



INTERNSHIP REPORT

Meta-learning algorithms for Few-Shot Computer Vision

April - July 2019

Etienne BENNEQUIN

Under the supervision of Antoine Toubhans, PhD



Abstract

Few-Shot Learning is the challenge of training a model with only a small amount of data. Many solutions to this problem use meta-learning algorithms, *i.e.* algorithms that learn to learn. By sampling few-shot tasks from a larger dataset, we can teach these algorithms to solve new, unseen tasks.

This document reports my work on meta-learning algorithms for Few-Shot Computer Vision. This work was done during my internship at Sicara, a French company building image recognition solutions for businesses. It contains:

1. an extensive review of the state-of-the-art in few-shot computer vision;
2. a benchmark of meta-learning algorithms for few-shot image classification;
3. the introduction to a novel meta-learning algorithm for few-shot object detection, which is still in development.

Special thanks

I would like to thank everyone at Sicara for their help on so many levels. In particular, thanks to Antoine Toubhans, Elisa Negra, Laurent Montier, Tanguy Marchand and Theodore Aouad for their daily support.

Contents

1	Introduction	4
2	Review	5
2.1	Few-Shot classification problem	5
2.2	Meta-learning paradigm	7
2.3	Meta-learning algorithms	8
2.3.1	Gradient-based meta-learning	8
2.3.2	Metric Learning	10
2.4	Few-Shot Image classification benchmarks	12
2.5	Few-Shot object detection	12
3	Contributions	14
3.1	Overview	14
3.2	Meta-learning algorithms for Few-Shot image classification	15
3.2.1	Implementation	15
3.2.2	Reproducing the results	18
3.2.3	Effects of label noise in the support set at evaluation time	21
3.2.4	Future work	24
3.3	MAML for Few-Shot object detection	25
3.3.1	The Few-Shot Object Detection problem	25
3.3.2	YOLOMAML	26
3.3.3	Implementation	26
3.3.4	First results and investigations	29
3.3.5	Future work	31
4	Conclusion	32

1 Introduction

In 1980, Kunihiko Fukushima developed the first convolutional neural networks. Since then, thanks to increasing computing capabilities and huge efforts from the machine learning community, deep learning algorithms have never ceased to improve their performances on tasks related to computer vision. In 2015, Kaiming He and his team at Microsoft reported that their model performed better than humans at classifying images from ImageNet [1]. At that point, one could argue that computers became better than people at harnessing billions of images to solve a specific task.

However, in real world applications, it is not always possible to build a dataset with that many images. Sometimes we need to classify images with only one or two examples per class. For this kind of tasks, machine learning algorithms are still far from human performance.

This problem of learning from few examples is called few-shot learning.

For a few years now, the few-shot learning problem has drawn a lot of attention in the research community, and a lot of elegant solutions have been developed. An increasing part of them use meta-learning, which can be defined in this case as *learning to learn*.

During my internship at Sicara, I focused on meta-learning algorithms to solve few-shot computer vision tasks, both for image classification and object detection. I compared the performance of four distinct meta-learning algorithms in few-shot classification tasks. I also started the development of a novel meta-learning model for few-shot object detection.

The first section is an extensive review of state-of-the-art solutions for solving few-shot image classification and few-shot image detection. It starts with the definition of the few-shot learning problem.

Then I will expose my contributions. The first part of it is a benchmark of state-of-the-art algorithms for few-shot image classification on several settings and datasets. The second part introduces the YOLOMAML, a novel solution for few-shot object detection. This algorithm is still in development.

This report shares details about the research process and the implementation of the algorithms and experiments. I hope this information about the issues raised during my work and my attempts at solving them will be useful for anyone who will work on meta-learning algorithms in the future.

2 Review

2.1 Few-Shot classification problem

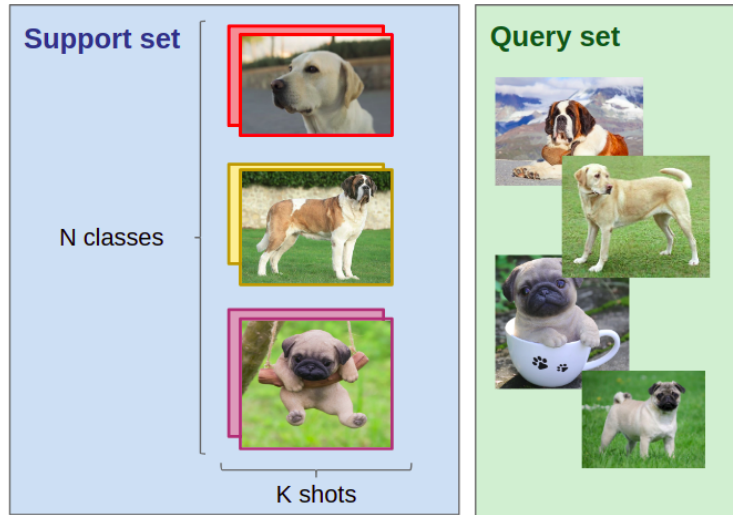


Figure 1: A 3-way 2-shot classification problem. Images from the query set would need to be classified in $\{ \text{Labrador, Saint-Bernard, Pug} \}$.

We define the N -way K -shot image classification problem as follows. Given:

1. a support set composed of:
 - N class labels;
 - For each class label, K labeled images;
2. Q query images;

we want to classify the query images among the N classes. The $N \times K$ images in the support set are the only examples available for these classes.

When K is small (typically $K < 10$), we talk about *few-shot image classification* (or *one-shot* in the case where $K = 1$). The problem in this case is that we fail to provide enough images of each class to solve the classification problem with a standard deep neural network, which usually require thousands of images. Note that this problem is different from semi or weakly supervised learning, since the data is fully labeled. The problem here is not the scarcity of labels, but the scarcity of training data.

A visual example of a few-shot classification problem is shown in Figure 1.

The Few-Shot Learning problem (which includes few-shot image classification) has drawn a lot of attention in the past few years. Many different ways of solving this problem have been imagined. They all have in common that they use additional information from a large *base-dataset*. The classes in the base-dataset are different from the ones in the support set of

the few-shot task we ultimately want to solve. For instance, if the target task is classifying images as Labrador, Saint-Bernard or Pug (Figure 1), the base-dataset can be composed of many other dog breeds. Here I provide an overview of these solutions.

Memory-augmented networks Santoro *et al.* (2016) [2] had the idea that new images from previously unseen classes could be classified by using stored information about previous image classification. Their model uses a Recurrent Neural Networks that learns both how to store and how to retrieve relevant information from past data. Other methods exploit the idea of extending neural networks with external memory [3] [4].

Metric learning Koch *et al.* (2015) [5] proposed the Siamese Neural Networks to solve few-shot image classification. Their model is composed of two convolutional neural networks with shared weights (the *legs*), that compute embeddings (*i.e.* features vectors) for their input images, and one *head* that compares the respective output of each leg. At training time (on the large *base-dataset*), the network receives couples of images as input, predicts whether they belong or not to the same class, and is trained upon the accuracy of this prediction. Ultimately, when evaluated on a few-shot classification class (see Figure 1), each query image is compared to every images in the support set, and is assigned to the class that is considered the closest (using for instance *k*-Nearest Neighbours).

This algorithm achieved interesting results on few-shot image classification. However, the task upon which it was trained (comparison of two images) differed from the task upon which it was evaluated (classification).

Vinyals *et al.* (2016) [6] considered that this was a drawback and proposed a slightly different version of this algorithm, inside of the *meta-learning framework* (see the definition of this framework in section 2.2). Their Matching Networks also classify query images by comparing their embedding to the embeddings computed from support set images, but the difference is that their training objective is image classification as well. They outperform Siamese Networks, thus validating their assumption.

Later works aim at improving this algorithm [7] [8]. They will be presented with more details in section 2.3.2.

Gradient-based meta-learners Other algorithms inside of the meta-learning framework learn an efficient way to fine-tune a convolutional neural network on the support set in order to accurately classify the query set. Finn *et al.* (2017) [9] developed a Model-Agnostic Meta-Learner (MAML) which tries to learn the best parameters for the CNN’s initialization in order to achieve good accuracy on the query set after only a few gradient descents on the support set. The Meta-SGD developed by Li *et al.* (2017) [10] goes further: in addition to the initialization parameters, this algorithm learns for each parameter a learning rate and an update direction. Ravi & Larochelle (2016) [11] proposed a Long-Short-Term-Memory network where the cell state (*i.e.* the variable supposed to carry long-term memory in a LSTM) is the parameters of the CNN. This allows to execute a *learned gradient descent*, where all the hyper-parameters of the CNN’s training are actually trained parameters of the LSTM.

Still inside the meta-learning framework, which they considered as a sequence-to-sequence problem, Mishra *et al.* (2018 [12]) combine temporal convolutions with causal attention to create their Simple Neural AttentIve Learner (SNAIL). Finally, Garcia & Bruna [13] proposed to use graph neural networks as an extension of all meta-learning algorithms for few-shot learning.

Data generation An other option to solve the problem of having too few examples to learn from is to generate additional data. Hariharan & Girshick (2017) [14] augmented metric learning algorithm with hallucinated feature vectors which were added to the feature vectors extracted from real images. Antoniou *et al.* (2017) [15] applied Generative Adversarial Networks to Few-Shot data augmentation: their GAN are able to take as input an image from a previously unseen class to generate new images which belong in the same class. Wang *et al.* (2018) [16] proposed a meta-learned imaginary data generator which can be trained in an end-to-end fashion with a meta-learning classification algorithm.

Among this plethora of solutions, I decided to focus on meta-learning algorithms, which currently achieve state of the art results in few-shot image classification, in addition to exploiting a conceptually fascinating paradigm. The next section proposes a formulation of this paradigm.

2.2 Meta-learning paradigm

Thrun & Pratt (1998) [17] stated that, given a task, an algorithm is *learning* “if its performance at the task improves with experience”, while, given a family of tasks, an algorithm is *learning to learn* if “its performance at each task improves with experience **and** with the number of tasks”. We will refer to the last one as a *meta-learning algorithm*. Formally, if we want to solve a task \mathcal{T}_{test} , the meta-learning algorithm will be trained on a batch of training tasks $\{\mathcal{T}_i\}$. The training experience gained by the algorithm from its attempts at solving these tasks will be used to solve the ultimate task \mathcal{T}_{test} .

I will now formalize the meta-learning framework applied to the few-shot classification problem described in section 2.1. A visualization is available in Figure 3.

To solve a N -way K -shot classification problem named \mathcal{T}_{test} , we have at our disposal a large meta-training set \mathcal{D} . The meta-training procedure will consist of a finite number of episodes.

An episode is composed of a classification task \mathcal{T}_i that is similar to the classification task \mathcal{T}_{test} we ultimately want to solve: from \mathcal{D} we sample N classes and K support-set images for each class, along with Q query images. Note that the classes of \mathcal{T}_i are entirely disjoint from the classes of \mathcal{T}_{test} (*i.e.* the classes of \mathcal{T}_{test} do not appear in the meta-training set \mathcal{D} , although they have to be similar for the algorithm to be efficient). At the end of each episode, the parameters of our model will be trained to maximize the accuracy of the classification of the Q query images (typically by backpropagating a classification loss such as negative log-probability). Thus our model learns across tasks the ability to solve an unseen classification task.

Formally, where a standard learning classification algorithm will learn a mapping $image \mapsto label$, the meta-learning algorithm typically learns a mapping $supportset \mapsto (query \mapsto label)$.

The efficiency of our meta-learning algorithm is ultimately measured on its accuracy on the target classification task \mathcal{T}_{test} .

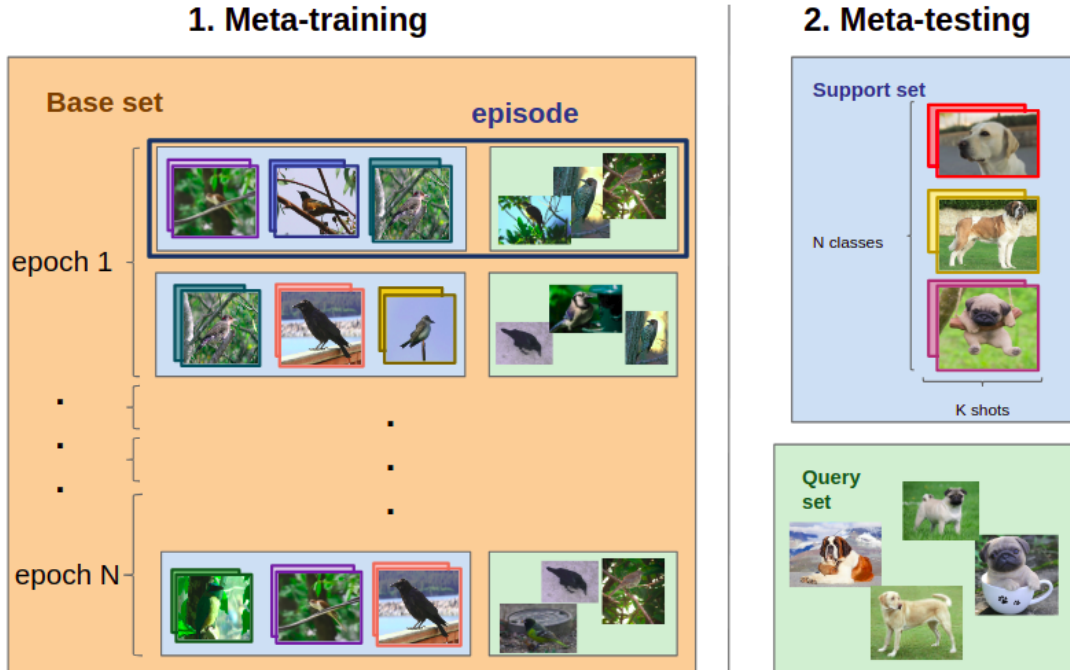


Figure 2: To solve a few-shot image classification task \mathcal{T}_{test} defined by a support set and a query set (on the right), we use a meta-training set \mathcal{D} (on the left) from which we sample episodes in the form of tasks \mathcal{T}_i similar to \mathcal{T}_{test} .

2.3 Meta-learning algorithms

Recently, several meta-learning algorithms for solving few-shot image classification are published every year. The majority of these algorithm can be labeled as either a metric learning algorithm or as a gradient-based meta-learner. Both kind are presented in this section.

2.3.1 Gradient-based meta-learning

In this setting, we distinguish the meta-learner, which is the model that learns across episodes, and a second model, the base-learner, which is instantiated and trained inside an episode by the meta-learner.

Let us consider an episode of meta-training, with a classification task \mathcal{T}_d which is defined by a support set of $N * K$ labeled images and a query set of Q images. The base-learner model, typically a CNN classifier, will be initialized, then trained on the support set (*e.g.* the base-training set). The algorithm used to train the base-learner is defined by the meta-learner model. The base-learner model is then applied to predict the classification of the Q query images. The meta-learner's parameters are trained at the end of the episode from the loss resulting from the classification error.

From this point, algorithms differ on their choice of meta-model.

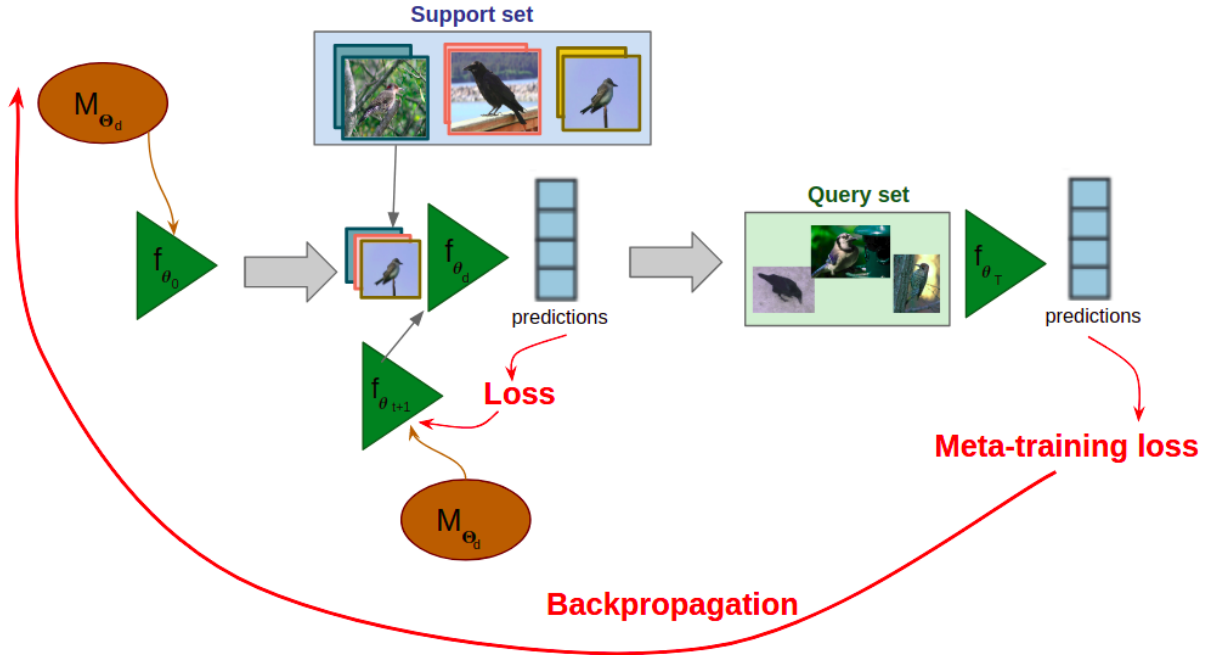


Figure 3: The d^{th} episode of meta-training, which follows this process: (1) the support set and the query set are sampled from the meta-training set; (2) the base-model f_{θ} is initialized by the meta-model M_{Θ_d} ; (3) the parameters of the base-model are fine-tuned on the support set (the fine-tuning process depends on M_{Θ_d}); (4) after T updates, the base-model is evaluated on the query set; (5) the parameters Θ of the meta-model are updated by backpropagating the loss resulting from the base-model’s predictions on the query set.

Meta-LSTM (2016) Ravi & Larochelle [11] decided to use a Long-Short-Term-Memory network [18]: the parameters θ of the base-learner f_{θ} are represented by the cell state of the LSTM, which leads to the update rule $\theta_t = f_t \odot \theta_{t-1} + i_t \odot \tilde{c}_t$ where f_t and i_t are respectively the forget gate and the input gate of the LSTM, and \tilde{c}_t is an input. We can see the update rule as an extension of the backpropagation, since with $f_t = 1$, i_t the learning rate and $\tilde{c}_t = -\nabla_{\theta_{t-1}} \mathcal{L}_t$ we obtain the standard backpropagation. Hence this model learns how to efficiently operate gradient descents on the base-model from the support set, in order to make this base-model more accurate on the query set.

Model-Agnostic Meta-Learning (2017) Finn *et al.* [9] proposed an algorithm that learns how to initiate the parameters of the base-model, but does not intervene in the base-model’s parameters update. Here, the meta-learner creates a copy of itself at the beginning of each episode, and this copy (the base-model) is fine-tuned on the support set, then makes predictions on the query set. The loss computed from these predictions is used to update the parameters of the meta-model (hence, the initialization parameters for the next episodes will be different). The algorithm as described by Finn *et al.* is shown in Figure 4.

The main feature of this method is that it is conceived to be agnostic of the base-model,

Algorithm 1 Model-Agnostic Meta-Learning

Require: $p(\mathcal{T})$: distribution over tasks
Require: α, β : step size hyperparameters

- 1: randomly initialize θ
- 2: **while** not done **do**
- 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
- 4: **for all** \mathcal{T}_i **do**
- 5: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ with respect to K examples
- 6: Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
- 7: **end for**
- 8: Update $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$
- 9: **end while**

Figure 4: Overview of the MAML algorithm with one gradient update on the support set (credits to [9])

which means that it can virtually be applied to any machine learning algorithm. Finn *et al.* tested it on supervised regression and classification, and on reinforcement learning tasks, but it could be used to solve many other problems necessitating fast adaptation of a Deep Neural Network, for instance for few-shot object detection (see section 3.3).

2.3.2 Metric Learning

In section 2.1, I presented the Siamese Networks algorithm [5], which was a first attempt at solving few-shot classification using metric learning, *i.e.* learning a distance function over objects (some algorithms actually learn a similarity function, but they are nonetheless referred to as metric learning algorithms).

As such, metric learning algorithms learn to compare data instances. In the case of few-shot classification, they classify query instances depending on their similarity to support set instances. When dealing with images, most algorithm train a convolutional neural network to output for each image an embedding vector. This embedding is then compared to embeddings of other images to predict a classification.

Matching Networks (2016) As explained in section 2.1, Siamese Networks train their CNN in a discrimination task (*are these two instances from the same class?*) but the algorithm is tested on a classification task (*to which class does this instance belong?*). This issue of task shift between training and testing time is solved by Vinyals *et al.* [6]. They proposed the Matching Networks, which is the first example of a metric learning algorithm inside the meta-learning framework.

To solve a few-shot image classification task, they use a large meta-training set from which they sample episodes (see Figure 3). For each episode, they apply the following procedure:

1. Each image (support set and query set) is fed to a CNN that outputs as many embeddings;
2. Each query image is classified using the softmax of the cosine distance from its embedding to the embeddings of support set images;

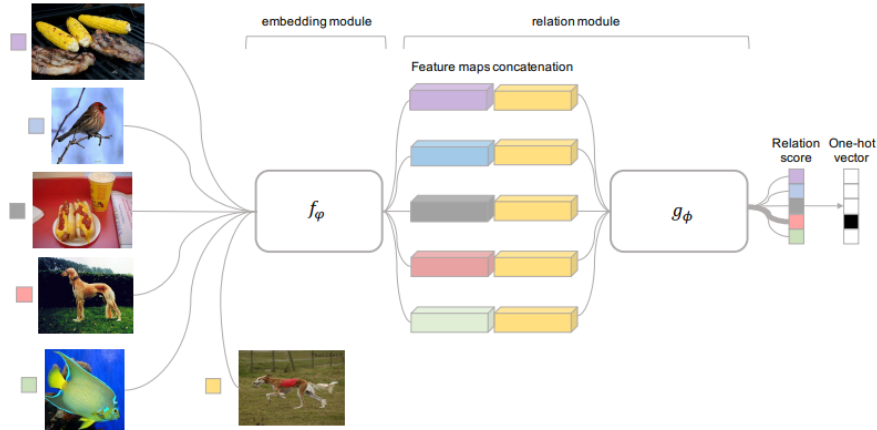


Figure 5: Relation Network architecture for a 5-way 1-shot problem with one query example (credits to [8]). Note that they chose to represent the final output with a one-hot vector obtained by a *max* function on the relation scores, but that during training time we need to use a *softmax* to make the network differentiable.

3. The cross-entropy loss on the resulting classification is backpropagated through the CNN;

This way, the Matching Networks learn to compute a representation of images that allows to classify them with no specific prior knowledge on the classes, simply by comparing them to a few instances of these classes. Since considered classes are different in every episode, Matching Networks are expected to compute features of the images that are relevant to discriminate between classes, whereas a standard classification learning algorithm is expected to learn the features that are specific to each class.

It is to be noted that Vinyals *et al.* also proposed to augment their algorithm with a Full Context Embedding process: the embedding of each image depends on the embeddings of the others thanks to bidirectional LSTM. They expect that this better exploit all the available knowledge on the episode. This process slightly improved the performance of their algorithm on the miniImageNet benchmark, but also demands a longer computing time.

Prototypical Networks (2017) Building on Matching Networks, Snell *et al.* [7] proposed Prototypical Networks. The process is essentially the same (although Full Context Embeddings are not used), but a query image is not compared to the embeddings of every images of the support set. Instead, the embeddings of the support set images that are from the same class are averaged to form a class prototype. The query image is then compared only to these prototypes. It is to be noted that when we only have one example per class in the support set (One-Shot Learning setting) the Prototypical Networks are equivalent to the Matching Networks. They obtained better results than Matching Networks on the miniImageNet benchmark, and expose that part of this improvement must be credited to their choice of distance metric: they notice that their algorithm and Matching Networks both perform better using Euclidean distance than when using cosine distance.

Relation Network (2018) Sung *et al.* [8] built on Prototypical Networks to develop the Relation Network. The difference is that the distance function on embeddings is no longer arbitrarily defined in advance, but learned by the algorithm (see Figure 5): a *relation module* is put on top of the *embedding module* (which is the part that computes embeddings and class prototypes from the input images). This relation module is fed the concatenation of the embedding of a query image with each class prototype, and for each couple outputs a relation score. Applying a softmax to the relation scores, we obtain a prediction.

2.4 Few-Shot Image classification benchmarks

Algorithms intended to solve the few-shot learning problem are usually tested on two datasets: Omniglot and miniImageNet.

Omniglot Lake *et al.* (2011) [19] introduced the Omniglot dataset. It is composed of 1623 characters from 50 distinct alphabets. Each one of these characters is a class and contains 20 samples drawn by distinct people. Each data instance is not only a 28x28x1 image, but also contains information about how it was drawn: how many strokes, and the starting and ending point of each stroke (see Figure 6). Although Lake *et al.* primarily used Omniglot for few-shot learning of visual concepts from their subparts [20], the dataset as a set of 28x28 one-channel images is used as a MNIST-like benchmark for few-shot image classification. Most algorithm now achieve a 98%-or-better accuracy on this dataset on most use cases [8].



Figure 6: Two different visualizations of a same instance of the Omniglot dataset. On the left, we can see how the character was drawn. On the right, we see a 28x28 one-channel image. Credits to [20]

miniImageNet Vinyals *et al.* [6] proposed to use a part of ImageNet as a new, more challenging benchmark for few-shot image classification. Their dataset consist of 100 classes, each containing 600 3-channel images. The commonly used train/validation/evaluation split of this dataset [11] separates it in three subsets of respectively 64, 16 and 20 classes. This way, we ensure that the algorithm is evaluated on classes that were not seen during training.

2.5 Few-Shot object detection

Although research in few-shot object detection is currently less advanced than in few-shot classification, some solutions to this problem have been proposed in the last few months. First,

we will go over the existing solutions for standard object detection, then we will learn about the recent efforts in developing algorithms for few-shot object detection.

Object detection Algorithms for object detection can be separated in two categories: single-stage detectors and the R-CNN family (two-stage detectors). Single-stage detectors aim at performing fast detection while algorithms like R-CNN are more accurate.

R-CNN [21] uses a first network to determine regions of interest in an image, and then a second network to classify the content of each region of interest. Fast R-CNN [22] and Faster R-CNN [23] improved the algorithm's efficiency by reducing redundant computations and the number of regions of interest. Mask R-CNN [24] uses the same principle as R-CNN but performs image segmentation.

Single-stage detectors perform object detection on an image in a single forward pass through a CNN: the bounding-box and the label of each object are predicted concurrently. Leading single-stage detectors are the SSD (for Single-Shot Detector) [25], RetinaNet [26] and YOLO (for You Only Look Once) [27].

YOLO went through two incremental improvements since its creation in 2016. Its last version, YOLOv3, contains three output layers. Each one is responsible for predicting respectively large, medium-size and small objects. For each output layer, three *anchors* are set as hyper-parameters of the model. An anchor is like a "default bounding box", and YOLOv3 actually predicts deformations to these anchors, rather than predicting a bounding box from scratch. The network is mostly composed of residual blocks [1]. In particular, the backbone of the model is a Darknet53, a 53-layer residual network pre-trained on ImageNet. A visualization of the YOLOv3 architecture is available in Figure 7.

Few-Shot Object Detectors To the best of my knowledge, the first few-shot detector was proposed in late 2018 by Kang *et al.* [28]. Their algorithm combines a standard single-stage detector with an other auxiliary network. This second model is responsible for reweighting the features outputted by the feature extractor of the model (in YOLOv3, this would be the output of the Darknet53). The goal of this reweighting is to give more importance to features related to the specific few-shot detection task being solved (the intuition is that the relevant features for detection depends of the type of object to detect). The reweighting model is trained in a meta-learning set-up (see section 2.2): at each meta-training episode, a few-shot detection task is sampled (for instance: detecting dogs and cats, with a few annotated examples of dogs and cats), and the training objective is the precision of the detector.

More recently, Fu *et al.* proposed the Meta-SSD [29]. They apply Li *et al.*'s Meta-SGD [10] to Liu *et al.*'s Single-Shot Detector. They end up with a fully meta-trainable object detector.

Concurrently, Wang *et al.* [30] developed a novel framework around the Faster R-CNN. The resulting algorithm can adapt to a new task with few labeled examples.

Previous works already tackled few-shot object detection [31] [32], although they considered a slightly different problem: they defined few-shot as few labeled images per category, but also used a large pool of unlabeled data.

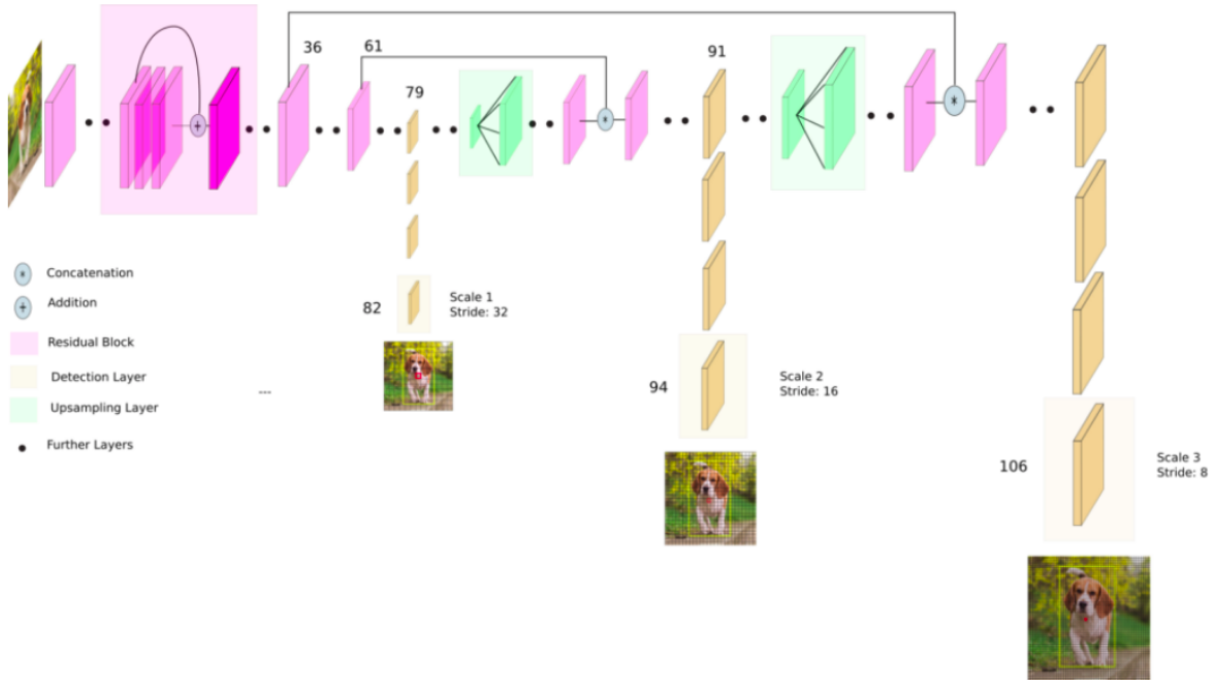


Figure 7: Architecture of YOLOv3. Credits to Ayoosh Kathuria.

3 Contributions

3.1 Overview

Sicara is a company which develops computer vision solutions based on machine learning algorithms for businesses. However, it is common that the amount of data made available by the customer is not large enough to effectively train a standard convolutional neural network. Also, we often need to harness this data with a very short lead time. Therefore, a company like Sicara needs an efficient and ready-to-use meta-learning algorithm for few shot learning problems related to computer vision.

I was in charge of the first step of this process, which is benchmarking several state-of-the-art algorithms, identifying the strengths and weaknesses of each algorithm, its performance on different kinds of datasets, and overall their relevance depending on the task that needs solving.

During this work on meta-learning algorithms, we decided to focus on the Model Agnostic Meta-Learner [9] and to switch from the few-shot image classification problem to the few-shot object detection problem, which had until then attracted less attention in the research community than few-shot classification. Our idea is to apply MAML to the YOLOv3 object detector in order to obtain an algorithm capable of detecting new classes of objects with little time and only a few examples.

In this section, I will first explain my work on meta-learning algorithms for few-shot image classification, then I will detail my progress so far in developing a novel algorithm: the YOLOMAML.

3.2 Meta-learning algorithms for Few-Shot image classification

I compared the performance of four meta-learning algorithms, using two datasets: miniImageNet (see section 2.4) and Caltech-UCSD Birds 200 (CUB) [33], which is the dataset containing 6,033 pictures of birds from 200 different classes. The four algorithms are the following:

- Matching Networks [6]
- Prototypical Networks [7]
- Relation Network [8]
- Model Agnostic Meta-Learner [9]

The primary intention was to conduct extensive experiments on these algorithms with variations on both their settings, the target tasks and the training strategy, in order to obtain a fine understanding of how these algorithms behave and how to best harness their abilities. I also intended to include other promising algorithms, such as the Simple Neural Attentive Learner [12] or the Meta-LSTM [11]. However, since we decided halfway through the benchmark to focus on the exciting opportunity of developing a novel meta-learning object detector, there wasn't enough time to go through the all set of experiments. Hence, my contribution for a deeper understanding of meta-learning consists in:

1. a documented implementation of meta-learning algorithms for few-shot classification tasks, with a focus on allowing future researchers in the field to easily launch new experiments, in a clear and reproducible way;
2. the reproduction of the results presented by Chen *et al.* [34], bringing the exposition of the challenges that we face when benchmarking meta-learning algorithms;
3. a study on the impact of label noise in the support set at evaluation time;

In this subpart I will present these contributions with more details, both on the results and on the process of obtaining these results.

3.2.1 Implementation

Chen *et al.* [34] published in April 2019 a first unified comparison of meta-learning algorithms for few-shot image classification, and made their source code available¹. For us, this code in PyTorch presents two main advantages:

1. It proposes a unified implementation of Matching Networks, Prototypical Networks, Relation Network, MAML and two baseline methods for comparison. This allows the experimenter to fairly compare algorithms.

¹<https://github.com/wyharveychen/CloserLookFewShot>

2. It contains a relatively consistent framework for the treatment of the several datasets (Omniglot, EMNIST [35], miniImageNet and CUB), from the description of the train / validation / evaluation split using `json` to the sampling of the data in the form of episodes for few-shot image classification, which uses the code from Hariharan *et al.* [14]².

For these reasons, I used this code as a (very advanced) starting point for my implementation. I identified three main issues:

1. The original code was very scarcely documented, which makes it difficult to understand, and even more difficult to modify, since it was not always clear what a chunk of code did, or what a variable represented.
2. Some experiment parameters were defined inside the code and therefore not easily customizable when launching an experiment, nor monitorable after the experiments, affecting the reproducibility of the experiments.
3. Some chunks of code were duplicated in several places in the project.

The main goal of my work on this code was to make it easily accessible, allowing future researcher to understand the way these algorithms work in practice, and to quickly be able to launch their own experiments. This goal was achieved by:

- cleaning the code and removing all duplicates;
- extensively document every class and function with the knowledge gained during my work on the code;
- integrate two useful tools for conducting experiments:
 - `pipeline` is an internal library at Sicara which allows to configure experiments with a YAML file: this file describes the different steps of the experiment and explicitly indicates all necessary parameters of the experiment;
 - `Polyaxon` is an open-source platform for conducting machine learning experiments; its main features (for our usage) are (1) an intuitive dashboard for keeping track of all passed, current and programmed experiments, with for each one the YAML configuration file, along with all logs and outputs, (2) the possibility to launch groups of experiments with varying parameters, and (3) a Tensorboard integrated to the platform.

The structure of the implementation is shown in Figure 8. The code can be divided in five categories, detailed below.³

²<https://github.com/facebookresearch/low-shot-shrink-hallucinate>

³The code is available at <https://github.com/ebennequin/FewShotVision>. Note that this version does not use the `pipeline`, which is private to Sicara.

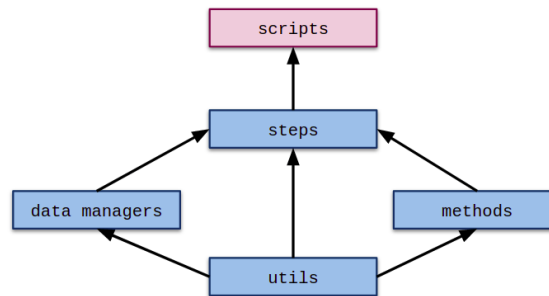


Figure 8: Structure of my code to conduct experiments on meta-learning algorithms for few-shot image classification

scripts These are the files that must be executed to launch the experiments. I used YAML files for compatibility with the `pipeline` library, but standard Python scripts could be used just as well (and are actually used in the publicly available repository). They describe the flow between the different steps (which steps uses which step outputs) and contain all the values parameterizing the experiment:

- dataset to work on (ex: miniImageNet);
- method to work with (ex: Matching Networks);
- backbone CNN of the method (ex: Resnet18);
- parameters of the N -way k -shot classification task with q queries per class (with N allowed to be different at training and evaluation time);
- whether to perform data augmentation on the meta-training set;
- number of meta-training epochs;
- number of episodes (*i.e.* classification tasks) per meta-training epoch;
- optimizer (ex: Adam);
- learning rate;
- which state of the model to keep for evaluation (the model trained on all the epochs, or the model that achieve the best validation accuracy);
- number of few-shot classification task to evaluate the model on;

steps They are called by the scripts and use the parameters explicited in it. One example of step is `MethodTraining`, which is responsible for the training of the model.

data managers They define the `SetDataset` and `EpisodicBatchSampler` classes, which respectively extend the PyTorch base classes `Dataset` and `Sampler` and are used to build a `DataLoader` that loads the data in the shape of few-shot classification task (*i.e.* a support set and a query set, instead of regular batches of arbitrary size).

methods Each file in this category defines a class corresponding to one meta-learning algorithm (ex: Prototypical Networks). Every class contains three essential methods:

- `set_forward(episode)`: takes as input an episode composed of a support set and a query set, and outputs the predictions of the model for the query set.
- `train_loop()`: executes one meta-training epoch on the meta-training set.
- `eval_loop()`: evaluates the model on few-shot classification tasks sampled from the evaluation set.

utils These files contain all the utilities used in the rest of the code.

3.2.2 Reproducing the results

The first thing to do with this reimplementaion was to validate it by reproducing the results reported by Chen *et al.* [34]. This unexpectedly granted us with interesting new knowledge.

I experimented on the CUB dataset for a shorter running time. I reproduced Chen *et al.*'s experiments in the 5-way 1-shot and 5-way 5-shot settings, for Matching Networks, Prototypical Networks, Baseline and Baseline++ (see Figure 9). I purposefully omitted MAML for this part, since this algorithm's training takes about five times longer than the others' (see Table 1). Relation Network is also omitted because its process is essentially similar to Prototypical Networks.

	CUB		miniImageNet	
	1 shot	5 shots	1 shot	5 shot
Baseline	1h10	1h00	10h05	10h07
Baseline++	56mn	51mn	10h25	10h28
MatchingNet	6h41	4h21	7h51	6h23
ProtoNet	6h38	5h07	7h40	6h08
MAML	28h05	22h28	31h22	25h22

Table 1: Running time of several algorithms depending on the setting and dataset. This is the running time of the whole process, from training to evaluation.

The parameters of the experiments follow those described by Chen *et al.*, *i.e.* a 4-layer CNN as a backbone, an Adam optimizer with a learning rate of 10^{-3} , 100 episodes per epoch and data augmentation on the training set. The baselines are trained for 200 epochs on CUB, and for 400 epochs on miniImageNet. The other algorithms are trained for 600 epochs in the 1-shot setting, and on 400 epochs in the 5-shots setting. We keep the state of the model that had the

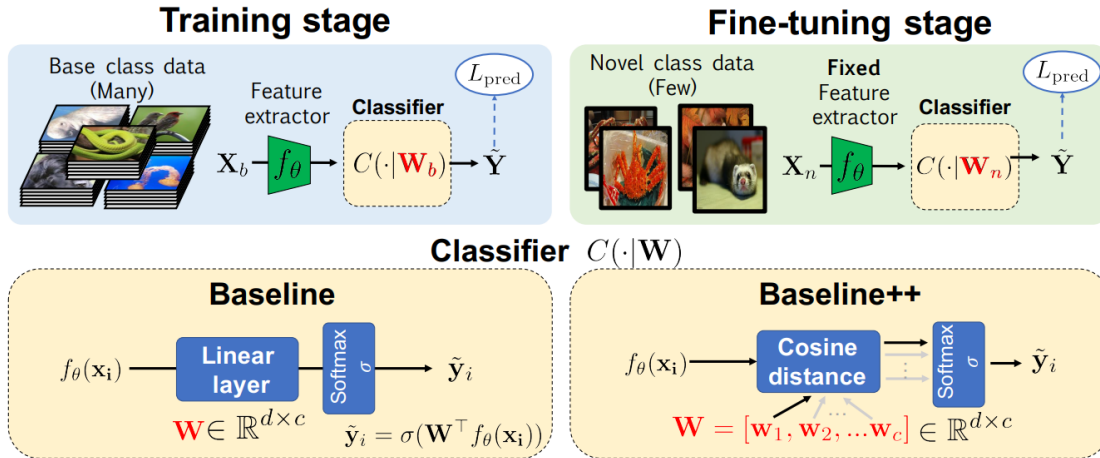


Figure 9: Baseline and Baseline++ few-shot classification methods. Both algorithms are pre-trained on the meta-training set. When evaluated on a few-shot classification task, the feature extractor f_θ is frozen and the classifier \mathcal{C} is fine-tuned on the support set before being applied to the query set. In Baseline++, the classifier is not a standard fully connected layer, but computes the cosine distance between its weights and the input features vector. Both algorithm are used to compare the meta-learning algorithms to non-meta-learning methods. This figure is credited to [34].

best accuracy on the validation set, and evaluate it on 600 few-shot classification tasks sampled from the evaluation set.

The results of these experiments are reported in Table 2. 6 out of 8 experiments gave results out of the 95% confidence interval reported by Chen et al, with a difference up to 6% in the case of 1-shot Baseline++. Our results fall below the confidence interval in 4 cases and above the confidence interval in 2 cases.

	our reimplementaion		Chen <i>et al.</i> 's	
	1 shot	5 shots	1 shot	5 shot
Baseline	46.57 \pm 0.73	68.36 \pm 0.66	47.12 \pm 0.74	64.16 \pm 0.71
Baseline++	53.71 \pm 0.82	75.09 \pm 0.62	60.53 \pm 0.83	79.34 \pm 0.61
MatchingNet	58.43 \pm 0.85	75.52 \pm 0.71	61.16 \pm 0.89	72.86 \pm 0.70
ProtoNet	50.96 \pm 0.90	75.48 \pm 0.69	51.31 \pm 0.91	70.77 \pm 0.69

Table 2: Comparison of the results of our reimplementaion compared to the results reported by Chen *et al.*, on the CUB dataset with a 5-way classification task. Our results are shown in bold when they are out of the 95% confidence interval reported by Chen *et al.*

A fair assumption was that my implementation was to blame for this incapacity to reproduce the original paper’s results. To verify it, I reproduced the experiments of Chen *et al.* using their original implementation. The results are shown in Table 3. In most cases, they are out of the 95% confidence interval reported in [34].

	CUB		miniImageNet	
	1 shot	5 shots	1 shot	5 shot
Baseline++	61.36 \pm 0.92	77.53 \pm 0.64	48.05 \pm 0.76	67.01 \pm 0.67
MatchingNet	59.55 \pm 0.89	75.63 \pm 0.72	48.43 \pm 0.77	62.26 \pm 0.70
ProtoNet	50.28 \pm 0.90	75.83 \pm 0.67	43.89 \pm 0.73	65.55 \pm 0.73
MAML	54.57 \pm 0.99	75.51 \pm 0.73	43.92 \pm 0.77	62.96 \pm 0.72

Table 3: Reproduction of the results of [34] on both CUB and miniImageNet, using the implementation provided with the paper. Our results are shown in bold when they are out of the 95% confidence interval reported in [34].

From there, my assumption was that the incertitude on the results didn’t come solely from the sampling of the evaluation tasks, but also from the training. I proceeded to verify this assumption. I relaunched the first experiment 8 times for Prototypical Networks and evaluated the 8 resulting model on the exact same classification tasks. The results are shown in Table 4. We can see that the accuracy can go from 74.20% to 76.04% on the same set of tasks. This validates that two same models trained with the same hyperparameters may obtain different accuracies on the same evaluation tasks.

76.04 \pm 0.71
74.77 \pm 0.71
75.45 \pm 0.71
75.54 \pm 0.70
75.14 \pm 0.70
74.90 \pm 0.71
74.91 \pm 0.71
74.20 \pm 0.70

Table 4: Accuracy of the Prototypical Network on the same set of evaluation tasks on the CUB dataset in the 5-way 5-shot setting, after 8 independent training processes.

From this work on the reproduction of the results reported by Chen *et al.*, we can retain two main take-aways:

1. The results obtained with a instance of *meta-training + evaluation* cannot be reproduced, although the `numpy` random seed is systematically set to 10 at the beginning of the process. I learned that setting the `numpy` random seed is not enough to fix the randomness of a training using `PyTorch`. I found that the following lines succeed in doing so:

```
np.random.seed(numpy_random_seed)
torch.manual_seed(torch_random_seed)
torch.backends.cudnn.deterministic = True
torch.backends.cudnn.benchmark = False
```

The third and fourth lines are only necessary when using the CuDNN backend⁴.

2. On the same set of evaluation tasks, the accuracy of a model can vary with an amplitude of up to 2% due to randomness in the training. This amplitude is similar to the reported difference in accuracy between algorithms and higher than the confidence intervals usually reported when evaluating meta-learning algorithms [34] [8] [9] [11] [7] [36]. I argue that a reported difference of a few percents in accuracy between two meta-learning algorithms on a set of classification tasks cannot be considered as a relevant comparator of these algorithms. It would be ideal to get an exact measure of the uncertainty by launching a sufficient number of trainings, but the necessary computing time for this operation is prohibitive (see Table 1).

3.2.3 Effects of label noise in the support set at evaluation time

In practice, meta-learning algorithms can be used this way:

1. The model is trained once and for all by the model’s designer on a large dataset (with a possibility to update when new labeled examples become available);
2. When faced with a novel few-shot classification task, the user feeds a few labeled examples to the model, and then is able to apply it to the query images.

As the model’s designer and the model’s user can be different entities, and as the source of the support set for the novel task may be different from the source of the meta-training data, the designer may not be able to control the quality of the data in the novel task. This is why the model’s robustness to noisy data in the support set is an important issue.

In this subsection, we address the issue of label noise (*i.e.* data instances assigned with the wrong label) in the support set of the evaluation classification task. To simulate this noise, we use label swaps: given an integer M , for each classification task, we execute M label swaps on the support set of the classification task. Here is the process of one label swap:

1. Sample uniformly at random two labels l_1, l_2 among the N labels of the support set
2. For each label l_x , select uniformly at random one image i_{l_x} among the K images in the support set associated with this label
3. Assign label l_1 to image i_{l_2} and label l_2 to image i_{l_1}

Note that even though one label swap changes the label of two images, M label swaps do not necessarily cause $2M$ falsely labeled images, since swapped images are sampled with replacement (in the following, you will see that most models reach an accuracy of 35% even after 10 label swaps were applied on the $5 \text{ labels} \times 5 \text{ images}$ support set, which would be hard to explain if 80% of the support set had false labels).

Also, label swaps are not a perfect simulation: in real cases, the fact that an image supposed to have the label l_1 was falsely labeled with l_2 does not mean that an other image supposed to

⁴<https://pytorch.org/docs/stable/notes/randomness.html>

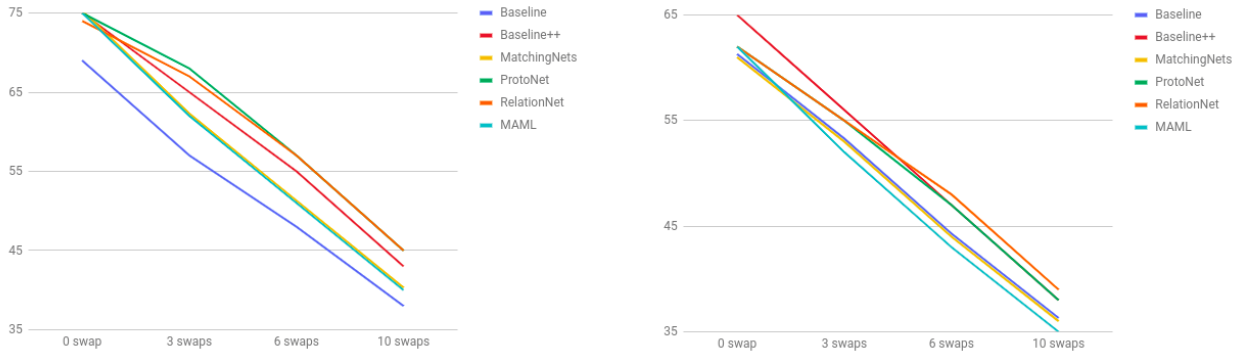


Figure 10: Accuracy of the methods for different number of label swaps in the support set of each classification task. Left: CUB. Right: miniImageNet.

have the label l_2 was falsely labeled with l_1 . However, this solution ensures that the support set is still balanced even after applying the artificial label noise (in a 25-images dataset, if one label has one example less than an other label, the dataset becomes noticeably unbalanced). Therefore, we know that the results will not come from errors in labelling, and not from an unbalanced dataset.

Measuring the effects of label noise in the support set at evaluation time First, we need to measure the effect of label noise on the model’s accuracy. I experimented both on CUB and miniImageNet, with the algorithms Baseline, Baseline++, Matching Networks, Prototypical Networks, Relation Network and MAML. All models were trained on 400 epochs, with the Adam optimizer and a learning rate of 10^{-3} . Meta-learning algorithms (*i.e.* all but Baseline and Baseline++) were trained on 5-way 5-shot classification tasks. No artificial label noise was added to the training set.

The models were then evaluated on 5-way 5-shot classification tasks on four different settings corresponding to four different number of label swaps in each classification task (0, 3, 6 and 10). I reported for each setting the mean of the accuracy on 600 tasks. Note that all models (here and in the following of this subsection) are evaluated on the same tasks. To be consistent with my remarks in section 3.2.2, the results are reported with a precision of 1%.

The results are shown in Figure 10. We observe that all algorithms endure a serious drop in accuracy on the query set when the label noise in the support set increases, which was expected. We notice that Prototypical Networks and Relation Network are slightly less impacted. This could be explained by the fact that both algorithms use the mean of the features vectors for each class, which reduces the impact of extreme values.

10-way training Snell *et al.* [7] showed that, when evaluating metric learning algorithms on N -way K -shot classification tasks, the models trained on N' -way K -shot classification tasks with $N' > N$ performed better than the models trained on N -way K -shot classification tasks (the intuition being that a model trained on more difficult tasks will generalize better to new tasks, or, in French, “*qui peut le plus peut le moins*”). I tested whether this trick also made the



Figure 11: Accuracy of the methods for different number of label swaps in the support set of each classification task, with a 5-way training and a 10-way training. Left: CUB. Right: miniImageNet.

model more robust to label noise.

I conducted the same experiment as the one described in the previous paragraph, with the exception that the training was done on 10-way 5-shot classification tasks (instead of 5-way 5-shot). This experiment was done only on metric learning algorithms (*i.e.* Matching Networks, Prototypical Networks, Relation Networks). Indeed, MAML does not allow to change the number of labels in the classification tasks, since the architecture of the CNN (ending in a N -filter linear layer) needs to stay the same.

The results are shown in Figure 11. They confirm that using a higher number of labels per classification task during training increases the accuracy of the model. However, this doesn't seem to have any effect on the robustness to label noise.

Simulating label noise during meta-training Coming from the idea that training and testing conditions must match, I assumed that incorporating artificial label noise in the support set of the classification tasks on which the models are meta-trained could increase their robustness to label noise at evaluation time. The following experiment tests this assumption. Label swaps are introduced in the classification tasks composing the meta-training, in the same way that they were applied to the classification tasks at evaluation time in the previous experiments. This results in three set-ups, respectively referred to as 0, 3 and 10-swap training:

1. Same experiment as in the first paragraph of this section, only on miniImageNet (not CUB)
2. Same, but in each episode of the meta-training, 3 label swaps are applied to the support set
3. Same, but in each episode of the meta-training, 10 label swaps are applied to the support set

Note that we do not experiment on the baselines, since they are not meta-learning algorithm and thus do not solve classification task during training. The results of this experiment are

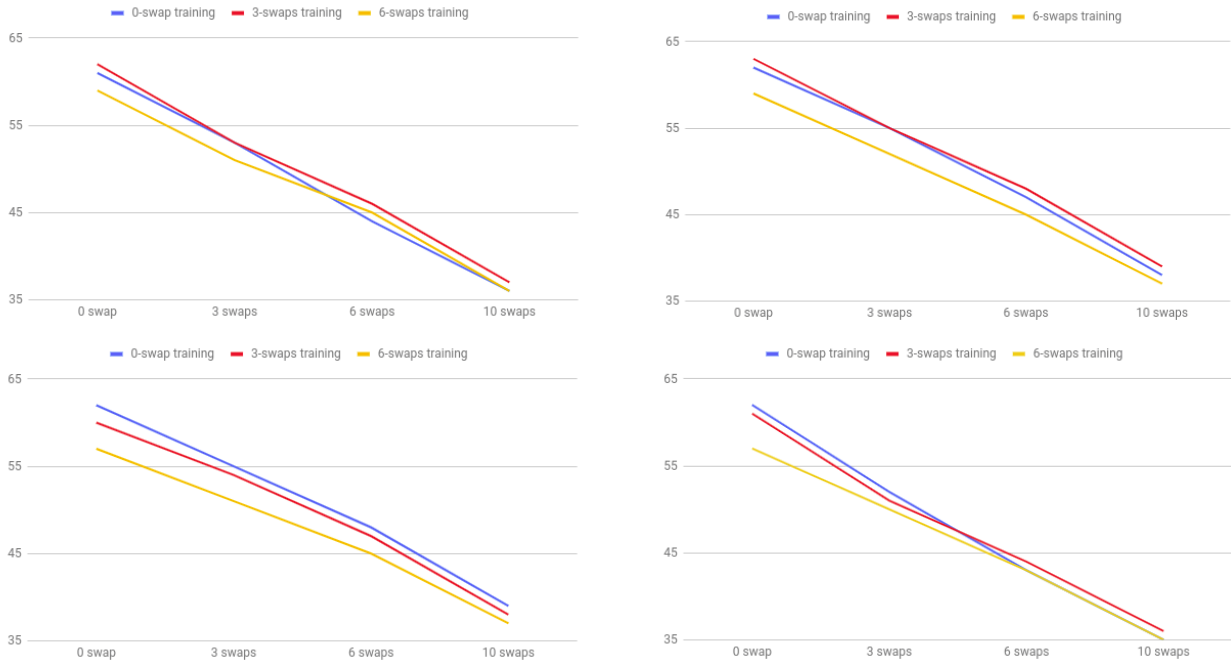


Figure 12: From left to right, top to bottom: Matching Networks, Prototypical Networks, Relation Network, MAML. For each method, accuracy on a model trained with three strategies, for different number of label swaps in the support set at evaluation time.

shown in Figure 12. We see that adding label swaps during meta-training causes a drop in accuracy when the model is evaluated on correctly labeled classification tasks. The difference is less obvious when the number of label swaps in evaluation tasks increases. Based on these experiments, there is no reason to believe that introducing artificial label noise during the meta-training makes meta-learning algorithms more robust to label noise in novel classification tasks.

3.2.4 Future work

In addition to the choice of the meta-learning algorithm, there are many possible ways to improve its performance with minor design choices, such as hyperparameter tuning, or, in the case of Prototypical Networks and their derivatives, the choice of the prototype function. The mean function could be replaced for instance by a "leaky" median (in a way that leaves the function differentiable).

However, we saw that these algorithms only differ by a small margin. It would be interesting to explore different ways to improve performance at few-shot classification. One way could be to compare the performance of meta-learning algorithms depending on the "shape" of the meta-training dataset. Would a dataset with 100 different classes and 500 examples per class allow better performance than a dataset with 50 classes and 1000 examples per class? My assumption is that it would, since it would allow the algorithm to better generalize to new classes, but this still needs to be proven.

Finally, in addition to the classification accuracy, it would be interesting to collect more intelligence about how meta-learning algorithm actually learn, for instance by studying the features representations, or using Explainable Machine Learning techniques, adapted to the meta-learning paradigm.

3.3 MAML for Few-Shot object detection

3.3.1 The Few-Shot Object Detection problem

We saw that in real world applications, we sometimes need to solve an image classification task with only few examples. The same problem is encountered in all other tasks composing the field of computer vision. Here, we tackle the Few-Shot Object Detection problem.

Here we define the object detection task as follow: given a list of object types and an input image, the goal is to detect all object belonging in the list. Detecting an object consists in:

1. localizing the object by drawing the smallest bounding box containing it;
2. classifying the object.

As such, object detection is the combination of a regression task and a classification task. An example is shown in Figure 13.



Figure 13: Consider the task of detecting objects belonging in $\{ \text{laptop, mug, notebook, lamp} \}$. Here the object detector detects and classifies a laptop, a mug and a notebook. It doesn't detect the plant and the pen since they are not part of the given task.

Following this, we define a N -way K -shot object detection task as follows. Given:

1. a support set composed of:
 - N class labels;
 - For each class, K labeled images containing at least one object belonging to this class;

2. Q query images;

we want to detect in the query images the objects belonging to one of the N given classes. The $N \times K$ images in the support set contain the only examples of object belonging to one of the N classes.

When K is small, we talk about *few-shot object detection*.

We can immediately spot a key difference with few-shot image classification: one image can contain multiple objects belonging to one or several of the N classes. Therefore, when solving a N -way K -shot detection tasks, the algorithm trains on *at least* K example objects for each class. During a N -way K -shot classification tasks, the algorithms sees *exactly* K examples for each class. Note that this can become a challenge: in this configuration, the support set may be unbalanced between classes. As such, this formalization of the few-shot object detection problem leaves room for improvement. It was chosen because it is a rather straightforward setup, which is also convenient to implement, as we will see in section 3.3.3.

3.3.2 YOLOMAML

To solve the few-shot object detection problem, we had the idea of applying the Model-Agnostic Meta-Learning algorithm [9] to the YOLO [37] detector. We call it YOLOMAML for lack of a better name.

As presented in section 2.3.1, MAML can be applied to a wide variety of deep neural networks to solve many few-shot tasks. Finn *et al.* considered few-shot classification and regression as well as reinforcement learning. It could as well be applied to a standard detector to solve few-shot object detection.

YOLOv3 is already used on other projects at Sicara. Our expertise on this detector motivated our choice to use it. Also, it prevents the advantage of being a single-stage detector. It appeared easier to apply MAML to YOLO than to a variant of R-CNN.

At the same time, Fu *et al.* proposed the Meta-SSD [29]. It applies the Meta-SGD [10] (a variant of MAML which additionally meta-learns hyper-parameters of the base model) to the Single-Shot Detector [25]. Fu *et al.* presented promising results. Although Meta-SSD and YOLOMAML are very similar, I argue that it is relevant to continue to work on YOLOMAML, in order to:

1. confirm or challenge the interesting results of Fu *et al.*, with a similar algorithm and on a wider variety of datasets;
2. disclose the challenges of developing such an algorithm.

YOLOMAML is a straightforward application of the MAML algorithm to the YOLO detector. The algorithm is shown in Algorithm 1.

3.3.3 Implementation

I decided to re-use the structure of the MAML algorithm from my work on Image Classification. For the YOLO model, I used the implementation from Erik Linder-Norén⁵, which is mostly a

⁵<https://github.com/eriklindernoren/PyTorch-YOLOv3>

Algorithm 1 YOLOMAML

Require: α, β , respectively the inner loop and outer loop learning rate

Require: `n_episodes` the number of few-shot detection tasks considered before each meta-gradient descent

Require: `number_of_updates_per_task` the number of inner loop gradient descents in each few-shot detection task

```

1: initialize the parameters  $\theta$  of the YOLO detector  $f_\theta$ 
2: while not done do
3:   sample n_episodes detection tasks  $\mathcal{T}_i$ , where each task is defined by a support set  $S_i = \{x_j^S, l_j^S\}$  and a query set  $Q_i = x_j^Q, l_j^Q$ 
4:   for  $\mathcal{T}_i$  in  $\{\mathcal{T}_i\}$  do
5:      $\theta_0 \leftarrow \theta$ 
6:     for  $t < \text{number\_of\_updates\_per\_task}$  do
7:       compute the gradient of the loss of the YOLO model  $f_{\theta_t}$  on the support set:
          $\nabla_{\theta_t} \mathcal{L}(f_{\theta_t}(\{x_j^S\}), \{l_j^S\})$ 
8:       update  $\theta_{t+1} \leftarrow \theta_t - \alpha \nabla_{\theta_t} \mathcal{L}(f_{\theta_t}(\{x_j^S\}), \{l_j^S\})$ 
9:     end for
10:    compute the gradient of the loss of the YOLO model  $f_{\theta_{\text{number\_of\_updates\_per\_task}}}$  on the query set relative to initial parameters  $\theta$ :  $\nabla_{\theta} \mathcal{L}(f_{\text{number\_of\_updates\_per\_task}}(\{x_j^Q\}), \{l_j^Q\})$ 
11:  end for
12:  update  $\theta \leftarrow \theta - \beta \sum_{\mathcal{T}_i \in \{\mathcal{T}_i\}} \nabla_{\theta} \mathcal{L}(f_{\text{number\_of\_updates\_per\_task}}(\{x_j^Q\}), \{l_j^Q\})$ 
13: end while

```

PyTorch reimplementation of Joseph Redmon’s original C implementation⁶. It contains two main parts:

- Data processing from raw images and labels to an iterable `Dataloader`. The class `ListDataset` is responsible for this process.
- The definition, training and induction of the YOLO algorithm, mostly handled by the class `Darknet`.
 - It creates the YOLO algorithm as a sequence of PyTorch `Module` objects, from a configuration file customizable by the user.
 - It allows to load pre-trained parameters for part or all of the network.
 - It defines the forward pass of the model and the loss computation.

The experiences in few-shot object detection were made on the COCO 2014 dataset [38].

I had to work on three main levels of the implementation to allow complementarity between YOLO and MAML:

- model initialization;

⁶<https://github.com/pjreddie/darknet>

- fast adaptation of weights in convolutional layers
- data processing in the form of few-shot detection episodes

Model initialization YOLOv3 in its standard form contains more than 8 millions parameters. Thus a full meta-training of it with MAML (which involves second order gradient computation) is prohibitive in terms of memory. Therefore:

1. Instead of the standard YOLOv3 neural network, I used a custom Deep Tiny YOLO. The backbone of the model is the Tiny Darknet⁷. On top of it, I added two output blocks (instead of three in the regular YOLOv3). The full configuration file of this network is available in the repository⁸ in `detection/configs/deep-tiny-yolo-5-way.cfg`.
2. I initialized the backbone with parameters trained on ImageNet, then froze those layers.

This way, there were only five trainable convolutional blocks left in the network. This allows to train the YOLOMAML on a standard GPU in a few hours. Note that there exists a Tiny YOLO, but there is no available backbone pre-trained on ImageNet for this network, which motivated my choice of a new custom network.

Fast adaptation The core idea of MAML is to update the trainable parameters on each new task, while training the initialization parameters across tasks. For this, we need to store the updated parameters during a task, as well as the initialization parameters. A solution for this is to add to each parameter a field `fast` which stores the updated parameters. In our implementation (inherited from [34]), this is handled by `Linear_fw`, `Conv2d_fw` and `BatchNorm2d_fw` which respectively extend the `nn.Linear`, `nn.Conv2d` and `nn.BatchNorm2d` PyTorch objects. I modified the construction of the `Darknet` objects so that they use these custom layers instead of the regular layers.

Data processing As in few-shot image classification, we can sample a N -way K -shot detection task with Q queries per class by first sampling N classes. Then, for each class, we sample $K + Q$ images which contain at least one box corresponding to this class. The difference in detection is that we then need to eliminate from the labels the boxes that belong to a class that does not belong to the detection task. There would be the same problem with multi-label classification.

To solve this problem, I created an extension of the standard PyTorch `Sampler` object: `DetectionTaskSampler`. In addition to returning the indices of the data instances to the `DataLoader`, it returns the indices of the sampled classes. This information is processed in `ListDataset` to feed the model proper few-shot detection task with no reference to classes outside the task.

⁷<https://pjreddie.com/darknet/tiny-darknet/>

⁸<https://github.com/ebennequin/FewShotVision>

3.3.4 First results and investigations

My attempts to build a working few-shot object detector are to this day unsuccessful. In this section, I will expose my observations and attempts to find the source(s) of the problem.

I launched a first experiment with a Deep Tiny YOLO initialized as explained in the previous section. It is trained on 3-way 5-shot object detection tasks on the COCO dataset. It uses an Adam optimizer with a learning rate of 10^{-3} (both in the inner loop and outer loop). It is trained for 10 000 epochs, each epoch corresponding to one gradient descent on the average loss on 4 episodes. During each episode, the model is allowed two updates on the support set before performing detection on the query set.

The loss is quickly converging (see Figure 14) but at inference time, the model is unable to perform successful detections (with a F1-score staying below 10^{-3}). Extensive hyperparameter tuning has been performed with no sensible improvement on the results.

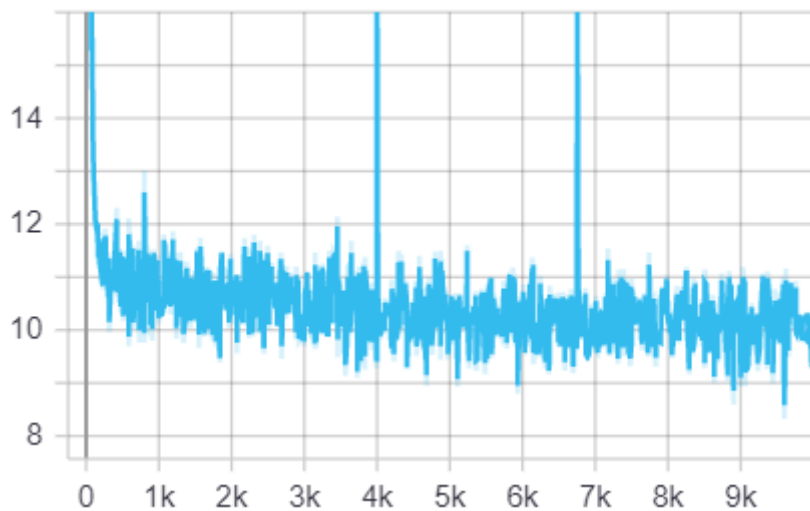


Figure 14: Total loss per epoch of YOLOMAML. Each point is the average total loss of the model on the query set of the episodes of one epoch.

To ensure that these disappointing performance was not due to my reimplemention of YOLO, I trained the Deep Tiny YOLO without MAML, in the same settings, for 40 epochs. Although this training is not optimal, the model is still able to perform relevant detections, which is not the case for YOLOMAML (see Figure 15).

The YOLOv3 algorithm aggregates three losses on three different parts of the predictions:

1. the shape and position of the bounding box of predicted objects, using Mean Square Error;
2. the objectness confidence (how sure is the model that there is truly an object in the predicted bounding box) using Binary Cross Entropy;
3. the classification accuracy on each predicted box, using Cross Entropy.



Figure 15: Object detection by the models YOLOMAML (left column) and YOLO (right column).

Figure 16 shows the evolution of these different parts of the loss. Loss due to objectness confidence has been further divided into two parts : the loss on boxes that contain an object in the ground truth, and the loss on boxes that do not contain an object in the ground truth.

We can see that the loss due to the classification and to the shape and position of the bounding box do not evolve during training. The no-object-confidence loss drops in the first thousand epochs before stagnating, while the yes-object-confidence rises to a critical amount before stagnating.

Figure 17 shows the same data for the training of YOLO. We can see that in this case, the yes-object-confidence drops after a peak in the first epochs. All parts of the loss decrease during the training, except the no-object-confidence, which reaches a floor value which is relatively small compared to the other parts.

Considering this, it is fair to assume that the bottleneck in training YOLOMAML is the prediction of the objectness confidence.

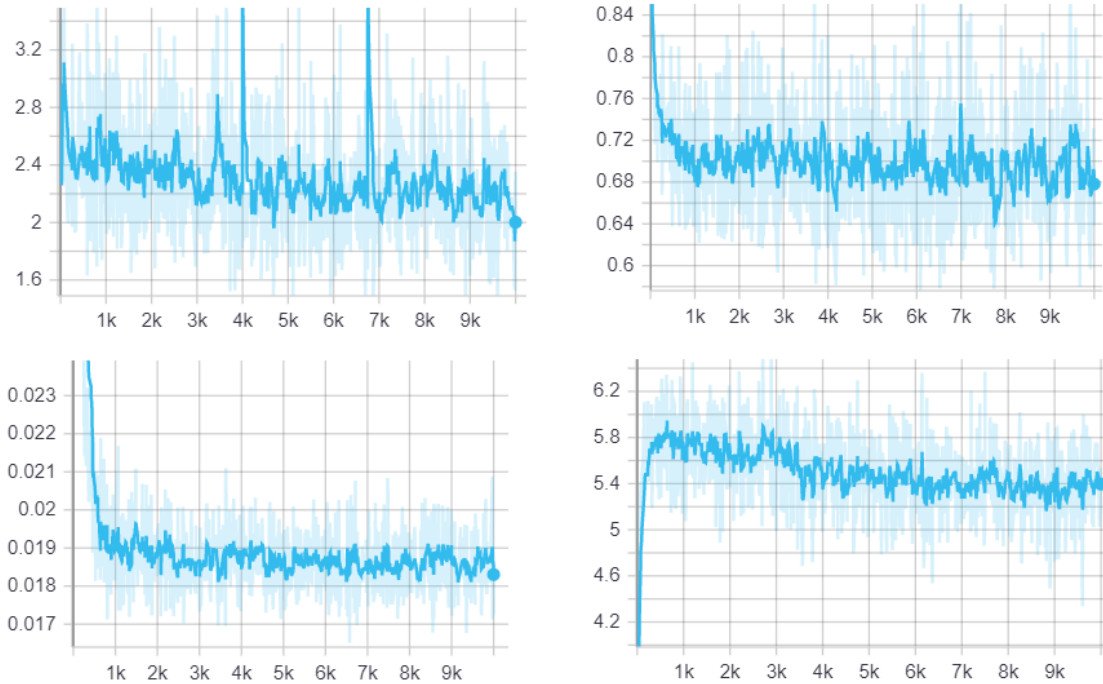


Figure 16: Evolution of the four parts of the loss of YOLOMAML during the same training as in Figure 14. Up-left: bounding box loss. Up-right: classification loss. Bottom-left: objectness confidence loss for boxes with no ground truth object. Bottom-right: objectness confidence loss for boxes with a ground truth object. Exponential moving average has been used to clearly show the patterns.

3.3.5 Future work

Unfortunately I did not have enough time to develop a working version of YOLOMAML. At this point I believe the answer resides in the prediction of the objectness confidence, but it is likely that other issues may rise when this one is solved.

An other direction of future work would be to constitute a dataset adapted to few-shot detection. Other works [28] [29] propose a split of the PASCAL VOC dataset adapted to few-shot detection. However, PASCAL VOC contains only 25 classes, while COCO contains 80 classes. I believe this makes COCO more adapted to meta-learning, which is entangled with the idea of learning to generalize to new classes.

Finally, a drawback of a (working) YOLOMAML would be that it does not allow way change, *i.e.* that a model trained on N -way few-shot detection tasks cannot be applied to a N' -way few-shot detection tasks. Solving this problem would be a useful improvement for YOLOMAML.

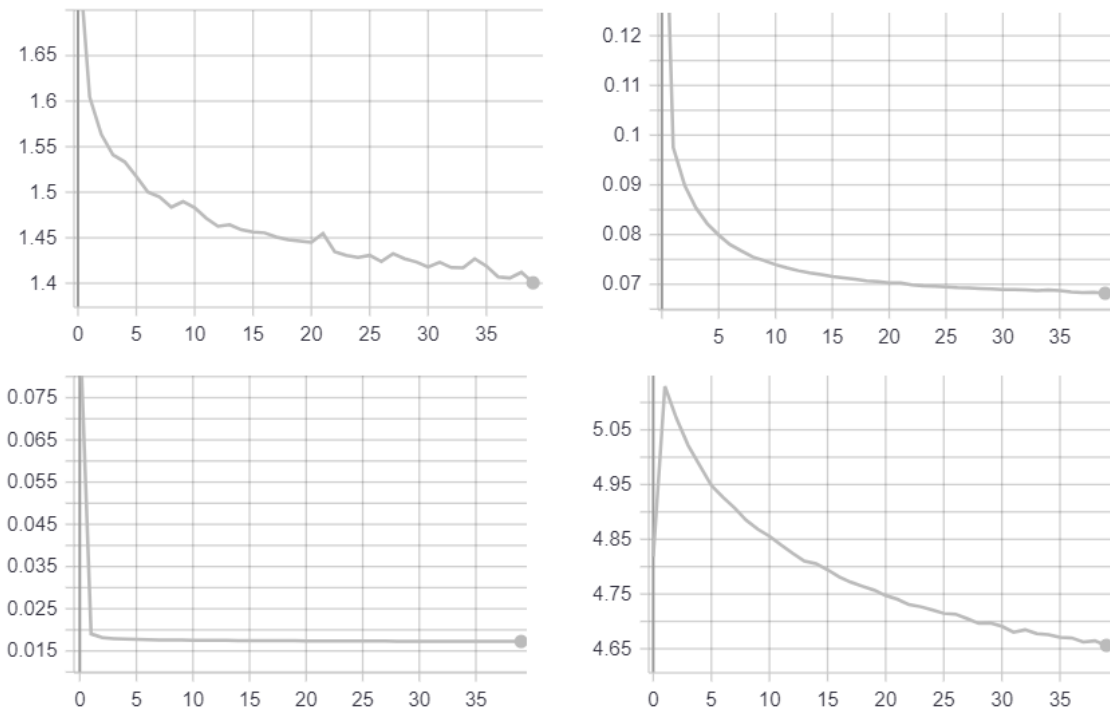


Figure 17: Evolution of the four parts of the loss of YOLO. Up-left: bounding box loss. Up-right: classification loss. Bottom-left: objectness confidence loss for boxes with no ground truth object. Bottom-right: objectness confidence loss for boxes with a ground truth object.

4 Conclusion

Advanced research in Few-Shot Learning is still young. Until now, only few works have tackled the few-shot object detection problem, for which there is yet no agreed upon benchmark (like mini-ImageNet for few-shot classification). However, solving this problem would be a very important step in the field of computer vision. Using meta-learning algorithms, we could have the ability to learn to detect new, unseen objects with only a few examples and a few minutes.

I am disappointed that I was not able to make YOLOMAML work during my internship at Sicara. However, I strongly believe that it is important to keep looking for new ways of solving few-shot object detection, and I intend to keep working on this.

References

- [1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [2] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *International conference on machine learning*, pages 1842–1850, 2016.
- [3] Tsendsuren Munkhdalai and Hong Yu. Meta networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2554–2563. JMLR. org, 2017.
- [4] Pablo Sprechmann, Siddhant M Jayakumar, Jack W Rae, Alexander Pritzel, Adria Puigdomenech Badia, Benigno Uria, Oriol Vinyals, Demis Hassabis, Razvan Pascanu, and Charles Blundell. Memory-based parameter adaptation. *arXiv preprint arXiv:1802.10542*, 2018.
- [5] Gregory Koch, Richard Zemel, and Ruslan Salakhutdinov. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, volume 2, 2015.
- [6] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, koray kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 3630–3638. Curran Associates, Inc., 2016.
- [7] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4077–4087. Curran Associates, Inc., 2017.
- [8] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip H.S. Torr, and Timothy M. Hospedales. Learning to compare: Relation network for few-shot learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [9] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1126–1135, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.
- [10] Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-sgd: Learning to learn quickly for few-shot learning. *arXiv preprint arXiv:1707.09835*, 2017.
- [11] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. 2016.

- [12] Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive meta-learner. In *International Conference on Learning Representations*, 2018.
- [13] Victor Garcia Satorras and Joan Bruna Estrach. Few-shot learning with graph neural networks. In *International Conference on Learning Representations*, 2018.
- [14] Bharath Hariharan and Ross Girshick. Low-shot visual recognition by shrinking and hallucinating features. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3018–3027, 2017.
- [15] Antreas Antoniou, Amos Storkey, and Harrison Edwards. Data augmentation generative adversarial networks. *arXiv preprint arXiv:1711.04340*, 2017.
- [16] Yu-Xiong Wang, Ross Girshick, Martial Hebert, and Bharath Hariharan. Low-shot learning from imaginary data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7278–7286, 2018.
- [17] Sebastian Thrun and Lorien Pratt. Learning to learn. 1998.
- [18] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [19] Brenden Lake, Ruslan Salakhutdinov, Jason Gross, and Joshua Tenenbaum. One shot learning of simple visual concepts. In *Proceedings of the annual meeting of the cognitive science society*, volume 33, 2011.
- [20] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- [21] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 580–587, June 2014.
- [22] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.
- [23] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [24] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [25] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.

- [26] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [27] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [28] Bingyi Kang, Zhuang Liu, Xin Wang, Fisher Yu, Jiashi Feng, and Trevor Darrell. Few-shot object detection via feature reweighting. *arXiv preprint arXiv:1812.01866*, 2018.
- [29] Kun Fu, Tengfei Zhang, Yue Zhang, Menglong Yan, Zhonghan Chang, Zhengyuan Zhang, and Xian Sun. Meta-ssd: Towards fast adaptation for few-shot object detection with meta-learning. *IEEE Access*, 7:77597–77606, 2019.
- [30] Tao Wang, Xiaopeng Zhang, Li Yuan, and Jiashi Feng. Few-shot adaptive faster R-CNN. *CoRR*, abs/1903.09372, 2019.
- [31] Hao Chen, Yali Wang, Guoyou Wang, and Yu Qiao. Lstd: A low-shot transfer detector for object detection. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [32] Xuanyi Dong, Liang Zheng, Fan Ma, Yi Yang, and Deyu Meng. Few-shot object detection. *ArXiv*, abs/1706.08249, 2017.
- [33] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-UCSD Birds 200. Technical Report CNS-TR-2010-001, California Institute of Technology, 2010.
- [34] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. *arXiv preprint arXiv:1904.04232*, 2019.
- [35] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. Emnist: an extension of mnist to handwritten letters. *arXiv preprint arXiv:1702.05373*, 2017.
- [36] Luca Bertinetto, Joao F Henriques, Philip HS Torr, and Andrea Vedaldi. Meta-learning with differentiable closed-form solvers. *arXiv preprint arXiv:1805.08136*, 2018.
- [37] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- [38] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.