



# Convolutional neural networks for enhanced classification mechanisms of metamodels

Phuong T. Nguyen<sup>a</sup>, Davide Di Ruscio<sup>a,\*</sup>, Alfonso Pierantonio<sup>a</sup>, Juri Di Rocco<sup>a</sup>, Ludovico Iovino<sup>b</sup>

<sup>a</sup> Università degli studi dell'Aquila, 67100 L'Aquila, Italy

<sup>b</sup> Gran Sasso Science Institute, Italy

## ARTICLE INFO

### Article history:

Received 16 June 2020

Received in revised form 1 October 2020

Accepted 2 November 2020

Available online 8 November 2020

## ABSTRACT

Conventional wisdom on Model-Driven Engineering suggests that metamodels are crucial elements for modeling environments consisting of graphical editors, transformations, code generators, and analysis tools. Software repositories are commonly used in practice for locating existing artifacts provided that a classification procedure is available. However, the manual classification of metamodel in repositories produces results that are influenced by the subjectivity of human perception besides being tedious and prone to errors. Therefore, automated techniques for classifying metamodels stored in repositories are highly desirable and stringent. In this work, we propose MEMOCNN as a novel approach to classification of metamodels. In particular, we consider metamodels as data points and classify them using supervised learning techniques. A convolutional neural network has been built to learn from labeled data, and use the trained weights to group unlabeled metamodels. A comprehensive experimental evaluation proves that the proposal effectively categorizes input data and outperforms a state-of-the-art baseline.

© 2020 Elsevier Inc. All rights reserved.

## 1. Introduction

Model-Driven Engineering (Schmidt, 2006) (MDE) helps tame the complexity of nowadays large-scale software development by using models to capture relevant knowledge of the problem domain. Models are typically specified in terms of concepts formalized in metamodels, and they are manipulated via automated transformations to achieve superior automation, whether it be refactoring, simulation, or code generation. As in other knowledge-intensive domains, repositories are large depots of unstructured resources that permit to store experience (in terms of developed artifacts) and to make it available for later use. While the relevance of software repositories is often recognized, the difficulties in building effective, scalable, and useful tools are often underestimated (White et al., 2015). Leveraging the informative structure within a repository, which comprehends insights, know-how, and generally, accumulated and distributed knowledge, requires the ability to locate and make existing artifacts amenable to reuse, analysis, and documentation (among others) (Di Rocco et al., 2015). Over the last decade, several model

repositories have been proposed by both academia and industry, e.g., Hein et al. (2009), Karasneh and Chaudron (2013), Koegel and Helming (2010), Kutsche et al. (2008), just to mention a few. However, unless repositories allow modelers to easily discover and retrieve the artifacts that satisfy their needs, the potential benefits related to the availability of such resources can be significantly compromised. Clustering and classification approaches have already been proposed to categorize the metamodel repositories automatically (Basciani et al., 2016; Nguyen et al., 2019). In contrast, manual procedures to classify artifacts are typically time-consuming, prone to errors, and might lead to accuracy problems. Clustering techniques have been employed in order to group similar metamodels in an unsupervised manner (Basciani et al., 2016). Unfortunately, the performance of clustering techniques heavily relies on the chosen similarity measure, and two main problems that hinder its effectiveness are namely: (i) timing performance; and (ii) the identification of the appropriate number of clusters (Nguyen et al., 2019).

In recent years, Machine Learning (ML) gained an interest in the scientific community as it promises to yield unprecedented results and performance compared with conventional approaches (Goodfellow et al., 2016). The core principle of an ML algorithm is that it attempts to simulate the learning activities of human beings (Portugal et al., 2015). As a result, ML systems are capable of autonomously extracting meaningful patterns from real-world examples, like humans (Domingos, 2012; Wuest et al.,

\* Corresponding author.

E-mail addresses: [phuong.nguyen@univaq.it](mailto:phuong.nguyen@univaq.it) (P.T. Nguyen), [davide.diruscio@univaq.it](mailto:davide.diruscio@univaq.it) (D. Di Ruscio), [alfonso.pierantonio@univaq.it](mailto:alfonso.pierantonio@univaq.it) (A. Pierantonio), [juri.dirocco@univaq.it](mailto:juri.dirocco@univaq.it) (J. Di Rocco), [ludovico.iovino@gssi.it](mailto:ludovico.iovino@gssi.it) (L. Iovino).

2016). Recently, the AURORA tool (Nguyen et al., 2019) has been proposed to classify metamodel repositories using neural networks automatically. By learning from a training set consisting of labeled metamodels, it classifies incoming unlabeled ones based on well-established techniques. Even though the approach has been validated on a dataset consisting of more than 500 metamodels, it has some performance issues as detailed later in this paper.

In this work, metamodels are considered as data points and they are classified by means of supervised learning techniques. In particular, we exploit convolutional neural networks (CNNs) (Krizhevsky et al., 2012) to learn from labeled data, and use the trained weights to group unlabeled metamodels. Such networks have been successfully employed in various domains, such as imaging applications (Zhao et al., 2019; He et al., 2015; Nguyen et al., 2016), and natural language processing (Belinkov and Glass, 2019). By means of supervised learning, we overcome the main limitation of different unsupervised clustering techniques where it is necessary to specify the number of clusters in advance. Meanwhile, with supervised learning the number of clusters depends on the training data, and it has been specified ex-ante in the labeled data. In our experiment, we used a dataset with nine categories, however such a number can easily be changed if there are more categories previously classified. In this respect, there is no need to dictate any number in advance with our approach.

The rationale behind the selection of a CNN is as follows. To choose a suitable network for classifying metamodels, we performed an investigation into various techniques, namely recurrent neural networks (RNNs), long short-term memory (LSTM) neural networks (Hochreiter and Schmidhuber, 1997), and CNNs. As a matter of fact, RNNs and LSTMs are suitable for working with time series data such as weather forecast or pedestrian trajectory (Zhao et al., 2019). We anticipate that they can be used to solve other tasks in MDE, e.g., recommendation of model repairs, or for assisting the specification of metamodels. For classification, we select CNNs as they have been designed to successfully work with images. The recent development of disruptive deep neural networks which are based on a CNN has enabled a large number of applications in different domains. Thus, the deployment of a CNN to classify metamodels allows us to take advantage of a well-founded background related to many technological breakthroughs, which have been proven through a series of work in image classification.

By mimicking the human visual system (Rawat and Wang, 2017), the approach (named MEMOCNN hereafter) considerably improves the accuracy of AURORA by outperforming the baseline for different quality metrics as demonstrated by the experiments. Moreover, the analysis reveals that MEMOCNN becomes more effective and efficient when there is more data available for training. In this respect, the contributions of the work can be summarized as follows: (i) a representation scheme to transform metamodels into a CNN-processable format; (ii) a convolution neural network for automatic classification of metamodel repositories; (iii) an empirical evaluation with different experimental configurations to study the proposed approach's performance; and finally (iv) a MEMOCNN prototype together with the meta-data used in our evaluation is released to the benefit of future research.<sup>1</sup>

**Outline.** The structure of the paper is as follows. Section 2 presents a background for our work. The proposed approach is illustrated in Section 3. The next section explains the evaluation settings used to study the approach's performance. The experimental results are reported and analyzed in Section 5. In

Section 6, we discuss the outcomes of the experiments and highlight the probable threats to validity. Section 7 reviews the related work and associates them with our approach. Finally, Section 8 outlines possible future work and draws some conclusions.

## 2. Background

Section 2.1 provides an introduction to metamodels. Afterward, Section 2.2 describes the problem of classification as a machine learning task. We introduce classification of metamodels in Section 2.3. Finally, a brief introduction to CNNs is provided in Section 2.4, highlighting the most relevant aspects.

### 2.1. A brief introduction to metamodels

A metamodel defines the abstract concepts of a domain where concepts, as well as the relationships among them, are expressed by modeling infrastructure (Schmidt, 2006). The Eclipse Modeling Framework (EMF)<sup>2</sup> provides modelers with an infrastructure to design, define and edit models. EMF includes *Ecore* as a meta-metamodel for describing metamodels. Fig. 1 depicts a simplification of the *Ecore* metamodel. All the represented meta classes are named elements. *EPackage* is composed of subpackages as well as classes, while *EClass* consists of *EStructuralFeatures*, i.e., *EAttributes* and *EReferences*. Moreover, *EClass* can inherit structural features from other classes. *EAttribute* and *EReference* inherit the lowerbound and upperbound attributes in order to define the cardinalities of structural features. *EReference* links the container to *EClass* by *eReferenceType* reference. Moreover, containment is a boolean attribute, which allows one to specify if the linked *EClass* is contained, e.g., to trigger cascade-deletion, or not in the container being modeled. *EAttributes* are typed as *EDataType* instance by *eAttributeType* reference. The encoding schema that we are going to present in Section 3.2 relies on the *Ecore* simplification depicted in Fig. 1.

### 2.2. Classification as a supervised learning task

Classification belongs to supervised learning where training data is used to identify recurrent patterns, which in turn are used to classify new instances. In other words, classification relies on the existence of predefined classes of objects and aims to assign a new item to a class. In this respect, labeled training data is needed to guide the learning process which is conducted to refine the constituent weights and biases, by introducing all the training instances. Meanwhile, testing is to predict a label for an input instance by means of the weights and biases previously obtained from learning. Several machine learning algorithms can be exploited for classification, e.g., Decision Trees, Random Forest, or Support Vector Machines (Kotsiantis et al., 2006). Among others, neural networks have been widely used for classification (Zhang, 2000), due to their effectiveness.

Classification offers a wide range of applications, covering different domains. The most prominent use cases include speech recognition (Nassif et al., 2019), biometric identification (Sundarajan and Woodard, 2018), document classification (Sebastiani, 2002), to name but a few. In image processing, the extraction of low-level features from images and videos have been extensively investigated (Hartmann and Lienhart, 2002), and recently the focus has been shifted to automatically classify images and videos into semantically meaningful categories. For instance, HTML documents have been classified into a hierarchy of categories (Ceci et al., 2003) to support Web documents sharing in

<sup>1</sup> <https://github.com/MDEGroup/memoCNN/>

<sup>2</sup> <https://www.eclipse.org/modeling/emf/>

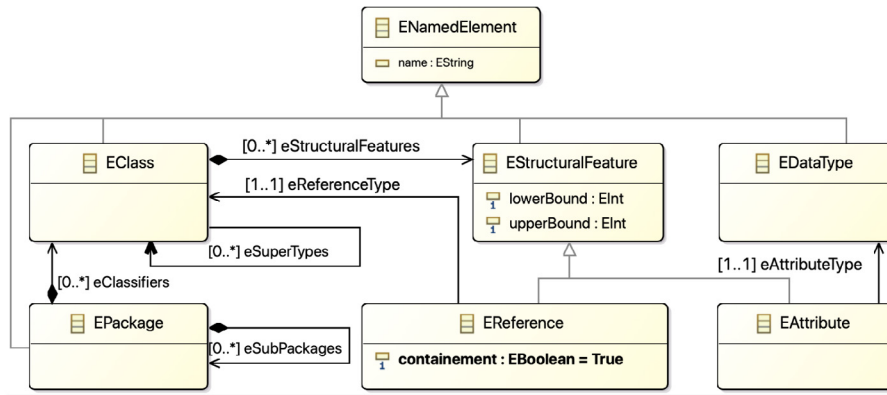


Fig. 1. An excerpt of the ecore metamodel.

a groups of users with common interests. Classification is also useful when the intent is to automatically and efficiently group images as static diagrams, e.g., UML diagrams. The automatic classification of images has been also successfully and extensively applied on medical images (Shi et al., 2009) to automate a process of pattern recognition related to some diseases. In recent years, the research community has been increasingly interested in categorizing large volumes of unstructured text, given the proliferation of social networks (Hartmann et al., 2019).

### 2.3. Classification of metamodels

Open-source software platforms typically provide developers with guidance to locate software artifacts to be reused through various built-in facilities. Unfortunately, finding relevant information in a repository presents technical difficulties as the stored artifacts need to be accurately classified before they can be detected. Until recently, such classification has been operated manually by annotating the artifacts with metadata whose granularity is adapted to the level of accuracy dictated by the problem at hand. Like other hand-operated activities, the manual classification of items stored in repositories is tedious, time-consuming, and prone to error. To this end, clustering has demonstrated to be a workable solution to the problem (Basciani et al., 2016). However, clustering techniques suffer from a low efficiency as well as the difficulty in selecting an appropriate number of clusters (Nguyen et al., 2019).

AURORA represents the first approach in the domain of automated classification of metamodels with machine learning techniques (Nguyen et al., 2019). The system works on the premise that the prescribed information related to metamodels and their categories can be used to train a supervised classifier. The approach is based on a feed-forward neural network, and it learns from a training data consisting of labeled metamodels and classifies unknown ones. An empirical evaluation (Nguyen et al., 2019) on a dataset of manually classified metamodels showed that the system obtains a considerably high prediction performance.

Nevertheless, AURORA has its own drawbacks. Though it performs well on the given dataset, which consists of a small number of metamodels, i.e., 555 (see Section 4.2), it is less efficient on larger datasets as it is the case with feed-forward neural networks (Driss et al., 2017). As a consequence, with larger metamodels, the network performs slowly, and the learning process could take longer to converge, i.e., the accuracy can only be increased after several epochs. In this respect, convolutional neural networks (CNN) may be a natural choice to obtain a scalable classification procedure. Indeed, such kind of networks is characterized by a degree of accuracy that improves with larger dataset populations. In the next subsection, we introduce the technical background of CNNs as a base for further discussions.

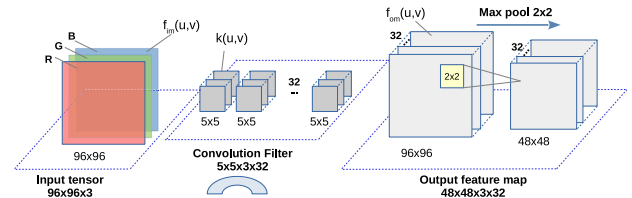


Fig. 2. Filtering and pooling.

### 2.4. Convolutional neural networks

Ordinary neural networks connect all the neurons of a layer to the next layer. In this way, the number of weights and biases increases exponentially with the number of layers. Such a structure is not suitable for working with images since it adds up to a huge number of parameters. CNNs (Rawat and Wang, 2017) have been specifically designed for working on images, and they attempt to capture the spatial and temporal structures, which are intrinsic features of an image. Furthermore, CNNs also minimize the number of parameters by using convolutional filters. A *filter* or *kernel* is a small square matrix used to capture specific features from an input image, e.g., nodes and edges. In practice, several kernels can be used to extract various types of features from images. Given a convolutional layer, each perceptron gets input from a certain part from the previous layer, which is called a *receptive field*. The convolution operation is performed by sliding the filter along the width and height of a *feature map*, which is either the input image, or the result of the convolution operation. An output feature map of one layer becomes the input feature map of the succeeding layer. A CNN consists of three types of layers, i.e., convolutional, pooling, and fully-connected layers, whose functionalities are explained below:

- **Convolution layer:** it performs most of the computation within a CNN to extract important features of an input image;
- **Pooling layer:** it scales down a feature map by taking the maximum value within a window, preferably a square one, to reduce the number of parameters (Rawat and Wang, 2017). The layer also extracts rotational and positional invariant features from a feature map;
- **Fully-connected layer:** the layer works as an ordinary perceptron, each of its neurons is fully connected to the previous layer.

For illustration purpose only, we take an example in Fig. 2 where there is an excerpt of a CNN. A tensor of size  $96 \times 96 \times 3$

represents an input image, and the number 3 corresponds to the three color channels in images, i.e., Red (R), Green (G), and Blue (B). A 4D filter of size  $5 \times 5 \times 3 \times 32$  is convolved with the input feature map to produce an output feature map of size  $96 \times 96 \times 1 \times 32$ . To elaborate on the computation, we consider only a slice of the 4D filter, corresponding to a matrix  $k(u, v)$ . Given an input feature map  $f_{im}(u, v)$ , the output feature map  $f_{om}(u, v)$  is created by convolving  $k(u, v)$  with  $f_{im}(u, v)$ :

$$f_{om}(u, v) = (f_{im} * k)(u, v) = \sum_i \sum_j f_{im}(i, j) \cdot k(u - i, v - j) \quad (1)$$

Pooling is performed by means of Maxpool  $2 \times 2$  to reduce the resulted feature map's width and height to a half, yielding a new feature map of size  $48 \times 48 \times 1 \times 32$ .

Furthermore, in a CNN there are the following features:

- **Dropout**: it is a parameter used to distribute the learned representation across all the neurons. It is an effective measure to minimize overfitting (Srivastava et al., 2014);
- **Softmax**: the last fully-connected layer of a CNN normally uses Softmax as the activation function and it converts a set of real numbers to probabilities which sum to 1.0 (Rawat and Wang, 2017). Given  $C$  classes, and  $y_k$  is the output of the  $k^{th}$  neuron, the final prediction is the class that gets the maximum probability, i.e.,  $\hat{y} = \operatorname{argmax} p_k, k \in 1..C$ , where  $p_k$  is computed below.

$$p_k = \frac{\exp(y_k)}{\sum_{k=1}^C \exp(y_k)} \quad (2)$$

- **Rectified Linear Units (ReLU)**: convolutional layers use ReLU as the activation function, which returns 0 given a negative input, and returns the input itself if it is larger than 0,  $f(x) = \max(0, x)$ .

### 3. MEMOCNN: A convolutional neural network to classify metamodels

In an attempt to simulate humans' cognition towards the relationship between metamodels and their categories, we conceptualized the usage of a well-defined deep learning technique, i.e., convolutional neural networks. In this section, we present MEMOCNN, an approach to classification of metamodels exploiting a Convolutional Neural Network. Section 3.1 presents the conceived architecture to realize MEMOCNN. An explanatory example showing how metamodels are encoded is presented in Section 3.2. Finally, Section 3.3 illustrates a concrete configuration of a MEMOCNN network.

#### 3.1. Architecture

Fig. 3 shows the main components of the MEMOCNN architecture supporting two main phases, i.e., *training* and *deployment*. Metamodels can be collected from different sources, and in the scope of this paper we make use of GITHUB as the main data source, since the platform offers a large number of metamodels repositories. Afterward, expert modelers are asked to inspect the metamodels and manually classify them into independent groups based on their experience. By means of labeled data including metamodels and labels assigned by humans, the training process is in charge of identifying the right *weights* to configure the CNN CLASSIFIER component. Finally, the weights obtained by CNN CLASSIFIER are then utilized to classify metamodels given as input to the deployment process.

The TRAINING phase starts with collecting metamodels from public repositories like GITHUB by means of the CRAWLER component ①. NLP PROCESSOR ② applies Natural Language Processing

steps (i.e., stop words removal, stemming, and lemmatization) on all the named elements, which are contained in the collected metamodels. DATA ENCODER ③ takes as input the manipulated metamodels and encodes them as matrixes enabling the subsequent phases of the approach. To identify the properties to be included in the encoding, we followed a *feature engineering* process that started with considering only terms of the metamodels (uni-gram), we subsequently added containment information (bi-gram), up to including all the missing metamodel properties, including relation cardinalities, typing information, etc. (n-gram). In particular, with the uni-gram encoding scheme, only the terms contained in the metamodels being processed are considered. All the aspects related to metamodel structure are discarded. With the bi-gram scheme, information about containment relations between named elements is also encoded. The n-gram scheme fully encodes the input metamodels by enriching the bi-gram scheme by including all the properties of the represented structural features (e.g., typing information, relations cardinalities, etc.). The whole approach has been evaluated by considering the different encoding schemes individually.

The proposed architecture allows for both offline and online training ④. By the former, one can train the network on a PC or a laptop installed with TensorFlow and Keras. By the latter, the code is executed directly on Google Colaboratory<sup>3</sup> connecting to metadata stored into Google Drive.<sup>4</sup> The final outcome of the training process is a set of weights and biases, which then can be used to deploy the network.

The DEPLOYMENT phase shown in Fig. 3 refers to the adoption of a previously trained and deployed CNN CLASSIFIER ⑦ to automatically classify input metamodels. To this end, as also done in the training phase, the metamodels of interest have to be manipulated by means of NLP steps ⑤ and properly represented by means of one of the encoding scheme previously discussed ⑥. MEMOCNN has been implemented on top of TensorFlow and Keras,<sup>5</sup> two well-founded Machine Learning frameworks. To facilitate future research, we made available online the MEMOCNN tool and the metadata used in the evaluation.

#### 3.2. Metamodel encoding

We propose a method to transform a metamodel into a picture by making use of its modeling elements such as metaclasses, packages, and attributes. In other words, we attempt to *picturize* a metamodel, exploiting its constituent features. Starting from a set of metamodels, we construct a set of terms and use it to convert all the metamodels into vectors.

To illustrate how the data processing works, we consider the *Library* metamodel in Fig. 4. This metamodel can be used for defining models, where a *Library* can be composed of a set of *Books* and *Writers*. A writer can edit books, and then a book can be written by a set of writers. A writer is defined in terms of its *firstname* and *lastname*, whereas a book is characterized by its *title* and the number of *pages*.

The lexical terms in the metamodel are extracted and encoded as prescribed by the above encoding schemes. Table 1 summarizes the outcome by presenting all lexical terms in the metamodel, and the corresponding encodings for both the raw and normalized feature vectors. For the sake of presentation,

<sup>3</sup> <https://colab.research.google.com/>

<sup>4</sup> <https://drive.google.com/>

<sup>5</sup> Keras has been integrated into TensorFlow and provides to developers a higher abstract level, and thus a more user-friendly programming interface. Meanwhile, TensorFlow offers a greater flexibility in managing a network's settings. We exploited them both to build MEMOCNN, thereby earning different degrees of granularity for the experimental configurations.



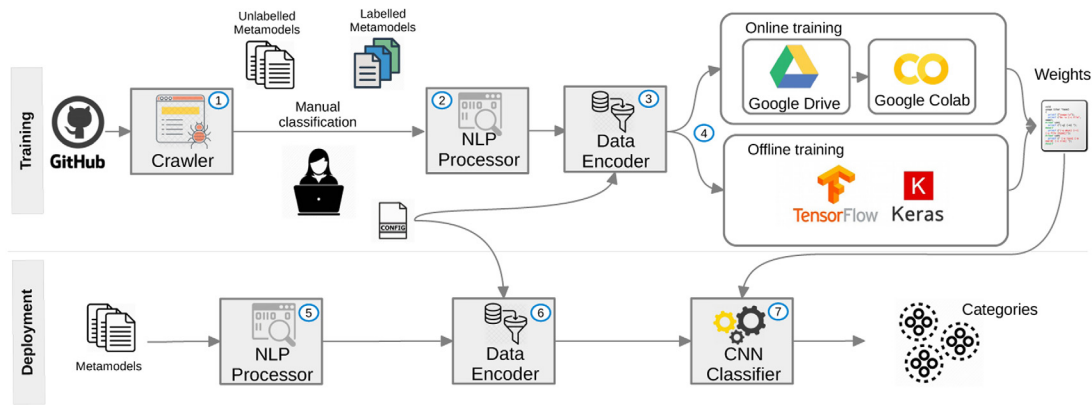


Fig. 3. MEMOCNN's system architecture.

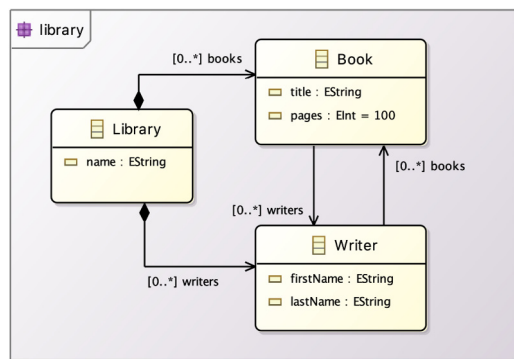


Fig. 4. The Library metamodel.

each metamodel element name is given with an identifier and its metamodel type. For instance, the package name *library* is (i) individually encoded in uni-grams; (ii) in bi-grams as part of the terms of the metaclasses #2, #6 and #10; and (iii) in n-grams as part of the terms of the structural features #3, #4, #5, #7, #8, #9, #11, #12 and #13. In addition, the encodings related to normalized feature vectors require the term to be subject to stemming, lemmatization, and stop words removal. Thus, the term *books* is translated into *book* and comprises both the metaclass *Book* and the feature *books*. Moreover, additional information about cardinalities, e.g., 0.1 and whether the node is subject to containment, e.g., TRUE, is appended.

With the aim of giving a flavour of the NLP steps that are applied by NLP PROCESSOR depicted in Fig. 3 (see ② and ⑤), the *normalized feature vectors* shown in Table 1 represent the output of stemming, lemmatization, and stop words removal tasks applied on the *raw feature vectors*. For Step ⑤, we make use of the *nlTK* python library<sup>6</sup> to perform the required NLP steps. First, a term is tokenized into subterms by using as delimiters any non-alphabetical character and blank spaces. After that, the *nlTK Wordnet Lemmatizer* module is invoked to squeeze the inflected subterms to the root English term. Then, the *Porter stemmer* algorithm (Porter et al., 1980) reduces the subterms by removing the common English endings from them. Afterward, English stop words which form the list of subterms are also removed. Finally, the resulting subterms are concatenated to create features.

By analyzing all the input metamodels, the corresponding terms are collected and used as indexes of vectors each encoding a corresponding metamodel. To illustrate how the encoding

paradigm works, we take an explanatory example in Fig. 5(a). For the sake of presentation, we consider a corpus with a total of 16 terms.<sup>7</sup> In the considered example, there are 13 of them related to the bi-gram encoding of the *Library* metamodel previously discussed. In Fig. 5(a), the blank cells represent no occurrences of the related terms, or in other words the corresponding value is equal to 0. Nevertheless, we refrain from using the one-hot vector encoding scheme, since such a representation neglects the importance of each constituent term. Thus, in our encoding scheme, the cells with a dot contain real numbers which quantify the occurrences of the corresponding terms in the encoded metamodel using the TF-IDF function. In particular, we normalized such values with respect to the most occurred terms in the dataset by *min-max feature scaling* (Juszczak et al., 2002). This allows us to keep a notionally common scale among the metamodel features representations.

Starting from the vector in Fig. 5(a), we serialize it into a matrix, preferably a square one as shown in Fig. 5(b). By traversing a vector from its beginning to the end, each time we select four consecutive entries to populate a row for the matrix. In this sense, the size of the vector needs to be a square number. So if the size is lower than the closest square number, then we pad the missing elements with zeros to constitute a square matrix. With this representation, we attempt to capture features by employing a unified structure across all the metamodels: each cell in the matrix represents a specific term, and adjacent cells with a dot form patterns like nodes and edges, as we find in images. In this respect, we hypothesize that metamodels within a same categories would share common pictorial patterns.

In real usage, even when the dataset becomes bigger, we can always apply the same process to parse and represent the data. A set of training metamodels is parsed using all the aforementioned steps, resulting in a corpus of the terms that appear in the dataset. Different from the previous work about AURORA (Nguyen et al., 2019), in this work we do not use any cut-off value but make use of all the parsed terms. Given a large dataset, each metamodel is then serialized into a vector, albeit with several entries. The final matrix will be populated accordingly, which is much bigger compared to the example one in Fig. 5(b). In this way, we have a matrix of  $M \times M$  size, where  $M$  can be hundreds, or thousands, like in image classification with high resolution images.

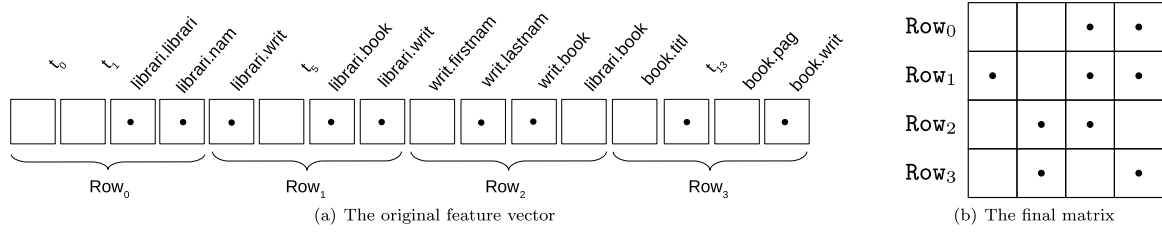
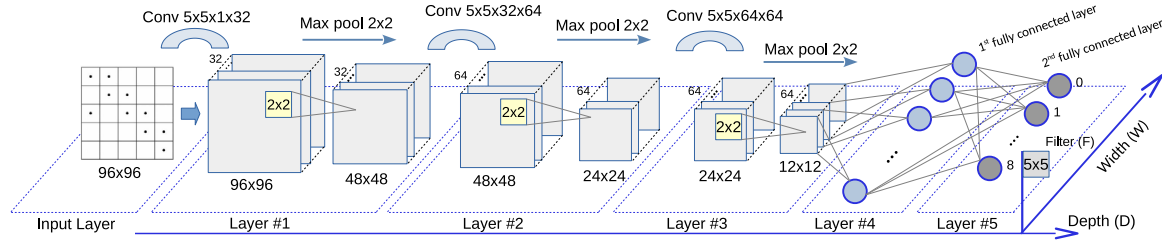
To demonstrate how the resulting pictures look like in practice, we exploit a real dataset whose characteristics are described in detail in Section 4.2. We represent them in the bitmap format following the proposed representation (cf. Figs. 5(a)) and 5(b)). We select for each class a metamodel and show the corresponding images, which are all of size  $148 \times 148$ , in Figs. 7(a)–7(i).

<sup>6</sup> <https://www.nltk.org/>

<sup>7</sup> In practice, the number of terms is much larger.

**Table 1**The extraction of the *Library* metamodel.

Metamodel element			Raw feature vector			Normalized feature vector		
Type	id	Name	uni-gram	bi-gram	n-gram	uni-gram	bi-gram	n-gram
package	#1	library	library	–	–	library*	–	–
metaclass	#2	Library	Library	library.Library	–	library*	librari.librari	–
attribute	#3	name	name	Library.name	Library.name.EString.0.1	name	librari.name	librari.name.estring.0.1
reference	#4	writers	writers	Library.writers	Library.writers.Writer.false.0.-1	writ**	librari.writ	librari.writ.writ.false.0.-1
reference	#5	books	books	Library.books	Library.books.Book.false.0.-1	book***	librari.book	librari.book.book.false.0.-1
metaclass	#6	Writer	Writer	library.Writer	–	writ**	librari.writ	–
attribute	#7	firstName	firstName	Writer.firstName	Writer.firstName.EString.0.1	firstnam	writ.firstName	writ.firstName.estring.0.1
attribute	#8	lastName	lastName	Writer.lastName	Writer.lastName.EString.0.1	lastnam	writ.lastName	writ.lastName.estring.0.1
reference	#9	books	books	Writer.books	Writer.books.Book.false.0.-1	writ.book	writ.book	writ.book.book.false.0.-1
metaclass	#10	Book	Book	library.Book	–	book***	librari.book	–
attribute	#11	title	title	Book.title	Book.title.EString.0.1	titl	book.titl	book.titl.estring.0.1
attribute	#12	pages	pages	Book.pages	Book.pages.EInt.0.1	page	book.pag	book.pag.eint.0.1
reference	#13	writers	writers	Book.writers	Book.writers.Writer.false.0.-1	writ**	book.writ	book.pag.eint.0.1

**Fig. 5.** The encoding of the *Library* metamodel.**Fig. 6.** A MEMOCNN network with five layers.

From the figures, it is evident that each class has a different distribution, and thus the density of dots across an image. For instance, Class 2 and Class 5 are much denser compared to Class 6 or Class 7. As a matter of fact, CNNs are highly suitable for capturing this type of features, e.g., nodes or edges (Krizhevsky et al., 2012), since they use the convolution operation to filter and retain features from input images. In this respect, we suppose that our proposed representation scheme will facilitate the recognition of metamodels. We are going to validate this hypothesis with concrete experiments in Section 5.

### 3.3. An example of MEMOCNN network

Fig. 6 illustrates an example of a MEMOCNN network. The input layer corresponds to an input metamodel which is represented as a 2D image of size  $96 \times 96$ . There are three convolutional layers, i.e., Layer #1, #2, #3 and two fully connected (FC) layers, i.e., Layer #4 and Layer #5. While the number of neurons of the 1st FC layer can be experimentally configured, the number of neurons of the 2nd FC layer is set to the number of categories. For each transition from one layer to the next one, a convolution filter is applied on the input feature map, resulting in an output feature map. For example, by Layer #3, a filter of size  $5 \times 5 \times 64 \times 64$  is convolved with the input feature map, and this yields an output feature map of size  $24 \times 24 \times 1 \times 64$ . Within the layer, pooling is conducted to reduce the dimension, producing an output feature map of size  $12 \times 12 \times 1 \times 64$ . From the left to the

right side of Fig. 6, the size of the feature maps decreases, whilst the number of filters increases. The network was the result of a series of experiments on the given dataset. We suppose that such a configuration is subject to change in practice, depending on the input data. For a larger dataset, we expect a deeper and/or wider network, which is generally the case with image datasets (Duong et al., 2020).

During the experiments conducted to evaluate MEMOCNN, we changed the number of layers as well as the number of filters to create different configurations and they are going to be clarified in the next section.

## 4. Evaluation

This section presents the experimental evaluation conducted to analyze the performance of MEMOCNN. Section 4.1 introduces the research questions being addressed in this paper. The dataset and the metrics used to evaluate the proposed approach are illustrated in Sections 4.2 and 4.3, respectively. Afterward, Section 4.4 highlights the experimental configurations.

### 4.1. Research questions

By using a series of experiments on a dataset that has been manually classified, we aim to answer the following research questions:

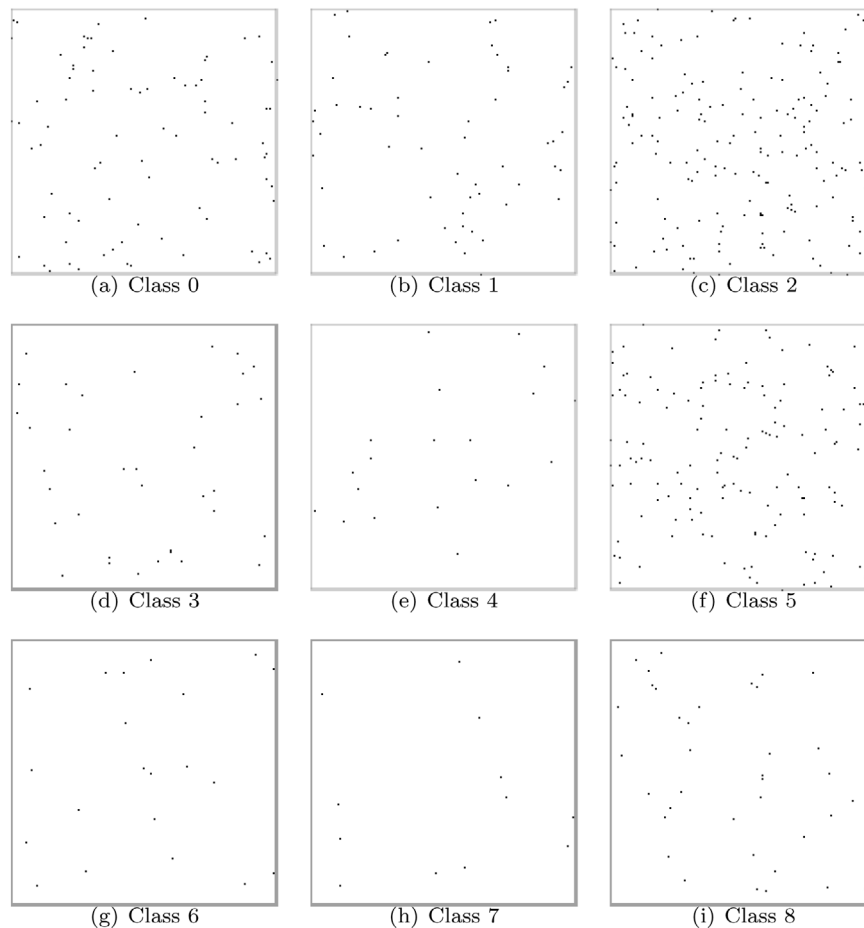


Fig. 7. Examples of the metamodels after being transformed into images.

- *RQ<sub>1</sub>: How effective is MEMOCNN at classifying the considered categories?* It is of great importance to know if MEMOCNN can accurately classify metamodels into their correct categories.
- *RQ<sub>2</sub>: Which network configuration contributes to a better performance?* We investigate whether the approach yields scalable performance by analyzing how a more complex network helps enhance the prediction accuracy on the given dataset.
- *RQ<sub>3</sub>: How does MEMOCNN compare with AURORA?* Finally, using a common dataset, we evaluate MEMOCNN against AURORA to ascertain whether the proposed approach gives a better prediction performance compared to the baseline.

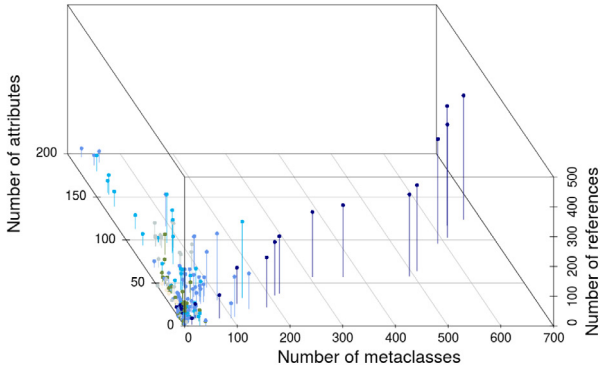
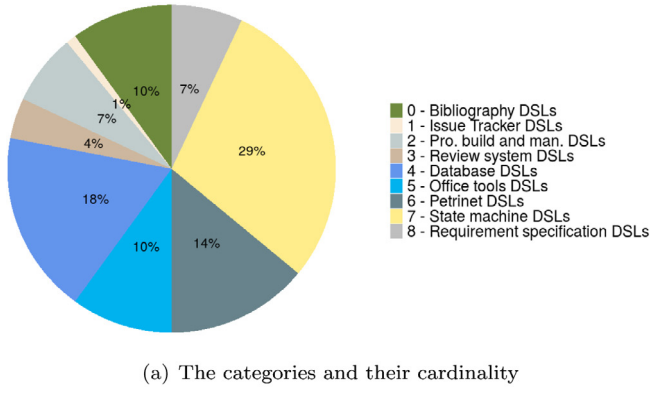
#### 4.2. Dataset

In the evaluation, we made use of the same dataset that was exploited to evaluate AURORA (Babur, 2019). The dataset consists of metamodels related to domain specific languages (DSLs) collected from GITHUB. All the metamodels were inspected and manually classified by an expert modeler based on his experience (Babur, 2019). The final dataset consists of 555 metamodels distributed in nine independent categories. A summary of the categories and their cardinality is provided in Fig. 8(a). The most populous category is *State machine DSLs* with 159 metamodels, accounting for 29% of the total amount. Meanwhile, the smallest category is *Issue Tracker DSLs* with only seven metamodels. Such a category may render the classification difficult since there is a limited amount of data for training. The distribution of metamodels with respect to the number of metaclasses, attributes, and

references is depicted in Fig. 8(b), where most of the metamodels contain a low number of metaclasses as well as attributes and references. The metamodels are concentrated on the bottom left corner of the tensor. Only those of the category *Review system DSLs* contain a large number of metaclasses, and they are scattered across the 3D space. In summary, the dataset exhibits a variety of categories with different characteristics and cardinalities. To identify possible cloned metamodels in the dataset, we exploited the *match-based* similarity function, which has been devised in one of our previous work (Basciani et al., 2016). In particular, the function takes two metamodels as input, i.e.,  $mm_1$  and  $mm_2$ , and it returns 1 if every element of  $mm_1$  exactly matches with one element of  $mm_2$ , and vice versa. By computing the similarity among all metamodels pairs, we detected 60 clones, accounting for 10.8% of the dataset. In practice, modeling repositories commonly include many model clones as it has been pointed out in a recent work (Lopes et al., 2017), where the authors conducted an empirical study on 4.5 millions of GITHUB repositories to identify code file duplicates. The study reveals that more than 40% of the code files are duplicated. Due to this reason, we believe that maintaining a controlled number of metamodel clones in the dataset helps resemble a real repository scenario. As a result, we decided to keep cloned metamodels in the dataset.

#### 4.3. Metrics

As seen in the dataset (Section 4.2), there are nine categories in total, and in the scope of this work we deal with multi-class (or multinomial) classification (Aly, 2005), i.e., each instance



**Fig. 8.** A summary of the dataset.

of the data needs to be classified into one of the nine categories. To thoroughly study the performance of our approach, we employed cross-validation (Kohavi, 1995) for evaluating the classification outcomes. For each of the categories mentioned in Section 4.2, we measured the relevance of the ground-truth data, i.e., the metamodels and their corresponding labels, with the classified results obtained by running the systems, i.e., AURORA and MEMOCNN. We are interested in understanding how well the produced classes match with the ground-truth data, considering the fact that the dataset is imbalance. Thus, *balanced accuracy*, *precision*, *recall*, *F<sub>1</sub>-score*, *ROC*, and *AUC* are utilized to study the systems' performance as detailed below. First, given a category, the following definitions are considered:

- *True positive (TP)*: the metamodels being correctly classified;
- *False positive (FP)*: the classified metamodels which actually do not match with the ground-truth data;
- *False negative (FN)*: the metamodels that should have been classified since they belong to the ground-truth data, but they are not;
- *True negative (TN)*: the metamodels that are not classified and they also do not belong to the ground-truth data.

Then, the metrics are defined as follows:

#### 4.3.1. Precision and recall

Precision and recall are used to quantify the level of relevance of the classified items to the ground-truth data as given below. Precision measures the fraction of correctly classified items to the total number of items, while Recall (or *True positive rate*, TPR) is

defined as the fraction of actual positive cases that are correctly classified.

$$P = \frac{TP}{TP + FP} \quad (3)$$

$$R = TPR = \frac{TP}{TP + FN} \quad (4)$$

#### 4.3.2. F<sub>1</sub>-score

It is a harmonic average of precision and recall.

$$F_1 = \frac{2 \cdot P \cdot R}{P + R} \quad (5)$$

In the context of multi-class classification, macro-averaged F<sub>1</sub>-score (macro-F<sub>1</sub> for short) is computed as the average of the F<sub>1</sub>-scores of all categories.

#### 4.3.3. False positive rate (FPR)

It is the ratio of the falsely classified items to the total number of items.

$$FPR = \frac{FP}{TN + FP} \quad (6)$$

#### 4.3.4. Balanced accuracy

As seen in Section 4.2, the dataset is imbalanced and thus we use balance accuracy to measure the performance as defined below:

$$Accuracy = \frac{TPR + FPR}{2} \times 100\% \quad (7)$$

#### 4.3.5. ROC curve and AUC

A receiver operating characteristic (ROC) curve represents the relationship between FPR (false positive rate) and TPR (true positive rate, Recall) (Fawcett, 2006). Such a curve is sketched in a 2D space where the x-Axis and the y-Axis correspond to FRP and TPR, respectively. In this respect, an ROC curve spans from (0, 0) to (1, 1) in the space. Furthermore, the area under the ROC curve (AUC) is used as an explicit indication of how good a classifier can be. A *dummy* classifier, i.e., the one that randomly assigns a label to metamodels, would have an AUC value of 0.5. In contrast, a perfect classifier would have an AUC value of 1.0. In this sense, an ROC curve being closer to the upper left corner accounts for a better prediction performance, compared to another one that is closer to the lower right corner in the 2D space.






Since ROC curves are dedicated for binary classification, in our work, we transform the multi-class problem into binary classification using the following method. Given one among the nine considered categories, i.e.,  $C_i$ , we consider it as the positive class, while all the remaining eight categories are considered as a single class, or negative. That means, if an instance is correctly classified into its category, we count it as a true positive. In contrast, if it is classified into any of the negative classes, then it is counted as a false negative. In this way, we are able to derive an ROC as like for binary classification.

#### 4.4. Configurations

In order to thoroughly study the performance of MEMOCNN, we considered different experimental settings as given in Table 2. The configurations are built by customizing various parameters, i.e., encoding schemes, the number of layers (depth), the number of filters (width), and kernel's size. Each change in the parameters aims to investigate possible effects on the final classification. By each configuration, the first layer's input tensor corresponds



**Table 2**  
Experimental configurations.

Configuration		C <sub>1</sub>	C <sub>2</sub>	C <sub>3</sub>	C <sub>4</sub>	C <sub>5</sub>	C <sub>6</sub>
ML framework		 Keras	 Keras	 Keras	 Keras	 Keras	TensorFlow
Encoding		uni-gram	bi-gram	n-gram	uni-gram	uni-gram	uni-gram
Layer #1	Input tensor	94 × 94 × 1	148 × 148 × 1	184 × 184 × 1	96 × 96 × 1	96 × 96 × 1	96 × 96 × 1
	Conv. kernel	3 × 3 × 1 × 16	3 × 3 × 1 × 16	3 × 3 × 1 × 16	3 × 3 × 1 × 16	3 × 3 × 1 × 16	5 × 5 × 1 × 32
	Pooling	Maxpool 2 × 2	Maxpool 2 × 2	Maxpool 2 × 2	Maxpool 2 × 2	Maxpool 2 × 2	Maxpool 2 × 2
	Dropout	0.5	0.5	0.5	0.5	0.5	0.5
	Act. Function	ReLU	ReLU	ReLU	ReLU	ReLU	ReLU
Layer #2	Input tensor				48 × 48 × 1 × 16	48 × 48 × 1 × 16	48 × 48 × 1 × 64
	Conv. kernel				3 × 3 × 1 × 32	3 × 3 × 1 × 32	5 × 5 × 32 × 64
	Pooling				Maxpool 2 × 2	Maxpool 2 × 2	Maxpool 2 × 2
	Dropout				0.5	0.5	0.5
	Act. Function				ReLU	ReLU	ReLU
Layer #3	Input tensor					24 × 24 × 1 × 32	24 × 24 × 1 × 64
	Conv. kernel					3 × 3 × 1 × 64	5 × 5 × 64 × 64
	Pooling					Maxpool 2 × 2	Maxpool 2 × 2
	Dropout					0.6	0.6
	Act. Function					ReLU	ReLU
1 <sup>st</sup> FC layer	# of nodes	32	32	32	128	128	256
	Act. Function	ReLU	ReLU	ReLU	ReLU	ReLU	ReLU
2 <sup>nd</sup> FC layer	# of nodes	9	9	9	9	9	9
	Act. Function	Softmax	Softmax	Softmax	Softmax	Softmax	Softmax
Number of epochs		15/30	15/30	15/30	15	15	15
Loss function [40]		Cross Entropy	Cross Entropy	Cross Entropy	Cross Entropy	Cross Entropy	Cross Entropy
Optimizer [41]		Adam	Adam	Adam	Adam	Adam	Adam

to the input *image*, and its size may vary depending on the encoding scheme. Three configurations namely C<sub>1</sub>, C<sub>2</sub>, C<sub>3</sub> are used to address RQ<sub>1</sub> and RQ<sub>3</sub>. Four configurations, namely C<sub>1</sub>, C<sub>4</sub>, C<sub>5</sub>, and C<sub>6</sub>, are used to answer RQ<sub>2</sub>. To study the performance of MEMOCNN in comparison to AURORA, we set up only one layer for MEMOCNN, and this is the reason why there is no information for Layer #2 and Layer #3 in Table 2, i.e., the corresponding cells are left blank. The number of nodes for the 2nd fully-connected layer corresponds to the number of categories, as shown in Fig. 8(a), and this explains why the number is always 9 throughout the table.

All the configurations have been built using Keras except C<sub>6</sub> that has been made with TensorFlow since the framework allows for more control over the network. It is worth noting that the Dropout value has been empirically set to 0.5 or 0.6 (Srivastava et al., 2014). By all, except the 2nd fully-connected layer, the activation function ReLU is used. The Softmax function is selected for the output layer as it is suitable for classifying inputs into multiple categories (see Section 2.4). Moreover, we used the cross entropy loss function (Zhang and Sabuncu, 2018) and the Adam optimizer (Kingma and Ba, 2015) throughout the experiments as these two functions have been widely exploited in various classification scenarios. We are going to give more details for each configuration in Section 5.

## 5. Experimental results

Although MEMOCNN supports both offline and online training (see Fig. 3), in the scope of this paper, we opt for the latter since Google offers a robust platform that helps accelerate the computation with MEMOCNN. In particular, we uploaded the related metadata to Google Drive to feed the network, which is then executed on Google Colab. This section reports and analyzes the results in detail by answering the research questions in Section 4.1.

### 5.1. RQ<sub>1</sub>: How effective is MEMOCNN at classifying the considered categories?

First, we ran experiments following the ten-fold cross-validation technique, using three configurations C<sub>1</sub>, C<sub>2</sub>, C<sub>3</sub>, and the number of epochs was set to 15 (cf. Table 2). We computed class-wise average precision, recall, and F<sub>1</sub>-score and show them in Table 3. The results in Table 3 demonstrate that all the considered encoding schemes help gain a considerably high precision. For example, the highest precision is 1.000, obtained for Class 1 by both C<sub>1</sub> and C<sub>2</sub>. The lowest precision is 0.736 by Class 4. Concerning recall, we can see that uni-gram brings the best prediction for Class 1 and 2. With respect to macro-F<sub>1</sub> and balanced accuracy, we also see that C<sub>1</sub> yields the best performance. Altogether, it is evident that using uni-gram as the encoding scheme (C<sub>1</sub>) helps obtain the best classification performance compared to bi-gram (C<sub>2</sub>) and n-gram (C<sub>3</sub>).

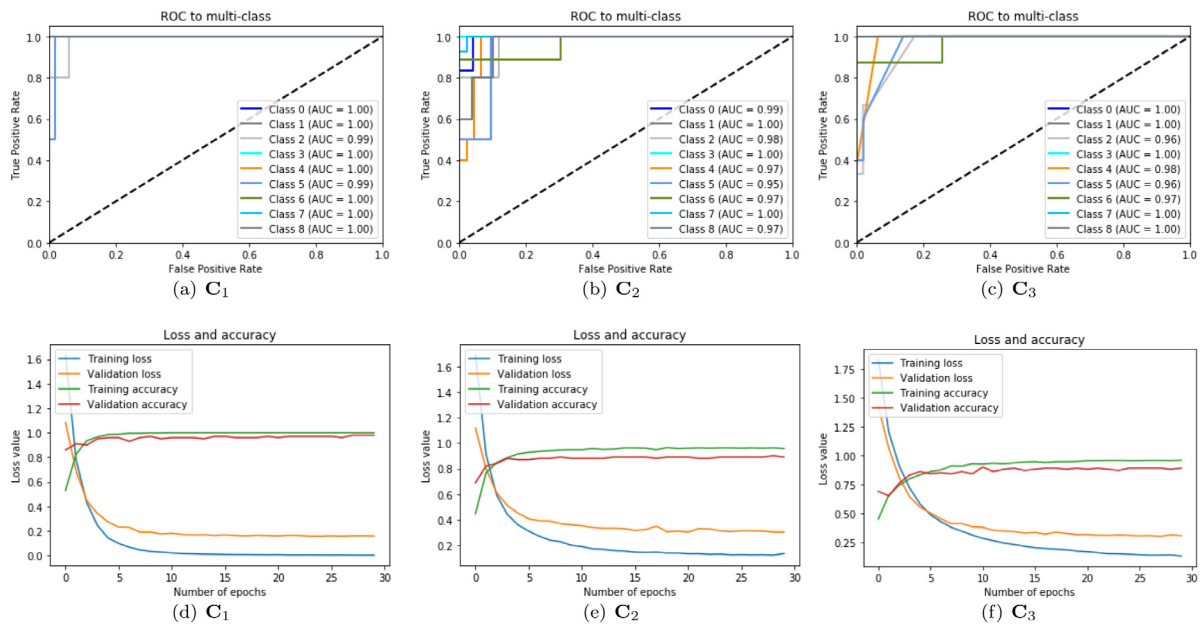
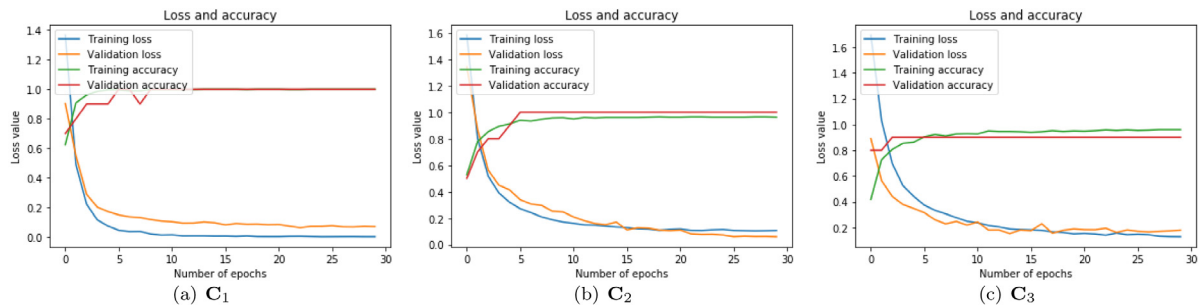
Afterward, we performed additional experiments to study how well MEMOCNN classifies the metamodels into the classes presented in Fig. 8. Furthermore, to better observe the change in performance, the number of epochs was increased to 30.

For this evaluation, ROC and AUC are employed to measure the outcomes as they give a greater insight into each classified category (Davis and Goadrich, 2006). In Fig. 9, 9(b), and 9(c), we depict the ROC curves together with their corresponding AUC values obtained from the experiments for all the considered categories. For uni-gram, we see that most of the ROC curves are close to the upper left corner (0,1), implying that the tool can properly classify the metamodels into their original categories. Similarly, it is evident that using bi-gram and n-gram to encode the metamodels also brings benefits to the classification, though a bit worse compared to using uni-gram. In particular, MEMOCNN gets AUC values ranging from 0.95 to 1.0 for all the categories, indicating a high match rate.

Given an ML model, it is essential to understand whether it is either *underfitting* or *overfitting*, since none of them is useful for the final deployment. A model is considered to be underfitting if

**Table 3**  
Precision, recall, and F<sub>1</sub>-score.

		Class (see Fig. 8 for greater detail)								
	Metric	0	1	2	3	4	5	6	7	8
C <sub>1</sub> (uni-gram)	Precision	0.924	1.000	0.950	0.875	0.904	0.968	0.988	0.937	0.868
	Recall	1.000	1.000	0.780	0.843	0.990	0.857	0.969	0.969	0.901
	F <sub>1</sub>	0.960	1.000	0.860	0.857	0.945	0.909	0.978	0.952	0.884
	macro-F <sub>1</sub>					0.927				
	Accuracy					0.897				
C <sub>2</sub> (bi-gram)	Precision	0.883	1.000	0.923	0.983	0.736	0.936	0.963	0.867	0.885
	Recall	0.852	0.682	0.670	0.921	0.913	0.738	0.827	0.899	0.676
	F <sub>1</sub>	0.867	0.810	0.748	0.951	0.815	0.805	0.889	0.882	0.767
	macro-F <sub>1</sub>					0.842				
	Accuracy					0.782				
C <sub>3</sub> (n-gram)	Precision	0.900	0.958	0.885	0.851	0.837	0.936	0.964	0.891	0.852
	Recall	0.859	0.817	0.791	0.901	0.892	0.822	0.890	0.722	0.663
	F <sub>1</sub>	0.879	0.882	0.835	0.875	0.836	0.875	0.925	0.797	0.740
	macro-F <sub>1</sub>					0.853				
	Accuracy					0.821				

**Fig. 9.** Receiver operating characteristic, loss and accuracy curves (training set = 80%, validation set = 20%).**Fig. 10.** Accuracy and loss (training set = 95%, validation set = 5%).

it fails to learn from data. In contrast, if it learns too well on the training data but performs poorly on the validation data, then it becomes overfitting. To combat overfitting, the common practice is to introduce regularization, for example by adding L1 and L2. However, in most cases, we realized that these methods did not seem to work as expected, e.g., we still witnessed overfitting here and there. We anticipate that this happened due to the small dataset used in our evaluation. It is our assumption that

the countermeasures to mitigate overfitting take effect only when there is a larger number of metamodels.

To investigate if MEMOCNN can be *good fit*, i.e., it is neither underfitting nor overfitting, starting from the data divided by the ten-fold cross-validation technique, we further split each training fold into two parts: 80% for training and 20% for validation. The accuracy and loss recorded during the experiments are shown in Fig. 9(d), 9(e), and 9(f). We see that while the training accuracy

and validation accuracy are almost concordant with each other, the training loss and validation loss are not. A common pattern for the validation plots is that after a dramatic fall, they start to decrease slowly once a certain threshold has been reached, e.g., after three epochs. Furthermore there is always a gap between the plot of validation loss and that of the training loss. In other words, the model appears to be slightly overfitting by the current configuration. We assume that this happens, again, due to the limited amount of data used for training (Rawat and Wang, 2017). Thus, to further examine the model, we increased the amount of data for training by specifying the following ratio: 95% of the data is used for training and only 5% is used for validation. The resulted learning curves for this setting are depicted in Fig. 10, 10(b), and 10(c). In these figures, a clear pattern can be seen: The curves for validation loss and training loss stay close and they are nearly identical. In this respect, we come to the conclusion that with more data used for training, the model becomes less overfitting and almost good fit.

The dataset used in the present paper is considerably small and imbalanced. While we see that some types of neural networks perform well on small dataset, e.g., feed-forward neural networks, we believe that the advantages of a CNN for the classification of metamodels can be better showcased, only if there is a fairly large as well as balanced dataset. In this sense, it is important to experiment on a dataset with more metamodels for each category, and we consider the issue as our future work.

**Answer to RQ<sub>1</sub>:** MEMOCNN effectively classifies the testing metamodels into their original categories. When it comes to a larger dataset for training, the system turns to be good fit.

5.2. RQ<sub>2</sub>: Which network configuration contributes to a better performance?

A CNN can be characterized with different dimensions, as shown in Fig. 6: Depth (D) is the number of layers, including the fully connected ones; Width (W) corresponds to the number of filters (kernels) for each convolutional layer. In recent years, several attempts have been made to increase the prediction performance of a CNN by extending it in different directions (Tan and Le, 2019a).

The research question is particularly relevant for investigating the influence of the network configurations. In particular, we study whether an increase in the number of layers, as well as the number of filters, helps improve the overall accuracy on the given dataset. To this end, four configurations are considered, namely C<sub>1</sub>, C<sub>4</sub>, C<sub>5</sub>, and C<sub>6</sub>, as shown in Table 2. For the sake of simplicity, in these configurations, we used uni-gram as the encoding scheme, as it has demonstrated itself, through RQ<sub>1</sub>, to obtain the best performance among others. As mentioned before, C<sub>1</sub> consists of only one convolutional layer. By C<sub>4</sub>, there are two convolutional layers, and the width of the first and the second layer is 16 and 32, respectively. Compared to C<sub>4</sub>, C<sub>5</sub> has one more layer with width 64. Finally, C<sub>6</sub> is one layer deeper than C<sub>5</sub>, and it is also wider by all the convolutional layers. By the first layer, there are 32 filters, while the second and the third layer consist of 64 filters. For the 2nd fully-connected layer, we used twice the number of nodes compared to C<sub>5</sub>, i.e., 256 instead of 128 nodes. Furthermore, the size of the filters is also larger than that of the other configurations: C<sub>6</sub> employs filters of size  $5 \times 5$ , instead of size  $3 \times 3$  as by C<sub>4</sub> and C<sub>5</sub>. Altogether, it is evident that the network by C<sub>6</sub> is not only *wider* and *deeper*, but also *bigger* (in filters' size) in comparison with the other experimental settings.

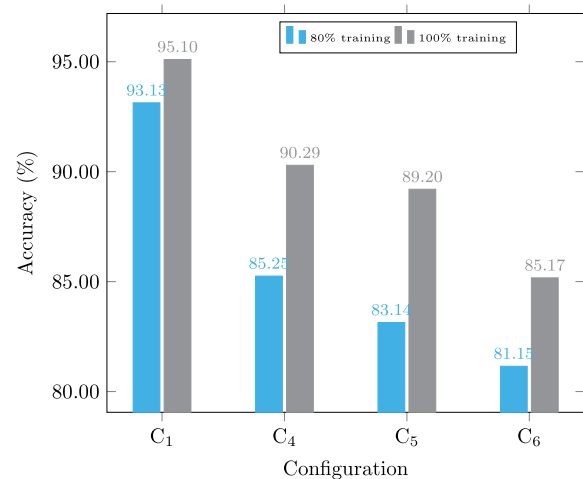


Fig. 11. Accuracy for C<sub>1</sub>, C<sub>4</sub>, C<sub>5</sub>, and C<sub>6</sub>.

To further study RQ<sub>1</sub>, different amounts of data are considered. First, we used 80% of the original training data for training and 20% for validation. Second, 100% of each training fold is used for training, no data is exploited for validation.

The final accuracies we got for the considered configurations are shown in Fig. 11. Overall, the figure suggests that using a more complex network causes a reduction in accuracy. For example, with C<sub>4</sub>, the accuracy is 90.29%, which is higher than that of C<sub>5</sub> and C<sub>6</sub>, i.e., 89.20%, and 85.17%, respectively. This is interesting since we expect that a broader and deeper network would bring in a performance gain (Tan and Le, 2019a). However, the results in this experiment clearly show that the converse is true. We are going to give some explanations for such an outcome in Section 6.1.

The results reported in Fig. 11 further support our findings in RQ<sub>1</sub>: MEMOCNN is more effective when it is fed with more training data. In Fig. 11, given each configuration among C<sub>1</sub>, C<sub>4</sub>, C<sub>5</sub>, and C<sub>6</sub>, when we compare the first column with the second one, we see that using more data to feed MEMOCNN induces a higher accuracy. Take as an example, for C<sub>5</sub>, the accuracy increases from 83.14% to 89.20% when more metamodels are used to train MEMOCNN. Furthermore, we spot a remarkable pattern: complex network structures favor more data. In particular, when we move from the left to the right side of Fig. 11, the gap between the performance of the two bars representing different levels of input data becomes more significant. For instance, by C<sub>1</sub>, the accuracies are 93.13% and 95.10%, meanwhile by C<sub>6</sub>, the corresponding scores are 81.15% and 85.17%. This indicates that deeper and broader networks suffer a setback when they are fed with a small amount of data. Altogether, we conclude that complex networks are supposed to achieve a better performance if we incorporate more data for the training process.

**Answer to RQ<sub>2</sub>:** An increase either in the network's width or depth reduces MEMOCNN's prediction accuracy on the given dataset. Given that a small amount of data is used to train complex networks, such a reduction becomes more evident.

### 5.3. RQ<sub>3</sub>: How does MEMOCNN compare with AURORA?

To aim for a reliable comparison with AURORA, we exploited the replication package which is already available online.<sup>8</sup> Furthermore, we ran the experiments with AURORA using the settings presented in the original paper (Nguyen et al., 2019) that yield the best performance. For MEMOCNN, three configurations have been employed, i.e., C<sub>1</sub>, C<sub>2</sub>, and C<sub>3</sub> as specified in Table 2. Moreover, the ten-fold cross-validation technique is applied on both systems: the dataset is divided into ten equal parts, and the validation is done ten times: each time one fold is used for testing, while the other nine folds are used for training.

The accuracy as well as the precision, recall, and macro-F<sub>1</sub> scores obtained by both systems for different encoding schemes, i.e., uni-gram, bi-gram, and n-gram are shown in Table 4. The table demonstrates that the baseline achieves a good performance as it has been previously claimed (Nguyen et al., 2019). In particular, AURORA has an accuracy being larger than 82% for all encoding schemes. Furthermore, our experiments also confirm that using uni-gram to feed AURORA yields the best accuracy, compared to using bi-gram and n-gram. In comparison with the baseline, MEMOCNN achieves a better accuracy with respect to all three encoding schemes. When encoding with uni-gram, our proposed approach gets an accuracy of 93.72% while the corresponding score of AURORA is 92.25%. Similarly, for other encodings, MEMOCNN outperforms AURORA in terms of accuracy.

We consider the precision, recall, macro-F<sub>1</sub> scores by both AURORA and MEMOCNN as shown in Table 4. With C<sub>1</sub>, MEMOCNN can provide better precision and recall, for example MEMOCNN gets a recall of 0.95 compared to 0.89 by AURORA. This also applies to the outcomes obtained by using bi-gram to provide input. As a consequence, our proposed approach gets a better sensitivity than AURORA does. The final macro-F<sub>1</sub> scores demonstrate that MEMOCNN achieves superior performance than that of the baseline. By C<sub>3</sub>, MEMOCNN outperforms AURORA by all metrics. For example, MEMOCNN gets 0.87, 0.84, and 0.86 as precision, recall, and macro-F<sub>1</sub>, respectively; while the corresponding scores by AURORA are 0.75, 0.62, and 0.69, respectively.

The results obtained by our approach appear to be encouraging at first sight. However, by carefully examining the outcomes presented in Table 4, we noticed that the difference in performance between the two approaches is not significantly big. For instance, with accuracy, in comparison with AURORA, MEMOCNN gets more than 1% and less than 2% performance gain typically, e.g., 82.23% compared to 81.20% (with n-gram), or 85.11% compared to 83.12% (with bi-gram). The same trend can also be witnessed with other quality indexes, i.e., precision, recall, and macro-F<sub>1</sub>. It is necessary to ascertain the cause of this outcome, and thus we came across the following question: “Is MEMOCNN only capable of narrowly surpassing AURORA, or is it because of the dataset?” If the former were the case, then the novelty of MEMOCNN would become questionable, as no-one is willing to sacrifice the simplicity and compactness of AURORA, as shown in Section 2.3, to get a moderate gain in performance. However, our intuition is that the outcome happens due to the considerably small dataset, and thus there is limited data available for training. In fact, a CNN generally consumes a large amount of data to tune its internal parameters (Driss et al., 2017; Yamashita et al., 2018). Furthermore, as we have shown in RQ<sub>1</sub> and RQ<sub>2</sub>, MEMOCNN favors a larger dataset for obtaining a better performance.

Thus, to validate our hypothesis, we conducted additional experiments by specifying different ratios of testing to training data through hold-out validation (Borovicka et al., 2012). Moreover, to simplify the comparison, we exploited only uni-gram

to provide input for both systems since the encoding facilitates the best performance among the other encoding schemes as we already showed by means of RQ<sub>1</sub> and RQ<sub>2</sub>. In particular, we increased by 10% for training data per each experiment, resulting in five different settings, as given in Table 5 which reports the results we got for different hold-out settings. While in general, AURORA gains an increase in performance, albeit small, when we change the amount of training data, it is clear that MEMOCNN outperforms AURORA by all the considered configurations. The most significant difference in accuracy is seen when we set 40% of the dataset as testing and the remaining 60% as training, i.e., AURORA gets 84.23%, and MEMOCNN gets 91.38% correctly classified metamodels. Furthermore, MEMOCNN obtains better accuracy if there is more data available for training. Starting from 89.22%, the accuracy gradually grows up to 96.10% when 10% of the data is used for testing, and 90% is used for training.

By measuring the execution time for both systems, we realized that MEMOCNN is more efficient in timing compared AURORA, i.e., 82 s to 90 s for the same unit of input data. However, this is only a marginal difference, and this happens due to the fact that we have a limited amount of data. Generally, feed-forward neural networks will be more time consuming if we run them on a larger dataset, i.e., containing more metamodels. This is the reason why no significant difference between the timing performance of AURORA and MEMOCNN is seen. Meanwhile, CNNs can maintain a linear increase in timing for a denser dataset, thanks to its convolution operations, which attempt to reduce the number of pairwise multiplications. Thus, we suppose that the difference in timing will be more evident given a dataset with a large number of metamodels.

**Answer to RQ<sub>3</sub>:** MEMOCNN outperforms AURORA with respect to accuracy, precision, recall, macro F<sub>1</sub>, and timing. Furthermore, MEMOCNN improves its prediction performance when being trained with more data.

## 6. Discussions

The lessons learned from the experiments are discussed together with future developments in Section 6.1. The next section discusses the threats to validity.

### 6.1. Lessons learned and provisions

The experience of preparing the experimental evaluation and the corresponding results provided us with invaluable insights. The most remarkable lessons include:

▷ *Data is the cure.* By considering RQ<sub>1</sub> and RQ<sub>2</sub> together, we realized that data augmentation helps tackle two different issues: low accuracy and overfitting. MEMOCNN is less overfitting and more accurate if being fed with more training data. Our work also confirms that adding data is an effective measure to combat overfitting (Yamashita et al., 2018).

▷ *Deeper and wider networks can be useful down the road.* The results in RQ<sub>2</sub> show that on the considered dataset (see Section 4.2), MEMOCNN configurations related to complex networks do not bring any improvements in performance. This lack of progress is understandable at the light of the following arguments. Deepening allows a network to approximate the target function with increased nonlinearity, thereby getting better feature representations. Nevertheless, it also adds up to the complexity of the network, which makes the network more difficult to optimize and susceptible to overfitting. For example, given a training meta-model, the model memorizes well the target class, and thus it fails

<sup>8</sup> <https://github.com/MDEGroup/AURORA>



**Table 4**  
Precision, recall, F<sub>1</sub>-score, and accuracy.

	System	Configuration		
		C <sub>1</sub>	C <sub>2</sub>	C <sub>3</sub>
Precision	AURORA	0.92	0.84	0.75
	MEMOCNN	<b>0.95</b>	<b>0.88</b>	<b>0.87</b>
Recall	AURORA	0.89	0.75	0.62
	MEMOCNN	<b>0.95</b>	<b>0.85</b>	<b>0.84</b>
macro-F <sub>1</sub>	AURORA	0.92	0.81	0.69
	MEMOCNN	<b>0.95</b>	<b>0.86</b>	<b>0.86</b>
Accuracy (%)	AURORA	92.25	83.12	81.20
	MEMOCNN	<b>93.72</b>	<b>85.11</b>	<b>82.23</b>

**Table 5**  
Accuracy for various hold-out settings (%).

System	Testing-Training				
	50%–50%	40%–60%	30%–70%	20%–80%	10%–90%
AURORA	85.20	84.23	89.25	89.14	92.52
MEMOCNN	<b>89.32</b>	<b>91.38</b>	<b>92.15</b>	<b>94.50</b>	<b>96.10</b>

to generalize from data. So a question arises at the moment: “Are deeper and wider networks not useful at all?” The experimental results in Fig. 11 give evidence that complex networks suffer a setback given a small amount of data used for training. Moreover, as it has been found out by related studies (Driss et al., 2017; Yamashita et al., 2018), CNNs favor more data due to the massive number of parameters. This capability essentially means that deeper/wider network structures can be beneficial under certain circumstances. In summary, the answer to the above question is: “Yes, deeper and wider networks are useful to a large dataset for enhancing prediction accuracy.”

▷ *Combination of AURORA and MEMOCNN.* To the best of our knowledge, the proposed approach is the first attempt that represents metamodels using matrices and exploits a convolutional neural network to classify them. Though MEMOCNN obtains a better performance compared to AURORA, we do not undermine the elegance of the baseline. The tool is suitable for working with a small dataset, and it obtains an excellent performance. In this respect, we suppose that a combination of both AURORA and MEMOCNN would be a complementary system.

Apart from the lessons, based on the current results, we envisage that there are the following future developments.

▷ *Application to classify documents.* The approach proposed by Kim (2014) works on sentences; however it may not apply to long documents since the representation matrix becomes enormous, which in turn places a burden on the computational performance. The fact that MEMOCNN performs well on metamodels implies that we can use it to classify documents. First, uni-gram, bi-gram, and n-gram are used to parse documents, as it is done with metamodels, then MEMOCNN can be used to classify the metadata. In a broader sense, our proposed approach can also be applied to group Stack Overflow posts into categories.

▷ *Incorporating more input dimensions.* During the evaluation, we made use of one-channel input tensors. However, we could think of extending it to more channels by incorporating various features. Referring to the Computer Vision domain, images are typically represented using three color channels, namely R, G, B (see Fig. 2). We anticipate the usage of a 3D input tensor by making use of the three encoding schemes together, i.e., uni-gram, bi-gram, and n-gram. The final aim is to better capture the features of a metamodel by taking into consideration different representation schemes.

▷ *Transfer learning.* For CNNs, it is crucial to collect enough data for training. However, labeled data can only be obtained through a manual process, which is both time-consuming and prone to error. To this end, transfer learning has been devised as a practical

solution to export the knowledge extracted from a mature source domain to a novice target domain (Weiss et al., 2016). Transfer learning facilitates the re-use of convolution weights from a model trained on large datasets for training new models with a limited amount of labeled items. Compared to using only random weights (as done with MEMOCNN in this paper), transfer learning is still useful even when the target domain is quite different from the original one where the weights are obtained (Huang et al., 2017). We suppose that the application of transfer learning to MEMOCNN is beneficial to the final classification since our representation scheme allows for the transformation of a metamodel into a picture-like format. This hypothesis needs to be properly validated, and we consider it as future work. Furthermore, once we have a larger dataset, we would be interested in embedding an auto-encoder, as this would possibly enhance the classification performance.

## 6.2. Threats to validity

In this section a discussion of validity threats, which might harm the experimental evaluation is given.

**Internal validity.** Such threats are the internal factors that may affect the outcomes. A probable threat that could be seen through the results obtained for the categories with a considerably low number of items for training and testing. While the threat is partly mitigated by some of the categories in the given dataset as they contain a significant number of metamodels, we believe that the threat can be completely eliminated only if there is a fairly large as well as balanced dataset. In this sense, it is necessary to consider a larger dataset with more metamodels for each category, and we consider the issue as our future work. Another threat is related to the coverage of metamodel constraints. The definition of a metamodel is often endowed with constraints that impose additional properties to the conforming models. In order to consider constraints in the proposed classification technique, the metamodel encoding should encompass constraints as well. However, while it should not pose too many difficulties, it must be investigated whether (i) it is possible to distinguish, for instance, two versions of the same metamodel, each associated with a different constraint, and (ii) it makes sense at all because a constraint might not have enough informative content to the end of the classification.

**External validity.** The main threat to external validity pertains to the generalizability of our findings, i.e., if they are still valid outside the scope of this study. We confronted the issue by evaluating AURORA and MEMOCNN using different validation scenarios,

aiming to simulate an actual usage of the systems. Moreover, we also performed different trials of the same experiments to examine if there is any significant difference among the outcomes obtained by each trial. Overall, we have seen that the obtained results are consistent with those presented in this paper.

**Construct validity.** The threats are related to the simulated setting used to evaluate the tools. The evaluation has been conducted on a training set and a test set, modeling a real application scenario. Depending on the purpose, we used two different settings, namely hold-out validation and ten-fold cross-validation. Furthermore, to aim for a reliable comparison of MEMOCNN with AURORA, we made use of the original implementation of AURORA without touching its internal design.

**Conclusion validity.** Validity threats may be related to whether the experiment methodology is related to the outcome, or there may also be other parameters that may have influenced it. The evaluation metrics, i.e., accuracy, precision, recall,  $F_1$ -score, ROC, and AUC might pose a threat to conclusion validity. The same metrics were employed to evaluate both systems as a mitigation measure.

## 7. Related work

The primary technical advantages claimed by the approach presented in this paper are related to the adoption of ML techniques in MDE. Thus, we review some of the most notable related work on the adoption of such cognitive approaches in the context of MDE, and the categorization of models. Afterward, we discuss some prominent applications of CNNs, as well as the attempts made to improve CNNs' prediction performance by dimension extension and transfer learning.

### 7.1. Machine learning in MDE

The MDE community has made considerable progress in recent years as regards adopting various ML algorithms to solve different issues (Babur et al., 2018; Burgueño et al., 2019). Kusmenko et al. (2019) implement *MontiAnna*, a holistic deep learning modeling framework where deep neural networks can be designed, trained, and integrated into Component & Connector systems. Afterward, *MontiAnna* has been extended using Reinforcement Learning (RL) (Gatto et al., 2019). By exploiting RL algorithms, two systems (Barriga et al., 2018, 2019) can repair errors in a set of broken models, which have been introduced by conflicting changes. Various studies (Dolques et al., 2010; Burgueño et al., 2019) use ML techniques to automatically infer model transformation rules from sets of source and target models. Meta-learning is a technique that aims at using ML itself to automatically learn the most appropriate algorithms and parameters for an ML problem. Hartmann et al. (2019) demonstrate the parallelism between metamodeling and meta-learning, and they envision a notation for describing meta-learning processes.

Learning algorithms have been successfully applied to solve different problems in the MDE scenario. Nevertheless, the adoption of cognitive techniques in this domain does not seem to have received attention being commensurate with its potential. In this respect, MEMOCNN proved to advance the existing approaches in the classification of metamodels.

### 7.2. Categorization of metamodels

Classification and clustering techniques have been employed to categorize metamodels for a diversity of purposes, not excluded reuse (Basciani et al., 2016; Babur et al., 2016; Babur and Cleophas, 2017). In particular, an agglomerative hierarchical clustering technique has been employed to organize metamodels

in repositories (Basciani et al., 2016) automatically. Hierarchical clustering techniques are the main engine to perform metamodel comparison, analysis, and visualization (Babur et al., 2016; Babur and Cleophas, 2017). In these approaches, metamodel named elements are converted into lexical terms and represented as a vector space model. This encoding enables the usage of hierarchical clustering algorithms to group similar objects and represent it as a dendrogram. The approaches mentioned above (Babur et al., 2016; Babur and Cleophas, 2017) share the same intent of this paper even though they are based on unsupervised learning techniques, whereas our work follows the supervised learning strategy. Lopez et al. (2002) present a domain-oriented approach for software requirements reuse. First, requirements captured from semiformal diagrams are injected into models, which are then analyzed to check for any probable quality issues. Finally, a mechanism to cluster the requirements is applied to allow the domain requirements patterns to be identified and reused. *Reuser* (Robinson and Woo, 2004) is a tool used to automatically retrieve related UML artifacts and propose them to the modeler. A search algorithm is applied to find similar artifacts by classifying them into a concept lattice. Then it compares the input artifact query against the concept lattice, to retrieve the closest set of matching artifacts.

Jiang et al. (2004) discuss different characteristics of UML metamodel extension mechanisms according to a four-level classification. Each level of metamodel extension has different features in (contrasting) aspects such as readability, expression capability, use scope, and tool support. The paper aims at providing modelers with a reliable theoretical base for selecting the right level to extend the metamodel to find the right trade-off of the above aspects. A classification model for UML stereotypes has been developed (Berner et al., 1999), where the artifacts are analyzed according to their potential to alter the syntax and semantics of the base language to be able to control their application practice. Feature-based criteria for classifying approaches according to the type of spatial relation involved have been widely discussed (Bottoni and Grau, 2004).

UML class diagrams are commonly used to specify the designs of systems, which then can be used to guide the construction of software (Osman et al., 2018). There are two main types of UML diagrams: (i) FwCD diagrams are hand-made as part of the forward-looking development process; (ii) RECD diagrams are reverse engineered from the source code. Osman et al. (2018) propose a supervised machine learning algorithm to build an automated classifier for classifying UML diagrams. The approach can be used to classify a diagram into either an FwCD or an RECD, and it has been evaluated using a test set of 999 class diagrams obtained from open source projects. By running experiments with different machine learning algorithms, the authors concluded that Random Forest is the most suitable algorithm for their classification tasks.

To our knowledge, AURORA (Nguyen et al., 2019) is the first effort to classify metamodels using a feed-forward neural network. Despite its simplicity, the tool is efficient, and on a considerably small dataset, it classifies the metamodels, obtaining high accuracy. In this work, we improved AURORA by employing a CNN, which is supposed to learn better on a large dataset.

### 7.3. Application and extension of CNNs

CNNs were originally designed to work with images, and thus they have been widely used to solve various recognition tasks. To name but a few, CNNs have been applied to detect fruits (Duong et al., 2020; Gatica et al., 2013), to classify images (Krizhevsky et al., 2012; Rawat and Wang, 2017), allowing for the transfer of well pre-trained hyper-parameters to a new application. The

success of CNNs in Computer Vision is a source of inspiration for applications in other domains. For instance, CNNs have been exploited to classify DNA sequences (Aoki and Sakakibara, 2018), sentences (Kim, 2014), or to process natural language (Belinkov and Glass, 2019). To the best of our knowledge MEMOCNN is the first attempt to bring CNNs into a new domain, i.e., MDE (Schmidt, 2006). By developing MEMOCNN on top of a convolutional neural network, we inherit a well-founded background of knowledge to solve the problem of metamodel classification. This also paves the way for deploying transfer learning, which allows for increasing both efficiency and effectiveness.

As seen in Fig. 6, a CNN can be characterized with width, depth, resolution, and filter size. So far, various approaches have been proposed to boost up the performance of a CNN by extending it in these dimensions. Two of the most recent studies of this type are EfficientNet and MixNet. The EfficientNet family (Tan and Le, 2019a) improves performance by extending the original network in depth, width, and resolution. It decreases the learning rate of the network, given that the number of layers reaches a high value. By relying on width scaling, it is possible to capture features more definitely, easing the training process. However, the scalability of the approach is questionable because of reduced precision over high-level features. The memory cost represents the main hindrance as the network achieves a higher accuracy with a high-quality resolution. Like this, the compound scaling method does not expose the limitations of traditional techniques that tend to act on a single dimension at once. MixNet (Tan and Le, 2019b) is a family of neural networks that combines multiple kernel sizes in a single convolution, thus allowing for recognizing various types of patterns from input images. MixNet has demonstrated itself to outperform some state-of-the-art CNNs significantly. As we already pointed out in RQ<sub>3</sub>, and Section 6.1, the expansion of MEMOCNN in different dimensions may bring benefits. However, the performance depends much on the input data: given a small amount of training data, enlarging the network turns out to be counterproductive.

#### 7.4. Transfer learning

The technique has been widely used by imaging applications to improve the learning performance of CNNs (Weiss et al., 2016; Huang et al., 2017). Several studies adopt weights pre-trained using the ImageNet dataset, which consists of a huge number of images covering various categories (Russakovsky et al., 2015). Some notable studies that make use of ImageNet are: the work by Razavian et al. to classify images (Razavian et al., 2014), or the one by Karpathy and Fei-Fei (2017) to generate natural language descriptions for images. The accuracy obtained by employing ImageNet is better than that by using random learning. We assume that transfer learning would be beneficial to the training of MEMOCNN. Compared to using only random weights, as we have done with MEMOCNN in this paper, transfer learning is still helpful, also when the target domain is quite different from the original one where the weights have been obtained (Huang et al., 2017).

## 8. Conclusions and future work

MEMOCNN has been proposed as a novel approach to the classification of metamodels by adapting a deep learning technique that has gained considerable traction in other domains. The comparison with AURORA, a state-of-the-art baseline, on a manually classified dataset demonstrates that MEMOCNN obtains a better performance with respect to different quality metrics. Our future plan is to investigate whether transfer learning is beneficial to MEMOCNN by adopting pre-trained weights using

the ImageNet dataset. Moreover, we believe that the performance of memoCNN can be better judged when there is more data for training. Thus, we will invest time and efforts to populate a larger dataset, aiming to empower the classification capability. We plan to collect metamodels from GitHub using the dedicated API, and then manually classify them by involving expert modelers to build a training set. Last but not least, we are going to deploy the framework to other domains, such as the classification of documents and Stack Overflow posts to support software development activities.

## CRedit authorship contribution statement

**Phuong T. Nguyen:** Conceptualization, Methodology, Software, Writing - review & editing. **Davide Di Ruscio:** Conceptualization, Writing - original draft, Writing - review & editing, Supervision. **Alfonso Pierantonio:** Conceptualization, Writing - original draft, Writing - review & editing, Supervision. **Juri Di Rocco:** Writing - review & editing, Data curation, Visualization. **Ludovico Iovino:** Writing - review & editing, Validation, Visualization.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgments

The research described in this paper has been carried out as part of the CROSSMINER Project, which has received funding from the European Union's Horizon 2020 Research and Innovation Programme under Grant 732223. We thank the anonymous reviewers for their valuable comments and suggestions that helped us improve the manuscript.

## References

- Aly, M., 2005. Survey on multiclass classification methods. *Neural Netw.* 1–9.
- Aoki, G., Sakakibara, Y., 2018. Convolutional neural networks for classification of alignments of non-coding RNA sequences. *Bioinformatics* 34 (13), i237–i244. <http://dx.doi.org/10.1093/bioinformatics/bty228>, URL <https://academic.oup.com/bioinformatics/article-pdf/34/13/i237/25098204/bty228.pdf>.
- Babur, Ö., 2019. A labeled ecore metamodel dataset for domain clustering. <http://dx.doi.org/10.5281/zenodo.2585456>.
- Babur, Ö., Chaudron, M.R., Cleophas, L., Di Ruscio, D., Kolovos, D., 2018. AMMoRe 2018: First international workshop on analytics and mining of model repositories. In: *CEUR Workshop Proceedings*, vol. 2245, pp. 778–779.
- Babur, Ö., Cleophas, L., 2017. Using n-grams for the automated clustering of structural models. In: Steffen, B., Baier, C., van den Brand, M., Eder, J., Hinchey, M., Margaria, T. (Eds.), *SOFSEM 2017: Theory and Practice of Computer Science*. Springer International Publishing, Cham, pp. 510–524.
- Babur, Ö., Cleophas, L., van den Brand, M., 2016. Hierarchical clustering of meta-models for comparative analysis and visualization. In: Wąsowski, A., Lönn, H. (Eds.), *Modelling Foundations and Applications*. Springer International Publishing, Cham, pp. 3–18.
- Barriga, A., Rutle, A., Heldal, R., 2018. Automatic model repair using reinforcement learning. In: *Proceedings of Workshops Co-Located with MODELS 2018*. Copenhagen, Denmark, October, 14, 2018, pp. 781–786.
- Barriga, A., Rutle, A., Heldal, R., 2019. Personalized and automatic model repairing using reinforcement learning. In: *2019 ACM/IEEE 22nd International Conference on Model Driven Engineering Languages and Systems Companion. MODELS-C. IEEE*, pp. 175–181.
- Basciani, F., Di Rocco, J., Di Ruscio, D., Iovino, L., Pierantonio, A., 2016. Automated clustering of metamodel repositories. In: Nurcan, S., Soffer, P., Bajec, M., Eder, J. (Eds.), *Advanced Information Systems Engineering*. Springer International Publishing, Cham, pp. 342–358.
- Belinkov, Y., Glass, J., 2019. Analysis methods in neural language processing: A survey. *Trans. Assoc. Comput. Linguist.* 7, 49–72. [http://dx.doi.org/10.1162/tacl\\_a\\_00254](http://dx.doi.org/10.1162/tacl_a_00254), URL <https://www.aclweb.org/anthology/Q19-1004>.



- Berner, S., Glinz, M., Joos, S., 1999. A classification of stereotypes for object-oriented modeling languages. In: France, R., Rumpe, B. (Eds.), *UML'99 – The Unified Modeling Language*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 249–264.
- Borovicka, T., Jirina, Jr., M., Jirina, M., Kordik, P., 2012. *Selecting Representative Data Sets*. INTECH Open Access Publisher.
- Bottoni, P., Grau, A., 2004. A suite of metamodels as a basis for a classification of visual languages. In: 2004 IEEE Symposium on Visual Languages – Human Centric Computing. pp. 83–90. <http://dx.doi.org/10.1109/VLHCC.2004.5>.
- Burgueño, L., Burdusel, A., Gérard, S., Wimmer, M., 2019. MDE Intelligence19: First International Workshop on Artificial Intelligence and Model-Driven Engineering.
- Burgueño, L., Cabot, J., Gérard, S., 2019. An LSTM-based neural network architecture for model transformations. In: 2019 ACM/IEEE 22nd International Conference on Model Driven Engineering Languages and Systems. MODELS. pp. 294–299. <http://dx.doi.org/10.1109/MODELS.2019.00013>.
- Ceci, M., Esposito, F., Lapi, M., Malerba, D., 2003. Automated classification of web documents into a hierarchy of categories. In: *Intelligent Information Processing and Web Mining*. Springer, pp. 59–68.
- Davis, J., Goadrich, M., 2006. The relationship between precision-recall and ROC curves. In: *ICML '06: Proceedings of the 23rd International Conference on Machine Learning*. ACM, New York, NY, USA, pp. 233–240. <http://dx.doi.org/10.1145/1143844.1143874>, URL <http://portal.acm.org/citation.cfm?id=1143874>.
- Di Rocco, J., Di Ruscio, D., Iovino, L., Pierantonio, A., 2015. Collaborative repositories in model-driven engineering [software technology]. *IEEE Softw.* 32 (3), 28–34.
- Dolques, X., Huchard, M., Nebut, C., Reitz, P., 2010. Learning transformation rules from transformation examples: An approach based on relational concept analysis. In: 2010 14th IEEE International Enterprise Distributed Object Computing Conference Workshops. pp. 27–32.
- Domingos, P., 2012. A few useful things to know about machine learning. *Commun. ACM* 55 (10), 78–87. <http://dx.doi.org/10.1145/2347736.2347755>.
- Driss, S.B., Soua, M., Kachouri, R., Akil, M., 2017. A comparison study between MLP and convolutional neural network models for character recognition. In: Kehtarnavaz, N., Carlssohn, M.F. (Eds.), *Real-Time Image and Video Processing 2017*, Vol. 10223. International Society for Optics and Photonics, SPIE, pp. 32–42. <http://dx.doi.org/10.1117/12.2262589>.
- Duong, L.T., Nguyen, P.T., Di Sipio, C., Di Ruscio, D., 2020. Automated fruit recognition using EfficientNet and MixNet. *Comput. Electron. Agric.* 171, 105326. <http://dx.doi.org/10.1016/j.compag.2020.105326>, URL <http://www.sciencedirect.com/science/article/pii/S1068169919319787>.
- Fawcett, T., 2006. An introduction to ROC analysis. *Pattern Recognit. Lett.* 27 (8), 861–874. <http://dx.doi.org/10.1016/j.patrec.2005.10.010>.
- Gatica, G., Best, S., Ceroni, J., Lefranc, G., 2013. Olive fruits recognition using neural networks. In: *First International Conference on Information Technology and Quantitative Management*. Procedia Comput. Sci. 17, 412–419. <http://dx.doi.org/10.1016/j.procs.2013.05.053>, URL <http://www.sciencedirect.com/science/article/pii/S1877050913001865>.
- Gatto, N., Kusmenko, E., Rumpe, B., 2019. Modeling deep reinforcement learning based architectures for cyber-physical systems. In: 2019 ACM/IEEE 22nd International Conference on Model Driven Engineering Languages and Systems Companion. MODELS-C. IEEE, pp. 196–202.
- Goodfellow, I., Bengio, Y., Courville, A., 2016. *Deep Learning*. The MIT Press.
- Hartmann, J., Huppertz, J., Schamp, C., Heitmann, M., 2019. Comparing automated text classification methods. *Int. J. Res. Mark.* 36 (1), 20–38.
- Hartmann, A., Lienhart, R., 2002. Automatic classification of images on the web. 4676, pp. 31–40. <http://dx.doi.org/10.1117/12.451108>.
- Hartmann, T., Moawad, A., Schockaert, C., Fouquet, F., Le Traon, Y., 2019. Meta-modelling meta-learning. In: 2019 ACM/IEEE 22nd International Conference on Model Driven Engineering Languages and Systems. MODELS. pp. 300–305. <http://dx.doi.org/10.1109/MODELS.2019.00014>.
- He, K., Zhang, X., Ren, S., Sun, J., 2015. Deep residual learning for image recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition. CVPR. pp. 770–778.
- Hein, C., Ritter, T., Wagner, M., 2009. Model-driven tool integration with modelbus. In: *Workshop Future Trends of Model-Driven Development*. pp. 50–52.
- Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural Comput.* 9 (8), 1735–1780. <http://dx.doi.org/10.1162/neco.1997.9.8.1735>.
- Huang, Z., Pan, Z., Lei, B., 2017. Transfer learning with deep convolutional neural network for SAR target classification with limited labeled data. *Remote Sens.* 9 (9), <http://dx.doi.org/10.3390/rs9090907>, URL <https://www.mdpi.com/2072-4292/9/9/907>.
- Jiang, Y., Shao, W., Zhang, L., Ma, Z., Meng, X., Ma, H., 2004. On the classification of uml's meta model extension mechanism. In: Baar, T., Strohmeier, A., Moreira, A., Mellor, S.J. (Eds.), *The Unified Modeling Language. Modeling Languages and Applications*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 54–68.
- Juszcak, P., Tax, D., Duin, R.P., 2002. Feature scaling in support vector data description. In: *Proc. Asci. Citeseer*, pp. 95–102.
- Karasneh, B., Chaudron, M.R., 2013. Online Img2UML repository: An online repository for UML Models. In: *EESMOD@ MoDELS*. pp. 61–66.
- Karpathy, A., Fei-Fei, L., 2017. Deep visual-semantic alignments for generating image descriptions. *IEEE Trans. Pattern Anal. Mach. Intell.* 39 (4), 664–676. <http://dx.doi.org/10.1109/TPAMI.2016.2598339>.
- Kim, Y., 2014. Convolutional neural networks for sentence classification. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*. EMNLP. Association for Computational Linguistics, Doha, Qatar, pp. 1746–1751. <http://dx.doi.org/10.3115/v1/D14-1181>, URL <https://www.aclweb.org/anthology/D14-1181>.
- Kingma, D.P., Ba, J., 2015. Adam: A method for stochastic optimization. In: Bengio, Y., LeCun, Y. (Eds.), *3rd International Conference on Learning Representations*. ICLR 2015. San Diego, CA, USA, May 7–9, 2015, Conference Track Proceedings. <http://arxiv.org/abs/1412.6980>.
- Koegel, M., Helmig, J., 2010. EMFStore: A model repository for EMF models. In: 2010 ACM/IEEE 32nd International Conference on Software Engineering, Vol. 2. IEEE, pp. 307–308.
- Kohavi, R., 1995. A study of cross-validation and bootstrap for accuracy estimation and model selection. In: *14th International Joint Conference on Artificial Intelligence*. Morgan Kaufmann Publishers Inc., San Francisco, pp. 1137–1143.
- Kotsiantis, S.B., Zaharakis, I.D., Pintelas, P.E., 2006. Machine learning: A review of classification and combining techniques. *Artif. Intell. Rev.* 26 (3), 159–190. <http://dx.doi.org/10.1007/s10462-007-9052-3>.
- Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks. In: Pereira, F., Burges, C.J.C., Bottou, L., Weinberger, K.Q. (Eds.), *Advances in Neural Information Processing Systems 25*. Curran Associates, Inc., pp. 1097–1105.
- Kusmenko, E., Nickels, S., Pavlitskaya, S., Rumpe, B., Timmermanns, T., 2019. Modeling and training of neural processing systems. In: 2019 ACM/IEEE 22nd International Conference on Model Driven Engineering Languages and Systems. MODELS. IEEE, pp. 283–293.
- Kutsche, R., Milanovic, N., Bauhoff, G., Baum, T., Carlsburg, M., Kumpe, D., Widiker, J., 2008. Bizycle: Model-based interoperability platform for software and data integration. In: *Proceedings of the MDTPI at ECMDA*, vol. 430.
- Lopes, C.V., Maj, P., Martins, P., Saini, V., Yang, D., Zitny, J., Sajjani, H., Vitek, J., 2017. DéjàVu: A map of code duplicates on GitHub. *Proc. ACM Program. Lang.* 1 (OOPSLA), 1–28.
- Lopez, O., Laguna, M.A., García, F.J., 2002. Reuse based analysis and clustering of requirements diagrams. In: *Pre-Proceedings of the Eighth International Workshop on Requirements Engineering: Foundation for Software Quality*. REFSQ'02, pp. 71–82.
- Nassif, A.B., Shahin, I., Attili, I., Azzeh, M., Shaalan, K., 2019. Speech recognition using deep neural networks: A systematic review. *IEEE Access* 7, 19143–19165.
- Nguyen, P.T., Di Rocco, J., Di Ruscio, D., Pierantonio, A., Iovino, L., 2019. Automated classification of metamodel repositories: A machine learning approach. In: 2019 ACM/IEEE 22nd International Conference on Model Driven Engineering Languages and Systems. MODELS. pp. 272–282. <http://dx.doi.org/10.1109/MODELS.2019.00011>.
- Nguyen, A., Kanoulas, D., Caldwell, D.G., Tsagarakis, N.G., 2016. Detecting object affordances with convolutional neural networks. In: 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems. IROS. pp. 2765–2770. <http://dx.doi.org/10.1109/IROS.2016.7759429>.
- Osman, M.H., Ho-Quang, T., Chaudron, M., 2018. An automated approach for classifying reverse-engineered and forward-engineered UML class diagrams. In: 2018 44th Euromicro Conference on Software Engineering and Advanced Applications. SEAA. pp. 396–399.
- Porter, M.F. et al., 1980. An algorithm for suffix stripping. *Program* 14 (3), 130–137.
- Portugal, I., Alencar, P.S.C., Cowan, D.D., 2015. The use of machine learning algorithms in recommender systems: A systematic review. *CoRR abs/1511.05263*, [arXiv:1511.05263](http://arxiv.org/abs/1511.05263).
- Rawat, W., Wang, Z., 2017. Deep convolutional neural networks for image classification: A comprehensive review. *Neural Comput.* 29 (9), 2352–2449. [http://dx.doi.org/10.1162/neco\\_a.00990](http://dx.doi.org/10.1162/neco_a.00990).
- Razavian, A.S., Azizpour, H., Sullivan, J., Carlsson, S., 2014. CNN features off-the-shelf: An astounding baseline for recognition. *CoRR abs/1403.6382*, URL <http://dblp.uni-trier.de/db/journals/corr/corr1403.html#RazavianASC14>.
- Robinson, W.N., Woo, H.G., 2004. Finding reusable UML sequence diagrams automatically. *IEEE Softw.* 21 (5), 60–67. <http://dx.doi.org/10.1109/MS.2004.1331304>.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L., 2015. ImageNet large scale visual recognition challenge. *Int. J. Comput. Vision* 115 (3), 211–252. <http://dx.doi.org/10.1007/s11263-015-0816-y>.
- Schmidt, D.C., 2006. Guest editor's introduction: Model-driven engineering. *Computer* 39 (2), 25–31. <http://dx.doi.org/10.1109/MC.2006.58>.
- Sebastiani, F., 2002. Machine learning in automated text categorization. *ACM Comput. Surv.* 34 (1), 1–47. <http://dx.doi.org/10.1145/505282.505283>.
- Shi, Z., He, L., Suzuki, K., Nakamura, T., Itoh, H., 2009. Survey on neural networks used for medical image processing. *Int. J. Comput. Sci.* 3 (1), 86.



- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R., 2014. Dropout: A simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* 15 (1), 1929–1958.
- Sundararajan, K., Woodard, D.L., 2018. Deep learning for biometrics: A survey. *ACM Comput. Surv.* 51 (3), <http://dx.doi.org/10.1145/3190618>.
- Tan, M., Le, Q., 2019a. EfficientNet: Rethinking model scaling for convolutional neural networks. In: Chaudhuri, K., Salakhutdinov, R. (Eds.), *Proceedings of the 36th International Conference on Machine Learning*. In: *Proceedings of Machine Learning Research*, vol. 97, PMLR, Long Beach, California, USA, pp. 6105–6114. URL <https://proceedings.mlr.press/v97/tan19a.html>.
- Tan, M., Le, Q.V., 2019b. MixConv: Mixed depthwise convolutional kernels. *CoRR abs/1907.09595*, [arXiv:1907.09595](https://arxiv.org/abs/1907.09595).
- Weiss, K., Khoshgoftaar, T., Wang, D., 2016. A survey of transfer learning. *J. Big Data* 3, <http://dx.doi.org/10.1186/s40537-016-0043-6>.
- White, M., Vendome, C., Linares-Vásquez, M., Poshyanyk, D., 2015. Toward deep learning software repositories. In: *2015 IEEE/ACM 12th Working Conference on Mining Software Repositories*. IEEE, pp. 334–345.
- Wuest, T., Weimer, D., Irgens, C., Thoben, K.-D., 2016. Machine learning in manufacturing: Advantages, challenges, and applications. *Prod. Manuf. Res.* 4 (1), 23–45. <http://dx.doi.org/10.1080/21693277.2016.1192517>.
- Yamashita, R., Nishio, M., Do, R.K.G., Togashi, K., 2018. Convolutional neural networks: An overview and application in radiology. *Insights Imaging* 9 (4), 611–629. <http://dx.doi.org/10.1007/s13244-018-0639-9>.
- Zhang, G.P., 2000. Neural networks for classification: A survey. *IEEE Trans. Syst. Man Cybern. C* 30 (4), 451–462.
- Zhang, Z., Sabuncu, M., 2018. Generalized cross entropy loss for training deep neural networks with noisy labels. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems* 31. Curran Associates, Inc., pp. 8778–8788.
- Zhao, Z., Zheng, P., Xu, S., Wu, X., 2019. Object detection with deep learning: A review. *IEEE Trans. Neural Netw. Learn. Syst.* 30 (11), 3212–3232. <http://dx.doi.org/10.1109/TNNLS.2018.2876865>.



**Phuong T. Nguyen** is a postdoctoral researcher at Università degli Studi dell'Aquila, Italy. He obtained a Ph.D. in Computer Science from the University of Jena, Germany. Since the graduation, Phuong has worked as a university teaching and research assistant in Vietnam and Italy. His research interests include Computer Networks, Semantic Web, Recommender Systems, and Machine Learning. Recently, he has been working to develop recommender systems in Software Engineering for mining open source code repositories.



**Davide Di Ruscio** is Associate Professor at the DISIM - Università degli Studi dell'Aquila, Italy. His main research interests are related to several aspects of Software Engineering, Open Source Software, and Model Driven Engineering (MDE) including domain specific modeling languages, model transformation, model differencing, and model evolution. He has published more than 130 papers in various journals, conferences and workshops on such topics. He is a member of the steering committee of the International Conference on Model Transformation (ICMT), of the Software

Language Engineering (SLE) conference, of the Seminar Series on Advanced Techniques & Tools for Software Evolution (SATTOSSE), of the Workshop on Modelling in Software Engineering at ICSE (MiSE) and of the International Workshop on Robotics Software Engineering (RoSE). He is in the editorial board of the International Journal on Software and Systems Modeling (SoSyM), of the Journal of Object Technology, and of the IET Software journal. More information is available at <http://people.disim.univaq.it/diruscio>.



**Alfonso Pierantonio** is Professor in Computer Science at Università degli Studi dell'Aquila, Italy. His current research interests include Model-Driven Engineering with a specific emphasis on model management analytics and coevolution. Alfonso is involved in several scientific projects, including EU H2020 Lowcomote and EU H2020 Typhon; moreover, he is principal co-investigator of the Italian Railways funded project ERMES. He has been general chair of STAF 2015, PC chair of ECMFA 2018, and designated general chair of MoDELS 2022 besides being on the organizing, steering, and program committees of many international conferences. Alfonso is the Editor-in-Chief of the Journal of Object Technology (JOT) and a member of the editorial board of the Journal on Software and Systems Modeling (SoSyM). He has published more than 140 papers in international journals, conferences, and workshops.



**Juri Di Rocco** is a postdoctoral researcher at the Department of Information Engineering Computer Science and Mathematics, Università degli Studi dell'Aquila, Italy. Previously, he obtained the Ph.D. degree from Università degli Studi dell'Aquila in the MDEGroup research team with Alfonso Pierantonio and Davide Di Ruscio. He is interested in all aspects of software language engineering. Main research interests are related to several aspects of Model Driven Engineering (MDE) including domain specific modeling languages, model transformation, model differencing, modeling repositories and mining techniques.



**Ludovico Iovino** is currently Assistant Professor at the GSSI - Gran Sasso Science Institute, L'Aquila - in the Computer Science scientific area. His interests include Model Driven Engineering (MDE), Model Transformations, Metamodel Evolution, code generation and software quality evaluation. Currently he is working on model-based artifacts and issues related to the metamodel evolution problem. He has been included in program committees of numerous conferences and in the local organization of the STAF conference 2015 and models and evolution workshop from 2 years. He is part of different academic projects related to Model Repositories, model migration tools and Eclipse Plugins.