

promoting access to White Rose research papers



Universities of Leeds, Sheffield and York
<http://eprints.whiterose.ac.uk/>

This is an author produced version of a paper published in **Ecological Modelling**.

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/4129/>

Published paper

Parry, H.R. and Evans, A.J. (2008) *A comparative analysis of parallel processing and super-individual methods for improving the computational performance of a large individual-based model*, Ecological Modelling, Volume 214 (2-4), 141 – 152.

Manuscript Number: ECOMOD1606R1

Title: A comparative analysis of parallel processing and super-individual methods for improving the computational performance of a large individual-based model

Article Type: Research Paper

Section/Category:

Keywords: Agent-based modelling; Individual-based modelling; Parallel computing; Super-individuals

Corresponding Author: Dr Hazel Ruth Parry, Ph.D.

Corresponding Author's Institution: Central Science Laboratory

First Author: Hazel Ruth Parry, Ph.D.

Order of Authors: Hazel Ruth Parry, Ph.D.; Andrew J Evans, Ph. D.

Manuscript Region of Origin: UNITED KINGDOM

1 A comparative analysis of parallel processing
2 and super-individual methods for improving
3 the computational performance of a large
4 individual-based model

Hazel R. Parry ^{a,b,*}, Andrew J. Evans ^b

^a*Central Science Laboratory, Sand Hutton, York, YO41 1LZ, England*

Abstract

Individual-based modelling approaches are being used to simulate larger complex spatial systems in ecology and in other fields of research. Several novel model development issues now face researchers: in particular how to simulate large numbers of individuals with high levels of complexity, given finite computing resources. A case study of a spatially-explicit simulation of aphid population dynamics was used to assess two strategies for coping with a large number of individuals: the use of ‘super-individuals’ and parallel computing. Parallelisation of the model maintained the model structure and thus the simulation results were comparable to the original model. However, the super-individual implementation of the model caused significant changes to the model dynamics, both spatially and temporally. When super-individuals represented more than around 10 individuals it became evident that aggregate statistics generated from a super-individual model can hide more detailed deviations from an individual-level model. Improvements in memory use and model speed were perceived with both approaches. For the parallel approach, significant speed-up was only achieved when more than five processors were used and memory availability was only increased once five or more processors were used. The super-individual approach has potential to improve model speed and memory use dramatically, however this paper cautions the use of this approach for a density-dependent spatially-explicit model, unless individual variability is better taken into account.

Key words: Agent-based modelling, Individual-based modelling, Parallel computing, Super-individuals

* Corresponding author. Address: Central Science Laboratory, Sand Hutton, York,
YO41 1LZ, England. Tel.: +44 1904 462724; Fax: +44 1904 462111.
Email address: h.parry@cs1.gov.uk (Hazel R. Parry).

29 **1 Introduction**

30 A desire to better understand and inter-link the complex dynamic structures
31 of ecosystems, along with self-organisation, emergence of spatial and temporal
32 patterns and apparent unpredictability, has prompted a shift in the general
33 approach to ecological modelling today (Grimm and Railsback, 2005; Parrott
34 and Kok, 2000). Following trends in other fields of research, from social science
35 (Gilbert and Troitzsch, 1999) to fluvial sediment transport (Schmeeckle and
36 Nelson, 2003), there has been a shift away from procedural, equation-based
37 models to object-based simulations. These include individual-based models
38 (IBMs), cellular automata and multi-agent simulation (MAS). Such models
39 are concerned with modelling variation among individuals in a population,
40 and the interaction between individuals (DeAngelis and Gross, 1992; Grimm,
41 1999; Grimm and Railsback, 2005; Grimm et al., 1999; Huston et al., 1988;
42 Judson, 1994; Uchmański and Grimm, 1996). IBM is closely related to multi-
43 agent simulation. MAS has arisen from artificial intelligence (AI) research and
44 is used widely in other fields such as social science and computing (Gilbert
45 and Troitzsch, 1999).

46 Object-based approaches have been successfully implemented to model a range
47 of ecological systems (for a review see Grimm, 1999; Grimm and Railsback,
48 2005). They have the potential to further understanding of the local processes
49 that influence regional species population dynamics spatially and temporally,
50 enabling better understanding of how individual local-level and field-scale in-
51 teractions result in larger scale population distributions. However, some of
52 the potential of MAS and IBM methods is constrained by the demands that
53 may be placed on computing power. For realistic scenarios, it may be nec-

54 essary to simulate large numbers of individuals. There may also be added
55 complexity, such as in models where interactions or agent-density are impor-
56 tant, but populations are sparse (for example insect populations, Parry et al.,
57 2006b, 2004), or where agents are memory-heavy because they are complex
58 (e.g. forest dynamics, Verzelen et al., 2006), or multiple types of agent are
59 used, such as in models of competition or predator-prey models (e.g. Hos-
60 seini, 2006). Haefner (1992: pp.156-157), with some foresight, identified future
61 developments in ecological individual-based models that would benefit from
62 advanced computing as: multi-species models; models of large numbers of in-
63 dividuals within a population; models with greater realism in the behavioural
64 and physiological mechanisms of movement; and models of individuals with
65 ‘additional individual states’ (e.g. genetic variation).

66 The key limitations imposed by computer hardware are: (1) the number of
67 calculations that can be performed in a reasonable time (controlled by pro-
68 cessing power); (2) the number of agents that can be modelled (controlled
69 by memory). Relationships were determined between increasing numbers of
70 initial agents and the memory and simulation speed of a simple agent model
71 (described in section 2) run on a single 2.80 GHz Intel Xeon processor 2097 MB
72 RAM machine. Once the model is running, processor use of memory is nearly
73 linear and using an equation derived from the curve we can predict that at a
74 maximum available memory capacity of 1.5GB RAM on the single machine,
75 the theoretical limit to the initial number of agents is approximately 7,500,000.
76 However, at this limit, the simulation is calculated to take approximately 1
77 million seconds (12 days) to run (calculated from the simulation speed curve
78 using a quadratic function). This may be an under-estimate, as there is a
79 slight processing overhead for dealing with additional memory blocks, which

80 would result in a less linear relationship over time. The potential number of
81 replicates of a stochastic simulation are affected by this, so for example the
82 use of Monte Carlo techniques would no longer be possible.

83 There are a number of solutions to the problem of large numbers of individu-
84 als in an individual- or agent-based simulation (table 1). Solutions may range
85 from hardware investment (such as obtaining a more powerful computer) to
86 computational solutions, such as changes in the software design (e.g. paralleli-
87 sation) or changes in the the model structure. This paper evaluates and com-
88 pares two such solutions to this problem. The first is a computational solution
89 that requires access to networked hardware: to parallel program the model
90 software to work across a network of powerful computers, so splitting the pro-
91 cessing/data load. The second is a mathematical solution, where the model
92 itself is altered so that individuals are aggregated into ‘super-individuals’ (af-
93 ter Scheffer et al., 1995). The two methodologies are applied to a case study of
94 a spatially-explicit, individual-based simulation of aphid population dynam-
95 ics in agricultural landscapes (Parry et al., 2006b). The comparability of the
96 model results with the original model are first determined, then the two ap-
97 proaches are evaluated in terms of improved model efficiency (memory use
98 and speed).

99 A key advantage of parallel programming is that it maintains the strengths
100 of an individual-based approach whilst potentially increasing the number of
101 agents that can be simulated, as opposed to the super-individual approach
102 where the key interactions in the model are altered. However, parallelisation
103 is a complex solution, and although the agent interactions are unchanged sig-
104 nificant restructuring of the model software is needed. Haefner (1992) outlined
105 the potential applications of parallel computing to individual-based simula-

106 tions in ecology, but also pointed out the need for ecological modellers to
107 improve their technical knowledge. Few examples of parallel simulations exist
108 in the ecological literature to date. Some examples of note are a parallel sim-
109 ulation of a school of fish by Lorek and Sonnenschein (1995) and several in
110 relation to the ATLSS project (<http://atlss.org/>), which include a parallel
111 individual-based model of Everglades deer ecology by Abbott et al. (1997) and
112 a parallel spatially-explicit fish model (ALFISH) by Wang et al. (2004). Other
113 agent simulation examples can be found outside ecology in the use of parallel
114 agents for reducing genetic algorithm search times (Lefley and McKew, 2004)
115 and performing large scale traffic simulations (Dupuis and Chopard, 2001).

116 The simplicity of the super-individual approach makes it attractive, particu-
117 larly as it does not require complex programming and powerful computer sys-
118 tems to implement. It maintains the philosophy and integrity of an individual-
119 based approach without reverting to a population model to deal with large
120 numbers of individuals. However, implementations of this approach to date
121 are primarily not spatially-explicit.

122 **2 Application**

123 The results presented in this paper relate to a simplified version of a spatially
124 explicit individual-based simulation model of aphid population dynamics in
125 agricultural landscapes (Parry, 2006; Parry et al., 2006b, 2004). A spatially-
126 explicit IBM of aphid populations was constructed to assess the impact of
127 variation in agronomic practices in time and space. These practices included
128 crop introduction and configuration, pesticide spray application, matrix habi-
129 tat availability and fragmentation. The impacts that these have upon aphid

130 populations were observed, including both regional and local population dy-
131 namics as well as individual movement paths. A key limitation of the model
132 was the restriction on the number of insect agents that could be modelled.
133 The simplified model was used to explore and evaluate the options for cop-
134 ing with large numbers and complexity in the full model. The simple model
135 moves aphid agents (*Rhopalosiphum padi*) randomly from cell to cell around
136 a uniform landscape and local agent density is recorded. Aphids reproduce
137 parthogenetically with winged and non-winged morphs produced. Density de-
138 termines the proportion of alate (winged) morphs that are born at each iter-
139 ation. The simulation is begun with a population of alate agents originating
140 from a central cell in a 50×50 cell landscape, where each cell is 25×25 m. The
141 wind is set to a constant speed of 8 km h^{-1} and a constant westerly direction.
142 In the full version of the model there are a number of variables (some of which
143 are density dependent), realistic immigration across a region and a more com-
144 plex environment. The complexity of the full version of the model increases
145 computational demands beyond those demonstrated here and a large number
146 of agents (several million) were required for the simulation to be realistic at
147 the landscape scale.

148 Initial populations of 10, 100, 1,000, 10,000, 100,000 and 500,000 aphids were
149 used, originating from a single central cell. Each simulation was run thirty
150 times and an average taken to represent the total population trend over time
151 (as several parameters in the model are stochastic). While the model was
152 allowed to run for 120 days, spatial comparisons were made after 2, 20 and
153 40 days by creating surfaces that show the mean density in each cell over the
154 thirty runs. Temporal comparisons were made of the population dynamics at
155 the central cell.

156 **3 Parallel computing**

157 *3.1 Implementation*

158 In order to address computing problems where a model is hindered by data
159 requirements far larger than can be accommodated at any individual pro-
160 cessing element, parallel solutions are often implemented. The combined or
161 ‘virtual shared’ RAM of several computers is used to cope with the amount of
162 data and processing needed, using a Sequential-Algorithm Multiple-Data ap-
163 proach (SAMD), where the same algorithm is applied to different data items
164 on different processors (“nodes”). The scale of the problem for each individ-
165 ual computer is therefore reduced (often speeding the model up), or more
166 resources are made available (allowing for larger models). To parallelise the
167 model software, sections are run on each node and then the nodes periodi-
168 cally communicate together to share results. This requires somewhat complex
169 communication strategies to make the physically distributed systems act as a
170 single unit.

171 Key to an efficient parallel model is minimising the inter-processor communi-
172 cations; it is these that take valuable time (Pacheco, 1997). Because of this,
173 in models with static agents that only interact locally it makes sense to divide
174 the environment and agents between processors. Conversely, in models with
175 roaming agents, it makes sense to divide up the agents and have a copy of the
176 whole environment on each processor. Hardest to deal with are the situations
177 where agents roam the environment, but also interact with each other. In such
178 cases dividing up either the agents or the environment results in an increase
179 in inter-processor messages or agent transfer.

180 In this case study, parallelisation essentially splits the agents in the simulation
181 between a number of processors (nodes), each containing information on the
182 environment and total agent densities. Direct agent interaction (to determine
183 morphology) is mediated by density in the model. Thus, aphids can be split
184 between processors because it is not necessary for each processor to know the
185 exact position of all the aphids, just the density within each section of the
186 environment on the other machines. This information can be collated at a
187 single ‘control’ node and the total densities broadcast to each ‘worker’ node.

188 The initial model was created using the agent-based simulation toolkit Repast
189 (<http://repast.sourceforge.net>). The toolkit was implemented in paral-
190 lel by running the Repast interface on the control node (including the GUI
191 etc.), while the rest of the model code is run independently on worker-nodes,
192 synchronised by the control node using message passing (Parry, 2006; Parry
193 et al., 2006a). Agents are established on the worker nodes, coordinated by
194 the control node. In coding the parallel version of the model software, the
195 same code is placed on each processor but different sections of code are run
196 dependent on the node ID. The code on the control node controls the model
197 input, output and program flow, using the standard Repast methods of ‘setup’,
198 ‘buildModel’, ‘preStep’, ‘step’ and ‘postStep’ to structure the code and to ini-
199 tiate the simulation steps (figure 1). For example, when the method preStep
200 is run, the control node (node zero) is programmed to send out messages to
201 the other nodes to invoke agent methods associated with preStep (the model
202 is strongly synchronised). The updated agents pass density information back
203 to the control node when needed (figure 1). It was expected that the speed of
204 the simulation would increase with the number of nodes used, compensating
205 for any minor time delay caused by the timing control of the simulation from

206 the control node. A similar strategy was employed by Lorek and Sonnenschein
207 (1995) for a non-Repast simulation, which was found to increase simulation
208 speed as well as enable the size of the simulation to increase.

209 *3.1.1 Message passing*

210 The agent model was parallelised using a Message-Passing Interface (MPI)
211 for Java, MPIJava (<http://www.hpjava.org>), run on a 30-node distributed
212 memory parallel computer known as a Beowulf cluster. Message passing (MP)
213 is the principle manner by which Beowulf clusters are linked. MPIJava uses the
214 open-source native MPI ‘LAM’ (<http://www.lam-mpi.org/>). Further details
215 on the methods used to incorporate the MPI into the model are given in
216 Parry (2006); Parry et al. (2006a). The particular Beowulf cluster used for the
217 simulations presented here was a dedicated cluster of thirty machines (nodes),
218 where each node has dual 2.66 GHz Intel Xeon processors with 1280 MB of
219 DDR memory and 40GB 7200rpm internal IDE disks running over a switched
220 GB network. Although the results presented in this paper refer to simulations
221 conducted on a dedicated Beowulf cluster, the principles of parallelising the
222 model for a multi-core machine or a non-dedicated cluster would be very
223 similar. The model presented here has since been adapted to run on an Intranet
224 cluster of non-dedicated PCs and on a multi-core processor system without
225 re-coding the parallelisation, only altering the MPI commands in the code to
226 work with a customised MPI. Non-dedicated clusters of machines with mixed
227 specifications may however introduce problems of network unreliability and
228 performance bottlenecking on the slowest machines. Such issues are explored
229 further in Parry (in press).

230 3.1.2 *Data mapping and load balancing*

231 Data must be evenly distributed between the nodes, known as ‘load balancing’
232 (Pacheco, 1997). In this model, a balanced load was calculated using a form
233 of ‘block mapping’ (Pacheco, 1997). Agents were split evenly across the sys-
234 tem, whilst each node contained full environmental information. The agents
235 remained on the same node throughout the simulation, thus maintaining a
236 balanced load while being able to roam the environment. For all the parallel
237 simulations, it was found that the maximum memory used by each node was so
238 similar that 95% confidence limits derived from the standard error evaluated
239 to ± 0.00 in all cases. This shows that the distribution of individuals across
240 the worker nodes was highly efficient, and the load very well balanced.

241 4 Super-individuals

242 4.1 *Implementation*

243 The super-individual approach to modelling large populations on an individ-
244 ual basis was proposed by Scheffer et al. (1995), comparable to the earlier
245 ‘generalised individuals’ of Metz and de Roos (1992). A super-individual ap-
246 proach ‘allows zooming from a real individual-by-individual model to a cohort
247 representation or ultimately an all-animals-are-equal view without changing
248 the model formulation’ (Scheffer et al., 1995: pp. 161). The simple idea is that
249 individuals in a population can be grouped together into ‘super-individuals’,
250 thus reducing the number of objects to simulate and therefore reducing the
251 memory and processing power required (figure 2). For populations such as
252 aphids where there are high reproductive and mortality rates leading to large

253 juvenile populations, this approach can be very useful (Grimm and Railsback,
254 2005). It is possible to use the approach to test the effects of grouping individ-
255 uals, and also to examine the degree to which individual behaviour explains
256 the observed phenomena. A similar approach is used in physical models such
257 as the lattice models of fluid dynamics, particle modelling and Lagrangian
258 modelling (e.g. Woods and Barkmann, 1994).

259 *4.1.1 Combining individuals into a single super-individual*

260 Although Scheffer et al. (1995) state that no changes to the model formula-
261 tion are required for the super-individual approach, there are some significant
262 changes to the model structure that potentially influence the model results. To
263 convert the individual-based model to super-individuals, individuals originat-
264 ing at the same spatial location (cell) were split by initial age and morphology
265 (whether they have wings or not) into super-individuals. Each super-individual
266 represented a fixed number of individuals throughout the course of the simu-
267 lation.

268 *4.1.2 Adding individual immigrants to super-individuals*

269 Initial immigrants were added as super-individuals of the same scale factor
270 and, as in the unmodified model, these were of uniform age and morphology
271 (adult alates (i.e. winged)).

272 *4.1.3 Mortality of individuals/super-individuals*

273 Estimating the mortality of super-individuals can be done in a number of
274 ways, all of which are prone to error. The three main approaches are given by

275 Grimm and Railsback (2005: pp. 267) (figure 3):

276 N = the number of individuals represented by the super-individual (i.e. the
277 scale factor). N_0 = the number of individuals represented by the super-individual
278 at the start of the simulation.

- 279 (1) The number of super-individuals remains constant, and mortality reduces
280 N .
- 281 (2) N is kept relatively constant, by mortality reducing N until super-individuals
282 are recombined when N falls below $N_0/2$.
- 283 (3) Assume that an entire super-individual dies when subject to mortality.

284 Both approaches 1 and 2 require dynamic updating of the number of individ-
285 uals represented by the super-individual, but in this way they do maintain
286 more of the original variability of the model. However, significant errors, par-
287 ticularly spatial errors, would be introduced as individuals are re-grouped,
288 and the process would be computationally intensive. Reducing the number of
289 super-individuals in approaches 2 and 3 has computational advantages (the
290 number of super-individuals to iterate is minimised and individual variability
291 is less important so calculations are less complex).

292 Approach 3 was chosen: super-individuals are subject to the same probability
293 of mortality as individuals and when the super-individual dies all individuals
294 represented by the super-individual die. This approach was chosen because
295 the variability between individuals of the model (particularly age) meant that
296 approach 2 (recombining individuals) was problematic. Approach 1 (main-
297 taining a constant number of super-individuals) would also be problematic to
298 implement as the constant updating and variability of N would be computa-
299 tionally intensive, particularly as the density of individuals is important to a

300 number of model processes. Approach 3 was therefore considered to be the
301 most computationally efficient, although the least biologically realistic as it
302 suggests that mortality affects equally a group of aphids of uniform age and
303 morph in a particular cell (discretization of mortality). Another potential is-
304 sue with approach 3 is that it may require a lower value of N than the other
305 approaches to avoid excessive discretization of mortality. This paper assesses
306 whether this is the case.

307 4.1.4 *Changes to the model structure*

308 The construction of a super-individual simulation involved very little alter-
309 ation of the model structure (for details of this structure see Parry et al.,
310 2006b, 2004). A variable was added to record the number of individuals all
311 super-individuals actually represent. Equations that were dependent on den-
312 sity (such as morphology determination) were altered so that the density val-
313 ues were related to the real number of individuals in the simulation, not the
314 number of super-individuals (see equation 1). This was because the proportion
315 of alates produced is in relation to the density of individuals.

316 Morph determination is represented by the equation:

$$317 \quad ALPROP = \frac{0.002 + 0.991}{(1 + EXP(-0.076 \times (DENSITY - 67.416)))} \quad (1)$$

318 where ALPROP = the proportion of newly laid nymphs that will become alate
319 and DENSITY = the total number of individual aphids per plant.

320

321 5 Evaluation

322 The parallel version of the model produced extremely similar results to the
323 non-parallel model, as expected (no changes were made to the model structure,
324 only to the software). Variability between the original model and the parallel
325 model was only due to the model's stochasticity. However, the super-individual
326 model did alter the model structure, therefore some variation was expected
327 in the output between the super-individual model and the original model.
328 This variability is presented first below, then a comparison is made of the
329 improvement in performance in terms of model speed and memory use for
330 both the parallel and super-individual approach in relation to the original
331 model.

332 5.1 *Super-individual temporal and spatial replication of the individual-based* 333 *simulation*

334 Movement of super-individuals followed the same rules as that of individuals,
335 however this produced spatial clustering of the populations. To test the super-
336 individual model, populations of 100, 1,000, 10,000 and 100,000 and 500,000
337 individuals were represented by varying numbers of super-individuals (Table
338 2). Results are compared to the original individual-based model, both tem-
339 porally and spatially, in the following sections. Results from simulations with
340 10,000 individuals are given in more detail as an example to demonstrate the
341 effects of combining individuals.

342 5.1.1 *Temporal*

343 Overall, for simulations of fewer than 10,000 individuals the super-individual
344 simulations produced population densities that were much lower than the
345 individual-based model equivalent (figure 4). For 10,000 individuals, densi-
346 ties only become significantly lower at the second population peak, and the
347 super-individual simulations also reach this peak earlier. This can be related
348 to the spatial results (below), where it is only after this point in time that
349 it is evident that differing spatial distributions and densities are beginning to
350 emerge. The only case where the super-individual simulation falls within the
351 95% confidence limits of the original model for the duration of the simulation
352 period is the simulation of 10,000 individuals with 1,000 super-individuals
353 (scale factor 10), figure 4. The percentage error between the temporal results
354 for all the super-individual simulations and the individual-based simulations is
355 shown graphically in figure 5. This confirms that super-individual simulations
356 of 10,000 aphids and above with low scale factors may be acceptable. This
357 also shows that when a large number of individuals are represented by very
358 few super-individuals (in this case 10 super-individuals) the error is greatest.
359 Excessive discretization of mortality is therefore evident (suggested in section
360 4.1), resulting in a need to reduce the scale factor for results to better represent
361 the individual-based model.

362 5.1.2 *Spatial*

363 Clustering is evident in the spatial distribution. The super-individuals are
364 contained in fewer cells, closer to the origin, than the individual-based simu-
365 lation. This is illustrated for 10,000 individuals by figure 6. The distribution

366 better replicates the unmodified model when the number of super-individuals
367 is maximised and the individuals they represent minimised, due to the assump-
368 tion that when mortality occurs, the whole super-individual dies. Only when
369 the number of individuals within the super-individual (N) is minimised in a
370 large population of super-individuals can this be overcome (Grimm and Rails-
371 back, 2005). However, even when this is the case, for 10,000 individuals with
372 1,000 super-individuals (scale factor 10) (figure 4) this still does not produce
373 a similar spatial distribution pattern, despite giving a satisfactory temporal
374 result. This suggests that errors in spatial distribution may be hidden in super-
375 individual models validated temporally. The super-individual patterns are in
376 fact most comparable to the patterns of the individuals for the same number,
377 e.g. 10 super-individuals compares well with the distribution of 10 individuals,
378 the difference is the density at each cell. This is the expected result when the
379 local redistribution of (super)individuals is the main process determining the
380 spatial distribution, despite density affecting morphology.

381 *5.2 Speed*

382 Super-individuals always improve the model speed (figure 7). The speed im-
383 provement is enormous for the largest simulations, where 500,000 individuals
384 simulated with super-individuals using a scale factor of 100,000 increases the
385 model speed by over 500 times the original speed. However, it was shown above
386 that only large simulations with a low scale factor (10-100) may benefit from
387 the super-individual approach, thus for these scale factors an improvement
388 in model speed of approximately 10,000-30,000% (100-300 times) the original
389 speed would result for simulations of 100,000 to 500,000 individuals.

390 Adding more processors does not necessarily increase the model speed. Fig-
391 ure 7 shows that for simulations run on two nodes (one control node, one
392 worker node) the simulation takes longer to run in parallel compared to the
393 non-parallel model. Message passing time delay and the modified structure of
394 the code are responsible. As the number of nodes used increases, the speed
395 improvement depends on the number of agents simulated. The largest im-
396 provement in comparison to the non-parallel model is when more than 500,000
397 agents are run across twenty-five nodes, although the parallel model is slower
398 by comparison for lower numbers of individuals. However, when only five nodes
399 are used the relationship is more complex: for 100,000 agents five nodes are
400 faster than the non-parallel model, but for 500,000 the non-parallel model is
401 faster. This is perhaps due to the balance between communication time in-
402 creasing as the number of nodes increases versus the decrease in time expected
403 by increasing the number of nodes. Overall, these results seem to suggest that
404 when memory is sufficient on a single processor, it is unlikely to ever be effi-
405 cient to parallelise the code.

406 5.3 *Memory usage*

407 Super-individuals always reduce the memory requirements of the simulation
408 (figure 8). The memory requirements for a simulation of super-individuals has
409 a similar memory requirement to that of an individual-based simulation with
410 the same number of agents. For simulations of 100,000 agents this can reduce
411 the memory requirement to less than 10% of the memory required for the
412 individual-based simulation with a scale factor of 10,000, and for simulations
413 of 500,000 agents this may be reduced to around 1% with the same scale

414 factor.

415 The mean maximum memory usage by each worker node in the parallel simu-
416 lations is significantly lower than the non-parallel model, for simulations using
417 more than two nodes (figure 8). The two node simulation used more memory
418 on the worker node than the non-parallel model when the simulation had
419 100,000 agents or above. This is probably due to the memory saved due to the
420 separation of the GUI onto the control node being over-ridden by the slight
421 additional memory requirements introduced by the density calculations. How-
422 ever, when 5 and 25 nodes are used, the memory requirements on each node
423 are very much reduced, below that of the super-individual approach in some
424 cases. The super-individual approach uses the least memory for 500,000 indi-
425 viduals, apart from when only a scale factor of 10 is used (then the 25 node
426 parallel simulation is more memory efficient).

427 **6 Discussion**

428 The parallel model produced identical results to the initial model, as this
429 modifies only the model software and not the model itself. However, the super-
430 individual approach did not produce identical results to the initial model,
431 especially when assessed spatially. The similarity between the super-individual
432 results and the initial, unmodified model varied according the number of real
433 individuals that the super-individual was representing, and the number of
434 individuals simulated. The super-individual approach can only be considered
435 in situations where the number of individuals is high and the number of real
436 individuals represented by each super-individual is low (i.e. a low scale factor).

437 For the super-individual approach, within-cell density peaks vary temporally
438 between simulations run with different super-individual sizes. This is due to
439 the differences in emigration and movement patterns as a result of the size
440 of the super-individuals, as well as the method used to represent mortality.
441 Excessive discretization of mortality is evident as it is assumed that an en-
442 tire super-individual dies when subject to mortality. Further assessment of
443 the model (Parry, 2006) shows that regionally, the total population density is
444 similar between the different super-individual configurations and the unmodi-
445 fied model, but as shown in figure 6 there is a clear difference in the dispersal
446 patterns. Overall, the evidence indicates that the variability is such that the
447 super-individual approach is not suitable for the spatially-explicit simulation
448 of the aphid model, as presented here. Indeed, although the aphid model
449 is more strongly density dependent than most ecological models, most are
450 to some degree density dependent, rendering super-individual models prob-
451 lematic for spatially informative work. Modifications to the approach could
452 make it a possibility for future work. Experimentation with the other rules
453 for super-individual mortality suggested in section 4.1 would be a first step.
454 Other possible modifications include:

- 455 (1) Weighted kernels around a central ‘super-individual’, so that a more re-
456 alistic dispersal pattern is achieved.
- 457 (2) Relocation of a percentage of the super-individuals from a cell, without
458 actual population redistribution.
- 459 (3) Cell population model with individual migration.

460 However, re-distribution of individuals could significantly increase run-time,
461 adds complexity to the simulations and may take more memory than the
462 individual-based approach. This would also rely on a non-naturalistic model

463 of dispersion. Most movement in the model is a short distance each day, so
464 there will be constant shifting from super-individual to individual or creation
465 of dispersal kernels.

466 Further investigation may also indicate that spatial heterogeneity may have a
467 strong impact on the accuracy of the super-individual approach. The simula-
468 tions presented here were conducted in a neutral landscape, but if the model
469 were run in a heterogeneous landscape the interactions of the individuals with
470 the landscape may create model feedback that might further affect the accu-
471 racy of the super-individual results, both spatially and temporally.

472 Although the parallel solution appears to be more appropriate, in order to en-
473 sure it is optimised for agent simulations the balance between the advantage
474 of increasing the memory availability and the cost of communication between
475 nodes must be assessed in relation to the number of individuals simulated.
476 When the number of individuals is low, parallel simulations take longer (fig-
477 ure 7) and are less efficient (figure 8) than a non-parallel model run on a
478 single node. Increasing the number of nodes can reduce the demands on each
479 individual node, but time to communicate between processors may also be
480 increased (depending on the way in which the model is parallelised).

481 For the model presented here, estimates of the maximum number of agents
482 that can be simulated for varying numbers of nodes (table 3) and the maxi-
483 mum number of agents for a given super-individual scale factor (table 3) were
484 calculated with 1GB RAM, based upon information in figure 8. For the parallel
485 version, when only two nodes are used the non-parallel simulation is estimated
486 to have a higher maximum agent capacity per worker node, because space is
487 not being used by the message passing code. However, from five nodes and

488 up there is a higher maximum agent capacity for the parallel version than the
489 non-parallel model. The maximum agent capacity of 25 nodes is very high, at
490 nearly 100 million. This is approximately ten times the number of agents that
491 can be run within a reasonable time in the individual-based model. For low
492 numbers of nodes run times are huge: for two and five nodes the run times are
493 estimated to be 13 days and 47 days respectively. This would be expected to
494 increase with the complexity of the simulation. For any given model there will
495 be a threshold below which parallelisation is not efficient. Investigating this
496 threshold is likely to be a matter of iterative development as demonstrated
497 here, starting with a stripped-down model containing just the basic message-
498 passing elements. Passing of agents between nodes is processor intensive, and
499 therefore should be minimised. In this model, only the environment object
500 and information on the number of agents to create on each node are passed
501 from the control node to each of the nodes, and only density information is
502 returned to the control node for redistribution and display (see figure 1).

503 For the super-individuals (table 4) the relationship between the maximum
504 number of individuals that can be simulated and the scale factor is very simple.
505 The run time remains the same, as the maximum number of super-individuals
506 or individuals is constant. The maximum number of individuals that can be
507 simulated by this approach therefore depends purely on the scale factor used.
508 It would therefore not be unrealistic to assume this approach may potentially
509 enable the simulation of very large numbers of individuals indeed, in excess
510 of $7E^{11}$ if a scale factor of 1,000,000 is used, for example (assuming this scale
511 factor may be acceptable).

512 7 Conclusion

513 In order to address the limitations on the number of agents imposed by the
514 processing power and memory available, two solutions have been tested: par-
515 allel processing and super-individuals. The parallel approach involved signifi-
516 cant recoding of the model software, but no changes were made to the model
517 structure and performance increased significantly, enabling the simulation of
518 at least ten times more agents. The parallel model produced results that are
519 comparable to the initial, non-parallel model, leading to the conclusion that for
520 the simulation of very large populations the parallel model is a good solution.

521 Although initially far simpler to implement, the super-individual approach is
522 inappropriate for spatial simulations in the form presented here. However,
523 it may be possible to use this approach if the model were to be signifi-
524 cantly altered, by using another approach to simulate super-individual mor-
525 tality, or a super-individual model merged with an individual-based model,
526 where dispersal can be simulated by switching from a super-individual to an
527 individual-based model when necessary. There is a high risk that the complex-
528 ity of switching between model or implementing retrospective re-distribution
529 of agents could introduce significant error and put high demands on the pro-
530 cessor and/or memory, which are already limited. Overall, the results pre-
531 sented here indicate that the super-individual approach is inappropriate to
532 the spatially-explicit simulation of populations with density-dependent func-
533 tions or interactive agents, unless individual variability is better taken into
534 account. If this can be achieved satisfactorily, it has been demonstrated that
535 the super-individual approach may lead to very large reductions in computa-
536 tional demands.

537 8 Acknowledgements

538 This work is part of a PhD funded by the Central Science Laboratory (De-
539 fra Seedcorn). Thanks to Daniel Parry for his advice on Java programming
540 and Phil Northing and two anonymous reviewers for their comments on the
541 manuscript. The authors are members of the Multi-Agent Systems and Simu-
542 lation Research Group, School of Geography, University of Leeds.

543 References

- 544 Abbott, C. A., Berry, M. W., Comiskey, E. J., Gross, L. J., Luh, H.-K., 1997.
545 Parallel individual-based modeling of everglades deer ecology. *IEEE Com-
546 putational Science and Engineering* 4 (4), 60–78.
- 547 DeAngelis, D., Gross, L. (Eds.), 1992. *Individual-based models and approaches
548 in ecology: populations, communities and ecosystems*. Routledge, Chapman
549 and Hall, New York, 544 pp.
- 550 Dupuis, A., Chopard, B., 2001. *Parallel simulation of traffic in Geneva using
551 cellular automata*. Nova Science Publishers, Inc., Commack, NY, USA.
- 552 Gilbert, N., Troitzsch, K. G., 1999. *Simulation for the Social Scientist*. Open
553 University Press, Buckingham, 288 pp.
- 554 Grimm, V., 1999. Ten years of individual-based modelling in ecology: what
555 have we learned and what could we learn in the future? *Ecological Modelling*
556 115, 129–148.
- 557 Grimm, V., Railsback, S. F., 2005. *Individual-based Modeling and Ecology*.
558 *Princeton Series in Theoretical and Computational Biology*. Princeton Uni-
559 versity Press, Princeton, 480 pp.
- 560 Grimm, V., Wyszomirski, T., Aikman, D., Uchmański, J., 1999. Individual-

561 based modelling and ecological theory: synthesis of a workshop. *Ecological*
562 *Modelling* 115, 275–282.

563 Haefner, J. W., 1992. Parallel computers and individual-based models: An
564 overview. In: DeAngelis, D. L., Gross, L. J. (Eds.), *Individual-based mod-*
565 *els and approaches in ecology: populations, communities and ecosystems.*
566 *Routledge, Chapman and Hall, New York, Ch. 7, pp. 126–164.*

567 Hosseini, P. R., 2006. Pattern formation and individual-based models: The
568 importance of understanding individual-based movement. *Ecological Mod-*
569 *elling* 194 (4), 357–371.

570 Huston, M., DeAngelis, D., Post, W., 1988. New computer models unify eco-
571 logical theory. *BioScience* 38 (10), 682–691.

572 Judson, O., 1994. The rise of the individual-based model in ecology. *Trends in*
573 *Ecology and Evolution* 9, 9–14.

574 Lefley, M., McKew, I. D., 2004. Can a parallel agent approach to genetic algo-
575 rithms reduce search times. In: Loffi, A., Garobaldi, J. (Eds.), *Applications*
576 *and Science in Soft Computing.* New York: Cambridge University Press, pp.
577 69–74.

578 Lorek, H., Sonnenschein, M., 1995. Using parallel computers to simulate
579 individual-oriented models in ecology: a case study. In: *Proceedings: ESM*
580 *'95 European Simulation Multiconference, Prague, June 1995.*

581 Metz, J. A. J., de Roos, A. M., 1992. The role of physiologically structured
582 population models within a general individual based model perspective.
583 In: DeAngelis, D. L., Gross, L. J. (Eds.), *Individual Based Models and*
584 *Approaches in Ecology: Concepts and Models.* Routledge, Chapman and
585 Hall, New York, pp. 88–111.

586 Pacheco, P. S., 1997. *Parallel Programming with MPI.* Morgan Kauffman Pub-

587 lishers, San Francisco, CA.

588 Parrott, L., Kok, R., 2000. Incorporating complexity in ecosystem modelling.
589 Complexity International 7, 1–19.
590 URL <http://www.csu.edu.au/ci/>

591 Parry, H. R., in press. Agent Based Modelling, Large Scale Simulations.
592 In: Meyers, R. A. (Ed.) Encyclopedia of Complexity and System Science.
593 Springer-Verlag, GmbH Berlin Heidelberg.

594 Parry, H. R., October 2006. Effects of land management upon species popu-
595 lation dynamics: A spatially explicit, individual-based model. Ph.D. thesis,
596 University of Leeds.

597 Parry, H. R., Evans, A. J., Heppenstall, A. J., 2006a. Millions of agents: Par-
598 allel simulations with the Repast agent-based toolkit. In: Trappl, R. (Ed.),
599 Cybernetics and Systems 2006, Proceedings of the 18th European Meet-
600 ing on Cybernetics and Systems Research. Austrian Society for Cybernetic
601 Studies, Vienna.
602 URL <http://www.lintar.disco.unimib.it/ABModSim/>

603 Parry, H. R., Evans, A. J., Morgan, D., 2006b. Aphid population response to
604 agricultural landscape change: A spatially explicit, individual-based model.
605 Ecological Modelling 199, 451–463.

606 Parry, H. R., Evans, A. J., Morgan, D., June 2004. Aphid population dy-
607 namics in agricultural landscapes: An agent-based simulation model. In:
608 Pahl-Wostl, C., Schmidt, S., Jakeman, T. (Eds.), iEMSs 2004 International
609 Congress: “Complexity and Integrated Resources Management”. Interna-
610 tional Environmental Modelling and Software Society, Osnabrueck, Ger-
611 many.

612 Scheffer, M., Baveco, J. M., DeAngelis, D. L., Rose, K. A., van Nes, E. H.,
613 1995. Super-individuals: a simple solution for modelling large populations

614 on an individual basis. *Ecological Modelling* 80, 161–170.

615 Schmeckle, M. W., Nelson, J. M., 2003. Direct numerical simulation of bed-
616 load transport using a local, dynamic boundary condition. *Sedimentology*
617 50, 279–301.

618 Uchmański, J., Grimm, V., 1996. Individual-based modelling in ecology: what
619 makes the difference? *Trends in Ecology and Evolution* 11 (10), 437–441.

620 Verzelen, N., Picard, N., Gourlet-Fleury, S., 2006. Approximating spatial in-
621 teractions in a model of forest dynamics as a means of understanding spatial
622 patterns. *Ecological Complexity* 3 (3), 209–218.

623 Wang, D., Gross, L., Carr, E., Berry, M., 2004. Design and implementation of
624 a parallel fish model for South Florida. In: *Proceedings of the 37th Annual*
625 *Hawaii International Conference on System Sciences (HICSS'04)*.

626 URL [http://csdl2.computer.org/comp/proceedings/hicss/2004/
627 2056/09/205690282c.pdf](http://csdl2.computer.org/comp/proceedings/hicss/2004/2056/09/205690282c.pdf)

628 Woods, J. D., Barkmann, W., 1994. Simulating plankton ecosystem by the La-
629 grangian Ensemble Method. *Philosophical Transactions of the Royal Society*
630 *London B* 343, 27–31.

631 **List of Tables**

632	1	Possible solutions when faced with a large number of	
633		individuals to model.	31
634	2	Table to show the construction of the tested super-individuals:	
635		individuals, super-individuals and the number of individuals	
636		each super-individual represents	32
637	3	The maximum number of agents that can be simulated for 2,	
638		5 and 25 processors, and the associated estimated run time	33
639	4	The maximum number of agents that can be simulated	
640		when the super-individual scale factor (number of individuals	
641		represented by each super-individual) is 10, 100, 1,000, 10,000	
642		and 100,000 and the associated estimated run time.	34

Solution	Pros	Cons
Invest in an extremely powerful computer.	No changes to model code or structure.	High cost.
Invest in an extremely powerful computer network and reprogram the model in parallel.	Makes available high levels of memory and processing power. Model remains the same.	High cost. Advanced computing skills required for restructuring of model software.
Super-individuals	Relatively simple solution. Little change to model formulation.	Reprogramming of model and altered structure and interactions. Untested in spatial context.
Reduce the number of individuals in order for model to run.	No reprogramming of model.	Unrealistic population. Alters model behaviour.
Revert to a population based modelling approach.	Could potentially handle any number of individuals.	Lose insights from IBM. Potentially unsuitable for the particular research questions of the study. Construction of entirely new model.

Table 1

Possible solutions when faced with a large number of individuals to model.

Number of individuals	Number of super-individuals	Number of individuals represented by each super-individual ('scale factor')
100	10	10
1,000	10	100
	100	10
10,000	10	1,000
	100	100
	1,000	10
100,000	10	10,000
	100	1,000
	1,000	100
	10,000	10
500,000	50	10,000
	500	1,000
	5,000	100
	50,000	10

Table 2

Table to show the construction of the tested super-individuals: individuals, super-individuals and the number of individuals each super-individual represents

Number of Nodes	Maximum number of agents	Estimated run time of simulation (seconds)
IBM	$7.49E^6$	$1.12E^6$
2	$3.33E^6$	$1.11E^6$
5	$1.43E^7$	$4.06E^6$
25	$1.00E^8$	$2.20E^4$

Table 3

The maximum number of agents that can be simulated for 2, 5 and 25 processors, and the associated estimated run time

Scale factor	Maximum number of individuals	Maximum number of super-individuals (agents)	Estimated run time of simulation (seconds)
IBM	$7.49E^6$	-	$1.12E^6$
10	$7.49E^7$	$7.49E^6$	$1.12E^6$
100	$7.49E^8$	$7.49E^6$	$1.12E^6$
1,000	$7.49E^9$	$7.49E^6$	$1.12E^6$
10,000	$7.49E^{10}$	$7.49E^6$	$1.12E^6$
100,000	$7.49E^{11}$	$7.49E^6$	$1.12E^6$

Table 4

The maximum number of agents that can be simulated when the super-individual scale factor (number of individuals represented by each super-individual) is 10, 100, 1,000, 10,000 and 100,000 and the associated estimated run time.

643 **List of Figures**

644	1	Flow chart illustrating the operation of rules at each stage of a	
645		model run for a simple Repast model, and the role of message	
646		passing to control the program flow between node 0 and the	
647		other nodes.	36
648	2	Super-individuals: Grouping of individuals into single objects	
649		that represent the collective	37
650	3	The three main approaches to estimating the mortality of	
651		super-individuals: (a) The number of super-individuals remains	
652		constant, and mortality reduces the number of individuals (N)	
653		represented by the super-individual. (b) N is kept relatively	
654		constant, by mortality reducing N then super-individuals are	
655		recombined when N falls below $N_0/2$. (c) Assume that an	
656		entire super-individual dies when subject to mortality.	38
657	4	10,000 individuals: comparison between individual-based	
658		simulation, 1,000 super-individual simulation (each represents	
659		10 individuals), 100 super-individual simulation (each	
660		represents 100 individuals) and 10 super-individual simulation	
661		(each represents 1,000 individuals), showing 95% confidence	
662		limits derived from the standard error.	39
663	5	Comparison of the mean (absolute) percentage error between	
664		the super-individual simulations and the individual-based	
665		simulation, at $t = 40$.	40
666	6	Spatial density distributions for individual-based versus	
667		super-individual simulations (10,000 aphids) at (a) 2 days	
668		(b) 20 days and (c) 40 days. The distribution further from	
669		the central cell is influenced by the constant westerly wind	
670		direction to result in a linear movement pattern.	41
671	7	Plot of the percentage speed up from the individual-based	
672		(non-parallel) model against number of agents modelled:	
673		comparison between parallel simulations using 2, 5 and 25	
674		nodes and super-individuals of scale factor 10, 100, 1,000,	
675		10,000, 100,000 and 500,000	42
676	8	Plot of the mean maximum memory used in a simulation	
677		run against number of agents for the model, for different	
678		numbers of nodes (memory per node) and scale factors for	
679		super-individuals	43

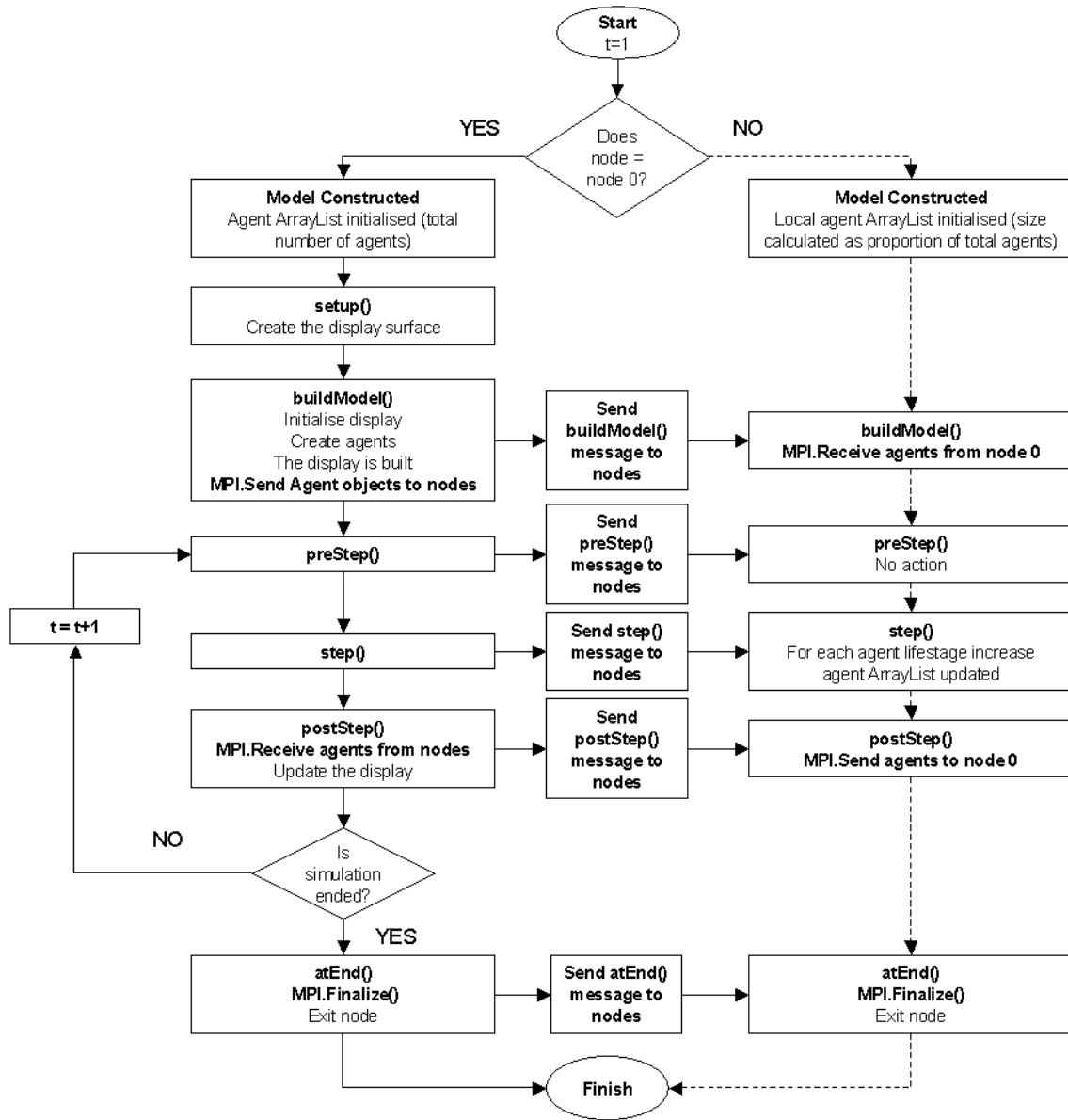


Fig. 1. Flow chart illustrating the operation of rules at each stage of a model run for a simple Repast model, and the role of message passing to control the program flow between node 0 and the other nodes.

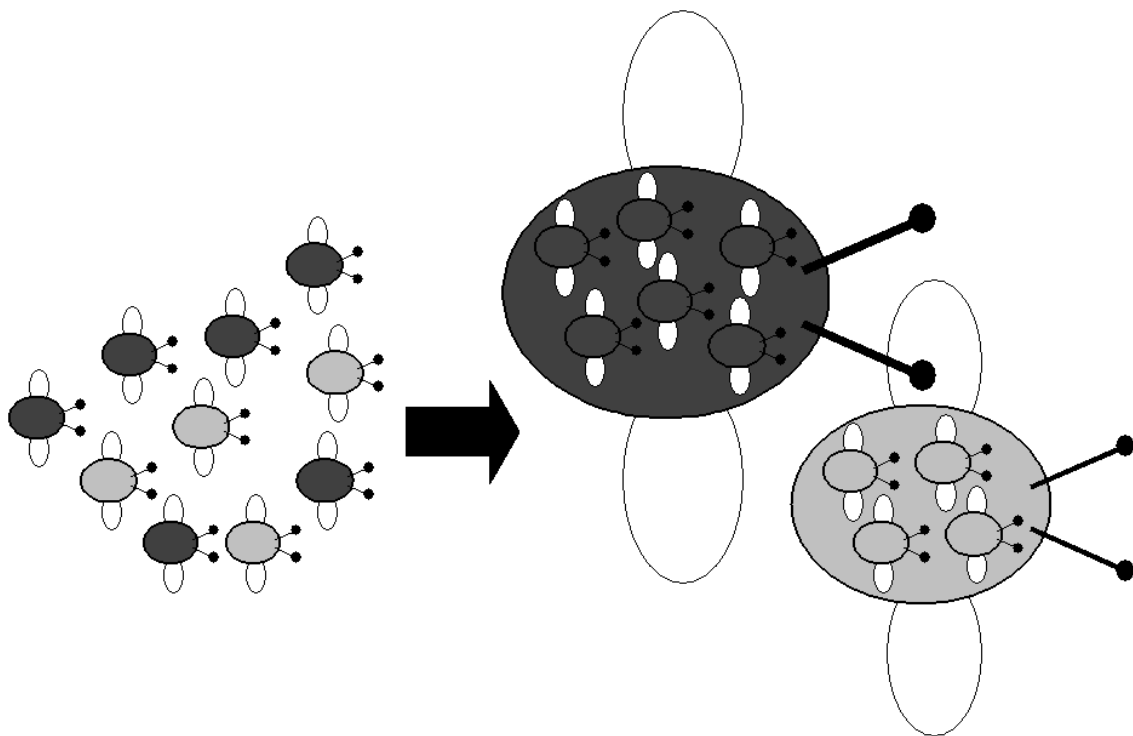


Fig. 2. Super-individuals: Grouping of individuals into single objects that represent the collective

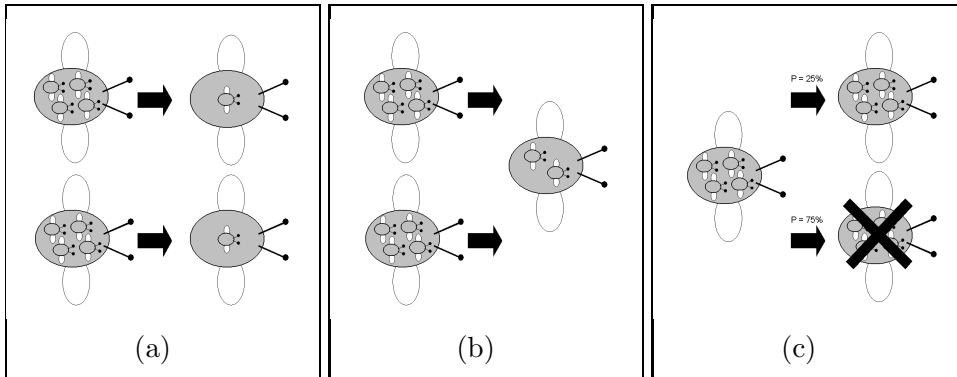


Fig. 3. The three main approaches to estimating the mortality of super-individuals: (a) The number of super-individuals remains constant, and mortality reduces the number of individuals (N) represented by the super-individual. (b) N is kept relatively constant, by mortality reducing N then super-individuals are recombined when N falls below $N_0/2$. (c) Assume that an entire super-individual dies when subject to mortality.

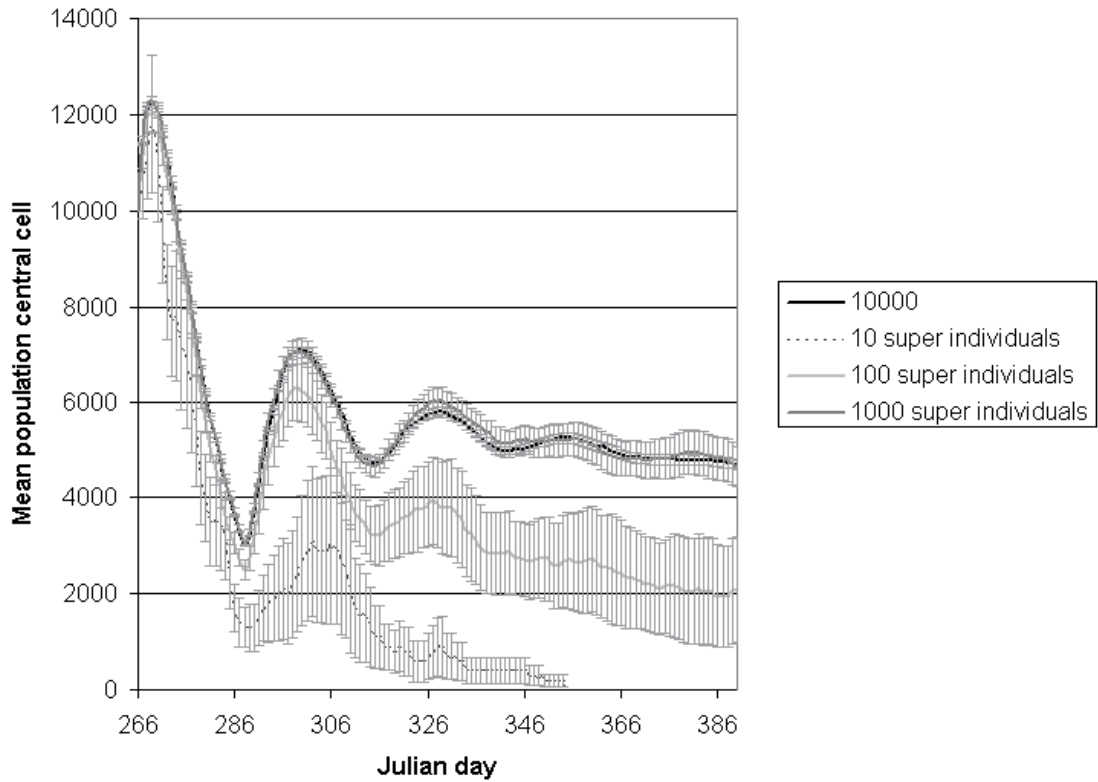


Fig. 4. 10,000 individuals: comparison between individual-based simulation, 1,000 super-individual simulation (each represents 10 individuals), 100 super-individual simulation (each represents 100 individuals) and 10 super-individual simulation (each represents 1,000 individuals), showing 95% confidence limits derived from the standard error.

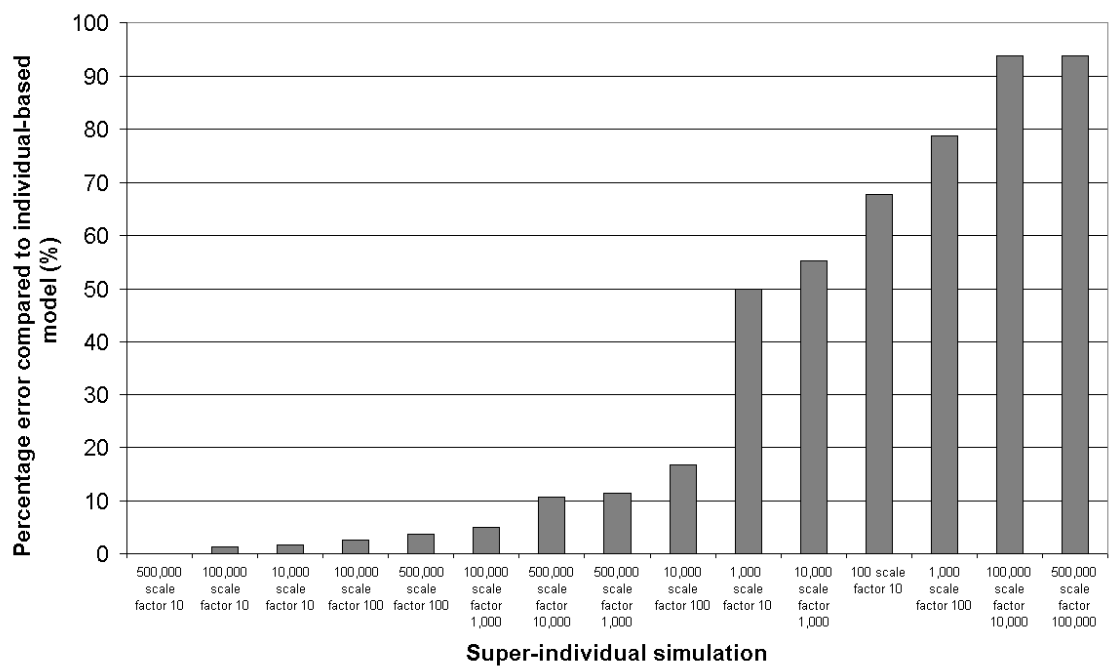
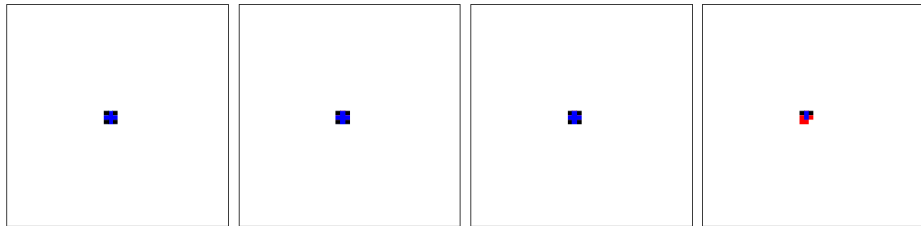
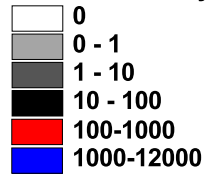
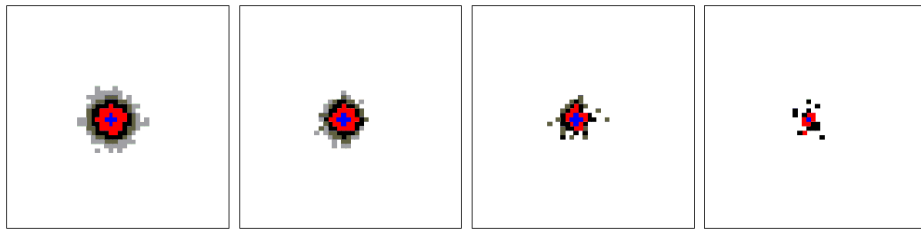


Fig. 5. Comparison of the mean (absolute) percentage error between the super-individual simulations and the individual-based simulation, at $t = 40$.

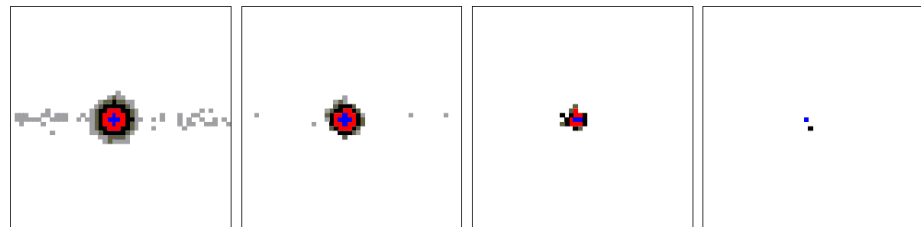
Mean cell density



(a) 10,000 individuals, density at 2 days: (l-r) Individual-based simulation, super-individual simulation scale factor 10, 100 and 1,000



(b) 10,000 individuals, density at 20 days: (l-r) Individual-based simulation, super-individual simulation scale factor 10, 100 and 1,000



(c) 10,000 individuals, density at 40 days: (l-r) Individual-based simulation, super-individual simulation scale factor 10, 100 and 1,000

Fig. 6. Spatial density distributions for individual-based versus super-individual simulations (10,000 aphids) at (a) 2 days (b) 20 days and (c) 40 days. The distribution further from the central cell is influenced by the constant westerly wind direction to result in a linear movement pattern.

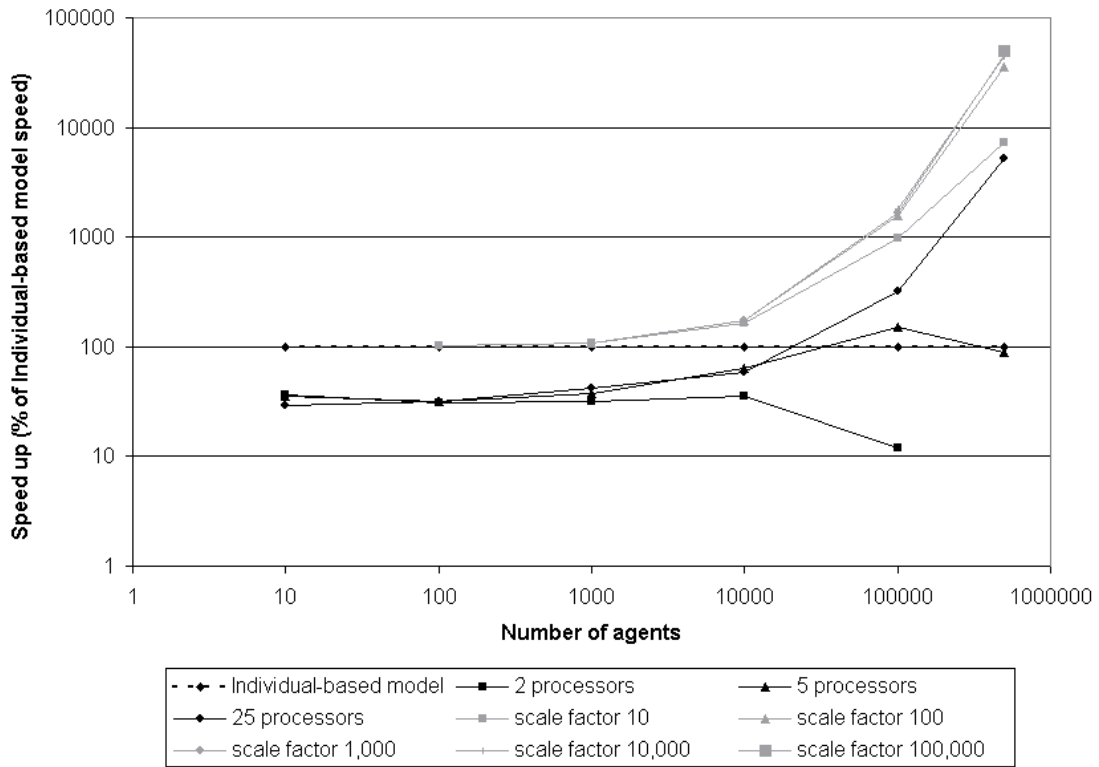


Fig. 7. Plot of the percentage speed up from the individual-based (non-parallel) model against number of agents modelled: comparison between parallel simulations using 2, 5 and 25 nodes and super-individuals of scale factor 10, 100, 1,000, 10,000, 100,000 and 500,000

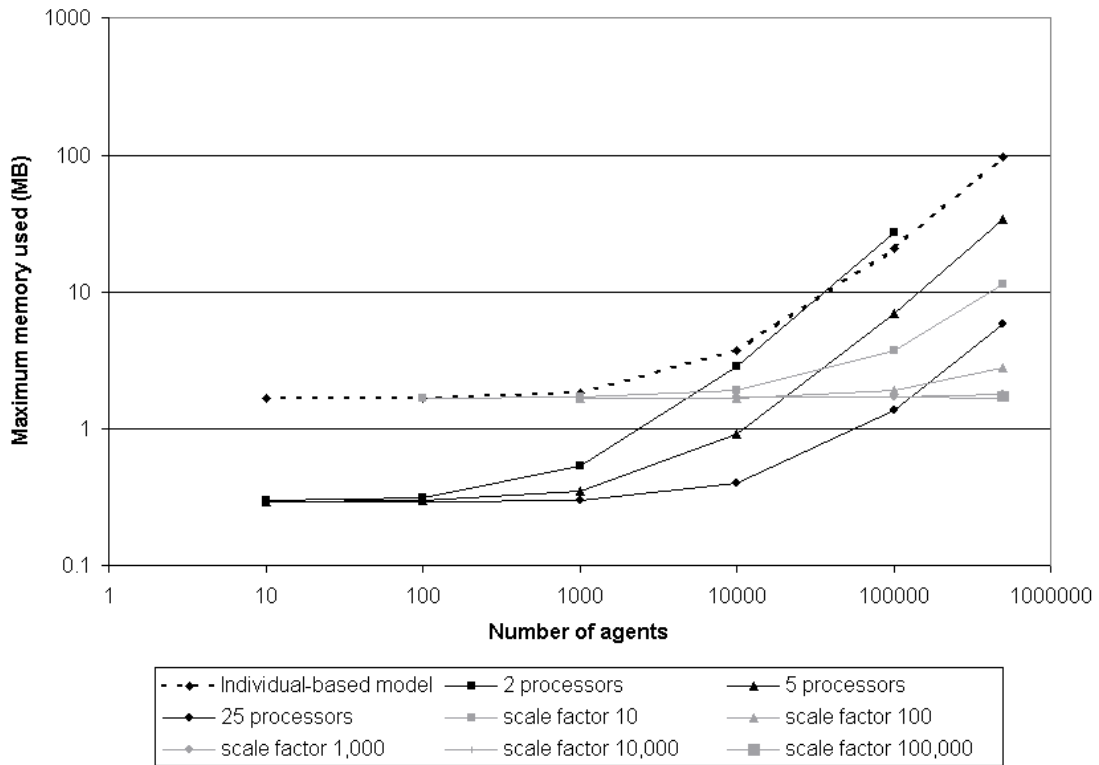


Fig. 8. Plot of the mean maximum memory used in a simulation run against number of agents for the model, for different numbers of nodes (memory per node) and scale factors for super-individuals