

# Lights and Shadows in Evolutionary Deep Learning: Taxonomy, Critical Methodological Analysis, Cases of Study, Learned Lessons, Recommendations and Challenges

Aritz D. Martinez<sup>a</sup>, Javier Del Ser<sup>a,b,\*</sup>, Esther Villar-Rodriguez<sup>a</sup>, Eneko Osaba<sup>a</sup>, Javier Poyatos<sup>c</sup>,  
Siham Tabik<sup>c</sup>, Daniel Molina<sup>c</sup>, Francisco Herrera<sup>c</sup>

<sup>a</sup>TECNALIA, Basque Research & Technology Alliance (BRTA), 48160 Derio, Spain

<sup>b</sup>University of the Basque Country (UPV/EHU), 48013 Bilbao, Spain

<sup>c</sup>DaSCI Andalusian Institute of Data Science and Computational Intelligence, University of Granada, 18071 Granada, Spain

---

## Abstract

Much has been said about the fusion of bio-inspired optimization algorithms and Deep Learning models for several purposes: from the discovery of network topologies and hyperparametric configurations with improved performance for a given task, to the optimization of the models parameters as a replacement for gradient-based solvers. Indeed, the literature is rich in proposals showcasing the application of assorted nature-inspired approaches for these tasks. In this work we comprehensively review and critically examine contributions made so far based on three axes, each addressing a fundamental question in this research avenue: a) optimization and taxonomy (*Why?*), including a historical perspective, definitions of optimization problems in Deep Learning, and a taxonomy associated with an in-depth analysis of the literature, b) critical methodological analysis (*How?*), which together with two case studies, allows us to address learned lessons and recommendations for good practices following the analysis of the literature, and c) challenges and new directions of research (*What can be done, and what for?*). In summary, three axes – optimization and taxonomy, critical analysis, and challenges – which outline a complete vision of a merger of two technologies drawing up an exciting future for this area of fusion research.

*Keywords:* Deep Learning, Neuroevolution, Evolutionary Computation, Swarm Intelligence.

---

## 1. Introduction

Nowadays there is overall consensus on the capital importance gained by Deep Learning in the Artificial Intelligence field [1]. Initial results of Deep Learning date back to the late 80's, stepping on a history of preceding achievements in neural computation [2]. However, it was not until years later when advances in high-performance computing, new achievements in neural network training [3], and the availability of massive datasets paved the way for the renowned success of this family of learning models. Nowadays, plenty of application areas have harnessed the superior modeling capabilities of Deep Learning models, from natural language processing [4], speech and audio processing [5, 6], social network analysis [7] or autonomous driving [8], to mention a few. As a result, Deep Learning models such as Convolutional Neural Networks (CNN, [9]), Recurrent Neural Networks (RNN, [10]) or Generative Adversarial Networks (GAN, [11]) prevail in many tasks, including image classification, time series forecasting or visual object generation.

---

\*Corresponding author. TECNALIA, Basque Research & Technology Alliance (BRTA), P. Tecnológico, Ed. 700. 48170 Derio (Bizkaia), Spain. E-mail: javier.delser@tecnalia.com (Javier Del Ser).

There are several properties of Deep Learning models that make them outperform traditional *shallow learning* methods. Among them, Deep Learning models can automatically learn hierarchical features from raw data, so that features organized in higher levels of the hierarchy are composed by a combination of simpler lower-level features. As a result of this capability, features with minimal human effort and domain knowledge can be learned and fused together for a manifold of tasks, such as classification, regression or representation learning [12]. Furthermore, Deep Learning models comprise a large number of parameters to represent such hierarchical features, which are adjusted (*trained*) as per the task under consideration. In addition, Deep Learning approaches can model highly non-linear mappings between their inputs and outputs [13, 14]. Finally, decisions issued by these black-box models can be explained to non-expert users, making these black-box models of practical use in domains where explainability is a must [15].

In the Artificial Intelligence field we can find many evidences of the potential of the fusion of different technologies to tackle complex tasks. Deep Learning is not an exception to this statement: the fact that the architectural design, hyper-parameter tuning and training of Deep Learning can be formulated as optimization problems has motivated a long story between these models and the field of bio-inspired optimization, particularly Evolutionary Computation and Swarm Intelligence methods. The number of layers, their dimension and type of neurons, intermediate processing elements and other structural parameters can span a large solution space demanding search heuristics for their efficient exploration. Similarly, hyper-parameter tuning in Deep Learning models can be approached via heuristic wrappers, whereas their training process is essentially the minimization of a task-dependent loss function with respect to all the trainable parameters.

The purpose of this manuscript is to perform a thorough assessment of the potential of meta-heuristic algorithms for Deep Learning, with a proper understanding of the current state-of-the-art of this research area. It is supported by an exhaustive critical examination of the recent literature falling in this intersection, and a profound reflection, informed with empirical results, on the lights and shadows of this research area. The contributions of this study can be summarized as follows:

- We perform a brief hindsight on the historical confluence of both research areas so as to frame the importance of our study.
- We mathematically define concepts and notions on bio-inspired optimization and Deep Learning that help the reader follow the rest of the overview.
- We present a taxonomy that allows categorizing every proposal reported so far according to three criteria: a) the Deep Learning model under consideration; b) the optimization goal for which the bio-inspired algorithm is devised, distinguishing among architectural design, hyper-parameter configuration and model training; and c) the search mechanism followed by the bio-inspired solver(s) in use, either Evolutionary Computation, Swarm Intelligence or hybrid methods.
- Based on this taxonomy, we perform a detailed examination of the literature belonging to each category, pausing at milestones that supposed a genuine advance in the field, as well as identifying poor practices and misconceptions that should be avoided for the benefit of the community.
- We design and discuss on two experiments focused on two cases of study aimed at eliciting empirical evidence on the performance of bio-inspired optimization algorithms when applied to the topological design, hyper-parameter tuning and training of Deep Learning models.
- We provide a series of lessons and methodological recommendations learned from the literature analysis and the experiments, which should establish the minimum requirements to be met by future studies on the fusion of bio-inspired optimization and Deep Learning.

- A prospect is made towards the future of this research area by identifying several challenges that remain insufficiently addressed to date, and by suggesting research directions that could bring effective solutions to such challenges.

In summary, we comprehensively review and critically examine contributions made in this research avenue based on three axes: a) *optimization and taxonomy*, which comprises a historical perspective on this fusion of technologies, a clear definition of the optimization problems in Deep Learning, and a taxonomy associated to an in-depth analysis of the literature; b) *critical analysis*, informed by two case studies, which altogether elicit a number of lessons learned and recommendations for good practices, and c) *challenges* that motivate new directions of research for the near future. These three axes aim to provide a clear response to four important questions related to Evolutionary Deep Learning<sup>1</sup>, which are represented in Figure 1:

- Why are bio-inspired algorithms of interest for the optimization of Deep Learning models?
- How should research studies falling in the intersection between bio-inspired optimization and Deep Learning be made?
- What can be done in future investigations on this topic?
- What should future research efforts be conducted for?

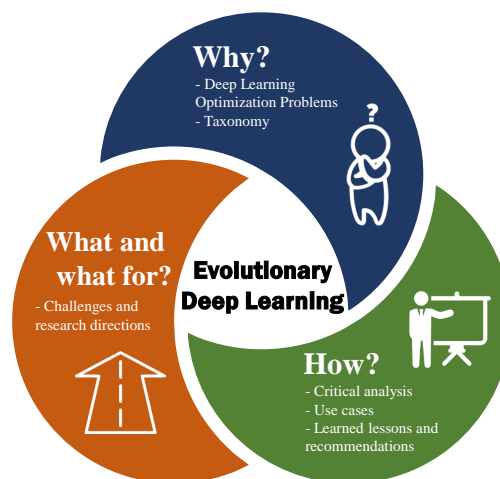


Figure 1: Diagram depicting the three axes and four fundamental questions on Evolutionary Deep Learning tackled in the overview, along with the specific aspects that contributes to each question.

The structure of the manuscript conforms to the above three axes: first, Section 2 briefly overviews the historical connection between Deep Learning and bio-inspired optimization. Next, Section 3 defines mathematically the optimization problems associated to Deep Learning models. Section 4 presents our taxonomy and an analysis of the literature falling on each of its categories. Section 5 exposes methodological caveats resulting from our critical literature study. Sections 6 and 7 present the two designed cases of study, and discuss the results obtained therefrom. Section 8 enumerates learned lessons and prescribes

<sup>1</sup>Throughout the survey we embrace the term *Evolutionary Deep Learning* to refer to the use of bio-inspired algorithms for solving optimization problems related to Deep Learning, no matter if they belong to Evolutionary Computation or to Swarm Intelligence.

good practices to be followed by prospective studies. Section 9 outlines several challenges and research directions that should drive future research efforts of the interested audience. Finally, Section 10 point out the final outline and conclusions of the overview. Additionally, Appendix A revisits several Deep Learning models, whereas Appendix B introduces the reader to the field of bio-inspired optimization, placing emphasis on Evolutionary Computation and Swarm Intelligence.

## 2. Historical Tour on Evolutionary Neural Networks and Evolutionary Deep Learning

Despite its relative youth, the current momentum of the synergy between Deep Learning and bio-inspired optimization is founded on a set of historical milestones that suggested the scientific community to combine these two branches of Artificial Intelligence. We herein revisit briefly this profitable background that led into the literature mainstream that motivates the current study. Figure 2 summarizes graphically such milestones, arranging them in a timeline along with the number of related publications reported in the last few years (data retrieved from the Scopus database, with the search terms indicated in the caption of the figure).

Although timid attempts at evolving neural networks with bio-inspired solvers had been reported in the late 90s [16], it was not until 2002 when Stanley and Miikkulainen settled a major breakthrough in the research community with their seminal work “Evolving neural networks through augmenting topologies” [17]. The NEAT approach proposed in this work allowed connection and layer types of an artificial neural network architecture (ANN) to be optimized by means of a meta-heuristic algorithm towards a progressively better precision of the evolved ANN model for a given task. NEAT embraces the main workflow of population-based meta-heuristics, particularly genetic algorithms: a population of encoded candidates is generated representing several network architectures, from which new candidate architectures are produced and evaluated on a given task (in the original work, a Reinforcement Learning task). After all candidates in the population have been evaluated, mutation and crossover operators are applied, generating a new population by means of combining network architectures, generating new layers or varying their hyper-parameters. This iterative search process is stopped when a stopping criterion set beforehand is met. As just stated, this algorithm was mainly proposed to solve reinforcement learning tasks, keeping in mind the profit of applying meta-heuristics to such environments, such as getting interesting behaviours and not falling in local optima in expense of precision.

Shortly after its first publication, NEAT’s unprecedented results spurred a flurry of new extensions and variants, not only in terms of new tasks and applications, but also in what refers to its core algorithmic components. Regarding the latter, the acknowledged importance of using good network encoding strategies soon became a major research goal in this literature strand, given the huge search space spawned by the evolution of architecture and weights of ANNs. In its original version, NEAT encoded candidates using a Direct encoding strategy, i.e. networks’ hidden units, connections and parameters (phenotype) were *directly* represented as an array representing each point of the network (genotype). However, in 2009 Stanley et al. proposed the so-called Hypercube-based NEAT, or HyperNEAT, [18], which relies on a generative encoding strategy to evolve large-scale neural networks using geometric regularities of the task domain and compositional pattern producing networks (namely, an ANN variant comprising multiple potentially heterogeneous activation functions that can be evolved via genetic algorithms). Besides the optimization of the activation function of each neuron in the network, HyperNEAT also proposed to utilize an indirect encoding to represent the networks to be evolved, inheriting other concepts from preceding NEAT versions such as speciation and historical marking. A new NEAT approach was developed years later for CNNs, which was coined as CoDeepNeat [19]. Following the NEAT design principles, CoDeepNeat was proved to excel at evolving layers, parameters, topology and hyper-parameters of CNNs. To this end, it uses an indirect encoding approach so that the network information can be encoded as rules or processes for creating individuals, ultimately yielding a reduced representation of the search space that can be explored more efficiently by the bio-inspired algorithm in use.

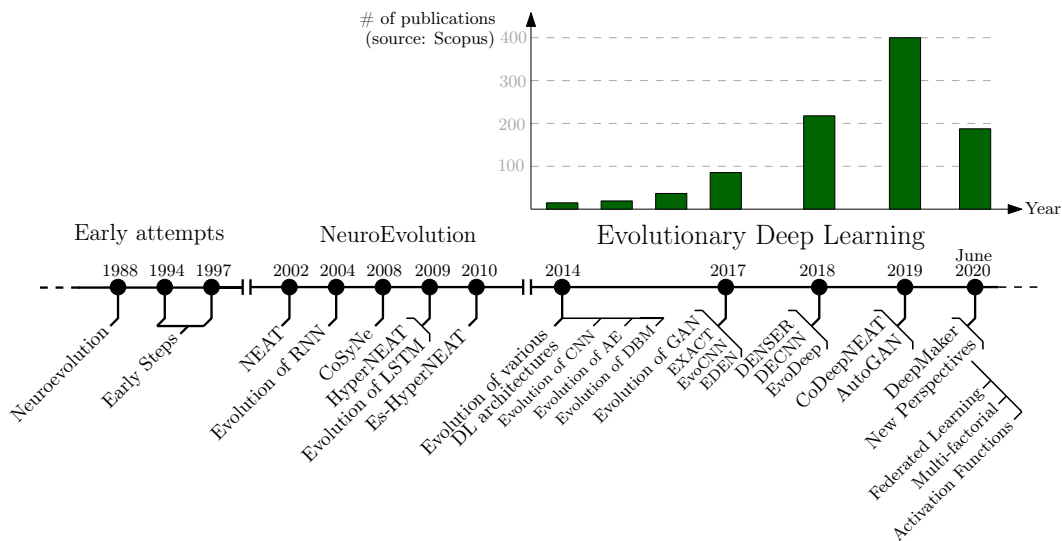


Figure 2: Timeline with main milestones in the history of Evolutionary Machine Learning, and a bar diagram showing the number of publications reported in this research area during the period 2013-June 2020. Data retrieved from Scopus by submitting the query (EVOLUTIONARY OR SWARM INTELLIGENCE) AND DEEP LEARNING.

Since its first appearance, neuroevolution approaches (with NEAT at their forefront) have been applied to multitude of tasks and problems. A major fraction of them relate to reinforcement learning, such as car controllers [38] or first-person agents control [39]. Interestingly, years after these applications were reported, the community shifted its research focus towards bio-inspired algorithms as an efficient replacement to train Deep Reinforcement Learning methods [40], showing up competitive results. Bio-inspired optimization was also applied to other Deep Learning tasks at the time, particularly for CNNs [41], and showcased in a diversity of applications including the classification of epileptic EEGs [42], time series forecasting [43] or scheduling deicing tasks in airports [44].

As a result of these findings, a growing literature corpus started to explore the potential of neuroevolution strategies for the topology and structural hyper-parameter optimization of different Deep Learning models, including AE, DBM and GANs. Naturally, the research community soon flowed into a further use of bio-inspired algorithms: the optimization of the trainable parameters of Deep Learning models [45]. An early approach was explored in [46], where a genetic algorithm is used for architecture optimization, and differential evolution for weight optimization. The main scientific interest of this work and those published thereafter is to assess whether meta-heuristics can avoid the convergence limitations of gradient descent methods when facing strongly non-convex search spaces as those characterizing the problem of training complex Deep Learning architectures. Nonetheless, most agreed on the dilated computation time and enormous computational resources needed by meta-heuristic solvers, which outweigh in practice the eventual convergence gains reported in experimental studies. As a result, gradient back-propagation approaches have dominated as the Deep Learning training solvers of choice over the years. Recently, this tendency is again emerging, stimulated by advances in highly parallel computing paradigms (including GPU and distributed asynchronous computation), prospecting new ways to train Deep Learning models with bio-inspired algorithms appearing frequently as promising alternatives.

Although it is still an incipient research field, some huge companies like Google, Facebook or Uber have glimpsed the power of this hybridization, and have started investing massively in software frameworks towards optimizing their Deep Learning models with meta-heuristics. The research community has also joined this momentum by open-sourcing software packages for this same purpose. EvoDeep [22],

Table 1: Software frameworks to evolve Deep Learning models using meta-heuristics and other search strategies.

Reference	Name	Year	Deep Learning models	Optimization domains	Optimization algorithm	Programming Language
[20]	EvoCNN	2017	CNN	Topology, structural hyper-parameters, weight initialization	Genetic algorithm	Python
[21]	DNF	2018	CNN	Trainable parameters (weights and biases)	Particle Swarm Optimization, Global-best Harmony Search and Differential Evolution	Python, TensorFlow
[22]	EvoDeep	2018	CNN	Topology, structural hyper-parameters, training hyper-parameters	Evolutionary Algorithm with specific mutation and crossover operators	Python, Keras, TensorFlow
[23]	DENSER	2019	CNN	Topology, structural hyper-parameters, training hyper-parameters	Genetic Algorithm and Dynamic Structured Grammatical Evolution	Python
[24]	Google Cloud AutoML	2019	CNN, RNN	Topology, structural hyper-parameters, training hyper-parameters	Transfer learning + neural architecture search	Proprietary
[25]	AutoKeras	2019	CNN, RNN	Topology	Hyperband, Random, Greedy	Python, Keras
[26]	ENAS	2018	CNN, RNN	Topology	Policy gradient-based subgraph selection	Python, Tensorflow
[27]	LEAF	2019	CNN, RNN	Topology, structural hyper-parameters, training hyper-parameters	CoDeepNEAT	Python, Pytorch
[28]	EvoAAA	2020	AE	Topology	Genetic Algorithm, Evolution Strategy, Differential Evolution	R
[29]	Ludwig	2019	CNN, LSTM, RNN, Fully-Connected	Structural and training hyper-parameters	Tree of Parzen Estimators	Python, TensorFlow
[30]	Keras-CoDeepNEAT	2020	CNN, RNN	Topology, structural hyper-parameters, training hyper-parameters	CoDeepNEAT	Python, Keras
[31]	MetaQNN	2017	CNN	Topology, structural hyper-parameters	Reinforcement Learning	Python, Caffe
[32]	DEvol	2017	CNN	Topology, structural hyper-parameters	Genetic Algorithm	Python, Keras
[33, 34]	AI Labs Neuroevolution Algorithms	2017-2018	CNN	Trainable parameters (weights and biases)	Genetic Algorithm + Novelty Search	Python, TensorFlow
[35]	CGP-CNN	2017	CNN	Topology, structural hyper-parameters	Cartesian Genetic Programming	Python, Chainer
[36]	AdaNet	2017	Architectures that can be represented as a directed acyclic graph	Topology, structural hyper-parameters, training hyper-parameters, trainable parameters	AdaNet Algorithm	Python, TensorFlow
[37]	Auto-Pytorch	2018	CNN (+ shallow learning)	Topology, structural hyper-parameters, training hyper-parameters	Random-forest-based Bayesian optimization	Python, Pytorch

AutoKeras [52] or Google Cloud AutoML are noteworthy platforms used for autonomously optimizing Deep Learning models, which we collect in Table 1 together with other alternatives from the literature.

The complexity of these problems, given by the cardinality of their search spaces and/or the large number of variables to be optimized, has stimulated a ever-growing corpus of literature that lasts to date. Advances held within this area have been reviewed in a number of surveys on this topic, listed in Table 2. However, our critical inspection of the achievements in this area reported over the years reveals poor methodological practices, unsolved technical caveats and research challenges that deserve a detailed analysis of where we currently stand in this effervescent area.

### 3. Deep Learning: Fundamentals and Optimization Problems

Before delving into the rest of this work, it is first convenient to settle the mathematical formulation of the optimization problems that lie at the core of this review. In this way, we establish what we understand by the different criteria in which the subsequent literature study is organized.

Table 2: Recent overviews on evolutionary Deep Learning and related topics.

Survey	Period	# reviewed works	Taxonomy	Coverage (DNN models/tasks)	Empirical study	Lessons learned and challenges
[47]	1987-2016	~ 15	Yes (optimization domain of the neural network*)	CNN, RNN, RL, DBN	No	Relevance of data quality. Evolutionary techniques good at exploration and exploitation but no single method for all optimization tasks.
[48]	2014-2018	~ 50	No (temporal analysis)	CNN, RL	No	High computational resources are required. Special emphasis on ensembles, transfer learning, multiobjective and modular evolutionary approaches.
[49]	2011-2018	~ 15	Yes (NN-based or GP-based and optimization problem(Architecture, Training, Multi-objective))	CNN	No	Lack of mathematical foundations, computational costs, scalability, poor generalization ability of the evolved model, lack of interpretability
[50]	2011-2019	~ 90	Yes (Evolutionary/Swarm Intelligence and Deep Learning model)	CNN, DBN, RNN, AE	No	Lack of rigurocity by the community. Costs of implementation, run time and overfitting.
[51]	2012-2019	~ 20	Yes (meta-heuristic and Deep Learning Architecture)	CNN, DBM, DBN	No	Time, more efforts on enhancing convergence speed and complexity (meta optimization)
<b>Ours</b>	2014-2020	~ 160	Yes (Deep Learning model/task and optimization problem)	CNN, AE, DBM, DBN, RNN, GAN, RL	Yes	See Sections 8, 9 and 10

Note: The column “# reviewed works” only takes into account the number of papers related to the optimization of Deep Learning problems. Any other non-related reference has been filtered out and not considered in the reported quantities.

\*: weights, architecture + weights, input layer, node, learning algorithm, combination of domains.

A Deep Learning model can be seen as a black-box optimizer where some parameters can be manually selected so that the model behaves in a different way. In fact, almost all parameters that can be tuned in a Deep Learning model can be treated as a task to be optimized. Therefore, depending on the parameters to be solved, we can differentiate different optimization problems. This being said, a Deep Learning model can be conceived mathematically as a composition of  $N$  different functions (*layers*) that maps its input  $\mathbf{x}_n \in \mathcal{X}_n$  to an output  $\mathbf{y}_n = f_n(\mathbf{x}_n; \mathbf{W}_n; T_n, \boldsymbol{\theta}_n) \equiv f_{n, \mathbf{W}_n}^{T_n, \boldsymbol{\theta}_n}(\mathbf{x}_n)$ , where:

- $T_n \in \mathcal{T}$  denotes the *type of layer*, with  $\mathcal{T}$  denoting the set of possible layer types (e.g. convolutional, LSTM, GRU).
- $\boldsymbol{\theta}_n$  is the vector of *structural hyper-parameters* of the layer. The specific parameters in this vector depend on the type  $T_n$  of the layer (e.g.  $\boldsymbol{\theta}_n$  will specify the sizes of the convolutional filters only if  $T_n = \text{convolutional}$ ).
- $\mathbf{W}_n$  denotes the *trainable parameters* (weights/filter coefficients and biases) of layer  $n$ , whose type and cardinality depend on  $T_n$  and  $\boldsymbol{\theta}_n$ . For instance, if  $\mathbf{x}_n$  represents RGB images (3 channels),  $T_n = \text{convolutional}$  and  $\boldsymbol{\theta}_n$  establishes that layer  $n$  comprises five  $3 \times 3$  convolutional filters,  $\mathbf{W}_n$  will comprise  $3 \times 3 \times 5 \times 3$  weights and 5 biases, yielding a total of  $|\mathbf{W}_n| = 140$  trainable parameters.

It is important to highlight, at this point, that the values of the trainable parameters  $\mathbf{W}_n$  must be learned by the model to efficiently perform a given task.

For the sake of simplicity in subsequent derivations, we will assume that we deal with a supervised learning task in which we assume a training dataset  $\mathcal{D}_{tr} = \{(\mathbf{x}_1^m, \mathbf{y}_N^m)\}_{m=1}^{M_{tr}}$ , with  $\mathbf{y}_N^m \in \mathcal{Y}$  denoting the supervised output of input  $\mathbf{x}_1 \in \mathcal{X}_1$ , and  $M_{tr}$  representing the number of training instances. The trainable parameters of a Deep Learning model are learned from  $\mathcal{D}_{tr}$  by means of a *training* algorithm  $\{\mathbf{W}_n\}_{n=1}^N = \text{ALG}(\mathcal{D}_{tr}, \{T_n, \boldsymbol{\theta}_n\}_{n=1}^N; \boldsymbol{\vartheta})$ , where we refer to  $\boldsymbol{\vartheta}$  as the set of *training hyper-parameters* of the training algorithm. In general, the training algorithm is driven by the minimization of a task-dependent loss function  $L(\hat{\mathbf{y}}_N^m, \mathbf{y}_N^m)$  that provides a measure of error between the supervision  $\mathbf{y}_N^m$  of input  $\mathbf{x}_1^m \in \mathcal{D}$  and the corresponding output of the Deep Learning model:

$$\hat{\mathbf{y}}_N^m = F(\mathbf{x}_1^m; \{\mathbf{W}_n\}_{n=1}^N; \{T_n\}_{n=1}^N, \{\boldsymbol{\theta}_n\}_{n=1}^N) \doteq f_{N, \mathbf{W}_N}^{T_N, \boldsymbol{\theta}_N} \circ f_{N-1, \mathbf{W}_{N-1}}^{T_{N-1}, \boldsymbol{\theta}_{N-1}} \circ \dots \circ f_{1, \mathbf{W}_1}^{T_1, \boldsymbol{\theta}_1}(\mathbf{x}_1^m), \quad (1)$$

where  $\circ$  denotes composition of functions, i.e.  $f \circ g(x) = f(g(x))$ . Such a measure of loss computed for every training instance can be averaged to yield a numerical estimation of the performance of the DL model when approximating the supervised instances in  $\mathcal{D}_{tr}$ :

$$L(F; \mathcal{D}_{tr}) = \frac{1}{M_{tr}} \sum_{m=1}^M L(F(\mathbf{x}_1^m; \{\mathbf{W}_n\}_{n=1}^N; \{T_n\}_{n=1}^N, \{\boldsymbol{\theta}_n\}_{n=1}^N), \mathbf{y}_N^m). \quad (2)$$

With this notation in mind, we define the following optimization problems that underlie the construction of Deep Learning models:

**Problem 1.** (*Topological Optimization*) Given a learning task defined on a training dataset  $\mathcal{D}_{tr}$ , the topological optimization of a DL model refers to the search for the topology of the DL model that best solves the task at hand, wherein topology involves the discovery of the optimal number of layers  $N$  and their types  $\{T_n\}_{n=1}^N$ . This problem assumes fixed values for  $\{\boldsymbol{\theta}_n\}_{n=1}^N$  (e.g. standard values), and relies on a training algorithm  $ALG(\mathcal{D}_{tr}, \{T_n, \boldsymbol{\theta}_n\}_{n=1}^N; \boldsymbol{\vartheta})$  to optimize the trainable parameters  $\{\mathbf{W}_n\}$ . Mathematically:

$$\min_{N, \{T_n\}_{n=1}^N} L(F; \mathcal{D}_{tr}), \quad (3)$$

where the dependence of the aggregate loss function with respect to  $N$  and  $\{T_n\}_{n=1}^N$  comes through (2), and  $\{\mathbf{W}_n\}_{n=1}^N$  are optimized by means of  $ALG(\mathcal{D}_{tr}, \{T_n, \boldsymbol{\theta}_n\}_{n=1}^N; \boldsymbol{\vartheta})$ .

Topology optimization is rarely conceived in isolation with respect to the rest of variables that define a DL model. Instead, topology is often optimized along with the values of their structural hyper-parameters. However, we define this second problem separately so as to allow for a fine-grained literature analysis:

**Problem 2.** (*Structural Hyper-parameter Optimization*) Given a learning task defined on a training dataset  $\mathcal{D}_{tr}$ , and a fixed topology of the DL model ( $N$  and  $\{T_n\}_{n=1}^N$ ), the optimization of the structural hyper-parameters of the DL model aims to find the best value of  $\boldsymbol{\theta}_n$  (structural hyper-parameters) for each of their compounding layers. Mathematically:

$$\min_{\{\boldsymbol{\theta}_n\}_{n=1}^N} L(F; \mathcal{D}_{tr}), \quad (4)$$

where the dependence of the aggregate loss function with respect to variables  $\{\boldsymbol{\theta}_n\}_{n=1}^N$  comes through (2), and  $\{\mathbf{W}_n\}_{n=1}^N$  are optimized by means of  $ALG(\mathcal{D}_{tr}, \{T_n, \boldsymbol{\theta}_n\}_{n=1}^N; \boldsymbol{\vartheta})$ .

Finally, the third optimization problem that can be formulated is the training process itself, which aims at finding the values of the parameters  $\{\mathbf{W}_n\}_{n=1}^N$  that minimizes the loss in (2). This is indeed the purpose of  $ALG(\mathcal{D}_{tr}, \{T_n, \boldsymbol{\theta}_n\}_{n=1}^N; \boldsymbol{\vartheta})$ . However, we note at this point that two different formulations of this problem can be made depending on whether variables to be optimized include the set of *training hyper-parameters*  $\boldsymbol{\vartheta}$ :

**Problem 3.** (*Training Hyper-parameter Optimization*) Given a learning task defined on a training dataset  $\mathcal{D}_{tr}$ , a fixed topology of the DL model ( $N$  and  $\{T_n\}_{n=1}^N$ ), fixed values of their structural hyper-parameters  $\boldsymbol{\theta}_n$ , and a training algorithm  $ALG(\mathcal{D}_{tr}, \{T_n, \boldsymbol{\theta}_n\}_{n=1}^N; \boldsymbol{\vartheta})$ , the training hyper-parameter optimization problem of a DL model aims to find the best value of  $\boldsymbol{\vartheta}$  (training hyper-parameters) as:

$$\min_{\boldsymbol{\vartheta}} L(F; \mathcal{D}_{tr}), \quad (5)$$

where the dependence of the aggregate loss function with respect to  $\boldsymbol{\vartheta}$  comes through the application of  $ALG(\mathcal{D}_{tr}, \{T_n, \boldsymbol{\theta}_n\}_{n=1}^N; \boldsymbol{\vartheta})$  to solve for  $\{\mathbf{W}_n\}_{n=1}^N$  as per (2).



**Problem 4.** (Trainable Parameter Optimization) Given a learning task defined on a training dataset  $\mathcal{D}_{tr}$ , a fixed topology of the DL model ( $N$  and  $\{T_n\}_{n=1}^N$ ), and fixed values of their structural hyper-parameters  $\theta_n$ , the trainable parameter optimization problem of a DL model seeks the best value of  $\{\mathbf{W}_n\}_{n=1}^N$  (trainable parameters) as:

$$\min_{\{\mathbf{W}_n\}_{n=1}^N} L(F; \mathcal{D}_{tr}), \quad (6)$$

for which an optimization (training) algorithm  $ALG(\mathcal{D}_{tr}, \{T_n, \theta_n\}_{n=1}^N; \vartheta)$  is utilized.

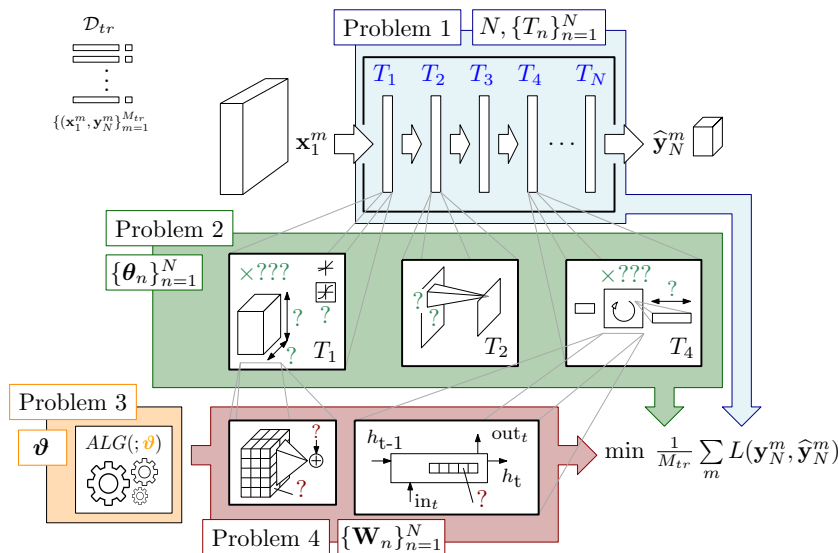


Figure 3: Optimization problems in Deep Learning for a generic model comprising, among others, a convolutional layer, a max-pooling layer and a recurrent layer.

A visual summary of the above problems is sketched in Figure 3. The above 4 optimization problems can represent the majority of contributions so far elaborating on new algorithms to address them efficiently. However, several practical remarks must be made on these definitions:

- First of all, it is important to highlight that other measures are often used in the objectives of these problems to replace the aggregate loss  $L(F; \mathcal{D}_{tr})$ , particularly in Problems 1, 2 and 3. This is the case of cross-validated task-dependent performance metrics (e.g. accuracy or  $F_1$  for classification problems), or even the same measures computed over a validation holdout so as to reduce the computational burden of repeatedly evaluating different solution candidates. In Problem 4, however, differentiable loss functions like the ones typically used by gradient-based back-propagation algorithms are often used disregarding the training algorithm adopted.
- Since the evaluation of the aggregate loss function in Expression (2) depends on all optimization variables considered in Problems 1 (topology), 2 (structural hyper-parameters) and 3 (training hyper-parameters), it is often the case that studies in the literature consider several problems jointly (e.g. joint topology and structural hyper-parameter optimization). However, as stated previously we find it very convenient to explicitly differentiate among all problems to set clear the optimization domain in which each contributes to the general understanding of the field.
- Given the scope of our study it is relevant to emphasize that the goal of the reviewed literature strand is to devise an algorithm that, when solving for the variables of each of the above problems, produces

lower loss values than other solvers with respect to the dataset and task at hand. However, the design target (the sought *optimization algorithm*) varies depending on the problem under consideration:

- In Problems 1, 2 and 3, the design target is an optimization algorithm that solves for  $N$ ,  $\{T_n\}_{n=1}^N$  (Problem 1),  $\{\theta_n\}_{n=1}^N$  (Problem 2) and  $\vartheta$  (Problem 3). This algorithm under target operates as a wrapper of the overall model, working in parallel to the training algorithm  $ALG(\mathcal{D}_{tr}, \{T_n, \theta_n\}_{n=1}^N; \vartheta)$ . As the latter falls out from the target of of the design process, the training algorithm in use is often set to a naive gradient back-propagation solver (e.g. stochastic gradient descent, Adam and the like).
- In Problem 4, the design target is the optimization algorithm  $ALG(\mathcal{D}_{tr}, \{T_n, \theta_n\}_{n=1}^N; \vartheta)$  itself solving for the trainable parameters  $\{\mathbf{W}_n\}_{n=1}^N$ , hence there is only one single optimization algorithm.

#### 4. Taxonomy and Literature Review

In light of the past history between bio-inspired optimization and Deep Learning, a need arises for properly organizing contributions so far in a taxonomy that covers which problems are addressed, which Deep Learning models are involved, and which bio-inspired algorithms are in use. In this section we perform this analysis, centering the discussion around a taxonomy that sorts the literature according to the three aforementioned criteria.

The main purpose of this taxonomy (Subsection 4.1) and the literature analysis made over each of its categories (Subsections 4.2, 4.3 and 4.4) is to highlight those areas where the community has so far placed most research efforts. This literature analysis settles a firm stepping stone towards a critical discussion of poor methodological practices and points of improvement observed in related contributions to date: the *shadows* in which this field is held nowadays. Such a discussion will be held in Section 5.

##### 4.1. Taxonomy

As has been stated in Section 3, we distinguish among four main optimization tasks: topological optimization (Problem 1), structural hyper-parameter optimization (Problem 2), training hyper-parameter optimization (Problem 3) and trainable parameter optimization (Problem 4). Our taxonomy gathers Problems 2 and 3 under the general hyper-parameter tuning category, discriminating between them in a lower level of the taxonomy.

The main reason for this special arrangement of the taxonomy is to highlight that as per the reviewed literature (>160 references), there is little explicit distinction between structural hyper-parameter and training hyper-parameter in related contributions. Our thorough examination of this corpus has discriminated interesting research opportunities in the extrapolation of studies and frameworks, from training to structural hyper-parameter tuning and vice versa. When it comes to Problem 4, a distinction is made between i) bio-inspired algorithms that do not incorporate any problem-specific knowledge in their design; and ii) bio-inspired solvers that are hybridized with local search solvers or combined with gradient back-propagation techniques. [175, 176, 177]

Figure 4 depicts graphically the taxonomy considered for the literature analysis. In the first level we consider the type of optimization problem under consideration (topology, hyper-parameter and trainable parameter optimization), followed by contributions sorted as per the Deep Learning model (Appendix A) and kind of bio-inspired solver (Appendix B) under choice. In what follows we analyze in depth on the contributions classified within each of these categories.

##### 4.2. Topology optimization

It is widely acknowledged that the topology or architecture of Deep Learning models have a direct impact on their performance. For this reason researchers have traditionally striven to develop automated methods for generating topologically small yet well-performing network architectures. In this context, the set of algorithms gathered under the *Neuroevolution* (NE) label aim at progressively augmenting

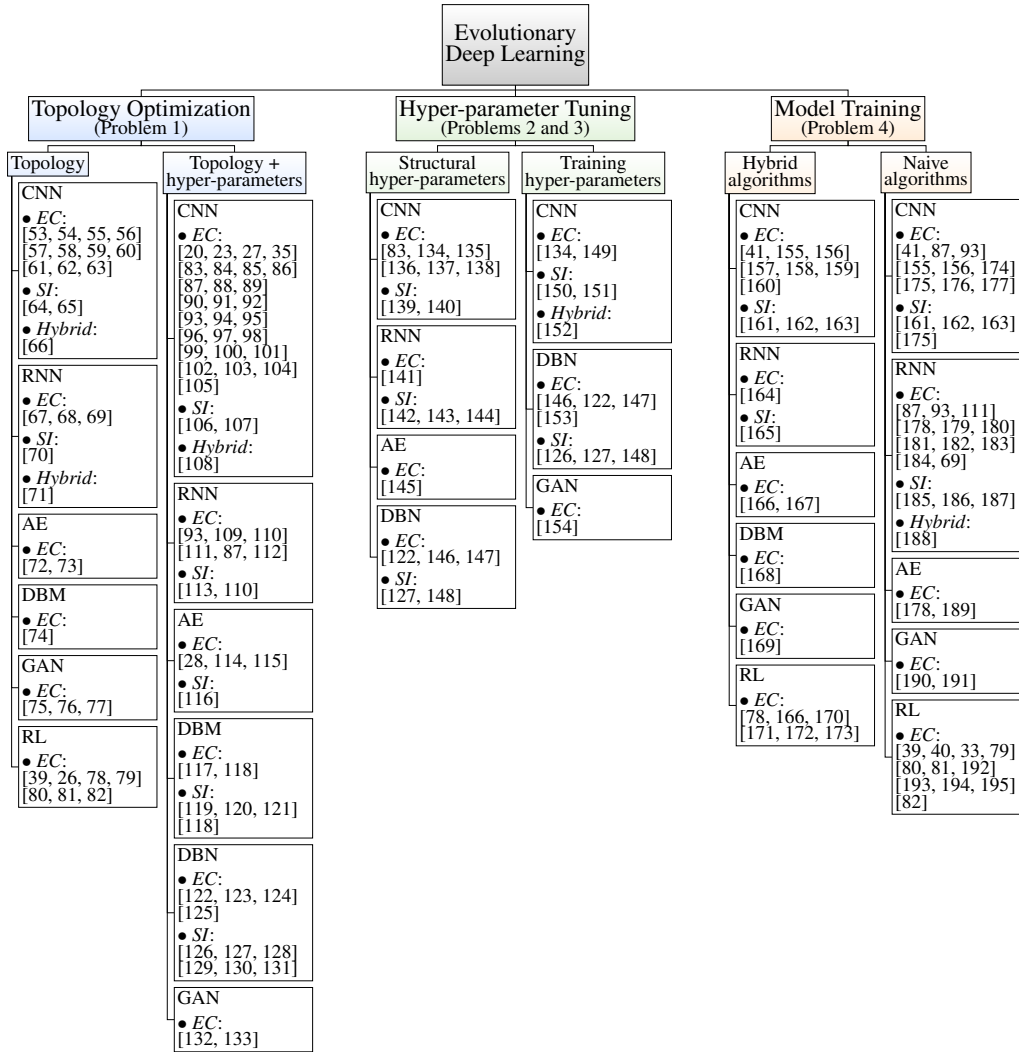


Figure 4: Taxonomy of the reviewed literature on Evolutionary Computation and Swarm Intelligence algorithms applied to the optimization of Deep Learning models. The taxonomy is structured by the domain of the Deep Learning model under focus (topology, hyper-parameters and trainable parameters), further discriminated by the specific Deep Learning model under consideration and the type of bio-inspired algorithm in use (*EC*: Evolutionary Computation; *SI*: Swarm Intelligence; *Hybrid*: a mixture of both).

the complexity of neural network topologies to attain increasingly better generalization properties while keeping its complexity to its minimum required. Originally applied to ANN models, different NE variants have been applied in the last few years to optimize Deep Learning models, not only in terms of their topology, but also jointly with their structural and training hyper-parameters (e.g. kernel size, activation function, dropout and learning rate). In some few cases, trainable parameters have also been considered in the set of variables to be optimized via NE [41]. Furthermore, since they resort to evolutionary algorithms at its core, NE approaches have stimulated over the years a manifold of other bio-inspired approaches, in a way to assess whether the same optimization problem can be tackled more effectively with alternative search strategies and operators.

All in all, in terms of topological optimization CNNs are arguably the most targeted Deep Learn-

ing models to date. CNNs' topology optimization is faced by scientific community in two ways; layer by layer or by blocks. In layer-wise optimization, hyper-parameters are fixed and networks are fully evolved using bio-inspired solvers, such as GA [55] and customized versions of Evolutionary Algorithms [53]. In this last work, the so-called AmoebaNet-A model settled a state-of-the-art landmark score on the ImageNet dataset (83.9% accuracy), including comparisons to other search strategies (random search and reinforcement learning). Another approach proposed in [65] resorts to Particle Swarm Optimization (PSO) – a popular Swarm Intelligence solver – to optimize a block formed by dense layers. Once optimized, this block is stacked along with convolutional and pooling layers configured with fixed hyper-parameters, and ultimately used to address a image classification task. This work exemplifies a research trend focused on optimizing the topology of certain parts of the entire Deep Learning architecture, in an attempt at reducing the cardinality of the search space and speeding up the search process, at the cost of being much less exploratory in terms of network configurations than other counterparts. There is another important matter to be taken in account when performing topological optimization: the encoding of solutions, which impacts directly on the dimensionality of the search space. Actually, initial improvements of NE approaches were achieved thanks to novel network encoding strategies, which allowed for an easier exploration and less computational cost than preceding alternatives.

Another strand of literature has elaborated on more complex problem formulations by jointly addressing the optimization of the topology of the network along with its hyper-parameters. Again, CNNs have become central in related studies. An illustrative work is the one in [20], where an Evolutionary Algorithm is used with different mutation operators operating on topological variables, structural and training hyper-parameters such as the filter size, the convolution stride, learning rate or the insertion/removal of convolutional layers, among others. Studies in this field tend to be similar to each other in terms of the complexity of the optimization problem under consideration. Thus, new proposals are usually made by customizing the operators (mutator, selector) of optimization algorithms or by developing custom encoding strategies, as in [106] where a PSO variant is introduced based on an IPv4 based codification scheme with varying length.

Despite the predominance of CNNs in topological optimization, RNNs (and in particular, LSTM networks) have also been a subject of study in this research area. In [67] a Differential Evolution (DE) solver is proposed for achieving this purpose and efficiently undertaking a wind forecasting regression task. It is relevant to observe that when both architecture and hyper-parameters are evolved for RNNs, certain hyper-parameters are recurrently considered in related studies, such as learning rate, dropping frequency factor [113] or batch size [110, 109]. In general terms, the aforementioned DE appears to be the most applied meta-heuristics in LSTM. A few exceptions can be found, such as [113] (Bat Algorithm), [70] (Ant Colony Optimization) and [110], where a comparison is made between DE, PSO and SA, concluding that DE reaches better performance levels. Hybrid approaches have been also explored for the optimization of RNNs, as in [71] where the architecture (i.e. connection pattern) is optimized by means of a hybrid PSO-GA solver. An interesting point arises when inspecting in detail this set of studies: the creation of custom objective functions to allocate different (usually conflicting) optimization goals. The work in [67] is an example of how a customized objective function can yield topologically optimized network designs that achieve a balance between performance and model complexity, being the latter of particular interest for the deployment of the model in resource-constrained embedded devices.

Other Deep Learning models have grasped a remarkable attention in topological optimization. AEs have been optimized topologically, often along with structural hyper-parameters in those cases where convolutional layers are involved. In their original formulation, AEs are composed by stacked dense layers (*encoder*) producing a low-dimensional representation of the input, which is reconstructed by another set of stacked dense layers (*decoder*). A contribution from 2015 [73] presented a way to generate promising AE architectures by mutating candidates by using a customized Evolutionary Algorithm, whose mutation operator is based on the reconstruction error achieved by the decoder. Also in this vein, a mini-batch variant training method was proposed (*evo-batches*) aimed at reducing the computational cost when a large

number of candidate networks have to be evaluated in large datasets. Decoder and encoder topologies of studies related to AEs are often assumed to be symmetrical [73]. However, in [72] a more flexible architecture was proposed, where the decoder is evolved along with the encoder and does not have to mimic its architecture. In this reference several operations were applied to topological variables, such as layer addition (random number of neurons), layer removal, application of Gaussian perturbation to the number of neurons, or layer swapping. To wrap up the activity noted in AEs, we highlight the work in [116] and [114], where PSO and GA are respectively used to evolve topology and structural hyper-parameters of AEs comprising convolutional layers.

We proceed forward with our analysis by pausing at DBMs, whose architecture can be very similar to DBNs in terms of the structural hyper-parameters involved in the optimization process. Given this similarity and the relative scarcity of studies related to these models, we analyze them jointly in what follows. The majority of works related to the topological optimization of DBMs and DBNs pay special attention to the process of optimizing both architecture and some hyper-parameters. Moreover, most of them rely on Swarm intelligence algorithms, such as PSO or ACO. An exception can be found in [126] where Artificial Bee Colony (ABC) is used to optimize DBNs' structure, learning rate, momentum and weight decay. The results are compared to those yielded by other bio-inspired solvers: Firefly Algorithm (FA), CS and Harmony Search (HS). FA for 2- and 3-layered DBN, and ACO for single-layer DBN, resulted to yield the best performing network architectures for the image reconstruction task under consideration. The rest of contributions consider a combination of all or some of learning rate, momentum or weight decay hyper-parameters, focusing the application of the optimized model to different practical problems such as traffic flow prediction [127] or the detection of turbine failures [129]. The work in [120] introduces a novel way to optimize structure and hyper-parameters of DBMs, and compares the performance of DBMs for image classification when optimized with different flavors of the PSO solver, random search and several HS variants. They concluded that bio-inspired techniques are suitable to optimize DBMs, beating Random Search in all considered datasets. Nevertheless, network topologies optimized via PSO and HS solvers scored similar performance levels. This last observation connects directly with one of the points remarked in our critical analysis of Section 5.

In terms of algorithmic variants, some hyper-heuristic techniques like [122] proposed an approach to optimise DBNs' structural hyper-parameters, i.e. number of hidden units, along with non-structural hyper-parameters like learning rate, and some hyper-parameters related to the heuristic algorithm (number of epochs or iterations). Hyper-heuristics have also been used to optimize CNNs, [149] presents a method to select the best heuristics, where batch size, number of epochs, neurons on the fully connected layer, dropout and learning rates, rho and epsilon factors are evolved.

There are also a few works dealing with the topological optimization of GANs. In [132], a meta-heuristic approach to evolve GANs' discriminator and generator is introduced. Specifically, a GA is used to evolve the architecture, activation functions of each layer and initialization mode in both generator and discriminator. Furthermore, an optimization of the loss function is done, taking them from a bunch of well-known formulations. Besides, training hyper-parameters are also mentioned in this work as potentially evolvable variables (yet not optimized in practice), such as the gradient-based solver used to learn the trainable parameters, the batch size and the number of epochs. Shortly thereafter, Costa et al. [76] proposed an approach to optimize the architecture and parameters of both the generator and discriminator modules of a GAN. The approach was based on DeepNEAT and adapted to the context of GAN optimization. Linear, convolutional and transpose convolutional layers were directly mapped to a phenotype – an array of genes – representing the final network. In all layers the activation function was evolved, and in the case of convolutional and transpose convolutional layers, the output channels were also considered.

Finally, in the field of RL, the tendency observed in our literature analysis is to use NE approaches to optimize both topology and trainable parameter (weights) of the neural network mapping the output of the environments to the actions to be taken by the agent. Commonly, NEAT is used for this purpose [80, 79, 78, 81], which becomes in charge of optimizing the neural network involved in Deep RL approaches.

A real-time adaptation of NEAT was used in [39] to evolve agents for the NERO videogame, placing an emphasis on the need for efficient workarounds to alleviate the complexity of neuro-evolution methods.

On a summarizing note, the literature on bio-inspired algorithms applied to architecture optimization has a long history departing from NE, which was originally applied to evolve ANN architectures. Since then, many other meta-heuristics have been applied to optimize architecture and hyper-parameters of Deep Learning models. Given that networks can have variable-length topologies, a good solution encoding strategy is essential to lessen computational costs and the time of execution without hindering the representability of all network configurations. Remarkably, modern bio-inspired solvers such as FA, BA and CS have been lately used with competitive results with respect to classical solvers (EA, PSO and ACO).

#### 4.3. Hyper-parameter optimization

Arguably, one of the optimization tasks where bio-inspired methods have been traditionally applied within the Machine Learning field is hyper-parameter tuning. It is well-known that hyper-parameter tuning usually yields better performance levels than off-the-shelf model configuration. When shifting the focus towards the hyper-parametric optimization of Deep Learning models, two major fields are spotted. On the one hand, literature focused on optimizing parameters related to the training algorithms, such as *learning rate*, *batch size* or *momentum*. On the other hand, architectural hyper-parameters, which are layer-type-dependent, e.g. filter size, number of kernels, activation functions or stride size in CNNs. Actually, given the high number of structural hyper-parameters of convolutional layers, CNNs have protruded as one of the most explored Deep Learning models for hyper-parameter optimization, in many cases jointly with training hyper-parameters such as the learning rate, momentum or weight decay of gradient solvers. Although there are some contributions focused exclusively on the optimization of structural hyper-parameters, the mainstream is to jointly address the optimization of topology and hyper-parameters, as the literature on topological optimization examined in the previous subsection has clearly revealed.

Let us start from 2016, when [150] proposed a set of bio-inspired meta-heuristics (BA, FA and PSO) to optimize the aforementioned hyper-parameters in a CNN used for Parkinson disease identification from image data. Likewise, [134] proposed a DE-based approach to optimize filter sizes, number of neurons in the fully-connected network end, weight initialization policy and dropout rate for sentiment analysis. Comparisons were made to GA and PSO, achieving better results in terms of accuracy and computational efficiency. The optimization of dropout probability is other common approach that some authors have tackled using different solvers, including CS, BA, FA and PSO [151] or hybrid GA and Tabu Search algorithms [152]. In this latter work random search and Bayesian optimization were proven to perform worse than the proposed hybrid meta-heuristic algorithm over the considered image classification datasets. Batch size and learning rate were also regarded as optimization variables.

Moving on to RNNs, very few papers are focused only on hyper-parameter optimization, conforming to the general trend observed in the analyzed literature. Dropout optimization is tackled for LSTM networks in [142] and [143], where ACO is used for engines vibration prediction. The connections between neurons are activated/deactivated to accomplish the task, which is very similar to the approach carried out in [144]. However, the aspect to be highlighted in this last reference is that a PSO is utilized to optimize the output connections of an Echo State Network (ESN), a randomization-based RNN model belonging to the family of Reservoir Computing models. We will later revolve on the possibilities that we have found in the extrapolation of advances in bio-inspired optimization for Deep Learning models to other less studied models.

#### 4.4. Trainable parameter optimization

In recent years the advent and progressive maturity of new parallel and distributed computing paradigms have reignited the interest of the scientific community in applying bio-inspired optimization algorithms to train Deep Learning models. Although each model type is different in terms of topology, they share

some disadvantages resulting from the adoption of gradient back-propagation training methods, such as gradient vanishing/exploding and proneness to get stuck in local optima. Evidently, the complexity of the problem grows up as the number of parameters involved in the optimization process increase, yielding largely non-convex search landscapes. These acknowledged issues have been extensively studied by the community by proposing different workarounds. Nonetheless, an increasing trend towards the use of bio-inspired solvers for this purpose can be noticed in recent literature, as their search operators do not rely on gradient back-propagation anyhow, and therefore avoid its drawbacks effectively. Based on this rationale, we now delve into how the community has adapted different bio-inspired solvers for training Deep Learning models.

A first examination of the literature exposes two main tendencies followed by the community, which are reflected in the second level of the corresponding taxonomy branch of Figure 4:

- Approaches that combine bio-inspired solvers with traditional training algorithms, which aim to overcome the disadvantages we have just introduced. Almost the entirety of studies adopting this hybrid design strategy are focused on CNNs, implementing the aforementioned combination in many different ways. A straightforward way to overcome falling in local optima is to evolve an initial set of values for the trainable parameters (weights, bias) that sets the gradient back-propagation solver on a promising path towards the global minimum of the loss function. In [162] this approach is adopted for training a CNN using the ABC algorithm. Other works [155] combine GA and SGD: GA evolves new candidates through its search operators, but the fitness function is evaluated after some training epochs of stochastic gradient back-propagation (SGD). Similarly, in [163] PSO is used to evolve the trainable parameters of the last layer of a CNN, while the parameters of the rest of the layers are learned via SGD. Comparisons with the CNN trained exclusively with SGD rendered an enhanced convergence speed and final accuracy on image classification tasks. Last but not least, in [165] the CS algorithm is used to train RNNs following two strategies: one trained using only this solver, and the other combining CS and gradient back-propagation. A benchmark comparison to networks trained with conventional gradient back-propagation and different variants of the ABC algorithm discovered that all models trained using bio-inspired optimization techniques performed better than the RNN trained with gradient back-propagation.
- Approaches in which training is performed completely using bio-inspired optimization methods. Most references embracing this second strategy deal with RNNs and CNNs, and differ from each other mostly in the search algorithm being considered. All in all, a common approach is to evolve the trainable parameters of the model (via the search operators of the bio-inspired solver at hand), and evaluate it in terms of loss value or any other performance estimator linked to the task at hand. In this line, in [174] and [176] two SA-based solvers are proposed and assessed for optimizing the parameters of a CNN, achieving better performance scores and better convergence speed than the same model trained via gradient back-propagation. LSTM network training has also been tackled by using different bio-inspired optimization techniques. A good exponent is the work in [179], where HS, Gray Wolf Optimizer (GWO), Sine Cosine Algorithm (SCA) and Ant Lion Optimization (ALO) were compared to each other when used to learn the trainable parameters of different LSTM model configurations. We emphasize that despite the diversity of methods considered in this work, no comparisons to traditional training solvers were reported, uncovering one of the critical points discussed in Section 5.

Before proceeding with this critical analysis, we briefly comment on Deep RL models. Bio-inspired algorithms have been lately postulated as efficient alternatives to solve several optimization problems underlying these models. A first approach is to train parts of the architecture via evolutionary algorithms, and the rest using SGD. This is indeed what the study in [166] proposes: a Covariance Matrix Adaptation Evolution Strategy (CMA-ES) is used to evolve weights for the behavior-generating network, while the rest of the Deep RL architecture (a convolutional AE) is trained by using a gradient based method. This

work continued the research line started years before in [192], where CMA-ES was used to train simpler RL networks. All in all, Evolutionary Algorithms have largely demonstrated to be efficient solvers to train Deep RL models, as shown in renowned studies such as [33, 40] (with GA), [33] (GA with Novelty Search) and [194] (DE and Novelty Search). Recent works [193] have also explored the capabilities of multi-task optimization evolving multiple related RL tasks at the same time, taking advantage of the transfer of genetic information between tasks. In other works, NEAT is used to optimize the trainable parameters and the topology of a Deep RL model [80, 79]. On the other hand, in [170, 171] an hybrid EA algorithm is proposed, where a population of networks is trained and periodically evolved, exploiting Lamarckian transfer capabilities. This same procedure is used in [172] as a way to inject information about the gradient in the population of individuals maintained by the Evolutionary Algorithm during the search. Other hybridization techniques consist of the use of different networks in the same architecture, where some of them are trained via SGD, and others evolved with evolutionary operators. This is the case of [173], where the parameters of a RNN used for determining the actions of an agent are evolved using Cooperative Synapse Neuroevolution (CoSyNE), an evolutionary algorithm that enforces subpopulations at the level of a single trainable parameter. This work built upon the findings in [196] and extrapolated them to complex Deep RL models, showcasing how evolutionary algorithms can evolve small networks capable of reaching competitive performance levels.

## 5. Critical Methodological Analysis

The above corpus of reviewed literature sheds evidence on the vibrant activity of the intersection between bio-inspired optimization and Deep Learning. So far the community has reported interesting findings in what refers to hyper-parameter optimization, topological search and small/medium-sized network training. Notwithstanding this noted activity, our critical analysis of these contributions has disclosed a number of poor practices and methodological shortages that should be underlined to set them down in black and white. We now discuss briefly about these issues, settling the necessary rationale for the experimental part of this survey:

- The lack of benchmark datasets/tasks to validate new advances (Subsection 5.1).
- The unrealistic scales of the Deep Learning models optimized via bio-inspired methods (Subsection 5.2).
- The need for good methodological practices when comparing among different solvers used for Deep Learning optimization (Subsection 5.3).
- The limited utility of software implementations (Subsection 5.4).
- The existence of metaphor-based publication series (Subsection 5.5).

### 5.1. Lack of benchmark datasets/tasks

A major problem observed in the literature is the heterogeneity of datasets used to validate new algorithmic approximations for the optimization problems under analysis. Even when the task is clearly defined (e.g. image classification or time series forecasting), the possibility to compare the results obtained by different studies becomes unfeasible since the considered datasets are not the same.

For the community to gain verifiable evidence about the claimed gains of upcoming proposals, consensus should be reached about the datasets/tasks that should be utilized for comparison purposes in the future. Unfortunately, the diversity of datasets/tasks over which some of the new contributions are assessed seem to go in the opposite direction, calling into question whether the reported performance improvements can be extrapolated to other learning problems.



### 5.2. *Unrealistic complexity of Deep Learning models*

Besides the heterogeneity of datasets/tasks discussed above, it is often the case that the Deep Learning models under consideration do not meet the complexity levels of the state of the art for the task under consideration. This is particularly concerning in works related to model training (Problem 4), where the cardinality of the search space faced by the meta-heuristic solver is in the order of several thousands to millions of optimization variables (trainable parameters). For instance, a recent work on image classification using the well-known MNIST dataset has recently established a new record in accuracy (0.1% test error rate) with a model comprising 57.02 millions of trainable parameters [197]. However, most of the reviewed literature on training via bio-inspired solvers rarely considers Deep Learning models that surpass a few thousands of trainable parameters.

This issue again calls for a major reflection on whether research advances are missing the real challenge underneath the use of bio-inspired algorithms in such large search spaces (scalability, exploitation of the correlation among decision variables), to instead focus on minor aspects of doubtful scientific contribution.

### 5.3. *Comparison methodology*

Even if addressing and effectively solving the preceding two issues, several methodological aspects still remain often overseen when comparing among different solvers for a given task/dataset/optimization problem scenario:

- **Baseline schemes:** our literature analysis revealed that a fraction of contributions discussed on extensive experiments with several new meta-heuristic algorithms for a given optimization problem, without including in the benchmark standard solvers utilized in the past for the same problem. This, again, is particularly worrying in regards to Problem 4 (model training): comparisons should compulsorily include gradient back-propagation based solvers widely used for the same purpose (e.g. SGD, Adam). Overlooking the analysis of whether bio-inspired algorithms perform competitively with respect to established solvers for the same purpose is counterproductive for the potentiality of this research area.
- **Objective function(s):** we noticed that relevant divergences emerge in how the optimization algorithms proposed over the years are guided when attempting to solve a given optimization problem. For instance, a common practice is to reserve a validation data subset over which a measure of performance related to the task at hand is computed (e.g. accuracy in image classification). This measure is used as the objective function guiding the search of the proposed optimization algorithm. However, depending on whether this validation subset is kept fixed or shuffled, a partition bias might affect the generalization capabilities of the evolved network, specially when dealing with small datasets.

On a similar reasoning, when dealing with imbalanced datasets the standard definition of accuracy is known to be not adequate to quantify the performance of the model in the minority class, and could exacerbate further the aforementioned problems. In what refers to Problem 4 (trainable parameter optimization), this issue becomes even more serious because derivatives of the loss function are not needed any longer, hence widening the portfolio of possible objective functions. To date, there is no clear answer whether differentiable loss functions should be selected at the objective function of bio-inspired optimization algorithms used for model training, or, instead, alternative task-dependent objectives should be formulated.

- **Parameter tuning of solvers:** additional issues arise in how different solvers are compared to each other for a given optimization problem. To begin with, it is often the case that no evidence is provided about the optimality of the parameters controlling the behavior of the search algorithm itself. In the research community working in bio-inspired optimization, it is largely accepted that a good parameter tuning is crucial for ensuring fair comparisons among algorithms [198]. Since the objective function evaluation of candidates in problems related to Deep Learning is usually costly in terms of computational effort,

the parameters of the bio-inspired algorithm are often set equal to values retrieved from past works, or conforming to common practice. This unfairly biases the discussion, usually leaving unclear whether the reported performance gaps are incidental. We acknowledge that the research record noted around the usage of hyper-heuristics [122, 149] avoids to an extent this comparison bias, but the problem still prevails in most works proposing new bio-inspired methods.

- Assessing the impact of randomness on the results: another methodological aspect that has not been properly considered in most related literature is the fact that several sources of randomness can collide into a optimization problem. For instance, in Problems 1, 2 and 3 not only the search algorithm comprises a number of stochastic operators, but the training algorithm in use can also induce randomness in the obtained results. For instance, the values of trainable parameters optimized by the SGD solver depends on the composition of the mini-batches over which gradient estimates are computed. Such mini-batches comprise a number of examples from the training set, which is usually shuffled between successive epochs. This source of randomness could justify training the optimized network several times (*runs*) and aggregating the results for a more reliable fitness computation. Otherwise, this should be conceived as an additional factor motivating a proper statistical assessment of the significance of performance gaps between different optimization algorithms, adding to the randomness induced by their search operators.

Surprisingly, only a few exceptions have embraced the usage of statistical tests for this purpose, leaving most of the experiments reported in this field dubious and inconclusive. Furthermore, experiments with several datasets, tasks and optimization algorithms should embrace the methodological practices for multiple comparisons deeply rooted on the scientific community, such as critical distance plots [199] or Bayesian tests [200].

- Reproducibility of results: although this is a claim that emerges in almost any field of research, the need for reproducibility becomes particularly pressing in this area. Reasons go far beyond the verification of the contribution reported in emerging studies: the community can expedite the achievement of novel advances in the field if the software and experimental results of preceding works are shared within the community. The current status in this matter is not that concerning, as many open-source software frameworks exist with the functionalities required to develop new approaches and experiments. Unfortunately, many contributions published lately still do not provide any access to the software implementing their proposals. When given, it is also frequent to see that the source code is made ad-hoc for the problem at hand, thereby dilating the time needed for building new proposals on top of existing ones.

#### 5.4. *Software implementations of limited practical utility*

Nowadays, Deep Learning architectures scoring competitively in tasks defined over real-world datasets are usually composed by millions of trainable parameters. When addressing Problem 4 (trainable parameter optimization), the huge search space faced by the optimization algorithm under consideration requires intelligent means to exploit the relationships existing among the optimization variables. Actually, gradient based solvers adopt this strategy by the back-propagation of the gradients throughout all layers compounding the Deep Learning model, which gives rise to an implicit mechanism to exploit the correlations between the optimization variables. In this context, there is an entire research area dedicated to the large-scale global optimization, plenty of algorithmic proposals where synergies among optimization variables are exploited by assorted means. Nonetheless, many works still revolve on the naive application of standard bio-inspired solvers, neglecting such interactions between variables.

Further along this line, we note that the vibrant activity of the area is not in accordance with the ultimate goal for which bio-inspired algorithms are being proposed for training Deep Learning models.

As told by the currently available implementations in the literature, the computational cost of population-based meta-heuristics is enormous, and yields far longer training times than off-the-shelf gradient back-propagation approaches. Even if the software implementation of algorithmic proposals reported to date may have been restricted to experimental settings, a complexity analysis should have been performed to accounting for both their benefits and drawbacks, so that the community can delimit realistic boundaries for the practical utility of these advances.

### 5.5. *Metaphor-based publication series*

Last but not least, we emphasize on the claims of recent studies about the justification of new meta-heuristic algorithms just by the biological metaphor in which it is allegedly inspired [201]. In our literature review we identified several publication series in which the same optimization problem were tackled by different bio-inspired solvers within a short time period. By no means these contributions provide any scientific value for the general knowledge of the field, even if publication workflows run at different speeds.

Disregarding the reasons for these practices, results with optimization algorithms inspired by different biological behaviors and phenomena should not be shattered over different short-elapsd publications. They can provide much more valuable insights when presented and discussed together.

## 6. **Case of study I: Architecture and Hyper-parameter Optimization of Deep Learning Models with Bio-inspired Algorithms**

In this first case of study we focus on the topology and structural hyper-parameter optimization of Deep Learning models, which has been an open challenge in the last few years. In particular, this experimental study focuses on finding the best CNN architecture along with their structural hyper-parameters for solving image classification tasks. We focus the case of study on two recent frameworks:

- EvoDeep, proposed in [22] as an evolutionary algorithm specially designed to optimize both hyper-parameters and architecture of a Deep Learning model, selecting the type and size of layers compounding the evolved architecture.
- AutoKeras [25], which is another modern AutoML systems built on top of the renowned Keras Python library, that is able to automatically generate highly-optimized neural network.

Taking into account the recent activity noted in this area, we herein compare and discuss on the performance of these two consolidated frameworks for topology and hyper-parameter optimization of Deep Learning models. The comparison is made using three important measures: accuracy, time and model complexity.

To this end, in what follows we describe EvoDeep (Subsection 6.1) and AutoKeras (Subsection 6.2), present the designed experimental setup (Subsection 6.3). The obtained results are discussed in Subsection 6.4, analyzing them in terms of accuracy, computation time and model complexity.

### 6.1. *EvoDeep: Evolutionary Computation for Deep Neural Network topology and hyper-parameter tuning*

EvoDeep is a framework based on an evolutionary algorithm that follows a  $(\lambda + \mu)$  strategy, where  $\lambda$  indicates the number of new individuals produced at each generation, and  $\mu$  represents the number of individuals selected for the next generation. Each individual of the population is a network architecture with its respective hyper-parameters. The fitness for each individual is the accuracy of the neural network when solving a classification problem. At every generation, the recombination and mutation operators are applied to generate a new individual for the next generation.

EvoDeep finds increasingly better neural network architectures for CNN using elementary convolutional and fully connected layers. Specifically, EvoDeep evolves a network by using only the original data and an set of permitted layers (i.e. `Convolution2D`, `Flatten`, `MaxPooling2D`, `Reshape`, `Dense` and `Dropout` as per Keras notation). The number of layers represents the depth of the evolved network, and its search range is set to  $[3, 20]$  with a step size of 1 layer. Moreover, the training hyper-parameters can be also specified in its configuration. These parameters are as follows:

- Optimizer, chosen among Adam, SGD, Rmsprop, Adagrad, Adamax or Nadam.
- Number of epochs: an integer value in the range  $[2, 20]$ , with a step size of 2 epochs.
- Batch size: this value has a range of  $[100, 5000]$  with a step size of 100 examples.

Another characteristic of EvoDeep is the fact that it requires the data in a specific manner. A two-dimensional matrix is the input for the algorithm. The number of row matches the number of examples of the database and the number of columns is the size of the images. With size of the image we mean the product of width, height and channels (three channels for RGB and one for grayscale images). Unfortunately, EvoDeep is more oriented towards optimizing grayscale rather than RGB images.

EvoDeep allows the user to specify the parameters of its evolutionary strategy. Table 3 shows the parameters and values used in this case of study.

Table 3: Parameter configuration set for the EvoDeep framework.

Parameter	$\lambda$	$\mu$	$cxpb, p_c$	$mutpb, p_m$	$newpb$	$ngen, n_{gen}$
<b>Description</b>	Number of newly produced individuals per generation	Number of selected individuals for the next population	Crossover probability	Mutation probability	Probability of adding a new layer to the network	Number of generations
<b>Value</b>	10	5	0.5	0.5	0.5	20

## 6.2. AutoKeras: an AutoML reference

AutoKeras is a software that also allows to finding the best neural network architecture for a given data set and task [25], offering several search engines for this purpose. In our study, the version of AutoKeras used is 1.0.2, which features the following search methods:

- Random: it performs a random search of the models in relation to their depth and layer type.
- Greedy: it groups the parameters into several categories. For each category, the tuner uses a greedy strategy to generate new values for the hyper-parameters and to generate new values for one of the categories of hyper-parameters. It uses the best trial so far for the rest of hyper-parameter values.
- Hyperband: departing from a given model, this bandit-based algorithm searches the best hyper-parameter values for this model by running several random configurations for a scheduled number of iterations, using earlier results to retain good candidate configurations that are evaluated for longer runs.
- Bayesian optimization: the search space is explored using morphing. This optimization method is based on the three basis of Bayesian optimization: update, generation and observation.
- Task-specific: AutoKeras tries with a task-specific tuner of the model, and evaluates the most commonly used models for the task at hand, which in our case is image recognition. This is the default configuration. AutoKeras has two trial blocks: Vanilla and ResNet. The best model is taken from one of these two.

### 6.3. Experimental setup

In our study we have chosen 4 diverse and representative data sets according to their complexity: Horses or Humans (HORSEHUMAN [202]); the `vangogh2photo` dataset from the so-called `cycle_gan` repository (VANGOGH [203]), MNIST [204] and CIFAR-10 [205]. On the one hand, Horses or Humans and Van Gogh or Photo have two classes, but Van Gogh or Photo is unbalanced (only 400 images belonging to the minority class). On the other hand CIFAR-10 and MNIST are more complex databases in terms of number of examples and classes: both databases have 10 different classes. Moreover, in our experiment we have considered both color (RGB) and grayscale versions of these datasets so as to assess the influence of the color space in the complexity and performance of the models evolved by the compared frameworks.

Table 4: Utilized datasets in Case of Study I.

Dataset	Shape	# classes	# Instances (train/test)	Characteristics
HORSEHUMAN	$300 \times 300 \times 3$	2	1,027 / 256	Balanced dataset, binary classification, RGB
HORSEHUMAN-G	$300 \times 300 \times 1$	2	1,027 / 256	Balanced dataset, binary classification, grayscale
VANGOGH	$256 \times 256 \times 3$	2	6,687 / 1,151	Imbalanced dataset, binary classification, RGB
VANGOGH-G	$256 \times 256 \times 1$	2	6,687 / 1,151	Imbalanced dataset, binary classification, grayscale
CIFAR-10	$32 \times 32 \times 3$	10	50,000 / 10,000	Balanced dataset, multi-class classification, RGB
CIFAR-10-G	$32 \times 32 \times 1$	10	50,000 / 10,000	Balanced dataset, multi-class classification, grayscale
MNIST	$28 \times 28 \times 1$	10	60,000 / 10,000	Balanced dataset, multi-class classification, grayscale

In order to account for the statistical behavior of the frameworks under comparison, we have carried out 5 independent runs of AutoKeras, EvoDeep and random models, from which we have chosen the best model in terms of its accuracy in validation. Random models are obtained using EvoDeep without the evolutionary process. The number of tested individuals is the number of evaluations that EvoDeep would make should the search be done via its evolutionary strategy. The training phase has been done by splitting the training set as follows: 80% for training and 20% for validation. After that, the model is trained again with the whole training set, and evaluated on the test set.

### 6.4. Results and discussion

Table 5 shows the obtained accuracy scores for models evolved with EvoDeep, AutoKeras and random search over the color and grayscale datasets under consideration. The best results for each dataset are highlighted in bold. On the one hand, it is straightforward to observe that EvoDeep outperforms random models over both train and test datasets. However, AutoKeras still offers a better performance than EvoDeep. AutoKeras features the best test results in each dataset. To sum up, the results of EvoDeep are close to AutoKeras, which it is one of the best AutoML tools in this area. In Vangogh or Photo and MNIST the difference between them is less than 1% in test and 2% in Horses or Humans. Therefore, EvoDeep shows a great performance on these databases.

When shifting the scope towards the results obtained for the grayscale databases, similar conclusions can be drawn. In relation to the HORSEHUMAN dataset, the results issued by EvoDeep are similar to the previous ones in terms of accuracy over test, whereas AutoKeras improves its performance by 3%. The results for CIFAR-10-G are worse when compared to the color ones. The exception in this trend is VANGOGH-G, due to the bad results obtained by the random models and EvoDeep. The accuracy in test decreases approximately 17% in EvoDeep and 27% in random models. By contrast, the results obtained for this dataset by AutoKeras are similar to the ones obtained for its colored counterpart.

Table 5: Results in terms of accuracy for color datasets and algorithms (in %).

Color datasets	HORSEHUMAN		VANGOGH		CIFAR-10			
	Train	Test	Train	Test	Train	Test		
<b>Random</b>	89.19	89.06	98.10	92.96	73.06	53.63		
<b>EvoDeep</b>	97.07	89.40	100.00	98.87	99.99	60.55		
<b>AutoKeras</b>	<b>100.00</b>	<b>91.40</b>	<b>100.00</b>	<b>99.48</b>	<b>95.28</b>	<b>73.26</b>		
Grayscale datasets	HORSEHUMAN-G		VANGOGH-G		CIFAR-10-G		MNIST	
	Train	Test	Train	Test	Train	Test	Train	Test
<b>Random</b>	94.06	89.84	94.02	65.24	73.44	49.35	99.87	98.40
<b>EvoDeep</b>	100.00	89.84	96.66	81.58	85.52	56.79	100.00	98.69
<b>AutoKeras</b>	<b>100.00</b>	<b>94.53</b>	<b>100.00</b>	<b>99.30</b>	<b>94.42</b>	<b>71.67</b>	<b>99.99</b>	<b>99.40</b>

After the comparison of EvoDeep, AutoKeras and random models in terms of accuracy, we now examine in Table 6 the execution time of the previous results, in minutes. A first inspection of the results in this table reveals that AutoKeras has the best performance with lower times than EvoDeep. We note that EvoDeep stops if no improvement is made over 5 consecutive generations. This is what occurs in MNIST: EvoDeep needs more computation time than AutoKeras and random models, but its accuracy results lie in between. This fact makes EvoDeep a reliable software because, even though it needs more computation time, the quality of the models in terms of predictive accuracy are close to those of AutoKeras in most cases. When it comes to grayscale datasets, in general the computation time is lower than the time taken by the experiments with color datasets. Runtimes of AutoKeras are almost the same except for HORSEHUMAN-G, in which the time is approximately 2.6 times faster than that of HORSEHUMAN. EvoDeep and random models require less computation time in all the cases when compared to the colored ones. For example, VANGOGH-G is around 2.25 times faster than the previous experiments.

Table 6: Results in terms of time (minutes) for all datasets and frameworks.

Color	HORSEHUMAN	VANGOGH	CIFAR-10-G	
<b>Random</b>	<b>5.5</b>	<b>11</b>	<b>92</b>	
<b>EvoDeep</b>	22.0	50	322	
<b>AutoKeras</b>	13.5	14	110	
Grayscale	HORSEHUMAN-G	VANGOGH-G	CIFAR-10	MNIST
<b>Random</b>	<b>3.08</b>	<b>7.14</b>	<b>43.44</b>	<b>46</b>
<b>EvoDeep</b>	13.17	22.25	295.99	230
<b>AutoKeras</b>	5.19	13.98	109.98	262

We now proceed in our analysis with Table 7, where we compare random, EvoDeep and AutoKeras in terms of the complexity of the models produced over the search (in terms of the number of layers and the number of trainable parameters of the best model). As proven by the results included in this table, random model produces very similar network architectures across the colored datasets: they all comprise 5 layers with varying size (between  $10^3$  and  $2 \cdot 10^6$  trainable parameters). EvoDeep has better results than AutoKeras in terms of number of layers for 3 datasets, and in terms of the number of trainable parameters for 2 datasets (HORSEHUMAN and MNIST). The best model for HORSEHUMAN is achieved by EvoDeep,

comprising a simple neural architecture with 3 layers and approximately  $1.9 \cdot 10^6$  parameters. AutoKeras’ model for the VANGOGH dataset has 9 layers and fewer trainable parameters in comparison to the models produced by the other frameworks. EvoDeep produces a model with only 5 layers for CIFAR-10, but AutoKeras’ model has less parameters. Finally, for the MNIST dataset, EvoDeep discovers a better model than AutoKeras in terms of model complexity: fewer layers and parameters. To sum up, in terms of model complexity EvoDeep can be declared to perform better than AutoKeras, at least in some of the considered colored datasets.

Table 7: Results in terms of model complexity for all datasets and algorithms.

Color datasets	HORSEHUMAN		VANGOGH		CIFAR-10	
	Layers	Parameters	Layers	Parameters	Layers	Parameters
<b>Random</b>	5	<b>179,922</b>	<b>5</b>	1,158,402	<b>5</b>	2,270,100
<b>EvoDeep</b>	<b>3</b>	1,895,012	9	18,756,012	<b>5</b>	10,821,230
<b>AutoKeras</b>	194	23,566,856	10	<b>31,944</b>	10	<b>144,849</b>

Grayscale datasets	HORSEHUMAN-G		VANGOGH-G		CIFAR-10-G		MNIST	
	Layers	Parameters	Layers	Parameters	Layers	Parameters	Layers	Parameters
<b>Random</b>	<b>4</b>	853,322	<b>3</b>	318,372	<b>5</b>	239,650	<b>5</b>	<b>211,245</b>
<b>EvoDeep</b>	7	<b>129,182</b>	6	783,352	12	1,883,220	8	5,409,980
<b>AutoKeras</b>	194	23,560,580	10	<b>31,364</b>	10	<b>144,269</b>	194	23,579,021

When analyzing the complexity of models discovered for the grayscale datasets, the results in Table 7 unveil a huge improvement in the number of parameters in random models and EvoDeep. Results by AutoKeras remain very similar to those for the colored datasets, with minimal differences in terms of the number of parameters. If we observe the number of layers, random models have almost the same number of layers as for the colored ones. However, more remarkable changes are noticed for EvoDeep. The models discovered for the HORSEHUMAN-G dataset has more layers, but significantly less parameters. The same statement can be said for CIFAR-10-G. The results for the VANGOGH-G dataset are the exception: the best model has fewer layers and fewer parameters. These observations support the aforementioned claims on the complexity of models produced by AutoKeras with respect to EvoDeep and Random.

The previous experiments have shown some differences between the color and the grayscale databases. In fact, when focusing the analysis on the best model found by EvoDeep across all datasets, remarkable differences arise between the test accuracy and the complexity of such models corresponding to color and grayscale datasets. As summarized in Table 8, accuracy results are similar in both cases. The accuracy achieved by the best EvoDeep model over CIFAR-10 is higher than that of its grayscale version (CIFAR-10-G). The case with the largest difference in terms of accuracy is VANGOGH, which has a 17% gap. In terms of complexity, the best model of each grayscale database has fewer neurons than the corresponding colored dataset. In HORSEHUMAN-G, the model has approximately 14.69 times fewer neurons than that of HORSEHUMAN. For VANGOGH and CIFAR-10, this ratio gets close to 24 times and 5.75 times, respectively.

Table 8: Results in terms of accuracy and complexity for the best model encountered by EvoDeep for both color and grayscale datasets.

	HORSEHUMAN	HORSEHUMAN-G	VANGOGH	VANGOGH-G	CIFAR-10	CIFAR-10-G	MNIST
<b>Test accuracy (%)</b>	89.45	89.94	98.87	81.58	60.55	56.79	98.69
<b># of neurons</b>	1,895	129	18,756	783	10,821	1,883	5,409

We summarize now the main conclusions drawn from this discussion, leaving a further elaboration on the general lessons learned in regards to topology and hyper-parameter optimization for Section 8:

- EvoDeep has similar results to AutoKeras in the color databases HORSEHUMAN, VANGOGH and CIFAR-10. Nevertheless, the computation time that EvoDeep requires is much higher than the one taken AutoKeras during its search. If we consider the grayscale datasets, the difference between AutoKeras and EvoDeep increases. In particular, accuracy gaps over the VANGOGH dataset is particularly large: 81.58% in VANGOGH-G and 98.87% in VANGOGH. In terms of accuracy, EvoDeep performs better with the colored databases.
- If we take a closer look at the results in terms of accuracy and model complexity, EvoDeep needs less computation time in the grayscale datasets. Furthermore, this statement also holds in terms of model complexity. Although there are some models that comprise more layers when dealing with grayscale datasets, the number of total trainable parameters for all the models is much lower. All these facts contribute to a better overall performance of EvoDeep with grayscale databases.

To summarize, although AutoKeras has better performance in terms of accuracy in all datasets, EvoDeep gives competitive results and requires less computation time for simple datasets (in our case, MNIST). For several datasets, EvoDeep produces models with fewer parameters, while AutoKeras yield very complex models that may be unsuitable for application scenarios with stringent memory restrictions. In other datasets, the simplicity of the layers supported by EvoDeep enforces a higher number of neurons (and parameters) than AutoKeras for the same performance level.

The empirical results reported in this first case of study underscore the potential of Evolutionary Algorithms for the topological and hyper-parameter optimization of CNN networks, suggesting several possible improvements to frameworks appearing in the literature in forthcoming years. First, frameworks should incorporate sophisticated layers at their core, so that they become eligible for the evolutionary algorithm in use and ultimately lead to a reduced overall complexity of the optimized models. The overly complex models encountered by EvoDeep in some of the considered image classification datasets is a clear evidence that, for the sake of fair comparisons, all frameworks should ensure that the search spaces explored by their optimization engines are comparable to each other as well. Another possible improvement is to formulate topological and hyper-parameter optimization as a multi-objective problem, embracing as conflicting objectives the accuracy of the model, and a measure of its complexity (layers, parameters, computation time, etc). We believe that such an output could allow the community to examine the behavior of new frameworks and solvers in regards to the performance and complexity of their produced solutions, making their output more flexible to implement the discovered models by taking into account both objectives. We will later revolve on these research directions on Sections 8 and 9.

## 7. Case of Study II: Training Deep Learning Models with Bio-inspired Algorithms

This second case of study aims to shed light on the performance of bio-inspired optimization algorithms when applied to model training (*trainable parameter optimization* as per our nomenclature in this survey). The ultimate goal is to check whether bio-inspired algorithms can be a competitive alternative to gradient-based methods when undertaking image classification tasks over well-known datasets, ensuring that Deep Learning architectures of realistic complexity are in use. To this end, we design an experimental setup to provide an informed response to the following research questions (RQ):

- RQ1: Should bio-inspired solvers exploit the layered structure of Deep Learning models during the search?
- RQ2: Do bio-inspired solvers perform competitively with respect to gradient-based solver for trainable parameter optimization, in terms of predictive accuracy and computational efficiency?



In the remainder of this second case of study, Subsection 7.1 introduces the evolutionary algorithm selected for the experiments, underscoring several changes made to its original definition to make it better suited to the optimization of trainable parameters. Subsection 7.2 provides further details on the experimental setup, including the Deep Learning models and datasets under consideration. Subsections 7.3 and 7.4 discuss in depth on the results obtained from the experiments, with a focus on RQ1 and RQ2, respectively.

### 7.1. SHADE-ILS: A reference evolutionary algorithm for large-scale global optimization

Given the large number of variables to be optimized, we select SHADE with Iterative Local Search (SHADE-ILS) as the evolutionary algorithm selected for our experiments. SHADE-ILS is a renowned large-scale global optimization algorithm that has won recent international competitions in the field [206]. In particular, SHADE-ILS resorts to the global exploration capability of an adaptive variant of the Differential Evolution algorithm (SHADE), which is helped by two local search methods – namely, a limited-memory version of the BroydenFletcherGoldfarbShanno (BFGS) algorithm, and multiple trajectory search – that improve the candidate solutions encountered during the search. At every iteration, SHADE applies evolutionary operators on a population of individuals, followed by the application of one of the local search methods to the best individual of the population. The selection of which local search to apply is driven by the best expected relative improvement of each of the considered local search operator, which is given by the results produced by each local search alternative during recent generations. Moreover, SHADE-ILS incorporates a restart mechanism to avoid stagnation. These key algorithmic aspects and the excellent results scored in benchmarks and competitions make SHADE-ILS one of the baseline algorithms for large-scale global optimization problems, just like the one addressed in this second case of study.

Several changes have been done to the original SHADE-ILS algorithm to make it better suited to the optimization of the trainable parameters of a Deep Learning model. The objective function to be minimized over the search is the same than that used by gradient back-propagation approaches for the same model, namely, binary cross-entropy for binary classification and categorical cross-entropy for multi-class classification. Both are measured for all examples belonging to the training set. Secondly, we take into account the layer-wise structure of Deep Learning models by devising several scheduling strategies. Each strategy establishes the criteria, order and number of generations for which the optimization variables (trainable parameters) belonging to each layer are evolved via the global search operators and local search techniques of SHADE-ILS. The design of the above strategies is motivated by RQ1, where the focus is placed on whether the layered structure of Deep Learning models should be exploited anyhow by the bio-inspired solver.

Specifically, the strategies considered in the simulations discussed later are as follows (see Figure 5 for a visual explanation):

- FULL-SHADE-ILS: all trainable parameters are optimized jointly, without considering the layered structure of the model to be trained. This is actually the strategy followed by most contributions to the literature dealing with the application of bio-inspired optimization techniques for Deep Learning models. As we will later show, this strategy only works for network architectures of relatively limited size.
- DOWN-SHADE-ILS: trainable parameters are optimized starting from those belonging to the first layer of the network. Such parameters are evolved via SHADE-ILS for a certain number of generations, keeping the values of the remaining parameters fixed to their Glorot-based initialized values. The optimization schedule is repeated in order from the first to the last layer of the network for a maximum number of epochs.
- UP-SHADE-ILS: this schedule is similar to the previous DOWN-SHADE-ILS strategy, but departing from the last layer of the network, and proceeding upwards until the first layer of the network.

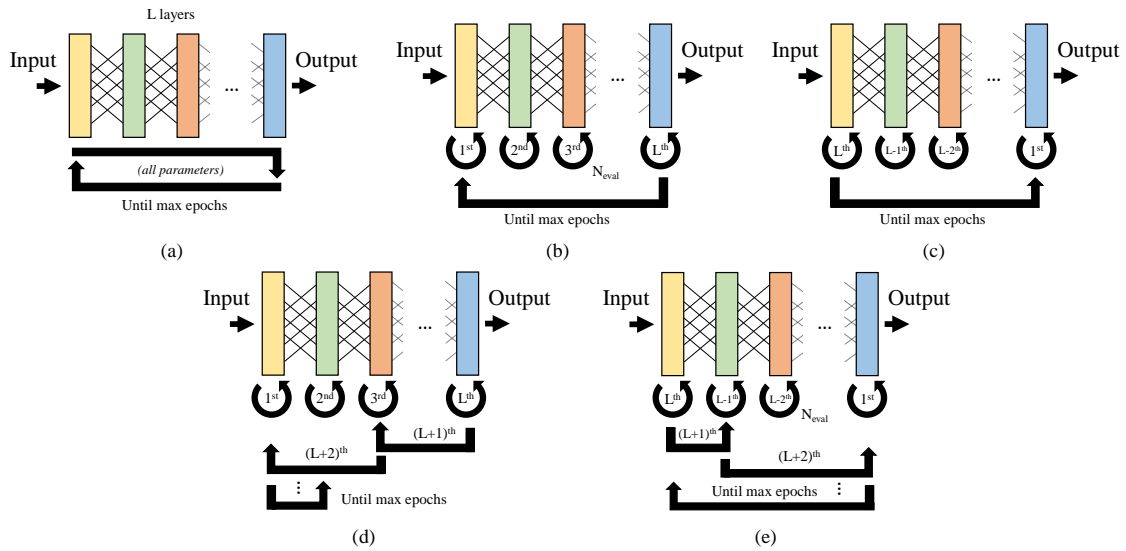


Figure 5: Diagram showing the different scheduling strategies proposed for optimizing Deep Learning models with the SHADE-ILS algorithm: (a) FULL-SHADE-ILS; (b) DOWN-SHADE-ILS; (c) UP-SHADE-ILS; (d) A-DOWN-SHADE-ILS; (e) A-UP-SHADE-ILS.

- **A-DOWN-SHADE-ILS**: an automated variant in which after a first iteration of the DOWN-SHADE-ILS optimization strategy, a relative improvement ratio of the network predictive accuracy is computed for every layer optimization step. This ratio is used to select which layer to optimize in subsequent epochs, so that layers whose optimization yielded larger improvements in the last epoch are more likely to be selected for optimization.
- **A-UP-SHADE-ILS**: this last strategy is similar to A-DOWN-SHADE-ILS, the difference being that the first layer-wise application of SHADE-ILS is done from the last to the first layer of the network (namely, under the DOWN-SHADE-ILS strategy). Once this initial stage is completed, SHADE-ILS proceeds analogously to the previous strategy, automatically selecting the layer with most potential margin of improvement.

When addressing RQ2, it is important to stress on the complexity of ensuring a fair comparison between gradient back-propagation solvers and bio-inspired optimization algorithms. The main reason is their essentially different search behavior. On one hand, gradient back-propagation approaches maintain a single solution to the problem, which is enhanced over epochs by exploit the mathematical relationships between the optimization variables (trainable parameters) through the tailored computation of their gradients. However, the search operator of gradient back-propagation techniques is simple, yet effective (gradients are personalized for every single variable) and efficient (gradient computations are highly parallelizable). By contrast, most bio-inspired algorithms rely on a population of individuals, which are evolved jointly by means of a series of search operators, so that the best individuals survive between generations. In summary, comparisons between both approaches should be fair not only in terms of accuracy performance, but also in terms of computational complexity.

For ensuring fairness in these terms, a straightforward decision is to define how *epoch* and *generation* relate to each other. First of all, it is clear that both concepts indicate when a full update of the network's parameters is complete: in gradient back-propagation, an epoch implies the application of a number of gradient updates to the whole network parameters. The number of updates depends on the size of the training set and the chosen batch size. Given that an epoch is defined as the optimization of all parameters

composing the model, in SHADE-ILS we will establish that an *epoch* corresponds to the optimization of all layers of the network under any of the strategies described previously. If we assume  $N_{eval}$  evaluations of the network per layer and  $L$  layers, an epoch for SHADE-ILS will comprise  $L \cdot N_{eval}$  total evaluations of the network per epoch. Each evaluation of the network involves predicting the entire  $N_{train}$ -sized training set and computing the loss function. As a result, in general the number of evaluated training instances per epoch differs between SHADE-ILS ( $N_{train} \cdot N_{eval} \cdot L$ ) and gradient back-propagation ( $N_{train}$ ). Nevertheless, we proceed forward with the experiments disregarding this issue, and analyze whether SHADE-ILS, even if endowed with more computational budget per epoch, can beat the accuracy of networks optimized via Adam, one of the most renowned gradient back-propagation solvers.

### 7.2. Experimental setup

The above two questions are tackled by considering 6 datasets for image classification: Hand gesture recognition (HANDS, [207]), Blood Cells classification dataset (BCCD, [208]), MNIST [204], Fashion MNIST (F-MNIST) [209], GTSRB [210] and CIFAR-10 [205]. Details of these datasets are given in Table 9. Given the amount of computational resources required to complete the experiments, we consider a subset of the examples for each dataset. For the same reason and except for the WBC dataset, images have been converted to grayscale to reduce the number of channels of the input image.

Table 9: Datasets and models utilized for the second case of study.  $C2D_{N,x \times y}$  denotes a convolutional layer with  $N$  filters of  $n$  rows and  $m$  cols;  $D_N$  represents a fully-connected (*dense*) layer with  $N$  output neurons;  $\boxplus$  is a  $2 \times 2$  max pooling layer;  $\odot$  is a  $2 \times 2$  average pooling layer; and  $Drop_p$  is a dropout layer with rate  $p$ . Layers enclosed within  $(\cdot)^L$  are concatenated  $L$  times.

Dataset	Shape	# classes	# Instances ( $N_{train}/N_{test}$ )	# trainable parameters	Network topology and structural hyper-parameters
HANDS	$30 \times 40 \times 1$	10	10,000 / 10,000	3,854	$C2D_{8,4 \times 4} - \boxplus - (C2D_{16,2 \times 2} - \boxplus)^2 - D_{20} - D_{10}$
BCCD	$30 \times 40 \times 3$	2	17,000 / 5,416	9,065	$C2D_{30,3 \times 3} - \boxplus - (C2D_{16,3 \times 3} - \boxplus)^2 - D_{16} - Drop_{0.7} - D_1$
MNIST	$28 \times 28 \times 1$	10	10,000 / 5,000	19,063	$C2D_{28,3 \times 3} - \boxplus - C2D_{14,3 \times 3} - \boxplus - C2D_{7,2 \times 2} - \boxplus - D_{128} - Drop_{0.2} - D_{80} - Drop_{0.3} - D_{10}$
F-MNIST	$28 \times 28 \times 1$	10	10,000 / 5,000	36,188	$C2D_{64,4 \times 4} - Drop_{0.25} - \odot - C2D_{16,4 \times 4} - Drop_{0.25} - \odot - Drop_{0.15} - D_{70} - D_{10}$
GTSRB	$32 \times 32 \times 1$	43	20,000 / 10,000	83,999	$C2D_{6,3 \times 3} - \odot - C2D_{16,3 \times 3} - \odot - D_{120} - D_{84} - D_{43}$
CIFAR-10-G	$32 \times 32 \times 1$	10	10,000 / 5,000	1,658,570	$C2D_{32,3 \times 3} - Drop_{0.1} - C2D_{64,5 \times 5} - Drop_{0.2} - D_{128} - Drop_{0.3} - D_{10}$

For each of these datasets a fixed Deep Learning architecture is considered, featuring a realistic level of complexity (given by its number of trainable parameters), and rendering a good prediction performance when trained via gradient back-propagation. The layer type and structural hyper-parameters of every layer compounding such models are also specified in the table. Thus, the aim of this section is not to evolve very precise, state-of-the-art networks, but to assess the limitations faced by Evolutionary Algorithms when used for training these deep neural networks. Table 10 summarizes the training hyper-parameters utilized for all image classification tasks under study.

Table 10: Training hyper-parameters used in our experiments.

Dataset	Adam		SHADE-ILS		Initializer	Epochs
	Batch size	Learning rate	Population size	$N_{eval}$		
HANDS	128	0.01	10	200	Glorot	20
BCCD	64	0.02	10	200	Glorot	20
MNIST	512	0.01	10	200	Glorot	20
F-MNIST	512	0.01	10	200	Glorot	20
GTSRB	64	0.02	10	200	Glorot	22
CIFAR-10-G	256	0.01	10	200	Glorot	30

### 7.3. Addressing RQ1: On the importance of the layered neural structure in the design of SHADE-ILS

Once the experimental setup has been described, we start our discussion by addressing RQ1, namely, a quantitative analysis of the impact and viability of exploiting the layered structure of Deep Learning

models by SHADE-ILS, similarly to what gradient-based solvers do when back-propagating the gradients. To this end, we evaluate the aforementioned scheduling strategies devised for SHADE-ILS over the datasets and Deep Learning models under consideration, and compare its performance to that of a naive application of SHADE-ILS that does not take into account any structure of the problem.

Table 11 summarizes the results obtained in regards to RQ1. Specifically, the average loss and accuracy measured over train and test subsets are reported for every dataset and SHADE-ILS schedule. Values are averaged over 5 independent runs of every (dataset,schedule) combination. In addition, results corresponding to the Adam gradient-based solver are also included as a reference. The best results among those yielded by SHADE-ILS are highlighted in bold.

Table 11: Average accuracy/loss (over 5 independent runs) corresponding to different schedules of the SHADE-ILS algorithm over the datasets under consideration. Results corresponding to the Adam gradient-based solver are also included as a reference.

	Adam		FULL SHADE-ILS		DOWN SHADE-ILS		UP SHADE-ILS		A-DOWN SHADE-ILS		A-UP SHADE-ILS	
	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test
	Acc	Acc	Acc	Acc	Acc	Acc	Acc	Acc	Acc	Acc	Acc	Acc
	Loss	Loss	Loss	Loss	Loss	Loss	Loss	Loss	Loss	Loss	Loss	Loss
HANDS	0.9983	0.9932	<b>0.9708</b>	<b>0.9553</b>	0.7038	0.6910	0.7150	0.7072	0.6492	0.6470	0.5253	0.5189
	0.0082	0.0296	<b>0.1082</b>	<b>0.1589</b>	2.6608	2.6894	0.8194	0.8705	3.5419	3.5750	1.3179	1.3533
BCCD	0.9611	0.9815	<b>0.8558</b>	<b>0.8489</b>	0.7499	0.7400	0.7527	0.7439	0.7852	0.7768	0.8212	0.8087
	0.1314	0.0566	<b>0.3321</b>	<b>0.3399</b>	0.5306	0.5377	0.5227	0.5322	0.4861	0.4962	0.4149	0.4301
MNIST	0.9480	0.9534	0.9524	0.9342	0.9747	0.9504	0.9748	0.9490	0.9719	0.9452	<b>0.9772</b>	<b>0.9508</b>
	0.1672	0.1534	0.1590	0.2326	0.0851	0.1685	0.0856	0.1798	0.0938	0.1828	<b>0.0764</b>	<b>0.1627</b>
F-MNIST	0.9636	0.9672	0.9265	0.9196	0.9489	0.9384	0.9422	0.9311	<b>0.9561</b>	<b>0.9447</b>	0.9435	0.9336
	0.1190	0.1100	0.2539	0.2940	0.1773	0.2157	0.2077	0.2407	<b>0.1576</b>	<b>0.1966</b>	0.2012	0.2317
GTSRB	0.7437	0.7046	0.2329	0.2355	0.3636	0.3473	0.3596	0.3504	<b>0.3956</b>	<b>0.3868</b>	0.3204	0.3133
	0.6957	0.7807	3.0994	3.1196	2.4532	2.4931	2.4304	2.4672	<b>2.2916</b>	<b>2.3323</b>	2.6405	2.6622
CIFAR-10-G	0.8347	0.6542	0.2628	0.2602	0.3696	0.3612	0.3708	0.3642	<b>0.3806</b>	<b>0.3790</b>	0.3765	0.3681
	0.4519	1.2418	2.0952	2.1087	1.7878	1.8023	1.7737	1.7867	<b>1.7483</b>	<b>1.7619</b>	1.7612	1.7787

Several interesting observations can be made after inspecting the above table. To begin with, we see that for relatively small-sized Deep Learning models (i.e. those for the HANDS and BCCD datasets), FULL-SHADE-ILS suffices for obtaining good scores, even superior than those rendered by its scheduled SHADE-ILS counterparts. This goes in line with our examination of the related literature, in which many contributions deal with model training using naive bio-inspired solvers, without taking into account the structure of the network. The above results confirm that when the number of network parameters is low, a powerful optimization algorithm can effectively (albeit not efficiently) find optimal values for the task at hand.

However, when increasing the complexity of the network, the trend changes, and the exploitation of the layered structure of the Deep Learning model becomes essential to maintain a good performance. This is specially remarkable in the GTSRB and CIFAR-10-G datasets, in which the accuracy values of FULL-SHADE-ILS degrade severely with respect to those attained by A-DOWN-SHADE-ILS. For networks of moderate size (MNIST, F-MNIST), accuracy differences between the scheduled and non-scheduled versions of SHADE-ILS become negligible. Nevertheless, in terms of loss metric values the difference results to be larger, showing evidence that SHADE-ILS performs better in terms of optimized loss when endowed with the automated schedule mechanism.

When comparing the accuracy and loss values measured over train and test, a quick glimpse at the table confirms that in general, the trained models do not overfit excessively. We pause briefly at the

case of the MNIST dataset, arguably one of the most utilized databases in this field. If we compare the average loss values achieved by the Adam optimizer (0.1672 over train, 0.1534 over test) to those of A-UP-SHADE-ILS (0.0764 over train, 0.1627 in test), one can state that A-UP-SHADE-ILS achieves lower loss values than Adam, thereby concluding that this scheduled SHADE-ILS variant is a better optimization algorithm for model training than Adam. However, we must bear in mind that the final goal of predictive modeling is to provide models that *generalize nicely*, namely, models that perform as expected when predicting unseen data instances. When placing our attention in the losses measured over the test set, networks evolved by A-UP-SHADE-ILS for the MNIST dataset seems to overfit, thereby providing lower accuracy scores than expected. A similar conclusion can be drawn in other datasets (e.g. F-MNIST), yet at a lower extent than MNIST. This leads to an interesting insight on the influence of overfitting that we will elaborate in depth in Section 8.

We conclude the discussion on this first set of results by emphasizing on the low scores obtained by all SHADE-ILS variants when the complexity of the network is very high (models corresponding to the GTSRB and CIFAR-10-G datasets). In these cases the large gaps to the accuracy and loss scores of the network optimized by the Adam solver indicate without doubts that this meta-heuristic algorithm is of no practical use for this level of complexity. Given that SHADE-ILS is specially tailored to deal with high-dimensional optimization problems, it is fair to conclude that Evolutionary Computation and Swarm Intelligence methods are still far from being a realistic replacement for gradient-based solver. Instead, these empirical findings should drive the interest of the community towards hybridizing bio-inspired algorithms with gradient-based information and/or solvers. We will later revolve on this postulated research line.

#### 7.4. Addressing RQ2: comparing bio-inspired optimization algorithms to gradient-based solvers

Experiments discussed in the previous subsection have concentrated on the performance comparison between different layer-wise optimization schedules of the SHADE-ILS solver. In our analysis of the results shown in Table 11 we also highlighted the large gap between the gradient-based Adam solver and the best performing SHADE-ILS schedule, specially for Deep Learning models of moderate-to-high levels of complexity. Although this remark provides a partial answer to RQ2, in this second part of the study we delve into the convergence of the training process when undertaken via SHADE-ILS and Adam, aiming to discern the reasons for these identified performance gaps.

This being said, we focus on the results corresponding to MNIST and F-MNIST, which are illustrative of the conclusions that can be drawn from the overall set of performed experiments. The dual plots in Figures 6.a to 6.d depict the loss/accuracy convergence plots measured over train and test subsets corresponding to Adam and the scheduled A-UP-SHADE-ILS variant of the SHADE-ILS algorithm. Plotted lines depict the average loss/accuracy over epochs computed over 5 experiments, whereas the shaded overlay areas denote their standard deviation. Net loss values are indicated in the left axis of every plot, whereas accuracy score values are indicated in the right axis.

Our discussion departs from Figures 6.a and 6.b, corresponding to the convergence plots over the MNIST dataset over train and test subsets, respectively. It is straightforward to observe that when measured over the train dataset (Figure 6.a, both loss and accuracy scores of the A-UP-SHADE-ILS approach are better than those rendered by Adam over all epochs. This fact is conclusive in regards to the comparable or superior performance of SHADE-ILS when optimizing the trainable parameters of relatively small-sized networks. However, when shifting the focus on the test set (which reflects the generalization capability the evolved Deep Learning models), the curves plotted in Figure 6 reveal that A-UP-SHADE-ILS yields trainable parameter values that are slightly overfitted, thus generalizing worse than the model optimized via the Adam solver. This observation suggests that the apparently worse performance of Adam as an optimization algorithm is actually an advantage (a sort of *implicitly regularization* mechanism) that yields models of better generalization properties.

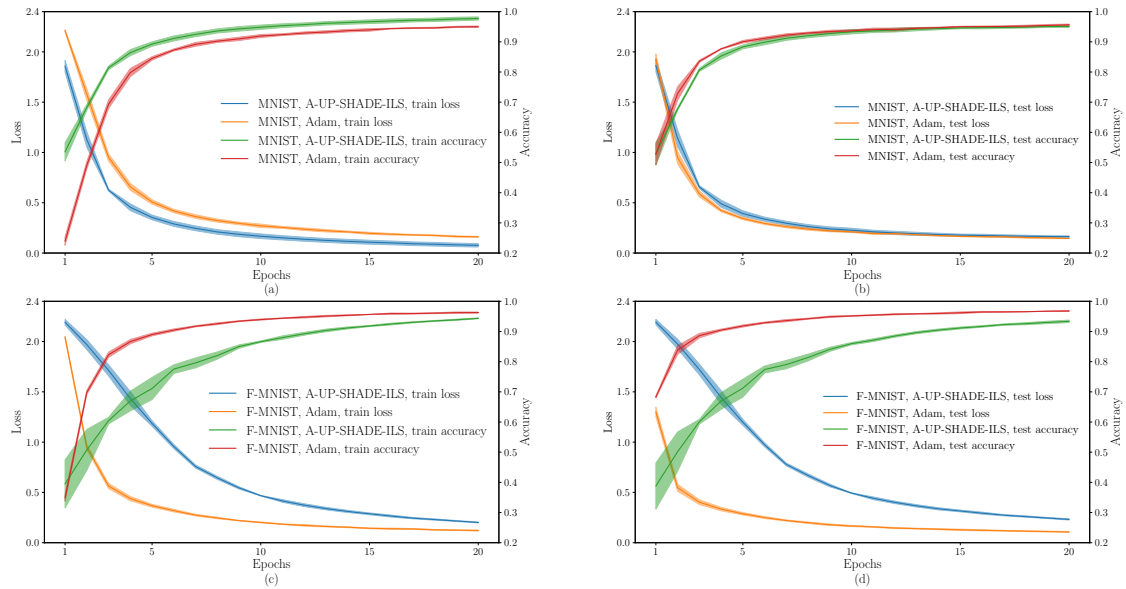


Figure 6: Accuracy and loss Convergence plots of Adam and A-UP-SHADE-ILS corresponding to (a) MNIST, measured over train set; (b) MNIST, measured over test set; (c) F-MNIST, measured over train set; (d) F-MNIST, measured over test set.

When increasing the complexity of the network to be evolved, Figures 6.c and 6.d illustrate a rather different behavior of the convergence plots. At this point we recall that the model selected to deal with the F-MNIST image classification problem has 36,188 trainable parameters, almost twice the complexity of the model designed for the MNIST dataset (19,063 trainable parameters). The convergence curves in these plots show that although the accuracy and loss values of the networks evolved with both solvers get close to each other after all epochs are completed, A-UP-SHADE-ILS perform steadily worse than Adam over all intermediate epochs. In light of the results for the rest of datasets with models of larger complexity (Table 11), the case of the F-MNIST dataset must be understood as an inflection point, beyond which SHADE-ILS fails to perform competitively with respect to gradient-based methods. Furthermore, this worse performance occurs even if SHADE-ILS is allowed to execute more network evaluations per epoch.

These experiments and our conclusions drawn therefrom reinforce even further our belief that for the time being, current bio-inspired optimization algorithms do not constitute a feasible replacement for gradient-based solvers. In the next section we collect and summarize the lessons learned through our literature analysis and experiments, and prescribe several good practices and recommendations that should be followed to achieve significant advances in evolutionary Deep Learning.

## 8. Learned Lessons, Good Practices and Recommendations

As it follows from our experiments and the performed literature analysis, lights and shadows still remain in the application of Evolutionary Computation and Swarm Intelligence algorithms to the diverse optimization problems arising from Deep Learning.

We now summarize the lessons learned throughout this study, classifying them depending on the optimization domain in which each lessons is most recurrently noted:

- General lessons and recommended practices that hold for all Deep Learning optimization problems (Subsection 8.1).

- Lessons and recommendations suited for studies related to topology optimization (Subsection 8.2).
- Lessons and good practices related to structural and training hyper-parameter optimization (Subsection 8.3).
- Lessons and recommendations for the optimization of trainable parameters (Subsection 8.4).

### 8.1. General lessons

The first lesson we summarize at this point is in close accordance with the general need for more methodological principles in meta-heuristic research across all the application scenarios where these solvers are applied nowadays. Unfortunately, the optimization problems tackled in this study are not an exception to this claim. The guidelines and procedures to be followed to reach solid and conclusive studies in the use of meta-heuristics are known to the community, specially in regards to the identification of the novel aspects of newly emerging algorithms and the proper design of comparison benchmarks. In this latter research stage, the assessment of the statistical significance is a must when dealing with problems related to Deep Learning topology and/or hyper-parameter optimization. Since several sources of uncertainty may coexist in the same problem statement (e.g. the operators of the search algorithm, the initialization of weights, and the stochasticity of the gradient-based training algorithm), reporting on the statistical significance via hypothesis testing should be considered a necessary step. Despite not an exclusive recommendation of the research area under study, leaving all code and results available in public software repositories for the research community is of utmost necessity, due to the dilated computation times usually required to run experiments with Deep Learning models.

Another general lessons that stems from our study is that the goal in Evolutionary Deep Learning is to yield models of improved generalization properties, namely, to find models that perform *better* when fed with unseen data. To an extent, in our literature review we noted that many contributions in the recent past dismisses this target goal and conclude that the solver at hand performs better since it achieves a more optimal objective value than gradient-based methods. Statements alike should be avoided in prospective studies, as there is no practical value in a model that generalizes worse *in the wild*, e.g. when predicting new data instances.

Another missing point in past contributions is a quantitative evaluation of the complexity of the solver, not only a mere indication of the quality of its produced output. We have shown in the cases of study that important complexity gaps exist between the solvers under study. Neglecting to inspect this important aspect of optimization solvers yields biased conclusions about the practical value of new proposals. Convergence plots like the ones depicted in the case of study II can provide a hint about the relationship between performance (accuracy) and complexity (number of epochs) of the solvers under comparison. Along these plots, a clear definition of an *epoch* should be provided, so as to establish a reference and ensure fair complexity comparisons.

Furthermore, other optimization objectives beyond predictive performance should also be considered in the future, such as the complexity of the optimized model in topology and/or structural hyper-parameter optimization. There are many reasons for considering such objectives in addition to predictive performance, ranging from an easier deployment of the optimized models in constrained computation hardware (e.g. cell phones) or a potentially more interpretable network structure [15]. However, most studies seem to focus just on the predictive performance of the optimized model, leaving unanswered relevant matters for model deployment such as the tradeoff between model's performance and complexity.

Finally, we advocate for a global consensus on the dimensions of realistic benchmark dataset and models. A Deep Learning model is not just a layered structure of perceptrons, but rather a hierarchical composite of neural layers of different nature. It is their capability to extract increasingly specialized features from large-dimensional data what should bestow the label *Deep Learning*. Unfortunately, we have observed many works where *Deep Learning* refers to a neural network with very few trainable parameters and fed with already handcrafted features. The same can be said about the datasets used for validation,

which in many studies lack the complexity that could argue the adoption of Deep Learning models. The community should agree on the minimum characteristics (task difficulty, diversity of datasets, model complexity) of experimental setups devised for reaching meaningful results and valuable conclusions.

### 8.2. *Topology optimization*

We now focus on the lessons learned from the literature that has so far tackled the optimization of the topology of Deep Learning models. First of all, in the first case of study we highlighted that the EvoDeep framework featured a diversity of layer types lower than that of AutoKeras, which restricts the search domain of the evolutionary algorithm that runs at its core. This fact can imprint a subtle yet impacting bias in prospective comparisons between new topology optimization frameworks. Such comparisons should ensure that the counterparts in the benchmark should utilize the same number and diversity of layer types when optimizing Deep Learning models for a given task. Otherwise, there is no certainty whether gaps found among the compared frameworks are due to the differences between their optimization algorithms or to the fact that some of them are in unfair disadvantage. Similar recommendations can be issued in regards to the formulation of the objective to be optimized, which should be set the same among frameworks. In what refers to the optimization algorithm, we have shown that bio-inspired solvers (in particular, evolutionary algorithms) are still far from the performance offered by other ad-hoc methods. Actually, the default optimization approach for the AutoKeras framework is a parameter-wise greedy strategy, and no Evolutionary Computation nor Swarm Intelligence methods are considered whatsoever. This observation implicitly suggests that bio-inspired solvers are still far from performing competitively, as our results in the first case of study have clearly shown. However, it is known that in other application domains, greedy search methods usually fall into local optima unless proper algorithmic countermeasures are included along the search. This opens up an opportunity to hybridize such greedy methods with global search heuristics or, alternatively, to design ad-hoc search algorithms that incorporate any of the ingredients featured by such greedy methods.

Another aspect in which much research effort has been invested is in the design of variable-size encoding strategies for the representation of evolved network architectures. This has been a subject of intense research for years since the advent of the first neuro-evolution frameworks, in particular the adaptation of compositional pattern producing networks made by Hyper-NEAT to represent and evolve neural networks. Despite the numerous works on neural representation strategies published thereafter to date, in many cases the selected encoding approach does not account for the validity of the sequence of neural layers it represents. To this end, EvoDeep resorts to finite state machines to model all possible transitions between layers, followed by a two-part encoding to represent global (training) hyper-parameters and layers' types and structural hyper-parameters, respectively. It is our belief that a key aspect for an efficient heuristic search is to tightly couple the solution encoding strategy to the design of the search operators.

Finally, we dedicate some words of reflection about the level of granularity at which topology optimization should be performed. While in some recent works Deep Learning models are optimized at the level of the computation graph, namely, without assuming any particular type of neural layer. This extreme is interesting should the goal be to transcend conventional layer types and seek more diverse neural structures. At the other end of the scale, other contributions have capitalized on the mixture of pretrained modules, aiming to enhance the structure of the Deep Learning model while leveraging, at the same time, high-level knowledge acquired in other related tasks. To the best of our knowledge, there is no clear consensus on whether high-level or low-level topological optimization strategies are more promising for Deep Learning topology optimization.

### 8.3. *Structural/training hyper-parameter optimization*

When it comes to the optimization of the hyper-parameters of the Deep Learning model, a first recommendation to elicit is to clearly specify the hyper-parameters to be optimized, as well as their search



ranges. Echoing the aforementioned need for fairness in comparison benchmarks, and following our conclusions drawn from the Case of Study I (Section 6), it is crucial to guarantee that the compared solvers explore search spaces of equal complexity so as to remove any bias due to differences in this matter. A good practice for this purpose is to include a table listing each parameter with its corresponding search range, so that fairness can be guaranteed and reproducibility eased for the interested audience.

Another recommendation in hyper-parameter optimization is to estimate the impact of different hyper-parameter values in terms of the predictive accuracy and overall complexity of the optimized model. By reporting on the correspondence between different hyper-parameters values and the overall performance and complexity of the model, the community can discern potentially good search ranges for such hyper-parameters, thereby reducing the time needed to perform new hyper-parametric optimization tasks.

Finally, our taxonomy of the existing literature shown in Figure 4 revealed that the number of works simultaneously tackling structural and hyper-parameter optimization surpassed those dealing only with structural or training hyper-parameter optimization in isolation. This fact calls into question whether the results obtained so far for the latter cases are conclusive. The selected topology for the Deep Learning model affects directly the search space of a structural hyper-parameter optimization problem defined over it, so it remains uncertain whether the conclusions drawn for the particular network topology under choice can be extrapolated to any other network topology. This is why we suggest increasing the number of experiments with different network topologies and datasets when performing hyper-parameter optimization. Otherwise, the claims delivered by prospective studies can be in doubt due to the lack of enough empirical evidence.

#### *8.4. Trainable parameter optimization*

To end with, the optimization of the trainable parameters of Deep Learning models is arguably the one grasping most interest from the community in recent times. Several learned lessons and recommendations can be issued in this regard.

First of all, the research community should come to an agreement and understand that currently, we are far from fully replacing gradient-based solvers with bio-inspired optimization algorithms. The results discussed in the Case of Study II (Section 7) buttress this statement with solid findings: one of the most renowned and competitive algorithms for large-scale global optimization (SHADE-ILS) has not been able to perform better than the gradient-based Adam solver. Besides, when increasing the complexity of the Deep Learning model, the performance of SHADE-ILS degrades severely even if granted more computational budget (number of loss evaluations) than its gradient-based counterpart. In summary, for the time being bio-inspired optimization algorithms cannot rival the computational efficiency and the quality of solutions produced by gradient-based methods.

The main reason for this conclusion can arguably be found in the structure of the Deep Learning model, which establishes relationships among the trainable parameters of consecutive layers that *should* be exploited by the optimization algorithm. This is actually what gradient back-propagation realizes in a clever yet computationally efficient fashion, even though creating other known issues (e.g. gradient vanishing). In our experiments we noticed that when adapting the search behavior of SHADE-ILS to the layered structure of the network, simple layer-wise schedules of this algorithm yields remarkable performance boosts, specially for networks of relatively small size. This suggests that more sophisticated hybridization strategies should be further investigated to embed problem-specific knowledge (namely, the structure of Deep Learning models or gradient information) within the search procedure of bio-inspired solvers.

Notwithstanding this noted performance gap, gradient-based solvers restrict the spectrum of loss functions to those for which derivatives can be computed. However, bio-inspired algorithms do not impose any requirement on the objective function to be optimized, nor do they require it to be differentiable. This fact could tilt the scale towards the use of bio-inspired algorithms in singular learning tasks that re-

quire a tailored definition of the objective function, as in cases with severe class imbalance or multilabel classification.

When it comes to computational efficiency, the higher complexity of population-based meta-heuristics when compared to that of gradient-based solvers should stimulate more parallel and distributed implementations of Evolutionary Computation and Swarm Intelligence methods. There are modern programming languages and frameworks that can be utilized to accelerate bio-inspired search algorithms, even if still lagging behind the typical runtimes of gradient-based techniques. Recent works aim indeed at this direction, reviewing implementations available so far and prescribing recommendations and guidelines for the implementation of meta-heuristics in GPU [211, 212] and asynchronous distributed computing architectures [213]. Experiments with large datasets and realistic Deep Learning models should capitalize on already available software packages that ease the seamless deployment of meta-heuristics in massively parallel computing hardware, such as *jMetalPy* [214] (Apache Spark and Dask) and *libCudaOptimize* [215] (CUDA for GPU). Interestingly, Tensorflow (the computation engine that underlies well-known software libraries for Deep Learning models) also provides a naive implementation of Differential Evolution as one of its functionalities [216]. Parallel, federated or distributed computation frameworks for Deep Learning models are also spreading fast [62, 63, 217]. Definitely new studies should leverage the availability of these tools to undertake experiments at realistic complexity scales.

We end with our learned lessons on trainable parameter optimization by emphasizing several good methodological practices that should be followed in prospective studies. First, the accuracy achieved by the optimized models over the test set should be informed jointly with the usual objective function statistics reported in experiments with bio-inspired meta-heuristics. This is particularly relevant in trainable parameter optimization to assess whether performance gaps identified between solvers do not come along with a penalty in the generalization of the evolved model. Furthermore, conventional gradient-based solvers should be always included in the benchmark, even if their lower complexity makes the comparison unfair in such terms. Finally, we recommend the use of convergence plots such as the ones depicted in Figures 6.a to 6.d as a visual tool to examine the relative differences between algorithms over epochs. This information can be very valuable for the deployment of the solver(s) in real hardware, as well as for the detection of overfitting issues like the one identified in our experiments.

## 9. Challenges and Research Directions

Our exhaustive literature review and performed experiments have unveiled several promising facts and unsolved caveats of Evolutionary Computation and Swarm Intelligence algorithms when used for addressing optimization problems related to Deep Learning. Nonetheless, many proposals have been contributed to date for assorted learning tasks, not only supervised and unsupervised learning, but also other paradigms relying on Deep Learning models (e.g. deep reinforcement learning). Despite this noted activity, several research niches still remain uncharted or insufficiently addressed in this fusion of technologies.

In this section we summarize several challenges and research directions that should be under the target of future efforts conducted in this area. Such challenges are schematically depicted in Figure 7, and contribute to the last two questions targeted by this overview: what can be done in future investigations on the confluence between bio-inspired optimization and Deep Learning, and what should future research efforts be conducted for?

### 9.1. Large-scale optimization for model training

The design of bio-inspired algorithms capable of efficiently tackling large-scale optimization problems seems to be one of the critical points that require further developments to train Deep Learning models of realistic complexity levels. This is the reason for the selection of SHADE-ILS as the search algorithm

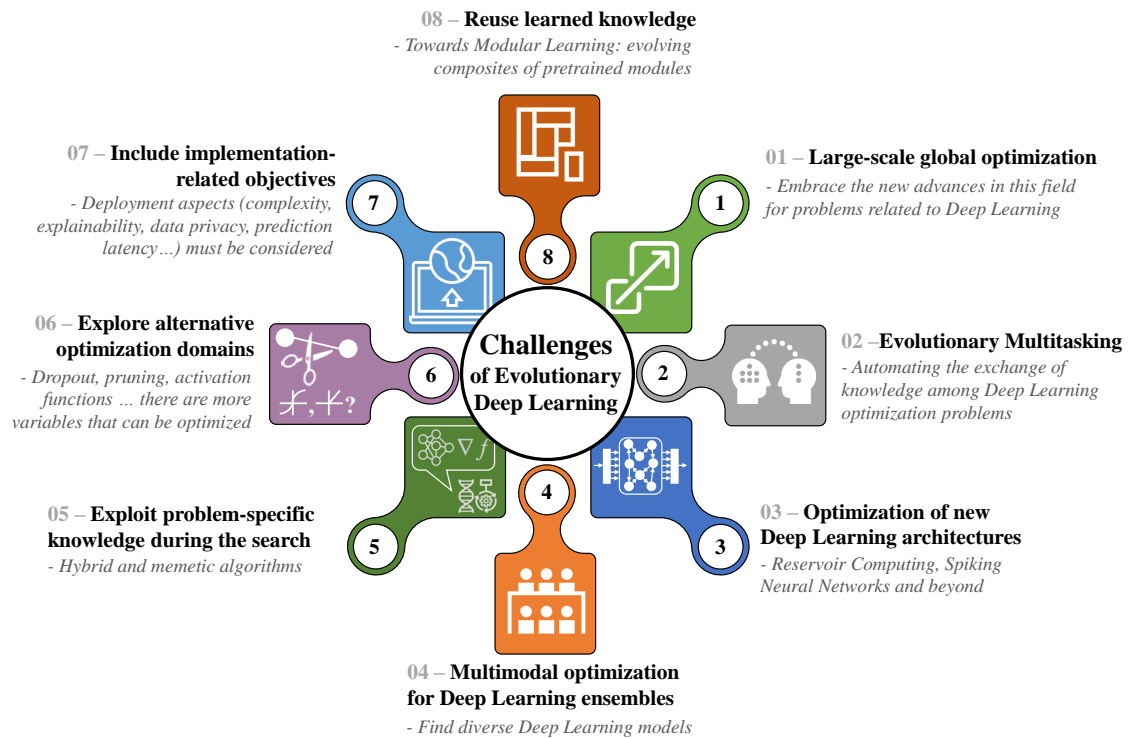


Figure 7: Challenges and research directions envisioned for Evolutionary Computation and Swarm Intelligence for the optimization of Deep Learning models.

in the second case of study. Indeed, SHADE-ILS remains nowadays as one of the most competitive proposals for large-scale global optimization, and is regularly considered as a baseline for competitions and benchmarks.

However, as in other research areas related to meta-heuristics, many advances in large-scale global optimization are regularly contributed to the community, featuring sophisticated ways to infer and exploit the correlation between variables during the search process (*interaction learning*). Improving this feature in large-scale solvers is often the main target of new proposals, either in an implicit fashion (as in Estimation of Distribution Algorithms, and Bayesian Optimization) or explicitly via grouping, statistical correlation-based methods, decomposition or other assorted means [218].

Unfortunately, our analysis has revealed that most works related to trainable parameter optimization have resorted to off-the-shelf variants of bio-inspired solvers. Consequently, no consideration is made about the interactions between variables (weights, biases) that are known to occur due to the neural connections throughout the multiple neural layers. This motivates a closer look to be taken at new advances in large-scale global optimization, for both single- and multi-objective optimization problems [219]. Given the upsurge of Deep Learning problems in which more than one objective is established [54, 57, 101, 220, 221], the use of multi-objective solvers for large-scale optimization seems to be a natural choice.

## 9.2. Evolutionary multitasking

An interesting research area has revolved lately around the design of evolutionary multitasking algorithms capable of simultaneously addressing several optimization problems within a unique search process that exploits the complementarities and synergies existing among such problems [222]. The

challenge in this area is to develop intelligent optimization methods that not only promote the exchange of knowledge among candidate solutions corresponding to related problems, but also prevents the convergence of the search from being affected by *counterproductive* knowledge transfers among unrelated tasks.

When framed within the current study, the adoption of large-scale evolutionary multitasking to optimize simultaneously different Deep Learning models can boost even further the possibilities foreseen for the intersection between Transfer Learning and bio-inspired optimization. For instance, the transfer of pretrained modules between tasks can be conceived as a crossover strategy between networks partially evolved for undertaking different tasks. Similarly, the exchange of the parameters values between Deep Learning models can be also automated via evolutionary multitasking towards evolving behavioral policies for different reinforcement learning tasks [193]. Evolutionary multitasking has been also used to achieve modular network topologies [223]. The relative youth and promising results shown by evolutionary multitasking techniques are a sign of objective evidence that optimization problems related to Deep Learning should be explored via these techniques, e.g. by leveraging the straightforward exchange of knowledge among networks allowed by their hierarchically layered structure. Furthermore, developments in multi-objective evolutionary multi-tasking [224, 225, 226] open up further opportunities towards considering other objectives beyond accuracy of relevance for Deep Learning, such as the complexity of evolved topologies.

### 9.3. Optimization of new Deep Learning architectures

Most of the reviewed literature on bio-inspired algorithms for Deep Learning has focused on traditional forms of neural computation, including convolutional filters and recurrent units. However, this major activity has set apart the optimization of other neural network flavors, for which *Deep* (multi-layered) versions have been proposed over the years. Such alternative deep architectures have fewer optimization variables, hence favoring the use of naive bio-inspired solvers for the different problems that can be formulated on them.

One of such neural families is *Reservoir Computing*, which comprises a number of recurrent neural networks where only the parameters of the output layer (the readout layer) are learned. The parameters of the rest of recurrent neurons (the *reservoir* are randomly initialized subject to some stability constraints, and kept fixed while the readout layer is trained [227]. Some works have been reported in the last couple of years dealing with the optimization of Reservoir Computing models, such as the composition of the reservoir, connectivity and hierarchical structure of Echo State Networks via Genetic Algorithms [228], or the structural hyper-parameter optimization of Liquid State Machines [229, 230] and Echo State Networks [231] using an adapted version of the Covariance Matrix Adaptation Evolution Strategy (CMA-ES) solver. The relatively recent advent of Deep versions of Reservoir Computing models [232] unfolds an interesting research playground over which to propose new bio-inspired solvers for topology and hyper-parameter optimization.

Despite more scarcely, optimization problems related other families of neural computation models have also been approached via Evolutionary Computation and Swarm Intelligence. The most remarkable case is the family of Spiking Neural Networks, in which topology and structural hyper-parameters have been addressed by means of different bio-inspired solvers [233, 234]. Training of synapses in spiking neural architectures has been also tackled in [235, 236]. We fully concur with the prospects outlined in the recent overview on training methods for Spiking Neural Networks [237]: efficient large-scale training methods should be investigated for new variants of these models, such as spiking deep belief networks and spiking convolutional neural networks.

### 9.4. Multimodal optimization for Deep Learning ensembles

Another interesting research path stems from the adoption of niching methods used in bio-inspired algorithms for multi-modal problems for the construction of Deep Learning ensembles (also referred to

as *committees*. Indeed, the evolved population of candidate networks can be employed to retain near-optimal yet diverse Deep Learning model configurations. Such a diversity can emerge from different evolved topologies and/or values of their (hyper-)parameters. Such retained network configurations can be assembled into a committee, allowing for a robust fusion of their issued decisions.

Work in this direction has been published recently in [238], where a niching mechanism is used to penalize individuals that are more similar to others in the population are penalized. Network configurations remaining in the population after the search are then combined together via majority voting, showing a significant improvement in performance with respect to the use of a single evolved Deep Learning model. Multi-modal optimization methods can also be combined with diversity induction techniques (e.g. Novelty Search [239]) to promote an efficient exploration and discovery of multiple global optima over strongly multi-modal search spaces, just like the ones known to characterize Deep Learning optimization problems.

### 9.5. *Exploitation of problem-specific knowledge during the search*

In light of our experiments, we believe that the research community working on bio-inspired optimization should *go memetic* and exploit the rich structural properties of neural networks when designing new algorithmic approaches for any of the problems related to Deep Learning. A diversity of strategies can be followed for this purpose, from simplest layer-wise search schedules as the ones proposed in Case of Study I, to more sophisticated means like iterating between meta-heuristic and gradient-based search, or the use of similarity measures between neural layers [240] for controlling the behavior of search operators in topology and/or structural hyper-parameter optimization.

By blending together the global search capabilities of bio-inspired meta-heuristics and the problem-specific knowledge embedded in e.g. back-propagated gradients, a major performance improvement can be obtained, effectively overcoming known issues such as gradient vanishing or slow convergence. This is actually the approach followed in [21], yet validated with small network sizes. Another problem-specific aspect hybridized with bio-inspired algorithms can be found in [34]. In this work, the need for inducing curiosity in reinforcement learning models (specially in environments with sparse rewards) is realized by not driving the search of the learning agent with the reward objective, but rather with a measure of the *behavioral novelty* of the resulting policy. Other examples of hybrid methods are [241, 242, 160, 158].

This can be conceived as another example of the importance of considering the particularities of the learning problem in the design of bio-inspired algorithms. We firmly advocate for more proposals along this direction when addressing optimization problems in Deep Learning.

### 9.6. *Exploration of alternative optimization domains*

The taxonomy around which our literature analysis has been organized considers three possible optimization domains on which a problem related to Deep Learning can be formulated: 1) topological variables, namely, the number and type of layers that compose the model; 2) hyper-parameters, which establish the details of the layers (structural hyper-parameters) and the optimization algorithm chosen for training the model (training hyper-parameters); and 3) trainable parameters (weights and biases). This threefold categorization collectively reflects most contributions reported to date in this research area.

However, there are more optimization domains related to Deep Learning that can be tackled with bio-inspired solvers. One of them is pruning, e.g. the selective removal of connections among neurons for different purposes, from regularization against overfitting to a lower computational burden of gradient back-propagation solvers. Different pruning strategies can be developed to select which connections to drop at every layer of the Deep Learning model, which have been reviewed in recent comprehensive surveys on this topic [243, 244]. To the best of our knowledge, the use of Evolutionary Computation and Swarm Intelligence for neural pruning has been only done at the level of optimizing dropout rates. When turning the focus on more fine-grained pruning strategies, large-scale optimization algorithms can

be used to determine the subset of neural connections that must be discarded to achieve a good trade-off between generalization performance and network compactness. Recent findings in the application of genetic algorithms to convolutional channel selection [245] should be followed by further studies evaluating the scales at which network pruning with meta-heuristics can be realized.

Other optimization domains for which Evolutionary Computation and Swarm Intelligence methods can be applied include the tailored design of activation functions [136], or the fusion of decisions issued within Deep Learning ensembles [89]. Definitely, these problems and other ones still to be proposed lay a magnificent panorama for bio-inspired optimization.

### *9.7. Inclusion of multiple implementation-related objectives*

When evolving Deep Learning models, objectives and constraints related to the application scenario on which they are to be deployed should be considered in the optimization problem. For instance, many software libraries and embedded electronic chips can be found nowadays in the market for the implementation and execution of neural network models in constrained computing devices, such as Internet of Things (IoT) sensors [246]. Likewise, real-world applications such as autonomous vehicular driving, remote precision surgery or wearable sensors restrict severely which Deep Learning models can be rolled out in their equipment. This suggests that a closer look should be taken at the complexity of evolved Deep Learning models during the optimization of their topology and hyper-parameters, among other domains.

A similar elaboration can be also made in regards to the progressive maturity of new paradigms such as Federated Learning [247] and Edge Computing [248]. In these paradigms not only local models lack the amount of computational resources needed to run complex Deep Learning models, but also additional objectives are imposed. Efficient (incremental) training algorithms, energy consumed by the model [249], data privacy preservation [250, 251], robustness against adversarial attacks [252], or the explainability and accountability of decisions [15] are some illustrative examples of the eventual confluence of multiple implementation-related objectives in Deep Learning optimization. However, these factors are rarely taken into account in current approaches for evolving Deep Learning models. Most of them rely just on accuracy or any other measure of predictive performance.

This being said, Deep Learning models should be evolved by considering together several objectives and constraints as the ones exemplified above. To this end, we foresee that bio-inspired algorithms for multi-objective optimization [253] can play a differential role in the future of Deep Learning. This branch of bio-inspired optimization, along with techniques devised to handle constraints during the search [254], can catalyze the practical deployment of Deep Learning models by addressing the improvement of the model's generalization performance together with implementation-related objectives and constraints.

### *9.8. Reuse of learned knowledge: towards modular learning*

The reutilization of pretrained modules between models corresponding to related learning tasks is at the core of Transfer Learning [255], whose most straightforward strategy is to reuse parts of a model developed for a task as the starting point for a model devised for another task. Such parts are often conceived as structural fragments of the network, particularly those capturing high-level features from the input to the Deep Learning model. Features corresponding to those parts are more easily reusable among tasks due to their high-level nature. In image classification, for instance, features learned by the first layers of the network are borders and broad shapes that could be of help for many different tasks. Therefore, it is intuitive to think that the output of such layers can be of help in other tasks for which annotated data instances are scarce.

The availability of Deep Learning models trained on different datasets and the effectiveness of just transferring layers between models corresponding to different tasks suggest a very interesting research path: expanding further the search space of topology and structural hyper-parameter search to also consider pretrained modules. The inclusion of such modules in the alphabet of possible layers could yield a major boost of Deep Learning models, specially for those cases with few labeled data. Furthermore,

the flexibility of bio-inspired algorithms when designing the encoding strategy that represents networks during the search could allow achieving finer levels of granularity in the knowledge imported from such pretrained modules, to the scales of convolutional filters or recurrent units. There is a great opportunity to bring together transfer learning and topology/structural hyper-parameter optimization around the same goal: to evolve and discover Deep Learning models of superior performance.

## 10. Conclusions and Outlook

This paper has presented a comprehensive and critical review on the use of Evolutionary Computation and Swarm Intelligence approaches to the topological, hyper-parametric and/or trainable parameter optimization of Deep Learning models. As we have previously indicated, the paper is focused on three axes: a) definition of optimization problems in Deep Learning and taxonomy; b) a critical methodological analysis of the related literature and two cases of study, allowing to prescribe learned lessons and recommendations for good practices; and c) an enumeration of challenges and new directions of research. We highlight the aforementioned two cases of study, providing factual results on the performance of bio-inspired optimization algorithms when applied to the architectural design, hyper-parameter tuning and training of Deep Learning models.

Our elaborations made throughout these three axes have yielded informed conclusions and insights about the four fundamental questions posed in the introduction, which round up the critical review of the field targeted in this overview. We synthesize below our responses to such questions, in the form of reflections stemming from the aforementioned axes:

### 1. **Why** are bio-inspired algorithms of interest for the optimization of Deep Learning models?

The increased scales and diversity of neural layers of modern Deep Learning approaches have lately reactivated the global interest in Deep Learning optimization with bio-inspired algorithms, as a means to automate efficiently the processes of designing their topology, tuning their hyper-parameters and learning their parameters. Such processes can be formulated as complex optimization problems, motivating the adoption of bio-inspired algorithms for solving them efficiently. Furthermore, the renowned global search capability of Evolutionary Computation and Swarm Intelligence methods makes them a suitable choice to deal with complex search spaces as those characterizing Deep Learning problems. Finally, the flexibility of bio-inspired solvers to be hybridized with problem-specific search methods is another reason supporting the hypothesis that Deep Learning optimization can largely benefit from them.

In conclusion, we observe solid grounds for this synergistic fusion of technologies, which has so far stimulated the community to tackle optimization problems in Deep Learning using bio-inspired algorithms. However, we recognize that this fusion has not yet achieved results that are truly a step forward in terms of quality and objective achievement. This is still a lost race of bio-inspired optimization algorithms with respect to gradient-based solvers, which remain as the horse at the head of the race.

### 2. **How** should research studies falling in the intersection between bio-inspired optimization and Deep Learning be made?

Our discussion on the results obtained in our cases of study suggest that bio-inspired algorithms can be used for topological and/or hyper-parameter optimization, yet performing worse than other methods when comparisons are fair in terms of search space and complexity. Furthermore, our experiments have also revealed that even competitive bio-inspired solvers for large-scale global optimization are outperformed by conventional gradient-based solvers for trainable parameter optimization. These results, along with several issues detected in the literature (most notably, the unrealistic scales of the

evolved model and the dataset/task under consideration), support our claims that there is a large space for improvement in this research area.

On a prescriptive note, we have identified several learned lessons and recommendations that trace *how* research should be done to reach solid conclusions and sound achievements. We highlight below the most relevant ones:

- Good methodological practices when designing experiments and benchmarks between different solvers, including realistic datasets and models, fairness in terms of computational complexity, assessment of the significance between performance gaps and reported performance scores over test instances, among others.
- A closer attention at encoding strategies for topology optimization that account for the validity of the composition of layers that they represent.
- A clear definition of the variable search ranges in structural and training hyper-parameter optimization, so that differences emerging between solvers can be attributed exclusively to their search efficiency.
- The exploitation of problem-specific knowledge in trainable parameter optimization: the interactions between trainable parameters imposed by the hierarchical structure of neural connections should be exploited further by bio-inspired solvers for them to step out from the shadows of their application to the training of Deep Learning models.

### 3. **What** can be done in future investigations on this topic?

In this regard we have underscored the need for overcoming the computational inefficiency observed in current bio-inspired optimization algorithms with respect to gradient-based solvers. It is known that these latter solvers have also their own drawbacks: they require differentiable loss formulations, they are sensitive to vanishing and exploding gradients and they are prone to local optima in non-convex search spaces. There are well-founded reasons why bio-inspired solvers can be a firm alternative to gradient-based methods, but more efficient designs and/or better performing implementations of bio-inspired approaches should be under active investigation in the future for them to become a practical choice for Deep Learning model training.

We have also stressed on other subareas in Evolutionary Computation and Swarm Intelligence of utmost interest for their application to Deep Learning optimization. Large-scale global optimization for training Deep Learning models of realistic complexity, or multi-modal optimization for the automated construction of Deep Learning ensembles, are just a few examples of the myriad of research niches in bio-inspired optimization that have not been explored yet. Evolutionary multitasking is another interesting direction to follow when approaching several Deep Learning optimization problems at the same time, mostly when such problems are related to each other and share a significant degree of overlap in their solutions (as occurs in transfer learning). Finally, optimization problems formulated on alternative Deep Learning models seem to have received less attention to date, and unleash further opportunities for bio-inspired optimization.

An additional research direction in this fusion of technologies has been identified in the existence of other variables and domains in which optimization problems can be formulated. Ensemble construction, network pruning or selective dropout strategies can also be described as optimization problems, favoring the adoption of bio-inspired optimization algorithms for their efficient solving.

### 4. **What** should future research efforts be conducted **for**?

There is a strong incentive for which Deep Learning optimization should be approached via bio-inspired algorithms in the future: the consideration of additional objectives and constraints linked to the application scenario on which the model is to be deployed. Aspects such as the available



computational resources, the need for periodically updating the model, or the time taken by the model for issuing a prediction should be considered during the design of the model for it to be of practical value.

Furthermore, new paradigms such as Edge Computing, explainable Artificial Intelligence and Federated Learning have underpinned the need for taking into account other objectives beyond the accuracy of the model. Aspects such as data privacy preservation, the explainability of decisions issued by the model, the non-stationary nature of data at the edge, or the complexity of the learning algorithm can be formulated as additional objectives for the design of the model, impacting on its topology, hyper-parameter values and other optimization variables.

These arguments, combined with the flexibility of bio-inspired algorithms to deal with multiple conflicting objectives, can be a primary *what for?* driver for adopting them to evolve Deep Learning models in practical settings. Specific scenarios and contexts that require ad-hoc designs to be evaluated with multiple objectives can and should also open the door to evolutionary Deep Learning models towards meeting imposed goals in terms of efficiency, performance and other application-related objectives.

All in all, we hope that the material and insights given in this work serve as a reference material for readers willing to arrive at this research area with a clear and thorough understanding of its recent past, current status, and potential in years to come. There are promising evidences that Evolutionary Computation and Swarm Intelligence can tackle optimization problems related to Deep Learning, but we firmly believe that they are not close to maturity, nor do they justify yet the replacement of other solvers used for the same purposes. Nevertheless, this is the role of research itself: to build upon the shadows of knowledge and bring light through scientific achievements. This survey has just lit a candle to illuminate this path through the field of Evolutionary Deep Learning.

## Acknowledgments

Aritz D. Martinez, Esther Villar-Rodriguez, Eneko Osaba and Javier Del Ser acknowledge the funding support received from the Basque Government through the EMAITEK and ELKARTEK programs (3KIA project, KK-2020/00049), and the Spanish *Centro para el Desarrollo Tecnológico Industrial* (CDTI, Ministry of Science and Innovation) through the *Red Cervera* Programme (AI4ES project). Javier Del Ser also acknowledges funding support from the Consolidated Research Group MATHMODE (IT1294-19) granted by the Department of Education of the Basque Government. Siham Tabik, Daniel Molina and Francisco Herrera would like to thank the Spanish Government for its funding support (SMART-DaSCI project, TIN2017-89517-P), as well as the BBVA Foundation through its *Ayudas Fundación BBVA a Equipos de Investigación Científica* 2018 call (DeepSCOP project).

## Appendix A. Deep Learning Models

The spectrum of Deep Learning models is certainly huge nowadays, with new learning variants for different tasks proposed continuously by the community. In what follows we provide a brief introduction of the DL models in which Evolutionary Algorithms and Swarm Intelligence methods have been mostly used for addressing the above problems.

### *Deep Boltzman Machines and Deep Belief Nets*

Deep Boltzman Machines (DBMs) and Deep Belief Nets (DBNs) are considered as generative probabilistic models comprised by a stacked hierarchy of Restricted Boltzmann Machine (RBM) layers. Unlike classical Boltzmann Machines, RBMs have no intra-layer connection between nodes, and comprise two

layers of fully-connected neurons (visible and hidden), in charge of learning the probability distribution of a set of binary-valued inputs. DBMs are constructed as a series of RBMs stacked on top of each other, whereas DBNs are hybrid models whose interactions are composed of indirect connections at the top layers (RBM) and downward directed belief nets between the lower ones.

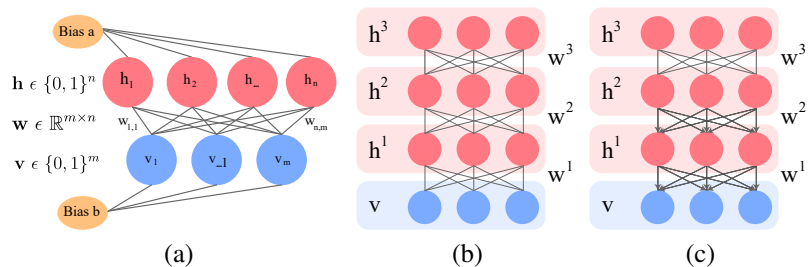


Figure A.1: Architectures of (a) a Restricted Boltzmann Machine; (b) a Deep Boltzmann Machine; and (c) a Deep Belief Net.

The learning process of DBM and DBN is conducted by a greedy layer-by-layer pre-training approach, followed by a discriminative fine-tuning process performed by a back-propagation algorithm according to labelled data provided in the training set (supervised learning). This procedure modifies the learned layerwise parameters to ensure the appropriate learning of the labels/desired outputs.

Both DBMs and DBNs have restricted structural topologies. Therefore, the optimization tasks that have been central for these models relate to the proper selection of structural (number of hidden units) and training hyper-parameters (training epochs and learning rate, among others). Efficiently tuning these hyper-parameters has been a subject of intense scientific research over the history of these models, as we show in our literature study.

#### Autoencoders

Autoencoders (AEs) are composed by two different typically mirror frameworks, named encoder and decoder, and are deemed as the inspiration behind  $G$  models. Briefly explained, the principal goal of AEs is to learn an encoded representation of their input data. In order to achieve that, the reconstruction error between the input data and the generated output is minimized by making the encoder create condensed low-dimensional data abstractions such that, when decoded, the original input is reconstructed with high fidelity. This process assist the model to learn critical and latent features present in the data. Thus, AEs are composed by four core parts:

- *Encoder*: in charge of projecting high-dimensional data down to a subspace of encoded representation.
- *Bottleneck*: the most compressed representation of the input data.
- *Decoder*: the model learns the process for the reconstruction of the data.
- *Reconstruction Loss*: responsible for assessing the adequacy of the solutions offered by the decoder, measuring how close the output is regarding the input data.

It is remarkable that a wide variety of AEs can be found along the literature, such as Convolutional, Variational or Denoising AEs, among many other approaches. Having each of them their specific characteristics and optimizable hyper-parameters. In all of them, the main structural hyper-parameter is the size of the bottleneck, also known as latent space size. Nevertheless, the fact that they have two separated networks to optimize, makes AEs attractive for topology optimization (Problem 1). Although the topological architectures of encoder and decoder parts are mirrored, they can be optimized on their own, as well as their connection patters. In terms of training, the AE is trained using gradient back-propagation.

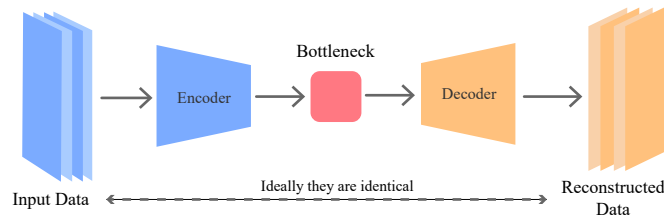


Figure A.2: Architecture of a basic Autoencoder.

As such, a fully or partially substitution of the algorithm can be performed (Problems 3 and 4 as defined in Section 3).

### Convolutional Neural Networks

Commonly employed to deal with pattern extraction from images and videos, Convolutional Neural Networks (CNNs) imitate the behavior of the visual cortex when processing images ([256, 257]) by laying out, in their simplest architectures, a hierarchical set of convolutional and pooling layers. Convolutional layers allow for a spatial analysis of their input data by applying a sliding dot product between the input tensor and a set of filters (also referred to as *kernels*), which transforms the input tensor into a feature map. Pooling layers perform a spatial downscaling on its input tensor, making the model less sensitive to small variations and achieving better generalization properties of the overall model.

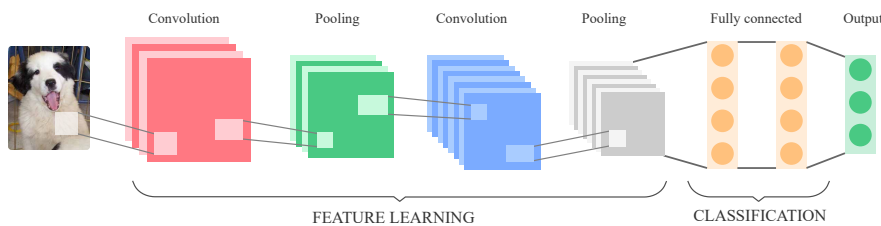


Figure A.3: Architecture of a Convolutional Neural Network.

CNNs unleash a rich background for optimization due to the variety of layers composing their architecture. Indeed, there are multiple structural hyper-parameters related to those layers, such as *kernel size*, *padding size* or *number of filters* in convolutional layers, or *pool size* and *stride number* in pooling layers, among others. CNNs are also suitable for topological optimization tasks, with the selection of layer types or the optimization of connection patterns among layers as the ones mostly addressed in the literature. The number of neurons/filters to allocate at each layer is considered as structural hyper-parameter optimization. Finally, the training optimization task related to CNNs is to substitute the conventionally used gradient back-propagation based solver by a meta-heuristic algorithm.

### Recurrent Neural Networks

Recurrent Neural Networks (RNNs) allow to learn from the relationship between items in sequence data, and are therefore typically used in speech and handwriting recognition, time series forecasting, natural language processing and video processing. RNNs are constructed as rolled neural networks with an internal state called *memory*, which allows processing inputs of variable length, as well as previous predictions to be used as inputs for subsequent predictions steps over the input sequence. Designed as blocks of memory cells, they internally comprise a repeated neural network comprising two core elements: a hidden state and gates. Information flows through the cells, with gates regulating which

portion of the input information should be thrown away or kept (*forget gate*), retained in the cell’s memory (*input gate*), and revealed to the cell’s output (*output gate*).

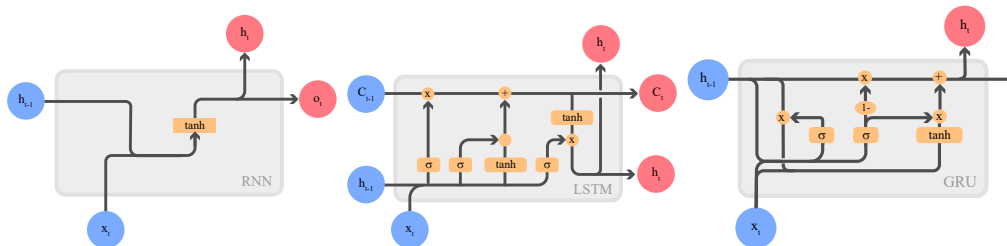


Figure A.4: RNN, LSTM, GRU architectures.

This is indeed the working mechanism underneath two of the most renowned RNN variants to date: Long Short-Term Memory units (LSTMs) and Gated Recurrent Units (GRUs). These two RNN models emerged as a workaround to the vanishing gradient problem faced by traditional RNNs. The most remarkable difference between simple RNNs and LSTMs/GRUs is that the latter incorporate a mechanism able to keep information along the units (temporally), letting the output at  $t - 1$  influence on the result at instant  $t$ . On the other hand, the main difference between GRUs and LSTMs is that LSTMs provides more flexibility in the control of the information flowing through the cell, by virtue of its internal cell state. On the contrary, GRU cells provide less control as they couple forget and input gates, but as a result have less parameters to be trained. Therefore, GRUs can be seen as a simplified version of LSTM cells.

The performance of recurrent neural networks is closely related to the number of units taking part in the recursive connection. Because of their restricted architecture, the main task for this networks is the optimization of the number of hidden units (topology optimization, problem 1 in Section 3). As noted in the literature analysis, dropout and recurrent dropout rates are also often optimized in recurrent architectures (structural hyper-parameter optimization, problem 2).

### Generative Adversarial Networks

Generative Adversarial Networks (GANs) were introduced by Goodfellow et al. in [11]. They are essentially two neural networks trying to beat each other, consequently progressing as per their opposite respective goals: a discriminator network, whose main objective is the classification of the input data and the evaluation of its authenticity; and a generator network, in charge of producing realistic synthetic data. The generator network attempts at “fooling” the discriminator network by generating new data instances increasingly closer to the genuine ones, whereas the discriminator aims at detecting whether the output of the discriminator is real or fake, progressively becoming more competent in this detection task. A schematic diagram of a typical GAN is depicted in Figure A.5.a.

Like CNNs, GAN architectures present a vast variety of optimization problems. Nevertheless, all of them are very similar to those of CNNs and RNNs described previously in this section. Depending on the nature of the data under consideration, both discriminator and generator may include convolutional and/or recurrent layers, thereby exposing the same need for optimizing the topology, hyper-parameters and/or trainable parameters of CNNs and RNNs.

### Deep Reinforcement Learning

Deep Reinforcement Learning (DRL) is considered as an special case due to the amount of contributions in the field. It is not considered a model architecture, but rather a different learning paradigm in which the output of the Deep Learning model defines how an agent should interact with an environment. In general terms, the modules involved in a DRL approach are:

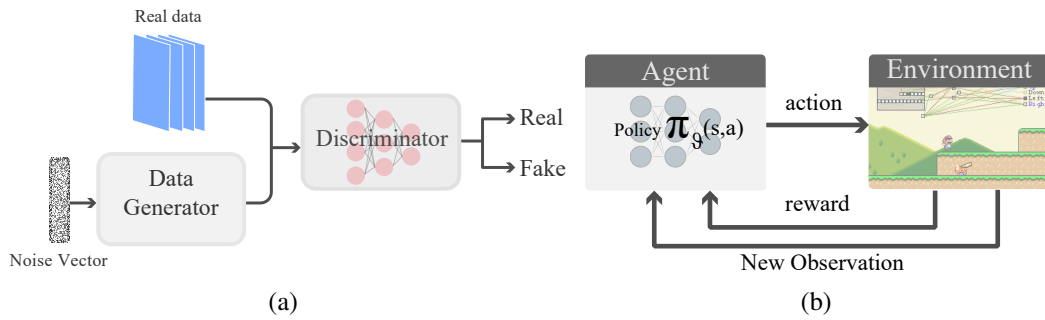


Figure A.5: General architecture of (a) a Generative Adversarial Network; and (b) a Deep Reinforcement Learning model.

- *Environment*, namely, a scenario/asset/system modeled in a way that a human can interact with them through several variables. Upon the application of certain actions, the environment feeds the agent with observations, and evaluates the actions taken by the agent returning a value of reward.
- *Agent*, which embodies a policy that maps observations to actions. In the specific formulation of Deep Reinforcement Learning, the policy itself is intrinsically encoded in the parameters of a Deep Learning model, which is designed to take the action that maximizes the reward given a specific observation from the environment.

The general architecture of this kind of approaches is illustrated in Figure A.5.b. Within the field of DRL many type of approaches are covered, from MLP to RNN, Dueling or Multiagent environments, which have been demonstrated to get good results when they are combined with meta-heuristic algorithms [39, 40, 38]. Besides the discussed architectural and topological hyper-parameters that depend on the model that is adopted in the DQL approach (e.g. CNN, RNN or LSTM), the environment information is also susceptible to be manipulated to decrease the complexity, or to try to make the learning procedure computationally affordable.

## Appendix B. Bio-inspired Optimization: Evolutionary Computation and Swarm Intelligence

The need for search algorithms capable of efficiently dealing with the optimization problems arising from Deep Learning has stimulated an upsurge of literature proposing different solvers for this purpose. We now provide a brief overview of the optimization research area, with a focus on meta-heuristic algorithms that are inspired by biological sources of inspiration. For a more detailed overview of developments and prospects in this research area we refer to recent comprehensive reviews in [201, 258].

As shown in the taxonomy of Figure B.1, optimization methods can be first grouped in three categories: exact methods, heuristics and meta-heuristics. Exact methods are those that always solve an optimization problem to optimality, either by exploring the entire space of solutions or by taking advantage of specific characteristics of the problem at hand (e.g. linearity, convexity). On the other hand, a heuristic search algorithm addresses a given optimization problem by resorting to knowledge related to the domain where the problem is formulated. By exploiting this domain-specific information in its search algorithm, a heuristic explores the space of feasible solutions efficiently, intensifying the search around the most promising areas as per the objective(s) under consideration. Finally, the third category corresponds to meta-heuristics, which lie at the core of this study.

Briefly explained, a meta-heuristic optimization algorithm solves a problem using only general information and knowledge common to a wide variety of problems with similar characteristics [259]. Meta-heuristic algorithms explore the solution space by progressively learning how candidate solutions should be modified towards optimality, with the aim of reaching increasingly promising results disregarding the

characteristics of the problem being tackled. Given their self-learning nature and their abstraction from the problem itself, meta-heuristic approaches are well-suited to deal with real-world problems featuring complex search spaces and even non-analytically defined objectives/constraints. This is in fact the reason why meta-heuristics have taken a prominent role when addressing the optimization problems underneath Deep Learning.

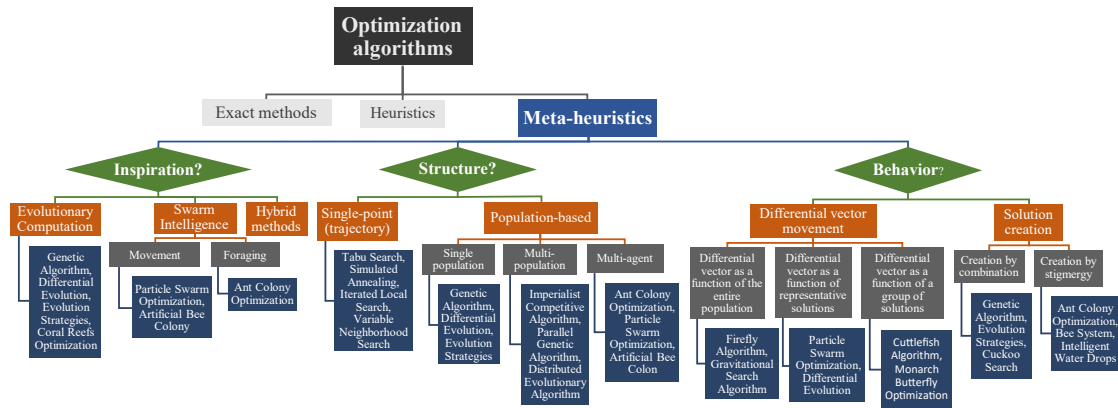


Figure B.1: Taxonomy of optimization algorithms, with a focus on meta-heuristic optimization algorithms as per the different criteria under which they can be classified. Some examples of algorithms are also given. For a further insight into nature and bio-inspired optimization approaches we encourage the reader to the extensive review in [260].

Deeper into the taxonomy of Figure B.1, meta-heuristics can be further divided into different groups depending on several criteria. To begin with, we can distinguish between 1) single-point (also referred to as trajectory-based) meta-heuristic methods, which rely on the progressive improvement of a single solution to the problem by exploring its neighborhood under a set of movement operators (as in e.g. Tabu Search, TS [261] or Simulated Annealing, SA [262]); and 2) population-based techniques, which maintain a set of possible solutions of the problem that interact with each other towards producing new solutions of increased quality (e.g. Genetic Algorithm, GA [263, 264], Ant Colony Optimization, ACO [265] or Particle Swarm Optimization, PSO [266]). This last category can be broken down in 2.1) multi-population techniques, where the population is divided into different that evolve separately, exchanging information periodically (for instance, the Imperialist Competitive Algorithm [267]); 2.2) multi-agent methods, whose population is composed by multiple diverse agents with different roles that interact with each other towards optimality (e.g. Artificial Bee Colony, ABC [268]); and 2.3) single-population approaches, such as the aforementioned GA. At the same time, meta-heuristics can also be divided as per its search behavior, yielding A) differential vector movement based methods, which rely on the computation of a differential vector to move from a reference solution towards a new candidate; and B) solution creation based methods, which generate new solutions to explore the search space instead of moving existing ones. Other criteria can be adopted for organizing the enormous corpus of literature related to optimization meta-heuristics, such as the support of the optimization variables of problems that can be tackled by the meta-heuristic at hand (discrete/continuous/mixed), the scope of the search (global/local), or the stochastic/deterministic nature of their search operators, among others.

Among these criteria, the inspiration underneath the search algorithm itself has sprung a vast area of research widely known as bio-inspired optimization [269]. Over the last decades, a manifold of behavioral patterns observed in biological systems have been emulated to yield intelligent algorithms capable of mimicking the learning and adaptation capabilities of such biological systems to address complex computational problems. Therefore, a bio-inspired meta-heuristic algorithm can be categorized as such if its main search strategy gets partially or fully inspired by biological phenomena, such as the evo-

lution of species, the echolocation of bats or the foraging behavior of ant colonies. A plethora of inspiring metaphors can be found nowadays in contributions dealing with new bio-inspired optimization algorithms, not without an ongoing controversy on the value of the metaphor itself for the novelty and scientific soundness of the reported methods [201].

Leaving such disputes aside, a research trend that has so far endured over the years is the hybridization of bio-inspired algorithms with problem-specific local search methods. The main reason behind this practice is to exploit the advantages of bio-inspired solvers, and to overcome their disadvantages when dealing with problems for which ad-hoc heuristics can be developed and inserted into the overall search process. Arguably, Memetic Algorithms [270] capitalize on this principle, with many application domains having so far harnessed this synergy between global and local search algorithms [271]. As highlighted in the prospective part of this survey, the incorporation of domain-based knowledge in the design of bio-inspired optimization algorithms can play a central role in the future when it comes to the intersection between Deep Learning and bio-inspired optimization.

## References

- [1] Y. LeCun, Y. Bengio, G. Hinton, Deep learning, *Nature* 521 (7553) (2015) 436–444.
- [2] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, L. D. Jackel, Back-propagation applied to handwritten zip code recognition, *Neural computation* 1 (4) (1989) 541–551.
- [3] G. E. Hinton, S. Osindero, Y.-W. Teh, A fast learning algorithm for deep belief nets, *Neural computation* 18 (7) (2006) 1527–1554.
- [4] T. Young, D. Hazarika, S. Poria, E. Cambria, Recent trends in deep learning based natural language processing, *IEEE Computational Intelligence Magazine* 13 (3) (2018) 55–75.
- [5] M. Kolb, Z.-H. Tan, J. Jensen, M. Kolb, Z.-H. Tan, J. Jensen, Speech intelligibility potential of general and specialized deep neural network based speech enhancement systems, *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)* 25 (1) (2017) 153–167.
- [6] Y. Zhang, W. Chan, N. Jaitly, Very deep convolutional networks for end-to-end speech recognition, in: *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2017, pp. 4845–4849.
- [7] S. Pal, Y. Dong, B. Thapa, N. V. Chawla, A. Swami, R. Ramanathan, Deep learning for network analysis: Problems, approaches and challenges, in: *MILCOM 2016-2016 IEEE Military Communications Conference*, IEEE, 2016, pp. 588–593.
- [8] S. Grigorescu, B. Trasnea, T. Cocias, G. Macesanu, A survey of deep learning techniques for autonomous driving, *Journal of Field Robotics* (2019).
- [9] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [10] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, Y. Bengio, Learning phrase representations using rnn encoder-decoder for statistical machine translation, *arXiv preprint arXiv:1406.1078* (2014).
- [11] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: *Advances in neural information processing systems*, 2014, pp. 2672–2680.

- [12] M. M. Najafabadi, F. Villanustre, T. M. Khoshgoftaar, N. Seliya, R. Wald, E. Muharemagic, Deep learning applications and challenges in big data analytics, *Journal of Big Data* 2 (1) (2015) 1.
- [13] K. Yu, L. Jia, Y. Chen, W. Xu, Deep learning: yesterday, today, and tomorrow, *Journal of computer Research and Development* 50 (9) (2013) 1799–1804.
- [14] S. Fong, S. Deb, X.-s. Yang, How meta-heuristic algorithms contribute to deep learning in the hype of big data analytics, in: *Progress in Intelligent Computing Techniques: Theory, Practice, and Applications*, Springer, 2018, pp. 3–25.
- [15] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins, R. Chatila, F. Herrera, Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI, *Information Fusion* 58 (2020) 82–115.
- [16] X. Yao, Y. Liu, A new evolutionary system for evolving artificial neural networks, *IEEE transactions on neural networks* 8 (3) (1997) 694–713.
- [17] K. O. Stanley, R. Miikkulainen, Evolving neural networks through augmenting topologies, *Evolutionary computation* 10 (2) (2002) 99–127.
- [18] K. O. Stanley, D. B. D’Ambrosio, J. Gauci, A hypercube-based encoding for evolving large-scale neural networks, *Artificial life* 15 (2) (2009) 185–212.
- [19] R. Miikkulainen, J. Z. Liang, E. Meyerson, A. Rawal, D. Fink, O. Francon, B. Raju, H. Shahrzad, A. Navruzyan, N. Duffy, B. Hodjat, Evolving deep neural networks, *ArXiv abs/1703.00548* (2017).
- [20] E. Real, S. Moore, A. Selle, S. Saxena, Y. L. Suematsu, J. Tan, Q. V. Le, A. Kurakin, Large-scale evolution of image classifiers, in: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, JMLR. org, 2017, pp. 2902–2911.
- [21] J. Muñoz-Ordóñez, C. Cobos, M. Mendoza, E. Herrera-Viedma, F. Herrera, S. Tabik, Framework for the training of deep neural networks in tensorflow using metaheuristics, in: *International Conference on Intelligent Data Engineering and Automated Learning*, Springer, 2018, pp. 801–811.
- [22] A. Martín, R. Lara-Cabrera, F. Fuentes-Hurtado, V. Naranjo, D. Camacho, Evodeep: a new evolutionary approach for automatic deep neural networks parametrisation, *Journal of Parallel and Distributed Computing* 117 (2018) 180–191.
- [23] F. Assunção, N. Lourenço, P. Machado, B. Ribeiro, Denser: deep evolutionary network structured representation, *Genetic Programming and Evolvable Machines* 20 (1) (2019) 5–35.
- [24] Google, Google Cloud AutoML.  
URL <https://cloud.google.com/automl/>
- [25] H. Jin, Q. Song, X. Hu, Auto-keras: An efficient neural architecture search system, in: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 1946–1956.
- [26] H. Pham, M. Y. Guan, B. Zoph, Q. V. Le, J. Dean, Efficient neural architecture search via parameter sharing, *arXiv preprint arXiv:1802.03268* (2018).
- [27] J. Liang, E. Meyerson, B. Hodjat, D. Fink, K. Mutch, R. Miikkulainen, Evolutionary neural automl for deep learning, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, 2019, pp. 401–409.



- [28] F. Charte, A. J. Rivera, F. Martínez, M. J. del Jesus, EvoAAA: An evolutionary methodology for automated neural autoencoder architecture search, *Integrated Computer-Aided Engineering (Preprint)* (2020) 1–21.
- [29] P. Molino, Y. Dudin, S. S. Miryala, Ludwig: a type-based declarative deep learning toolbox (2019). [arXiv:arXiv:1909.07930](https://arxiv.org/abs/1909.07930).
- [30] J. da Silveira Bohrer, B. I. Grisci, M. Dorn, Neuroevolution of neural network architectures using codeepeat and keras (2020). [arXiv:2002.04634](https://arxiv.org/abs/2002.04634).
- [31] B. Baker, O. Gupta, N. Naik, R. Raskar, Designing neural network architectures using reinforcement learning, *International Conference on Learning Representations* (2017).
- [32] J. Davison, Devol-deep neural network evolution (2017).
- [33] F. P. Such, V. Madhavan, E. Conti, J. Lehman, K. O. Stanley, J. Clune, Deep neuroevolution: Genetic algorithms are a competitive alternative for training deep neural networks for reinforcement learning, *arXiv preprint arXiv:1712.06567* (2017).
- [34] E. Conti, V. Madhavan, F. P. Such, J. Lehman, K. Stanley, J. Clune, Improving exploration in evolution strategies for deep reinforcement learning via a population of novelty-seeking agents, in: *Advances in neural information processing systems*, 2018, pp. 5027–5038.
- [35] M. Sukanuma, S. Shirakawa, T. Nagao, A genetic programming approach to designing convolutional neural network architectures, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, 2017, pp. 497–504.
- [36] C. Cortes, X. Gonzalvo, V. Kuznetsov, M. Mohri, S. Yang, Adanet: Adaptive structural learning of artificial neural networks, in: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, JMLR. org, 2017, pp. 874–883.
- [37] H. Mendoza, A. Klein, M. Feurer, J. T. Springenberg, M. Urban, M. Burkart, M. Dippel, M. Lindauer, F. Hutter, Towards automatically-tuned deep neural networks, in: F. Hutter, L. Kotthoff, J. Vanschoren (Eds.), *AutoML: Methods, Systems, Challenges*, Springer, 2018, Ch. 7, pp. 141–156, to appear.
- [38] L. Cardamone, D. Loiacono, P. L. Lanzi, Evolving competitive car controllers for racing games with neuroevolution, in: *Proceedings of the 11th Annual conference on Genetic and evolutionary computation*, ACM, 2009, pp. 1179–1186.
- [39] K. O. Stanley, B. D. Bryant, R. Miikkulainen, Real-time neuroevolution in the nero video game, *IEEE transactions on evolutionary computation* 9 (6) (2005) 653–668.
- [40] F. P. Such, V. Madhavan, E. Conti, J. Lehman, K. O. Stanley, J. Clune, Deep neuroevolution: Genetic algorithms are a competitive alternative for training deep neural networks for reinforcement learning, *arXiv preprint arXiv:1712.06567* (2017).
- [41] P. Verbancsics, J. Harguess, Generative neuroevolution for deep learning, *arXiv preprint arXiv:1312.5355* (2013).
- [42] X.-Z. Gao, J. Wang, J. M. Tanskanen, R. Bie, P. Guo, Bp neural networks with harmony search method-based training for epileptic eeg signal classification, in: *2012 Eighth International Conference on Computational Intelligence and Security*, IEEE, 2012, pp. 252–257.

- [43] J. P. Donate, X. Li, G. G. Sánchez, A. S. de Miguel, Time series forecasting by evolving artificial neural networks with genetic algorithms, differential evolution and estimation of distribution algorithm, *Neural Computing and Applications* 22 (1) (2013) 11–20.
- [44] X. Mao, A. Ter Mors, N. Roos, C. Witteveen, Using neuro-evolution in aircraft deicing scheduling, *Adaptive and Learning Agents and Multi-Agent Systems* (2007) 138–145.
- [45] G. Morse, K. O. Stanley, Simple evolutionary optimization can rival stochastic gradient descent in neural networks, in: *Proceedings of the Genetic and Evolutionary Computation Conference 2016*, ACM, 2016, pp. 477–484.
- [46] K. Mason, J. Duggan, E. Howley, Neural network topology and weight optimization through neuro differential evolution, in: *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, ACM, 2017, pp. 213–214.
- [47] V. K. Ojha, A. Abraham, V. Snášel, Metaheuristic design of feedforward neural networks: A review of two decades of research, *Engineering Applications of Artificial Intelligence* 60 (2017) 97–116.
- [48] A. Baldominos, Y. Saez, P. Isasi, On the automated, evolutionary design of neural networks: past, present, and future, *Neural Computing and Applications* (2019) 1–27.
- [49] H. Al-Sahaf, Y. Bi, Q. Chen, A. Lensen, Y. Mei, Y. Sun, B. Tran, B. Xue, M. Zhang, A survey on evolutionary machine learning, *Journal of the Royal Society of New Zealand* 49 (2) (2019) 205–228.
- [50] A. Darwish, A. E. Hassanien, S. Das, A survey of swarm and evolutionary computing approaches for deep learning, *Artificial Intelligence Review* 53 (3) (2020) 1767–1812.
- [51] H. Chiroma, A. Y. Gital, N. Rana, M. A. Shafii, A. N. Muhammad, A. Y. Umar, A. I. Abubakar, Nature inspired meta-heuristic algorithms for deep learning: Recent progress and novel perspective, in: *Science and Information Conference*, Springer, 2019, pp. 59–70.
- [52] H. Jin, Q. Song, X. Hu, Efficient neural architecture search with network morphism, *arXiv preprint arXiv:1806.10282* (2018).
- [53] E. Real, A. Aggarwal, Y. Huang, Q. V. Le, Regularized evolution for image classifier architecture search, *arXiv preprint arXiv:1802.01548* (2018).
- [54] Y.-H. Kim, B. Reddy, S. Yun, C. Seo, Nemo: Neuro-evolution with multiobjective optimization of deep neural network for speed and accuracy, in: *JMLR: Workshop and Conference Proceedings*, Vol. 1, 2017, pp. 1–8.
- [55] L. Xie, A. Yuille, Genetic cnn, in: *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1379–1388.
- [56] M. Suganuma, S. Shirakawa, T. Nagao, A genetic programming approach to designing convolutional neural network architectures, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, ACM, 2017, pp. 497–504.
- [57] Z. Lu, I. Whalen, V. Boddeti, Y. Dhebar, K. Deb, E. Goodman, W. Banzhaf, Nsga-net: a multi-objective genetic algorithm for neural architecture search, *arXiv preprint arXiv:1810.03522* (2018).
- [58] P. R. Lorenzo, J. Nalepa, Memetic evolution of deep neural networks, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, 2018, pp. 505–512.

- [59] Z. Chen, Y. Zhou, Z. Huang, Auto-creation of effective neural network architecture by evolutionary algorithm and resnet for image classification, in: 2019 IEEE International Conference on Systems, Man and Cybernetics (SMC), IEEE, 2019, pp. 3895–3900.
- [60] B. Evans, H. Al-Sahaf, B. Xue, M. Zhang, Evolutionary deep learning: A genetic programming approach to image classification, in: 2018 IEEE Congress on Evolutionary Computation (CEC), IEEE, 2018, pp. 1–6.
- [61] M. J. Shafiee, A. Mishra, A. Wong, Deep learning with darwin: Evolutionary synthesis of deep neural networks, *Neural Processing Letters* 48 (1) (2018) 603–613.
- [62] H. Zhu, Y. Jin, Real-time federated evolutionary neural architecture search, *arXiv preprint arXiv:2003.02793* (2020).
- [63] T. Desell, Large scale evolution of convolutional neural networks using volunteer computing, in: *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, 2017, pp. 127–128.
- [64] B. Wang, Y. Sun, B. Xue, M. Zhang, Evolving deep neural networks by multi-objective particle swarm optimization for image classification, *Genetic and Evolutionary Computation conference (GECCO 2019)* (2019).
- [65] B. Wang, B. Xue, M. Zhang, Particle swarm optimisation for evolving deep neural networks for image classification by evolving and stacking transferable blocks, *arXiv preprint arXiv:1907.12659* (2019).
- [66] B. Wang, Y. Sun, B. Xue, M. Zhang, A hybrid ga-pso method for evolving architecture and short connections of deep convolutional neural networks, in: *Pacific Rim International Conference on Artificial Intelligence*, Springer, 2019, pp. 650–663.
- [67] Y.-L. Hu, L. Chen, A nonlinear hybrid wind speed forecasting model using lstm network, hysteretic elm and differential evolution algorithm, *Energy conversion and management* 173 (2018) 123–142.
- [68] A. Rawal, R. Miikkulainen, From nodes to networks: Evolving recurrent neural networks, *arXiv preprint arXiv:1803.04439* (2018).
- [69] P. J. Angeline, G. M. Saunders, J. B. Pollack, An evolutionary algorithm that constructs recurrent neural networks, *IEEE Transactions on Neural Networks* 5 (1) (1994) 54–65.
- [70] T. Desell, S. Clachar, J. Higgins, B. Wild, Evolving deep recurrent neural networks using ant colony optimization, in: *European Conference on Evolutionary Computation in Combinatorial Optimization*, Springer, 2015, pp. 86–98.
- [71] C.-F. Juang, A hybrid of genetic algorithm and particle swarm optimization for recurrent network design, *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 34 (2) (2004) 997–1006.
- [72] F. Assuncao, D. Sereno, N. Lourenco, P. Machado, B. Ribeiro, Automatic evolution of autoencoders for compressed representations, in: 2018 IEEE Congress on Evolutionary Computation (CEC), IEEE, 2018, pp. 1–8.
- [73] S. Lander, Y. Shang, Evoae—a new evolutionary method for training autoencoders for deep learning networks, in: 2015 IEEE 39th Annual Computer Software and Applications Conference, Vol. 2, IEEE, 2015, pp. 790–795.

- [74] K. Liu, L. M. Zhang, Y. W. Sun, Deep boltzmann machines aided design based on genetic algorithms, in: *Applied Mechanics and Materials*, Vol. 568, Trans Tech Publ, 2014, pp. 848–851.
- [75] K. N. Mehta, Neuroevolutionary training of deep convolutional generative adversarial networks (2019).
- [76] V. Costa, N. Lourenço, P. Machado, Coevolution of generative adversarial networks, in: *International Conference on the Applications of Evolutionary Computation (Part of EvoStar)*, Springer, 2019, pp. 473–487.
- [77] V. Costa, N. Lourenço, J. Correia, P. Machado, Coegan: evaluating the coevolution effect in generative adversarial networks, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, 2019, pp. 374–382.
- [78] A. P. Poulsen, M. Thorhauge, M. H. Funch, S. Risi, DIne: A hybridization of deep learning and neuroevolution for visual control, in: *2017 IEEE Conference on Computational Intelligence and Games (CIG)*, IEEE, 2017, pp. 256–263.
- [79] S. Pham, K. Zhang, T. Phan, J. Ding, C. L. Dancy, Playing snes games with neuroevolution of augmenting topologies, in: *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [80] M. Hausknecht, J. Lehman, R. Miikkulainen, P. Stone, A neuroevolution approach to general atari game playing, *IEEE Transactions on Computational Intelligence and AI in Games* 6 (4) (2014) 355–366.
- [81] K. O. Stanley, R. Miikkulainen, Efficient reinforcement learning through evolving neural network topologies, in: *Proceedings of the 4th Annual Conference on Genetic and Evolutionary Computation*, Morgan Kaufmann Publishers Inc., 2002, pp. 569–577.
- [82] A. Etor, J. Ceberio, A. Pérez, E. Irurozki, Neural architecture search for time series classification, *The Genetic and Evolutionary Computation Conference*, 2020.
- [83] S. Fujino, N. Mori, K. Matsumoto, Deep convolutional networks for human sketches by means of the evolutionary deep learning, in: *2017 Joint 17th World Congress of International Fuzzy Systems Association and 9th International Conference on Soft Computing and Intelligent Systems (IFSA-SCIS)*, IEEE, 2017, pp. 1–5.
- [84] A. Baldominos, Y. Saez, P. Isasi, Evolutionary convolutional neural networks: An application to handwriting recognition, *Neurocomputing* 283 (2018) 38–52.
- [85] N. N. Ali Bakhshi, Stephan Chalup, Fast evolution of cnn architecture for image classification, in: N. N. Hitoshi Iba (Ed.), *Deep Neural Evolution*, Springer Singapore, 2020, Ch. 8, pp. 209–229.
- [86] E. Dufourq, B. A. Bassett, Eden: Evolutionary deep networks for efficient machine learning, in: *2017 Pattern Recognition Association of South Africa and Robotics and Mechatronics (PRASA-RobMech)*, IEEE, 2017, pp. 110–115.
- [87] R. Akut, S. Kulkarni, Neuroevolution: Using genetic algorithm for optimal design of deep learning models., in: *2019 IEEE International Conference on Electrical, Computer and Communication Technologies (ICECCT)*, IEEE, 2019, pp. 1–6.
- [88] F. Assunção, N. Lourenço, P. Machado, B. Ribeiro, Fast denser: Efficient deep neuroevolution, in: *European Conference on Genetic Programming*, Springer, 2019, pp. 197–212.

- [89] E. Bochinski, T. Senst, T. Sikora, Hyper-parameter optimization for convolutional neural network committees based on evolutionary algorithms, in: 2017 IEEE International Conference on Image Processing (ICIP), IEEE, 2017, pp. 3924–3928.
- [90] J. Prellberg, O. Kramer, Lamarckian evolution of convolutional neural networks, in: International Conference on Parallel Problem Solving from Nature, Springer, 2018, pp. 424–435.
- [91] Y. Sun, B. Xue, M. Zhang, G. G. Yen, Automatically evolving cnn architectures based on blocks, arXiv preprint arXiv:1810.11875 (2018).
- [92] S. Zhang, Y. Chen, X. Huang, Y. Cai, Text classification of public feedbacks using convolutional neural network based on differential evolution algorithm, International Journal of Computers Communications & Control 14 (1) (2019) 124–134.
- [93] R. Miiikkulainen, J. Liang, E. Meyerson, A. Rawal, D. Fink, O. Francon, B. Raju, H. Shahrzad, A. Navruzyan, N. Duffy, et al., Evolving deep neural networks, in: Artificial Intelligence in the Age of Neural Networks and Brain Computing, Elsevier, 2019, pp. 293–312.
- [94] T. Elsken, J.-H. Metzen, F. Hutter, Simple and efficient architecture search for convolutional neural networks, arXiv preprint arXiv:1711.04528 (2017).
- [95] B. Ma, X. Li, Y. Xia, Y. Zhang, Autonomous deep learning: A genetic dcnn designer for image classification, Neurocomputing 379 (2020) 152–161.
- [96] T. N. Masanori Sukanuma, Shinichi Shirakawa, Designing convolutional neural network architectures using cartesian genetic programming, in: N. N. Hitoshi Iba (Ed.), Deep Neural Evolution, Springer Singapore, 2020, Ch. 7, pp. 185–208.
- [97] X. Gu, Z. Meng, Y. Liang, D. Xu, H. Huang, X. Han, C. Wu, Esae: Evolutionary strategy-based architecture evolution, in: International Conference on Bio-Inspired Computing: Theories and Applications, Springer, 2019, pp. 193–208.
- [98] M. Sukanuma, M. Kobayashi, S. Shirakawa, T. Nagao, Evolution of deep convolutional neural networks using cartesian genetic programming, Evolutionary Computation 28 (1) (2020) 141–163.
- [99] H. Zhu, Y. Jin, Multi-objective evolutionary federated learning, IEEE transactions on neural networks and learning systems (2019).
- [100] M. Loni, S. Sinaei, A. Zoljodi, M. Daneshtalab, M. Sjödin, Deepmaker: A multi-objective optimization framework for deep neural networks in embedded systems, Microprocessors and Microsystems (2020) 102989.
- [101] Z. Lu, I. Whalen, Y. Dhebar, K. Deb, E. Goodman, W. Banzhaf, V. N. Boddeti, Multi-criterion evolutionary design of deep convolutional neural networks, arXiv preprint arXiv:1912.01369 (2019).
- [102] Y. Sun, B. Xue, M. Zhang, G. G. Yen, Evolving deep convolutional neural networks for image classification, IEEE Transactions on Evolutionary Computation (2019).
- [103] Y. Sun, B. Xue, M. Zhang, G. G. Yen, Completely automated cnn architecture design based on blocks, IEEE transactions on neural networks and learning systems 31 (4) (2019) 1242–1254.
- [104] H. Rakhshani, H. Ismail, L. Idoumghar, G. Forestier, J. Lepagnot, J. Weber, M. Brvilliers, P.-A. Muller, Neural architecture search for time series classification, 2020 International Joint Conference on Neural Networks (IJCNN), 2020.

- [105] Z. Lu, K. Deb, E. Goodman, W. Banzhaf, V. N. Boddeti, Nsganetv2: Evolutionary multi-objective surrogate-assisted neural architecture search, arXiv preprint arXiv:2007.10396 (2020).
- [106] B. Wang, Y. Sun, B. Xue, M. Zhang, Evolving deep convolutional neural networks by variable-length particle swarm optimization for image classification, in: 2018 IEEE Congress on Evolutionary Computation (CEC), IEEE, 2018, pp. 1–8.
- [107] M. Z. Bin Wang, Bing Xue, Particle swarm optimization for evolving deep convolutional neural networks for image classification: Single- and multi-objective approaches, in: N. N. Hitoshi Iba (Ed.), *Deep Neural Evolution*, Springer Singapore, 2020, Ch. 6, pp. 161–190.
- [108] B. Wang, Y. Sun, B. Xue, M. Zhang, A hybrid differential evolution approach to designing deep convolutional neural networks for image classification, in: *Australasian Joint Conference on Artificial Intelligence*, Springer, 2018, pp. 237–250.
- [109] L. Peng, S. Liu, R. Liu, L. Wang, Effective long short-term memory with differential evolution algorithm for electricity price prediction, *Energy* 162 (2018) 1301–1314.
- [110] B. Nakisa, M. N. Rastgoo, A. Rakotonirainy, F. Maire, V. Chandran, Long short term memory hyperparameter optimization for a neural network based emotion recognition framework, *IEEE Access* 6 (2018) 49325–49338.
- [111] A. Rawal, R. Miikkulainen, Evolving deep lstm-based memory networks using an information maximization objective, in: *Proceedings of the Genetic and Evolutionary Computation Conference 2016*, 2016, pp. 501–508.
- [112] V. C. Lobo Neto, L. A. Passos, J. P. Papa, Evolving long short-term memory networks, in: V. V. Krzhizhanovskaya, G. Závodszy, M. H. Lees, J. J. Dongarra, P. M. A. Slood, S. Brissos, J. Teixeira (Eds.), *Computational Science – ICCS 2020*, Springer International Publishing, 2020, pp. 337–350.
- [113] P. Bento, J. Pombo, S. Mariano, M. do Rosário Calado, Short-term load forecasting using optimized lstm networks via improved bat algorithm, in: *2018 International Conference on Intelligent Systems (IS)*, IEEE, 2018, pp. 351–357.
- [114] M. van Knippenberg, V. Menkovski, S. Consoli, Evolutionary construction of convolutional neural networks, in: *International Conference on Machine Learning, Optimization, and Data Science*, Springer, 2018, pp. 293–304.
- [115] F. Charte, A. J. Rivera, F. Martínez, M. J. del Jesus, Automating autoencoder architecture configuration: An evolutionary approach, in: *International Work-Conference on the Interplay Between Natural and Artificial Computation*, Springer, 2019, pp. 339–349.
- [116] Y. Sun, B. Xue, M. Zhang, G. G. Yen, A particle swarm optimization-based flexible convolutional autoencoder for image classification, *IEEE transactions on neural networks and learning systems* (2018).
- [117] J. P. Papa, G. H. Rosa, A. N. Marana, W. Scheirer, D. D. Cox, Model selection for discriminative restricted boltzmann machines through meta-heuristic techniques, *Journal of Computational Science* 9 (2015) 14–18.
- [118] L. A. Passos, J. P. Papa, A metaheuristic-driven approach to fine-tune deep boltzmann machines, *Applied Soft Computing* (2019) 105717.

- [119] T. Kuremoto, S. Kimura, K. Kobayashi, M. Obayashi, Time series forecasting using restricted boltzmann machine, in: *International Conference on Intelligent Computing*, Springer, 2012, pp. 17–22.
- [120] L. A. Passos, D. R. Rodrigues, J. P. Papa, Fine tuning deep boltzmann machines through meta-heuristic approaches, in: *2018 IEEE 12th International Symposium on Applied Computational Intelligence and Informatics (SACI)*, IEEE, 2018, pp. 000419–000424.
- [121] J. Wang, K. Wang, Y. Wang, Z. Huang, R. Xue, Deep boltzmann machine based condition prediction for smart manufacturing, *Journal of Ambient Intelligence and Humanized Computing* 10 (3) (2019) 851–861.
- [122] N. R. Sabar, A. Turky, A. Song, A. Sattar, Optimising deep belief networks by hyper-heuristic approach, in: *2017 IEEE Congress on Evolutionary Computation (CEC)*, IEEE, 2017, pp. 2738–2745.
- [123] D. Hossain, G. Capi, M. Jindai, Evolution of deep belief neural network parameters for robot object recognition and grasping, *Procedia Computer Science* 105 (C) (2017) 153–158.
- [124] G. H. de Rosa, J. P. Papa, Soft-tempering deep belief networks parameters through genetic programming (2019).
- [125] D. R. M. R. Leandro Aparecido Passos, Gustavo Henrique de Rosa, J. P. Papa, On the assessment of nature-inspired meta-heuristic optimization techniques to fine-tune deep belief networks, in: N. N. Hitoshi Iba (Ed.), *Deep Neural Evolution*, Springer Singapore, 2020, Ch. 3, pp. 67–96.
- [126] M.-H. Horng, Fine-tuning parameters of deep belief networks using artificial bee colony algorithm, *DEStech Transactions on Computer Science and Engineering* (2017).
- [127] L. Li, L. Qin, X. Qu, J. Zhang, Y. Wang, B. Ran, Day-ahead traffic flow forecasting based on a deep belief network optimized by the multi-objective particle swarm algorithm, *Knowledge-Based Systems* (2019).
- [128] S. Goudarzi, M. Kama, M. Anisi, S. Soleymani, F. Doctor, Self-organizing traffic flow prediction with an optimized deep belief network for internet of vehicles, *Sensors* 18 (10) (2018) 3459.
- [129] M. Ma, C. Sun, X. Chen, Discriminative deep belief networks with ant colony optimization for health status assessment of machine, *IEEE Transactions on Instrumentation and Measurement* 66 (12) (2017) 3115–3125.
- [130] M. O. K. K. S. M. Takashi Kuremoto, Takaomi Hirata, Search heuristics for the optimization of dbn for time series forecasting, in: N. N. Hitoshi Iba (Ed.), *Deep Neural Evolution*, Springer Singapore, 2020, Ch. 5, pp. 131–152.
- [131] D. Rodrigues, X.-S. Yang, J. Papa, Fine-tuning deep belief networks using cuckoo search, in: *Bio-Inspired Computation and Applications in Image Processing*, Elsevier, 2016, pp. 47–59.
- [132] U. Garciarena, R. Santana, A. Mendiburu, Evolved gans for generating pareto set approximations, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, ACM, 2018, pp. 434–441.
- [133] Y. Lu, B. Kakillioglu, S. Velipasalar, Autonomously and simultaneously refining deep neural network parameters by a bi-generative adversarial network aided genetic algorithm, *arXiv preprint arXiv:1809.10244* (2018).

- [134] A. Dahou, M. A. Elaziz, J. Zhou, S. Xiong, Arabic sentiment classification using convolutional neural network and differential evolution algorithm, *Computational Intelligence and Neuroscience* 2019 (2019).
- [135] S. R. Young, D. C. Rose, T. P. Karnowski, S.-H. Lim, R. M. Patton, Optimizing deep learning hyper-parameters through an evolutionary algorithm, in: *Proceedings of the Workshop on Machine Learning in High-Performance Computing Environments*, ACM, 2015, p. 4.
- [136] G. Bingham, W. Macke, R. Miikkulainen, Evolutionary optimization of deep learning activation functions, *arXiv preprint arXiv:2002.07224* (2020).
- [137] J. Kim, S. Cho, Evolutionary optimization of hyperparameters in deep learning models, in: *2019 IEEE Congress on Evolutionary Computation (CEC)*, 2019, pp. 831–837.
- [138] S. Gonzalez, R. Miikkulainen, Improved training speed, accuracy, and data utilization through loss function optimization (2020). *arXiv:1905.11528*.
- [139] P. R. Lorenzo, J. Nalepa, M. Kawulok, L. S. Ramos, J. R. Pastor, Particle swarm optimization for hyper-parameter selection in deep neural networks, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, ACM, 2017, pp. 481–488.
- [140] T. Yamasaki, T. Honma, K. Aizawa, Efficient optimization of convolutional neural networks using particle swarm optimization, in: *2017 IEEE Third International Conference on Multimedia Big Data (BigMM)*, IEEE, 2017, pp. 70–73.
- [141] P. Ortego, A. Diez-Olivan, J. Del Ser, F. Veiga, M. Penalva, B. Sierra, Evolutionary lstm-fcn networks for pattern classification in industrial processes, *Swarm and Evolutionary Computation* 54 (2020) 100650.
- [142] A. ElSaid, F. El Jamiy, J. Higgins, B. Wild, T. Desell, Optimizing long short-term memory recurrent neural networks using ant colony optimization to predict turbine engine vibration, *Applied Soft Computing* 73 (2018) 969–991.
- [143] A. ElSaid, F. E. Jamiy, J. Higgins, B. Wild, T. Desell, Using ant colony optimization to optimize long short-term memory recurrent neural networks, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, ACM, 2018, pp. 13–20.
- [144] H. Wang, X. Yan, Optimizing the echo state network with a binary particle swarm optimization algorithm, *Knowledge-Based Systems* 86 (2015) 182–193.
- [145] T. Silhan, S. Oehmcke, O. Kramer, Evolution of stacked autoencoders, in: *2019 IEEE Congress on Evolutionary Computation (CEC)*, IEEE, 2019, pp. 823–830.
- [146] J. P. Papa, W. Scheirer, D. D. Cox, Fine-tuning deep belief networks using harmony search, *Applied Soft Computing* 46 (2016) 875–885.
- [147] J. P. Papa, G. H. Rosa, D. R. Pereira, X.-S. Yang, Quaternion-based deep belief networks fine-tuning, *Applied Soft Computing* 60 (2017) 328–335.
- [148] G. Rosa, J. Papa, K. Costa, L. Passos, C. Pereira, X.-S. Yang, Learning parameters in deep belief networks through firefly algorithm, in: *IAPR Workshop on Artificial Neural Networks in Pattern Recognition*, Springer, 2016, pp. 138–149.



- [149] M. ul Hassan, N. R. Sabar, A. Song, Optimising deep learning by hyper-heuristic approach for classifying good quality images, in: *International Conference on Computational Science*, Springer, 2018, pp. 528–539.
- [150] C. R. Pereira, D. R. Pereira, J. P. Papa, G. H. Rosa, X.-S. Yang, Convolutional neural networks applied for parkinsons disease identification, in: *Machine Learning for Health Informatics*, Springer, 2016, pp. 377–390.
- [151] G. H. De Rosa, J. P. Papa, X.-S. Yang, Handling dropout probability estimation in convolution neural networks using meta-heuristics, *Soft Computing* (2018) 1–10.
- [152] B. Guo, J. Hu, W. Wu, Q. Peng, F. Wu, The tabu\_genetic algorithm: A novel method for hyper-parameter optimization of learning algorithms, *Electronics* 8 (5) (2019) 579.
- [153] A. R. Ismail, O. A. Mohammad, et al., Evolutionary deep belief networks with bootstrap sampling for imbalanced class datasets, *International Journal of Advances in Intelligent Informatics* 5 (2) (2019) 123–136.
- [154] M. Jaderberg, V. Dalibard, S. Osindero, W. M. Czarnecki, J. Donahue, A. Razavi, O. Vinyals, T. Green, I. Dunning, K. Simonyan, et al., Population based training of neural networks, arXiv preprint arXiv:1711.09846 (2017).
- [155] K. Pawelczyk, M. Kawulok, J. Nalepa, Genetically-trained deep neural networks, in: *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, 2018, pp. 63–64.
- [156] L. Rere, M. I. Fanany, A. M. Arymurthy, Metaheuristic algorithms for convolution neural network, *Computational intelligence and neuroscience* 2016 (2016).
- [157] L.-O. Fedorovici, R.-E. Precup, F. Dragan, R.-C. David, C. Purcaru, Embedding gravitational search algorithms in convolutional neural networks for ocr applications, in: *2012 7th IEEE International Symposium on Applied Computational Intelligence and Informatics (SACI)*, IEEE, 2012, pp. 125–130.
- [158] A. Martn Garca, V. Vargas Yun, P. A. Gutierrez, D. Camacho, C. Martnez, Optimising convolutional neural networks using a hybrid statistically-driven coral reef optimisation algorithm, *Applied Soft Computing* 90 (2020). doi:10.1016/j.asoc.2020.106144.
- [159] J. Zhang, F. B. Gouza, Gadam: genetic-evolutionary adam for deep neural network optimization, arXiv preprint arXiv:1805.07500 (2018).
- [160] X. Cui, W. Zhang, Z. Tüske, M. Picheny, Evolutionary stochastic gradient descent for optimization of deep neural networks, in: *Advances in neural information processing systems*, 2018, pp. 6048–6058.
- [161] D. Zang, J. Ding, J. Cheng, D. Zhang, K. Tang, A hybrid learning algorithm for the optimization of convolutional neural network, in: *International Conference on Intelligent Computing*, Springer, 2017, pp. 694–705.
- [162] A. Banharsakun, Towards improving the convolutional neural networks for deep learning using the distributed artificial bee colony method, *International Journal of Machine Learning and Cybernetics* (2018) 1–11.
- [163] M. H. Khalifa, M. Ammar, W. Ouarda, A. M. Alimi, Particle swarm optimization for deep learning of convolution neural network, in: *2017 Sudan Conference on Computer Science and Information Technology (SCCSIT)*, IEEE, 2017, pp. 1–5.

- [164] Y. Li, Z. Zhu, D. Kong, H. Han, Y. Zhao, Ea-lstm: Evolutionary attention-based lstm for time series prediction, *Knowledge-Based Systems* 181 (2019) 104785.
- [165] N. M. Nawi, A. Khan, M. Rehman, H. Chiroma, T. Herawan, Weight optimization in recurrent neural networks with hybrid metaheuristic cuckoo search techniques for data classification, *Mathematical Problems in Engineering* 2015 (2015).
- [166] S. Alvernaz, J. Togelius, Autoencoder-augmented neuroevolution for visual doom playing, in: *2017 IEEE Conference on Computational Intelligence and Games (CIG)*, IEEE, 2017, pp. 1–8.
- [167] O. E. David, I. Greental, Genetic algorithms for evolving deep neural networks, in: *Proceedings of the Companion Publication of the 2014 Annual Conference on Genetic and Evolutionary Computation*, 2014, pp. 1451–1452.
- [168] E. Levy, O. E. David, N. S. Netanyahu, Genetic algorithms and deep learning for automatic painter classification, in: *proceedings of the 2014 Annual Conference on Genetic and Evolutionary Computation*, 2014, pp. 1143–1150.
- [169] A. Al-Dujaili, T. Schmiedlechner, U.-M. O’Reilly, et al., Towards distributed coevolutionary gans, *arXiv preprint arXiv:1807.08194* (2018).
- [170] S. Khadka, K. Tumer, Evolutionary reinforcement learning, *arXiv preprint arXiv:1805.07917* (2018).
- [171] S. Khadka, S. Majumdar, T. Nassar, Z. Dwiell, E. Tumer, S. Miret, Y. Liu, K. Tumer, Collaborative evolutionary reinforcement learning, *arXiv preprint arXiv:1905.00976* (2019).
- [172] S. Khadka, K. Tumer, Evolution-guided policy gradient in reinforcement learning, in: S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, R. Garnett (Eds.), *Advances in Neural Information Processing Systems* 31, Curran Associates, Inc., 2018, pp. 1188–1200.
- [173] J. Koutník, J. Schmidhuber, F. Gomez, Evolving deep unsupervised convolutional networks for vision-based reinforcement learning, in: *Proceedings of the 2014 Annual Conference on Genetic and Evolutionary Computation*, ACM, 2014, pp. 541–548.
- [174] L. R. Rere, M. I. Fanany, A. M. Arymurthy, Simulated annealing algorithm for deep learning, *Procedia Computer Science* 72 (2015) 137–144.
- [175] C. A. de Pinho Pinheiro, N. Nedjah, L. de Macedo Mourelle, Detection and classification of pulmonary nodules using deep learning and swarm intelligence, *Multimedia Tools and Applications* (2019) 1–29.
- [176] V. Ayumi, L. R. Rere, M. I. Fanany, A. M. Arymurthy, Optimization of convolutional neural network using microcanonical annealing algorithm, in: *2016 International Conference on Advanced Computer Science and Information Systems (ICACSIS)*, IEEE, 2016, pp. 506–511.
- [177] G. Rosa, J. Papa, A. Marana, W. Scheirer, D. Cox, Fine-tuning convolutional neural networks using harmony search, in: *Iberoamerican Congress on Pattern Recognition*, Springer, 2015, pp. 683–690.
- [178] S. Risi, K. O. Stanley, Deep neuroevolution of recurrent and discrete world models, *arXiv preprint arXiv:1906.08857* (2019).
- [179] T. A. Rashid, P. Fattah, D. K. Awla, Using accuracy measure for improving the training of lstm with metaheuristic algorithms, *Procedia Computer Science* 140 (2018) 324–333.

- [180] T. A. Rashid, M. K. Hassan, M. Mohammadi, K. Fraser, Improvement of variant adaptable lstm trained with metaheuristic algorithms for healthcare analysis, in: *Advanced Classification Techniques for Healthcare Analysis*, IGI Global, 2019, pp. 111–131.
- [181] N. Van Hoorn, J. Togelius, J. Schmidhuber, Hierarchical controller learning in a first-person shooter, in: *2009 IEEE symposium on computational intelligence and games*, IEEE, 2009, pp. 294–301.
- [182] C. A. Duchanoy, M. A. Moreno-Armendáriz, L. Urbina, C. A. Cruz-Villar, H. Calvo, J. d. J. Rubio, A novel recurrent neural network soft sensor via a differential evolution training algorithm for the tire contact patch, *Neurocomputing* 235 (2017) 71–82.
- [183] B. Jana, S. Mitra, S. Acharyaa, Reconstruction of gene regulatory network using recurrent neural network model: A harmony search approach, in: *Soft Computing and Signal Processing*, Springer, 2019, pp. 129–138.
- [184] S. Biswas, S. Acharyya, A bi-objective rnn model to reconstruct gene regulatory network: a modified multi-objective simulated annealing approach, *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)* 15 (6) (2018) 2053–2059.
- [185] A. M. Ibrahim, N. H. El-Amary, Particle swarm optimization trained recurrent neural network for voltage instability prediction, *Journal of Electrical Systems and Information Technology* 5 (2) (2018) 216–228.
- [186] H. Hisashi, Deep boltzmann machine for evolutionary agents of mario ai, in: *2014 IEEE Congress on Evolutionary Computation (CEC)*, 2014, pp. 36–41.
- [187] N. M. Nawi, A. Khan, M. Rehman, H. Chiroma, T. Herawan, Weight optimization in recurrent neural networks with hybrid metaheuristic cuckoo search techniques for data classification, *Mathematical Problems in Engineering* 2015 (2015).
- [188] C.-F. Juang, Y.-C. Chang, I.-F. Chung, Optimization of recurrent neural networks using evolutionary group-based particle swarm optimization for hexapod robot gait generation, *Hybrid Metaheuristics: Research And Applications* 84 (2018) 227.
- [189] Q. Song, Y.-J. Zheng, Y. Xue, W.-G. Sheng, M.-R. Zhao, An evolutionary deep neural network for predicting morbidity of gastrointestinal infections by food contamination, *Neurocomputing* 226 (2017) 16–22.
- [190] C. Wang, C. Xu, X. Yao, D. Tao, Evolutionary generative adversarial networks, *IEEE Transactions on Evolutionary Computation* 23 (6) (2019) 921–934.
- [191] J. Toutouh, E. Hemberg, U.-M. O’Reilly, Spatial evolutionary generative adversarial networks, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, 2019, pp. 472–480.
- [192] C. Igel, Neuroevolution for reinforcement learning using evolution strategies, in: *The 2003 Congress on Evolutionary Computation*, 2003. CEC’03., Vol. 4, IEEE, 2003, pp. 2588–2595.
- [193] A. D. Martinez, E. Osaba, J. Del Ser, F. Herrera, Simultaneously evolving deep reinforcement learning models using multifactorial optimization, *arXiv preprint arXiv:2002.12133* (2020).
- [194] K. Mason, J. Duggan, E. Howley, Maze navigation using neural networks evolved with novelty search and differential evolution, in: *Adaptive and Learning Agents Workshop (at ICML-AAMAS 2018)*, 2018.

- [195] P. Chrabaszcz, I. Loshchilov, F. Hutter, Back to basics: Benchmarking canonical evolution strategies for playing atari, arXiv preprint arXiv:1802.08842 (2018).
- [196] F. Gomez, J. Schmidhuber, R. Miikkulainen, Accelerated neural evolution through cooperatively coevolved synapses, *Journal of Machine Learning Research* 9 (May) (2008) 937–965.
- [197] S. Tabik, R. F. Alvear-Sandoval, M. M. Ruiz, J.-L. Sancho-Gómez, A. R. Figueiras-Vidal, F. Herrera, Mnist-net10: A heterogeneous deep networks fusion based on the degree of certainty to reach 0.1% error rate. ensembles overview and proposal, *Information Fusion* 62 (2020) 1–8.
- [198] A. LaTorre, D. Molina, E. Osaba, J. Del Ser, F. Herrera, Fairness in bio-inspired optimization research: A prescription of methodological guidelines for comparing meta-heuristics, arXiv preprint arXiv:2004.09969 (2020).
- [199] J. Demšar, Statistical comparisons of classifiers over multiple data sets, *Journal of Machine learning research* 7 (Jan) (2006) 1–30.
- [200] J. Carrasco, S. García, M. Rueda, S. Das, F. Herrera, Recent trends in the use of statistical tests for comparing swarm and evolutionary computing algorithms: Practical guidelines and a critical review, *Swarm and Evolutionary Computation* 54 (2020) 100665.
- [201] J. Del Ser, E. Osaba, D. Molina, X.-S. Yang, S. Salcedo-Sanz, D. Camacho, S. Das, P. N. Suganthan, C. A. C. Coello, F. Herrera, Bio-inspired computation: Where we stand and what’s next, *Swarm and Evolutionary Computation* 48 (2019) 220–250.
- [202] L. Moroney, Horses or humans dataset (2019).  
URL <http://laurencemoroney.com/horses-or-humans-dataset>
- [203] J. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired image-to-image translation using cycle-consistent adversarial networks, *CoRR* abs/1703.10593 (2017). arXiv:1703.10593.  
URL <http://arxiv.org/abs/1703.10593>
- [204] Y. LeCun, C. Cortes, MNIST handwritten digit database (2010).
- [205] A. Krizhevsky, V. Nair, G. Hinton, Cifar-10 (canadian institute for advanced research).  
URL <http://www.cs.toronto.edu/~kriz/cifar.html>
- [206] D. Molina, A. LaTorre, F. Herrera, Shade with iterative local search for large-scale global optimization, in: *2018 IEEE Congress on Evolutionary Computation (CEC)*, IEEE, 2018, pp. 1–8.
- [207] T. Mantecón, C. R. del Blanco, F. Jaureguizar, N. García, Hand gesture recognition using infrared imagery provided by leap motion controller, 2016.
- [208] Blood cell classification dataset, [https://github.com/Shenggan/BCCD\\_Dataset](https://github.com/Shenggan/BCCD_Dataset) (2019).
- [209] H. Xiao, K. Rasul, R. Vollgraf, Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, arXiv preprint arXiv:1708.07747 (2017).
- [210] J. Stallkamp, M. Schlipsing, J. Salmen, C. Igel, Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition, *Neural Networks* (0) (2012) –.
- [211] M. Essaid, L. Idoumghar, J. Lepagnot, M. Brévilliers, Gpu parallelization strategies for meta-heuristics: a survey, *International Journal of Parallel, Emergent and Distributed Systems* 34 (5) (2019) 497–522.

- [212] Y. Tan, K. Ding, A survey on gpu-based implementation of swarm intelligence algorithms, *IEEE transactions on cybernetics* 46 (9) (2015) 2028–2041.
- [213] G. Schryen, Parallel computational optimization in operations research: A new integrative framework, literature review and research directions, *European Journal of Operational Research* (2019).
- [214] A. Benitez-Hidalgo, A. J. Nebro, J. Garcia-Nieto, I. Oregi, J. Del Ser, jmetalpy: A python framework for multi-objective optimization with metaheuristics, *Swarm and Evolutionary Computation* 51 (2019) 100598.
- [215] Y. S. Nashed, R. Ugolotti, P. Mesejo, S. Cagnoni, libcudaoptimize: an open source library of gpu-based metaheuristics, in: *Proceedings of the 14th annual conference companion on Genetic and evolutionary computation*, 2012, pp. 117–124.
- [216] M. e. a. Abadi, Tensorflow: A system for large-scale machine learning, in: *12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16)*, 2016, pp. 265–283.
- [217] P. Balaprakash, M. Salim, T. Uram, V. Vishwanath, S. Wild, Deephyper: Asynchronous hyperparameter search for deep neural networks, in: *2018 IEEE 25th International Conference on High Performance Computing (HiPC)*, IEEE, 2018, pp. 42–51.
- [218] S. Mahdavi, M. E. Shiri, S. Rahnamayan, Metaheuristics in large-scale global continues optimization: A survey, *Information Sciences* 295 (2015) 407–428.
- [219] J.-H. Yi, L.-N. Xing, G.-G. Wang, J. Dong, A. V. Vasilakos, A. H. Alavi, L. Wang, Behavior of crossover operators in nsga-iii for large-scale optimization problems, *Information Sciences* 509 (2020) 470–487.
- [220] J. Liang, E. Meyerson, R. Miikkulainen, Evolutionary architecture search for deep multitask networks, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, 2018, pp. 466–473.
- [221] Y. Li, S. Fang, X. Bai, L. Jiao, N. Marturi, Parallel design of sparse deep belief network with multi-objective optimization, *Information Sciences* (2020).
- [222] Y.-S. Ong, A. Gupta, Evolutionary multitasking: a computer science view of cognitive multitasking, *Cognitive Computation* 8 (2) (2016) 125–142.
- [223] R. Chandra, A. Gupta, Y.-S. Ong, C.-K. Goh, Evolutionary multi-task learning for modular knowledge representation in neural networks, *Neural Processing Letters* 47 (3) (2018) 993–1009.
- [224] A. Gupta, Y.-S. Ong, L. Feng, K. C. Tan, Multiobjective multifactorial optimization in evolutionary multitasking, *IEEE transactions on cybernetics* 47 (7) (2016) 1652–1665.
- [225] S. Yao, Z. Dong, X. Wang, L. Ren, A multiobjective multifactorial optimization algorithm based on decomposition and dynamic resource allocation strategy, *Information Sciences* 511 (2020) 18–35.
- [226] K. K. Bali, A. Gupta, Y.-S. Ong, P. S. Tan, Cognizant multitasking in multiobjective multifactorial evolution: Mo-mfea-ii, *IEEE Transactions on Cybernetics* (2020).
- [227] M. Lukoševičius, H. Jaeger, Reservoir computing approaches to recurrent neural network training, *Computer Science Review* 3 (3) (2009) 127–149.
- [228] M. Dale, Neuroevolution of hierarchical reservoir computers, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, 2018, pp. 410–417.

- [229] Y. Zhou, Y. Jin, J. Ding, Evolutionary optimization of liquid state machines for robust learning, in: *International Symposium on Neural Networks*, Springer, 2019, pp. 389–398.
- [230] Y. Zhou, Y. Jin, J. Ding, Surrogate-assisted evolutionary search of spiking neural architectures in liquid state machines, *Neurocomputing* (2020).
- [231] K. Liu, J. Zhang, Nonlinear process modelling using echo state networks optimised by covariance matrix adaption evolutionary strategy, *Computers & Chemical Engineering* (2020) 106730.
- [232] C. Gallicchio, A. Micheli, L. Pedrelli, Deep reservoir computing: A critical experimental analysis, *Neurocomputing* 268 (2017) 87–99.
- [233] R. A. Vazquez, Training spiking neural models using cuckoo search algorithm, in: *2011 IEEE Congress of Evolutionary Computation (CEC)*, IEEE, 2011, pp. 679–686.
- [234] C. D. Schuman, J. S. Plank, A. Disney, J. Reynolds, An evolutionary optimization framework for neural networks and neuromorphic architectures, in: *2016 International Joint Conference on Neural Networks (IJCNN)*, IEEE, 2016, pp. 145–154.
- [235] R. A. Vazquez, B. A. Garro, Training spiking neural models using artificial bee colony, *Computational intelligence and neuroscience 2015* (2015).
- [236] R. Carino-Escobar, J. Cantillo-Negrete, R. A. Vazquez, J. Gutierrez-Martinez, Spiking neural networks trained with particle swarm optimization for motor imagery classification, in: *International Conference on Swarm Intelligence*, Springer, 2016, pp. 245–252.
- [237] X. Wang, X. Lin, X. Dang, Supervised learning in spiking neural networks: A review of algorithms and evaluations, *Neural Networks* (2020).
- [238] A. Baldominos, Y. Saez, P. Isasi, Hybridizing evolutionary computation and deep neural networks: an approach to handwriting recognition using committees and transfer learning, *Complexity* 2019 (2019).
- [239] A. D. Martinez, E. Osaba, I. Oregi, I. Fister, I. Fister, J. D. Ser, Hybridizing differential evolution and novelty search for multimodal optimization problems, in: *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, 2019, pp. 1980–1989.
- [240] S. Kornblith, M. Norouzi, H. Lee, G. Hinton, Similarity of neural network representations revisited, *arXiv preprint arXiv:1905.00414* (2019).
- [241] A. Pourchot, O. Sigaud, Cem-rl: Combining evolutionary and gradient-based methods for policy search, *arXiv preprint arXiv:1810.01222* (2018).
- [242] K. Maziarz, M. Tan, A. Khorlin, K.-Y. S. Chang, S. Jastrzebski, Q. de Laroussilhe, A. Gesmundo, Evolutionary-neural hybrid agents for architecture search, *arXiv preprint arXiv:1811.09828* (2018).
- [243] D. Blalock, J. J. G. Ortiz, J. Frankle, J. Gutttag, What is the state of neural network pruning?, *arXiv preprint arXiv:2003.03033* (2020).
- [244] A. Labach, H. Salehinejad, S. Valaee, Survey of dropout methods for deep neural networks, *arXiv preprint arXiv:1904.13310* (2019).
- [245] Z. Wang, F. Li, G. Shi, X. Xie, F. Wang, Network pruning using sparse learning and genetic algorithm, *Neurocomputing* (2020).

- [246] M. Mohammadi, A. Al-Fuqaha, S. Sorour, M. Guizani, Deep learning for iot big data and streaming analytics: A survey, *IEEE Communications Surveys & Tutorials* 20 (4) (2018) 2923–2960.
- [247] Q. Yang, Y. Liu, T. Chen, Y. Tong, Federated machine learning: Concept and applications, *ACM Transactions on Intelligent Systems and Technology (TIST)* 10 (2) (2019) 1–19.
- [248] J. Chen, X. Ran, Deep learning with edge computing: A review, *Proceedings of the IEEE* 107 (8) (2019) 1655–1674.
- [249] E. García-Martín, C. F. Rodrigues, G. Riley, H. Grahn, Estimation of energy consumption in machine learning, *Journal of Parallel and Distributed Computing* 134 (2019) 75–88.
- [250] M. Nasr, R. Shokri, A. Houmansadr, Comprehensive privacy analysis of deep learning: Stand-alone and federated learning under passive and active white-box inference attacks, *arXiv preprint arXiv:1812.00910* (2018).
- [251] N. Rodríguez-Barroso, G. Stipcich, D. Jiménez-López, J. A. Ruiz-Millán, E. Martínez-Cámara, G. González-Seco, M. Luzón, M. Á. Veganzones, F. Herrera, Federated learning and differential privacy: Software tools analysis, the sherpa. ai fl framework and methodological guidelines for preserving data privacy, *Information Fusion*, to appear (2020).
- [252] A. N. Bhagoji, S. Chakraborty, P. Mittal, S. Calo, Analyzing federated learning through an adversarial lens, in: *International Conference on Machine Learning*, 2019, pp. 634–643.
- [253] C. A. C. Coello, G. B. Lamont, D. A. Van Veldhuizen, et al., *Evolutionary algorithms for solving multi-objective problems*, Vol. 5, Springer, 2007.
- [254] E. Mezura-Montes, C. A. C. Coello, Constraint-handling in nature-inspired numerical optimization: past, present and future, *Swarm and Evolutionary Computation* 1 (4) (2011) 173–194.
- [255] K. Weiss, T. M. Khoshgoftaar, D. Wang, A survey of transfer learning, *Journal of Big data* 3 (1) (2016) 9.
- [256] C. Enroth-Cugell, J. G. Robson, The contrast sensitivity of retinal ganglion cells of the cat, *The Journal of physiology* 187 (3) (1966) 517–552.
- [257] S. Hochstein, R. Shapley, Quantitative analysis of retinal ganglion cell classifications., *The Journal of physiology* 262 (2) (1976) 237–264.
- [258] D. Molina, J. Poyatos, J. Del Ser, S. García, A. Hussain, F. Herrera, Comprehensive taxonomies of nature- and bio-inspired optimization: Inspiration versus algorithmic behavior, critical analysis and recommendations, *Cognitive Computation* (2020).
- [259] I. BoussaïD, J. Lepagnot, P. Siarry, A survey on optimization metaheuristics, *Information sciences* 237 (2013) 82–117.
- [260] D. Molina, J. Poyatos, J. Del Ser, S. García, A. Hussain, F. Herrera, Comprehensive taxonomies of nature-and bio-inspired optimization: Inspiration versus algorithmic behavior, critical analysis and recommendations, *arXiv preprint arXiv:2002.08136* (2020).
- [261] F. Glover, M. Laguna, Tabu search, in: *Handbook of combinatorial optimization*, Springer, 1998, pp. 2093–2229.
- [262] S. Kirkpatrick, C. D. Gelatt, M. P. Vecchi, Optimization by simulated annealing, *science* 220 (4598) (1983) 671–680.

- [263] D. Goldberg, Genetic algorithms in search, optimization, and machine learning, Addison-Wesley Professional, 1989.
- [264] K. De Jong, Analysis of the behavior of a class of genetic adaptive systems, Ph.D. thesis, University of Michigan, Michigan, USA (1975).
- [265] M. Dorigo, M. Birattari, Ant colony optimization, Springer, 2010.
- [266] J. Kennedy, R. Eberhart, et al., Particle swarm optimization, in: Proceedings of IEEE international conference on neural networks, Vol. 4, Perth, Australia, 1995, pp. 1942–1948.
- [267] E. Atashpaz-Gargari, C. Lucas, Imperialist competitive algorithm: an algorithm for optimization inspired by imperialistic competition, in: 2007 IEEE congress on evolutionary computation, IEEE, 2007, pp. 4661–4667.
- [268] D. Karaboga, B. Basturk, A powerful and efficient algorithm for numerical function optimization: artificial bee colony (abc) algorithm, Journal of global optimization 39 (3) (2007) 459–471.
- [269] X.-S. Yang, Z. Cui, R. Xiao, A. H. Gandomi, M. Karamanoglu, Swarm intelligence and bio-inspired computation: theory and applications, Newnes, 2013.
- [270] P. Moscato, et al., On evolution, search, optimization, genetic algorithms and martial arts: Towards memetic algorithms, Caltech concurrent computation program, C3P Report 826 (1989) 1989.
- [271] F. Neri, C. Cotta, Memetic algorithms and memetic computing optimization: A literature review, Swarm and Evolutionary Computation 2 (2012) 1–14.