Made available by Hasselt University Library in https://documentserver.uhasselt.be

A systematic analysis of deep learning algorithms in high-dimensional data regimes of limited size Peer-reviewed author version

Jaxy, Simon; Nowe, Ann & LIBIN, Pieter (2024) A systematic analysis of deep learning algorithms in high-dimensional data regimes of limited size. In: 2024 IEEE 36TH International conference on tools with artificial intelligence, ICTAI, IEEE COMPUTER SOC, p. 515 -523.

DOI: 10.1109/ICTAI62512.2024.00079 Handle: http://hdl.handle.net/1942/46203

A systematic analysis of deep learning algorithms in high-dimensional data regimes of limited size

1st Simon Jaxy Artificial Intelligence Lab Vrije Universiteit Brussel (VUB) Brussels, Belgium simon.jaxy@vub.be 2nd Ann Nowé Artificial Intelligence Lab Vrije Universiteit Brussel (VUB) Brussels, Belgium ann.nowe@vub.be 3th Pieter Libin Artificial Intelligence Lab, Data Science Institute Vrije Universiteit Brussel (VUB), Universiteit Hasselt (UH) Brussels, Belgium, Hasselt, Belgium pieter.libin@vub.be

Abstract—There is a substantial demand for deep learning methods that can work with limited, high-dimensional, and noisy datasets. Nonetheless, current research mostly neglects this area, especially in the absence of prior expert knowledge or knowledge transfer. In this work, we bridge this gap by studying the performance of deep learning methods on the true data distribution in a limited, high-dimensional, and noisy data setting. To this end, we conduct a systematic evaluation that reduces the available training data while retaining the challenging properties mentioned above. Furthermore, we extensively search the space of hyperparameters and compare state-of-the-art architectures and models built and trained from scratch to advocate for the use of multi-objective tuning strategies. Our experiments highlight the lack of performative deep learning models in current literature and investigate the impact of training hyperparameters. We analyze the complexity of the models and demonstrate the advantage of choosing models tuned under multi-objective criteria in lower data regimes to reduce the likelihood to overfit. Lastly, we demonstrate the importance of selecting a proper inductive bias given a limited-sized dataset. Given our results, we conclude that tuning models using a multi-objective criterion results in simpler yet competitive models when reducing the number of data points.

Index Terms—Limited Data, Deep Learning, Multi Objective Optimization, Overfit

I. INTRODUCTION

In recent years, deep learning celebrated unprecedented success ranging from the rise of large language models [1]–[3] to the discovery of new materials [4], [5] and image generation [6], [7]. Pushing models further toward the boundary of computation enables a sheer unforeseeable range of possibilities [8], provided that the amount of data matches the demand of the steadily growing models and their hunger for data [9].

However, looking at the other end, where data is scarce but high-dimensional and noisy, we are still facing substantial obstacles. Research has long abstained from investigating deep learning models in these circumstances, although countless real-world applications require performative models to drive development and research. Rare diseases [10], [11] is a particular case of such an application where data can have very high-dimensional, noisy measurements and is limited by definition. Likewise, the discovery of new molecules to improve on existing drugs is a fundamentally low data problem as many newly proposed molecules are incompatible or even toxic [12]. Furthermore, archaeological discoveries rely on limited, highly complex data [13], [14]. Deep learning has made tremendous progress in the last few years, particularly regarding methodologies tailored for data in large quantities. In this work, we address settings where the data is limited, high-dimensional, and noisy by nature.

In recent work, Banerjee et al. [10] reflect on the state and applicability of deep learning methods for investigating rare diseases. To facilitate learning, they propose to increase the amount of data by combining data sets, injecting prior knowledge, or transferring the weights of a deep learning method trained on a related domain. However, in newly emerging fields such as diagnosing new rare diseases, simply merging datasets falls out of the question as we do not have relatable data. Therefore, using transfer learning to enable learning in high-dimensional yet scarce data environments is not an option since transferring weights from one domain to another might not always benefit performance [15].

We explore the limitations of deep learning when confronted with limited yet high-dimensional and noisy data without introducing prior knowledge or transferring weights from another domain. To do so, we conduct a systematic evaluation¹ that mirrors the challenging data contexts and enables us to analyze deep learning and a deep Gaussian Processbased model in a controlled framework, by systematically decreasing the available data, starting from the entire data set, until only one percent of the original data remains. We then task the algorithms with predicting on the complete test data to evaluate their abilities on the true data distribution. Furthermore, we extensively investigate the spaces of possible architectures and hyperparameters of several deep learning techniques to provide insights and guidelines to train models in this challenging set up. Specifically, we question the tradeoff between performance and complexity in a multi-objective context. With our work, we demonstrate the capability of deep learning models to learn in this challenging framework.

The rest of the paper is structured as follows. We provide

¹https://github.com/simomoxy/limited_data.git

an overview of related work in section II. In sections III and IV, we present our evaluation and deep learning models in detail before explaining our experimental procedures. Finally, we display and discuss our results in sections V and VI respectively.

II. RELATED WORK

Training small and efficient models that function correctly in adverse circumstances remains an important challenge. Recently, the research community has recognized the lack of resources invested in this branch of deep learning, and part of it is shifting its focus away to reduce model complexity [16]–[18]. However, popular techniques to enable stable training in small data regimes, such as transfer learning [19], [20] or knowledge distillation [21]–[23], are not applicable to our problem framework since we are concerned with studying the behavior of deep learning methods for entirely new domains. Banerjee et al. [10] investigate the current state of machine learning in the context of rare diseases, where the authors propose several techniques to improve the dataset, such as harmonically combining data sets and reducing class imbalance with decision tree-based methods, as well as augmenting models, e.g., with knowledge graphs. Furthermore, applications of these techniques are summarized and analyzed. In our study, we take a general perspective on this difficult-to-train context, where the domain of rare diseases constitutes a specific case instead of focusing entirely on this field. Moreover, compared to Banerjee et al., we do so without prior knowledge or the transfer of weights. Dou et al. [24] provide a broad review of machine learning methods facing small data challenges in molecular sciences. Our work differs since we create a synthetic evaluation in which we systematically control the quantity of information available to the learner. Additionally, we limit our exploration of deep learning models and abstain from other machine learning techniques as we seek to process high-dimensional data, a domain in which deep learning methods have been excelling in recent years, lifting the requirement to carefully design feature extraction techniques as it is frequent for traditional machine learning models. Brigato et al. [16] identify the need to find models that can work under small data availability. In their pioneering study, the authors reveal the benefit of using models with lower than state-of-the-art complexity by investigating convolutional neural networks under different image classification evaluations. Contrary to their work, we do not consider using sophisticated data augmentation techniques that need to be to be hand-designed and require sufficient prior knowledge [16]. In contrast to other works on limited data [16], [17], where the authors only present a few models per study, we expand our examinations to consider a broader range of deep learning models. As we search the space of deep learning architectures exhaustively, our work is related to Neural Architecture Search (NAS). NAS extends hyperparameter optimization by additionally searching architectural parameters [18], [25], [26]. It is concerned with finding the optimal network architecture without relying on a researcher's

prior experience and freeing the process of required intuition by reducing the necessity of human intervention [26]. Our work builds on the principles of neural architecture search by considering a model's complexity and performance as two criteria for finding optimal architecture and tuning hyperparameters. We optimize hyperparameters using Bayesian and bandit-based search [27] to find suitable architectures and tuning hyperparameters. When performing our search concerning performance and complexity, we are interested in finding Pareto-optimal models resulting from the trade-off between these conflicting criteria. The Pareto frontier marks the boundary where we achieve an optimal trade-off between performance and model size such that we cannot find a model further optimized in one aspect without losing in the other [18].

III. METHODS

We base our study on three pillars: III-A the evaluation, III-B the search strategies, III-C the models. We now describe them starting with the evaluation.

A. Evaluation

In our experiments, we want to shed light on the performance and limitations of modern deep-learning techniques when faced with highly complex and noisy data under limited availability. To gain insights into the impact of the size of the dataset, we aim to show the limitations of such techniques when we systematically reduce the available data. More specifically, we consider the PTB-XL dataset [28], consisting of 21.837, 12-lead electrocardiogram recordings of 18.885 patients, having 71 labels in total [29] as our base data set, denoted as \mathcal{X} . The dataset contains an unbalanced label distribution, with possibly multiple labels per data point and noisy, high-dimensional time series data. We split \mathcal{X} into a training \mathcal{T} , validation \mathcal{V} , and test set \mathcal{U} in accordance with the stratified folds, as proposed in [28]. We use a downsampling rate, $\delta \in \{0.0, 0.2, 0.4, 0.6, 0.8, 0.85, 0.9, 0.95, 0.99\}$, to draw K subsets from \mathcal{T} uniformly at random. For each δ , the random subset fulfills the following condition

$$|\mathcal{T}_{\delta}^{i}| = \lfloor |\mathcal{T}|(1-\delta)\rfloor \tag{1}$$

with i = 1, ..., K, $|\cdot|$ denoting the size of a set and $\lfloor \cdot \rfloor$ the floor function. By training our models on the same fixed sets \mathcal{T}_{δ}^{i} , for all values of *i*, and thus, avoiding training our models on different sub-samples of the data, we ensure comparability across the training runs of the different models. Further, we consider the complete validation set to tune the parameters according to the true data distribution. The validation in the original data split has been ensured to have high-quality samples and a balanced label distribution [28]. Thus, our tuning data set concerns

$$\mathcal{D}^i_{\delta} = \{\mathcal{T}^i_{\delta}, \mathcal{V}\} \tag{2}$$

We determine the performances and limitations of modern deep learning models in a two-fold procedure, where we define performance in terms of macro-averaged area under the curve (macro AUC). To compute the macro AUC, each class AUC is calculated independently before averaging them. This accounts for a better judgment of the classifier than accuracy as it relieves the requirement to optimize for a threshold value [29]. Additionally, choosing micro-averaging, i.e., computing the average based on the combined true and false positive rates of each class, would lead to an overrepresentation of highly populated classes [28].

B. Search Strategies

The task of automatically finding the optimal hyperparameters, be it architectural or training settings, is an active field of research [27], [30]. In hyperparameter optimization frameworks, the aim is to find the optimal setting of hyperparameters without having a human-in-the-loop to eliminate possible biases and the need for prior experience regarding specific machine learning models. For this, trials are sampled according to user-specified performance metrics. For our purposes, we differentiate between bandit-based optimization via the Asynchronous Successive Halving Algorithm (ASHA) [31], and Bayesian optimization using the Multi-objective Tree Parsen Estimator (MTPE) [32].

1) Asynchronous Successive Halving Algorithm: ASHA [31] can be best understood as a best-arm identification problem in a multi-armed bandit set up. The algorithm samples hyperparameter configurations where each configuration corresponds to an arm. It aims to identify the best-performing arm and, where in this setting this arm corresponds to the best-performing hyperparameter configuration. Formally, given a set of hyperparameter configurations Θ , we want to find the single best-performing arm according to the evaluation function f,

maximize
$$f(\theta)$$

subject to $\theta \in \Theta$ (3)

In doing so, the algorithm allocates a uniform computational budget to a predetermined number of configurations. After each evaluation, poorly performing configurations are eliminated, and their computational budget is redistributed to the remaining trials. This procedure is called successive halving [33]. In ASHA, this elimination process occurs asynchronously by evaluating each arm as soon as possible instead of waiting for the remaining arms, repeating the process until only one trial is left. Due to its asynchronous nature, the algorithm efficiently enables parallelization and scalability.

2) Multi-objective Tree Parsen Estimator: MTPE [32] is a multi-objective hyperparameter optimization algorithm in the form of

maximize
$$f(\theta) = (f_1(\theta), ..., f_M(\theta))$$

subject to $\theta \in \Theta$ (4)

that, similar to ASHA, tries to find the best hyperparameter configuration, but instead of a single objective evaluation, it considers multiple objective functions. For this, a metric vector is constructed given an evaluation, $\zeta = f(\theta)$. Using this metric vector, containing the respective evaluations of each objective function, we can compare two hyperparameter settings by the concept of domination. Following² [32], a vector $\zeta \in \mathbb{R}^{M}$ dominates another vector $\zeta' \in \mathbb{R}^{M}$ if for all $i \in \{1, \dots, M\}$ it holds that $\zeta_i \geq \zeta_i$, and there exists $j \in \{1, \ldots, M\}$ such that $\zeta_j > \zeta_j$. We write domination as $\zeta \succ \zeta_j$, meaning that the metric vector ζ performs better in at least one metric $j \in \{1, \ldots, M\}$, while performing better or equal in the other metrics. Similarly, we define weak domination between two vectors when for all $i \in \{1, \ldots, M\}$ it holds that $\zeta_i \geq \zeta_i$, denoted as $\zeta \succeq \zeta'$. Further, a vector ζ dominates (or weakly dominates) a set of vectors $Z \subseteq \mathbb{R}^M$, denoted as $\zeta \succ Z$ (or $\zeta \succ Z$), if and only if for all $\zeta \prime \in Z$ it holds that $\zeta \succ \zeta \prime$ (or $\zeta \succeq \zeta l$). Likewise, set of vectors $Z \subseteq \mathbb{R}^M$ dominates (or weakly dominates) a vector $\zeta \in \mathbb{R}^M$, $Z \succ \zeta$ (or $Z \succeq \zeta$), if there exists $\zeta' \in Z$ such that $\zeta' \succ \zeta$ (or $\zeta' \succeq \zeta$ respectively). Following this, we define incomparability as $\zeta || Z$, i.e., neither $\zeta \succeq Z$ nor $\zeta \preceq Z$. With these concepts, we can define two densities, $l(\theta)$ and $g(\theta)$,

$$p(\theta|\zeta) = \begin{cases} l(\theta) \text{ if } \zeta \succ Z^* \lor \zeta || Z^* \\ g(\theta) \text{ if } \zeta \preceq Z^* \end{cases}$$
(5)

where Z^* is a set of vectors such that $p(\zeta \succ Z^* \lor \zeta || Z^*) = \gamma$. Here, Z^* contains trials that are inferior or incomparable to ζ , while $\gamma \in (0, 1)$ constitutes a threshold parameter that can be set by the user. Essentially, γ splits the hyperparameter configurations into two groups: those with good performance, measured by $l(\theta)$, and those with poor performance, measured by $g(\theta)$. This helps to distinguish between trials that improve over earlier ones and those that do not where new configurations are sampled by evaluating an acquisition function. In the case of MTPE, the corresponding acquisition function is the Expected Hypervolume Improvement (EHVI) given by

$$EHVI_{Z^*}(x) \propto \left(\gamma + \frac{g(\theta)}{l(\theta)}(1-\gamma)\right)^{-1}$$
 (6)

giving more emphasis to samples with a high probability under $l(\theta)$ and a low probability under $g(\theta)$. Subsequently, after evaluating the new hyperparameter configuration in terms of the objective functions $\{f_m\}_{m=1}^M$, the probability densities are updated, and a proceeding sample is chosen according to the acquisition function.

C. Models

We test a broad range of models to demonstrate the state of deep learning under limited data availability. We investigate the current best-performing models reported for the original PTB-XL dataset. Next, we question the performance of popular deep learning methods by extensively searching possible architectural settings.

²Please note that instead of following the original notation of [32], we decided to switch the signs to emphasize that we are dealing with a maximization problem.



Fig. 1: The number of Pareto optimal models decreases with higher down-sample rates.

1) State of the Art (SOTA): As a baseline, we consider the models reported in [34], [35] consisting of a state space model (S4) [34], [36], a one-dimensional XResnet model [37], and an LSTM-based model (CPC) [35]. Each of these models comes with a fixed architecture. First, we train the models on each \mathcal{T}^i_{δ} with their default training hyperparameters as described in [34], [35]. The training hyperparameters under investigation are the learning rate and weight decay as these have a crucial impact on a deep learning model's performance [27], and batch size, which has shown to directly correlate with the ability to generalize [38]–[40]. Next, we tune each of these hyperparameters using the bandit-based hyperparameter optimization ASHA to analyze their impact on the performance on smaller data samples. For our purposes, we do not tune the SOTA models on the multi-objective criterion as their model size does not change since we consider their architectures fixed.

2) Base models: In addition to the state-of-the-art models, we investigate a range of deep learning models, which we refer to as Base models, comprised of a convolutional neural network architecture (CNN), an LSTM-based architecture (LSTM), a transformer encoder (ENC), and a state-space model (S4) [36]. Furthermore, we explore the performance of deep Gaussian Process approximations using Random Fourier Features and convolutional layers (ConvRFF) [41].

We tune these models regarding architectural hyperparameters, such as the width and depth of the layers, activation functions, and dropout. Further, we tune them for the same tuning hyperparameters as the SOTA models using ASHA [31] to gain insights over the preferred architectures in low but high-dimensional data settings. Additionally, we tune the hyperparameters of the Base models using a multi-objective search procedure realized by the MTPE by maximizing for performance and minimizing for complexity.

IV. EXPERIMENTAL SET UP

We implement our experiments using PyTorch [42] and PyTorch Lightning [43]. We use RayTune [44] to implement ASHA and Optuna [45] for the MTPE. Our tuning is set up such that each sample of architectural or tuning hyperparameters constitutes a trial. An experiment consists of 100 trials per

TABLE I: An overview of the experiments. Each SOTA model has precisely one configuration. The number of Pareto trials varies for each model and each down-sample rate. In total, we conduct 1772 training runs.

Model	Configurations	Pareto Trials	Full Runs
S4	900	216	351
CPC	450	-	90
XResnet	450	-	90
CNN	900	312	357
ConvRFF	900	248	293
LSTM	450	207	252
ENC	450	294	339
Total	4500	1277	1772

model, both SOTA, and Base, for each \mathcal{T}^{i}_{δ} . We run each trial for 10 epochs and set K = 5 to report the mean and variances per model and per down-sample rate. After completing the tuning, we select the best hyperparameter according to their respective performances in the single objective optimization and every model that is part of the Pareto front for the multi objective search. We then train each of the best-performing models for 100 epochs. Subsequently, the models are tested on the complete testing data as we judge their quality to capture the true data distribution. All experiments (training, tuning and testing) are run on Nvidia A100 and P100 GPUs.

V. RESULTS

We search the space of training and architectural hyperparameters using the single objective ASHA and the multiobjective TPE and present an overview of the associated experiments in Table I. Each model is tuned for every downsample rate for each of the five data samples and 100 trials per optimization strategy, resulting in $100|\delta|K = 450$ trials per optimization strategy. In the case of ASHA, we select the top-performing model of each data sample to train it on the respective data sample. In contrast, since the multi-objective search does not result in a single best model but rather in a Pareto front of model configurations, we train each member of the Pareto front. In total, we train 1.772 models over 100 epochs, ranging over 9 different down-sample rate settings and five data samples per down-sample rate.

Figure 1 displays the number of Pareto models found per model and down-sample rate. Notably, the number of Pareto optimal models decreases together with the data.

Performance

At first, we compare the performances on the test data, \mathcal{U} , of each model, trained for each \mathcal{T}_{δ}^{i} . We depict their results in Figure 2 starting with the SOTA models on the left, the best-performing retuned models via the single-objective ASHA, and lastly, the models found by the multi-objective Bayesian search. In the case of the single objective hyperparameter tuning, we report that the performance of the SOTA models is higher than that of the Base models. However, the higher the down-sample rate, the smaller the gap becomes, resulting in a marginal performance advantage for the S4 (0.68 ± 0.07) compared to the CNN (0.68 ± 0.03) for a down-sample



Fig. 2: Performance is measured as macro AUC for all models, starting from the SOTA models (left), single-objectively tuned SOTA models (middle-left), single-objectively tuned Base models (middle-left), and multi-objectively tuned models (right) across all down-sample rates. Surprisingly, the models are stable throughout higher down-sample rates where the S4 model is superior throughout most data subsets and optimization strategies. For higher down-sample rates, models that are trained from scratch perform equal or better than models with specific architectures.



Fig. 3: Loss of performance when comparing each down-sample rate to the average performance on the entire data set, measured as macro AUC for all models, starting from the SOTA models (left), single-objectively tuned SOTA models (middle-left), single-objectively tuned Base models (middle-left), and multi-objectively tuned models (right) across all down sampling rates. The performance remains close to the baseline up to a down-sample rate of 40%.

rate of 99%. In the multi-objective case, the CNN surpasses the S4 for a down-sample rate of 60%, 95%, and 99%. For a down-sample rate of 99%, the Gaussian process-based model ConvRFF significantly closes the gap between the two models based on point estimates and even beats the S4 model regarding macro AUC. Notably, the performance for all models remains stable upon a loss of 40 percent of the sample size, after which it starts to decline, as shown in Figure 3. Overall S4 is the superior model for most of the down-sample rates. Figure 3 shows the difference in the average performance of each model given a down-sample rate of 0%, i.e., the complete data set, depicting the SOTA models (left), the single objective models (middle), and the multi-objective models (right). Remarkably, when comparing all three configurations, we see that the multi-objectively tuned models remain stable even with a down-sample rate of up to 60%, consistently outperforming the others at all lower down-sample rates.

Complexity vs. Performance

The trade-off between complexity and performance for the multi-objectively optimized models is shown in Figures 4, where we depict the Pareto front of a CNN model trained on $\mathcal{T}_{0.85}^3$ after tuning it for 10 epochs. In the case of the CNN in Figure 4, we demonstrate that a model's complexity directly influences its performance in the early epochs. However, we remark that in some trials the multi-objective search does not necessarily result in a linear or convex shape of the Pareto front, which would be preferable. Figure 5 shows the performances of fully trained models for all subset of the data using a down-sample rate of 85%. From this figure, we can observe that, given a fixed down-sample rate, the model size can be reduced without severely impacting the performance by carefully investigating the Pareto front.



Fig. 4: Pareto front of a CNN tuned on $\mathcal{T}_{0.85}^3$ while being optimized for performance and complexity. The higher the macro AUC and the lower the model size the better the model is. We would expect the Pareto front (red) to be of linear or convex shape.



Fig. 5: Size vs. performance of fully trained models after being optimized on $\mathcal{T}_{0.85}^i$ for performance and complexity simultaneously. Displayed for all i = 1...5. For most models, the size can be reduced without harming the performance.

Training Hyperparameters

Next, we investigate the role of the training hyperparameters. More specifically, we study whether the learning rate, weight decay, dropout, or batch size is impacted by the reduction of data samples. We note that the weight decay has no particular tendency except for the Bayesian optimisation procedure that saw an increase from 0.004 ± 0.01 for a downsample rate of 0% to 0.02 ± 0.03 given only 1% of the data set. The optimization algorithms set the batch size at a slightly lower value compared to the original batch size of 32 as reported in [34], particularly for higher down-sample rates (26.11 ± 31.77 to 21.30 ± 24.52 for a down-sample rate of 0% and 99% averaged over all search results). Overall, all three searches slightly increased the learning rates in their configurations compared to the learning rate used on the entire data set (0.005 ± 0.008 to 0.018 ± 0.022 again for $\delta = 0$ and $\delta = 0.99$ averaged over all searches). Lastly, the multiobjective search resulted in higher dropout rates by roughly 15% (0.24 ± 0.19 to 0.40 ± 0.22 , $\delta = 0$ and $\delta = 0.99$ respectively).

Overfit

Subsequently, we inspect the models for their susceptibility to overfit when reducing the amount of training data as seen in Figure 6. Here, we determine overfit as the training loss minus the validation loss, and likewise training macro AUC minus the validation macro AUC. Surprisingly, as seen in Figure 6a, the SOTA models are well-tuned on average with respect to overfitting in both loss and performance, exhibiting just a slight overfit on the loss. The models found by ASHA show a larger average overfit with respect to the performance values. In contrast, the multi-objective models demonstrate a comparably marginal overfit in terms of averaged loss and macro AUC. Figure 6b displays single trajectories of the overfit and loss per model tuned with single-objective and multi-objective search, respectively. In both cases, the trajectories indicate that for a subset of models, the overfit diverges more strongly from the average case. Additionally, we present the overfit in terms of loss and macro AUC for each of the K sampled datasets while setting $\delta = 0.99$ in Figures 7 and 8, which similarly demonstrate a tight confidence interval, while single trajectories diverge from the average. Overall, the multi-objective search criterion results in models that are welltuned on average.

VI. DISCUSSION

Our study compared the performance of state-of-the-art models and basic architectures via a systematic evaluation that tests a model's ability to deal with limited, highly dimensional, and noisy data. Furthermore, we extensively explored the space of architectural and training hyperparameters.

Each model suffers a performance loss when reducing the amount of data points. However, some models, such as the transformer encoder, are more affected than others (e.g., the state space model S4).

We found that the models are surprisingly robust to overfitting. Overall, we state that the models that are tuned with respect to performance and complexity are well-tuned (i.e., no overnor underfit) compared to their single-objective counterparts. Training hyperparameters, such as the learning rate and batch size, tend to be impacted by a reduced number of data points, corroborating earlier findings in the literature [39],



(a) Mean overfit in terms of loss and macro AUC with 95% confidence interval.



(b) Single trajectories.

Fig. 6: Overfit in terms of loss and macro AUC, calculated by subtracting the validation value from the training value, averaged over all down-sample rates, and displayed for all epochs using the average and a 95% confidence interval (a) and single trajectories (b). In the case of the loss, we normalized the training and validation values before subtracting the latter from the former. The plots compare single-objectively (left) and multi-objectively (right) tuned models.

particularly for the lower batch sizes [38], [40]. Weight decay only plays a marginal role in the single objective search, but tends to be increased when considering performance and complexity. Similarly, dropout becomes more important as the number of data points decreases, emphasizing the need to balance performance and model complexity. This finding is in accordance with the reports made by Brigato et al. [46]. We postulate that the superior performance of the state space model and the convolutional neural network-based models arise due to their inductive biases, specific for a given data type. This difference is noticeable in the single-objective

models depicted in both of the mid-columns of Figure 2,

where the remaining models fall off quickly for higher downsample rates. The ability of the state space model to capture long-range signals [36] appears to be beneficial to classify electrocardiogram signals, as demonstrated in the original PTB-XL evaluation [34]. Thus, one way of augmenting the inferior models would be to expand their inductive biases, e.g., as proposed by Yin et al., who introduce a convolutional bias into their transformer architecture [47].

We highlighted the importance of balancing performance and model size by advocating for a multi-objective model selection when confronted with limited yet high-dimensional and noisy data. Our experiments demonstrated that the performance of our multi-objectively tuned models remains more stable when compared with single-objectively tuned models, and is robust against overfitting. Furthermore, we have shown that, in low data regimes, a model's complexity can be significantly reduced without sacrificing its performance, as seen in Figure 5. When ready-made state-of-the-art models are not yet available in the corresponding literature, we propose tuning the architectural and training hyperparameters of a model from scratch, as it can result in a performance close to the best-performing models in the respective domain, especially in lower data set sizes. Albeit the accuracy may not be the best possible, the reduction in overfit is a strong argument in favor. Additionally, the resulting Pareto front allows for a case-sensitive model selection in which we favor one objective over the other.

In certain cases, however, the Pareto front can have a suboptimal shape, resulting in a set of models forming a concave curve. By averaging over these models, we obtain a new model that dominates the set of Pareto optimal models [48], meaning that it performs better in terms of performance while being of smaller complexity than the models of the Pareto front found to be optimal. Concretely, this means that if the Pareto front exhibits a concave shape, the overfit curves presented in Figure 6a overestimate the performance of the models. This is confirmed by analyzing Figures 6b-8, from which it is evident that some model trajectories diverge stronger than the average, e.g., the transformer-based model. However, the ConvRFF and S4 models are well-represented by their respective averages. Overall, our results indicate that the trajectories exhibit a reduction in overfitting, albeit with the possibility of degenerate cases.

Given the plethora of different hyperparameter settings, our work only explored the tip of the iceberg, nonetheless focused on the parameters that are most relevant for model creation. However, we used efficient search strategies to navigate the space of possible settings.

Lastly, we want to address why we have chosen to tune our models on the entire validation set for every down-sample rate instead of reducing the amount of validation similar to the training data. With this study, we want to reflect on the state of current deep learning methods when faced with reduced data in a high-dimensional and noisy setting in terms of capturing the true data distribution. While downsampling the validation data would be tailored specific to a case, we chose to emphasize



Fig. 7: Overfit in terms of loss with 95% confidence interval (upper) and the corresponding single trajectories (lower, calculated by subtracting the validation value from the training value, displayed for all epochs as single trajectories. The plots compare single-objectively (left) and multi-objectively (right) tuned models. The multi-objectively tuned models show smaller confidence intervals with a mean closer to zero while having single diverging trajectories.



Fig. 8: Overfit in terms of macro AUC with 95% confidence interval (upper) and the corresponding single trajectories (lower), calculated by subtracting the validation value from the training value, displayed for all epochs as single trajectories. The plots compare single-objectively (left) and multi-objectively (right) tuned models. The multi-objectively tuned models show smaller confidence intervals with a mean closer to zero while having single diverging trajectories.

the general capabilities of the deep learning methods across limited data samples. Using the entire validation data allows us to optimize the models according to the true data distributions and thus gives a better reflection of the model performances for low data regimes.

VII. CONCLUSION

In this study, we have analyzed the current state of deep learning methods in limited, high-dimensional, and noisy data settings. Specifically, we have introduced a novel approach that allows us to investigate models when reducing the data in a systematic manner. Furthermore, we have extensively searched the space of architectural and training hyperparameters using single-objective and multi-objective search. Our experiments revealed the importance of model selection according to multiobjective criteria to yield well-tuned models, competitive with state of the art methods. Moreover, we demonstrated the impact of the inductive bias of a model in limited, highdimensional and noisy data sets.

VIII. ACKNOWLEDGEMENT

S.J. gratefully acknowledges support from Fonds Wetenschappelijk Onderzoek (FWO) via FWO PhD Fellowship strategic basic research, Belgium 1SHHV24N. P.J.K.L. wishes to express gratitude for the support received from FWO via postdoctoral fellowship, Belgium 1242021N and the research council of the Vrije Universiteit Brussel (OZR-VUB via grant number OZR3863BOF). This research was supported by funding from the Flemish Government under the "Onderzoeksprogramma Artificiële Intelligentie (AI) Vlaanderen" program and through the IMAGIca project by the Interdisciplinary Research Program of the Vrije Universiteit Brussel (reference IRP8_b). Lastly, we want to thank the HPC administration and support service of Vrije Universiteit Brussel that helped tremendously during the experimental phase, Bart Bogaerts for providing us with essential feedback and guidance throughout the development of this research and finally, Bram Silue, Denis Steckelmacher, and Samuele Pollaci for proofreading.

REFERENCES

- T. Le Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić, D. Hesslow, R. Castagné, A. S. Luccioni, F. Yvon, M. Gallé *et al.*, "Bloom: A 176bparameter open-access multilingual language model," 2022.
- [2] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, "Llama: Open and efficient foundation language models," 2023.

- [3] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg, H. Nori, H. Palangi, M. T. Ribeiro, and Y. Zhang, "Sparks of artificial general intelligence: Early experiments with gpt-4," 2023.
- [4] A. Merchant, S. Batzner, S. S. Schoenholz, M. Aykol, G. Cheon, and E. D. Cubuk, "Scaling deep learning for materials discovery," *Nature*, vol. 624, no. 7990, pp. 80–85, 2023.
- [5] C. Chen, D. T. Nguyen, S. J. Lee, N. A. Baker, A. S. Karakoti, L. Lauw, C. Owen, K. T. Mueller, B. A. Bilodeau, V. Murugesan, and M. Troyer, "Accelerating computational materials discovery with artificial intelligence and cloud high-performance computing: from largescale screening to experimental validation," 2024.
- [6] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, "Zero-shot text-to-image generation," 2021.
- [7] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "Highresolution image synthesis with latent diffusion models," 2022.
- [8] K. Grace, H. Stewart, J. F. Sandkühler, S. Thomas, B. Weinstein-Raun, and J. Brauner, "Thousands of ai authors on the future of ai," 2024.
- [9] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901.
- [10] J. Banerjee, J. N. Taroni, R. J. Allaway, D. V. Prasad, J. Guinney, and C. Greene, "Machine learning in rare disease," *Nature Methods*, pp. 1– 12, 2023.
- [11] J. Schaefer, M. Lehne, J. Schepers, F. Prasser, and S. Thun, "The use of machine learning in rare diseases: a scoping review," *Orphanet journal* of rare diseases, vol. 15, pp. 1–10, 2020.
- [12] H. Altae-Tran, B. Ramsundar, A. S. Pappu, and V. Pande, "Low data drug discovery with one-shot learning," ACS Central Science, vol. 3, no. 4, pp. 283–293, 2017.
- [13] A. Karamitrou, F. Sturt, P. Bogiatzis, and D. Beresford-Jones, "Towards the use of artificial intelligence deep learning networks for detection of archaeological sites," *Surface Topography: Metrology and Properties*, vol. 10, no. 4, p. 044001, oct 2022.
- [14] A. H. Jamil, F. Yakub, A. Azizan, S. A. Roslan, S. A. Zaki, and S. A. Ahmad, "A review on deep learning application for detection of archaeological structures," *Journal of Advanced Research in Applied Sciences and Engineering Technology*, vol. 26, no. 1, p. 7–14, Jan. 2022.
- [15] W. Zhang, L. Deng, L. Zhang, and D. Wu, "A survey on negative transfer," *IEEE/CAA Journal of Automatica Sinica*, vol. 10, no. 2, pp. 305–329, 2023.
- [16] L. Brigato and L. Iocchi, "A close look at deep learning with small data," in 2020 25th International Conference on Pattern Recognition (ICPR). IEEE, 2021, pp. 2490–2497.
- [17] S. Greydanus and D. Kobak, "Scaling down deep learning with mnist-1d," 2024.
- [18] G. Menghani, "Efficient deep learning: A survey on making deep learning models smaller," *Faster, and Better. arXiv*, vol. 2106, 2021.
- [19] K. Weiss, T. M. Khoshgoftaar, and D. Wang, "A survey of transfer learning," *Journal of Big data*, vol. 3, no. 1, pp. 1–40, 2016.
- [20] R. Zhang, J. Han, A. Zhou, X. Hu, S. Yan, P. Lu, H. Li, P. Gao, and Y. Qiao, "Llama-adapter: Efficient fine-tuning of language models with zero-init attention," *arXiv preprint arXiv:2303.16199*, 2023.
- [21] P. Kaliamoorthi, A. Siddhant, E. Li, and M. Johnson, "Distilling large language models into tiny and effective students using pqrnn," *CoRR*, vol. abs/2101.08890, 2021.
- [22] W. Wang, F. Wei, L. Dong, H. Bao, N. Yang, and M. Zhou, "Minilm: Deep self-attention distillation for task-agnostic compression of pretrained transformers," 2020.
- [23] Z. Sun, H. Yu, X. Song, R. Liu, Y. Yang, and D. Zhou, "Mobilebert: a compact task-agnostic BERT for resource-limited devices," *CoRR*, vol. abs/2004.02984, 2020.
- [24] B. Dou, Z. Zhu, E. Merkurjev, L. Ke, L. Chen, J. Jiang, Y. Zhu, J. Liu, B. Zhang, and G.-W. Wei, "Machine learning methods for small data challenges in molecular science," *Chemical Reviews*, vol. 123, no. 13, pp. 8736–8780, 2023.

- [25] T. Elsken, J. H. Metzen, and F. Hutter, "Neural architecture search: A survey," *The Journal of Machine Learning Research*, vol. 20, no. 1, pp. 1997–2017, 2019.
- [26] P. Ren, Y. Xiao, X. Chang, P.-Y. Huang, Z. Li, X. Chen, and X. Wang, "A comprehensive survey of neural architecture search: Challenges and solutions," ACM Computing Surveys (CSUR), vol. 54, no. 4, pp. 1–34, 2021.
- [27] M. Feurer and F. Hutter, "Hyperparameter optimization," Automated machine learning: Methods, systems, challenges, pp. 3–33, 2019.
- [28] P. Wagner, N. Strodthoff, R.-D. Bousseljot, D. Kreiseler, F. I. Lunze, W. Samek, and T. Schaeffter, "Ptb-xl, a large publicly available electrocardiography dataset," *Scientific data*, vol. 7, no. 1, p. 154, 2020.
- [29] N. Strodthoff, P. Wagner, T. Schaeffter, and W. Samek, "Deep learning for ecg analysis: Benchmarks and insights from ptb-xl," *IEEE Journal* of Biomedical and Health Informatics, vol. 25, no. 5, pp. 1519–1528, 2021.
- [30] T. Yu and H. Zhu, "Hyper-parameter optimization: A review of algorithms and applications," *CoRR*, vol. abs/2003.05689, 2020.
- [31] L. Li, K. Jamieson, A. Rostamizadeh, E. Gonina, J. Ben-Tzur, M. Hardt, B. Recht, and A. Talwalkar, "A system for massively parallel hyperparameter tuning," *Proceedings of Machine Learning and Systems*, vol. 2, pp. 230–246, 2020.
- [32] Y. Ozaki, Y. Tanigaki, S. Watanabe, and M. Onishi, "Multiobjective treestructured parzen estimator for computationally expensive optimization problems," in *Proceedings of the 2020 Genetic and Evolutionary Computation Conference*, ser. GECCO '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 533–541.
- [33] K. Jamieson and A. Talwalkar, "Non-stochastic best arm identification and hyperparameter optimization," in *Artificial intelligence and statistics*. PMLR, 2016, pp. 240–248.
- [34] T. Mehari and N. Strodthoff, "Advancing the state-of-the-art for ecg analysis through structured state space models," *arXiv preprint* arXiv:2211.07579, 2022.
- [35] —, "Self-supervised representation learning from 12-lead ecg data," *Computers in biology and medicine*, vol. 141, p. 105114, 2022.
- [36] A. Gu, K. Goel, and C. Ré, "Efficiently modeling long sequences with structured state spaces," arXiv preprint arXiv:2111.00396, 2021.
- [37] T. He, Z. Zhang, H. Zhang, Z. Zhang, J. Xie, and M. Li, "Bag of tricks for image classification with convolutional neural networks," *CoRR*, vol. abs/1812.01187, 2018.
- [38] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang, "On large-batch training for deep learning: Generalization gap and sharp minima," *CoRR*, vol. abs/1609.04836, 2016.
- [39] F. He, T. Liu, and D. Tao, "Control batch size and learning rate to generalize well: Theoretical and empirical evidence," Advances in neural information processing systems, vol. 32, 2019.
- [40] I. Kandel and M. Castelli, "The effect of batch size on the generalizability of the convolutional neural networks on a histopathology dataset," *ICT express*, vol. 6, no. 4, pp. 312–315, 2020.
 [41] T. Wang, L. Xu, and J. Li, "Sdcrkl-gp: Scalable deep convolutional
- [41] T. Wang, L. Xu, and J. Li, "Sdcrkl-gp: Scalable deep convolutional random kernel learning in gaussian process for image recognition," *Neurocomputing*, vol. 456, pp. 288–298, 2021.
- [42] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Z. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, highperformance deep learning library," *CoRR*, vol. abs/1912.01703, 2019.
- [43] W. Falcon and The PyTorch Lightning team, "PyTorch Lightning," Mar. 2019.
- [44] R. Liaw, E. Liang, R. Nishihara, P. Moritz, J. E. Gonzalez, and I. Stoica, "Tune: A research platform for distributed model selection and training," arXiv preprint arXiv:1807.05118, 2018.
- [45] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A next-generation hyperparameter optimization framework," 2019.
- [46] L. Brigato and L. Iocchi, "A close look at deep learning with small data," in 2020 25th International Conference on Pattern Recognition (ICPR), 2021, pp. 2490–2497.
- [47] M. Yin, Z. Chang, and Y. Wang, "Adaptive hybrid vision transformer for small datasets," in 2023 IEEE 35th International Conference on Tools with Artificial Intelligence (ICTAI). IEEE, 2023, pp. 873–880.
- [48] F. Felten, E.-G. Talbi, and G. Danoy, "Multi-objective reinforcement learning based on decomposition: A taxonomy and framework," 2024. [Online]. Available: https://arxiv.org/abs/2311.12495