



AARHUS UNIVERSITY

Faculty of Natural Sciences

Department of Molecular Biology and Genetics

**Martyna Stankevičiūtė**

# **The relationship between evolutionary constraint and genetic diversity across primate genomes**

**Master's Thesis**

Master's degree programme in Bioinformatics

*Supervisor:* **Juraj Bergman**

Tenure-track Assistant Professor,

Department of Molecular Biology and Genetics – BiRC,

Aarhus University

Section for Bioinformatics and Computational Biology – BiRC



*Aarhus 2026*

# Contents

<b>Acknowledgments</b> .....	<b>3</b>
<b>Abstract</b> .....	<b>4</b>
<b>1 Introduction</b> .....	<b>5</b>
1.1 Genetic Diversity .....	6
1.2 Evolutionary Constraint and Purifying selection .....	8
1.3 Measuring Constraint: Conservation Scores .....	9
1.4 Relationship of Constraint and Diversity .....	11
1.5 Aim .....	12
1.6 Objectives .....	13
<b>2 Materials and Methods</b> .....	<b>13</b>
2.1 Data .....	14
2.1.1 Human Reference genome .....	14
2.1.2 Annotation data of <i>hs1</i> genome .....	14
2.1.3 Multiple Genome Alignment .....	15
2.1.4 Phylogenetic tree .....	15
2.1.5 Diversity data .....	15
2.2 Methods .....	16
2.2.1 Parametric conservation scoring using PHAST .....	16
2.2.1.1 Neutral substitution model estimation .....	17
2.2.1.2 PhastCons: neighborhood-aware conservation scoring .....	17
2.2.1.3 PhyloP: per-site conservation scoring .....	18
2.2.2 Non-parametric Conservation scoring (GPN-Star) .....	19
2.2.2.1 Multiple sequence alignment processing .....	21
2.2.2.2 Training dataset construction .....	21
2.2.2.3 Model configuration and hyper-parameter selection .....	22
2.3 Data integration and processing .....	23
2.3.1 Quality control of PHAST conservation scores .....	23
2.3.2 Integration with heterozygosity data .....	23
2.4 Statistical analysis and linear modelling .....	24

---

2.4.1 Correlation analysis .....	24
2.4.2 Categories of conservation scores .....	24
2.4.3 Transformation of heterozygosity .....	25
2.4.4 Linear modelling .....	25
2.5 Result Visualisation .....	27
<b>3 Results .....</b>	<b>27</b>
3.1 Evolutionary constraint across the primate genome .....	28
3.2 Heterozygosity across primate species .....	34
3.3 Heterozygosity and evolutionary conservation relationship .....	36
3.4 Regression models .....	40
3.5 GPN-star model training .....	44
<b>4 Discussion .....</b>	<b>46</b>
4.1 Limitations and future directions .....	49
<b>5 Conclusions .....</b>	<b>51</b>
<b>6 Declaration of use of GAI tools .....</b>	<b>53</b>
<b>Bibliography .....</b>	<b>54</b>
<b>Appendix .....</b>	<b>63</b>
Supplementary Figures .....	64
Supplementary Tables .....	66

## Acknowledgments

I want to express my gratitude to all people who help me along the during this study. Thank you to all the amazing people from BiRC at Aarhus Univeristy, the professors and fellow students. A huge thank you for people who maintain the local HPC cluster “GenomeDK” which i used for all the computational tasks. I want to thank my closest 3 friends who i met during my studies: Antonia, Laurids and Sebastian - for always reminding to take a break and boosting the mood. I am grateful for my family and Mykolas for always waiting for me to comeback home. Finally, my biggest thank you goes to the people behind the Primate Diversity Panel: for greatest help in the most important moments of the thesis to Bjarke Pedersen and for the best supervision and constant support to my supervisor Juraj Bergman.

## Abstract

Genomic regions under strong purifying selection are expected to harbor less standing genetic variation. Testing this relationship systematically across the primate order requires both a high-resolution map of genomic constraint and diversity estimates spanning a phylogenetically broad sample of species.

This study combines a 49-species whole-genome alignment with heterozygosity estimates from 219 non-human primate species to ask whether evolutionary constraint, quantified at 100 kb resolution, predicts patterns of genetic diversity across the primate genome. PhastCons and PhyloP conservation scores were derived from the whole-genome alignment using the PHAST package and combined with per-window heterozygosity estimates, calculated as heterozygous sites per callable site. Heterozygosity was systematically reduced in more conserved genomic windows, a negative relationship observed consistently across all 219 species individually despite differences in baseline diversity. Linear mixed models confirmed this relationship even after accounting for species-specific diversity, which explained 27% of total variance independently of local constraint. The X chromosome showed lower diversity than autosomes beyond its elevated conservation.

In parallel, a phylogeny-aware DNA language model GPN-star was trained, providing a foundation for future inference of per-site conservation scores using a non-parametric approach.

Together, the results demonstrate that purifying selection leaves a consistent and recoverable imprint on primate genetic diversity across 100 kb genomic windows, detected within all 219 species regardless of their baseline genomic diversity.