

A parallel solution with GPU technology to predict Energy Consumption in Spatially Distributed Buildings using Evolutionary Optimization and Artificial Neural Networks

J.R.S. Iruela^a, L.G.B. Ruiz^a, M.C. Pegalajar^a, M.I. Capel^b

^a*Department of Computer Science and Artificial Intelligence, University of Granada, Spain*

^b*Department of Software Engineering, University of Granada, Spain*

Abstract

Today all governments talk about climate change and its consequences. One of the ways to tackle this problem is by studying the energy consumption of the buildings around us. The study of energy consumption may give us relevant information to make better decisions, and thus reduce costs and pollution. However, ANN training models, in order to achieve those goals, has a high computational cost in terms of time. To solve that problem, this paper presents a GPU-based parallel implementation of NSGA-II to train ANNs whose evaluation has also been implemented in a parallel GPU scheme. Our methodology is designed to predict the energy consumption of a series of public buildings, and thus, to model consumption, save energy and improve the energy efficiency of these buildings without compromising their performance obtaining the prediction in a very short period of time. We compared the sequential implementation of the evolutionary algorithm NSGA-II with our new version developed in parallel and the parallel implementation gets better results in much faster execution time.

Keywords: Energy consumption forecasting, Artificial Neural Networks, GPU, Evolutionary Algorithm

1. Introduction

- 2 In the last 50 years, the planet has degraded more than in 100 centuries [1].
- 3 According to the United Nations, 60% of all-natural resources are exploited

4 in a non-sustainable manner. It becomes, therefore, necessary to carry out
5 efficient use of the energy resources for the time to come. Moreover, pre-
6 dicting energy consumption (EC) in a set of buildings associated with an
7 institution is a difficult task to perform Feng et al. [2] that large companies
8 and administrations around the world are aiming today Liu et al. [3], Park
9 et al. [4], and to solve it is of paramount importance for the achievement of
10 efficient EC by 2020. One of the most important areas in which to tackle this
11 problem is in the local EC (EC) forecasting, which will allow us to anticipate
12 future events and, in this way, propitiate to make wise decisions regarding
13 energy savings. In this realm, to analyze the EC collected by sensors at an
14 individual level, over delimited zones or buildings, can help to reduce energy
15 costs and environmental impact created by energy production. To make,
16 however, predictions of EC within a time range, e.g., a forecast temporal
17 horizon (time of prediction of the model) larger than 1 hour is not usually
18 useful, and it is actually an extremely difficult task to carry out because of
19 the massive amount of data to be processed within an established time-span.

20 In this work we propose to use ANN to predict EC in buildings that, to
21 the best of our knowledge, has not been tried before in conjunction with a
22 multi-objective function optimization over a population, which is carried out
23 with the well-known NSGA-II algorithm to optimize both models. In order
24 to obtain the necessary high performance of the application for making useful
25 EC predictions, we propose an implementation based on GPU to leverage all
26 the performance capabilities of that implementation of the models can bring
27 to us. The models were tested by using different buildings of the University
28 of Granada, and the obtained results, by making a daily prediction in an
29 approximate time of 60 seconds, are shown.

30 Several studies provided solutions to the problem of predicting EC in
31 buildings by using Evolutionary Algorithm (EA) and ANN [5, 6]. However,
32 the main drawback of those methods lies in their non-dependable time re-
33 sponse. So far, some approximations have been proposed to solve this prob-
34 lem [5]. As a consequence, there is still a lot of work to do in this research line
35 whereas some interesting GPU implementations of EAs have been proposed
36 in [7, 8, 9, 10, 11] in which different data structures, configurations, data size
37 and complexity were studied to solve the problem. The main issue in these
38 solutions lies in the fact that all of them only parallelize the repetitive process
39 of the evolutionary algorithm needed to find the fitness of the population,
40 without going any further in terms of the parallel structure of the sub-tasks
41 of that algorithm. To fill up this gap, we propose here to advance in the

42 design of the EA by taking full advantage of parallelism not only in the par-
43 allelization of the EA intrinsic scheme but also at the level of the evaluation
44 of each individual of the mentioned population. Thus, by the evaluation of
45 the ANN associated with each individual, we designed a parallel solution to
46 deal with the fitness function calculation of the population of the EA.

47 Although, as mentioned above, the problem of parallelization of EAs has
48 been addressed multiple times with success, however the proposed solutions
49 up today present the hurdle of ignoring time restrictions. The function that
50 usually takes the longest time to be calculated is the fitness function, which
51 is the weak point regarding the performance of all the previously mentioned
52 parallelization proposals of these algorithms [10]. In order to achieve the
53 above requirements and solve the problem of EC prediction, a new design of
54 a parallel multi-objective evolutionary algorithm of the NSGA-II is presented
55 here. By an original assignment of blocks to threads of the GPU that consists
56 of starting a processing structure of 2 levels of GPU-threads, we obtained that
57 the evaluation with genetic operators of each individual of the population was
58 independently performed. Thus, the iterative loop performed in each ANN to
59 evaluate all data is done in parallel without safety violation (race conditions,
60 non-mutual exclusion, . . .) that might hinder one another evaluation.

61 The main point to solve the mentioned problems is how to cope with the
62 high time cost of calculations. To do so, parallel computing was adopted as a
63 fundamental tool for speeding up the creation of new prediction models [12].
64 However, there are two different methodologies within the area of parallelism
65 to follow. The first one, based on the parallelization of the algorithmic parts
66 of the calculation only at the CPU level, which is expensive in terms of
67 computation resources and limited by the small acceleration factor of the
68 calculations. The second approach is aimed at the deployment of software at
69 the GPU level, and thus the entire algorithm and data are programmed as
70 CUDA kernels on the GPU. This approach is usually more economical than
71 the former one and has better performance since it can have up to 6000 cores
72 available (e.g., NVIDIA Tesla graphics) that are efficient for the execution
73 of matrix operations. For this reason, we selected this parallelism in GPU,
74 which proved to work successfully on modelling EC in our University.

75 Our proposal has proven to be an excellent solution to deal the high
76 time cost needed to obtain good forecasting models, achieving substantial
77 improvements in terms of time and an enormous difference between sequential
78 and parallel evaluation of the population. The rest of the work has been
79 divided into the following sections: Section 2 presents a brief description

80 of the GPU architecture. Section 3 describes the proposed algorithms to
81 solve the power consumption prediction problem. Section 4 explains how the
82 proposed multi-objective evolutionary algorithm has been parallelized and
83 implemented on GPUs with the goal of reducing the run-time of the entire
84 algorithm. Section 5 presents the results of the experiments, which were
85 obtained from the UGR data-set and the implementation of the developed
86 algorithm. Section 6 concludes with a summary of the main conclusions
87 obtained and future research work.

88 **2. Related work**

89 Studying energy consumption can help to reduce energy costs and the envi-
90 ronmental impact created by its production. For this reason, to improve
91 energy efficiency in buildings is a concern in many situations nowadays.
92 When studying energy consumption we are mainly interested in obtaining
93 behavioural patterns [13] to detect anomalies [14, 15] and to make predic-
94 tions of energy demand [16], as well as to acquire consumption profiles of
95 different buildings[17].

96 In the related literature there are many techniques to predict energy
97 consumption in different scenarios, such as ARIMA[18], Grey Models[19],
98 Regression Tree[20], Support Vector Machine[21], Artificial Neural Networks
99 (ANN) [22] or even combination of different techniques to optimize solutions,
100 such as, neural fuzzy systems[23], regression trees using clustering methods
101 [24], ANN and Evolutionary Algorithm [25], among others Wu et al. [26], Fan
102 et al. [27]. Particularly, ANN has proven to be a powerful technique to pre-
103 dict energy consumption. Several kinds of ANN can be found in literature to
104 solve these problems. In [28] an Artificial Neural Network Inverse is utilized
105 to optimize multiple variables in an absorption heat transformer in order to
106 provide a method for optimizing energy usage. An study of local perceptron
107 networks can be found in [29] where authors discuss several approaches to
108 carry out forecasting on time series. They demonstrate that ANNs many of-
109 ten outperform traditional prediction techniques. For a deeper analysis and
110 a comprehensive study of artificial neural networks, a fundamental reference
111 is [30].

112 Nevertheless, all those techniques lack of a method to optimize and im-
113 prove their results, and therefore, they are highly improvable very often.
114 Hence, other techniques must be used to deal with this drawback, and these
115 are Evolutionary Algorithms (EA)-based ones Krzywanski et al. [31]. The

116 EAs bring the possibility of optimizing models in many senses, e.g., sim-
117 plifying its structure, dealing local minimum, optimizing relations among
118 parameters. Several studies have provided solutions to the problem of pre-
119 dicting energy consumption in a series of buildings using EA, with which
120 sequential solutions are proposed to obtain accurate results. However, the
121 cost due to the execution time of EA is often seen as very high to get results
122 on time [5], especially in those systems which presents time criticality. On
123 the other hand, EA has proven to be very suitable to obtain its parallelization
124 and has successfully solved numerous problems. In [9] a flight planner was
125 developed with EA, the results are compared in the study with the sequential
126 version. With the parallel algorithm, the execution time was reduced by a
127 factor of 290x compared to sequential execution. A similar study with differ-
128 ent techniques is carried out in [7] where the problem of energy consumption
129 forecasting of a hydroelectric power plant was solved in real-time successfully
130 by using a parallel GPU implementation of a evolutionary algorithm. In that
131 study, results were shown competitive compared to the serial version. And
132 several approaches were proposed in many other studies [8, 10], where dif-
133 ferent proposals of parallel-based evolutionary algorithms were investigated,
134 and so distinct different data structures, configurations of the algorithm, data
135 size and complexity of the problem, have been compared over recent years
136 [11].

137 Since in our approach we do not only look for accuracy but also intend
138 to speed up the algorithm execution, we focus the research on parallelizing
139 a multiobjective evolutionary algorithm with neural networks in order to
140 obtain the results in the shortest period. In order to achieve these objectives
141 of precision and acceleration, the execution of the evolutionary algorithm
142 has been parallelized along with the evaluation of individuals through neural
143 networks.

144 **3. Preliminary concepts**

145 GPU was originally designed to be used in video games or image rendering
146 [32] as it presents a high degree of parallelism that is required in these kinds
147 of problems. But as of today it has been assumed in many fields of research,
148 and it has yielded the technical term known as GPGPU (General-Purpose
149 computation on GPU) [32, 33]. Note that this section is intended to introduce
150 some essential concepts concerning the GPU, so the reader can skip this
151 section if he is an expert in this field, or otherwise he should not since the

152 association between our solution and the structure of the GPU could not be
 153 easy to follow.

154 Therefore, the GPU’s particular structure is the first thing we must
 155 point out. Thus, figure 1 compares the architectures of the CPU and GPU.
 156 Whereas CPU reserves more space for control units and their associated stor-
 157 age, and thus the remaining space is intended for the logical arithmetic units
 158 (ALU)(see Figure 1a), on the contrary the GPU allocates more space to ALU
 159 and less space for control and storage units in order to reach a higher degree
 160 of explicit parallelism and throughput(see Figure 1b).

161 CPU is biased to process serial instructions that use a lot of memory, while
 162 a GPU is a better choice for processing parallel instructions that use much less
 163 memory. Another difference is that CPU has few powerful processing cores,
 164 while GPU has thousands of these. The GPU architecture works better with
 165 highly parallelizable applications while the CPU gets more computational
 166 performance by parallelizing applications with longer sequential code [34].
 167 This point is important as we leverage these features in the representation of
 168 our solution (the ANN) as it shows a high level of potential parallelism [7],
 169 which will be detailed in the next sections.

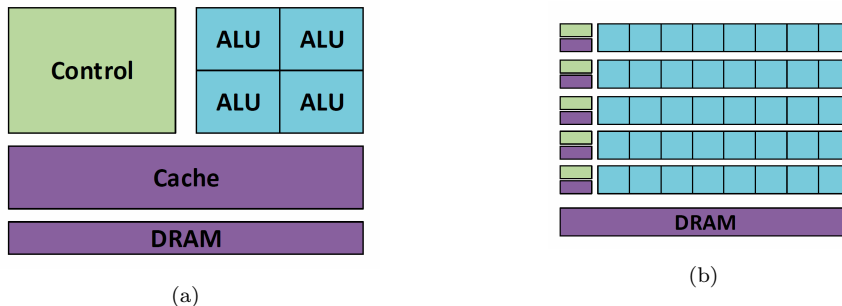


Figure 1: CPU (a) and GPU (b) architecture.

170 3.1. GPU notations

171 The architecture and software that support the use of GPUs allow us to use
 172 them for general-purpose calculations. There are different frameworks with
 173 which it is possible to develop software and generate code for the GPU. The
 174 two most important frameworks are OpenCL [35] and CUDA [36], both of
 175 which are platforms that allow the GPU to be used for high-level program-
 176 ming. OpenCL is a cross-platform programming language with no hardware

177 constraints that is intended to be used on heterogeneous platforms, whereas
178 CUDA is a parallel computing platform that includes a compiler and tools
179 for programming algorithms in NVIDIA GPUs. The fact that CUDA was
180 developed exclusively to work with NVIDIA graphics cards gives it a clear
181 advantage over OpenCL by obtaining better performance results [37]. On
182 the other hand, OpenCL is supported by more software applications than
183 CUDA actually is.

184 In the CUDA environment, CPU is referred as the *host*, and GPU as the
185 *device*. The set of instructions running on GPUs are structured into functions
186 called kernels. These kernels will be the most critical point in our work as
187 the success of our solution regarding performance improvement depends on
188 the optimal implementation of them. They will be the operations that our
189 algorithm will carry out, i.e., the creation of the solutions, its evaluation,
190 etcetera.

191 A kernel can be launched from either the CPU or GPU. The advantage
192 of calling a kernel from the GPU is that it reduces the bottleneck caused by
193 the necessary communication between GPU and CPU memories, and thus
194 we take advantage of this feature and therefore try to keep data into the
195 GPU as long as possible to speed up the computations. The kernels are
196 executed by the different threads available on the GPU. As can be seen
197 in the figure 2 that shows the internal structure of the GPU, threads are
198 organized in blocks, which have a 3D structure that defines the number of
199 threads in each dimension. At the same time, blocks are organized in a 2D
200 grid. When a kernel is launched, it runs as a grid of parallel threads. The
201 threads in a grid are organized in 2 hierarchy levels but only to reference
202 them in the program. At the top level, each grid is made up of one or more
203 blocks of threads. All blocks have the same number of threads. Each block
204 is unequivocally identified by two coordinates assigned by CUDA, for more
205 detail we refer to [36].

206 NVIDIA uses "Single Instruction Multiple Threads" (SIMT) as the run-
207 time model. The threads of each block are executed in sets of 32 threads
208 called *warps*. All threads in the same warp must execute the same instruc-
209 tion at the same time instant. When some threads belonging to the same
210 warp need to execute a different instruction, a divergence occurs, and the
211 group of threads of the same warp is executed sequentially, thus no simulta-
212 neously. In order to minimize convergence, it is advisable to avoid instruc-
213 tions that produce thread bifurcation. However, in our solution, this would
214 be required as some functions are needed for it. For instance, the mutation

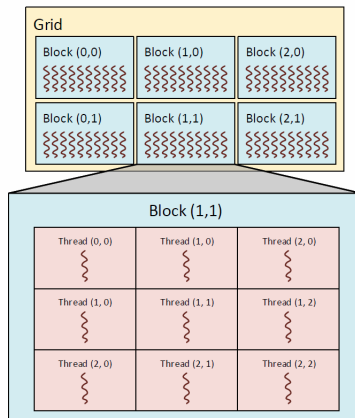


Figure 2: Structure block and threads.

215 operator, because some individuals will change their values and others will
 216 not, depending on a probability. All of this will be discussed in detail later
 217 on.

218 CUDA defines its own memory hierarchy, which includes, among the dif-
 219 ferent memories shared and non-shared by threads, there exist global memory
 220 and shared memory. Global memory can be accessed and modified from the
 221 host and the device, therefore all the threads can access to the global mem-
 222 ory. Shared memory is faster than global memory but its capacity is smaller.
 223 All the threads of the same block can access the shared memory. All data
 224 stored in this memory are lost when the block ends its execution. Data that
 225 are frequently used by the same block must be moved to the block's shared
 226 memory, and so the performance offered by the GPU can be improved.

227 4. Methodology

228 In this work, we combine two machine learning techniques to address the
 229 problem of EC forecasting applied to prevent energy waste in buildings. The
 230 first technique is artificial neural networks (ANN), which are used to model
 231 and predict energy usage. The second one is the multi-objective evolutionary
 232 algorithm NSGA-II deployed in our study to provide a procedure to obtain
 233 optimal forecasting models, i.e., the neural network with the least number of
 234 hidden neurons and which makes the least error in the set of examples. Note
 235 that in our solution, we will not use any classic machine learning algorithm,

236 but we will use the same EA to modify the weights of the ANNs so we carry
237 out the two task we pursue at once. In this way, we minimize the error
238 and the complexity of the model simultaneously. Since the cost of deploying
239 these techniques for EC forecasting is very high in terms of time, the models
240 supporting them were designed for implementation on a GPU framework.
241 With this fact in mind, in the following subsections, we will go over these
242 concepts.

243 *4.1. Artificial neural networks*

244 ANNs are algorithm based on brain functioning which are widely known for
245 their good results, and used in tasks such as process control, optimization,
246 pattern recognition, prediction, etc. [38, 39, 40].

247 One of the most recognized models, used in multivariate regression prob-
248 lems, are the neural networks called feed-forward. Its topology can be ob-
249 served in 3 and, as this figure shows, the ANN connect their input neurons
250 to the hidden neurons, and then these to the output neurons. This type of
251 ANN uses the values of the input neurons and the weights assigned to each
252 link, which connect them to the hidden neurons, to assign values to the states
253 of the network. Once the states of the hidden neurons and the weights asso-
254 ciated with the output layer are calculated, the value for each of the output
255 neurons will be finally obtained.

256 On the other hand, Back-propagation (BP) is one of the most popular
257 algorithms used for training feed-forward neural networks. BP is a method
258 based on gradient descent and may be subject to convergence at a premature
259 local optimal solution. As a consequence, the solutions found by this tech-
260 nique depend on the fair initial randomness and on the lack of that fairness
261 during the calculus of solutions, which is not easy to prevent [41]. Therefore,
262 ANNs can make errors caused by multiple local optimal solutions and the
263 application of the BP method may yield results from a local solution that are
264 not the global optimal [42]. In addition to this, the use of ANN, like many
265 other techniques, implies determining several criteria before its execution,
266 such as to obtain the number of neurons and the selection of the training
267 procedure. These decisions can be made by an expert through trial and er-
268 ror, as usual, on account of high computing time of carrying out the entire
269 procedure, it would be better to reduce the number of trials. We must keep
270 in mind that making these decisions is generally a difficult task, since the
271 change of a single parameter, among the possible configurations of an ANN,

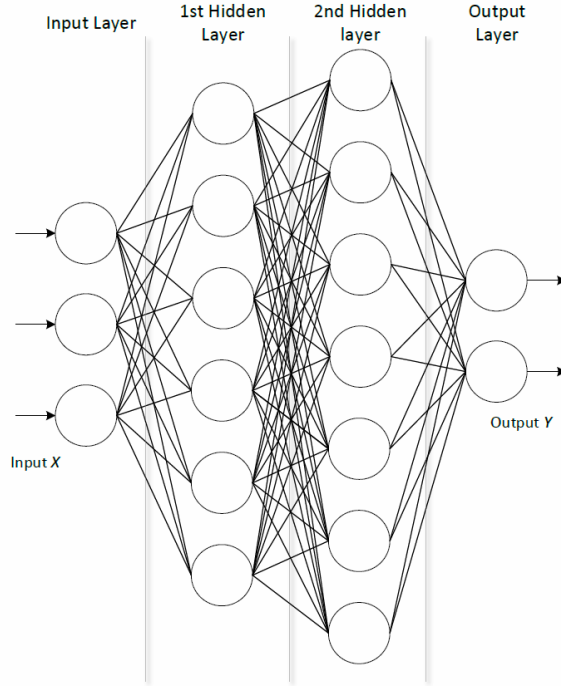


Figure 3: Typology of a feed-forward ANN.

272 may lead to a substantial negative effect on the performance of the learn-
 273 ing algorithm, which makes it difficult to find the optimal global and how it
 274 affects the performance of the neural network. As a solution to those draw-
 275 backs, we propose to deploy an evolutionary algorithm capable of carrying
 276 out these tasks at the same time. In other words, we utilize the evolutionary
 277 algorithm not only to adjust some of ANN's parameters, such as the number
 278 of neurons, but also to combine the information of the different solutions to
 279 learn about the data and thus to provide an implicit training mechanism,
 280 thus fitting in this way the different weights of the models. Indeed, this ap-
 281 proach has shown its potential to solve many problems like this one up to
 282 now [43][44].

283 4.2. Evolutionary Algorithms and Multi-objective Evolutionary Algorithms

284 Evolutionary algorithms (EA) are a global search method based on a pop-
 285 ulation of individuals (associated with sub-optimal solutions), inspired by
 286 the natural mechanisms of the genetic evolution of biological species. On

287 the whole, a EA employs three operators: selection, mutation and crossover.
288 Those procedures are used to generate new solutions (called individuals) that
289 will lead to exploration of new search points. EA may become very effective
290 in a wide variety of problems, however, their execution time may become a
291 limiting factor for deploying it when programming the solution of some large
292 problems. This fact is due to the great number of operations which should
293 be done to achieve the final solution. One of the most expensive, in terms
294 of computation time, of such procedures is the evaluation of the population,
295 also called calculus of the *fitness* function. This is because all individuals
296 must be evaluated many times to estimate the accuracy of their associated
297 solutions. Fortunately, evaluations can be carried out independently for each
298 individual in the population, and it makes the EA a candidate method to be
299 parallelized with high performance [45].

300 In the specific literature, EA are usually designed to solve a problem
301 with a single objective, e.g., error, saving, space, etc. However, it is not
302 uncommon in real problem solving to seek solutions that require finding
303 values of parameters according to multiple criteria [46]. However, it is normal
304 that the optimal way to combine the different objectives is not known or it
305 is difficult to do so. For these problems with multiple objectives, there is
306 not always a single solution that can be considered the best one, but a set of
307 solutions that represent the best compromises between the different criteria.
308 This set is called the optimum Pareto set and its representation in the target
309 space is called the Pareto front [47, 48].

310 On the other hand, EAs are recognized for their great effectiveness in
311 solving difficult optimization problems [6] but addressing them requires large
312 amounts of computing resources as every solution has to be evaluated many
313 times. In this vein, it becomes necessary to take advantage of parallelization,
314 which has proven to be the fastest and most effective approach [49] [50].
315 For these reasons, our design, developed by using this technology, will be
316 explained in the sequel.

317 *4.3. Multi-objective evolutionary algorithms to optimize ANN*

318 In this work, the multi-objective evolutionary algorithm NSGA-II is used to
319 combine the knowledge of every solution so as to optimize the topology of the
320 networks and to train the weights of the ANNs. By doing so, one can obtain
321 the least error with the optimum number of neurons [51]. This algorithm sets
322 the ground to develop a software component of the method based on fair-
323 randomness of solutions choice that allows us to better explore the search

324 space, making it easier not to fall into local optimal solutions but to improve
 325 solution searching by performing several procedures properly optimized.

326 In the following section, we will detail and adapt all the components of
 327 the NSGA-II algorithm in order to obtain an optimal ANN with the final
 328 aim of predicting EC in buildings.

329 4.3.1. Coding of individuals

330 In the calculation of individuals two parts are differentiated: a first one that
 331 will represent the number of hidden neurons in the sought networks and a
 332 second part that will represent the set of associated weights to the links bound
 333 to these neurons; note that an appropriate value of the weights will provide
 334 better forecasting in the model. These two parts can be seen represented
 335 graphically in figure 4. This figure shows a diagram with the coding used
 336 in this study, where N is the number of hidden neurons that make up the
 337 network and W_{ij} and V_{ij} are the different weights of the network. The gene
 338 associated with the number of hidden neurons will be an integer value and
 339 the genes associated with the set of network's weights will be coded as real
 340 values.

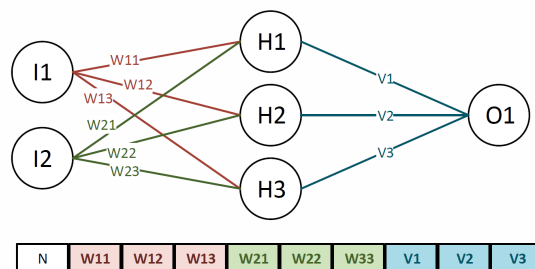


Figure 4: Data structure.

341 4.3.2. Objectives to optimize

342 In this study, the main objective is to obtain the optimum neural network.
 343 To do so, the following objectives have been taken into account:

344 **The prediction error:** In our problem, for an ANN, we can say the lower
 345 the error is, the better solution is found. Therefore, to minimize the error
 346 obtained in the training set will be one of the objectives of our method.

347 In equation 1 we see the definition of MSE, this refers to the first objective,
 348 where T is the number of instances available for training, $Y(t)$ is the expected
 349 output and $O(t)$ is the output calculated by the network at time t .

$$\text{Min} \{f_1(s)\} = \min \left\{ \frac{1}{T} \sum_{t=1}^T (Y(t) - O(t))^2 \right\} \quad (1)$$

350 **Number of hidden neurons:** This objective allows us to compare networks
 351 of different sizes with the intention of finding networks with fewer hidden
 352 neurons. In function 2 you can see the definition of this objective, where $h(s)$
 353 is the number of hidden neurons for the network s .

$$\text{Min} \{f_2(s)\} = \min \{h(s)\} \quad (2)$$

354 Eventually, the NSGA-II algorithm will focus on finding individuals that
 355 minimize these objectives, i.e., solutions that have the lowest prediction error
 356 and obtained with the simplest model architecture. It is common to find a
 357 better solution if the model is more complex, therefore discovering a simpler
 358 model while keeping the error level of a more complex one is not a trivial
 359 task [43][44][52]. Next sections will detail the whole scheme of the proposed
 360 algorithm.

361 4.3.3. Initialization of the population

362 Each of the individuals who are part of the population will represent a possi-
 363 ble, non-optimal, solution to the problem. Each solution represents a distinct
 364 ANN. The generation of the individuals is made randomly, in order to explore
 365 the most of the search space. To generate the population, in our problem,
 366 each individual will be assigned a random number of hidden neurons between
 367 a previously declared minimum and a maximum. Therefore, in this study,
 368 4 and 32 respectively were established as these values. Once the *length* of
 369 the individual has been defined, its structure will be completed by randomly
 370 assigning weights between a previously defined range [-30, 30].

371 4.3.4. Fitness function

372 The fitness function is used to know how good is a solution. In our case
 373 and, since it is a classic evolutionary algorithm, it is normal to choose as
 374 fitness function the error returned by the neural network evaluated with the
 375 training set. In that case, as it is a multi-objective algorithm, fitness will be

376 determined by 2 objectives, the error obtained by ANN and the number of
377 hidden neurons.

378 *4.3.5. Selection*

379 Selection is the first genetic operator to be used in a EA. This operator,
380 based on the fitness of each individual, is used to improve the search for
381 parent chromosomes. The main objective of the selection is to promote and
382 make survive parents with better fitness. There are several strategies to carry
383 out the selection, some of the most used are tournament selection, random
384 selection and roulette selection.

385 In this work, the selection by tournament has been used, in which two
386 parents are chosen at random among all the population and these are con-
387 fronted by means of a comparison of their fitness.

388 *4.3.6. Crossover*

389 The crossing is a genetic operator that is in charge of generating new children
390 from the parents. This operator simulates the mating and genetic reproduc-
391 tion of nature. For getting so, this operator selects two parents, and then
392 two things can happen:

- 393 1. Both parents have the same number of hidden neurons; in this case, two
394 children are generated of the same length of the parents, and the genes
395 of the children will be obtained by crossing the ones of both parents.
- 396 2. The parents have a different number of hidden neurons; in this case,
397 two children will also be generated. The first child will have the same
398 number of hidden neurons as the first parent and the second child will
399 have the same number of hidden neurons as the second parent. The
400 genes of the children will be generated by crossing the genes of both
401 parents. When the number of genes is different, then these will be
402 completed with the longest parent's genes.

403 The heuristic Wright's crossover [53] has been used in this work, and its
404 application has yielded good results, indeed [43, 54]. Figure 5 shows graph-
405 ically the heuristic Wright's crossover where two parents are crossed based
406 on a crossing probability if the children do not cross because the probability
407 of crossing is not satisfied, the parents replace the children in the next gen-
408 eration. If they can finally be crossed, the children will be obtained in the
409 following way:

410 Let a_i and b_i be the parents to be crossed and r is a random real number
 411 (float) belonging to $[0,1]$ then a child is obtained by the following equation:

$$h_1 = r \times (b_i - a_i) + a_i \quad (3)$$

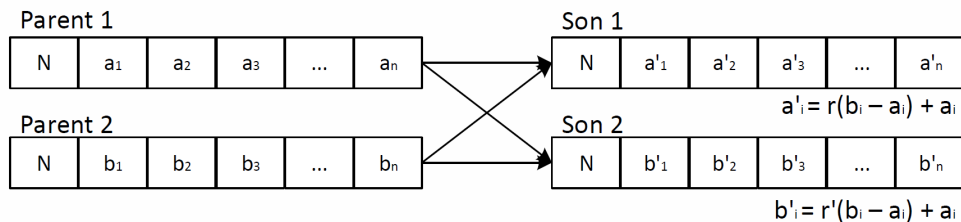


Figure 5: Heuristic Wright crossover.

412 4.3.7. Mutation

413 The mutation operator is applied in EA to preserve the diversity of the
 414 population. This operator creates new individuals by slightly modifying or
 415 drastically the genes of existing individuals. Thus, in our work, two types of
 416 mutation can be distinguished: structural and genetic, respectively. Figure
 417 6 shows a graphical representation of both mutations. While the structural
 418 mutation modifies the topology of the network (neuron numbers) and causes
 419 a drastic change in the individual, the genetic mutation modifies one or
 420 several genes corresponding to the weights of the network. When the size of
 421 the network is changed, the chromosome must be restructured by expanding
 422 or decreasing its size; this benefits the exploration of the search space.

423 In the following sub-section, the implementation of the proposed NSGA-II
 424 algorithm is detailed.

425 4.4. NSGA-II description

426 The NSGA-II (Elitist Non-Dominated Sorting Genetic Algorithm) was pro-
 427 posed by Deb et al., in 2002. Since this problem requires to work with several
 428 targets to optimize, it is necessary to introduce the concept of dominance.
 429 Therefore, prior to explaining the algorithm, a series of definitions must be
 430 given to properly understand how NSGA-II works:

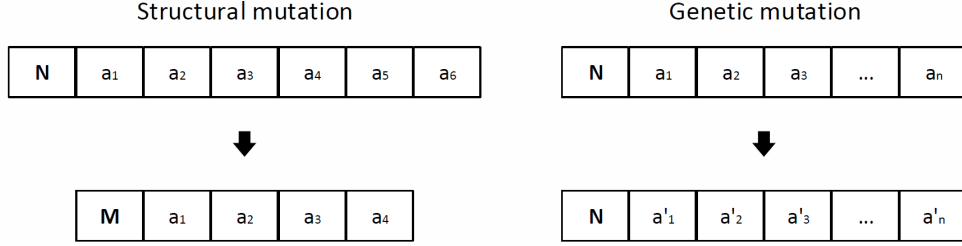


Figure 6: Structural and genetic mutation.

431 **Definition 4.1.** Solution A dominates solution B if the following two con-
 432 ditions are met:

- 433 1. If $A \leq B$ for all M objectives and
- 434 2. If $A < B$ for at least one of the M objectives.

435 In the event that one of the above conditions is not met, solution A
 436 would not dominate solution B. For example, if a ANN A has better error
 437 and less neurons than another ANN B, then A dominates B; otherwise, it
 438 doesn't. Considering a population of N solutions each with M values of target
 439 functions, the following procedure is used to find the non-dominated set of
 440 solutions [55]:

- 441 1. Do $i = 1$.
- 442 2. For all $j \neq i$, compare solutions x^i and x^j to determine dominance,
 443 using the 2 conditions cited above.
- 444 3. If for any j , x^i is dominated by x^j , mark x^i as dominated. Increase i
 445 by one and go to step 2.
- 446 4. If all solutions ($i = j$) in the set have been considered, go to step 5;
 447 otherwise increase i by one and go to step 2.
- 448 5. All solutions that are not marked as dominated are non-dominated
 449 solutions.

450 By applying the above definition, we can divide the population of indi-
 451 viduals into different groups according to their dominance number. The first
 452 frontier will be formed by individuals who are not dominated by any other
 453 individual in the population and then, the dominance number is 0. The sec-
 454 ond and following frontiers will be calculated by performing an iterative loop

455 where individuals who are not dominated by other individuals of the same
 456 group are singled out and the rest are inserted into a new front. This process
 457 will be repeated until all individuals in the population are organized into
 458 fronts. Note that the first front is a set that only contains the individuals
 459 non-dominated by others, and thus the set of P solutions associated to these
 460 individuals is known as the Pareto front [51].

461 In order to avoid the selection of two *too similar* parents, NSGA-II uses
 462 the individual's attribute known as *distance crowding*, which allows to pre-
 463 serve the diversity of the population. This procedure guarantees the diversity
 464 of the population within the same Pareto front by blending all the objectives.
 465 In this way, when the population converges towards the optimal Pareto front,
 466 the algorithm ensures that the solutions are distant enough to acknowledge
 467 dissimilarity from each other [55].

468 Having defined the main concepts necessary for NSGA-II [34], we will
 469 now explain how it works, firstly by showing the pseudocode and then a
 470 brief explanation of it:

```

1 generationsCount = 0;
2 currentPopulation = generateRandomIndividuals(N);
3 evaluateFitness(currentPopulation);
4 nonDominatedSort(currentPopulation);
5 while generationsCount < X do
6   while newPopulation.count < X do
7     parent1, parent2 = select(currentPopulation);
8     child1, child2 = crossover(parent1, parent2);
471 9     mutate(child1, child2);
10    newPopulation.Add(child1, child2);
   end
11   nextGeneration = currentPopulation + newPopulation;
12   nonDominatedSort(nextGeneration);
13   currentGeneration = selectBestN(nextGeneration);
14   generationsCount = generationsCount + 1;
end

```

Algorithm 1: NSGA-II.

472 Initially, the NSGA-II algorithm creates a random population of P_0 par-
 473 ents, on line 2. The population is sorted (within different fronts) on line 4
 474 according to their number of non-dominance, once they have been evaluated
 475 (see lines 2-3). Using the selection, crossover and mutation operators, the

476 descendant population Q_0 will be generated, on lines 7-9. The population of
477 P_t parents of size N is used to create the descendant population Q_t of size
478 N , line 10. The two populations are then joined to form R_t of size $2N$, line
479 11. Once the populations have joined, a non-dominated operator is used to
480 classify the R_t population on different Pareto fronts and the new population
481 is generated according to that procedure, as shown on lines 12-13. First, the
482 individuals belonging to the best non-dominated front F_1 will be added, then
483 individuals from the second front F_2 , third front F_3 , \dots , so on a so forth. The
484 adding of individuals to different fronts continues until obtaining N fronts,
485 where N is the number of solutions.

486 For a more detailed description of the implementation of the algorithm
487 see [56].

488 **5. Parallel Multiobjective Evolutionary Algorithms for ANN opti-** 489 **mization**

490 So far, the general scheme of the implemented EA has been introduced as
491 well as many necessary concepts for understanding our solutions. In this
492 section, the bulk of our contribution is presented. This section details the
493 strategy that has been carried out to parallelize the multi-objective evolu-
494 tionary algorithm NSGA-II and also to calculate in parallel the computation
495 error that all the population of individuals also makes. These two approaches
496 that we followed provided a substantial increase in speed with regard to their
497 sequential versions. Logically, parallel versions follow a different development
498 strategy than the serial versions, since the design of an algorithm for being
499 deployed in a CPU is not valid to be run on a GPU directly, because mem-
500 ory and thread management in a GPU differ from the CPU ones, yet the
501 algorithm intention shares the same philosophy.

502 In figure 7 we can see the parallelization of the algorithm's general scheme.
503 We can observe that two sections are differentiated; the sequential section,
504 where the load of the dataset is carried out and the definition of parameters
505 of the algorithm. That section is also where the selection and ordering of the
506 best individuals that have been obtained after executing the whole process
507 (main loop of algorithm 1) is carried out to finally check the stop criterion.
508 The parallel section is made up of the evaluation of the individuals and all the
509 functions belonging to the evolutionary algorithms NSGA-II. In the following
510 sections we will design the entire algorithm for its optimization on the GPU.

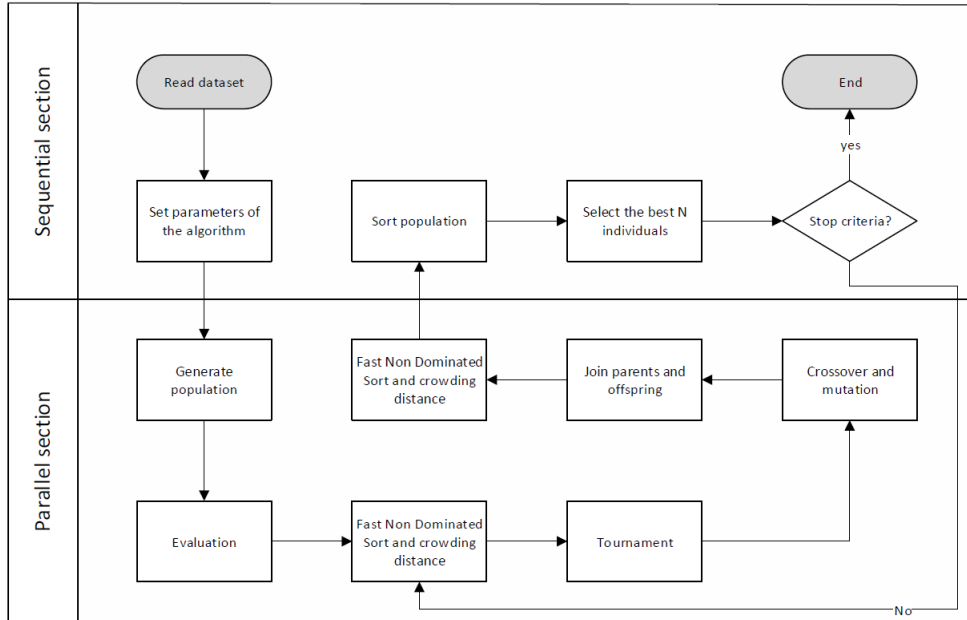


Figure 7: Flowchart of kernels.

511 *5.1. GPU-based Parallel Strategy for evaluations of ANNs in the population*

512 As we have mentioned in previous sections, the ANN used in our study
 513 and the one to be evaluated are the feed-forward neural networks. The
 514 evaluations of each of the networks are independent of each other, so we
 515 can consider launching them on the GPU at the same time, in parallel. In
 516 order to evaluate all individuals at once, we developed a kernel in which each
 517 individual is assigned to a specific block, as shown in figure 8. This block
 518 is called a thread-block in CUDA, which is a programming abstraction that
 519 represents, with CUDA Toolkit 10, a group of up to 1024 threads. Threads
 520 in the same block run on the same stream processor and can communicate
 521 with each other via shared memory or atomic operations.

522 Therefore, as in the kernel designed, we associated one individual to each
 523 block of the grid for parallel computation with GPU. The number of threads
 524 per block will be equal to the maximum number of neurons present in the
 525 initial parameters of the algorithm $\langle\langle Pop_size, N_max \rangle\rangle$. To map the
 526 grid's blocks into the individuals of NSGA-II, we used a global float matrix

527 in which each row corresponds to an individual. Each individual will be
 528 composed of the network weights, the number of hidden neurons and the
 529 error obtained after executing the evaluation.

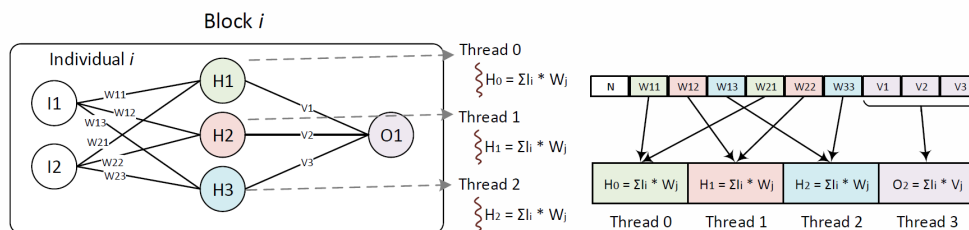


Figure 8: Parallel Strategy for Evaluating the ANN.

530 *5.2. Parallel Strategy for GPU Evaluation of NSGA-II Operators*

531 The strategy followed to parallelize the functions corresponding to the NSGA-
 532 II algorithm are detailed below. The GPU parallelization design is simpler
 533 than the one carried out in the ANN evaluation of the population's individ-
 534 uals of the previous subsection. Now, only a single thread-block of the GPU
 535 will be used to execute all individual actions. All kernels corresponding to
 536 this part are executed individually and serially. They have been executed
 537 with a block size equal to 1 whose number of threads is equal to the pop-
 538 ulation size $\ll 1, pop_size \gg$. The reason of this is because of the low
 539 cost in terms of space of those operations and the advantage of using a faster
 540 memory, i.e., in the same block all elements may access to a small space of
 541 memory with substantially lower latency than uncached global memory. As
 542 we can see in the figure 9, the data will be read from shared memory and
 543 global memory. In this way, each one of the operations associated to the
 544 functions of the algorithm is carried out in parallel by each population indi-
 545 vidual with the consequent advantage of being able to share memory of quick
 546 access, and thus obtaining a greater gain in the execution of each one of the
 547 functions of the evolutionary algorithm. These data will be processed by the
 548 different threads of the same block. Depending on the kernel to execute, the
 549 threads of the unique block will perform one task or another in order to get
 550 the following results, for example, for f_{GI} will apply the first function that
 551 will result in R_{GI} , which translates into an array with randomly initialized
 552 individuals:

- 553 • **Generation of individuals (GI):** A matrix will be obtained with a
554 new, randomly generated, population.
- 555 • **Non_dominated_sorting (NDS):** It will return the population sorted
556 by fronts.
- 557 • **Crowding distance (CD):** It will return, in a similar way to NDS,
558 the crowding distance of each individual with respect to its neighbors
559 that belong to its same front.
- 560 • **Selection operator (SO):** A set of the best parents will be obtained.
- 561 • **Crossover and mutation operators (CMO):** It will return a new
562 population of children from the parent selection carried out in SO.

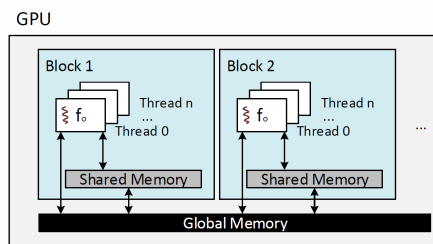


Figure 9: Memory management and design for each kernel f_i leveraging shared memory.

563 GI: A kernel has been created in which each thread randomly generates
564 an individual. The xoroshiro128+ [57] algorithm has been used to create
565 random numbers in the GPU. These new values will be stored in a float
566 matrix where each row will correspond to an individual, and the columns
567 will be the weights of the network that will later be evaluated in order to
568 obtain an error. As many threads as individuals in the population will be
569 used.

570 NDS: In this kernel, each thread is in charge of evaluating the dominance
571 of each individual with respect to the others. To this end, the measures of
572 goodness of the individuals have been loaded into shared memory in order to
573 accelerate access to these data. Once the dominance of each of the individuals
574 has been calculated, a single thread will be in charge of creation of the
575 different fronts. The data structure used to classify individuals in fronts

576 has been shown in figure 12, where the elements of a first array of integer
 577 type contains the indices of individuals sorted by fronts and the components
 578 of a second array of integer type holds the size of each front. To execute this
 579 kernel, as many threads as individuals in the population will be assigned to
 580 the thread block.

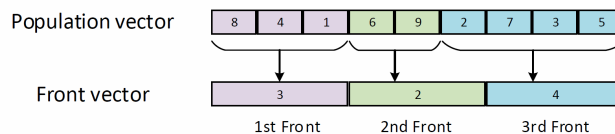


Figure 10: Data structure for fronts.

581 CD: Each of the threads is in charge of assigning the crowding distance
 582 between individuals of a specific front. To do this, each thread is responsible
 583 for obtaining the solutions with maximum and minimum fitness values, re-
 584 spectively, in each front, along with the crowding distance. As many threads
 585 as fronts will be used, so the kernel will be launched with a number of threads
 586 equal to the size of the population, which is meant to address the extreme
 587 case of each individual belongs to a single front. Two vectors have been
 588 created in shared memory of the integer and float types, respectively, which
 589 contain the measures of goodness. Finally, the distance between individuals
 590 of the same front is calculated and stored in a data structure similar to the
 591 one described in figure 12.

592 SO: In this kernel each of the threads is responsible for conducting a
 593 tournament between two individuals from which a winner is obtained. In
 594 order to optimize the speed of access to these data, shared memory area has
 595 been used to store the information of the fronts and the crowding distance.
 596 Two float and integer vectors have been created for this purpose.

597 CMO: Unlike the other kernels, this kernel's execution initiates with a
 598 number of threads per block equal to $pop_size/2$. This is because each thread
 599 is responsible for making a cross and a mutation through which two new
 600 individuals, also called children, will be generated. This kernel works with
 601 two float matrices that store parents and new children, both stored in global
 602 memory due to their large size. In addition, several variables stored in the
 603 registers of each thread are created to know the size of each individual, fitness
 604 and probabilities of crossing and mutation.

Data (Bytes)		
Global memory	Constant memory	Shared memory (per block)
8114×10^3	65535	4915
Other GeForce characteristics		
Registers per block	Warp size	Maximum Threads per MP
65536	32	2048

Table 1: Features of the NVIDIA Geforce GTX 1080.

605 6. Experiments and results

606 This section shows the experiments and results carried out to solve the prob-
607 lem of EC prediction in a set of distributed buildings of the University of
608 Granada. Developed models analyze historical data and are trained by fol-
609 lowing our parallel-optimized design of the NSGA-II evolutionary algorithm.
610 Four different facilities have been used in this study whose dataset has about
611 four years hourly data. Although the UGR has more buildings available, the
612 rest of them show a very similar energy expenditure to one of those selected.
613 In summary, each data set is made up of 35,000 instances each one thus
614 a building provides hourly information about its consumption, i.e., about
615 35,000 samples per building. Most of the pre-processing consisted in compil-
616 ing consumption by days; in this way, we will work with the past 24 hours
617 as time window for predicting.

618 All measurements were performed using a graphics card NVIDIA Geforce
619 GTX 1080, which has a total of 2560 cores, a maximum of 1024 threads per
620 block and the memory hierarchy that is depicted in Table 1. In this work,
621 Python was used due to its versatility and its broad compatibility with other
622 packages related to Machine Learning methods and algorithms. Numba is
623 utilized to work with CUDA, a Python compiler supported by CUDA that
624 can compile the code for its execution in GPU.

625 NSGA-II need to set different parameters before algorithm’s execution,
626 all of which were previously tested, and by modifying them one can observe
627 the algorithm yields more exploratory results or more convergence of results
628 depending on its execution behavior. In order not to hamper the final aim
629 of this contribution, some of these parameters were skipped as we want to
630 focus on the significant improvement in terms of execution time of our imple-

631 mentation of the algorithm. Once this has been done, we decided to choose
 632 as fundamental parameters for the NSGA-II algorithm the ones shown in
 633 Table 2 only.

Probability	
Crossing	0.9
Tournament's winner	0.9
Structural mutation	0.5
<i>Genetic mutation</i>	0.1

Table 2: Chosen parameters for NSGA-II implementation

634 Although several population sizes were tested for the NSGA-II algorithm
 635 implementation, its size was set to 160 individuals owing to the results that
 636 will be explained in the following paragraphs. The ANN instances deployed
 637 in our study were set as to have a range of neurons between 4 and 32. The
 638 *energy consumption* (EC) data have been normalized between 0 and 1 to
 639 have the same range of values for each input to the ANN, thus ensuring that
 640 the model will not give more weight to those attributes which values belong
 to a wider range. The data will be reconstructed as shown in Figure 11.

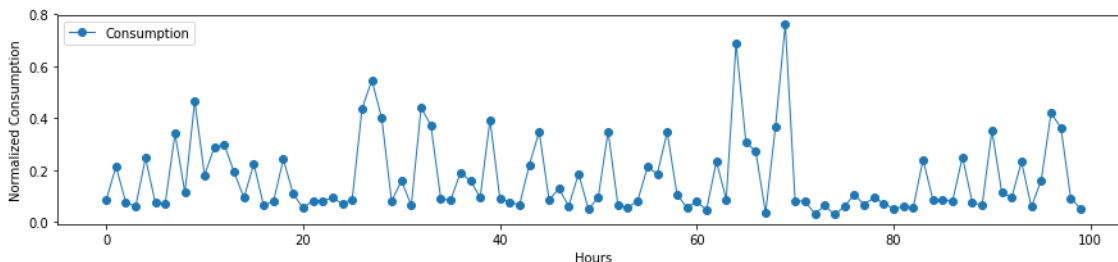


Figure 11: Normalized consumption over a 100-hour period.

641
 642 The NSGA-II developed in this work allows us to obtain a set of optimal
 643 solutions, which are those that form the Pareto front. On figure XX is shown
 644 a visual example of the Pareto front obtained for a particular building of our
 645 study. In this graph, the X axis shows the error committed (MSE) and the Y
 646 axis the number of hidden neurons, each one of the graph-points represents
 647 a non-dominated optimal solution and the entire set of points represent the
 648 Pareto front

649 On this front of Pareto we will be see the solutions that have obtained the
650 least error in the prediction of energy consumption with a minimum network
size

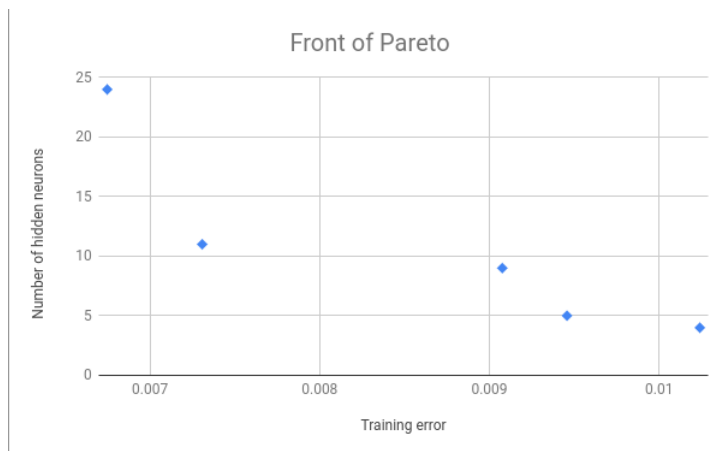


Figure 12: The Pareto optimal front provided by the NSGA-II

651

652 We run ten executions for each algorithm and building, so that we could
653 perform a statistical analysis of the results. Each dataset was randomly di-
654 vided into training data (70%) and test data (30%), to prevent over-training.

655 After the parameters are established, the first experiment is intended to
656 compare both implementations of the NSGA-II, the sequential and our paral-
657 lel proposal. Notice that both implementations are exactly the same. The
658 sequential design is performed by launching one block and just one thread
659 in the GPU. Table 3 illustrates the speed up achieved by our design. In this
660 table, all kernels decrease their execution time in the parallel version of the
661 evolutionary algorithm with respect to the sequential one. Some functions
662 need more computational cost than others, for example, the generation of
663 individuals (GI) and the cross-mutation operator (CMO) take similar time
664 because their computations are focused on simple operations over the popu-
665 lation. In the first case, GI is responsible for the creation of random numbers
666 to generate all individuals. In the second case, CMO not only creates some
667 random values but also modifies some genes of individuals to perform the
668 mutation and crossover computations. For this reason, the implementation
669 of CMO operator shows more speedup than that of the GI. Another example
670 of an effortless function is the selection operator (SO) as it just needs to
671 decide which individuals are selected for crossing. Differently to execution

672 times shown by those functions, the non-dominated sort operator (NDS) is
 673 the most time-consuming function. In this case, more barriers that synchronise
 674 and cause contention of threads arise in the NDS function because of
 675 the high dependency among fronts to build each one of them, and this is why
 676 the lowest speed-up is obtained here. Along with the NDS, the crowding distance
 677 (CD) is the basis of the NSGA-II algorithm. The CD, however, may
 678 be optimized in a better way than NDS because CD calculates the *distance*
 679 matrix among solutions, and this task may be done separately. This procedure
 680 takes considerable time as all individuals must be computed each other,
 681 for this reason, the third-best speed-up is attained here. Finally, the most
 682 consuming function is often associated with the fitness function (FF), due
 683 to its iterative behaviour. This function used to be launched over and over
 684 during the algorithm. As a consequence, the emphasis has been placed on
 685 looking for a good parallel implementation of FF in our proposal. Thus, in
 686 this table one may observe a great difference in execution times of all functions
 687 in its sequential version. The same happens if we observe the parallel
 688 version results. Nevertheless, the optimized design of the FF function empowers
 689 the algorithm and gets a surprising speed-up up to 787.90. What's
 690 more, it enables to decrease the execution time of the algorithm from several
 691 hours to some minutes.

Kernel Name	Sequential	Parallel	Speedup
GI	17.85	0.259	68.92
NDS	479.67	62.6	7.66
CD	1026.50	17.5	58.66
SO	3.70	0.122	30.33
CMO	17.56	1.02	17.22
FF	104350	132.44	787.90

Table 3: Time (ms) of each kernel, generation of individuals (GI), non-dominated sort and crowding distance functions (NDS and CD, respectively), selection operator (SO), crossover and mutation operators (CMO) and the fitness function (FF).

692 Some experiments were launched using different population size. Those
 693 results are gathered in table 4. The execution time in nearly all cases is the
 694 same, although it is observed that it increases a little as the population grows.
 695 There is almost no difference because there is still some memory in the GPU
 696 which may be used. For this reason, from 200 individuals on it takes longer.

697 In addition to this, the mean squared error achieved is slightly better as the
698 population increases. However, the implementation run typically reaches a
699 point between 160 and 200 individuals beyond which the algorithm does not
700 provide any further improvement. This is on account of the amount of the
701 information that individuals may provide, and therefore the amount of useful
702 knowledge each individual can share.

Population Size	Time (s)	MSE	Memory Usage (MB)	GPU Usage
80	22.989	0.0122	143	92%
160	19.792	0.0098	145	95%
200	20.129	0.0094	147	98%
250	24.345	0.0099	151	100%
300	25.932	0.0105	160	100%

Table 4: Scalability of the designed algorithm using different population size.

703 Similarly, table 5 shows the average execution time of each individual
704 with a different number of individuals if all of them got the same number
705 of threads. The first column shows the number of individuals, the second
706 column the minimum time, the third column the maximum time and the last
707 column the time average. That table provides information about how each
708 individual moderately increases their time cost as the number of individuals
709 increases. This is due to the increase in data that are loaded into the device as
710 the number of individuals increases. The amount of data is relatively small, as
711 a consequence, all data are located in shared memory except the information
712 related to each individual, and this fact explains this little variability.

Individuals	Min	Max	Avg
80	0.8249	1.1721	0.9286
160	0.8443	1.1491	0.9492
200	0.8713	1.0945	0.9620
250	0.8717	1.0833	0.9684
300	0.8802	1.0643	0.9820

Table 5: Population size - Time (sec) execution of an individual.

713 The next experiment is intended to compare the accuracy of both solu-
714 tions. In table 6 we can see the experimentation for sequential and parallel
715 models. Columns 2,3,4 and 5 show the experimentation for the buildings
716 B1, B2, B3 and B4, respectively. Column 1 describes the metric used. For
717 sequential and parallel experiments we show the average fitness, which con-
718 tains the average MSE of 10 executions of each algorithm; the best fitness
719 and worst fitness with the minimum and maximum MSE obtained with 10
720 runs and the standard deviation. The number of evaluations has been fixed
721 at 5000 and a total of 150 individuals. As we can see in the table 6, the re-
722 sults obtained by both implementations are quite alike. However, the parallel
723 version shows lower error values on average.

724 For a better analysis of the results in table 6, we have included boxplots
725 of the MSE distribution. Figure 13 contains the boxplots of the MSE for the
726 different buildings in sequential and parallel models. The median value in
727 boxplots is highlighted with an orange line. The whiskers plot illustrates in
728 all buildings present a mean error closer to the Q1 in sequential, and their
729 parallel versions get closer to the Q3. Just in the fourth instance, these results
730 slightly change due to the higher complexity of the behaviour presented by
731 that building. This fact give us an idea that all results are quite similar among
732 them, and their differences are caused by the randomness in the algorithm.
733 In nearly all cases, the algorithm reaches a set of solutions comparable in
734 every experiment. However, in the situation that the building presents a
735 consumption harder to predict, then their solutions vary in a higher range as
736 it is not so easy to find the optimal model, yet note that this range is quite
737 small in terms of error as illustrated in the following figures.

738 The implementation carried out in parallel maintains the quality, under-
739 stood as the lower possible MSE of solutions obtained with the algorithm,
740 with respect to the solutions found in sequential implementation as shown in
741 table 6. With the intention of validating the quality of the results obtained,
742 a statistical test has been used. First, the normality of the errors was verified
743 with the Shapiro test. Since all error distributions follow a normal distribu-
744 tion, a parametric test was carried out using the T-test with 99% confidence
745 to compare the results in the sequential and parallel version for each of the
746 buildings. Each cell contains the p-value resulting from the t-test. A value
747 < 0.01 means that there are significant differences between the samples com-
748 pared in the test, i.e., the sequential and parallel versions of the algorithm
749 in our case, and a value > 0.01 means that there are no differences between
750 the two versions. We may conclude from this table that both versions show

	Measures	B1	B2	B3	B4
Sequential	Average Fitness	0.0168	0.0118	0.0182	0.0174
	Best Fitness	0.0140	0.0079	0.0130	0.0165
	Worst Fitness	0.0197	0.0145	0.0237	0.0245
	Standard Deviation	± 0.0018	± 0.0022	± 0.0031	± 0.0035
Parallel	Average	0.0153	0.0095	0.0152	0.0162
	Best Fitness	0.0117	0.0061	0.0107	0.0130
	Worst Fitness	0.0176	0.0122	0.0191	0.0218
	SD	± 0.0017	± 0.0018	± 0.0025	± 0.0027
t-Test		0.0821	0.0227	0.0305	0.0329

Table 6: Average MSE of sequential and parallel implementations.

751 similar behaviour in terms of accuracy. There are no significant differences
752 except in the fourth case, caused by the variability in the solutions obtained
753 experiments that correspond to the sequential implementation. Note that
754 Figure 13 shows a greater box in that building. This is because the be-
755 haviour of that building is harder to model. Consequently, finding a model
756 with good accuracy is more complicated.

757 Table 7 shows the percentage of GPU utilization of each of the kernels.
758 Starting from the bottom upwards, the GI is the function that makes less
759 use of the device, followed by the SO and the CMO. This is due to their
760 computations are notably simpler than the rest. After them, CD is placed
761 in the third place, this kernel must be launched every time all individuals
762 are evaluated. It takes longer than the previous functions because it must
763 calculate distances among all individuals. The function that needs the sec-
764 ond more time to get their results is the NDS as this kernel has to perform a
765 sorting algorithm according to the dominance of the solutions and this kind
766 of algorithms often takes a great amount of time. Finally, as expected, the
767 core of the algorithm is mainly focused on the FF. This kernel takes the
768 longest time on account of the data that need to process and the number
769 of iterations that need to be done to evaluate all individuals with all data.
770 Nearly half of the time required to execute the algorithm is spent in FF.
771 For this reason, our efforts were focused on improving and optimizing this
772 part of the evolutionary algorithm by associating one individual per block.
773 Thus, the number of threads was optimized according to the maximum num-

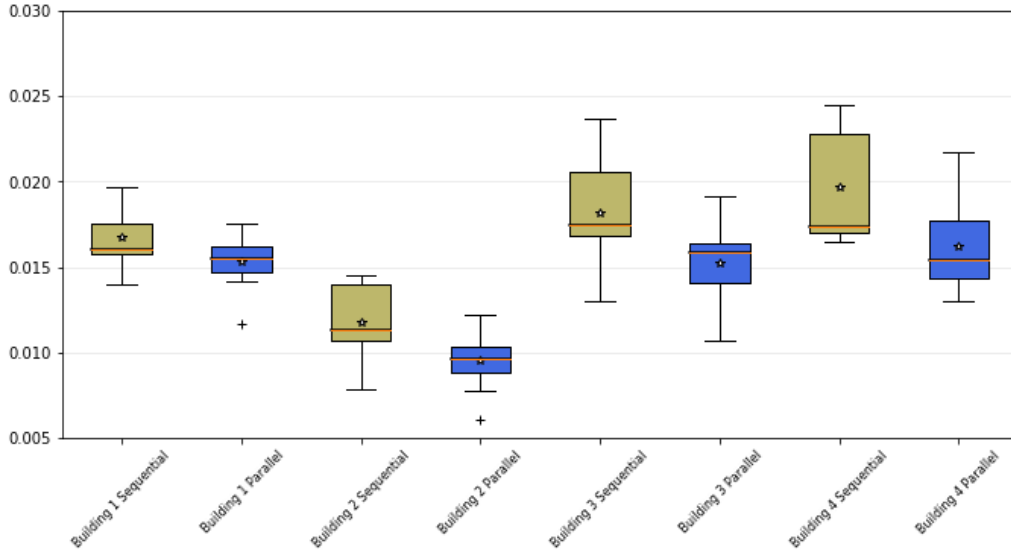


Figure 13: Boxplots of the error rate (MSE) for buildings.

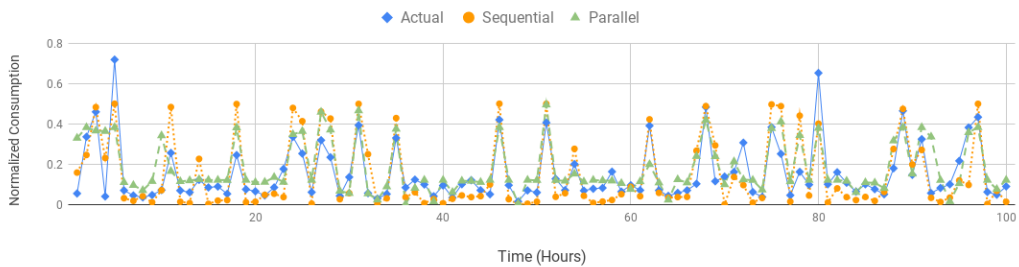
774 ber of neurons and therefore, all computations can be executed in parallel,
 775 achieving the excellent speed-up obtained. This fact demonstrates that our
 776 implementation might be considered a significant contribution in the current
 777 state-of-art since not only improves the parallel design of a hybrid solution of
 778 the NSGA-II with ANN but is also showing better figures regarding speed up
 779 in particularly complex functions of the algorithm. Additionally, it lessens
 780 the MSE with respect to other NSGA-II implementations on GPU.

Kernel Name	Usage Rate
FF	48.11%
NDS	39.55%
CD	11.46%
CMO	0.23%
SO	0.07%
GI	0.01%

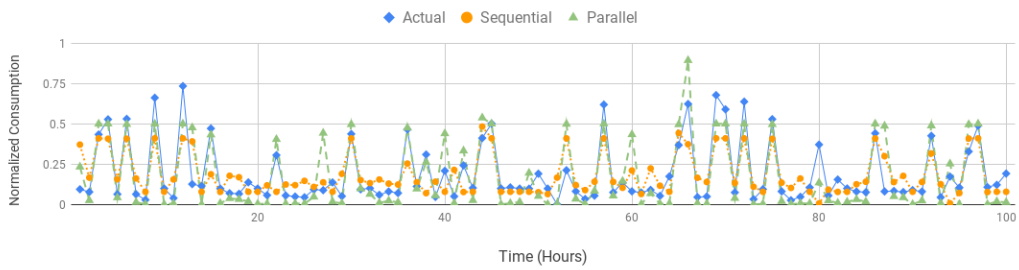
Table 7: Kernel profiling.

781 Finally, figure 14 shows a graphical representation of the forecast of EC
 782 for a total of 100 hours of a set of test data. In these plots, the horizontal

783 axis represents the hour that is predicted and the vertical axis represents
784 the normalized consumption for that hour. Figure 14a shows the best result
785 obtained and figure 14b shows the worst result of the parallel and sequential
786 versions. As we can see in the figures the predicted results fit well to the real
787 values in both cases. The results of both versions have been put together to
788 visually verify how the fit does not differ considerably from each other.



(a)



(b)

Figure 14: Example of the (a) best and worst (b) predictions of energy consumption by both parallel and sequential models.

789 *6.1. Performance comparison between NSGA-II and Backpropagation.*

790 After showing the good results obtained in the performance of our parallel
791 NSGA-II algorithm in comparison to its sequential version, an experiment
792 has been carried out to compare the algorithm developed in this paper with
793 the well-known Backpropagation (BP) one. BP is one of the most popular
794 local search algorithms used for neural network training.

795 Given the long execution time that the BP algorithm implementation
796 takes, a total of 9 runs of the algorithm were launched for each of the buildings
797 included in this study. Each of the executions has been performed to a total
798 of 50 iterations. The number of neurons in the Nsga-II experiment has been
799 set between a range of 4 and 32, however the range chosen for BP algorithm
800 execution was coarser (by increasing in two from 2 to 32), since the average
801 execution time of the BP took too long. The average time taken to run the
802 parallel NSGA-II and then obtain a set of optimal solutions is 25.932 s., while
803 the BP had an average execution time of 10 hours to obtain the results with
804 respect to all the different number of neurons aforementioned. In table 8 we
805 can see the performance comparison for both algorithms. This table shows,
806 for each building, the best model obtained after conducting the experiment.
807 The first column shows the identifier of each building. In column 2 and 5 we
808 can see the number of hidden neurons, which the best model obtains for the
809 different buildings with the NSGA-II and the BP respectively. In columns 3
810 and 6 it can be seen that the error given by the NSGA-II and BP for the B2
811 building is almost equal, and it is slightly better in the rest of the buildings
812 as for the BP results. Finally, columns 4 and 7 show the execution time of
813 the NSGA-II and BP, it can be seen that our algorithm is much faster than
814 the BP. Although the error of the models is important, in this work we have
815 focused more in improving the execution time sacrificing some tenths of error
816 since, since in some problems is most relevant to obtain faster and a little
817 less precise models.

818 **7. Conclusion**

819 This paper presents a parallel implementation of NSGA-II developed in GPU
820 to train an ANN whose evaluation has also been implemented in parallel in
821 GPU to solve the problem of predicting EC. The objective of our proposal is
822 to reduce the execution time of the different functions that compose the algo-
823 rithm. Our implementation is able to obtain an optimal neural network with
824 the least number of neurons necessary to learn the set of examples and also

Building	NSGA-II neurons	NSGA-II MSE	NSGA-II time (seconds)	BP neurons	BP MSE	BP time (seconds)
B1	4	0.0117	26.998	14	0.0044	7821.44
B2	4	0.0061	25.965	4	0.0049	595.40
B3	4	0.0107	25.441	18	0.0054	12200
B4	4	0.0130	25.899	18	0.0045	13089.30

Table 8: NSGA-II and BP performance comparison.

825 provides the weights trained making a minimum error for a specific building,
826 and thus taking less than 60 seconds with a total of 20000 evaluations and
827 an average MSE of 0.0093, achieving in some functions of the algorithm a
828 speedup of 788 compared with the sequential version.

829 In the experimentation section, different results of the behaviour of the
830 algorithms developed in this work have been shown. Our proposal obtains
831 good results in time and accuracy with respect to the sequential version on
832 account of the good utilization of the resources of the GPU as they are the
833 deployment of memory of fast access. It is worth mentioning that 100%
834 of GPU use has been achieved, for which it has been necessary an elaborate
835 CUDA development work and extensive experimentation. On the other hand,
836 the use of the GPU limits the parameters of the algorithm since it cuts down
837 the size of the population of individuals that can be used due to constraints
838 by GPU resources availability. Finally, the version developed in parallel has
839 shown to obtain better results in execution time without loss of precision.
840 This assertion has been verified with statistical tests, which have shown that
841 there are no significant differences between the accuracy obtained in the serial
842 version and the parallel version, even though there is a great difference in
843 the performance of an advantageous way by the parallel version.

844 In future work we would like to explore and solve our current device lim-
845 its by deploying our implementation on a GPU with more computational
846 resources, as the Nvidia Tesla with the intention of knowing if a higher pop-
847 ulation size and a higher number of executions will produce better results
848 without increasing too much the execution time of the entire algorithm. We
849 are also confident that CUDA’s future capabilities will be improved by help-
850 ing us to scale our problem. Besides, as a continuation of this work, several
851 lines of research remain open and in which it is possible to continue working.
852 These lines have arisen during the development of this work and we hope to

853 be able to work in the near future. One of these future lines is to develop
854 more complex networks for the evaluation of the models, intending to achieve
855 an improvement in the precision of the results. On the other side, we would
856 like to analyze the aggregation of other data sources to reinforce the results
857 obtained. These new data sources could be measurements of temperature,
858 humidity, luminosity and occupancy.

859 **Acknowledgement**

860 This work has been supported by the project TIN201564776-C3-1-R.

861 **References**

- 862 [1] S. Belt Ibérica, En los últimos 50 años, el plan-
863 eta se ha degradado más que en 100 siglos,
864 <http://www.belt.es/noticias/2005/abril/26/planeta.htm>, 2005.
- 865 [2] Y. Feng, D. Gong, Q. Zhang, S. Jiang, L. Zhao, N. Cui, Evaluation of
866 temperature-based machine learning and empirical models for predicting
867 daily global solar radiation, *Energy Conversion and Management* 198
868 (2019) 111780.
- 869 [3] Y. Liu, Y. Zhou, D. Wang, Y. Wang, Y. Li, Y. Zhu, Classification of
870 solar radiation zones and general models for estimating the daily global
871 solar radiation on horizontal surfaces in china, *Energy Conversion and*
872 *Management* 154 (2017) 168 – 179.
- 873 [4] J.-K. Park, A. Das, J.-H. Park, A new approach to estimate the spatial
874 distribution of solar radiation using topographic factor and sunshine
875 duration in south korea, *Energy Conversion and Management* 101 (2015)
876 30 – 39.
- 877 [5] L. G. B. Ruiz, M. I. Capel, M. C. Pegalajar, Parallel memetic algorithm
878 for training recurrent neural networks for the energy efficiency problem,
879 *Applied Soft Computing* 76 (2019) 356–368.
- 880 [6] L. G. B. Ruiz, R. Rueda, M. P. Cuéllar, M. C. Pegalajar, Energy
881 consumption forecasting based on elman neural networks with evolutive
882 optimization, *Expert Systems with Applications* 92 (2018) 380–389.

- 883 [7] L. B. d. Oliveira, C. G. Marcelino, A. Milanés, P. E. M. Almeida, L. M.
884 Carvalho, A successful parallel implementation of NSGA-II on GPU for
885 the energy dispatch problem on hydroelectric power plants, in: 2016
886 IEEE Congress on Evolutionary Computation (CEC), pp. 4305–4312.
- 887 [8] J. Bergstra, O. Breuleux, F. Bastien, P. Lamblin, R. Pascanu, G. Des-
888 jardins, J. P. Turian, D. Warde-Farley, Y. Bengio, Theano : A CPU
889 and GPU Math Compiler in Python.
- 890 [9] V. Roberge, M. Tarbouchi, G. Labonté, Fast Genetic Algorithm Path
891 Planner for Fixed-Wing Military UAV Using GPU, IEEE Transactions
892 on Aerospace and Electronic Systems 54 (2018) 2105–2117.
- 893 [10] R. S. Sinha, S. Singh, S. Singh, V. K. Banga, Accelerating Genetic
894 Algorithm Using General Purpose GPU and CUDA.
- 895 [11] D. D’Agostino, G. Pasquale, I. Merelli, A Fine-Grained CUDA Im-
896 plementation of the Multi-objective Evolutionary Approach NSGA-II:
897 Potential Impact for Computational and Systems Biology Applications,
898 in: C. DI Serio, P. Liò, A. Nonis, R. Tagliaferri (Eds.), Computational
899 Intelligence Methods for Bioinformatics and Biostatistics, Lecture Notes
900 in Computer Science, Springer International Publishing, 2015, pp. 273–
901 284.
- 902 [12] I. M. Coelho, V. N. Coelho, E. J. d. S. Luz, L. S. Ochi, F. G. Guimarães,
903 E. Rios, A GPU deep learning metaheuristic based model for time series
904 forecasting, Applied Energy 201 (2017) 412–418.
- 905 [13] O. Guerra Santin, Behavioural Patterns and User Profiles related to
906 energy consumption for heating, Energy and Buildings 43 (2011) 2662–
907 2672.
- 908 [14] D. B. Araya, K. Grolinger, H. F. ElYamany, M. A. M. Capretz, G. Bit-
909 suamlak, An ensemble learning framework for anomaly detection in
910 building energy consumption, Energy and Buildings 144 (2017) 191–
911 206.
- 912 [15] Y. Zhang, W. Chen, J. Black, Anomaly detection in premise energy
913 consumption data, in: 2011 IEEE Power and Energy Society General
914 Meeting, pp. 1–8.

- 915 [16] B. B. Ekici, U. T. Aksoy, Prediction of building energy consumption by
916 using artificial neural networks, *Advances in Engineering Software* 40
917 (2009) 356–362.
- 918 [17] S. Heiple, D. J. Sailor, Using building energy simulation and geospatial
919 modeling techniques to determine high resolution building sector energy
920 consumption profiles, *Energy and Buildings* 40 (2008) 1426–1436.
- 921 [18] S. Barak, S. S. Sadegh, Forecasting energy consumption using ensemble
922 ARIMA–ANFIS hybrid algorithm, *International Journal of Electrical
923 Power & Energy Systems* 82 (2016) 92–104.
- 924 [19] Y.-S. Lee, L.-I. Tong, Forecasting energy consumption using a grey
925 model improved by incorporating genetic programming, *Energy Con-
926 version and Management* 52 (2011) 147–152.
- 927 [20] J. Fan, X. Wang, L. Wu, F. Zhang, H. Bai, X. Lu, Y. Xiang, New
928 combined models for estimating daily global solar radiation based on
929 sunshine duration in humid regions: A case study in south china, *Energy
930 Conversion and Management* 156 (2018) 618 – 625.
- 931 [21] B. Dong, C. Cao, S. E. Lee, Applying support vector machines to predict
932 building energy consumption in tropical region, *Energy and Buildings*
933 37 (2005) 545–553.
- 934 [22] L. Ekonomou, Greek long-term energy consumption prediction using
935 artificial neural networks, *Energy* 35 (2010) 512–517.
- 936 [23] K. Li, H. Su, J. Chu, Forecasting building energy consumption using
937 neural networks and hybrid neuro-fuzzy system: A comparative study,
938 *Energy and Buildings* 43 (2011) 2893–2899.
- 939 [24] B. Choubin, G. Zehtabian, A. Azareh, E. Rafiei-Sardooi, F. Sajedi-
940 Hosseini, Kişi, Precipitation forecasting using classification and regres-
941 sion trees (CART) model: a comparative study of different approaches,
942 *Environmental Earth Sciences* 77 (2018) 314.
- 943 [25] A. Azadeh, S. F. Ghaderi, S. Tarverdian, M. Saberi, Integration of
944 artificial neural networks and genetic algorithm to predict electrical en-
945 ergy consumption, *Applied Mathematics and Computation* 186 (2007)
946 1731–1741.

- 947 [26] L. Wu, G. Huang, J. Fan, F. Zhang, X. Wang, W. Zeng, Potential
948 of kernel-based nonlinear extension of arps decline model and gradient
949 boosting with categorical features support for predicting daily global
950 solar radiation in humid regions, *Energy Conversion and Management*
951 183 (2019) 280 – 295.
- 952 [27] J. Fan, X. Wang, L. Wu, H. Zhou, F. Zhang, X. Yu, X. Lu, Y. Xiang,
953 Comparison of support vector machine and extreme gradient boosting
954 for predicting daily global solar radiation using temperature and precip-
955 itation in humid subtropical climates: A case study in china, *Energy*
956 *Conversion and Management* 164 (2018) 102 – 111.
- 957 [28] R. A. Conde-Gutiérrez, U. Cruz-Jacobo, A. Huicochea, S. R. Casolco,
958 J. A. Hernández, Optimal multivariable conditions in the operation of an
959 absorption heat transformer with energy recycling solved by the genetic
960 algorithm in artificial neural network inverse, *Applied Soft Computing*
961 72 (2018) 218–234.
- 962 [29] Y. Dong, J. Wang, Z. Guo, Research and application of local perceptron
963 neural network in highway rectifier for time series forecasting, *Applied*
964 *Soft Computing* 64 (2018) 656–673.
- 965 [30] G. S. Georgiou, P. Christodoulides, S. A. Kalogirou, Implementing ar-
966 tificial neural networks in energy building applications — A review, in:
967 2018 IEEE International Energy Conference (ENERGYCON), pp. 1–6.
- 968 [31] J. Krzywanski, K. Grabowska, F. Herman, P. Pyrka, M. Sosnowski,
969 T. Prauzner, W. Nowak, Optimization of a three-bed adsorption chiller
970 by genetic algorithms and neural networks, *Energy Conversion and*
971 *Management* 153 (2017) 313 – 322.
- 972 [32] M. Heitzler, J. C. Lam, J. Hackl, B. T. Adey, L. Hurni, GPU-Accelerated
973 Rendering Methods to Visually Analyze Large-Scale Disaster Simulation
974 Data, *Journal of Geovisualization and Spatial Analysis* 1 (2017) 3.
- 975 [33] S. S. Stone, J. P. Haldar, S. C. Tsao, B. Sutton, Z.-P. Liang, et al.,
976 Accelerating advanced mri reconstructions on gpus, *Journal of parallel*
977 *and distributed computing* 68 (2008) 1307–1318.
- 978 [34] A. Syberfeldt, T. Ekblom, A comparative evaluation of the GPU vs.
979 the CPU for parallelization of evolutionary algorithms through multiple

- 980 independent runs, *International Journal of Computer Science & Infor-*
981 *mation Technology (IJCSIT)* 9 (2017) 1–14.
- 982 [35] OpenCL - The open standard for parallel programming of heterogeneous
983 systems, <https://www.khronos.org/ocl/>, 2013.
- 984 [36] J. Sanders, E. Kandrot, *CUDA by Example: An Introduction to*
985 *General-Purpose GPU Programming*, Addison-Wesley Professional, 1st
986 edition, 2010.
- 987 [37] OpenCL vs. CUDA: Which Has Better Application Support,
988 <https://create.pro/blog/opengl-vs-cuda/>, 2017.
- 989 [38] C. H. Dagli, *Artificial Neural Networks for Intelligent Manufactur-*
990 *ing*, Springer Science & Business Media, 2012. Google-Books-ID:
991 K4ftCAAQBAJ.
- 992 [39] S. Samarasinghe, *Neural Networks for Applied Sciences and Engineer-*
993 *ing : From Fundamentals to Complex Pattern Recognition*, Auerbach
994 Publications, 2016.
- 995 [40] D. Roffman, G. Hart, M. Girardi, C. J. Ko, J. Deng, Predicting non-
996 melanoma skin cancer via a multi-parameterized artificial neural net-
997 work, *Scientific Reports* 8 (2018) 1701.
- 998 [41] J. Škutová, *Weights initialization methods for mlp neural networks*,
999 *Vysoká škola báňská-Technická Univerzita Ostrava* (2008).
- 1000 [42] L. G. B. Ruiz, M. Cuéllar, M. Delgado, M. C. Pegalajar, An application
1001 of non-linear autoregressive neural networks to predict energy consump-
1002 tion in public buildings, *Energies* 9 (2016) 684.
- 1003 [43] M. Delgado, M. C. Pegalajar, A multiobjective genetic algorithm for
1004 obtaining the optimal size of a recurrent neural network for grammatical
1005 inference, *Pattern Recognition* 38 (2005) 1444–1456.
- 1006 [44] A. Blanco, M. Delgado, M. C. Pegalajar, A real-coded genetic algorithm
1007 for training recurrent neural networks, *Neural Networks* 14 (2001) 93–
1008 105.

- 1009 [45] P. Pospichal, J. Jaros, J. Schwarz, Parallel Genetic Algorithm on the
1010 CUDA Architecture, in: C. Di Chio, S. Cagnoni, C. Cotta, M. Ebner,
1011 A. Ekárt, A. I. Esparcia-Alcazar, C.-K. Goh, J. J. Merelo, F. Neri,
1012 M. Preuß, J. Togelius, G. N. Yannakakis (Eds.), Applications of Evo-
1013 lutionary Computation, Lecture Notes in Computer Science, Springer
1014 Berlin Heidelberg, 2010, pp. 442–451.
- 1015 [46] E. Zitzler, K. Deb, L. Thiele, Comparison of Multiobjective Evolution-
1016 ary Algorithms: Empirical Results, *Evol. Comput.* 8 (2000) 173–195.
- 1017 [47] C. Von Lüken, A. Hermosilla, B. Barán, Algoritmos evolutivos para
1018 optimización multiobjetivo: un estudio comparativo en un ambiente par-
1019 alelo asíncrono, X Congreso Argentino de Ciencias de la Computación
1020 (2004).
- 1021 [48] C. A. Coello Coello Coello, A Short Tutorial on Evolutionary Mul-
1022 tiobjective Optimization, in: E. Zitzler, L. Thiele, K. Deb, C. A.
1023 Coello Coello, D. Corne (Eds.), Evolutionary Multi-Criterion Optimiza-
1024 tion, Lecture Notes in Computer Science, Springer Berlin Heidelberg,
1025 2001, pp. 21–40.
- 1026 [49] G. Wu, J. L. Greathouse, A. Lyashevsky, N. Jayasena, D. Chiou,
1027 GPGPU performance and power estimation using machine learning, in:
1028 2015 IEEE 21st International Symposium on High Performance Com-
1029 puter Architecture (HPCA), pp. 564–576.
- 1030 [50] M. C. Pegalajar, M. Capel, L. G. B. Ruiz, M. Delgado, A parallel ap-
1031 proach to intelligent data analysis for efficient energy management in
1032 distributed facilities (pioneer), in: European Space Projects: Develop-
1033 ments, Implementations and Impacts in a Changing World - Volume 1:
1034 EPS Porto 2017,, INSTICC, SciTePress, 2017, pp. 28–49.
- 1035 [51] K. Deb, A. Pratap, S. Agarwal, T. Meyarivan, A Fast and Elitist Mul-
1036 tiobjective Genetic Algorithm: NSGA-II, *Trans. Evol. Comp* 6 (2002)
1037 182–197.
- 1038 [52] M. Delgado, M. P. Cuellar, M. C. Pegalajar, Multiobjective hybrid opti-
1039 mization and training of recurrent neural networks, *IEEE Transactions*
1040 *on Systems, Man, and Cybernetics, Part B (Cybernetics)* 38 (2008) 381–
1041 403.

- 1042 [53] A. H. Wright, Genetic Algorithms for Real Parameter Optimization,
1043 in: Foundations of Genetic Algorithms, Morgan Kaufmann, 1991, pp.
1044 205–218.
- 1045 [54] M. Delgado, M. P. Cuellar, M. C. Pegalajar, Multiobjective Hybrid
1046 Optimization and Training of Recurrent Neural Networks, IEEE Trans-
1047 actions on Systems, Man, and Cybernetics, Part B (Cybernetics) 38
1048 (2008) 381–403.
- 1049 [55] C. A. C. Flórez, R. A. Bolaños, A. M. Cabrera, Algoritmo Multiobjetivo
1050 Nsga-Ii Aplicado Al Problema De La Mochila., Scientia Et Technica XIV
1051 (2008) 206–211.
- 1052 [56] K. Deb, D. Kalyanmoy, Multi-Objective Optimization Using Evolution-
1053 ary Algorithms, John Wiley & Sons, Inc., New York, NY, USA, 2001.
- 1054 [57] Xoshiro, xoroshiro generators and the PRNG shootout,
1055 <http://xoshiro.di.unimi.it/>, 2018.