



iDiv

German Centre for Integrative Biodiversity Research (iDiv)
Halle-Jena-Leipzig



HELMHOLTZ
Zentrum für Umweltforschung

Fast code

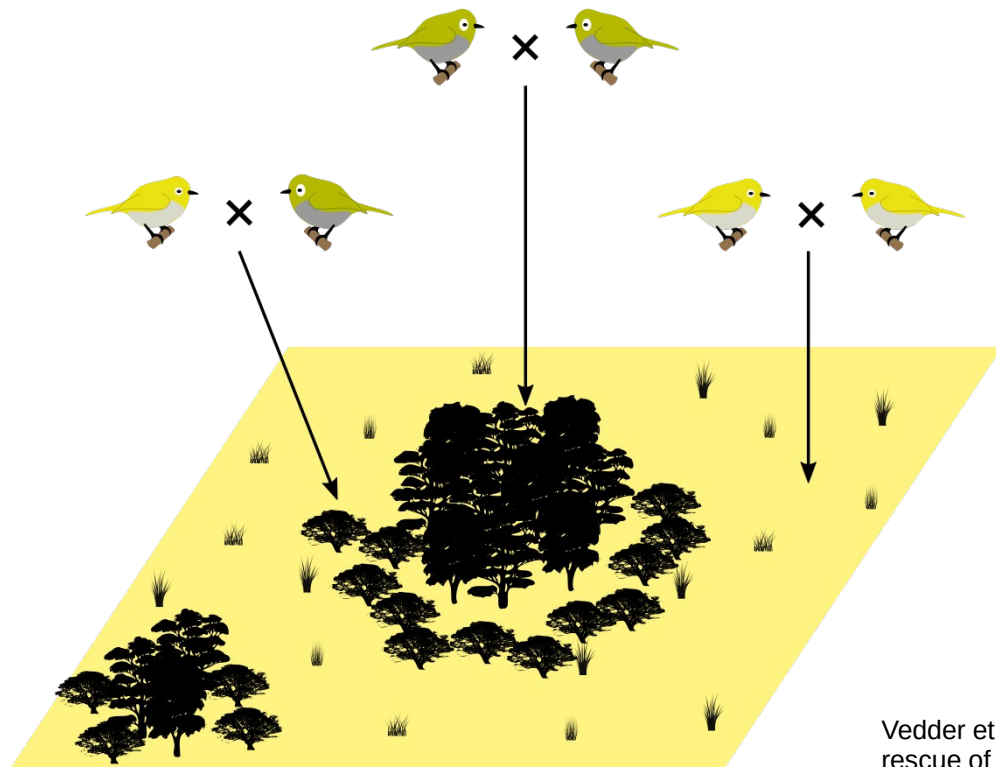
Performance optimisation in Julia

Daniel Vedder

iDiv is a research centre of the

DFG Deutsche
Forschungsgemeinschaft

Houston, we have a problem...



measured

350 000 individuals
x 500 years
= 17 hours
→ 10.3M IndYears/h

calculated

350 000 individuals
x 2M years
= 70 000 hours
→ 8 years!

Vedder et al. (2022). Hybridization may aid evolutionary rescue of an endangered East African passerine. *Evolutionary Applications*, 15(7), 1177–1188.

<https://doi.org/10.1111/eva.13440>

Before we start

“We should forget about small inefficiencies, say about 97% of the time; premature optimization is the root of all evil.”

Donald Knuth

Before we start

Reasons to optimise:

- Make a problem computationally tractable
- Speed up the modelling cycle
(*code change* → *simulation* → *data analysis* → *code change* ...)
- It's fun, and teaches you a lot

Reasons *not* to optimise:

- Optimisation takes time (lots!)
- Code readability suffers
- Outcome is uncertain, performance gains are hard to predict and sometimes counterintuitive

Measure twice, cut once

- Knuth (1971): 4% of the software consume >50% of the time
Boehm (1987): 20% of the software consume 80% of the time
- Do you want to optimise CPU or RAM usage (i.e. speed or memory)?
- Therefore:
 - get the program working *first*
 - then, profile to find the relevant bottleneck
 - *then* (and only then), optimise
 - finally, profile again to verify it worked

Profiling in Julia

```
1#!/usr/bin/env julia
2# A very thin wrapper to start a GeMM simulation. If you need something fancier,
3# have a look at `rungemmparallel.jl`, or import the GeMM module into your own
4# wrapper script.
5
6thisDir = joinpath(pwd(), "src")
7any(path -> path == thisDir, LOAD_PATH) || push!(LOAD_PATH, thisDir)
8using Pkg
9Pkg.activate(".")
10using GeMM
11
12using Profile
13Profile.clear()
14
15rm("results/taita_test", recursive=true, force=true)
16
17@profile rungemm("zosterops.config")
18
19open("profile_flat.txt", "w") do s
20    Profile.print(IOContext(s, :displaysize=>(300,145)), format=:flat, mincount=10, sortedby=:count)
21end
22
23open("profile_tree.txt", "w") do s
24    Profile.print(IOContext(s, :displaysize=>(300,300)), mincount=10)
25end
26
```

Profiling in Julia

1	Count	Overhead	File	Line	Function
2	=====	=====	=====	=====	=====
202	417	0	.../Studium/Masterarbeit/model/src/scheduling.jl	13	simulate! (::Array{GeMM.Patch,1}, ::Int64, ::Int64)
203	482	0	.../Studium/Masterarbeit/model/src/scheduling.jl	25	simulate! (::Array{GeMM.Patch,1}, ::Int64, ::Int64)
204	539	1	...beit/Studium/Masterarbeit/model/src/output.jl	357	(::GeMM.var"#logprint#44"{String,Bool})(::String, ::Bool)
205	545	0	...beit/Studium/Masterarbeit/model/src/output.jl	363	simlog (::String, ::Char, ::String, ::Bool)
206	548	9	...tudium/Masterarbeit/model/src/reproduction.jl	135	createoffspring (::Int64, ::GeMM.Individual, ::GeMM.Individual, ::Bool)
207	552	0	...beit/Studium/Masterarbeit/model/src/output.jl	350	logprint
208	556	0	...beit/Studium/Masterarbeit/model/src/output.jl	348	simlog
209	791	3	...t/Studium/Masterarbeit/model/src/zosterops.jl	251	zreproduce! (::GeMM.Patch)
210	850	0	@Base/io.jl	323	open
211	862	1	.../Studium/Masterarbeit/model/src/scheduling.jl	75	macro expansion
212	1096	0	@Base/threadingconstructs.jl	81	(::GeMM.var"#70#threadsfor_fun#117"{Array{GeMM.Patch,1}})(::Bool)
213	1098	0	@Base/threadingconstructs.jl	48	(::GeMM.var"#70#threadsfor_fun#117"{Array{GeMM.Patch,1}})(::Bool)
214	1129	0	...dium/Masterarbeit/model/src/run_simulation.jl	22	runsim (::String, ::Int64)
215	1256	0	@Base/timing.jl	174	macro expansion
216	1257	0	...dium/Masterarbeit/model/src/run_simulation.jl	36	rungemm (::String, ::Int64)
217	1257	0	@Base/client.jl	457	include (::String)
218	1257	0	@Base/boot.jl	331	eval (::Module, ::Any)
219	1257	0	.../usr/share/julia/stdlib/v1.5/REPL/src/REPL.jl	134	eval_user_input (::Any, ::REPL.REPLBackend)
220	1257	0	.../usr/share/julia/stdlib/v1.5/REPL/src/REPL.jl	195	repl_backend_loop (::REPL.REPLBackend)
221	1257	0	.../usr/share/julia/stdlib/v1.5/REPL/src/REPL.jl	180	start_repl_backend (::REPL.REPLBackend, ::Any)
222	1257	0	.../usr/share/julia/stdlib/v1.5/REPL/src/REPL.jl	292	run_repl (::REPL.AbstractREPL, ::Any; backend_on_current_task :: Bool)
223	1257	0	.../usr/share/julia/stdlib/v1.5/REPL/src/REPL.jl	288	run_repl (::REPL.AbstractREPL, ::Any)
224	1257	0	@Base/client.jl	399	(::Base.var"#807#809"{Bool,Bool,Bool,Bool})(::Module)
225	1257	0	@Base/essentials.jl	710	#invokelatest#1
226	1257	0	@Base/essentials.jl	709	invokelatest
227	1257	0	@Base/client.jl	383	run_main_repl (::Bool, ::Bool, ::Bool, ::Bool, ::Bool)
228	1257	0	@Base/client.jl	313	exec_options (::Base.JLOptions)
229	1257	0	@Base/client.jl	506	_start()
230	Total snapshots: 3491				

Profiling in Julia

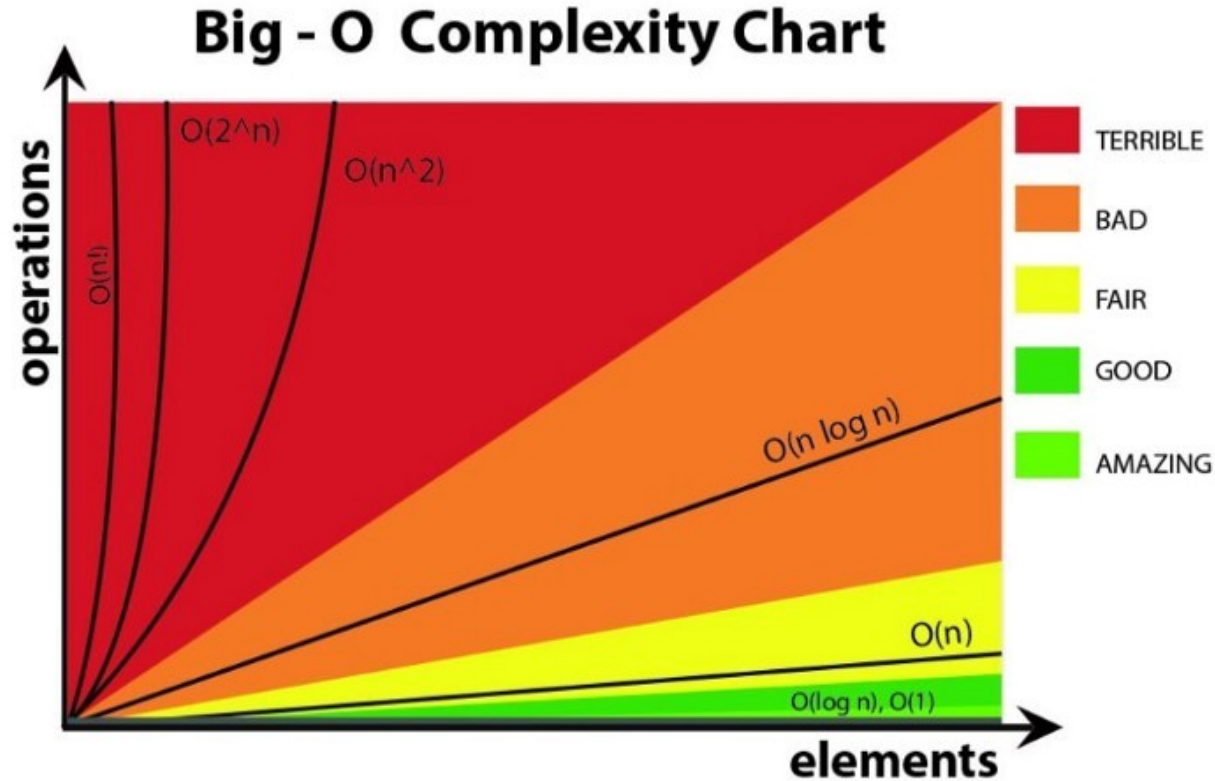
```
1 Overhead | [+additional indent] Count File:Line; Function
2 =====
```

```
245 37| 47 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/zosterops.jl:247; zreproduce! (::GeMM.Patch)
246 3| 791 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/zosterops.jl:251; zreproduce! (::GeMM.Patch)
247 | 66 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/reproduction.jl:130; createoffspring (::Int64, ::GeMM.Individual, ::GeMM.Individual, ::Bool)
248 | 12 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/genetics.jl:17; meiosis (::Array{GeMM.Chromosome,1}, ::Bool, ::String, ::Bool)
249 12| | 12 @Base/boot.jl:406; Array
250 8| 43 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/genetics.jl:21; meiosis (::Array{GeMM.Chromosome,1}, ::Bool, ::String, ::Bool)
251 29| | 29 @Base/Base.jl:33; getproperty
252 2| 55 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/reproduction.jl:131; createoffspring (::Int64, ::GeMM.Individual, ::GeMM.Individual, ::Bool)
253 4| 35 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/genetics.jl:21; meiosis (::Array{GeMM.Chromosome,1}, ::Bool, ::String, ::Bool)
254 17| | 17 @Base/Base.jl:33; getproperty
255 6| 16 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/reproduction.jl:134; createoffspring (::Int64, ::GeMM.Individual, ::GeMM.Individual, ::Bool)
256 9| 548 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/reproduction.jl:135; createoffspring (::Int64, ::GeMM.Individual, ::GeMM.Individual, ::Bool)
257 1| 44 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/genetics.jl:67; gettraitdictfast (::Array{GeMM.Chromosome,1}, ::Array{String,1})
258 | 43 @Base/dict.jl:90; Dict{String,Float64}()
259 36| | 36 @Base/boot.jl:406; Array
260 | 12 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/genetics.jl:69; gettraitdictfast (::Array{GeMM.Chromosome,1}, ::Array{String,1})
261 | 12 @Base/boot.jl:420; Array
262 12| | 12 @Base/boot.jl:406; Array
263 148| 341 /media/DATA/Arbeit/Studium/Masterarbeit/model/src/genetics.jl:71; gettraitdictfast (::Array{GeMM.Chromosome,1}, ::Array{String,1})
264 19| | 19 @Base/Base.jl:33; getproperty
265 154| 155 @Base/Base.jl:33; getproperty (::GeMM.Trait, ::Symbol)
266 12| 12 @Base/array.jl:809; getindex (::Array{GeMM.Trait,1}, ::Int64)
```


Runtime classes

- In order to get the result of a computation (e.g. sorting), how often must an operation (e.g. comparing two numbers) be performed on each element?
- How does the runtime depend on the number of elements n ?
 - $O(n)$ = worst-case runtime
 - $\Omega(n)$ = best-case runtime
 - $\Theta(n)$ = average-case runtime

Runtime classes



Runtime classes

Examples from IBM development:

- Find a random individual in a patch: $O(1)$
- Find the biggest individual in a patch: $O(n)$
- Find a suitable mate for all individuals in the patch: $O(n^2)$

Dealing with bad runtimes:

- Reduce the runtime class (e.g.: finding the biggest individual in a patch is $O(n)$, but doing so in a size-sorted patch is $O(1)$)
→ algorithms and data structures
- Do it less often (reduce n)

Code tuning

Some common miscreants:

- disk I/O
- nested loops
- array allocation

Strategies:

- buffer output
- structure data
- reduce allocations
- cache calculations
- parallelise
- compress

Buffer output

```
252 """
253     recordlineages(w)
254
255 Save the abundance of each lineage per patch. (Low-detail data recording function.)
256 """
257 function recordlineages(world::Array{Patch,1}, timestep::Int)
258     #XXX despite being low-detail, calling this frequently still means a lot of I/O
259     if !isfile(joinpath(setting("dest"), "lineages.log"))
260         simlog("t,X,Y,lineage,abundance,temp,prec", 'i', "lineages.log", true)
261     end
262     datastring = ""
263     for p in world
264         for l in unique(map(x -> x.lineage, p.community))
265             datastring *= string(timestep)*","*string(p.location[1])*","*string(p.location[2])*
266             ","*string(l)*","*string(length(findall(x -> x.lineage == l, p.community)))*","*
267             string(p.temp)*","*string(p.prec)*"\n"
268         end
269     end
270     simlog(datastring, 'i', "lineages.log", true)
271 end
```

first construct a complete string, then write it out

Structure data

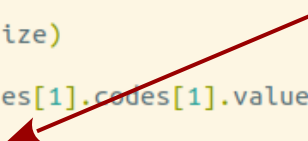
```
28 """
29     gettraitdict(chromosomes, traitnames)
30
31 Convert a genome (an array of chromosomes) into a dict of traits and their values.
32 """
33 function gettraitdict(chrms::Array{Chromosome, 1}, traitnames::Array{String, 1})
34     #TODO can this be made more efficient? It's called really often...
35     traitdict = Dict{String, Float64}()
36     traits = Array{Trait, 1}()
37     nchrms = 0
38     ngenes = 0
39     for chr in chrms
40         nchrms += 1
41         for gene in chr.genes
42             ngenes += 1
43             append!(traits, gene.codes)
44         end
45     end
46     for traitidx in eachindex(traitnames)
47         wantedtraits = skipmissing(map(x -> x.value, filter(x -> x.nameindex == traitidx, traits)))
48         traitdict[traitnames[traitidx]] = mean(wantedtraits)
49         traitdict[traitnames[traitidx] * "sd"] = std(wantedtraits)
50     end
51     traitdict["ngenes"] = ngenes
52     traitdict["nlnkgunits"] = nchrms
53     traitdict
54 end
```

$O(n) = n^2$

Structure data

```
56 """
57     gettraitdictfast(chromosomes, traitnames)
58
59 Convert a genome (an array of chromosomes) into a dict of traits and their values.
60 This is an optimised version that can be run if `degpleiotropy` is 0 and `linkage` is "none".
61 """
62 function gettraitdictfast(chrms::Array{Chromosome, 1}, traitnames::Array{String, 1})
63     # Makes use of the fact that with `degpleiotropy == 1` and `linkage == "none"`,
64     # there is exactly one trait per chromosome (one gene per chromosome and one trait per gene),
65     # and the chromosomes are arranged in trait-order.
66     genomesize = length(chrms)
67     traitdict = Dict{String, Float64}()
68     haploidlength = Int(genomesize/2)
69     values = Array{Float64,}(undef, genomesize)
70     for c in eachindex(chrms)
71         @inbounds values[c] = chrms[c].genes[1].codes[1].value
72     end
73     for traitidx in eachindex(traitnames)
74         wantedtraits = (values[traitidx], values[traitidx+haploidlength])
75         traitdict[traitnames[traitidx]] = mean(wantedtraits)
76         traitdict[traitnames[traitidx] * "sd"] = std(wantedtraits)
77     end
78     traitdict["ngenes"] = genomesize
79     traitdict["nlnkgunits"] = genomesize
80     traitdict
81 end
```

$O(n) = n$



Reduce allocations

```
129 """
130     zdisperse!(bird, world, location)
131
132 Dispersal of a single bird. Birds look patches with a suitable
133 habitat and a free territory or available mate. (Cf. Aben et al. 2016)
134 """
135 function zdisperse!(bird::Individual, world::Array{Patch,1}, location::Tuple{Int, Int})
136     # keep track of where we've been and calculate the max dispersal distance
137     x, y = location
138     route = [location]
139     !(bird.traits["dispshape"] > 0) && @goto failure # Logistics() requires  $\theta > \text{zero}(\theta)$ 
140     maxdist = rand(Logistic(bird.traits["dispmean"], bird.traits["dispshape"]))
141     while maxdist > 0
142         # calculate the best habitat patch in the surroundings (i.e. the closest to AGC optimum)
143         target = [(x-1, y-1), (x, y-1), (x+1, y-1),
144                 (x-1, y),      (x+1, y),
145                 (x-1, y+1), (x, y+1), (x+1, y+1)]
146         bestdest = nothing
147         bestfit = nothing
148         for t in target
149             (t in route) && continue
150             possdest = coordinate(t[1], t[2], world)
151             isnothing(possdest) && continue
152             patchfit = abs(possdest.prec - bird.traits["precopt"])
153             if isnothing(bestdest) || patchfit < bestfit
154                 bestdest, bestfit = possdest, patchfit
155             end
156         end
157         (isnothing(bestdest)) && @goto failure
158         x, y = bestdest.location
159         # check if the patch is within the bird's AGC range and has free space
160         if (bestfit <= bird.traits["prectol"] && length(bestdest.community) < setting("cellsize"))
161             # can we settle here?
```

array allocation

Reduce allocations

```
129 """
130     zdisperse!(bird, world, location)
131
132 Dispersal of a single bird. Birds look patches with a suitable
133 habitat and a free territory or available mate. (Cf. Aben et al. 2016)
134 """
135 function zdisperse!(bird::Individual, world::Array{Patch,1}, location::Tuple{Int, Int})
136     # keep track of where we've been and calculate the max dispersal distance
137     x, y = location
138     route = [location]
139     !(bird.traits["dispshape"] > 0) && @goto failure # Logistics() requires θ > zero(θ)
140     maxdist = rand(Logistic(bird.traits["dispmean"], bird.traits["dispshape"]))
141     while maxdist > 0
142         # calculate the best habitat patch in the surroundings (i.e. the closest to AGC optimum)
143         bestdest = nothing
144         bestfit = nothing
145         for xnew in (x-1):(x+1)
146             for ynew in (y-1):(y+1)
147                 ((xnew, ynew) in route) && continue
148                 possdest = coordinate(xnew, ynew, world) #XXX expensive?
149                 isnothing(possdest) && continue
150                 patchfit = abs(possdest.prec - bird.traits["precpot"])
151                 if isnothing(bestdest) || patchfit < bestfit
152                     bestdest, bestfit = possdest, patchfit
153                 end
154             end
155         end
156         (isnothing(bestdest)) && @goto failure
157         x, y = bestdest.location
158         # check if the patch is within the bird's AGC range and has free space
159         if (bestfit <= bird.traits["prectol"] && length(bestdest.community) < setting("cellsize"))
160             # can we settle here?
```

no array allocation
(but deeply nested loops)

Cache calculations

```
129 """
130     zdisperse!(bird, world, location)
131
132 Dispersal of a single bird. Birds look patches with a suitable
133 habitat and a free territory or available mate. (Cf. Aben et al. 2010)
134 """
135 function zdisperse!(bird::Individual, world::Array{Patch,1})
136     # keep track of where we've been and calculate the maximum distance
137     x, y = location
138     route = [location]
139     !(bird.traits["dispshape"] > 0) && @goto failure # Logically impossible
140     maxdist = rand(Logistic(bird.traits["dispmean"], bird.traits["dispmax"]))
141     while maxdist > 0
142         # calculate the best habitat patch in the surrounding area
143         bestdest = nothing
144         bestfit = nothing
145         for xnew in (x-1):(x+1)
146             for ynew in (y-1):(y+1)
147                 ((xnew, ynew) in route) && continue
148                 possdest = coordinate(xnew, ynew, world) #XXX expensive?
149                 isnothing(possdest) && continue
150                 patchfit = abs(possdest.prec - bird.traits["precopt"])
151                 if isnothing(bestdest) || patchfit < bestfit
152                     bestdest, bestfit = possdest, patchfit
153                 end
154             end
155         end
156         (isnothing(bestdest)) && @goto failure
157         x, y = bestdest.location
158         # check if the patch is within the bird's AGC range and has free space
159         if (bestfit <= bird.traits["prectol"] && length(bestdest.community) < setting("cellsize"))
160             # can we settle here?
```

```
259 let width = 0, height = 0
260     """
261         coordinate(x, y, world)
262
263     A utility function to perform a fast look-up for the patch at coordinate x/y.
264     Important: this assumes a rectangular world with coordinates in row-major order!
265     Returns the index of the desired patch.
266     """
267     global function coordinate(x::Int, y::Int, world::Array{Patch,1})
268         if iszero(width)
269             width = maximum(p -> p.location[1], world)
270             height = maximum(p -> p.location[2], world)
271         end
272         (x <= 0 || y <= 0 || x > width || y > height) && return
273         i = ((y-1) * width) + x
274         return i
275     end
276 end
277
```

Cache calculations

```
259 let width = 0, height = 0
260     ""
261     coordinate(x, y, world)
262
263     A utility function to perform a fast look-up for the patch at coordinate x/y.
264     Important: this assumes a rectangular world with coordinates in row-major order!
265     Returns the index of the desired patch.
266     ""
267     global function coordinate(x::Int, y::Int, world::Array{Patch,1})
268         if iszero(width)
269             width = maximum(p -> p.location[1], world)
270             height = maximum(p -> p.location[2], world)
271         end
272         (x <= 0 || y <= 0 || x > width || y > height) && return
273         i = ((y-1) * width) + x
274         return i
275     end
276 end
277
278 ""
279 findneighbours(world)
280
281 Construct a list of neighbours for each patch in the world, for faster lookup
282 later on. (Must be called during initialisation.)
283 ""
284 function findneighbours!(world::Array{Patch,1})
285     for patch in world
286         for x in (patch.location[1]-1):(patch.location[1]+1)
287             for y in (patch.location[2]-1):(patch.location[2]+1)
288                 ((x,y) == patch.location) && continue
289                 neighbour = coordinate(x,y,world)
290                 (isnothing(neighbour)) && continue
291                 push!(patch.neighbours, neighbour)
292             end
293         end
294     end
295 end
```

```
139 ""
140     zdisperse!(bird, patch, world)
141
142 Dispersal of a single bird. Birds look patches with a suitable
143 habitat and a free territory or available mate. (Cf. Aben et al. 2016)
144 ""
145 function zdisperse!(bird::Individual, patch::Patch, world::Array{Patch,1})
146     # keep track of where we've been and calculate the max dispersal distance
147     route = [patch.location]
148     !(bird.traits["dispshape"] > 0) && @goto failure # Logistics() requires 0 > zero(0)
149     maxdist = rand(Logistic(bird.traits["dispmean"], bird.traits["dispshape"]))
150     while maxdist > 0
151         # calculate the best habitat patch in the surroundings (i.e. the closest to AGC optimum)
152         bestdest = nothing
153         bestfit = nothing
154         for pid in patch.neighbours
155             neighbour = world[pid]
156             (neighbour.location in route) && continue
157             neighbourfit = abs(neighbour.prec - bird.traits["precopt"])
158             if isnothing(bestdest) || neighbourfit < bestfit
159                 bestdest, bestfit = neighbour, neighbourfit
160             end
161         end
162         (isnothing(bestdest)) && @goto failure
163         # check if the patch is within the bird's AGC range and has free space
164         if (bestfit <= bird.traits["prectol"] && length(bestdest.community) < bestdest.capacity)
165             # can we settle here?
166             partner = zfindmate(bestdest.community, bird)
167             if !isnothing(partner)
168                 # if we've found a partner
```

single loop

50% total runtime reduction!

Parallelise

```
66 """
67     zosteropsexperiment(world)
68
69 The annual update procedure for the Zosterops experiments, this time for bird populations.
70 """
71 function zosteropsexperiment(world::Array{Patch,1})
72     # The first four processes are patch-internal and can therefore be parallelised
73     # Note: multithreading requires calling Julia with the -p parameter
74     Threads.@threads for patch in world
75         establish!(patch)
76         survive!(patch)
77         zreproduce!(patch)
78         if setting("mutate")
79             mutate!(patch)
80         end
81     end
82     zdisperse!(world)
83 end
```

Compression

GeMM has a very high RAM usage

– how do we reduce this?

→ idea: compress genetic sequences

20 bytes

4 bytes

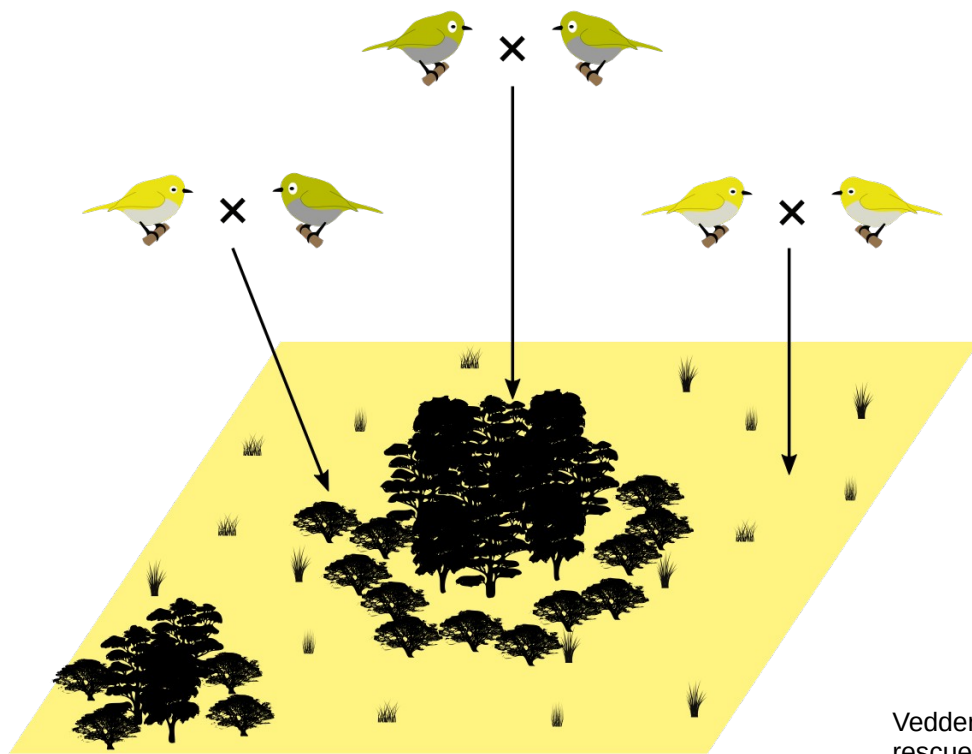
sequence: acgtaatgccccacatttggga

binary: 100101110111100100111
110101101101100101100
111111111110110100

decimal: 682176507920252852

```
171 """
172     seq2num(sequence)
173
174 Convert a DNA base sequence (a string) into binary and then into an integer.
175 This saves memory.
176 """
177 function seq2num(sequence::String)
178     num::Int64 = 0 # Int64 allows for max length of 21bp
179     for b in eachindex(sequence)
180         if sequence[end+1-b] == 'a'
181             num += 2^(3*(b-1)) * 4 # b'100'
182         elseif sequence[end+1-b] == 'c'
183             num += 2^(3*(b-1)) * 5 # b'101'
184         elseif sequence[end+1-b] == 'g'
185             num += 2^(3*(b-1)) * 6 # b'110'
186         elseif sequence[end+1-b] == 't'
187             num += 2^(3*(b-1)) * 7 # b'111'
188         end
189     end
190     num
191 end
192
193 """
194     num2seq(n)
195
196 Convert an integer into binary and then into a DNA base sequence string.
197 """
198 function num2seq(n::Integer)
199     bases = "acgt"
200     binary = string(n, base = 2)
201     sequence = ""
202     for i in 1:3:(length(binary) - 2)
203         sequence *= string(bases[parse{Int, binary[i:(i + 2)], base = 2} - 3]]
204     end
205     sequence
206 end
```


Houston, we've solved the problem!



measured

350 000 individuals
 x 500 years
 = 175 000 hours
 → → 154M Ind/h

calculated

350 000 individuals
 x 200 years
 = 70 000 hours
 → → 18 years!

Vedder et al. (2022). Hybridization may aid evolutionary rescue of an endangered East African passerine. *Evolutionary Applications*, 15(7), 1177–1188.

<https://doi.org/10.1111/eva.13440>

Thank you for your attention!

Any questions?