

“R and SLiM walk into a bar...”

– introduction to *slendr*, plus a few tips & tricks
discovered during its development



www.slendr.net

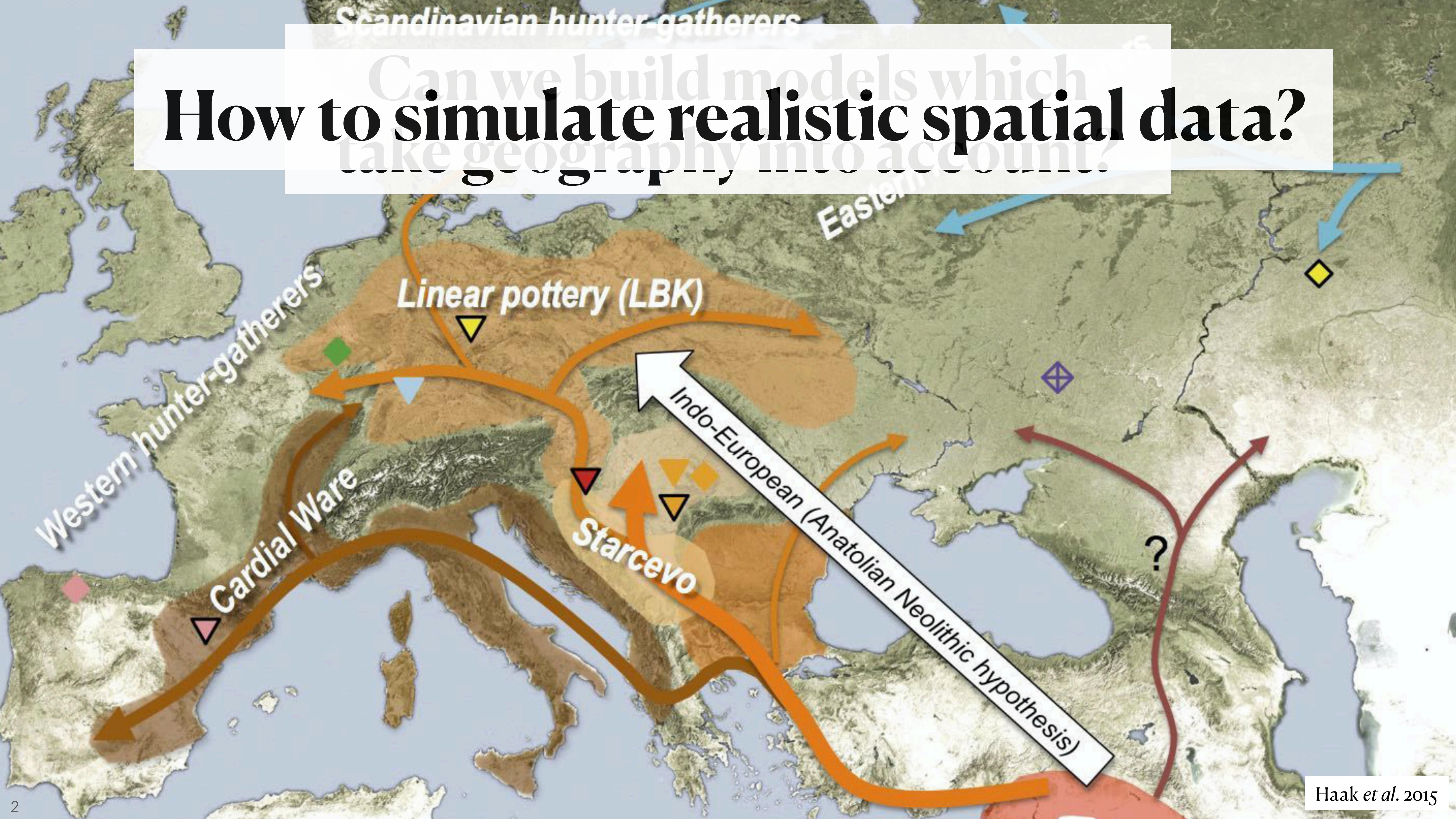
Martin Petr

mp@bodkan.net

 @fleventy5

Globe Institute
University of Copenhagen

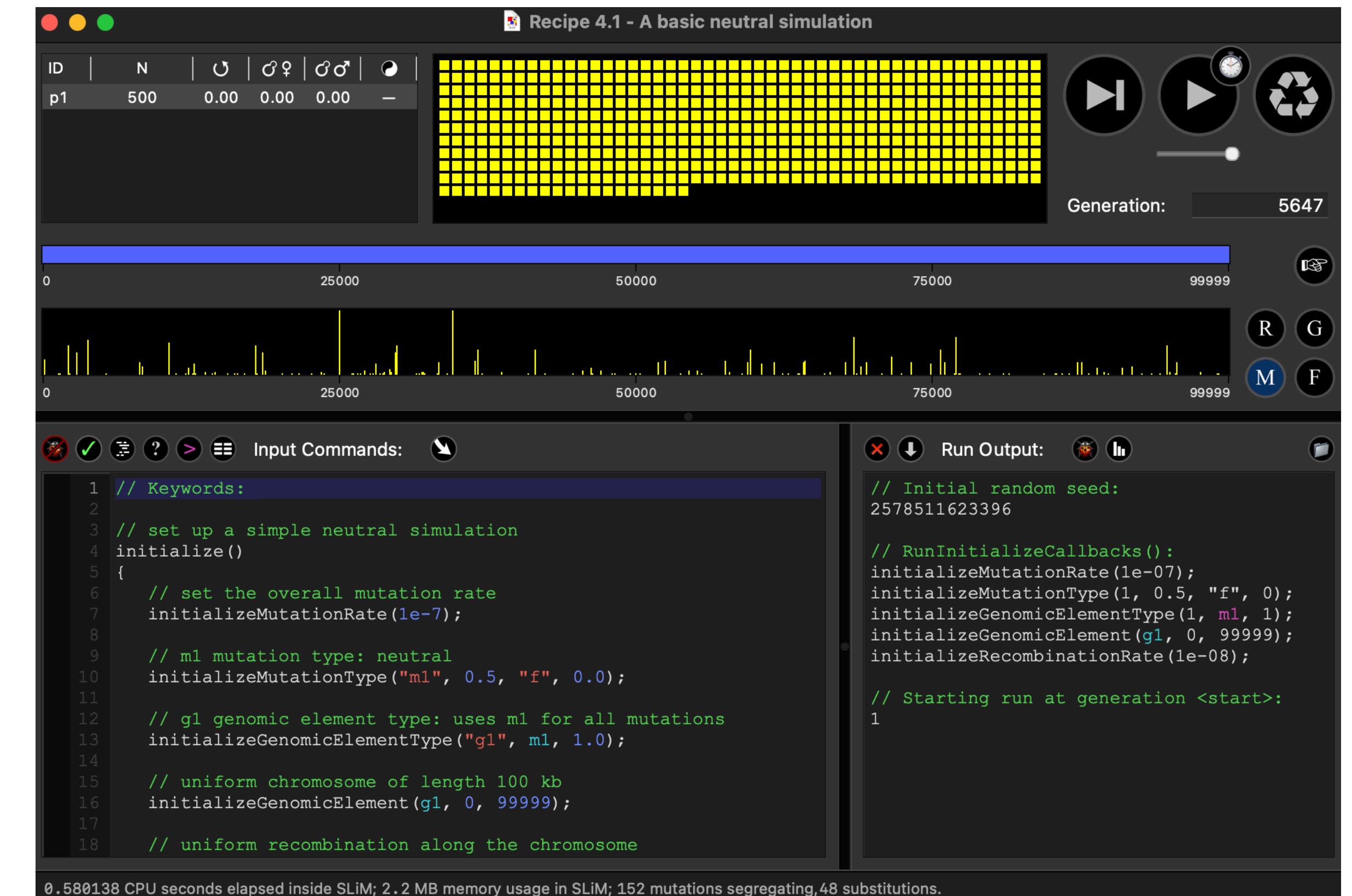
How to simulate realistic spatial data? take geography into account?



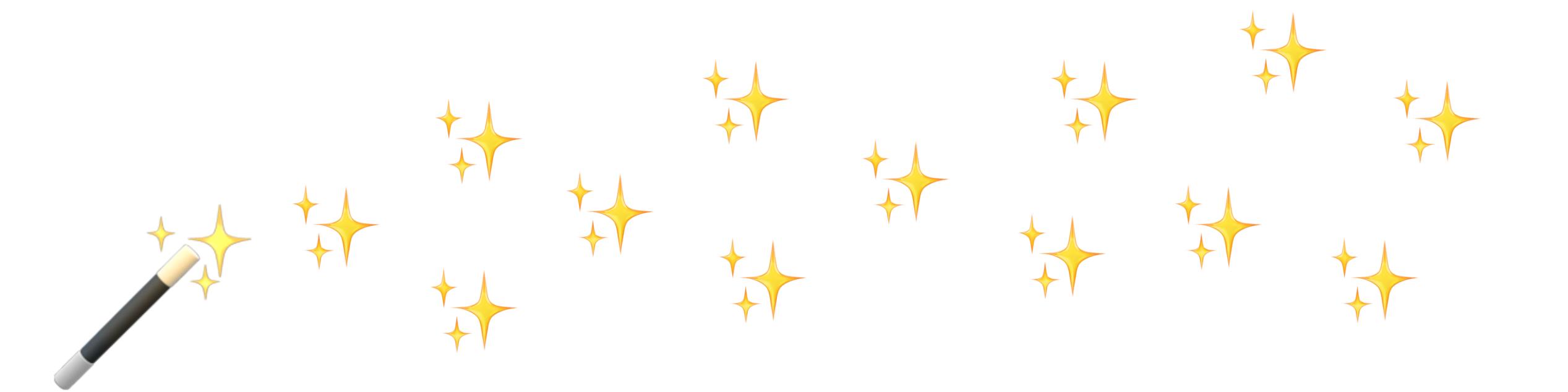
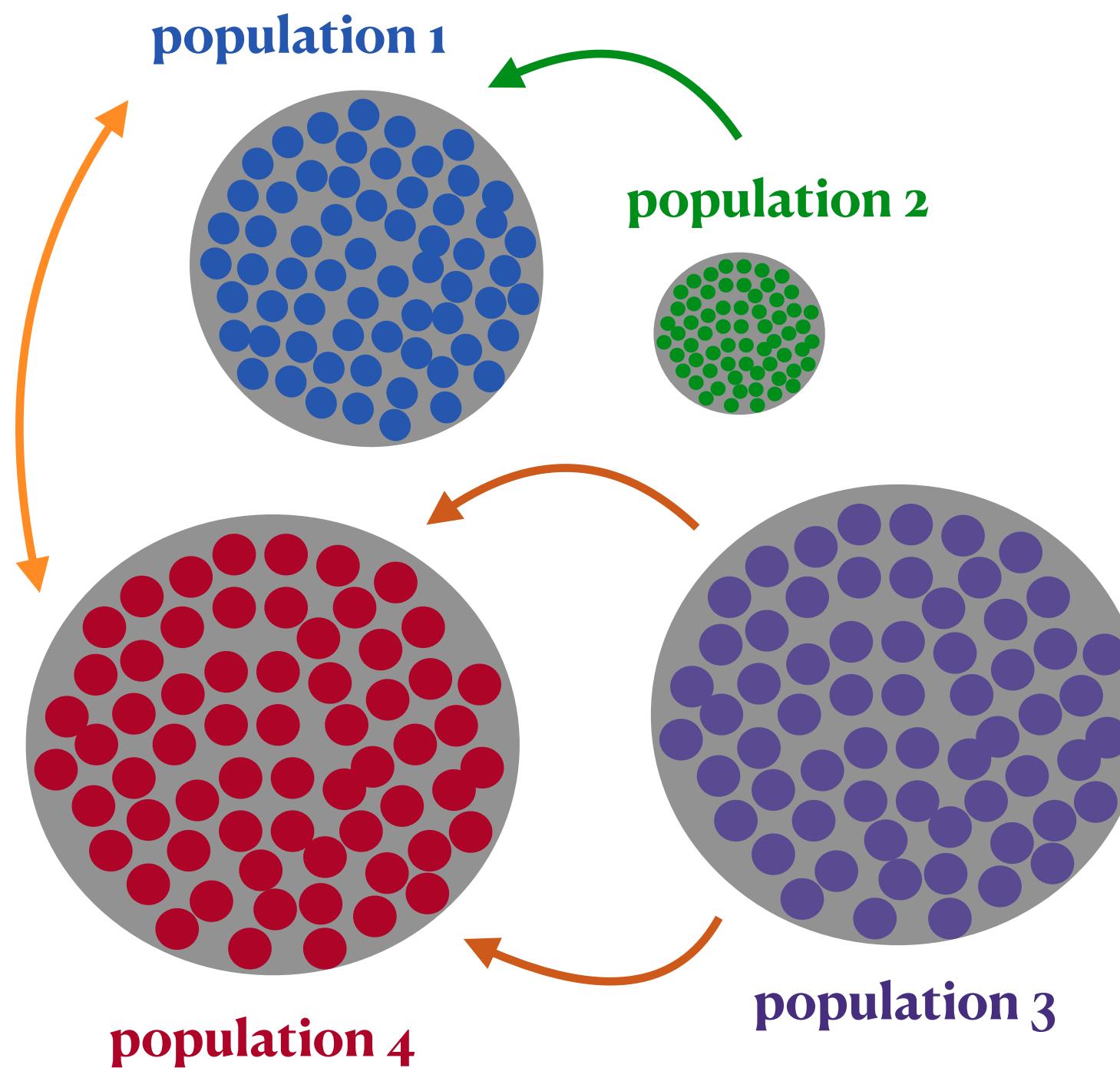
SLiM as the obvious first building block!

- extremely popular and powerful
- custom programming language
- complex demographic models
- custom mating schemes

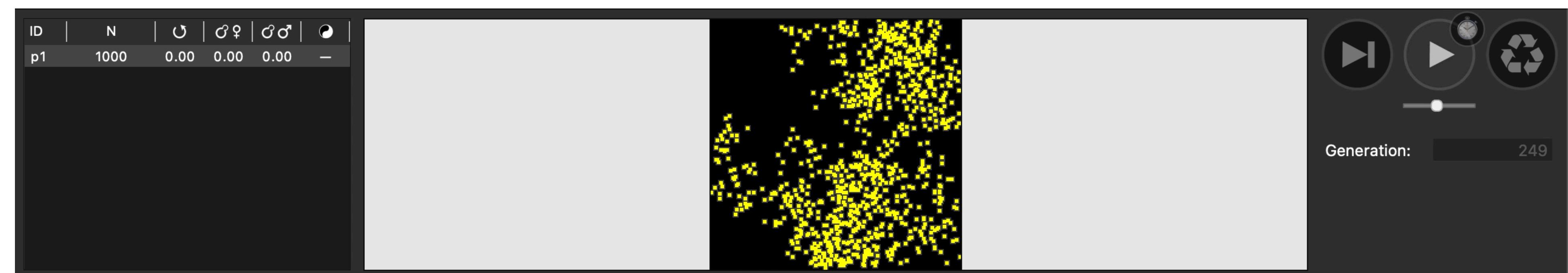
... also continuous space!



Random-mating simulation:

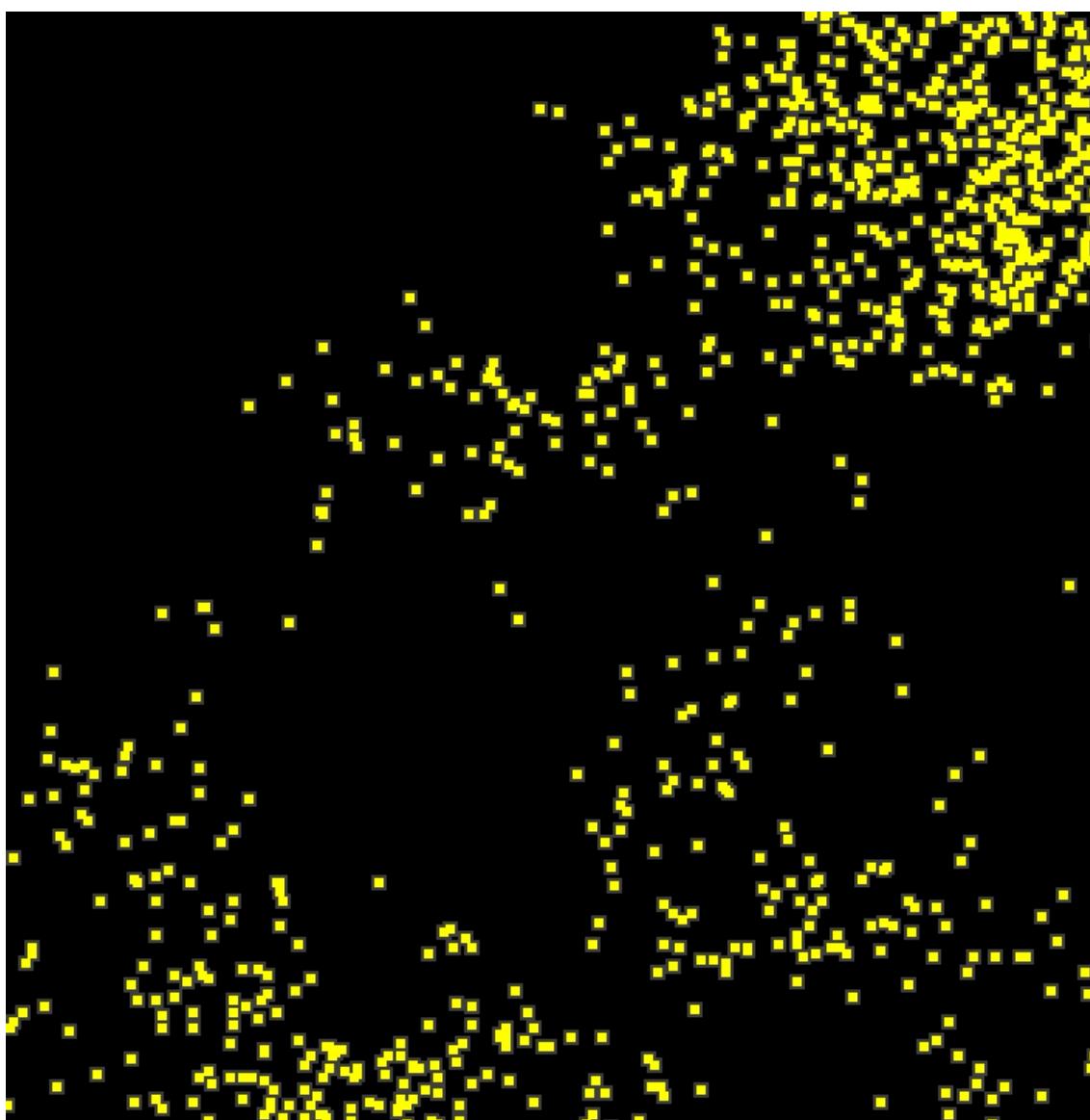


```
initializeSLiMOptions(dimentionality="xy");
```

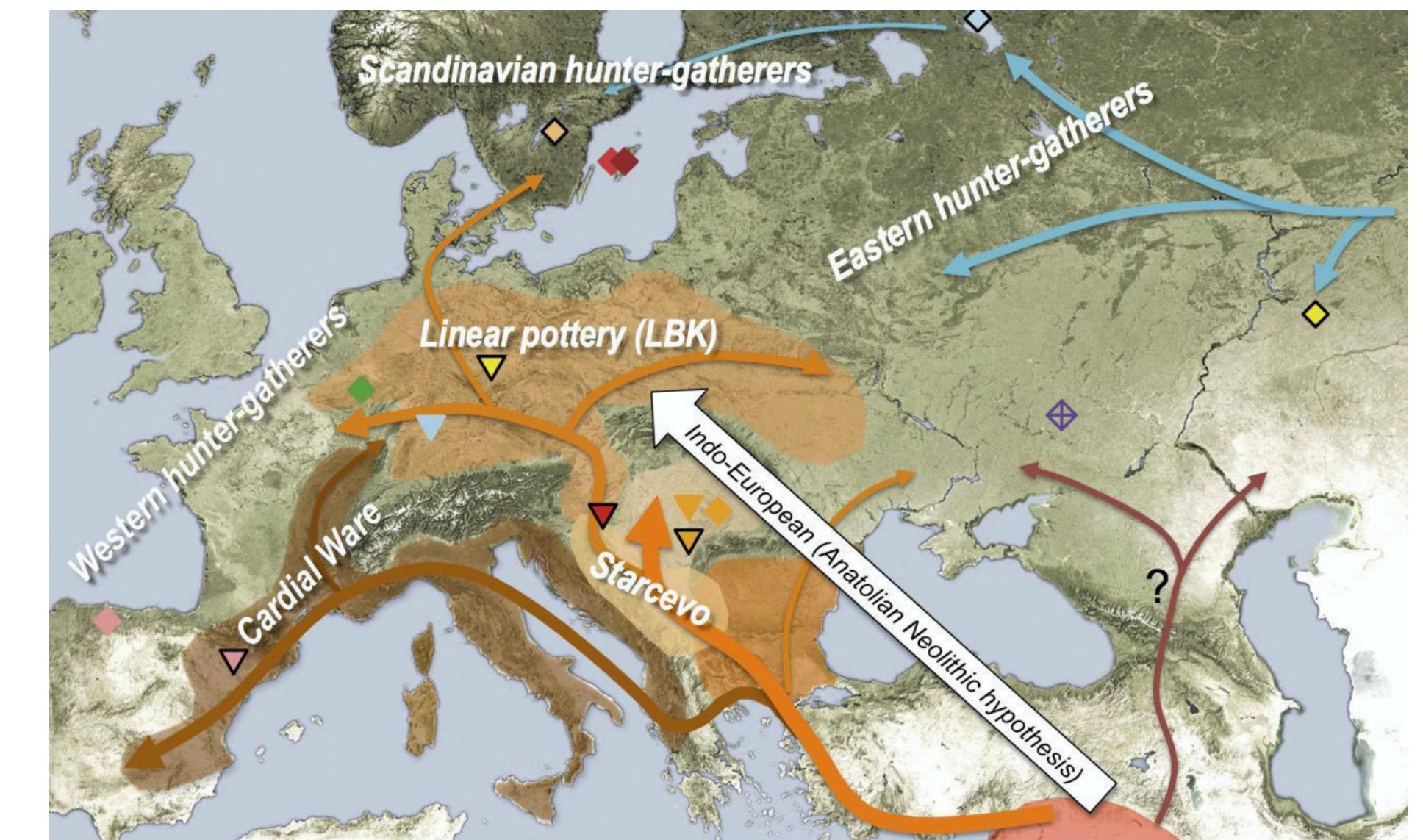


Each individual has a coordinate on a two-dimensional plane!

How can we get from here...

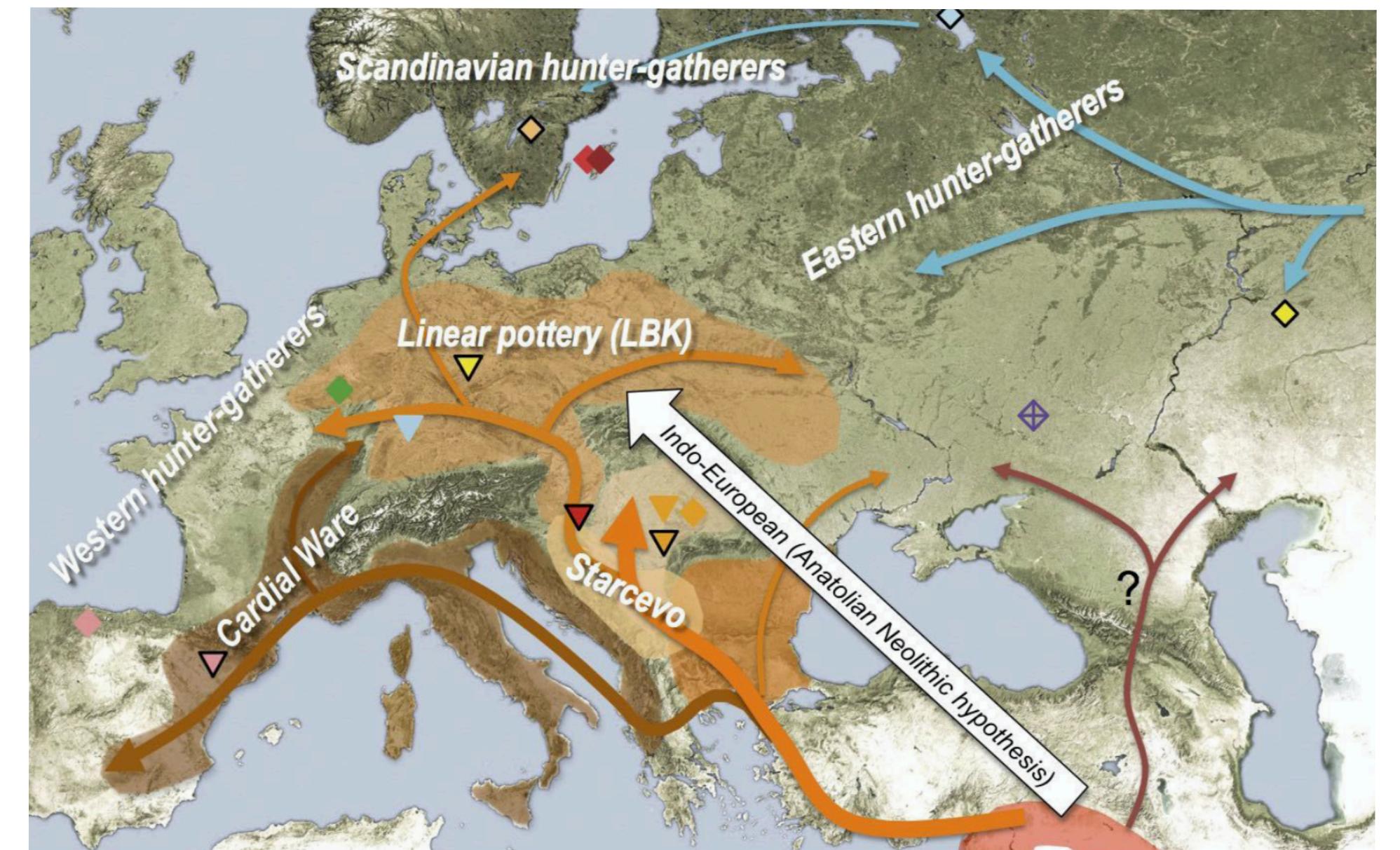


...to here?



Design goal: a programmable framework for spatial popgen

- create a "mini-language" for formal definition of
 - population ranges, movement, expansion, ...
 - no geospatial know-how required!
- use SLiM as a “low-level” engine
- implement in 
- build a research tool, but also a means of exploration



R package *slendr*



www.slendr.net

Model definition & simulation

Model components

<code>population()</code>	Define a population and its original spatial range
<code>world()</code>	Define a world map for all spatial operations
<code>geneflow()</code>	Define a geneflow event
<code>region()</code>	Define a geographic region

Population dynamics

<code>move()</code>	Move the population to a new location in a given amount of time
<code>expand()</code>	Expand the population range
<code>shrink()</code>	Shrink the population range
<code>boundary()</code>	Update the population range
<code>resize()</code>	Resize the population size
<code>dispersal()</code>	Change dispersal parameters

Manipulation of spatial objects

<code>join()</code>	Merge two spatial <code>slendr</code> objects into one
<code>overlap()</code>	Generate the overlap of two <code>slendr</code> objects
<code>subtract()</code>	Generate the difference between two <code>slendr</code> objects
<code>reproject()</code>	Reproject coordinates between coordinate systems
<code>distance()</code>	Calculate the distance between a pair of spatial boundaries
<code>area()</code>	Calculate the area covered by the given <code>slendr</code> object
<code>dimensions()</code>	Return the dimensions of the world map

Model visualization and diagnostics

<code>plot(<slendr>)</code>	Plot <code>slendr</code> geographic features on a map
<code>plot_graph()</code>	Plot geneflow graph based on given model configuration
<code>plot_ancestry()</code>	Plot simulated ancestry proportions
<code>animate()</code>	Animate the simulated population dynamics
<code>explore()</code>	Open an interactive browser of the spatial model

Compiling and running spatial models

<code>compile()</code>	Compile the spatial demographic model
<code>sampling()</code>	Define sampling events at specified times for a given set of populations
<code>slim()</code>	Run a <code>slendr</code> model as a SLiM script
<code>read()</code>	Read a previously serialized model configuration
<code>script()</code>	Substitute variables in a template SLiM script

Tree-sequence analysis

Tree sequence loading and processing

<code>ts_load()</code>	Load a tree sequence file produced by a given model
<code>ts_recapitiate()</code>	Recapitiate the tree sequence
<code>ts_simplify()</code>	Simplify the tree sequence down to a given set of individuals
<code>ts_mutate()</code>	Add mutations to the given tree sequence
<code>ts_coalesced()</code>	Check that all trees in the tree sequence are fully coalesced

Extracting genotypes from tree sequences

<code>ts_genotypes()</code>	Extract genotype table from the tree sequence
<code>ts_eigenstrat()</code>	Extract genotypes from the tree sequence in the EIGENSTRAT format
<code>ts_vcf()</code>	Save genotypes from the tree sequence as a VCF file

Accessing tree sequence components

<code>ts_data()</code>	Extract combined annotated table of individuals and nodes
<code>ts_ancestors()</code>	Infer spatio-temporal ancestral history for given nodes/individuals
<code>ts_individuals()</code> <code>ts_edges()</code> <code>ts_nodes()</code>	Get the table of individuals/nodes/edges from the tree sequence
<code>ts_samples()</code>	Extract names and times of individuals scheduled for sampling
<code>ts_tree()</code>	Get a tree from a given tree sequence
<code>ts_draw()</code>	Plot a graphical representation of a single tree

Tree sequence statistics

<code>ts_f2()</code> <code>ts_f3()</code> <code>ts_f4()</code> <code>ts_f4ratio()</code>	Calculate the f2, f3, f4, and f4-ratio statistics
<code>ts_afs()</code>	Compute the allele frequency spectrum (AFS)
<code>ts fst()</code> <code>ts_divergence()</code>	Calculate pairwise statistics between sets of individuals
<code>ts_diversity()</code>	Calculate diversity in given sets of individuals
<code>ts_tajima()</code>	Calculate Tajima's D for given sets of individuals
<code>ts_segregating()</code>	Calculate the density of segregating sites for the given sets of individuals

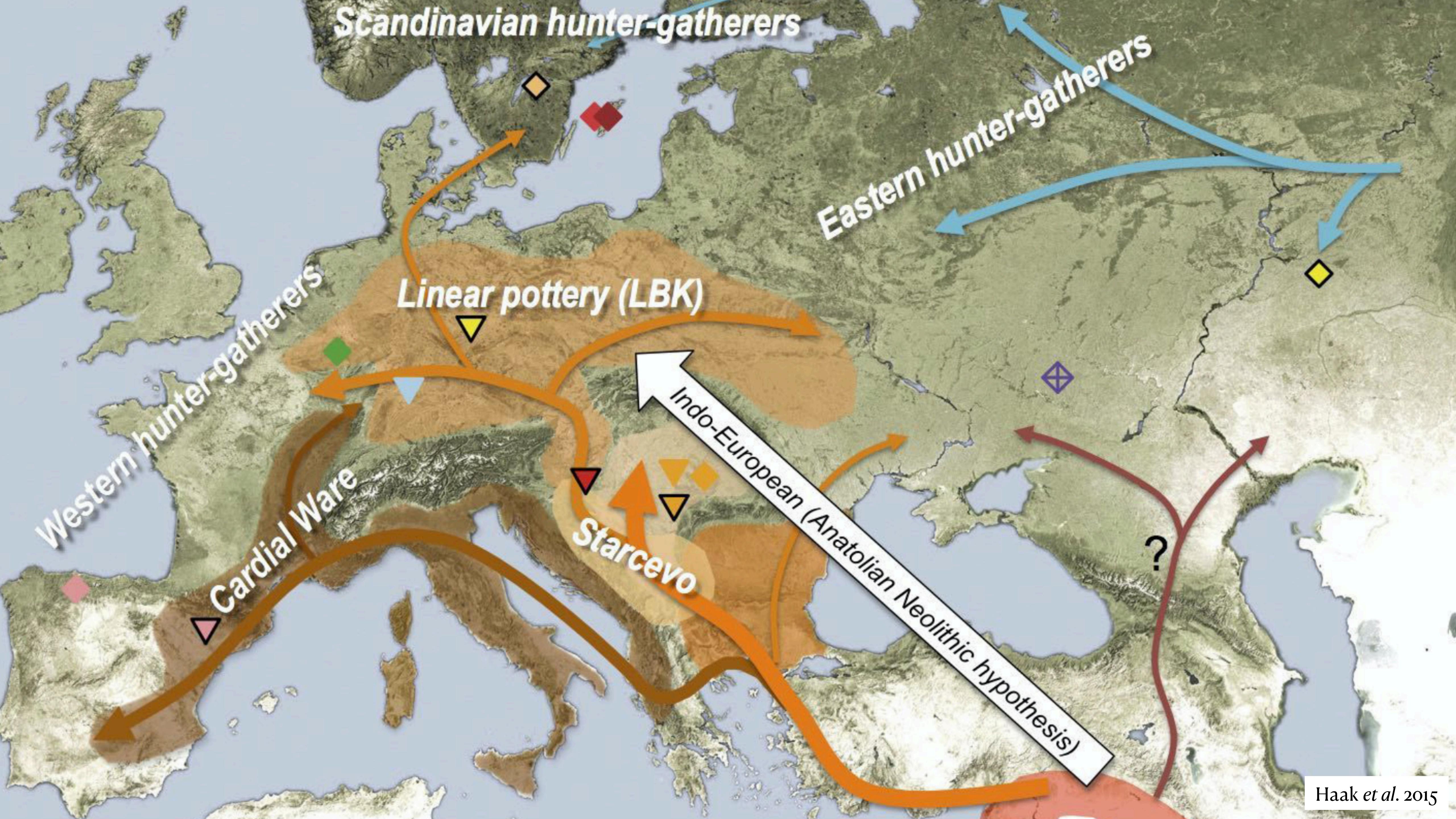
slendr workflow demo

interactive RStudio session in the browser:

github.com/bodkan/slendr

click on  launch binder

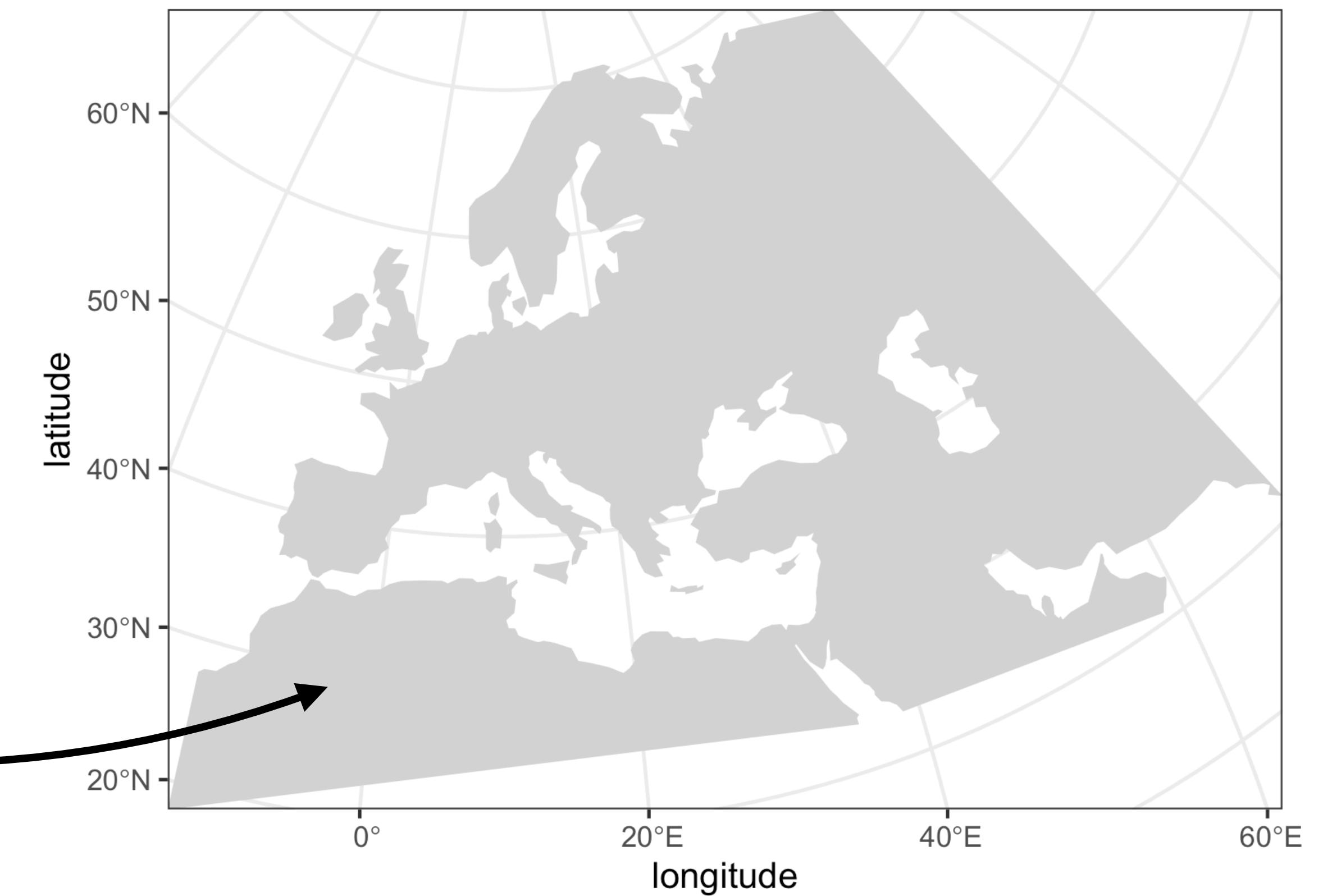
(no installation needed!)



1. Define a spatial context (“real world()”)

(section of the world which will be occupied by populations)

```
map <- world()  
xrange = c(-13, 70),          # min-max longitude  
yrange = c(18, 65),          # min-max latitude  
landscape = "naturalearth",  # source of data  
crs = "EPSG:3035"           # coordinate reference system (CRS)
```



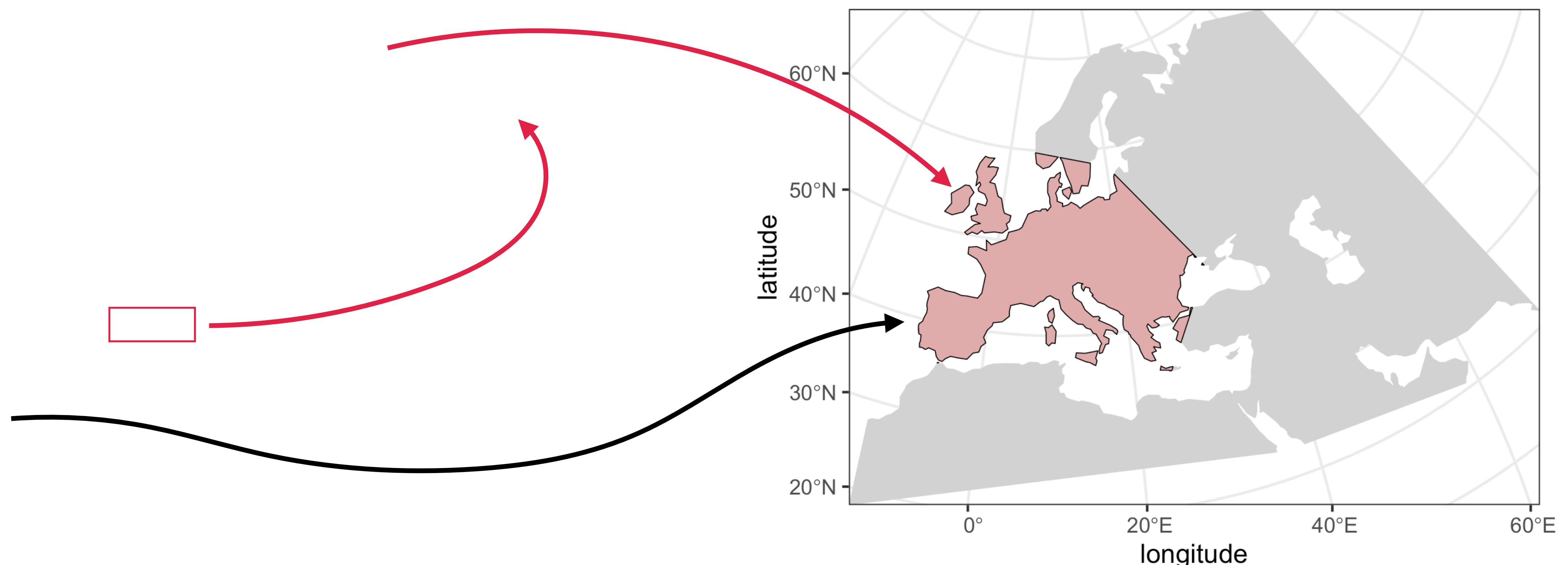
(using freely available *Natural Earth* cartographic data)

2. population() (and its boundaries)

Everything is stored in a vectorised form!

```
> eur$geometry[[1]]
```

```
POLYGON ((2676851 1520500, 2568327 3290814, 3489173 4188384, 4910944 3919000,  
6104864 2710417, 5695390 1461722, 5240847 1387400, 4321000 1657812, 2965317 15  
70227, 2676851 1520500))
```



3. Population displacement

```
yam <- population( # Yamnaya steppe population  
  name = "YAM", time = 7000, N = 500, remove = 2500, map = map,  
  polygon = list(c(26, 50), c(38, 49), c(48, 50),  
    c(48, 56), c(38, 59), c(26, 56))  
)
```

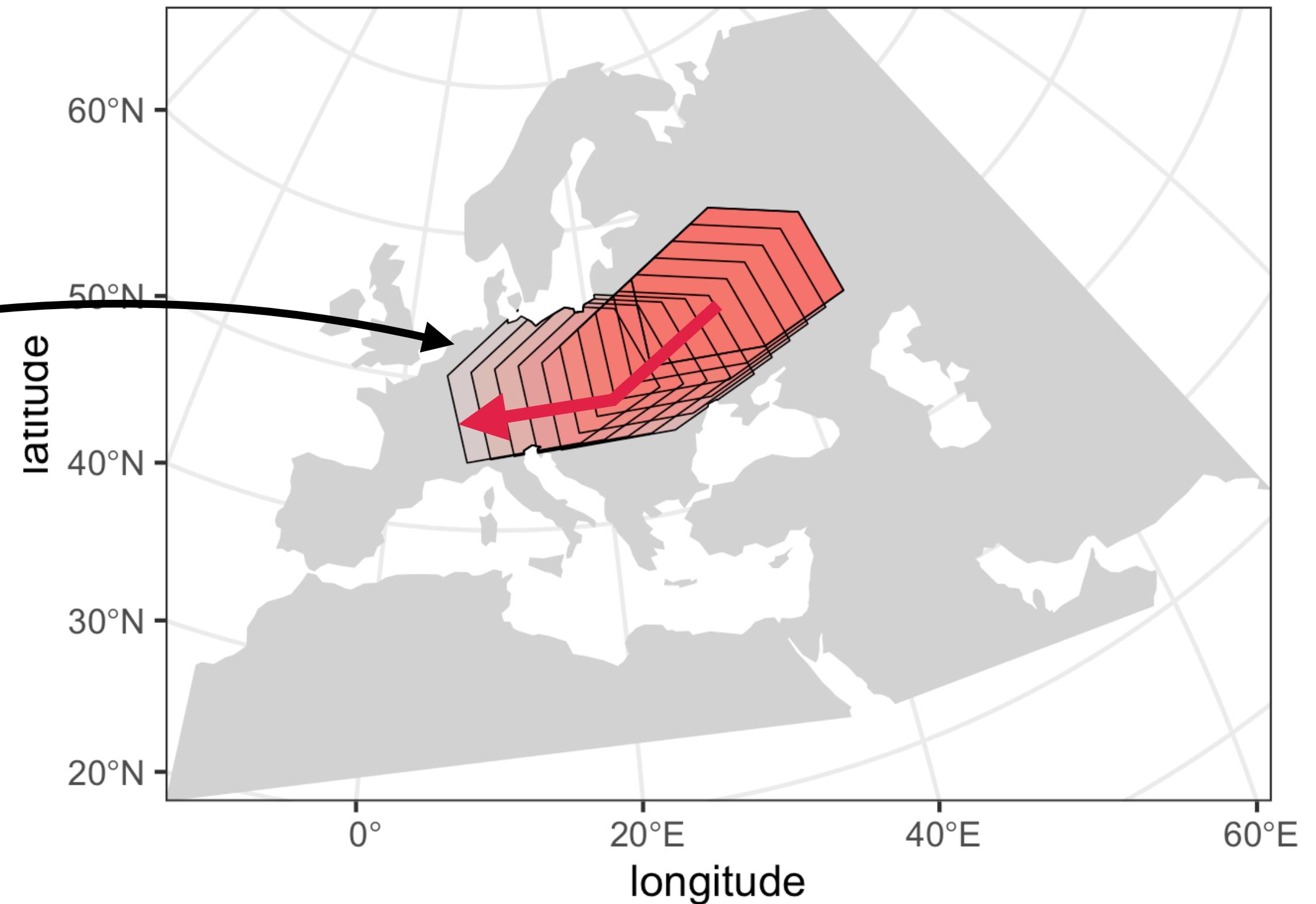
```
> yam  
slendr 'population' object
```

```
-----  
name: YAM  
habitat: terrestrial
```

```
number of spatial maps: 11  
map: internal coordinate reference system EPSG 3035  
scheduled removal at time 2500
```

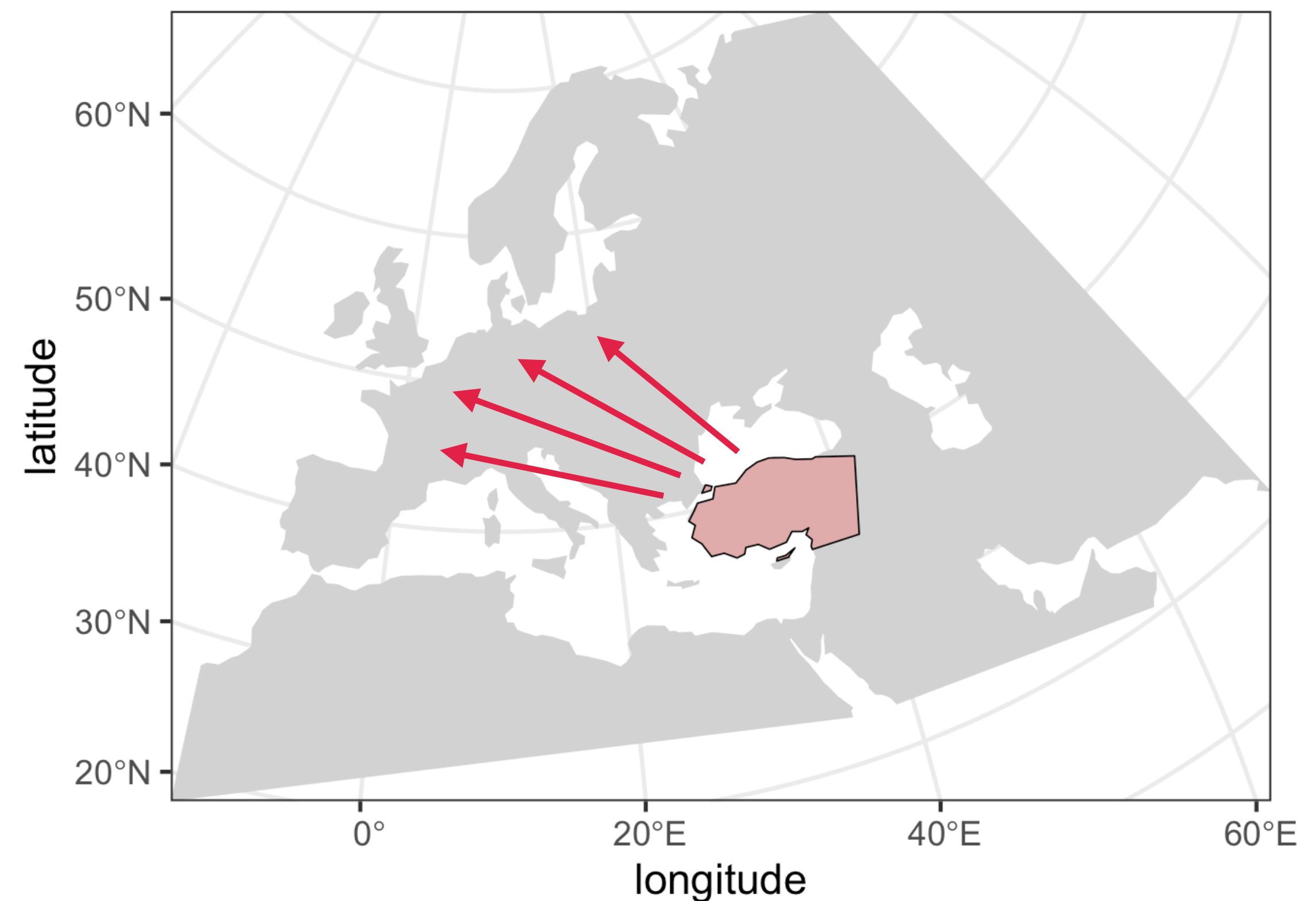
population history overview:

- time 7000: created as an ancestral population ($N = 500$)
- time 5000-3000: movement across a landscape



4. Population range expansion

```
ana <- population( # Anatolian farmers  
  name = "ANA", time = 28000, N = 3000, map = map,  
  center = c(34, 38), radius = 500e3, polygon = anatolia  
)
```

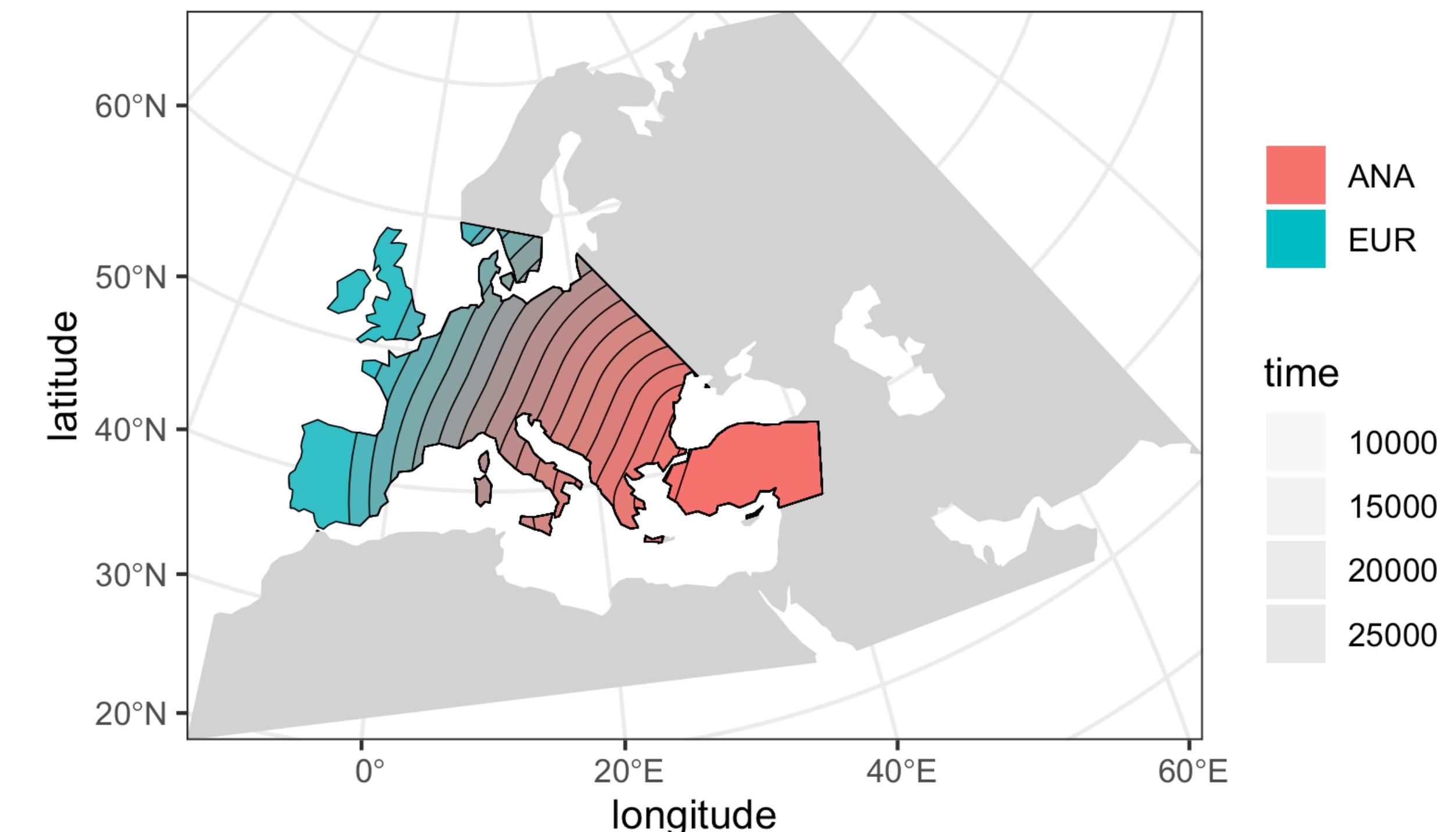


5. gene_flow() events

```
eur <- population(  
  name = "EUR", time = 25000, N = 2000,  
  map = map, polygon = europe  
)
```

```
ana <- population( # Anatolian farmers  
  name = "ANA", time = 28000, N = 3000, map = map,  
  center = c(34, 38), radius = 500e3, polygon = anatolia  
) %>%  
  expand_range( # expand the range by 2.500 km  
    by = 2500e3, start = 10000, end = 7000,  
    polygon = join(europe, anatolia),  
)
```

```
gf <- gene_flow(from = ana, to = eur,  
  rate = 0.5, start = 8000, end = 6000)
```



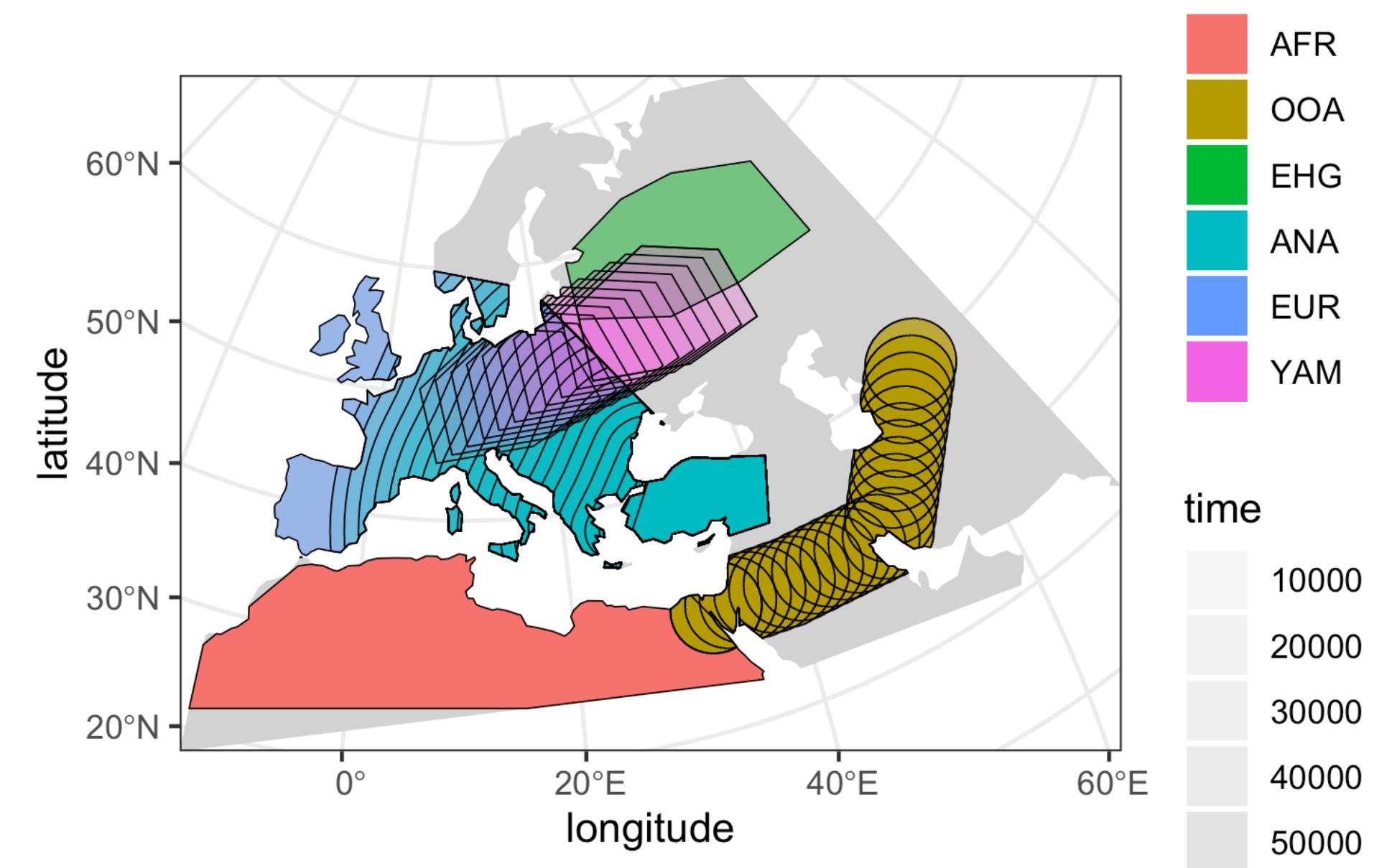
(gene_flow() performs sanity checks to make sure the event makes sense as programmed)

6. Model compilation

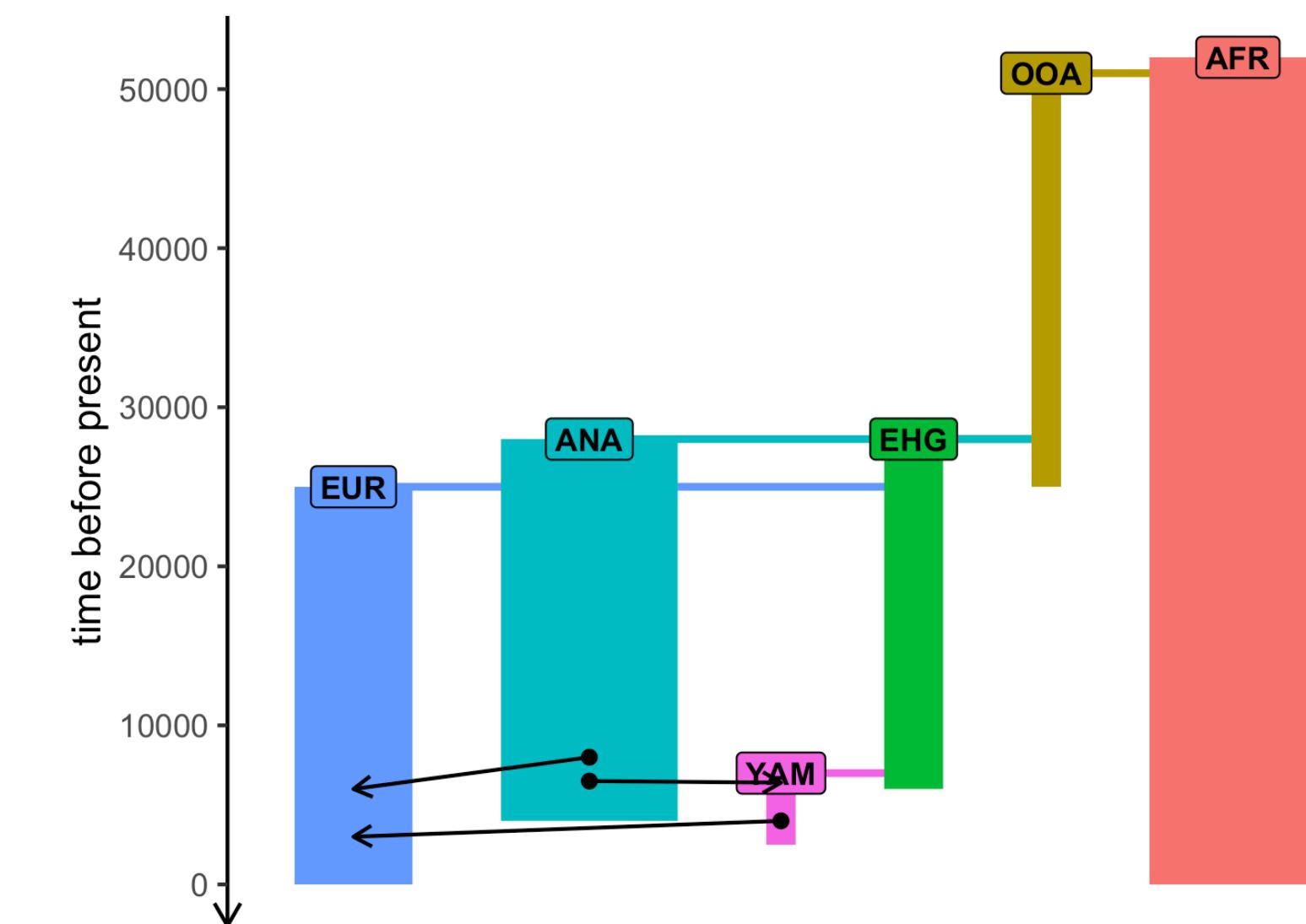
we created these R objects earlier

```
model <- compile_model(  
  populations = list(afr, ooa, ehg, eur, ana, yam),  
  gene_flow = gf,  
  generation_time = 30,  
  resolution = 10e3,           # resolution in meters per pixel  
  competition = 130e3, mating = 100e3, # spatial interaction parameters  
  dispersal = 70e3            # how far can offspring 'migrate'  
)
```

`plot_map(model)` – “compressed” spatial dynamics



`plot_model(model)` – embedded demographic model



How is a **slendr** model represented in the background?

`compile_model()` generates these files on disk:

1.png	PNG image
2.png	PNG image
3.png	PNG image
4.png	PNG image
5.png	PNG image
[...]	
57.png	PNG image
58.png	PNG image
59.png	PNG image
60.png	PNG image
checksums.tsv	Plain Text
description.txt	Plain Text
direction.txt	Plain Text
dispersals.tsv	Plain Text
geneflow.tsv	Plain Text
generation_time.txt	Plain Text
length.txt	Plain Text
maps.tsv	Plain Text
orig_length.txt	Plain Text
populations.tsv	Plain Text
ranges.rds	R Data File
resolution.txt	Plain Text
script.py	Python Source
script.slim	SLiM model

a bunch of
image files

a bunch of
text files

a Python script
a SLiM script

```
model <- compile_model()
populations = list(afr, ooa, ehg, eur, ana, yam),
gene_flow = gf,
generation_time = 30,
resolution = 10e3, # resolution in meters per pixel
competition = 130e3, mating = 100e3, # spatial interaction parameters
dispersal = 70e3 # how far can offspring 'migrate'
```

```
> model
slendr 'model' object
-----
populations: AFR, OOA, EHG, ANA, EUR, YAM
geneflow events: 3
generation time: 30
time direction: backward
total running length: 52000 model time units
model type: spatial
- number of spatial maps: 60
- resolution: 10000 distance units per pixel

configuration files in: /private/var/folders/70/b_q2zdh116b9pfg29p
03sx60000gn/T/RtmpZwjanI/filed3d867bb9977_slendr_model
```

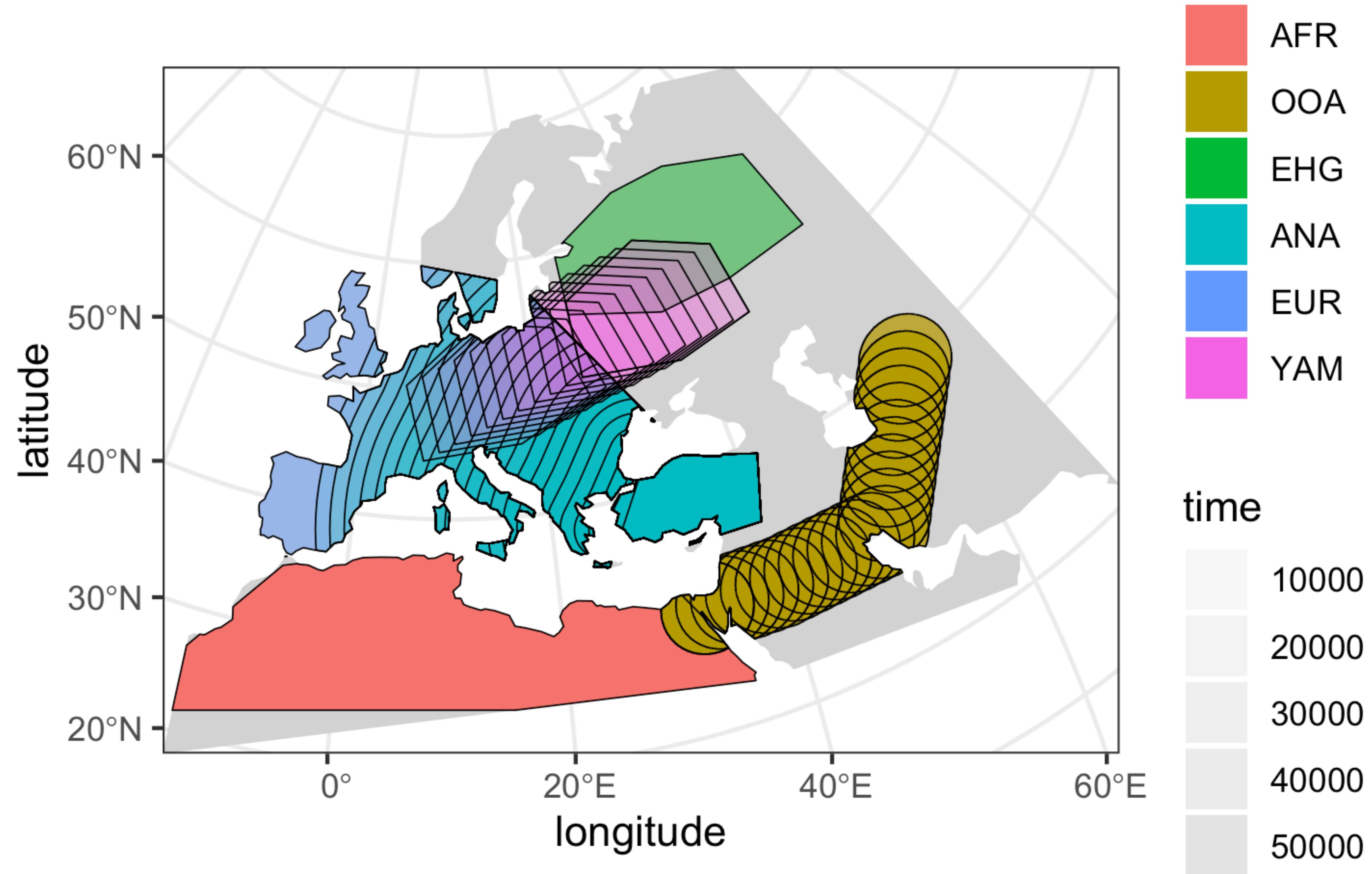
7. Simulation of the model

[slim() or msprime()]

```
ts <- slim(  
    model, # model object compiled earlier  
    sequence_length = 10e6,    # simulate 10Mb of sequence ...  
    recombination_rate = 1e-8 # ... with this uniform recombination rate  
)
```

The `slim()` function executes a built-in SLiM script designed for `slendr` demographic models.

The result `ts` is a tree sequence.



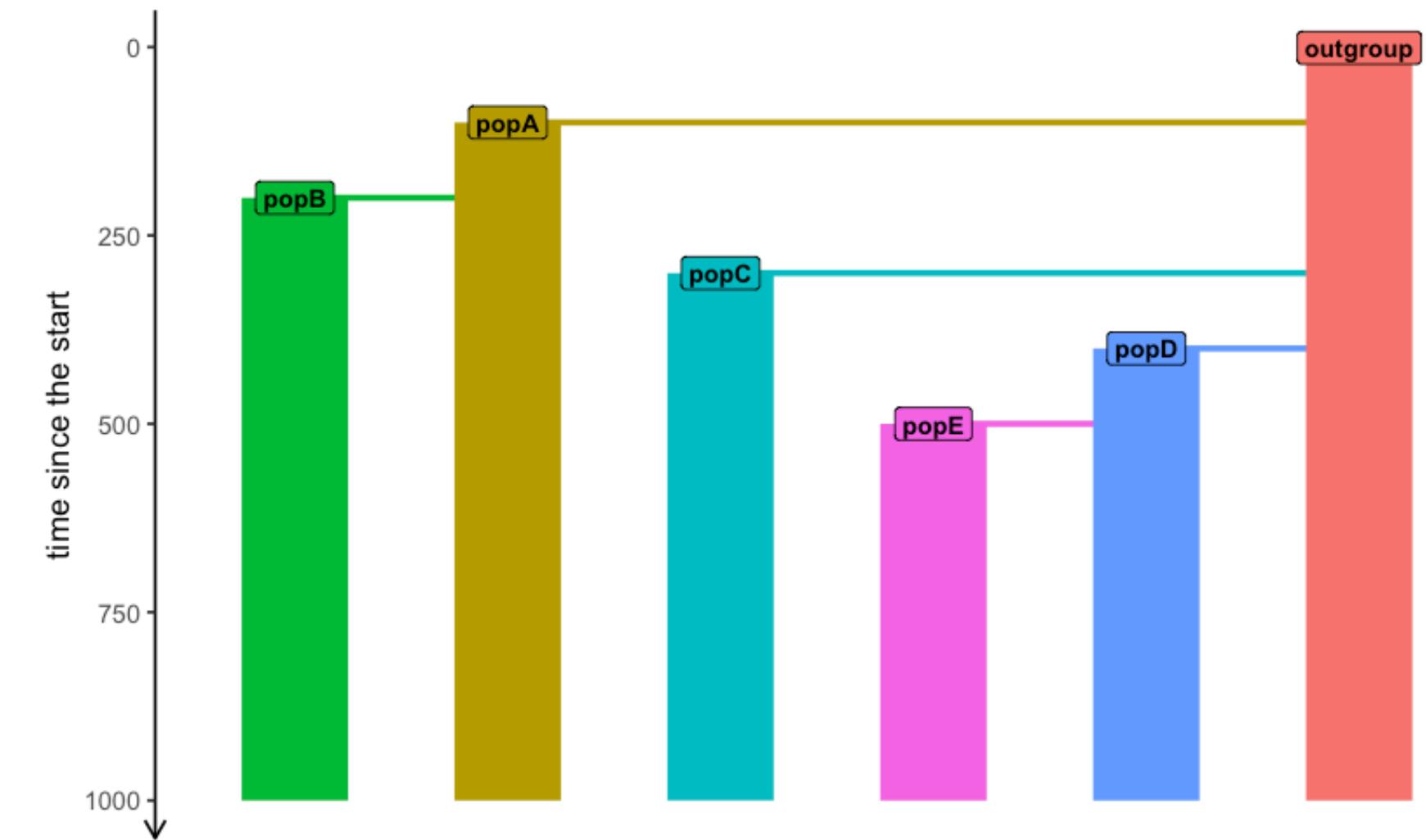
We know the genome and exact location & time of everyone



Generated by a formally defined spatio-temporal model

How is a *slendr* model executed in the background?

```
1 outgroup <- population("outgroup", time = 1, N = 1000)
2 popA <- population("popA", time = 100, N = 1000, parent = outgroup)
3 popB <- population("popB", time = 200, N = 1000, parent = popA)
4 popC <- population("popC", time = 300, N = 1000, parent = outgroup)
5 popD <- population("popD", time = 400, N = 1000, parent = outgroup)
6 popE <- population("popE", time = 500, N = 1000, parent = popD)
7
8 model <- compile_model(
9   populations = list(outgroup, popA, popB, popC, popD, popE),
10  generation_time = 1, simulation_length = 1000
11 )
```



(we will focus on non-spatial *slendr* models now, for simplicity)

The most basic SLiM equivalent of this model

```
1 initialize() {
2     initializeMutationRate(0);
3     initializeMutationType("m0", 0.5, "f", 0.0);
4     initializeGenomicElementType("g1", m0, 1.0);
5     initializeGenomicElement(g1, 0, 999999);
6     initializeRecombinationRate(1e-8);
7 }
8
9 late() {
10    sim.addSubpop("p0", 1000);
11    p0.name = "outgroup";
12 }
13
14 100 late() {
15    sim.addSubpopSplit("p1", 1000, p0);
16    p1.name = "popA";
17 }
18
19 200 late() {
20    sim.addSubpopSplit("p2", 1000, p1);
21    p2.name = "popB";
22 }
23
24 300 late() {
25    sim.addSubpopSplit("p3", 1000, p0);
26    p3.name = "popC";
27 }
28
29 400 late() {
30    sim.addSubpopSplit("p4", 1000, p0);
31    p4.name = "popD";
32 }
33
34 500 late() {
35    sim.addSubpopSplit("p5", 1000, p4);
36    p5.name = "popE";
37 }
38
39 1000 late() {
40    sim.simulationFinished();
41 }
```

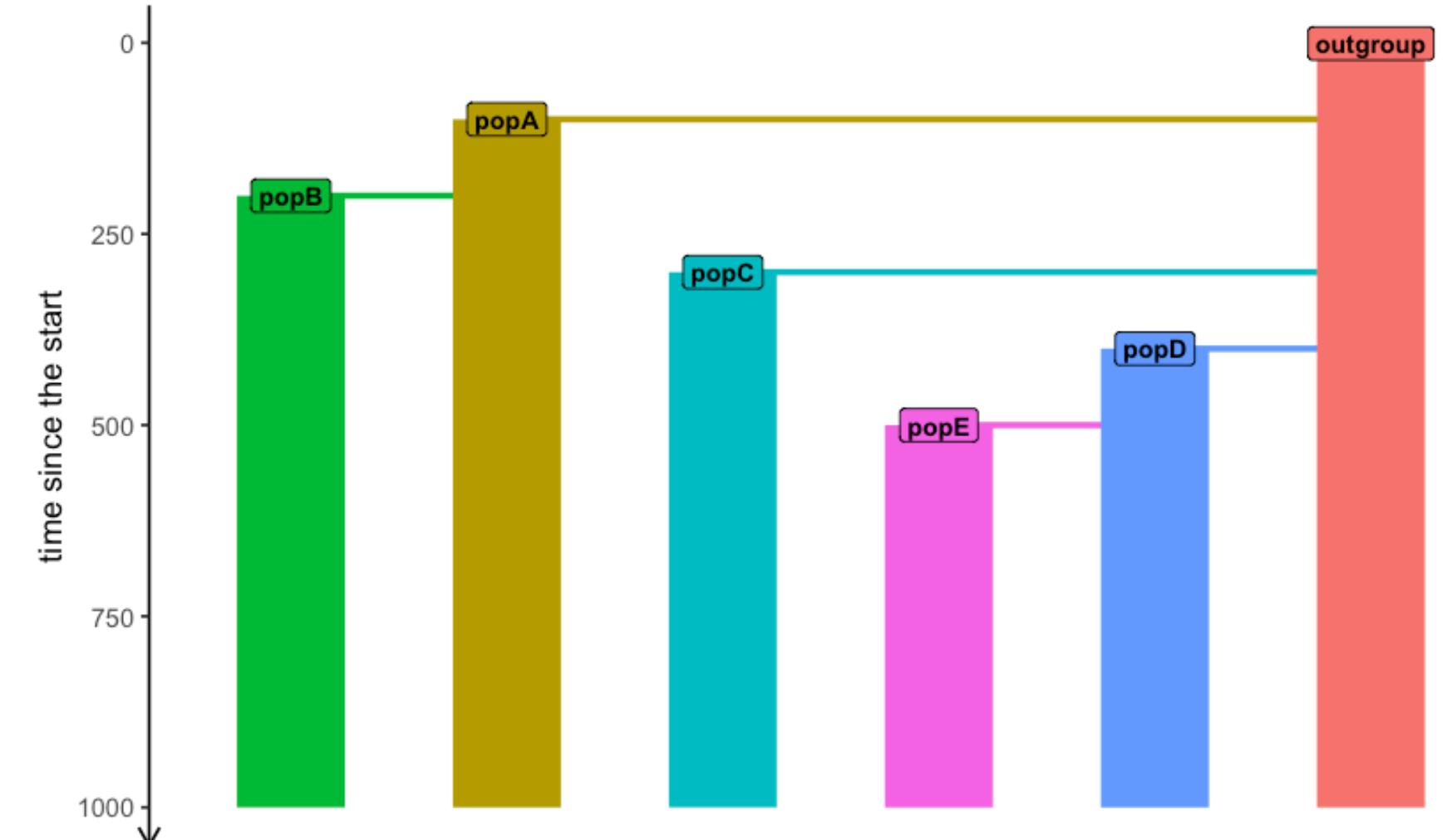
Each demographic event
(here just population splits)
has its own script block

But there are many many other demographic events we might want to simulate:

- population resizes
- gene-flow events
- spatial dynamics (expansions, contractions, displacements, ...)

Each of them would likely need a dedicated script block.

We might even want to parametrize these events.



This can make for a long and complex SLiM script, which is OK if we know exactly what sort of model to simulate – but *slendr* doesn't “know” this.

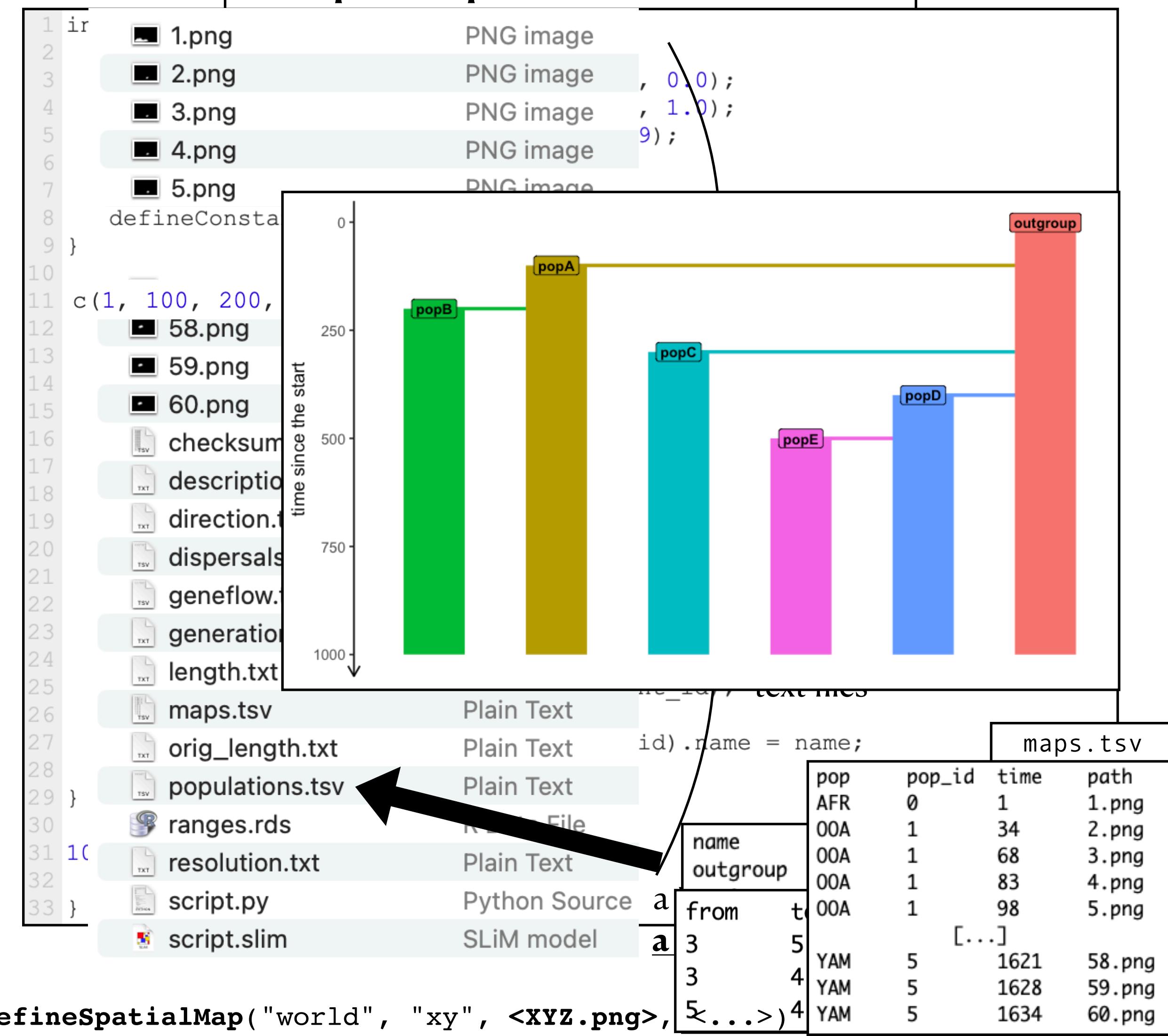
slendr solves this by turning scripts/models “inside out”

```

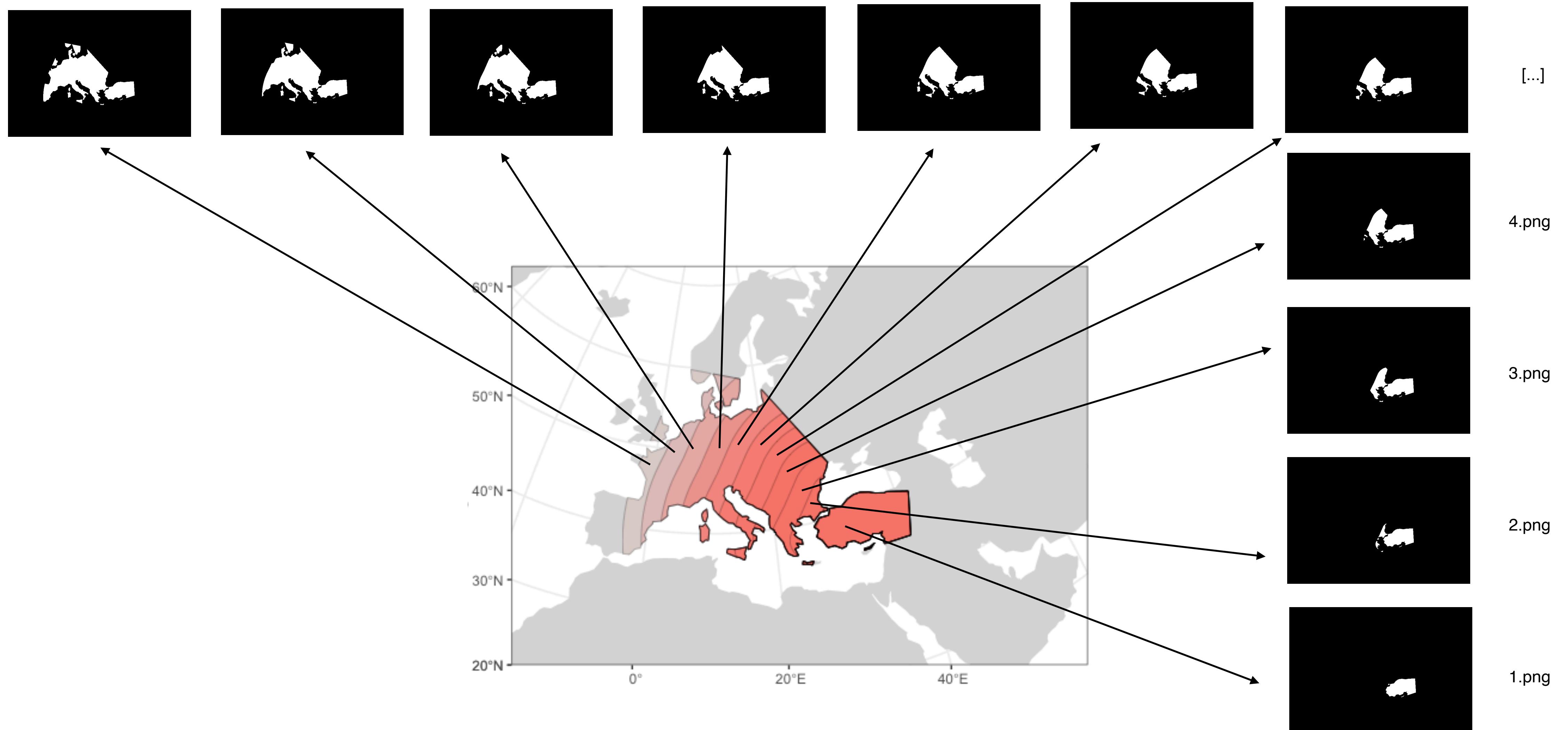
1 initialize() {
2     initializeMutationRate(0);
3     initializeMutationType("m0", 0.5, "f", 0.0);
4     initializeGenomicElementType("g1", m0, 1.0);
5     initializeGenomicElement(g1, 0, 999999);
6     initializeRecombinationRate(1e-8);
7 }
8
9 late() {
10    sim.addSubpop("p0", 1000);
11    p0.name = "outgroup";
12}
13
14 100 late() {
15    sim.addSubpopSplit("p1", 1000, p0);
16    p1.name = "popA";
17}
18
19 200 late() {
20    sim.addSubpopSplit("p2", 1000, p1);
21    p2.name = "popB";
22}
23
24 300 late() {
25    sim.addSubpopSplit("p3", 1000, p0);
26    p3.name = "popC";
27}
28
29 400 late() {
30    sim.addSubpopSplit("p4", 1000, p0);
31    p4.name = "popD";
32}
33
34 500 late() {
35    sim.addSubpopSplit("p5", 1000, p4);
36    p5.name = "popE";
37}
38
39 1000 late() {
40    sim.simulationFinished();
41}

```

`compile_mode()` generates these files on disk. A simplified equivalent code inside **slendr**:

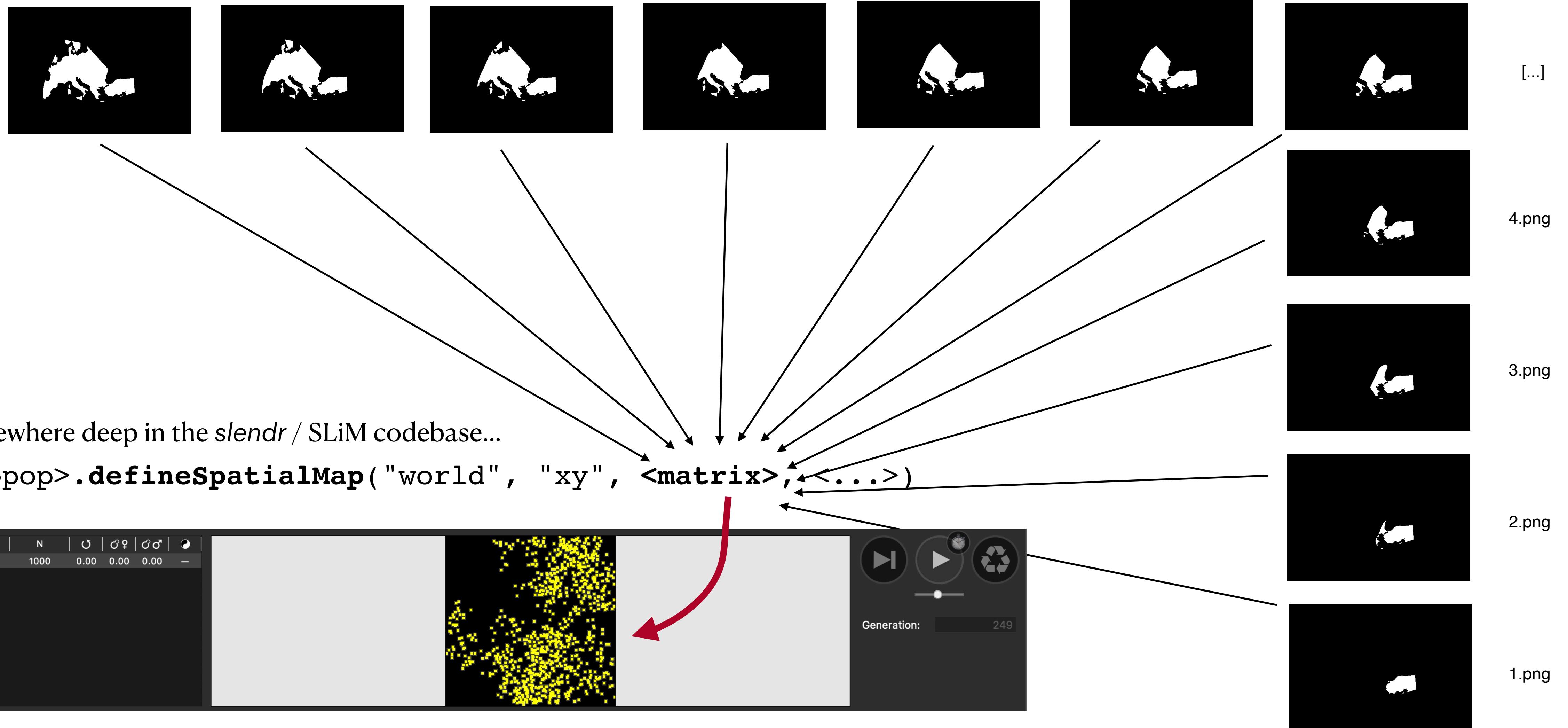


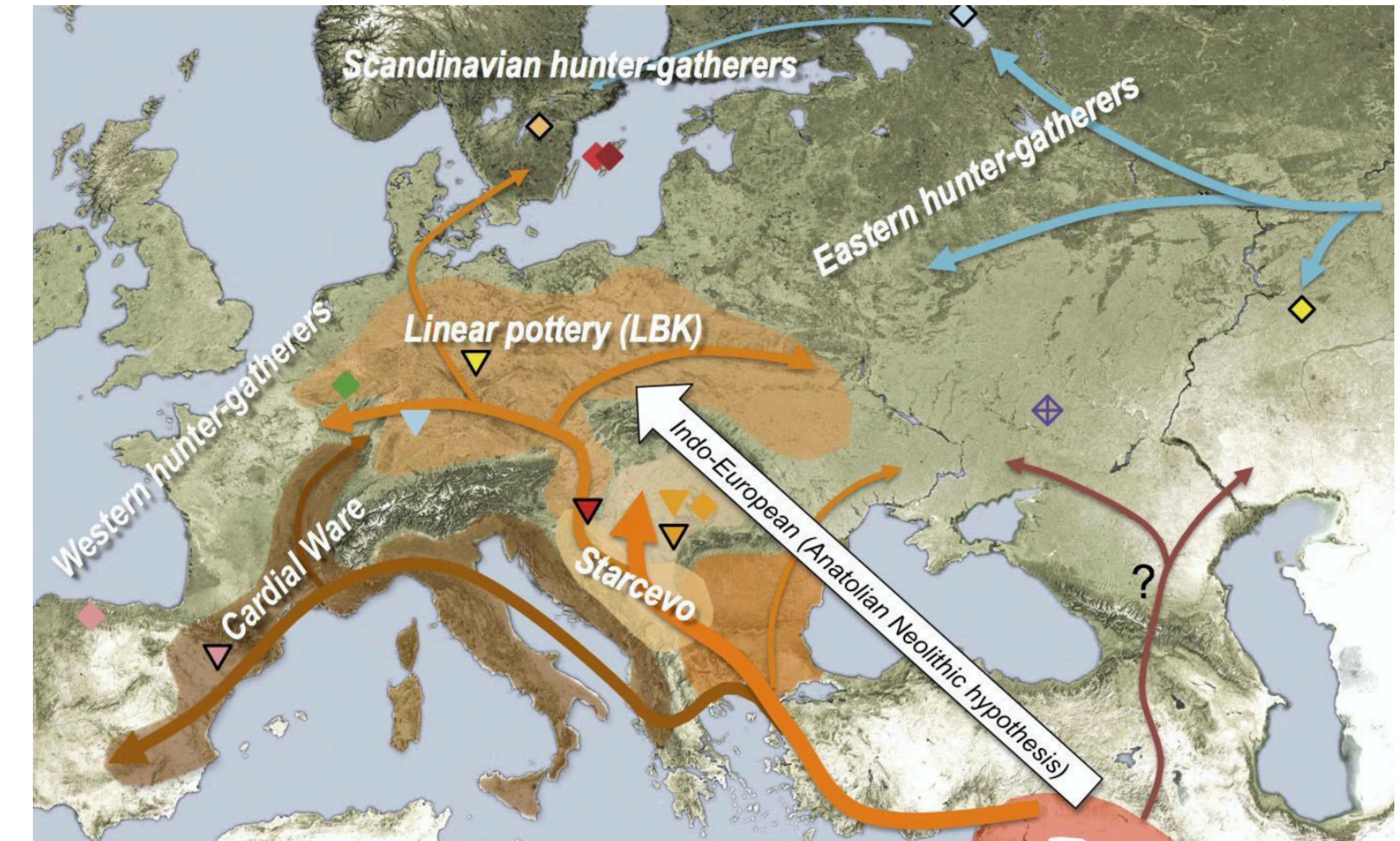
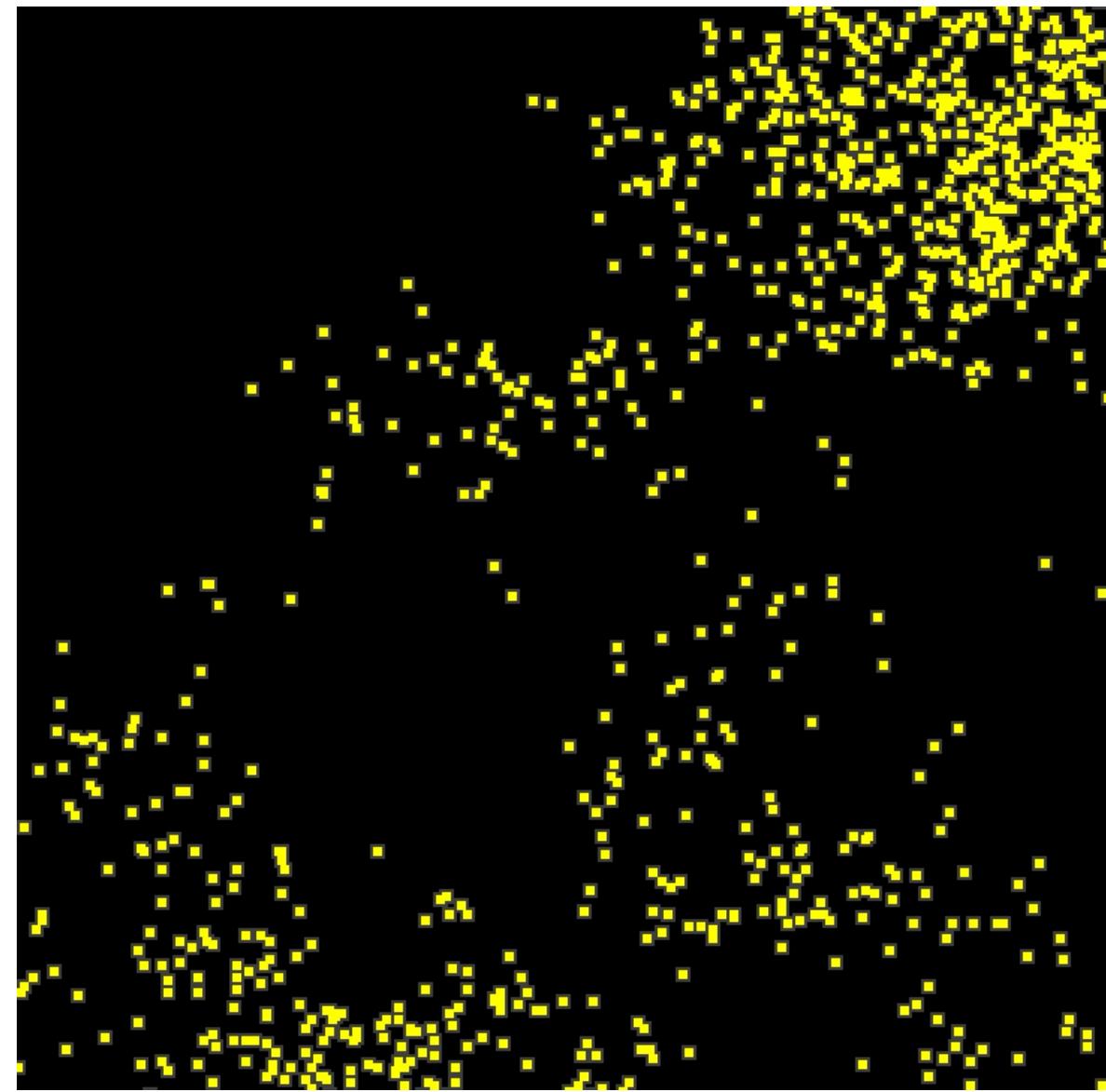
Vector → raster conversion during compilation



■ pixel →

□ pixel →





Somewhere deep in the *slendr* / SLiM codebase...

```
<subpop>.defineSpatialMap( "world", "xy", <matrix>, <...>)
```



Model definition & simulation

Model components

<code>population()</code>	Define a population and its original spatial range
<code>world()</code>	Define a world map for all spatial operations
<code>geneflow()</code>	Define a geneflow event
<code>region()</code>	Define a geographic region

Population dynamics

<code>move()</code>	Move the population to a new location in a given amount of time
<code>expand()</code>	Expand the population range
<code>shrink()</code>	Shrink the population range
<code>boundary()</code>	Update the population range
<code>resize()</code>	Resize the population size
<code>dispersal()</code>	Change dispersal parameters

Manipulation of spatial objects

<code>join()</code>	Merge two spatial <code>slendr</code> objects into one
<code>overlap()</code>	Generate the overlap of two <code>slendr</code> objects
<code>subtract()</code>	Generate the difference between two <code>slendr</code> objects
<code>reproject()</code>	Reproject coordinates between coordinate systems
<code>distance()</code>	Calculate the distance between a pair of spatial boundaries
<code>area()</code>	Calculate the area covered by the given <code>slendr</code> object
<code>dimensions()</code>	Return the dimensions of the world map

Model visualization and diagnostics

<code>plot(<slendr>)</code>	Plot <code>slendr</code> geographic features on a map
<code>plot_graph()</code>	Plot geneflow graph based on given model configuration
<code>plot_ancestry()</code>	Plot simulated ancestry proportions
<code>animate()</code>	Animate the simulated population dynamics
<code>explore()</code>	Open an interactive browser of the spatial model

Compiling and running spatial models

<code>compile()</code>	Compile the spatial demographic model
<code>sampling()</code>	Define sampling events at specified times for a given set of populations
<code>slim()</code>	Run a <code>slendr</code> model as a SLiM script
<code>read()</code>	Read a previously serialized model configuration
<code>script()</code>	Substitute variables in a template SLiM script

Tree-sequence analysis

Tree sequence loading and processing

<code>ts_load()</code>	Load a tree sequence file produced by a given model
<code>ts_recapitiate()</code>	Recapitiate the tree sequence
<code>ts_simplify()</code>	Simplify the tree sequence down to a given set of individuals
<code>ts_mutate()</code>	Add mutations to the given tree sequence
<code>ts_coalesced()</code>	Check that all trees in the tree sequence are fully coalesced

Extracting genotypes from tree sequences

<code>ts_genotypes()</code>	Extract genotype table from the tree sequence
<code>ts_eigenstrat()</code>	Extract genotypes from the tree sequence in the EIGENSTRAT format
<code>ts_vcf()</code>	Save genotypes from the tree sequence as a VCF file

Accessing tree sequence components

<code>ts_data()</code>	Extract combined annotated table of individuals and nodes
<code>ts_ancestors()</code>	Infer spatio-temporal ancestral history for given nodes/individuals
<code>ts_individuals()</code> <code>ts_edges()</code> <code>ts_nodes()</code>	Get the table of individuals/nodes/edges from the tree sequence
<code>ts_samples()</code>	Extract names and times of individuals scheduled for sampling
<code>ts_tree()</code>	Get a tree from a given tree sequence
<code>ts_draw()</code>	Plot a graphical representation of a single tree

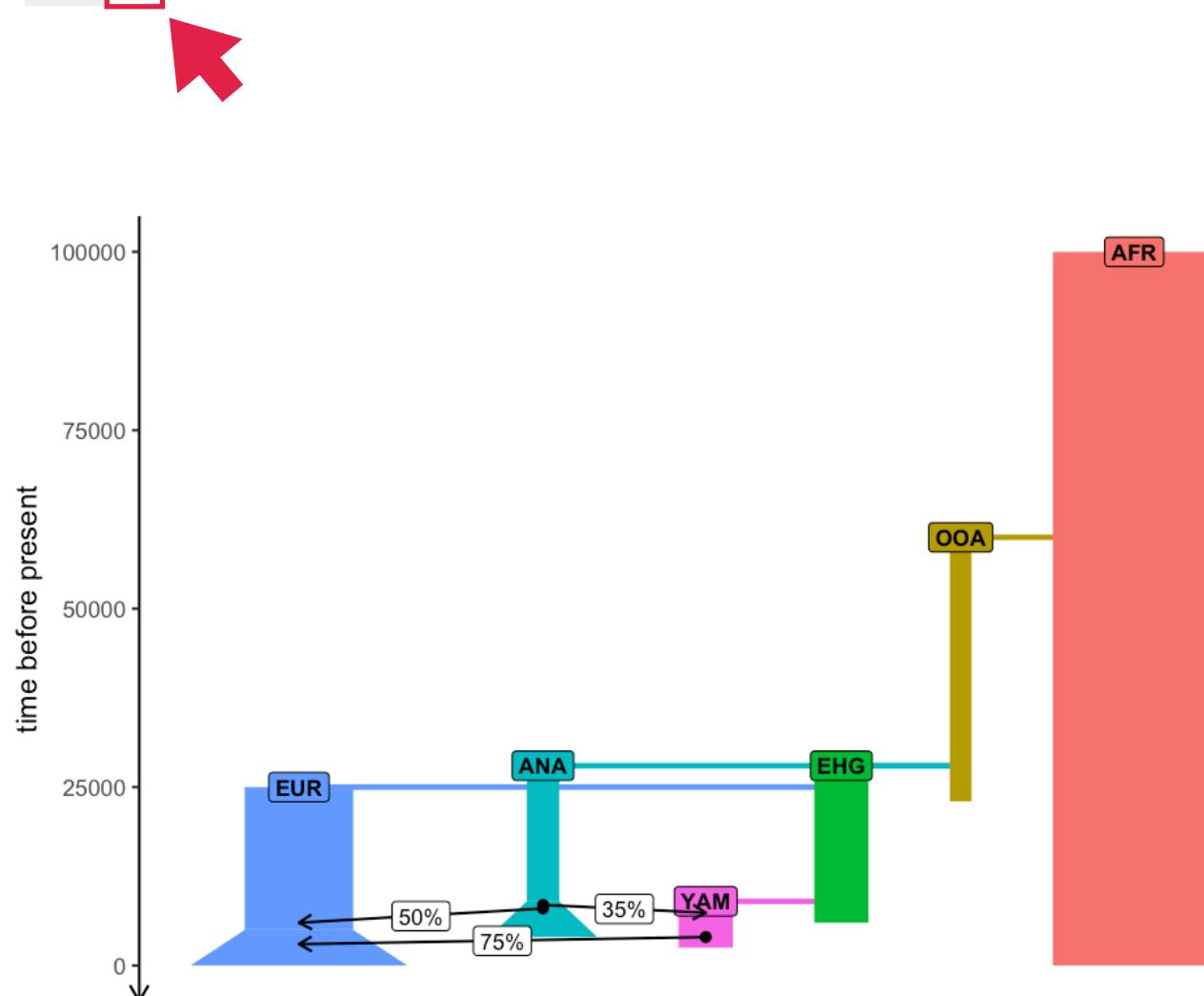
Tree sequence statistics

<code>ts_f2()</code> <code>ts_f3()</code> <code>ts_f4()</code> <code>ts_f4ratio()</code>	Calculate the f2, f3, f4, and f4-ratio statistics
<code>ts_afs()</code>	Compute the allele frequency spectrum (AFS)
<code>ts fst()</code> <code>ts_divergence()</code>	Calculate pairwise statistics between sets of individuals
<code>ts_diversity()</code>	Calculate diversity in given sets of individuals
<code>ts_tajima()</code>	Calculate Tajima's D for given sets of individuals
<code>ts_segregating()</code>	Calculate the density of segregating sites for the given sets of individuals

slendr's tree-sequence analysis library

(all internally backed by **tskit** and **pyslim**)

```
1 library(slendr); init_env()
2
3 afr <- population("AFR", time = 100000, N = 15000)
4 ooa <- population("OOA", parent = afr, time = 60000, N = 2000)
5 ehg <- population("EHG", parent = ooa, time = 28000, N = 5000)
6 eur <- population("EUR", parent = ehg, time = 25000, N = 10000) %>%
7   resize(N = 20000, how = "exponential", time = 5000, end = 0)
8 ana <- population("ANA", time = 28000, N = 3000, parent = ooa) %>%
9   resize(N = 10000, how = "exponential", time = 9000, end = 4000)
10 yam <- population("YAM", time = 9000, N = 5000, parent = ehg)
11
12 gf <- list(
13   gene_flow(from = ana, to = yam, rate = 0.35, start = 8500, end = 7400),
14   gene_flow(from = ana, to = eur, rate = 0.5, start = 8000, end = 6000),
15   gene_flow(from = yam, to = eur, rate = 0.75, start = 4000, end = 3000)
16 )
17
18 model <- compile_model(
19   populations = list(afr, ooa, ehg, eur, ana, yam),
20   gene_flow = gf, generation_time = 30
21 )
22
23 ts <- slim(model, sequence_length = 50e6, recombination_rate = 1e-8)
```



Recapitate the simulated tree sequence and overlay mutations:

```
ts <-
  ts %>%
  ts_recapitiate(Ne = 15000, recombination_rate = 1e-8) %>%
  ts_mutate(mutation_rate = 1e-8)
```

... do the same but also simplify to only a few defined individuals:

```
ts <-
  ts %>%
  ts_recapitiate(Ne = 15000, recombination_rate = 1e-8) %>%
  ts_mutate(mutation_rate = 1e-8) %>%
  ts_simplify(c("AFR_1", "OOA_1", "EHG_1", "ANA_1", "EUR_1", "YAM_1"))
```

Compute an f_3 admixture test:

```
> ts_f3(ts, A = "EUR_1", B = "EHG_1", C = "ANA_1")
# A tibble: 1 × 4
  A     B     C     f3
  <chr> <chr> <chr>   <dbl>
1 EUR_1 EHG_1 ANA_1 -0.0000104
```

Compute allele-frequency spectrum across all individuals:

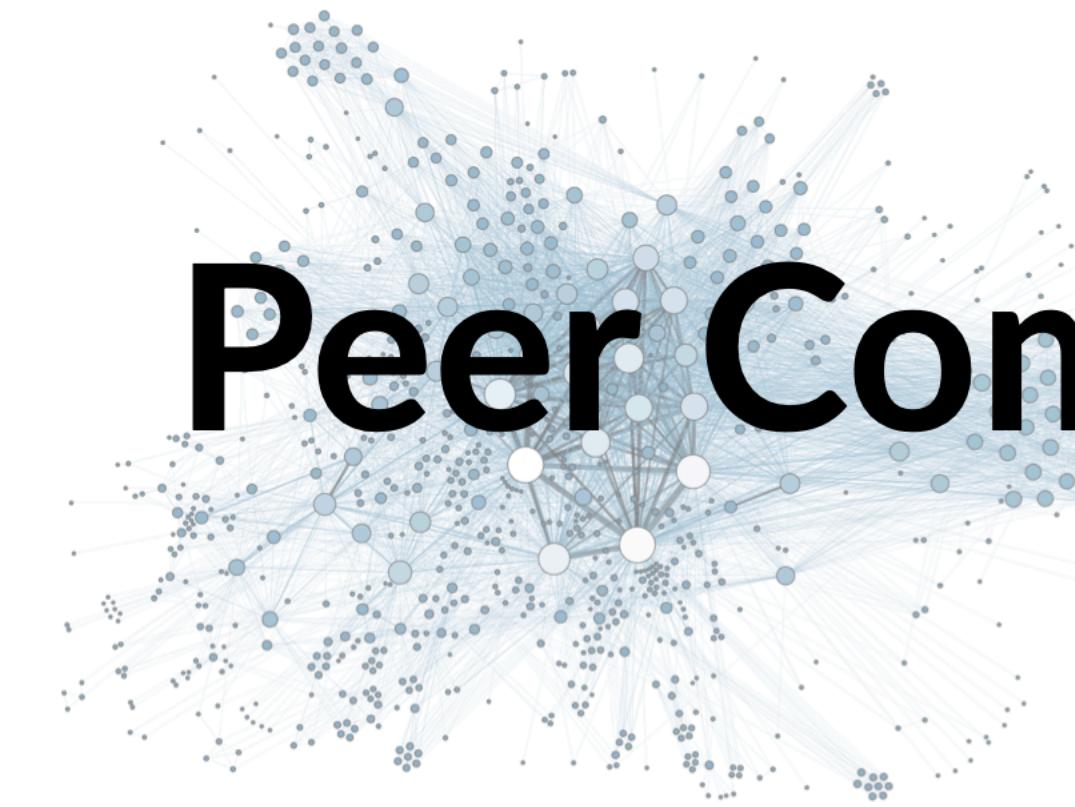
```
> ts_afs(ts)
[1] 0 25718 10869 8127 6638 5529 4678 4070 3749 3356 3335 3519
```

F_{ST} between AFR and EUR populations:

```
> ts fst(ts, list(afr = "AFR_1",
+                   eur = c("EHG_1", "ANA_1", "EUR_1")))
# A tibble: 1 × 3
  x     y     Fst
  <chr> <chr>  <dbl>
1 afr   eur   0.0764
```



www.slendr.net



PeerCommunity Journal

Section: Evolutionary Biology

RESEARCH ARTICLE
2023

slendr: a framework for spatio-temporal population genomic simulations on geographic landscapes

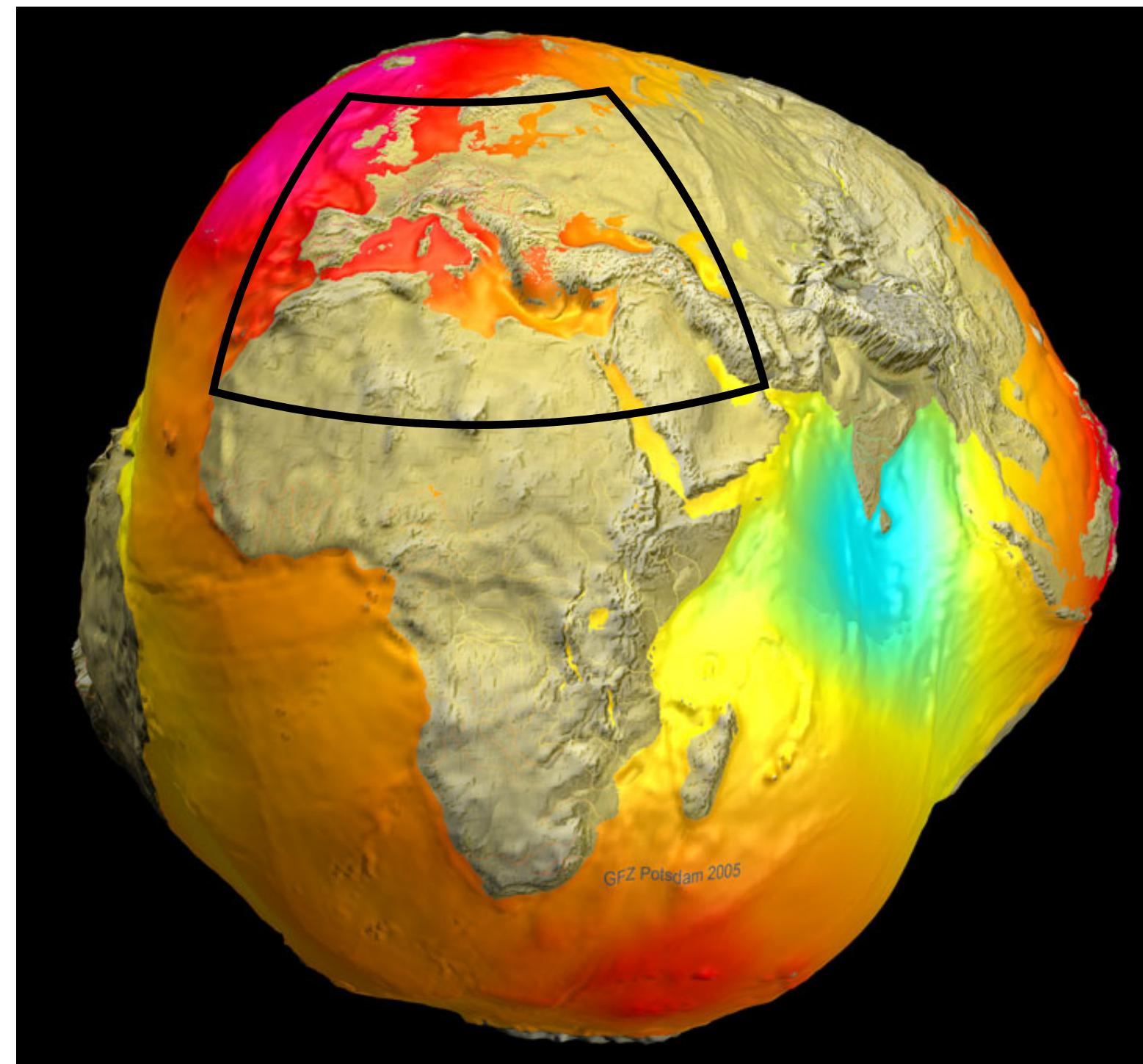
Martin Petr ,^{1,2} Benjamin C. Haller ,³ Peter L. Ralph ,⁴ and Fernando Racimo ,^{1,2}

Contact: mp@bodkan.net

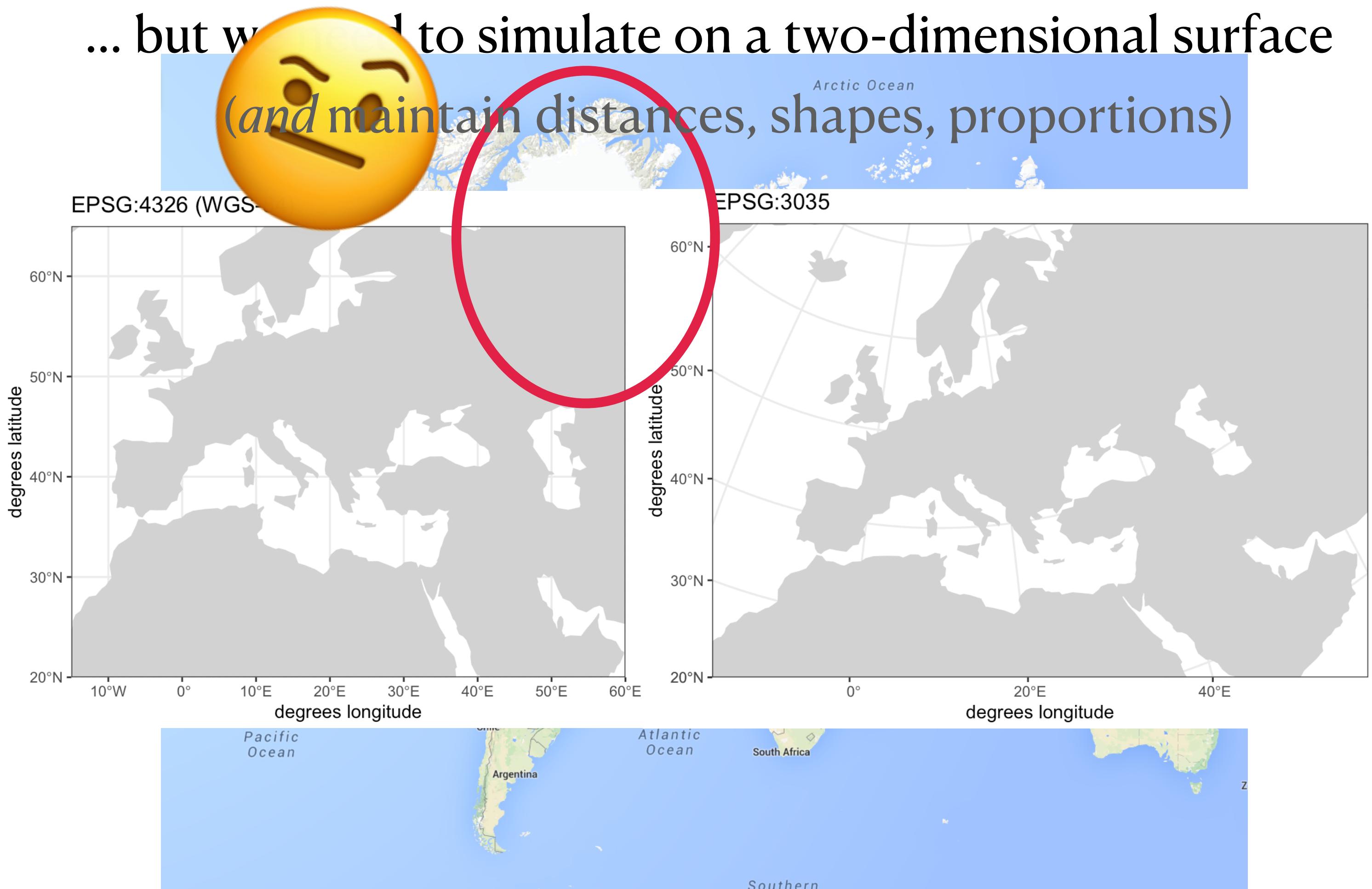
Slides will be at: <https://bodkan.net/talks>

Why an explicit *Coordinate Reference System (CRS)*?

Earth is a ~~potato~~ ~~sphere~~ ellipsoid...



... but we need to simulate on a two-dimensional surface
(and maintain distances, shapes, proportions)

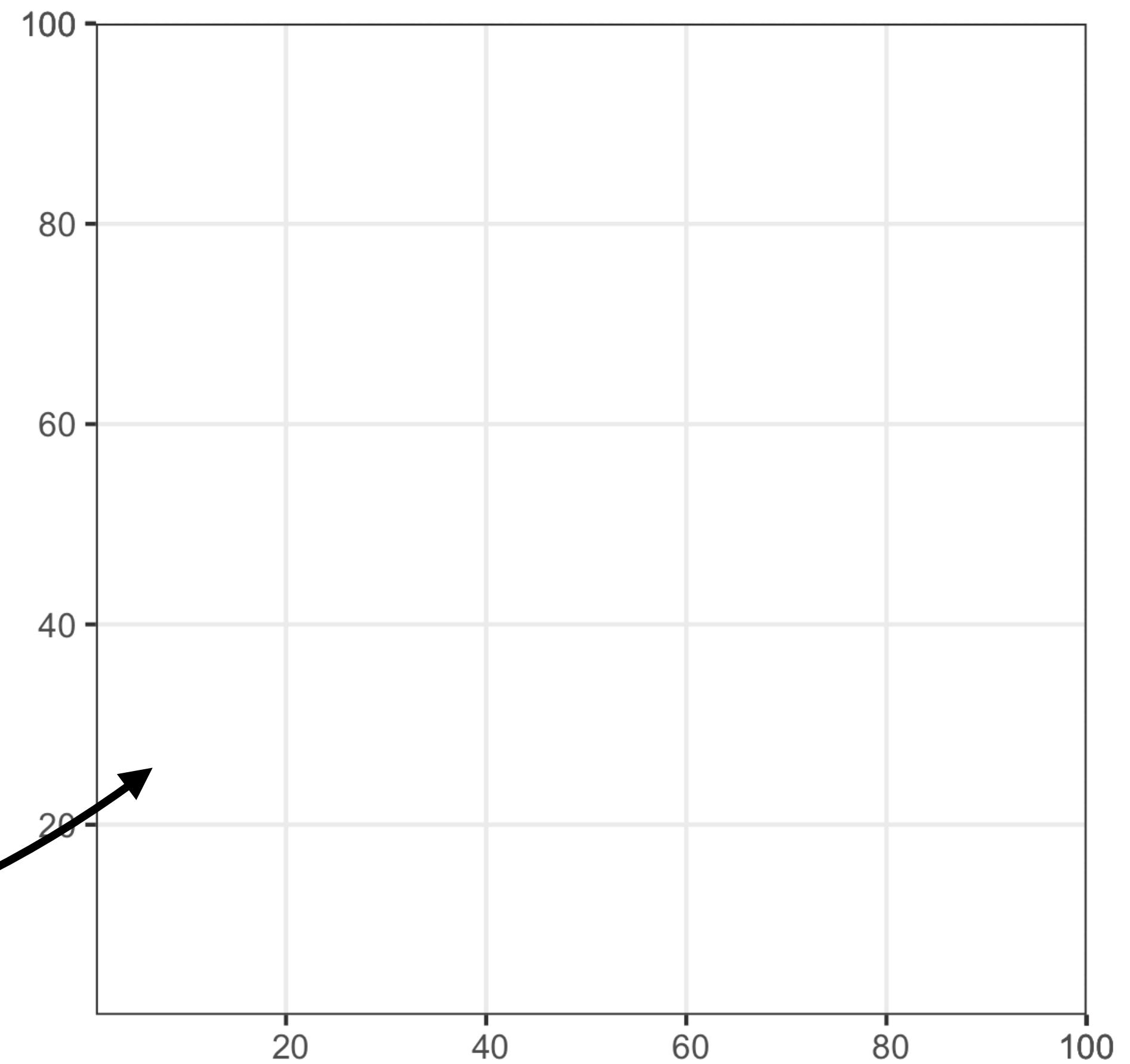


Astronomy Picture of the Day: 2014 December 15

1. Define a spatial context (“emptyworld()”)

(section of the world which will be occupied by populations)

```
map <- world(  
  xrange = c(1, 100),  
  yrange = c(1, 100),  
  landscape = "blank")
```



1. Define a spatial context (“custom world()”)

(section of the world which will be occupied by populations)

```
# three island regions...
island1 <- region(center = c(70, 50), radius = 15)
island2 <- region(center = c(40, 75), radius = 10)
island3 <- region(polygon = list(c(20, 20), c(40, 15), c(50, 20),
                                c(50, 35), c(40, 40), c(20, 35)))
# ... combined in a single map
map <- world(
  xrange = c(1, 100),
  yrange = c(1, 100),
  landscape = do.call(rbind, list(island1, island2, island3))
)
```

