

Simulation scalability of large brain neuronal networks thanks to time asynchrony

Cyrille Mascart^{a,1}, Gilles Scarella^{b,a}, Patricia Reynaud-Bouret^b, and Alexandre Muzy^a

^aUniversité Côte d'Azur, CNRS, I3S, France.; ^bUniversité Côte d'Azur, CNRS, LJAD, France.

This manuscript was compiled on June 11, 2021

We present here a new algorithm based on a random model for simulating efficiently large brain neuronal networks. Model parameters (mean firing rate, number of neurons, synaptic connection probability and postsynaptic duration) are easy to calibrate further on real data experiments. Based on time asynchrony assumption, both computational and memory complexities are proved to be theoretically linear with the number of neurons. These results are experimentally validated by sequential simulations of millions of neurons and billions of synapses in few minutes on a single processor desktop computer.

Brain neuronal networks | Random models | Point processes | Stochastic simulation | Discrete event systems | Time asynchrony | Hawkes processes.

There are more and more vast research projects, whose aim is to simulate brain areas or even complete brains to better understand the way it works. Let us cite for instance: the Human Brain Project (1) in Europe, the Brain Mapping by Integrated Neurotechnologies for Disease Studies (Brain/MINDS) (2) in Japan or the Brain Initiative (3) in the United-States. Several approaches are feasible. There is the biochemical approach (4), which is doomed for systems as complex as the brain. A more biophysical approach has been investigated, see for instance (5), where cortical barrels have been successfully simulated, but are limited to about 10^5 neurons. However, the human brain contains about 10^{11} neurons whereas a small monkey, like marmosets (2), has already 6×10^8 neurons (6) and a bigger monkey, like a macaque, has 6×10^9 neurons (6).

To simulate such huge networks, models reduction have to be made. In particular, a neuron has no more physical shape and is just represented by a point in a network with a certain voltage. Hodgkin-Huxley equations (7) are able to reproduce the physical shape if it is combined to other differential equations, representing the dynamic of ion channels, but the complexity of these coupled equations that form a chaotic system (8), makes the system quite difficult to simulate for huge networks. If ion channels dynamic is neglected, the simplest model of voltage is the Integrate-and-Fire model (9). With such models, it has been possible on supercomputers to simulate a human-scale cerebellar network reaching about 68×10^9 neurons (10).

However there is another point of view, which might allow us to simulate such massive networks with simplified models. Indeed, one can use much more random models to reproduce the essential dynamics of the neurons: their firing pattern. The randomization of not only the connectivity graph but also the dynamics on the graph is making the model closer to the data at hand and explain to a certain extent their variability. The introduction of randomness is not new and has been done in many models including Hodgkin-Huxley (11) and Leaky

Integrate-and-fire (LIF for short) (12).

Here we want to focus on particular random models - point processes (13) - which have a particular property: time asynchrony, that is the inability of the model to have two spikes that are produced exactly at the same time by two different neurons. This includes in particular Hawkes models and variants such as GLM, Wold processes, Galves-Löcherbach models, and even some random LIF models with random or soft threshold, all of them having been used to fit real data (13–19).

This property, which is well known in mathematics (14, 20), combined with graph sparsity lead us to propose a new algorithm in (21). The computational complexity of this new algorithm has been computed. Thanks to time asynchrony and to the computational activity tracking of firing neurons (22), we have shown in particular that if the graph is sparse, the complexity cost of the computation of a new point in the system is linear in the number of neurons. However the memory burden was too high to reach networks of 10^8 neurons.

In a preliminary work (23) focusing on the mathematical aspects of mean-field limits of LIFs, we formalized a way to deal with this memory aspect without putting it into practice: the main point is to not keep in memory the whole network, but to regenerate it when need be. Recently, the same idea, under the name of procedural connectivity, has been applied with success on LIF models in (24): using GPU-based parallel programming, and without time asynchrony, the authors have been able to simulate a network of 4×10^6 neurons and 24×10^9 synapses on a desktop GPU computer.

But, as we show in the present article, the gain of procedural connectivity is even huger when combined with time

Significance Statement

Time asynchrony refers to models that prevent any two neurons to fire at the exact same time (up to the numerical precision). Brain models that are time asynchronous have a huge numerical advantage on other ones: they can be easily simulated without huge parallelization or use of GPU. Indeed, this work proposes to exploit the time asynchrony property to derive new simulation algorithms, whose computational and memory complexity can be estimated beforehand. Application on realistic set of parameters show that they can indeed easily reach the size of a small monkey brain, on a usual single processor desktop computer.

Please provide details of author contributions here.

Please declare any competing interests here.

¹ A.O. (Author One) contributed equally to this work with A.T. (Author Two) (remove if not applicable).

² To whom correspondence should be addressed. E-mail: author.twoemail.com

69 asynchrony. Indeed, classical parallel programming usually
70 uses a discrete simulation time and computes for all neurons
71 (or synapses) what happens at each time step in a parallel
72 fashion (even if spikes can be communicated in between two
73 time steps (10)). At each time step, each process correspond-
74 ing to a different neuron has to wait for the calculations of
75 all the other processes to know what needs to be updated
76 before computing the next step. With time asynchrony, we
77 can leverage discrete-event programming (25–28) to track the
78 whole system in time by jumps: from one spike in the net-
79 work to another spike in the network. Since a very small
80 percentage of a brain is firing during a given unit of time
81 (29), the gain we have is tremendous in terms of computations.
82 Thanks to the procedural connectivity, the memory needed
83 to access for the computation of each new spike is also con-
84 trolled. Hence, thanks to procedural connectivity combined
85 with time asynchrony, we propose a new algorithm for time
86 asynchronous models, running sequentially on a single pro-
87 cessor, thus simulating a realistic network of 10^8 neurons for
88 which computational complexity as well as memory costs can
89 be controlled beforehand.

90 Results

91 **Time synchrony for parallel simulation.** Brain simulation of
92 large networks is usually done in parallel based on simulation
93 synchronization. Of course, this depends on both the mathe-
94 matical model at hand and the simulation algorithm. However,
95 for most models, differential equations are used to derive the
96 time course of the membrane voltage for each individual neu-
97 rons. These equations are (approximately) solved by usual
98 discrete-time numerical schemes (*cf.* Figure 1a) (24).

99 In this kind of implementation, when one presynaptic neu-
100 ron fires at a time t , *i.e.*, emits a spike (red dots in Figure
101 1a), the synaptic transmission to post-synaptic neurons is
102 done at the next time step $t + \Delta t$ (orange dots in Figure 1a).
103 Between two synaptic transmissions, the membrane potential
104 of a neuron evolves independently of the other neurons and
105 can be computed in parallel (green dots in Figure 1a). How-
106 ever, since one does not know when a spike will be emitted
107 in the network in advance, the membrane potential of all the
108 neurons are classically computed in a synchronous way to be
109 able to eventually transmit spikes, at each time step Δt of the
110 algorithm.

111 **Time asynchrony for sequential simulation.** As said in the in-
112 troduction, point processes models of neuronal network may
113 guarantee time asynchrony if they have a stochastic intensity.
114 Such processes include Hawkes processes, GLM approaches,
115 Wold processes or Galvès Löcherbach models in continuous
116 time (13, 15–18). Most of these models have proved their
117 efficiency in terms of goodness-of-fit with respect to real spike
118 train data (13–17, 30). Often, the intensity in these models
119 can be informally interpreted as a function of the membrane
120 voltage and for more evidence, we refer the reader to (31),
121 where the spike train of a motor neuron has been shown to
122 be adequately modeled by a point process whose stochastic
123 intensity is a function of the membrane voltage.

124 These point processes models differ from classical LIF,
125 mainly because the higher the intensity (or the membrane
126 potential) of a given neuron is, the more likely it is that the
127 neuron fires, but this is never for sure. In classical LIF models,

128 the neurons fire when their membrane potential reaches a fixed
129 threshold. Therefore it may happen that if a presynaptic neu-
130 ron fires, and if the corresponding postsynaptic neurons have
131 a potential close to the threshold, then all the postsynaptic
132 neurons fire at the exact same time. This phenomenon can
133 be massive (32, 33): this corresponds to the mathematical
134 phenomenon of blow-up, which happens for some mean-field
135 limits of such models. In this case, no time asynchrony is
136 possible but such phenomenon is completely unrealistic from
137 a biological point of view. There are LIF models with random
138 or soft threshold (19, 23) which might not have this problem
139 and which may also satisfy time asynchrony.

140 In (21), we proposed a discrete-event algorithm to simulate
141 point processes with stochastic intensities. This algorithm is
142 based on the theory of local independence graph (34), which
143 is the directed neuronal network in our present case.

144 The algorithm works as follows (see Figure 1b). The spike
145 events happen in continuous time in the system (up to the
146 numerical precision). Once a spike on a particular presynaptic
147 neuron happens (red dots in Figure 1b), the postsynaptic
148 neurons are updated (orange dots in Figure 1b). The presy-
149 naptic and the post synaptic neurons compute their intensities
150 (assimilated to membrane potential) and forecast its evolution
151 (green arrows in Figure 1b) if nothing in between occurs in
152 the system. They are therefore able to forecast their potential
153 next spike (gray dots in Figure 1b). The algorithm maintains
154 a scheduler containing all potential next spikes on all neurons
155 and decides that the next neuron to fire effectively is the one
156 corresponding to the minimum of these potential next spikes.
157 For more details, we refer to (21).

158 The gain comes from the fact that neurons that are not
159 firing a lot, do not require a lot of computation either. In
160 particular we do not have to update all neurons at each spike
161 but only the pre and post synaptic neurons that are involved in
162 the spiking event. This is the main difference with the parallel
163 simulation framework detailed above. The other difference
164 is that we can work with arbitrary precision, typically 10^{-15}
165 if necessary, without impeding the time complexity of the
166 algorithm.

167 Note that the whole algorithm is possible only because
168 two neurons in the network will not spike at the same time:
169 the whole concept is based on time asynchrony to be able to
170 jump from one spike in the system to the next spike in the
171 system. Of course, this is true only up to numerical precision:
172 if two potential next spikes (gray dots on Figure 1b) happen
173 at the exact same time with resolution 10^{-15} , by convention
174 the neuron with the smaller index is said to fire. But the
175 probability of such event is so small that this is not putting
176 the simulation in jeopardy.

177 Also note that this does not prevent neurons to eventually
178 synchronize frequently over a small time duration of a few
179 milliseconds, as defined for instance in (35) and the references
180 therein.

181 **Procedural connectivity.** One of the memory burden of both
182 methods (parallel and time asynchrony) comes from the fact
183 that a classic implementation stores the whole connectivity
184 graph, which is huge for brain scale models.

185 If the connectivity is the result of a random graph and
186 that each presynaptic neuron is randomly connected to its
187 postsynaptic neurons, one can store the random seed instead
188 of the result of the random attribution. Hence the whole

189 graph is never stored in full but only regenerated when need
 190 be (see Figure 1b). The random connectivity is regenerated
 191 at each spike taking advantage of the deterministic nature of
 192 the pseudo-random generator used in the simulation. Storing
 193 the generator initial seed, the seed of each neuron is computed
 194 based on initial seed value and neuron index (see Figure 1c).
 195 With this method, only the initial seed is stored in memory.
 196 Of course this dynamic regeneration at each spike has a cost
 197 in terms of time complexity, but this cost is negligible with
 198 respect to the other computations that need to be made and
 199 this saves memory.

200 This method has been evoked at first in (23) for time asyn-
 201 nchronous algorithms, without being put into practice, whereas
 202 this method has been already implemented with success on
 203 parallel programming with GPU (24).

204 **Computational and memory costs.** In (21), we obtained an
 205 accurate estimate of the complexity of the algorithm without
 206 procedural connectivity, for the simulation of linear Hawkes
 207 processes (*cf.* Equation 7 of (21)). This can be reused to
 208 compute the computational complexity of the same model,
 209 when the procedural connectivity step is added. Thus the
 210 overall time complexity of our algorithm is of the order of

$$211 \quad \mathcal{O}(T [Md^2\bar{m}^2 + \log(M)Md\bar{m} + Md\bar{m}]) \quad [1]$$

212 where M is the number of neurons, p the connection probabili-
 213 ty, $d = \lceil pM \rceil$ the average degree of the network, \bar{m} the average
 214 firing rate of the network and T the simulation duration. The
 215 last term in Eq. (1) corresponds to the computational cost
 216 of the procedural connectivity at each spike, which is indeed
 217 negligible with respect to other terms, as explained before.

218 Note that such computational costs depend in particular
 219 on the intensity shape of the underlying point process model.
 220 The linearity of the Hawkes processes makes it easy to derive,
 221 whereas this can be much more cumbersome with other models
 222 such as stochastic LIF, which needs to compute the distribution
 223 of the time at which the threshold is reached.

224 The maximal memory cost of the procedural connectivity,
 225 without the spike times storage, is of the order of $\mathcal{O}(d\omega) =$
 226 $\mathcal{O}(pM\omega)$, with ω the number of bits necessary for representing
 227 the index of a post-synaptic neuron.

228 The memory cost of a static storage of the whole graph is
 229 of the order of $\mathcal{O}(dM\omega) = \mathcal{O}(pM^2\omega)$. We do not include in
 230 this the memory costs for the storage of each spike of each
 231 neuron. However, this cost is the same whatever the method
 232 and it is of the order of $MT\bar{m}\epsilon$, where ϵ is the number of bits
 233 necessary to represent a spike, which depends on the numerical
 234 precision with which time is recorded. If d is thought to be a
 235 fixed parameter, the memory cost of our complete algorithm
 236 with procedural connectivity is thus

$$237 \quad \mathcal{O}(d\omega + MT\bar{m}\epsilon) \quad [2]$$

238 Conversely, the use of a static storage of the network is
 239 $\mathcal{O}(dM\omega + MT\bar{m}\epsilon)$. Note however that depending on what
 240 the program needs to return, we might not want to have the
 241 whole set of points but only summary statistics such as firing
 242 rates that will cost in memory much less than $\mathcal{O}(MT\bar{m}\epsilon)$.

243 **Choice of brain scale parameters.** Because of the precision of
 244 the actual measurements and the brain region and neuron

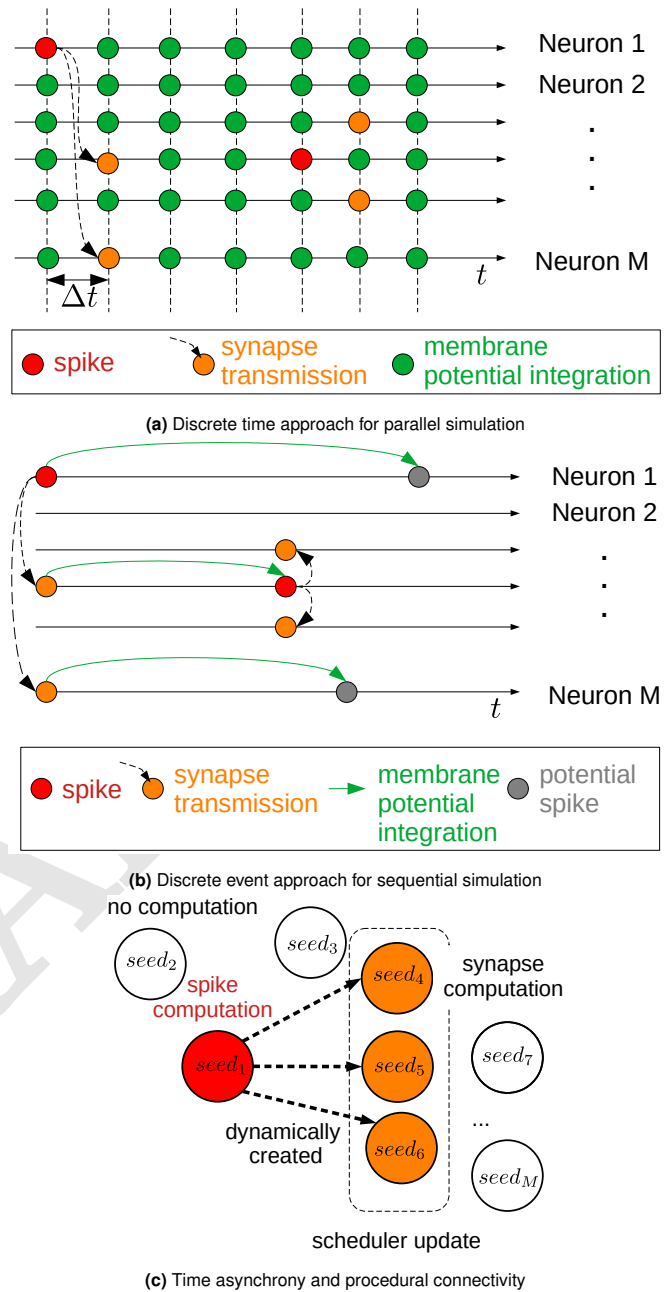


Fig. 1. Neuronal computations in time and in the network

245 variability, it is difficult to estimate quantitatively both physi-
 246 ological (number of synapses per neuron, etc.) and dynamic
 247 parameters (average firing rate, etc.) of neuronal networks
 248 in primates (6) and humans (36). Only rough estimates are
 249 available. The human brain being the more computationally
 250 intensive, we estimate here its main parameters for simulation.
 251 Our goal is indeed to show the algorithm scalability to simulate
 252 large networks with such parameters.

253 To our knowledge, the best documented region of the human
 254 brain is the (neo)cortex. Based on the structural statistics
 255 (number of neurons and synaptic connections) of neuronal
 256 networks in the (neo)cortex, we extrapolate here their repre-
 257 sentative parameter values at brain scale.

258 The firing rate of a neuron in the brain can be estimated

259 by the limited resources at its disposal, especially glucose.
 260 Measures of ATP consumption have shown (see (29)) that the
 261 firing rate of a neuron in human neocortex can be estimated
 262 around $0.16Hz$. Still based on ATP consumption, only 10% of
 263 the neurons in the neocortex can be active at the same time.
 264 So it seems coherent to choose an average of $0.16Hz$. These
 265 values can be extrapolated to the whole brain*, as follows.

The neocortex represents 80% of the volume of the brain (37) and consumes 44% of its energy (29). Considering that the energy consumed by the brain is proportional to the firing rate of the neurons, the power ratio then consists of

$$\frac{P_{cortex}}{P_{brain}} \sim \frac{V_{cortex}\bar{m}_{cortex}}{V_{brain}\bar{m}_{brain}},$$

266 with \bar{m}_{cortex} the mean firing rate of individual neurons in the
 267 neocortex (resp. in the brain) and V_{cortex} the volume of the
 268 neocortex (resp. in the brain). The average firing rate of
 269 the brain then consists of $\bar{m}_{brain} = 0.8 \times \frac{0.16}{0.44} = 0.29 Hz$ per
 270 neuron.

271 This average firing rate should not be confused with the fact
 272 that particular neurons can have a much larger firing rate.
 273 Particularly, groups of neurons synchronize together for
 274 achieving a particular cognitive task: this is the concept of
 275 neuronal assemblies (38). In an assembly, neurons can usually
 276 increase their rates to tens Hz (possibly $50Hz$) over a short
 277 duration. Therefore, we choose a firing rate in the brain where
 278 most of the neurons have a firing rate of $0.3Hz$ but some have
 279 a much higher firing rate (up to $50Hz$) using an heavy tailed
 280 distribution, see Materials and Methods and Table 2.

281 The average number of synaptic connections in human
 282 brains is hard to estimate and depends heavily on the neuron
 283 types and brain regions. For example, in the brain, it is assumed
 284 that the majority of neurons are cerebellum granule cells (39).
 285 In (40), the number of synaptic connections to granule neurons
 286 is estimated to an average of only 4 connections, matching those
 287 observed anatomically. On the other hand, Purkinje neurons can
 288 have up to 200,000 synapses on only one dendrite in the human
 289 brain (39). The approximate number of synapses in the cortex is
 290 0.6×10^{14} (41). Assuming that the volume of the cortex
 291 represents around 80% of the volume of the brain, the number
 292 of synapses in the brain is of order 10^{14} . Considering that the
 293 number of neurons in the human brain is of order 10^{11} (36), we
 294 find that the average number of synapses is about 1,000
 295 synapses per neuron[†]. The synaptic connection probability
 296 thus depends on the number of neurons M : $p_M = \frac{1,000}{M}$.

298 Finally, an action potential arriving on one pre-synaptic
 299 neuron produces an Excitatory PostSynaptic Potential (EPSP),
 300 or an Inhibitory PostSynaptic Potential (IPSP), in the postsynaptic
 301 neuron. The duration of these postsynaptic potentials is about
 302 $\tau = 20ms$ (39).

303 Therefore the parameters that we used in the simulation
 304 are indicated in Table 1. Notice that these parameters are
 305 generic and intuitive and can be taken easily into account in
 306 further studies, either at a biological or at theoretical model
 307 level.

* This calculus can be found on AI impact project webpage: <https://aiimpacts.org/rate-of-neuron-firing/> (lastly verified: 02/09/2021)

† Calculus on AI impact project webpage: <https://aiimpacts.org/scale-of-the-human-brain/> (lastly verified: 02/09/2021).

Simulation duration	$T = 5s$
Mean firing rate	$\bar{m} = 0.3Hz$
Number of neurons simulated	$M = \{10^5, 10^6, 10^7, 10^8\}$
Synaptic connection probability	$p_M = 1000/M$
Postsynaptic duration	$\tau = 20ms$

Table 1. Neuronal network parameters at human brain scale level.

308 **Software and hardware configurations.** The simulations have
 309 been run on a Symmetric shared Memory multiProcessor
 310 (SMP) computer equipped with Intel CascadeLake@2.6GHz
 311 processors[‡]. This kind of computer is used here to have access
 312 to larger memory capacities. At computational level, only
 313 one processor was used for the simulations. For small sizes
 314 of networks requiring small amounts of memory (*cf.* Figure
 315 3b), *e.g.* a network of 10^6 neurons with a total of 10^9 synaptic
 316 connections, this computer is equivalent to a simple desktop
 317 computer. The simulation of such networks takes only 25
 318 minutes for each biological second. This is of the order of
 319 the 4.13×10^6 neurons and 24.2×10^9 synapses simulated on
 320 GPUs (24), which takes about 15 minutes for each biological
 321 second. This GPU-based simulation was already running up
 322 to 35% faster than on 1024 supercomputer nodes (one rack of
 323 an IBM Blue Gene/Q) (42). Our simulation only requires a
 324 single usual processor and no GPU.

325 The implementation of the algorithm is written in C++
 326 (2011) programming language and compiled using g++ 9.3.0.

327 **Firing rate at network level.** Table 2 presents classical elemen-
 328 tary statistics on the simulated firing rates, whereas Figure 2
 329 presents the corresponding densities. As one can see in Section
 330 Material and Methods, the system is initialized with a lot of
 331 neurons whose spontaneous spiking activity is null. The sys-
 332 tem needs to warm up to have almost all neurons spiking. This
 333 explains why the density at $T = 5s$ is still rippled whereas, at
 334 $T = 50s$, it looks much smoother. This last case corresponds
 335 basically to the stationary version of the process. Indeed, as
 336 explained in Section Material and Methods, the parameters
 337 of the Hawkes model (in particular the spontaneous spiking
 338 activity) have been fixed to achieve a certain stationary distri-
 339 bution of the firing rates (with mean $0.3Hz$), which is heavy
 340 tailed to achieve records as large as $50 Hz$. As one can see
 341 (even if at $T = 5s$ the system is not warmed up yet with a
 342 lot of non spiking neurons), one can still achieve the desired
 343 average firing rate and extremal values. These basic statistics
 344 are not varying a lot with T (see Table 2). Note that the
 345 density plots are roughly the same for all configurations: with
 346 ripples at $T = 5s$ and smooth curve at $T = 50s$.

347 Our approach is particularly adapted to simulate precisely
 348 and efficiently a huge disparity in frequency distributions.
 349 Indeed, our simulation algorithm (22) allows focusing efficiently
 350 the computing resources on highly spiking neurons without
 351 computing anything for almost silent neurons (*cf.* Figures 1a
 352 and 1b). The last advantage of our approach is to be able to
 353 store time stamps with a precision of $10^{-15}s$.

354 **Execution times and memory usage.** The simulation execu-
 355 tion times are presented in Figure 3a for different sizes of

[‡] We used v100l and v100xl partitions on Joliot-Curie supercomputer at TGCC as a Fenix Infra-structure resource. Each node of v100l and v100xl has Intel CascadeLake@2.6GHz processors. A node on v100l is a dual-socket one with 2x18 cores, each core having a memory of 10 GBytes, so the total amount of available memory is 360 GBytes. A node on v100xl is a quad-socket one with 4x18 cores, each core having a memory of 41.5 GBytes, so the total amount of available memory is 3 TBytes.

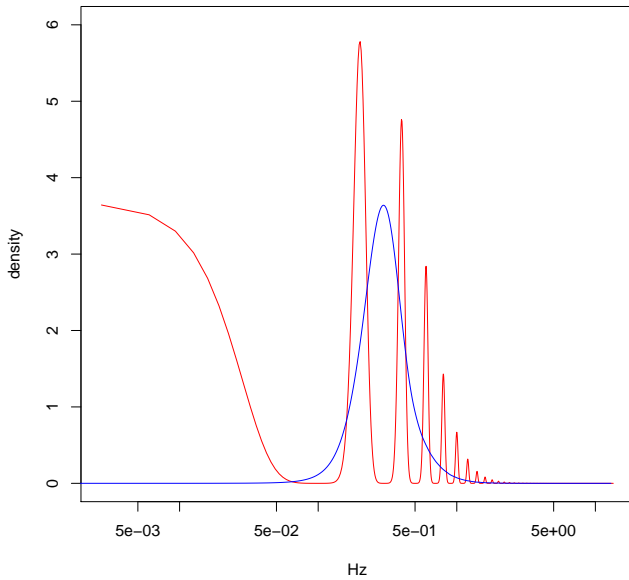


Fig. 2. Densities (on a logarithmic scale) of the simulated firing rates in the network with $M = 10^6$ neurons and $d = 1000$ post-synaptic connections in average. In red, for $T = 5$ s and in blue for $T = 50$ s. These densities are obtained with a Gaussian kernel estimator with bandwidth 0.02.

average number of post-synaptic connections per neuron has few impact on the memory. Indeed, within the network, only the post-synaptic connections receiving spikes are dynamically generated.

Material and Methods

Model. For a set of M neurons, we first design the graph of interaction by saying that neuron j influences neuron i if a Bernoulli variable $Z_{j \rightarrow i}$ of parameter p is non zero. The resulting network is an Erdős-Rényii graph.

Once the network is fixed, we design the spike apparition thanks to a Hawkes process, that is a point process whose intensity is given by

$$\lambda^i(t) = \nu_i + \sum_{j=1}^M \int_0^t h_{j \rightarrow i}(t - \tau) dN_{\tau}^j,$$

with dN^j the point measure associated to neuron j . In this formula, ν_i represents the spontaneous firing rate of the neuron i if the other neurons do not fire, whereas $h_{j \rightarrow i}$ is the interaction function, that is $h_{j \rightarrow i}(u)$ is the increase (if positive) or decrease (if negative) that the firing rate of neuron i suffers due to a spike on j , which happens u seconds before.

We are interested in a particular case of the Hawkes process where all the interaction functions are always the same when they are non null. More precisely, we set the interaction function

$$h_{j \rightarrow i} = Z_{j \rightarrow i} \theta h,$$

where h is a fixed positive interaction function of integral 1, θ is a tuning parameter that we need to calibrate to avoid explosion of the process. We also set $h_{i \rightarrow i} = 0$ (no self interaction). We take $h = 50 \mathbf{1}_{[0,0.02]}$: the interaction function is a constant and non zero only on a small interval of length 20ms, which corresponds to typical Post Synaptic Potentials in the brain.

Let us denote $H_{j \rightarrow i} = \int_0^{+\infty} h_{j \rightarrow i}(t) dt$ and $H = (H_{j \rightarrow i})_{i,j=1,\dots,M}$ is the corresponding matrix (line i corresponds to a triggered neuron, column j to a triggering neuron).

Note in particular that in this model, there are only excitatory neurons : if in the brain, there are inhibitory neurons, this will only reduce the number of points without changing the complexity. Moreover when all the interaction functions are non negative, we can easily understand the explosion condition.

Indeed, this Hawkes process explodes, that is, it produces an exponentially increasing number of point per unit of time (see (43)) if the spectral radius of H is larger than 1.

If (*Condition Stat*) the spectral radius is strictly smaller than 1 (44), then a stationary version exists and the corresponding vector of mean firing rates $m = (m_i)_{i=1,\dots,M}$ is given by

$$m = (I - H)^{-1} \nu. \quad [3]$$

Note also that if we start the simulation without points before 0 in (*Condition Stat*), the process is not stricto sensu stationary but it will converge to an equilibrium given by the stationary state (ergodic theorem) and that the number of points that will be produced is always smaller than the stationary version.

In the present case we want (i) to prevent explosion and (ii) to reach a certain vector m which is biologically realistic (average around 0.3 Hz, records around 50 Hz). Both of these calibrations can be done mathematically beforehand in the Hawkes model : we can guarantee the behavior of the whole

neural networks and different numbers of synaptic connections (called children). The experimental execution times obtained are in agreement with Eq. (1) which predicts, for instance, $O(10^{12})$ operations for $M = 10^7$ and $d = 10^3$. The curves are almost linear (with slopes around 1.1) with respect to the number of neurons, for different numbers of synaptic connections.

The total amount of memory used is displayed in Figure 3b. They are in agreement with the procedural memory complexity of Eq. (2) and also almost linear (with slopes slightly below 1). Note in particular that for $M = 10^7$, $d = 10^3$, $\omega = 32$ and $\epsilon = 64$ (leading to a 10^{-15} precision in time), the memory cost predicted by Eq. (2) is $O(10^{11})$ for the static implementation, whereas it is $O(10^9)$ for the procedural connectivity implementation. Besides notice that, as expected, increasing the

M	d	Average freq.	Freq. min.	Freq. max.	Freq. std.	Percentage of non spiking neuron
1e5	250	0.279 (0.279)	0 (0)	14 (13.94)	0.315 (0.222)	31.2 (0.01)
1e5	500	0.334 (0.333)	0 (0.02)	6.6 (5.76)	0.328 (0.218)	23.5 (0)
1e5	1000	0.399 (0.398)	0 (0.04)	10.4 (11.64)	0.345 (0.220)	16.9 (0)
1e6	250	0.267 (0.267)	0 (0)	19.2 (20.84)	0.308 (0.217)	33 (0.01)
1e6	500	0.322 (0.324)	0 (0)	38.4 (39.38)	0.329 (0.225)	25.1 (0.00)
1e6	1000	0.383 (0.387)	0 (0.02)	13.4 (12.9)	0.344 (0.223)	18.5 (0)
5e6	250	0.26 (0.261)	0 (0)	27.4 (28.5)	0.307 (0.217)	34.1 (0.02)
5e6	500	0.315 (0.316)	0 (0)	34.2 (34.1)	0.324 (0.220)	26 (0.00)
5e6	1000	0.377	0	19.8	0.342	19
1e7	250	0.258 (0.259)	0 (0)	23.2 (21.62)	0.306 (0.217)	34.4 (0.02)
1e7	500	0.311	0	21.6	0.322	26.4
1e7	1000	0.374	0	21.8	0.342	19.3
5e7	250	0.253	0	38.2	0.304	35.4
5e7	500	0.305	0	50.6	0.321	27.2
1e8	250	0.251	0	46.2	0.303	35.7

Table 2. Firing rates elementary statistics obtained by simulation for different sizes of neural networks and different numbers of synaptic connections and $T = 5$ s. The number between parentheses displays the results at $T = 50$ s for the less complex simulations

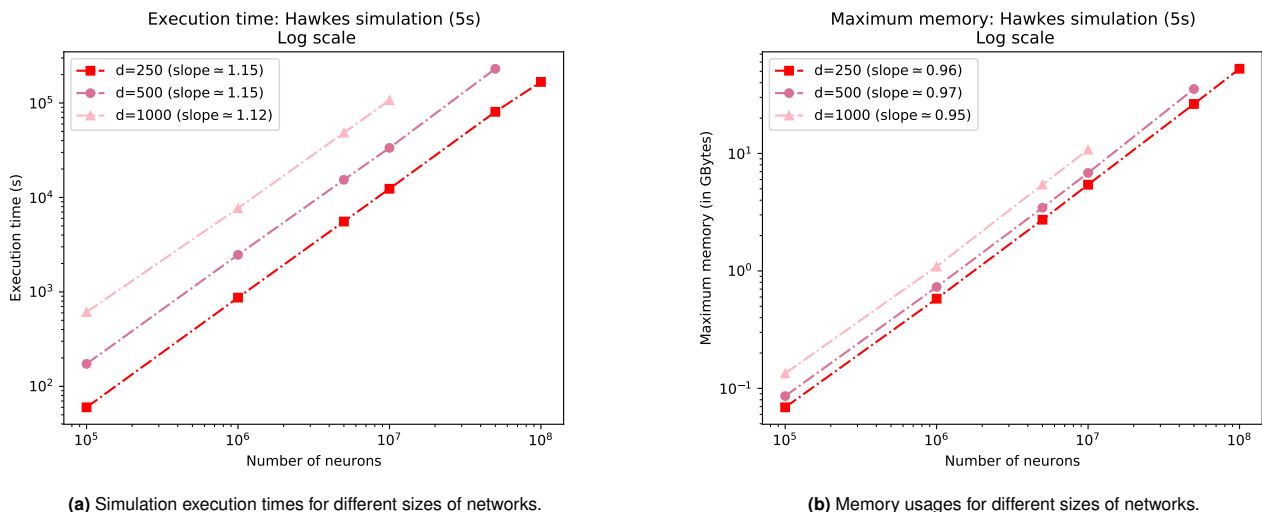


Fig. 3. Simulation execution times (3a) and memory usage (3b) for different sizes of networks.

418 system even before performing the simulation, whereas this
419 might be much more intricate for other models such as LIF.

420 **Choice of θ or how to avoid explosion.** Note that $H = \theta Z$,
421 with $Z = (Z_{j \rightarrow i})_{i,j=1,\dots,M}$. So if we can compute the largest
422 eigenvalue of Z or an upper bound, we can decide how to
423 choose θ .

We can use Gershgorin circles (45) to say that any complex
eigenvalue λ of Z satisfies (because the diagonal is null),

$$|\lambda| \leq \max_{i=1,\dots,M} \sum_{j \neq i} Z_{j \rightarrow i}.$$

424 Therefore the spectral radius is upper bounded by
425 $\max_{i=1,\dots,M} B_i$, where $B_i = \sum_{j \neq i} Z_{j \rightarrow i}$. This random quantity
426 can be computed for small networks but it is clearly too
427 intensive in our setting: indeed, with the procedural connectivity
428 implementation, it is always easy to access the children
429 ℓ of a given i , i.e. such that $i \rightarrow \ell$ is in the graph, but we need
430 to look at all the neurons in the graphs to find out the set of
431 parents j of i , i.e. such that $j \rightarrow i$ is in the graph. However,
432 probabilistic estimates might be computed mathematically.
433 Indeed B_i is just a sum of i.i.d. Bernoulli variables. So we
434 can apply Bernstein's inequality (46). This leads to, for all
435 positive x ,

$$\begin{aligned} &\forall i = 1, \dots, M, \\ &\mathbb{P}(B_i \geq (M-1)p + \sqrt{2(M-1)p(1-p)x} + x/3) \leq e^{-x}, \end{aligned}$$

and, by union bound, for the maximum

$$\mathbb{P}(\max_{i=1,\dots,M} B_i \geq (M-1)p + \sqrt{2(M-1)p(1-p)x} + x/3) \leq M e^{-x}.$$

Therefore let us fix a level α , say 1%, and take $x = \log(M) + \log(1/\alpha)$ in the previous equation. We obtain that with probability larger than $1 - \alpha$, the spectral radius of Z is upper bounded

$$\rho_{\max} = (M-1)p + \xi_\alpha$$

with

$$\xi_\alpha = \sqrt{2(M-1)p(1-p)[\log(M) + \log(1/\alpha)]} + [\log(M) + \log(1/\alpha)]/3$$

Note that ρ_{\max} is roughly $(M-1)p$, which is the largest
eigenvalue of $\mathbb{E}(Z)$. Finally it means that if we take $\theta <$
 $1/\rho_{\max}$, the process will not explode with probability larger
than $1 - \alpha$. In practice, to ensure a strong enough interaction,
we take $\theta = 0.9\rho_{\max}^{-1}$.

443 **Choice of ν_i or how to constraint the distribution of the mean**
444 **firing rates.** The first step consists in deciding for a target
445 distribution for the m_i . We have chosen to pick the m_i 's
446 independently as $0.1X$ where X is the absolute value of a
447 student variable with mean 3 and 4 degrees of freedom. The
448 choice of the student variable was driven by the wish of having
449 a moderate heavy tail, which will ensure records around 50
450 Hz and a mean around 0.3Hz.

The problem is that the m_i 's are not parameters of the
model, so we need to understand how to tune ν_i to get such
 m_i 's. Note that by inverting Eq. (3), we get that

$$(I - H)m = \nu$$

that is for all i

$$\nu_i = m_i - \theta \sum_{j \neq i} m_j Z_{j \rightarrow i},$$

451 which is very intuitive (47). Indeed the spontaneous rate that
452 we need to put is the mean firing rate m_i minus what can be
453 explained with the parents of i .

454 So in theory, the Hawkes model is very easy to tune for pre-
455 scribed firing rates since there is a linear relationship between
456 both. However, and for the same reasons as before, it might
457 be too computationally intensive to compute this explicitly.

458 One possible way is to again use concentration inequalities,
459 but this time on $\sum_{j \neq i} m_j Z_{j \rightarrow i}$ and not on B_i . However we
460 decided to do something simpler, which works well (as seen in
461 Figure 2).

462 Indeed $\sum_{j \neq i} m_j Z_{j \rightarrow i}$ is a sum of about $(M - 1)p \simeq 1000$
463 i.i.d variables with mean $\bar{m} = 0.3\text{Hz}$. Hence it should be close
464 to $\bar{m}B_i$. With the previous computations, we know already
465 that ν_i should therefore be larger than $m_i - \theta\rho_{\max}\bar{m}$.

With the previous choice of $\theta = 0.9\rho_{\max}^{-1}$, we have chosen to
take the positive part for the ν_i 's in the simulation, that is :

$$\nu_i = \max(m_i - 0.9\bar{m}, 0).$$

466 Therefore ν_i remains positive or null, which guarantees that
467 the Hawkes process stays linear. However, this also means
468 that a non negligible portion of the neurons start with a null
469 spontaneous firing rate, which explains the ripples of Figure 2.

470 With this choice, we cannot hope to have exactly the same
471 distribution as the desired m_i 's, but it conserves the same
472 heavy tail and roughly the same mean firing rate as the one
473 we wanted, as one can see on Table 2.

474 Discussion

475 Thanks to time asynchrony, we propose a new scalable algo-
476 rithm to the simulate spiking activity of neuronal networks.
477 We are able to generate roughly the same firing pattern as a
478 real brain for a range between 10^5 and 10^8 neurons, in a few
479 minutes on a single processor, most parameters being tuned
480 thanks to general considerations inferred from the literature.
481 Corresponding computational and memory complexities are
482 shown to be both linear.

483 At simulation level, whereas usual simulations are based
484 on the continuous variation of the electrical potentials of LIF
485 neurons, point processes lead to much more efficient simula-
486 tions. In particular, instead of computing the small continuous
487 variations for all neurons, only discrete spikes and their in-
488 teractions are simulated in the network. Between two spikes
489 no computations are done. Point processes also lead to time
490 asynchrony (two spikes cannot occur at the same time in the
491 network), which is a fundamental hypothesis for the algorithm
492 to work.

493 Combining the time asynchrony hypothesis with procedural
494 connectivity drastically reduces the memory consumption and
495 also, for the same network activity, reduces the computations
496 per spikes (*cf.* Figures 1a and 1b). In particular, complexities
497 (both theoretical and concrete) can be computed and proved
498 to be almost linear in the number of neurons, when Hawkes
499 processes are generated, leading to simulation scalability of
500 the whole approach without precision loss.

501 Both modeling and simulation results open many research
502 perspectives. We are currently developing new discrete event
503 algorithms that are able to simulate the spiking activity of
504 neurons embedded in potentially infinite neuronal networks
505 (48). This paves the path for simulation of parts of the brain
506 as an open physical system.

507 Furthermore, the minimal number of computations and
508 memory storage obtained here open new exciting perspectives
509 with respect to massive neuromorphic computers, by improving
510 the energy saving consumption of neuromorphic components
511 (49).

512 Finally, if we have proved that the simulation is doable,
513 the point process model used here can be calibrated further
514 on real data, by incorporating inhibition and more variability
515 in the interaction functions. Also, this model can be used for
516 reconstructing the functional connectivity of experimentally
517 recorded neurons (15, 47) to have access to more realistic

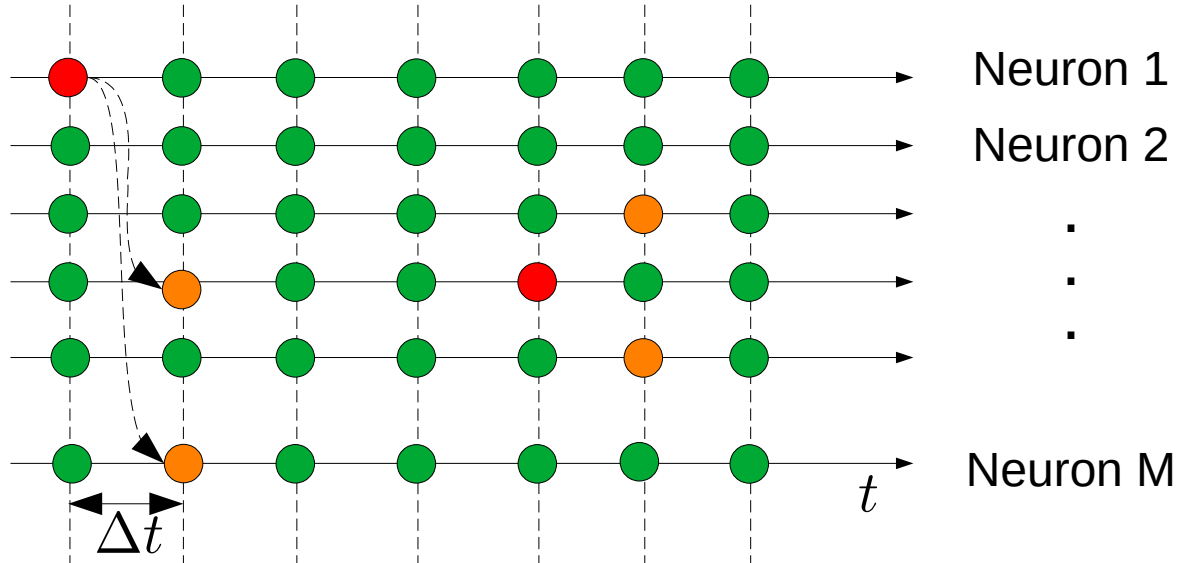
interaction functions. Besides, as only few computing resources
(one single processor) is used with a minimal memory amount,
this opens new possibilities to run in parallel many independent
replications of stochastic simulations of large networks. This
is particularly interesting for calibrating models based on real
data collections.

ACKNOWLEDGMENTS. This work is part of the project Hyper-
Brain from Human Brain Project (HBP) EBRAINS EU initiative.
The simulations were run on Fenix Infrastructure resources, which
are partially funded from the European Union's Horizon 2020 re-
search and innovation program through the ICEI project under the
grant agreement No. 800858. Our research was supported by the
French government, through CNRS, the UCA^{Jedi} and 3IA Côte
d'Azur Investissements d'Avenir managed by the National Research
Agency (ANR-15- IDEX-01 and ANR-19-P3IA-0002), directly by
the ANR project ChaMaNe (ANR-19-CE40-0024-02) and by the in-
terdisciplinary Institute for Modeling in Neuroscience and Cognition
(NeuroMod) of the Université Côte d'Azur.

1. K Amunts, et al., The human brain project—synergy between neuroscience, computing, infor-
matics, and brain-inspired technologies. *PLoS biology* **17**, e3000344 (2019).
2. MR Dando, Japan's brain/minds project (2020).
3. W Koroshetz, et al., The state of the nih brain initiative. *J. Neurosci.* **38**, 6427–6438 (2018).
4. RR Netz, WA Eaton, Estimating computational limits on theoretical descriptions of biological
cells. *PNAS* **118** (2021).
5. E Gal, et al., Rich cell-type-specific network topology in neocortical microcircuitry.
Nat. Neurosci. **20**, 1004–1013 (2017).
6. S Herculano-Houzel, CE Collins, P Wong, JH Kaas, Cellular scaling rules for primate brains.
Proc. Natl. Acad. Sci. **104**, 3562–3567 (2007).
7. DA McCormick, Y Shu, Y Yu, Hodgkin and huxley model—still standing? *Nature* **445**, E1–E2
(2007).
8. J Guckenheimer, RA Oliva, Chaos in the hodgkin–huxley model. *SIAM J. on Appl. Dyn. Syst.*
1, 105–114 (2002).
9. L Lapicque, Recherches quantitatives sur l'excitation électrique des nerfs traitée comme une
polarization. *J. Physiol. Pathology* **9**, 620–635 (1907).
10. H Yamaura, J Igarashi, T Yamazaki, Simulation of a human-scale cerebellar network model on
the k computer. *Front. neuroinformatics* **14**, 16 (2020).
11. RF Fox, Stochastic versions of the hodgkin-huxley equations. *Biophys. journal* **72**, 2068–
2074 (1997).
12. AA Lazar, EA Pnevmatikakis, Reconstruction of sensory stimuli encoded with integrate-
and-fire neurons with random thresholds. *EURASIP J. on Adv. Signal Process.* **2009**, 1–14
(2009).
13. W Truccolo, UT Eden, MR Fellows, JP Donoghue, EN Brown, A point process framework
for relating neural spiking activity to spiking history, neural ensemble, and extrinsic covariate
effects. *J. Neurophysiol.* **93**, 1074–1089 (2005).
14. EN Brown, R Barbieri, V Ventura, RE Kass, L Frank, The time-rescaling theorem and its
application to neural spike train data analysis. *Neural Comput.* **14**, 325–346 (2002).
15. P Reynaud-Bouret, V Rivoirard, F Grammont, C Tuleau-Malot, Goodness-of-fit tests and non-
parametric adaptive estimation for spike train analysis. *The J. Math. Neurosci.* **4**, 3 (2014).
16. J Pillow, et al., Spatio-temporal correlations and visual signalling in a complete neuronal
population. *Nature* **454**, 995–999 (2008).
17. C Pouzat, A Chaffiol, Automatic spike train analysis and report generation. an implementation
with r, r2html and star. *J. Neurosci. Methods* **181**, 119–144 (2009).
18. A Galves, E Löcherbach, Infinite systems of interacting chains with memory of variable
length—a stochastic model for biological neural nets. *J. Stat. Phys.* **151**, 896–921 (2013).
19. L Sacerdote, MT Giraudo, *Stochastic Biomathematical Models*. (Lecture Notes in Mathemat-
ics, Springer) Vol. 2058, pp. 99–148 (2013).
20. P Brémaud, *Point processes and queues*. (Springer-Verlag, New York-Berlin), (1981) Martin-
gale dynamics, Springer Series in Statistics.
21. C Mascart, A Muzy, P Reynaud-bouret, Efficient simulation of sparse graphs of point pro-
cesses (2020).
22. A Muzy, Exploiting activity for the modeling and simulation of dynamics and learning pro-
cesses in hierarchical (neurocognitive) systems. *Comput. Sci. Eng.* **21**, 84–93 (2019).
23. Grazieschi, Paolo, et al., Network of interacting neurons with random synaptic weights.
ESAIM: ProcS **65**, 445–475 (2019).
24. JC Knight, T Nowotny, Larger gpu-accelerated brain simulations with procedural connectivity.
Nat. Comput. Sci. **1**, 136–142 (2021).
25. A Muzy, BP Zeigler, F Grammont, Iterative specification as a modeling and simulation formal-
ism for i/o general systems. *IEEE Syst. J.* **12**, 2982–2993 (2017).
26. R Brette, et al., Simulation of networks of spiking neurons: a review of tools and strategies.
J. computational neuroscience **23**, 349–398 (2007).
27. M Rudolph, A Destexhe, Analytical integrate-and-fire neuron models with conductance-based
dynamics for event-driven simulation strategies. *Neural computation* **18**, 2146–2210 (2006).
28. JD Toublou, OD Faugeras, A markovian event-based framework for stochastic spiking neural
networks. *J. computational neuroscience* **31**, 485–507 (2011).
29. P Lennie, The cost of cortical computation. *Curr. Biol.* **13**, 493–497 (2003).
30. A Duarte, A Galves, E Löcherbach, G Ost, Estimating the interaction graph of stochastic
neural dynamics. *Bernoulli* **25**, 771–792 (2019).
31. P Jahn, RW Berg, Jr Hounsgaard, S Ditlevsen, Motoneuron membrane potentials follow a
time inhomogeneous jump diffusion process. *J. Comput. Neurosci.* **31**, 563–579 (2011).

- 597 32. MJ Cáceres, JA Carrillo, B Perthame, Analysis of nonlinear noisy integrate & fire neuron
598 models: blow-up and steady states. *J. Math. Neurosci.* **1**, Art. 7, 33 (2011).
- 599 33. F Delarue, J Inglis, S Rubenthaler, E Tanré, Global solvability of a networked integrate-and-
600 fire model of McKean-Vlasov type. *Ann. Appl. Probab.* **25**, 2096–2133 (2015).
- 601 34. V Didelez, Graphical models of markes point processes based on local independence.
602 *J.R. Stat. Soc. B* **70**, 245–264 (2008).
- 603 35. C Tuleau-Malot, A Rouis, F Grammont, P Reynaud-Bouret, Multiple Tests Based on a
604 Gaussian Approximation of the Unitary Events Method with Delayed Coincidence Count.
605 *Neural Comput.* **26**, 1408–1454 (2014).
- 606 36. S Herculano-Houzel, R Lent, Isotropic fractionator: A simple, rapid method for the quantifica-
607 tion of total cell and neuron numbers in the brain. *J. Neurosci.* **25**, 2518–2521 (2005).
- 608 37. H Stephan, H Frahm, G Baron, New and revised data on volumes of brain structures in
609 insectivores and primates. *Folia primatologica* **35**, 1–29 (1981).
- 610 38. F Grammont, A Riehle, Spike synchronization and firing rate in a population of motor cortical
611 neurons in relation to movement direction and reaction time. *Biol. cybernetics* **88**, 360–373
612 (2003).
- 613 39. GM Shepherd, *The synaptic organization of the brain*. (Oxford university press), (2004).
- 614 40. A Litwin-Kumar, KD Harris, R Axel, H Sompolinsky, L Abbott, Optimal degrees of synaptic
615 connectivity. *Neuron* **93**, 1153–1164.e7 (2017).
- 616 41. DA Drachman, Do we have brain to spare? *Neurology* **64**, 2004–2005 (2005).
- 617 42. MO Gewaltig, M Diesmann, Nest (neural simulation tool). *Scholarpedia* **2**, 1430 (2007).
- 618 43. S Delattre, N Fournier, M Hoffmann, Hawkes processes on large networks. *Ann. App. Probab.*
619 **26**, 216 – 261 (2016).
- 620 44. AG Hawkes, D Oakes, A cluster process representation of a self-exciting process.
621 *J. Appl. Probab.* **11**, 493–503 (1974).
- 622 45. R Varga, *Gershgorin and his circles*. (Springer-Verlag), (2004) Springer Series in Computa-
623 tional Mathematics.
- 624 46. S Boucheron, G Lugosi, P Massart, *Concentration inequalities*. (Oxford University Press,
625 Oxford), (2013) A nonasymptotic theory of independence, With a foreword by Michel Ledoux.
- 626 47. P Reynaud-Bouret, A Muzy, I Bethus, Towards a mathematical definition of functional con-
627 nectivity. (2021).
- 628 48. TC Phi, A Muzy, P Reynaud-Bouret, Event-Scheduling Algorithms with Kalikow Decomposi-
629 tion for Simulating Potentially Infinite Neuronal Networks. *SN Computer Science* **1** (2020).
- 630 49. LA Plana, et al., spinnlink: Fpga-based interconnect for the million-core spinnaker system.
631 *IEEE Access* **8**, 84918–84928 (2020).

DRAFT



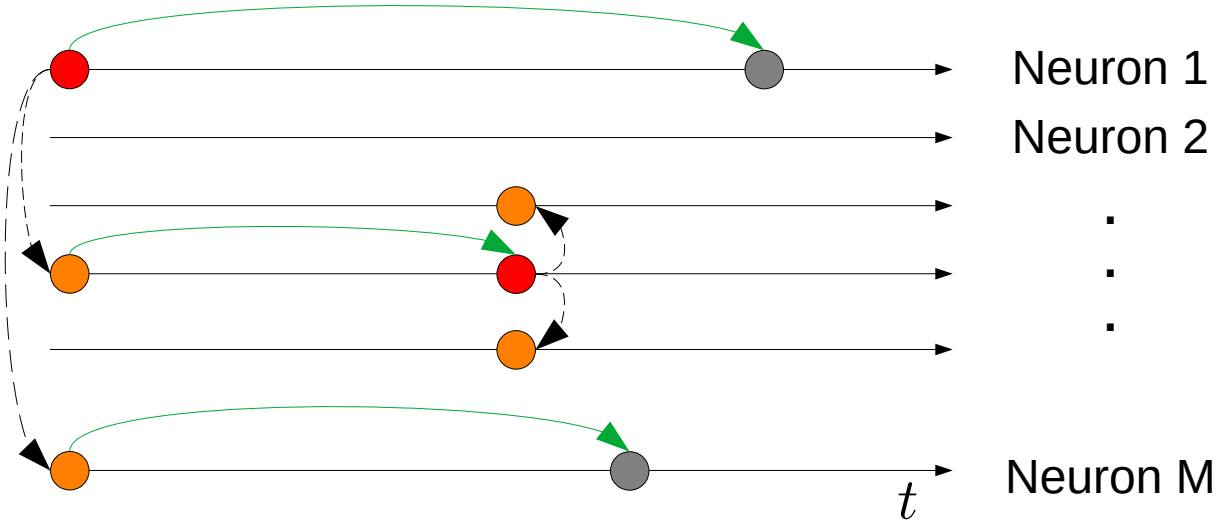
● spike



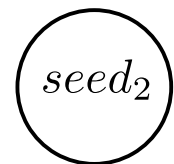
synapse
transmission



membrane
potential integration



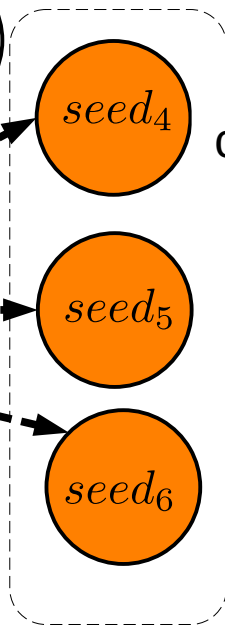
no computation



spike
computation



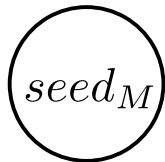
dynamically
created



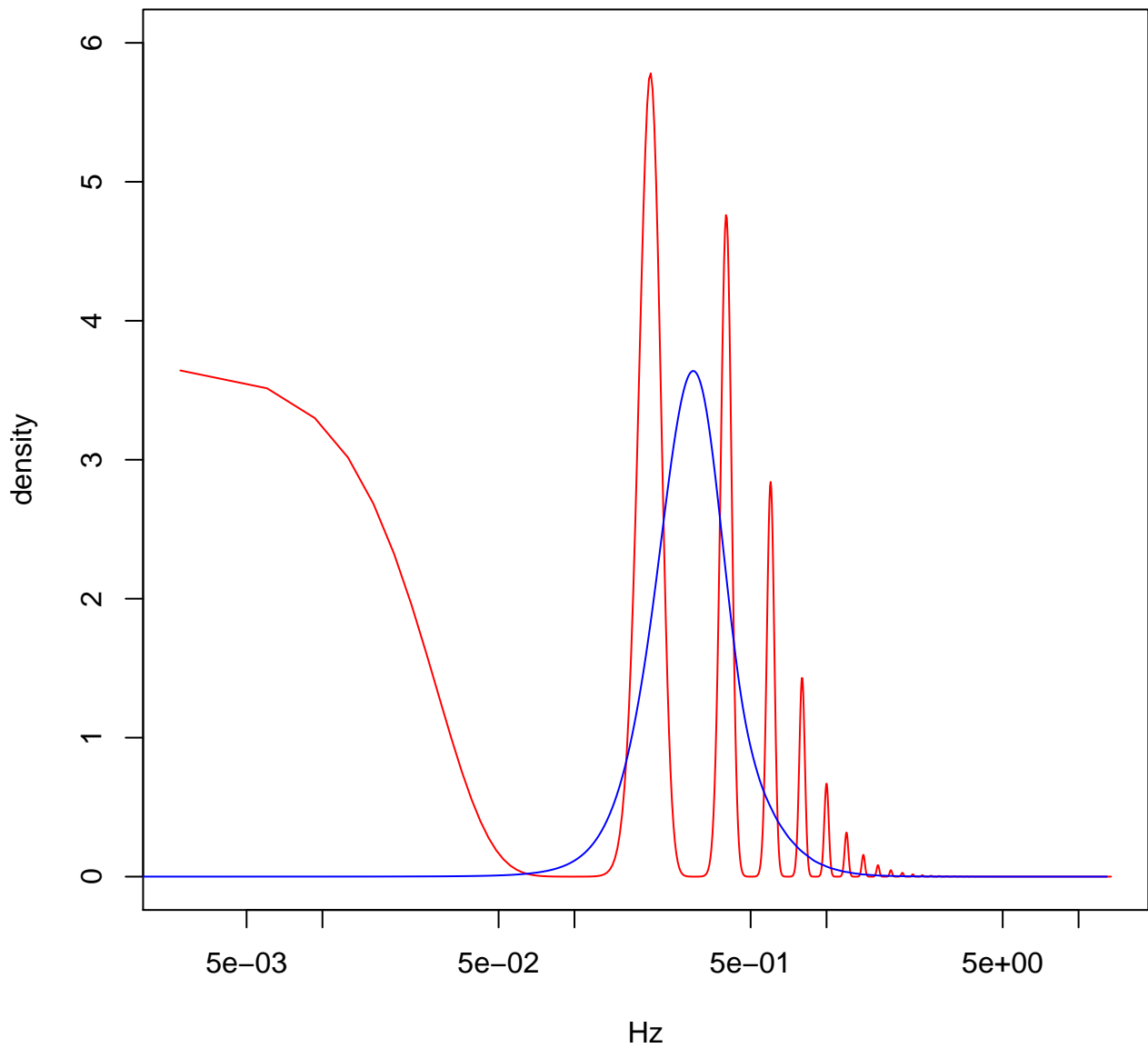
synapse
computation



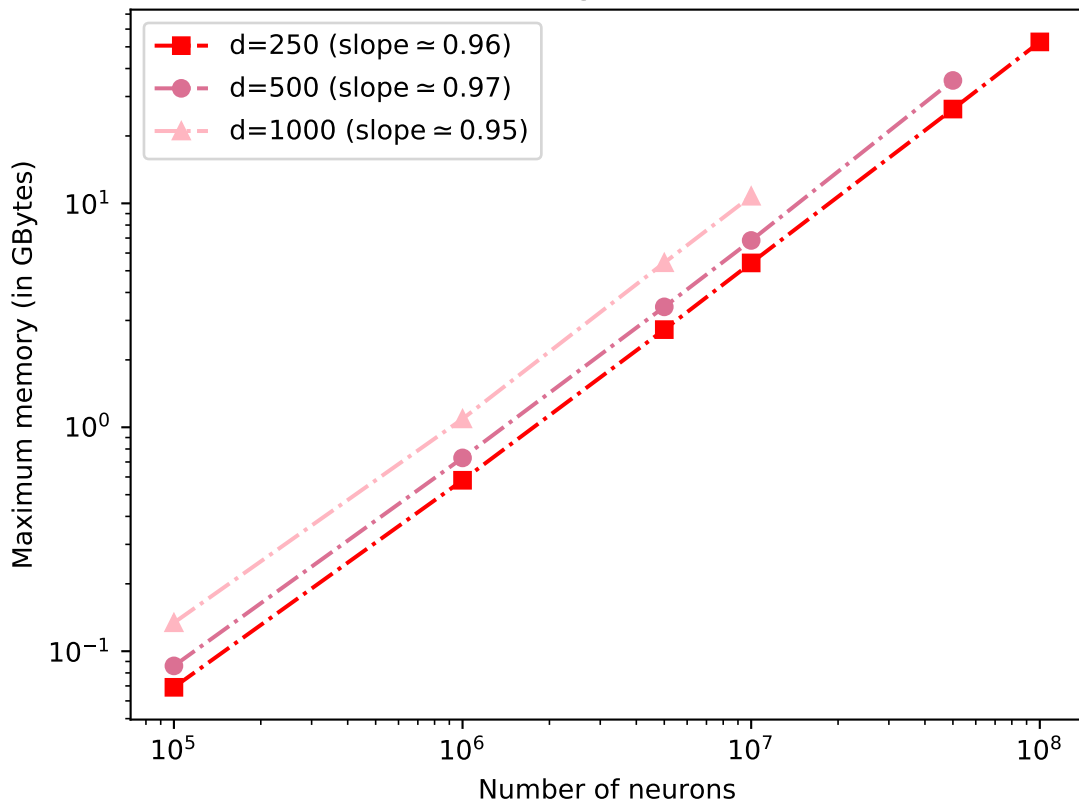
...



scheduler update



Maximum memory: Hawkes simulation (5s)
Log scale



Execution time: Hawkes simulation (5s)
Log scale

