

# What is computational biology?

David Welch

School of Computer Science, Centre for Computational Evolution, and Te  
Pūnaha Matatini

# What is biology?

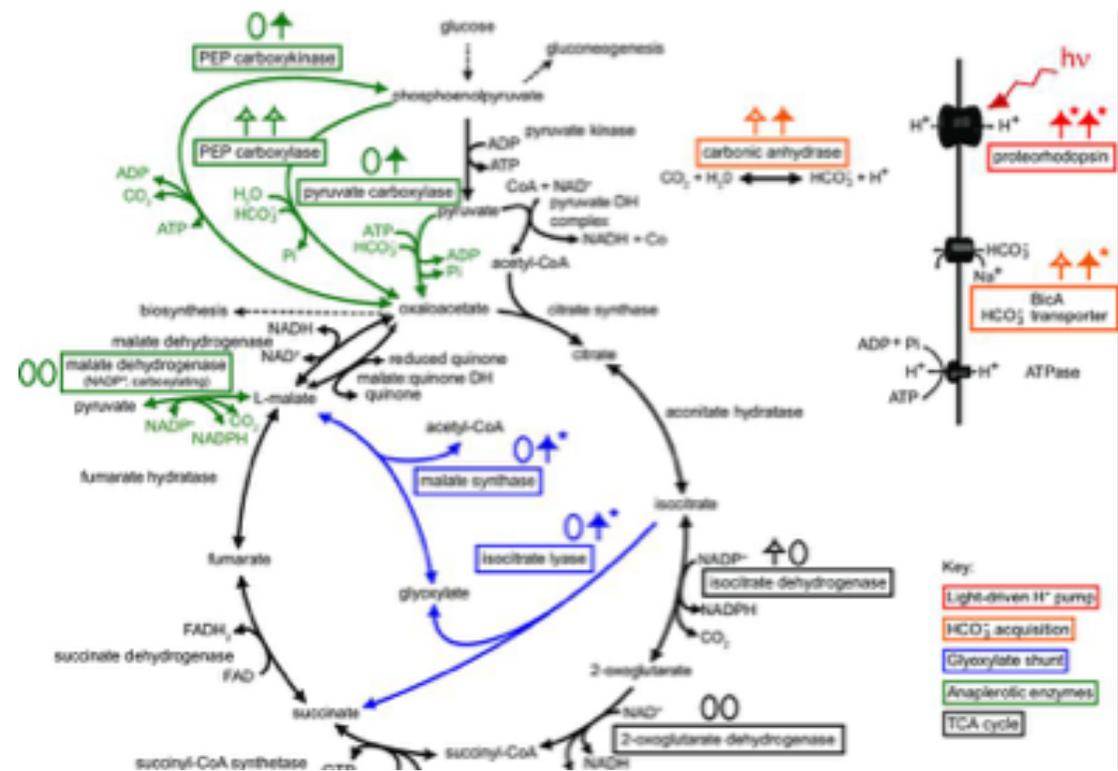
Biology is the study of life. This is a broad target!

- **individual organisms**
- **populations of organisms**
- **evolution of populations and species**
- **ecological systems** (interactions between diverse populations)
- **organism subsystems** (e.g. organs)
- **cell biology**
- **genetics**
- **metabolism**
- **ethology** (behaviour)
- **the origin of life** (how is life different from non-life, how did life emerge?)
- And more....!

# Biological systems differ from eg physical ones

Biological systems are **complex**

- lots of interacting parts
- parts are different from each other (**heterogeneity**)
- the parts interact **non-linearly**





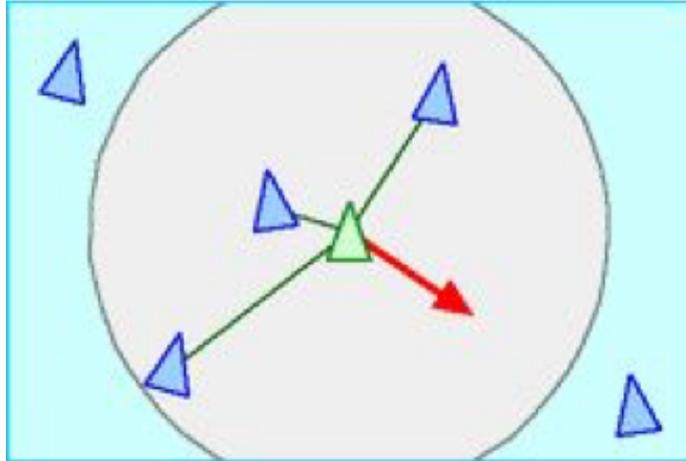
## 3 examples

- Boids: modelling, simulation, emergence
- Alignment and structural prediction: algorithms, gamification, heuristics
- Phylodynamics of COVID-19: inference, big data, visualisation

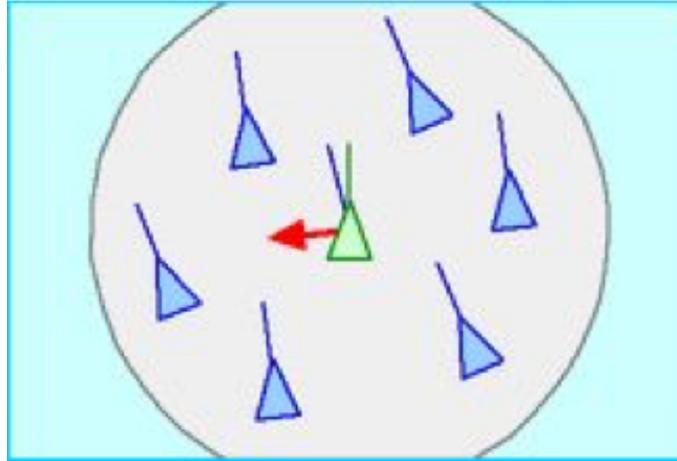
# Boids: a simple model of flocking

Each bird is following a few simple rules and the complex pattern simply emerges from them

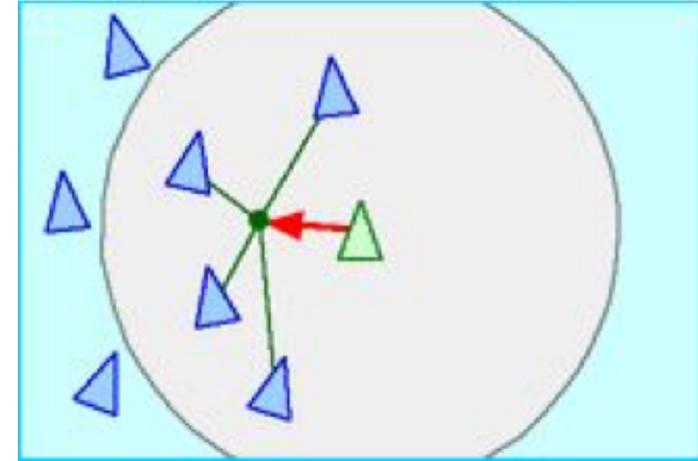
Separation



Alignment



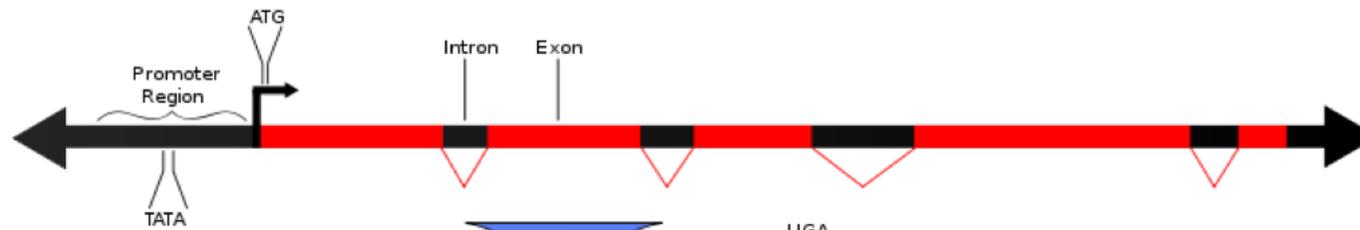
Cohesion



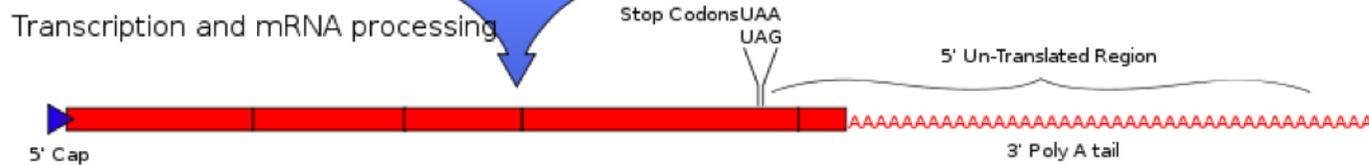
Images taken from Craig Reynold's website. <http://www.red3d.com/cwr/boids/>

# Central Dogma of Molecular Biology : Eukaryotic Mode

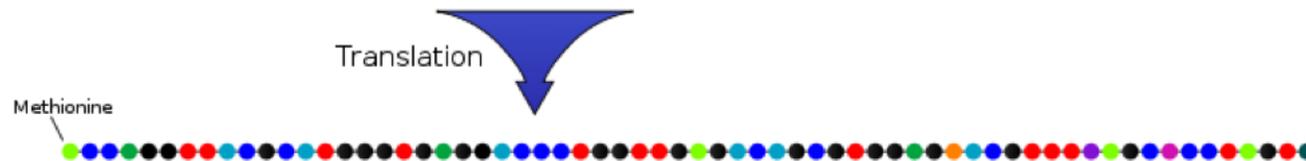
DNA



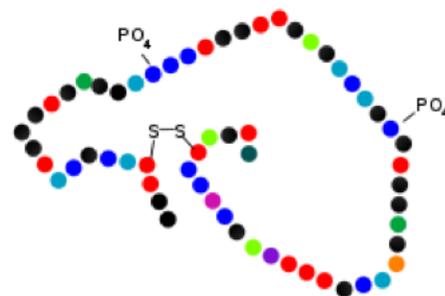
mRNA



Protein



Post-Translational Modification



# Alignments of Human Beta-Globin to Other Globins



Human beta-globin VHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGN  
 LTPEE VT LWGKVV VVGGGALGRLLVVYPWTQRFFESFGDLS PDA MGN  
 Ring-tailed lemur beta-globin TFLTPEENGHVTSLWGKVNVEKVGGEALGRLLVVYPWTQRFFESFGDLSSPDA MGN  
 PKVKAHGKKVLGAFSDGLAHLNLDKGTFFATLSELHCDKLHVDPENFRLLGNVLCVLAHHFGKEFTPPVQAAAYQKVVAGVANALAHKYH  
 PKVKAHGKKVL AFS GL HLDNLKGTFA LSELHC LHVDPENF LLGNVLV VLAHHFG F P QAA QKVV GVANALAHKYH  
 PKVKAHGKKVLSAFSEGLHHLNLDKGTFAQLSELHCVLHVDPENFKLLGNVLVIVLAHHFGNDFSPQTQAAFQKVVIGVANALAHKYH



Human beta-globin VHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNP  
 V T E SA LWGK N DE G AL R L VYPWTQR F FG LS P A MGNP  
 Goldfish beta-globin VEWTDAERSAIIIGLWGKLNDELGPQALARCLIVYPWTQRYFATFGNLSSPAAIMGNP  
 KVKAHGKKVLGAFSDGLAHLNLDKGTFFATLSELHCDKLHVDPENFRLLGNVLCVLAHHFG--KEFTPPVQAAAYQKVVAGVANALAHKYH  
 KV AHG V G DN K T A LS H KLHVDP NFRLL A FG F VQ A QK V AL YH  
 KVAAHGRTVMGGLERAIKNMDNIKATYAPLSVMHSEKLVHVDPNFRLLADCITVCAAMKFGPSGFNADVQEAWQKFLSVVVSALCRQYH



Human beta-globin VHLTPEEKSAVTALW----GKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVKA  
 L V W G N V G E L F F S P V  
 Bloodworm globin IV MGLSAAQRQVVASTWKDIAGSDNGAGVGKECFKFLSAHHDIAAVF-GFSGAS-----DPGVAD  
 HGKKVLGAFSDGLAHL-DNLKGTFFATLSELHCDK----LHVDPENFRLLGNVLCVLAHHFGKEFTPPVQAAAYQKVVAGVANALAHKYH  
 G KVL D HL D K K H E F LG L H G T A A AL  
 LGAKVLAQIGVAVSHLGDEGKMVAEMKAVGVRHKGYGYKHIAEYFEPLGASLLSAMEHRIGGKMTAAAKDAWAAAYADISGALISGLQ



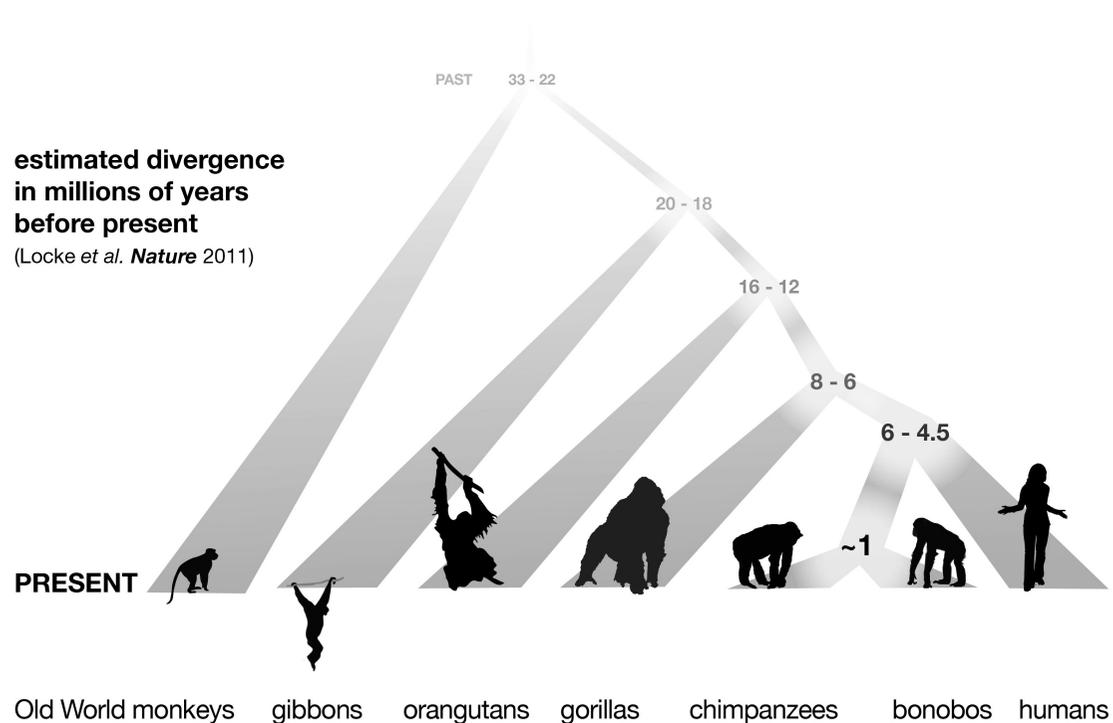
Human beta-globin VHLTPEEKSAVTALWG--KVVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPK  
 V T V K N L P F P NPK  
 Soybean leghemoglobin VAFTEKQDALVSSSFEAFKANIPQYSVVFYTSILEKAPAAKDLFSLANGVDPT----NPK  
 VKAHGKKVLGAFSDGLAHLNLDKGTFA--TLSELHCDKLHVDPENFRLLGNVLCVLAHHFGKEFTPPVQAAAYQKVVAGVANALAHKYH  
 H K D L A L H K DP F L G A A A  
 LTGHAEKLFALVRDSAGQLKASGTVVADAALGSVHAQKAVTDPQ--FVVVKEALLKTIKAAVGDKWSDELREWEVAYDELAAAIKKA--

# How to align multiple sequences?

C_aminophilum	AGCT.YCGCA	TGRAGCAGTG	TGAAAA.....	.....ACTCCGGT	GGTACAGGAT
C_colinum	AGTA..GGCA	TCTACAAGTT	GGAAAA.....	.....ACTGAGGT	GGTATAGGAG
C_lentocellum	GGTATTCGCT	TGATTATNAT	AGTAAA.....	.....GATTTATC	GCCATAGGAT
C_botulinum_D	TTTA.TGGCA	TCATACATAA	AATAATCAAA	.....GGAGCAATCC	GCTTTGAGAT
C_novyi_A	TTTA.CGGCA	T....CGTAG	AATAATCAAA	.....GGAGCAATCC	GCTTTGAGAT
C_gasigenes	AGTT.TCGCA	TGAAACA...	GC.AAATTA...	.....GGAGAAATCC	GCTATAAGAT
C_aurantibutyricum	A.NT.TCGCA	TGGAGCA...	AC.AATCAAA	.....GGAGCAAT.C	ACTATAAGAT
C_sp_C_quinii	AGTT.T.GCA	TGGGACA...	GC.AAATTA...	.....GGAGCAATCC	GCTATGAGAT
C_perfringens	AAGA.TGGCA	T.CATCA...	TTCAACCAAA	.....GGAGCAATCC	GCTATGAGAT
C_cadaveris	TTTT.CTGCA	TGGGAAA...	GTC.ATGAAA	.....GGAGCAATCC	GCTGTAAGAT
C_cellulovorans	ATTC.TCGCA	TGAGAGA...	.TGATCAAA	.....GGAGCAATCC	GCTATAAGAT
C_K21	TTGR.TCGCA	TGATCKAAAC	ATCAAAGGAT	..TTTTCTTTGGAAAATTC	ACTTTGAGAT
C_estertheticum	TTGA.TCGCA	TGATCTTAAC	ATCAAAGGAA	..TTT..TTTCGG..AATTC	ACTTTGAGAT
C_botulinum_A	AGAA.TCGCA	TGATTTTTCTT	ATCAAAGATT	..T.....ATT..	GCTTTGAGAT
C_sporogenes	AGAA.TCGCA	TGATTTTTCTT	ATCAAAGATT	..T.....ATT..	GCTTTGAGAT
C_argentinense	AAGG.TCGCA	TGACTTTTAT	ACCAAAGGAG	..T.....AATCC	GCTATGAGAT
C_subterminale	AAGG.TCGCA	TGACTTTTAT	ACCAAAGGAG	..T.....AATCC	GCTATGAGAT
C_tetanomorphum	TTTT.CCGCA	TGAAAAACTA	ATCAAAGGAG	..T.....AAT.C	GCTTTGAGAT
C_pasteurianum	AGTT.TCACA	TGGAGCTTTA	ATTAAGGAG	..T.....AATCC	GCTTTGAGAT
C_collagenovorans	TTGA.TCGCA	TGGTCGAAAT	ATTAAGGAG	..T.....AATCC	GCTTACAGAT
C_histolyticum	TTTA.ATGCA	TGTTAGAAAG	ATTAAGGAG	.....CAATCC	GCTTTGAGAT
C_tyrobutyricum	AGTT.TCACA	TGGAATTTGG	ATGAAAGGAG	..T.....AATTC	GCTTTGAGAT
C_tetani	GGTT.TCGCA	TGAAACTTTA	ACCAAAGGAG	..T.....AATCT	GCTTTGAGAT
C_barkeri	GACA.TCGCA	TGGTGTT...	.TTAATGAAA	.....ACTCCGGT	GCCATGAGAT
C_thermocellum	GGCA.TCGTC	CTGTTAT...	.CAAAGGAGA	.....AATCCGGT	...ATGAGAT
Pep_prevotii	AGTC.TCGCA	TGGNGTTATC	ATCAAAGA..	.....TTTATC	GGTGTAAAGAT
C_innocuum	ACGGAGCGCA	TGCTCTGTAT	ATTAAGCGC	CCTTCAAGGCGTGAAC....	...ATGGAT
S_ruminantium	AGTTTCCGCA	TGGGAGCTTG	ATTAAGATG	GCCTCTACTTGTAAGCTATC	GCTTTGCGAT

# A phylogeny shows relationships between sampled species, populations, or individuals

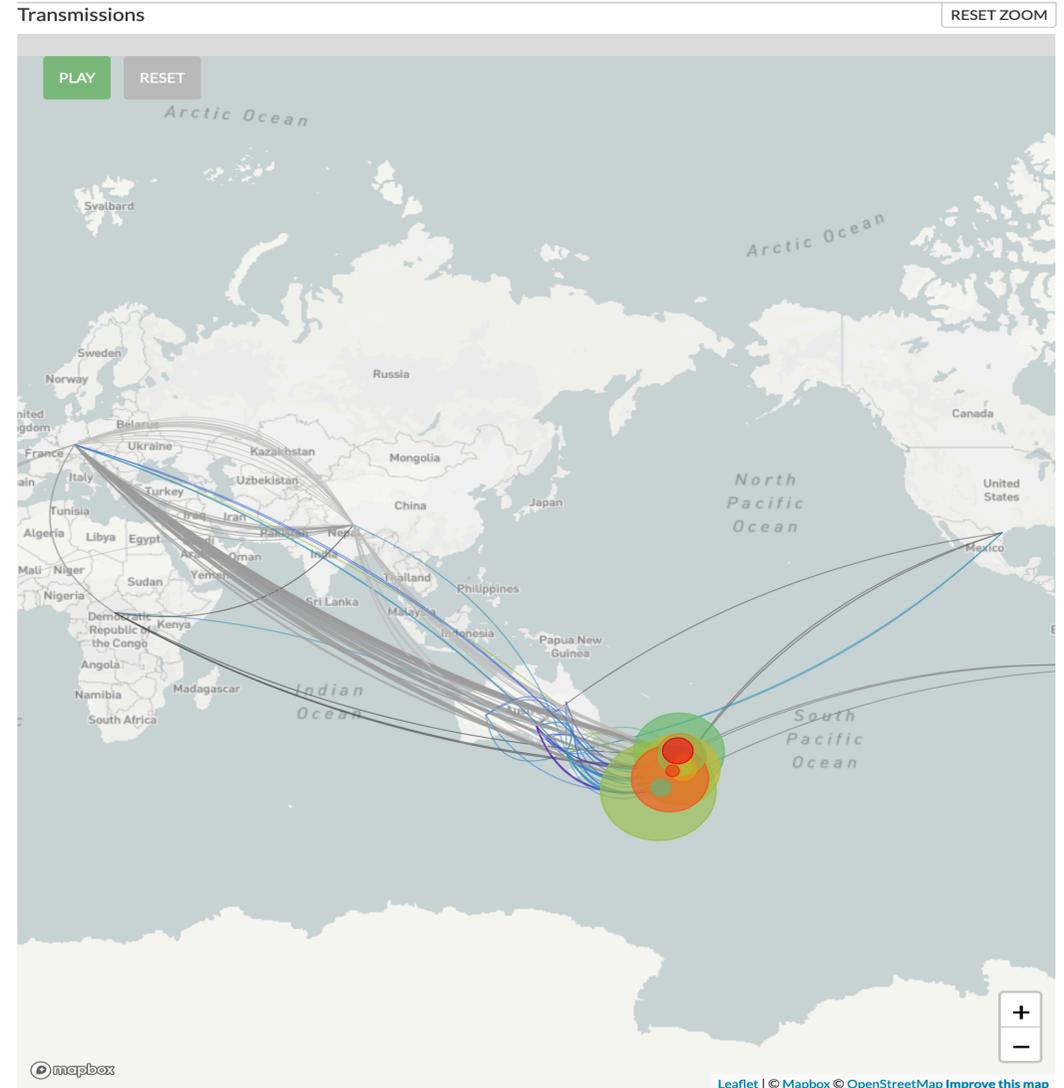
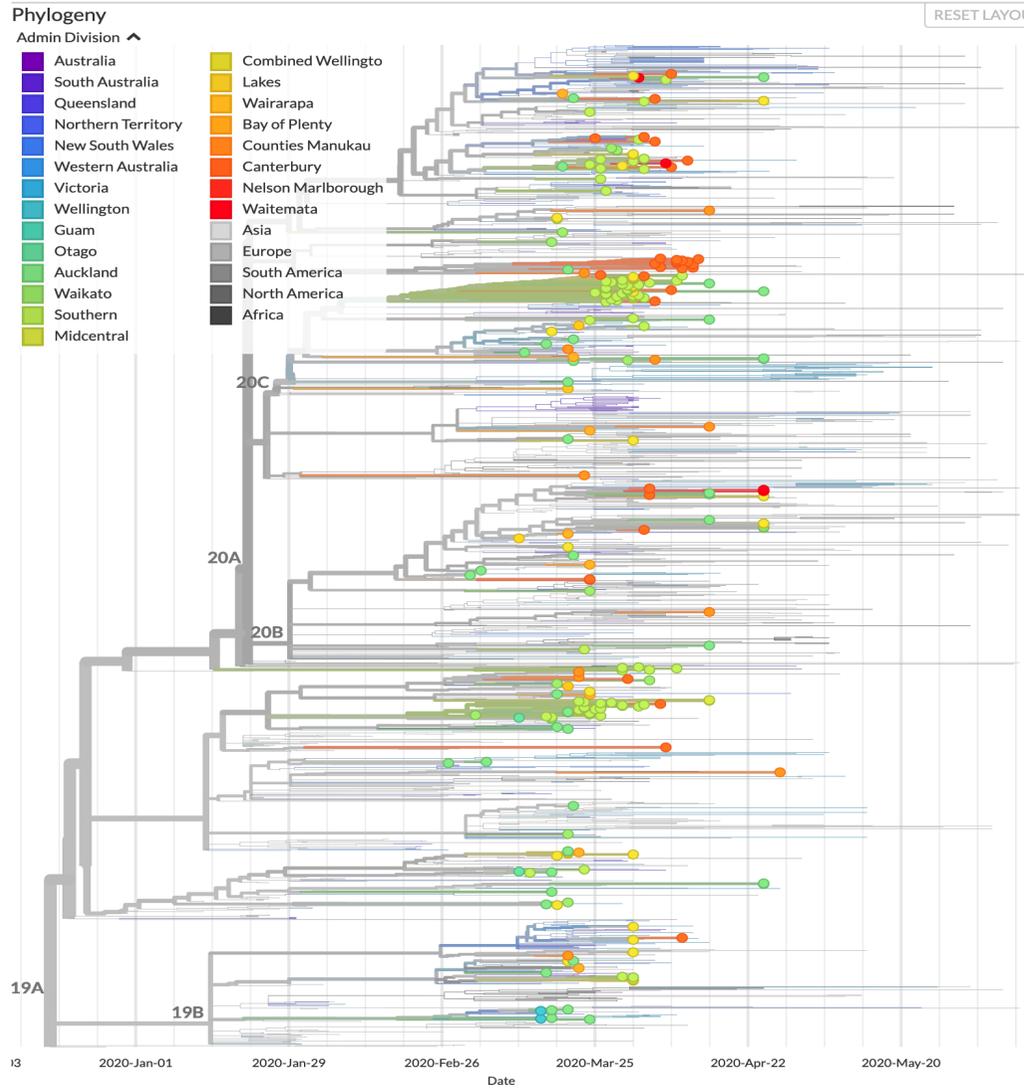
## Phylogenetic Relationships of Humans and Other Primates



# Phylodynamics of COVID-19

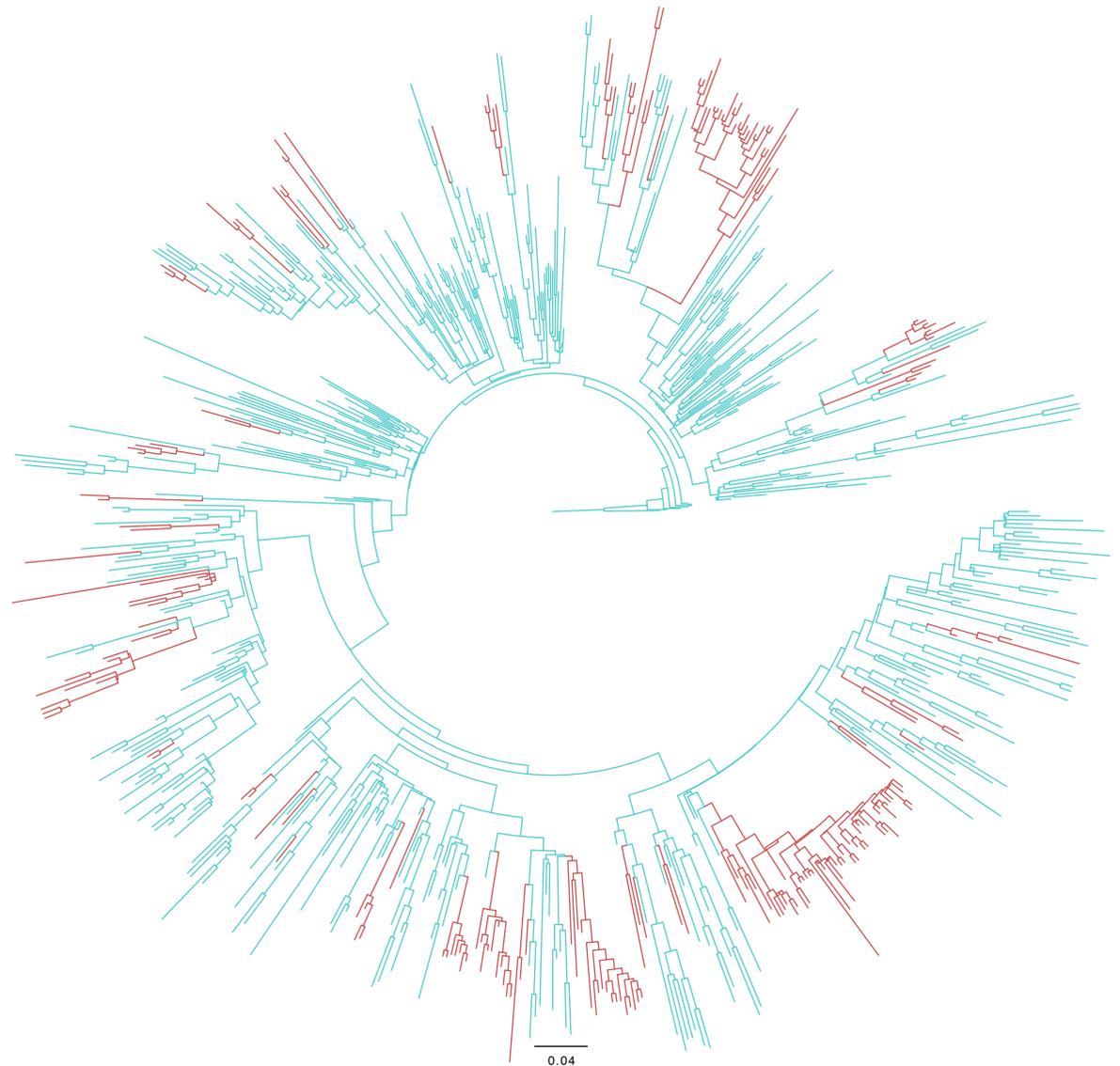
- Phylogenetics builds trees from genomes by modelling mutation processes, typically ignores population altogether.
- building a phylogenetic tree can tell us a lot about a virus: suggests local transmission chains, places sequences geographically, etc
- Phylodynamics accounts for the dynamics in the population of interest, incorporates information like
  - number of cases
  - Location and movement
  - Sampling dates

# Nextstrain.org provides an excellent global overview

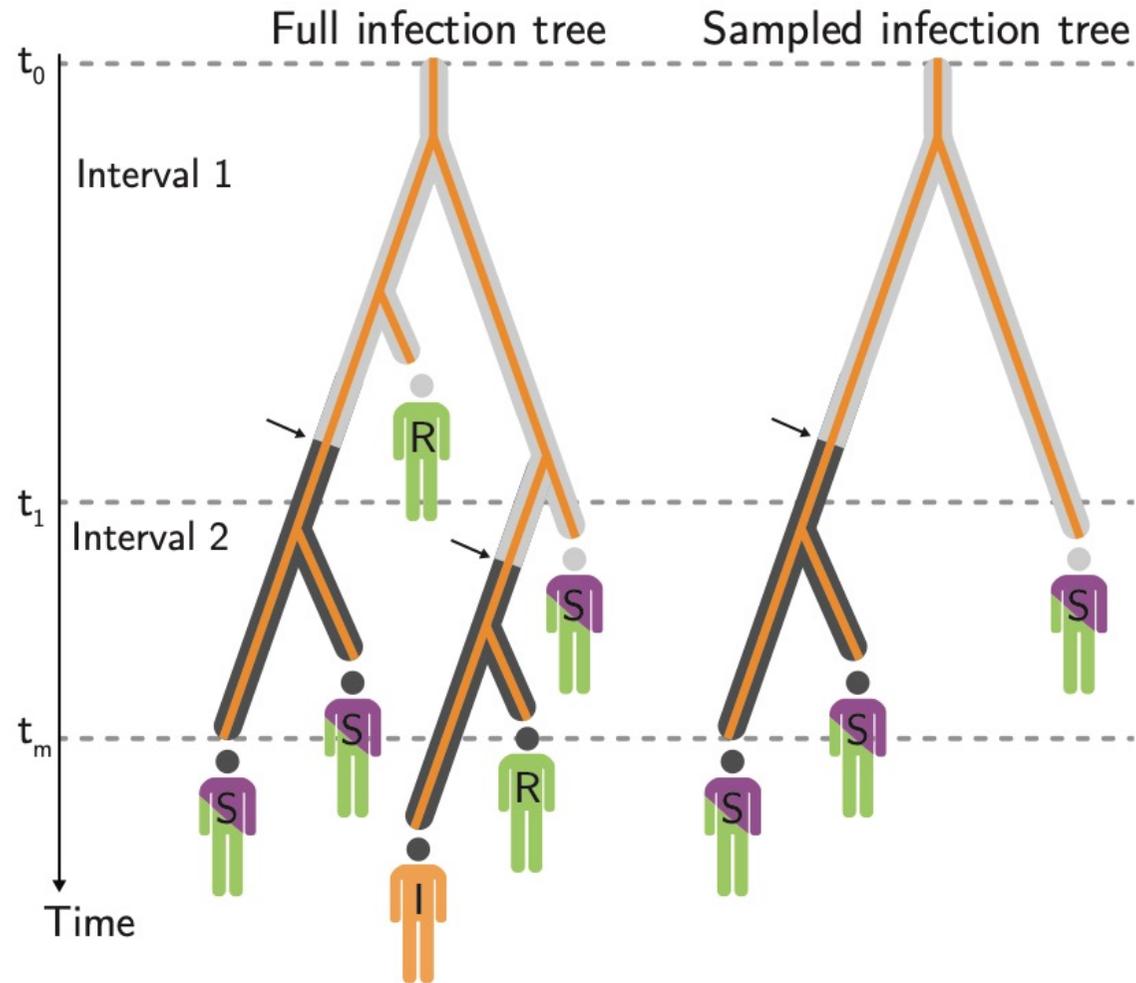


220 NZ  
samples, 500  
for rest of  
world, time  
increases  
towards  
leaves

location  
■ NZ  
■ RoW



# A basic phylodynamic model



Two locations: World and NZ

Parameters

Migration rates

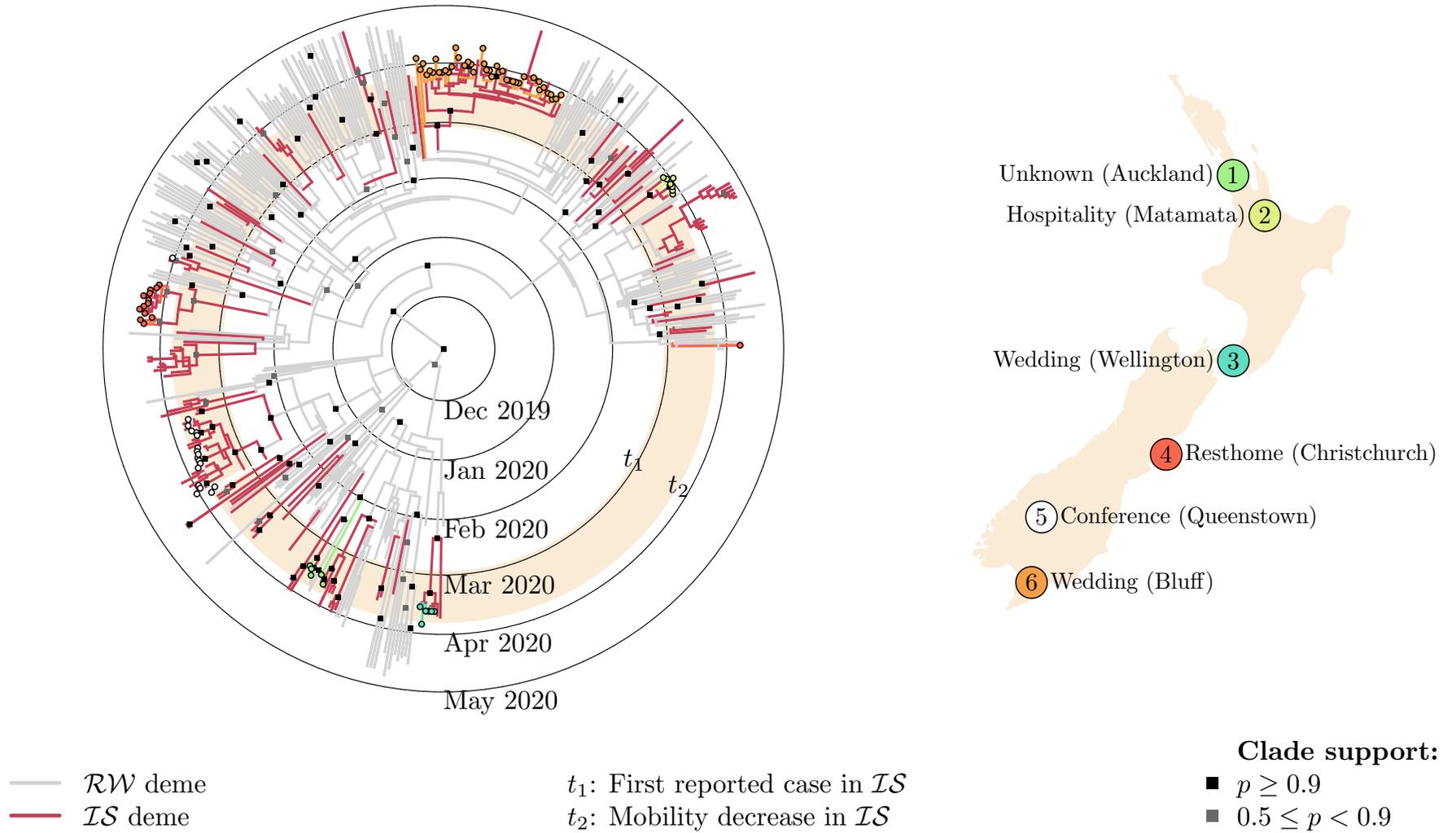
Transmission rates ("birth")

Removal rates ("death")

Sampling rates

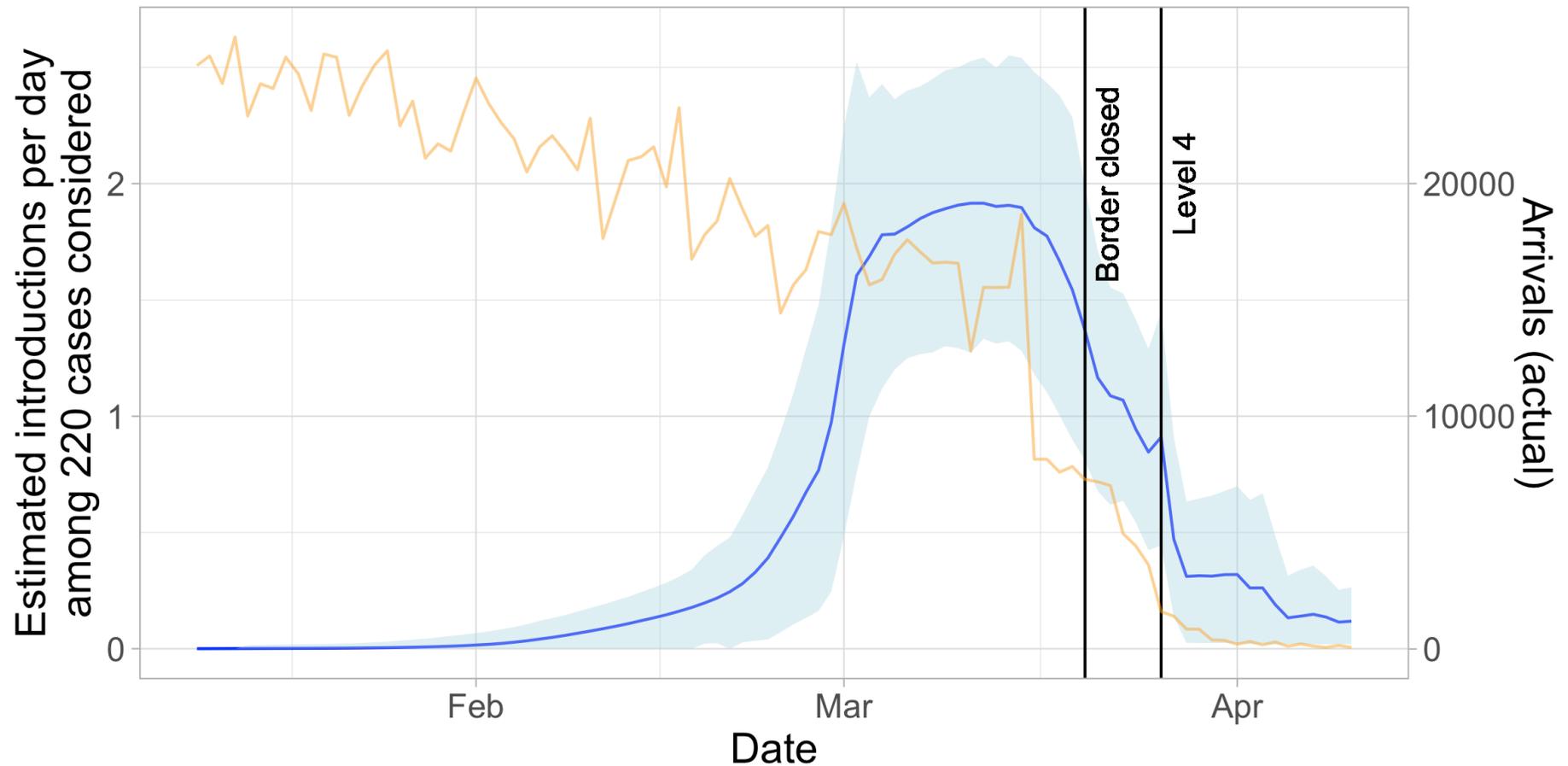
Parameters can differ in different intervals

# Major outbreaks in New Zealand

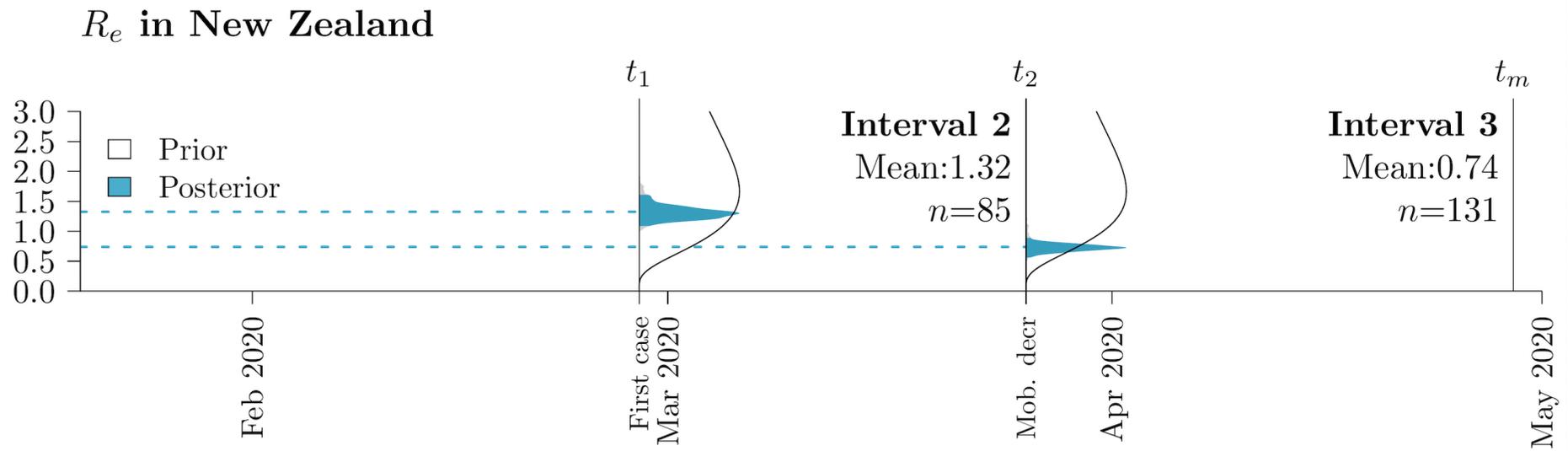




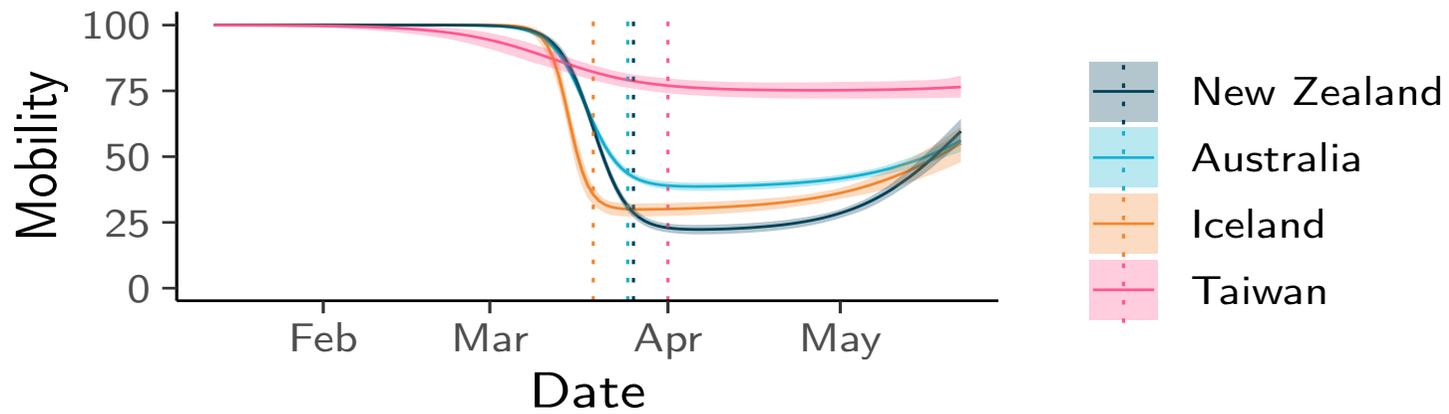
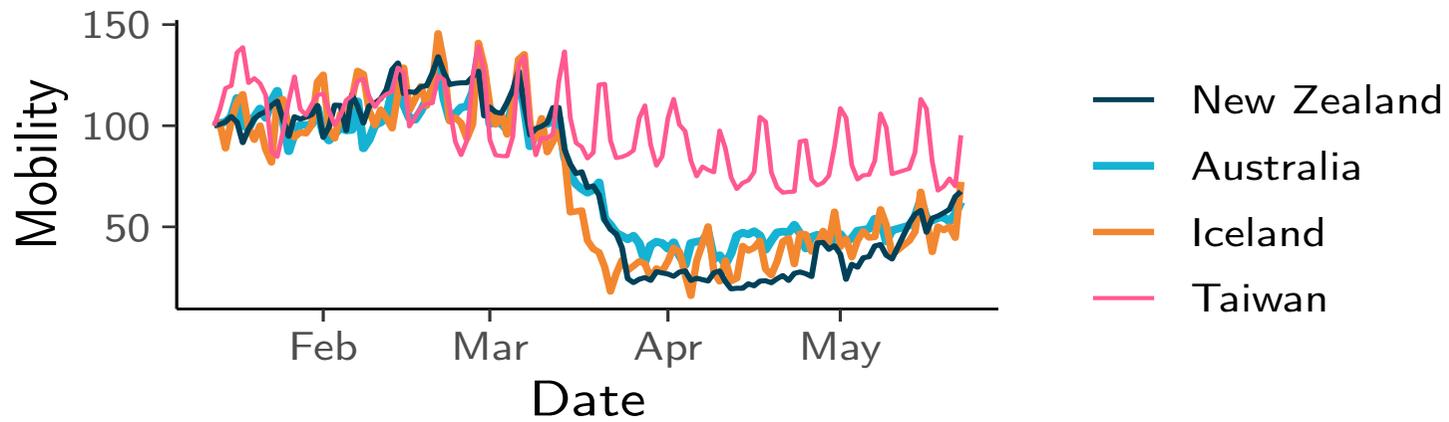
# When were cases introduced to New Zealand?



# Estimating the reproduction number



# Using mobility data to define intervals



# Challenges in computational biology

- More data being generated than we can handle...
- ...and data is noisy, often highly correlated
- High levels of cross disciplinary work required
- Models are easy to build and simulate but inference is hard
- Finding balance of simplicity and complexity is very hard
- How can we combine multiple sources of data?

# Links

- Murmuration: <https://www.youtube.com/watch?v=DmO4EIlgmd0>
- Boids: <https://www.youtube.com/watch?v=nbbd5uby0sY>
- Protein folding animation:  
<https://www.youtube.com/watch?v=jVOyaT56LEU>
- Foldit: <https://fold.it/portal/info/about>
- Multiple alignment game: <https://phylo.cs.mcgill.ca/play.php>
- Nextstrain COVID: <https://nextstrain.org/ncov/global>