

# What the heck is an SFS (site frequency spectrum)?

	Site 1
Sample 1	AA
Sample 2	AG
Sample 3	AA
Sample 4	AA
Sample 5	AA
Sample 6	AA
Sample 7	AA
Sample 8	AA

# What the heck is an SFS (site frequency spectrum)?

	Site 1
Sample 1	AA
Sample 2	AG
Sample 3	AA
Sample 4	AA
Sample 5	AA
Sample 6	AA
Sample 7	AA
Sample 8	AA

1

	Site 1	Site 2
Sample 1	AA	CC
Sample 2	AG	CC
Sample 3	AA	CC
Sample 4	AA	CC
Sample 5	AA	CC
Sample 6	AA	CC
Sample 7	AA	CC
Sample 8	AA	CC
	1	0

	Site 1	Site 2	Site 3	Site 4	Site 5	Site 6	Site 7	Site 8	Site 9	Site 10
Sample 1	AA	CC	TT	AA	AG	TT	TT	CC	TT	AA
Sample 2	AG	CC	CC	AA	AA	TT	TT	CC	TT	AA
Sample 3	AA	CC	TT	AA	AG	TT	TT	CC	TT	AA
Sample 4	AA	CC	TT	AA	AA	TT	TT	CC	TT	GG
Sample 5	AA	CC	CT	AA	AA	TT	TT	CC	TT	AA
Sample 6	AA	CC	TT	AA	AA	TT	TT	CC	TT	AA
Sample 7	AA	CC	TT	AA	AA	TT	CT	CC	TT	AA
Sample 8	AA	CC	TT	AA	AA	TT	TT	CC	TT	AA
	1	0	3	0	2	0	1	0	0	2

# A “folded” SFS summarizes the minor allele count (MAC)

Site 1	Site 2	Site 3	Site 4	Site 5	Site 6	Site 7	Site 8	Site 9	Site 10
1	0	3	0	2	0	1	0	0	2

# A “folded” SFS summarizes the minor allele count (MAC)

MAC	0	1	2	3	4	5	6	7	8
Count	5	2	2	1	0	0	0	0	0

Site 1  
1

Site 2  
0

Site 3  
3

Site 4  
0

Site 5  
2

Site 6  
0

Site 7  
1

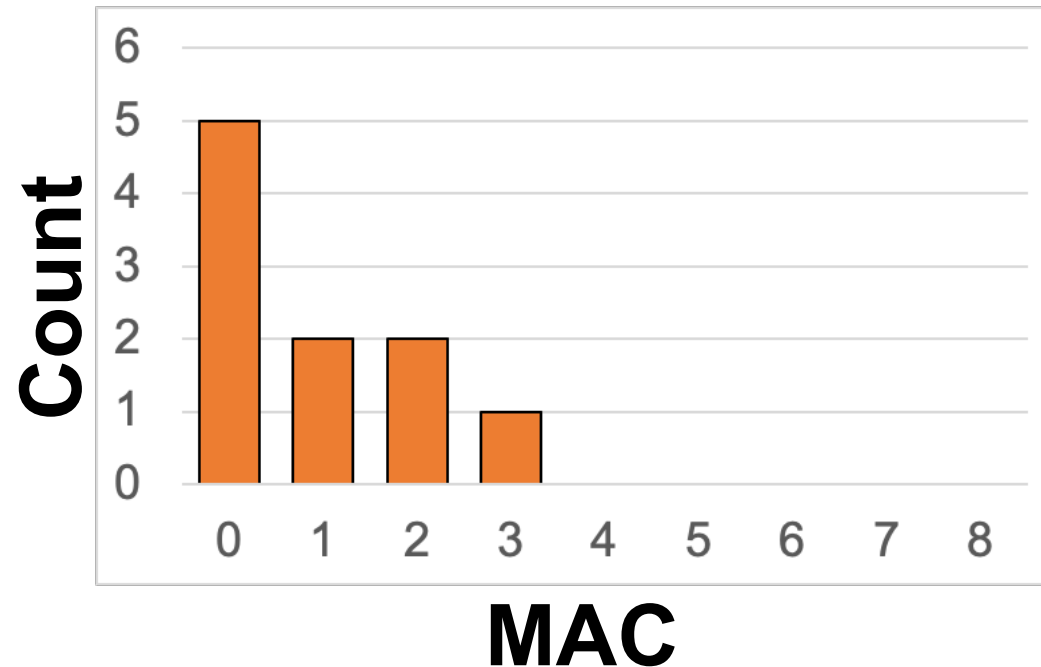
Site 8  
0

Site 9  
0

Site 10  
2

# A “folded” SFS summarizes the minor allele count (MAC)

MAC	0	1	2	3	4	5	6	7	8
Count	5	2	2	1	0	0	0	0	0



**An “unfolded” SFS summarizes  
the derived allele count**



Sample 1	AA	CC	TT	AA	AG	TT	TT	CC	TT	AA
Sample 2	AG	CC	CC	AA	AA	TT	TT	CC	TT	AA
Sample 3	AA	CC	TT	AA	AG	TT	TT	CC	TT	AA
Sample 4	AA	CC	TT	AA	AA	TT	TT	CC	TT	GG
Sample 5	AA	CC	CT	AA	AA	TT	TT	CC	TT	AA
Sample 6	AA	CC	TT	AA	AA	TT	TT	CC	TT	AA
Sample 7	AA	CC	TT	AA	AA	TT	CT	CC	TT	AA
Sample 8	AA	CC	TT	AA	AA	TT	TT	CC	TT	AA
	1	0	3	0	2	0	1	0	0	2

	A	C	C	A	A	C	T	C	T	A/G
Sample 1	AA	CC	TT	AA	AG	TT	TT	CC	TT	AA
Sample 2	AG	CC	CC	AA	AA	TT	TT	CC	TT	AA
Sample 3	AA	CC	TT	AA	AG	TT	TT	CC	TT	AA
Sample 4	AA	CC	TT	AA	AA	TT	TT	CC	TT	GG
Sample 5	AA	CC	CT	AA	AA	TT	TT	CC	TT	AA
Sample 6	AA	CC	TT	AA	AA	TT	TT	CC	TT	AA
Sample 7	AA	CC	TT	AA	AA	TT	CT	CC	TT	AA
Sample 8	AA	CC	TT	AA	AA	TT	TT	CC	TT	AA
	1	0	3	0	2	0	1	0	0	2

	A	C	C	A	A	C	T	C	T	A/G
Sample 1	AA	CC	TT	AA	AG	TT	TT	CC	TT	AA
Sample 2	AG	CC	CC	AA	AA	TT	TT	CC	TT	AA
Sample 3	AA	CC	TT	AA	AG	TT	TT	CC	TT	AA
Sample 4	AA	CC	TT	AA	AA	TT	TT	CC	TT	GG
Sample 5	AA	CC	CT	AA	AA	TT	TT	CC	TT	AA
Sample 6	AA	CC	TT	AA	AA	TT	TT	CC	TT	AA
Sample 7	AA	CC	TT	AA	AA	TT	CT	CC	TT	AA
Sample 8	AA	CC	TT	AA	AA	TT	TT	CC	TT	AA
	1	0	13	0	2	16	1	0	0	—

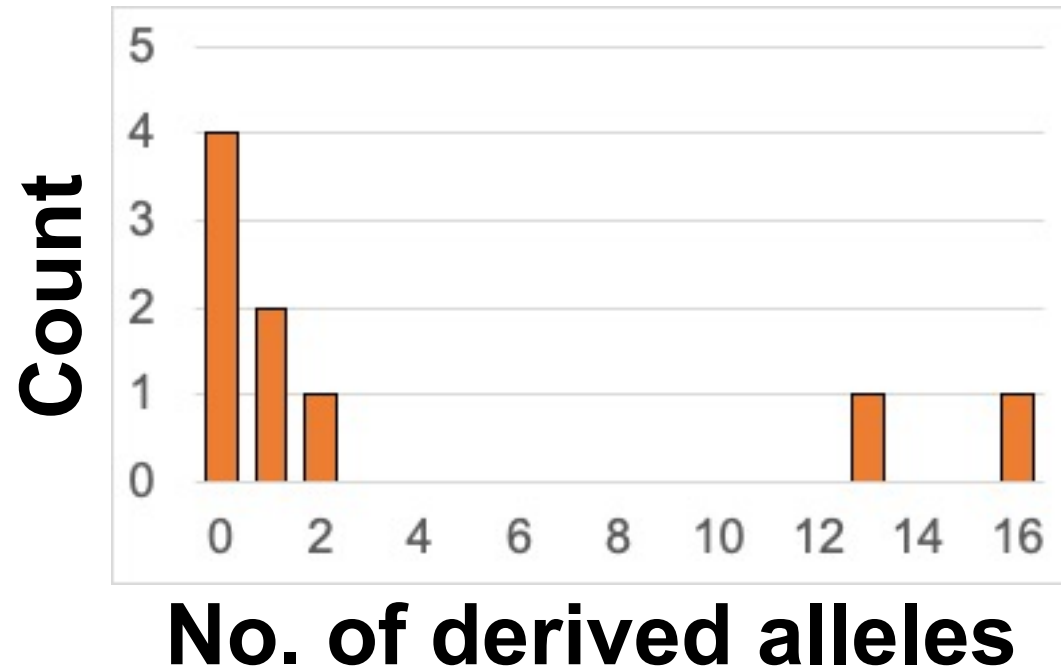
# An “unfolded” SFS summarizes the derived allele count

MAC	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
Count	4	2	1	0	0	0	0	0	0	0	0	0	0	1	0	0	1

Site 1	Site 2	Site 3	Site 4	Site 5	Site 6	Site 7	Site 8	Site 9	Site 10
1	0	13	0	2	16	1	0	0	—

# An “unfolded” SFS summarizes the derived allele count

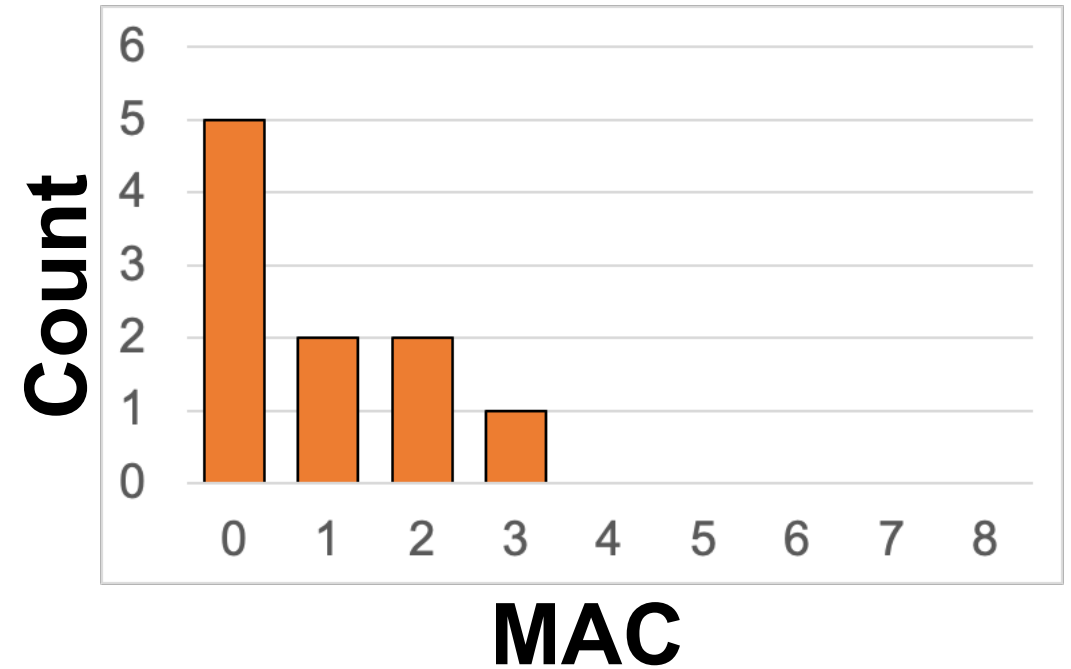
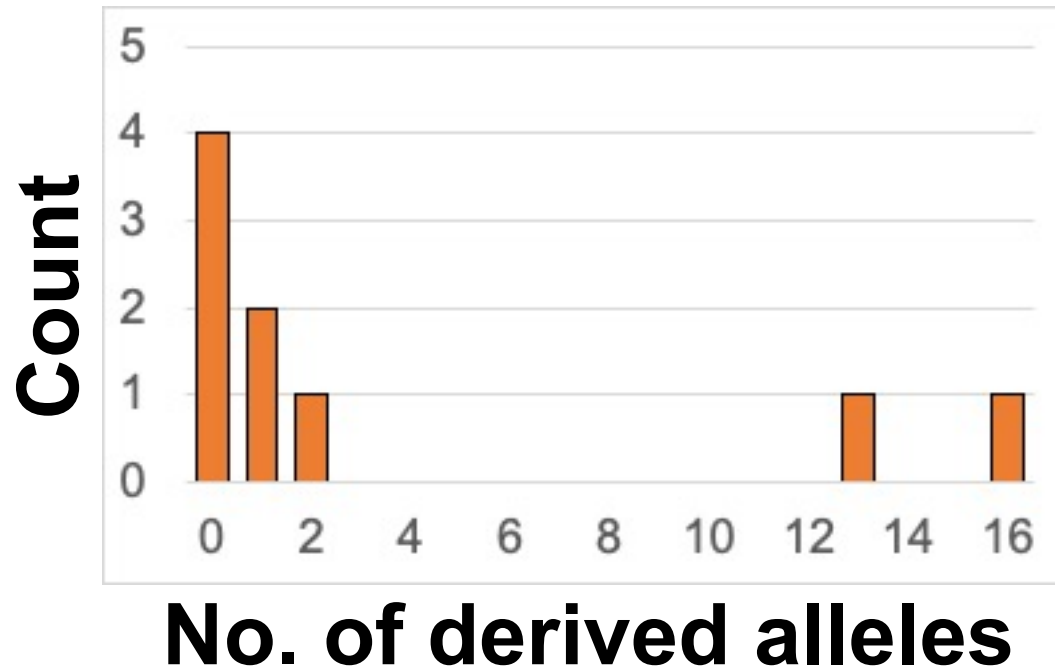
MAC	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
Count	4	2	1	0	0	0	0	0	0	0	0	0	0	1	0	0	1



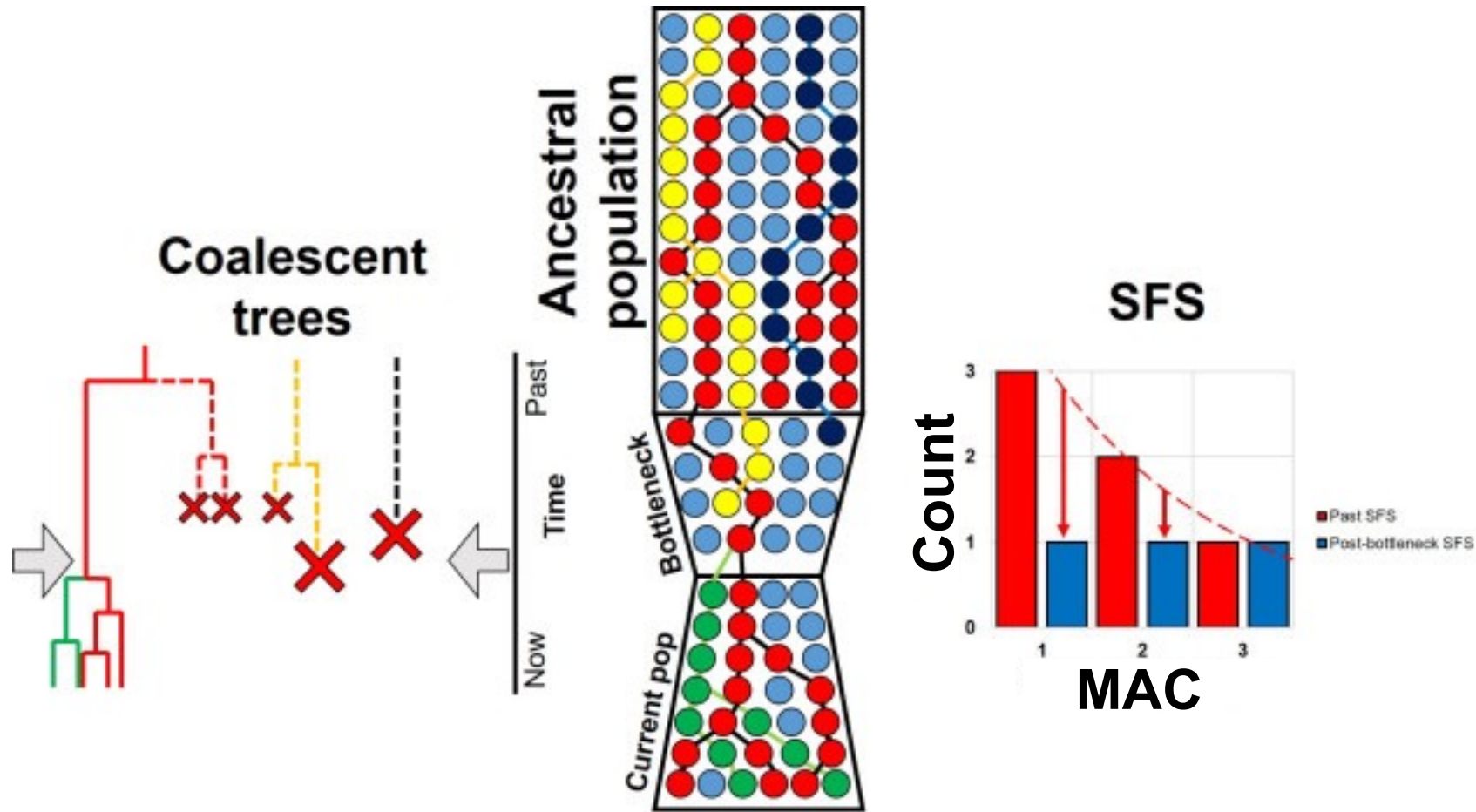
# Unfolded vs folded SFS

- $2N$  entries
- More information (more accurate inferences)
- Requires data from outgroup to estimate ancestral state

- $N$  entries
- Less information (less accurate inferences)
- Only requires data from the species of interest

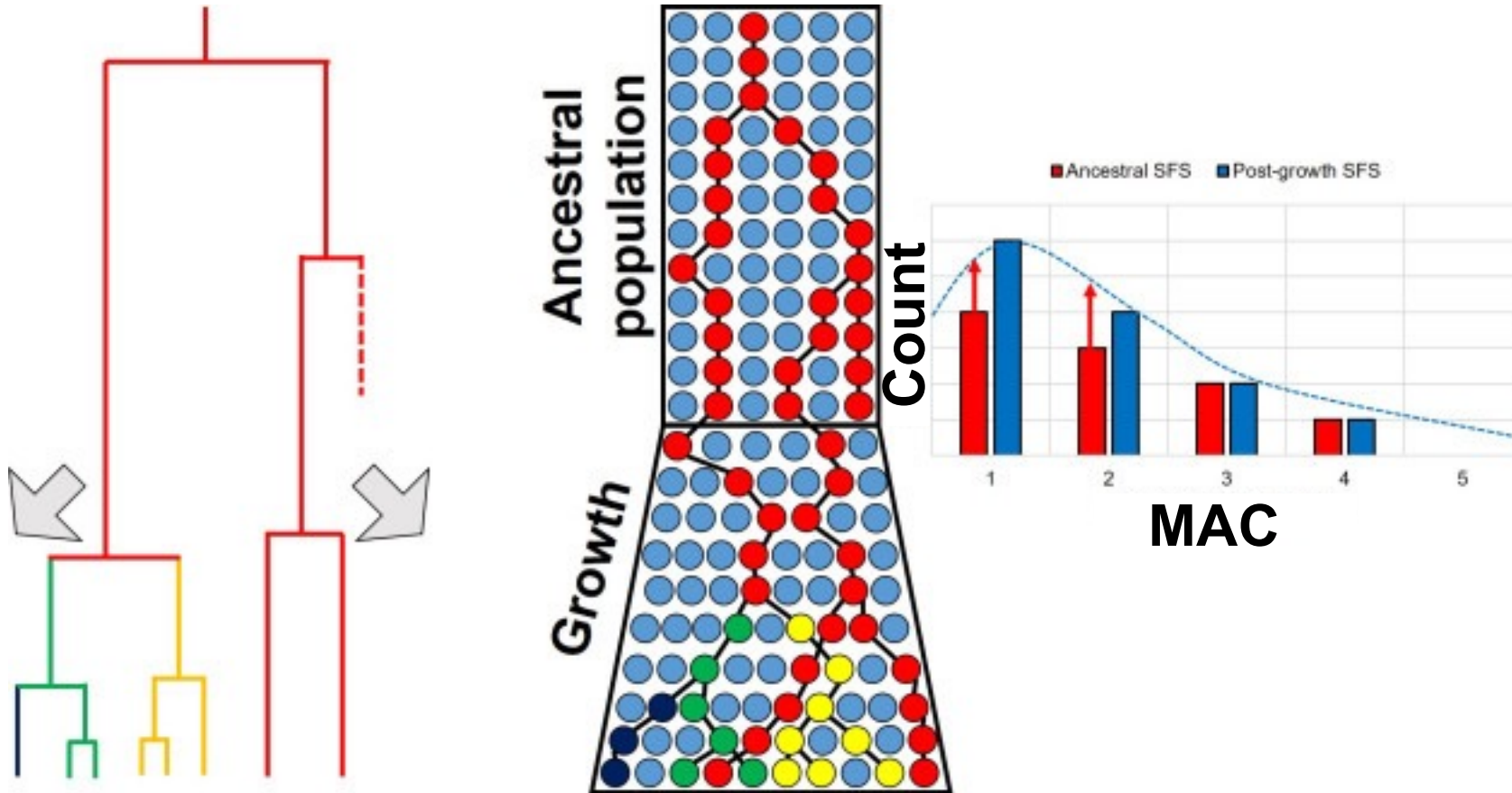


# Demography impacts the SFS



A bottleneck “evens” out the counts of different minor allele counts (e.g. reduces the number of singletons)

# Demography impacts the SFS



Growth/an expansion leads to an excess of singleton sites (a more uneven site frequency spectrum)



# Can use the SFS to infer demography

Bunch of approaches for doing  
this, including:

- Fastsimcoal2
- Dadi
- CubSFS

# Can use the SFS to infer demography

Bunch of approaches for doing this, including:

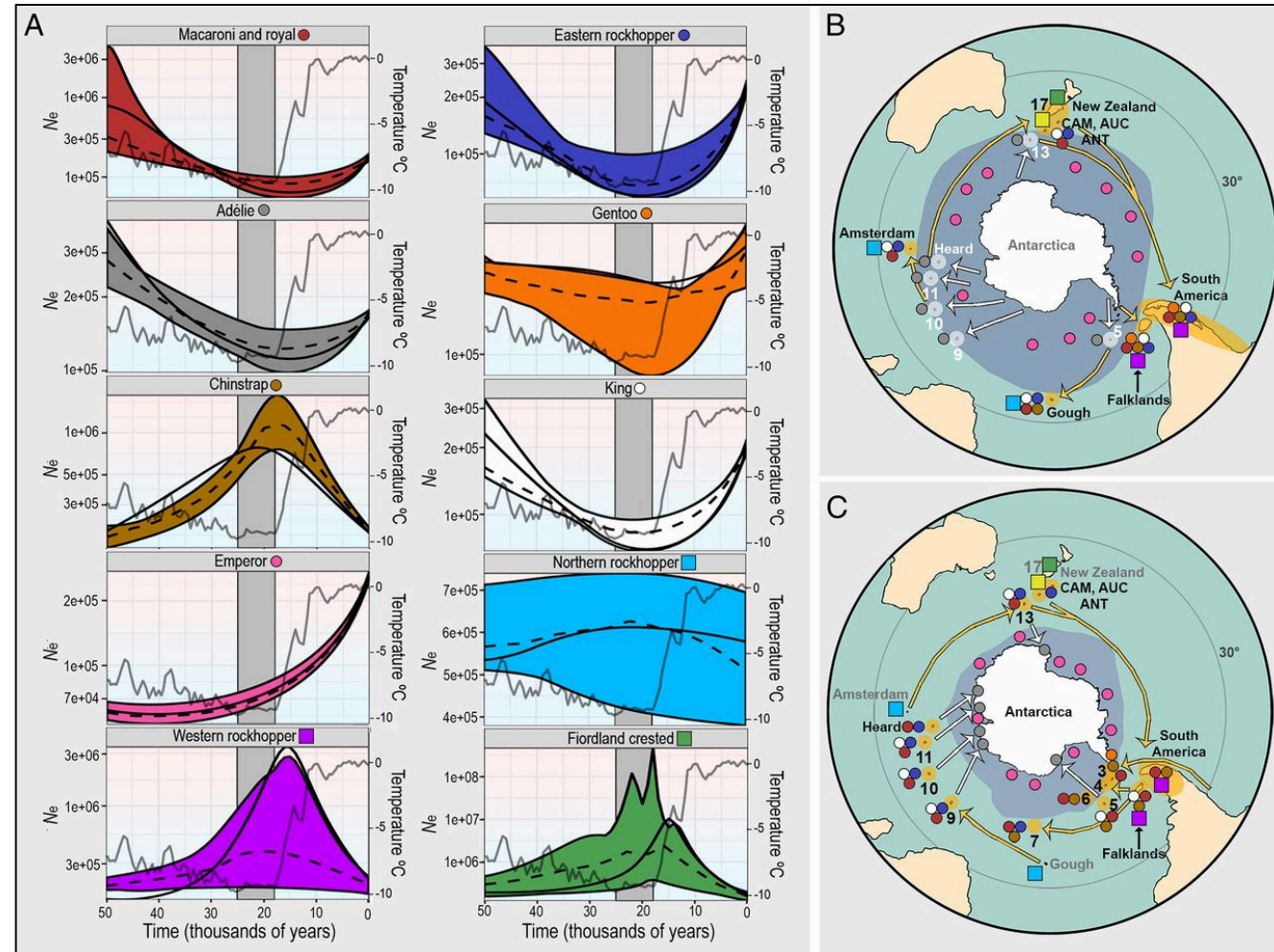
- Fastsimcoal2
- Dadi
- **CubSFS**

## Receding ice drove parallel expansions in Southern Ocean penguins

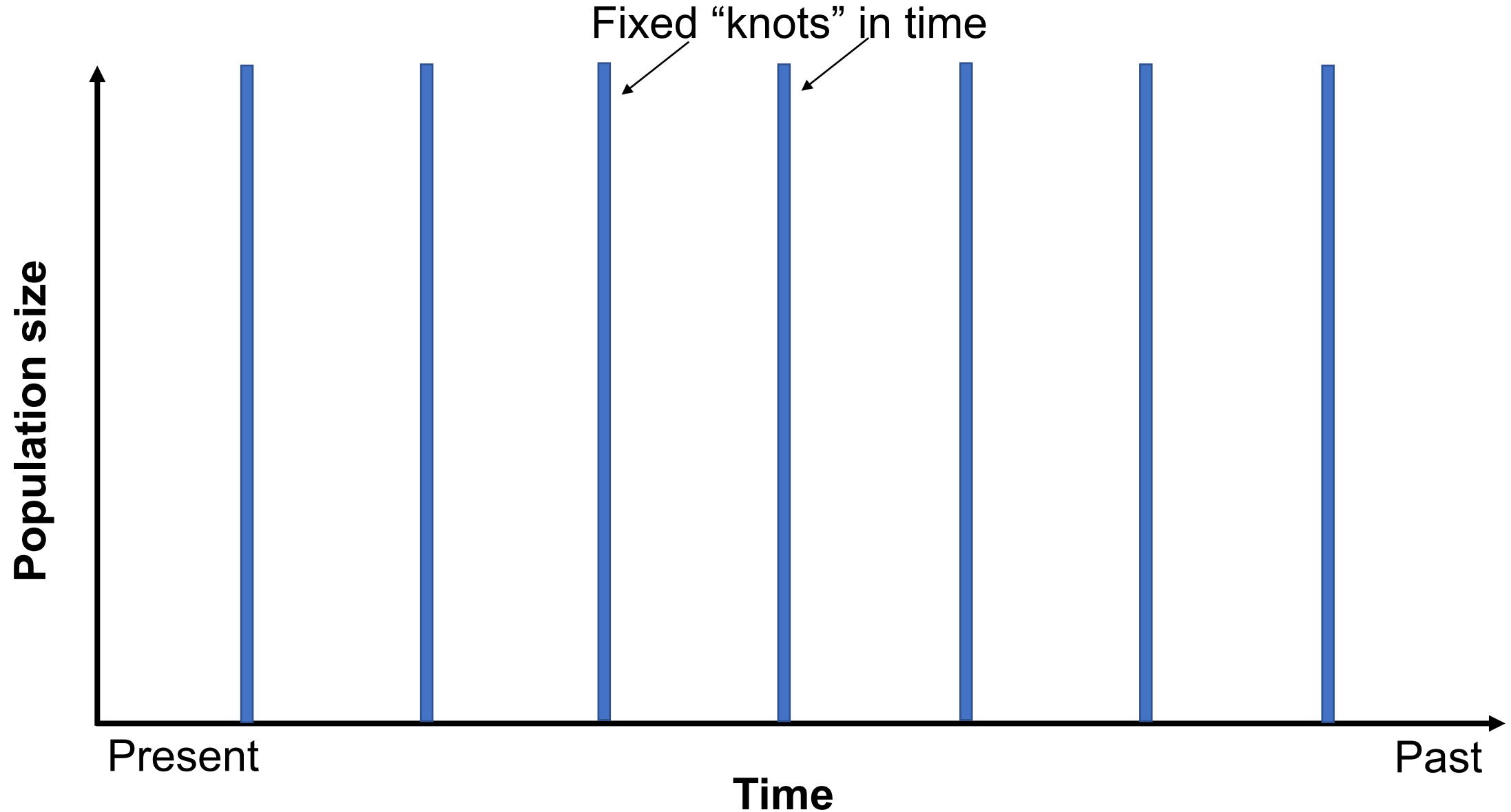
Theresa L. Cole<sup>a,b,1</sup>, Ludovic Dutoit<sup>a,2</sup>, Nicolas Dussex<sup>c,d,2</sup>, Tom Hart<sup>e,2</sup>, Alana Alexander<sup>d,2</sup>, Jane L. Younger<sup>f</sup>, Gemma V. Clucas<sup>g,h</sup>, María José Frugone<sup>i,j</sup>, Yves Cherel<sup>k</sup>, Richard Cuthbert<sup>l,m</sup>, Ursula Ellenberg<sup>n,o</sup>, Steven R. Fiddaman<sup>p</sup>, Johanna Hiscock<sup>q</sup>, David Houston<sup>r</sup>, Pierre Jouventin<sup>s</sup>, Thomas Mattern<sup>a</sup>, Gary Miller<sup>t,u</sup>, Colin Miskelly<sup>v</sup>, Paul Nolan<sup>w</sup>, Michael J. Polito<sup>x</sup>, Petra Quillfeldt<sup>y</sup>, Peter G. Ryan<sup>z</sup>, Adrian Smith<sup>p</sup>, Alan J. D. Tennyson<sup>y</sup>, David Thompson<sup>aa</sup>, Barbara Wienecke<sup>bb</sup>, Juliana A. Vianna<sup>cc</sup>, and Jonathan M. Waters<sup>a</sup>

26690–26696 | PNAS | December 26, 2019 | vol. 116 | no. 52

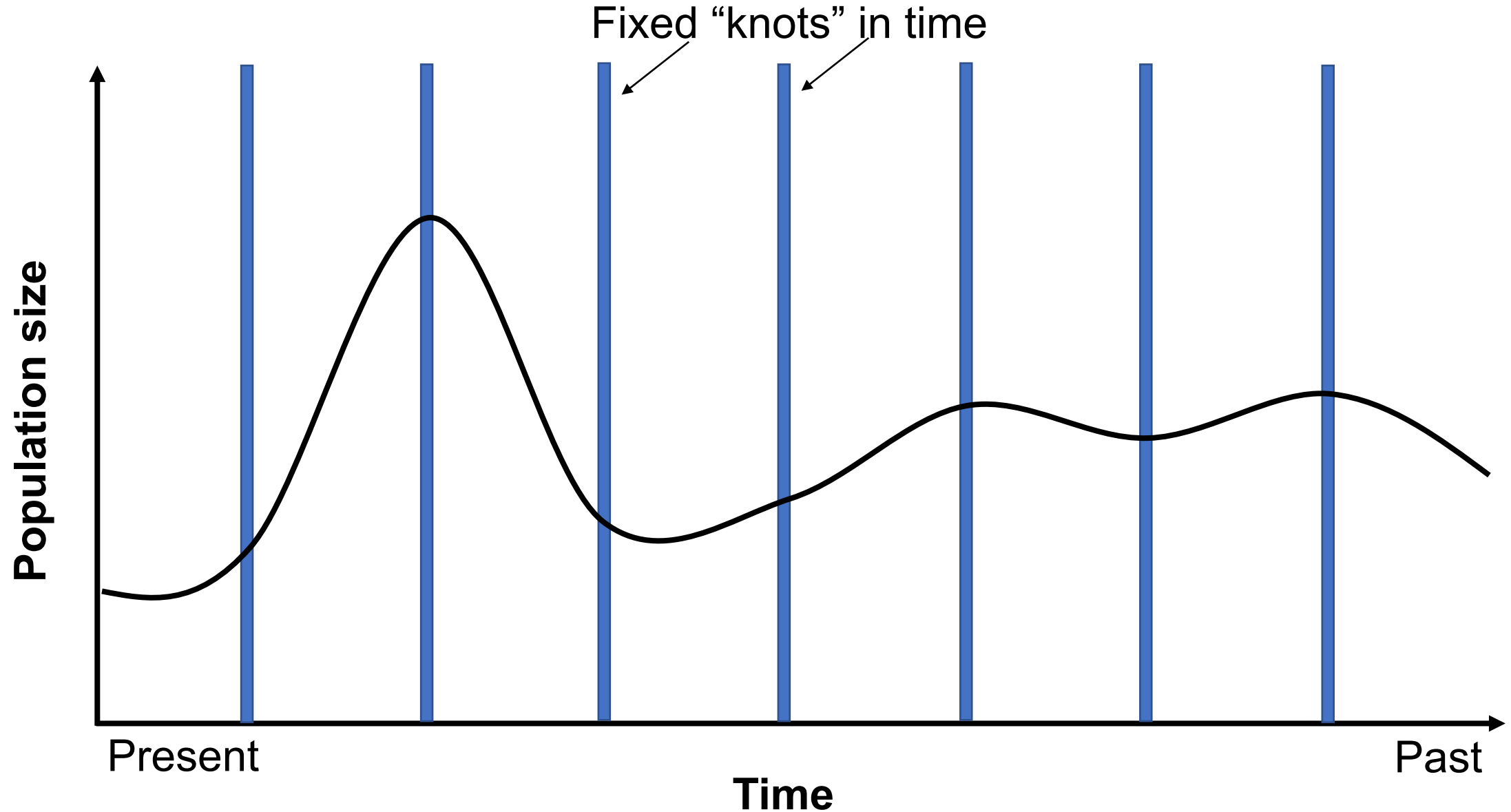
[www.pnas.org/cgi/doi/10.1073/pnas.1904048116](http://www.pnas.org/cgi/doi/10.1073/pnas.1904048116)



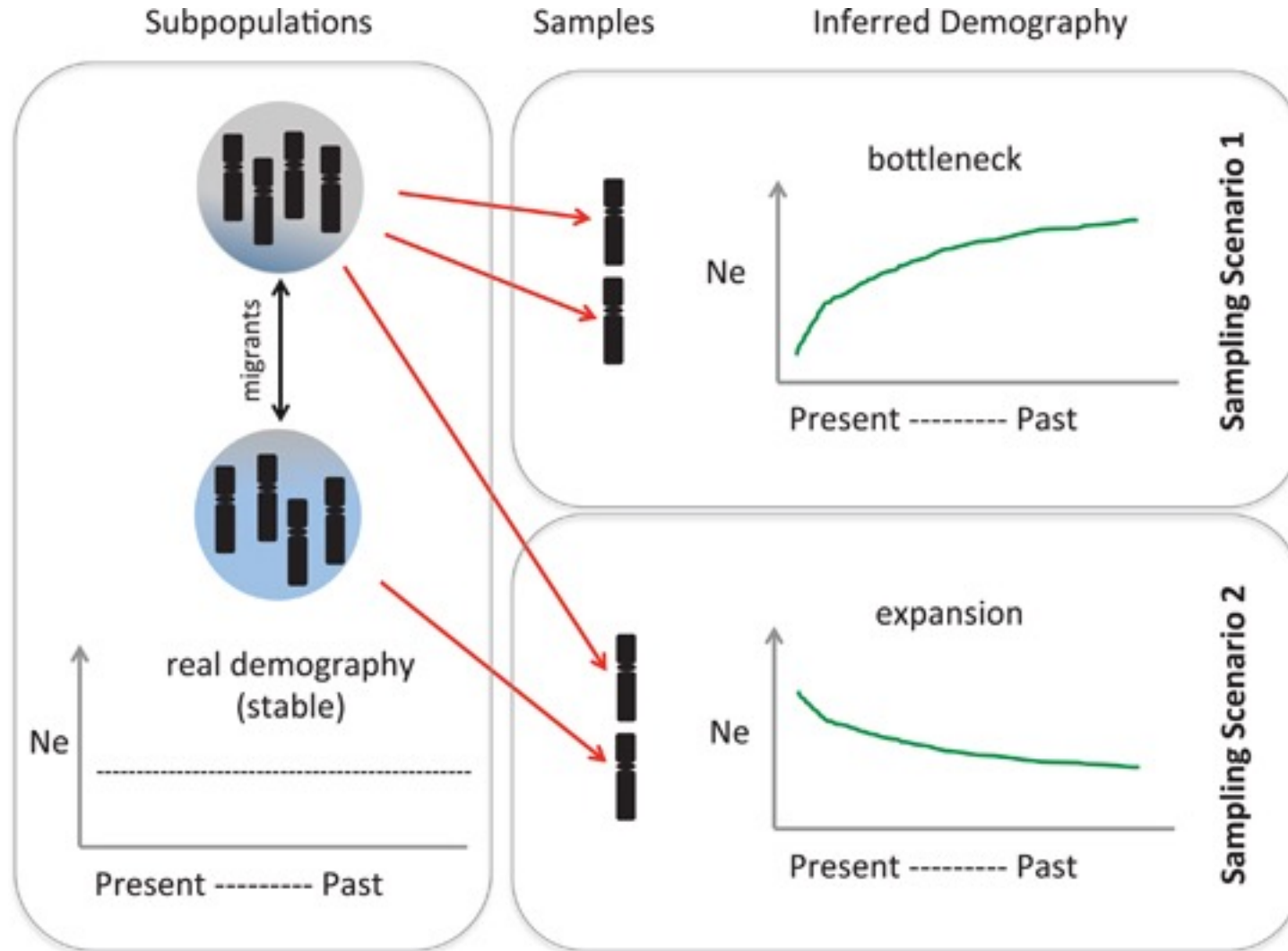
# CubSFS is quick...



# CubSFS is quick...



# Like many models, only valid if no population structure



# Using CubSFS to investigate historical demography of kanakana (lamprey)



# Using CubSFS to investigate historical demography of kanakana (lamprey)



Access to folded SFS data thanks to Allison Miller:

<https://gemmell-lab.otago.ac.nz/our-team/21-team/phd-students/153-allison-miller>

Lamprey community science:

[https://www.inaturalist.org/projects/lamps\\_for\\_champs\\_obs](https://www.inaturalist.org/projects/lamps_for_champs_obs)

<https://a3miller.wixsite.com/fishybites>