

Ancestral inference and encoding of bacterial nonsense mutations for phylogenetic genome-wide association studies and convergence analysis

Matthew P. Moore^{1,2}, Aditya K. Lankapalli¹, Xavier Didelot^{1,2}
¹ University of Warwick 2 Health Protection Research Unit in Gastrointestinal Infections

Background

Adaptation by loss of function has been observed in many human-adapted, pathogenic bacteria¹⁻³. In the *Shigella* spp., convergent loss of function is well characterised and the smaller genome size, compared with *E. coli*, is likely the result of genome streamlining⁴. Across *Shigella dysenteriae*, *Shigella flexneri* and *Shigella sonnei* parallel loss of function events have resulted in convergent metabolic capability reduction and accounted for 47% of within-species metabolic gene variation⁴. A study of *Salmonella enterica* serotype Typhimurium found that pseudogenisation was linked to host adaptation including, for example, in the *sseL* gene linked, via inactivation, with systemic dissemination. Further, pseudogenisation of *sseK2*, *sseK3*, *avrA* and *sseL* involved in

activating the pro-inflammatory response during infection were observed⁵.

Bacterial genome-wide association studies (GWAS) have successfully been used to make associative and correlative inferences about the genetic basis of phenotypes of interest⁶⁻⁸. Confounding any GWAS inferences are the clonality of bacterial population structures⁹, genetic content variability¹⁰, imperfect heritability¹¹ and homologous recombination¹². Corrections for confounding population structure have been developed based on phenotype rearrangement, clustering and dimension-reduction¹³⁻¹⁹. Methods have been developed that take advantage of

the population structure, providing phylogenetic approaches. Inferences can be made between pseudogenisation events at the gene level, though possibly resulting from many independent nonsense mutations, by collapsing to a binary genotype as with gene presence/absence data. Additionally, correlations may be inferred between an ancestral change in both genotype and phenotype. However, providing binary genotype data of LoF per gene, from independent but phylogenetically clustered nonsense mutations could incorrectly appear as a single ancestral event.

Methods

For the GWAS analysis a dataset of avian *E. coli* was selected, comprised of both infectious and asymptomatic carriage isolates from chickens, with some additional environmental isolates²⁰. Variants were called by Snippy²¹ against the reference genome (APEC078), providing SnpEff²² formatted variant call file (.vcf) output of variant impact on gene function. Variants categorised as 'HIGH' impact by SnpEff were included in the analysis of loss of function. Genome-wide associations were performed for each phylogroup, as in the original study: A (n=71), B1 (n=85), B2 (n=152) and G/ST-117 (n=220) with TreeWAS2. Core SNP trees were generated for each phylogroup. Core genome alignments were generated with snippy-multi²¹, IQ-TREE²³ with the GTR+I+G substitution model and 1000 ultrafast bootstraps (UFBoot2)²⁴ was used to generate an initial phylogeny and ClonalFrameML²⁵ was used to detect regions of homologous recombination and adjust branch lengths to non-recombinant SNPs only.

All shared mutations (same position&variant) were subjected to tests of ancestry. The first test, whether the nonsense mutation is monophyletic is performed with the ETE 3 Toolkit²⁶. The most recent common ancestor (MRCA) of all genomes with the mutation is determined and if they comprise all the children of that node, ancestry is asserted. If polyphyly is observed, ancestral state reconstruction (ASR) is conducted for every internal node back to the nonsense mutation MRCA with PastML²⁷. All variants at the reference position are provided as mutated reference codons unless there is no variation, whereby the reference codon is provided. Where nonsense mutations are nested within an ancestral nonsense mutation clade, the ancestral nonsense mutation is considered to be canonical. Whether mutations are canonical will determine the pattern of ancestry that is reported for that LoF event; ignoring those that are nested.

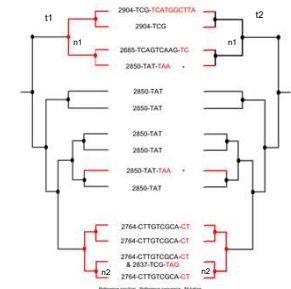


Figure 1. The imagined core genome phylogeny of 14 genomes (11&2) with some nonsense alleles in a single gene. In t1 branches are coloured red based on the ASR of binary nonsense allele status (t1) or retention of true status with Nonsense/Canonical (t2). MRCA node n2, is inferred to carry the ancestral-state of the frame-shifting deletion at reference position 2764. However, one genome in the clade also has a premature stop codon. The ancestral mutation is taken as canonical. (*) Shows a parallel mutation

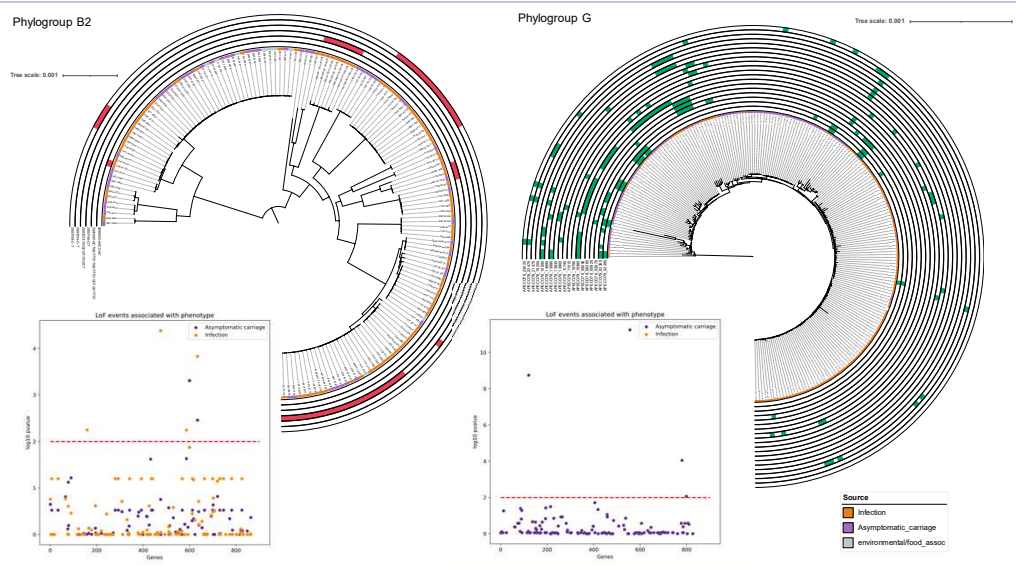


Figure 2. (Left) Core genome phylogeny of *E. coli* phylogroup B2 isolates. Colour rings correspond to phenotypes (inner) and *cirA* catecholase siderophore receptor gene nonsense mutations (outer rings). Each colour ring moving outwards corresponds to a different nonsense mutation (labelled). (Scatterplot, left) shows the p-values (log10) where nonsense mutations are associated with the 'Infection' or 'Asymptomatic carriage' phenotype. Above the red dashed line represents a p-value of <0.1. (Scatterplot, right) displays significant values for phylogroup G. (Right) Core genome phylogeny of *E. coli* phylogroup G isolates. Colour rings correspond to phenotypes (inner) and moving outwards, a selection of genes with the largest number of independent nonsense events. Green blocks show where a genome carries (any) inactivating mutation in the gene (labelled)

Table 1. Group B2 and G genes with >4 parallel inactivation events, hypothetical genes not included

Phylogroup	Locus ID	No. parallel nonsense events	Product
GroupG	APEC078_05430	21	Diguanylate cyclase <i>adpA</i>
	APEC078_02390	21	Predicted diacylglyceride-binding enzymes
	APEC078_02385	15	Transposase and inactivated derivatives
	APEC078_19150	13	G-TAT mismatch-specific DNA glycosylase <i>yggF</i>
	APEC078_09920	9	bifunctional UDP-sugar hydrolase/5'-nucleotidase periplasmic precursor <i>ushA</i>
	APEC078_22475	8	Beta-glucuronidyl-epiphyse-beta-glucuronidase-beta-lactamase <i>bgbB</i>
	APEC078_10385	8	4-aminobutylamine aminotransferase and related aminotransferases <i>gusE/gusC</i>
	APEC078_14690	8	Site-specific recombinases, DNA integrase <i>Pin</i> homologs <i>conJ2</i>
	APEC078_11515	7	Putative tail fiber chaperone, <i>Qin</i> phage <i>shoQ</i>
	APEC078_17410	6	CRISPR-associated helicase <i>Cas3</i>
	APEC078_23025	6	Mn-dependent Dnaase <i>yggX</i>
	APEC078_11740	6	Kraeonic dihydroxybenzoates, typically selenocysteine-containing <i>yggF</i>
	APEC078_09000	6	Glucose-1-phosphatase/inositol phosphatase
	APEC078_19255	6	Inner membrane protein <i>yggV</i> family
	APEC078_13800	6	Predicted transposon component
GroupB2	APEC078_15690	9	Cellulase M and related proteins <i>yggE</i>
	APEC078_06385	6	Glycosyltransferases, probably involved in cell wall biogenesis <i>mfbB</i>
	APEC078_13990	5	Nucleoside-diphosphate-sugar epimerases <i>bcfI</i>
	APEC078_10790	5	NAD-dependent aldehyde dehydrogenases <i>aldA</i>
	APEC078_13835	5	Hisdinal dehydrogenase
	APEC078_03975	5	Type II secretory pathway, ATPase <i>PilE/PilT</i> , pilus assembly pathway, ATPase <i>PilB</i>
	APEC078_21370	5	Arac-type DNA-binding domain-containing proteins <i>gadX</i>
	APEC078_08030	10	Nitroreductase <i>AraA</i>
	APEC078_14490	7	Calcitriol siderophore receptor <i>cirA</i>
	APEC078_02375	6	Transposase and inactivated derivatives
	APEC078_20405	6	Putative peptidoglycan-binding domain-containing protein <i>TSS3</i> <i>gpaA</i>
	APEC078_10515	5	Hydrolase <i>gpgZ</i>
	APEC078_09920	5	bifunctional UDP-sugar hydrolase/5'-nucleotidase periplasmic precursor <i>ushA</i>
	APEC078_14995	5	Type V secretory pathway, adhesin <i>aldA</i>
	APEC078_24655	5	Enterocytin <i>yggL</i>

Results

It was possible to automatically determine nonsense ancestry patterns for most genes with nonsense mutation(s) in >1 genome. For the largest two phylogroups, groupB2 and groupG: 97.95% and 98.57% of genes ancestries were inferred, respectively. Table 1 displays the genes with the largest number of parallel LoF events in groupG & B2 genomes. Group G is thought to be an avian-host specialising lineage and as such convergence on LoF may indicate adaptive trends both towards and away from infectivity/pathogenicity. As groupG is specialised to the avian host

many genes observed to have high levels of nonsense convergence in other groups are already missing or with a high proportion of inactivation in group G. Others are observed to be convergent across lineages such as *ushA*, which has convergently been inactivated across groups B1 (9 events), B2 (5 events) and G (9 events). An *ushA* homolog has additionally been reported to be adaptively inactivated in *Salmonella enterica*. Others such as *bgB* (8 inactivation events in groupG) is a known (strain-specific) pseudogene in *Shigella flexneri*.

Determination of the ancestral patterns also allowed a full, phylogenetic GWAS in order to infer genes, that by their inactivation, are associated with phenotypic switches. GroupG had 4 significant nonsense-genes (<.01), all associated with asymptomatic carriage. In group B2, 4 genes were observed to be associated, with infection with another 2 associated with asymptomatic carriage.

Discussion

Across all lineages the nonsense ancestry patterns were determined for >98% of genes with nonsense mutations in >1 genome (278-388 genes). However, genes were later filtered due to uncertainty at some nodes during ASR. Development is ongoing on how to treat these sites. Development is also ongoing into breaking ancestry due to differential patterns of homologous recombination. Finally, work is ongoing into determining the encoding of very-similar indels, compensatory mutations, reversals and determining the most appropriate ASR models to infer ancestral phenotypes. At present, runtime ranged from 8-21 minutes. However, the rate-limiting step of ASR with pastML is readily parallelisable

- Seakwiche, E. V., Hest, D. L. & Dijkshuis, D. F. Pathoadaptive mutations drive loss and variation in bacterial pathogens. *Frontiers in Microbiology* (2019). doi:10.3389/fmicb.2019.01948
- Seakwiche, E. V., Hest, D. L., Peto, T. E., Cook, D. W. & Wilson, D. J. Within-host evolution of bacterial pathogens. *Nature Reviews Microbiology* (2020). doi:10.1038/s41579-020-123-13
- Hest, D. J. et al. Bacterial Adaptation through Loss of Function. *PLoS Genet* (2020).
- Hest, D. J., Moore, M. P., Lankapalli, A., Wilson, D. J. & Didelot, X. Impact of insertion sequences on convergent evolution of *Shigella* species. *PLoS Genet* (2020).
- Brown, A. et al. Evolution of *Salmonella enterica* serotype typhimurium driven by convergent selection and niche adaptation. *PLoS Genet* (2020).
- Didelot, X. et al. Using genome-wide association studies: A new direction for bacteriology. *Genome Medicine* (2024).
- Chen, P. & Shapiro, B. J. The advent of genome-wide association studies for bacteria. *Curr Opin Microbiol* (2021). doi:10.1016/j.coi.2021.04.004
- Didelot, X., Lankapalli, A. & Didelot, X. Microbial genome-wide association studies: lessons from human GWAS. *Nature Reviews Genetics* (2016). doi:10.1038/nrg.2016.137
- Didelot, X., Lankapalli, A. & Didelot, X. Inference of homologous recombination structure using whole-genome sequencing. *Genetics* (2020). doi:10.1534/genetics.120.1221
- Madire, D., Dancic, C., Tettelin, M., Mangoni, V. & Rappelli, R. The microbial pan-genome. *Current Opinion in Genetics and Development* (2020). doi:10.1016/j.cog.2020.05.005
- Arm, A., M., & Didelot, X. Bayesian inference of the evolution of a phenotype distribution in a phylogenetic tree. *Genetics* (2016). doi:10.1534/genetics.125.19096
- Didelot, X. & Mankin, M. C. Impact of recombination on bacterial evolution. *Frontiers in Microbiology* (2020).
- Didelot, X., Dancic, C. & Didelot, X. Discontinuous analysis of principal components: A new method for the analysis of genetically structured populations. *BMC Genet* (2020).
- Wernert, L. A. et al. Genetic signatures of human and animal disease in the *Escherichia coli* pan-genome. *Nat. Rev. Genet.* (2015). doi:10.1038/nrg.2015.780
- Didelot, X. et al. Comprehensive Identification of Single Nucleotide Polymorphisms Associated with Bacterial Resistance within *Phenoxymethylpenicillin*. *PLoS Genet* (2014).
- Didelot, X. et al. The use of genome-wide association studies to determine the genetic basis of bacterial phenotypes. *PLoS Genet* (2014). doi:10.1371/journal.pgen.1003848
- Earl, S. K. et al. Identifying lineage effects when controlling for population structure improves power in human genome studies. *Nat. Microbiol* (2016). doi:10.1038/nrmicro1221
- Braydon, D., Babin, L., Schiffer, L. & Didelot, X. Rapid scoring of genes in microbial pan-genome-wide association studies with Scan. *Genome Biol.* (2020).
- Maguire, L. et al. Genome evolution and the emergence of pathogenicity in *Escherichia coli*. *Nat. Rev. Genet.* (2021). doi:10.1038/s41579-021-00886-6
- Seakwiche, E. V. et al. Available at: <http://www.victronics.com/software/scanpy/scanpy>
- Capelle, E. et al. A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff. SNPs in the genome of *Drosophila melanogaster* strain w1118; iso2-iso4; iso6-iso72; iso864. *BMC Bioinform.* (2015). doi:10.1186/s12859-015-0601-4
- Nguyen, L. T., Schmidt, H. A., von Haeseler, A., Minh, B. Q. & Truong, L. A. Efficient local bootstrap estimation for assessing the accuracy of maximum likelihood phylogenies. *Mol. Biol. Evol.* (2015). doi:10.1093/molbev/mst021
- Hoang, T. H., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum likelihood phylogenies. *Mol. Biol. Evol.* (2018). doi:10.1093/molbev/mty151
- Didelot, X. et al. Reconstructing the ancestral history of *Salmonella enterica* serovar *Paratyphi*. *PLoS Genet* (2015). doi:10.1371/journal.pgen.1004641
- Didelot, X. et al. Reconstructing the ancestral history of *Salmonella enterica* serovar *Paratyphi*. *PLoS Genet* (2015). doi:10.1371/journal.pgen.1004641
- Didelot, X. et al. Reconstructing the ancestral history of *Salmonella enterica* serovar *Paratyphi*. *PLoS Genet* (2015). doi:10.1371/journal.pgen.1004641
- Didelot, X. et al. Reconstructing the ancestral history of *Salmonella enterica* serovar *Paratyphi*. *PLoS Genet* (2015). doi:10.1371/journal.pgen.1004641