



THE UNIVERSITY  
of EDINBURGH



Biotechnology and  
Biological Sciences  
Research Council



THE ROYAL  
SOCIETY

# Storing and analysing a million genomes on a desktop computer

Gregor Gorjanc, Jana Obsteter, Gabriela Mafra Fortuna, Roger Ros-Freixedes, Martin Johnsson, Ivan Pocrnic

InterBull & EAAP

Lyon, 2023-03-24



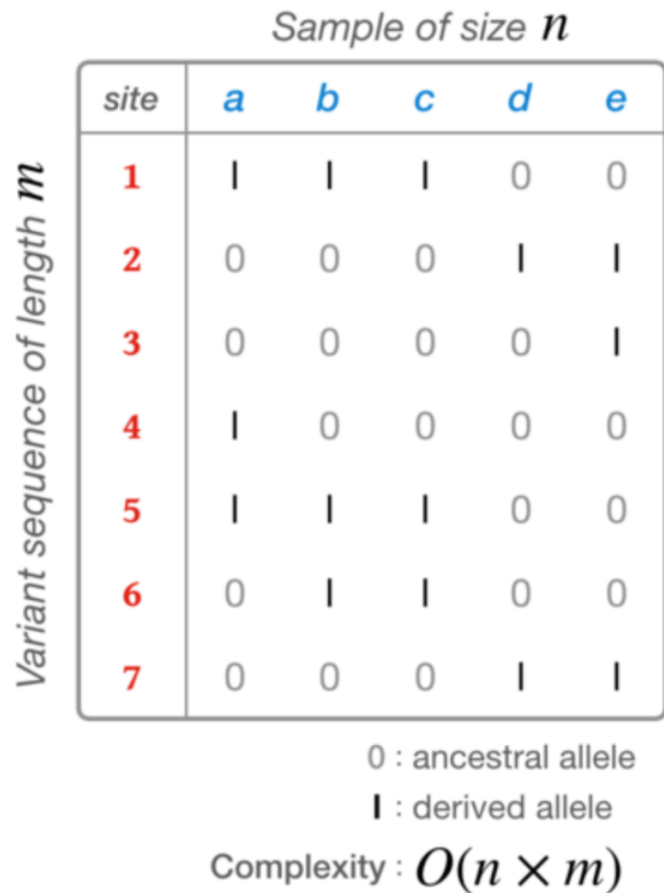
# Where have we been & where are we going?

- Started career with pedigree-based mixed models
- Rode on the excitement of introducing genomic selection
  - 2010s  $\sim 10^3$  individuals
  - 2015s  $\sim 10^5$  individuals      SNP array genotypes
  - 2020s  $\sim 10^6$  individuals
  - 2030s  $\sim 10^9$  individuals???
- Contributed to the whole-genome sequencing “craze”
  - 2015s  $\sim 10^3$  individuals
  - 2020s  $\sim 10^{4-6}$  individuals      Whole-genome sequences  
(with & without imputation)
  - 2030s  $\sim 10^{6-9}$  individuals???

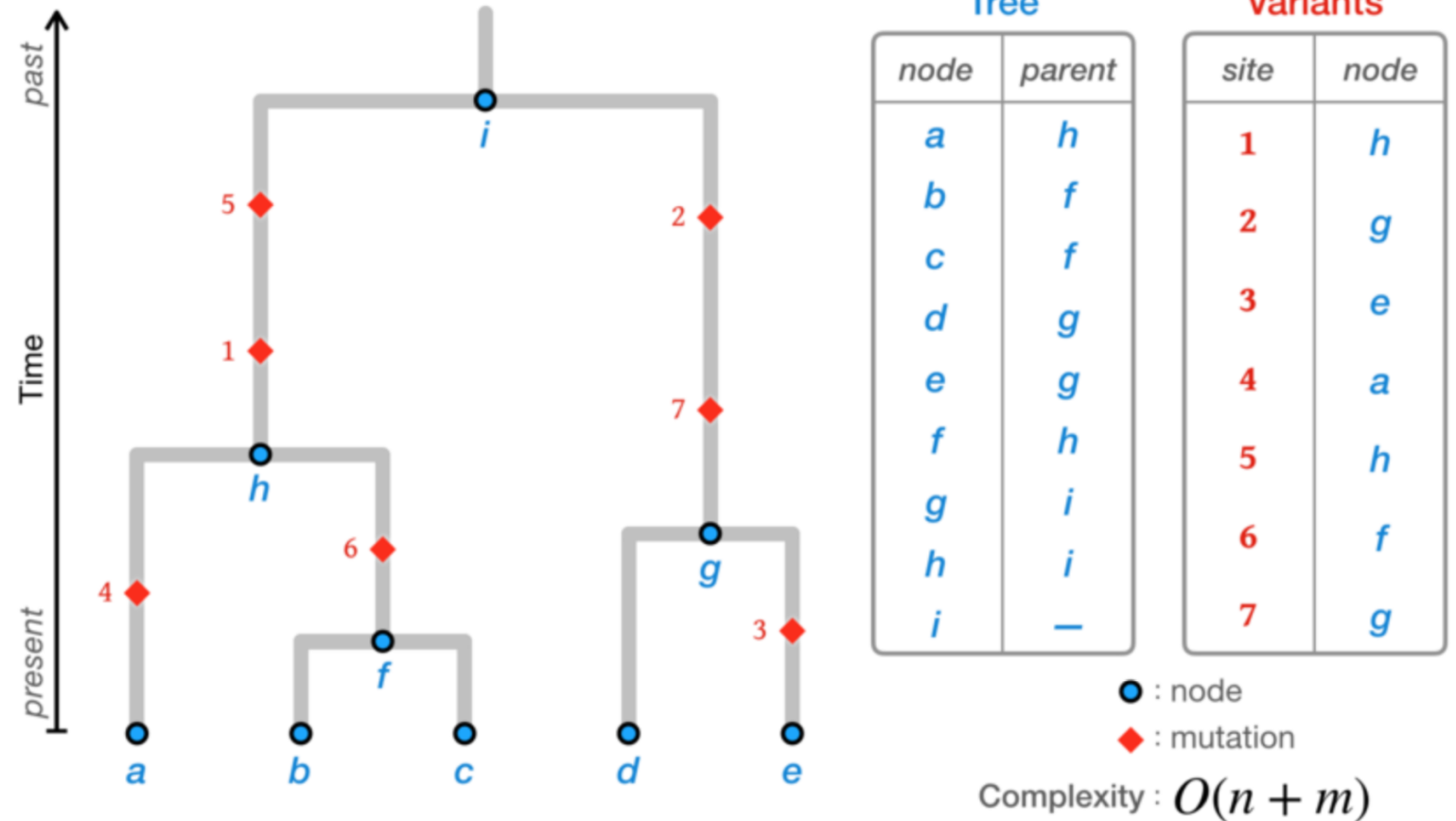
# Handling MEGA-SCALE through data generation process

Kelleher et al. (2019, Nature Genetics; ...)

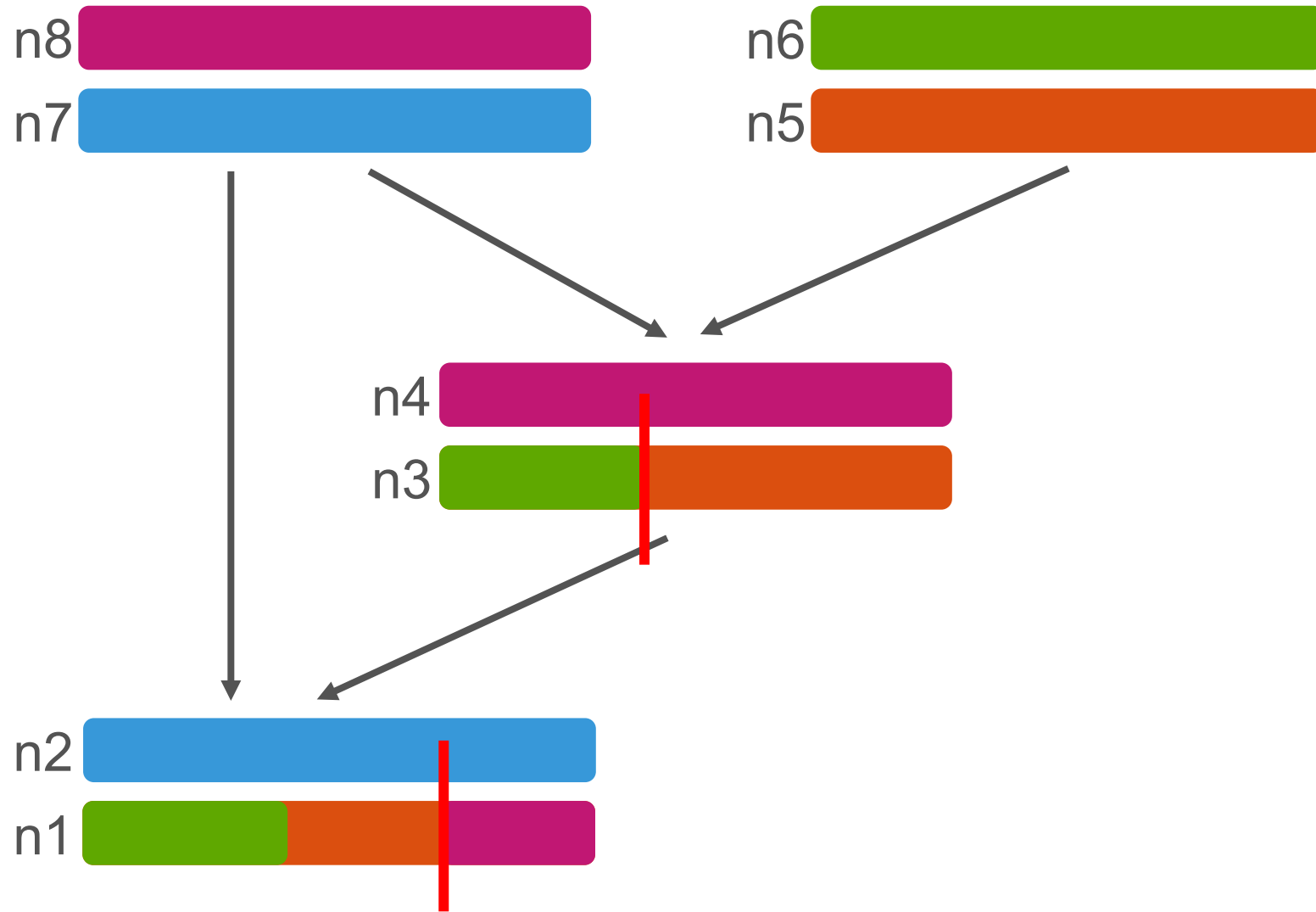
**A** Conventional data storage



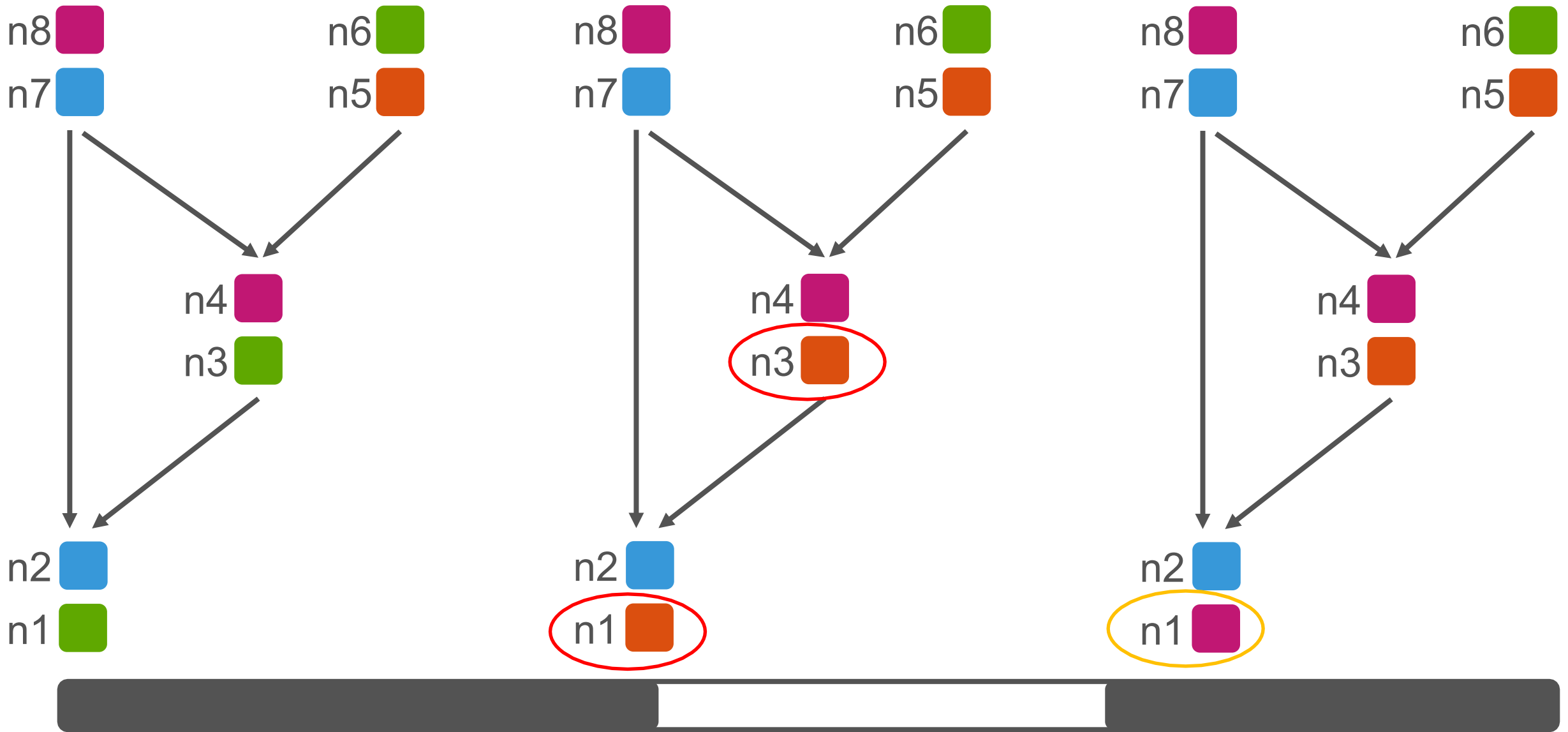
**B** Tree encoding



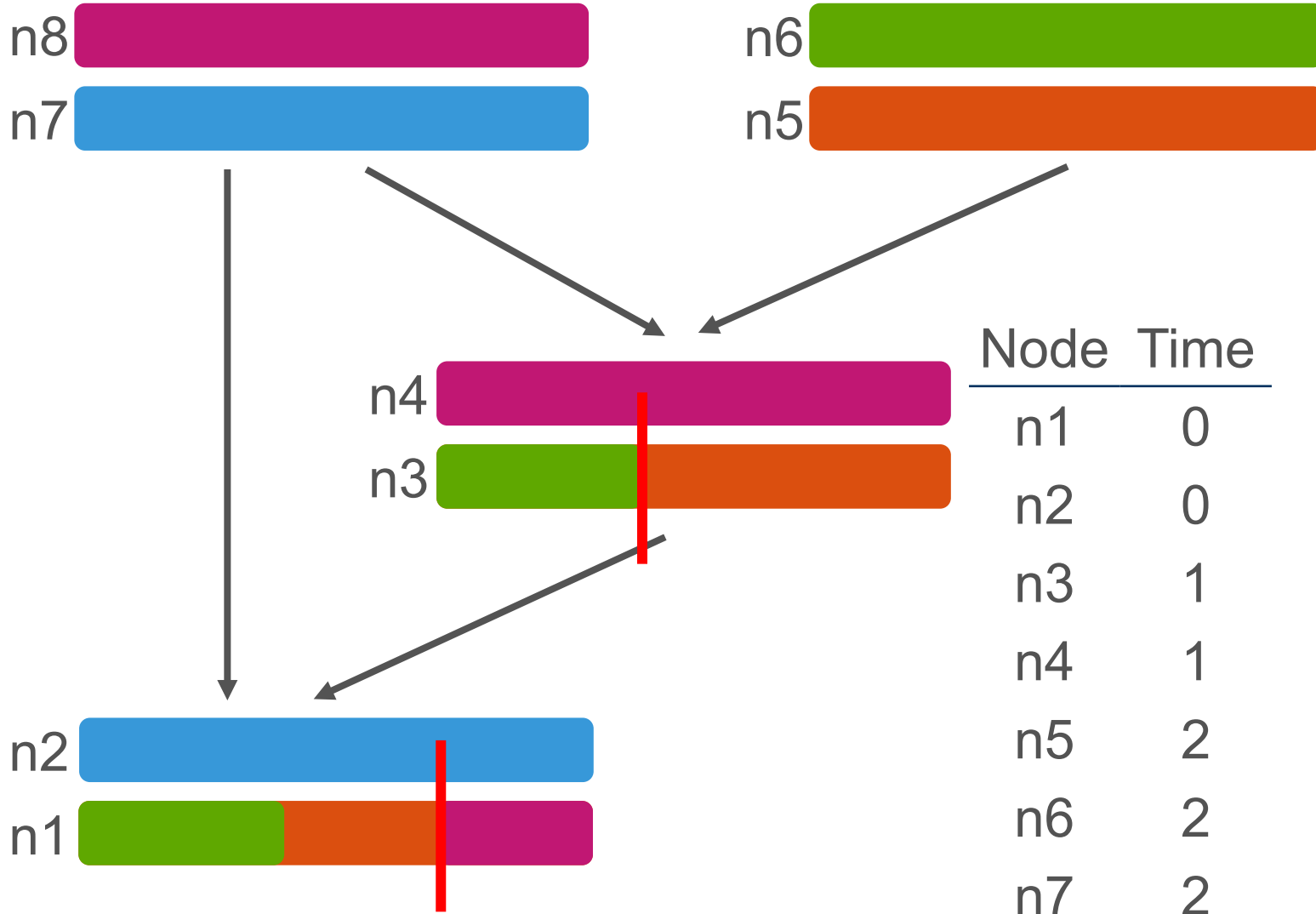
# Tracking chromosome segments



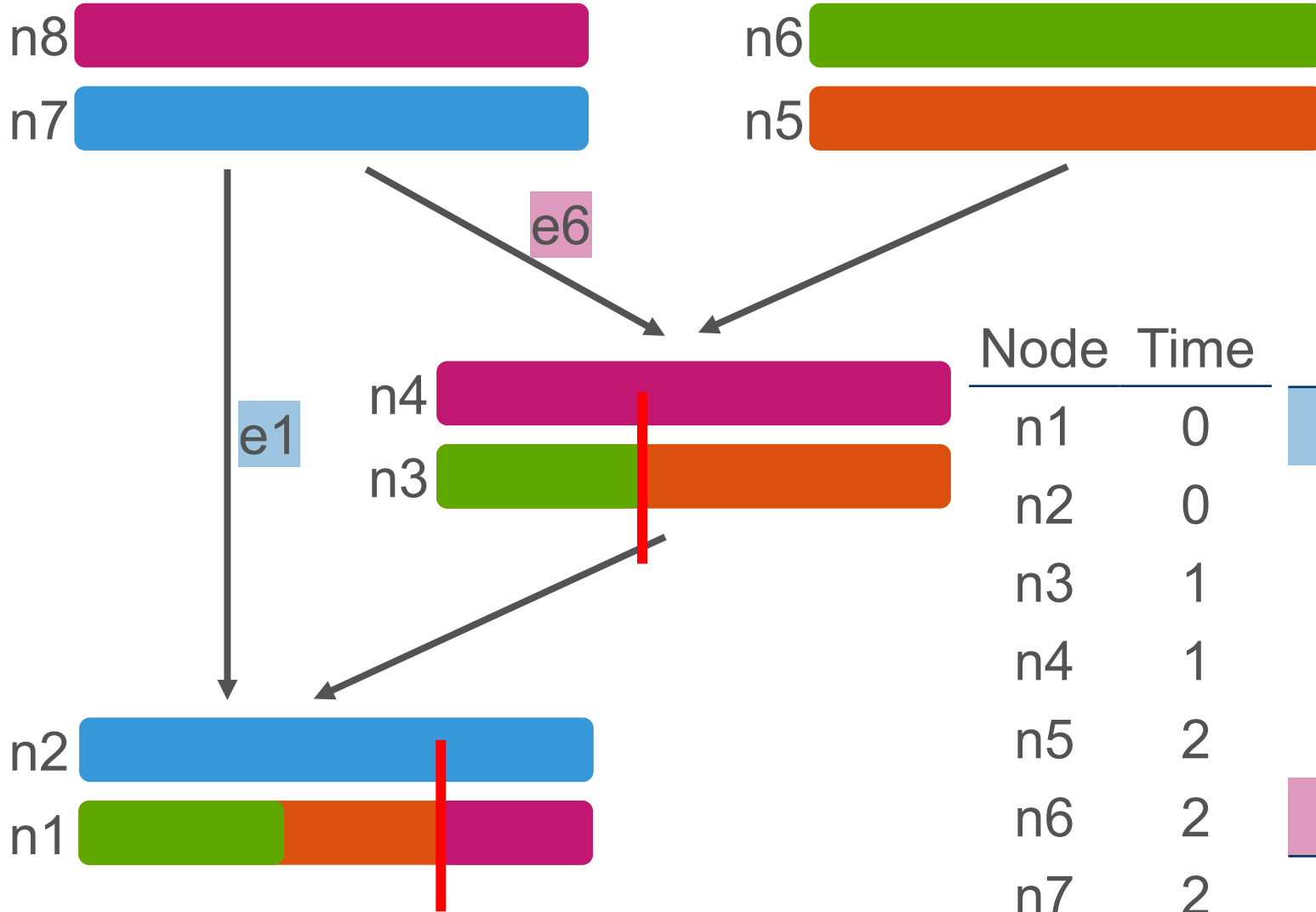
# Ancestral recombination graph & Local trees



# Tree sequence – Nodes

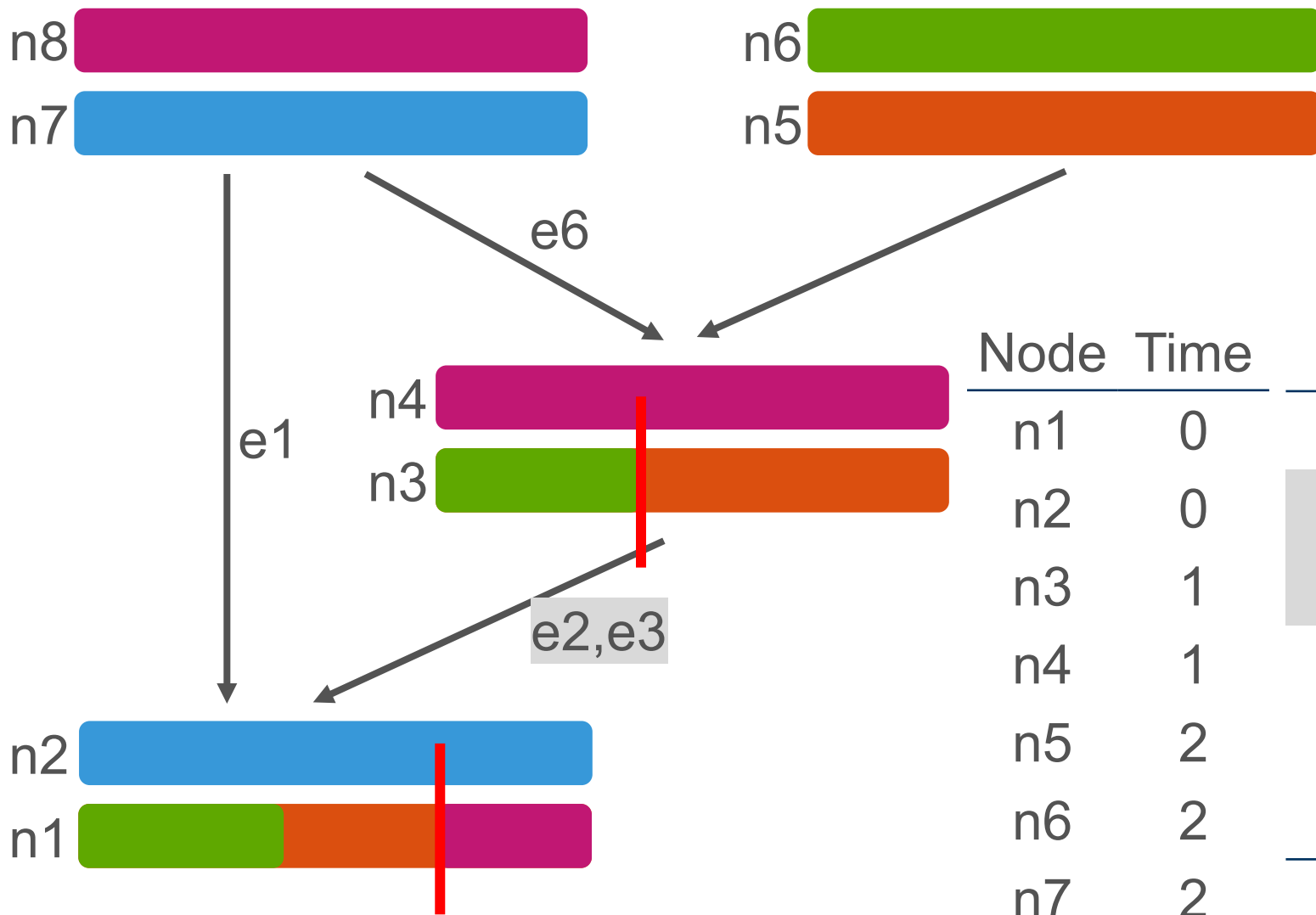


# Tree sequence – Nodes & Edges (pedigree)



Node	Time	Edge	Desc	Anc	Start	Stop
n1	0	e1	n2	n7	0	100
n2	0					
n3	1					
n4	1					
n5	2					
n6	2	e6	n4	n8	0	100
n7	2					
n8	2					

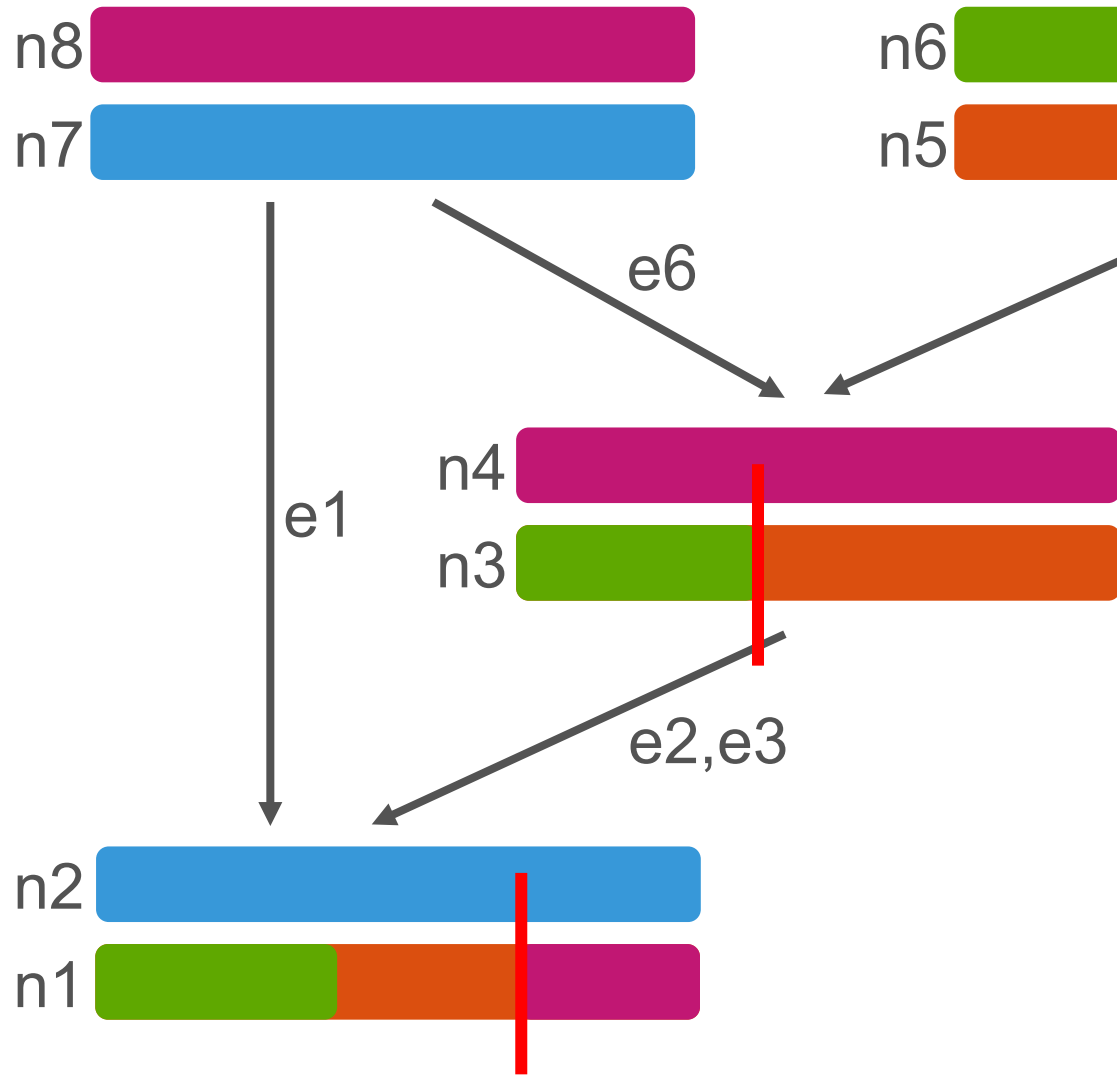
# Tree sequence – Nodes & Edges (pedigree)



Node	Time	Edge	Desc	Anc	Start	Stop
n1	0	e1	n2	n7	0	100
n2	0	e2	n1	n3	0	70
n3	1	e3	n1	n4	71	100
n4	1					
n5	2					
n6	2	e6	n4	n8	0	100
n7	2					
n8	2					



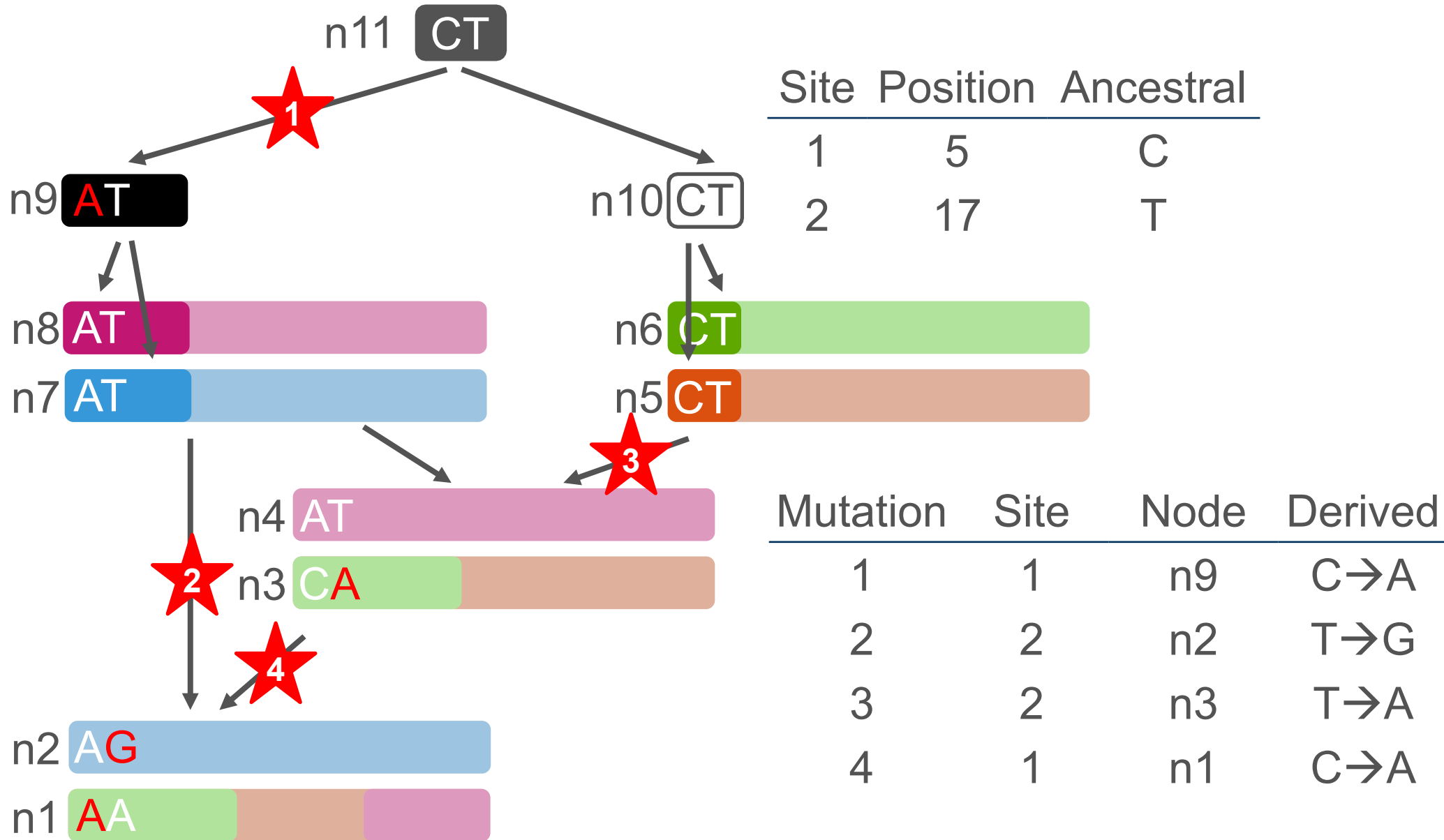
# Tree sequence – Nodes & Edges (pedigree)



Node	Time	Edge	Desc	Anc	Start	Stop
n1	0	e1	n2	n7	0	100
n2	0	e2	n1	n3	0	70
n3	1	e3	n1	n4	71	100
n4	1	e4	n3	n6	0	40
n5	2	e5	n3	n5	41	100
n6	2	e6	n4	n8	0	100
n7	2					
n8	2					



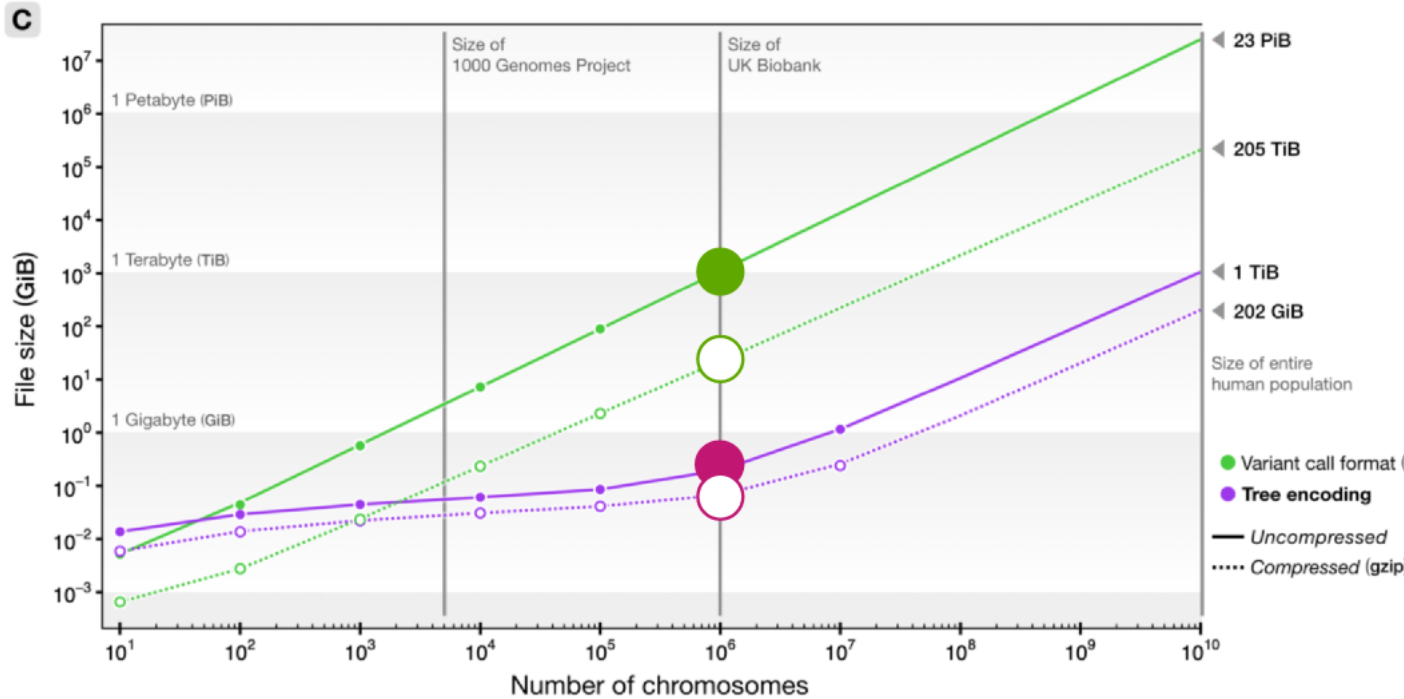
# Tree sequence – Nodes, Edges, Sites, & Mutations



# Computational power of tree sequences

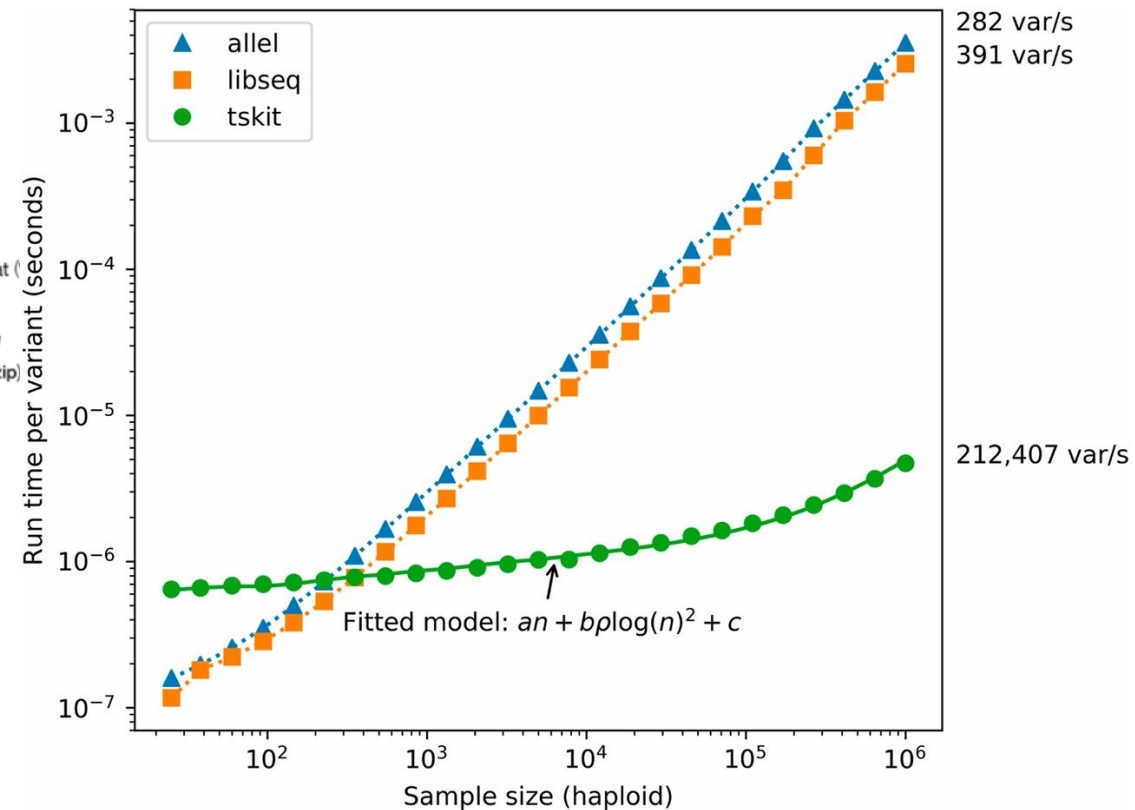
## STORAGE

Kelleher et al. (2019, Nature Genetics; ...)



## COMPUTE

Ralph et al. (2020, Genetics)



# Roslin-Genus/PIC ~1 million pig genomes project

- Pedigree & SNP array data
  - 9 lines with a total of ~450K pigs
  - ~15K-50K markers
- Whole-genome sequence
  - ~8K pigs (a mix of ~1x and ~30x)
  - ~46M variants passed quality control across lines
- Accurate imputation of whole-genomes
  - ~450K diploid pigs \* 2 = **~900,000 haploid genomes**
  - ~450K pigs \* ~46M sites \* 8 bytes /  $2^{40}$  = ~152 TiB of memory ☹️  
(2-bit storage → ~5TiB of memory)

# Roslin-Genus/PIC ~1 million pig genomes project

- Pedigree & SNP array data
  - 9 lines with a total of ~450K pigs
  - ~15K-50K markers

- Whole-genome sequence

- ~8K pigs (a mix of ~1x and ~30x)
- ~46M variants passed quality control across lines

Tables				Nodes	
Edges				ID	time
left	right	parent	child	0	0.0
0	20	4	0	1	0.0
0	20	4	1	2	0.0
0	10	4	2	3	0.0
0	10	5	3	4	2.0
0	10	5	4	5	3.0
10	20	4	6	6	1.0
10	20	6	2		
10	20	6	3		

Mutations			Sites			
ID	site	node	derived	ID	position	ancestral
0	1	3	T	0	2	C
1	2	2	G	1	4	A
2	4	4	T	2	5	C
3	6	6	G	3	7	G
4	8	2	T	4	8	C
				5	9	T
				6	12	T
				7	15	C
				8	18	G
				9	19	C

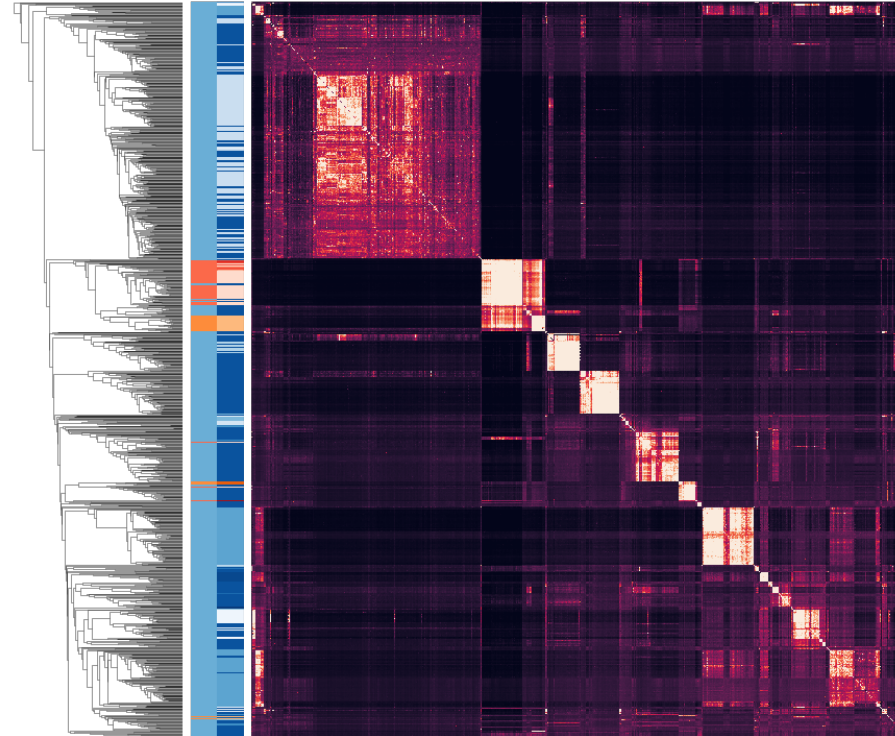
**~15+ GiB vs. ~152 TiB**  
**~99.9% “compression”!!!**

- Accurate imputation of whole-genomes

- ~450K diploid pigs \* 2 = **~900,000 haploid genomes**
- ~450K pigs \* ~46M sites \* 8 bytes / 2<sup>40</sup> = ~152 TiB of memory ☹️  
(2-bit storage → ~5TiB of memory)

# Real application: 1000 Bull Genomes data

- 2,716 samples & 157 groups
  - Bos taurus
  - Bos indicus
  - Crossbred & African
  - Bos taurus coreanae (nat. Korean)
  - Bos primigenius (auroch)
  - Bos grunniens (yak)
- 29 autosomal chromosomes with ~116M variants
- Shapelt (phase), tsinfer (infer tree sequence), & tskit (analyse)
- **VCFs: HUGE SIZE** → **Tree sequences: DECENT SIZE** → **"Compression": LARGE** → **Analysis: FAST**



Session 83  
Genetic diversity  
Thursday  
~ 16h 30m

# Conclusion

- MEGA-SCALE genomic datasets are here & growing
- Tree sequence data format to the rescue!?
  - PROS
    - Succinctly encodes the inheritance process
    - Combines pedigree, coalescent, phylogenetics (gene & species trees), segregation, recombination, gene conversion, mutations, IBS, IBD, ...
    - Significant storage reduction & fast analyses
    - Novel insights & modelling
  - CONS (on-going work)
    - Novel way of thinking
    - Ancestral alleles, inference from real data, and inputs still HUGE
    - Need to develop specialised algorithms



Join us @



THE UNIVERSITY  
of EDINBURGH



4-year PhD studentship available (start in autumn 2024)!!!







```
for (Year in 1:10) {
  Pop = randCross2(males = Sires,
                  females = Dams,
                  nCrosses = 750,
                  nProgeny = 100)

  Dams = selectInd(Pop,
                  nInd = 750,
                  sex = "F")

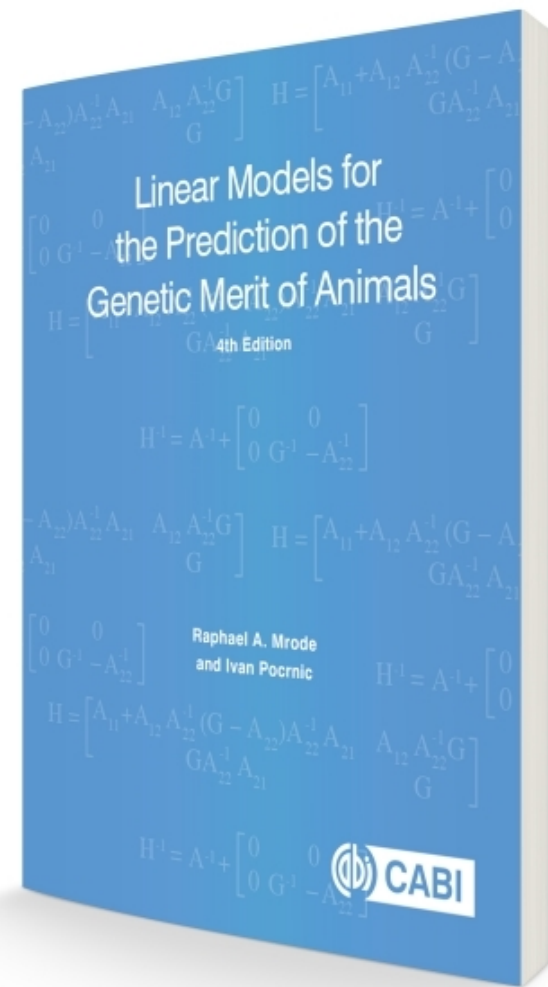
  Sires = selectInd(Pop,
                  nInd = 25,
                  sex = "M")
}
```



```
for (Year in 1:10) {
  Variety = selectInd(EYT, nInd = 1)
  EYT = selectInd(AYT, nInd = 10)
  AYT = selectInd(PYT, nInd = 50)
  PYT = selectInd(HDRW, nInd = 500)
  HDRW = madeDH(F1, nDH = 100)
  Parents = c(EYT, AYT)
  F1 = randCross(Parents, nCrosses = 100)
}
```

Free short online course

# Breeding Programme Modelling with AlphaSimR





THE UNIVERSITY  
of EDINBURGH



Biotechnology and  
Biological Sciences  
Research Council



THE ROYAL  
SOCIETY

# Storing and analysing a million genomes on a desktop computer

Gregor Gorjanc, Jana Obsteter, Gabriela Mafra Fortuna, Roger Ros-Freixedes, Martin Johnsson, Ivan Pocrnica

InterBull & EAAP

Lyon, 2023-03-24

