

# unikseq

**unique & conserved** sequence identification  
using **multiple, complete** genomes

René L Warren  
2022

# COMPARATIVE GENOMICS

## COMPARISON OF COMPLETE GENETIC M.

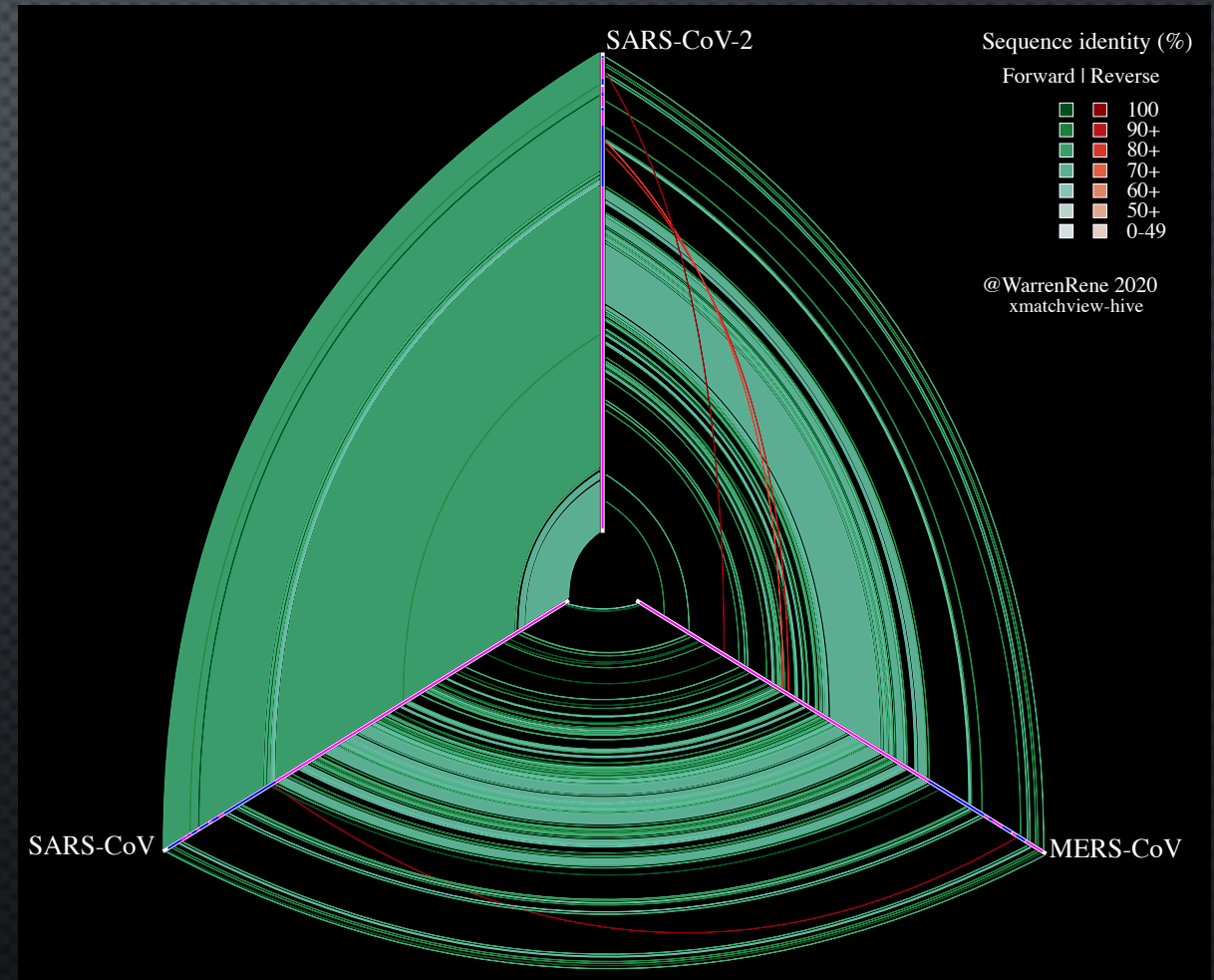
- SPECIES EVOLUTION / GENE FUNCTION
- BASE-LEVEL RESOLUTION

## VISUALIZED TO GAIN INSIGHTS

- REGION CONSERVATION / UNIQUENESS
- CO-LINEARITY / STRUCTURE

## BIOINFORMATICS AT THE HEART

- SEQUENCE ALIGNMENTS
- PAIRWISE / MULTIPLE SEQUENCE ALIGNMENTS (MSA)



# MSA

COMPUTATIONALLY MORE COMPLEX VS. PAIR-WISE ALIGNMENTS + COSTLY + NOT SCALABLE

## MSA KEY TO COMPARATIVE GENOMICS

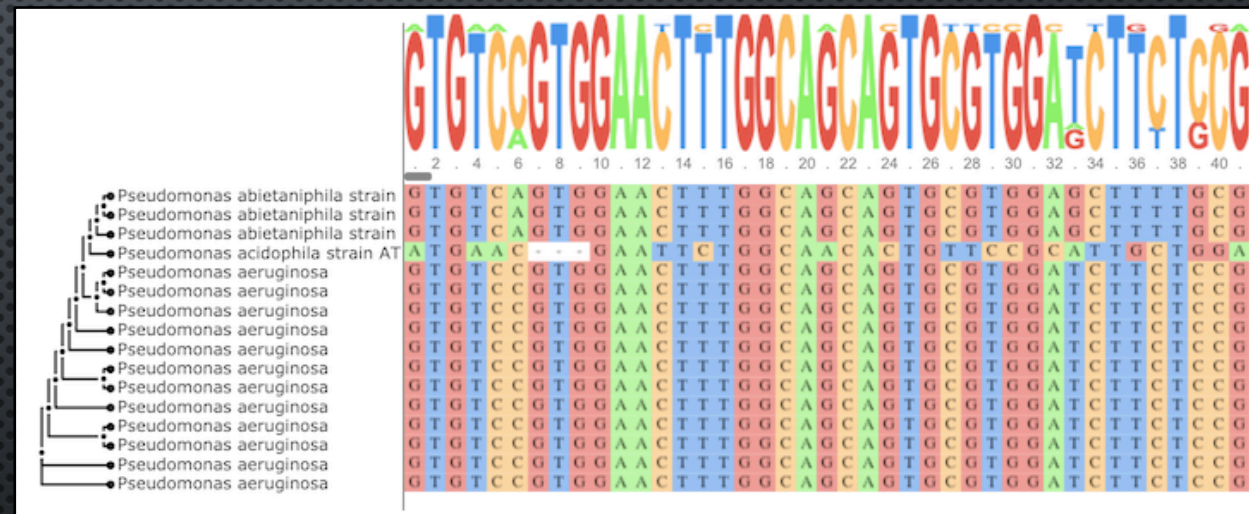
- CLUSTAL, MUSCLE, MAFFT
- USEFUL FOR PCR ASSAY DESIGN

## WITH PCR, NEED VISUAL INSPECTION

- IDENTIFY UNIQUE / CONSERVED REGIONS

## I/O ISSUES

- FINICKY INPUT : STRAND / SEQUENCE START
- OUTPUT INTERPRETABILITY ON KB-STRETCHES, >10s ENTRIES



Sequences compared with FastTree (Price 2009), Gblocks (Castresana 2002), Muscle (Edgar 2004) visualized with MSAviewer/PATRIC (<https://www.bv-brc.org/>)

# POLYMERASE CHAIN REACTION

## REVOLUTIONARY

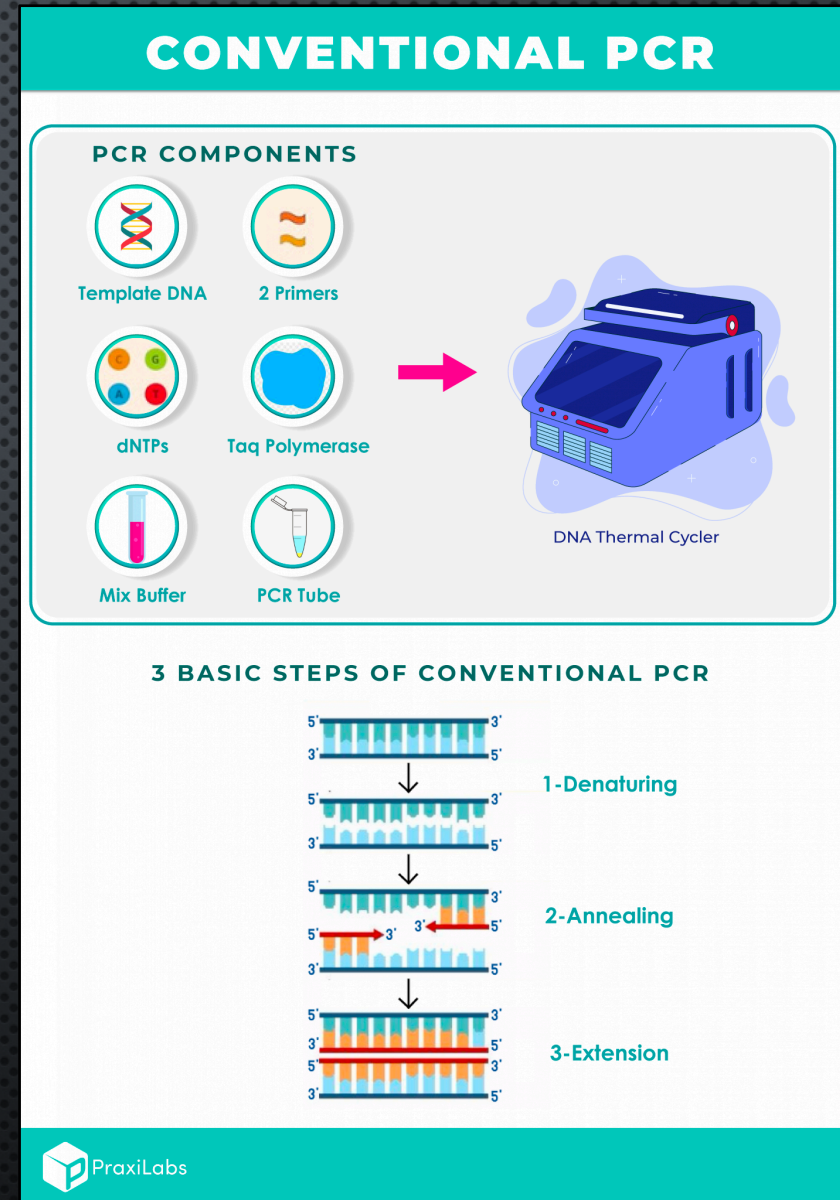
- MOLEC BIO LAB TECHNIQUES FOR
- NUCLEIC ACID AMPLIFICATION, DETECTION, QUANT.

## BROAD APPLICATIONS

- FORENSICS, RESEARCH, CLINICAL SETTINGS, BIOLOGY
- SPECIES MONITORING, CONSERVATION, ETC.

## EFFECTIVE PCR : SENSITIVE & SPECIFIC

- DESIGNING PCR PRIMERS (& QUANTITATIVE QPCR PROBE)
- CHALLENGING + LABOUR + TIME



# BIODIVERSITY LOSS

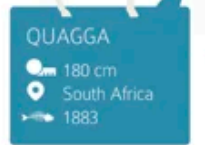
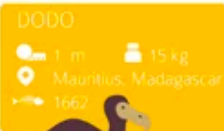
ANIMALS AND PLANTS DESTRUCTION

environmental infographic

- 226 species of mammals
- 181 species of birds
- 168 species of fish
- 77 species of reptiles
- 35 species of amphibians



## EXTINCT ANIMALS



## CAUSES



swamp drying



deforestation



steppes plowing



poaching



dams construction



hunting



grazing



urban construction



pollution



overfishing



animal skins



tropical fish trade



climate change



death from vehicles



wild animals trade

# iTrackDNA



non-destructive precision genomics for environmental impact tracking in a global climate change era



<https://itrackdna.ca/>



# iTrackDNA



- 4 YEARS \$12M LARGE SCALE APPLIED RESEARCH PROJECT
- TO FILL CRITICAL KNOWLEDGE GAPS & PROMOTE eDNA UPTAKE
- BUILD GENOMIC RESOURCES & eDNA CAPACITY (100s KITS)
- DEVELOP NATIONAL STANDARD FOR TARGETED eDNA ASSAYS
- PROMOTE TRAINING / CERTIFICATION
- INFLUENCE NATURAL RESOURCE MANAGEMENT / POLICY



# ENVIRONMENTAL DNA / RNA (eDNA / eRNA)

GENETIC MATERIAL EXTRACTED FROM ENVIRONMENTAL SAMPLES



RAPID  
ACCURATE  
COST-EFFECTIVE  
NON-INVASIVE  
ACTIONABLE



BARRIERS TO UPTAKE  
NEED STANDARDS  
NEED GENOMIC RESOURCES

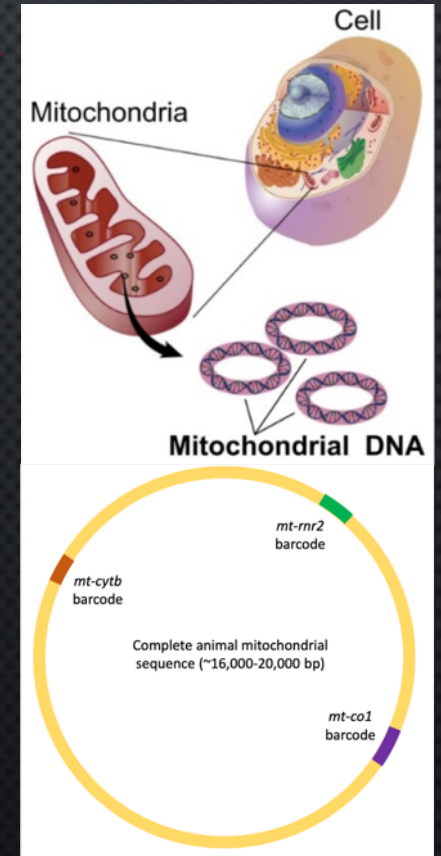
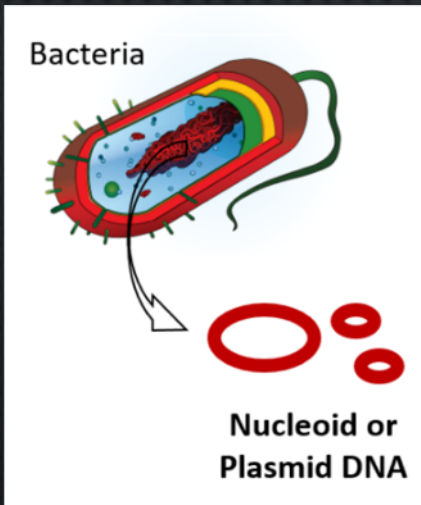




# eDNA

microbial

bulk sample  
macrobial  
(shed)



ANIMAL MITOGENOMES

ABUNDANT, TAXA SPECIFICITY & SEQUENCE CONSERVATION 9

# UTILITY OF eDNA

DETECTION / ENFORCEMENT OF TRAFFICKED SPECIES

*CITES* : CONVENTION ON INTERNATIONAL TRADE IN ENDANGERED SPECIES OF WILD FAUNA AND FLORA

*WAPTR* : WILD ANIMAL AND PLANT TRADE REGULATIONS

ENVIRONMENTAL/ECOSYSTEM MONITORING (E.G. IMPACT OF OIL SANDS, CLIMATE, ETC.)



*CITES*, wikipedia

# eDNA ASSAY

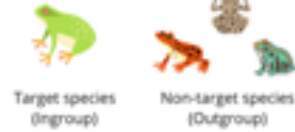
multi-step workflow

+collections of complete mitogenomes [sequences]

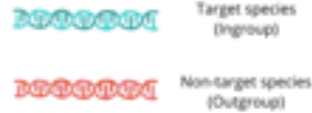
develop high-quality eDNA assays / kits

Environmental DNA

Relevant species identification



Mitogenome acquisition from public repositories or direct sequencing



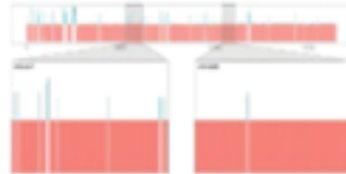
Multiple mitogenome sequence alignment and phylogenetic analysis



Unikseq

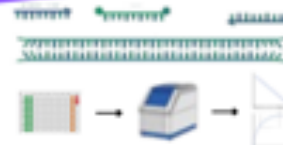


Parsing of mitochondrial sequences into k-mers



Identification of unique regions within the mitogenome

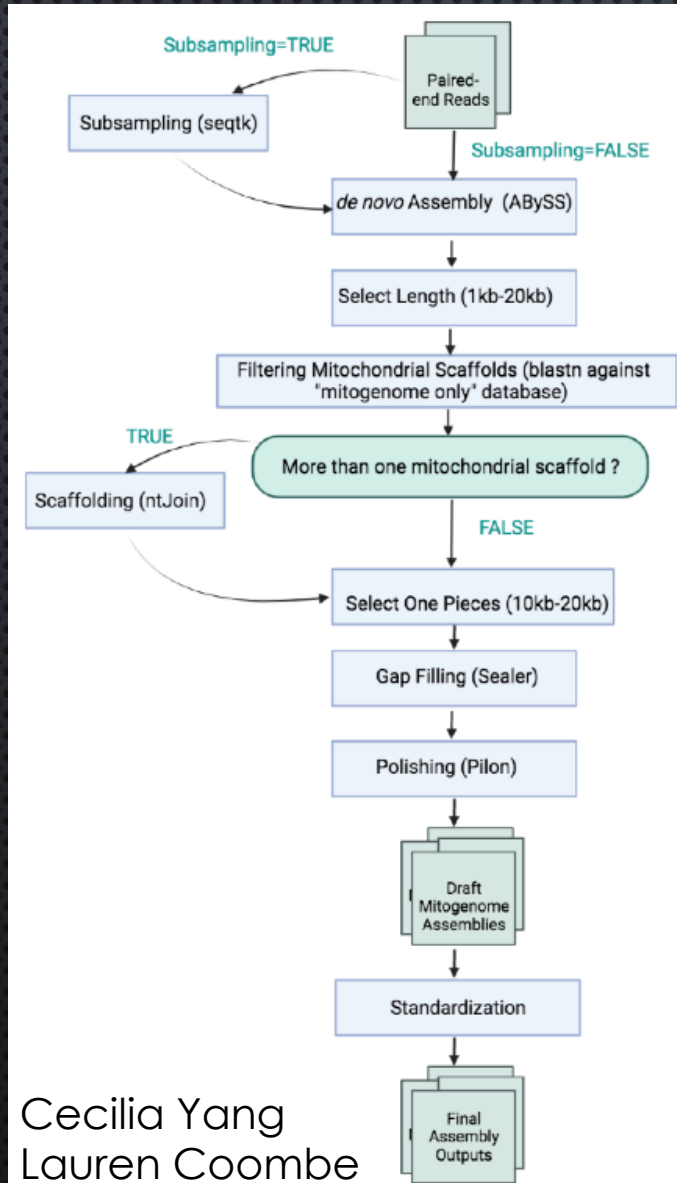
eDNA assay design and *in vitro* analysis



eDNA assay *in situ* (field) validation



# BUILDING MITOGENOME RESOURCES



**Sablefish (4)**

**White sharks (9)**

**Rockfish (57)**



**Thornyhead (4)**



**Greater sturgeons (9)**



**Longfin smelt (9)**



**Oyster (10)**



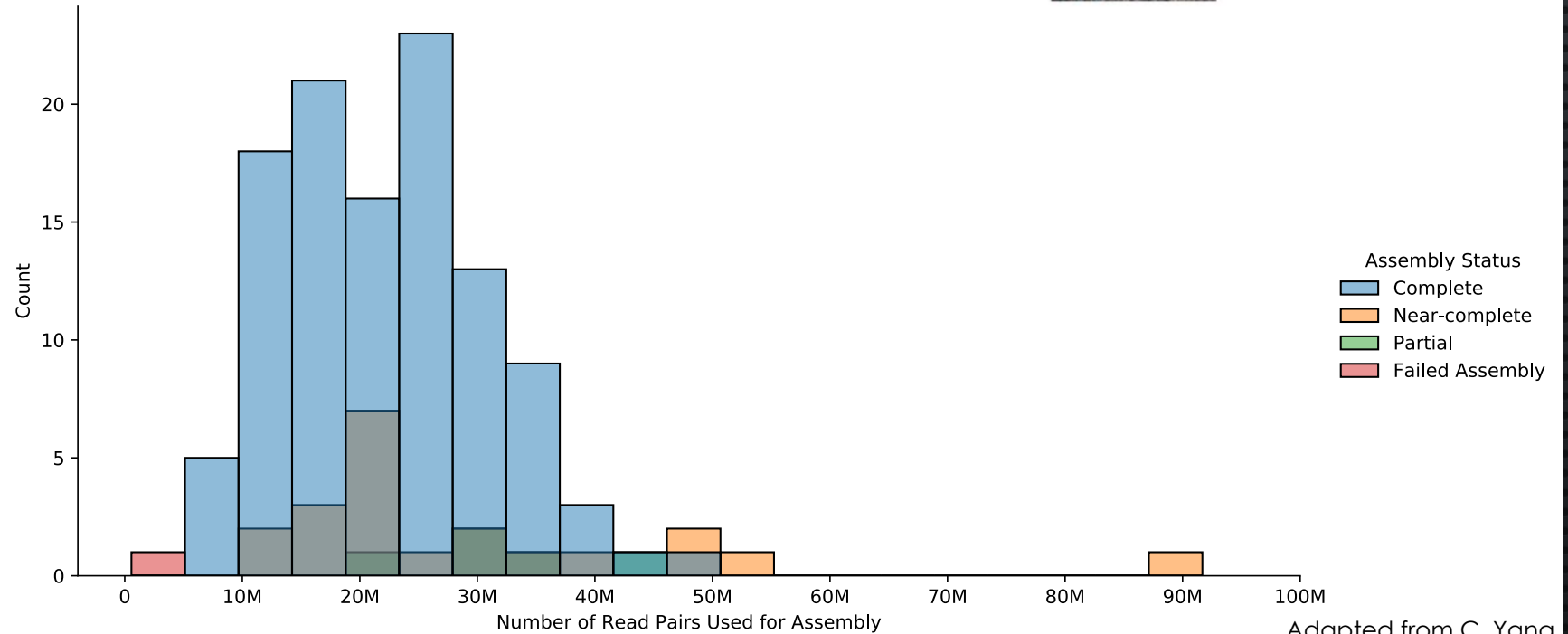
**Brown frog (2)**



**Long-tailed shrew (24)**



**Worm (9)**



Cecilia Yang  
Lauren Coombe

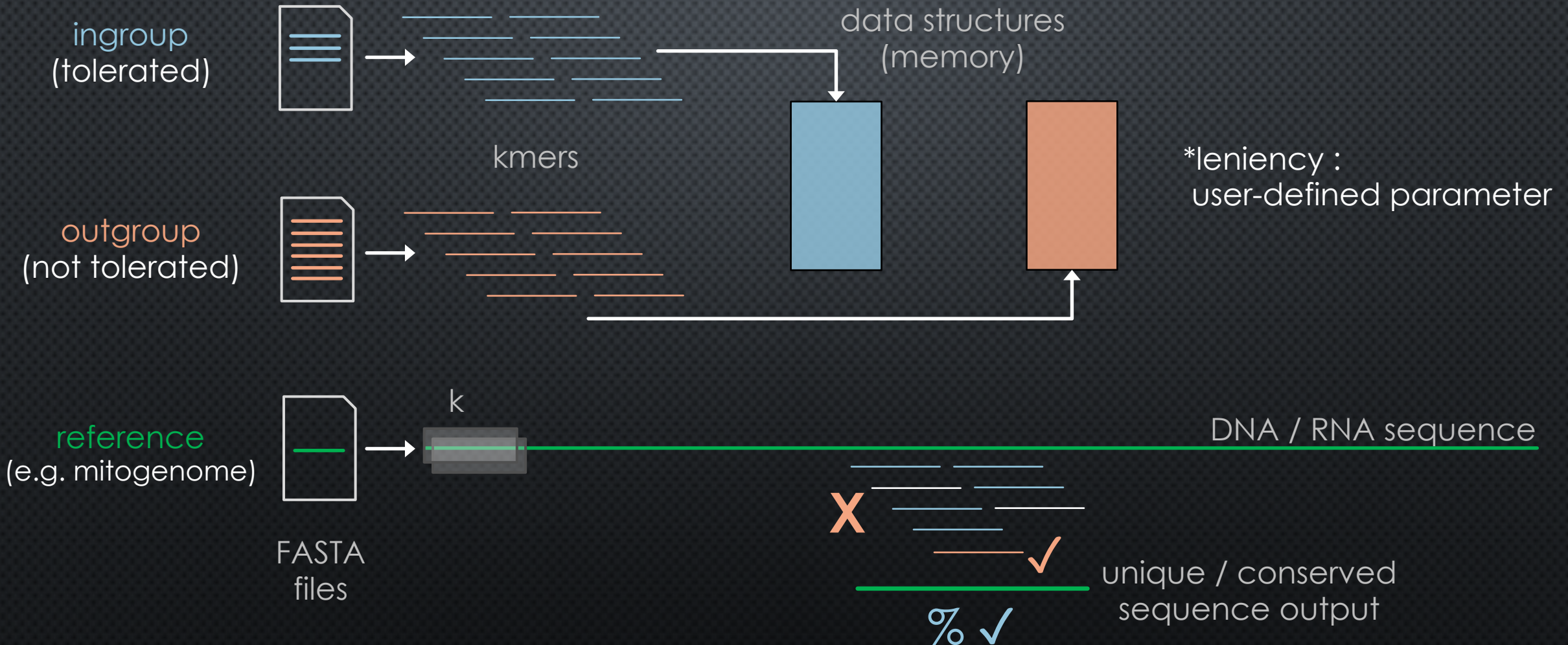
Adapted from C. Yang



unique / conserved region identification in  
DNA / RNA sequences using a kmer-based approach

# UNIQUE REGION IDENTIFICATION

IN A REFERENCE, TOLERATED IN INGROUP, NOT\* IN OUTGROUP



# GREAT WHITE SHARK : USE CASE



- VULNERABLE SPECIES

ENVIRONMENT MONITORING : CONSERVATION EFFORT INSIGHTS

- FIND MITOGENOME REGIONS UNIQUE &

HIGHLY (100%) CONSERVED IN GREAT WHITES

- SPECIFIC qPCR-BASED eDNA ASSAY DESIGN

121 great white shark mitogenome [isolates] (ingroup)

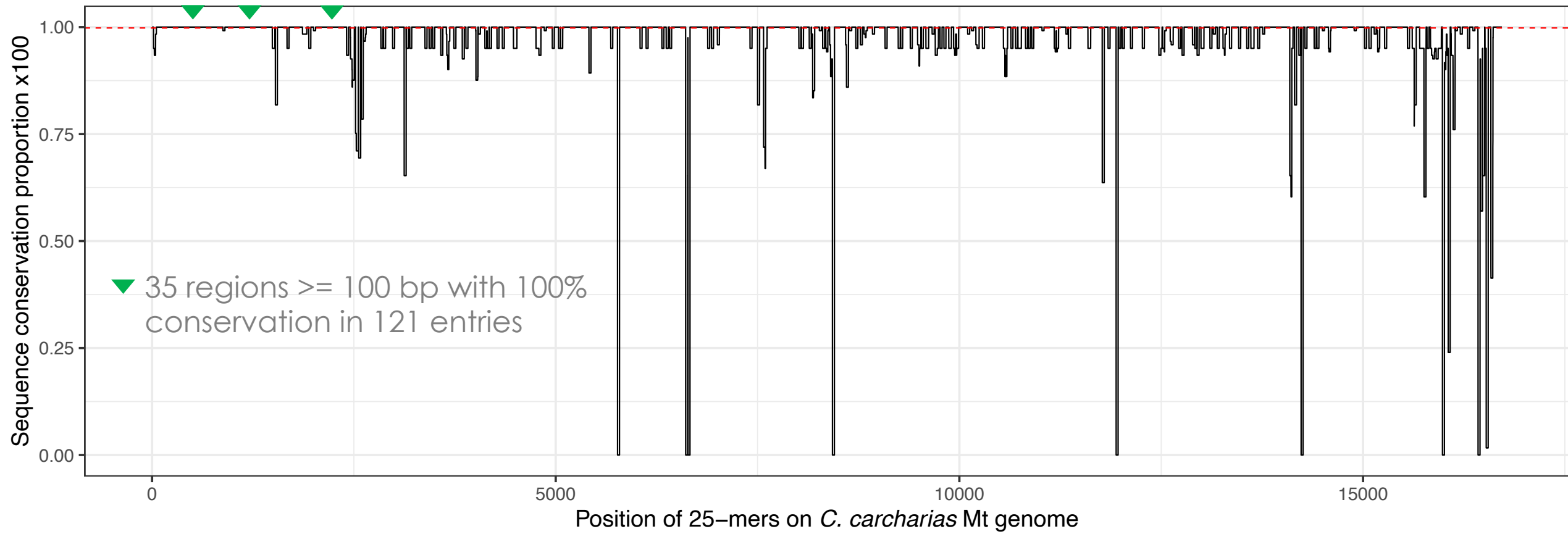
856 fish mitogenomes (outgroup)



# SEQUENCE CONSERVATION

## ACROSS 121 WHITE SHARK MTG (ISOLATES)

comparing 2,000,000 bp



clustal 24h52m

MUSCLE 28m27s

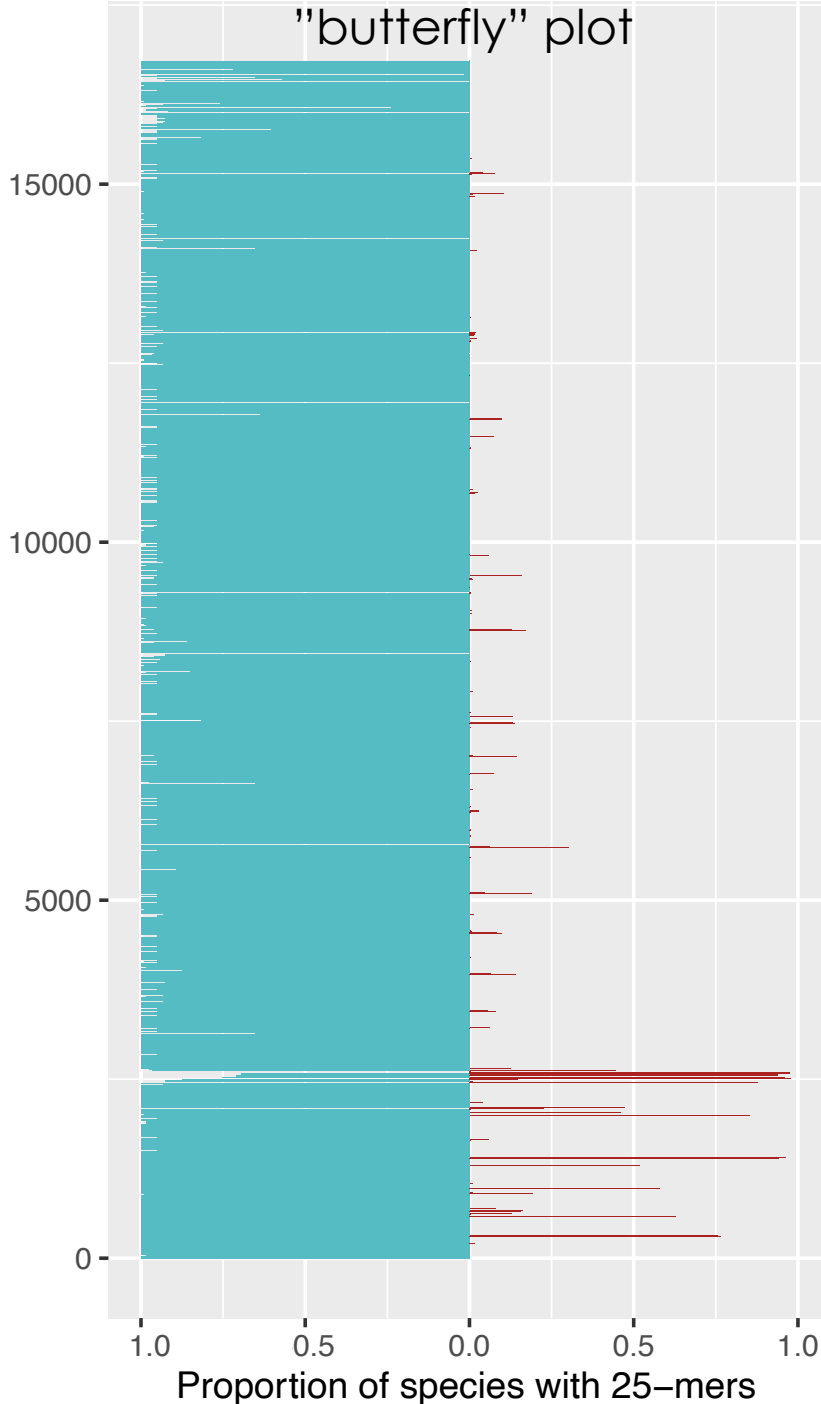
MAFFT 14s

**unikseq 3.2s**



Position of 25-mers on *C. carcharias* (CM030070.1) Mt genome

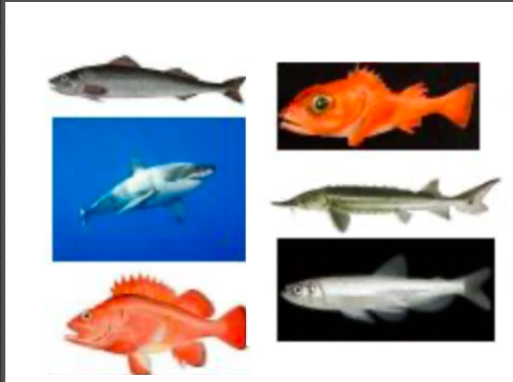
"butterfly" plot



factor(condition)  
teal ingroup-unique  
red outgroup  
(856 fishes)



VS



# UNIQUE REGIONS

- MANY REGIONS UNIQUE IN GREAT WHITE REF  
43 UNIQUE REGIONS (100 BP+), 12.9 KB
- IDEAL FOR DESIGNING **SPECIFIC** qPCR

left (blue) : sequence uniqueness / conservation  
right (red) : reference kmers in outgroup

# UNIQUE + CONSERVED REGIONS



>ch-CACA1-CM030070.1region11205-11314\_size110\_propspclN98.9\_propunivsOUT100.0\_avgOUTentries0.0

ccgtaagccacataggcctggtagcaggagcaatcctaatacaaacacacatgaagcttgcaggagcaattacactaataattgccacaggcctgatttcacccgcca

>ch-CACA1-CM030070.1region11317-11476\_size160\_propspclN98.9\_propunivsOUT100.0\_avgOUTentries0.0

ctgcttagccaaactaactacgagcgaatccacagccgaacaataactcctagcccaggcATACAAATTATTTTCCACTAATAGCTACCTGATGATTCCTTGCCATCCTAGCCAACCTGCCCCTCCC  
CCCTCTCCCAATTTATAGGAGAACTCCTT

>ch-CACA1-CM030070.1region11480-11719\_size240\_propspclN99.3\_propunivsOUT100.0\_avgOUTentries0.0

ATTACCTCCTTATTAACTGATCCAACACTGAACACTATTACCCTCTCGGGCCTTGGAGTATAATCACCGCTTCCTACTCCCCTTATATATTCCTAATAATTCAACGCGGACCAACccccctaccatatac  
ttatcattaaatccaagccacACACGAGAACACCTCCTCCTGAGCCTCCACCTCATGCCCGTCTACTTCTAATATTTAAACCAGAACTTATCTGAGGCTGAACACTCTGTGCT

>ch-CACA1-CM030070.1region11724-12333\_size610\_propspclN93.4\_propunivsOUT100.0\_avgOUTentries0.0

GTTAACCAAAAACATTAGATTGTGGTCTAAAAATAAAAGTAAAAccttttaaccaccgagagaggctggggacatgaagaactgctaattcttctcatcatggctcaaatccatgactcactcagcttctgaaa  
gatatacagtaatactattggcttaggaacccaaaacccctgggtgcaactccaagcaaaagctatgaatactatcttfaactcactctctcctaatcttggttatctctctctccactaataacctcattaaatcccaagaaact  
cactcctaactgggctctctcctacgcaaaaacagctgtaaaaatctccttctcattagccttaccactatccatttcttagaccaagggttagagtcataataccaactacaactgaatcaacattgggCCTTTTCTGA  
TATAATATAAGCTTCAAATTTGATACATACTCTGTTCTATTACCCCTGTGGCCCTCTACGTTACTTGATCTATCCTGAATTTGCCCTATGATACATACTCCGACCCAAAcatcaaccgcttcttt  
aaatatacttCTACTCTCCTAATCTCAATAATCATTTTAGTGACTGCCAAC

>ch-CACA1-CM030070.1region12339-12756\_size418\_propspclN98.4\_propunivsOUT100.0\_avgOUTentries0.0

ATTCCAACCTGTTTCATTGGCTGAGAGGGAGTTGGAATCATATCATTCTCCTCATTGGTTGATGACATAGTCGAACAGACGCCAACACAGCCGCCCTACAAGCTGTAATCTATAACCGAG  
TAGGAGATATTGGcctaattcttagcatagcttgactggctataaattaaactcctgagaatcaacaactatcattatccaagacataaacttaacctacctctcttggcttctgctagccgcagctggaaat  
ccgcacaattcggccttaccatgacttccatcggccatagaaggaccacaccagctcctcgccttactccactccagcacaatagttgtagccggtatcttctcctaattcgccttaccataatccacaataatca  
ttaattctaaca

>ch-CACA1-CM030070.1region12936-13134\_size199\_propspclN97.5\_propunivsOUT100.0\_avgOUTentries0.0

taaagcaatacttttctctgttccggatccattatccatagcctcaatgacgaacaagatatccgcaaaataggagggtcttcaaaactcctgccattcacctcctccttaactattggaagcctagccctcataggcatg  
cccttcttatcaggetcttttcaaaagatgccatcatcgaagccataaacact

>ch-CACA1-CM030070.1region13136-13248\_size113\_propspclN98.5\_propunivsOUT100.0\_avgOUTentries0.0

ctcacctcaacgcctgagcccttactcttaccctaattgcaacatcattacagccatttacagcctccgctcattttttcacattaataaattttccacggtttaattca

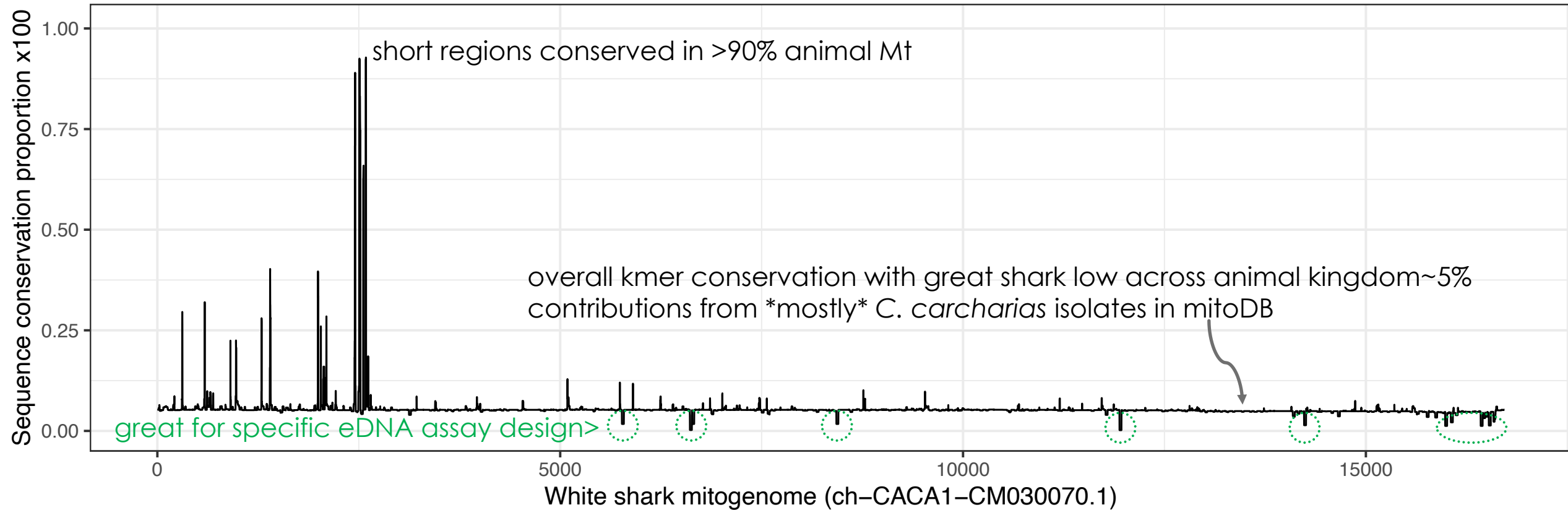
>ch-CACA1-CM030070.1region13267-14066\_size800\_propspclN98.9\_propunivsOUT100.0\_avgOUTentries0.0

aaacacccctaatattaaccccatcaaacgcctagcttacggaagtatcctgtctggcctcattatcacatccaacataatccctacaaaaccccaattatgactataattcccctactaaactctccgcccactaatt  
actattgctggccttctactagccttagaactagcaaacctaacctcagcttaaaataacccccaccctctatccccatcacttctcaaatatattgggatactttcccaaatatccaccgctcctgcccdaaatta  
acctatcctgagcccaacatgthccaccacctaattgaccaaacatggctctgaaaaaattggaccaaaaagtgccttattcaacaaacccctcctaattaaattatccaccaacctcaacaaggctatattaaagthta  
cctaacaactacttttctcacactagccttagccatactcactacattaacctaacccacagtaatgtccctcatgcttagacctCGAGTCAACTCCAACACCACAAACAAAGTCAATAATAACACCCACCC  
ACTTAAACTAACAACCACCCCCCATCCTCATAAAGCAAAGCCACCCCCACAAAATCCCCACGAGTTATCTCCATATTGCTCAACTCCTCTACCCCTGACCAATCCAACCTCAAATC  
ACTCTACCATAAAATAATTACCAACAAAAATAACTACTAAATAAAAACCAACATACAACAAAACAGATCAATTACCCCATGACTCAGGATATGGCTCAGCAGCAAGAGCTGCCG



# SEQUENCE CONSERVATION ACROSS 4,000 ANIMAL MITOGENOMES

comparing 65,000,000 bp



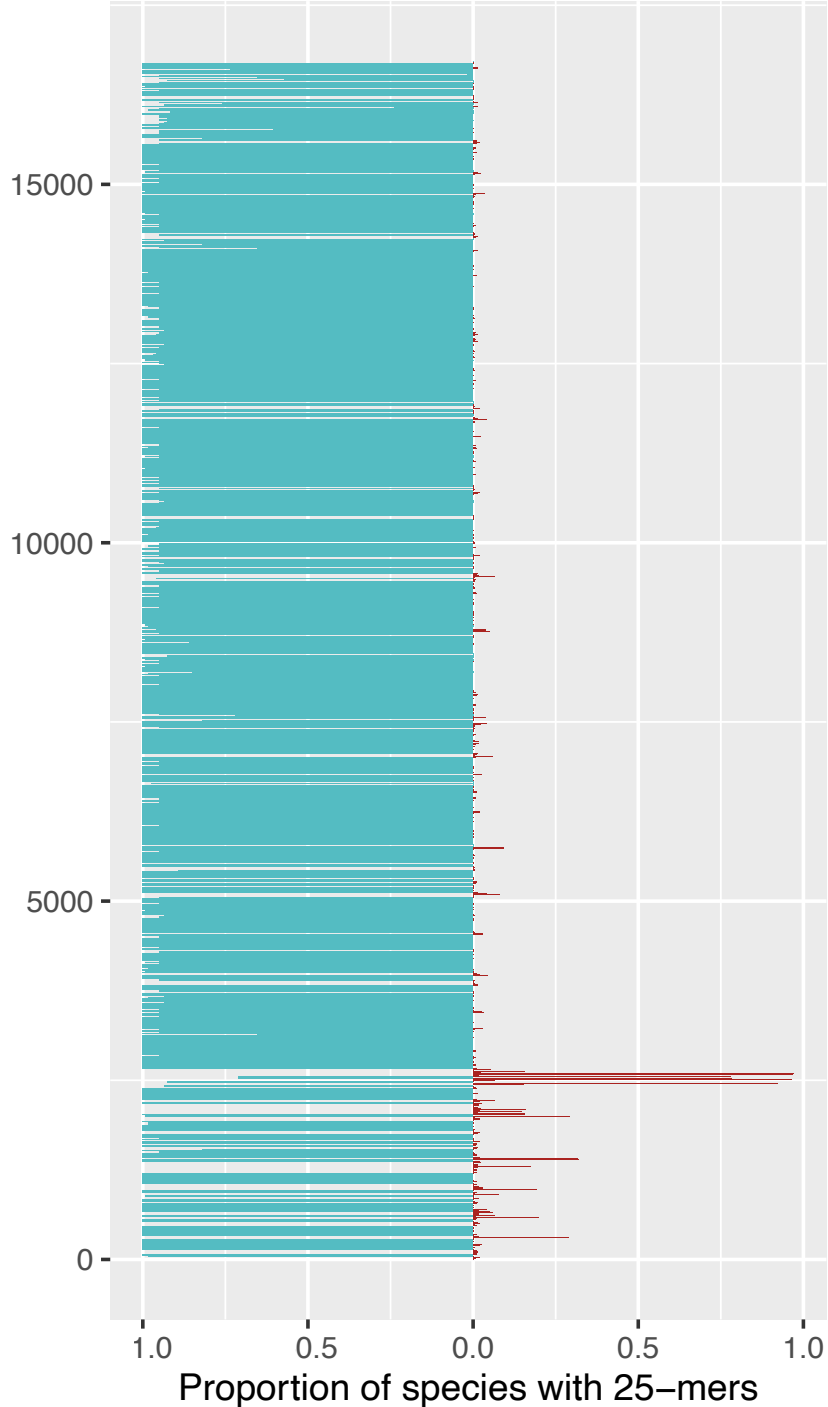
clustal >3weeks

MUSCLE error

MAFFT 9h50m, **4.4GB RAM**

**unikseq 2m, 9.0GB RAM**

Position of 25-mers on *C. carcharias* (CM030070.1) Mt genome



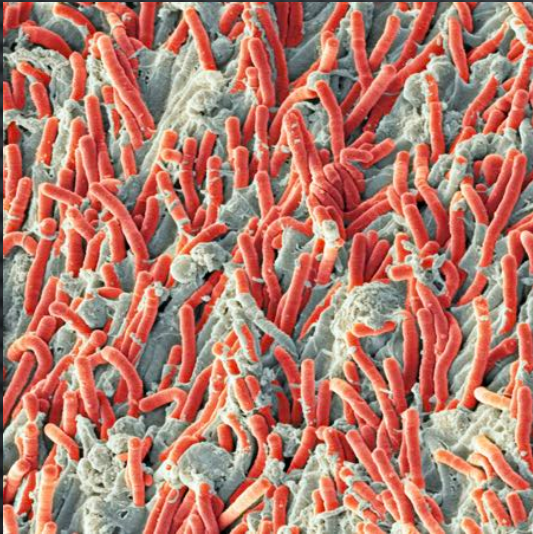
VS



# UNIQUE REGIONS

- LARGER OUTGROUP ~ RESTRICTIONS FOR UNIQUE REGION IDENTIFICATION
- 29 UNIQUE REGIONS (100 BP+), 5.2 KB  
3 : 100% AVG INGROUP CONSERVATION, 461 BP
- GREAT WHITE SHARK MITOGENOMES HARBOR MANY UNIQUE REGIONS CONSERVED ACROSS ISOLATES

# FUSOBACTERIUM NUCLEATUM



“Dental plaque bacteria hitch a ride in blood to fortify colon cancers” — *Cosmos mag.*

> [Genome Res.](#) 2012 Feb;22(2):299-306. doi: 10.1101/gr.126516.111. Epub 2011 Oct 18.

## Fusobacterium nucleatum infection is prevalent in human colorectal carcinoma

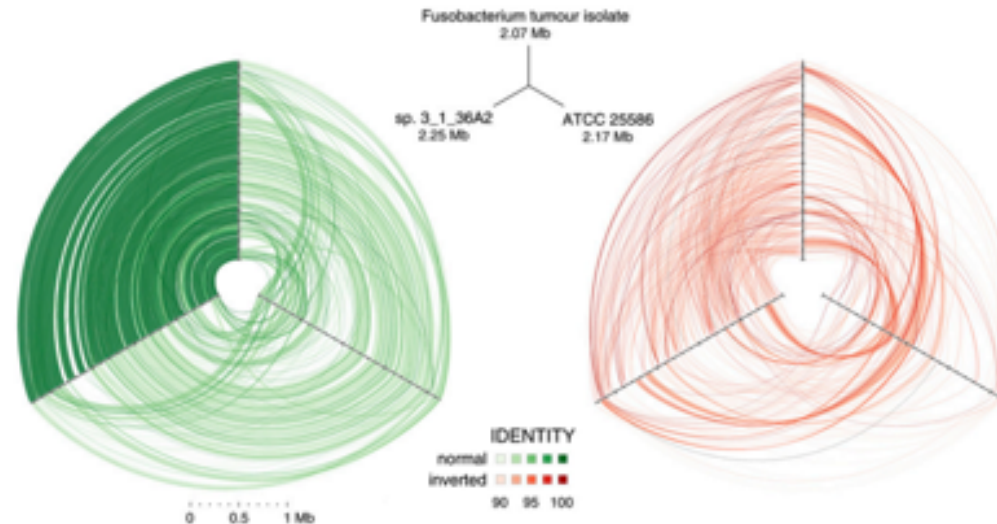
Mauro Castellarin<sup>1</sup>, René L Warren, J Douglas Freeman, Lisa Dreolini, Martin Krzywinski, Jaclyn Strauss, Rebecca Barnes, Peter Watson, Emma Allen-Vercoe, Richard A Moore, Robert A Holt

2.2 MBP CIRCULAR  
GENOME

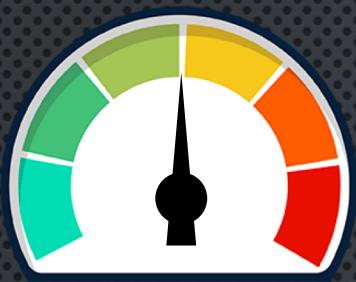
ASSOCIATED WITH  
CRC POLYPS

TUMOUR ISOLATE  
SEQUENCED 2012  
ASSEMBLED WITH  
SSAKE : **CC53**

SIMILAR TO *F. N.*  
VINCENTII AND ATCC  
(ORAL STRAIN)

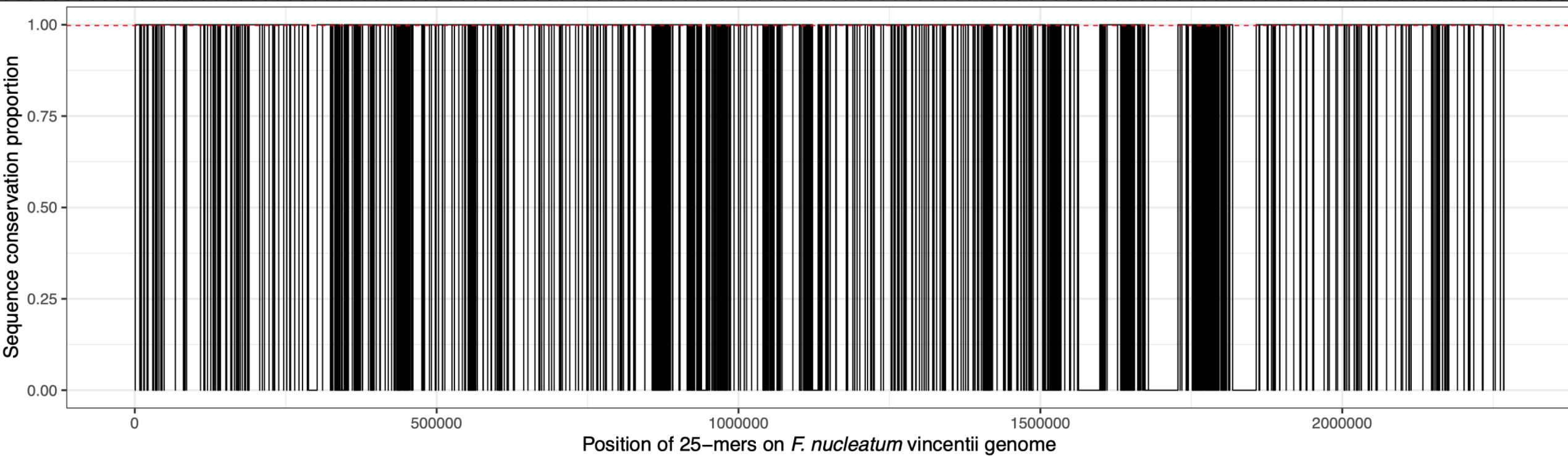


**Figure 3.** Hive plots showing alignment of three *Fusobacterium* genomes. Approximately 32 million high-quality WGS Illumina HiSeq reads ( $\geq 99$  consecutive Q30 bases) from *Fusobacterium* tumour isolate CC53 were assembled with SSAKE (v3.7, default options) into 379 contigs. The contigs were aligned using cross\_match (-minmatch 29 -minscore 59 -masklevel 101) to the complete *F. nucleatum* subsp. *nucleatum* ATCC 25586 genome and, independently, to the 12-contig HMP *Fusobacterium* sp. 3\_1\_36A2 assembly, respectively; and ordered/oriented based on the highest identity to the latter sequence. Three-way cross\_match (<http://www.phrap.org>) alignments between each *Fusobacterium* genome were performed and represented visually using hive plots (<http://www.hiveplot.com>). For each, the top, left, and right axes are proportional to genome size and represent the *Fusobacterium* tumor



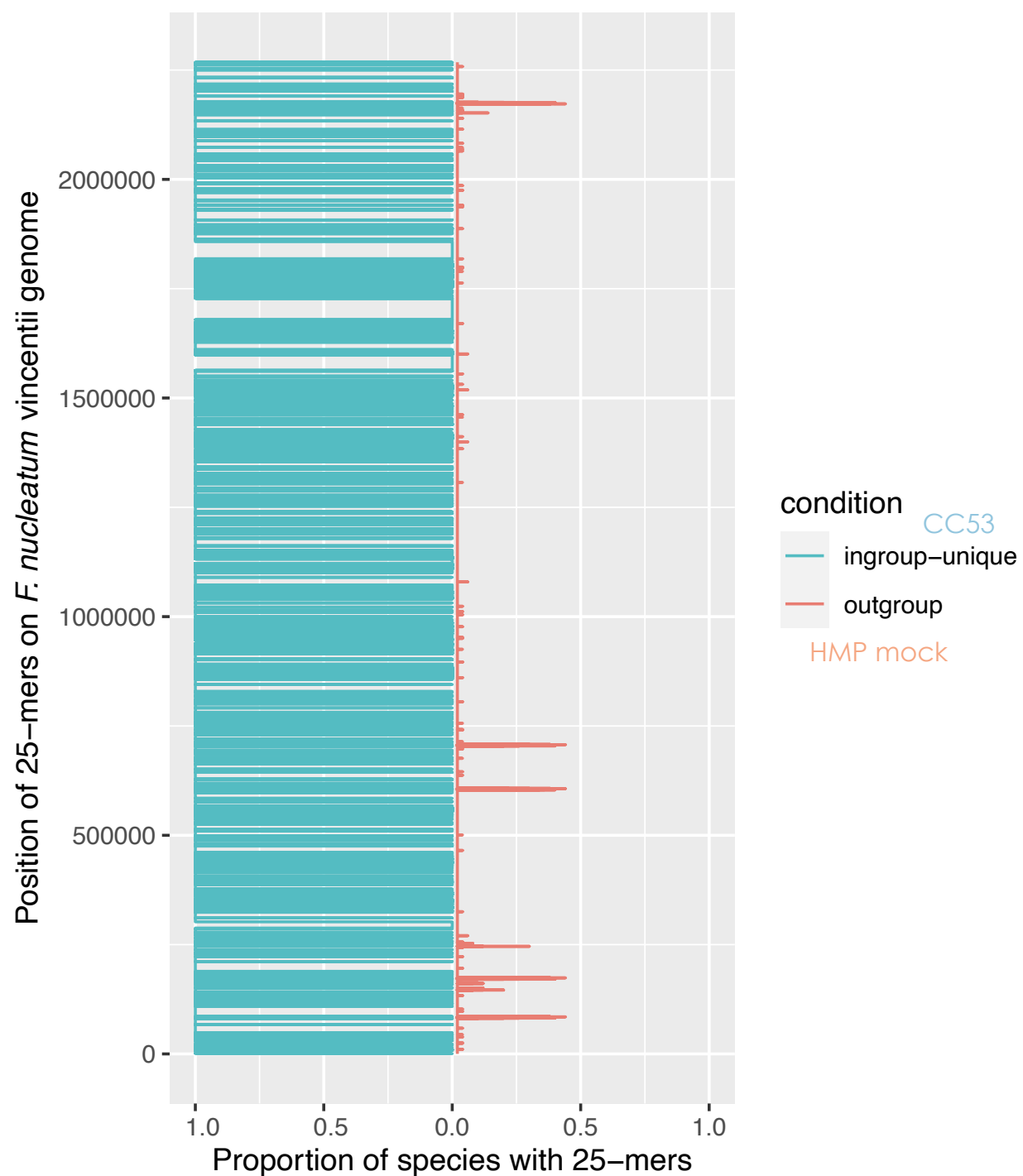
# CONSERVATION *FN. VINCENTII* VS CC53

comparing 90,000,000 bp



1 REFERENCE (*F. NUCLEATUM VINCENTII* FROM HMP)  
1 INGROUP ENTRY (*F. NUCLEATUM* CC53 CRC ISOLATE)  
50 OUTGROUP ENTRIES (HMP MOCK + *FUSOBACTERIUM* ATCC25586)

clustal >3weeks    MUSCLE error    MAFFT >3weeks    **unikseq 5m, 54GB RAM**



# UNIQUE REGIONS

## UNIQUE+CONSERVED WITH VINCENTII

n	Min	N50	Max	Sum
3,357	100	188	6,503	646,774

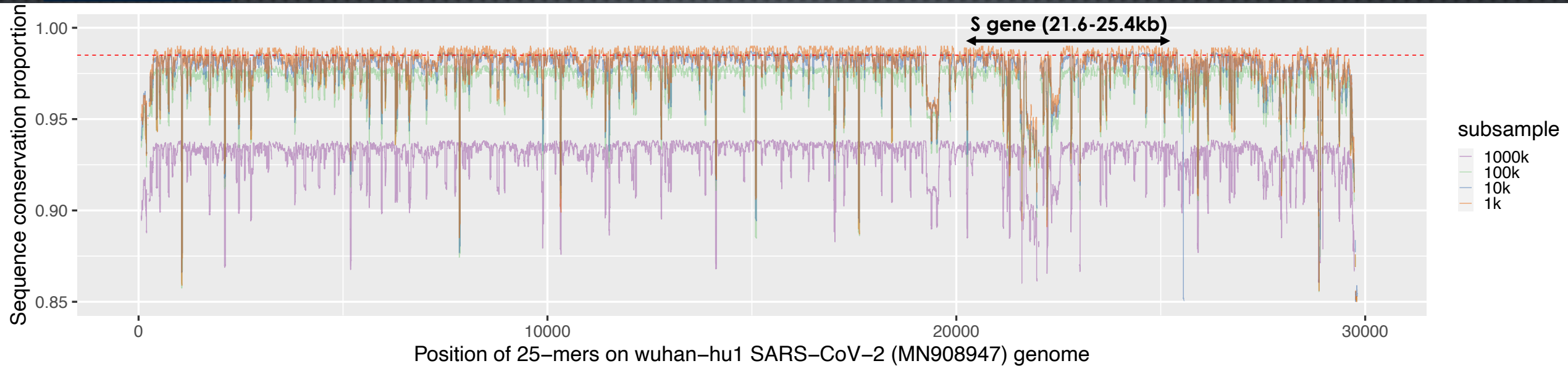
## UNIQUE TO CC53

n	Min	N50	Max	Sum
81	100	457	3,570	21,720



# 1 MILLION SARS-CoV-2 GENOMES

comparing 30,000,000,000 bp



1 CPU @ hpcg02

Subsample (k)	Run time (hh:mm:ss)	Memory (GB)
1	00:00:55	3.2
10	00:12:28	35.9
100	02:02:11	323.5
1000	11:27:02	1238.5

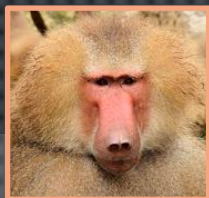
FURTHER BENEFIT FROM :  
1) COUNTING BLOOM FILTERS  
2) PARALLELISM



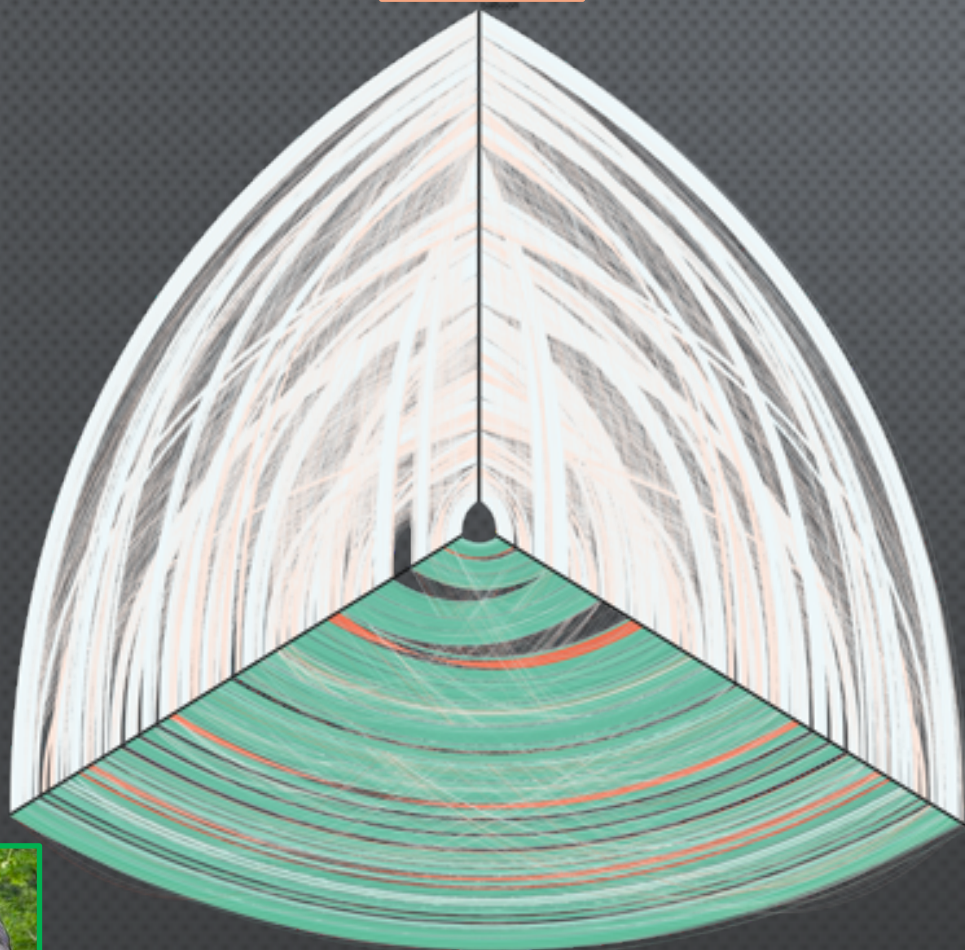


# BONOBO (& HUMAN) VS BABOON GENOMES

comparing 9,000,000,000 bp



USING TWO (INGROUP/OUTGROUP) BLOOM FILTERS



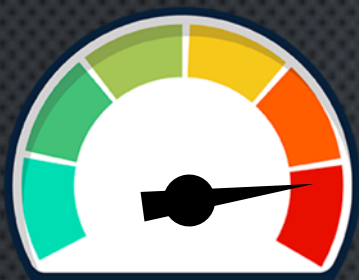
UNIQUE+CONSERVED

n	Min	N50	Max	Sum
415,629	100	123	913	52.9e6 (1.8%)

Bloom filters : 45' (16GB RAM) unikseq : 2h wall clock (21 GB RAM)

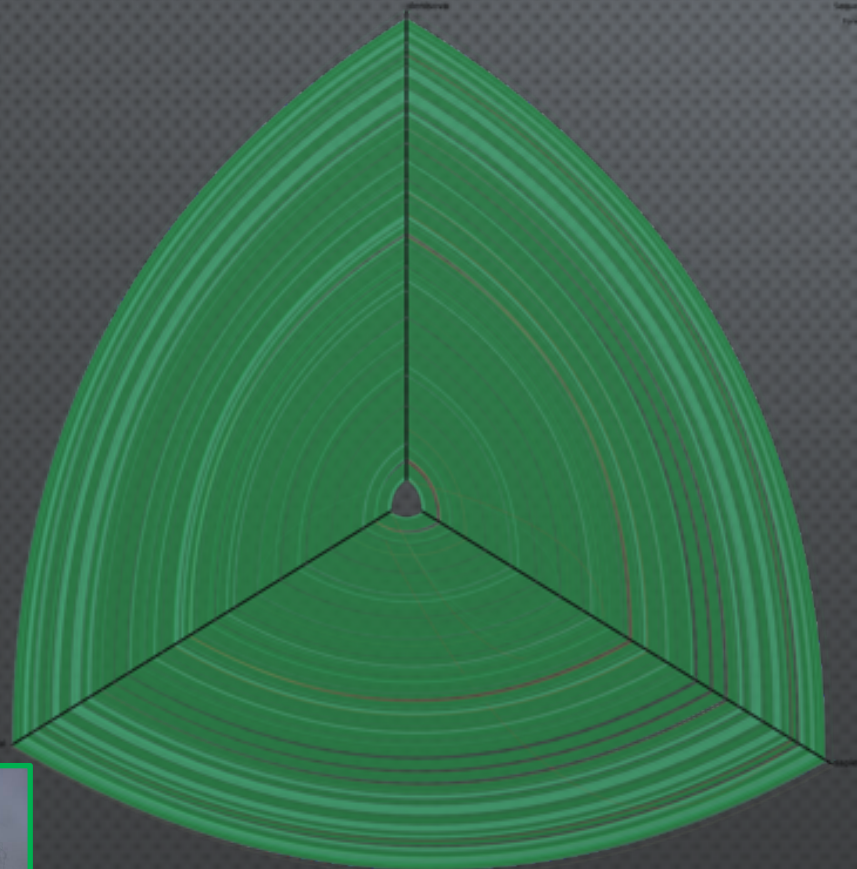
98.8%





# NEANDERTHAL (& DENISOVAN) VS HUMAN GENOMES

comparing 15,000,000,000 bp



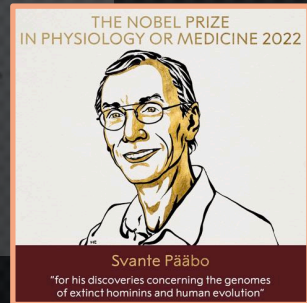
99.5%

## UNIQUE+CONSERVED (100%)

n	Min	N50	Max	Sum
6,314	100	682	5,407	2.5e6 (0.09%)

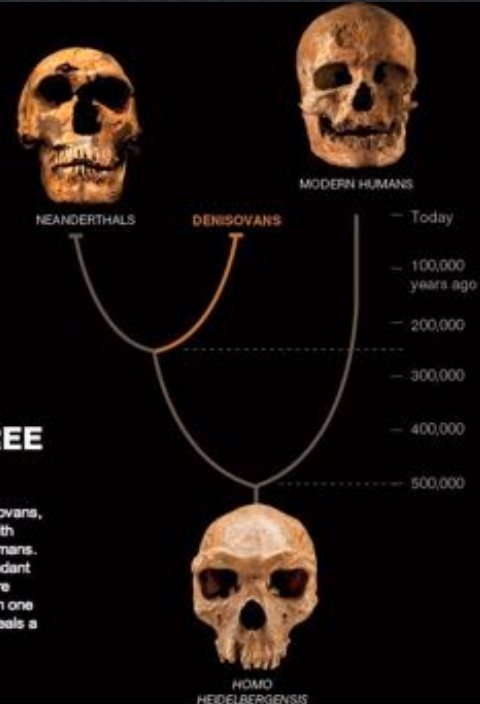
Bloom filters : 44' (16GB RAM) unikseq : 2h wall clock (23 GB RAM)

%	HG02055	HG01243	NA24385
<b>Neanderthal</b>	0.02	0.01	0.02
<b>Denisovan</b>	0.01	0.01	0.01

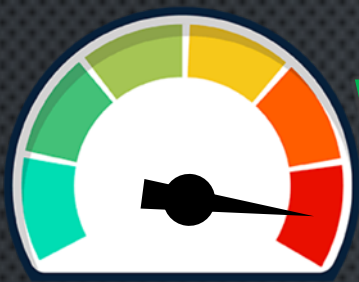


## A TALE OF THREE HUMANS

A third kind of human, called Denisovans, seems to have coexisted in Asia with Neanderthals and early modern humans. The latter two are known from abundant fossils and artifacts. Denisovans are defined so far only by the DNA from one bone chip and two teeth—but it reveals a new twist to the human story.



**THE FAMILY TREE**  
Neanderthals and Denisovans were closely related. DNA comparisons suggest that our ancestors diverged from theirs some 500,000 years ago.



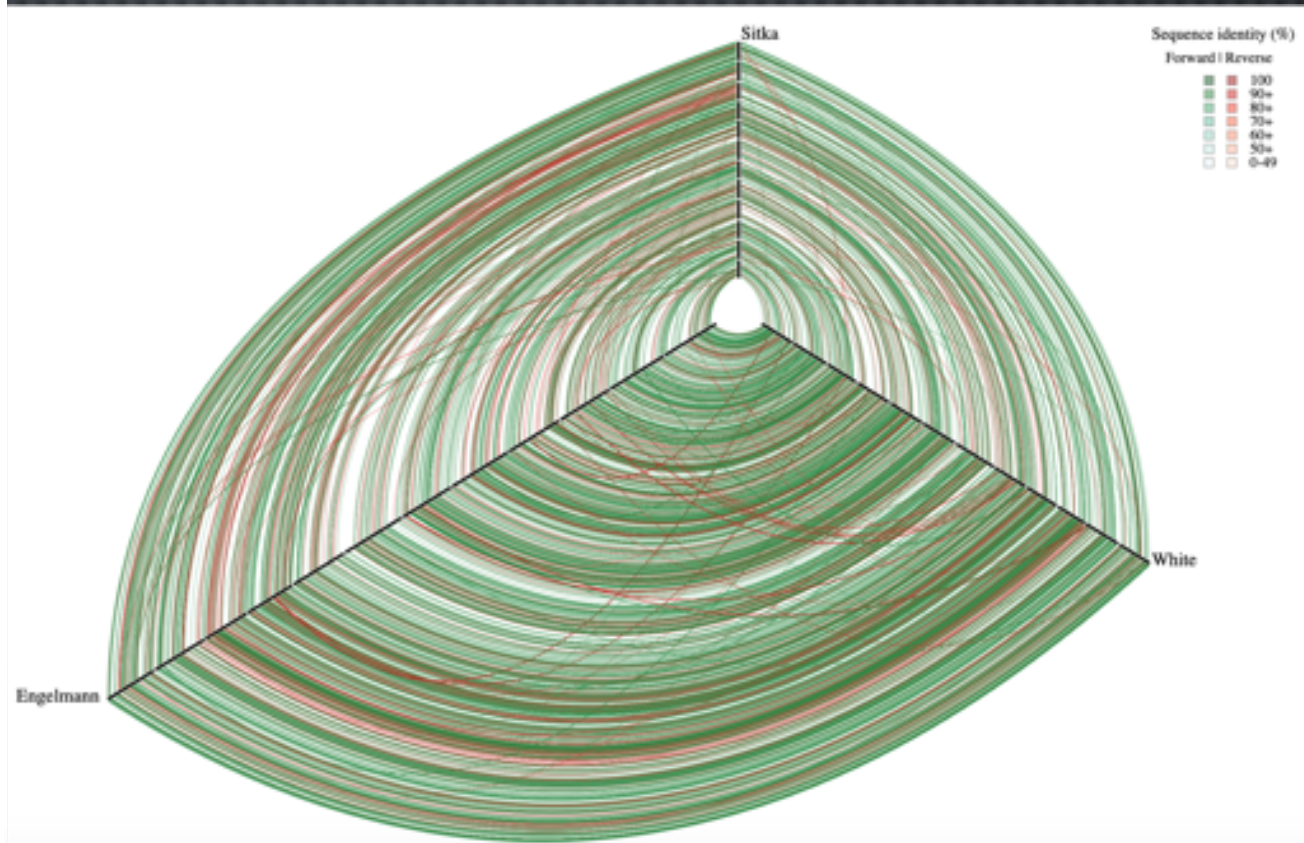
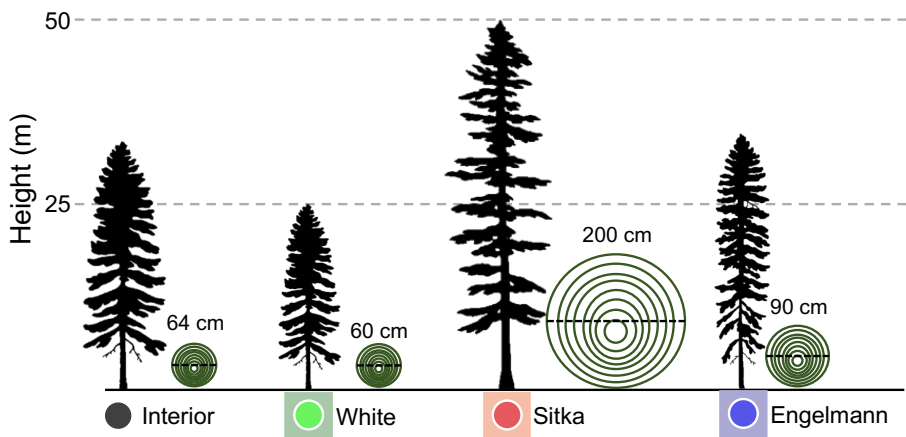
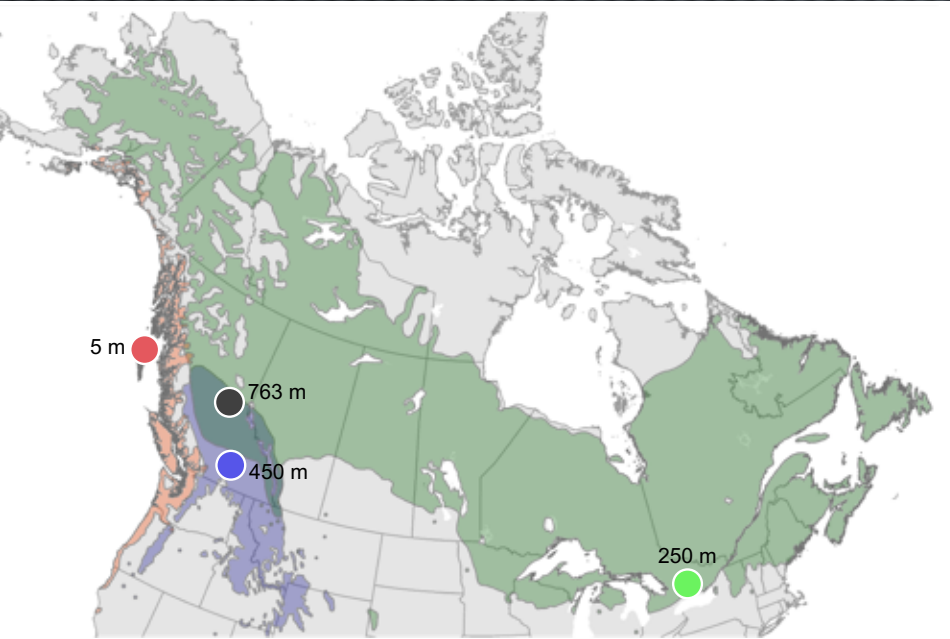
# WHITE (& BC INTERIOR) VS SITKA+ENGELMANN SPRUCE GENOMES

comparing 80,000,000,000 bp

UNIQUE+CONSERVED (100%)

n	Min	N50	Max	Sum
59,565	100	129	1,642	8.3e6 (0.04%)

Bloom filters : 9h (86GB RAM) unikseq : 14h07 (124GB RAM)



# SUMMARY

## NOVEL UTILITY FOR COMPARATIVE GENOMICS

- IDENTIFY UNIQUE / CONSERVED DNA / RNA REGIONS
- K-MER BASED : NO MSA / PAIR-WISE ALIGNMENTS REQUIRED
- FAST & SCALABLE
- NO STRICT INPUT : GENOMES / CONTIGS / READS

## BROAD UTILITY

- SPECIFIC eDNA ASSAY DESIGN
- CLINICAL / DIAGNOSTIC TESTS



# ACKNOWLEDGEMENTS

## GSC

Inanc Birol

Lauren Coombe

Cecilia Yang

## UVic

Caren Helbing

Michael Allison

Neha Acharya-Patel

Louie Lopez



GenomeCanada



Genome  
BritishColumbia

Leading • Investing • Connecting



GenomeQuébec



[github.com/bcgsc/unikseq](https://github.com/bcgsc/unikseq)

**QUESTIONS?**

# UNIKSEQ

## v1.2.2

- OUTPUTS CONSERVED REGIONS
- ROBUST TO DOS/WINDOWS AND UNIX-FORMATTED RNA/DNA FASTA FILES
- HANDLES TRADITIONAL MULTI-LINE FASTA RECORDS (SEQUENCE ON MULTI-LINES)
- VERSION CONTROLLED [HTTPS://GITHUB.COM/BCGSC/UNIKSEQ](https://github.com/bcgsc/unikseq)
- WORKS ON VARIED INPUT, NON-CONTIGUOUS GENOMES / READ SETS

Usage: ./unikseq v1.2.2

-----input files-----

-r reference FASTA (required)

-i ingroup FASTA (required)

-o outgroup FASTA (required)

-----kmer uniqueness filters-----

-k length (option, default: -k 25)

-l [leniency] min. non-unique consecutive kmers allowed in outgroup (option, default: -l 0)

-m max. [% entries] in outgroup tolerated to have a reference kmer (option, default: -m 0 % [original behaviour])

-----output filters-----

-t print only first t bases in tsv output (option, default: -t 25)

-c output conserved FASTA regions between reference and ingroup entries (option, -c 1==yes -c 0==no, [default, original unikseq behaviour])

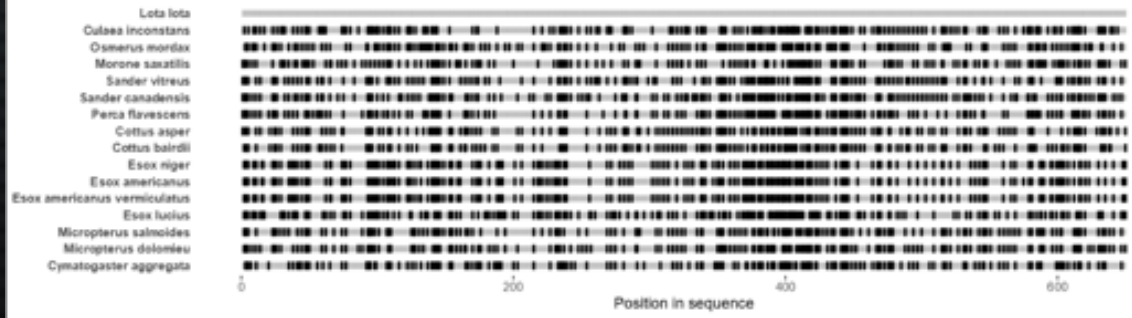
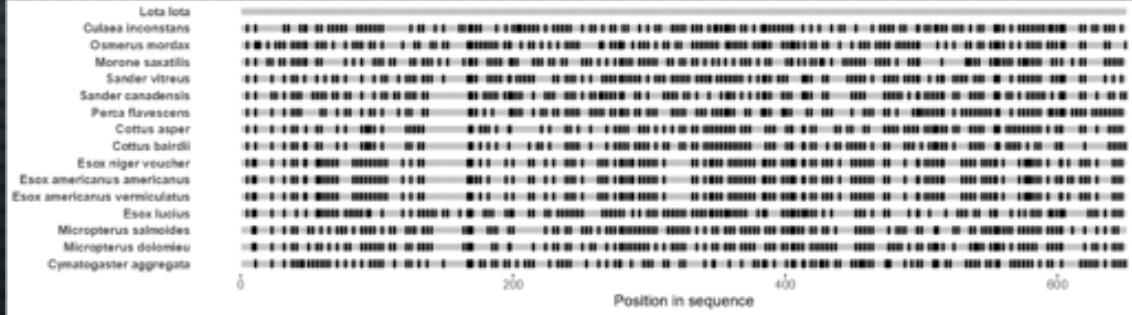
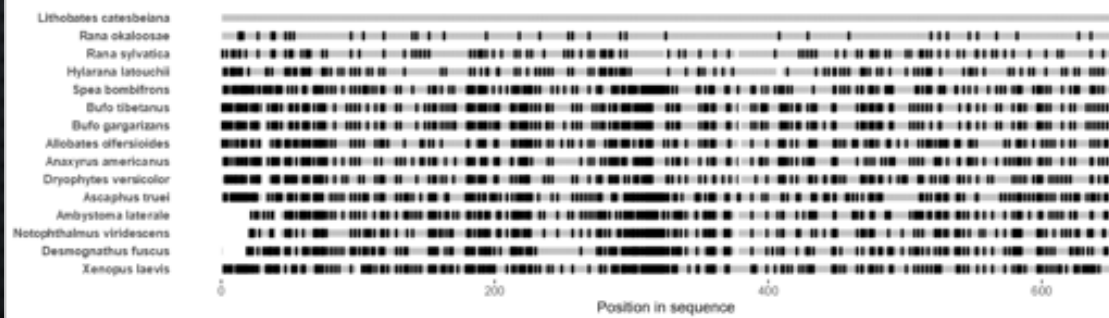
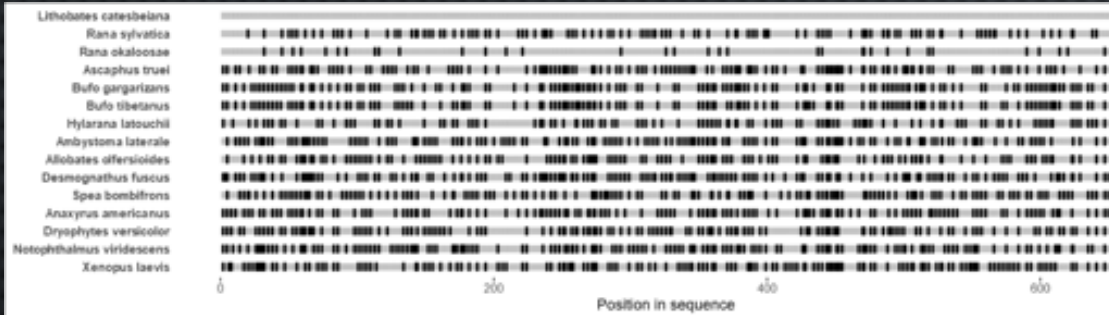
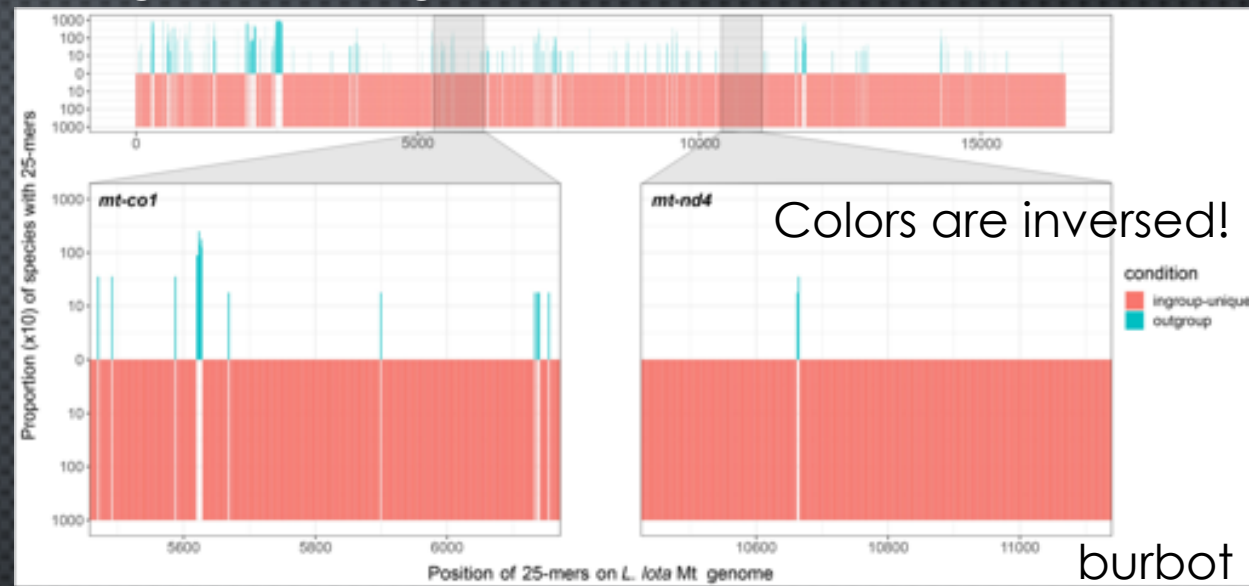
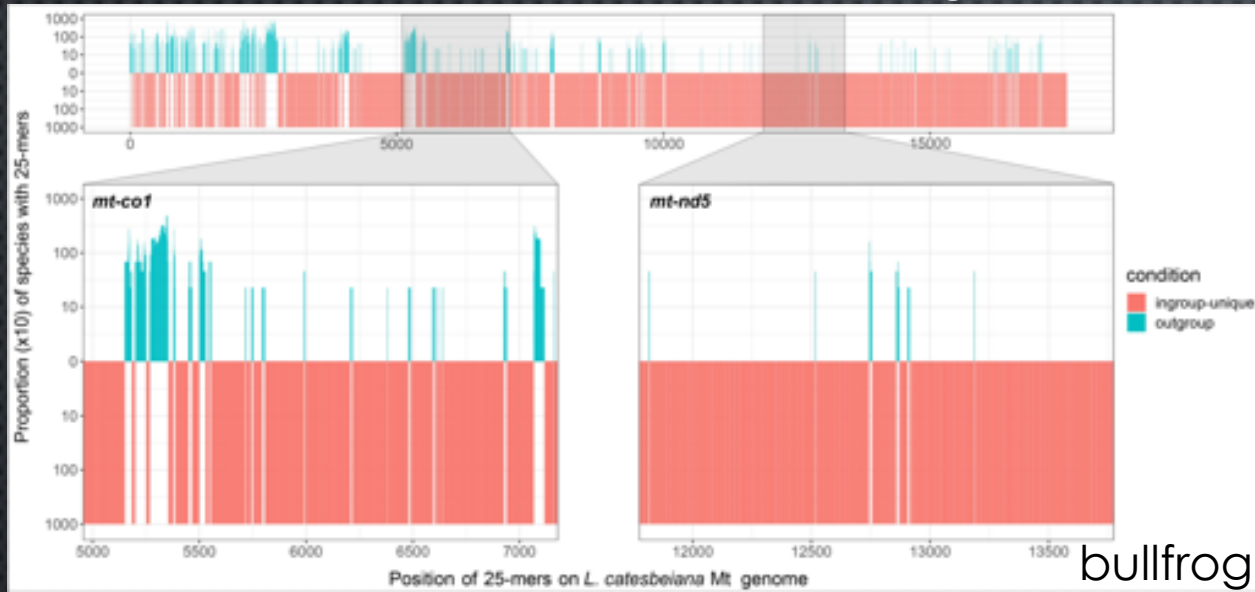
-s min. reference FASTA region [size] (bp) to output (option, default: -s 100 bp)

-p min. [-c 0:region average /-c 1: per position] proportion of ingroup entries (option, default: -p 25 %)

-u min. [% unique] kmers in regions (option, default: -u 90 %)

# eDNA ASSAY DESIGN : CASE STUDY

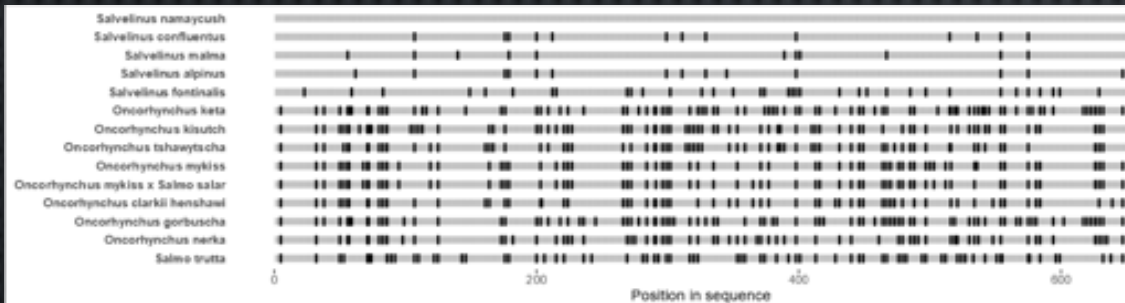
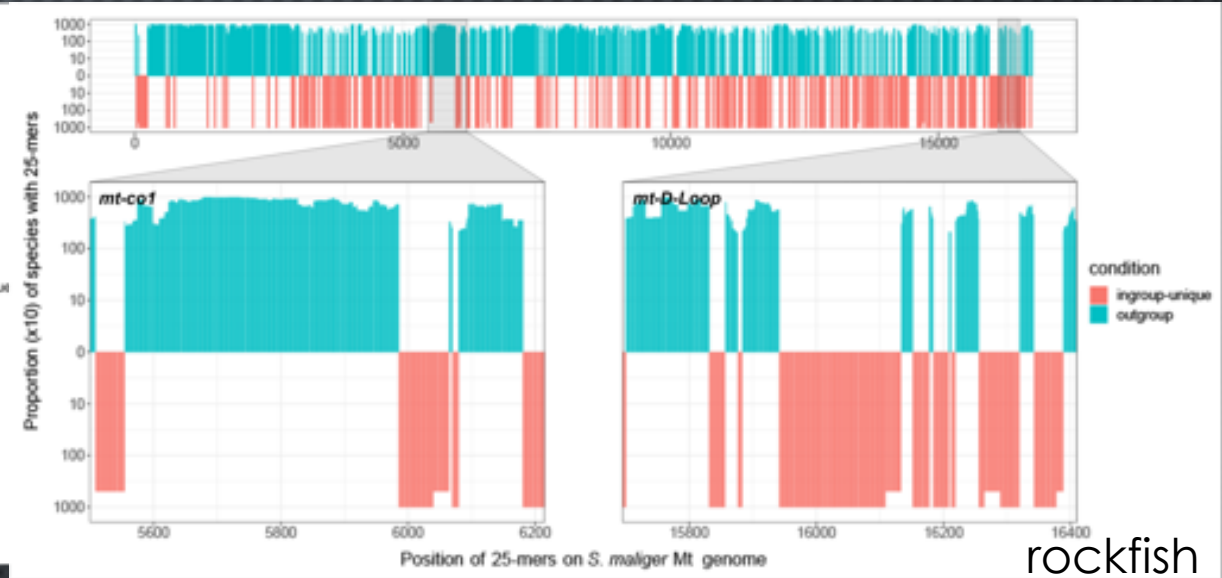
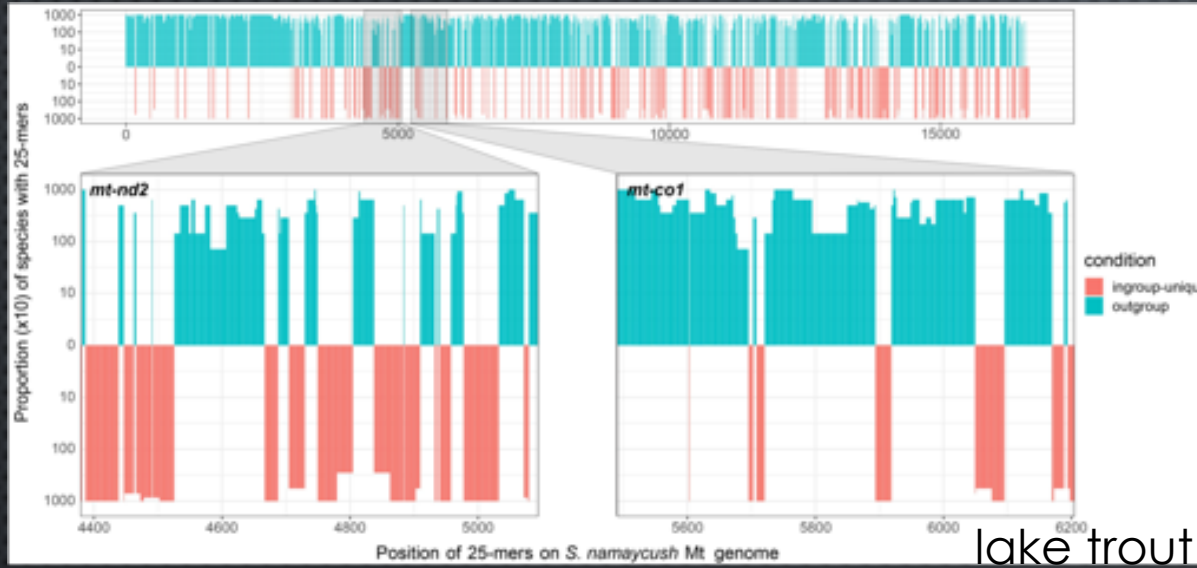
american bullfrog and burbot : strong assay design potential



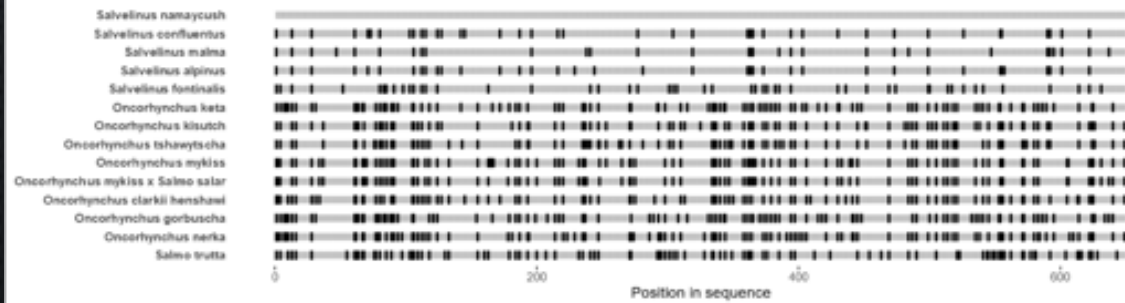


# eDNA ASSAY DESIGN : CASE STUDY

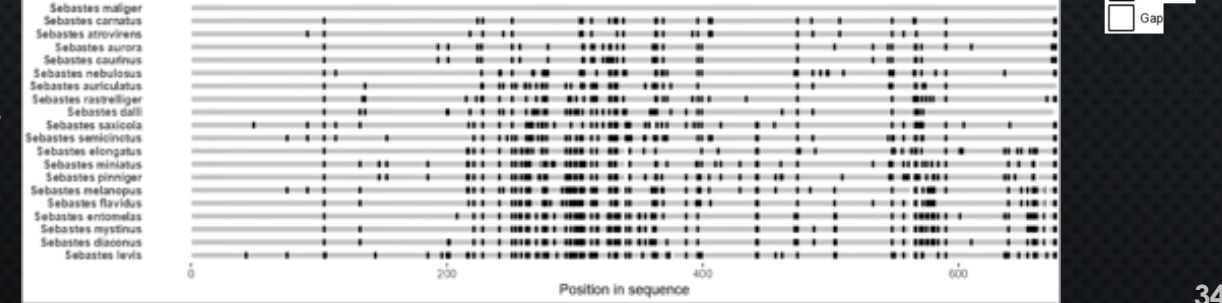
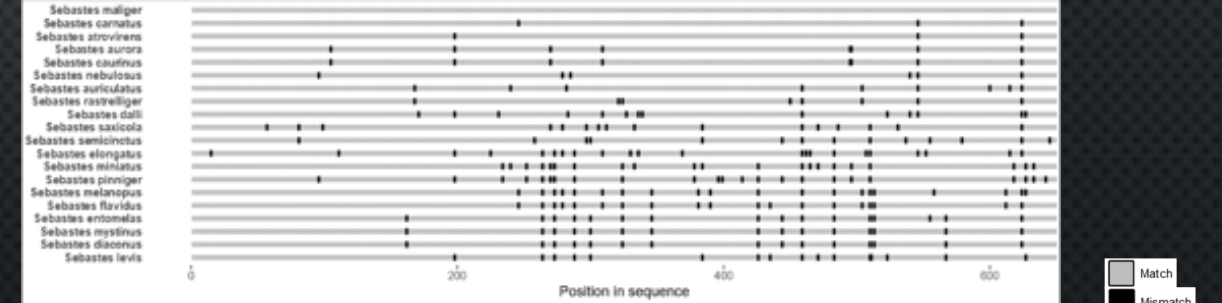
trout and rockfish : closely related taxa limit mitogenome qPCR design potential



mt-coi1 barcode



unikseq



Match  
Mismatch  
Gap