# Analysis of Multicast Traffic Induced by Spiking Neural Networks in the ANA System

**Anonymous Author**
Anonymous Institution

## Abstract

The *Anonymous Neuromimetic Architecture* is a massively parallel system devised to simulate very large-scale spiking neural networks in biological real-time. Its largest configuration incorporates up to 64K bespoke multicore System-on-Chips which are interconnected using a two-dimensional triangular torus network. Neural events (spikes) are modeled as packets that propagate through this fabric relying on a novel on-chip *multicast* router. This router has very modest characteristics when compared with off-chip routers found on datacenters. However, we confirm analytically that this multicast router is a *better solution* for the target application than more powerful point-to-point routers, such as Cray's SeaStar2 or those in IBM's Blue Gene/P. Our analysis shows that the resource savings obtained by the use of multicast more than compensate a lower bandwidth. Furthermore we provide and evaluate four different strategies to generate multicast trees, each of them providing different features and drawbacks.

## 1    Introduction

The Anonymous Neuromimetic Architecture system – hereafter, ANA – is a biologically-inspired massively parallel architecture based on custom-made multicore System-on-Chip (SoC). ANA is designed with the aim of simulating very large-scale spiking neural networks (up to $10^9$ neurons, roughly a 10% of the neurons in the human cortex) in biological real-time. Currently, some *test* chips with two cores and a fully functional router have been produced and successfully demonstrated running spiking neural models [1]. ANA chips, due to their low-power design can be used as control systems for robots providing them with *real-time stimulus-response* behavior [2, 3].

Biological spiking neural networks communicate by means of spike events which occur when a neuron is stimulated beyond a given threshold and fires. Spike signals are communicated to all connected neurons, with typical *fan-outs* of the order of 1000 [4]. These applications exhibit abundant parallelism, no explicit requirement to maintain coherent shared memory and are naturally resilient to failures; neurons may die (1 per second in the human brain [5]) and spikes may be missed, but the brain continues working properly. ANA takes advantage of these characteristics to deploy a well-balanced, low-power massively parallel architecture. The largest configuration houses 64K SoCs – creating a system with over $10^6$ computing cores – which are interconnected using a two-dimensional triangular torus (see Fig. 1). Neurons are modeled in software and their spikes generate packets that propagate through the on- and inter-chip communication fabric relying on custom-made on-chip *multicast* routers. The use of multicast routers alleviates the pressure suffered by the interconnection network because of the high connectivity of the simulated neural models. This paper shows analytical proof of these benefits measured both in terms of network bandwidth and number of neurons that can be supported by the system. However, constructing multicast trees from a source to a set of destinations is not *trivial*, and for this reason we describe sev-
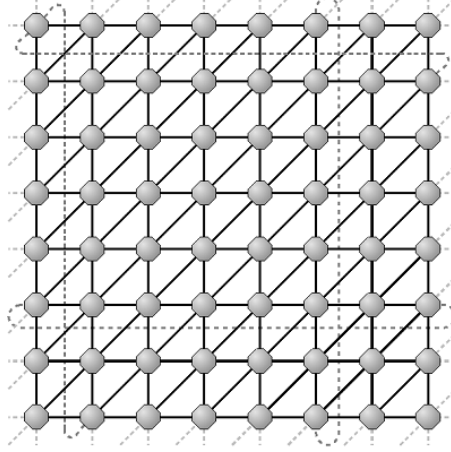
Figure 1. 8×8 triangular torus topology. Most peripheral wrap-around links not showed.

eral strategies to generate multicast trees. We also compare them considering their resource requirements.

The remaining of this paper is organized as follows. Section 2 describes the main characteristics of the system. Section 3 is devoted to define terminology and describe the four multicast strategies. Section 4 compares the proposed strategies. Section 5 provides an analytical evaluation of unicast and multicast alternatives, showing the benefits of using multicast. Finally, Section 6 closes the paper with some concluding remarks.

## 2    System architecture

Each ANA SoC (depicted in Fig. 2) contains 18 low-power ARM968 cores running an independent *event-driven* neural process. Events are generated by different modules: timer, vector interrupt controller (VIC), communication controller and DMA controller. Detailed simulations of the chip confirmed that each core can model in biological real-time up to around 1000 neurons [6]. The chip includes a SDRAM of 128MB mainly used to store synaptic information. Resources within the chip are connected by a custom-made self-timed Network-on-Chip (NoC) [7].

Each chip incorporates a router [8] whose primary role is to direct neural event packets to those chips and cores containing destination neurons. The router has 18 ports for the cores and six ports to communicate with six adjacent chips. Ports are hierarchically merged in three stages before using the actual routing engine. The router is able to handle a single packet at once, but works 8 times faster than transmission ports and is not expected to become a bottleneck – most of the time routers will be idle, and router delay will barely affect the pace at which packets are processed. The router supports point-to-point and multicast communications using small packets of 40 bits. We will see later how the multicast engine
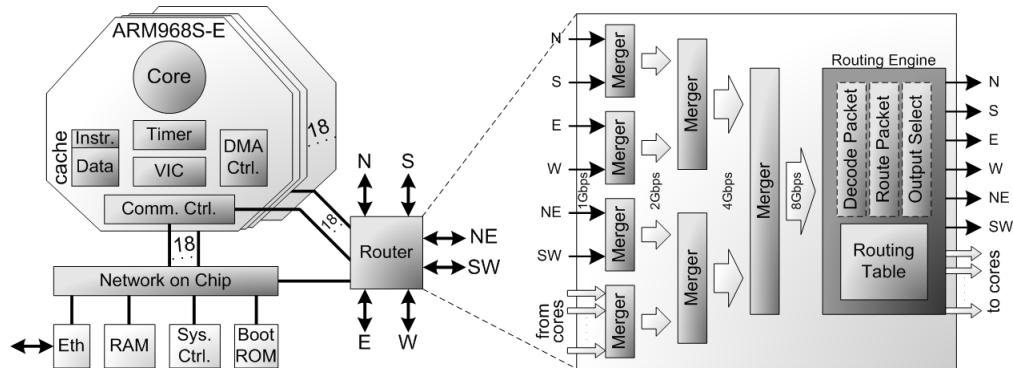


Figure 2. Schematic model of an ANA chip with all its components depicted.

2

66 reduces pressure at the injection ports, and also the number of packets traversing the net-
67 work, when compared to a point-to-point alternative. Following an Address Event Represen-
68 tation [9] a neural-event packet does not contain any information about its destination(s),
69 only the neuron that has fired. The information necessary to deliver a neural packet to all the
70 relevant cores and chips is compressed and distributed across the 1024-word routing table
71 within each router. To allow further compression, routing tables offer a *masked associative*
72 route look-up and routers are designed to perform a *default routing* – which requires no entry
73 in the routing table – by sending the packet to the port opposite to the one the packet comes
74 from, i.e. if a packet comes from the North it will be sent to the South.

75 Network *flow-control* is very simple. When a packet arrives to the router, one or more output
76 ports are selected and the packet is transmitted through them. If the packet cannot be for-
77 warded, the router will keep trying for a while and after a timeout the packet will be dropped
78 to avoid deadlock. Similarly, to avoid livelock situations, packets wandering around the net-
79 work for too much time are considered outdated and dropped. Emulating the behavior of
80 biological neural networks, dropped packets in ANA are not re-sent. As discussed, losing
81 neurons or signals does not impede the normal operation of the biological process; however
82 dropping level must be kept (very) low.

## 3      Definitions

84 Throughout this paper we will use the following conventions to represent the variables in our
85 study. The average distance between a source and its destinations is defined as $\bar{d}$ (measured
86 in hops). The fan-out – number of destinations – is defined as $F$ (nodes/spike). The link
87 bandwidth is represented as $B_L$ (bits/sec). The number of links per node is identified as $L$
88 (links/node). $n_n$ represents the number of neurons running in each node (neurons/node). $f_s$
89 represents the frequency at which neurons generate spikes (spikes/second). $p_s$ defines the
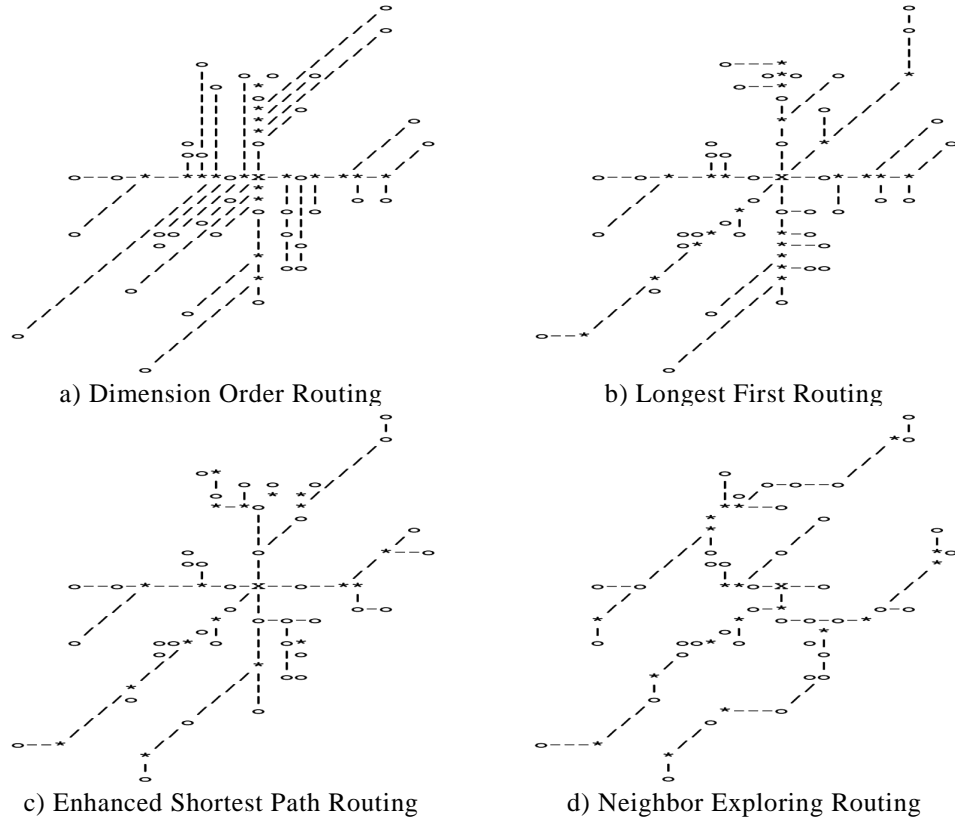90 packet size (bits).



a) Dimension Order Routing          b) Longest First Routing

c) Enhanced Shortest Path Routing        d) Neighbor Exploring Routing

Figure 3. Example of multicat trees. *legend*: 'x' source, 'o' destinations, '*' a entry is
required in routing table, '-' '|' '/' default route (no entry)

3

91  Next we describe four strategies to generate multicast paths. We did not find any previous
92  strategy in the literature that can be compared with the proposed strategies as multicast tree
93  generation is used often with rather different objectives (avoid blocking in multistage net-
94  works [10, 11] and topology discovery in IP networks [12, 13]). In the algorithmic defini-
95  tions, all the functions have comprehensible names and will not be discussed unless required
96  for a proper explanation of the algorithm. In all cases, `source` represents the node for which
97  the multicast tree is being generated, `dests` contains the destination set and `route` stores the
98  tree under construction. The operation of all the proposed algorithms is similar; they start
99  with an empty `route` and build the multicast tree by adding links as required. Examples of
100 routing trees generated with each strategy are shown in Fig. 3.

101 Dimension Order Routing (**DOR**) is the simplest way of generate a multicast tree: follow an
102 oblivious dimension order routing from the source to each destination (first move in X, then
103 in Y and then in the diagonal). Note that only shortest paths are considered and that most
104 communications will share the hops in dimension X or Y. An algorithmic definition of this
105 strategy is shown in Algorithm I. The `dor` function returns the collection of links crossed
106 when applying DOR from `source` to `d`.

107 Longest First Routing (**LFR**) is also oblivious using shortest paths. In this case, the routes
108 are generated travelling first in the dimension with more hops. The definition of this strategy
109 (Algorithm II) is similar to DOR's, but using a different routing function, `lf` which returns
110 the links traversed using longest first routing.

111 Enhanced Shortest Path Routing (**ESPR**) is also based on shortest path routes but is noticea-
112 bly more complex. Each destination, `d`, looks for the closest node, `c`, that it can reach using a
113 shortest path towards the `source` and adds the route from `c` to `d` (using longest first routing)
114 to the solution. If there is no node, then `c` will be the `source`. ESPR creates trees similar to
115 those produced by LFR but avoids creating parallel branches when several nodes are close
116 from each other. Algorithm III shows the pseudo-code definition of this strategy.

117 Finally, Neighbor Exploring Routing (**NER**) is the most complex of the presented strategies.
118 The destinations are sorted from the closest to the furthest. Then each destination, `d`, looks
119 for the closest route point, `p`, (the `source`, other destination or a used link) in its surround-
120 ings and adds the route between `d` and `p` to the solution. This way, as each node is connected
121 to its closest route point, the requirements in terms of network resources are drastically re-
122 duced. Algorithm IV shows the pseudo-code definition of this strategy. Function `sort` ar-
123 ranges the node-list in destinations by their distance to source.

| ```
route=∅
for d in dests
 r = dor(source, d)
 for l in r
  if l ∉ route
   add(l, route)
return route
``` | ```
route=∅
for d in dests
 r = lf(source, d)
 for l in r
  if l ∉ route
   add(l, route)
return route
``` | ```
route=∅
for d in dests
 c = closest(source, d)
 r = lf(c, d)
 for l in r
  if l ∉ route
   add(l, route)
return route
``` | ```
route=∅
sort(dests, source)
for d in dests
 p = surroundings(d)
 r = lf(p, d)
 for l in r
  if l ∉ route
   add(l, route)
return route
``` |
|---|---|---|---|
| Algorithm I. DOR strategy pseudo code | Algorithm II. LFR strategy pseudo code | Algorithm III. ESPR strategy pseudo code | Algorithm IV. NER strategy pseudo code |

## 4  Evaluation of multicast strategies

125 We have compared the four discussed strategies bearing in mind their resource requirements.
126 The figures of merit in the study are: (i) use of network resources, (ii) balanced use of re-
127 sources, (iii) entries in the routing tables and (iv) algorithm complexity – i.e. computing time
128 required to construct the trees. Fig. 4 shows the average resources utilization of $10^5$ runs for
129 a given configuration of $\bar{d} = 32$ and $F = 256$. Looking at the figure, it is clear that each strat-
130 egy offers different characteristics.

131 DOR requires a reduced number of entries in the routing tables and allows constructing mul-
132 ticast trees very quickly but, in turn, employs lots of network resources. Another negative
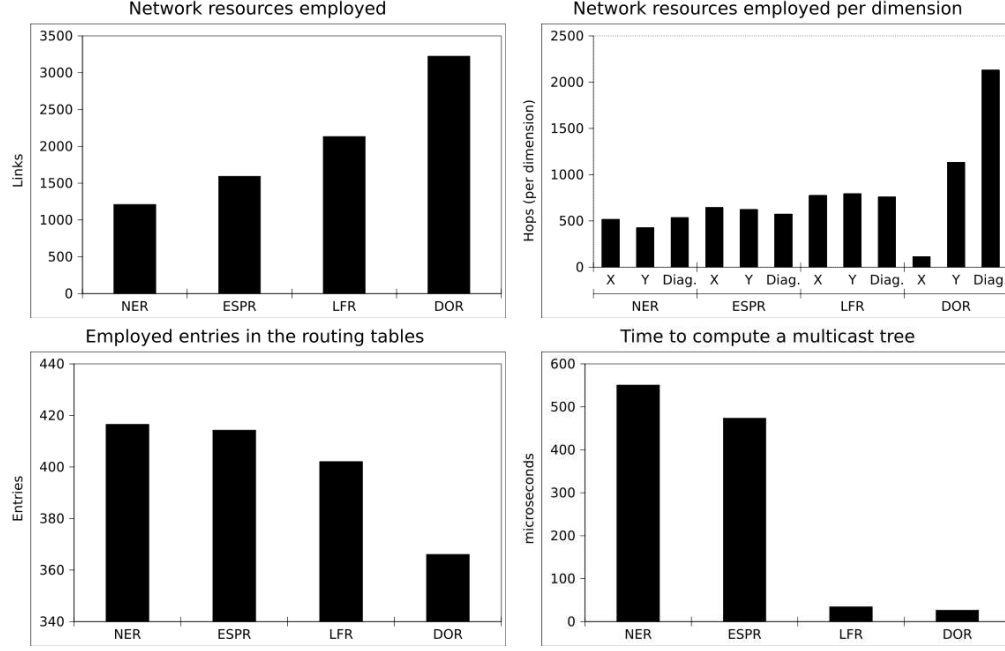133 aspect of DOR is that the utilization of the different dimensions of the topology is very un-

Figure 5. Per strategy resource requirements. $\bar{d} = 32$, $F = 256$ (average of $10^5$ runs).

balanced: the diagonal may become a communication bottleneck. These two negative charac-
teristics motivate the rejection of this strategy except in those cases in which reducing the
number of entries in the routing tables is critical.

LFR is almost as fast as DOR, but drastically reduces the requirements in terms of network
resources. This reduction is obtained at the cost of increasing the number of entries in the
routing tables. However the main advantage of LFR is that it produces the most balanced
utilization of the dimensions so none of them will become a bottleneck. In general, we can
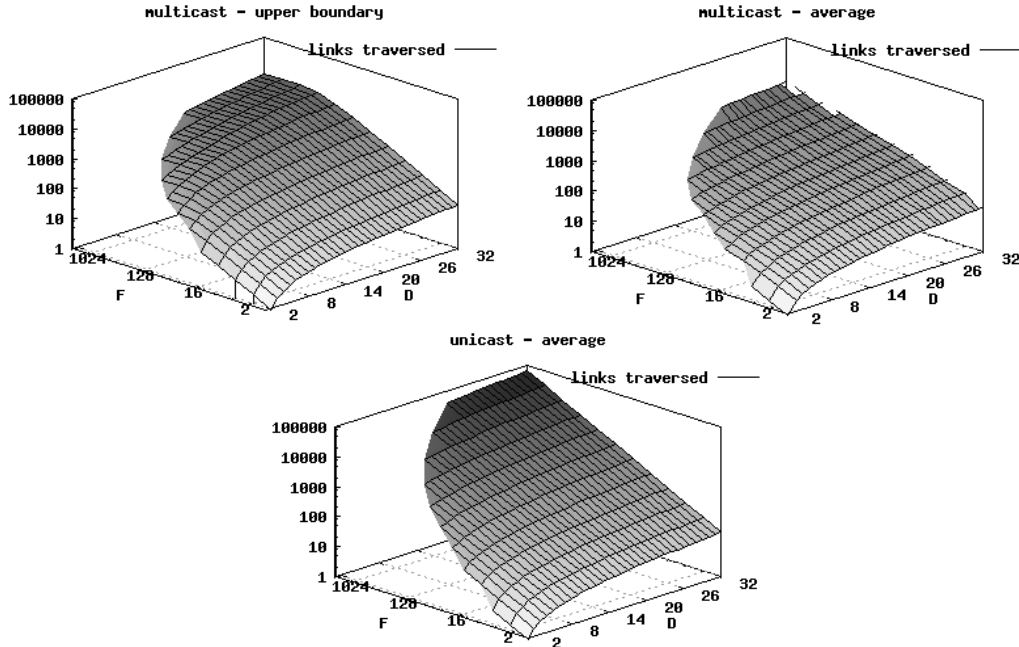state that this is a well-rounded alternative for generating multicast routes.



Figure 4. Network resources required to perform the communications.

142 ESPR reduces the network resource requirements when compared with DOR and LFR and
143 keeps balanced the per-dimension utilization. The costs to pay are increases of the number of
144 entries in the routing tables and of the time taken to produce the multicast trees.

145 Finally, NER produces the multicast trees using the lowest amount of network resources.
146 However results are produced very slowly and require the largest amount of entries in the
147 routing tables, a situation that can be exacerbated because routes are not straight, which
148 would complicate route masking. For this reason NER should only be used in those cases in
149 which reducing network utilization is *critical* and the number of table entries is not.

## 5    Motivation for the use of multicast

151 This section is devoted to formulate analytic models of network requirement to support neu-
152 ral communications, both for multicast and unicast approaches. We derive first the network
153 resources employed to communicate a spike to all connected neurons (*N*). It is measured as
154 the total number of traversed links (hops). We will derive it from $\bar{d}$ and $F$, both of them de-
155 pending only on the workload, not on the actual hardware. In the case of unicast, deriving
156 the average network requirement from these two parameters is straightforward; $F$ packets
157 sent to an average distance of $\bar{d}$ then, the *average* network requirement for unicast will be:

$$N_u = \bar{d} \cdot N \qquad (1)$$

159 In the case of multicast, it is not easy to derive an expression for the average network re-
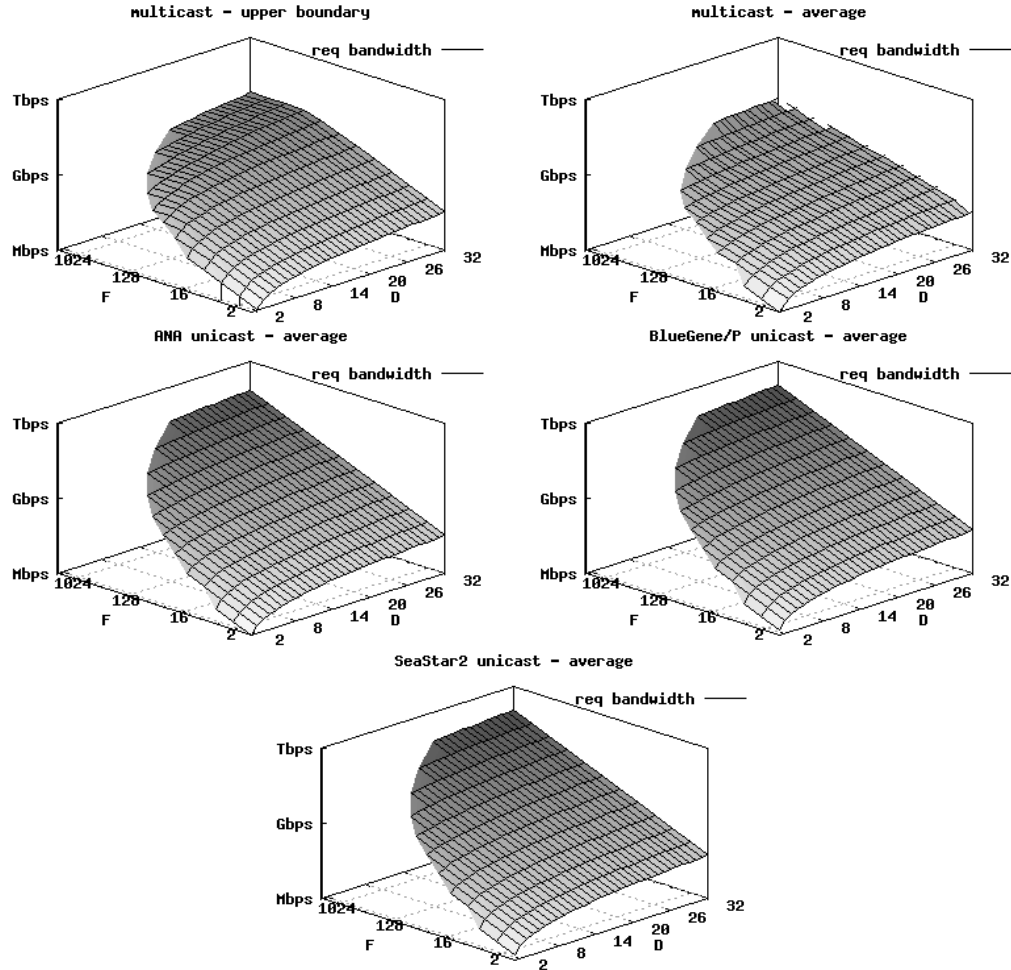160 quirement because this figure strongly depends on the employed strategy. Instead we can



Figure 6. Network bandwidth ($B_L$) required to support regular system operation.

derive an upper bound estimator of network requirement when using multicast traffic. This
estimator takes into account the worst configuration that could exist; this is, when routes to
the destinations share as less as possible of their paths. This estimator can be formulated as:

$$N_M = \sum_{i=1}^{i=2\cdot\bar{d}-1} \min\left(6\cdot i, \frac{F}{2}, 2\cdot 6\cdot(\bar{d}-i)\right) \tag{2}$$

Fig. 5 shows $N_u$ and $N_M$ with a sweep of the values of interest for the parameters $\bar{d}$ and $F$.
For the sake of completeness we also show $N_m$, the average network utilization when using
the previously defined ESPR strategy which has been computed as the average value of $10^5$
random runs. In the figure it is clear that $N_M$ is always lower or equal than $N_u$. This is obvi-
ous as the worst case for multicast is having completely disjoint routes, i.e. unicast behavior.

From (1) and (2) we can derive some other figures considering hardware characteristics.
More specifically in the remaining of this section we will consider five different configura-
tions: (i) worst-case and (ii) average multicast with ANA router (1Gbps, $p_s$ = 40 bits), (iii)
average unicast with ANA router, (iv) average unicast with the Blue Gene/P router (3.2Gbps,
$p_s$ = 64 bits) [14] and (v) average unicast with the SeaStar2 on Cray's XT4 (6Gbps, $p_s$ = 40
bits) [15]. Unless otherwise stated, the system will be modeled in a regular operational
status, i.e. $n_n$ = 18000, $f_s$ = 10Hz and $L$ = 6. We will derive next the network bandwidth re-
quired to support a given system configuration, assuming evenly distributed traffic.

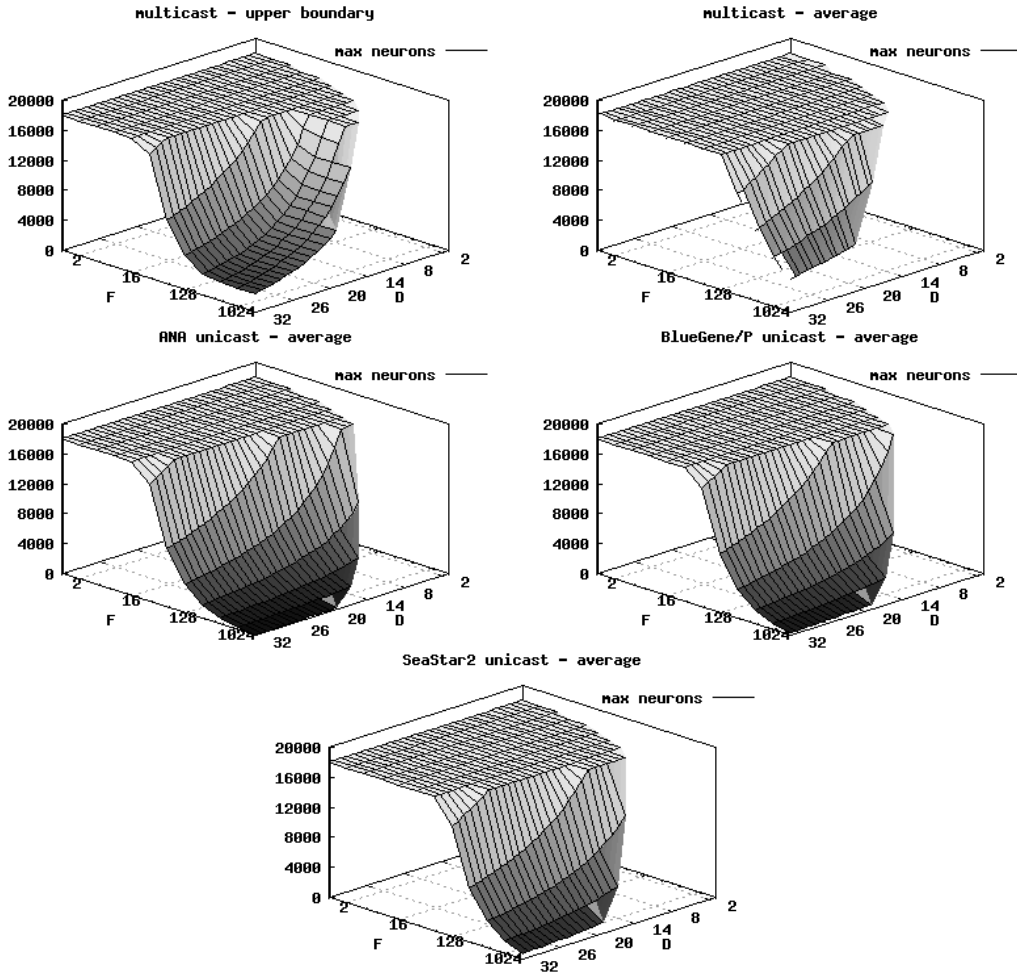$$B_L = \frac{N\cdot f_s\cdot n_n\cdot p_s}{L} \tag{3}$$



Figure 7. Number of neurons per core (firing at 10Hz) that can be supported by the system.
There is an upper bound of 18000 neurons, imposed by the cores..

179 From (1), (2), (3) and the described values for the different systems, we plot in Fig. 6 the
180 network bandwidth required to support the neural simulation. The *average* bandwidth re-
181 quired when using any *unicast* router can be more than one order of magnitude higher than
182 the *upper bound* of *multicast* router. We will close the section deriving an important figure of
183 merit, the number of neurons per-core that the system can support using each configuration.

$$184 \qquad n_n = \frac{B_L \cdot L}{N \cdot f_s \cdot p_s} \qquad\qquad (4)$$

185 From (1), (2), (4) and the discussed parameter values, we can compute the amount of neu-
186 rons that the system is able to execute in each node. For the sake of clarity we will limit the
187 amount of neurons that can be simulated to the maximum acceptable by the cores (18k). This
188 way we will be able to see when the network becomes a limiting factor. This figure is shown
189 in Fig. 7. In these plots it is evident that in the less demanding network configurations (low
190 number of destinations and/or low distances) all networks can support the execution of a
191 regular neural network simulation (plateaus in the plots). However as these figures increase
192 the networks start becoming the limiting factor (slopes). This happens more often with the
193 unicast routers, even when the Blue Gene and SeaStar2 have noticeably higher bandwidths
194 than ANA's multicast router. Furthermore, in those cases in which the network becomes the
195 limiting factor, the multicast alternative is able to support a larger number of neurons; in the
196 worst case unicast routers fall to support only *hundreds* of neurons, while multicast is still
197 able to support *thousands* of them.

## 6    Conclusions

199 In this paper we have highlighted the adequacy of multicast routing for the kind of applica-
200 tions executed on top of ANA in which each neuron must communicate its activation to
201 thousands of other neurons. Our analysis shows that a multicast approach requires noticeably
202 less network bandwidth than a unicast approach and therefore it is able to support the simu-
203 lation of a larger amount of neurons. More specifically we have shown how the multicast
204 router implemented in ANA can achieve better performance than more complex router de-
205 signs such as the SeaStar2 or Blue Gene's router, even though they provide noticeably higher
206 network bandwidth.

207 In addition we described four different strategies to construct multicast trees and measured
208 their resource consumption. Generating multicast trees using DOR reduces the utilization of
209 the entries within the routing tables (a precious resource as there are only 1024 entries per
210 chip), but demands the highest network resources. LFR is a well-rounded solution as the
211 multicast trees can be generated very quickly while keeping resource requirements low
212 enough. More sophisticated strategies which look for routes in their surroundings have been
213 also considered. ESPR searches for connections using always shortest path; in contrast NER
214 searches in all the surroundings even if no shortest path is used. These two strategies further
215 reduce the network requirements but they are two orders of magnitude slower to generate
216 than oblivious strategies.

217 These four strategies were described and implemented bearing in mind a two-dimensional
218 triangular torus topology, but we want to remark that they are general enough to be em-
219 ployed in other cube-like topologies such as the 3D tori implemented in *state-of-the-art* mas-
220 sively parallel processors such as IBM's Blue Gene or Cray's XT families. In fact, the adap-
221 tation would be straightforward simply by using 3D torus routing functions (DOR and LF).

### References

223 [1] *<Anonymized for blind review>*

224 [2] *<Anonymized for blind review>*

225 [3] T Elliott, N Shadbolt. "Developmental robotics: Manifesto and application". Philosophical Trans.
226 Royal Soc., vol. A, no. 361, 2003.

227 [4] P Dayan, LF Abbott. "Theoretical Neuroscience: Computational and Mathematical Modeling of Neural
228 Systems". The MIT Press. 2005.

229  [5] H Brody. "Cell counts in cerebral cortex and brainstem". Alzheimer disease senile dementia and related
230  disorders 7, 1978, pp. 345–351.

231  [6] *<Anonymized for blind review>*

232  [7] *<Anonymized for blind review>*

233  [8] *<Anonymized for blind review>*

234  [9] W Maas, CM Bishop. "Pulsed neural networks". The MIT Press. 1998.

235  [10] S Bhattacharya, G Elsesser, WT Tsai, DZ Du. "Multicasting in Generalized Multistage Intercon-
236  nection Networks". Journal of Parallel and Distributed Computing 22(1) 1994, pp. 80-95.

237  [11] S Coll, FJ Mora, J Duato, F Petrini. "Efficient and Scalable Hardware-Based Multicast in Fat-
238  Tree Networks". IEEE Transactions on Parallel and Distributed Systems 20(9), September 2009

239  [12] Y Breitbart, et al. "Topology discovery in heterogeneous IP networks: the NetInventory system".
240  IEEE/ACM Transactions on Networking 12(3), 2004, pp. 401-414.

241  [13] JK Shapiro, J Kurose, D Towsley, S Zabele."Topology discovery service for router-assisted multi-
242  cast transport". Procs. of IEEE Open Architectures and Network Programming, 2002, pp. 14-24 .

243  [14] S. Alam, et al. "Early evaluation of IBM Blue Gene/P". Proc. of the 2008 ACM/IEEE Conference
244  on Supercomputing, pp. 1-12, 15-21 Nov. 2008, Austin, TX, USA. DOI: 10.1109/SC.2008.5214725

245  [15] AR Sadaf, et al. "Cray XT4: an early evaluation for Petascale scientific simulation". Proc. of the
246  2007 ACM/IEEE Conference on Supercomputing, pp. 1-12, 10-16 Nov. 2007, Reno, NV, USA.