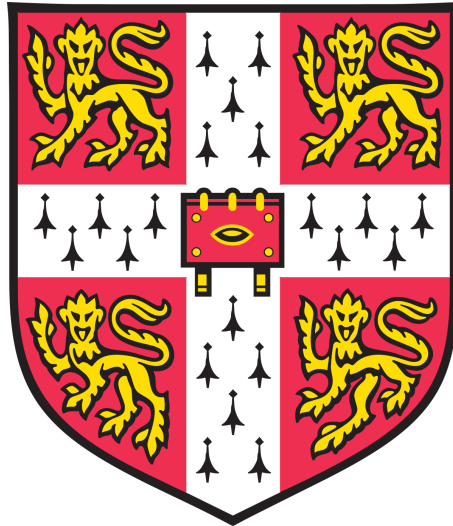


Integrative Analysis of the Human Gut Phageome

Using a Metagenomics Approach



Luis Fernando Camarillo Guerrero

Gonville & Caius College
University of Cambridge

This thesis is submitted for the degree of
Doctor of Philosophy

August 2020

Dedicated to my parents (Rosa María and Leopoldo) and sister (Marisol) for loving me and always supporting me in an unconditional manner from Day 1 of my life

Dedicada a mis padres (Leopoldo y Rosa María) y hermana (Marisol) por quererme y apoyarme de una manera incondicional desde mi primer día de vida.

Declaration

This thesis is the result of my own work and includes nothing which is the outcome of work done in collaboration except as declared in the Preface and specified in the text. It is not substantially the same as any that I have submitted, or, is being concurrently submitted for a degree or diploma or other qualification at the University of Cambridge or any other University or similar institution except as declared in the Preface and specified in the text. I further state that no substantial part of my thesis has already been submitted, or, is being concurrently submitted for any such degree, diploma or other qualification at the University of Cambridge or any other University or similar institution except as declared in the Preface and specified in the text. It does not exceed the prescribed word limit for the relevant Degree Committee

Luis Fernando Camarillo Guerrero

August 2020

Summary

Integrative Analysis of the Human Gut Phageome Using a Metagenomics Approach

Luis Fernando Camarillo Guerrero

Bacteriophages (or phages; viruses that infect bacteria and archaea) profoundly influence microbial communities. Given the impact of the gut microbiome composition and function on human health, there is a growing focus on phages that inhabit the gut ecosystem. However, the extent of viral diversity, biology, and worldwide epidemiology of gut phages remain largely unknown. In this thesis, I carry out a comprehensive genomic analysis of gut phages by harnessing the biggest collection of phage genomes, gut bacteria isolates, and human gut metagenomes.

I begin by introducing the Gut Phage Database (GPD) which is the largest genomic resource to date of human gut phage genomes and product of mining 28,060 faecal metagenomes and 2898 gut bacteria isolate genomes. I use machine learning to improve the quality of the predictions and investigate ways to organise the viral diversity in order to improve the characterisation of gut phages in downstream analyses.

Afterwards, I describe common functions and auxiliary metabolic genes encoded by human gut phages. I also highlight instances of hypervariable domains which may indicate the presence of phage receptor binding proteins. I then shift the focus to the analysis of two clades of gut phages, namely the Gubaphage and the *Picovirinae* subfamily. The Gubaphage is a novel phage clade uncovered in this work which is highly prevalent across the world. The *Picovirinae* clade was the most common predicted phage taxonomy in GPD. Host assignment allows me to study patterns of phage diversity across bacterial clades of the human gut and investigate their host range.

Finally, I analyse global patterns of the human gut phageome and its association with lifestyle and bacterial composition. I assess the idea of a core virome as well as in what degree my data agrees with this concept.

Acknowledgements

I would like to thank my supervisor Trevor Lawley for trusting me and allowing me to embark on a PhD project of my own. I'm aware of how fortunate I am to have had this freedom. Trevor, I am very grateful for your continuous guidance, and support throughout my PhD journey.

I would like to thank my colleagues from my lab team¹⁶² who I have been fortunate to learn a lot from, in particular Yan Shao and Hilary Browne. Yan not only became a great friend, he challenged me throughout my PhD with stimulating discussions about my project. Thanks Yan, I appreciate all your support and you intellectually pushed me to reach new heights in my knowledge about phages. Hilary, I appreciate your accessibility and the time you invested proofreading presentations/reports/thesis/manuscript, you give amazing advice.

I also would like to thank my manuscript collaborator Alexandre Almeida. Alex, your positive outlook was always refreshing. I learned a lot working along with someone of your scientific calibre. Thanks for all your advice and friendship.

I extend my gratitude to the Wellcome Sanger Institute for funding my PhD, Carl Anderson, Christina Hedberg-Delouka, Annabel Smith, and the Pathogen Informatics team.

I would like to acknowledge Darwin College which was the foundation of my social life in Cambridge. To all the people I met there, I appreciate your friendship and have fond memories of you. To the Darwinian Michael Schneider, who actually I met at Imperial College London, thanks for your friendship and for all the nights out at the Cambridge Bops. Honourable mention to the Darwin Salsa society, every week I was looking forward to go dancing on Monday at 8:00 pm. I would also like to thank the Mexicans living in Cambridge, I had so much fun hanging out with you all and really made me feel at home.

I thank Jose Manuel Aguilar Yañez for playing an integral role in my formation as a Scientist and always believing in me.

Last but not least, I thank my family to whom I dedicate this thesis. Your unconditional love and support were vital for me throughout this endeavour – *¡Lo logramos!*

Publications

Camarillo-Guerrero L.F., Almeida A., Rangel-Pineros G., Finn R.D., Lawley T. (2020). Massive expansion of human gut bacteriophage diversity. (In press). *Cell*.

Fung C., Tan S., Nakajima M., Skoog E.C., **Camarillo-Guerrero L.F.**, Klein J.A., Lawley T.D., Solnick J.V., Fukami T., Amieva M.R. "High-resolution mapping reveals that microniches in the gastric glands control *Helicobacter pylori* colonization of the stomach." *PLoS biology* 17.5 (2019): e3000231.

Contributions

This thesis is the result of my own work except:

- Metagenome assembly, sequence viral prediction with VirFinder and VirSorter, and dereplication at 95% sequence identity was carried out by Alexandre Almeida.
- Read mapping of GPD predictions to 28,060 metagenomes was carried out by Alexandre Almeida.
- Bacterial taxonomic assignment of gut isolates with the GTDB toolkit was carried out by Alexandre Almeida.
- The tool to assign a taxonomic rank to GPD predictions was developed by Guillermo Rangel Pineros.

Abbreviations

Acr	Anti-CRISPR
Abi	Abortive infection
AMG	Auxiliary Metabolic Gene
ANI	Average Nucleotide Identity
ARG	Antibiotic Resistance Gene
BACON	Bacteroidetes-Associated Carbohydrate-binding Often N-terminal
BAM	Bacteriophage Adherence to Mucus
CRISPR	Clustered Regularly Interspaced Short Palindromic Repeats
CTHR	Collagen Triple Helix Repeat
DNA	Deoxyribonucleic acid
FP	False Positive
GPD	Gut Phage Database
GTDB	Genome Taxonomy Database Toolkit
HGT	Horizontal Gene Transfer
HMM	Hidden Markov Model
IBD	Inflammatory Bowel Disease
ICTV	International Committee on Taxonomy of Viruses
ICE	Integrative and Conjugative Element
Ig	Immunoglobulin
ImmeDB	Intestinal Microbiome Mobile Elements Database
KEGG	Kyoto Encyclopaedia of Genes and Genomes
MCL	Markov Cluster
MGE	Mobile Genetic Element
ML	Machine Learning
NCBI	National Centre for Biotechnology information
OMV	Outer Membrane Vesicle
PCA	Principal Component Analysis
PC	Protein Cluster
PD	Parkinson's Disease
PICI	Phage-Inducible Chromosomal Islands
PtW	Piggyback-the-Winner

QC	Quality Control
RBP	Receptor Binding Protein
RNA	Ribonucleic acid
RT	Reverse Transcriptase
RM	Restriction Modification
SaPI	Staphylococcus Aureus Pathogenicity Islands
VC	Viral Cluster
VLP	Viral-Like Particle
VMR	Virus to Microbe Ratio

Figure 1.1. The lifecycles of bacteriophages.....	2
Figure 1.2. The diversity of phages	3
Figure 1.3. Bacterial anti-phage defences	7
Figure 1.4. Taxonomy of phages	10
Figure 1.5. Main bacteriophage morphological types detected in a faecal sample	13
Figure 1.6. Phage dynamics in the human gut.....	15
Figure 1.7. Inter- and intra-diversity and stability of the human gut phageome	18
Figure 3.1. Generation of the Gut Phage Database (GPD)	33
Figure 3.2 A machine learning approach to distinguish phages from ICEs.....	35
Figure 3.3. GPD taxonomy assignment and comparison against other gut phage databases .	37
Figure 3.4. Genome completeness of GPD	40
Figure 3.5. Clustering of phages into VCs	41
Figure 3.6. Distribution of genomes per VC and phylogenetic structure of GPD	43
Figure 3.7. DotBlast tool	46
Figure 3.8. HyperVir tool.....	48
Figure 3.9. vMatch inner workings and visualization of output matrix.....	51
Figure 4.1. Functions encoded by the human gut phageome	60
Figure 4.2. Protein clusters (PCs) encoded by gut phages	61
Figure 4.3. Hypervariable domains can narrow down protein function in phages.....	63
Figure 4.4. Investigation of the Gubaphage clade relationship to other crAss-like phages....	66
Figure 4.5. Expansion of the Picovirinae subfamily.....	68
Figure 4.6. Viral diversity across gut bacteria clades	73
Figure 4.7. Host range of gut phages	75
Figure 5.1. Rarefaction curves for viral richness.....	82
Figure 5.2. Human lifestyle is associated with global gut distribution of phageome types....	84
Figure 5.3. Phage carriage.....	86
Figure 5.4. Rank prevalence curve for VCs	87
Figure 5.5. Global gut phage clades and their bacterial hosts.....	90
Figure 5.6. Investigating the concept of a core-virome	93

Table of Contents

Declaration	i
Summary	ii
Acknowledgements.....	iii
Publications	iv
Contributions	v
Abbreviations.....	vi
List of Figures.....	viii
Table of Contents.....	ix
Chapter 1: Introduction.....	1
1.1 General overview of bacteriophages.....	1
1.1.1 The life cycle of bacteriophages.....	1
1.1.2 The outstanding diversity of bacteriophages.....	2
1.1.3 Phages: friends or foes of bacteria?.....	4
1.1.4 The arms-race between phage and bacteria.....	5
1.1.5 Evolutionary phage-host dynamics.....	7
1.1.6 Predator-prey dynamics	8
1.1.7 Taxonomy and the recent explosion of phage diversity.....	8
1.1.8 Prediction of phages from metagenomic sequences	11
1.2 Bacteriophages in the human gut.....	12
1.2.1 Discovery and isolation of faecal VLPs.....	12
1.2.2 Taxonomy of gut phages.....	12
1.2.3 The case of the crAssphage.....	13
1.2.4 Phage dynamics in the human gut.....	14
1.2.5 From lysogeny to the lytic cycle in the gut.....	15
1.2.6 Hosts and host ranges of gut phages.....	16
1.2.7 Commonly encoded genes	16
1.2.8 Stability, inter- and intra-diversity of the human gut phageome.....	17
1.2.9 Gut phages and human disease.....	19
1.2.10 Phage therapy	19
1.3 Thesis aims.....	20
Chapter 2: Methods	21
2.1 Chapter 3: The Gut Phage Database.....	21
2.1.1 Metagenome assembly.....	21
2.1.2 Viral sequence prediction.....	21
2.1.3 Sequence clustering.....	22
2.1.4 Quality control of GPD predictions.....	22
2.1.5 Genome completeness and contamination	23
2.1.6 Viral taxonomic assignment.....	23
2.1.7 Clustering of phages into VCs.....	23
2.1.8 Bioinformatics tools	24
2.2. Chapter 4: Biology of human gut phages.....	25
2.2.1 Detection of function in gut phages.....	25
2.2.2 Clustering of proteins into protein clusters (PCs).....	25
2.2.3 Phylogenetic analyses.....	25

2.2.4 Taxonomic assignment of bacterial genomes	25
2.2.5 Host assignment	26
2.2.6 Assessing viral diversity patterns	26
2.2.7 Host range analysis.....	26
2.3 Chapter 5: Global distribution and epidemiology of gut phages.....	28
2.3.1. Metagenomic read mapping.....	28
2.3.2 Dependency of phages detected and sample sequencing depth.....	28
2.3.3. Geographical distribution of metagenomic samples.....	28
2.3.4 Calculation of phage carriage.....	28
2.3.5 Detection of enterotypes targeted by VCs	29
2.3.6 Network of globally distributed phages.....	29
2.3.7 Core virome analyses.....	29
2.4 GPD resource and metadata	30
Chapter 3: The Gut Phage Database.....	31
3.1 Introduction and aims	31
3.2 Results and discussion	33
3.2.1 Construction of the gut phageome database (GPD).....	33
3.2.2 Decontamination using a machine learning approach	34
3.2.3 GPD significantly expands gut bacteriophage diversity	36
3.2.4 Genome completeness	38
3.2.5 Clustering of phages into VCs.....	40
3.2.6 Viral clusters reconstruct the phylogenetic structure of gut phages	41
3.2.7 Bioinformatics tools	44
3.2.8 Synteny analysis for viral genomes (dotBlast).....	44
3.2.9 Hypervariation analysis (hyperVir)	47
3.2.10 Exploring viral taxonomy through shared protein clusters (vMatch).....	49
3.3 Conclusions.....	52
Chapter 4: Function, phylogeny and host assignment of gut phages.....	54
4.1 Introduction and aims	54
4.2 Results and discussion	56
4.2.1 Functions encoded by gut phages.....	56
4.2.2 Protein clusters encoded by gut phages	60
4.2.3 Identification of hypervariation domains uncovers putative phage tropism determinants..	62
4.2.4 The Gubaphage represents a novel clade of gut phages	63
4.2.5 Expansion of the Picovirinae subfamily	66
4.2.6 Viral diversity across gut bacteria clades.....	70
4.2.7 Evaluating host range of gut phages.....	73
4.3 Conclusions.....	77
Chapter 5: Global distribution and epidemiology of gut phages	80
5.1 Introduction and aims	80
5.2 Results and discussion	81
5.2.1 Saturation curves for VCs.....	81
5.2.2 Human lifestyle associated with global gut distribution of phageome types	82
5.2.3 Phage carriage	85
5.2.4 Uncovering most prevalent phage in global human populations.....	86
5.2.5 Global distribution of 280 dominant human gut phages.....	88
5.2.6 Investigating the concept of a core-virome	90
5.3 Conclusions.....	94

Chapter 6: Summary and future work	96
6.1 Summary	96
6.1.1 Development of the GPD	96
6.1.2 Characterising phage functions and host range	97
6.1.3 Epidemiology of gut phages.....	99
6.2 Main findings of this work.....	100
6.3 Future work	101
References	103
Appendix 1. Predicted hosts of the crAss-like family	113
Appendix 2. Metadata of deeply sequenced samples.....	115