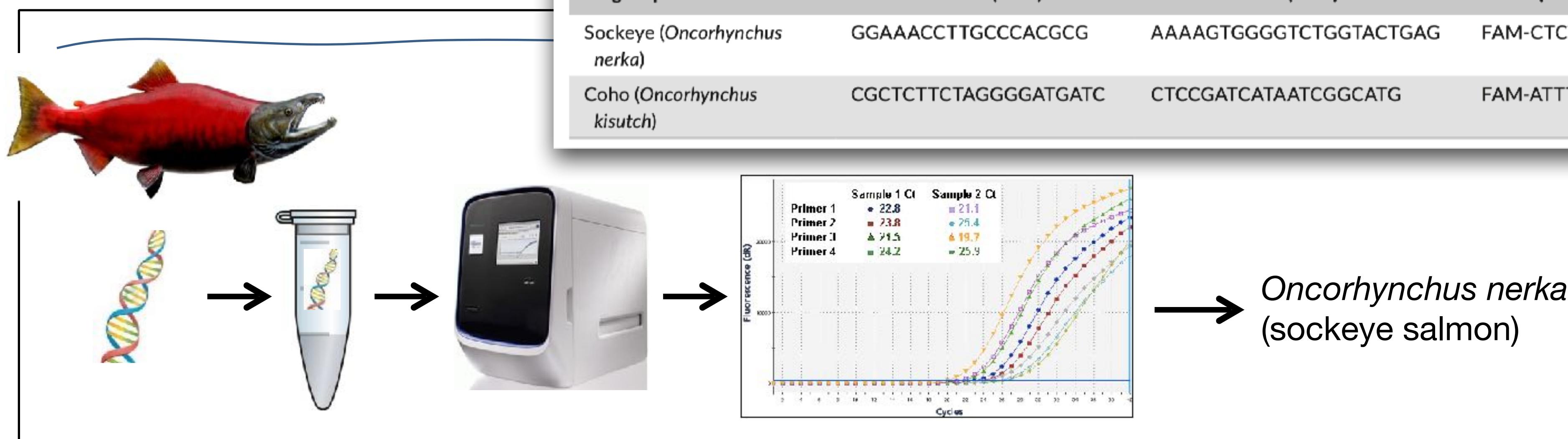
Validation data from Auke Creek Station in Juneau, Alaska

Daily counts of

- in-migrating Coho and Sockeye **adults**
- out-migrating Sockeye **smolts (juveniles)**

Does environmental DNA (eDNA) contain enough information to estimate salmon escapement sizes?

TABLE 1 Species-specific primers and probes used in this study (Rasmussen Hellberg et al., 2010)



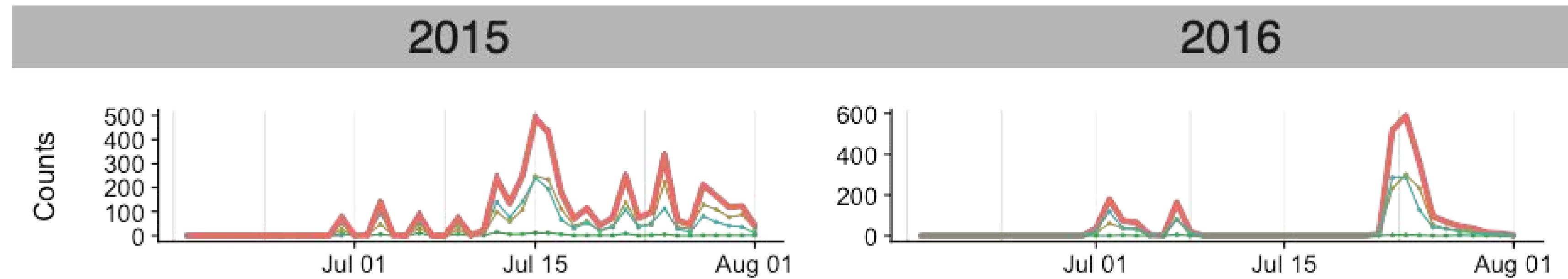
=



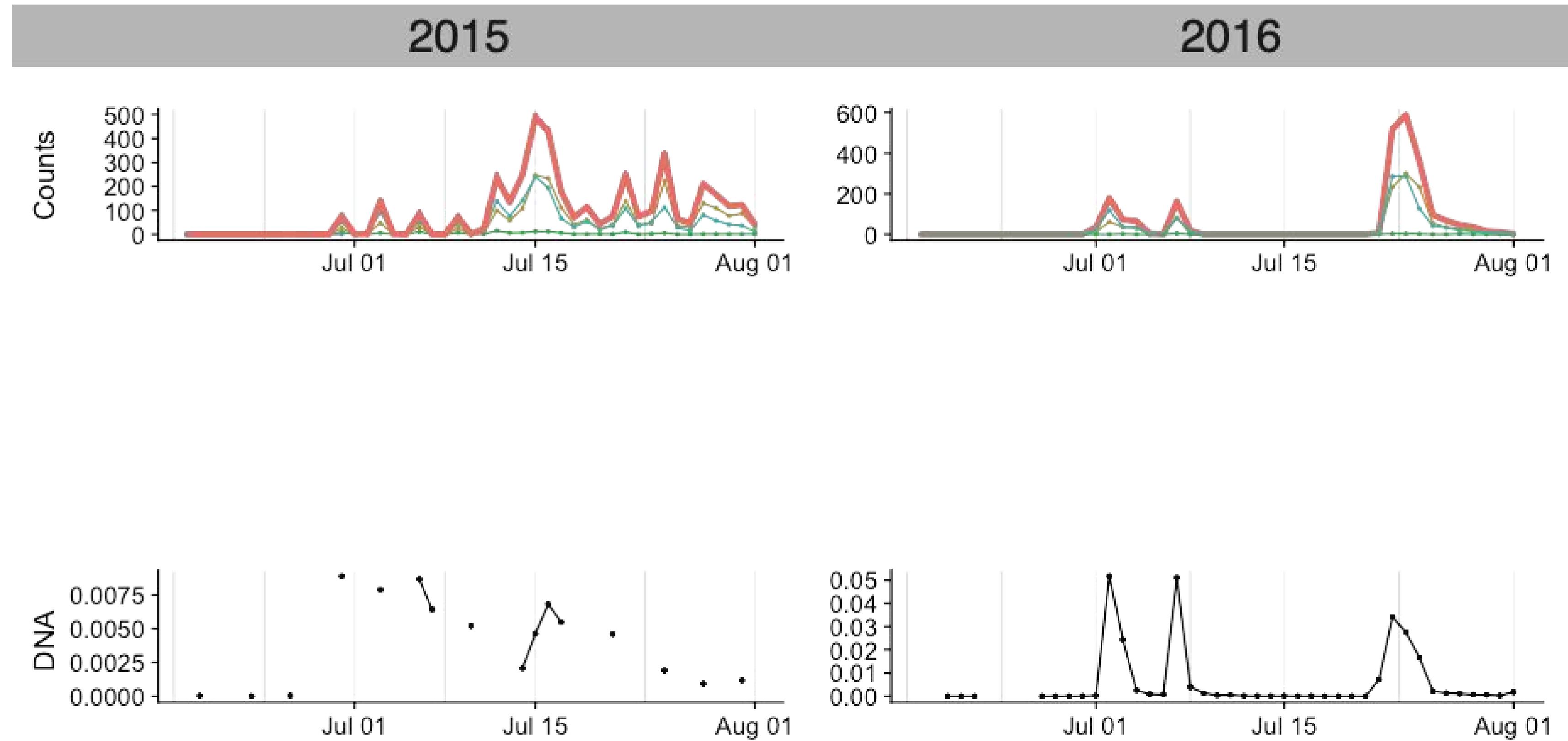
One person, one sample /day

Many people, all day

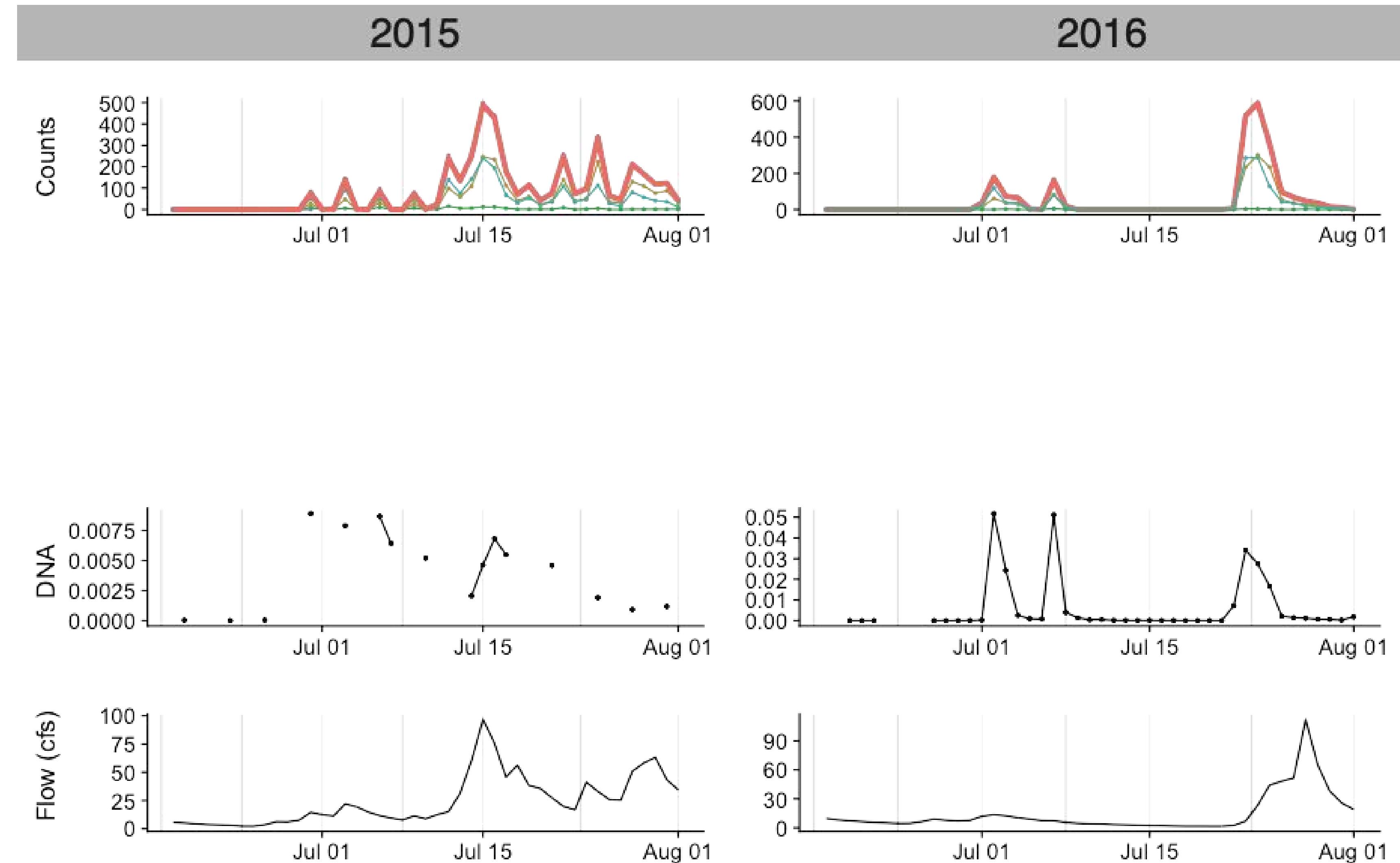
Sockeye in-migrating adults



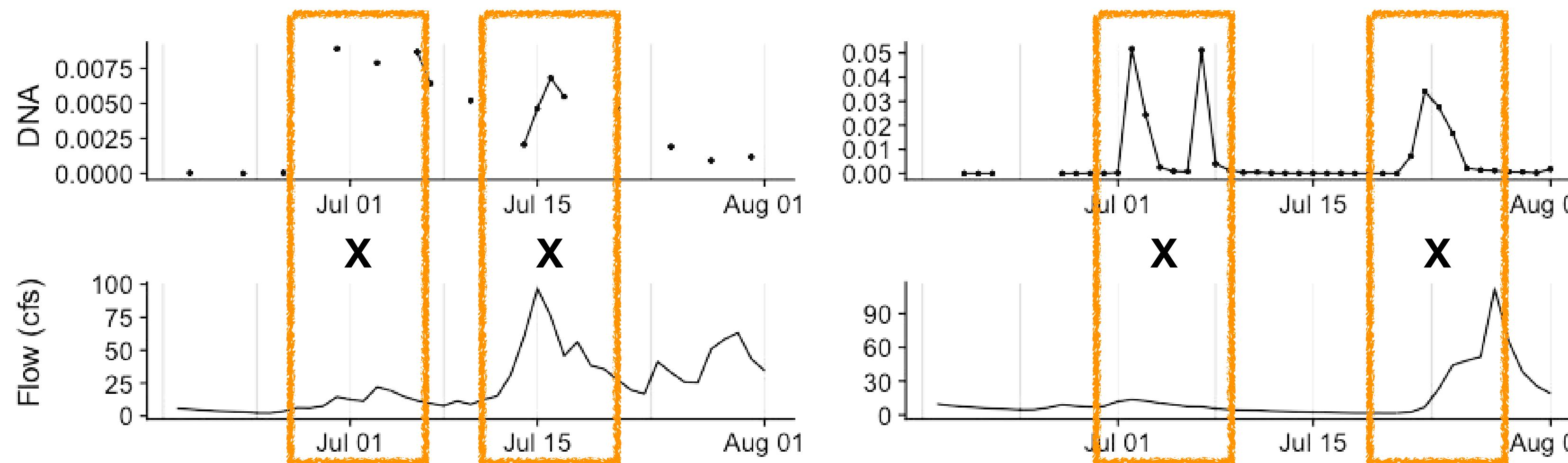
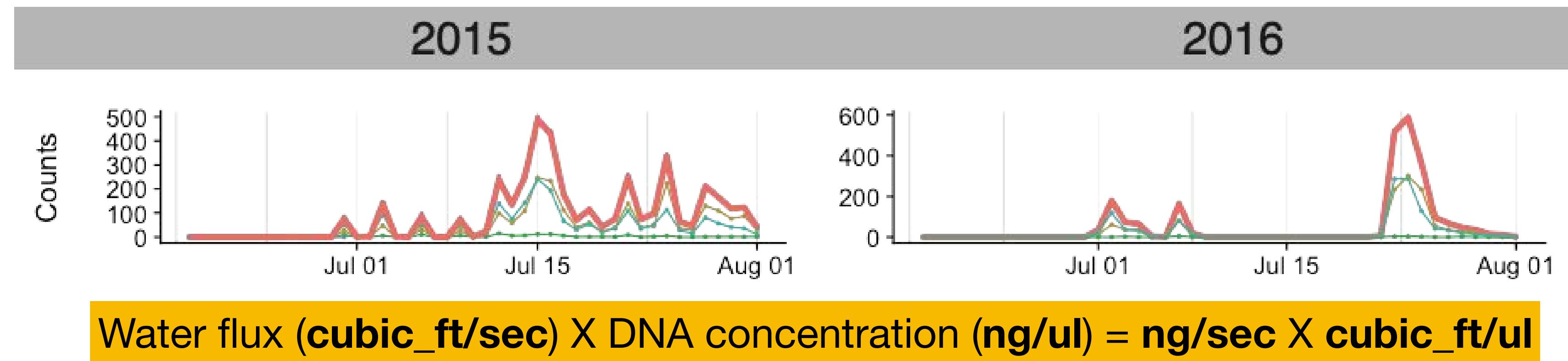
Sockeye in-migrating adults



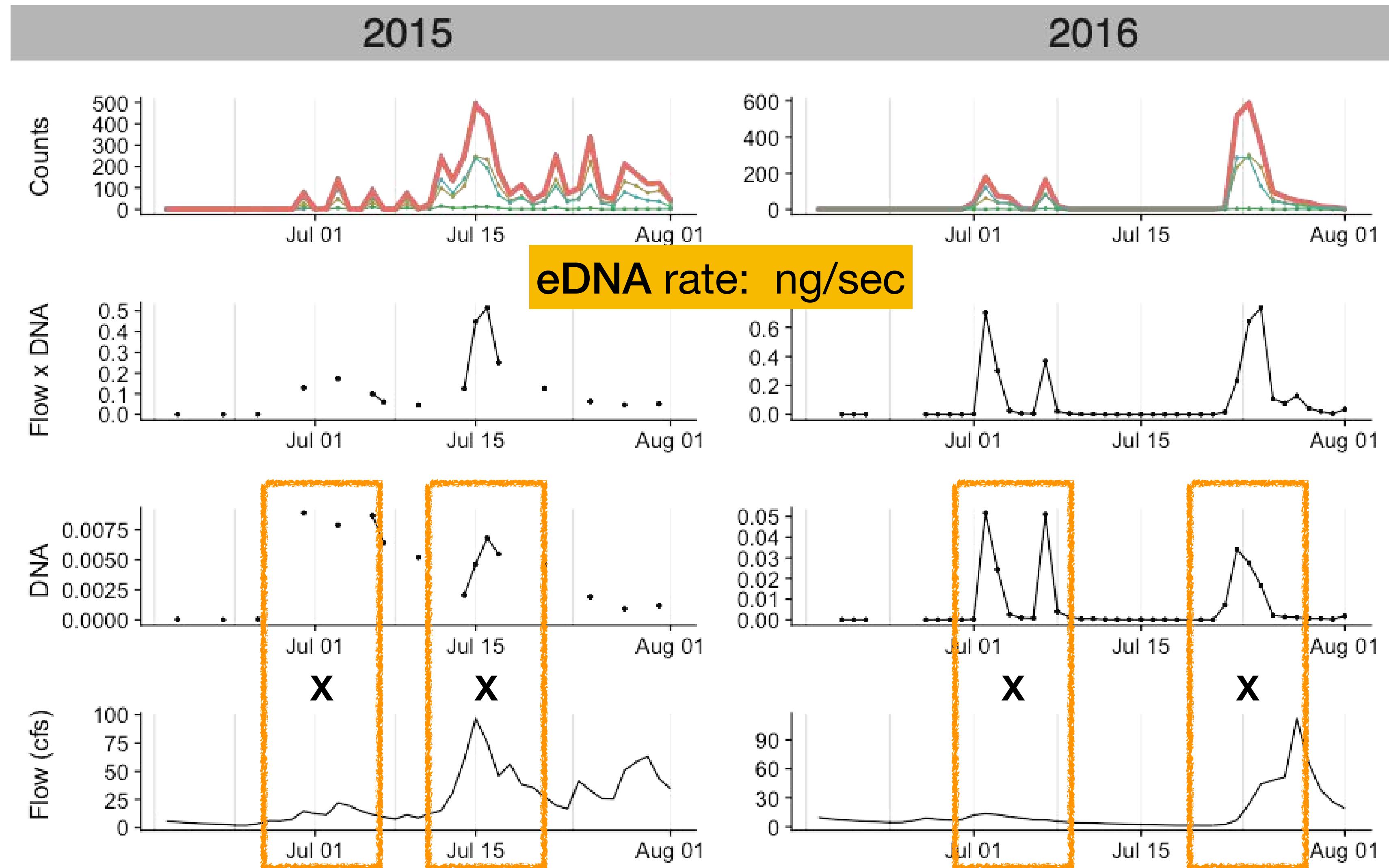
Sockeye in-migrating adults



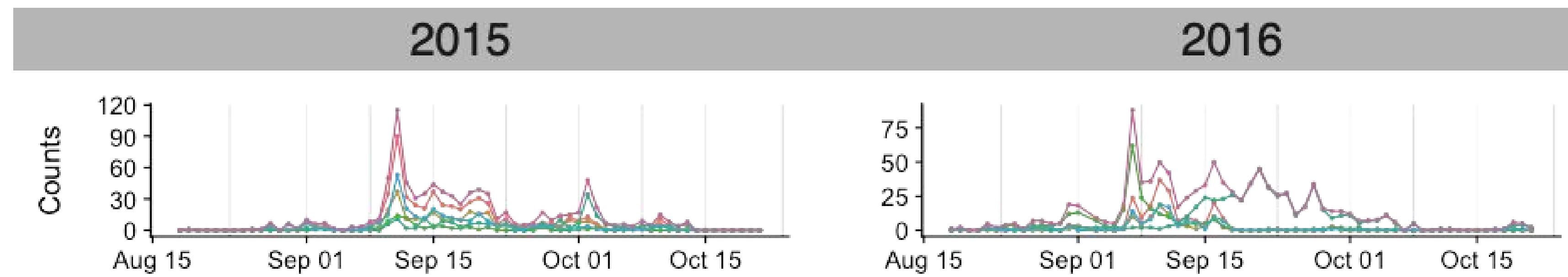
Sockeye in-migrating adults



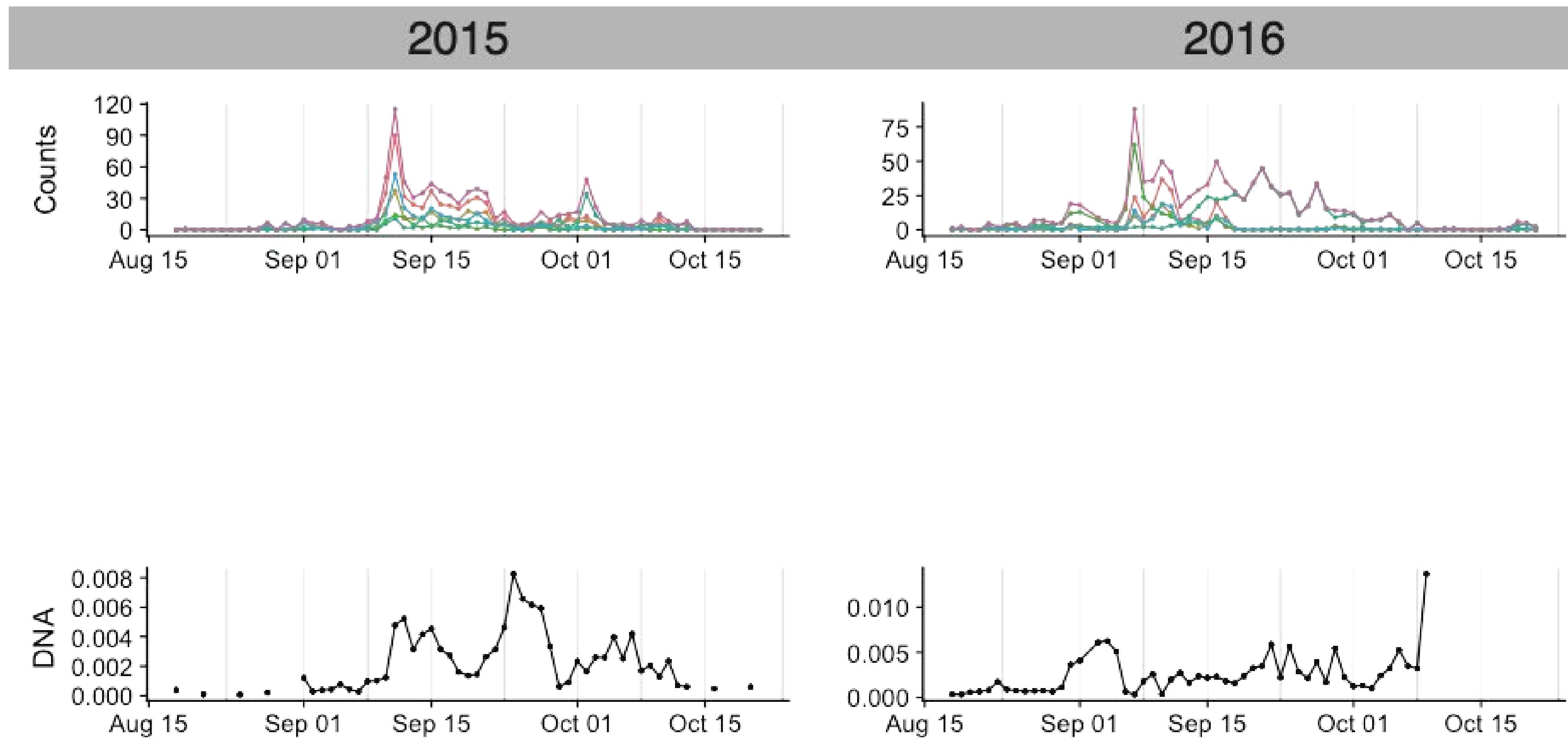
Sockeye in-migrating adults



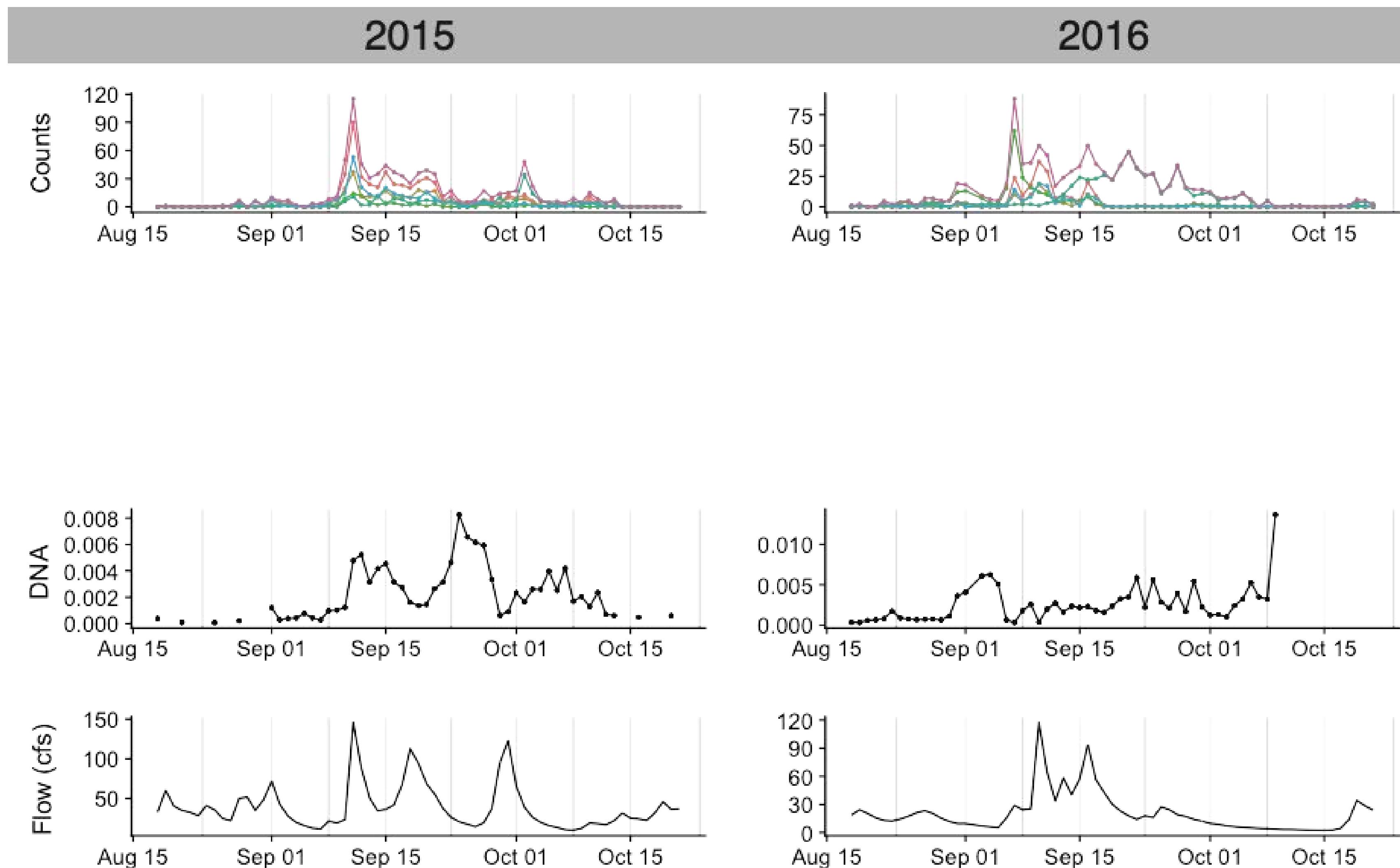
Coho in-migrating adults



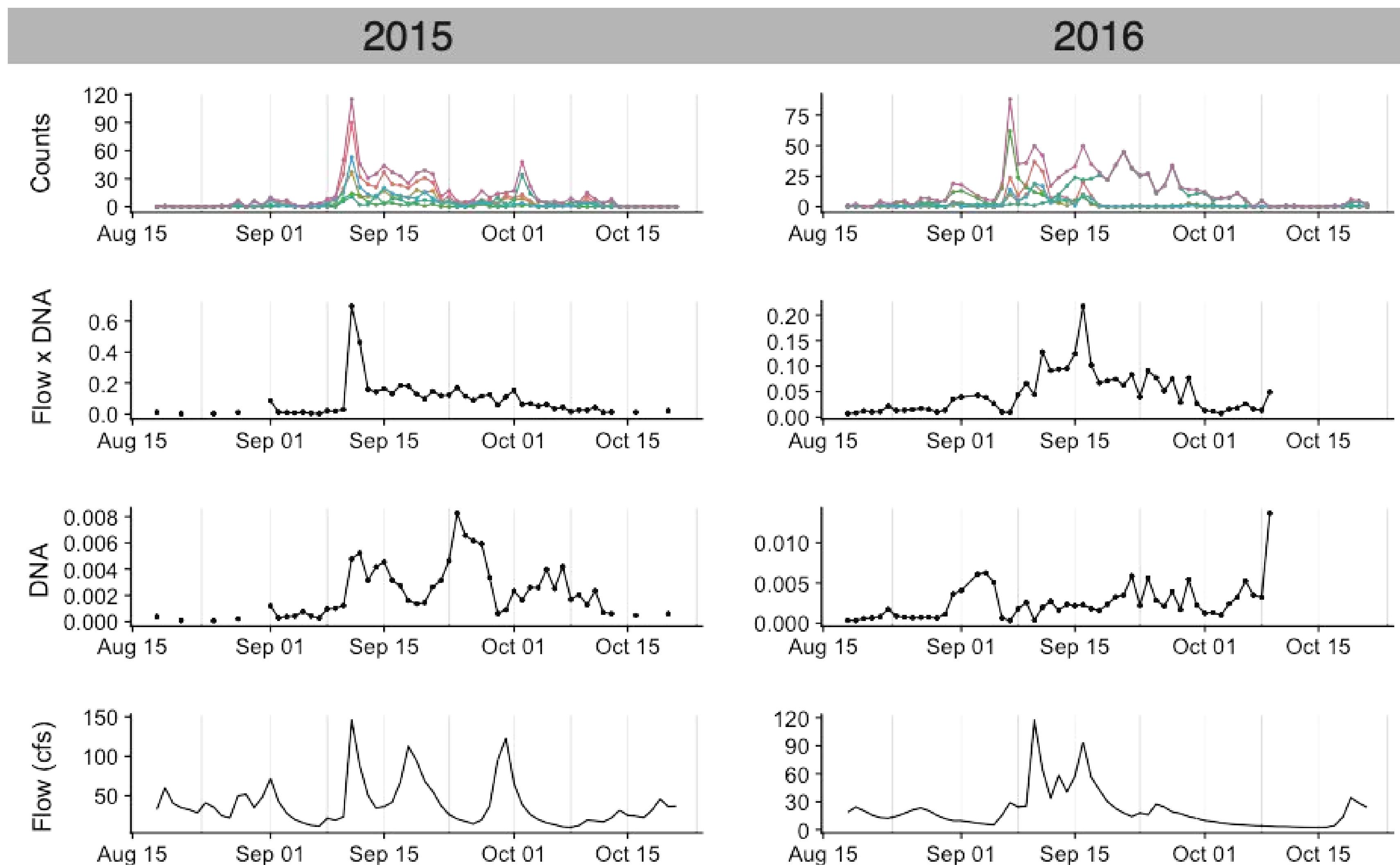
Coho in-migrating adults



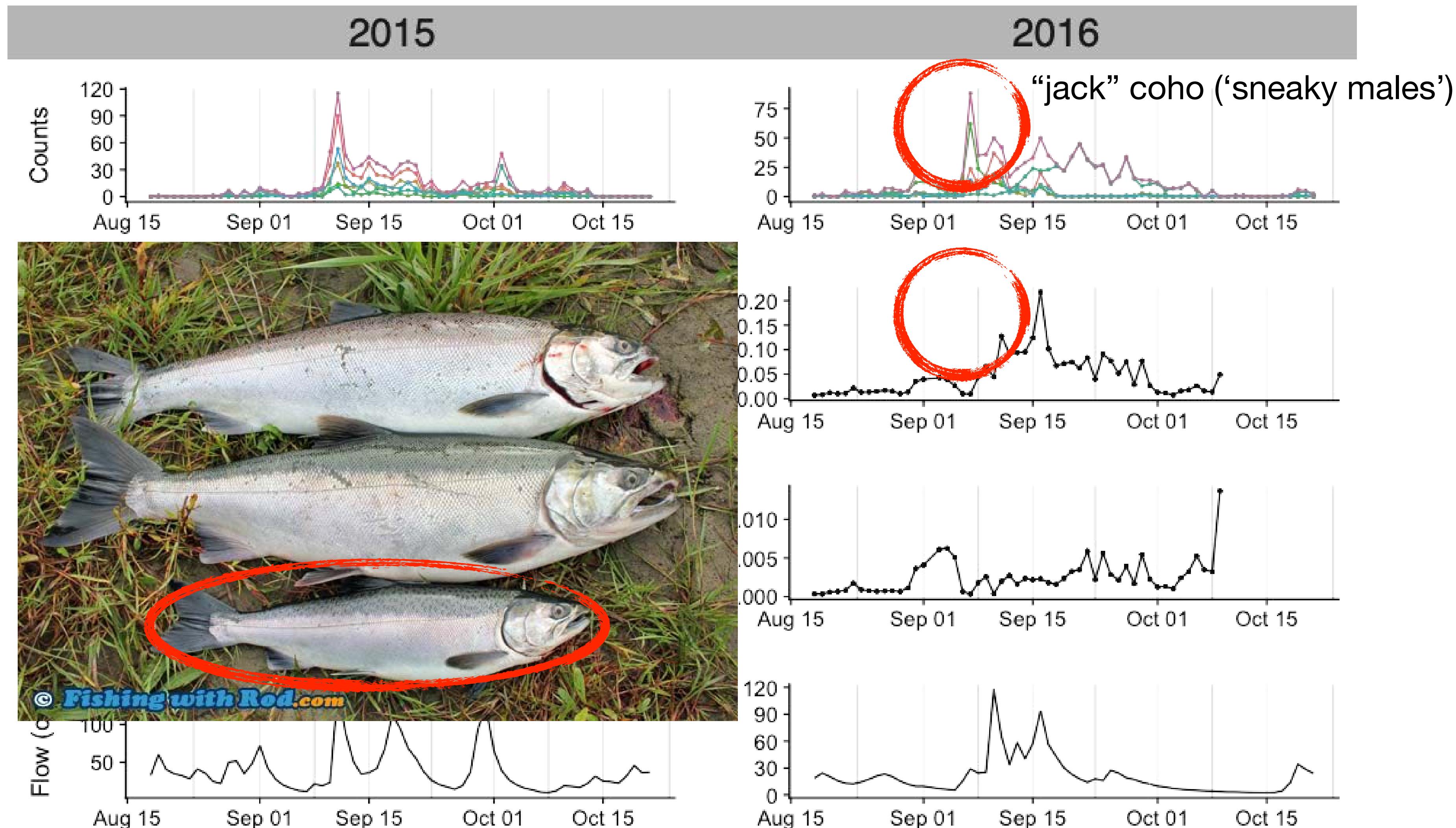
Coho in-migrating adults



Coho in-migrating adults

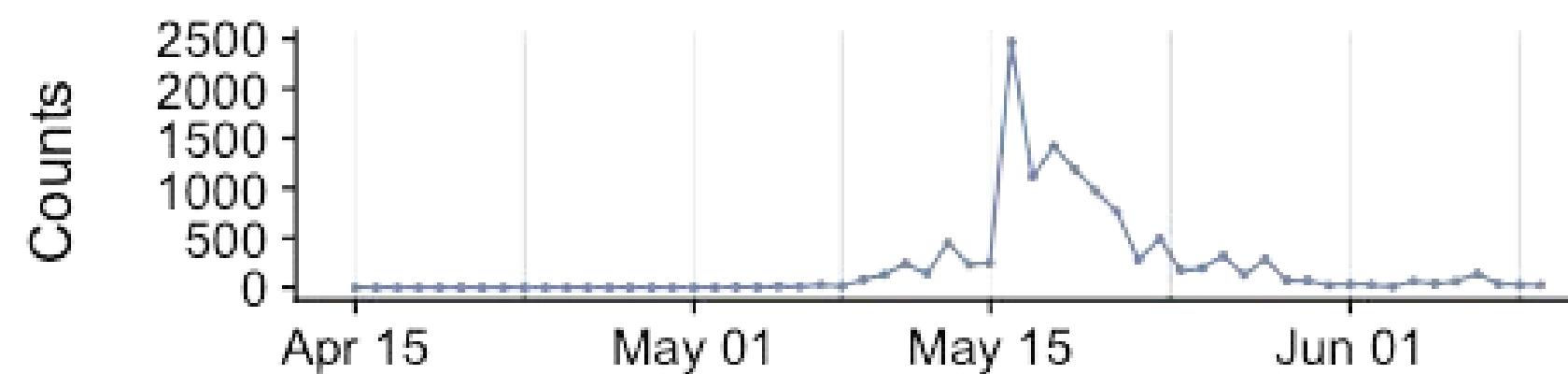


Coho in-migrating adults

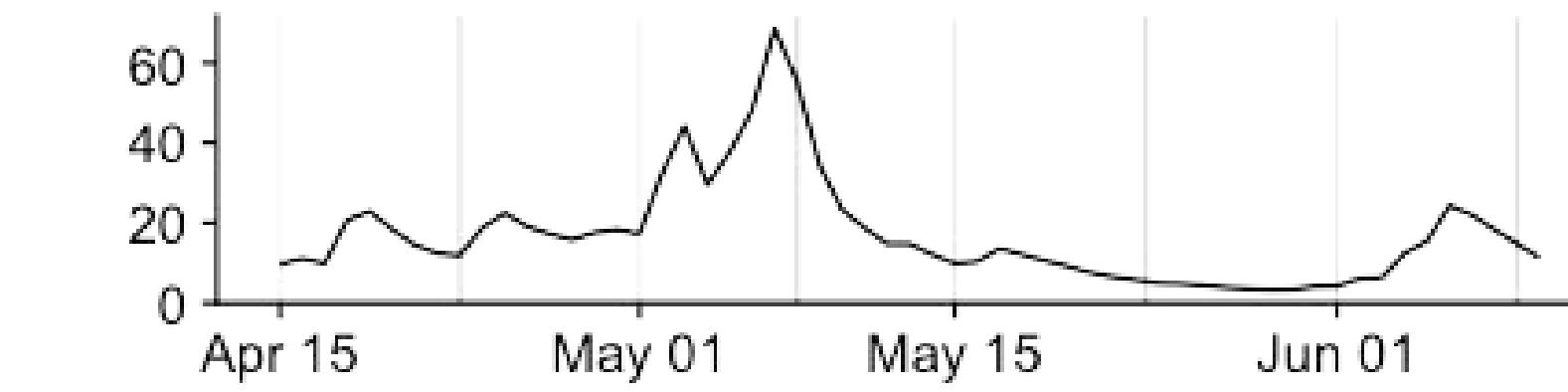
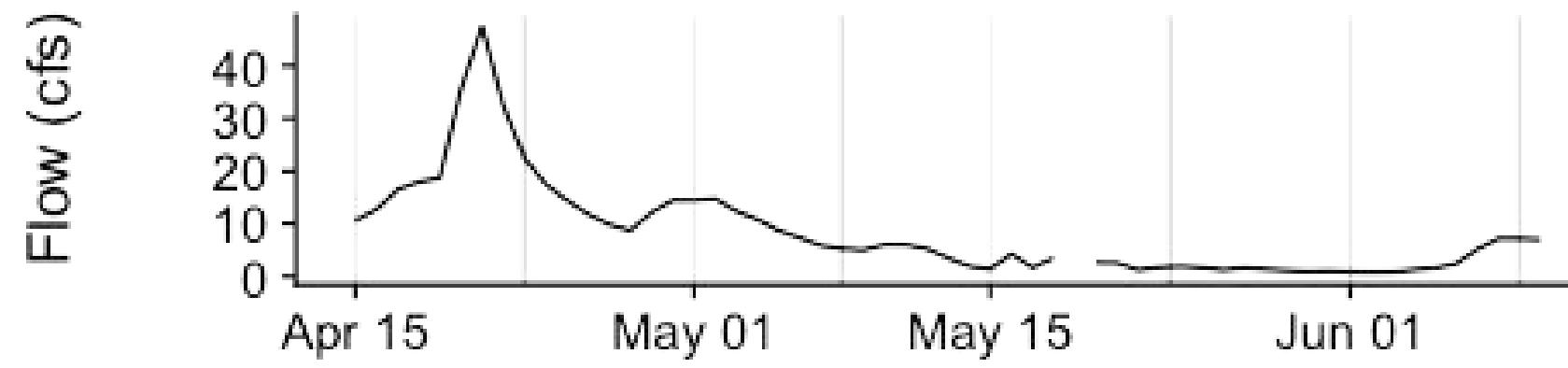
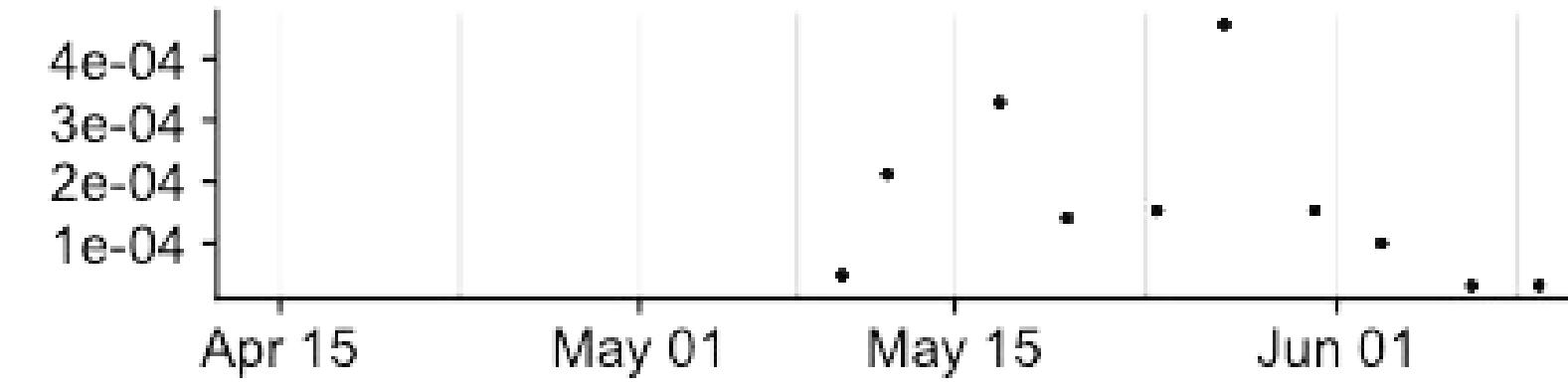
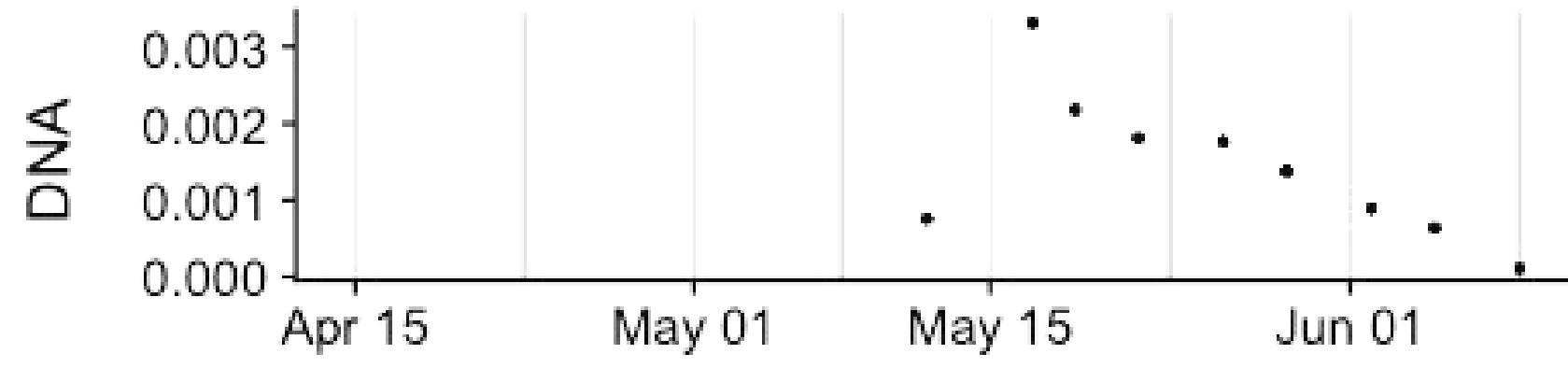
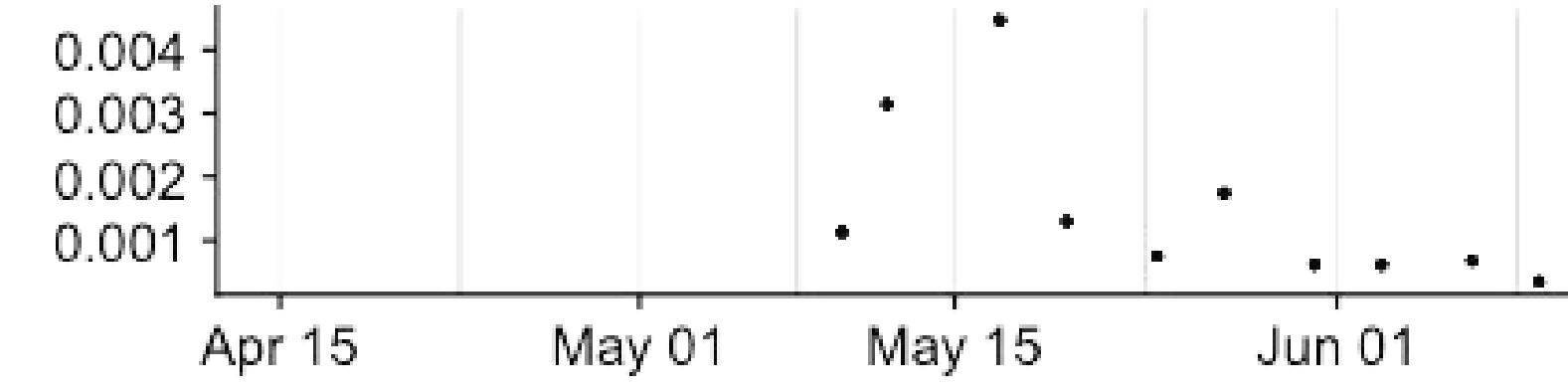
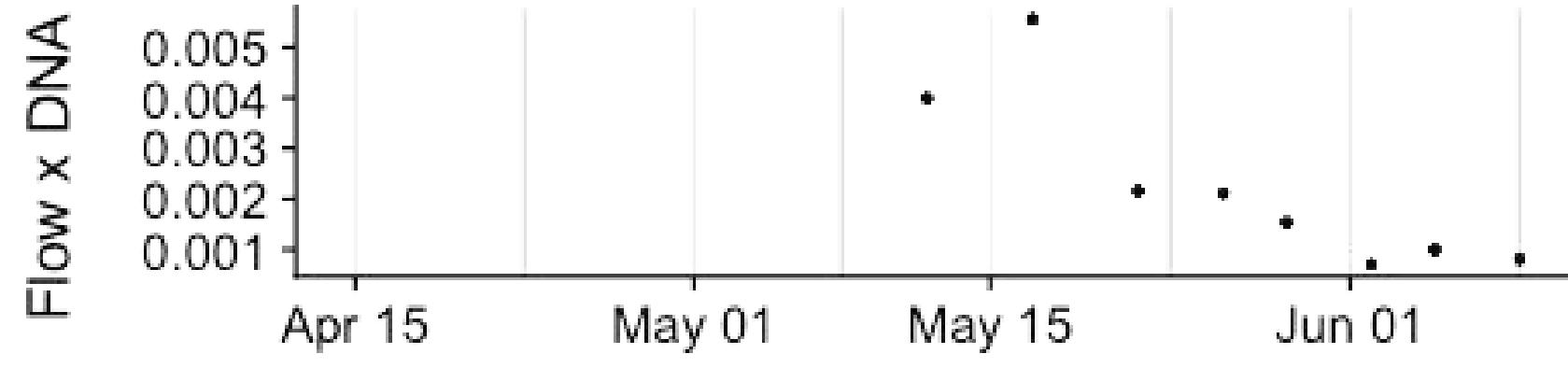
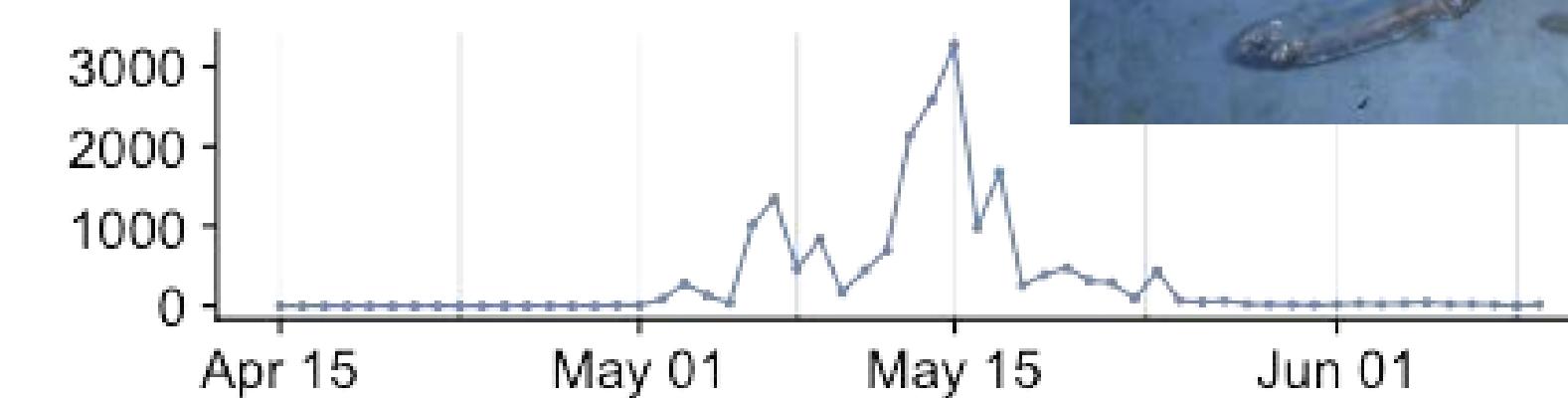


Sockeye out-migrating smolts (juveniles)

2015



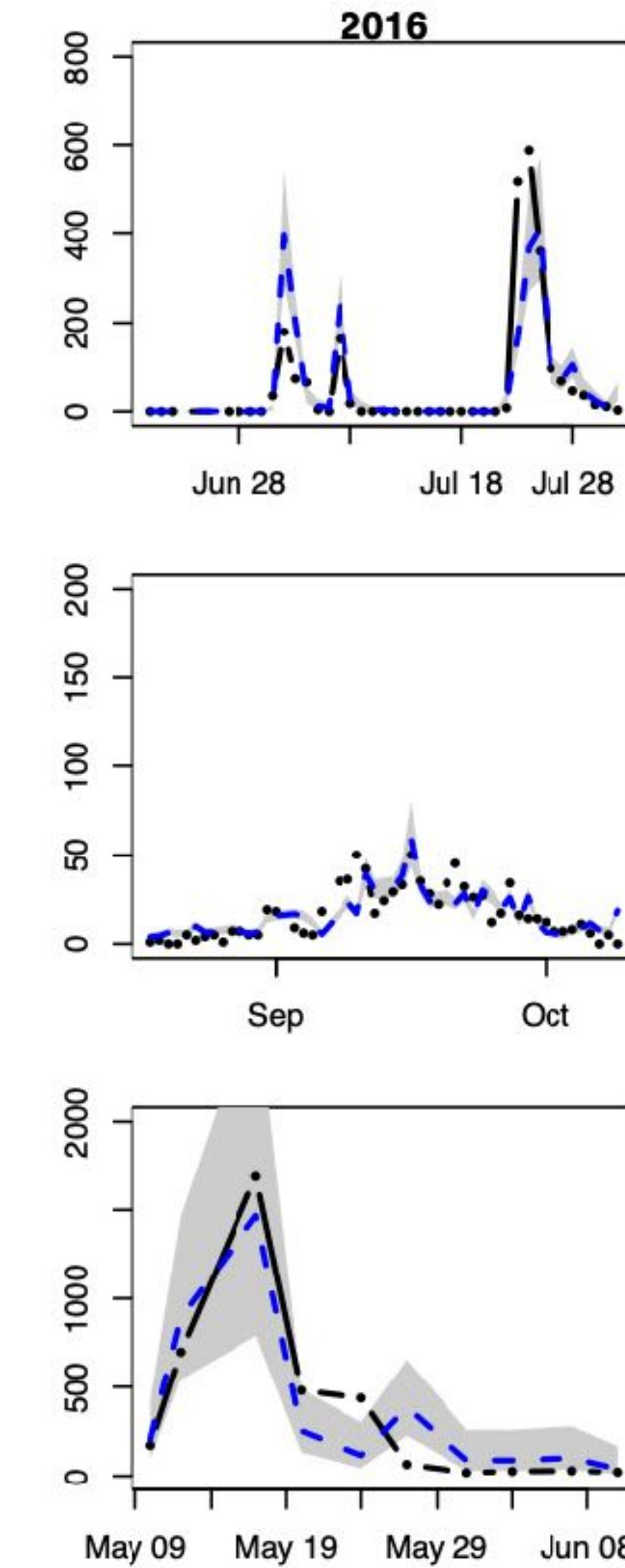
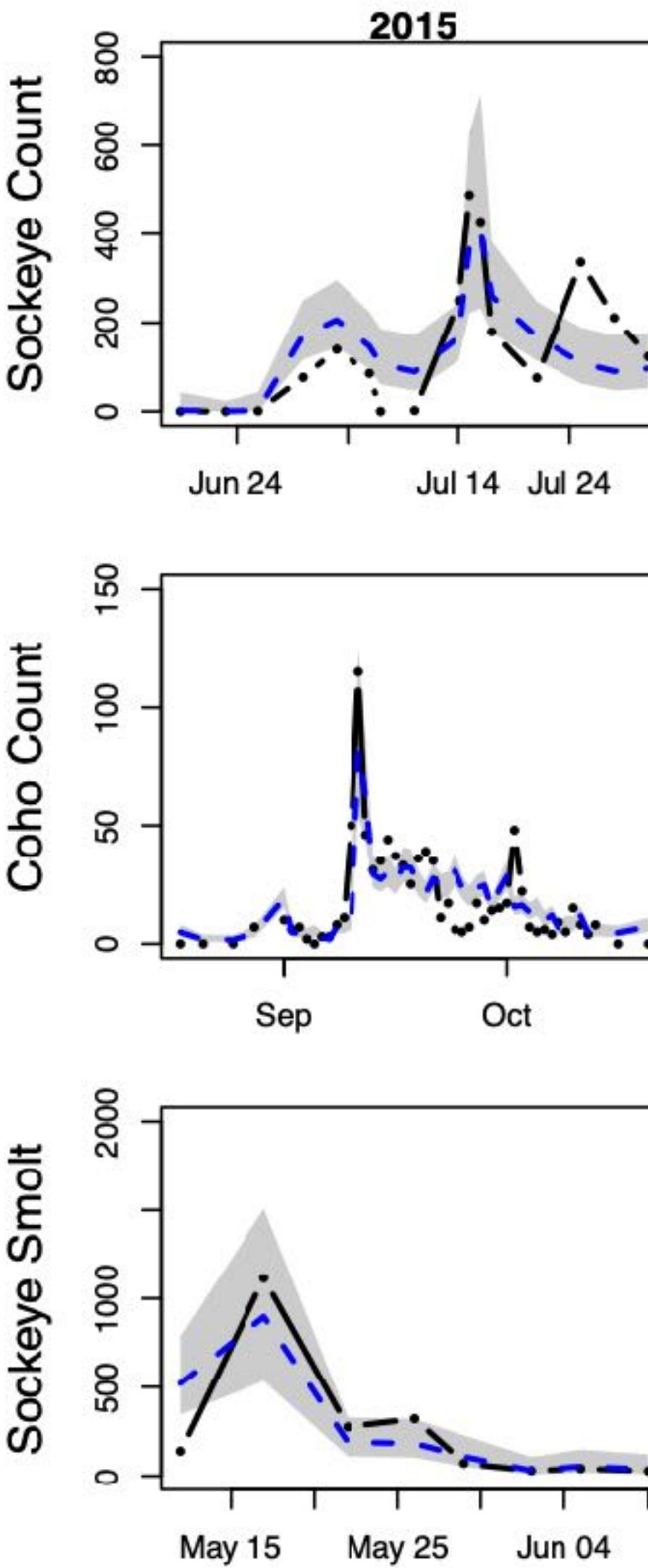
2016





Fitted models (blue) vs. Data (black)

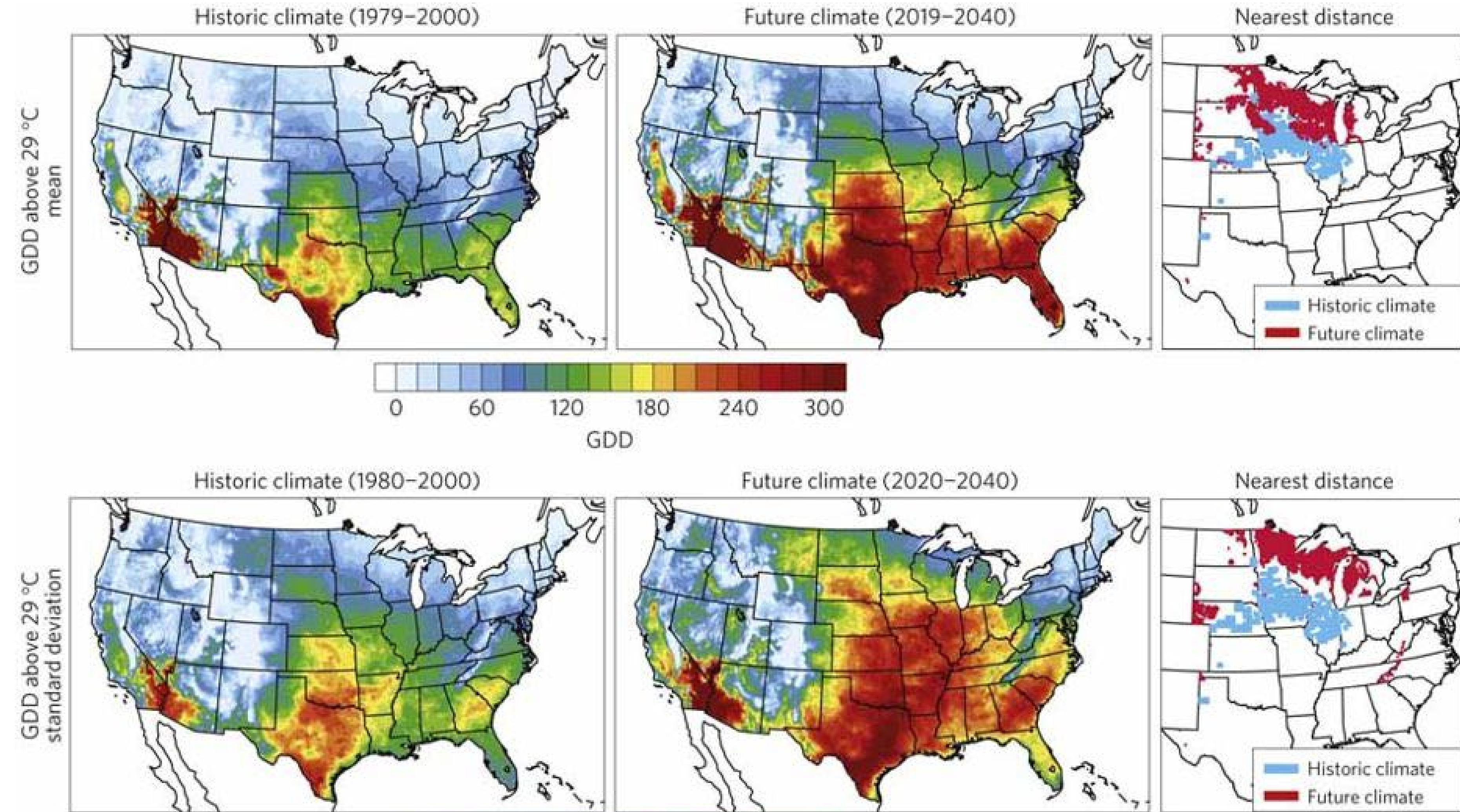
- `glm(log(count) ~ log(eDNA_rate), quasipoisson)`
- All highly significant
- Models not significantly different **between years**
- Small signal of salmon one day ago,
No effect of salmon two days ago



Methods to extract abundance information from DNA data

- Single-species quantitative PCR (qPCR): *control for eDNA transport dynamics*
- **Multiplexed *individual* barcoding (mBRAVE)**
- Mitogenomics and DNA spike-in (SPIKEPIPE)
- Metabarcoding and DNA spike-in (qSeq)
- Reverse metagenomics (RevMet)

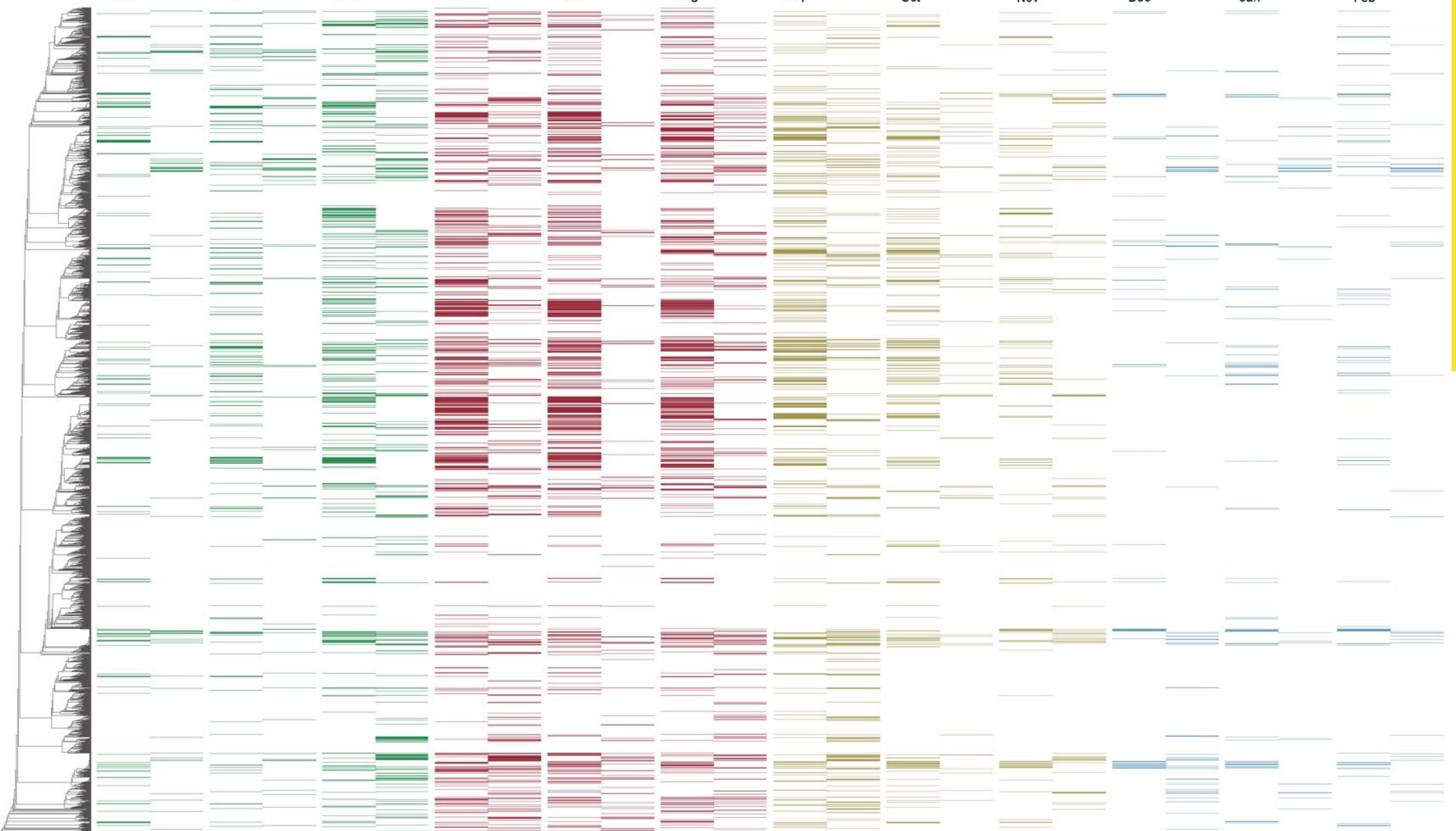
How much biodiversity will we lose because of climate change?



Changes in distribution due to changes in climate envelopes



South North

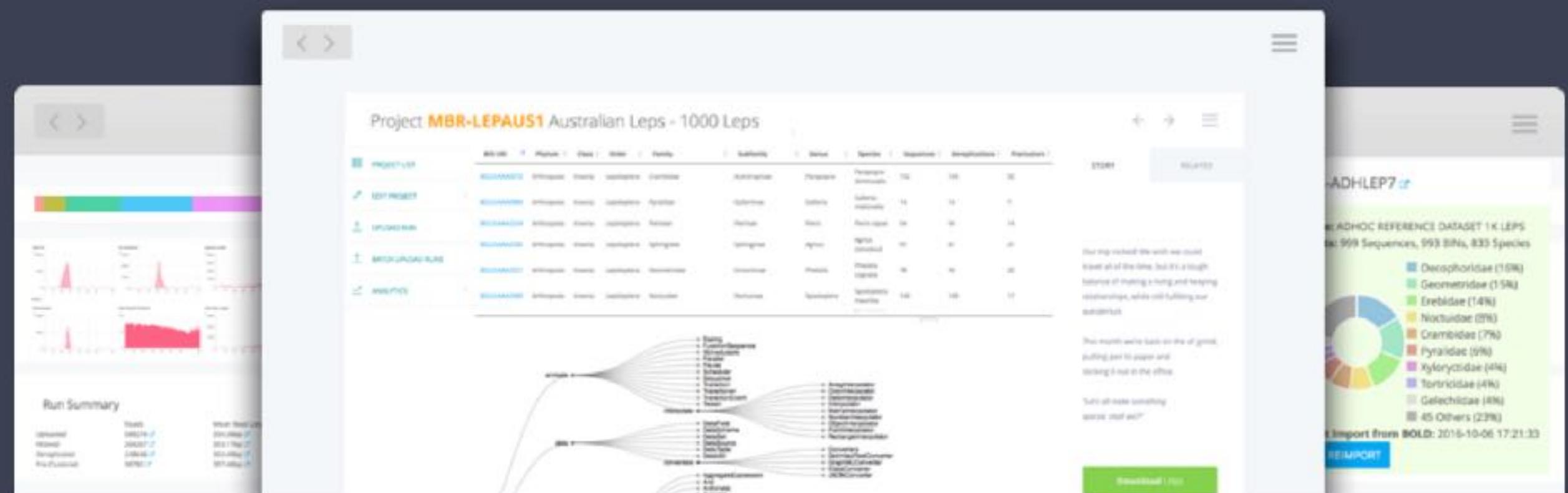


- 10 462 beetles
- collected biweekly over winter and summer
- alpine and subtropical sites
- individually DNA barcoded
- 2210 beetle BINS (“species”)
 - Barcode Identification Numbers



Multiplex Barcode Research And Visualization Environment

mBRAVE is a multi-user platform supporting the storage, validation, analysis, and publication of highly multiplexed projects based on high-throughput sequencing (HTS) instruments. This system builds on the [BOLD Platform](#) to support species identification and discovery for HTS data.



Login

Register

mBRAVE is currently accepting a limited number of new users on a weekly basis.

Storage

Analysis

Indexing

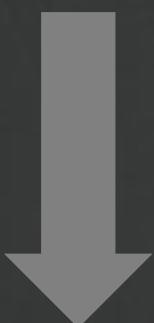
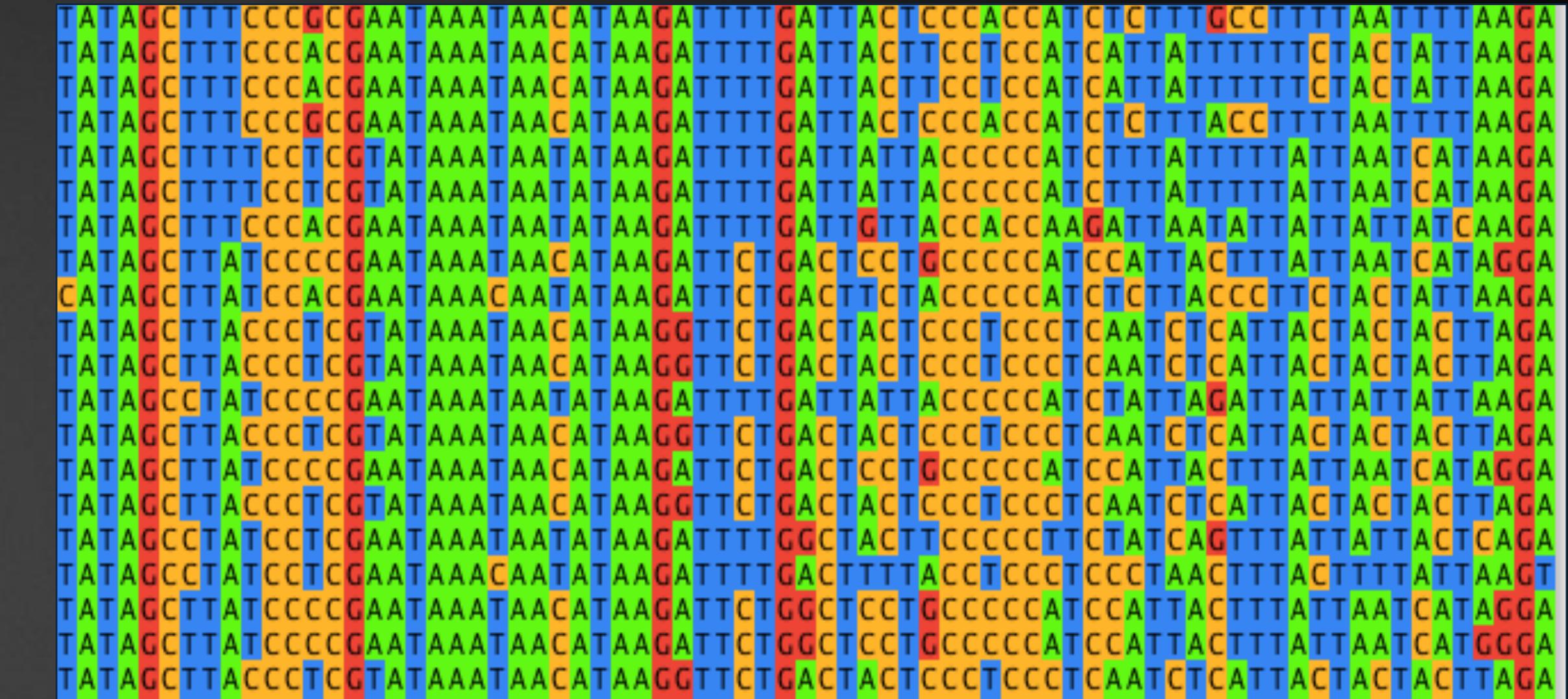
Discovery

Standards

mbraive.net

Massively parallel individual barcoding

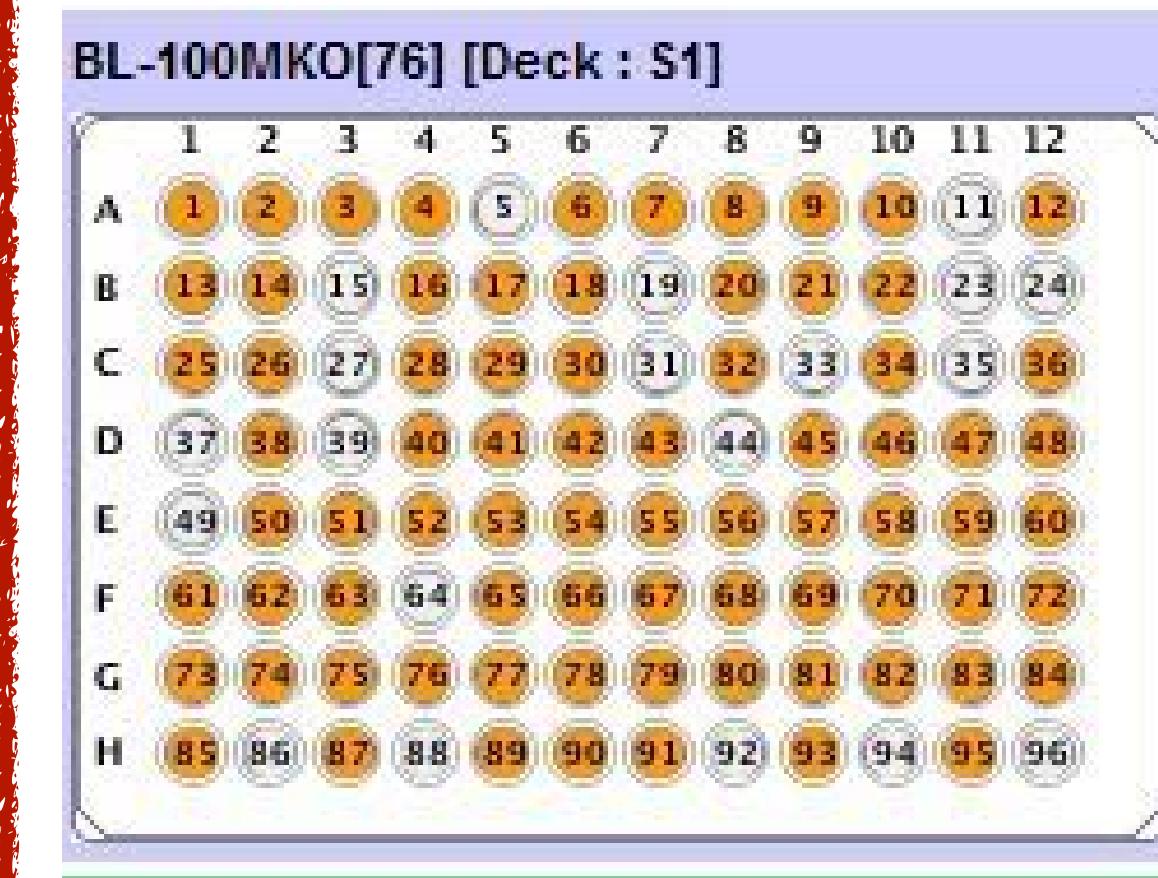
Multiplexed individual barcoding



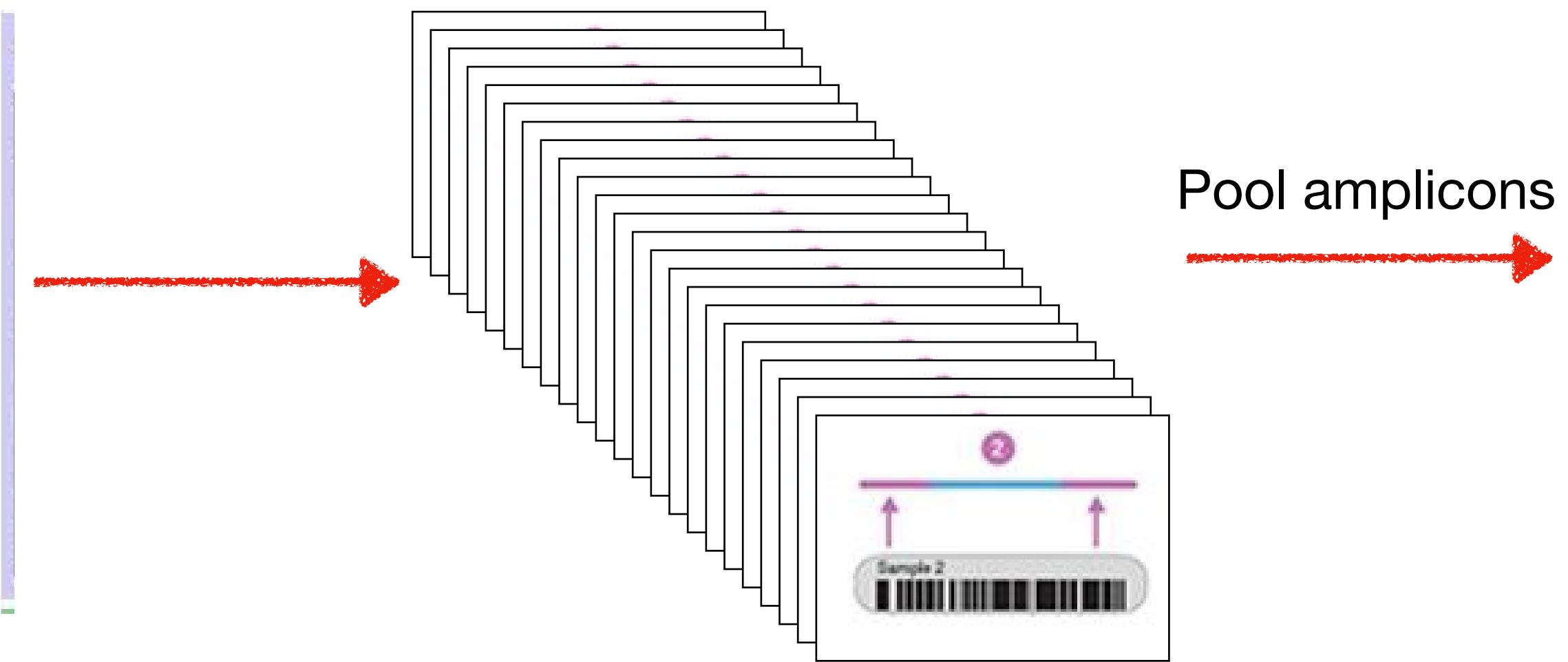
TATACTTATCCCCGAATAAAATAACATAAGATTGACTCCCCCATCATCAAATTATAAT

DNA from 96 beetles in 96 wells

\$\$ Add 96 COI primer pairs
with 96 different twin tags \$\$



Tag F Primer barcode sequence R Primer Tag



Pool amplicons

A consensus COI sequence for each beetle,
plus metadata, in fasta file: **count data!!**

```
>ANML_DLJ06_E12|run:ANML_DLJ06_dtd_0|contig_id:1|rep_count:1169|id_similarity:98.  
8950276243|c_count:1|cmnd:1.1|cmnd:0.1116338751|cmnd:None|p:None|c:None|o:None|f:None|g:  
None|s:None|otu:None|date:2019-08-13  
AACTCTTTATTTTATTGCTGGGGCTGGAAATGTTGTATACCTTCTAGTACTGATTGCGTTGAATTAAGAACACCTGGTA  
CTCTACTAGGTGACGCCAATATAATAGTAACTCGTACAGTCATGCATTATCATATAATTTCTTATAGTTACCTATTGTAATT  
GGT  
>ANML_DLJ04_C04|run:ANML_DLJ04_dtd_0|contig_id:1|rep_count:419|id_similarity:91.  
1602209945|c_count:3|cmnd:1.7|cmnd:0.2088385489|cmnd:3.3|p:Arthropoda|c:Insecta|o:None|f:  
None|g:None|s:None|otu:OTU14|date:2019-08-13  
AACACTATTTCTTTGGAGCTTGAGCAGGTATAGTTGGGACATCACTGAGGATTAATTCGAGCAGAATTAGGAACCCGGAT  
CTCTAATTGGTGTGACCAAATCTATAACGTAATTGTAACAGCCATGCTTTGTAATAATTTTTATAGTTACCTATTGTAATT  
GGT  
>ANML_2909_D06|run:ANML_2909_dtd_0|contig_id:1|rep_count:20|id_similarity:nan|c_count:2|  
cmnd:1.1|cmnd:0.255|cmnd:2.2|p:None|c:None|o:None|f:None|g:None|s:None|otu:OTU37|date:  
2019-08-01  
TACTTGTATTTCTTTGGTGTGAGCTGGAAATGAGTGTAGGAACTCTTAAAGCTCTAATTGCGAGAATTAGGGAACTCTGGTT  
CCCTAATTGGTGTGATCAGATTATAATGTAATTGTAACTGCACATGCTTTATTATAATTTTTATAGTTACCAATTGTT  
GGA  
>ANML_DLJ05_H08|run:ANML_DLJ05_dtd_0|contig_id:1|rep_count:6196|id_similarity:nan|c_count:  
6|cmnd:1.0|cmnd:0.103518399|cmnd:2.2|p:None|c:None|o:None|f:None|g:None|s:None|otu:OTU2|  
date:2019-08-14  
CACTCTTTATTTCTTTGGAGCATGGGCTGGAAATCTGGAAACATCTTAAGCTCTCATCGCGAGCACCTGGAAACCCGGTT  
CATTAACTGGAAACGATCAATTACAGTAATTGTTACTGCCAGCTTCTATAATTCTCATGTAATAACCAATCATATAATT  
GGT  
>ANML_2922_F09|run:ANML_2922_dtd_0|contig_id:1|rep_count:1342|id_similarity:nan|c_count:2|  
cmnd:1.7|cmnd:0.1236214685|cmnd:2.2|p:None|c:None|o:None|f:None|g:None|s:None|otu:OTU7|  
date:2019-08-02  
ATCCTTAACTCTTTGGTGTGATCAGGAATAGTGGAAACCTTAAAGATTACTAATTGCGCTAGAATTAGTAATTCTGGTT  
CTTAATTGGAGACGCCAATTTATAATGTAATTGTAACTGCCCATGCTTTATTATAATTCTTATAGTTACCTGTAATAATA  
GGA  
>ANML_DLJ01_F05|run:ANML_DLJ01_dtd_0|contig_id:3|rep_count:10|id_similarity:nan|c_count:3|  
cmnd:0.6|cmnd:0.24|cmnd:2.6|p:None|c:None|o:None|f:None|g:None|s:None|otu:OTU15|date:  
2019-08-13  
AACCTTTACTTCTTTGGAGCATGAAGAGGTATAATCGAACCTCCATAAGATTGATGATTGAAACAGAATTAGGAACAGCTGGT  
CTCTAATTGGAGATGACAATTATAATGTAATTGTAACAGCACATGCTTTATCATGTAATAACCCATCATGATA  
GGA
```

Demultiplex by tag

Keep the most
abundant read per tag
(= well = beetle), filter
out the rest, and
Assign a taxonomy to
using BOLD



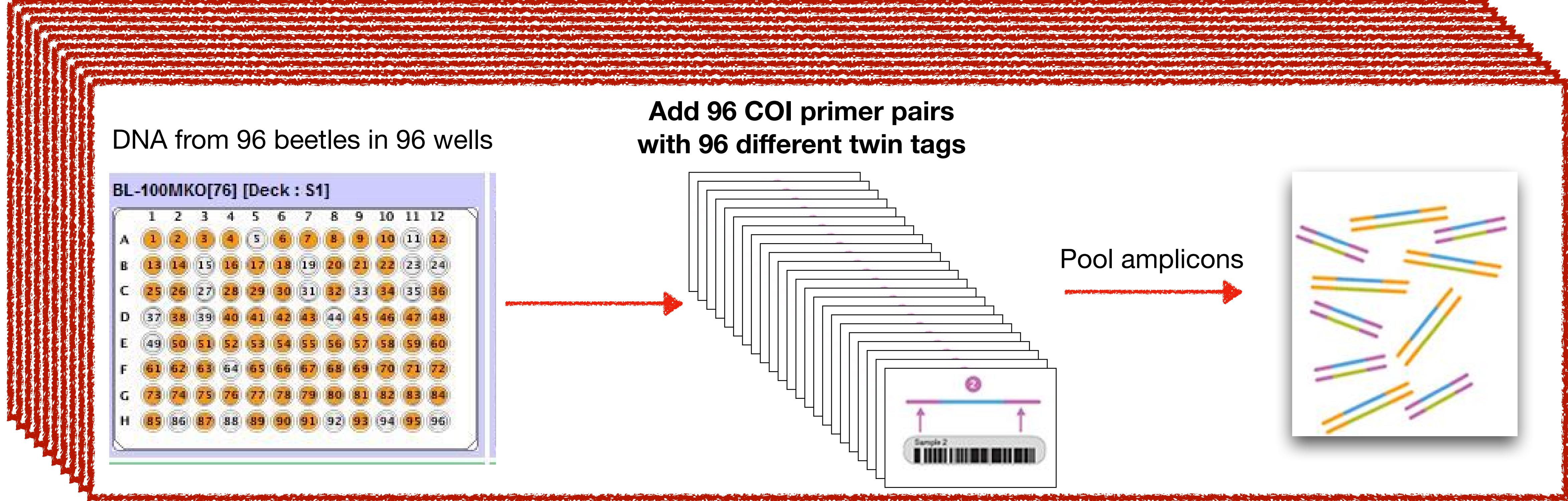
Upload fastq files
to mBRAVE

A	B	C	D
Forward Labels	Reverse Labels	Label	Group
0048_Asf_Ibc269	0096_Asr_Ibc239_rc	CONTROL_H12_BI0UG39505	BL-106R12
0048_Asf_Ibc269	0095_Asr_Ibc172_rc	CONTROL_H12_BI0UG39372	BL-106R12
0048_Asf_Ibc269	0094_Asr_Ibc151_rc	GMPAT212-18	BL-106R12
0048_Asf_Ibc269	0093_Asr_Ibc180_rc	GMPAT206-18	BL-106R12
0048_Asf_Ibc269	0092_Asr_Ibc311_rc	GMPAT211-18	BL-106R12
0048_Asf_Ibc269	0091_Asr_Ibc122_rc	GMPAT205-18	BL-106R12
0048_Asf_Ibc269	0090_Asr_Ibc269_rc	GMPAT210-18	BL-106R12
0048_Asf_Ibc269	0089_Asr_Ibc283_rc	GMPAT204-18	BL-106R12
0048_Asf_Ibc269	0088_Asr_Ibc266_rc	GMPAT209-18	BL-106R12
0048_Asf_Ibc269	0087_Asr_Ibc261_rc	GMPAT203-18	BL-106R12
0048_Asf_Ibc269	0086_Asr_Ibc2	GMPAT208-18	BL-106R12
0048_Asf_Ibc269	0085_Asr_Ibc261_rc	GMPAT202-18	BL-106R12
0048_Asf_Ibc269	0084_Asr_Ibc305_rc	GMPAT207-18	BL-106R12
0048_Asf_Ibc269	0083_Asr_Ibc310_rc	GMPAT201-18	BL-106R12

Upload tag
information to
mBRAVE



Illumina
sequence



A consensus COI sequence for each beetle, plus metadata, in fasta file: **count data!!**

```
>ANML_DLJ06_E12|run:ANML_DLJ06_dtd_0|contig_id:1|rep_count:1169|id_similarity:98.8950276243|c_count:1|cmxd:1.1|cmnd:0.1116338751|cnnd:None|p:None|c:None|o:None|f:None|g:None|s:None|otu:None|bold:AAE7930|date:2019-08-13
AACTCTTATTITATTITGGTGGCTGGGAATAGTTGGTATATCCCTTAGATACTGATTGAGTTGAATTAAGACACCTGGTA
CTCTACTAGGTGACGCCAATATATAATGTAATCGTACAGTCATGCATTATCATATAATTCTTATAGTTACCTATTGTAATT
GGT
>ANML_DLJ04_C04|run:ANML_DLJ04_dtd_0|contig_id:1|rep_count:419|id_similarity:91.1602209945|c_count:3|cmxd:1.7|cmnd:0.2088385489|cnnd:3.3|p:Arthropoda|c:Insecta|o:None|f:None|g:None|s:None|otu:OTU14|date:2019-08-13
AACTATATTCTTTGGAGCTTGAGCAGGTATGTTGGGACATCACTGAGAGTTACTAATTGAGCAGAATTAGGAACCCGGAT
CTCTAATTGGTGTGACCAAATCTATAACGTAATTGTAACAGCCATGCTTTGTAATAATTTTTATAGTTACCTATTGTAATT
GGT
>ANML_2909_D06|run:ANML_2909_dtd_0|contig_id:1|rep_count:20|id_similarity:nan|c_count:2|cmxd:1.1|cmnd:0.255|cnnd:2.2|p:None|c:None|o:None|f:None|g:None|s:None|otu:OTU37|date:2019-08-01
TACTTGTATTTCTTTGGTGGCTGGAGCTGGAAATGAGCTTAAAGCTCTAATTGAGCAGAATTAGGAACATCTGGTT
CCCTAATTGGTGTGATCAGATTATAATGTAATTGTAACTGCACATGCTTCAATTATAATTTTTATAGTTACCAATTGTTATT
GGA
>ANML_DLJ05_H08|run:ANML_DLJ05_dtd_0|contig_id:1|rep_count:6196|id_similarity:nan|c_count:6|cmxd:1.7|cmnd:0.103518399|cnnd:2.2|p:None|c:None|o:None|f:None|g:None|s:None|otu:OTU2|date:2019-08-14
CACTCTTATTITATTCTTGAGCATGGCTGGAAATCTGGACATCTTAAGCTCTCATCCGGCAGAACCTGGAAACCCGGTT
CATTAACTGGAAACATCAATTACAGTAATTGTTACTGCCAGCTTCTATAATTCTCATGTAATAACCAATCATATAATT
GGT
>ANML_2922_F09|run:ANML_2922_dtd_0|contig_id:1|rep_count:1342|id_similarity:nan|c_count:2|cmxd:1.7|cmnd:0.1236214685|cnnd:2.2|p:None|c:None|o:None|f:None|g:None|s:None|otu:OTU7|date:2019-08-02
ATCTTAACTCTTTGGTGTATGATCAGGAATAGTTGGAACCTTAAAGATTACTAATTGCTGAGATTAAAGTAATCTGGGT
CTTAAATTGGAGACGCCAATTTAATGTAATTGTAACTGCCCATGCTTATAATTCTTATAGTTACCTGTAATAATAA
GGA
>ANML_DLJ01_F05|run:ANML_DLJ01_dtd_0|contig_id:3|rep_count:10|id_similarity:nan|c_count:3|cmxd:0.6|cmnd:0.24|cnnd:2.6|p:None|c:None|o:None|f:None|g:None|s:None|otu:OTU15|date:2019-08-13
AACCCTTACTTCATTTGGAGCATGAAGAGGTATAATCGAACCTCCATAAGATTGATGATTGAGACAGAATTAGGAACAGCTGGT
CTCTAATTGGAGATGATCAAATTATAATGTAATTGTAACAGCACATGCTTTATCATGTAATAACCCATCATGATA
GGA
```

Demultiplex by library and tag

Keep the most abundant read per tag (= well = beetle), filter out the rest, and Assign a taxonomy to using BOLD



Methods to extract abundance information from DNA data

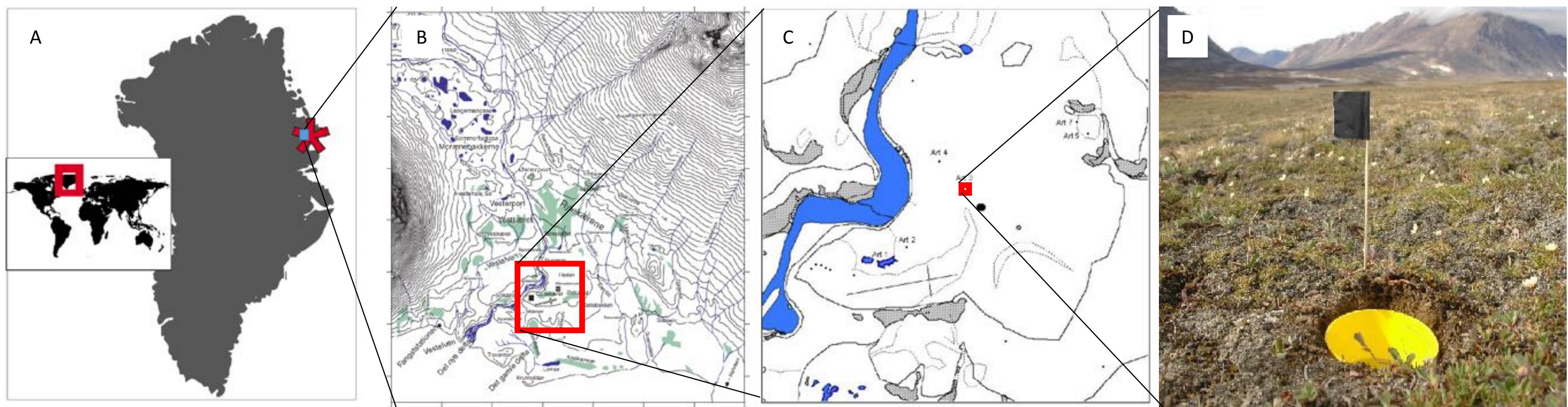
- Single-species quantitative PCR (qPCR)
- Multiplexed individual barcoding (mBRAVE): *individual count data in the 1000s*
- **Mitogenomics and DNA spike-in (SPIKEPIPE)**
- Metabarcoding and DNA spike-in (qSeq)
- Reverse metagenomics (RevMet)

Zackenberg Research Station, Greenland



SPIKEPIPE: A metagenomic pipeline for the accurate quantification of eukaryotic species occurrences and intraspecific abundance change using DNA barcodes or mitogenomes

Yinqui Ji^{1*} | Tea Huotari^{2*} | Tomas Roslin^{2,3} | Niels Martin Schmidt^{4,5} |
Jixin Wang¹ | Douglas W. Yu^{1,6,7} | Otso Ovaskainen^{8,9}



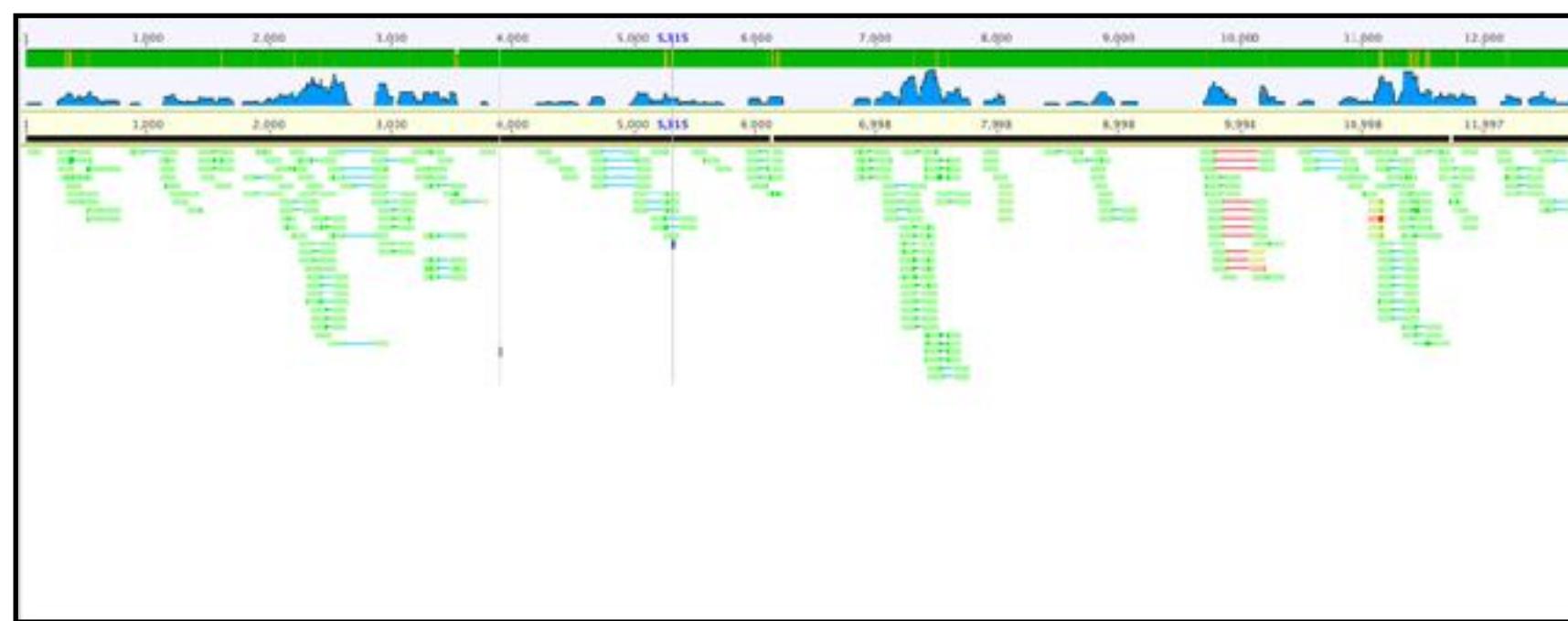


- entire aboveground arthropod community **~375 species**
- **>760,000 arthropods** collected in weekly samples and multiple pan traps from **1996-2013** (and ongoing)
- We assembled **308 mitogenome** sequences
- We shotgun-sequenced **~750 samples: 3 samples per week from 1997-2013** + technical replicates

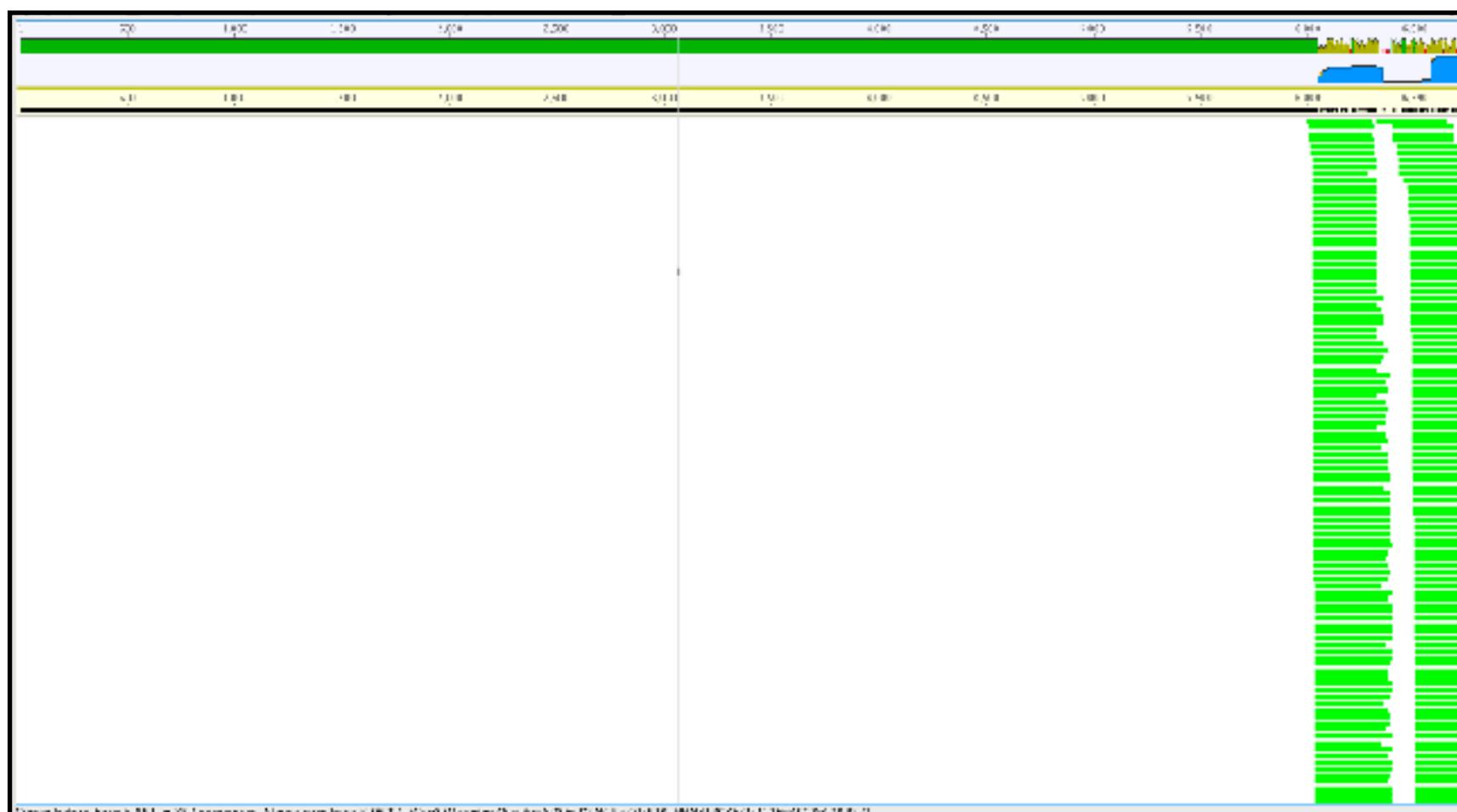
We mapped each sample's short reads to 308 mitogenomes



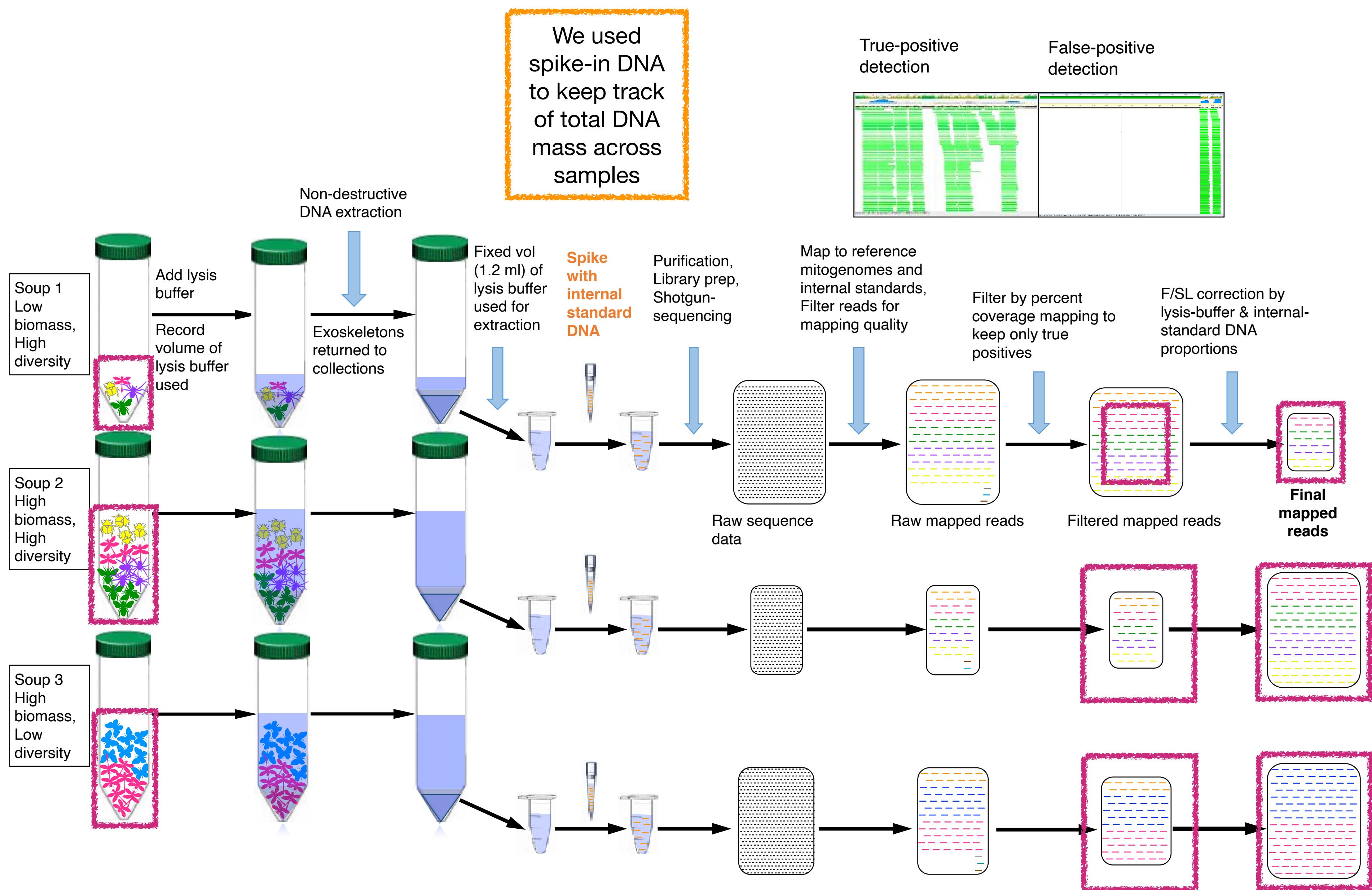
This mitogenome (a species) is confidently present in high biomass



This species is confidently present in low biomass



This species is not present, despite the many mapped reads (they mapped only to a conserved sequence)

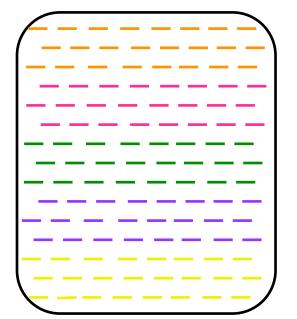


Spike-in
Soup 1
Low
biomass



stochasticity
随机性

big
fastq
dataset

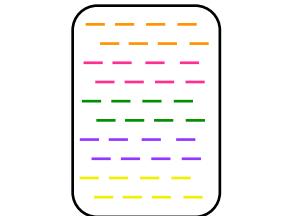


Spike-in
Soup 2
High
biomass



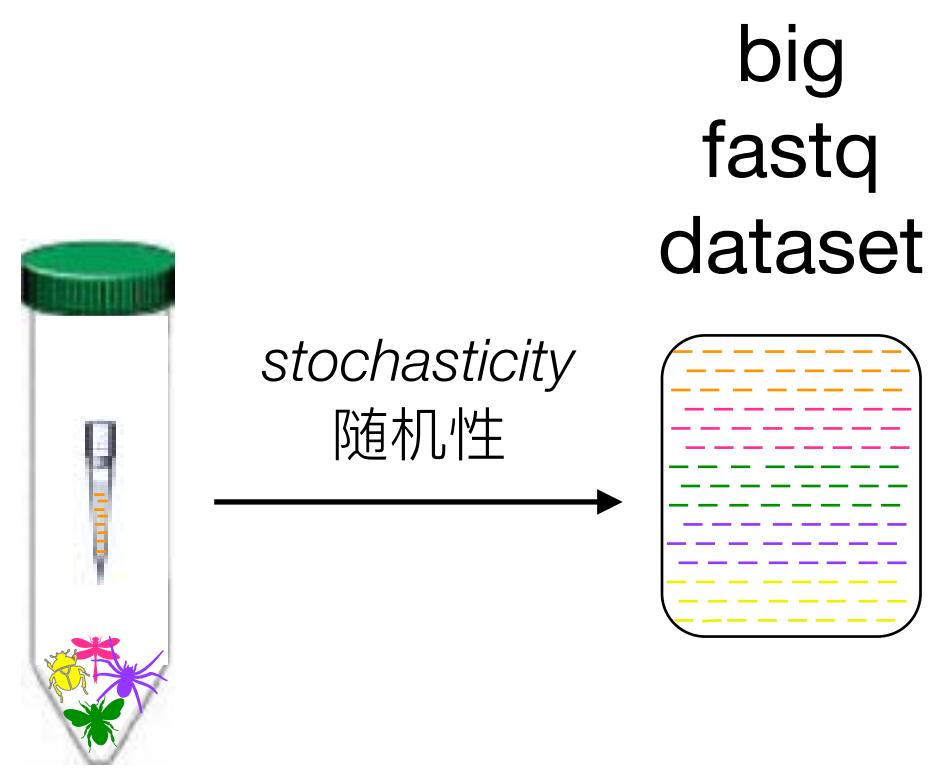
stochasticity
随机性

small
fastq
dataset

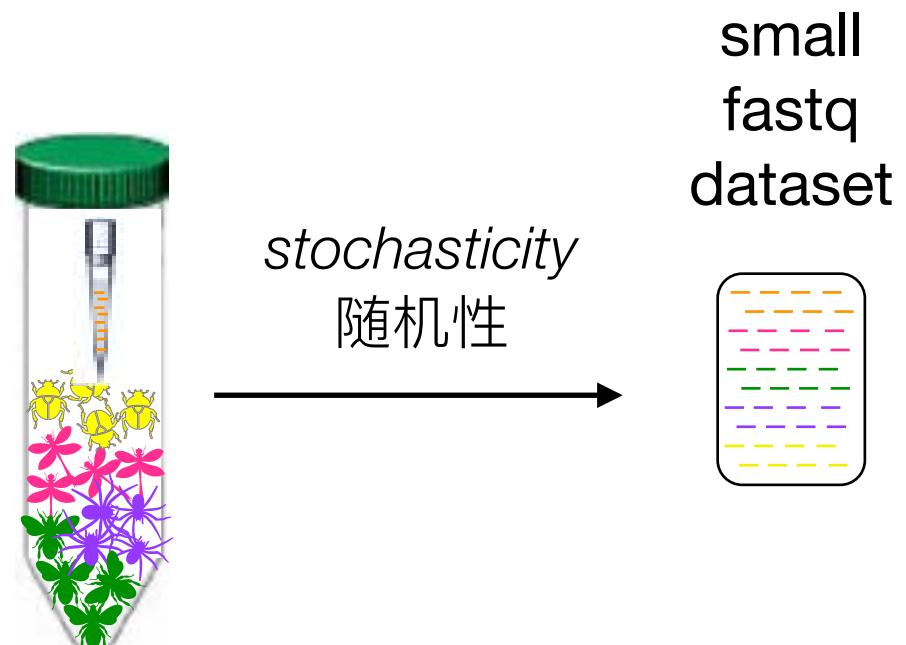


Sample 1	Reads_Raw
Taxon1	1000
Taxon2	1200
Taxon3	1600
Taxon4	2000
Taxon5	400
Taxon6	200
Sample 2	Reads_Raw
Taxon1	500
Taxon2	600
Taxon3	800
Taxon4	1000
Taxon5	200
Taxon6	100

Spike-in
Soup 1
Low
biomass

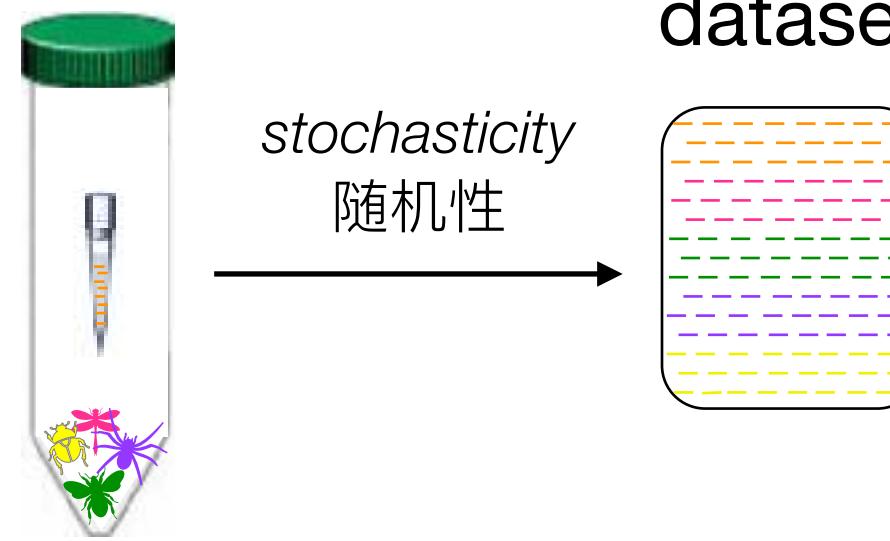


Spike-in
Soup 2
High
biomass

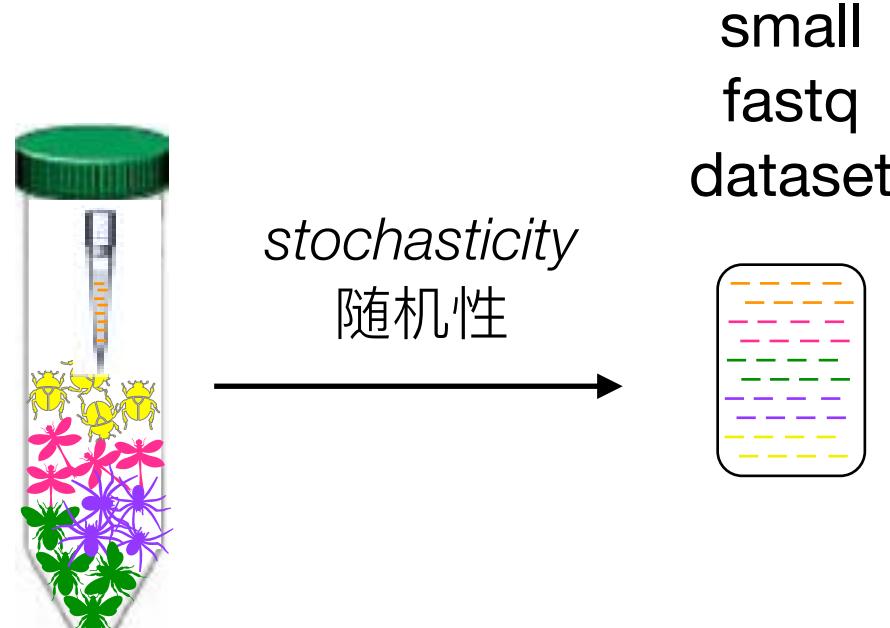


Sample	Reads_Raw
Sample 1	
OTU1	1000
OTU2	1200
OTU3	1600
OTU4	2000
OTU5	400
OTU6	200
OTU_SPIKE	40
Sample 2	
OTU1	500
OTU2	600
OTU3	800
OTU4	1000
OTU5	200
OTU6	100
OTU_SPIKE	10

Spike-in
Soup 1
Low
biomass



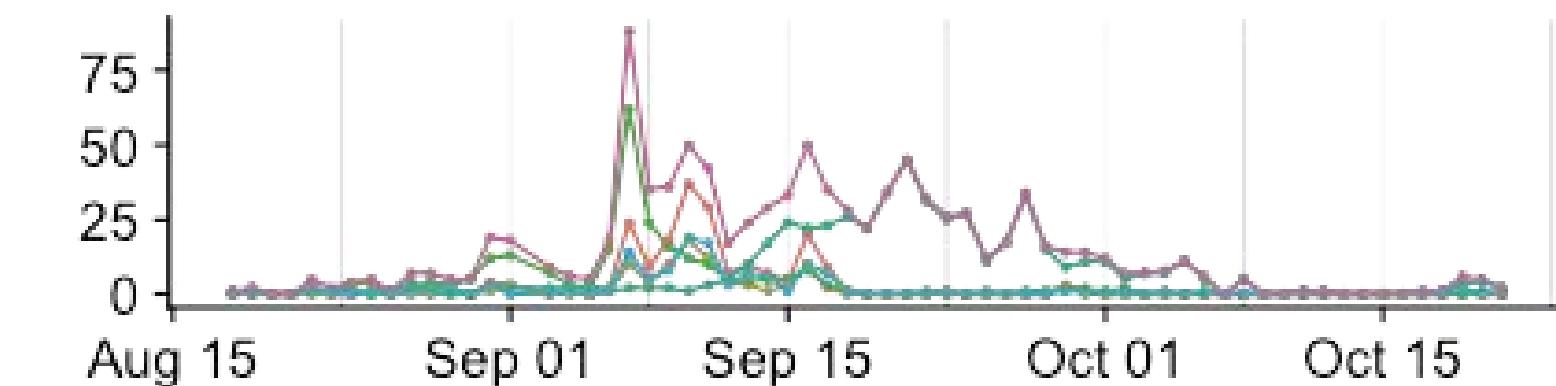
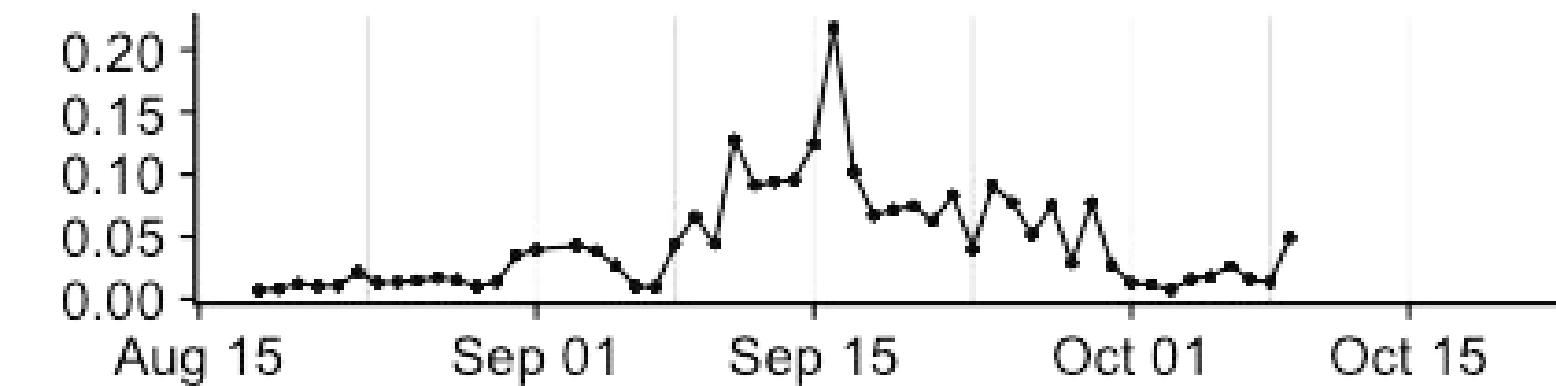
Spike-in
Soup 2
High
biomass



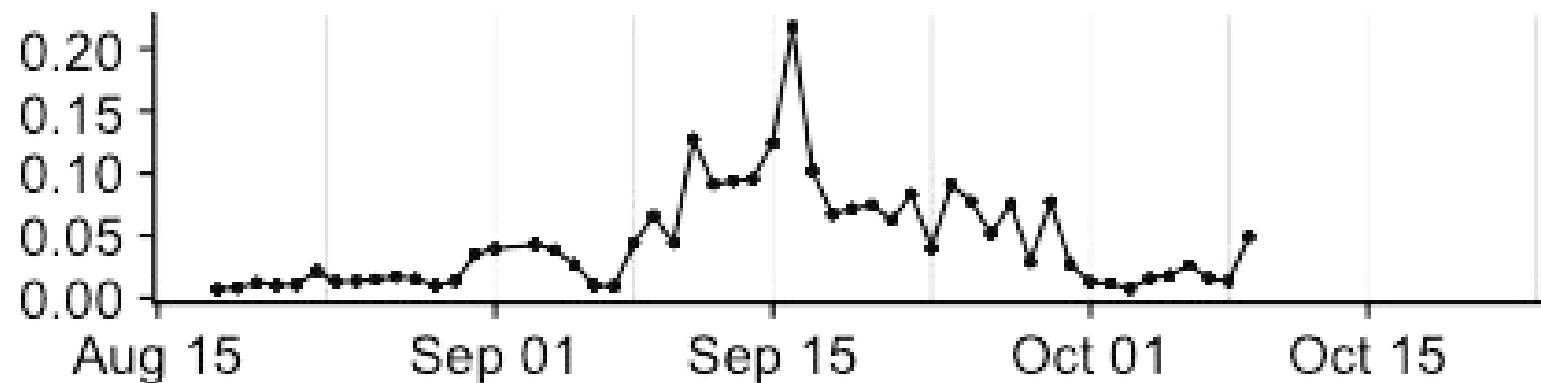
Sample 1	Reads_Raw	Reads_Corr	
OTU1	1000	25	$1000/40 = 25$
OTU2	1200	30	$1200/40 = 30$
OTU3	1600	40	$1600/40 = 40$
OTU4	2000	50	$2000/40 = 50$
OTU5	400	10	$400/40 = 10$
OTU6	200	5	$200/40 = 5$
OTU_SPIKE	40		
Sample 2	Reads_Raw	Reads_Corr	
OTU1	500	50	
OTU2	600	60	
OTU3	800	80	
OTU4	1000	100	
OTU5	200	20	
OTU6	100	10	
OTU_SPIKE	10		

Two kinds of quantification

- **Within-species quantification:**
“Species A is more abundant in these samples, less abundant in those samples.” Follow one species.
- *Across-species quantification:*
“Species A is more abundant than Species B in this sample” Compare two species.
- We tested quantification by making mock samples with 20 species, with different DNA quantities



How can within-species quantification be useful?



- **Within-species quantification:** “Species A is more abundant in these samples, less abundant in those samples.” Follow one species.

We can rescale within-species abundances to [0,1], known as quasiprobability

Potential importance of QP data, which preserves within-species frequency information

A. QP distribution data for one species

Species shows a preference for green habitat

0.1	1.0	0.1
0.1	0.7	
0.1	0.1	0.1

optimal env conditions

suboptimal env conditions

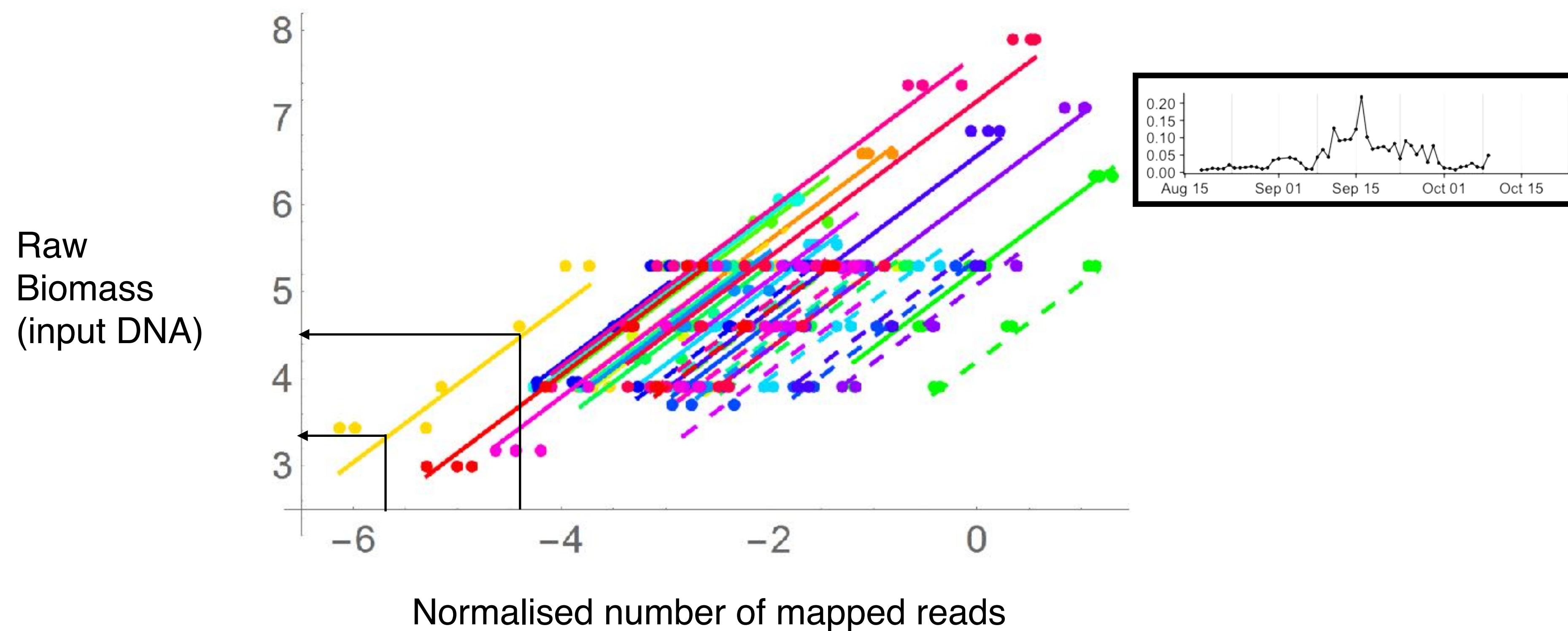
B. PA distribution data for one species

Species incorrectly shows no habitat preference

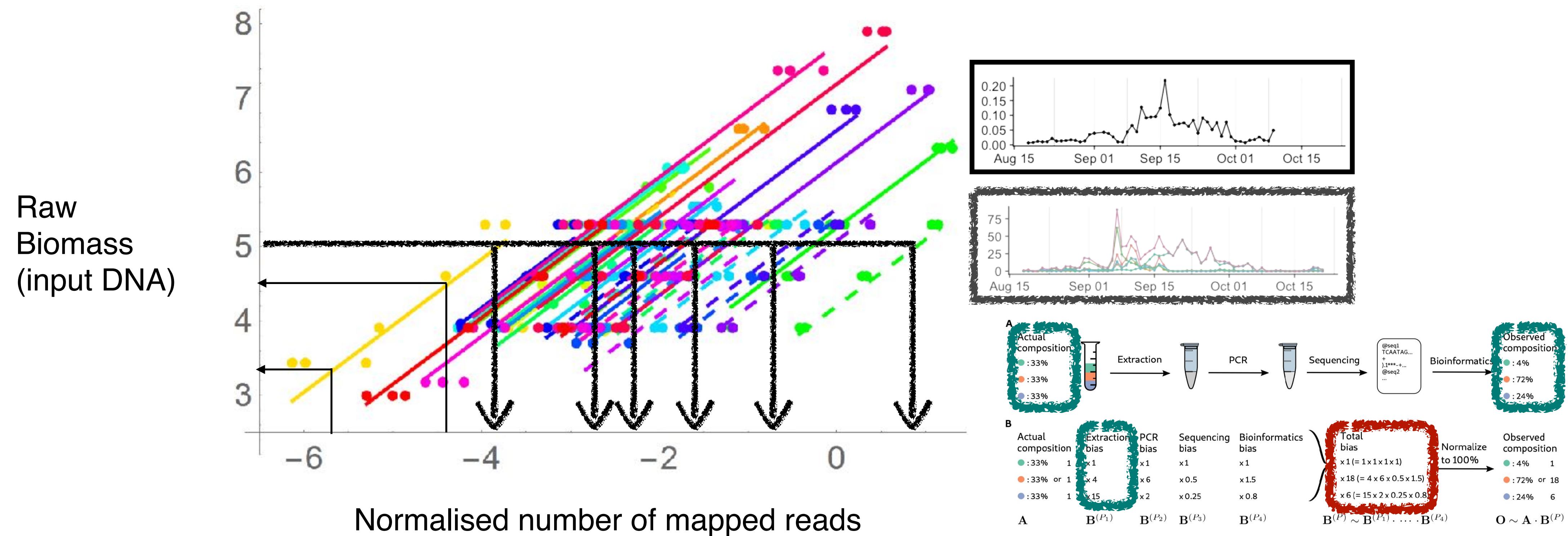
1	1	1
1	1	
1	1	1

Information loss with
Presence-Absence [0 or 1]

- Within-species quantification works well when you use a spike-in
 - Each colour is a different species. More mapped reads (x-axis) correctly predict more input DNA (y-axis)
 - Thus, we can observe how each species' population dynamics has changed from 1997 to 2013



- **Within-species quantification** works well. Each species is a different colour. More mapped reads (x-axis) correctly predicts more input DNA (y-axis). So we can observe how each species' population dynamics has changed from 1997 to 2013
- **Across-species quantification** does not work well. (we would need to have species correction factors)



	1. True OTU table					True rowSum
	OTU01	OTU02	OTU03	OTU04	spikeOTU	
Sample 1	10	20	0	0	20	30
Sample 2	0	100	20	50	20	170
Sample 3	40	40	5	50	20	135
Sample 4	60	0	30	100	20	190

Example from pa_vs_qp_tables_20210402.xlsx

Metabarcoding pipeline

$$500 = 10 * 5 * 10$$

2. Observed OTU table with row noise and species-specific biases

	OTU01	OTU02	OTU03	OTU04	spikeOTU	row noise
Sample 1	500	200	0	0	50	5
Sample 2	0	1000	300	1000	50	5
Sample 3	1200	240	45	600	30	3
Sample 4	600	0	90	400	10	1

species-specific biases

10 2 3 4 0.5

$$50 = 20 * 5 * 0.5$$

Laboratory stochasticity causes row noise
(e.g. from pooling amplicons for library prep)

Species-specific biases are caused by species properties

Actual
composition

● : 33% 1

● : 33% or 1

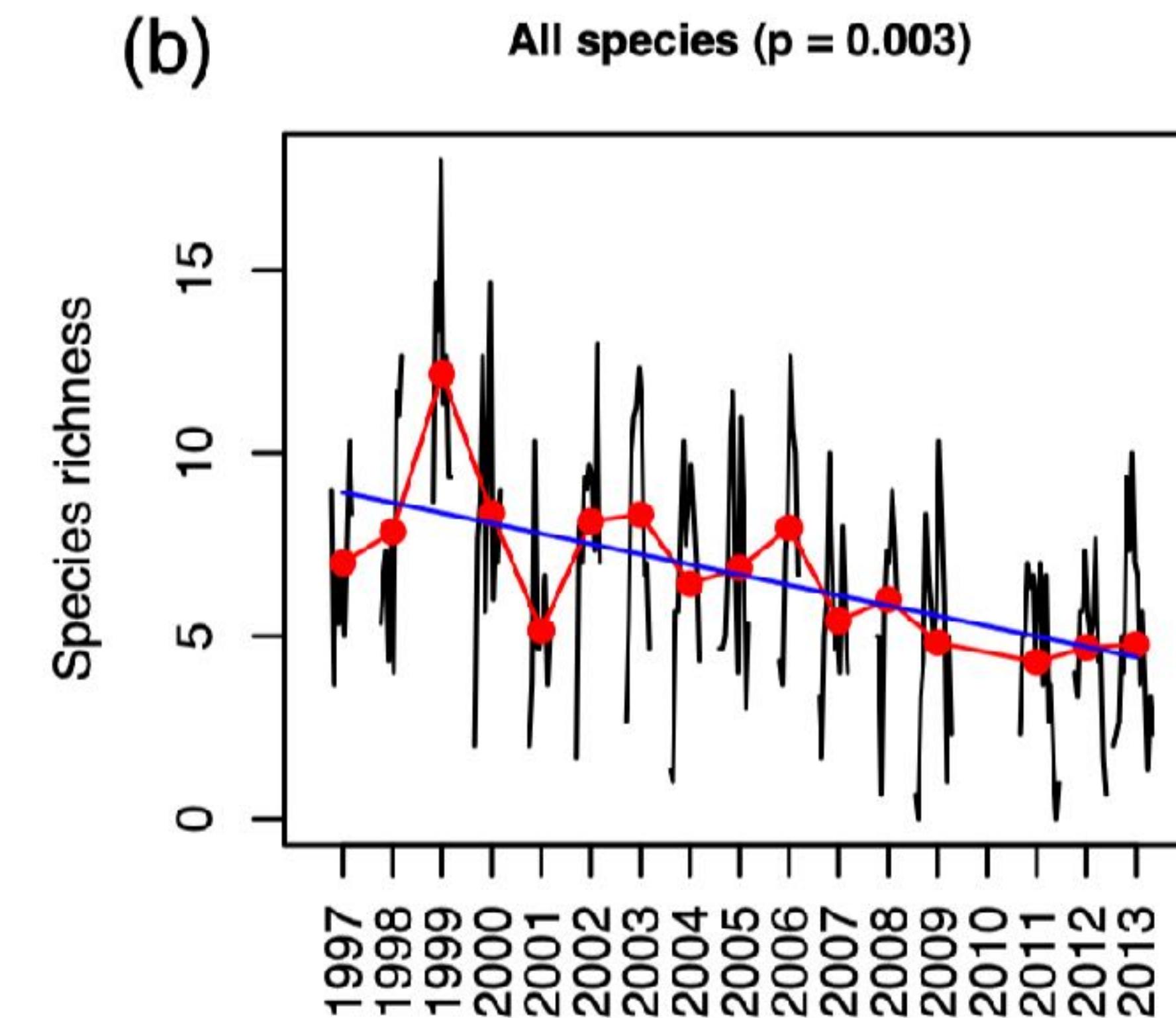
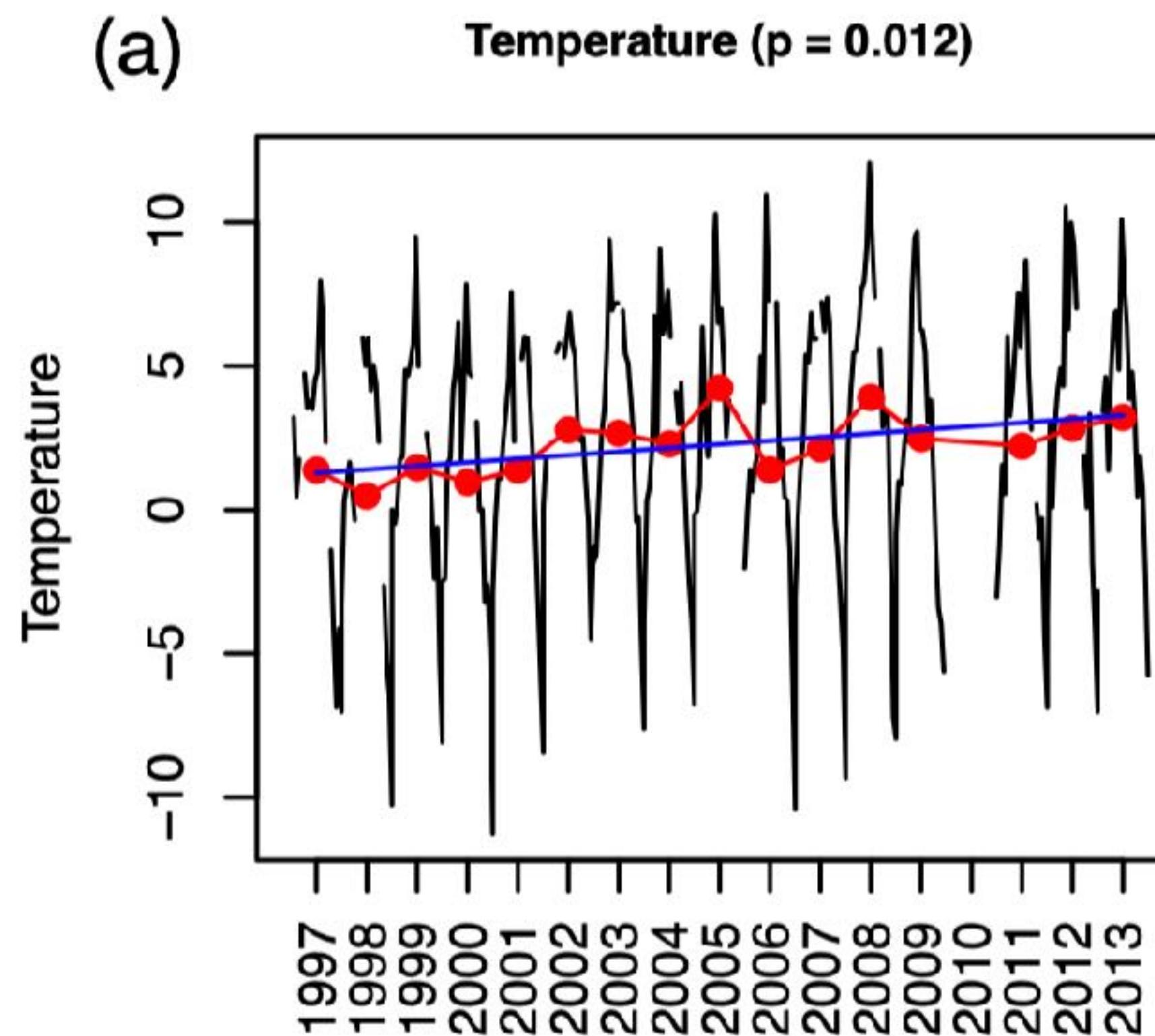
● : 33% 1

Extraction bias	x1
PCR bias	x1
Sequencing bias	x1
Bioinformatics bias	x1

Accounting for species interactions is necessary for predicting how arctic arthropod communities respond to climate change

Nerea Abrego, Tomas Roslin, Tea Huotari, Yinqi Ji, Niels Martin Schmidt, Jiaxin Wang, Douglas W. Yu and Otso Ovaskainen

Ecography
44: 1–12, 2021
doi: 10.1111/ecog.05547

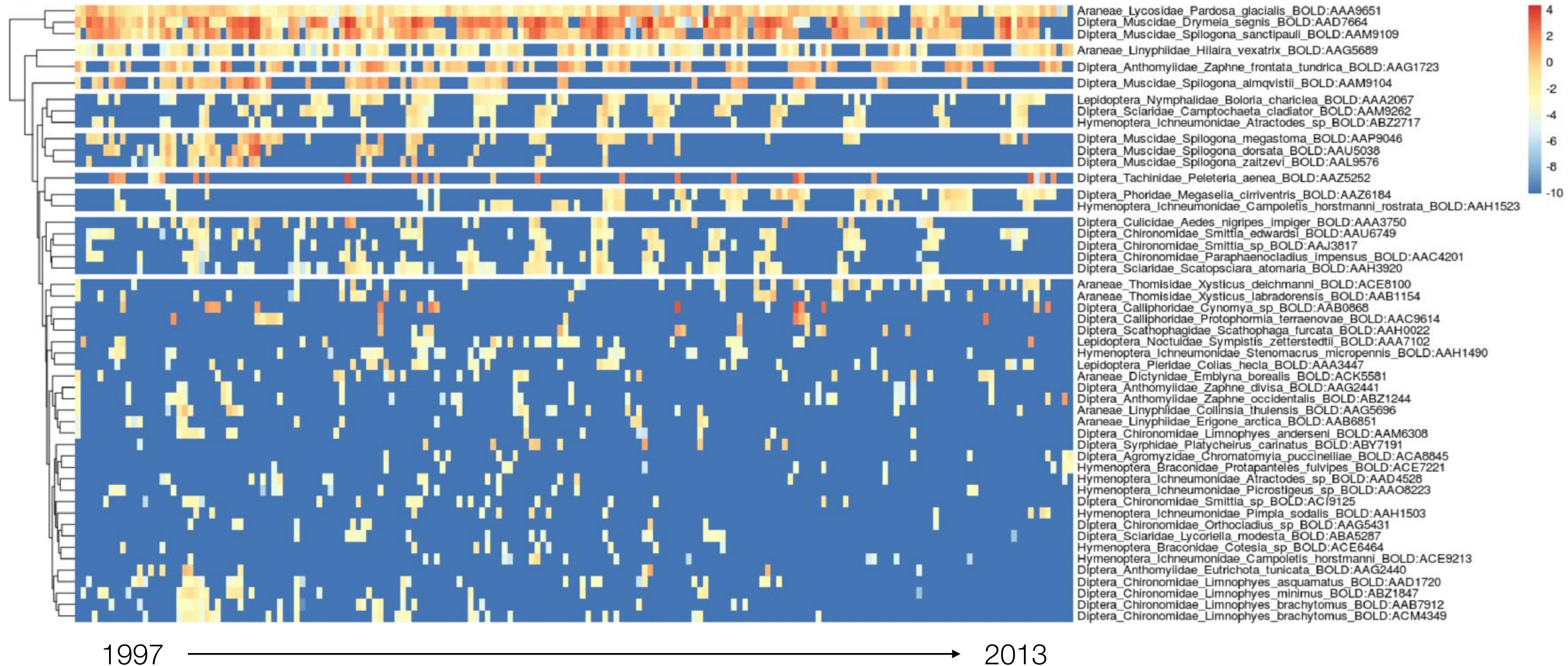


52 species in
542 trap-week samples

average summer temperature has increased by 2.0°C (from 1.28°C to 3.28°C)

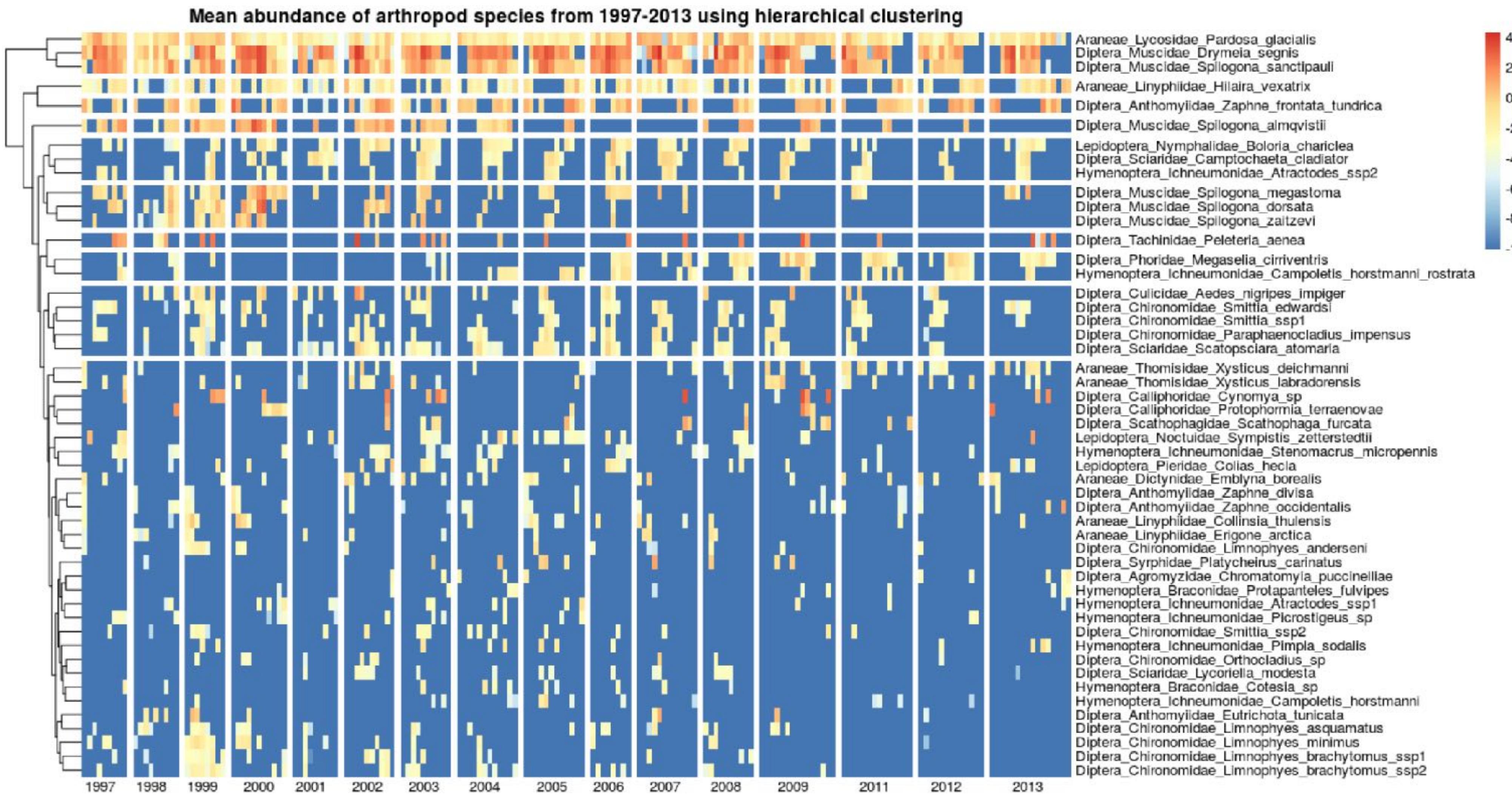
(observed) species richness has declined by half (8.9 to 4.4 species per trap-week, decline is mostly in the predator species)

What we get is a ***community time series***
(red = highest **within-species** abundance, blue = absence)



more blue cells (more absences) toward 2013

What we get is a ***community time series***
 (red = highest **within-species** abundance, blue = absence)

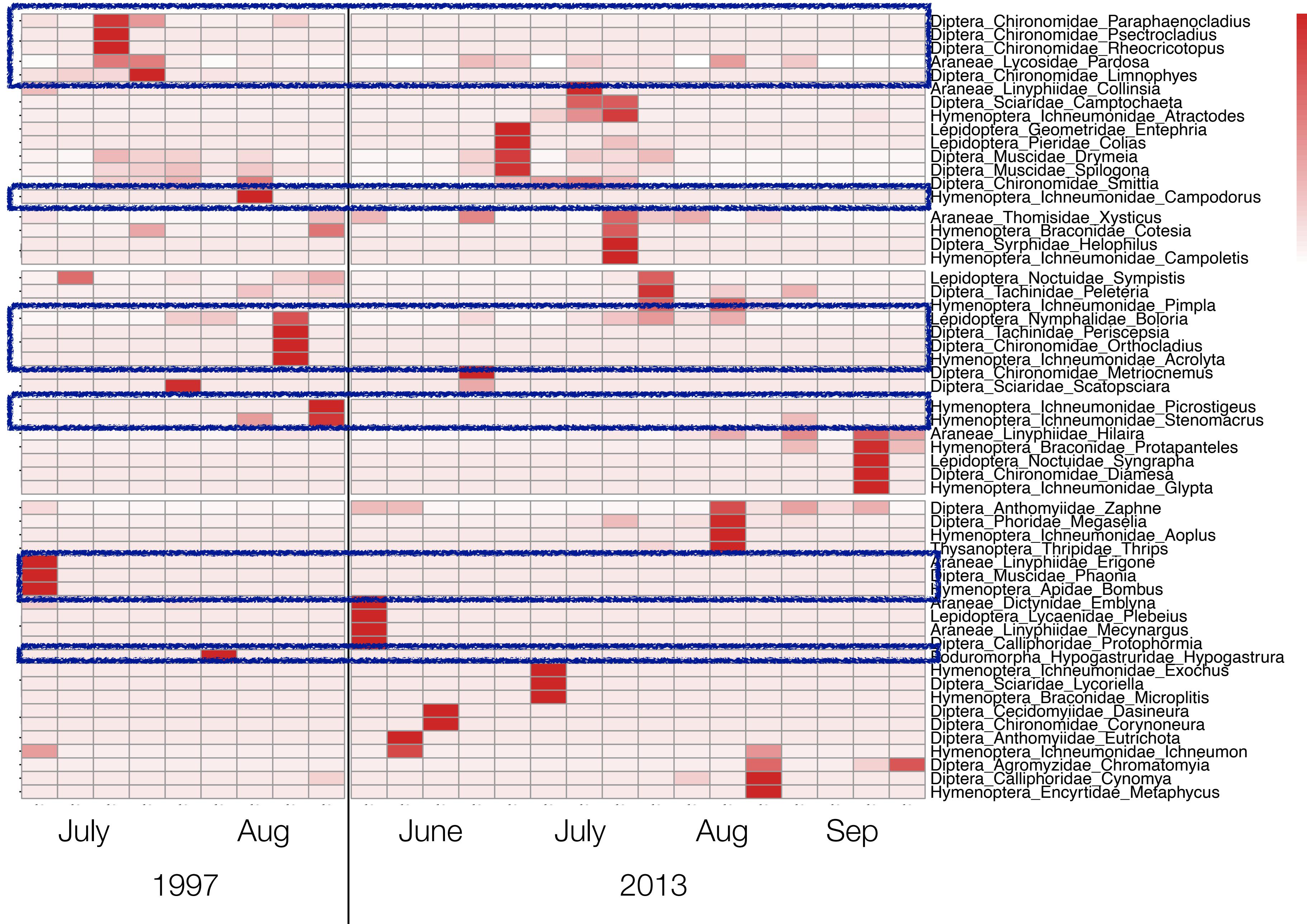


Summer has approx doubled in length in 16 years

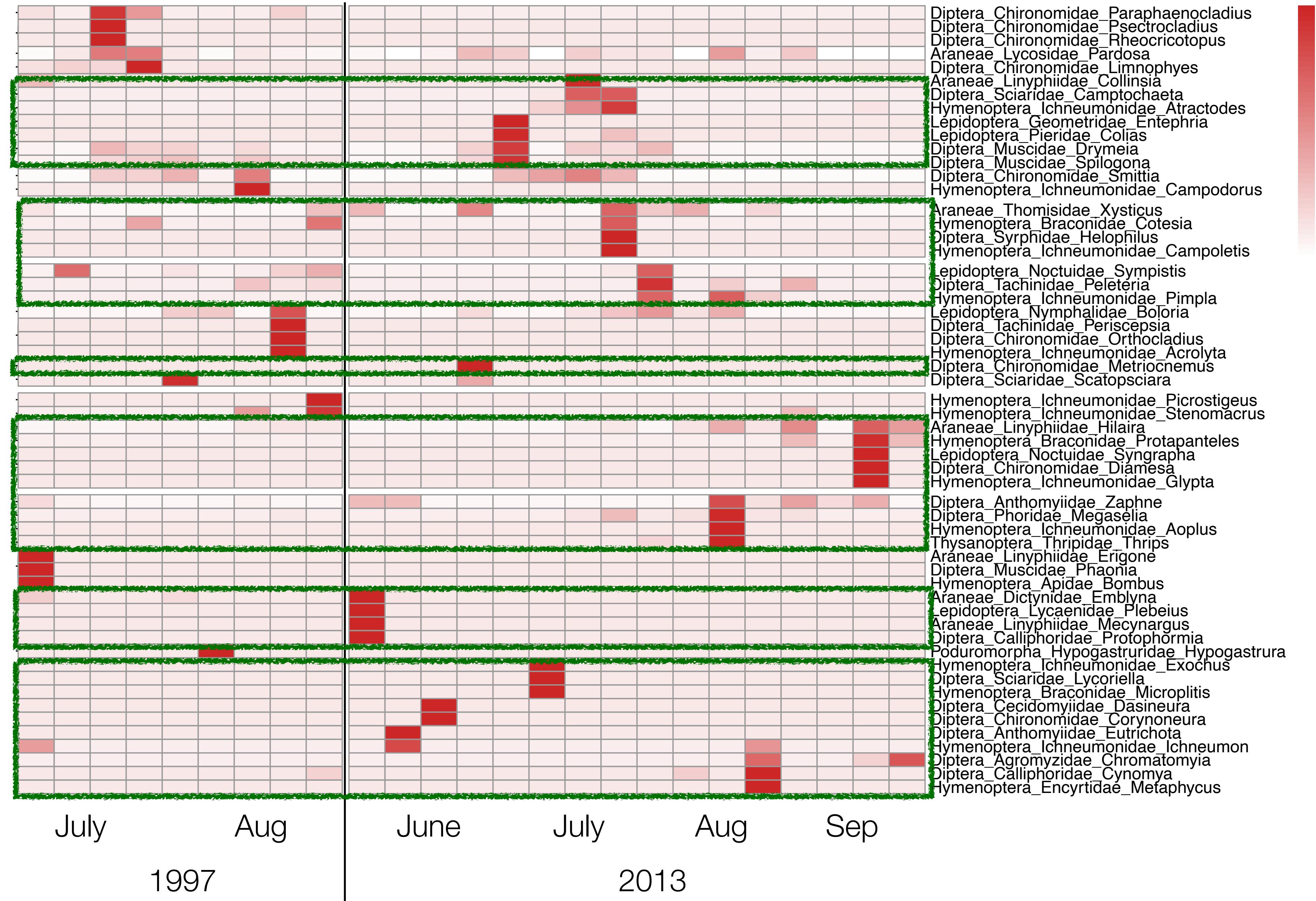
Let's look at 2 years of real data (dark red = more biomass of that species)



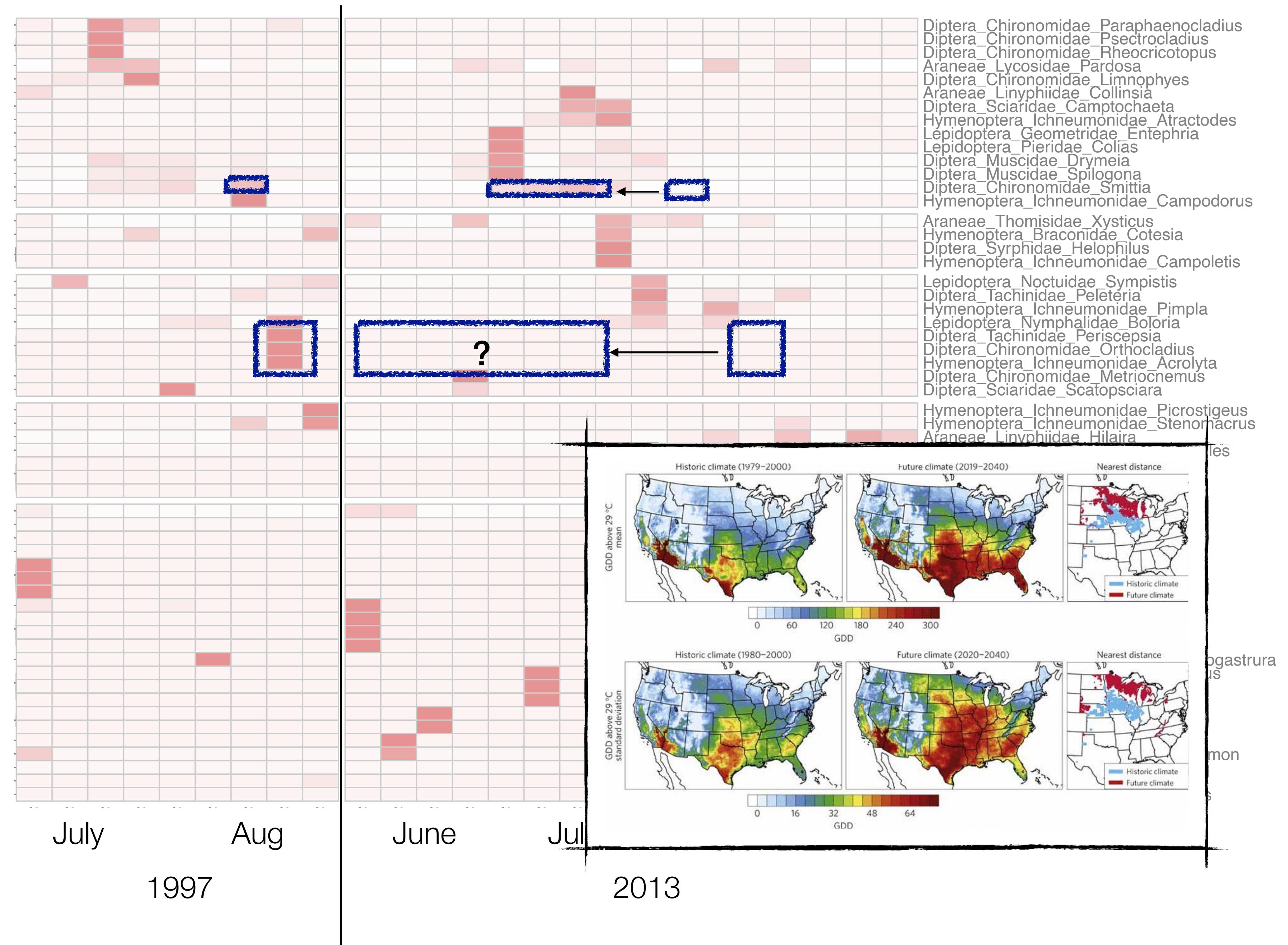
Many species that exhibited activity peaks in 1997, now no longer do



Many species that showed no activity peaks in 1997, now do have activity peaks!



We rarely see simple shifts in phenology (timing of peaks): the climate envelope idea

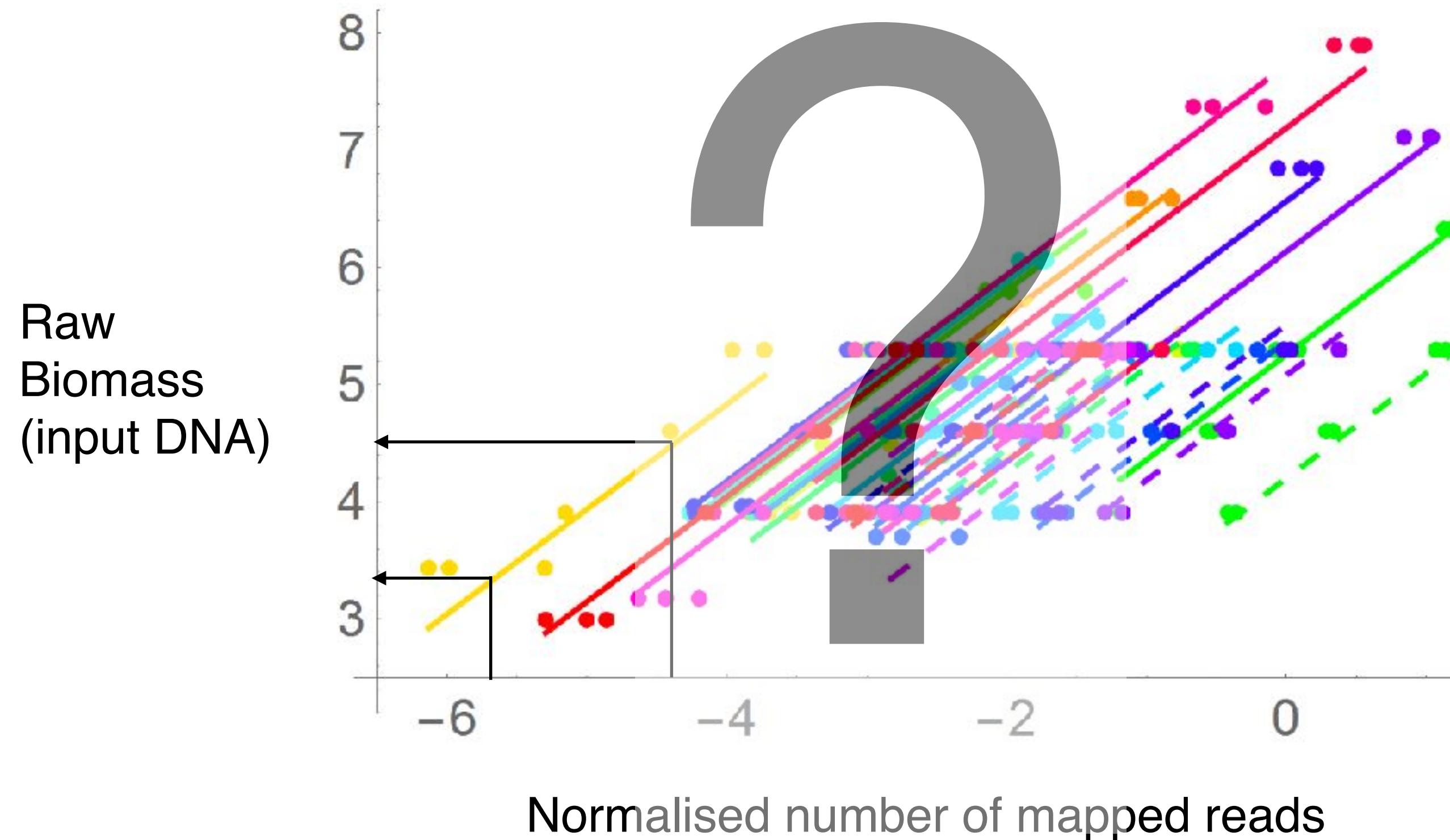
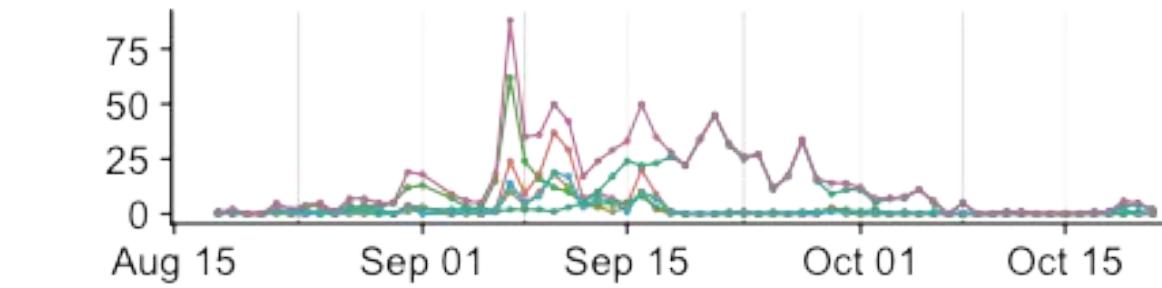
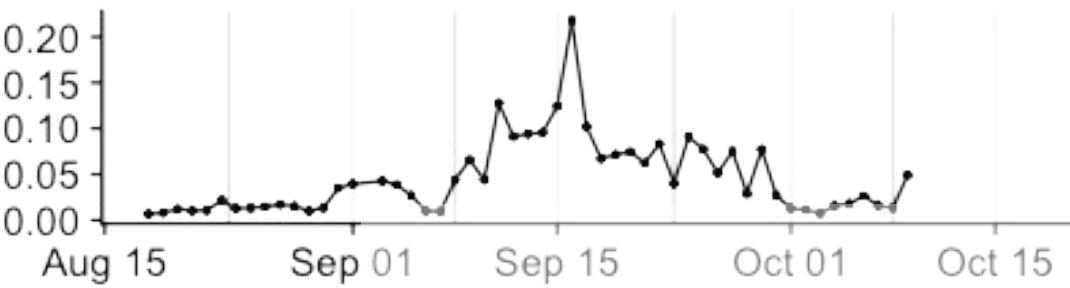


Methods to extract abundance information from DNA data

- Single-species quantitative PCR (qPCR)
- Multiplexed individual barcoding (mBRAVE)
- Mitogenomics and DNA spike-in (SPIKEPIPE)
- **Metabarcoding and DNA spike-in (qSeq)**
- Reverse metagenomics (RevMet)

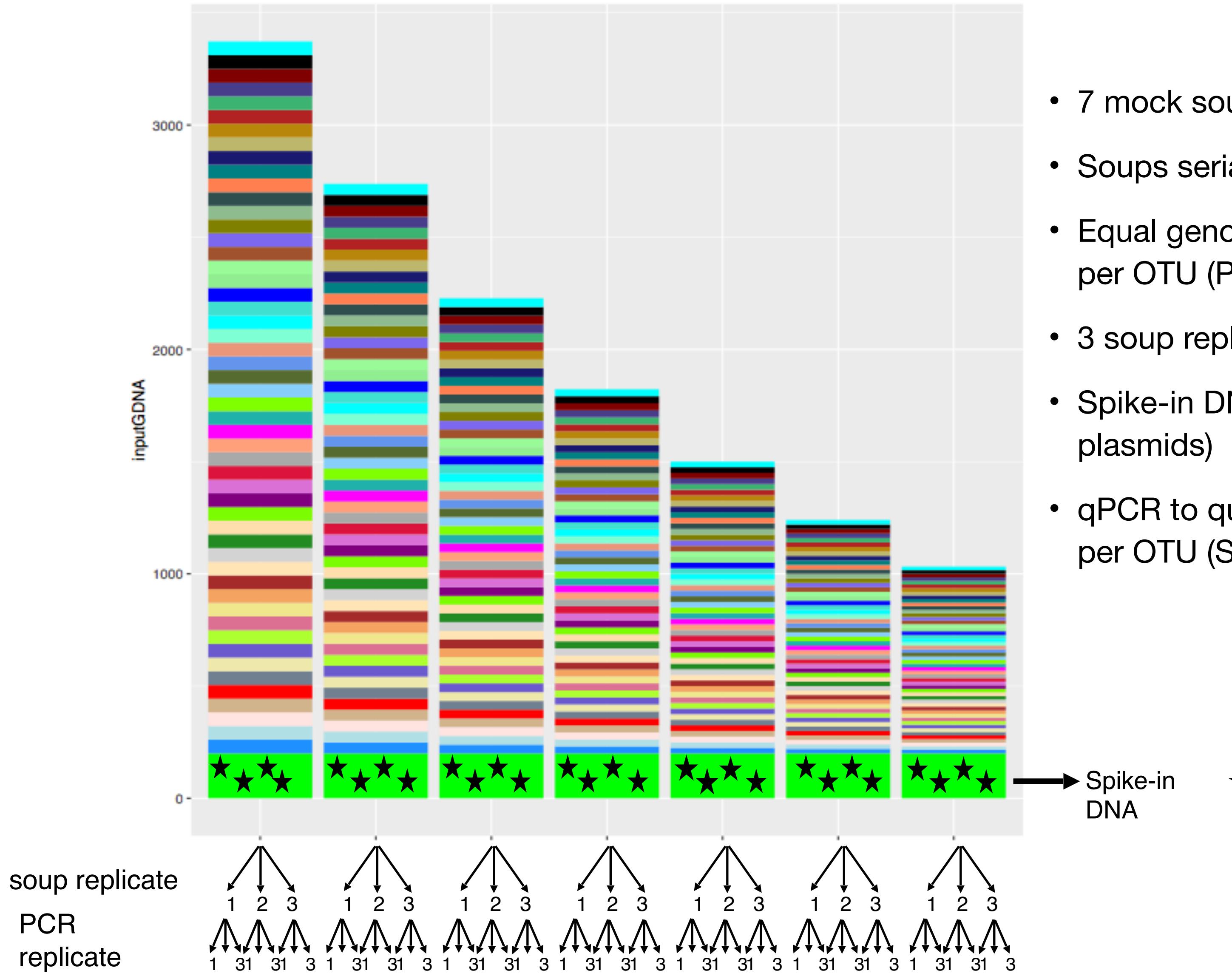
Is quantitative metabarcoding possible (qSeq)?

Mingjie Luo, Yinqi Ji, Yuanheng Li, Douglas Yu



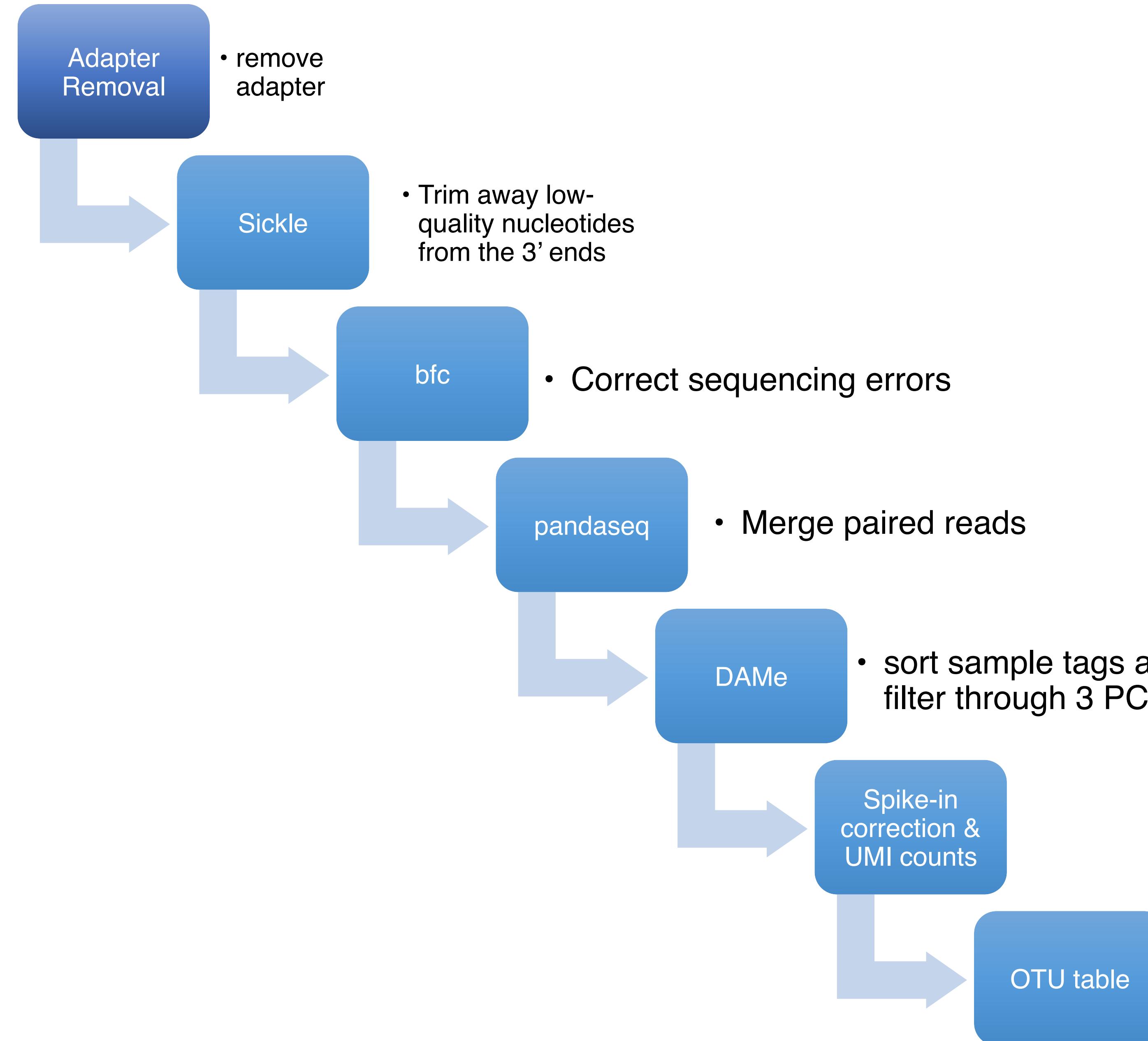
With metabarcoding, we could avoid the initial cost of creating a mitogenome dataset

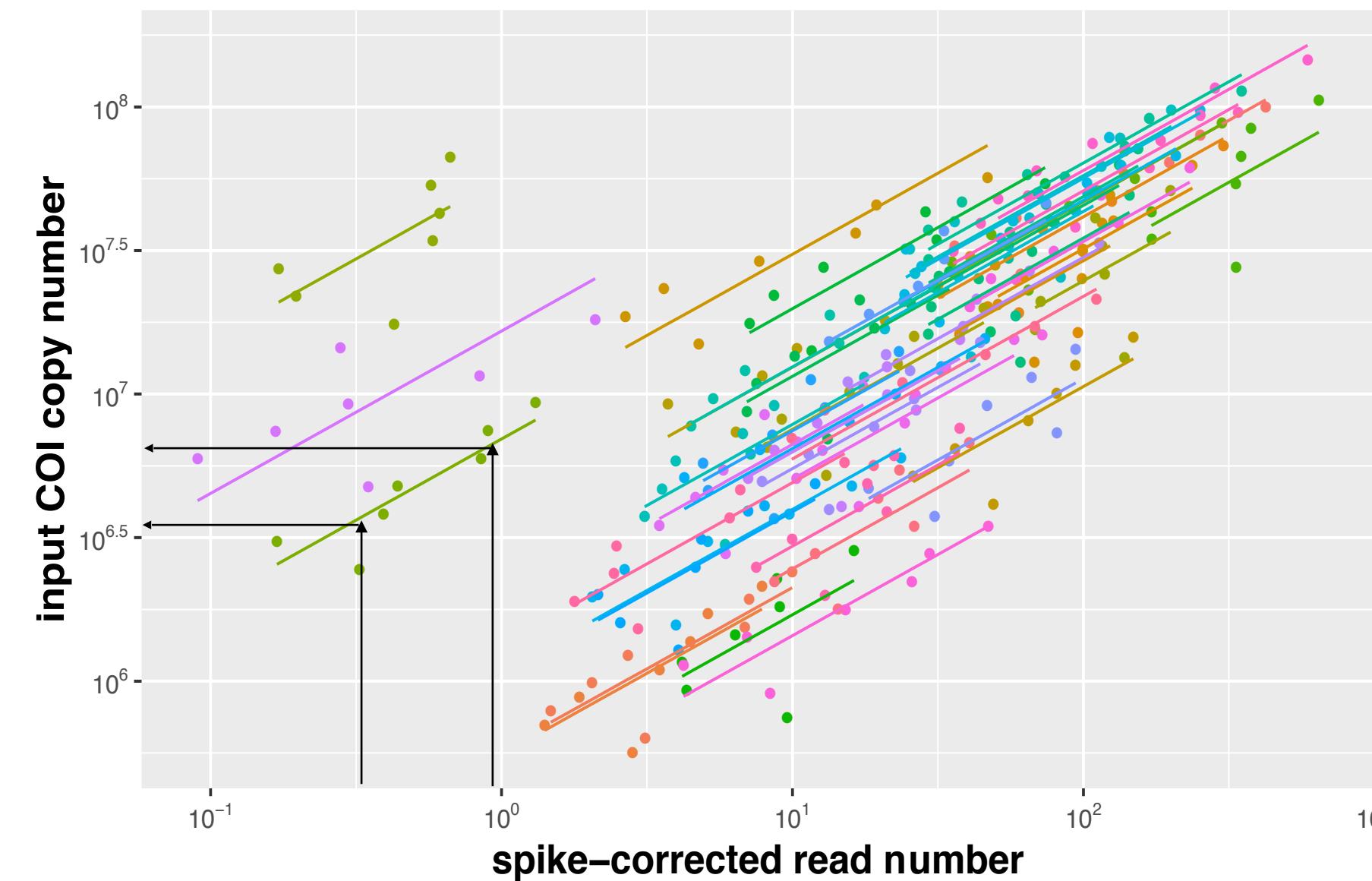
qSeq Experiment



- 7 mock soups, 52 insect OTUs.
- Soups serially diluted by 0.8X
- Equal genomic DNA concentration per OTU (PicoGreen)
- 3 soup replicates, 3 PCR replicates.
- Spike-in DNA (3 DNA barcodes in plasmids)
- qPCR to quantify COI copy number per OTU (SYBR green)

Bioinformatic Pipeline



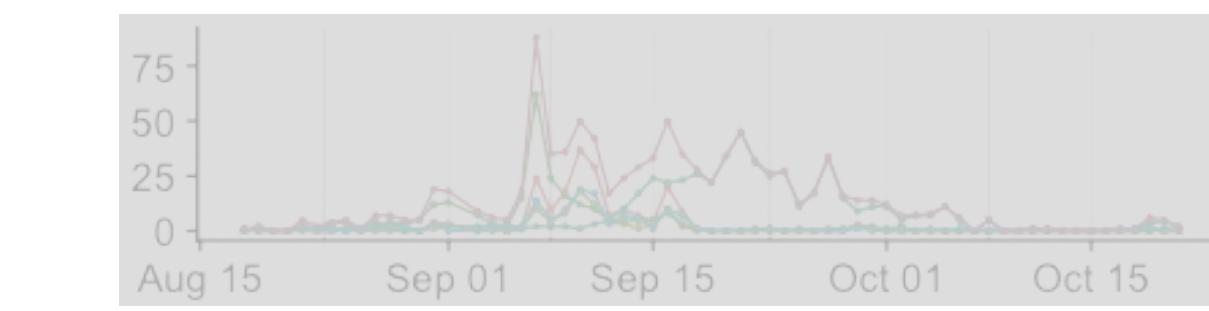
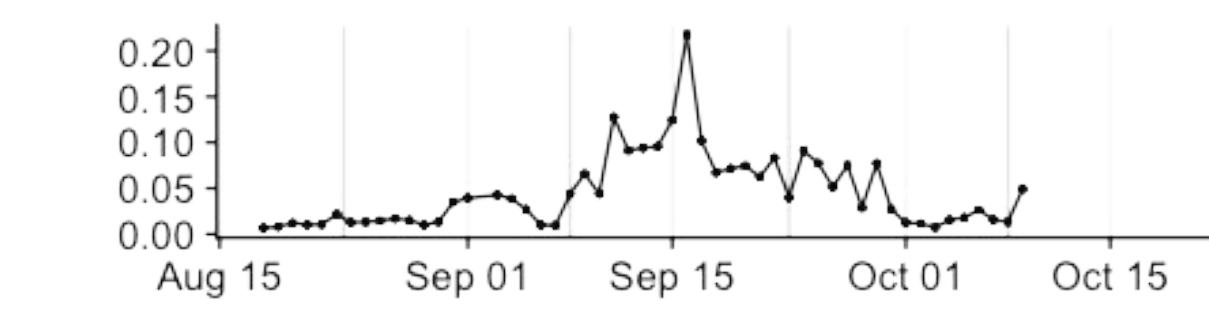
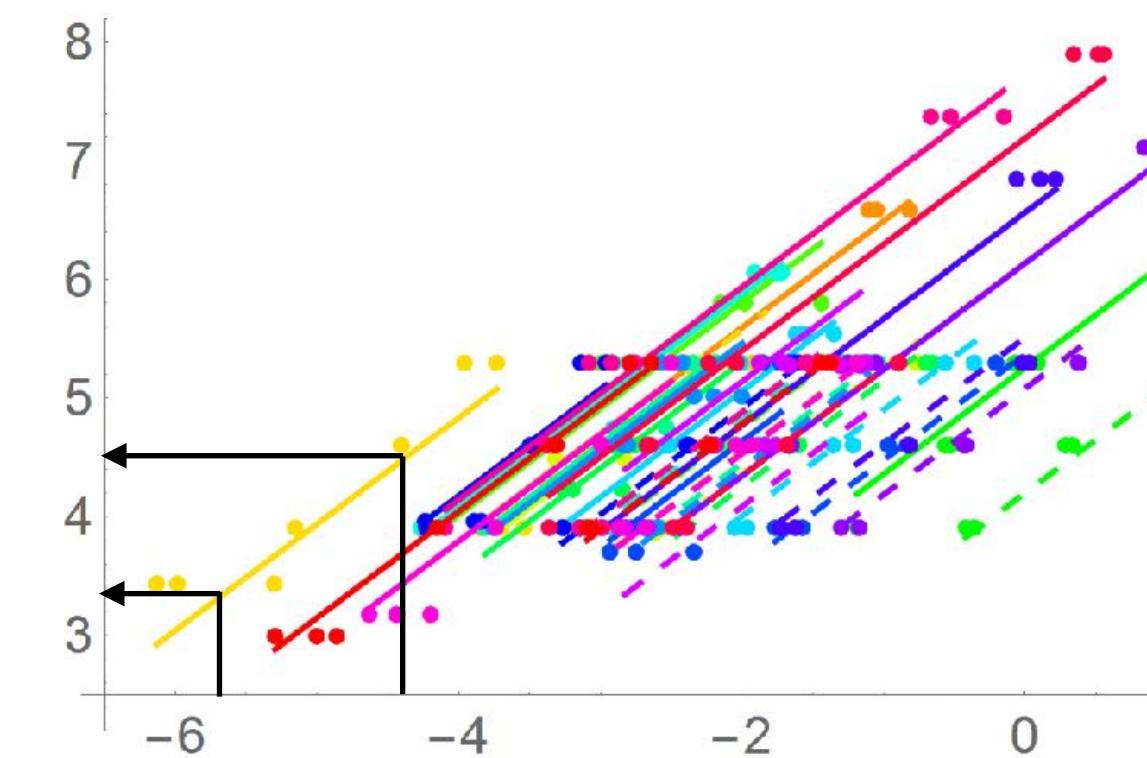
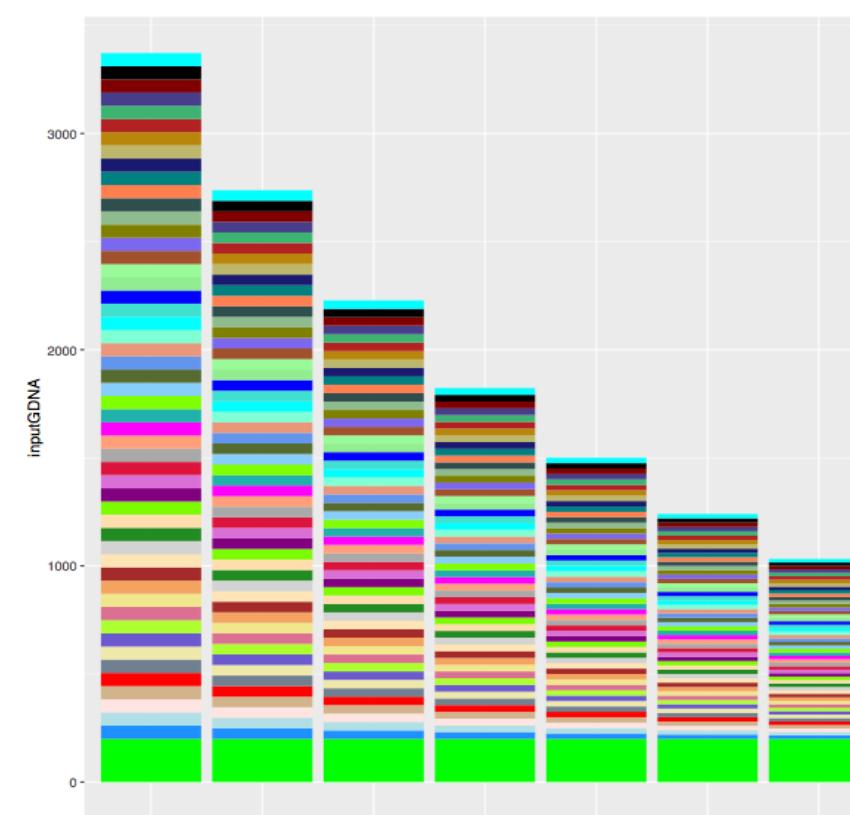


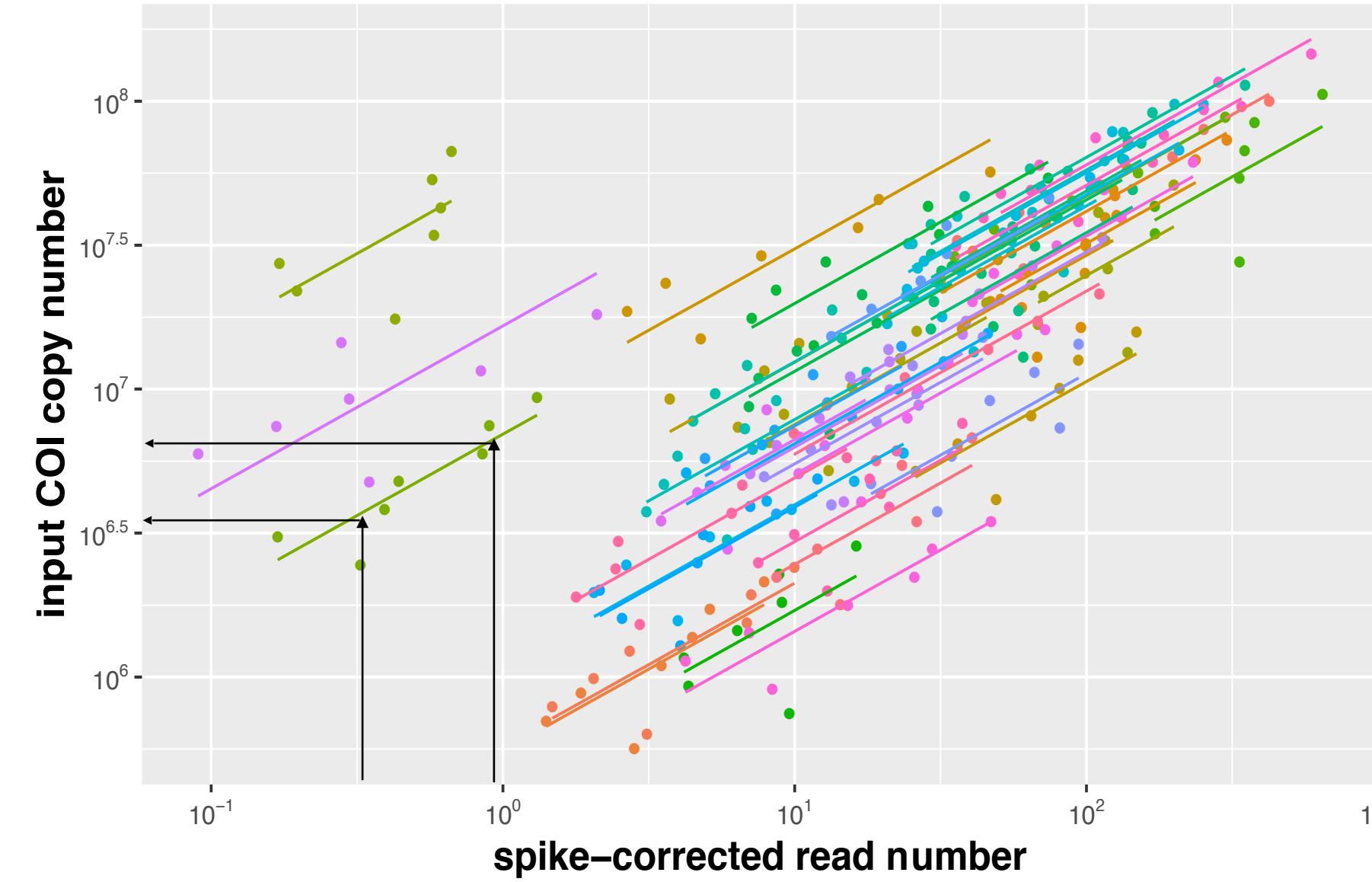
It works!

`input_COI_copy_num ~ readNum_spikeCorrected + otuID`

slope= 0.56, R²=0.95 (soup and PCR replicates averaged)

As with mitogenomics, we can estimate **within-species** change in abundance (DNA mass)



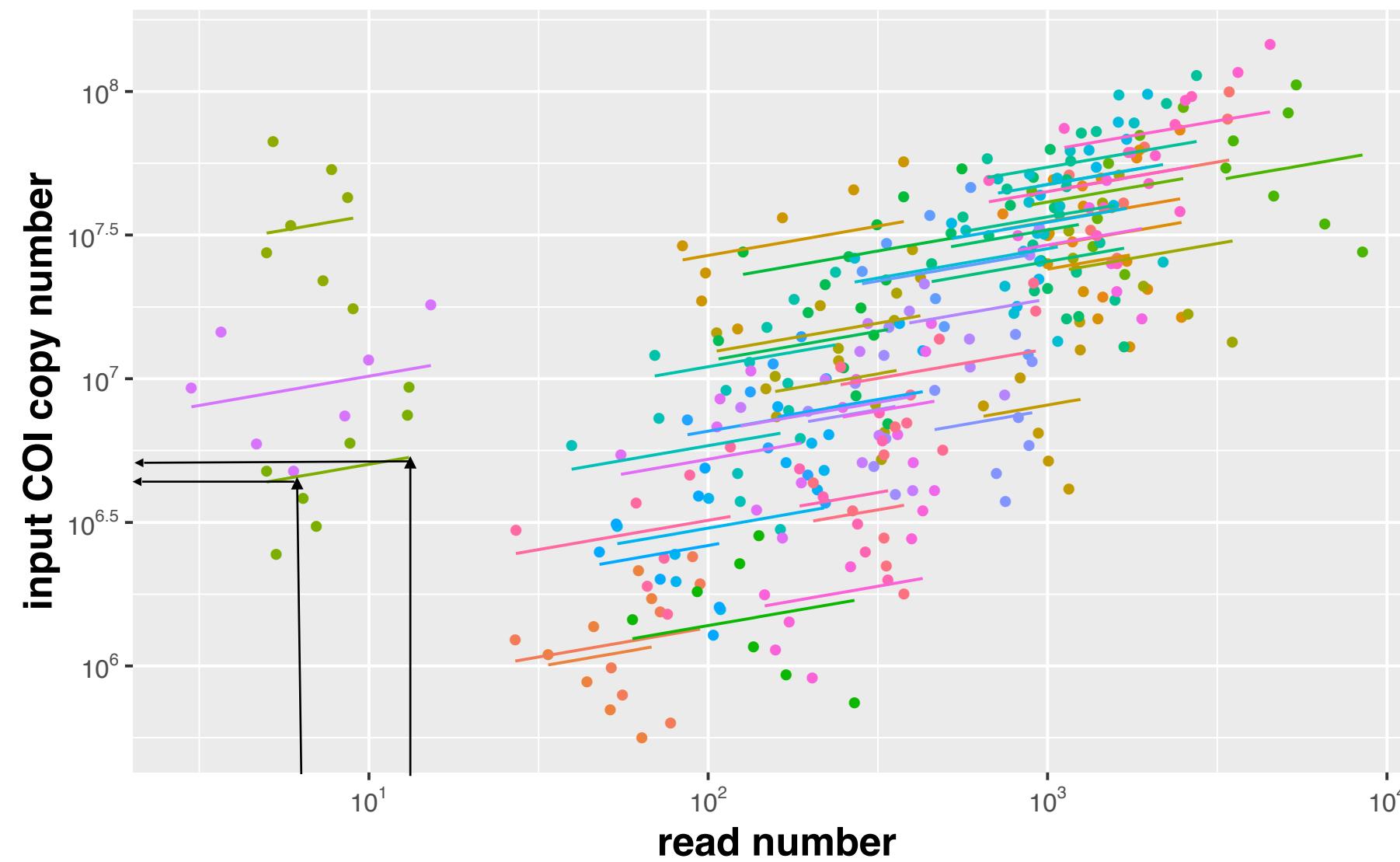


It works!

$\text{input_COI_copy_num} \sim \text{readNum_spikeCorrected} + \text{otuID}$

slope= 0.56, R²=0.95 (soup and PCR replicates averaged)

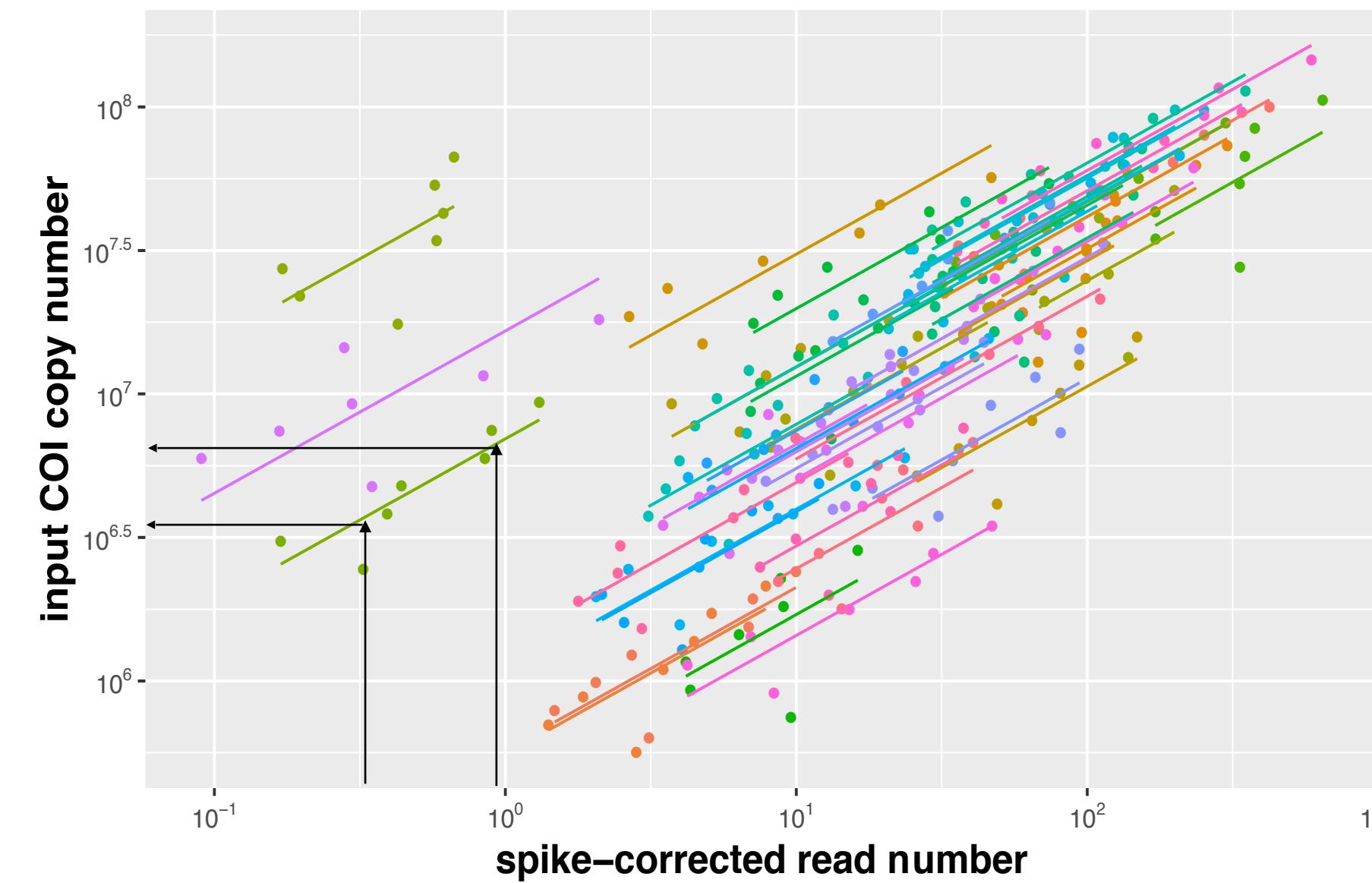
We can estimate **within-species** change in abundance (COI copy number)



Spike correction is necessary: increases slope and R²

$\text{input_COI_copy_num} \sim \text{readNum} + \text{otuID}$

slope= 0.2, R²=0.85 (soup and PCR replicates averaged)



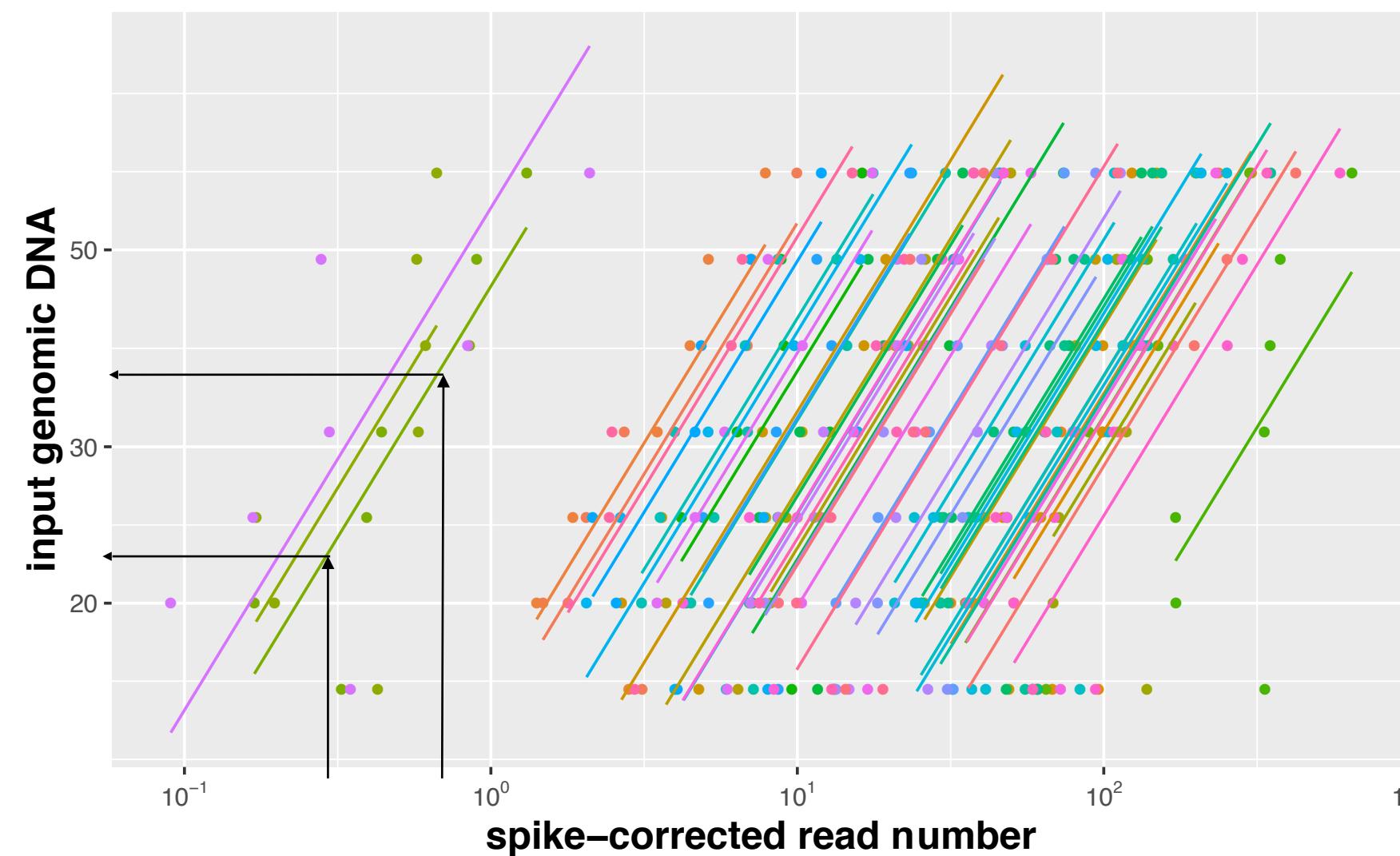
It works!

`input_COI_copy_num ~ readNum_spikeCorrected + otuID`

slope= 0.56, R²=0.95 (soup and PCR replicates averaged)

We can estimate **within-species** change in abundance (COI copy number)

Spike correction is necessary: increases slope and R²

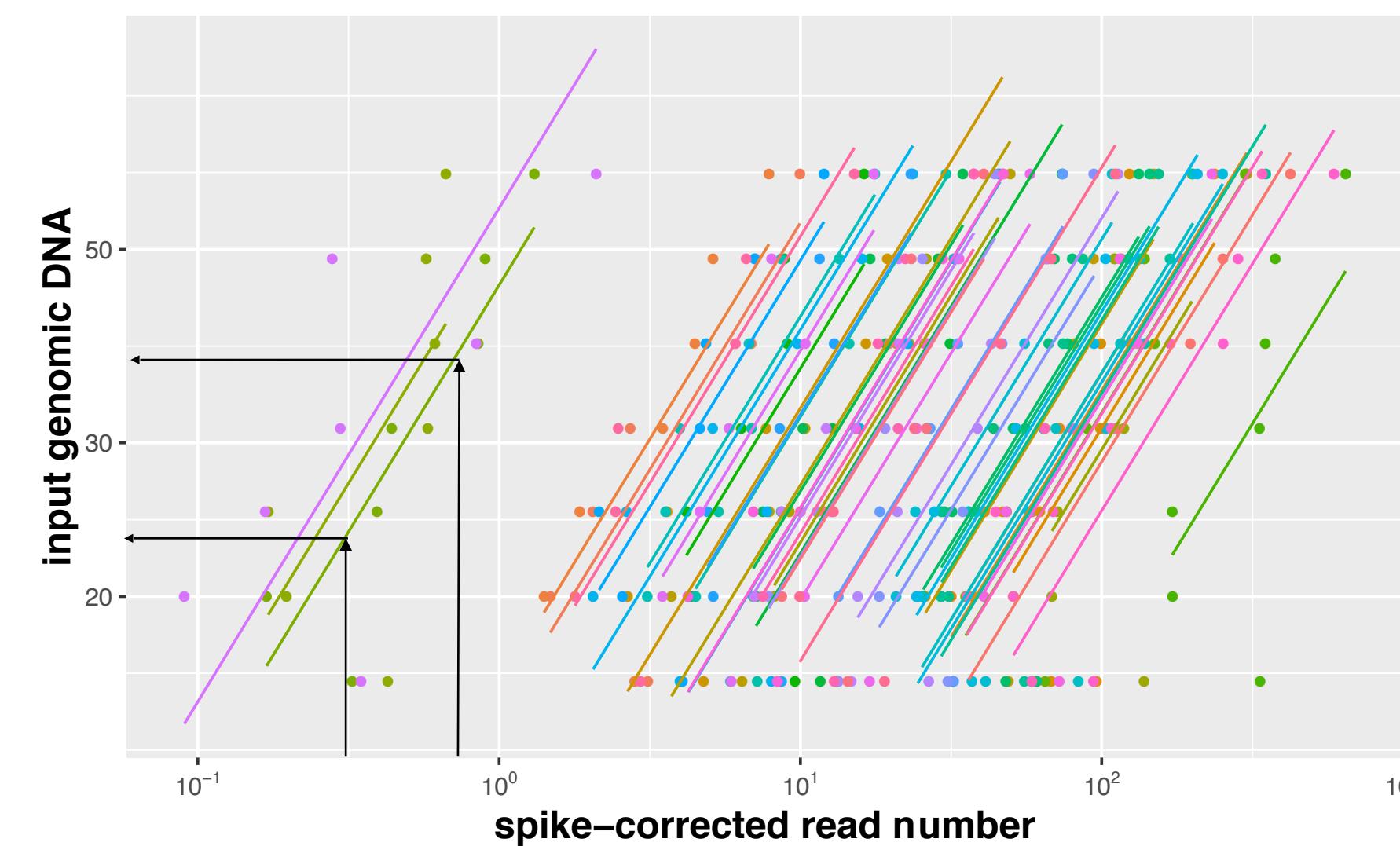


We can also predict input genomic DNA (our real goal)

`input_gDNA ~ readNum_SpikeCorr + otuID`

slope= 0.57, R²=0.62 (soup and PCR replicates averaged)

Note that the y-axis scale is different (ng, not copy number)



It works!

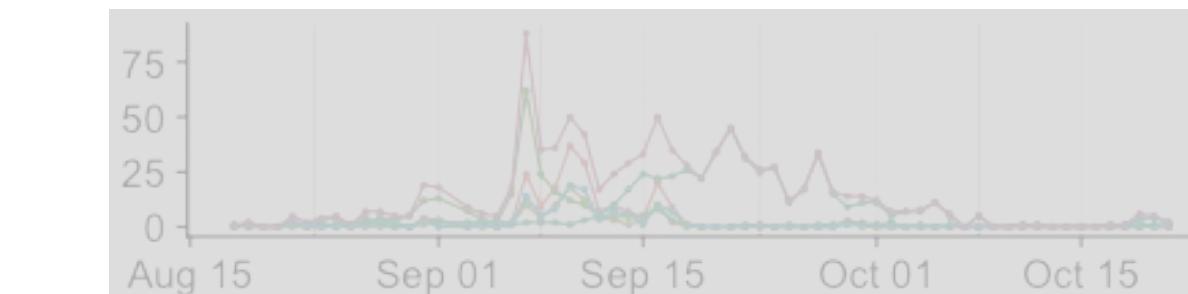
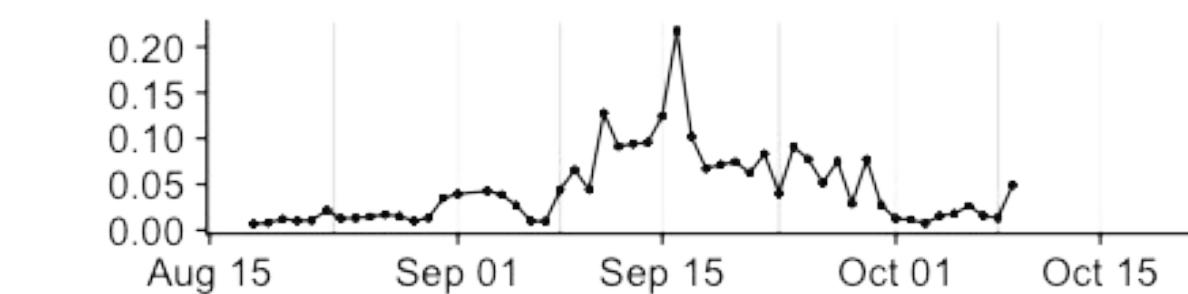
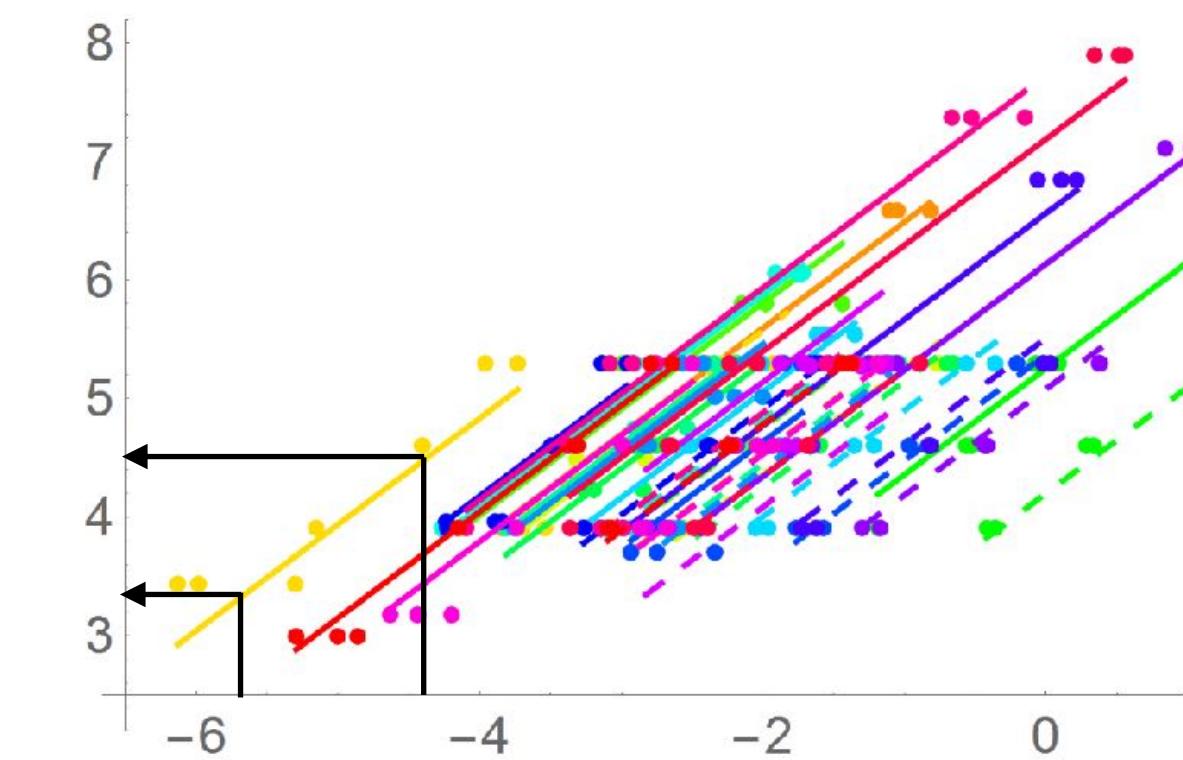
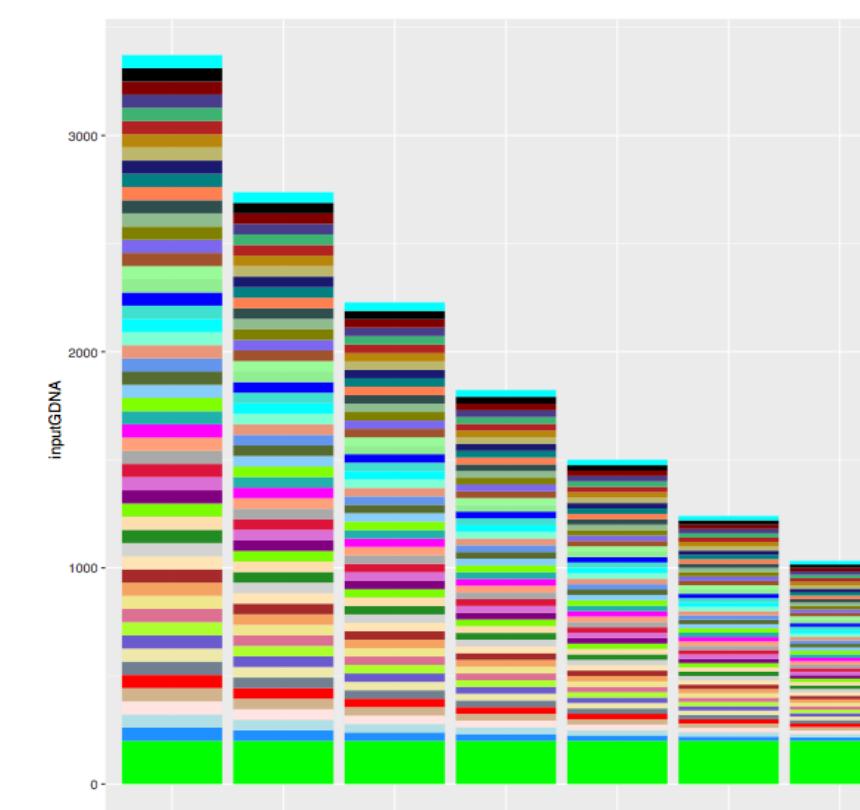
input_gDNA ~ readNum_SpikeCorr + otuID

slope= 0.57, R²=0.62 (soup and PCR replicates averaged)

We can estimate **within-species** change in abundance (as genomic DNA == our real goal)

Spike correction is necessary: increases slope and R²

We can indeed use metabarcoding to estimate within-species change in abundance (gDNA mass)



Methods to extract abundance information from DNA data

- Single-species quantitative PCR (qPCR)
- Multiplexed individual barcoding (mBRAVE)
- Mitogenomics and DNA spike-in (SPIKEPIPE)
- Metabarcoding and DNA spike-in (qSeq)
- **Reverse metagenomics (RevMet)**



We want to identify pollen to species **with across-species quantification** (which pollen species is more abundant in the honeybee's diet?)



DOI: 10.1111/2041-210X.13265

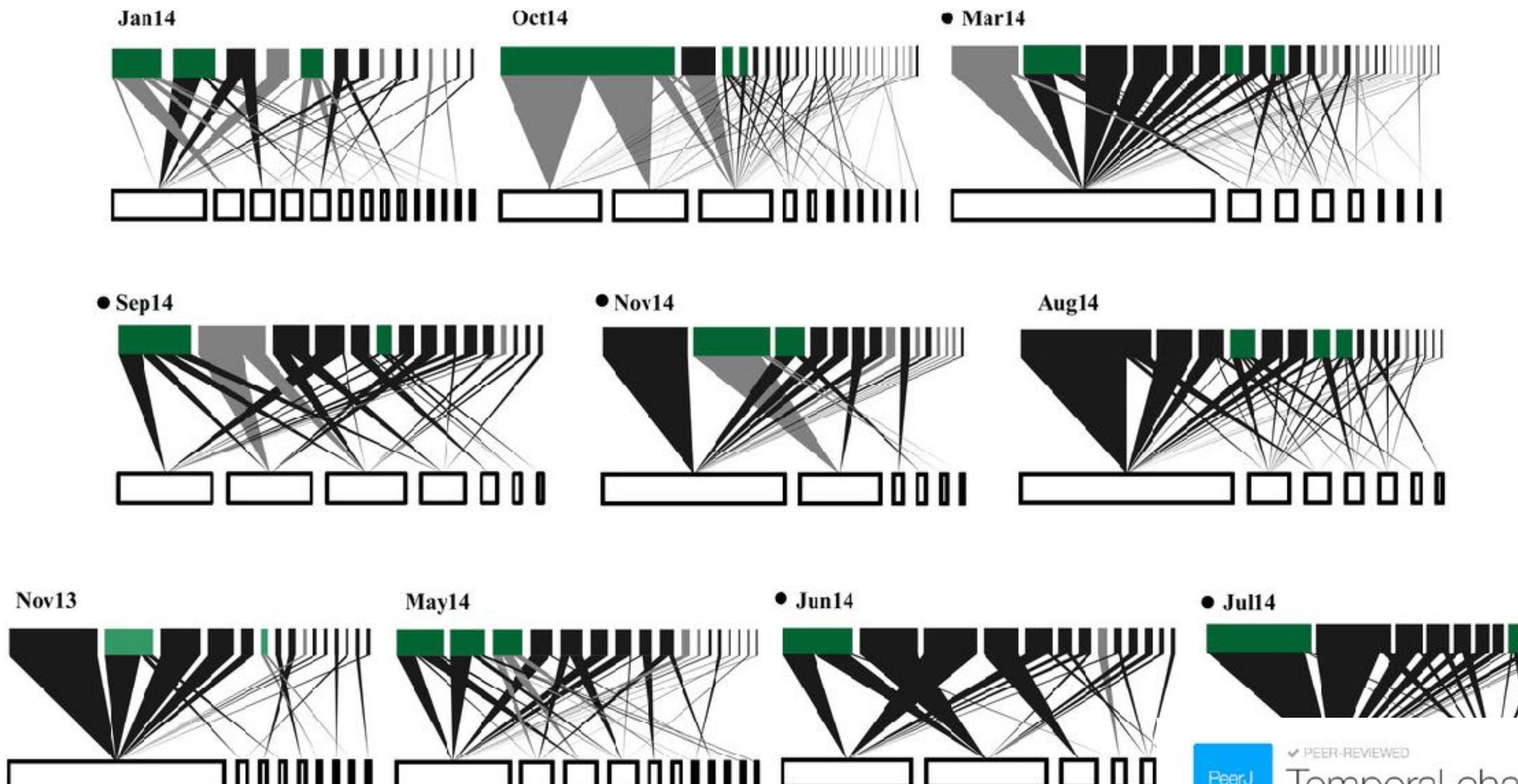
RESEARCH ARTICLE

Methods in Ecology and Evolution  BRITISH
ECOLOGICAL
SOCIETY

Semi-quantitative characterisation of mixed pollen samples using MinION sequencing and Reverse Metagenomics (RevMet)

Ned Peel^{1,2} | Lynn V. Dicks² | Matthew D. Clark^{1,3} | Darren Heavens¹ |
Lawrence Percival-Alwyn¹ | Chris Cooper⁴ | Richard G. Davies² | Richard M. Leggett¹ |
Douglas W. Yu^{2,5,6} 

Basic: Construct quantitative pollinator-plant networks



PeerJ PEER-REVIEWED
Temporal changes in the structure
of a plant-frugivore network are
influenced by bird migration and fruit
availability

Related
research

Biodiversity Biogeography Conservation Biology Ecology Zoology

Michelle Ramos-Robles¹, Ellen Andrensen², Cecilia Díaz-Castelazo¹

Published June 8, 2016 PubMed 27330852

Applied: Diet analysis to identify preferred plant species for bees: **need across-species quantification**



Crop plants require a pulse of pollination, which requires a large wild bee population



Preferred plants differ across bee species and might differ over time



Our goal – a protocol that:

- Does not need specialist molecular-lab skills
- Can be used on small amounts of pollen
- Uses the **nuclear genome**, not (just) the chloroplast genome
- Does **NOT require assembled genomes** for the reference database
- Is '**semi-quantitative**' (low vs. high biomass frequency)
- Can **identify pollen to plant species**

Plastid:Nuclear DNA ratio
might vary across pollen
grains and plant species

Ease of building references

We want the **dominant** pollen

MinION sequencer

- **long reads** (1000s of bp per read)
- ~ 85-95% accuracy rate
- low-input DNA kit (200 ng)



SQK-RBK001
Rapid Barcoding Kit
\$672.00

1 [Add to cart](#)

Simple and rapid library preparation, with barcoding for up to 12 gDNA samples

gDNA | 200 ng | <10 min | No PCR | Speed / simplicity

[More information](#)

A screenshot of a product listing for the SQK-RBK001 Rapid Barcoding Kit. The listing includes the product name, price (\$672.00), quantity selector, an 'Add to cart' button, a brief description of the kit's purpose, and several color-coded tags: gDNA (purple), 200 ng (orange), <10 min (teal), No PCR (blue), and Speed / simplicity (yellow). A 'More information' button is at the bottom.

MinION metagenomics



Query ('question') dataset

Sample database: each pollen sample is sequenced without PCR



Pollen sample

MinION sequencing



Reference dataset

Reference database:
each species is a 0.5X
'genome skim': a cloud of short reads

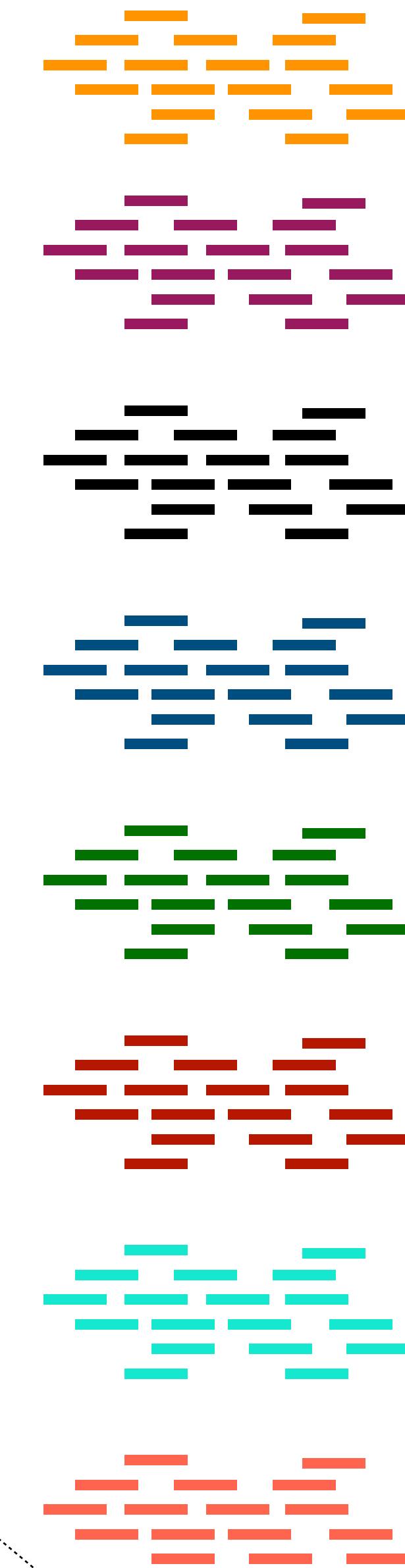


Illumina sequencing

100s to 1000s bp long



100-250 bp long



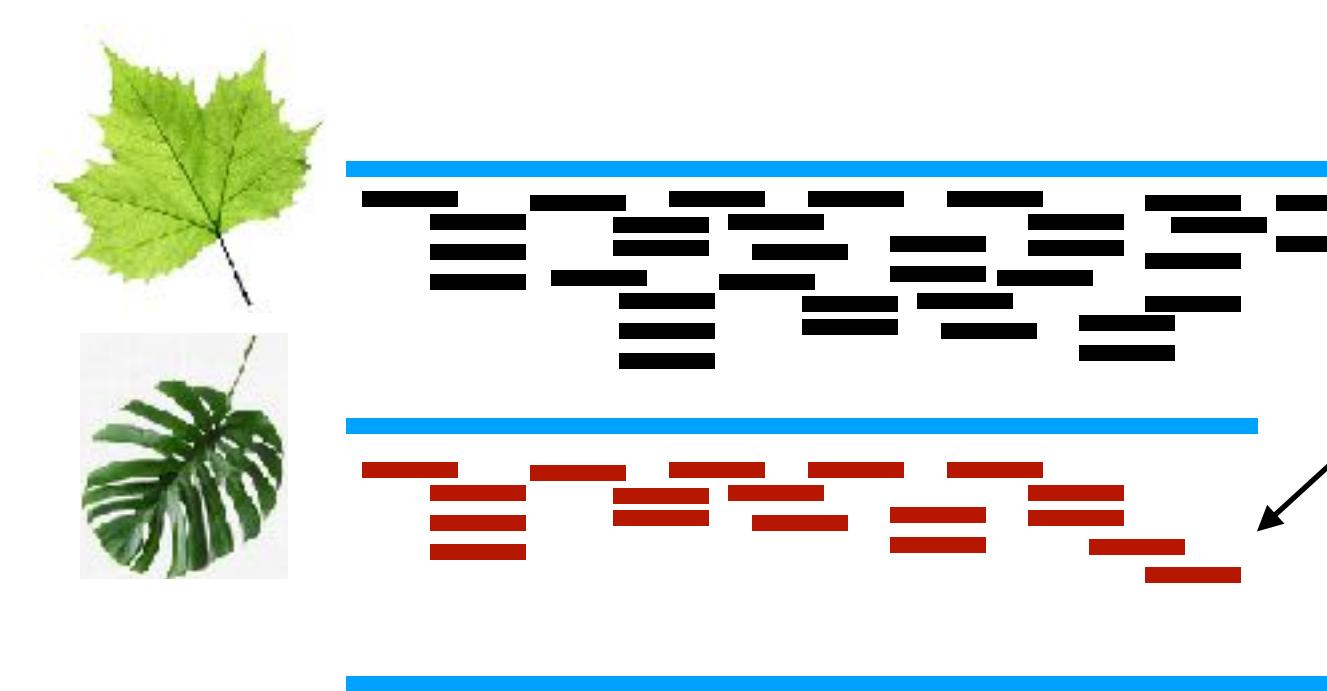
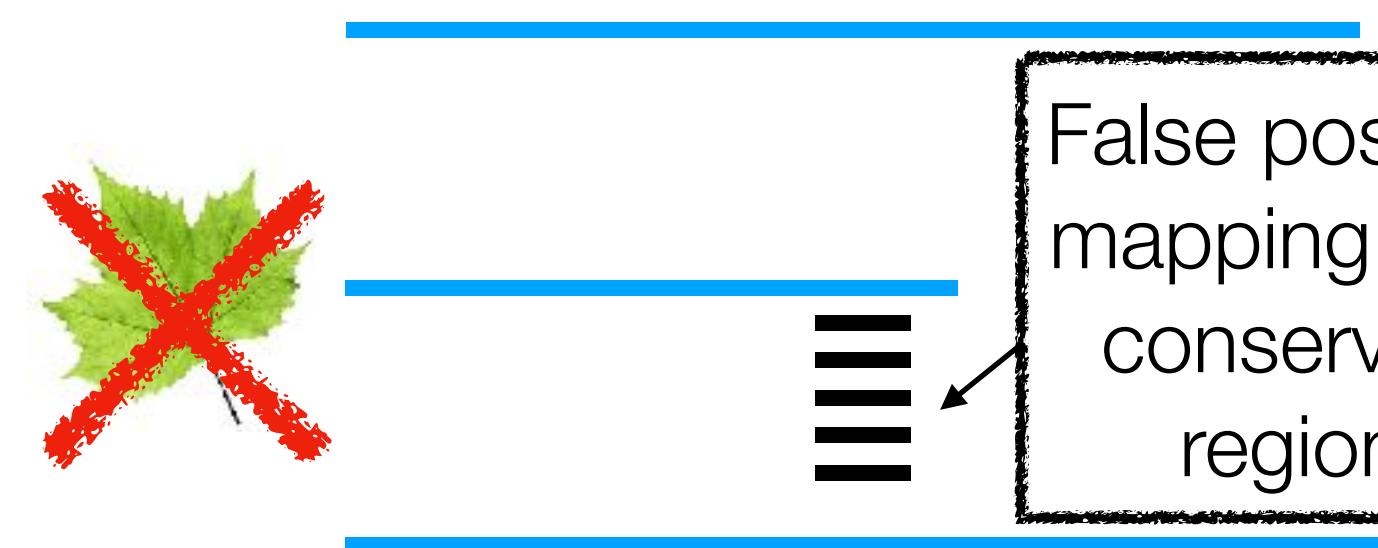
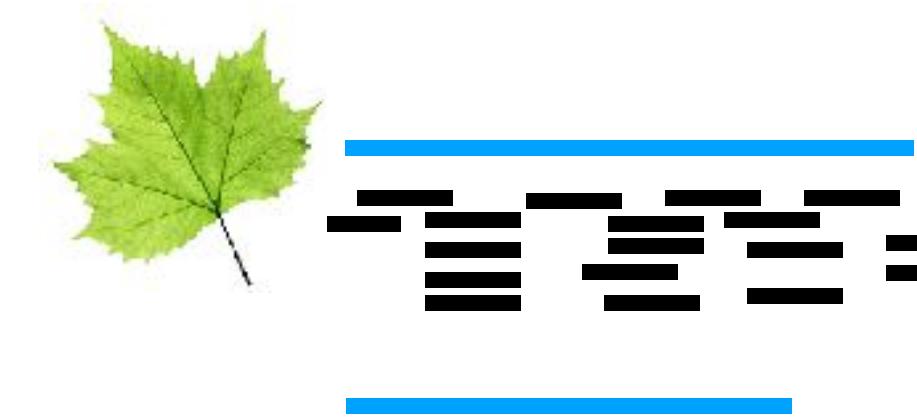
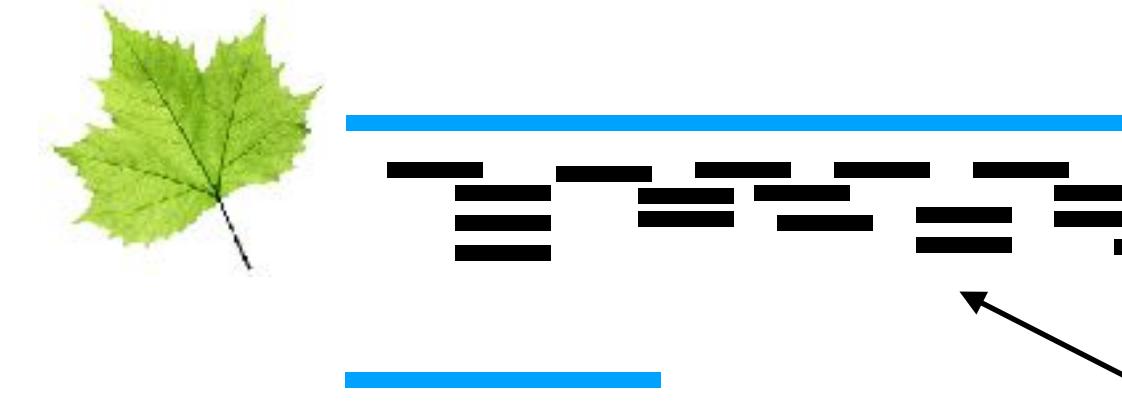
We map the references to the queries

Queries:
100s to 1000s bp long

References:
100-250 bp long



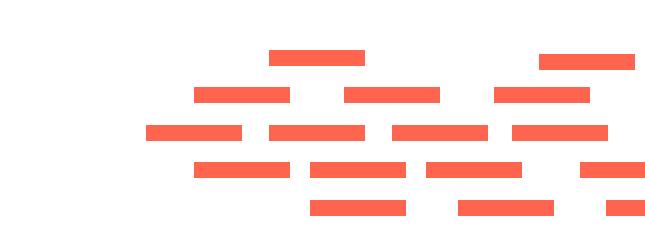
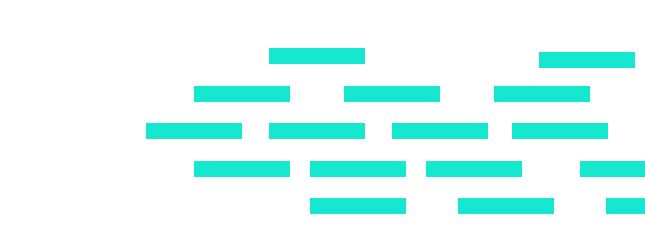
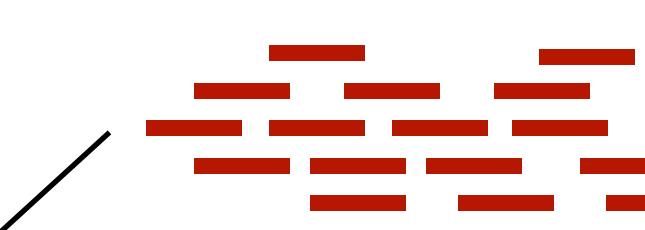
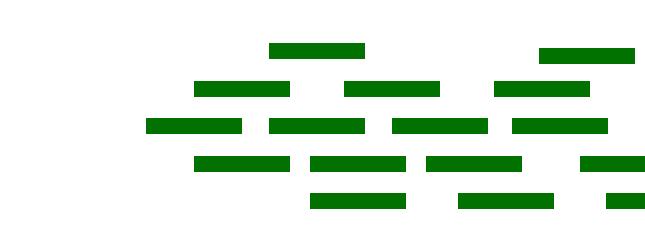
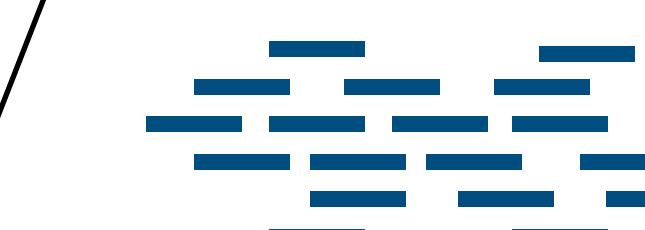
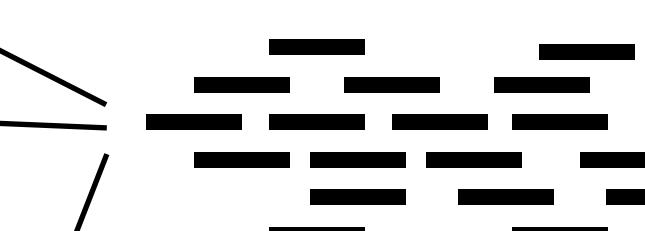
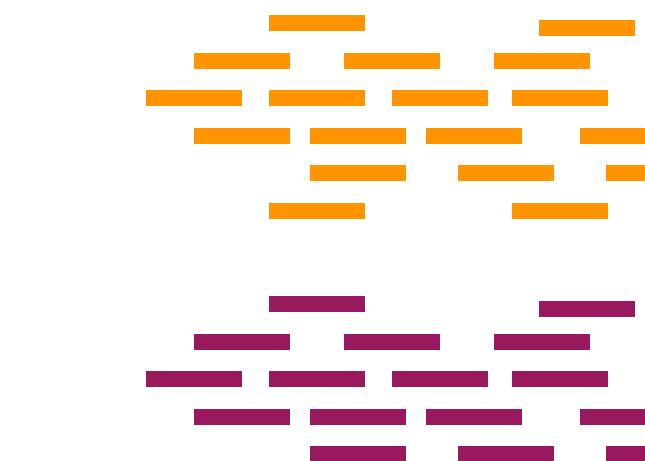
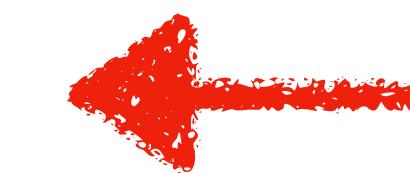
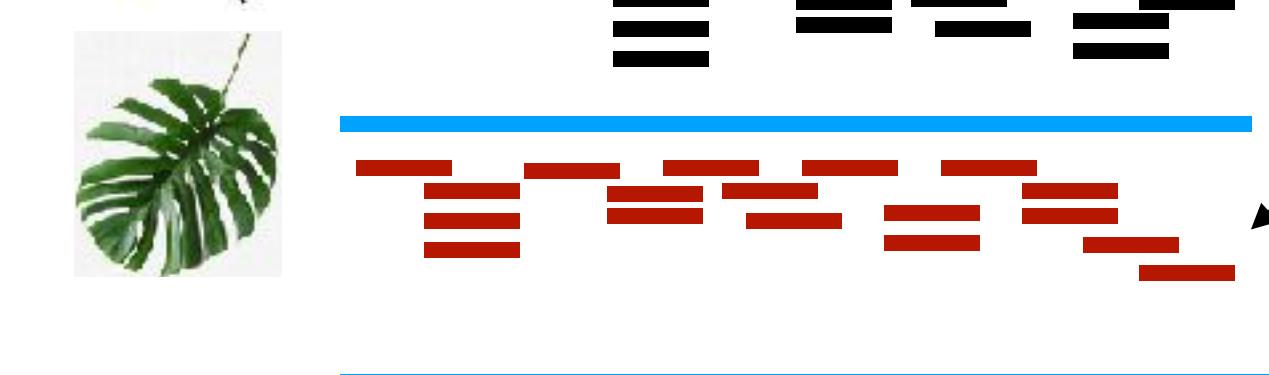
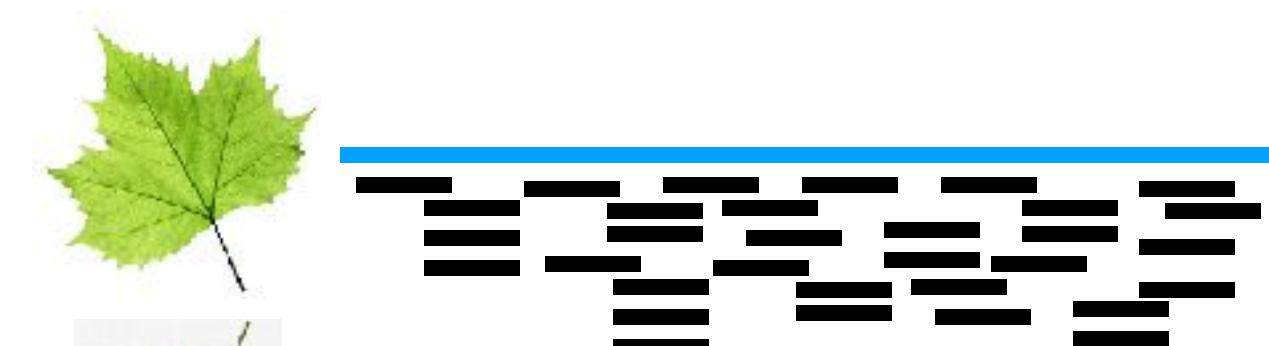
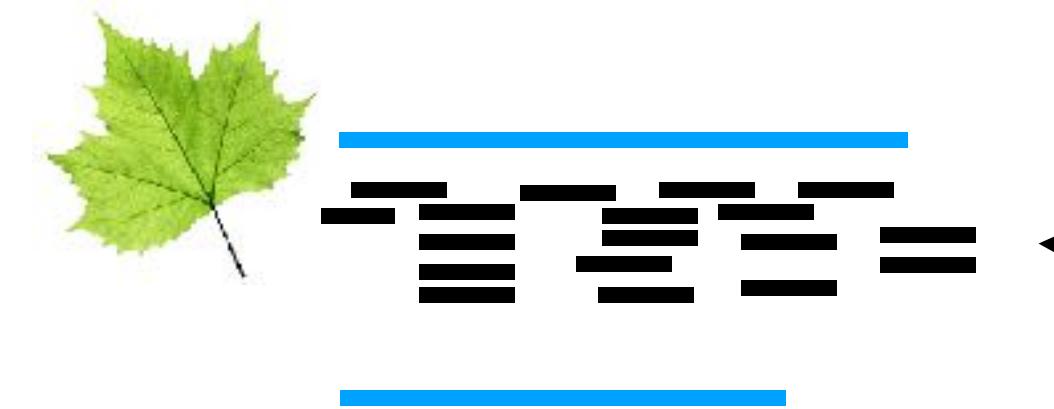
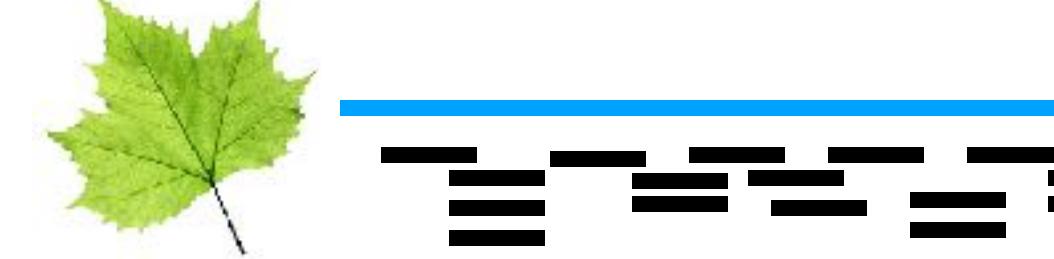
False positive mapping to a conserved region



We map the
references
to the
queries

Queries:
100s to 1000s bp long

References:
100-250 bp long



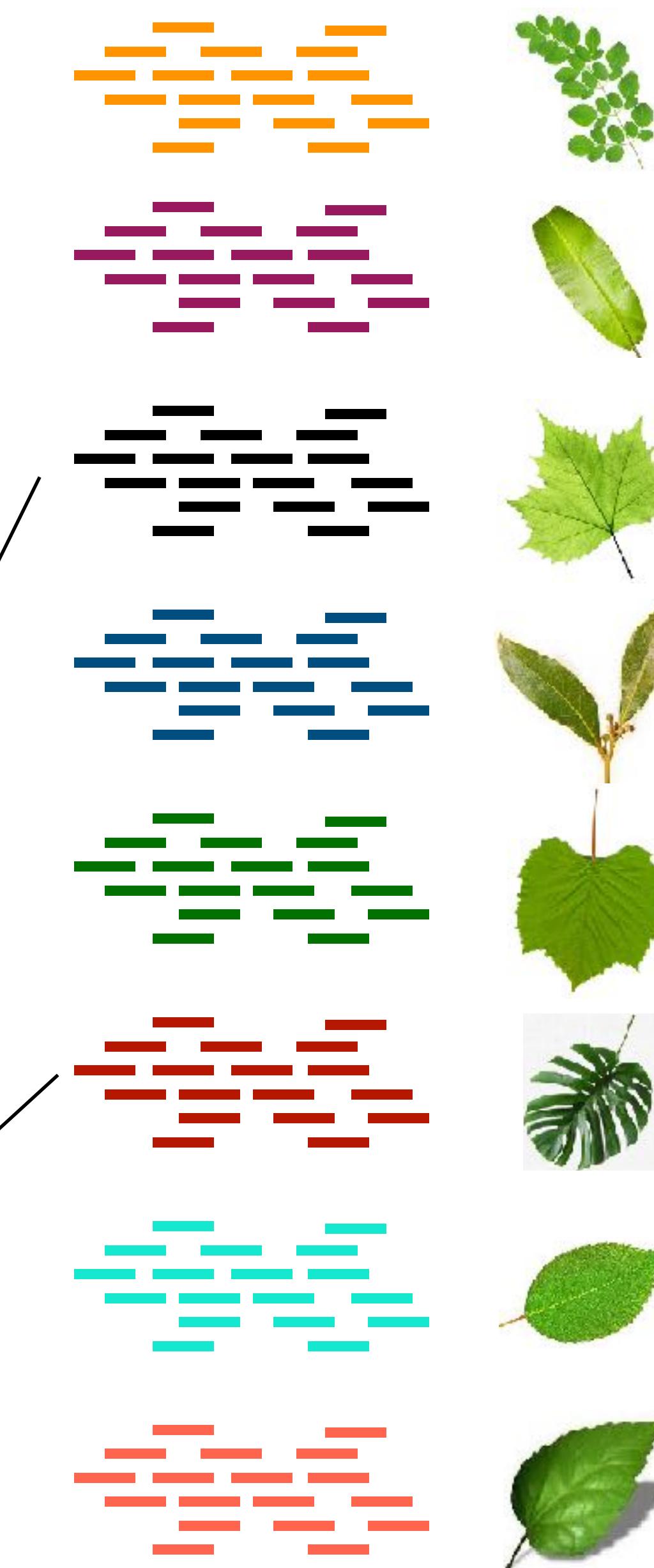
We map the
references
to the
queries



Queries:
100s to 1000s bp long



References:
100-250 bp long



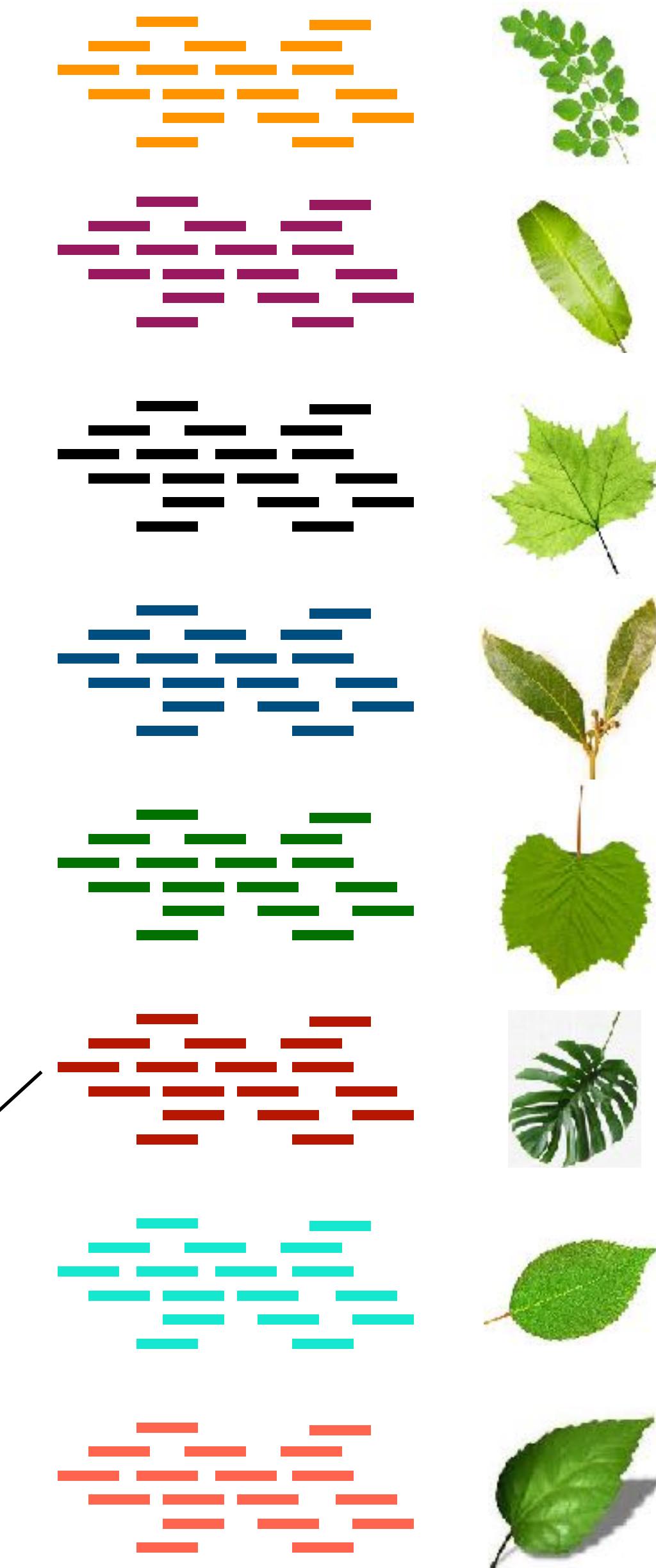
We map the
references
to the
queries



Queries:
100s to 1000s bp long



References:
100-250 bp long

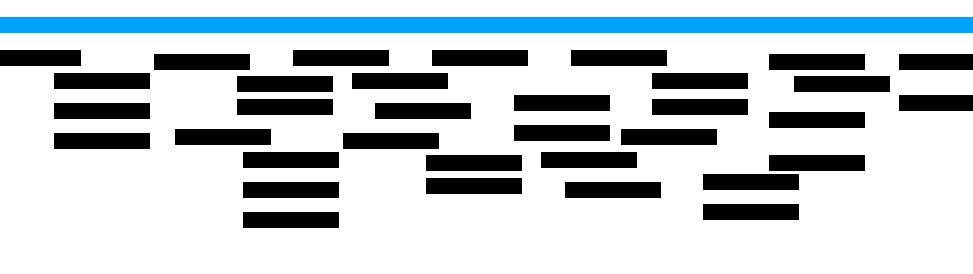


Our goals:

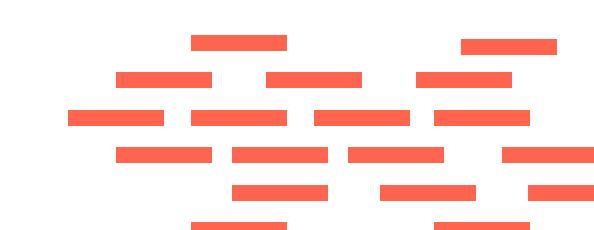
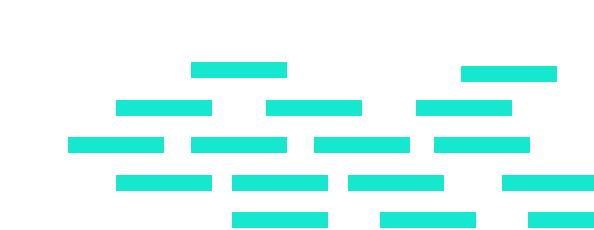
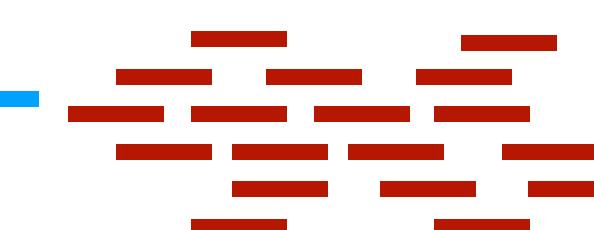
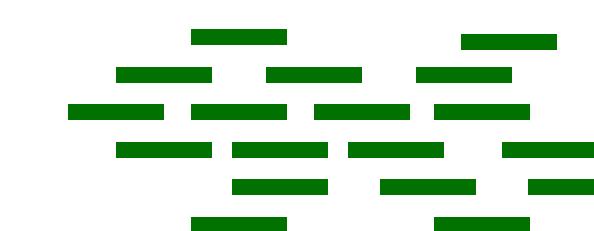
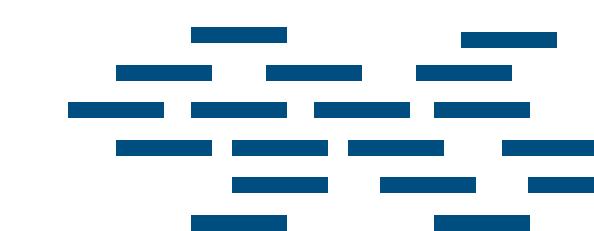
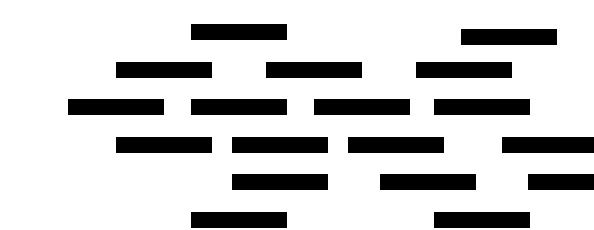
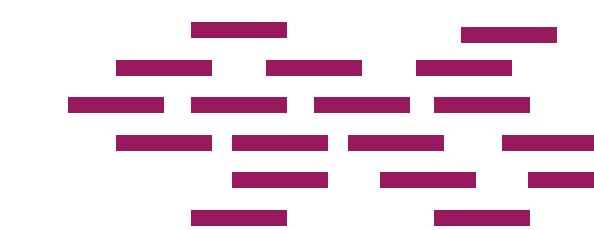
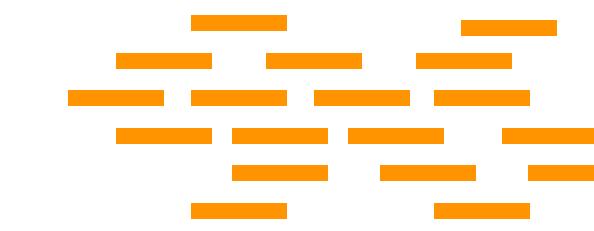
- Does not need specialist molecular-lab skills
- Can be used on small amounts of pollen
- Uses the **nuclear genome**, not (just) the chloroplast genome
- **Does NOT require assembled genomes** for the reference database
- Is '**semi-quantitative**' (low vs. high biomass frequency)
- Can **identify pollen to plant species**



each read is from one cell

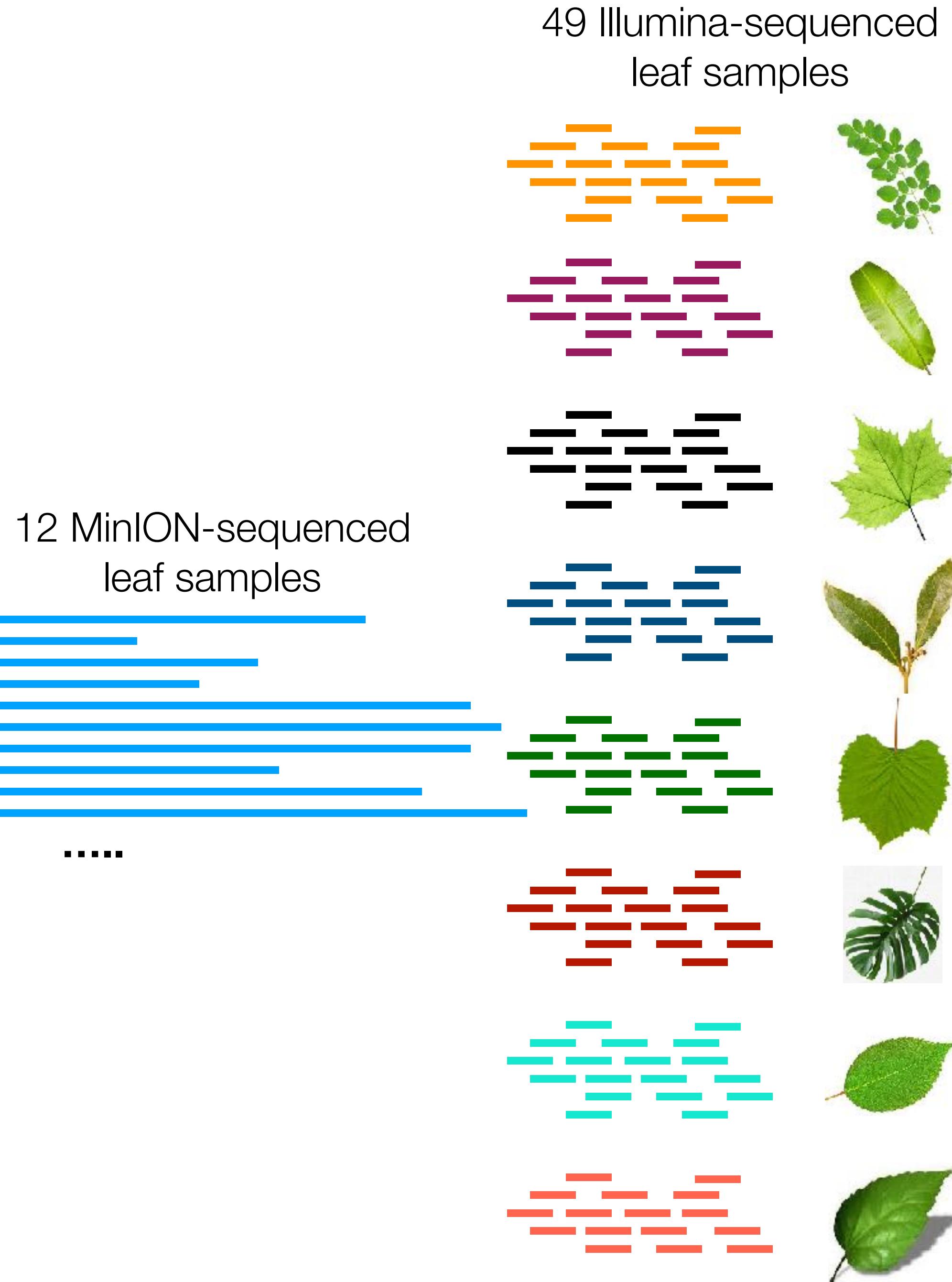


References:
100-250 bp long



Our test:

- 49 UK plant species **Illumina shotgun-sequenced**, 250 PE, 0.3-1.0X coverage
- 12 mock leaf samples of known composition, **MinION sequenced**
- **Map all Illumina skims to all MinION reads**
 - $49 \times 10 = 490$ combinations
 - *minimap2* → *samtools* → *R*
 - filter by depth and evenness of mapping coverage



12 Mock samples – DNA per species

Species_name_Latin	hml - 1	hml - 2	lmh - 1	lmh - 2	hml - 1	hml - 2	Hhl - 1	Hhl - 2	hhll - 1	hhll - 2	llhh - 1	llhh - 2
<i>Knautia arvensis</i>	100	100							100	100	1	1
<i>Galium verum</i>	100	100							100	100	1	1
<i>Crepis capillaris</i>	100	100							100	100	1	1
<i>Papaver somniferum</i>	10	10	1	1			1	1			100	100
<i>Anagallis arvensis</i>	10	10	1	1					1	1	100	100
<i>Sambucus nigra</i>	10	10	1	1			1000	1000	1	1	100	100
<i>Bryonia dioica</i>	1	1	10	10	10	10			1	1	100	100
<i>Ranunculus repens</i>	1	1	10	10	100	100			100	100		
<i>Lotus corniculatus</i>	1	1	10	10			100	100				
<i>Digitalis purpurea</i>	0	0	100	100	10	10						
<i>Leucanthemum</i>	0	0	100	100	100	100						
<i>Stachys sylvatica</i>	0	0	100	100	1	1	100	100	1	1	1	1

+ 37 more
species in
reference

Expectation is that ONLY Illumina reads from the included plant species should successfully map onto MinION reads in that sample

12 Mock samples – DNA per species

Species_name_Latin	hml - 1	hml - 2	lh - 1	lmh - 2	hml - 1	hml - 2	Hhl - 1	Hhl - 2	hhll - 1	hhll - 2	llhh - 1	llhh - 2
<i>Knautia arvensis</i>	100	100							100	100	1	1
<i>Galium verum</i>	100	100							100	100	1	1
<i>Crepis capillaris</i>	100	100							100	100	1	1
<i>Papaver somniferum</i>	10	10	1	1			1	1			100	100
<i>Anagallis arvensis</i>	10	10	1	1					1	1	100	100
<i>Sambucus nigra</i>	10	10	1	1			1000	1000	1	1	100	100
<i>Bryonia dioica</i>	1	1	10	10	10	10			1	1	100	100
<i>Ranunculus repens</i>	1	1	10	10	100	100			100	100		
<i>Lotus corniculatus</i>	1	1	10	10			100	100				
<i>Digitalis purpurea</i>	0	0	100	100	10	10						
<i>Leucanthemum</i>	0	0	100	100	100	100						
<i>Stachys sylvatica</i>	0	0	100	100	1	1	100	100	1	1	1	1

Species_name_Latin	hml - 1	hml - 2	lh - 1	lmh - 2	hml - 1	hml - 2	Hhl - 1	Hhl - 2	hhll - 1	hhll - 2	llhh - 1	llhh - 2
<i>Knautia arvensis</i>	286	267	0	0	1	0	0	0	139	223	4	2
<i>Galium verum</i>	127	79	1	1	0	0	1	0	28	40	1	1
<i>Crepis capillaris</i>	342	331	0	0	1	0	2	0	179	231	1	5
<i>Papaver somniferum</i>	13	10	1	0	0	0	1	0	0	0	127	88
<i>Anagallis arvensis</i>	20	24	3	1	0	0	2	0	1	2	370	286
<i>Sambucus nigra</i>	25	25	1	6	1	6	1225	783	2	3	423	298
<i>Bryonia dioica</i>	1	0	10	18	54	33	0	0	0	0	205	152
<i>Ranunculus repens</i>	3	3	21	31	543	516	1	1	180	239	2	2
<i>Lotus corniculatus</i>	0	0	3	12	0	0	44	30	0	0	0	0
<i>Digitalis purpurea</i>	0	0	42	70	19	9	0	0	0	0	0	0
<i>Leucanthemum</i>	1	0	78	143	229	232	0	0	0	0	1	0
<i>Stachys sylvatica</i>	0	1	299	494	7	6	209	144	5	6	5	3

12 Mock samples – DNA per species

Species_name_Latin	hml - 1	hml - 2	lh - 1	lmh - 2	hml - 1	hml - 2	Hhl - 1	Hhl - 2	hhll - 1	hhll - 2	llhh - 1	llhh - 2
<i>Knautia arvensis</i>	100	100							100	100	1	1
<i>Galium verum</i>	100	100							100	100	1	1
<i>Crepis capillaris</i>	100	100							100	100	1	1
<i>Papaver somniferum</i>	10	10	1	1			1	1			100	100
<i>Anagallis arvensis</i>	10	10	1	1					1	1	100	100
<i>Sambucus nigra</i>	10	10	1	1			1000	1000	1	1	100	100
<i>Bryonia dioica</i>	1	1	10	10	10	10			1	1	100	100
<i>Ranunculus repens</i>	1	1	10	10	100	100			100	100		
<i>Lotus corniculatus</i>	1	1	10	10			100	100				
<i>Digitalis purpurea</i>	0	0	100	100	10	10						
<i>Leucanthemum</i>	0	0										
<i>Stachys sylvatica</i>	0	0										

We can see which species are **present** and which are **absent**, and we can see which plant species have **higher biomass**

Species_name_Latin	hml - 1	hml - 2	lh	lh	lh	lh	lh	lh	lh	lh	lh	lh
<i>Knautia arvensis</i>	286	267										
<i>Galium verum</i>	127	79	1	1	0	0	1	0	28	40	1	1
<i>Crepis capillaris</i>	342	331	0	0	1	0	2	0	179	231	1	5
<i>Papaver somniferum</i>	13	10	1	0	0	0	1	0	0	0	127	88
<i>Anagallis arvensis</i>	20	24	3	1	0	0	2	0	1	2	370	286
<i>Sambucus nigra</i>	25	25	1	6	1	6	1225	783	2	3	423	298
<i>Bryonia dioica</i>	1	0	10	18	54	33	0	0	0	0	205	152
<i>Ranunculus repens</i>	3	3	21	31	543	516	1	1	180	239	2	2
<i>Lotus corniculatus</i>	0	0	3	12	0	0	44	30	0	0	0	0
<i>Digitalis purpurea</i>	0	0	42	70	19	9	0	0	0	0	0	0
<i>Leucanthemum</i>	1	0	78	143	229	232	0	0	0	1	0	0
<i>Stachys sylvatica</i>	0	1	299	494	7	6	209	144	5	6	5	3

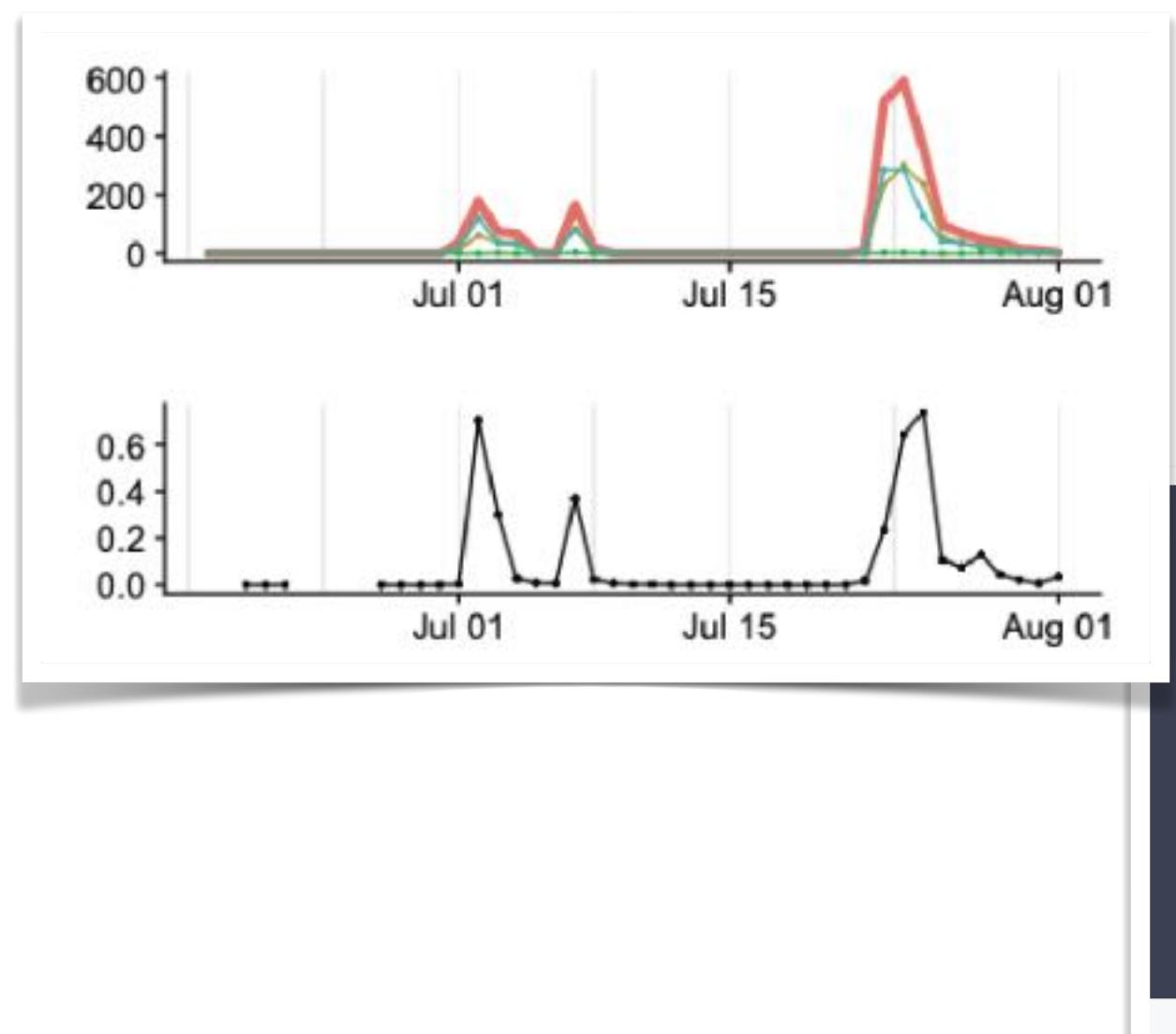
Our goals – a protocol that:

- Does not need specialist molecular-lab skills
- Can be used on small amounts of pollen
- Uses the **nuclear genome**, not (just) the chloroplast genome
- **Does NOT require assembled genomes** for the reference database
- Is '**semi-quantitative**' (low vs. high biomass frequency)
- Can **identify pollen to plant species**
- MinION library prep still complex but getting easier (and NO PCR!)
- ✓
- ✓
- ✓
- ~ Good at identifying dominant constituents, poorer at identifying rare components
- ~ Not great at differentiating closely related congeners (e.g. *Ranunculus acris* and *R. repens*, which also hard to differentiate using barcodes and morphology)

Summary

Single-species quantitative PCR (qPCR)

- **use eDNA rate** not concentration! ng/sec not ng/ μ l



mBRAVE

Multiplex Barcode Research And Visualization Environment

mBRAVE is a multi-user platform supporting the storage, validation, analysis, and publication of highly multiplexed projects based on high-throughput sequencing (HTS) instruments. This system builds on the BOLD Platform to support species identification and discovery for HTS data.

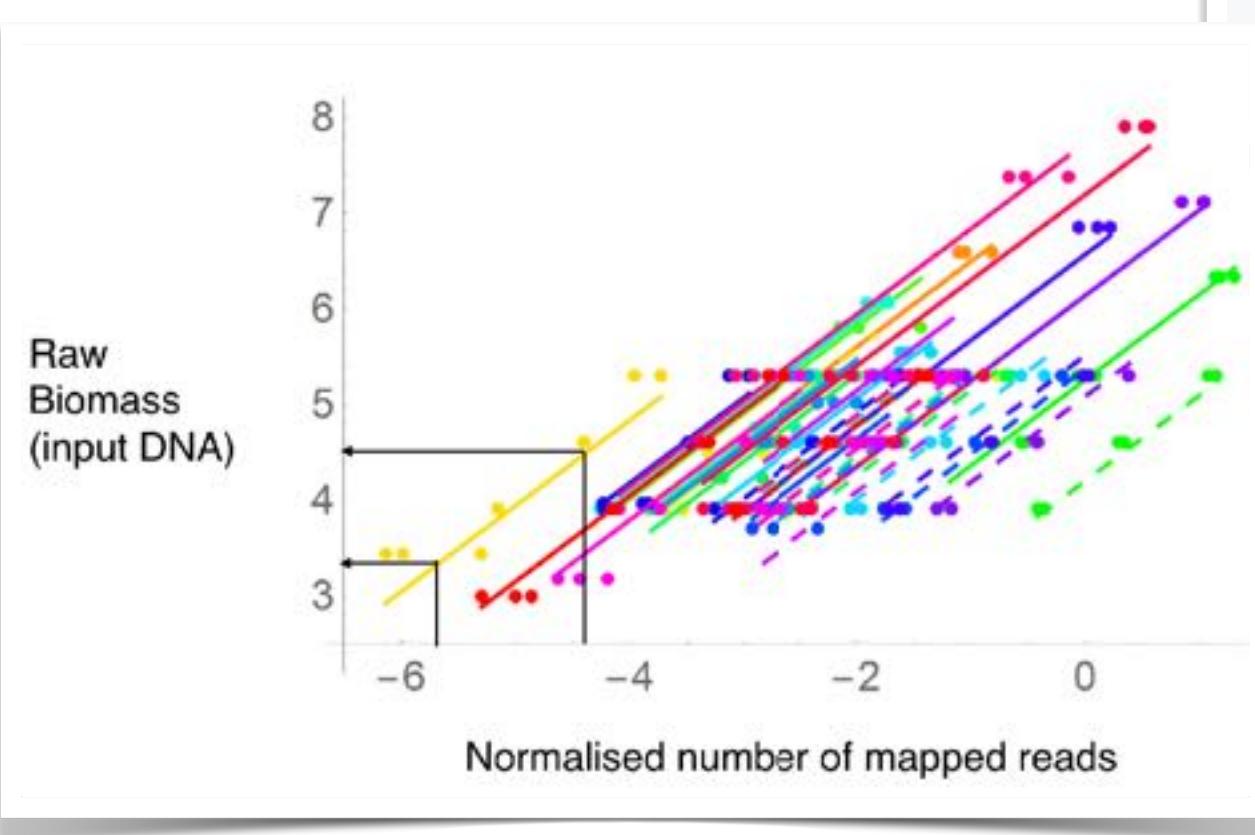


Login Register

mBRAVE is currently accepting a limited number of new users on a weekly basis.

Species_name_Latin	hm - 1	hm - 2	hm - 3
<i>Knautia arvensis</i>	100	100	
<i>Galium verum</i>	100	100	
<i>Crepis capillaris</i>	100	100	
<i>Papaver somniferum</i>	10	10	
<i>Anagallis arvensis</i>	10	10	
<i>Sambucus nigra</i>	10	10	
<i>Bryonia dioica</i>	1	1	
<i>Ranunculus repens</i>	1	1	
<i>Lotus corniculatus</i>	1	1	
<i>Digitalis purpurea</i>	0	0	
<i>Leucanthemum</i>	0	0	
<i>Stachys sylvatica</i>	0	0	

Species_name_Latin	hm - 1	hm - 2	hm - 3
<i>Knautia arvensis</i>	286	267	
<i>Galium verum</i>	127	79	
<i>Crepis capillaris</i>	342	331	
<i>Papaver somniferum</i>	13	10	
<i>Anagallis arvensis</i>	20	24	
<i>Sambucus nigra</i>	25	25	
<i>Bryonia dioica</i>	1	0	
<i>Ranunculus repens</i>	3	3	
<i>Lotus corniculatus</i>	0	0	
<i>Digitalis purpurea</i>	0	0	
<i>Leucanthemum</i>	1	0	
<i>Stachys sylvatica</i>	0	1	



Mitogenomics/Metabarcoding + DNA spike-in

(SPIKEPIPE, qSeq)

- **For within-species differences across samples** (e.g. time series)
- (difficult to get across-species diffs)

Reverse metagenomics (RevMet)

- **For within- and across-species differences within a sample** (e.g. diet analysis)