

## Which Design Decisions in AI-enabled Mobile Applications Contribute to Greener AI?

Roger Creus Castanyer · Silverio  
Martínez-Fernández · Xavier Franch

Received: date / Accepted: date

**Abstract Background:** The usage of complex artificial intelligence (AI) models demands expensive computational resources. While currently available high-performance computing environments can support such complexity, the deployment of AI models in mobile devices, which is an increasing trend, is challenging. Environments with low computational resources imply limitations in the design decisions during the AI-enabled software engineering life-cycle that balance the trade-off between the accuracy and the complexity of the mobile applications.

**Objective:** Our objective is to systematically assess the trade-off between accuracy and complexity when deploying complex AI models (e.g. neural networks) to mobile devices in pursuit of greener AI solutions. We aim to cover *(i)* the impact of the design decisions on the achievement of high-accuracy and low resource-consumption implementations; and *(ii)* the validation of profiling tools for systematically promoting greener AI.

**Method:** We implement neural networks in mobile applications to solve multiple image and text classification problems on a variety of benchmark datasets.

---

Roger Creus Castanyer  
Universitat Politècnica de Catalunya  
Barcelona, Spain  
E-mail: roger.creus@estudiantat.upc.edu  
ORCID: 0000-0003-1952-3357

Silverio Martínez-Fernández  
Universitat Politècnica de Catalunya  
Barcelona, Spain  
E-mail: silverio.martinez@upc.edu  
ORCID: 0000-0001-9928-133X

Xavier Franch  
Universitat Politècnica de Catalunya  
Barcelona, Spain  
E-mail: xavier.franch@upc.edu  
ORCID: 0000-0001-9733-8830

We then profile and model the accuracy, storage weight and time of CPU usage of the AI-enabled applications in operation with respect to their design decisions. Finally, we provide an open-source data repository following the EMSE open science practices and containing all the experimentation, analysis and reports in our study.

**Results:** We find that the number of parameters in the AI models makes the time of CPU usage scale exponentially in convolutional neural networks and logarithmically in fully-connected layers. We also see the storage weight scales linearly with the number of parameters, while the accuracy does not. For this reason, we argue that a good practice for practitioners is to start small and only increase the size of the AI models when their accuracy is low. We also find that Residual Networks (ResNets) and Transformers have a higher baseline cost in time of CPU usage than simple convolutional and recurrent neural networks. Finally, we find that the dataset used for experimentation affects both the scaling properties and accuracy of the AI models, hence showing that researchers must study the presented set of design decisions in each specific problem context.

**Conclusions:** We have depicted an underlying and existing relationship between the design of AI models and the performance of the applications that integrate these, and we motivate further work and extensions to better characterise this complex relationship.

**Keywords** AI-enabled Applications · Mobile Applications · Model Accuracy · Application Performance · Greener AI · Neural Networks

## 1 Introduction

Artificial Intelligence (AI) plays a key role in the world we live in. AI is about making machines mimic intelligent behaviour. Machine Learning (ML) sets the foundations of AI [1]. ML works over the most valuable and key source of knowledge that exists: data. In ML, machines take in data and learn patterns that would be difficult for humans to learn. ML is valuable in the sense that the data processing (or data classification, segmentation, representation) capabilities go far beyond than what humans can achieve. ML finds its extension in Deep Learning (DL), which introduces complex neural network (NN) models with the purpose to provide highly performing capabilities in an increasing number of challenging domains [2] (e.g. image-based and language-based contexts). Nowadays, NNs have achieved outstanding results and have outperformed humans in domains as diverse as machine translation, character and handwriting recognition, speech and facial recognition and videogames [3–6]. Key to the success of NNs is their capability of learning complex data representations in an unsupervised manner. In this way, NNs do not demand expert feature engineering tasks, which are a bottleneck of ML models performance. However, training NNs requires both large amounts of data and high computational resources. These two needs are not a stopper for NNs since we

are living the *Big Data, Big Compute, Big Models* revolution [7–9]. As a consequence, following the exponential evolution of the availability of data and computational power, the tendency of solving complex tasks with complex models is increasing [9–12]. This tendency does not take into account resource consumption regulations, and does not promote less data-intensive and lighter models, which are considerations that set the bases of greener AI. For example, recent language models like the GPT-3 have 175 billion parameters [9], and benchmarks for image classification have been trained for more than 6 months [13].

For *green AI* we understand AI research and practices that are more environmentally friendly and inclusive [14]. Green AI refers to the development and deployment of AI systems that minimize their negative impact on the environment and promote sustainable energy and resource consumption [15]. The development of greener AI can be achieved through various means such as optimizing the energy efficiency of hardware used for AI processing, reducing the carbon footprint of data centers that host AI systems, and developing algorithms that promote sustainable practices. In contrast, *red AI* refers to AI research that seeks to improve accuracy (or related measures) through the use of massive computational power while disregarding the environmental cost [14]. Ensuring efficient implementations is key to achieve greener AI. Some of the key practices of green AI are to measure and lower the environmental impact of AI [16], regulate the resource consumption of the AI-enabled applications [17] and balance the effectiveness and efficiency of the AI models [18, 19]. Central to these practices is the capability of measuring and monitoring the efficiency-related metrics of AI-enabled applications in operation, which allows developers to gain awareness of the performance of AI-enabled applications with respect to their complexity.

Deploying AI models to operate in environments with low computational resources constrains the design decisions in AI-enabled applications<sup>1</sup>. One of such low computational resources environments are mobile applications. With the growth of the mobile applications market, approaches on the integration of AI models in lightweight devices are becoming popular given their ability to provide great services to the user [20]. In the context of AI-enabled mobile applications, reducing energy consumption, time of CPU usage (i.e. in CPU or GPU) and storage weight is fundamental for deploying a user-friendly experience [21]. However, AI practitioners aim at providing high-accuracy results, which might involve deploying computationally demanding AI models, so balancing the efficiency and effectiveness of AI-enabled mobile applications is key to pursuing greener AI. In this paper, we approach green AI solutions by developing AI-enabled mobile applications whose focus is to maximize their accuracy while keeping their complexity under control.

In this work we use mobile devices as environments that support the deployment of AI models. Few frameworks are available that enable the end-to-end

---

<sup>1</sup> note that in this work, whenever we relate to AI models or AI components we mean NNs, since these are the ones that we study, develop and test in this paper

implementation of AI models into mobile applications, and such frameworks have a strong influence in the overall applications performance [22]. For this reason, there exists the need of a systematic approach for measuring the performance of an AI-enabled mobile application that takes into account the accuracy and complexity-related metrics.

With this, we aim to model and investigate the contribution of design decisions during the AI-enabled software engineering lifecycle to the overall mobile applications performance in terms of accuracy and complexity. Our objective is to report statistically significant relationships between the design decisions and their impact in order to build the bases for predictive tools that empower practitioners with consciousness of the performance of their AI-enabled applications.

This article is structured as follows. Section 2 introduces the Related Work; Section 3 describes the experimental design in detail (i.e. variables, hypotheses, research questions); Section 4 presents the experiment's results and findings; in Section 5 we discuss the implications of our work for AI engineering research and industry; in Section 6 we identify the threats to validity; and finally in Section 7 we present our conclusions and motivate future work on the topic.

## 2 Related work

Engineering AI-enabled applications that operate in edge devices is challenging [20, 23–25]. By edge devices, we understand any device that can host an AI-enabled application so that the AI models run in local machines (e.g. PCs, laptops, mobile phones or even IoT devices) and not in the cloud. Independently of the platform where the AI-enabled software is deployed, the development lifecycle for AI-enabled software, which is significantly different from the traditional software one [26], can be split into 4 phases. First, the *Data Management* which consists of the collection and processing of raw data. Second, the *Modelling* of the AI components which consists of adjusting different models to obtain the best-performing possible solution. Third, the *Development* of the environment where the models will be deployed and the integration of the AI models in it. Fourth, the *System Operation* of the AI-enabled applications. This lifecycle introduces concerning challenges that are not usual in traditional software engineering [27]. Data centrality is what increases the risk in many quality attributes (e.g. adaptability, scalability, explainability, privacy) [24]. As improving the sustainability of the aforementioned AI-based software engineering lifecycle has never been a key requirement historically, now there exists a high potential for improving the efficiency of the lifecycle and the related development tools [28].

### 2.1 Research to build Greener AI systems

In pursuit of greener AI, several works have put the focus on improving energy efficiency [29]. In the context of traditional mobile applications (without

AI components), it has been shown that different works deal with energy efficiency using optimizations of different nature [30]. For example, leveraging automatic refactoring tools to fix specific code smells in mobile applications has been shown to improve energy efficiency in practice [31–33]. Another approach to study and optimize energy efficiency consists of monitoring the application’s performance using profilers [34]. Monitoring energy consumption, which has been studied from the software (e.g. energy profilers) and hardware (e.g. power monitors) viewpoints, is not trivial [34,35]. In the majority of cases, when measuring energy consumption, one cannot isolate the software from the system where it is running. The latter can be achieved by making sure that no unnecessary applications or processes are running on the device<sup>2</sup>. Some examples of energy profilers are the Intel Power Gadget and Intel PowerLog, which provide detailed energy consumption profiled data (in Joules) plus CPU utilisation (in %) and power (in Watt). However, note that these tools are meant to run on a computer machine and not on a mobile device. For this reason, in this work we adopt a technology stack that allows practitioners to both deploy AI-enabled applications into mobile devices and monitor the CPU utilization directly in the edge devices using profilers.

More recently, researchers have found that software energy consumption can be modeled (and thus, estimated) accurately [34, 36]. The CPU usage information has been demonstrated to be very valuable for achieving accurate energy consumption estimations [36, 37]. In this way, the CPU usage, which is popular for being easily monitored, is a great representative of the energy cost of an application. In light of the results learned from these works, we include the time of CPU usage as one of the dependent variables of the study and define a set of design decisions and provide a novel interpretation of the relationship between these.

In the context of AI-enabled mobile applications, studies on the frameworks for developing the AI components have been proven to cause differences in the overall performance of the applications. Concretely, two popular frameworks which are Pytorch and Tensorflow have been found to carry differences in the accuracy [22] and also in green characteristics like the energy consumption [37]. In this work, we also use PyTorch and Tensorflow to implement the AI models.

Green AI solutions have also been studied from model-centric [33] and data-centric viewpoints. In this context, it has been observed empirically that optimizations on the datasets can be performed to reduce the energy consumption of the AI models without significant losses of accuracy [38]

Compared to the aforementioned works, we adopt a hybrid approach between the model-based and data-centric approaches for constraining resource consumption. In this context, our study is novel in exploring design decisions for both the AI models and the datasets used and how these affect the greenability (i.e., energy efficiency) of the AI-enabled mobile applications.

---

<sup>2</sup> <https://luiscruz.github.io/2021/07/20/measuring-energy.html>

### 3 Experimental Design

#### 3.1 Research Objectives

We define our goal and Research Question (RQ) following the GQM guidelines [39].

Overall, our goal is to analyse *the design decisions of AI-enabled applications* with the purpose to *assess their impact* with respect to *green characteristics* from the point of view of *a developer* in the context of *mobile applications with neural networks for image and language classification*.

We determine the units of analysis and formulate the RQs as follows. As aforementioned above, we work towards deploying green AI-based mobile applications by putting the focus on optimising their performance, defined as the trade-off between accuracy and complexity. We aim to *(i)* evaluate and quantify the implications of a set of key design decisions that we identify along the AI-enabled software lifecycle with respect to the overall application performance; and *(ii)* leverage profiling tools for measuring the performance and support the decision-making in the design of AI-enabled applications.

We establish research objectives, variables and hypotheses from which we perform an empirical study. In particular, we study the systematic variation of design decisions to measure their influence on the accuracy and complexity-related metrics in an experimental set-up. For this purpose, we implement and monitor AI-enabled applications in the image and language classification domains. We profile the performance of the applications by taking into account their time of CPU usage, storage weight and the accuracy in operation. Finally, we define linear regression equations for measuring the implications of the design decisions to the AI-enabled mobile applications performance. In this sense, our research objectives consist on finding statistically significant relations between dependent and independent variables. As no statistical method can show causation, we aim to provide a well-suited regression analysis for minimizing the error in the causal interpretation of our models [40].

With all this, we identify two design decisions that are related to the AI models, which are the number of parameters and the architecture type, and also one related to the data which is the dataset used for training the AI models. The number of parameters represents the total sum of weights and biases that operate inside the AI models (i.e. the NNs). The architecture type controls the degree of complexity introduced in the models and correlates with the amount of computational power provided. The dataset used for training conditions the quality and quantity of the data that feeds the models, hence affecting the models' performance (e.g. in terms of accuracy).

In this way, whenever we refer to the design decisions in this paper, we refer to this set of three variables that we have identified. Then, as stated above, we characterise the performance of the AI-enabled mobile application in terms of the accuracy and complexity trade-off and for that we measure the time of CPU usage, storage weight and the accuracy of the applications in operation. For these reasons, we establish the RQs as follows:

- RQ1 - What is the impact of the design decisions on the performance of AI-enabled applications in terms of accuracy and complexity?
  - RQ1.1 - How is the time of CPU usage of the AI models affected by their number of parameters, architecture type and dataset used?
  - RQ1.2 - How is the storage weight of the AI models affected by their number of parameters, architecture type and dataset used?
  - RQ1.3 - How is the accuracy of the AI models affected by their number of parameters, architecture type and dataset used?

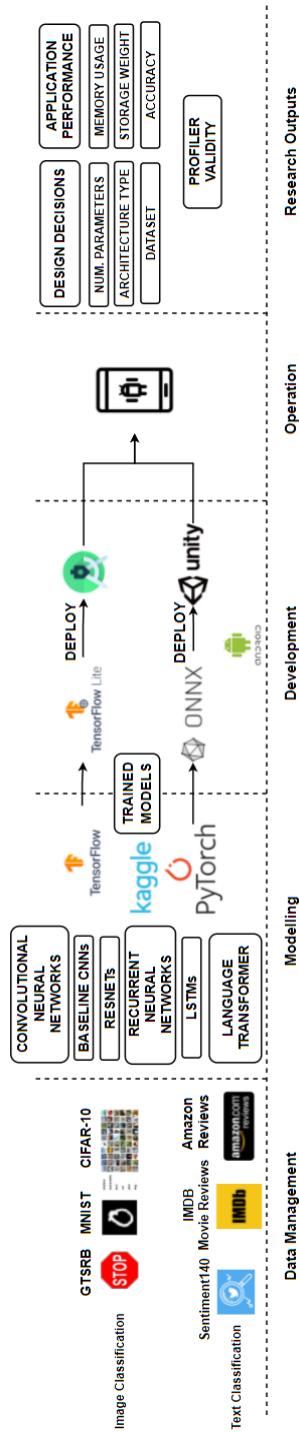
With *RQ1* we aim to validate our criterion for reasoning about high-performance implementations of AI models in AI-enabled software applications (i.e., which provide high-accuracy results while keeping complexity at a low level). We aim to provide a quantitative analysis of the importance and relation between design variables and performance metrics in terms of accuracy and complexity. Concretely, we aim to deploy a quantitative evaluation of the performance of the AI-enabled applications in operation. With this, our goal is to delineate the differences in performance caused by the AI models configuration (i.e. architecture type, number of parameters) and dataset used, and to validate them across the vision and language classification domains.

As we frame our study to the statistical modelling of the green characteristics from a set of design decisions, the validity of our hypotheses aims to be independent of whether we prove an existent relationship between our objects of study or we discard it. In this way, we deny any confirmation bias related to prior beliefs of ours on the possibly existing relations between the variables in the study. Moreover, we perform repeated experimentation, so the results in the profiled metrics and accuracy measurements are averaged across at least 25 examples which can be found in our public repository<sup>3</sup>. Finally, we discuss and perform the statistical analysis of the results in consensus between all the authors to generalize our conclusions and remove personal interpretation biases.

The development schema of the empirical study can be seen in Figure 1.

---

<sup>3</sup> <https://github.com/roger-creus/Which-Design-Decisions-in-AI-enabled-MobileApplications-Contribute-to-Greener-AI/tree/main>



**Fig. 1** Schema of our empirical study in the end-to-end software engineering lifecycle for AI-enabled applications. In the *(i)* Data Management phase we obtain and preprocess the data sources; in the *(ii)* Modelling phase we train the AI models for solving the defined tasks using the free NVIDIA TESLA P100 GPU hours offered weekly by Kaggle ; in the *(iii)* Development phase we convert and deploy the trained models in the applications that we build; in the *(iv)* Operation phase we make use of the applications and monitor their performance within the operation settings by means of profilers. Finally, we provide answers to the RQs.



## 3.2 Variables

In the following subsections, we respectively report the independent, dependent, and confounding variables of our experimental design.

### 3.2.1 Independent Variables

In this study, we define multiple independent (i.e. explanatory) variables in order to evaluate the contribution of the design decisions to the overall applications performance.

Regarding the data management phase of the software engineering lifecycle for AI-enabled applications, we experiment with a set of different datasets for each of the vision and language domains. In this way, we define the dataset (D) nominal variable indicating which dataset is used for training the AI models. Implicitly, each dataset provides different quantity and quality of data and this certainly affects the performance of the models. Concretely, the *dataset* takes values in a subset of the most popular benchmark datasets in the vision and language domains: the German Traffic Sign Recognition Benchmark (GTSRB) dataset, MNIST, and CIFAR for solving image classification, and Sentiment140, IMDB Movie Reviews and Amazon Reviews for text classification (see Section 3.4.1 for a detailed description of the datasets). We chose these datasets since they are well-known benchmarks which can be solved with NNs and at the same time carry significant differences in terms of the quality and quantity of the data (e.g. CIFAR contains 32x32 RGB images and MNIST contains 28x28 grayscale images). The presented datasets have many thousands of monthly downloads in HuggingFace <sup>4</sup>, and as mentioned above, we believe that all together they form a rich population of diverse datasets in terms of quality in the quantity of data.

Regarding the modelling phase of the AI components, we experiment with *(i)* different number of parameters (P), which is a numerical variable indicating the complexity of the AI models; and *(ii)* different architecture types (AT), encoded as a nominal variable characterizing different model architectures (in a subset of the SOTA architectures formed by the LSTM (Long Short-Term Memory), GRU (Gated Recurrent Unit), Transformer, CNN (Convolutional Neural Network), ResNet (Residual Networks) and MobileNets). The idea behind studying the *number of parameters* is to measure how scaling up the AI models in terms of computational power helps with achieving better accuracy and how it impacts the performance of the applications. Nowadays, many successes in the AI literature rely on very deep AI models trained in computational environments that are inaccessible to the vast majority of researchers. For this reason, we aim to understand the role of the scale of the AI models in the overall AI-enabled applications to help practitioners make the best use of it in any environment. With the *architecture type*, we aim to study how more complex architectures in both the vision and language domains (e.g.

---

<sup>4</sup> <https://huggingface.co/datasets?sort=downloads>

ResNets, Transformers) help achieving better accuracy or performance in the AI-enabled applications. These architectures are based on the simple convolutional and recurrent neural networks (that we also include in our study) but usually involve some novel operations (e.g. residual layers in ResNets, attention mechanism in Transformers) that are worth studying.

### 3.2.2 Dependent Variables

We evaluate the contribution of the design decisions with respect to the overall performance of the built applications in operation, which we study as a balance between the accuracy and the complexity. To define the performance of AI-enabled applications, we consider three dependent variables, the first two referring to the concept of complexity and the latter representing the accuracy.

- *Time of CPU usage* (M), which quantifies the resource consumption that the applications require to run (in ms).
- *Storage weight* (S), which quantifies the cost of deploying the AI models to the applications and storing them (in MB).
- *Accuracy* (A), which measures how capable is the application (and the AI model) of providing high-accuracy and satisfactory results in operation (% of correctly classified inputs across an evaluation dataset <sup>5</sup>).

Importantly, in our experiments we only report the time of CPU usage and do not study the impact of using GPUs or other modern dedicated AI-processing cores (e.g. Neural Processing Units). As can be seen later in the results section (e.g. Figure 17), all the AI model usage falls under the CPU module.

Regarding accuracy, we plan to measure the accuracy of the AI models when operating in the AI-enabled applications during the system operation phase of the AI-enabled software lifecycle. That is, we test the accuracy of the models with data outside the context of the training datasets to test the models' generalization capabilities. Concretely, we craft a test dataset for each domain (vision and language) containing data from the operation environment with a sufficiently large number of examples to test the applications systematically and provide statistically significant conclusions on the accuracy achieved. For the image domain we plan to report the top-3 accuracy (which considers a model output to be accurate if it contains the actual class in the top-3 predicted classes by the model when ranked by confidence), as we consider that the top-1 accuracy might be too demanding for large multi-class classification problems). For the language (i.e. text) domain we report the top-1 accuracy as text classification is generally a problem with a smaller number of classes.

### 3.2.3 Confounding Variables & others

We identify a confounding factor related to the experience and knowledge of the developer who builds the AI models and applications, since it conditions

<sup>5</sup> <https://github.com/roger-creus/Which-Design-Decisions-in-AI-enabled-MobileApplications-Contribute-to-Greener-AI/tree/main>

the overall performance of the AI-enabled applications in terms of their accuracy and complexity. A developer is responsible for the implementation of both multiple sophisticated AI models and mobile applications that integrate them. Moreover, he is also in charge of operating with the AI-enabled mobile applications, profiling and analyzing the applications' performance during the system operation phase. We consider that the more experience and knowledge the developer has, the less overhead in the AI-enabled applications. Ideally, as we aim to profile the performance of AI-enabled applications, we want to account only for the performance changes produced by our objects of study and not by deficient implementations. For mitigating the latter we provide as simple implementations as possible by only implementing the required functionalities to the application and by following official tutorials for the development of these.

**Table 1** The variables of the study

Class	Name and Abbreviation	Description	Scale	Operationalization
Independent	Number of Parameters (P)	The number of trainable parameters of the NNs	numerical	See Sections 3.2 and 3.3
	Architecture Type (AT)	The architecture that defines the NNs functionality	nominal	See Sections 3.2 and 3.3
	Dataset (D)	The dataset used for training the AI models	nominal	See Section 3.4.1
Dependent	Time of CPU usage (M)	The time that the AI-enabled applications need to answer a query	numerical	Profiled
	Storage Weight (S)	The cost of storing the AI models in the edge devices	numerical	Retrieved from the built applications
	Accuracy (A)	The precision of the AI models when answering queries in the edge devices	numerical	Measured in the system operation phase (see Section 3.2.2)
Confounding & others	Developer experience	The ability to implement high-quality software applications	-	-
	AI modelling framework	Technology used for modelling the AI components	-	Experiment with Pytorch and Tensorflow
	Export AI components	Technology used for exporting the AI components	-	Experiment with ONNX and Tensorflow Lite
	AI-enabled development framework	Technology used for developing the AI-enabled applications	-	Experiment with Unity3D and Android Studio

We also consider the technology stack used for deploying our experimentation a factor that conditions the results obtained. For that we experience with two different frameworks for modelling the AI components which are Pytorch and Tensorflow; we use two different technologies for exporting the trained AI models which are ONNX and Tensorflow Lite; and we also use two different frameworks for developing the AI-enabled mobile applications: Unity3d and Android studio. With this pipeline, we provide experiences with two different technology stacks that enable practitioners to deploy AI models in mobile applications. In this way, we can assess if there exists noticeable performance differences between the applications developed with different technologies and we can also suggest the most suitable tools and practices for developers willing to integrate and profile AI models in mobile applications.

### 3.3 Hypotheses and Analysis Operationalization

In this study we hypothesize that there is an underlying relation between the design decisions of AI-enabled software and their consequences and we aim to characterize it. Concretely, we study the accuracy (A) and the complexity (M, *storage weight*) as a model of the design decisions (*number of parameters, architecture type, dataset*), with separate hypothesis for the vision and language domains. The models that we consider are linear models between the design decisions (i.e. independent, explanatory variables) and the accuracy and complexity-related metrics (i.e. dependent, response variables). When fitting linear models we provide the following capabilities:

- Determining the goodness of fit by quantifying the amount of variability of the dependent variables that is captured (i.e. explained) by the set of independent variables.
- Using the fit as a predictive model for regressing accuracy and complexity responses for untested experimental conditions.
- Quantifying the strength of the relationship between the dependent and independent variables, and determining whether some independent variables may have no linear relationship with the responses at all, or identifying which subsets of independent variables may contain redundant information about the responses.
- Providing directional analysis on what independent variables contribute more to the target dependent variables.
- Determining whether the interactions between pairs of dependent variables significantly contribute to the independent variables.

With all this, we define the models as follows:

$$M = \beta_{P_1}P + \beta_{AT_1}AT + \beta_{D_1}D + \epsilon_1 \quad (1)$$

$$S = \beta_{P_2}P + \beta_{AT_2}AT + \beta_{D_2}D + \epsilon_2 \quad (2)$$

$$A = \beta_{P_3}P + \beta_{AT_3}AT + \beta_{D_3}D + \epsilon_3 \quad (3)$$

With equations 1, 2, 3 we define multiple linear regressions to explain the variability of the time of CPU usage, storage weight and accuracy with respect to the design decisions separately (1: M; 2: S; 3: A) [41]. For example, with  $\beta_{P_1}$  we study the impact of the number of parameters only to the time of CPU usage. Ideally, we want all the variables  $\epsilon_i$  to be of low magnitude, as it represents the unexplained variance of the dependent variables.

We fit each of the three multiple linear models with the least-squares estimation technique. In particular, we fit the coefficients  $\beta_i$  (for each of the design decisions) to minimize the squared error between model predictions and observed values (belonging to the data collected in our experimentation):

$$\hat{\beta}_i = \arg \min_{\beta_i} \sum_{n=1}^N (\beta_i * \mathbf{x}_n - \mathbf{y}_{n_i})^2 \quad (4)$$

where  $\hat{\beta}_i$  is the vector of fitted coefficients of the  $i$ -th model,  $N$  is the number of experiments that we perform,  $\mathbf{x}_n$  are the conditions of the  $n$ -th experiment, and  $\mathbf{y}_{n_i}$  is the observed response variable  $i$  that we profile in the  $n$ -th experiment (following the association of  $i = 1$ : M; 2: S; 3: A).

Then, we first test whether equations 1, 2, 3 are suitable models of the aforementioned relations. We do so with an F-test for each of the three multiple linear regressions. We define the hypotheses of the  $i$ -th F-test as follows, for  $i = 1$ : M; 2: *storage weight*; 3: *accuracy*.

$$\mathcal{H}_{0_i} : \beta_{P_i} = \beta_{AT_i} = \beta_{D_i} = 0 \quad (\text{Null Hypothesis})$$

$$\mathcal{H}_{1_i} : \exists \beta \in \{\beta_{P_i}, \beta_{AT_i}, \beta_{D_i}\} : \beta \neq 0 \quad (\text{Alternative Hypothesis})$$

In this sense, we first test whether there is a linear relation between each of the independent variables and the dependent ones. In the F-test, we compare the fitted model sum of squares against the residual sum of squares, and calculate the value of the F-statistic to derive the p-value associated to a level of confidence and to reject/support the null hypothesis [42].

Furthermore, we can evaluate how valuable the fitted models are with the adjusted  $R^2$  statistic [43]. In our multiple linear regression models, the  $R^2$  is the square of the multiple correlation coefficient, which measures how well a response variable can be predicted with a linear function of a set of independent variables.

We also formulate an hypothesis for each explanatory variable in each of the models of the independent variables ( $\mathcal{H}_{P_i}$ ,  $\mathcal{H}_{AT_i}$ ,  $\mathcal{H}_{D_i}$ ) that concerns the

statistical significance of each design decision in each of the specific models. Hence, we have three hypotheses for each model which sums to a total of another 9 hypotheses. With these, we aim to separately answer the research questions RQ1.1, RQ1.2 and RQ1.3. Each of these hypotheses can be tested with an adjusted t-test, which is an inferential statistic used to determine if there is a significant difference between the means of two groups [44]. Concretely, we aim to determine whether the coefficients encoding the contribution of each of the design decisions to the accuracy and complexity-related metrics are significantly different from zero. If the fitted coefficient corresponding to a design decision, e.g. *number of parameters* of the AI models, happens to be significantly different from zero and take a value  $y$  in a specific model  $i$ , we can state that if we fix all the other non-zero coefficients of model  $i$ , then for each change of 1 unit in the *number of parameters*, the response variable of model  $i$  changes  $y$  units.

$\mathcal{H}_{P_{0_i}} : \beta_{P_i} = 0$	(Null Hypothesis)
$\mathcal{H}_{P_{1_i}} : \beta_{P_i} \neq 0$	(Alternative Hypothesis)
$\mathcal{H}_{AT_{0_i}} : \beta_{AT_i} = 0$	(Null Hypothesis)
$\mathcal{H}_{AT_{1_i}} : \beta_{AT_i} \neq 0$	(Alternative Hypothesis)
$\mathcal{H}_{D_{0_i}} : \beta_{D_i} = 0$	(Null Hypothesis)
$\mathcal{H}_{D_{1_i}} : \beta_{D_i} \neq 0$	(Alternative Hypothesis)

These hypotheses test whether a specific independent variable (e.g. P) is significantly adding to the variability of the  $i$ -th variables of study given that the  $i$ -th model also considers other explanatory variables (e.g. *architecture type, dataset*).

With all this, we aim to (i) determine whether the accuracy and complexity of an AI-enabled application can be modelled by taking into account the design decisions; (ii) investigate the contribution of each design decision, analyse their interactions, and determine the most impactful design decisions on the overall performance of the applications.

For this reason, we explore a complete set of representative combinations of experimental settings. Concretely, we fix all the design decisions and iterate variations on a single one (and we do this for each of the design decisions) to consistently quantify the contribution of each of the design decisions.

### 3.4 Execution of the empirical study

In this section we explain how we have followed the empirical study design depicted above in Figure 1 and defined in our registered report [45]. Following the open science initiative of the EMSE journal [46], we share a replication package available at <https://doi.org/10.5281/zenodo.6523801>. The replication package enables to reproduce the experiments described in Section 3.4.2. Concretely, it consists of the source code for (i) developing the AI models (in

Python); (ii) the source code of the AI-enabled mobile applications (in C#); (iii) the source code for the statistical analysis (in R); and (iv) all the data needed for the aforementioned tasks. Deviations are also reported. With (i) we implement and train the AI models; with (ii) we develop the AI-enabled applications, build and deploy them to the mobile devices, and use the profilers to measure their performance. During operation, we also use the applications to query them with multiple inputs and measure their average accuracy and performance; and with (iii) we conduct the statistical analysis on the data generated during the repeated experiments in operation.

### 3.4.1 Data management: datasets

There exists a wide variety of NN models (e.g. recurrent, convolutional, Transformer-based) for tackling different problems. To provide generalist conclusions about AI-enabled applications in the vision and language domains, we aim to study a complete set of representative NN architectures. Also, we want to test the aforementioned architecture types in multiple benchmark datasets to provide a model and data-centric context. For this reason, we make use of the German Traffic Sign Recognition Benchmark (GTSRB) [47], MNIST [48] and Cifar-10 [49] datasets for solving image classification, and the Sentiment140 [50], Amazon Reviews [51] and IMDB Movie Reviews [52] for solving text classification.

For image classification, the GTSRB dataset is the output of an attempt to benchmark traffic sign recognition and it contains 51,840 images representing 43 unique traffic sign classes. MNIST is a benchmark dataset of handwritten digits (i.e. contains 10 unique classes belonging to the digits) containing 60,000 training examples, and a test set of 10,000 examples. Cifar-10 is a generic dataset for image classification containing 60,000 32x32 colour images in 10 classes, with 6,000 images per class.

For text classification, Sentiment140 contains 1,600,000 tweets annotated with either *positive* or *negative* labels, which encode the sentiment of the tweets. Then, the Amazon Reviews for Sentiment Analysis<sup>6</sup> contains 4 million reviews also tagged either as positive or negative, and the IMDB Movie reviews<sup>7</sup> dataset contains 50,000 reviews also for binary sentiment analysis.

### 3.4.2 Modelling

For image-based AI-enabled applications, we develop multiple Convolutional and ResNet-based neural networks, and for text-based ones we develop multiple LSTMs (Long Short-Term Memory), GRUs (Gated Recurrent Units) and Transformers. Then, we integrate them in Android mobile applications using Unity3D and Android Studio. For the statistical analysis of the accuracy and complexity relations we proceed as follows. First, we train the AI models and

---

<sup>6</sup> <https://www.kaggle.com/bittlingmayer/amazonreviews>

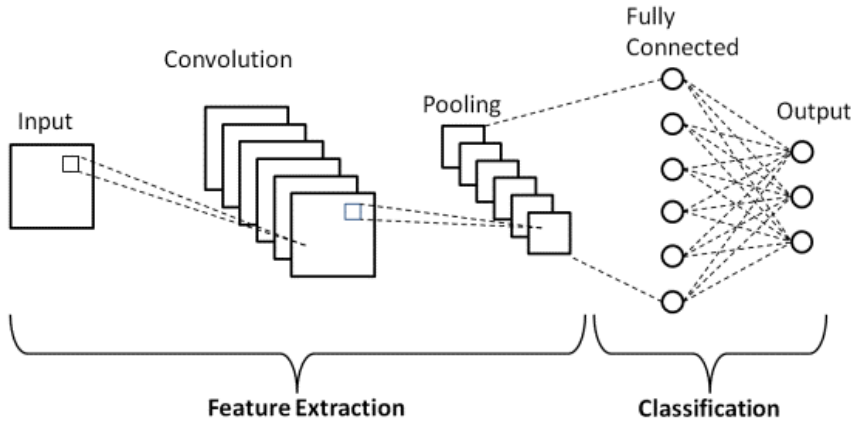
<sup>7</sup> <https://www.kaggle.com/lakshmi25npathi/imdb-dataset-of-50k-movie-reviews>

annotate their number of parameters (P), their architecture type (AT), their storage weight (S) (before being integrated in the mobile applications) and the dataset (D) used for training them. Secondly, we integrate the AI models in the mobile applications and profile the time of CPU usage (M) which we relate to the time consumption for answering a query and measure the accuracy in the AI-enabled system operation (A). Then, we export the profiled results and post-process them to aggregate all the operations that happen in the models (both in initialization and inference). Finally, we merge all the aforementioned information and craft an analysis dataset for the image and text domains (see the analysis datasets in the replication package provided).

Regarding the execution of our experiments, we perform some additional steps that we describe in the following. For image-based AI models, we add an extra categorical annotation to the Convolutional models (and not to the ResNet-based ones) indicating in which part of the model architecture the change in parameters (i.e. both addition or subtraction of weights and biases) happens. We find the need of making this distinction because we observe that parameters in the convolutional layers and in the fully connected layers behave differently, and hence can be studied separately. Still, note that the variable *number of parameters* indicates the total number of parameters and the extension just indicates which part of the model (convolutional, fully connected) contains the majority of the parameters. Hence we consider mainly convolutional CNNs (C-CNNs) and mainly fully-connected CNNs (FC-CNNs). For the ResNet-based models, we consider the five standard models which are the ResNet18, ResNet34, Resnet50, Resnet101 and Resnet 152. These have conceptually the same architecture and just vary in their complexity. We do not make a distinction between C and FC extensions for the ResNet-based models as they scale in a balanced way within these five models. In Figure 2 we show an instance of a CNN architecture. Whenever we relate to C-CNNs we mean that we increase the number of filters and layers in the *Feature Extraction* part, and when we relate to FC-CNNs we mean that we increase the number of neurons and layers in the *Classification* part.

For text-based AI models, we also encounter relevant differences in the behaviours of what we call the *embedding* and *network* parameters. The former are stored in the embedding layer, which is a look-up table, and are in charge of learning representations for the word tokens. There are  $tokens * embedding\ dimension$  embedding parameters in a text model, where *tokens* are the number of unique words in the dataset and *embedding dimension* is a hyperparameter to set. Then, the network parameters are in the neural network (e.g. in the recurrent and fully connected layers of the LSTM and in the attention and fully connected layers in the Transformers) and are in charge of learning the semantic relations between the words. Hence, we make the difference between EMB-Transformers, EMB-LSTMs, NN-Transformers and NN-LSTMs indicating which layers in the models dominates in numbers of parameters (in the same way we do for FC-CNNs and C-CNNs) and study them separately.





**Fig. 2** Example architecture of a CNN.

### 3.4.3 Development and Operation

As aforementioned in Section 3.4, during operation we query the AI-enabled mobile applications with evaluation sets containing 25 input examples each, which are all available in the replication package. We then measure the average accuracy and performance using the profilers and export the results for statistical analysis.

For our experiments, we have used a Xiaomi Poco X3 NFC mobile device to install and run the AI-enabled applications, which contains an Octa-core (2x2.3 GHz Kryo 470 Gold & 6x1.8 GHz Kryo 470 Silver) CPU.

During the experiments, we found several incompatibilities between the technologies that we use which block some paths defined in the scope of our study. First, we experience that the MobileNets architecture for images can be deployed in Android Studio but not in Unity due to an unsupported operation in the model [23]. Secondly, the GRU architecture cannot be imported neither in Android Studio nor in Unity, since it is a quite unpopular architecture compared to LSTMs and Transformers.

Finally, we have experienced that the profiler tools integrated in Android Studio and Unity work in a very different way and both make it quite difficult to export the profiled results into a file format that can be analysed in a simple way (i.e. a CSV). As we were incapable of doing so in Android Studio, we stick to an external package of Unity named *Profiler analyzer* which allows us to export the original outputs of the profiler (a binary file) into a CSV table that we can later analyze using the R programming language. Hence, the following results section is based on the models developed with Pytorch, exported with ONNX, and operated and profiled in the mobile applications built with Unity. The experimentation to carry out with this technology stack is available in the replication package provided.

## 4 Experiments' Results

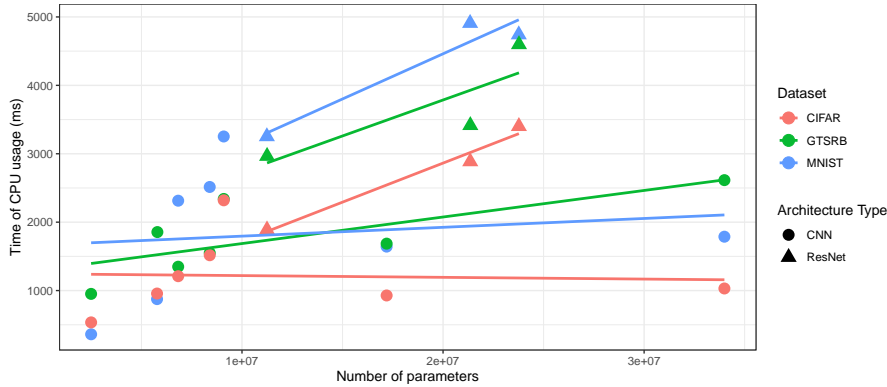
In this section, we answer our RQs by testing our hypotheses (see Section 3.3). In the following, we provide details on the experimentation and statistical analysis when operating with the AI-enabled mobile applications in each of the two domains of image and text. Concretely, we define a section for each RQ, and a subsection for each sub-RQ1 corresponding to the analysis of each dependent variable. In the sub-RQ1s we analyse the contribution of all the independent variables of study (i.e. *number of parameters*, *architecture type* and *dataset*) to each of the dependent ones (M, *storage weight* and A) for both domains of image and text. Note that we do this for all the datasets at the same time. The tables of results that gather all the evidence collected in the experiments are available in our public repository <sup>8</sup>.

4.1 RQ1.1: How is the time of CPU usage of the AI models affected by their number of parameters, architecture type and dataset used?

We first answer RQ1.1 which regards the relation between the the *time of CPU usage* and our independent variables of study. For that, we study the relation between the *number of parameters*, *architecture type*, the *dataset* and the *time of CPU usage* in the two domains of image and text.

### 4.1.1 Image-domain

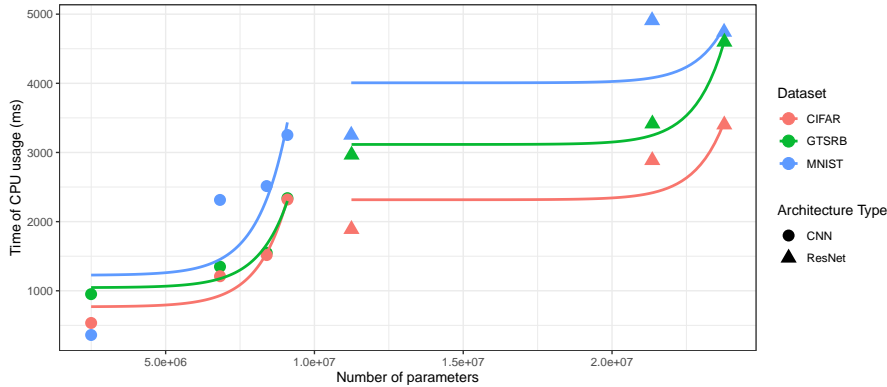
In Figure 3 we show the joint distribution between the *number of parameters* and the *time of CPU usage* for all different image datasets and architecture types.



**Fig. 3** Joint distribution of the dependent variables of study (*number of parameters*, *architecture type* and the *dataset*) and the *time of CPU usage*.

<sup>8</sup> <https://github.com/roger-creus/Which-Design-Decisions-in-AI-enabled-MobileApplications-Contribute-to-Greener-AI/tree/main>

Regarding the *architecture type*, Figure 3 shows that for all ResNets there is a positive linear relation indicating that the more parameters of the models are used, the higher time cost in operation. However, the distribution of the *time of CPU usage* shows a more complex shape for the CNNs (i.e. the CNN with higher time cost is not the one with more parameters). In the following we show the distribution of the *time of CPU usage* for C-CNNs and FC-CNNs separately to depict more precise and meaningful insights. Figure 4 shows the distribution of the *time of CPU usage* for C-CNNs and ResNets and Figure 5 shows the distribution of the *time of CPU usage* for FC-CNNs and ResNets.

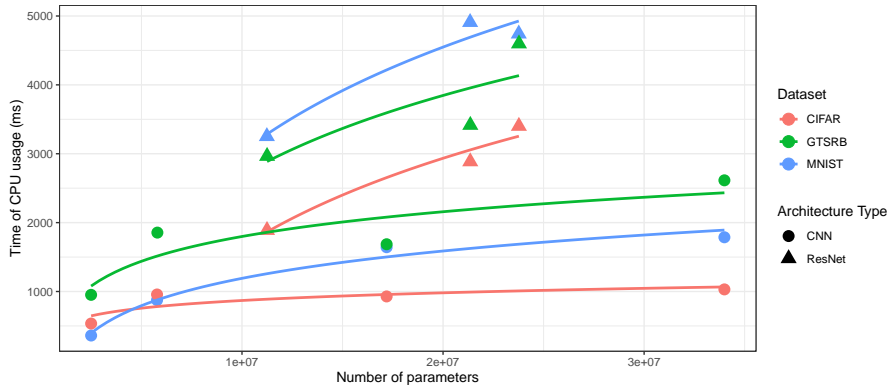


**Fig. 4** Joint distribution of the dependent variables of study (*number of parameters*, *architecture type* and the *dataset*) and the *time of CPU usage* of the C-CNNs.

In Figure 4 it can be seen that the joint distribution of convolutional parameters and the *time of CPU usage* is closer to a positive exponential form than to a linear one. This justifies that C-CNNs with a much smaller number of parameters than other FC-CNNs have a higher time cost in operation. This fact communicates that the convolution operation that happens in the first layers of all CNNs is significantly more time-consuming than the multiplications that happen in the FC layers.

Furthermore, in Figure 5 it can be seen that the joint distribution of fully-connected parameters and the *time of CPU usage* is more similar to a logarithmic distribution rather than a linear one. This fact might be due to the hardware components of the operation devices optimizing its resources as they deal with bigger models and hence with more complex operations (i.e. parallel computations).

So far we have identified positive relations between the *number of parameters* and the *time of CPU usage* for both ResNets and CNNs in all the image datasets, even though they have different forms (i.e. linear, exponential, logarithmic). However, we notice that the datasets used also have an impact in the *time of CPU usage* of the models, which might not be obvious if we assume that the *time of CPU usage* of a model only depends on its specification



**Fig. 5** Joint distribution of the dependent variables of study (*number of parameters*, *architecture type* and the *dataset*) and the *time of CPU usage* of the FC-CNNs.

(*number of parameters* and *architecture type*). Concretely, we identify that the differences in the *time of CPU usage* within the three different image datasets studied are due to the differences in the *number of parameters* that the *dataset* involves. Specifically, there are two dataset-specific factors that have an effect in the *number of parameters* and hence in the M:

- (i) The number of channels of the images in the dataset affects the size of the first convolutional input layer in the models. For example, as MNIST contains gray-scale images, these are represented as one-channel tensors and hence the first convolutional layer has a depth of 1. In contrast, as GTSRB and CIFAR contain RGB images they require a model with a first convolutional layer of depth 3.
- (ii) The number of classes in which we perform classification affects the number of output neurons in the last layer of the models and also the number of connections between the last two fully-connected layers. For instance, MNIST models contain 10 output neurons corresponding to each of the existing digits and GTSRB models contain 43 neurons belonging to the different traffic sign classes in the dataset. If we suppose that the second last fully connected layer contained 128 neurons, MNIST models would have  $128 * 10$  connections + 10 neurons = 1,290 parameters and GTSRB models would have  $128 * 43$  connections + 43 neurons = 5,547 parameters in the last two layers.

In the following, we fit Model 1 to describe the memory cost and test our hypotheses (see Section 3.3). Our hypotheses are based on testing what variables have a coefficient which fit is significantly different from zero.

In Table 2 it can be seen that the variable that makes the most significant difference in the *time of CPU usage* is the *architecture type*. According to the model fit, the CNNs used in this experiment imply a baseline time cost of 745ms and the ResNets of 2520ms. Also, each parameter adds an additional

**Table 2** Summary of the *time of CPU usage* Model fit in the image domain. The model reports  $R^2 = 91.54$  which defines a suitable goodness-of-fit.

Variable	Coefficient	Coefficient fit	P-value	Significance
P	$\beta_{P_1}$	2.79e-05	0.0761	.
AT-CNN	$\beta_{AT_{CNN_1}}$	7.45e+02	0.022	*
AT-ResNet	$\beta_{AT_{RN_1}}$	2.52e+03	2.62e-06	***
D-MNIST	$\beta_{D_{MNIST_1}}$	6.65e+02	0.051	.
D-GTSRB	$\beta_{D_{GTSRB_1}}$	8.98e+02	0.011	*

The model is described as:  $M = \beta_{P_1}P + \beta_{AT_1}AT + \beta_{D_1}D + \epsilon_1$

Note that  $\beta_{AT_1}$  and  $\beta_{D_1}$  have a separate coefficient for each *architecture type* and D respectively.

Significance levels are: No, Very little '.', Little '\*', Important '\*\*', Very important '\*\*\*'

2.79e-05ms. Additionally, experiments with GTSRB provide models that add a mean time cost of 665ms and in MNIST of 898ms.

For a more accurate estimate in the time cost that each parameter type adds, we fit the same Model 1 for C-CNNs and FC-CNNs separately, which provides the results presented in Tables 3 and 4.

**Table 3** Summary of the *time of CPU usage* model fit for ResNets and C-CNNs. The model reports  $R^2 = 96.54$  which defines a very suitable goodness-of-fit.

Variable	Coefficient	Coefficient fit	P-value	Significance
P	$\beta_{P_1}$	1.52e-04	5.84e-05	***
AT-CNN	$\beta_{AT_{CNN_1}}$	1.43e+02	0.63	No
AT-ResNet	$\beta_{AT_{RN_1}}$	1.84e+02	0.75	No
D-MNIST	$\beta_{D_{MNIST_1}}$	1.08e+03	0.01	.
D-GTSRB	$\beta_{D_{GTSRB_1}}$	4.87e+02	0.001	**

Significance levels are: No, Very little '.', Little '\*', Important '\*\*', Very important '\*\*\*'

**Table 4** Summary of the *time of CPU usage* model fit for ResNets and FC-CNNs. The model reports  $R^2 = 95.31$  which defines a very suitable goodness-of-fit.

Variable	Coefficient	Coefficient fit	P-value	Significance
P	$\beta_{P_1}$	4.24e-05	0.003	**
AT-CNN	$\beta_{AT_{CNN_1}}$	4.73e+01	0.878	No
AT-ResNet	$\beta_{AT_{RN_1}}$	2.17e+03	1.25e-05	***
D-MNIST	$\beta_{D_{MNIST_1}}$	8.49e+02	0.008	.
D-GTSRB	$\beta_{D_{GTSRB_1}}$	9.22e+02	0.013	**

Significance levels are: No, Very little '.', Little '\*', Important '\*\*', Very important '\*\*\*'

In Tables 3 and 4 it can be seen that the estimated time costs for convolutional and fully-connected parameters are 1.52e-04 and 4.24e-05 respectively, indicating that a convolutional parameter implies a time cost approximately

3.5 times bigger than a fully-connected parameter. Finally, we summarize our findings in the distribution of the *time of CPU usage* in image-based AI-enabled systems:

**Finding 1:** There exists a positive relation between the *number of parameters* and the *time of CPU usage*, meaning that an increase in the number of parameters of a CNN or a ResNet increases the time cost when operating with it.

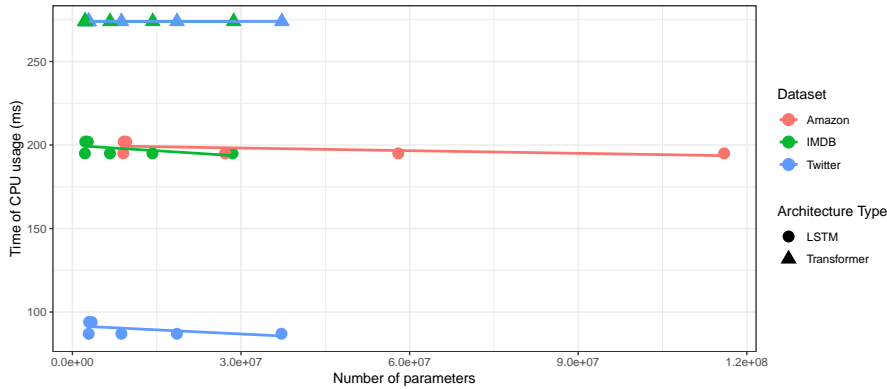
**Finding 2:** The *time of CPU usage* increases close to exponentially when there is an increase in the number of convolutional parameters of a CNN.

**Finding 3:** The *time of CPU usage* increases close to logarithmically when there is an increase in the number of fully-connected parameters of a CNN.

**Finding 4:** The image depth and number of classes in a dataset has an effect in the *number of parameters* of the AI models that fit it, and hence in their *time of CPU usage*.

#### 4.1.2 Text-domain

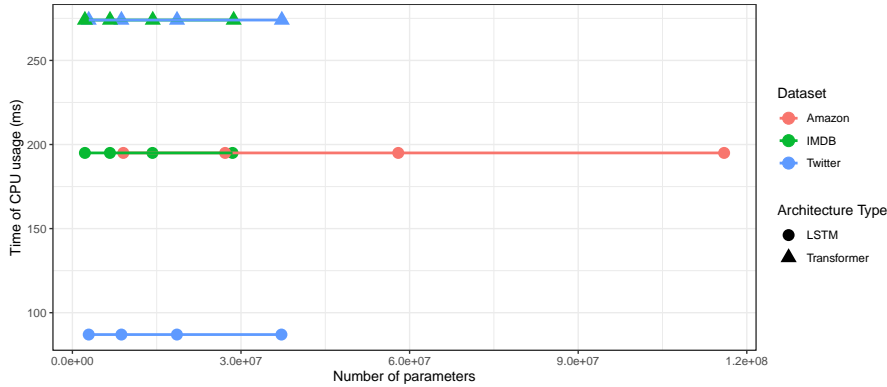
In Figure 6 we show the joint distribution between the *number of parameters* and the *time of CPU usage* for all different text datasets and architecture types.



**Fig. 6** Joint distribution of the *number of parameters*, *architecture type*, *dataset* and the *time of CPU usage*. The figure does not show evidence of a positive nor negative relation between the *time of CPU usage* and the *number of parameters*.

Figure 6 shows no evidence that increasing the number of parameters of a text model increases its time cost. However, it can be seen that Transformers have a higher time cost than LSTMs, and also that LSTMs in the Twitter dataset have a smaller *time of CPU usage* than in the other two datasets. Note that the models have a different number of parameters depending on the

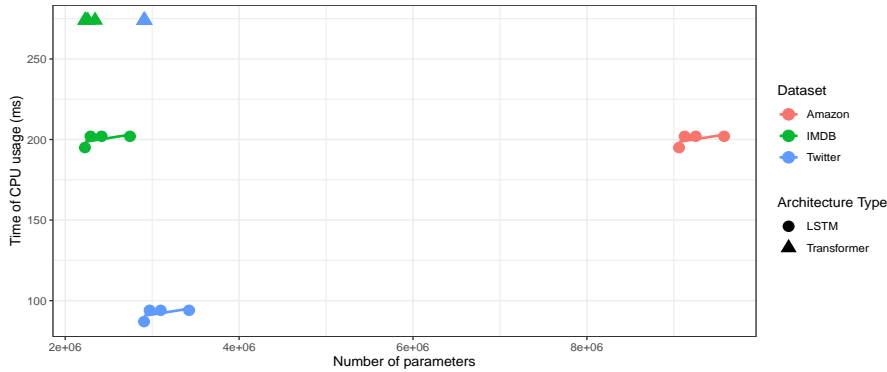
dataset that they are trained on. This is due to the size of the embedding layer depending on the number of the unique words in each dataset, being 937,500 in the Amazon dataset, 234,375 in the IMDB dataset and 225,452 in the Sentiment140 dataset. Figure 7 shows the distribution of the *time of CPU usage* for EMB-Transformers and EMB-LSTMs and Figure 8 shows the distribution of the *time of CPU usage* for NN-Transformers and NN-LSTMs.



**Fig. 7** Joint distribution of the *number of parameters*, *architecture type*, *dataset* and the *time of CPU usage* of the EMB-extended models.

In Figure 7 it can be seen that no matter what size the embedding layer has, the models need the same time to answer any query. As mentioned above, increasing the embedding parameters means increasing the size of a look-up table, and the fact that the time cost does not change is thanks to an efficient search algorithm on the look-up table which time complexity is independent of the size of the table. Recall that the text models are a combination of the embedding layer (look-up table) and fully-connected layers (as is explained in Section 3.4.2). In this way, a fully-connected layer consists of doing matrix multiplication, and hence the bigger the size of the matrix the higher the time cost. However, for the embedding layer, no complex data structure operation is required but only a search on the look-up table. For this reason, the fact that increasing the *number of parameters* does not increase the *time of CPU usage* is thanks to the software used (i.e. Pytorch, Tensorflow) internally providing a search algorithm on the embedding layer that scales greatly with the size of the look-up table.

In Figure 8 it can be seen that the increase in the number of network parameters in the text models causes what looks like a logarithmic increase in the time cost. This is very similar to what we also observe in the FC-CNNs, and it makes sense because all FC-CNNs, NN-Transformers and NN-LSTMs have fully connected layers that work in the same way. In Figure 8 it is difficult to see the latter for NN-Transformers and that is because we could not experiment with a wide range of number of network parameters since increasing them



**Fig. 8** Joint distribution of the *number of parameters*, *architecture type*, *dataset* and the *time of CPU usage* of the NN-extended models.

increases their memory cost drastically and we could not train them using our available computing environments (nor the one offered for free in Kaggle). That is, the memory cost of training such complex models is not accessible for researchers using standard computational environments but usually only available to big companies [14].

In the following, we fit the linear model 1 to describe the *time of CPU usage* in the text domain and test our hypotheses. Table 5 reports that increasing the number of parameters (both network and embedding parameters) does not have a significant effect over the *time of CPU usage*. Also, it shows that LSTMs carry a baseline time cost of 117ms and Transformers of 247ms. Then, compared to the Sentiment140 dataset, IMDB dataset carries a baseline time cost of 53.8ms and the Amazon dataset of 82.4ms (e.g. fixing the value of the *number of parameters*, a Transformer model in the amazon model has a time cost of  $247\text{ms} + 82.4\text{ms} = 329.4\text{ms}$  for each query).

**Table 5** Summary of the model fit for the *time of CPU usage* in the text domain. The model reports  $R^2 = 98.55$  which defines a very suitable goodness-of-fit.

Variable	Coefficient	Coefficient fit	P-value	Significance
P	$\beta_{P_1}$	-6.28e-08	0.79	No
AT-LSTM	$\beta_{AT_{LSTM}_1}$	1.17e+02	5.39e-14	***
AT-Transformer	$\beta_{AT_{Transformer}_1}$	2.47e+02	2e-16	***
D-IMDB	$\beta_{D_{IMDB}_1}$	5.38e+01	7.01e-06	***
D-Amazon	$\beta_{D_{Amazon}_1}$	8.24e+01	2.23e-06	***

Significance levels are: No, Very little '.', Little '\*', Important '\*\*', Very important '\*\*\*'



Finally, we summarize our findings in the distribution of the *time of CPU usage* in text-based AI-enabled systems:

**Finding 5:** Contrary to what happens in the image domain, there is no evidence that there is a significant relation between the *number of parameters* and the *time of CPU usage*. Even though we conclude that the latter holds for embedding parameters, we believe that a more extensive analysis over the network parameters could provide evidence that increasing the *number of parameters* also increases the memory and time cost.

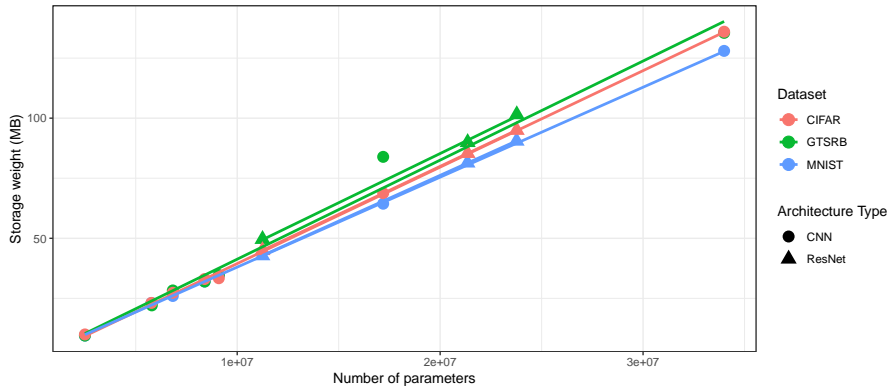
**Finding 6:** The number of unique words in a text dataset drastically conditions the number of parameters of the text models since it defines the size of the embedding layer. However, this does not have an impact on the *time of CPU usage*.

**Finding 7:** Transformers have a baseline time cost that is higher than the one in LSTMs.

4.2 RQ1.2: How is the storage weight of the AI models affected by their number of parameters, architecture type and dataset used?

#### 4.2.1 Image domain

In the following we study the relation between the dependent variables of study (*number of parameters*, *architecture type*, *dataset*) and the storage weight (S). In Figure 9 we show the joint distribution between the *number of parameters* and *storage weight* for both CNNs and ResNets models and all the image datasets.



**Fig. 9** Joint distribution of the *number of parameters*, *architecture type*, *dataset* and the *storage weight*. The figure clearly shows a positive linear relationship between the *number of parameters* and *storage weight* for all architectures and datasets.

In Figure 9 it can be seen that there is a very clear and positive linear relationship between the *number of parameters* and *storage weight* for all architectures and datasets of study. In this case, the *storage weight* distribution does not make a distinction between C-CNNs and FC-CNNs either. When fitting the distribution of *storage weight* with respect to the three variables of the *number of parameters*, *architecture type* and the *dataset* described in model 2, we obtain the results shown in Table 6. The linear fit reports an  $R^2$  value of 0.99 which indicates that with the selected independent variables we are capturing all the variance in the target variable. Furthermore, we can see that only the *number of parameters* is significantly causing variations in the distribution of *storage weight*. Hence, we can provide a valuable estimated coefficient for the *number of parameters* of value  $3.96e-06$ . With this, we can make predictions on the storage weigh (S) of a model by only knowing its number of parameters: e.g. a CNN (or a ResNet) with 3,545,342 parameters will approximately weigh  $3,545,342 * 3.96e-06 = 14.04\text{MB}$ .

**Table 6** Results of the linear fit between the *number of parameters* and the *architecture type*, and *storage weight*.

Variable	Coefficient	Coefficient fit	P-value	Significance
P	$\beta_{P_1}$	3.96e-06	<2-16	***
AT-ResNet	$\beta_{AT_{RN_2}}$	1.06e+00	0.5973	No
AT-CNN	$\beta_{AT_{CNN_2}}$	-1.97e-01	0.8929	No
D-GTSRB	$\beta_{D_{GTSRB_2}}$	2.99e+00	0.1646	No
D-MNIST	$\beta_{D_{MNIST_2}}$	-2.46e+00	0.1250	No

The model is described as:  $S = \beta_{P_2}P + \beta_{AT_2}AT + \beta_{D_2}D + \epsilon_2$

Note that  $\beta_{AT_2}$  and  $\beta_{D_2}$  have a separate coefficient for each *architecture type* and *dataset* respectively.

Significance levels are: No, Very little '.', Little '\*', Important '\*\*', Very important '\*\*\*'

We provide the following findings:

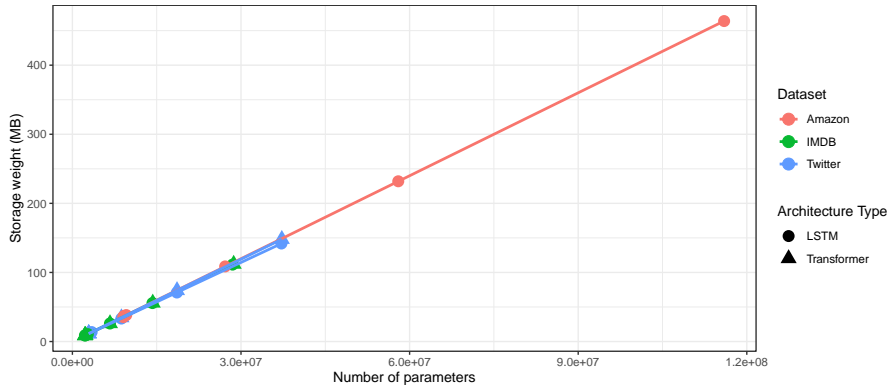
**Finding 8:** There exists a linear relationship between *storage weight* and the *number of parameters* described as  $S = 3.96e-06 * P$ .

**Finding 9:** No other variables significantly cause variations in the distribution of *storage weight*.

#### 4.2.2 Text-domain

In the following we study the relation between the dependent variables of the study (*number of parameters*, *architecture type*, *dataset*) and *storage weight* for text models. In Figure 10 we show the joint distribution between the *number of parameters* and *storage weight* for both Transformers and LSTMs and all the text datasets.

As it happened with the image models, there exists a very clear and positive linear relation between the *number of parameters* and *storage weight*, and it



**Fig. 10** Joint distribution of the *number of parameters*, *architecture type*, *dataset* and *storage weight*. The figure clearly shows a positive linear relationship between the *number of parameters* and *storage weight* for all architectures and datasets.

can be seen for all text models and datasets. When fitting the linear model 2 to the distribution of *storage weight* over the *number of parameters* we obtain the results in Table 7.

**Table 7** Results of the linear fit between the *number of parameters* and *storage weight*.

Variable	Coefficient	Coefficient fit	P-value	Significance
P	$\beta_{P_2}$	3.97e-06	2e-16	***
AT-LSTM	$\beta_{ATLSTM_2}$	2.37e+00	0.7265	No
AT-Transformer	$\beta_{ATTransformer_2}$	1.91+01	0.5918	No
D-IMDB	$\beta_{DIMDB_2}$	-6.14e-01	0.2739	No
D-Amazon	$\beta_{DAmazon_2}$	3.80e+00	0.3561	No

Significance levels are: No, Very little '.', Little '\*', Important '\*\*', Very important '\*\*\*'

Again, we obtain a R2 of 0.99 indicating a quasi-perfect goodness-of-fit. The linear model reports that each parameters (be a network or embedding parameter) adds an additional 3.97e-06MB to the model’s weight, which is a very similar result than the one obtained for image models, which was 3.96e-06MB.

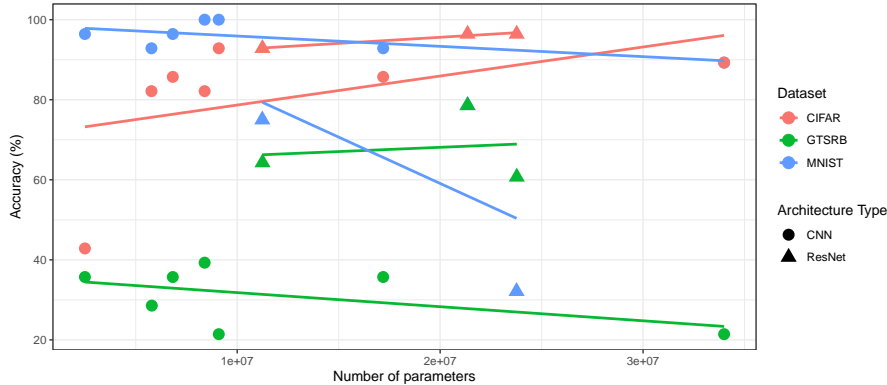
In this way, we report an additional finding relating both domains.

**Finding 10:** In the text domain there also exists a clear and linear relationship between *storage weight* and the *number of parameters* described as  $S = 3.97e-06 * P$ .

4.3 RQ1.3: How is the accuracy of the AI models affected by their number of parameters, architecture type and dataset used?

#### 4.3.1 Image-domain

Finally, we study the joint distribution of the *accuracy* and the different values of the *number of parameters*, *architecture type* and the *dataset*. The overall distribution for the image-domain is shown in Figure 11.



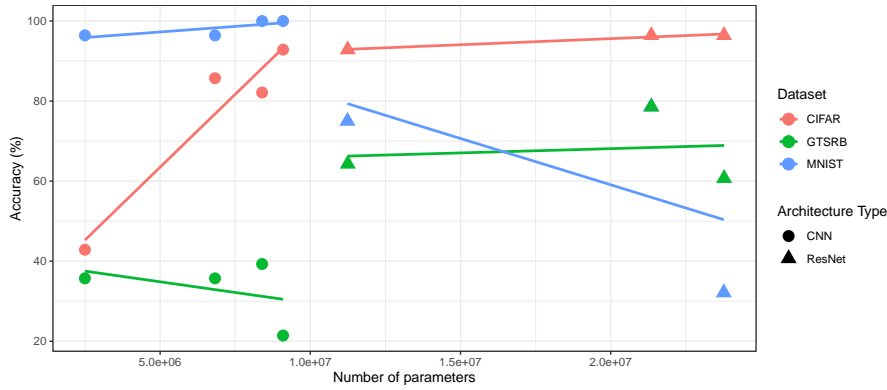
**Fig. 11** Joint distribution of the *number of parameters*, *architecture type*, *dataset* and *accuracy*. The figure shows a cloud of points which indicates that the distribution of *accuracy* might be more complex than linear with respect to the independent variables.

Inspecting Figure 11 we observe that increasing the *number of parameters* has a positive effect in the *accuracy* in CIFAR but a negative effect in GTSRB and MNIST. To obtain more granular insights about the relation between the *number of parameters*, *architecture type*, *dataset* and *accuracy* we show the latter separately C-CNNs and FC-CNNs. In the following, Figures 12 and 13 show the *accuracy* achieved by C-CNNs and FC-CNNs.

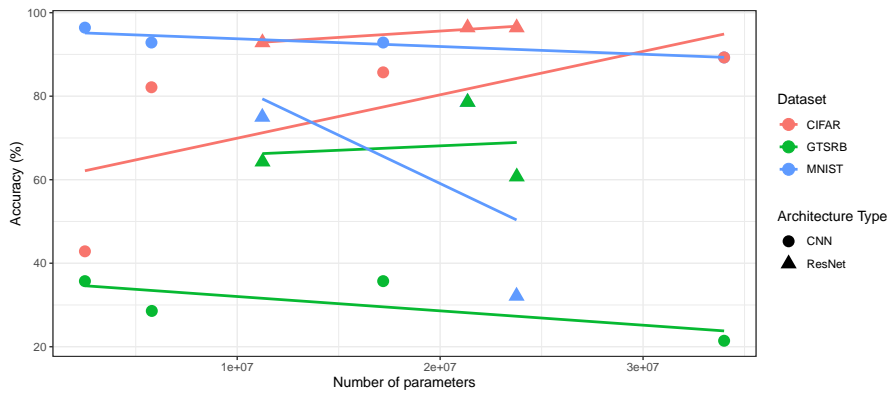
A clearly delineated insight emerges from Figures 12 and 13: the dataset has a significant effect over the accuracy that the AI models achieve. Concretely, we identify MNIST as the easiest dataset for achieving the highest accuracy when using CNNs. However, ResNets appear to be a too complex architecture type for fitting the MNIST dataset.

In addition, there is no common trend in terms of positive or negative relations between *accuracy* and the *number of parameters*, as the *dataset* seems to have a very significant effect over the distribution of *accuracy*. Specifically, the following results can be retrieved from the visualizations:

- (i) In MNIST, it seems that increasing the convolutional parameters of the CNNs helps the models extract better features from the images and hence obtain higher accuracy, but increasing the fully-connected parameters adds too much unnecessary complexity to the models and lowers their accuracy.



**Fig. 12** Joint distribution between the *number of parameters*, *dataset* and *accuracy* for C-CNNs and in all the image datasets.



**Fig. 13** Joint distribution between the *number of parameters*, *dataset* and *accuracy* for FC-CNNs and in all the image datasets.

Also, ResNets seem too complex models for MNIST as they obtain lower accuracy than basic CNNs.

- (ii) In GTSRB, simpler is better, as the smallest models are the ones performing better. This can be possible if the distribution of the images can be easily separable among the 43 existing classes in GTSRB and hence the more complex models look for too sophisticated patterns for classifying the data. Also, ResNets perform much better than CNNs in operation because the data in GTSRB is very different from the one used in operation for measuring the accuracy (i.e. GTSRB contains traffic signs from Germany and operation data consists of traffic signs from Spain), and ResNets offer better generalization than CNNs.
- (iii) CIFAR seems to be a too complex dataset for being overfitted and so the more complex the models, the better accuracy they achieve (i.e.

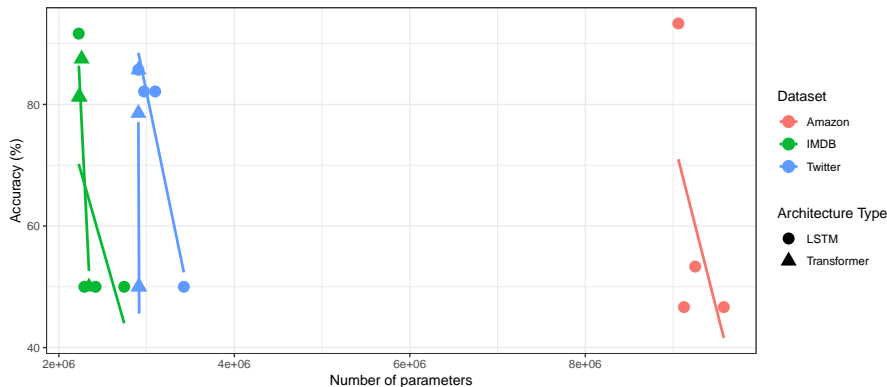
the bigger the *number of parameters*, the higher *accuracy*, and ResNets perform better than CNNs).

**Finding 11:** The *dataset* has a significant impact over *accuracy*. Increasing the *number of parameters* increases *accuracy* only when there is enough data in the *dataset*. If that is not the case, the distribution of *accuracy* is non-linear and one needs to search for the sweet spot of the *number of parameters* and the *architecture type*.

**Finding 12:** Similar to the *dataset*, the *architecture type* also has a significant effect over *accuracy*. Using more complex architectures (models that integrate more sophisticated pattern recognition operations, e.g. ResNets being more complex than CNNs) only has a positive effect over *accuracy* when fitting large amounts of data that need such adjustments to be fit. In the other cases, using simpler models like CNNs works best, so it is always recommended to start with the simpler fits and add complexity incrementally only when needed.

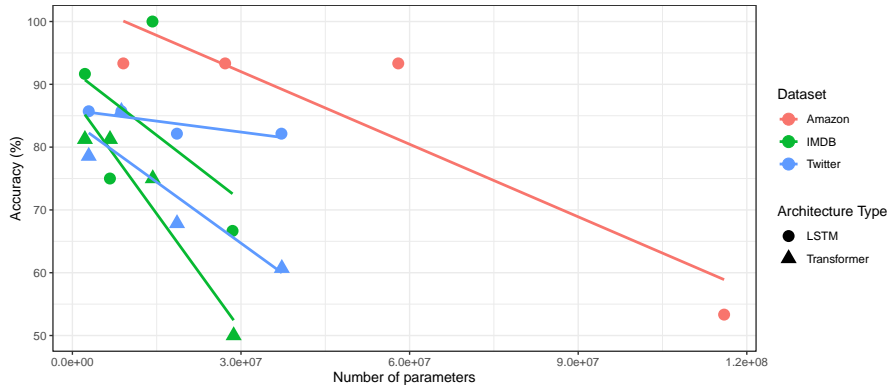
#### 4.3.2 Text-domain

In the following we show the distribution of *accuracy* for text models, considering the EMB and NN models separately.



**Fig. 14** Joint distribution between the *number of parameters*, *dataset*, ATs and *accuracy* for NN-models and all text datasets.

In Figure 14, we show the distribution of *accuracy* for NN-Transformers and NN-LSTMs and in Figure 15 for EMB-Transformers and EMB-LSTMs. In both cases, we observe that for the selected datasets, increasing the *number of parameters* (both network or embedding parameters) has a negative effect over *accuracy*. The latter means that the simpler models are capable of learning the best language representations with the smallest embedding spaces and the way to relate and operate with them with the simplest non-linear functions



**Fig. 15** Joint distribution between the *number of parameters*, *dataset*, *architecture type* and *accuracy* for EMB-models and all text datasets.

(learnt by the fully-connected layers). In general, no dataset seems easier to solve than the others because the models achieve a similar accuracy on all of them. Also, generally the LSTMs achieve a higher accuracy than Transformers in the datasets of study.

**Finding 13:** For the three text datasets of study, the simpler models always achieve the best results. Hence it is recommended to start modelling the data in a small embedding space and with a small number of parameters.

## 5 Discussion

In our work we have identified and depicted an underlying relationship between some design decisions of a set of AI models with respect to the overall performance of the mobile applications that integrate them. In this section, we discuss (i) the subjective experience on the validity of profiling tools to analyze the accuracy-complexity trade-off of AI-enabled mobile applications and (ii) the implications of our results on both the AI engineering research and industry.

### 5.1 Is profiling a suitable tool to systematically analyze the trade-off between accuracy and complexity of AI-enabled applications?

In the following, we report our subjective empirical experience when collecting data using profilers to help determine whether profiling tools can become a standard approach for systematically reasoning about the performance of AI-enabled applications. In particular, we focus on the Android and Unity profilers, which are the ones used in this work.

According to the Android Developers documentation, the Android Profiler tools provide real-time data to help one to understand how an app uses CPU, memory, network, and battery resources<sup>9</sup>. Furthermore, the Android Profiler allows to create sessions, which allows to record data of different runs and compare it. A drawback of the Android Profiler is that it is very sophisticated and that makes it quite difficult to use. It offers a lot of capabilities but are hard to configure. The good thing is that there are lots of official documentation which help to make the process easier. An example of usage of this tool consists of running an application, selecting which profiler module to view (CPU, memory, network, battery), selecting a time frame in the event timeline, and tracing the system calls. Regarding the latter, one can wrap a piece of code to be generated as a single system call in the profiler so the resources used are aggregated for all the operations that are wrapped (e.g. the *predict* wrapper can aggregate all the convolutional and fully connected layer operations happening during a forward-pass of a CNN). All this capabilities are suitable for the developer experience on Android Studio. However, for this work we needed to statistically analyze the results with regression models and for that we needed an available dataset or data table that we can export and iterate externally. Although we have persistently tried to find a way to export the Android Profiler results to a readable data format (e.g. a CSV file) we have found no support for that, and hence we have not been able to do anything with the profiled results but seeing them displayed in our screens. The only supported export operation saves the profiled results in a *.trace* file, which cannot be parsed into a format that can fit an statistical analysis with a software like R or Python. Figure 16 shows the Android Profiler when computing an image class prediction with a CNN.

Regarding the Unity Profiler, the official documentation defines it as a tool that gathers and displays data on the performance of an application in areas such as the CPU, memory, renderer, and audio. It's a useful tool to identify areas for performance improvement in the application, and iterate on those areas<sup>10</sup>. As Unity is a graphical engine, the integrated Profiler provides modules like the illumination module, the physics, the UI, or the rendering ones. Then, just like the Android Profiler, it also provides the CPU and the memory modules, plus a GPU profiler. However, it lacks the network and battery modules.

Compared to the Android Profiler, the one integrated in Unity is much easier to use, as one opens it and already sees the event timeline being updated in real-time. The user can click on a "Record" button to start tracing all system calls of the selected module. Then, it can stop recording and see the results in a sequential or hierarchical view.

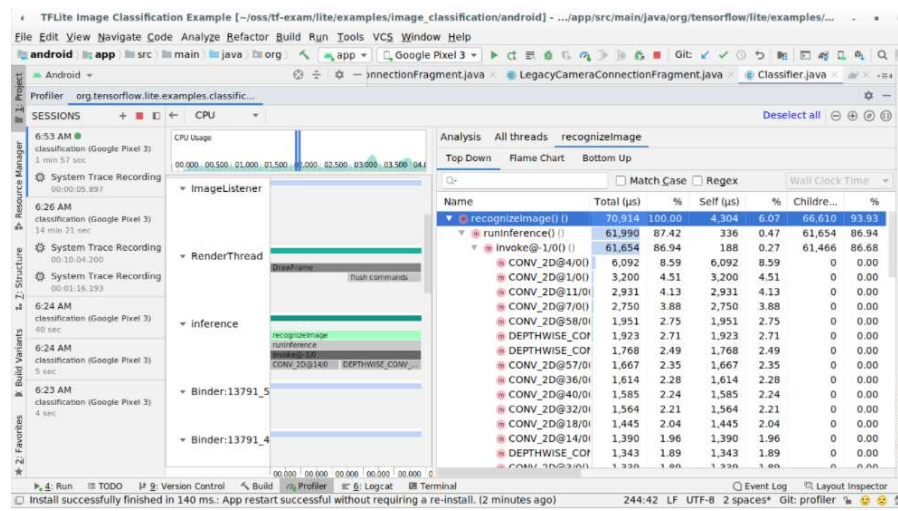
A great thing of the Unity profiler is that it can connect to an application run in an external device like a mobile phone. In this way, the profiler in the computer reports the resource usage of the external device, which is what we

---

<sup>9</sup> <https://developer.android.com/studio/profile/android-profiler>

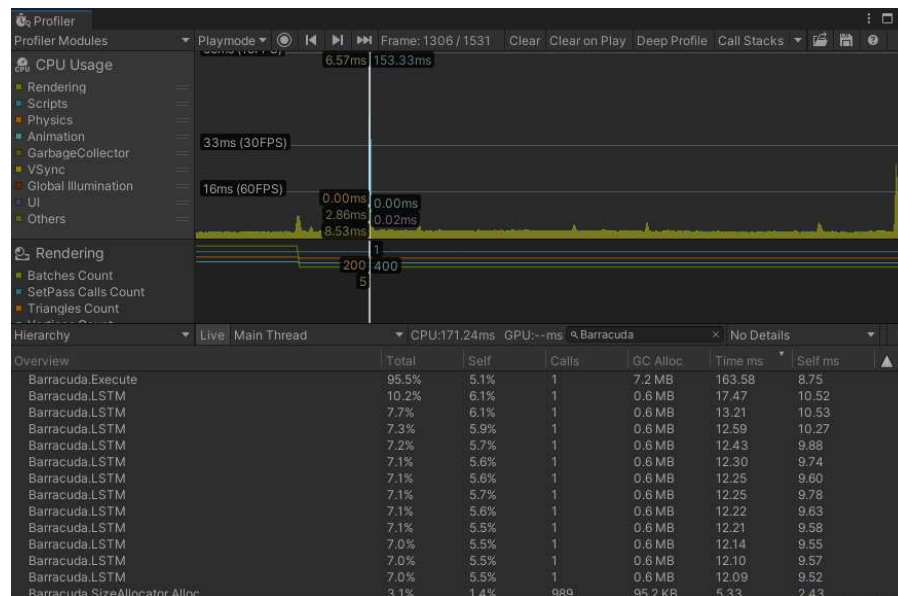
<sup>10</sup> <https://docs.unity3d.com/Manual/Profiler.html>





**Fig. 16** The Android Studio Profiler. It can be seen that each layer of the CNN is profiled separately, which is useful to get a better idea of which layers are the ones carrying more computations and resource usage.

needed in this work, where we study the importance of the resource management in mobile devices. Figure 17 shows the Unity3D Profiler when computing a text sentiment prediction with a LSTM.



**Fig. 17** The Unity3D Profiler. As in the Android Studio Profiler, each of the different LSTM layers is profiled separately.

Finally, just like with the Android profiler, there is no integrated support in the Unity profiler to export the profiled data into an externally readable format. However, there is a tool in Unity named the *Package Manager* which allows to import community and official packages to Unity, and we found a package named *Profiler Analyzer* which allows the user to export the profiled data into a CSV file, which is what we needed for this work.

For the availability to profile the performance of the mobile applications and to export the results into CSV files we mainly use Unity for analyzing the results of our experiments. However, we consider the Android profiler a very suitable tool for developers since it provides high-level insights of the performance of the applications.

Profilers were of critical relevance in this work, because without these tools we could not have any data at all to analyze. Hence, in this work, we argue that profiling tools are valid and useful to extract meaningful knowledge regarding the performance of AI-enabled applications. We consider that the two studied profilers are supported by rich documentation and are reasonably simple to use for practitioners, except for doing more complicated operations like exporting the results to specific formats. Finally, systematically comparing the capabilities of each profiler is out of the scope of this work, as it is to perform a more wide study using many other available and recent profiling tools (e.g. CodeCarbon <sup>11</sup>, Nvidia-smi <sup>12</sup>).

## 5.2 Implications for AI engineering research

Our work provides the AI engineering research community with an end-to-end experience on how to develop AI-based mobile applications, monitor their performance and interpret the profiling results. At the same time, we provide researchers with knowledge and awareness of the impact that the design decisions of AI models can have in the mobile applications' performance. Similar to other works like [37], we use Pytorch and Tensorflow as the DL frameworks for building a set of AI models with different specifications (in both their work and in ours we experiment with Transformers, CNNs and ResNets), and we also focus on monitoring the resource consumption of the AI models and describing statistically significant hypotheses. Compared to works like [38], where the path to achieve greener AI is studied from a data-centric view, we work in a model-centric approach in which we focus on design decisions on the AI models (e.g. *number of parameters* and *architecture type*). Nevertheless, we also include the *dataset* as a design decision that takes into account a data-based variable in the study of green AI. Compared to the aforementioned works, we target mobile devices, which are environments with low computational resources, and we formulate linear models that can uncover significant relations between the models' configuration and the mobile applications' performance.

---

<sup>11</sup> <https://codecarbon.io>

<sup>12</sup> <https://developer.nvidia.com/nvidia-system-management-interface>

Regarding the AI models that we experiment with (i.e. CNNs, ResNets, LSTMs, Transformers) we believe that the architecture and operations that these use are so different that it makes the *architecture type* to play a key role in the results we obtain. In this sense, we consider that the impact of some design decisions (e.g. *number of parameters* or *dataset*) can be significantly different depending only on the *architecture type*.

Comparing Android Studio and Unity, which are the frameworks that support the development of AI-enabled mobile applications that we use in our work, we consider that Unity comes with a more user-friendly experience for developers. From feeding data to the AI models to profiling the applications' performance, Unity provides the same capabilities as Android Studio but requiring less coding of technical components.

Finally, we provide the AI engineering research community with motivation to extend the degree in which this underlying relation can be characterised. One such example of the latter is the work by *Verdecchia et al.* on studying what variables influence the performance of AI-enabled software from a data-centric approach. Other paths that we motivate are extending the set of design decisions (identifying new ones), including more metrics to study the impact to the applications' performance, using other frameworks that support the implementation of AI models in applications or working in edge devices other than mobile devices.

### 5.3 Implications for AI engineering industry

Our work provides the AI engineering industry with a detailed analysis of the impact of only working towards maximizing the accuracy of AI models without taking their complexity into account. In our work, we have shown that the more complex models are responsible of higher usage of resources and this correlates with more energy consumption. The trend of pursuing maximum accuracy models disregarding any other cost is more popular in the industry sector because companies want to take the maximum profit out of their AI models, and they usually have enough resources for doing so. However, with green AI becoming more and more popular and the need to reduce carbon footprint more important, the AI engineering industry has the urge to switch to greener AI-based solutions. We are living in a moment where the planet requires greener AI and the customers are gaining consciousness of the relevance of a company's compromise to work towards greener AI [53].

For these reasons, with our work we validate that a very suitable and feasible methodology for implementing and deploying AI models is to first make use of the simpler architectures (e.g. CNNs for image and LSTMs for text), and with a small number of parameters. From the results obtained, we have experienced that high accuracy results can be often obtained with models with such simple configurations, depending on the dataset used (which affects the quality and quantity of the data) and the problem to solve. For this reason, we argue that a suitable methodology for developing AI-enabled software is

to start with the simplest solutions and add complexity only when needed. A valuable lesson learned that generalizes across the vision and language domains is that the distribution of the accuracy obtained with different models is difficult to frame with linear relationships, and hence we suggest tuning and evaluating the configurations of the AI models meticulously before considering more complex solutions.

In particular, we suggest following the practices adopted in this paper for implementing the AI models:

- Have a finite set of models with substantially different scales to quickly gain insights on how the number of parameters is affecting the usage time of CPU and helping achieve higher accuracy (learned from Findings 1,2,3).
- Keep the set of possible AI model scales small and have one for each specific architecture type of interest.
- Starting with the biggest AI models, keep discarding the ones that do not significantly achieve higher accuracy. These will be the least valuable in terms of the accuracy-performance trade-off (learned from Findings 11,12,13).
- Predict ahead the storage requirements of the model (following the formula uncovered in this manuscript) and see if it satisfies the application needs (learned from Findings 8,10).
- Test its time of CPU usage and real-time performance in operation and see if it satisfies the application needs.

## 6 Threats to validity

As with any empirical study, there might be limitations to our study design. Next, we report a few mitigation actions taken during the design of this experiment.

*Construct validity:* extent to which the theory constructs are correctly operated in the experiment. We build convolutional (i.e. baseline, ResNets), recurrent (i.e. LSTMs) and Transformer-based NNs in order to mitigate mono-operation bias.

*Conclusion validity:* ability to draw correct conclusion between treatment and outcome, we define and fit linear models to approximate the relation between design decisions and accuracy and complexity-related metrics. Hence, we estimate the models in order to minimize the error when fitting the collected experimental data with linear relations. In this way, the results of our linear models report that these are able to capture almost all the variability of the performance metrics with respect to linear relations between the design decisions. With that and together with the statistical tests on the significance of each of the design decisions we validate our conclusions. Regarding the validity of the AI-enabled applications themselves, their quality depends on the experience of the developers involved since that can influence their complexity. For mitigating the latter we provide as simple implementations as possible

by only implementing the required functionalities to the application and by following official tutorials for the development of these.

*Internal validity:* extent to which statements can be made about the causal effects of treatments on outcome. A wide variety of NNs are developed and deployed using different technology stacks, mitigating that the results of our analysis of the overall performance of the applications are caused accidentally rather than by the variables of study themselves. Also, the error in the causal interpretation of our conclusions is tied to the estimation capabilities of our analysis models, which we report to be high.

*External validity.* Our results are tied to the datasets and technology used for experimentation, that is, on training the AI models with Pytorch and Tensorflow, exporting them with ONNX and TensorflowLite, and developing the mobile applications in Android Studio and Unity3D. For the datasets, the results are tied to the image datasets of GTSRB, MNIST and CIFAR, and the text datasets of Sentiment140, IMDB Movie Reviews and Amazon Reviews. All the aforementioned determines the experimental settings to which the results obtained that regard the accuracy and complexity could be generalized.

Furthermore, we would like to remark that in this experimental work, we have followed the standards of empirical studies<sup>13</sup>, where at the core, we manipulate independent variables, study dependent variables and apply each treatment independently to several experimental units. Our work contains the standard essential attributes of empirical studies, for instance: (i) we state formal hypotheses; (ii) we describe the dependent variables and justify how they are measured; (iii) we describe the independent variables and how they are manipulated or measured; (iv) we design an appropriate protocol for studying the stated research questions and hypotheses; and, (v) we discuss the interpretations of the results. Finally, our work also follows a set of other desirable attributes: (vi) we provide a replication package containing all the source code, analysis scripts and data used and collected; (vii) we show visualizations of the observed data distributions; and, (viii) we average our results across data of different sources and modalities.

## 7 Conclusions and future work

Our hypotheses have always been oriented to characterise an existing underlying relationship between the design decisions of AI models and the overall performance of the applications that integrate them. We have obtained results that prove that this relation exists and we have observed evidence of how adding more complexity to our AI models (more complex architectures or higher number of parameters) causes the models to require a higher usage of resources. We have tested this relationship in an environment in which we can derive statistically significant relations (i.e. a substantially big group of models with different configurations has been developed, variations on the design

---

<sup>13</sup> <https://github.com/acmsigsoft/EmpiricalStandards/blob/master/docs/Experiments.md>

decisions have been iterated while keeping all other factors constant, and a total of 6 datasets in the vision and language domains have been used). Hence, we believe that our experimentation successfully targets the most popular AI models in the research community, and that provides meaningful knowledge to the vision and language-based communities. Finally, we have also provided experience on the end-to-end AI-enabled mobile applications development lifecycle, and have reported some of the issues that arise during such engineering process. During this lifecycle, we have also described our experience when using profiling tools to monitor the performance of AI-enabled mobile applications. With all this, we have provided important findings that arise from the results obtained.

Additionally, we have provided practitioners from the AI engineering research and industry with awareness of how a set of design decisions on the AI models contributes to the overall performance of the AI-enabled applications. From this point, we motivate future work to characterise the relation that we have uncovered in more detail. For example, future work can *(i)* identify a bigger set of design decisions during the AI-enabled software engineering lifecycle; *(ii)* identify other green characteristics that can be monitored and studied as a function of the design decisions; *(iii)* extend the range of AI model architectures to study; *(iv)* generalise the edge devices to others than mobile devices; *(v)* extend the technology stacks used for experimentation; and *(vi)* provide experiences with other tools to monitor the performance of the AI-enabled applications.

**Acknowledgements** This work has been partially supported by the GAISSA project (TED2021-130923B-I00, which is funded by MCIN/AEI/10.13039/501100011033 and by the European Union “NextGenerationEU”/PRTR); and, by the “Beatriz Galindo” Spanish Program (BEAGAL18/00064).

## Data Availability Statement

We provide a public repository that can be found here. The repository contains *(i)* the source code to train all the AI models in the study, which includes the training datasets; *(ii)* the source code of the AI-enabled mobile applications; *(iii)* the evaluation datasets used to profile the metrics in the study during operation; *(iv)* the profiled datasets containing the values of the profiled metrics (i.e. accuracy, time of CPU usage, storage weight profiled during operation); and *(v)* the source code to carry the statistical analysis of the profiled metrics.

## Conflicts of Interest

The authors declared that they have no conflict of interest.

## References

1. M. Mohri, A. Rostamizadeh, A. Talwalkar, *Foundations of machine learning* (MIT press, 2018)
2. Y. LeCun, Y. Bengio, G. Hinton, *nature* **521**(7553), 436 (2015)
3. C.C. Tappert, C.Y. Suen, T. Wakahara, *IEEE Transactions on pattern analysis and machine intelligence* **12**(8), 787 (1990)
4. D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, et al., *nature* **550**(7676), 354 (2017)
5. A.v.d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, K. Kavukcuoglu, arXiv preprint arXiv:1609.03499 (2016)
6. H. Mao, M. Cheung, J. She, in *Proceedings of the 25th ACM international conference on Multimedia* (2017), pp. 1183–1191
7. C. Byun, W. Arcand, D. Bestor, B. Bergeron, M. Hubbell, J. Kepner, A. McCabe, P. Michaleas, J. Mullen, D. O’Gwynn, et al., in *2012 IEEE Conference on High Performance Extreme Computing* (IEEE, 2012), pp. 1–6
8. K. Dowd, C. Severance, (2010)
9. T. Brown, B. Mann, N. Ryder, M. Subbiah, J.D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al., *Advances in neural information processing systems* **33**, 1877 (2020)
10. J. Rasley, S. Rajbhandari, O. Ruwase, Y. He, in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (2020), pp. 3505–3506
11. R. Zellers, A. Holtzman, H. Rashkin, Y. Bisk, A. Farhadi, F. Roesner, Y. Choi, arXiv preprint arXiv:1905.12616 (2019)
12. J. Devlin, M.W. Chang, K. Lee, K. Toutanova, arXiv preprint arXiv:1810.04805 (2018)
13. G. Hinton, O. Vinyals, J. Dean, et al., arXiv preprint arXiv:1503.02531 **2**(7) (2015)
14. R. Schwartz, J. Dodge, N.A. Smith, O. Etzioni, arXiv preprint arXiv:1907.10597 (2019)
15. C. Calero, M.Á. Moraga, M. Piattini, *Software Sustainability* pp. 1–15 (2021)
16. A. Lacoste, A. Luccioni, V. Schmidt, T. Dandres, arXiv preprint arXiv:1910.09700 (2019)
17. L. Cruz, R. Abreu, in *2019 IEEE/ACM 41st International Conference on Software Engineering: New Ideas and Emerging Results (ICSE-NIER)* (2019), pp. 101–104. DOI 10.1109/ICSE-NIER.2019.00034
18. M. Tan, Q. Le, in *International Conference on Machine Learning* (PMLR, 2019), pp. 6105–6114
19. V. Sanh, L. Debut, J. Chaumond, T. Wolf, arXiv preprint arXiv:1910.01108 (2019)
20. Z. Chen, Y. Cao, Y. Liu, H. Wang, T. Xie, X. Liu, in *Proceedings of the 28th ACM ESEC/FSE* (2020), pp. 750–762
21. J. Siebert, L. Joeckel, J. Heidrich, A. Trendowicz, K. Nakamichi, K. Ohashi, I. Namba, R. Yamamoto, M. Aoyama, *Software Quality Journal* (2021). DOI 10.1007/s11219-021-09557-y. URL <https://doi.org/10.1007/s11219-021-09557-y>
22. Q. Guo, S. Chen, X. Xie, L. Ma, Q. Hu, H. Liu, et al., in *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)* (IEEE, 2019), pp. 810–822
23. R.C. Castanyer, S. Martínez-Fernández, X. Franch, in *2021 IEEE/ACM 1st Workshop on AI Engineering-Software Engineering for AI (WAIN)* (IEEE, 2021), pp. 27–34
24. L.E. Lwakatere, A. Raj, I. Crnkovic, J. Bosch, H.H. Olsson, *Information and Software Technology* **127**, 106368 (2020)
25. S. Martínez-Fernández, J. Bogner, X. Franch, M. Oriol, J. Siebert, A. Trendowicz, A.M. Vollmer, S. Wagner, *ACM Trans. Softw. Eng. Methodol.* **31**(2) (2022). DOI 10.1145/3487043. URL <https://doi.org/10.1145/3487043>
26. L.E. Lwakatere, I. Crnkovic, J. Bosch, in *2020 SoftCOM* (2020), pp. 1–6. DOI 10.23919/SoftCOM50211.2020.9238323
27. L. Pons, I. Ozkaya, arXiv preprint arXiv:1911.02912 (2019)
28. C. Calero, M. Piattini, (2017)
29. R. Verdecchia, J. Sallou, L. Cruz, arXiv preprint arXiv:2301.11047 (2023)

30. L. Bao, D. Lo, X. Xia, X. Wang, C. Tian, in *2016 IEEE/ACM 13th Working Conference on Mining Software Repositories (MSR)* (IEEE, 2016), pp. 37–48
31. L. Cruz, R. Abreu, arXiv preprint arXiv:1803.05889 (2018)
32. A. Banerjee, A. Roychoudhury, in *Proceedings of the International Conference on Mobile Software Engineering and Systems* (2016), pp. 139–150
33. Y. Xu, S. Martínez-Fernández, M. Martínez, X. Franch, (2023)
34. D. Di Nucci, F. Palomba, A. Prota, A. Panichella, A. Zaidman, A. De Lucia, in *2017 IEEE 24th International Conference on Software Analysis, Evolution and Reengineering (SANER)* (2017), pp. 103–114. DOI 10.1109/SANER.2017.7884613
35. L. Cruz, R. Abreu, *Empirical Software Engineering* **24**(4), 2209 (2019)
36. S. Chowdhury, S. Borle, S. Romansky, A. Hindle, *Empirical Software Engineering* **24**(4), 1649 (2019)
37. S. Georgiou, M. Kechagia, T. Sharma, F. Sarro, Y. Zou, (ACM: Association for Computing Machinery, 2022)
38. R. Verdecchia, L. Cruz, J. Sallou, M. Lin, J. Wickenden, E. Hotellier, arXiv preprint arXiv:2204.02766 (2022)
39. V.R. Basili, G. Caldiera, D.H. Rombach. The Goal Question Metric Approach. *Encyclopedia of Software Engineering 1* (1994) (1994)
40. M.B. Miles, A.M. Huberman, *Qualitative data analysis: An expanded sourcebook* (sage, 1994)
41. T.L. Lai, H. Robbins, C.Z. Wei, *Journal of multivariate analysis* **9**(3), 343 (1979)
42. P. Pope, J. Webster, *Technometrics* **14**(2), 327 (1972)
43. J. Miles, Wiley StatsRef: Statistics Reference Online (2014)
44. Student, *Biometrika* pp. 1–25 (1908)
45. R.C. Castanyer, S. Martínez-Fernández, X. Franch, ESEM 2021 REGISTERED REPORT. Available on arXiv preprint arXiv:2109.15284 (2021). URL <https://arxiv.org/abs/2109.15284>
46. D. Méndez Fernández, M. Monperrus, R. Feldt, T. Zimmermann, *Empirical Software Engineering* **24**(3), 1057 (2019)
47. J. Stallkamp, M. Schlipsing, J. Salmen, C. Igel, *Neural networks* **32**, 323 (2012)
48. L. Deng, *IEEE Signal Processing Magazine* **29**(6), 141 (2012)
49. A. Krizhevsky, G. Hinton, et al., (2009)
50. A. Go, R. Bhayani, L. Huang, CS224N project report, Stanford **1**(12), 2009 (2009)
51. J. Ni, J. Li, J. McAuley, in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* (2019), pp. 188–197
52. A.L. Maas, R.E. Daly, P.T. Pham, D. Huang, A.Y. Ng, C. Potts, in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies* (Association for Computational Linguistics, Portland, Oregon, USA, 2011), pp. 142–150. URL <http://www.aclweb.org/anthology/P11-1015>
53. R. Verdecchia, P. Lago, C. Ebert, C. de Vries, *IEEE Software* **38**(6), 7 (2021). DOI 10.1109/MS.2021.3102254