

Ethan Hu 2024 Summer at  
Sym Geno Evo Lab

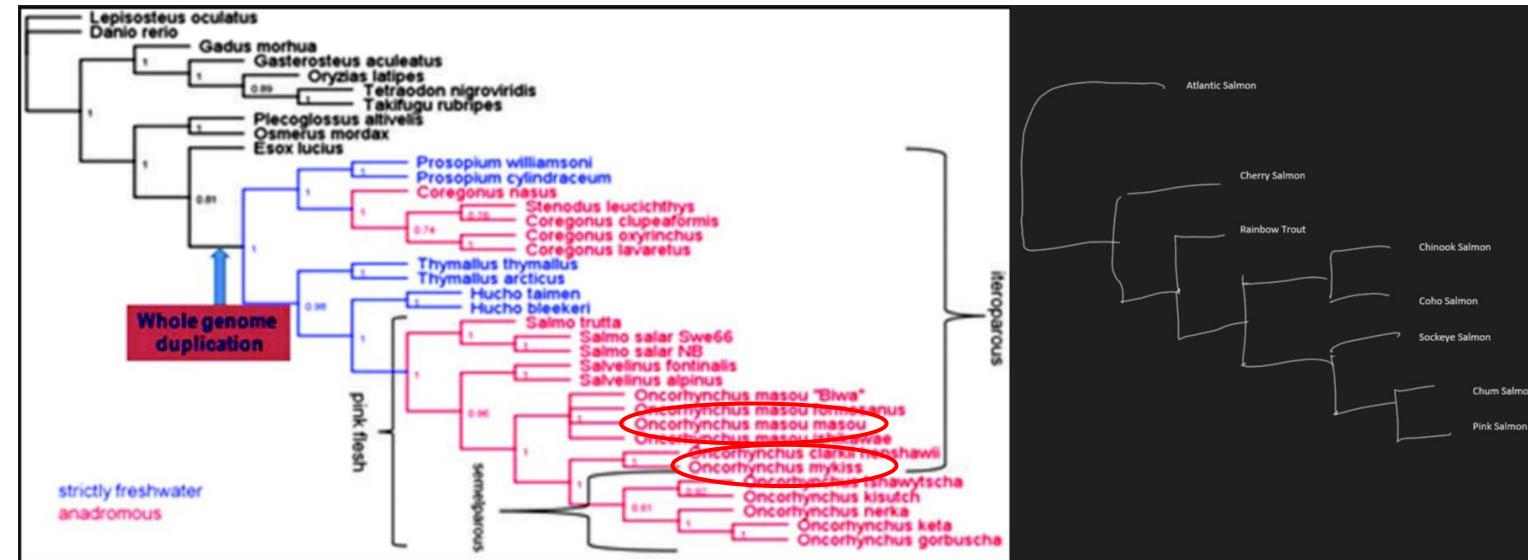
# Overview

- Summary of Phylogeny Tree and BRAKER installation
- Literature Reviews
  - Yabin Guo, An overview on the DNA nucleotide compositions across kingdoms
  - James C. Schnable, Genes and gene models, an important distinction
- Macrosynteny Plot
- Results
- Future Work

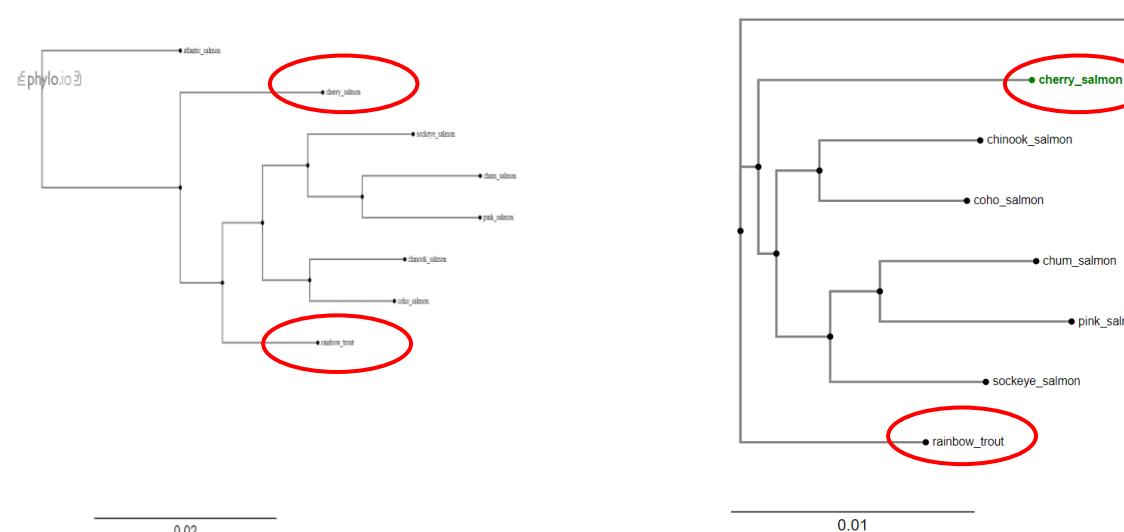
# Summary of Phylogeny Tree and BRAKER install

## Things of Interest:

- Phylogeny Tree pipeline doesn't work as well for closely related species
- BRAKER3 installation:
  - When installing perl in conda, had to use version 5.26.2, to accommodate for perl scalar util numeric
  - StringTie2, Bedtools, GFFRead, Diamond, SRA Toolkit, samtools all in Genemark-ETP tools
  - When using compleasm in BUSCO mode, delete the Darwin folder in sepp/tools/bundled



Davidson, William S. *Genome*, vol. 56, no. 10, Oct. 2013, pp. 548–550, doi:10.1139/gen-2013-0163.



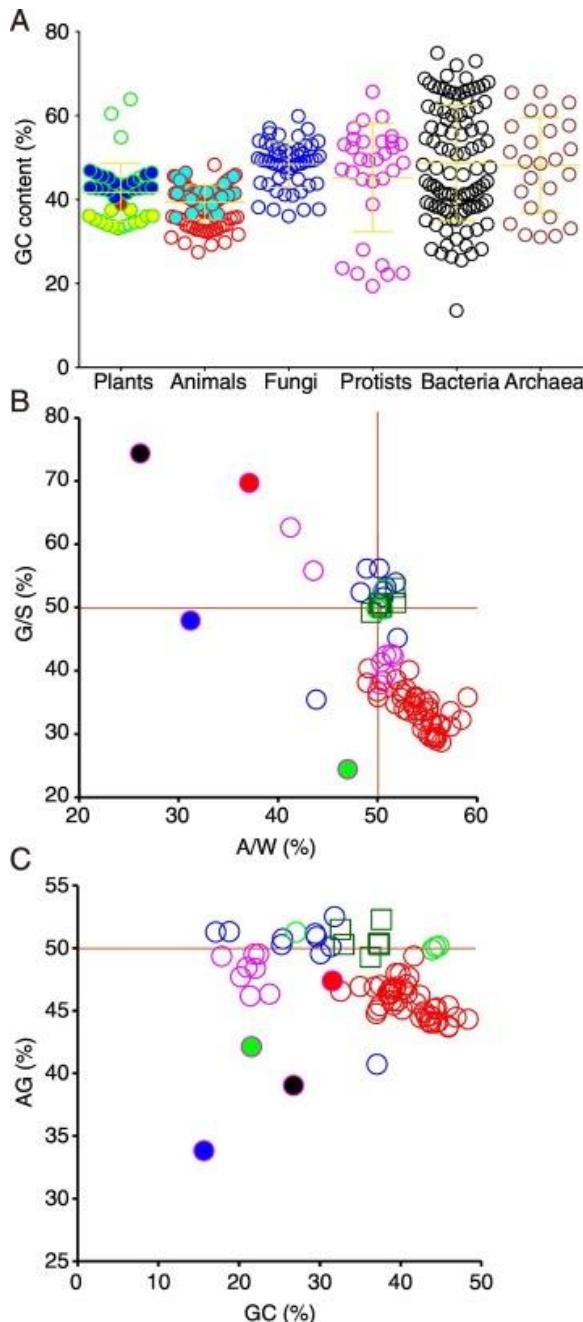
Orthofinder Tree

Full Pipeline

# Literature Review

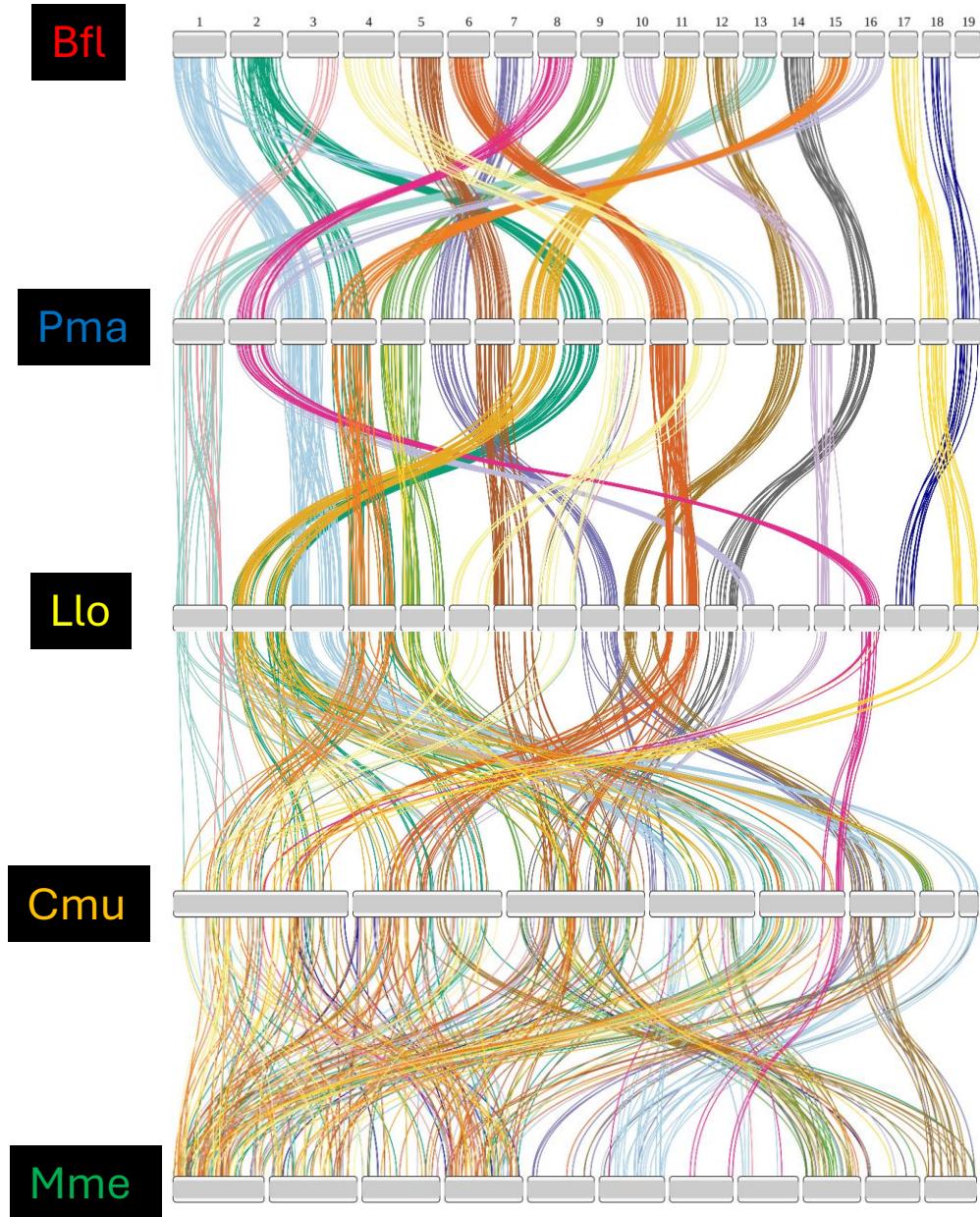
## An overview on the DNA nucleotide compositions across kingdoms

- Yabin Guo, Gene Reports
- GC content: Proportion of GC in DNA strand
  - Chargaff's first rule, # of G = # of C, # of A = # of T
- GC skew and AT skew:  $(G-C)/(G+C)$  and  $(A-T)/(A+T)$ , measures asymmetry of the nucleic acids
  - Chargaff's second rule,  $G/C \sim 1$ ,  $A/T \sim 1$ 
    - Substitution Rule
    - Symmetry is higher entropy
- Purine content: Proportion of AG in DNA strand
  - Szbalski's rule, DNA template strands have higher purine content <- politeness hypothesis
- Conclusion
  - Chargaff's second rule holds true
  - Some protists have asymmetric chromosomal strands
  - Asymmetric satellite DNA regions occurred in some animal genomes
  - Szbalski's rule is not universal, purine richness of coding regions due to energy requirements
  - Thermophilicity is correlated to both genome GC content and CDS AG% in archaea



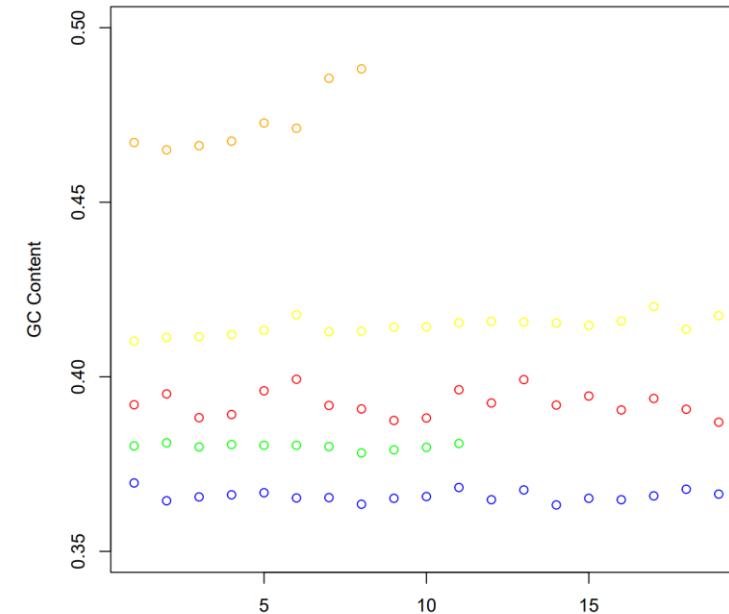
Category	Properties
GC	- Triple Hydrogen Bond
AT	- Double Hydrogen Bond
Purine (AG)	- Double Ring - Heavier
Pyrimidine (CT)	- Single Ring - Cheaper to produce

# Macrosynteny Plot

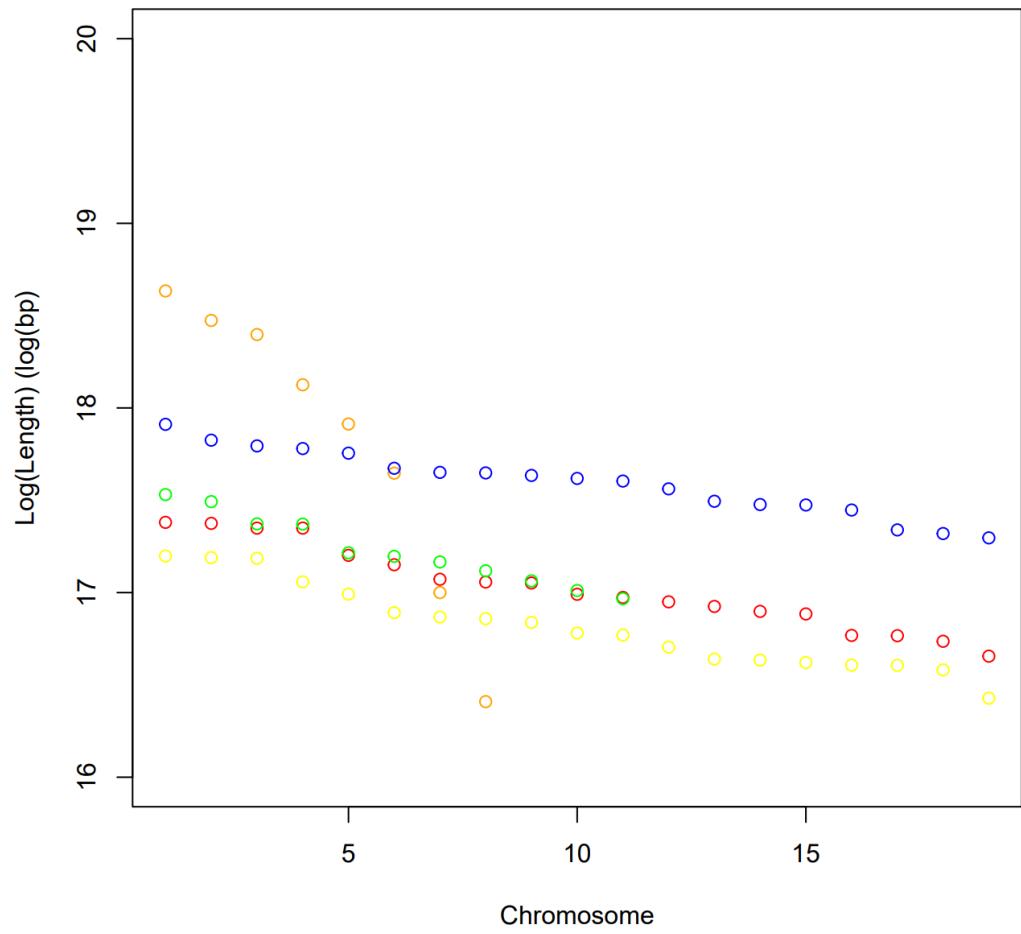


Name	Source
<i>B. floridae</i> , Bfl	<a href="#">Ferdinand</a> <a href="#">Marlétaz's lab</a> + SyntenFinder default
<i>P. maximus</i> , Pma	NCBI
<i>L. longissimus</i> , Llo	NCBI
<i>C. mucedo</i> , Cmu	GEO + Dryad
<i>M. membranacea</i> , Mme	NCBI + Dryad

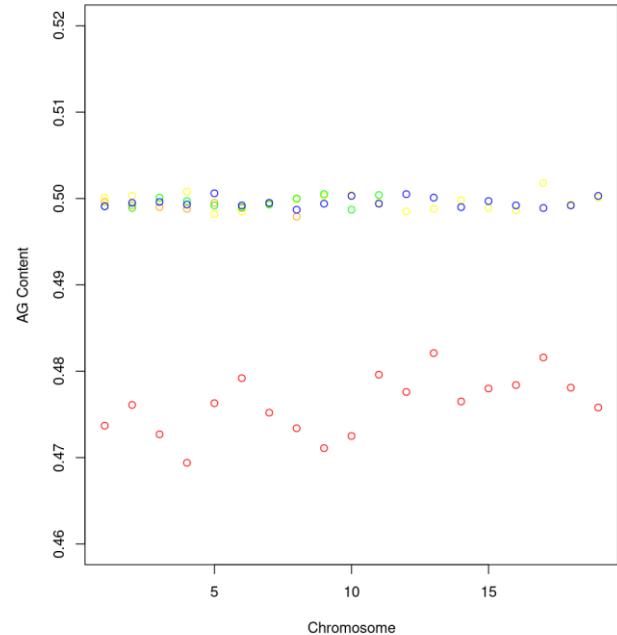
GC Content vs Chromosome



Sequence Length vs Chromosome

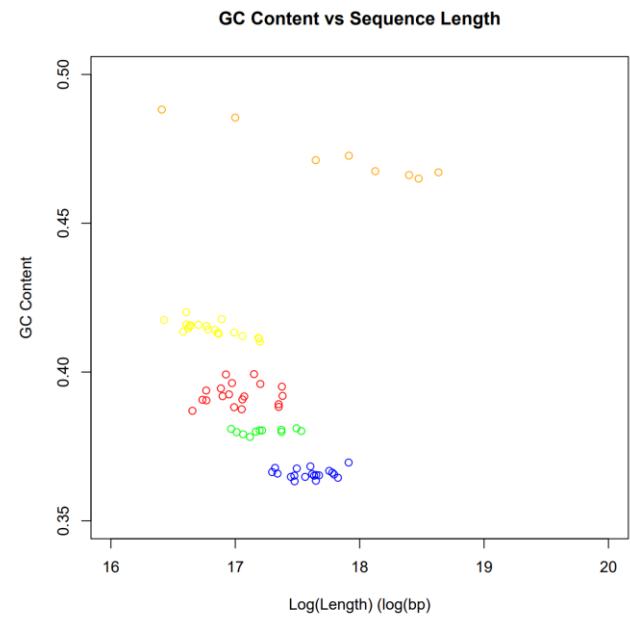
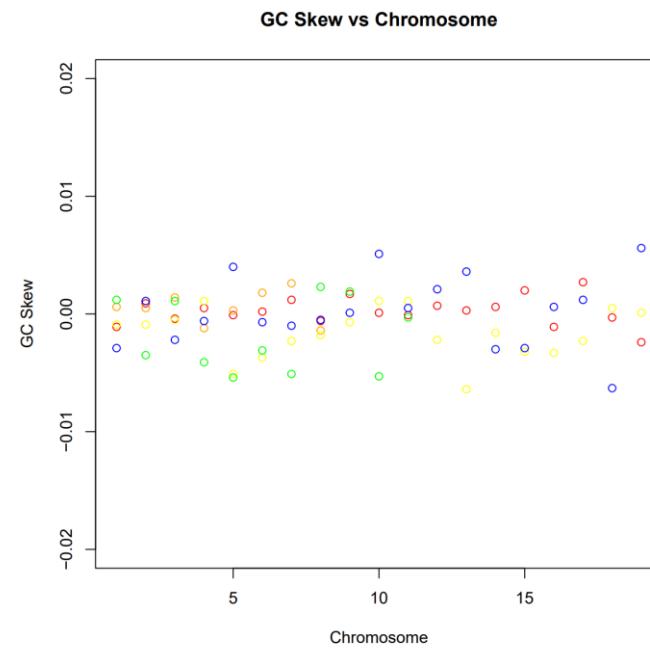
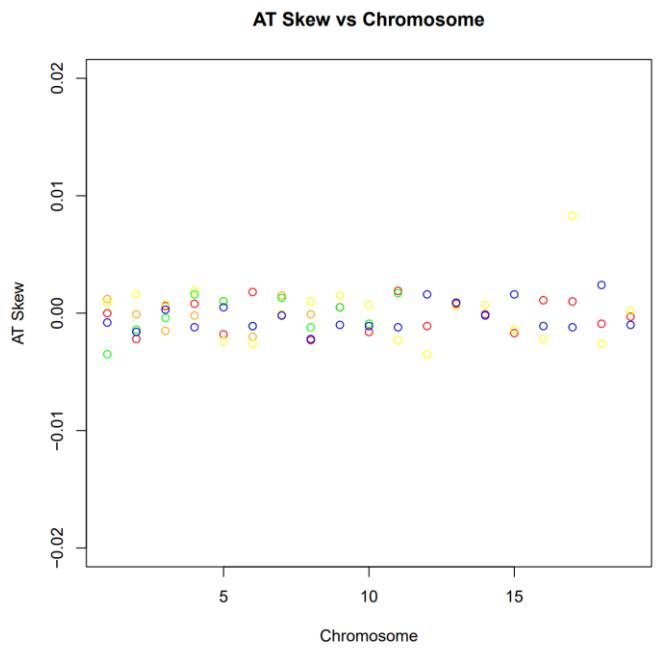


AG Content vs Chromosome



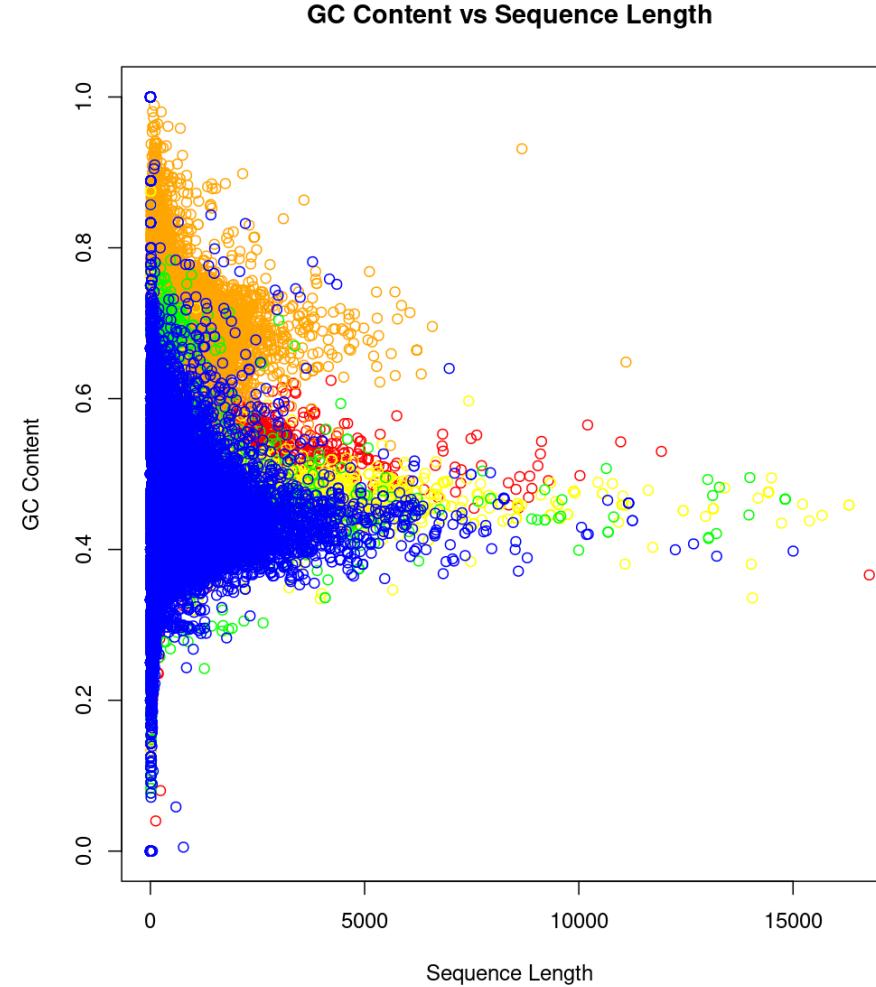
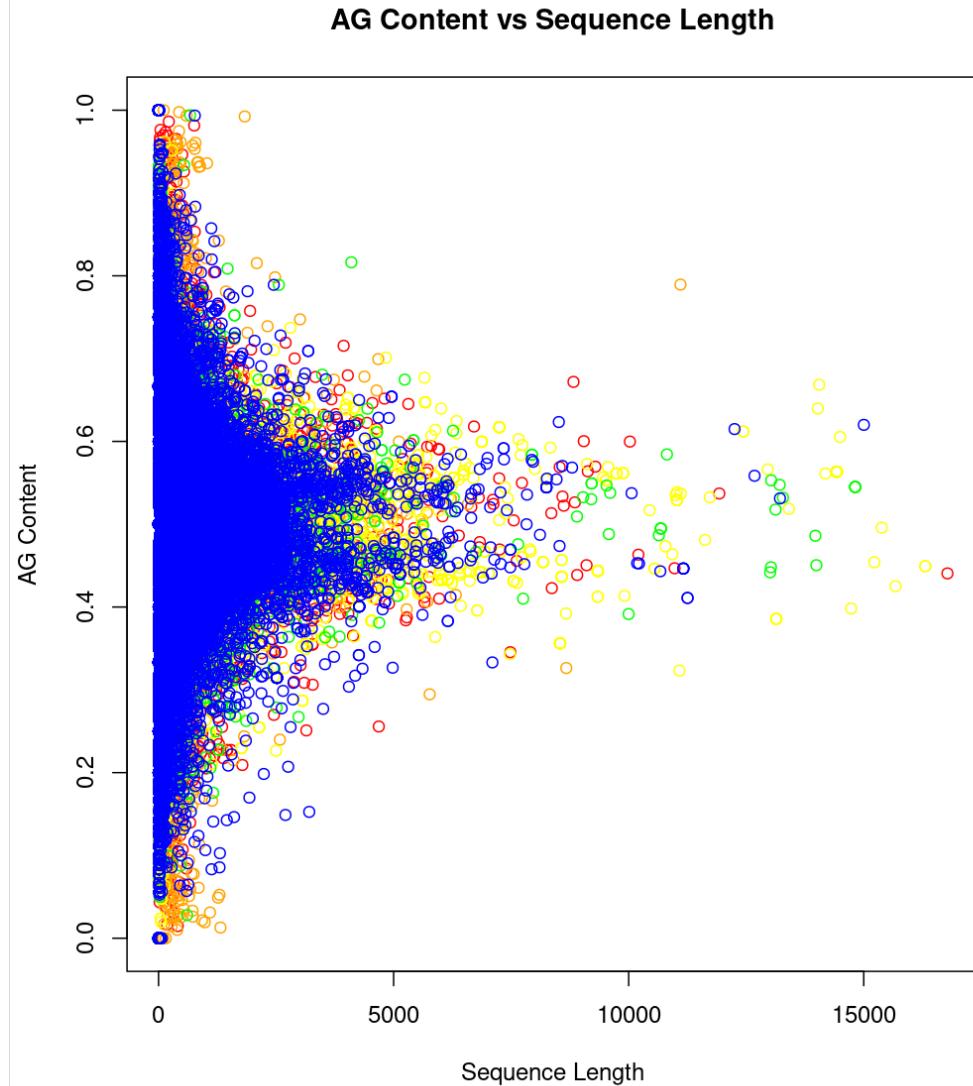
Bfl Pma Llo Cmu Mme

Bfl Pma Llo Cmu Mme

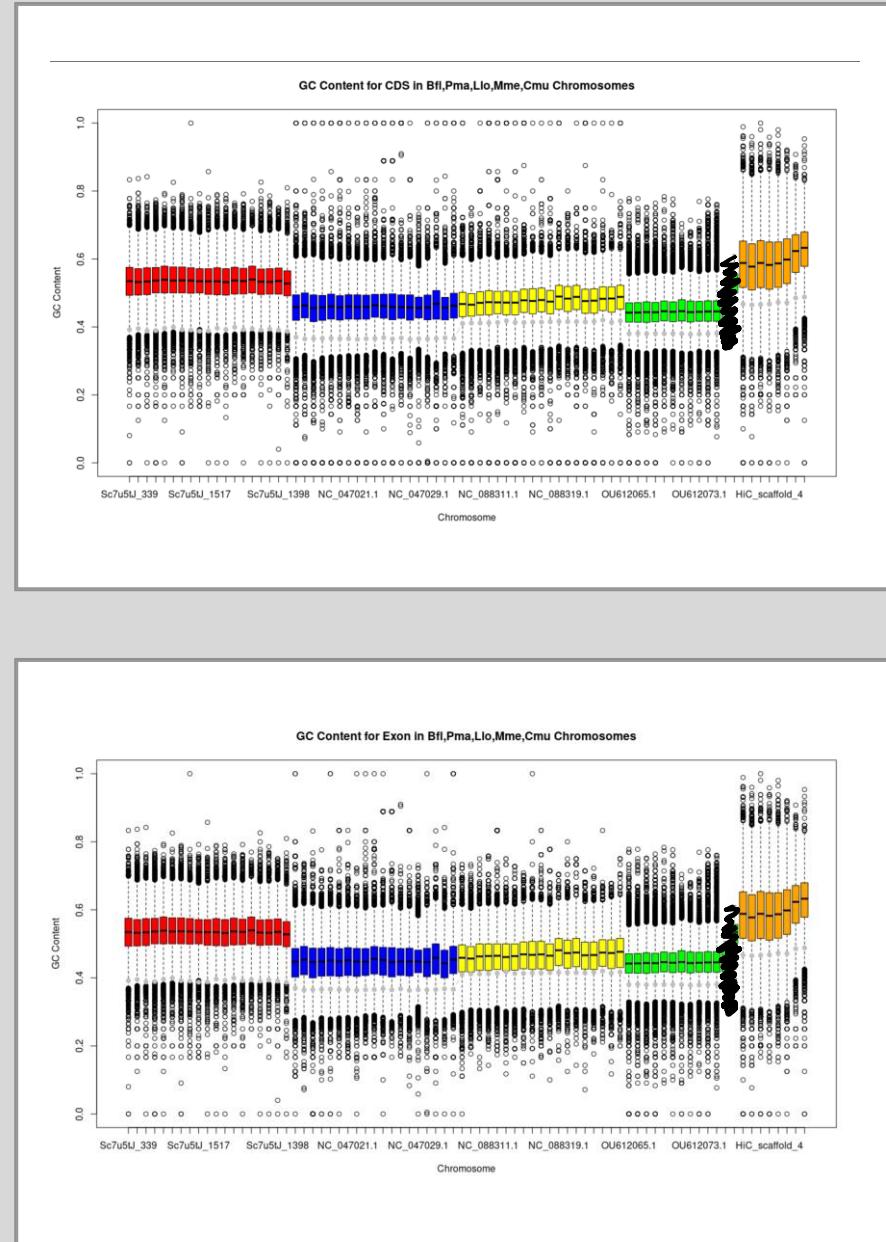
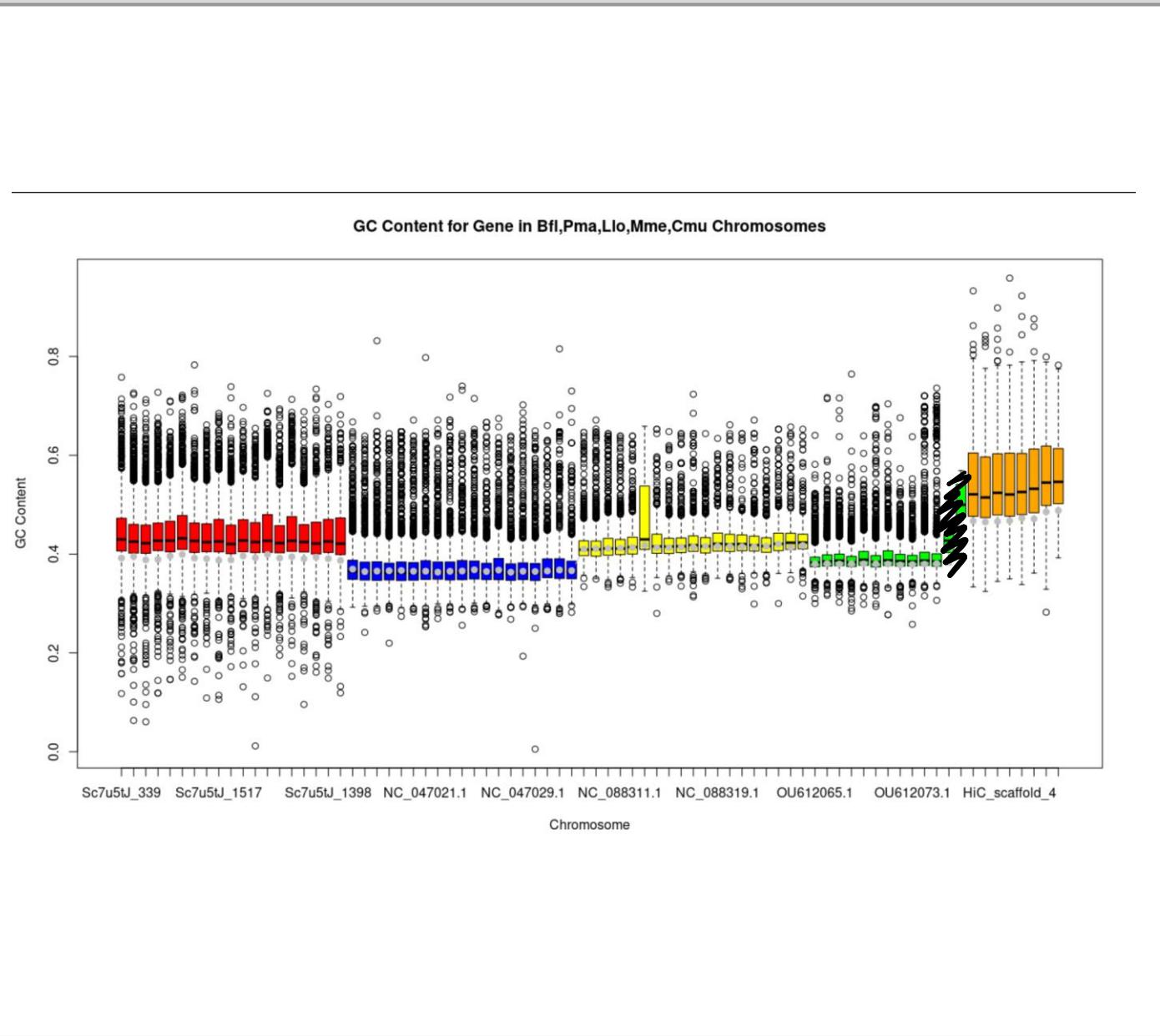


Bfl Pma Llo Cmu Mme

\*Only CDS entries



# Overview of Work:

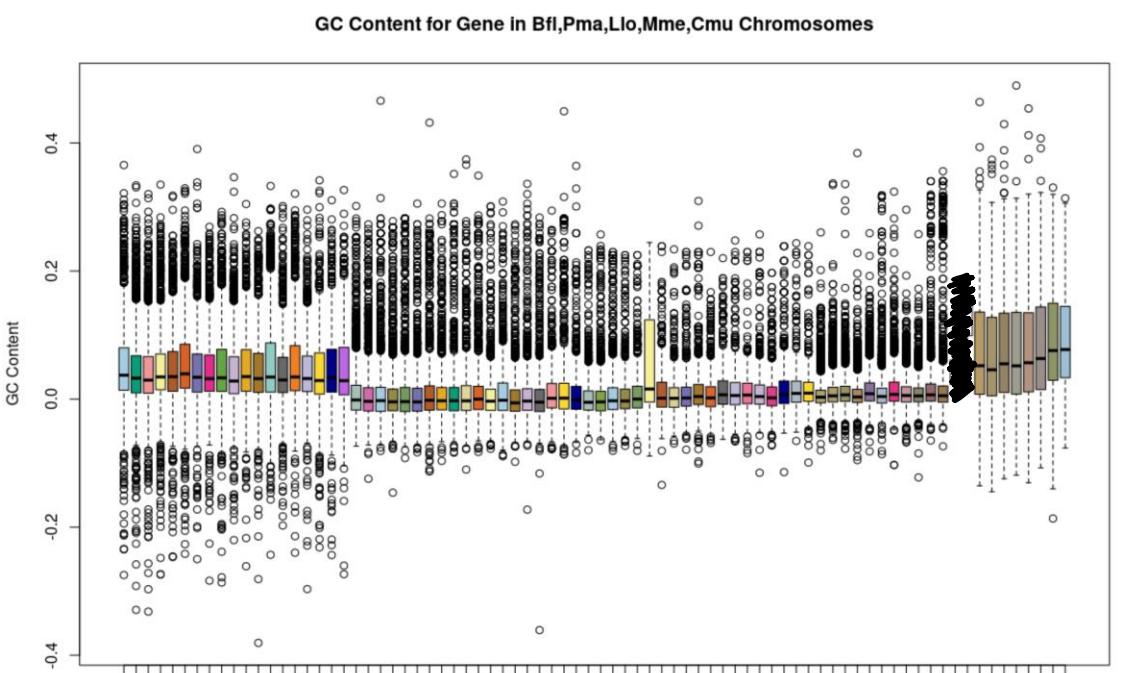
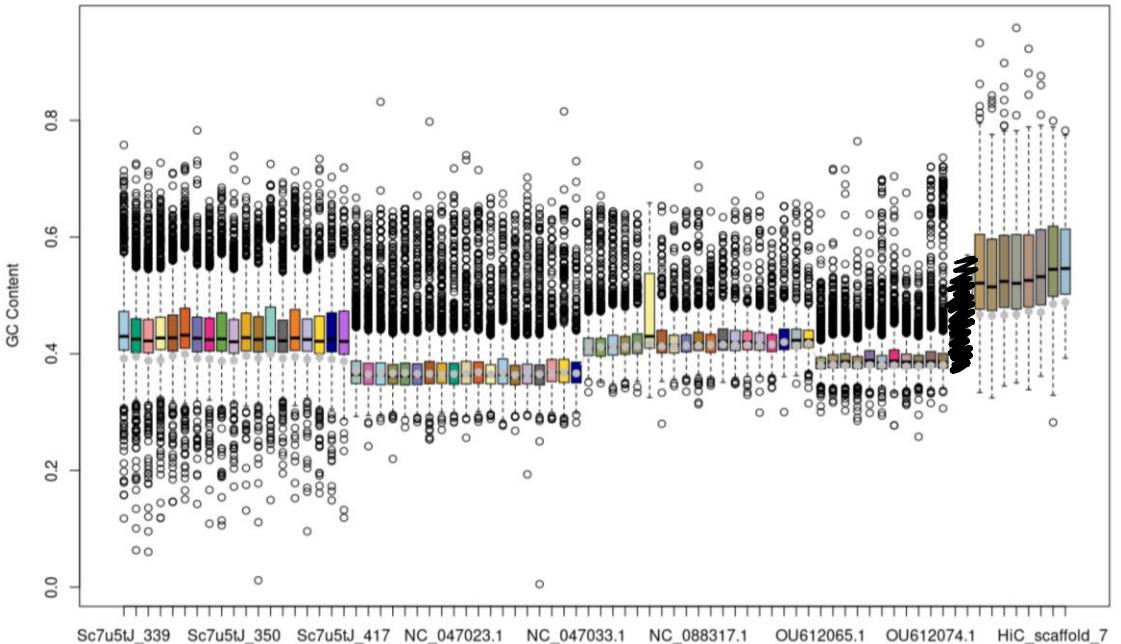


# Alg Tracing

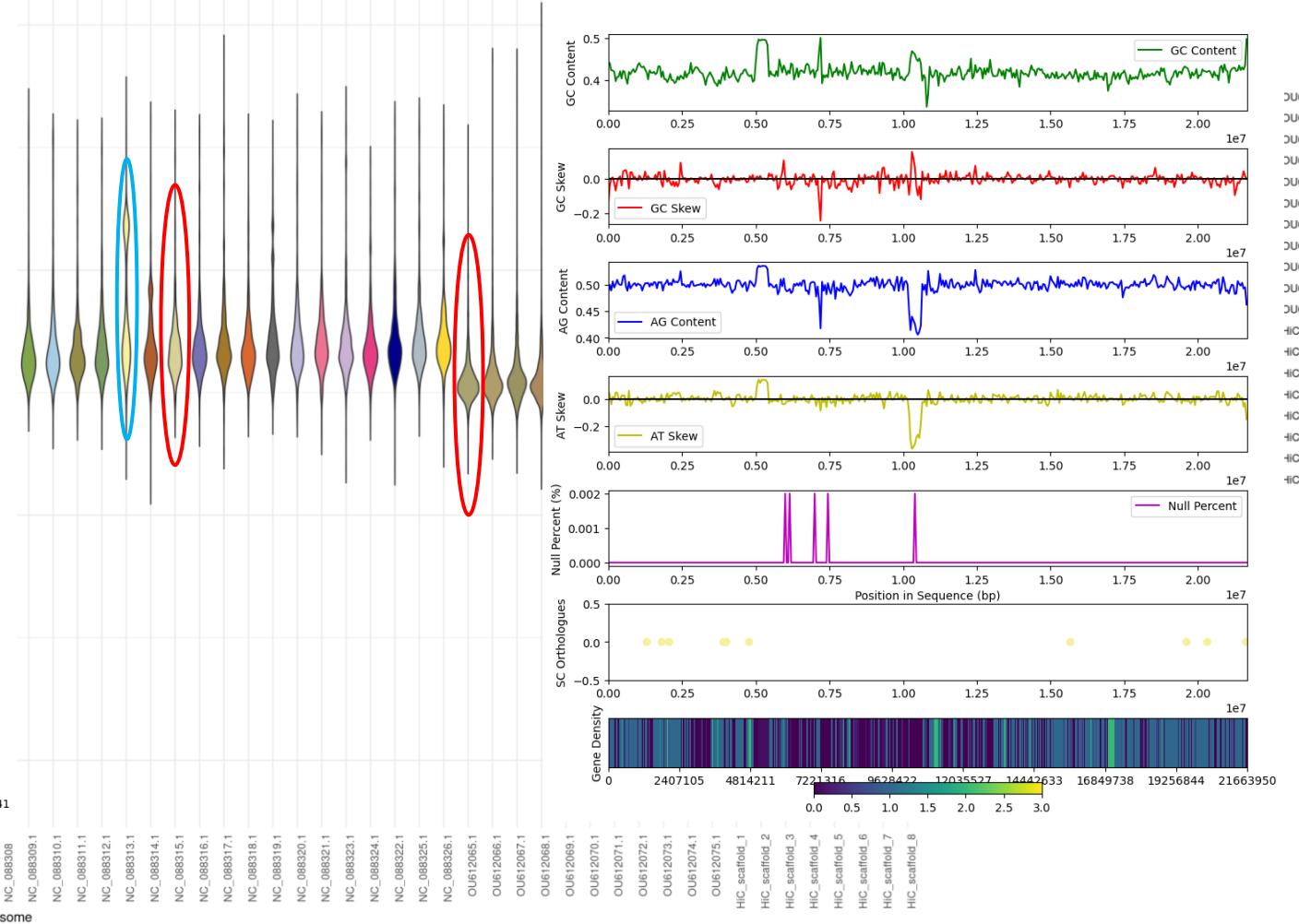
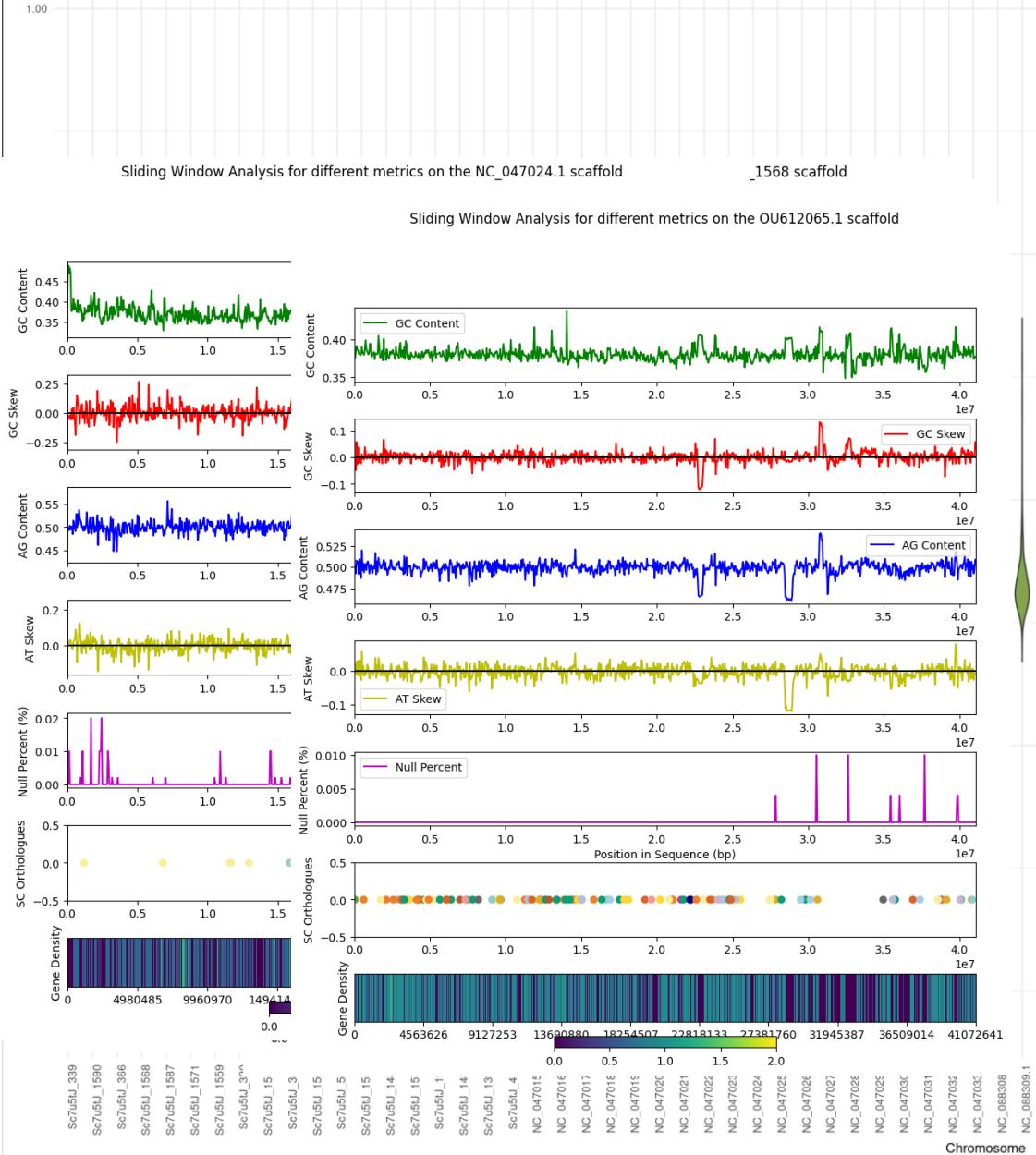
- Averaged the color of all the single copy orthologues in the chromosome

Top: Plotted chromosome average as grey line

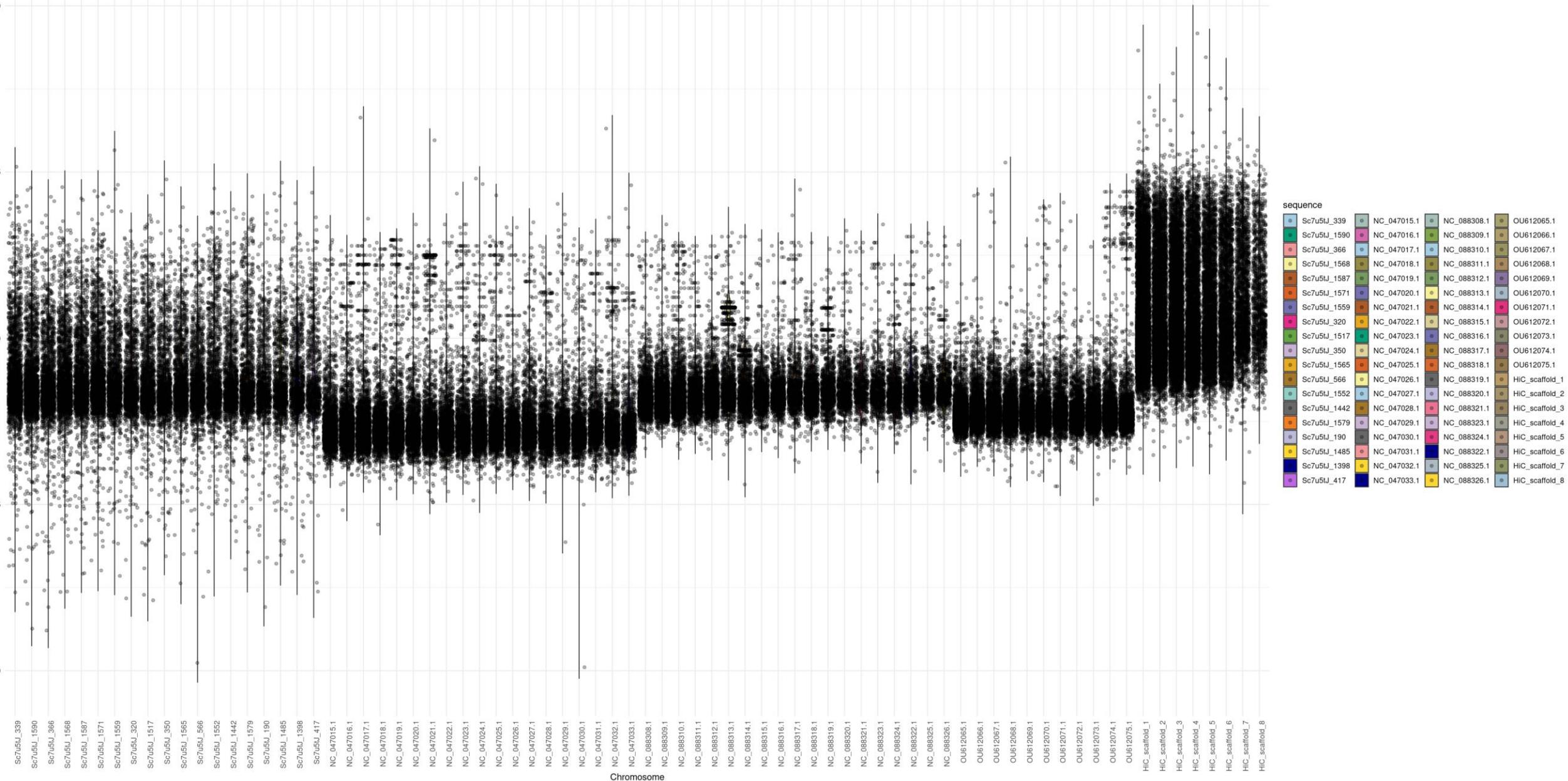
Bottom: Shifted each chromosome by avg GC content of species



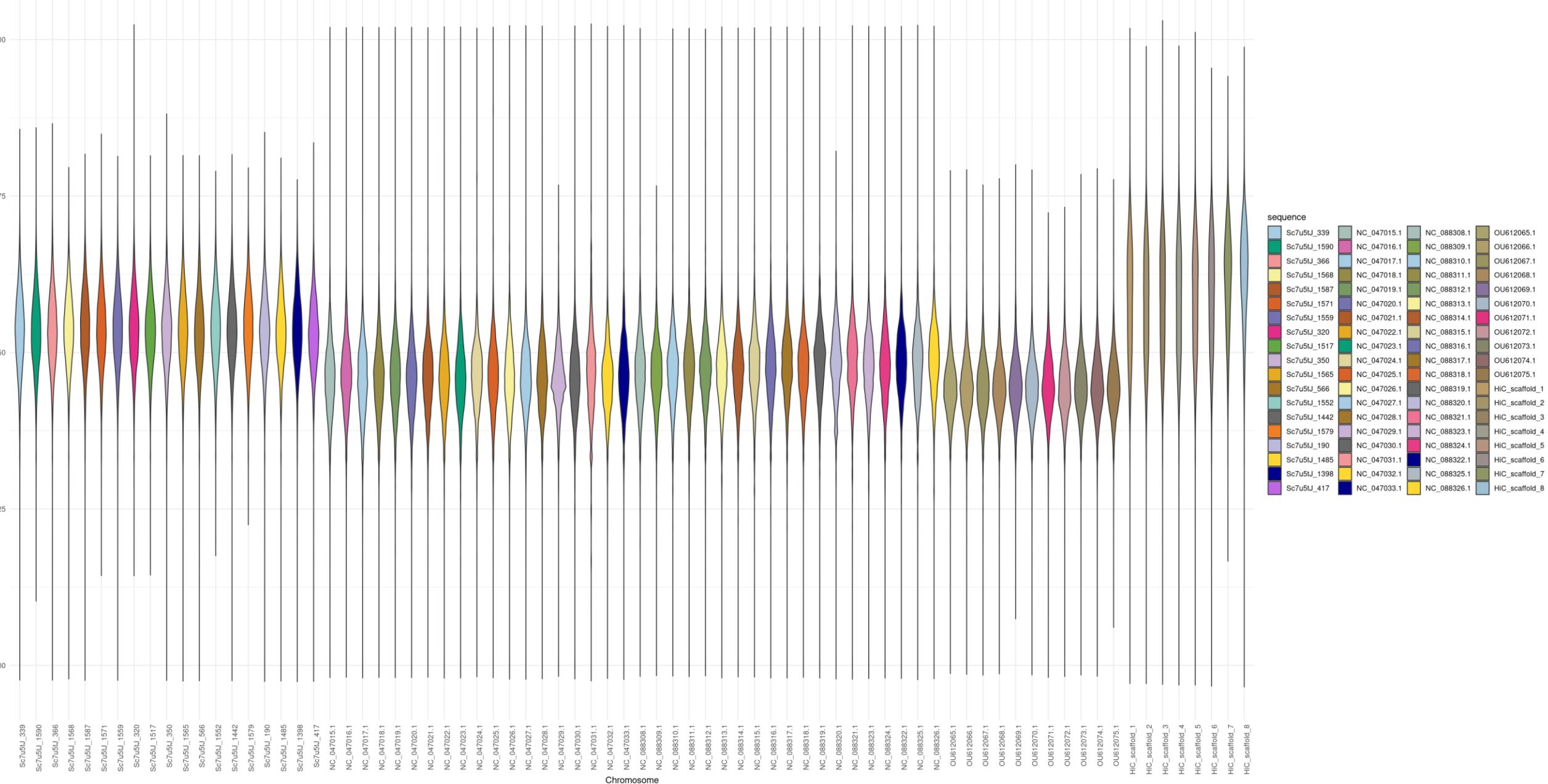
## GC Content for Genes in Bfl, Pma, Llo, Mme, Cmu Chromosomes



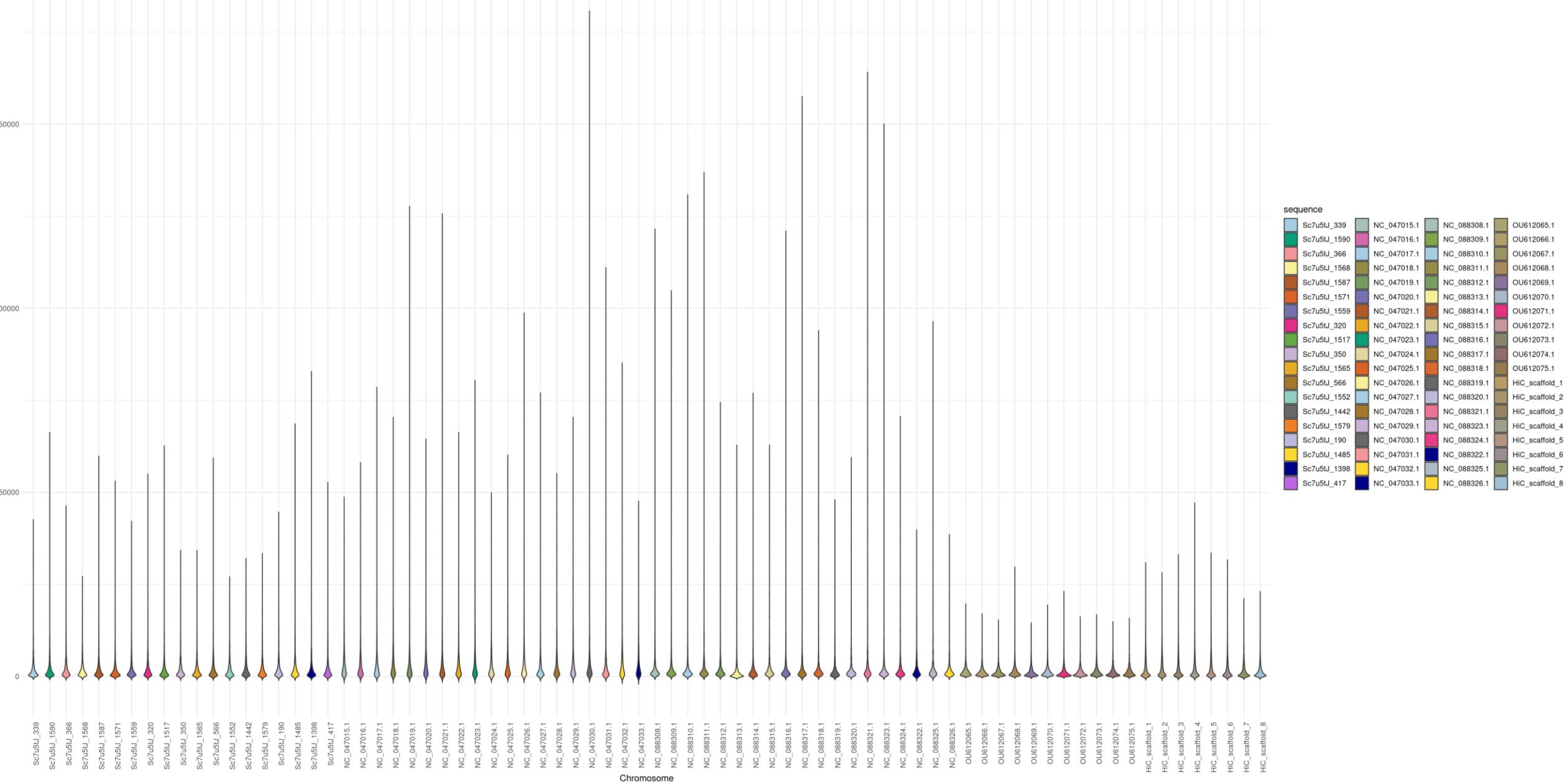
### GC Content for Genes in Bfl, Pma, Llo, Mme, Cmu Chromosomes



GC Content for CDS in Bfl, Pma, Llo, Mme, Cmu Chromosomes



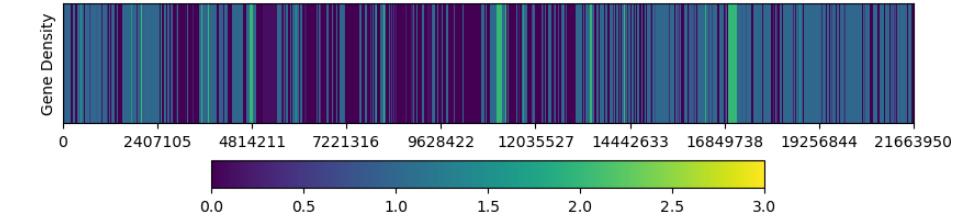
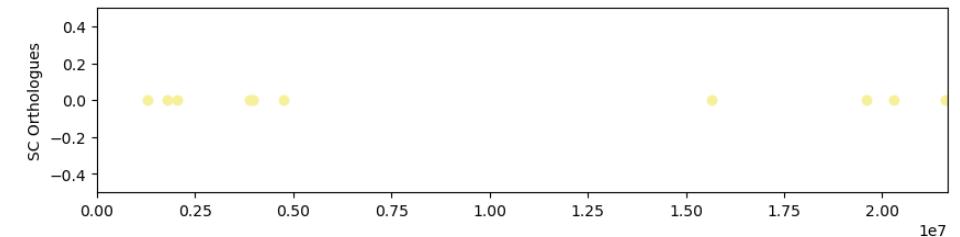
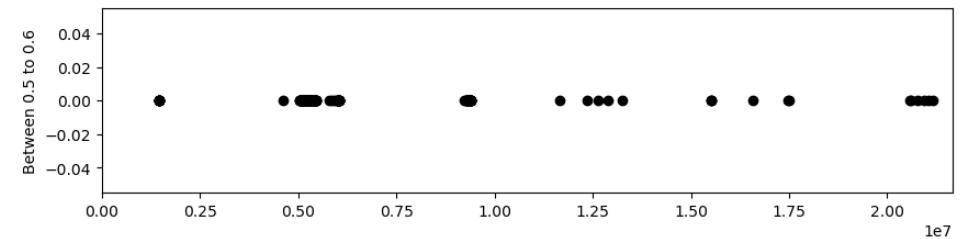
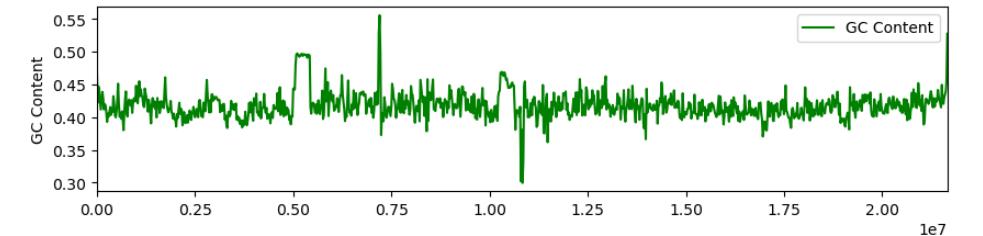
## Sequence Length for Genes in Bfl, Pma, Llo, Mme, Cmu Chromosomes



# Llo Biomodal Distribution

- High GC genes cluster
- Could be due to large centromeres

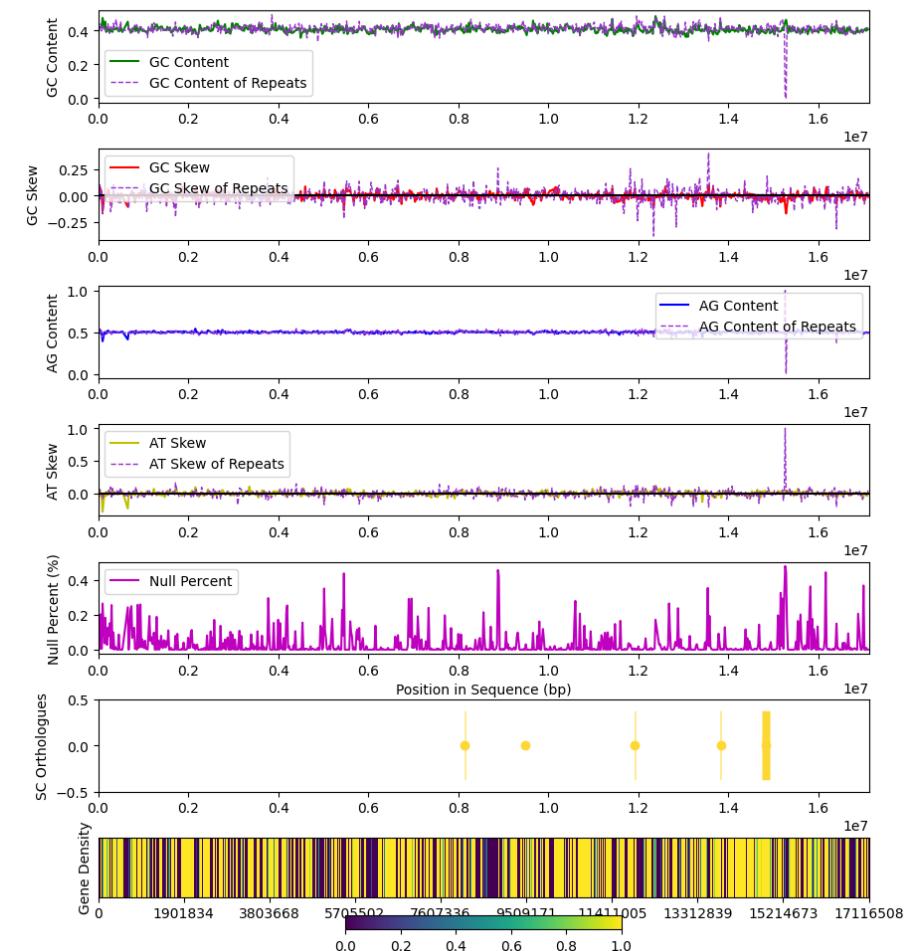
Sliding Window Analysis for different metrics on the NC\_088313.1 scaffold



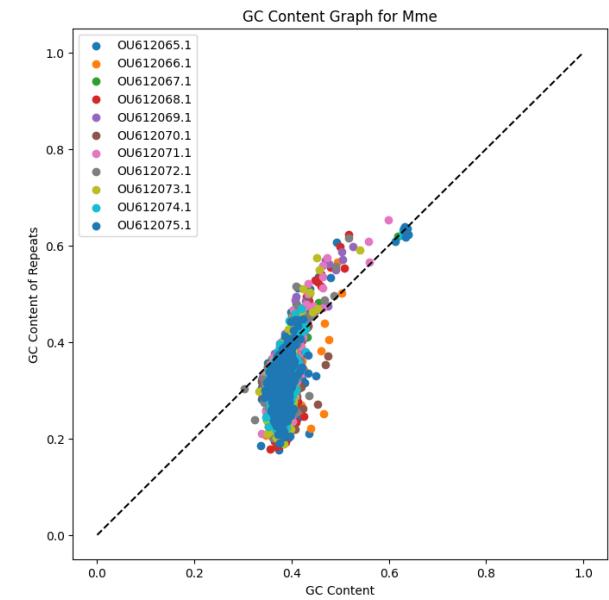
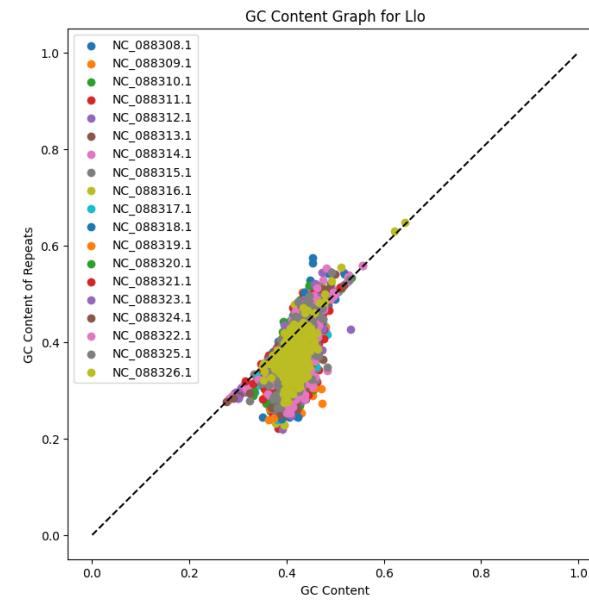
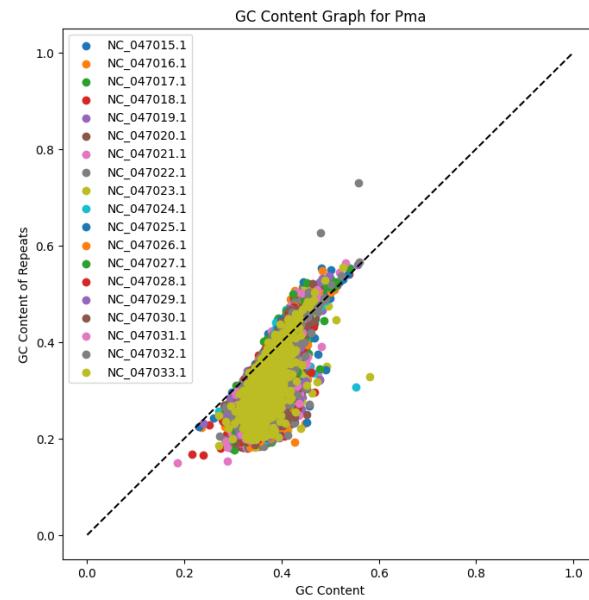
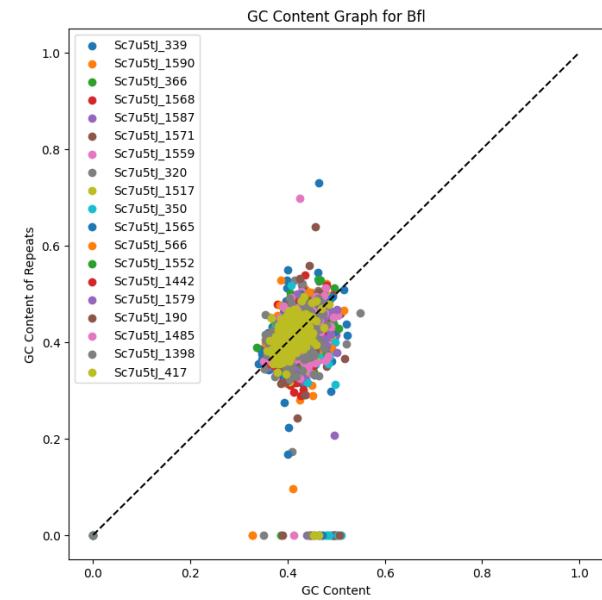
# Sliding Window Comparison of Repeats

- Repeat content is much more extreme

Sliding Window Analysis for different metrics on the Sc7u5tJ\_417 scaffold



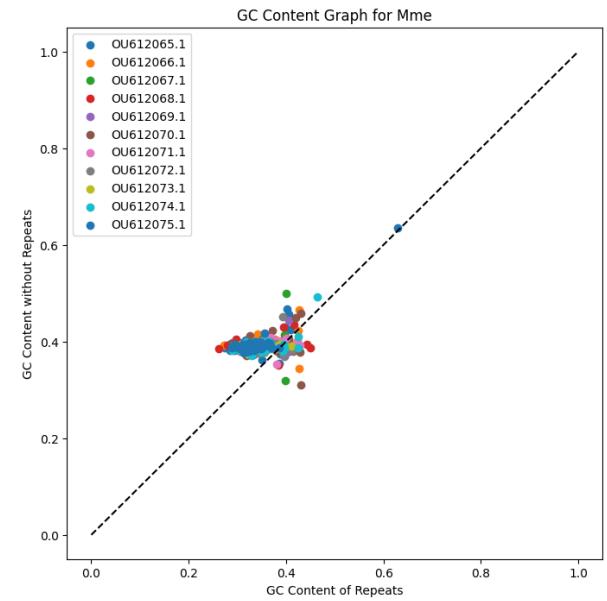
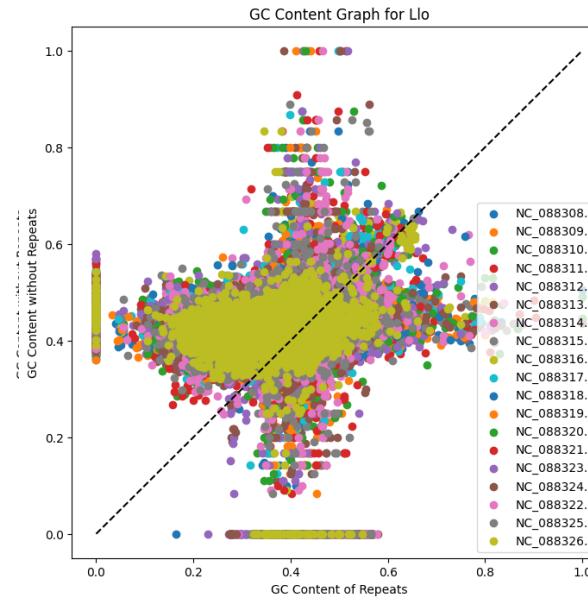
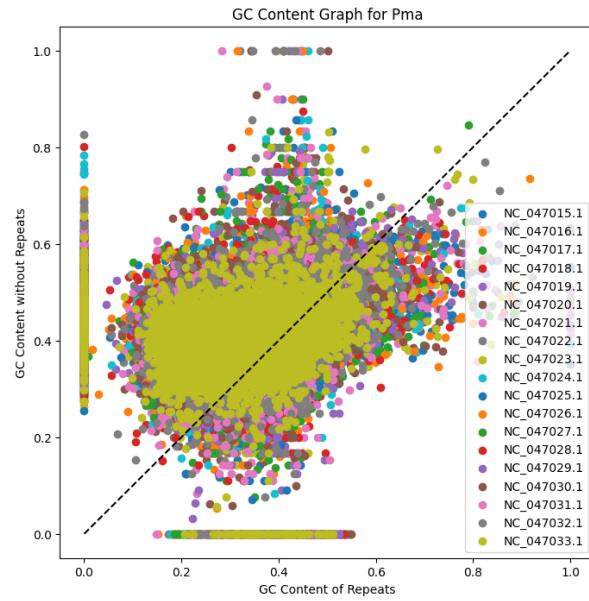
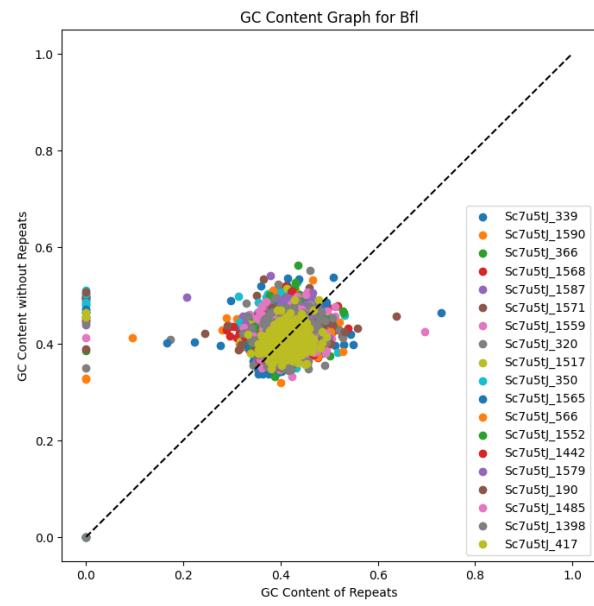
# Effects of Repeats

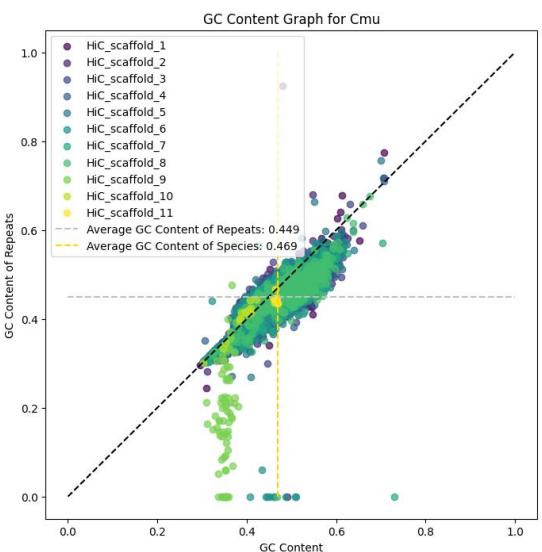
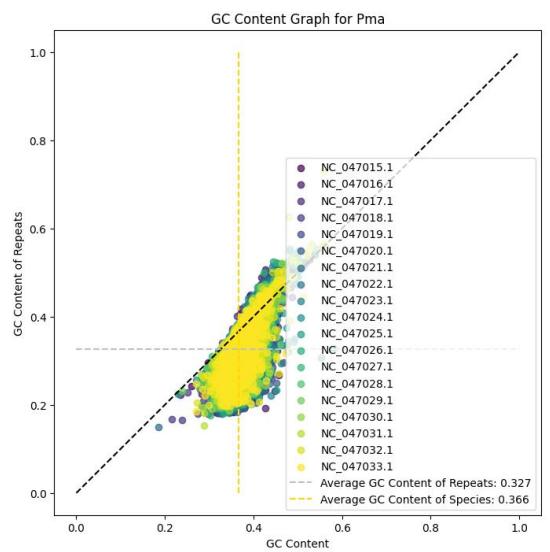
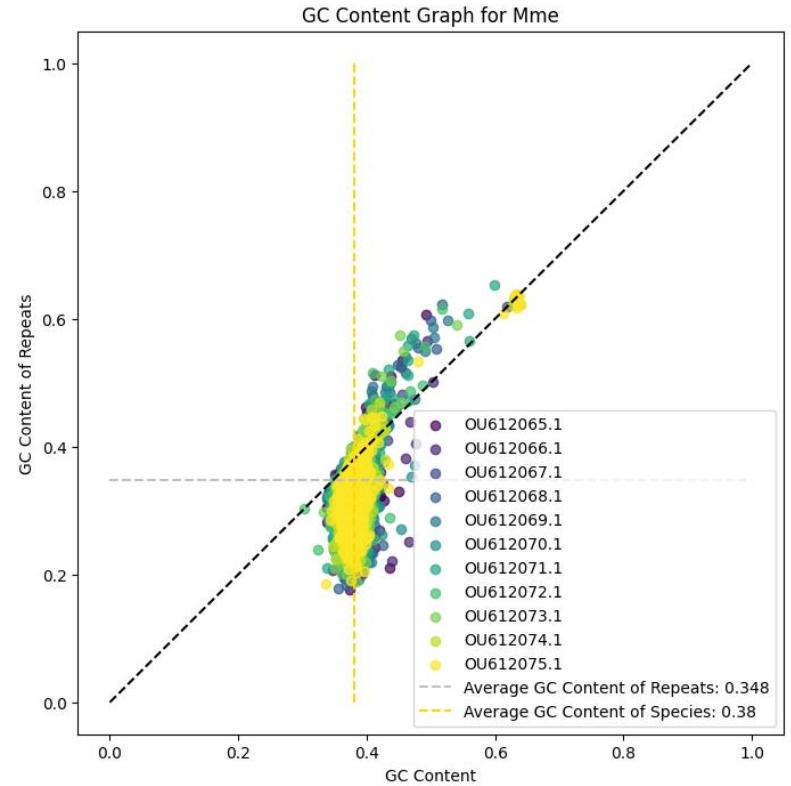
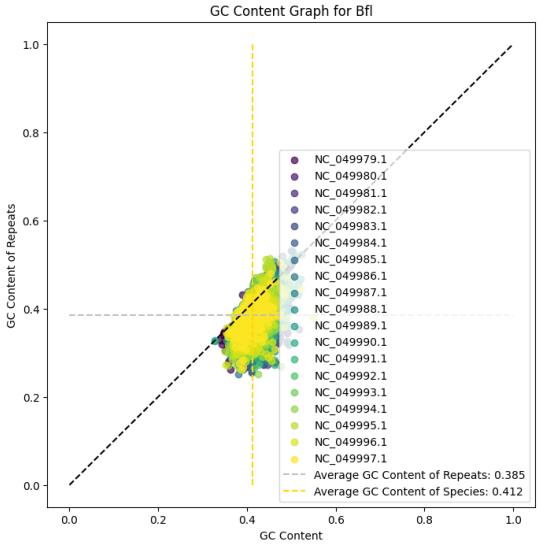
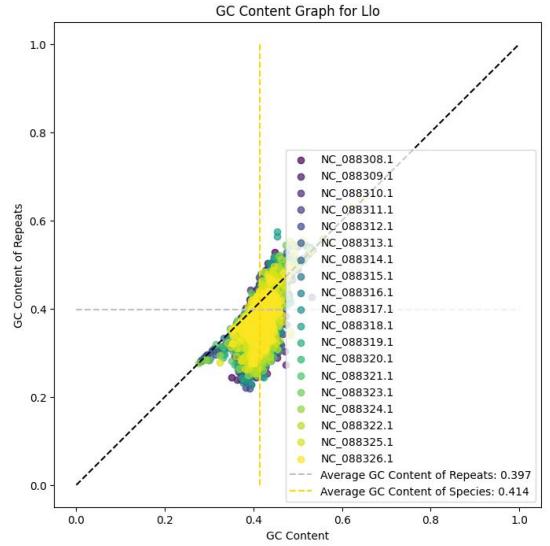


Weak positive correlation between GC content of window and GC content of repeats in window

# NR v. R

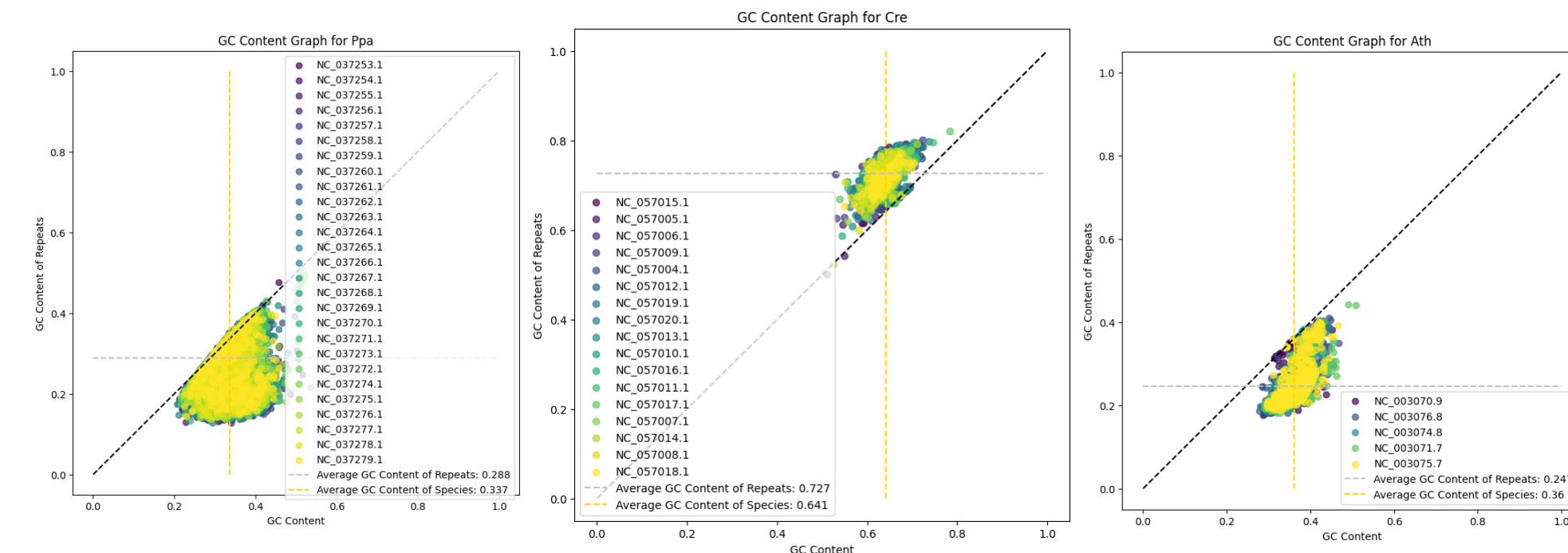
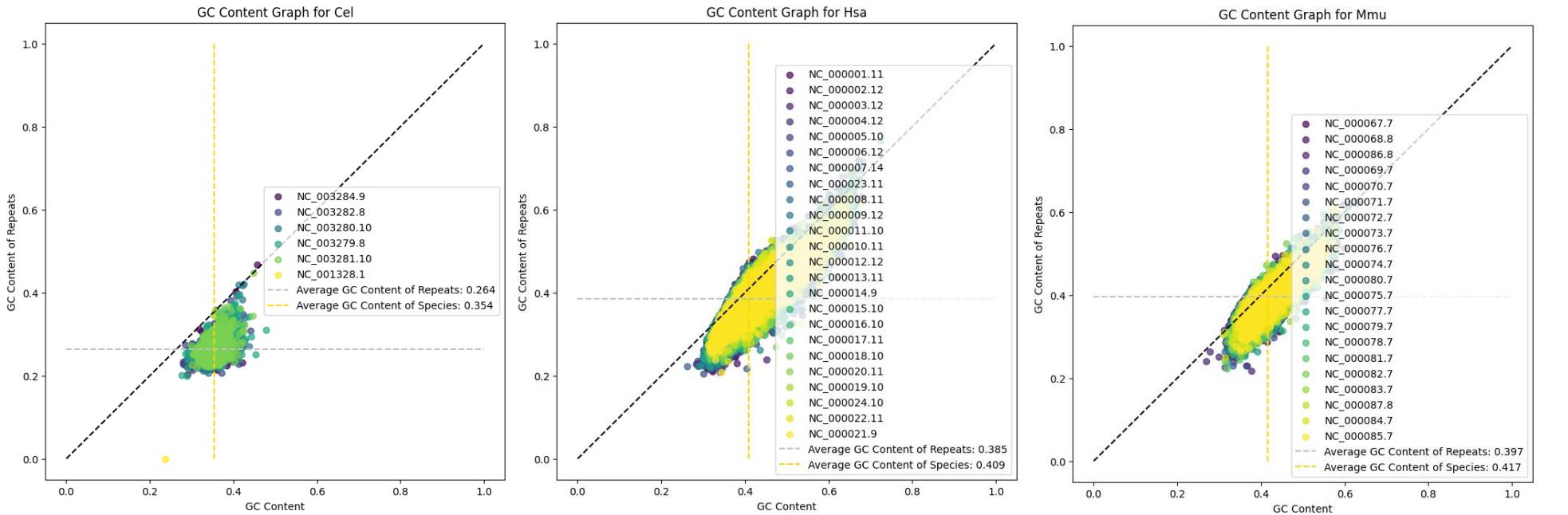
- No correlation
- When window size is too small, crosses appear -> sections of high variance due to small length, while large length has low variance





# Other Species

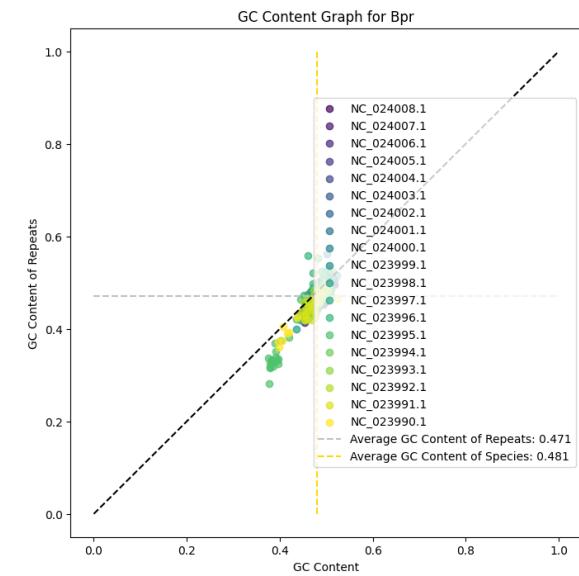
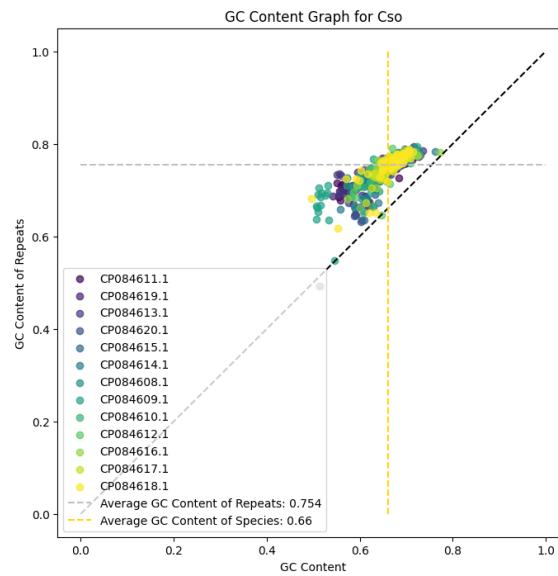
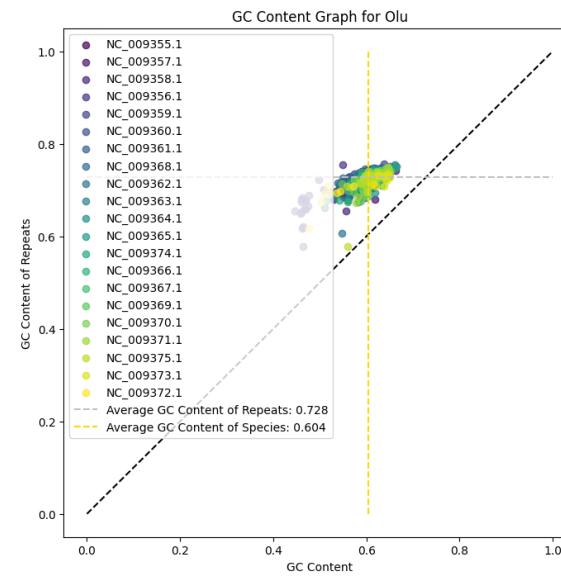
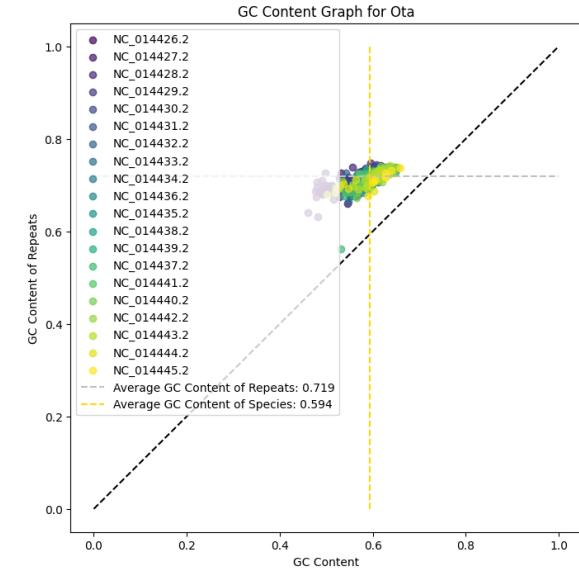
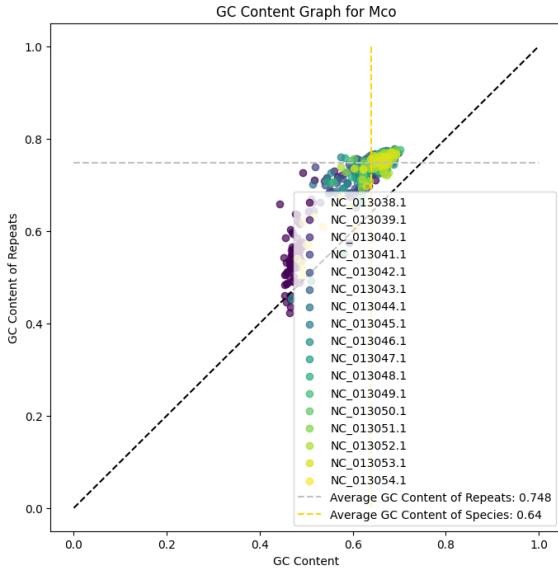
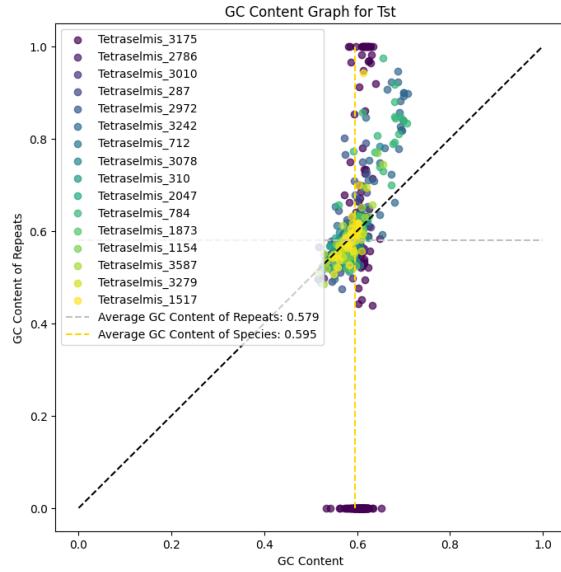
- *C. Elegans*: GCF\_000002985.6
- *H. Sapiens*: GCF\_000001405.40
- *M. Musculus*: GCF\_000001635.27
- *P. Patens*: GCF\_000002425.4
- *C. reinhardtii*: GCF\_000002595.2
- *A. Thaliana*: GCF\_000001735.4



# Added in other species

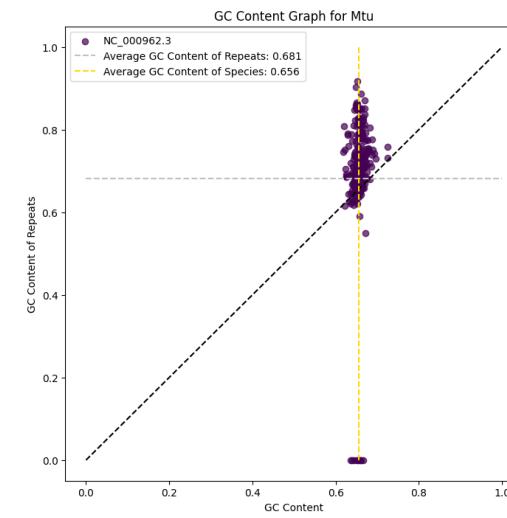
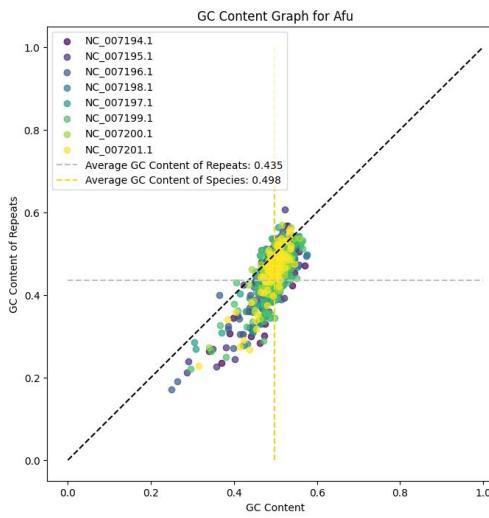
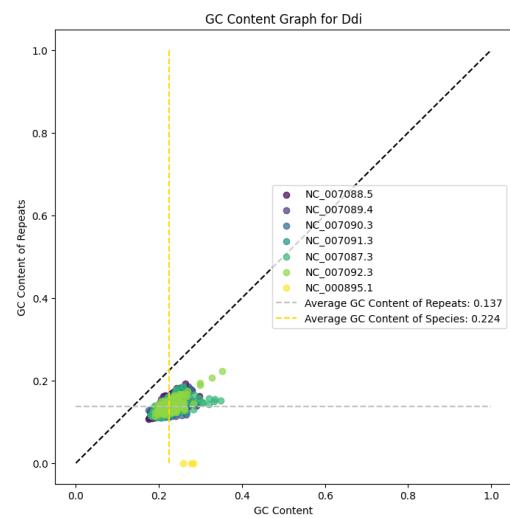
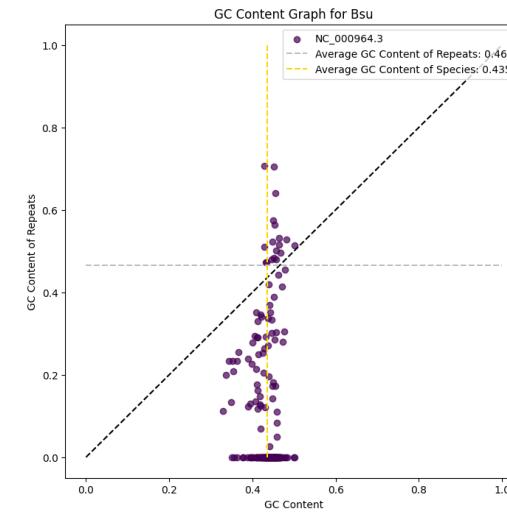
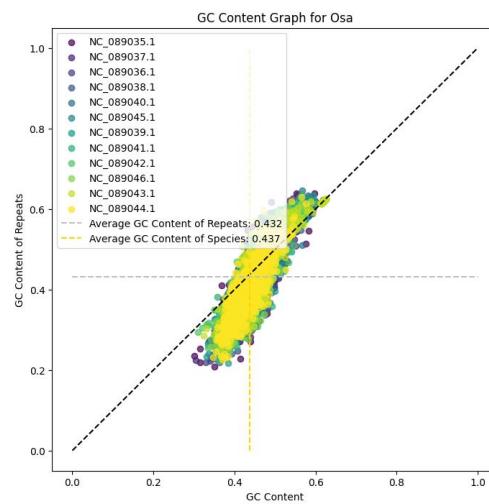
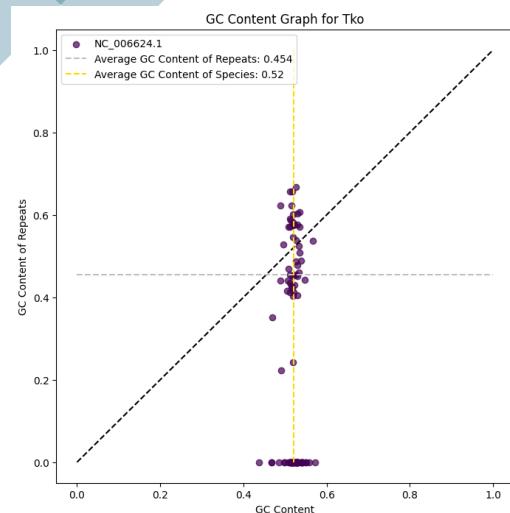
## Chlorophyta

- Micromonas commoda: GCF\_000090985.2
- Ostreococcus lucimarinus: GCF\_000092065.1
- Bathycoccus prasinos: GCF\_002220235.1
- Ostreococcus Tauri: GCF\_000214015.3
- Chlorella sorokiniana: GCA\_025917655.1
- Tetraselmis Striata: From lab



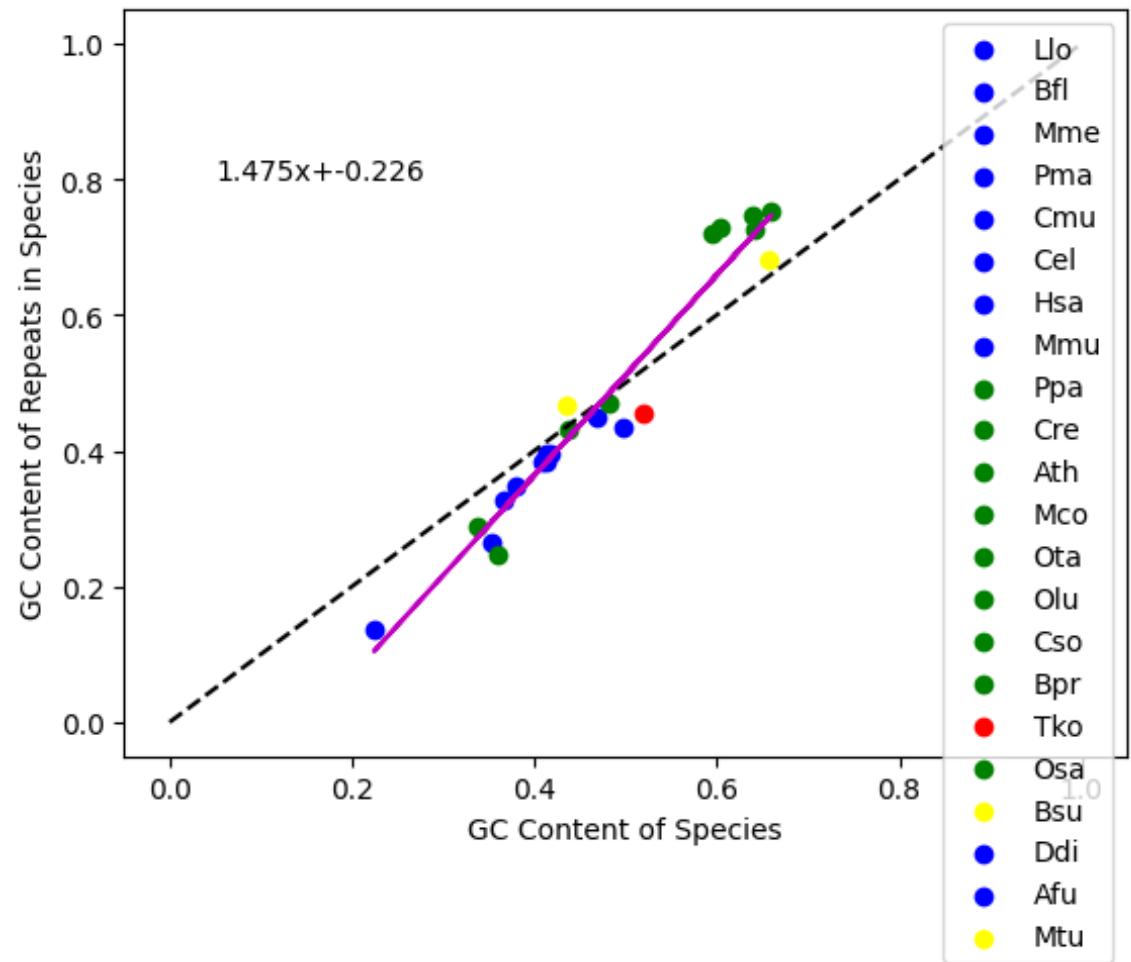
# More:

- Archaea: *Thermococcus kodakarensis* KOD1:  
[https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_000009965.1/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_000009965.1/)
- Amoebozoa: *Dictyostelium discoideum* AX4:  
[https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_000004695.1/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_000004695.1/)
- *Mycobacterium tuberculosis* H37Rv:  
[https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_000195955.2/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_000195955.2/)
- Streptomyces Coelicolor A3(2): [https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_008931305.1/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_008931305.1/)
- *Aspergillus fumigatus*: [https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_000002655.1/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_000002655.1/)
- *Oryza Sativa Japonica* Group: [https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_034140825.1/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_034140825.1/)
- *Bacillus subtilis* subsp. *subtilis* str. 168:  
[https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_000009045.1/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_000009045.1/)
- *Haloferax volcanii* DS2: [https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_000025685.1/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_000025685.1/)
- *Pyrococcus furiosus* DSM 3638:  
[https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_008245085.1/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_008245085.1/)



# Repeat Segments are more extreme?

- GC Content in repeats seem to increase more rapidly than the overall GC Content for all species
- Could be due to faster replication rate
  - if mutation bias exists
    - Lack of correlation in data from sliding window analysis contradicts this
  - if selective pressure
    - If so, why?
  - Would not be effected by biased gene conversion, as recombination
- Could be due to weaker constraints
  - Repeats aren't subject to as many selective constraints
- Error: [0.07691182 0.03715188]

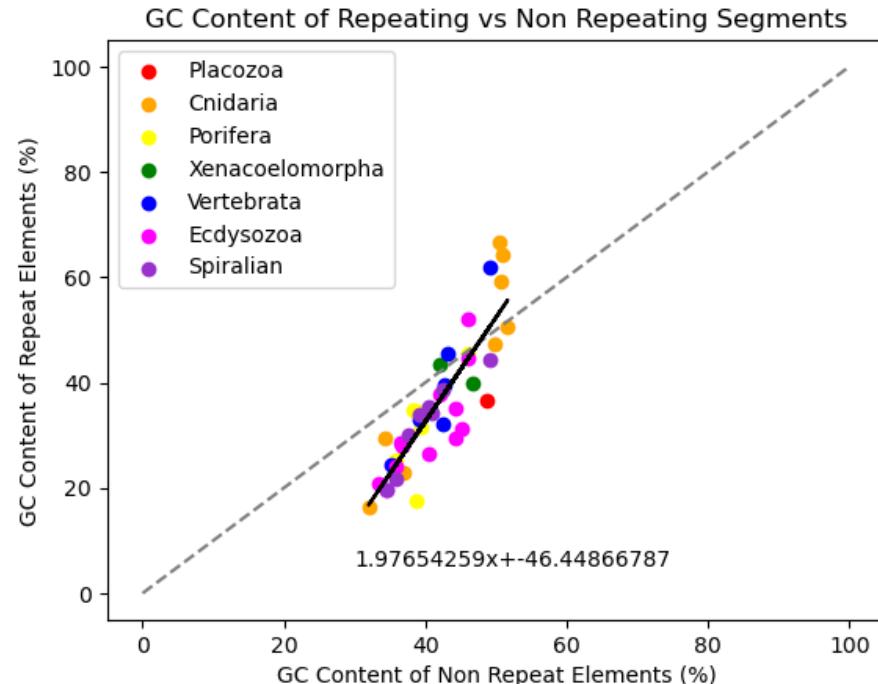


# Viruses\*

- Sars CoV 2:  
[https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_009858895.2/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_009858895.2/)
- HIV:  
[https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_000864765.1/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_000864765.1/)
- Japanese Encephalitis:  
[https://www.ncbi.nlm.nih.gov/datasets/genome/GCF\\_000862145.1/](https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_000862145.1/)

\*Can't mask, sequences too short

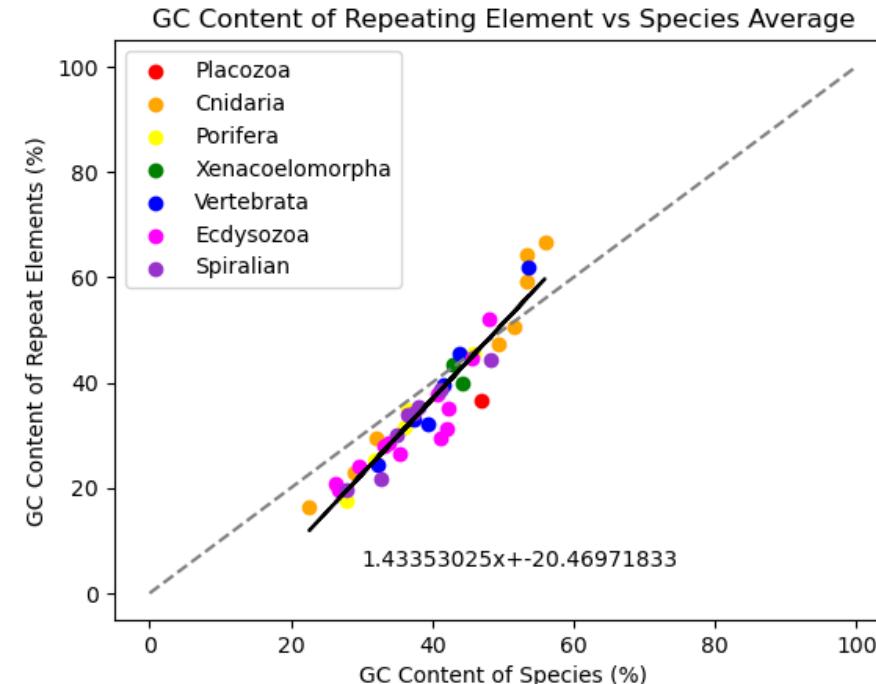
# Focusing on Eukaryotes



Error for R vs NR: [0.17515199 7.34616462]

$$\text{If } x = y, \frac{46.44866787}{1.97654359-1} = 47.56440563$$

$$\text{Error: } 47.56440563 \sqrt{\left(\frac{0.17515199}{1.97654259-1}\right)^2 + \left(\frac{7.34616462}{46.44866787}\right)^2} = 11.37408807$$

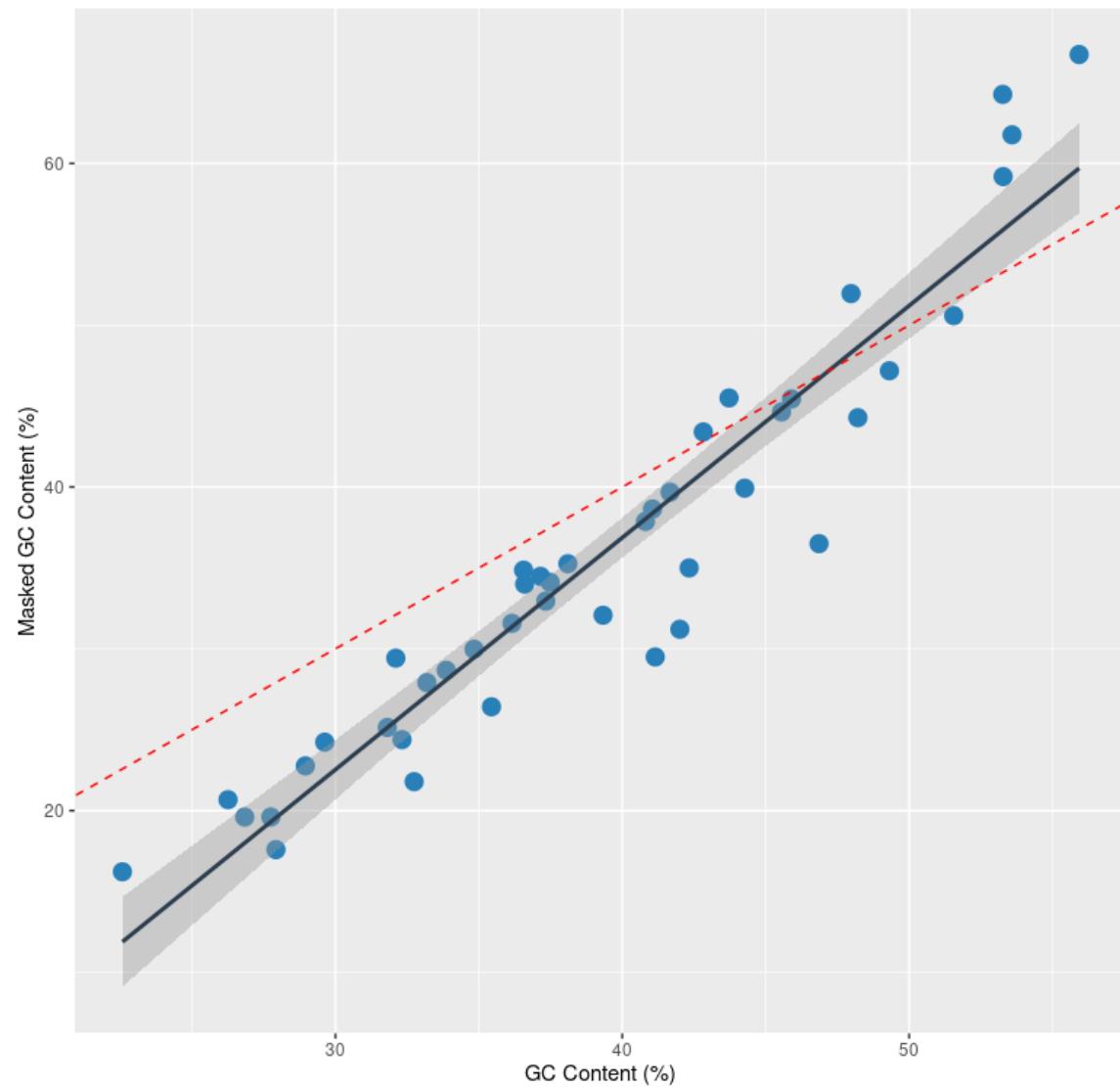


Error for R vs Species: [0.07408874 2.96851225]

$$\text{If } x = y, \frac{20.46971833}{1.43353025-1} = 47.21635533$$

$$\text{Error: } 47.21635533 \sqrt{\left(\frac{0.07408874}{1.43353025-1}\right)^2 + \left(\frac{2.96851225}{20.46971833}\right)^2} = 10.58281436$$

Masked GC Content vs. GC Content

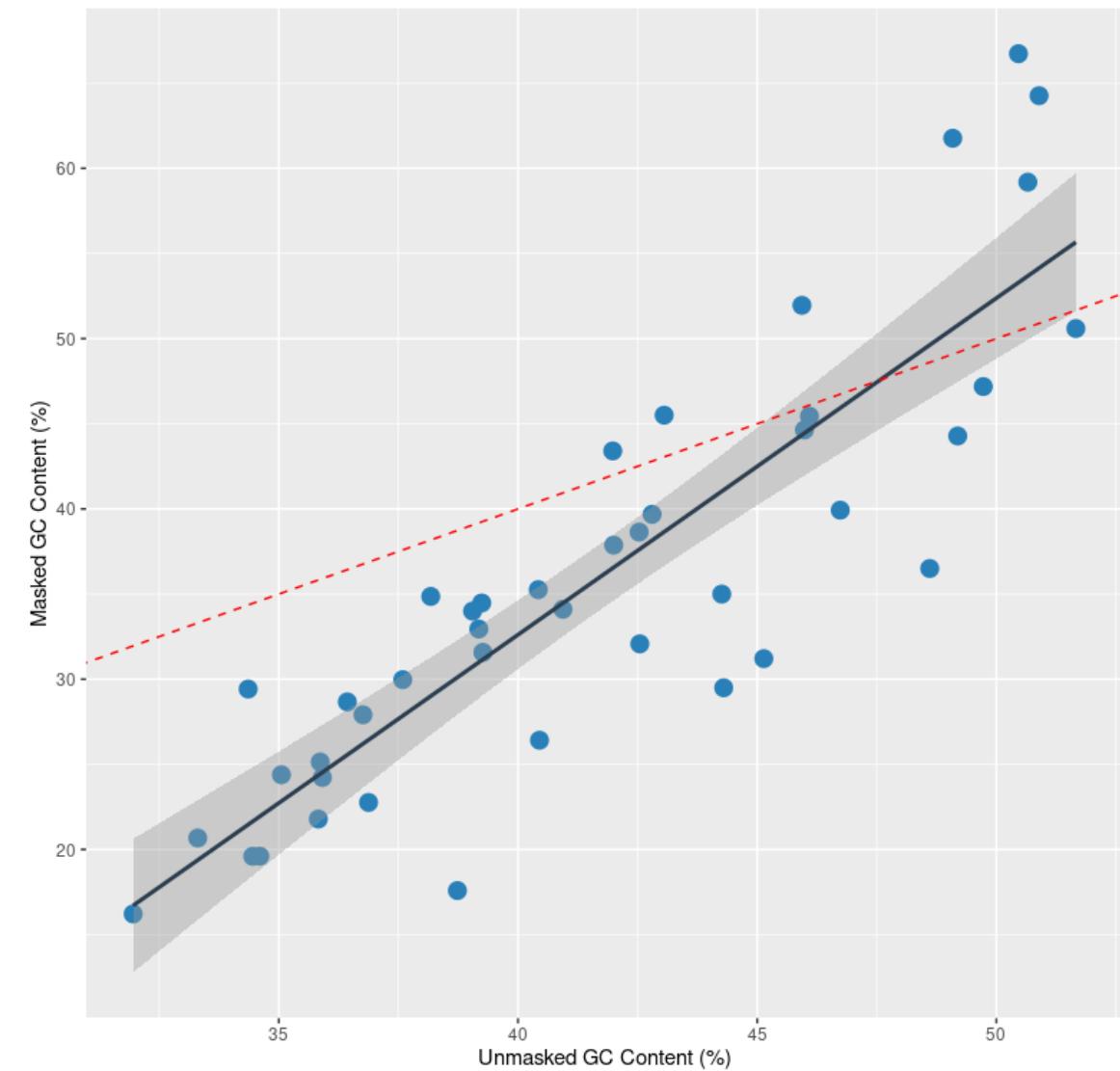


```

Call: lm(formula = masked_gc_content ~ gc_content, data = df)
Residuals: Min 1Q Median 3Q Max
-10.2002 -2.4081 0.4221 2.3565 8.3815
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) -20.46972 2.96851 -6.896 2.31e-08 ***
gc_content 1.43353 0.07409 19.349 < 2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 3.997 on 41 degrees of freedom
Multiple R-squared: 0.9013,
Adjusted R-squared: 0.8989
F-statistic: 374.4 on 1 and 41 DF, p-value: < 2.2e-16

```

Masked GC Content vs. Unmasked GC Content

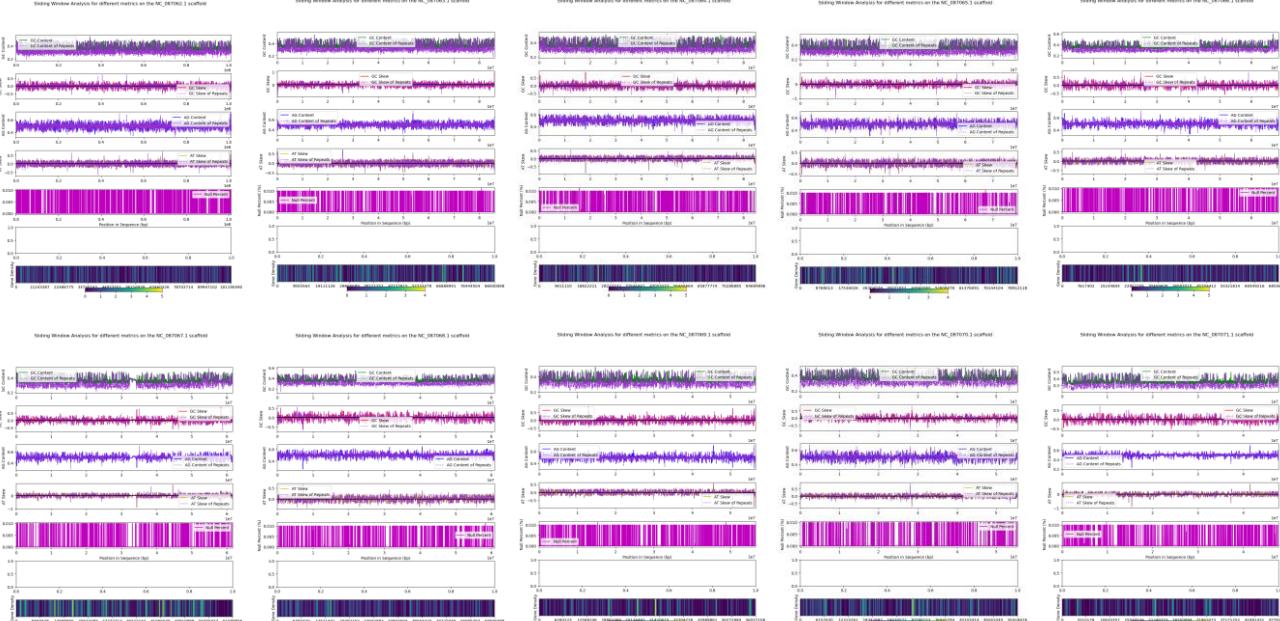


```

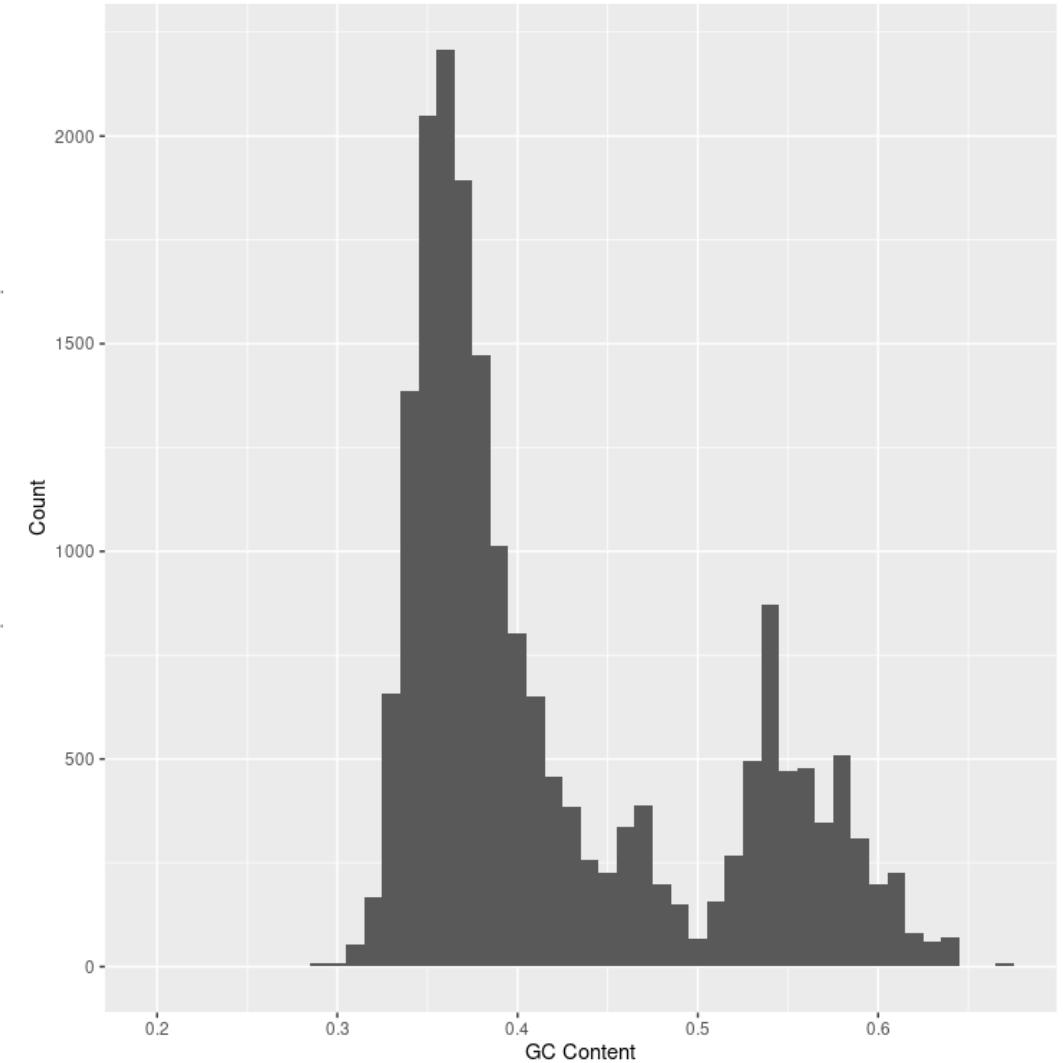
Call: lm(formula = masked_gc_content ~ unmasked_gc_content, data = df)
Residuals: Min 1Q Median 3Q Max
-13.1258 -4.1637 0.7848 3.1862 13.4383
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) -46.4487 7.3462 -6.323 1.50e-07
unmasked_gc_content 1.9765 0.1752 11.285 3.76e-14 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 6.278 on 41 degrees of freedom
Multiple R-squared: 0.7565,
Adjusted R-squared: 0.7505
F-statistic: 127.3 on 1 and 41 DF, p-value: 3.761e-14

```

# Histograms for detecting HGT?



## Histogram of GC Content in Genes *Symsagittifera roscoffensis*



## Initial Observations

- Compositional Domains exist in *S. roscoffensis*
  - Small subset of high GC genes seem to exist
  - Don't know if horizontal gene transfer would conserve GC content to an extent where you can use it as a test