

Alma Mater Studiorum Università di Bologna
Archivio istituzionale della ricerca

Lightweight Neural Architecture Search for Temporal Convolutional Networks at the Edge

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

Published Version:

Risso, M., Burrello, A., Conti, F., Lamberti, L., Chen, Y., Benini, L., et al. (2022). Lightweight Neural Architecture Search for Temporal Convolutional Networks at the Edge. IEEE TRANSACTIONS ON COMPUTERS, 0, 744-758 [10.1109/TC.2022.3177955].

Availability:

This version is available at: <https://hdl.handle.net/11585/904601> since: 2022-11-21

Published:

DOI: <http://doi.org/10.1109/TC.2022.3177955>

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>).
When citing, please refer to the published version.

(Article begins on next page)

This is the final peer-reviewed accepted manuscript of:

M. Risso *et al.*, "Lightweight Neural Architecture Search for Temporal Convolutional Networks at the Edge," in *IEEE Transactions on Computers*, vol. 72, no. 3, pp. 744-758, 1 March 2023.

The final published version is available online at:

<https://doi.org/10.1109/TC.2022.3177955>

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>)

When citing, please refer to the published version.

Lightweight Neural Architecture Search for Temporal Convolutional Networks at the Edge

Matteo Risso, Alessio Burrello, *Member, IEEE*, Francesco Conti, *Member, IEEE*,
Lorenzo Lamberti, *Member, IEEE*, Yukai Chen, *Member, IEEE*, Luca Benini, *Fellow, IEEE*,
Enrico Macii, *Fellow, IEEE*, Massimo Poncino, *Fellow, IEEE*, and Daniele Jahier Pagliari, *Member, IEEE*

Abstract—Neural Architecture Search (NAS) is quickly becoming the go-to approach to optimize the structure of Deep Learning (DL) models for complex tasks such as Image Classification or Object Detection. However, many other relevant applications of DL, especially at the edge, are based on time-series processing and require models with unique features, for which NAS is less explored. This work focuses in particular on Temporal Convolutional Networks (TCNs), a convolutional model for time-series processing that has recently emerged as a promising alternative to more complex recurrent architectures. We propose the first NAS tool that explicitly targets the optimization of the most peculiar architectural parameters of TCNs, namely dilation, receptive-field and number of features in each layer. The proposed approach searches for networks that offer good trade-offs between accuracy and number of parameters/operations, enabling an efficient deployment on embedded platforms. Moreover, its fundamental feature is that of being *lightweight* in terms of search complexity, making it usable even with limited hardware resources. We test the proposed NAS on four real-world, edge-relevant tasks, involving audio and bio-signals: (i) PPG-based Heart-Rate Monitoring, (ii) ECG-based Arrhythmia Detection, (iii) sEMG-based Hand-Gesture Recognition, and (iv) Keyword Spotting. Results show that, starting from a single seed network, our method is capable of obtaining a rich collection of Pareto optimal architectures, among which we obtain models with the same accuracy as the seed, and $15.9\text{-}152\times$ fewer parameters. Moreover, the NAS finds solutions that Pareto-dominate state-of-the-art hand-tuned models for 3 out of the 4 benchmarks, and are Pareto-optimal on the fourth (sEMG). Compared to three state-of-the-art NAS tools, ProxylessNAS, MorphNet and FBNetV2, our method explores a larger search space for TCNs (up to $10^{12}\times$) and obtains superior solutions, while requiring low GPU memory and search time. We deploy our NAS outputs on two distinct edge devices, the multicore GreenWaves Technology GAP8 IoT processor and the single-core STMicroelectronics STM32H7 microcontroller. With respect to the state-of-the-art hand-tuned models, we reduce latency and energy of up to $5.5\times$ and $3.8\times$ on the two targets respectively, without any accuracy loss.

Index Terms—Neural Architecture Search, Temporal Convolutional Networks, Deep Learning, Edge Computing, Energy Efficiency

1 INTRODUCTION

DEEP LEARNING (DL) models are at the core of many time-series processing applications. Notable examples are audio classification [1], bio-signals analysis [2], [3] and predictive maintenance [4], [5]. For many years, the state-of-the-art DL models for time-series-based tasks have been Recurrent Neural Networks (RNNs) [6]. Recently, however, new architectures have been proposed as viable alternatives to RNNs, such as Attention-based Transformers and Temporal Convolutional Networks (TCNs) [7]. The latter, in particular, are uni-dimensional Convolutional Neural Networks (CNNs) specialized for time series, which have been shown to provide an accuracy comparable to RNNs, while providing computational advantages, namely higher arithmetic intensity, smaller memory footprint and more

data reuse opportunities [7]. Thanks to these features, TCNs are particularly interesting for edge computing applications, where inference is directly executed on Internet of Things (IoT) edge devices, rather than on centralized servers in the cloud. On-device inference requires small and efficient DL models, compatible with the limited memory spaces and tight energy budgets of edge nodes, while avoiding the transmission of raw data to the cloud provides many advantages, such as better privacy, higher energy efficiency and more predictable response latency [8].

To meet such tight constraints, however, selecting an efficient model such as a TCN is just the first step. Next, it is paramount to optimize its architectural hyperparameters based on the task at hand, so that the resulting network occupies as low memory and performs as few operations as possible to reach the desired accuracy level. Nowadays, rather than manually, such architectural optimization is increasingly performed with automatic Neural Architecture Search (NAS) tools. A plethora of different NAS approaches have been proposed in the last few years, and several of these works have targeted edge devices [9], [10], [11], [12], [13]. However, to the best of our knowledge, none of them has focused specifically on models for time-series processing, nor specifically on TCNs, despite the unique features of these networks.

In fact, while TCNs share most of their key architectural

- M. Risso, Y. Chen, M. Poncino and D. Jahier Pagliari are with the Department of Control and Computer Engineering, Politecnico di Torino, 10129, Turin, Italy. E-mail: name.firstsurname@polito.it
- E. Macii is with the Interuniversity Department of Regional and Urban Studies and Planning, Politecnico di Torino, 10129, Turin, Italy. E-mail: enrico.macii@polito.it
- A. Burrello, F. Conti, L. Lamberti and L. Benini are with the Department of Electrical, Electronic, and Information Engineering, University of Bologna, 40136 Bologna, Italy. E-mail: name.surname@unibo.it
- L. Benini is also with the Department of Information Technology and Electrical Engineering at the ETH Zurich, 8092 Zurich, Switzerland. E-mail: lbenini@iis.ee.ethz.ch

features with standard CNNs, the peculiar 1D convolution operations at the heart of these networks increase the importance of some hyperparameters, such as the filters dilation and receptive field, resulting in a much larger variety in their values than what is common in 2D models for image-processing. Although there are NAS tools, originally designed for 2D CNNs, that can be easily extended to explore these parameters, they do so in a coarse-grain way, basically creating a different copy of all network layers for each architectural setting [13]. This approach results in highly memory- and time-consuming searches, requiring 100s of GPU hours even for relatively simple tasks, which in turn translate into large energy wastes and CO2 emissions. In contrast, lightweight NAS approaches explore a finer-grain space with lower complexity, but they achieve this result at the cost of focusing only on the key characteristics of a specific model type, e.g., the number of channels in each layer of 2D CNNs for computer vision [10], [11].

In Risso et al. [14], we proposed the first lightweight NAS explicitly designed for optimizing TCNs by tuning the dilation hyperparameter. In this work, we extend and complete [14] by including the optimization of the receptive field and of the number of channels of all convolutional layers in a TCN, as well as the number of neurons in Fully Connected layers, with a search time comparable to that of a single, standard training. Starting from a single seed model, our proposed tool, called *Pruning In Time (PIT)* can produce a rich set of Pareto optimal architectures in terms of number of operations/parameters versus accuracy. The following are the main contributions of our work:

- We frame the optimization of receptive field and dilation as a *structured weight pruning*, in which additional trainable masking parameters are added to different layer’s weights so that their binarized values encode valid settings of the architectural hyperparameters. These masks are then trained with a regularizer to reduce the model complexity as much as possible while preserving accuracy. While similar masking approaches already exist for optimizing the number of channels in a 2D convolutional layer [11], our work is the first to extend this approach to filter size and dilation.
- We consider two different regularizers, targeting respectively the reduction of the number of parameters and of the number of inference operations. This allows us to enlarge and enrich the collection of Pareto architectures found by our NAS.
- We test and validate PIT on four benchmarks relative to real-world time-series processing tasks where TCNs are commonly employed and for which a deployment on edge devices is relevant: (i) PPG-Based Heart-Rate Monitoring; (ii) ECG-based Arrhythmia Detection; (iii) sEMG-based Hand-Gesture Recognition; (iv) Keyword Spotting. Results show that PIT can find multiple Pareto-optimal architectures starting from a single seed network, achieving 15.9-152× parameter reduction while maintaining the same accuracy of the seed. PIT is also capable of either matching or surpassing the accuracy and computational cost of state-of-the-art hand-tuned networks.

Furthermore, our approach Pareto-dominates three popular NAS tools developed for computer vision, thanks to the exploration of a larger search space.

- We deploy some of the relevant Pareto-optimal solutions found for each task on two different edge devices, in order to measure their memory footprint, latency and energy consumption. The two considered platforms are the multicore GAP8 IoT processor [15] and the single-core STM32H7 MCU [16]. The deployment results show that, at iso-accuracy, solutions found by PIT reduce energy consumption and latency up to 5.45× on GAP8 and up to 3.83× on the STM32H7, compared to hand-tuned networks.

The code of PIT is released as open-source at: <https://github.com/EmbeddedML-EDAGroup/PIT>.

The rest of the paper is structured as follows. Section 2 provides the required background and surveys some of the most relevant NAS methods proposed in the literature. Section 3 presents the proposed methodology. Section 4 details the target benchmarks while Section 5 discussed the obtained results, and Section 6 concludes the paper.

2 BACKGROUND AND RELATED WORKS

2.1 Temporal Convolutional Networks

Temporal Convolutional Networks are 1-dimensional (1D) CNN variants that have recently gained significant traction for efficient time-series processing, obtaining state-of-the-art results in several tasks [17], [18], [19]. With respect to RNNs and their successive evolutions, such as the Long-Short Term Memory (LSTM) and Gated Recurrent Unit (GRU), TCNs are less affected by training-time issues, such as vanishing/exploding gradients and the large amount of training memory required by RNNs for long input sequences. Moreover, they also have computational advantages at inference time, since they share the better data locality and arithmetic intensity of standard CNNs, which makes them latency- and energy-efficient [7].

The main building blocks of TCNs are the same ones found in standard CNNs, i.e. Convolutional, Pooling and Fully Connected (FC) layers. However, the convolutional layers of a TCN are characterized by *causality* and *dilation*, two properties that make them suited for temporal inputs.

Causality enforces that the outputs of convolutions do not violate the natural cause-effect ordering of events. In practice, the outputs y_t of a TCN convolution only depend on a finite set of past inputs $x_{[t-F:t]}$, where t is a discrete index. *Dilation* is the mechanism used in TCNs for enlarging the receptive field of convolutions on the time axis, without requiring more trainable parameters and without increasing the number of operations required for inference. It is a fixed step d inserted between the input samples processed by each convolutional filter. Eq. 1 summarizes the 1D dilated convolution operation implemented by TCN layers:

$$y_t^m = \sum_{i=0}^{K-1} \sum_{l=0}^{C_{in}-1} x_{ts-d i}^l \cdot W_i^{l,m}, \forall m \in [0, C_{out}-1], \forall t \in [0, T-1] \quad (1)$$

where x and y are the input/output activations, T is the output sequence length, W the array of filter weights, C_{in}/C_{out}

the number of input/output channels, K the filter size and s the stride. We also define $F = d \cdot (K - 1) + 1$ the *receptive field* of the layer.

2.2 Neural Architecture Search

In recent years, several manually designed efficient and compact convolutional neural network architectures for edge devices have been proposed, including early MobileNets [20], ShuffleNets [21], EfficientNet [22], SqueezeNet [23], etc. While these models are very efficient, obtaining them required a long and time-consuming manual tuning of hyper-parameters, which has to be repeated from scratch when considering a different target task, or a different deployment target.

To solve this issue, many automated or semi-automated methods to optimize neural network architectures, easing the burden of designers, have been proposed. These approaches, generally denoted as Neural Architecture Search (NAS) algorithms, explore a large design space made of different combinations of layers and/or hyper-parameter values, selecting solutions that optimize a cost metric. The latter is often a function of both the accuracy of the network, and its computational cost (e.g., number of parameters or inference operations).

Table 1 qualitatively compares some of the most relevant works in this field, in terms of search time, memory requirements during training (Mem.), search space size, and possibility to vary the topology (number and type of layers) of the resulting NNs. For Time and Mem., smaller is better, whereas for Search Space, larger is better. Early NAS tools were based on Reinforcement Learning (RL) [9], [12], [24], [25] or Evolutionary Algorithms (EA) [26]. At each search iteration, these methods sample one or more architectures from the search space. Sampled networks are then trained to convergence to evaluate their accuracy (and possibly cost), which is then used to drive the next sampling. The repeated training in each iteration is the main drawback of these tools, for which a single search requires 1000s of GPU hours, even on relatively simple tasks. Accordingly, these methods are associated with large search time in Table 1. Memory occupation is low and comparable to a standard training, since each sampled architecture can be trained separately. The search space size is virtually unlimited, and these tools can easily support variable topologies. Notable exceptions are [9], which searches over a fixed convolutional topology of a variable number of layers without varying their type and the connections between them, and [24], that constrains its search space to a set of only 13 different layers per node.

To solve the search time issue of RL and EA methods, more recent *Differentiable* NAS (DNAS) approaches have proposed the so-called *supernets* [27]. Supernets are DNNs that include all possible alternative layers to be considered during the optimization. For instance, a single supernet layer might include multiple Convolutional layers with different kernel sizes, operating in parallel. The problem of choosing a specific architecture is then translated into the problem of choosing a *path* in the supernet [27]. The choice between the different paths is encoded with binary variables, jointly trained with the standard weights of the network using gradient-based learning. To search for accurate and efficient architectures, DNAS tools enhance the

TABLE 1
State-of-the-art NAS (Values: \uparrow = large, \nearrow = medium, \downarrow = small).

	Time	Mem.	Search Space	Topology
Reinforcement Learning				
Zoph et al. [9]	\uparrow	\downarrow	\nearrow	Variable*
MNASNET [12]	\uparrow	\downarrow	\uparrow	Variable
NASNET [24]	\uparrow	\downarrow	\nearrow	Variable
MetaQNN [25]	\uparrow	\downarrow	\uparrow	Variable
Evolutionary				
Real et al. [26]	\uparrow	\downarrow	\uparrow	Variable
DifferentiableNAS				
DARTS [27]	\nearrow	\uparrow	\downarrow	Variable
ProxylessNAS [13]	\nearrow	\nearrow	\nearrow	Variable
DmaskingNAS				
FBNetV2 [10]	\downarrow	\downarrow	\uparrow	Fixed
MorphNet [11]	\downarrow	\downarrow	\nearrow	Fixed
S-Path NAS [28]	\downarrow	\downarrow	\nearrow	Fixed
PIT (this work)	\downarrow	\downarrow	\uparrow	Fixed

* Depth only

normal training loss function with an additional differentiable regularization term that encodes the cost of the network. Typical cost metrics are the number of parameters and the number of Floating Point Operations (FLOPs) per inference [11]. Mathematically, DNAS tools search for:

$$\min_{W, \theta} \mathcal{L}(W; \theta) + \lambda \mathcal{R}(\theta) \quad (2)$$

where \mathcal{L} is the standard loss function, W is the set of standard trainable weights (e.g., convolutional filters), θ is the set of additional NAS-specific trainable parameters that encode the different paths in the supernet, \mathcal{R} is the regularization loss that measures the cost of the network and λ is a hand-tuned *regularization strength*, used to balance the two loss terms.

While DNAS algorithms are more efficient than early RL/EA-based solutions, training the entire supernet still requires huge computational resources both in terms of training time and memory occupation. This, in turn, translates in a reduction of the explored search space for practical DNASEs such as [27], which have to limit the search to few alternatives per layer, in order to keep the memory occupation under reasonable bounds. The authors of [13] have proposed ProxylessNAS, an advanced DNAS that reduces the memory requirements, keeping in memory at most two supernet paths for each batch of inputs. In ProxylessNAS, the normal weights and the additional parameters encoding supernet paths are trained and updated in an alternate manner. First, path parameters are frozen, and based on their current value, one sub-architecture of the supernet is stochastically sampled. Then, the weights of the sampled architecture are updated based on the training set. Second, the normal weights are frozen and the architectural parameters are trained on the validation set. This second phase updates two different paths at a time, sampling them from a multinomial distribution. In turn, this clever strategy allows ProxylessNAS to explore a significantly larger search space compared to other DNAS tools.

A further evolution in the direction of lightweight NAS is constituted by DMaskingNAS [10], fine-grain NAS [11]

and Single-Path NAS [28] approaches. In these solutions, the supernet is replaced by a single, usually large, architecture with a unique path. Optimized architectures are found as modifications of this initial *seed model*, obtained tuning hyper-parameters, such as the number of channels in each layer [11]. The key mechanism that enables this tuning within a normal training loop is the use of *trainable masks*, used to prune parts of the network. DMaskingNAS tools pursue the same DNAS objective of (2), where θ now represents the set of trainable masks. FBNet-V2 [10], for instance, uses a set of dedicated masks, each of which encodes a different number of output channels or a different spatial resolution, and is weighted with a trainable parameter. At the end of the search, the mask coupled with the largest parameter is used to determine the final architectural setting. Similarly, MorphNet [11] exploits as masking parameters the pre-existing multiplicative terms of batch normalization layers [29]. When these parameters assume a value lower than a threshold, the corresponding channels/feature maps from the preceding Convolutional layer are eliminated.

These approaches are more constrained than supernet-based ones in terms of NN topology. In fact, they do not allow to select between alternative layers (e.g., standard convolution versus depth-wise + point-wise convolution). On the other hand, they have two key advantages. First, they have much lower memory cost and search time, while still being able to find high-quality architectures. Crucially, the search time of a DMaskingNAS is comparable to a standard network training. Second, some DMaskingNASes (including our work) can explore the search space at a much finer grain. For example, MorphNet [11] can easily select between 1 and 32 output channels in a Convolutional layer with a granularity of 1, by starting from a 32 channels seed layer, and eliminating those corresponding to the smallest batch normalization multiplicative parameters. Obtaining the same result with a standard DNAS would require a very large supernet, with 32 parallel convolutional layers. The masking and supernet approaches can also be combined, to bypass the limitations of DMaskingNAS [30].

The NAS literature referenced above focuses almost exclusively on 2D-CNNs for computer vision. None of the existing approaches has been applied to time-series processing tasks, despite the fact that a large amount of edge-relevant real-world tasks deal with uni-dimensional time-dependent signals (e.g., bio-signals, audio, energy-traces, sensor readings from industrial machines, etc). Our work tries to fill this gap, by proposing a novel DMaskingNAS that targets the optimization of 1D networks. Moreover, the working principles of our tool are general, and could form the basis for a more general NAS, able to explore temporal hyper-parameters of arbitrary N-dimensional Convolutional layers (e.g., including also 3D-CNNs for spatio-temporal data processing), although this paper focuses exclusively on TCNs.

3 PROPOSED METHOD

We name our proposed tool *Pruning in Time (PIT)*, since it targets networks that process time-series, and the core mechanism of a DMaskingNAS is very similar to structured pruning [31]. PIT explores the architectures of convolutional and fully-connected (FC) layers, the two most compute-

TABLE 2
List of symbols used in the paper.

Symbol	Description
x	Input activations of a convolutional layer
y	Output activations of a convolutional layer
T	Output sequence length of a convolutional layer
C_{in}, C_{out}	Number of input/output channels of a conv. layer
W	Convolutional filter weights
K	Convolution filter size
s	Convolution stride
d	Convolution dilation
F	Convolution receptive field
\mathcal{L}	Task-specific loss function
\mathcal{R}	Regularization loss function
λ	Regularization strength
$\mathcal{S}, \hat{\mathcal{S}}$	Search space and sampled architecture
L_n	Generic convolutional/FC layer
N	Number of convolutional/FC layers
θ, Θ	Generic NAS architectural parameters and corresponding binary mask
α, Θ_A	NAS architectural parameters to optimize C_{out} and corresponding binary mask
β, Θ_B	NAS architectural parameters for F , and corresponding binary mask
$\gamma, \Gamma, \Theta_\Gamma$	NAS architectural parameters for d , intermediate binary mask elements and final binary mask
C_β, C_γ	Transformation matrices to generate Θ_B and Θ_Γ from β and γ .
$k(i)$	Index mapping function used to generate Θ_Γ from Γ

and memory-expensive operations present in TCNs. For each convolutional layer, PIT jointly explores the *number of channels* (C_{out}), the *receptive field* (F), and the *dilation* (d). Moreover, by tuning both F and d , it also indirectly affects the *filter size* K . To the best of our knowledge, no DMaskingNAS from literature has optimized the receptive field or the dilation, even for 2D-CNNs. Similarly, PIT can also optimize the number of output neurons of FC layers¹.

We provide an overview of the search space explored by our tool and of its general working principle in Section 3.1. Then, we detail the mechanisms used to generate differentiable masks for each considered hyper-parameter in Sections 3.1.1-3.1.4. Finally, the two cost regularizers used to drive the search and the overall training procedure are described in Section 3.2 and 3.3 respectively. Table 2 summarizes the main mathematical symbols used throughout the paper.

3.1 Search Space

As shown in Figure 1, PIT’s search space encompasses all sub-architectures derived from a seed TCN by tweaking the three aforementioned hyperparameters. In particular, PIT can decide to *reduce* C_{out} or F , and to *increase* d with respect to the seed, all of which have the effect of reducing the complexity and memory occupation of the layer.

To achieve this objective, each convolutional/FC layer of the seed is modified to become a function $L_n(W^{(n)}; \theta^{(n)})$ of its original weights tensor $W^{(n)}$ and of a new set of architectural parameters $\theta^{(n)}$. For a TCN with N layers, the search space of PIT is therefore defined by the set:

$$\mathcal{S} = \{L_n(W^{(n)}; \theta^{(n)})\}_{n=0}^{N-1} \quad (3)$$

1. This can be seen as a corner case of the C_{out} optimization, since FC layers are just a special case of 1D convolutions with $F = K = d = 1$ and C_{out} equal to the number of output neurons. Accordingly, the rest of this section describes PIT’s functionality for convolutions.

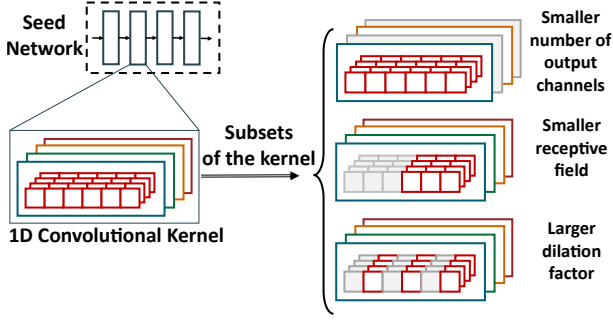


Fig. 1. Search space of PIT.

During the search, the elements of $\theta^{(n)}$ are properly combined to form a *binary mask* $\Theta^{(n)}$, which is used to *prune* a portion of the layer weights. In practice, an architecture \hat{S} is sampled from \mathcal{S} in each search iteration, by performing the Hadamard product between $W^{(n)}$ and $\Theta^{(n)}$, i.e., $\hat{S} = \{L_n(W^{(n)} \odot \Theta^{(n)})\}_{n=0}^{N-1}$. This eliminates the portions of $W^{(n)}$ that correspond to 0-valued mask elements, effectively letting the seed layer produce the same output that would be obtained with a smaller number of channels or receptive field, or with a larger dilation. The way in which $\Theta^{(n)}$ is generated from $\theta^{(n)}$ to produce this effect is the topic of Sections 3.1.1-3.1.3.

Having *binary masks* is required to either completely eliminate slices of $W^{(n)}$ (with value 0) or keep them untouched (with value 1) when sampling an architecture with the Hadamard product. In practice, this corresponds to sampling only *feasible* architectures (with integer C_{out} , F and d). To this end, we binarize $\Theta^{(n)}$ in the forward-pass of our search/training, applying an Heaviside step function with a fixed threshold $th = 0.5$.

At the same time, we also need to make the $\theta^{(n)} \rightarrow \Theta^{(n)}$ transformation differentiable, in order to embed the search into the standard gradient-based training of the network, learning contextually both the weights $W^{(n)}$ and the architectural parameters $\theta^{(n)}$. To cope with the Heaviside function derivation issues, i.e., derivative equal to 0 almost everywhere and not existent in δ , we follow the approach proposed in BinaryConnect [32], based on a Straight-Through Estimator (STE). Accordingly, during the backward-pass, the step function is simply replaced with an identity.

For notation simplicity, in the rest of the section we divide $\theta^{(n)}$ parameters in three groups: $\alpha^{(n)}$, used to tune the number of channels, $\beta^{(n)}$, which tune the receptive field, and $\gamma^{(n)}$, which affect the dilation factor. We also drop the superscript (n) when not needed. In PIT, each of these three groups of parameters is used to generate an *independent* binary mask, which can be then combined with the other two. Having independent masks for C_{out} , F and d , gives PIT the flexibility to optimize the three hyper-parameters either separately or jointly. At most, during a joint search, PIT explores:

$$|\mathcal{S}| \approx \prod_{n=0}^{N-1} (C_{out,seed}^{(n)} \cdot F_{seed}^{(n)} \cdot \lceil \log_2(F_{seed}^{(n)}) \rceil) \quad (4)$$

different solutions, where $C_{out,seed}$ and F_{seed} in (4) are those of the seed layers. The logarithmic term in (4) comes from

the fact that we only consider power-of-2 dilation factors, as detailed in Section 3.1.3. For a relatively small seed with $N = 8$, $F_{seed}^{(n)} = 17$, and $C_{out,seed}^{(n)} = 128 \forall n$, this corresponds to evaluating $\approx 10^{32}$ architectures in a single training.

3.1.1 Channels Search

To explore the number of channels in each convolutional layer, we take inspiration from [11]. In that work, the parameters of batch normalization (BN) layers [29] were transformed into binary masks to prune entire output channels and explore the space of all sub-layers with $C_{out} < C_{out,seed}$. However, requiring the presence of a BN layer after each convolution, although common in modern 2D-CNNs, still limits the applicability of the approach of [11]. Therefore, in PIT, we decouple the channel search from BN, and instead we exploit a dedicated trainable set of parameters α to zero-out entire filters from the W tensor of convolutional layers. PIT treats each output channel independently. So, it uses an α array of length $C_{out,seed}$, and it generates binary masks simply as:

$$\Theta_A = \mathcal{H}(|\alpha|) \quad (5)$$

where \mathcal{H} is the Heaviside binarization. Then, the layer function defined in (1) is modified to:

$$\tilde{y}_t^m = \sum_{i=0}^{K-1} \sum_{l=0}^{C_{in}-1} x_{ts-di}^l \cdot (\Theta_{A,m} \cdot W_i^{l,m}) \quad (6)$$

In practice, each binarized mask element is multiplied with all the weights of the *same convolutional filter*, i.e., with an entire slice of the weights tensor over the output channels axis. Each filter multiplied with a 0-mask effectively removes the corresponding output channel from the layer. Figure 2 depicts the application of Θ_A parameters to a simple layer with $C_{out,seed} = 4$.

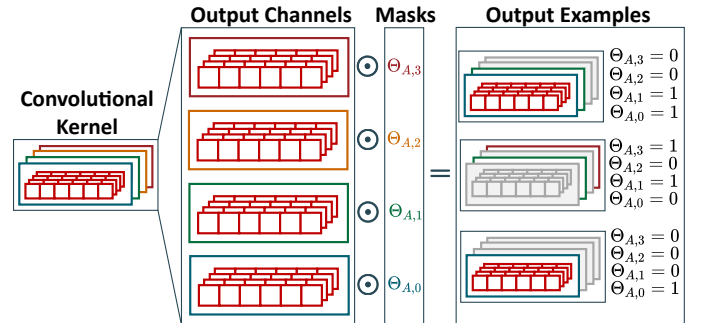


Fig. 2. Channels search example. Each $\Theta_{A,m} = 0$ zeroes-out the m -th convolutional filter, i.e., a slice of size $K \times C_{in}$ of the weights tensor W .

Noteworthy, besides reducing the number of channels, PIT can also *eliminate* entire layers from the network, if the latter includes skip-connections. In particular, if all the $\Theta_{A,m}$ of a convolutional layer are zeroed-out, then the inputs only flow through the skip connection, effectively reducing the number of the layers in the network by one. If skip connections are not present, instead, at least one output channel is always kept active to avoid breaking the network connectivity.

Our channels search scheme differs significantly from existing DMaskingNAS such as FBNetV2 [10], since we mask *weights tensors* rather than output activations. Fundamentally, as explained below, this makes our method easily extensible to the exploration of other hyper-parameters such as F and d , which would be much more difficult to optimize with an activations mask. Moreover, we use *independent* binarized parameters to mask each channel, rather than a set of predefined masks with an increasing number of trailing 0s, combined via Gumbel Softmax, as done in [10]. This, in turn, means that we can eliminate *any* combination of channels, not just the trailing ones.

3.1.2 Receptive Field Search

The second critical hyperparameter that we explore is the receptive field F , i.e., the range of input time-steps involved in a convolution. In standard convolutions, the receptive field is equal to the filter size ($F = K$). However, as detailed in Section 2.1, this no longer holds for TCNs when the dilation factor d is > 1 , and the general relation becomes: $F = (K - 1) \cdot d + 1$. As noted at the beginning of this section, by exploring both F and d , PIT also indirectly optimizes the filter dimension K .

The receptive field is explored using an array of additional trainable parameters β of length F_{seed} . Differently from the output channels, however, the β need to be further combined to define the corresponding binary differentiable masks. The reason is that, to “simulate” the effect of a smaller receptive field through masking, it is not sufficient to mask *any* set of time-slices in the weights tensor: in a causal TCN convolution, the receptive field extends exclusively in the past. Thus, the slices that should be eliminated are always the “oldest” ones, i.e., those that are multiplied with input time-steps that are farthest in the past. To do so, we derive elements of the binary masks Θ_B from β as:

$$\Theta_{B,i} = \mathcal{H} \left(\sum_{j=1}^{F_{seed}-i} |\beta_{F_{seed}-j}| \right) \quad (7)$$

Each $\Theta_{B,i}$ is then multiplied with a time-slice of the W tensor during the forward-pass, as shown in Figure 3. Therefore, when searching for the receptive field, (1) becomes:

$$y_t^m = \sum_{i=0}^{K-1} \sum_{l=0}^{C_{in}-1} x_{ts-di}^l \cdot (\Theta_{B,di} \odot W_i^{l,m}) \quad (8)$$

Thanks to the construction of (7), we have that if $i > j$, then $\Theta_{B,i} \leq \Theta_{B,j}$. In turn, this ensures that the first weight slices to be pruned are always the leftmost ones, as shown in the example on the right of Figure 3. Importantly, β_0 is always kept constant and equal to 1. This ensures that, once binarized, $\Theta_{B,0}$ is also always $= 1$, and consequently, that all convolutions take *at least one time-step* as input.

In practice, for efficiency reasons, we generate binary masks using the matrix transformation:

$$\Theta_B = \mathcal{H}(C_\beta \cdot |\beta|) \quad (9)$$

where C_β is a constant upper triangular matrix of 1s generated once at the beginning of a search, as shown on the left of Figure 4.

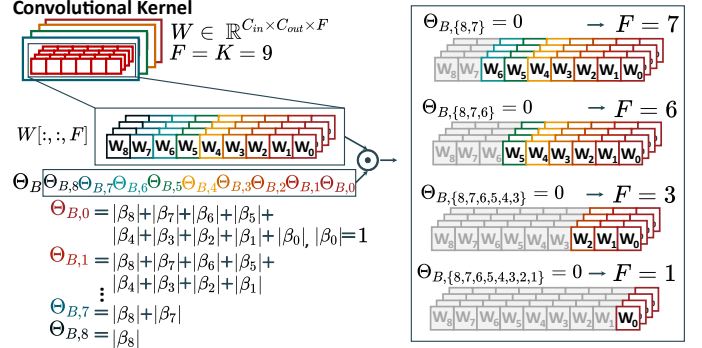


Fig. 3. Receptive field search example. Each $\Theta_{B,i} = 0$ eliminates the contribution of 1 input time-step from the convolution output, by zeroing out a time-slice of size $C_{out} \times C_{in}$ of the weights tensor W .

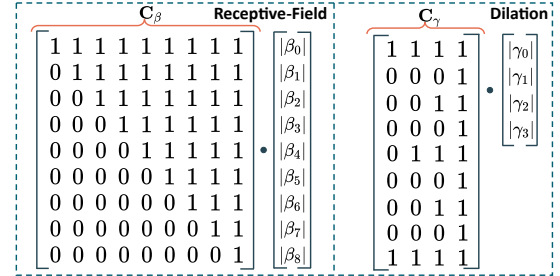


Fig. 4. Example of conversion between trainable architectural parameters β and γ and corresponding binary masks Θ_B and Θ_Γ , for a layer with $F_{seed} = 9$.

3.1.3 Dilation Search

Lastly, PIT also explores the dilation factor d . Similarly to the receptive field, also searching for dilation imposes some constraints on the portions of the weights tensor that should be pruned by our NAS. In particular, we need to ensure that only *regular* dilation factors are generated, i.e., that the time-steps gaps between consecutive convolution inputs are all equal for a given layer. For example, we do not want to obtain a layer that takes as input time-steps t , $t - 1$, $t - 3$, and $t - 10$, corresponding to gaps of 0, 1, and 6 time-steps respectively. In fact, such a layer would not be supported by most inference libraries, in particular those for edge devices [33], [34], which only implement regular dilation, as the latter enables more regular memory access patterns and better low-level optimizations.

Based on these observations, we follow an approach similar to the one described in Section 3.1.2. We start from an array of trainable parameters γ , which are then combined to compose differentiable binary masks². Our method only supports power-of-2 dilation factors which, besides being the most commonly used values, also simplify the generation of the masks. Thus, we have: $\text{len}(\gamma) = \lceil \log_2(F_{seed}) \rceil$.

In order to obtain the elements of Θ_Γ , we pass through

2. In our preliminary work of [14] we used a different mechanism to generate dilation masks as a *product* of γ elements instead of a sum. However, we found that this new approach is superior as it does not introduce nonlinear terms in the γ gradients.

Algorithm 1

```
1: for  $i \leftarrow 1, \dots, \text{Steps}_{\text{wu}}$  do #warmup loop
2:   Update  $W$  based on  $\nabla_W \mathcal{L}(W)$ 
3: end for
4: while not converged do #search loop
5:   Update  $W$  and  $\theta$  based on  $\nabla_{W,\theta}(\mathcal{L}(W; \theta) + \lambda \mathcal{R}(\theta))$ 
6: end while
7: for  $i \leftarrow 1, \dots, \text{Steps}_{\text{ft}}$  do #fine-tuning loop
8:   Update  $W$  based on  $\nabla_W \mathcal{L}(W)$ 
9: end for
```

see Section 3.3), $K_{\text{eff}}^{(n)}$ corresponds to the real filter size of the seed. In fact, each $\Theta_{B/T}$ is obtained as sum of a different number of γ (or β) elements. As a result, without normalization, the estimated cost would be higher than the real filter size. For instance, in a layer with $F_{\text{seed}} = 5$ and with all β/γ initialized at 1, without the denominators, we would have $K_{\text{eff}} = 33$, which is clearly incorrect. Conversely, with the denominators, we have $K_{\text{eff}} = 5 = F_{\text{seed}}$, which is correct, since the initialization of $\gamma = 1$ implicitly imposes $d = 1$.

3.2.2 OPs Regularizer

The second proposed regularizer R_{ops} estimates the number of operations required to perform an inference. Since the number of OPs of a 1D convolutional layer is $T \times C_{\text{in}} \times C_{\text{out}} \times K$, where T is the output sequence length defined in (1), the regularizer expression is simply:

$$\mathcal{R}_{\text{ops}} = \sum_{n=1}^N (\mathcal{R}_{\text{size}}^{(n)} \cdot T^{(n)}) \quad (17)$$

In practice, when targeting the reduction of the total OPs for inference, the only difference in the regularizer is that the cost of each layer is weighted by the output sequence length. This is particularly important in presence of layers such as pooling, strided convolution, etc., which significantly reduce T , and consequently the number of OPs for the downstream part of the network.

3.3 Training Procedure

Algorithm 1 summarizes the three main phases of a PIT architecture search. The first phase consists of Steps_{wu} iterations of warmup. At this stage of the algorithm, all θ parameters (i.e., α , β and γ) are initialized to 1 and frozen. Accordingly, all elements of the binary masks Θ are also binarized to 1. Therefore, warmup coincides with a normal training of the seed network, where the only objective is minimizing the task loss function \mathcal{L} . The number of warmup iterations is a user-defined parameter. In practice, in all our experiments, we warm up to convergence.

The second phase is where the actual NAS takes place. In the search loop, the model weights W and the architectural parameters θ are optimized simultaneously. Accordingly, the goal of this phase is to minimize the sum of the task-specific loss \mathcal{L} and of one of the two the regularization losses \mathcal{R} discussed in Section 3.2, weighted by the regularization strength λ . The duration of the search phase is controlled by an early-stop mechanism which monitors the value of \mathcal{L} on an unseen validation split of the target dataset, and stops the search when the latter does not improve for 20 epochs.

Finally, in the third and last phase the θ parameters, and corresponding Θ binary masks, are frozen to their latest values. This corresponds to sampling from the search space the architecture that PIT determined as optimal during the previous phase. Then, the weights W of the selected network are fine-tuned or re-trained from scratch, considering only the \mathcal{L} loss.

In order to obtain different Pareto points in the accuracy versus cost (size or OPs) space with PIT, it suffices to repeat Algorithm 1 changing the regularization strength λ . More precisely, the warmup phase can be performed just once, saving the final weights of the seed network. Overall, Algorithm 1 has a complexity that is comparable to a single TCN training. Moreover, the requirements in terms of GPU time and memory are greatly reduced with respect to a supernet-based DNAS. Therefore, obtaining 10s of Pareto points by changing λ still has a manageable cost, as shown for example by the results of Figure 10.

4 BENCHMARKS

We test PIT on four edge-relevant real-world benchmarks. We select diverse benchmarks to comprehensively evaluate the effectiveness of the proposed NAS. Specifically, we consider regression as well as classification tasks; the inputs analyzed are both raw or extracted features, and the TCNs employed as seed are based on different architectural styles.

4.1 PPG-based Heart-Rate Monitoring

The first benchmark deals with Heart-Rate (HR) monitoring on wrist-worn devices, using Photoplethysmography (PPG) sensors coupled with tri-axial accelerometers to mitigate the effect of motion artifacts [17], [35]. We target the PPG-Dalia [35] dataset, and the task is formulated as a *regression* of the HR value, whose ground truth is derived with ECG measurements. All results refer to the same input windowing and cross-validation scheme proposed in [35].

The seed network for this task is TEMPONet, a TCN originally proposed in [2] and later used for HR monitoring with state-of-the-art results in [17]. The network is composed of three *feature extraction* blocks and a final *regressor* module with three FC layers. Each feature extraction block is made of three convolutional layers with BatchNorm and ReLU activation, followed by an average pooling. The FC layers are also followed by BatchNorm and ReLU, and by a dropout layer with 50% rate. With respect to the original TEMPONet, our seed is obtained doubling the receptive field of all convolutions and setting the dilation to 1.

4.2 ECG-based Arrhythmia Detection

Our second benchmark deals with Electrocardiogram (ECG)-based arrhythmia detection, for wearable medical devices. We target the ECG5000 dataset [36], and the task consists in classifying the ECG signals in 5 classes: Normal, R-on-T Premature Ventricular Contraction, Premature Ventricular Contraction, Supraventricular Premature or Ectopic beat, and Unclassified Beat.

The reference TCN is ECGTCN, originally proposed in [37]. Differently from TEMPONet, ECGTCN is based on residual blocks. It has a first convolutional layer that

enlarges the number of input channels, followed by three modular blocks, each including two dilated convolutions with ReLU activation, BatchNorm and 50% dropout. The input and output feature maps of each block are then summed together. When the number of input and output channels differs, the residual path also includes a point-wise convolution (i.e., $K = 1$) in order to adapt the tensor sizes. The PIT seed is obtained from ECGTCN, setting the dilation of all layers to 1, while keeping the original receptive field.

4.3 sEMG-based Hand-Gesture Recognition

The third benchmark deals with hand-gesture recognition based on surface electromyography (sEMG) signals. For this task, we target the NinaPro DB1 dataset [38], which includes 52 heterogeneous gesture classes, using the same data pre-processing and augmentation described in [19].

The seed network is TCCNet, originally proposed in [19]. The architecture includes three feature extraction blocks, each composed of two dilated convolutions with ReLU and dropout (5% rate) and a residual branch with a point-wise convolution. The classifier includes an attention layer of the type described in [39] and a final FC layer with 53 output neurons (52 hand-gestures + 1 unknown class). The PIT seed is obtained simply setting the dilation to 1 in all layers.

4.4 Keyword Spotting

Our last benchmark is keyword spotting (KWS). We target the Speech Commands v2 dataset [40], following the pre-processing scheme proposed by the MLPerf Tiny benchmark suite [41] which produces 12 possible labels, including 10 words and two special classes for "unknown" and "silence".

As seed, we use the TCN presented in [18], called TC-ResNet14. The main difference with the other reference TCNs is that the original TC-ResNet14 did not use dilation, and the modular convolutional blocks alternate plain convolutions with strided convolution with $s = 2$. PIT's seed is obtained doubling the receptive field in each layer.

5 EXPERIMENTAL RESULTS

This section discusses the results obtained by PIT on the four aforementioned benchmarks. In particular, in Section 5.1, we present the global results of our NAS search in the accuracy versus number of parameters and accuracy versus number of OPs planes. In Section 5.2 we conduct ablation studies on one of the benchmarks, and in Section 5.3 we compare our approach with a state-of-the-art DNAS, ProxylessNAS [13], and with two state-of-the-art DMaskingNAS approaches, namely, MorphNet [11] and FBNetV2 [10]. Since the code for [10] is not publicly available, we re-implemented it based on the information provided in the paper. Finally, Section 5.4 presents the memory, latency and energy consumption results obtained deploying some of the networks found by PIT on two commercial edge devices.

PIT is written in Python (v3.6) and it is based on PyTorch (v1.7.1). All our training experiments and NAS searches are performed on a single NVIDIA TITAN Xp GPU with 12GB memory. The two deployment targets considered are: i) the multicore GAP-8 IoT processor by GreenWaves Technologies [15] and ii) the single-core STM32H7 MCU by

STMicroelectronics [16]. As inference software backend, we use the open-source layers library of [42] coupled with the tiling tool of [43] for GAP-8, and the CMSIS-NN library [44] for the STM32H7. All deployed networks are quantized to 8-bit, using PyTorch's built-in quantization algorithm.

5.1 Search Space Exploration

Figure 6 shows the results of applying PIT to the four benchmarks. The graphs report the TCNs accuracy (for classification tasks) or Mean Absolute Error (MAE, for regression tasks) on the x axis, and the number of parameters or OPs per inference on the y axis. The curves correspond to the outputs of PIT, where different points are obtained varying the regularization strength λ and considering both size and OPs regularizers. Moreover, each plot also reports the metrics of two additional TCNs. Black triangles correspond to the results obtained by the *hand-tuned* state-of-the-art TCNs directly taken from [17], [18], [19], [37], with the original number of channels, receptive fields, and dilation factors. Black squares, instead, indicate the metrics of the PIT *seeds*, i.e., the same networks modified as described in Section 4 (setting $d = 1$ everywhere, etc.) to enlarge the PIT search space.

The upper-left part of Figure 6 reports the results on the PPG-DaLia dataset for the PPG-based HR monitoring task. This is the only regression task considered, so the network performance is measured with the MAE, for which lower values are better. As shown by the graphs, starting from a single seed network, PIT is able to obtain a rich collection of Pareto-optimal architectures, spanning more than one order of magnitude both in terms of parameters (4.7k-78k) and OPs (0.27M-9.6M). Notably, PIT networks dominate in the Pareto sense both the seed architecture and the hand-tuned state-of-the-art TEMPONet. In particular, we obtain a similar MAE to the seed TCN (5.38 vs 5.40 BPM), with $120.0\times$ less parameters and $96.0\times$ less operations. Moreover, PIT also finds a new state-of-the-art deep learning model for this task, achieving a MAE of just 5.03 BPM while requiring only 53k parameters and 5.1M OPs, improving the best performing architecture proposed in [17]³ requiring $8.03\times$ and $5.42\times$ less parameters and OPs.

The Upper-right pair of charts shows the results obtained on the ECG5000 dataset for Arrhythmia Detection. PIT results span almost one order of magnitude in parameters (0.91k-5.36k) and OPs (50.3k-293.5k). Moreover, both the seed network and the hand-tuned one are Pareto-dominated. The best performing architecture found by our NAS improves the accuracy of the hand-tuned network (+1.03%) reducing both the number of parameters (-64.7%) and the FLOPs (-85.8%).

Lower-left part of Figure 6 shows the results obtained for the sEMG-based Hand-Gesture Recognition task on the NinaPro-DB1 dataset. The richness and diversity of the found architectures in terms of size and number of OPs are similar to the previous two benchmarks. However, while PIT results still dominate the seed, in this case the hand-tuned TCNNet sits on the Pareto front. Indeed, the PIT network that is nearest to the hand-tuned architecture on the

3. Note that this result is achieved without applying any additional post-processing as described in [17].

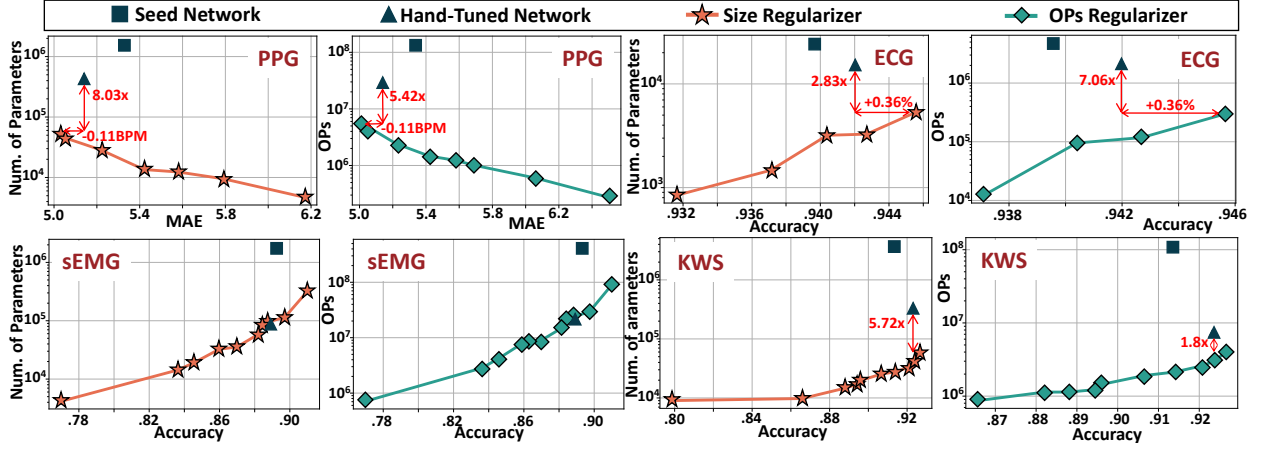


Fig. 6. Overall PIT Pareto fronts for the four target benchmarks, and comparison with seed and hand-tuned TCNs.

TABLE 3
Range of regularizer strength (λ) values for the four benchmarks.

Regularizer	PPG	ECG	sEMG	KWS
\mathcal{R}_{size}	1e-7 : 5e-4	5e-7 : 7.5e-3	1e-7 : 5e-6	5e-10 : 1e-5
\mathcal{R}_{ops}	1e-8 : 5e-5	5e-8 : 5e-4	5e-10 : 5e-8	1e-10 : 1e-6

curve, achieves a slightly lower accuracy (-0.47%) traded-off with a reduction of size (-3.33%). This result demonstrates the goodness of the original TCNNet proposed in [19] but, at the same time, it shows the good quality of the architectures found by our NAS, which despite starting from an oversized seed, is still able to produce optimized networks that closely resemble those tuned by experts.

Lastly, the lower-right part of Figure 6 shows the two Pareto fronts obtained on the Google Speech Commands dataset for Keyword Spotting. Once again, we largely outperform both the seed and the hand-tuned TCNs. Specifically, the most accurate PIT architecture slightly improves the accuracy of the hand-tuned network (+0.36%) while greatly reducing both the number of parameters (-82.53%) and FLOPs (-44.53%). Moreover, we obtain Pareto points that span 10k-98k parameters and 0.87M-3.98M OPs. It is important to note that the bad performance obtained by the seeds (black squares) for all four benchmarks is due to over-fitting, which in turn is caused by the large number of channels and receptive fields, and the absence of dilation.

Table 3 reports the range of regularizer strengths λ used on the four benchmarks to obtain these results. In general, λ should be set so that the two additive terms in the loss (\mathcal{L} and $\lambda\mathcal{R}$) assume comparable values at the beginning of a training. This ensures that PIT takes into account both accuracy and inference cost in its search, without degenerating to one of the two corner cases, i.e., accuracy-driven-only and cost-driven-only optimization. The corresponding values of λ vary for different tasks, as shown in the table. However, we found that a good rule of thumb, which works for all benchmarks, to identify the order of magnitude of the regularization strength is to start from $\lambda = 1/(\text{Seed Model Size})$. Then, based on the results of a PIT search with this initial value, one can decide to increase/decrease

λ to obtain smaller/more accurate TCNs respectively. By monitoring the loss in the initial epochs, it is also very easy to detect when the NAS is falling in one of the corner cases (one term much larger than the other) and stop the search immediately, without wasting training time.

5.2 Ablation Studies

This section analyzes the impact of some of the most important PIT parameters. Due to space limitations, we report the results of this study only for the PPG-based HR monitoring benchmark.

5.2.1 Hyper-parameters

Figure 7 analyzes the contribution of different hyper-parameters to the quality of results found by PIT. For this experiment, we use the \mathcal{R}_{size} regularizer and consider solutions in the MAE versus number of parameters space. We then repeat the NAS search 3 times. In each run, we freeze two of the three sets of architectural parameters (α , β and γ) to 1, letting PIT tune the third set. This gives us: i) the results of a search that only optimizes the number of channels in each layer (*Ch-Only*), performed on a TCN with maximal receptive field and $d = 1$, ii) the results of a receptive field-only search (*Rf-Only*), on a TCN with maximal C_{out} and $d = 1$, and iii) the results of a dilation-only search (*Dil-Only*) on a network with maximal F and C_{out} . The Pareto fronts obtained in each of these 3 conditions by varying the regularization strength λ are shown in the figure, together with the output of a complete search that optimizes all three hyper-parameters simultaneously (*All-in-One*).

The results clearly show that the main source of parameters reduction and performance improvement is the search along the channels dimension. This is probably due to the fact that the channels represent a large source of redundancy in hand-tuned TCNs, since their number is typically set using common heuristics, irrespective of the target task (e.g., C_{out} multiple of 32, progressively increasing along the depth of the network). However, Figure 7 also shows that optimizing *only* the number of channels is not sufficient, and that a combined optimization that also consider receptive field and dilation can yield Pareto-optimal networks across the entire MAE/parameters range.

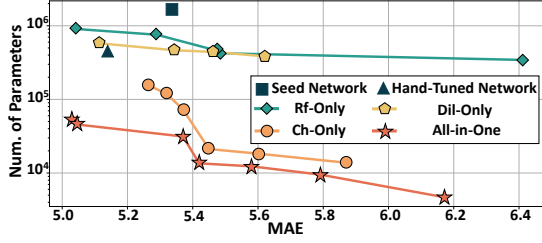


Fig. 7. Comparison between the results of PIT searches with different combinations of hyper-parameters for PPG-DaLia.

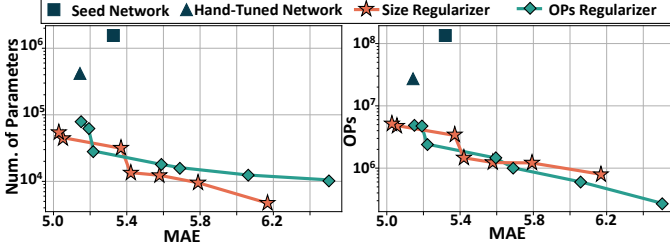


Fig. 8. Comparison of \mathcal{R}_{size} and \mathcal{R}_{ops} regularizers for PPG-DaLia.

5.2.2 Regularizers

Figure 8 compares the Pareto fronts obtained using the \mathcal{R}_{size} regularizer (with orange stars) and \mathcal{R}_{ops} regularizer (with green diamonds). Note that the PPG-based HR monitoring benchmark is the one for which the distinction between model size and number of OPs is most relevant, due to the presence of several layers (average pooling and strided convolution) that modify the activation array length T .

The Figure shows that, as expected, the majority of the Pareto points in the MAE versus number of parameters plane are produced when using the \mathcal{R}_{size} regularizer, with the few exceptions being due to local minima. Vice versa, the \mathcal{R}_{ops} regularizer tends to generate superior solutions in terms of MAE versus number of OPs.

5.3 Comparison with state-of-the-art NAS tools

Figure 9 compares the Pareto fronts obtained with PIT and three state-of-the-art NAS tools, namely ProxylessNAS [13] MorphNet [11] and FBNetV2 [10], on the HR monitoring benchmark. Results show that PIT outperforms all three across the entire design space, except for one MorphNet and one FBNetV2 point, that achieve a very low number of operations, although at the cost of a quite large MAE. The main reason for the superior results of PIT is the fact that our NAS explores a larger and finer-grain search space with respect to the baselines. For what concerns MorphNet and FBNetV2, this is partly due to the intrinsic nature of those tools, which cannot explore receptive field, nor dilation [10], [11]. Accordingly, F and d in their respective seeds have been set to the hand-tuned values of the state-of-the-art network. This different search starting point compared to PIT is the reason why, in the low-size/high-MAE regime, these tools find a single Pareto-optimal point.

For FBNetV2, we considered a *Coarse* search space, including 4 C_{out} alternatives per layer, uniformly spaced, i.e., $1/4 C_{out,seed}$, $1/2 C_{out,seed}$, $3/4 C_{out,seed}$ and $C_{out,seed}$ and a

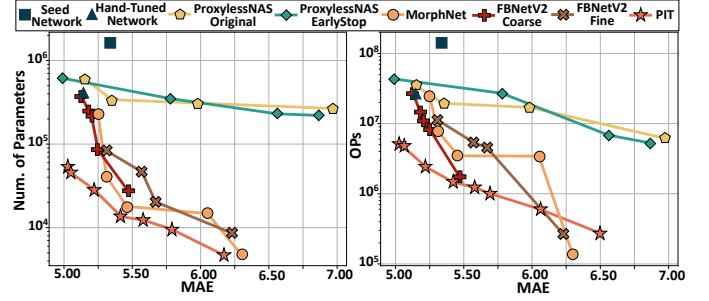


Fig. 9. Quality of results comparison between PIT and state-of-the-art NAS tools on the PPG-DaLia dataset.

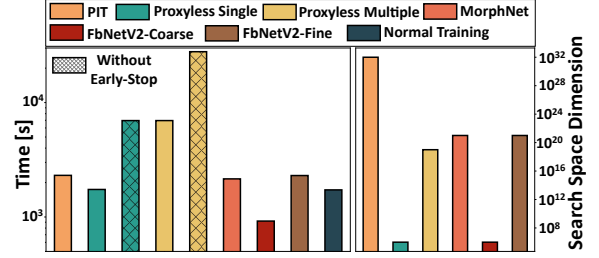


Fig. 10. Search space and time comparison between PIT and state-of-the-art NAS tools on the PPG-DaLia dataset.

Fine search space, which instead evaluates all C_{out} values with a granularity of 1. The latter is more similar to PIT, but the former achieves superior results in most cases. This is because FBNetV2 uses a pre-defined binary mask for each layer variant, combining them through a Gumbel softmax, as explained in Sec. 3.1.1. Experimentally, we found that with a too large number of masks, the search becomes unstable and yields sub-optimal results. In contrast, PIT does not have this limitation since it uses *independent* trainable masks that keep or eliminate an individual channel.

ProxylessNAS, being a super-net-based DNAS, would be virtually able to explore the entire PIT search space, as long as all the versions of layers to be explored are included in the super-net [13]. However, doing so would result in a too huge network, impossible to train due to memory and time requirements. In fact, as detailed in Section 3, PIT explores C_{out} and F with a granularity of 1, and for d , it considers all possible power-of-2 values. Therefore, each super-net node should include $C_{out,seed} \cdot F_{seed} \cdot \lceil \log_2(F_{seed}) \rceil$ different layers, connected in parallel. With the same parameters used for the example at the end of Section 3.1, this would correspond to ≈ 10000 different versions of *each* layer.

Therefore, we select a coarser-grain search space for ProxylessNAS, trying to make the comparison with PIT as fair as possible, while keeping the search space size similar to the one of the original paper [13]. To do so, we use the following procedure. First, we perform multiple ProxylessNAS searches on C_{out} , F and d separately, keeping the two not-optimized hyper-parameters at the seed values. In each of these searches, we consider 4 layers variants in each super-net node, uniformly sampling the PIT search space (in the same way described above for FBNetV2-Coarse). We then run ProxylessNAS multiple times with different regularization strengths. We identify, for every layer, the

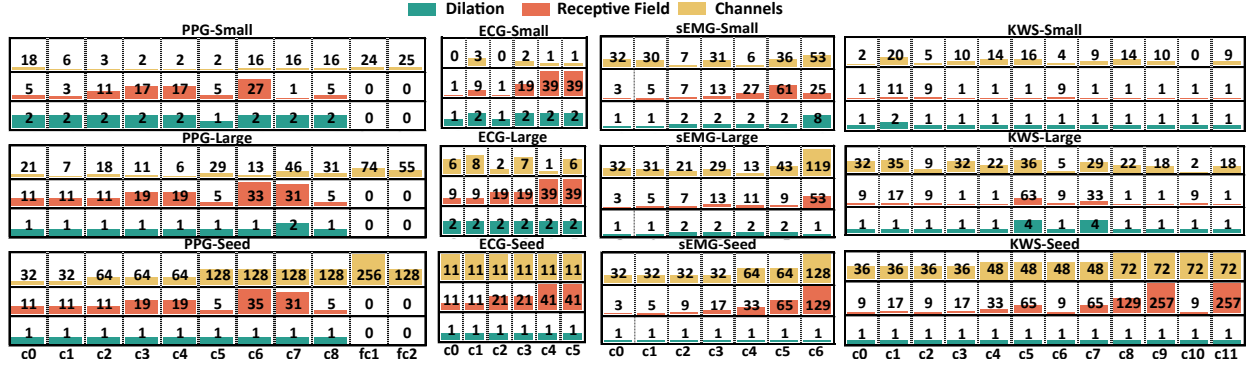


Fig. 11. Hyperparameters of the deployed PIT architectures and corresponding seed network for the four benchmarks.

two values of each hyper-parameter that have been chosen more frequently. The 2^3 possible combinations of the latter are used to generate the *combined* search space for ProxylessNAS, which, accordingly, includes 8 layer variants in each super-net node. The Pareto fronts of Figure 9 are obtained running ProxylessNAS multiple times on this combined search space, with different regularization strengths. We report both the results obtained with the training scheme proposed in the original paper (*Original* curve), which runs for a fixed number of epochs, and with the same early-stop mechanism employed for PIT (*EarlyStop* curve). As shown, the quality of results is similar in both cases.

Figure 10 compares the search space dimension and average execution time of the different tools on the HR monitoring benchmark. For reference, the execution time of a standard training of the seed network is also reported. All time results refer to the search phase only (without warm up), and are obtained on a single NVIDIA Titan XP GPU with a batch-size of 128. For ProxylessNAS, we report the results of the the initial single-hyper-parameter searches (*Proxyless-Single*) and of the final combined search (*Proxyless-Multiple*), both with and without early-stopping.

Our algorithm explores a $10^{26} \times 10^{12} \times$ larger search space than Proxyless-Single/-Multiple. Further, it is only $1.13 \times$ slower than Proxyless-Single with early-stopping, and $3.55 \times$ faster than the variant without early-stopping, while it is $3.0 \times / 14.22 \times$ faster compared to Proxyless-Multiple with/without the early-stopping training. With respect to MorphNet, we explore a $10^{11} \times$ larger search space at the cost of a small $1.07 \times$ increase in runtime. FBNetV2-Coarse is the fastest tool, converging in few search epochs. Whereas offering a $2.5 \times$ speedup with respect to PIT, the explored search space is 10^{26} smaller. Instead, FBNetV2-Fine explores a 10^{11} smaller space while requiring the same search time of the proposed approach. Lastly, PIT’s time-overhead with respect to a normal training is only 34%.

5.4 Embedded Deployment

This section analyzes the results obtained deploying two TCNs for each target benchmark on the GAP8 IoT processor (running at 100 MHz) and on the STM32H7 MCU (at 480 MHz). For each task, we deploy the best performing network in terms of MAE or accuracy (L). Moreover, we also select a small network that achieves a MAE drop < 1 BPM,

TABLE 4
Detailed deployment results for the four benchmarks.

Task	TCN	Perf. int8 (float32)	Mem. [kB]	GAP8		STM32	
				Lat. [ms]	En. [mJ]	Lat. [ms]	En. [mJ]
PPG	HT	5.01 (5.14) BPM	423	23.2	1.2	58.3	13.6
	S	5.71 (6.17) BPM	4.7	1.18	0.06	3.2	0.75
	L	5.01 (5.03) BPM	53.2	4.25	0.22	15.2	3.56
ECG	HT	94.2 (94.2) %	15.2	2.69	0.14	6.66	1.56
	S	92.84 (93.16) %	0.9	0.78	0.04	1.8	0.42
	L	94.13 (94.13) %	5.4	1.26	0.06	2.84	0.66
sEMG	HT	88.89 (88.87) %	88.8	61.0	3.11	291	68.1
	S	86.97 (86.98) %	35.4	39.6	2.02	169	39.5
	L	91.2 (90.99) %	317.8	238	12.1	960	225
KWS	HT	92 (92.31) %	323.4	13.4	0.68	30.7	7.17
	S	87 (86.58) %	9.8	1.40	0.07	2.66	0.62
	L	92.16 (92.64) %	56.5	3.74	0.19	10.6	2.48

or an accuracy drop $< 5\%$ with respect to the best performing one (S). For comparison, we also deploy the baseline hand-tuned architectures (HT). Table 4 reports the results in terms of performance (MAE or accuracy, depending on the dataset), memory footprint, inference latency and energy consumption, while Figure 11 shows the hyper-parameters selected by our NAS.

PIT finds competitive solutions for both hardware targets and for all 4 tasks, despite the large difference in complexity among them, testified by the more than two orders of magnitude span in memory, latency and energy consumption in the results of Table 4. For PPG-based HR monitoring, the L/S models achieve a $0/0.70$ BPM MAE increase with respect to the hand-tuned TEMPONet respectively, while resulting in a $8.03/90.8 \times$ lower memory footprint, and a $5.45/19.6 \times$ lower latency and energy consumption on GAP8. On the STM32 MCU, the latency and energy reduction of the two PIT outputs is $3.83/18.2 \times$. PIT’s L/S models for ECG processing, instead, achieve $+0.07\%/-1.36\%$ accuracy with respect to the hand-tuned ECGNET, with a $2.83/16.8 \times$ lower memory footprint, $2.13/3.44 \times$ latency and energy reduction on GAP8, and $2.34/3.7 \times$ on the STM32. For the sEMG gesture recognition task, the L/S models found by PIT obtain $+2.31\%/-1.92\%$ accuracy compared to TCCNet. In this case, the higher accuracy of the large model is paid with a $3.57 \times$ larger memory footprint, and a $3.85 \times$ latency and energy increase on GAP8 ($3.33 \times$ on the STM32H7), proving once again the goodness of the hand-tuned model for this task. The small TCN, instead,

results in a $2.51\times$ memory reduction, and $1.54\times$ and $1.72\times$ lower latency and energy on the two targets. Lastly, the L/S PIT outputs for KWS obtain $+0.16\%/-5\%$ accuracy with respect to TCResNet-14, with a $5.72/33.1\times$ lower memory footprint, $3.58/9.54\times$ lower energy and latency on GAP8, and $2.9/11.54\times$ lower energy and latency on STM32H7.

Figure 11 shows the high variability of hyper-parameters settings found by PIT, and the different optimization behaviours for different benchmarks. In general, comparing PIT outputs with the respective seeds, we can observe that the optimized networks found by our tool contradict several “rules of thumb” of manual DNN design, such as progressively increasing the number of channels and dilation for deeper layers. Accordingly, our NAS could also provide interesting insights for better TCN design. For instance, for the PPG benchmark, PIT finds solutions that are characterized by a large number of channels in the first and last layers, while keeping an overall high receptive field in the core of the network.

6 CONCLUSION

We have proposed PIT, a lightweight NAS tool for TCNs, able to explore a large, fine-grained search space of architectures with low GPU memory requirements. PIT is, to the best of our knowledge, the first DNN NAS tool explicitly designed for 1D convolutional networks, and the first to target the optimization of the receptive field and dilation of convolutional layers. With experiments of four real-world benchmarks, we have shown that PIT is able to find improved versions of state-of-the-art TCNs, with a memory compression of up to $8.03\times$ ($90.8\times$) and a latency and energy reduction of up to $5.45\times$ ($19.6\times$) without (with a reasonable) accuracy drop, when deployed on commercial edge devices. Our future work will focus on extending PIT principles to generic N-dimensional CNNs.

REFERENCES

- [1] W. He *et al.*, “Deep neural networks for multiple speaker detection and localization,” in *IEEE ICRA*, 2018, pp. 74–79.
- [2] M. Zanghieri *et al.*, “Robust real-time embedded emg recognition framework using temporal convolutional networks on a multicore iot processor,” *IEEE Trans. Biomed. Circuits Syst.*, 2019.
- [3] N. D. Truong *et al.*, “Convolutional neural networks for seizure prediction using intracranial and scalp electroencephalogram,” *Neural Networks*, vol. 105, pp. 104–111, 2018.
- [4] M. Azimi *et al.*, “Data-driven structural health monitoring and damage detection through deep learning: State-of-the-art review,” *Sensors*, vol. 20, no. 10, p. 2778, 2020.
- [5] T. Cerquitelli *et al.*, “Manufacturing as a data-driven practice: Methodologies, technologies, and tools,” *Proc. IEEE*, vol. 109, no. 4, pp. 399–422, 2021.
- [6] Z. C. Lipton *et al.*, “A critical review of recurrent neural networks for sequence learning,” *arXiv:1506.00019*, 2015.
- [7] S. Bai *et al.*, “An empirical evaluation of generic convolutional and recurrent networks for sequence modeling,” *arXiv:1803.01271*, 2018.
- [8] W. Shi *et al.*, “Edge computing: Vision and challenges,” *IEEE Internet Things J.*, vol. 3, no. 5, pp. 637–646, Oct 2016.
- [9] B. Zoph *et al.*, “Neural architecture search with reinforcement learning,” *arXiv preprint arXiv:1611.01578*, 2016.
- [10] A. Wan *et al.*, “Fbnetv2: Differentiable neural architecture search for spatial and channel dimensions,” in *Proc. IEEE/CVF CVPR*, 2020, pp. 12 965–12 974.
- [11] A. Gordon *et al.*, “Morphnet: Fast & simple resource-constrained structure learning of deep networks,” in *Proc. of the IEEE CVPR*, 2018, pp. 1586–1595.
- [12] M. Tan *et al.*, “Mnasnet: Platform-aware neural architecture search for mobile,” in *Proc. IEEE CVPR*, 2019, pp. 2820–2828.
- [13] H. Cai *et al.*, “Proxylessnas: Direct neural architecture search on target task and hardware,” *arXiv preprint arXiv:1812.00332*, 2018.
- [14] M. Risso *et al.*, “Pruning in time (pit): A light-weight network architecture optimizer for temporal convolutional networks,” in *Proc. 58th DAC*, 2021, pp. 1–6.
- [15] E. Flamand *et al.*, “Gap-8: A risc-v soc for ai at the edge of the iot,” in *Proc. IEEE 29th ASAP*. IEEE, 2018, pp. 1–4.
- [16] ST Microelectronics. STM32H7. [Online]. Available: <https://www.st.com/en/microcontrollers-microprocessors/stm32h7-series.html>
- [17] M. Risso *et al.*, “Robust and energy-efficient ppg-based heart-rate monitoring,” in *Proc. IEEE ISCAS*. IEEE, 2021, pp. 1–5.
- [18] S. Choi *et al.*, “Temporal convolution for real-time keyword spotting on mobile devices,” *arXiv preprint arXiv:1904.03814*, 2019.
- [19] P. Tsinganos *et al.*, “Improved gesture recognition based on semg signals and tcn,” in *Proc. IEEE ICASSP*. IEEE, 2019, pp. 1169–1173.
- [20] M. Sandler *et al.*, “Mobilenetv2: Inverted residuals and linear bottlenecks,” in *Proc. IEEE CVPR*, 2018, pp. 4510–4520.
- [21] N. Ma *et al.*, “Shufflenet v2: Practical guidelines for efficient cnn architecture design,” in *Proc. ECCV*, 2018, pp. 116–131.
- [22] M. Tan *et al.*, “Efficientnet: Rethinking model scaling for convolutional neural networks,” *ArXiv*, vol. abs/1905.11946, 2019.
- [23] F. N. Iandola *et al.*, “Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 1mb model size,” *ArXiv*, vol. abs/1602.07360, 2016.
- [24] B. Zoph *et al.*, “Learning transferable architectures for scalable image recognition,” *Proc. IEEE/CVF CVPR*, pp. 8697–8710, 2018.
- [25] B. Baker *et al.*, “Designing neural network architectures using reinforcement learning,” *ArXiv*, vol. abs/1611.02167, 2017.
- [26] E. Real *et al.*, “Large-scale evolution of image classifiers,” in *Proc. ICML*. PMLR, 2017, pp. 2902–2911.
- [27] H. Liu *et al.*, “Darts: Differentiable architecture search,” *arXiv:1806.09055*, 2019.
- [28] D. Stamoulis *et al.*, “Single-path mobile automl: Efficient convnet design and nas hyperparameter optimization,” *IEEE J. Sel. Topics Signal Process.*, vol. 14, no. 4, pp. 609–622, 2020.
- [29] S. Ioffe *et al.*, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *Proc. ICML*. PMLR, 2015, pp. 448–456.
- [30] S. R. Chaudhuri *et al.*, “Fine-Grained Stochastic Architecture Search,” *arXiv:2006.09581*, 2020.
- [31] J. Yu *et al.*, “Scalpel: Customizing dnn pruning to the underlying hardware parallelism,” *ACM SIGARCH Computer Architecture News*, vol. 45, no. 2, pp. 548–560, 2017.
- [32] M. Courbariaux *et al.*, “Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1,” *arXiv preprint arXiv:1602.02830*, 2016.
- [33] S. Microelectronics. (2017) X-cube-ai. [Online]. Available: <https://www.st.com/en/embedded-software/x-cube-ai.html>
- [34] G. Technologies. (2019) Gap8 nntool. [Online]. Available: <https://greenwaves-technologies.com/manuals/>
- [35] A. Reiss *et al.*, “Deep ppg: large-scale heart rate estimation with convolutional neural networks,” *Sensors*, vol. 19, no. 14, p. 3079, 2019.
- [36] Y. Chen *et al.*, “A general framework for never-ending learning from time series streams,” *Data Min. Knowl. Discov.*, vol. 29, no. 6, pp. 1622–1664, 2015.
- [37] T. M. Ingolfsson *et al.*, “Ecg-tcn: Wearable cardiac arrhythmia detection with a temporal convolutional network,” in *Proc. 3rd IEEE AICAS*. IEEE, 2021, pp. 1–4.
- [38] M. Atzori *et al.*, “Building the ninapro database: A resource for the biorobotics community,” in *Proc. 4th IEEE RAS & EMBS BioRob*. IEEE, 2012, pp. 1258–1265.
- [39] Z. Yang *et al.*, “Hierarchical attention networks for document classification,” in *Proc. 2016 NAACL*, 2016, pp. 1480–1489.
- [40] P. Warden, “Speech commands: A dataset for limited-vocabulary speech recognition,” *arXiv preprint arXiv:1804.03209*, 2018.
- [41] C. R. Banbury *et al.*, “Benchmarking tinyml systems: Challenges and direction,” *arXiv preprint arXiv:2003.04821*, 2020.
- [42] A. Burrello *et al.*, “Tcn mapping optimization for ultra-low power time-series edge inference,” in *Proc. IEEE/ACM ISLPED*, 2021, pp. 1–6.
- [43] A. Burrello *et al.*, “Dory: Automatic end-to-end deployment of real-world dnns on low-cost iot mcus,” *IEEE Trans. Comput.*, 2021.

[44] ARM. CMSIS-NN. [Online]. Available: https://arm-software.github.io/CMSIS_5/NN/html/index.html



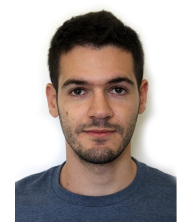
Matteo Rizzo received his B.Sc degree in Physical Engineering and M.Sc degree in Electronic Engineering at the Politecnico di Torino, Italy, in 2018 and 2020. He is currently working toward his Ph.D. degree at Politecnico di Torino, Italy. His research interests include Embedded Machine Learning and Energy-Efficient Embedded Systems.



Alessio Burrello received his B.Sc and M.Sc degree in Electronic Engineering at the Politecnico di Torino, Italy, in 2016 and 2018. He is currently working toward his Ph.D. degree at the Department of Electrical, Electronic and Information Technologies Engineering (DEI) of the University of Bologna, Italy. His research interests include parallel programming models for embedded systems, machine and deep learning, hardware oriented deep learning, and code optimization for multi-core systems.

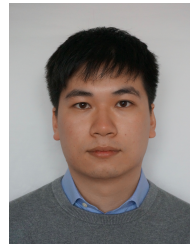


Francesco Conti received the Ph.D. degree in electronic engineering from the University of Bologna, Italy, in 2016. He is currently an Assistant Professor in the DEI Department of the University of Bologna. From 2016 to 2020, he held a research grant in the DEI department of University of Bologna and a position as post-doctoral researcher at the Integrated Systems Laboratory of ETH Zurich in the Digital Systems group. His research focuses on the development of deep learning based intelligence on top of ultra-low power, ultra-energy efficient programmable Systems-on-Chip. His research work has resulted in more than 60 publications in international conferences and journals and has been awarded several times, including the 2020 IEEE TCAS-I Darlington Best Paper Award.



Lorenzo Lamberti graduated with honors in both Bachelor's (2016) and Master's (2019) degree at the University of Bologna, Italy, where he is now pursuing a Ph.D. in Electronic Engineering in the Energy-Efficient Embedded Systems Laboratory. Previously, he has been an intern at the Fermi National Accelerator Laboratory of Chicago, US, and at the Datalogic Artificial Intelligence Laboratory in Pasadena, US. His research is currently targeted at autonomous navigation for nano-size unmanned aerial vehicles.

In particular, he focuses on neural architecture search, training and in-field deployment of Deep Neural Network-based navigation tasks on low-power MCUs.



Yukai Chen received the M.Sc. and Ph.D. degrees in computer engineering from Politecnico di Torino, Turin, Italy, in 2014 and 2018, respectively. He is currently a Postdoc researcher at Politecnico di Torino. His research interest includes computer-aided design for non-functional properties (power, thermal, reliability) modelling, simulation, and optimization of Cyber-Physical Systems, emphasizing low-power design, design automation and design space exploration.



He is a Fellow of the ACM and a member of the Academia Europaea.

Luca Benini is the Chair of Digital Circuits and Systems at ETH Zürich and a Full Professor at the University of Bologna. He has served as Chief Architect for the Platform2012 in STMicroelectronics, Grenoble. Dr. Benini's research interests are in energy-efficient systems and multi-core SoC design. He is also active in the area of energy-efficient smart sensors and sensor networks. He has published more than 1'000 papers in peer-reviewed international journals and conferences, five books and several book chapters.



mobility, clean and intelligent manufacturing. He is a Fellow of the IEEE.

Enrico Macii is a Full Professor of Computer Engineering with the Politecnico di Torino, Torino, Italy. He holds a Laurea degree in electrical engineering from the Politecnico di Torino, a Laurea degree in computer science from the Università di Torino, Turin, and a PhD degree in computer engineering from the Politecnico di Torino. His research interests are in the design of digital electronic circuits and systems, with a particular emphasis on low-power consumption aspects energy efficiency, sustainable urban mobility, clean and intelligent manufacturing. He is a Fellow of the IEEE.



Massimo Poncino is a Full Professor of Computer Engineering with the Politecnico di Torino, Torino, Italy. His current research interests include various aspects of design automation of digital systems, with emphasis on the modeling and optimization of energy-efficient systems. He received a PhD in computer engineering and a Dr.Eng. in electrical engineering from Politecnico di Torino. He is a Fellow of the IEEE.



Daniele Jahier Pagliari received the M.Sc. and Ph.D. degrees in computer engineering from the Politecnico di Torino, Turin, Italy, in 2014 and 2018, respectively. He is currently an Assistant Professor with the Politecnico di Torino. His research interests are in the computer-aided design and optimization of digital circuits and systems, with a particular focus on energy-efficiency aspects and on emerging applications, such as machine learning at the edge.