

---

# The AutoDL challenge 2019-2020 (automl\_freiburg video team)

---

Philipp Jankov Christopher Krolla Thomas Nierhoff

## Abstract

This documentation summarizes the progress of the automl\_freiburg video team on the AutoDL 2019 – 2020 challenge. Two methods are presented: The first method relies on transfer learning of a single, pretrained video classification network which is fine tuned during the challenge. As this approach has shown several limitations if the test datasets vary too much, a more elaborate method based on meta learning is developed: Multiple classification networks are evaluated on different training datasets a priori and the incumbent for every dataset is stored in a lookup table. During test time, the similarity between the test dataset and all training datasets is measured and the incumbent of the most similar training dataset is used for training/evaluation on the test dataset. The first method has been submitted and evaluated on the AutoCV2 challenge and resulted in a mediocre score. Four different approaches are presented and evaluated for the second method, resulting in scores on par with the best performing image and video classification team.

## 1. Introduction

Machine learning (ML) based on artificial neural networks is a fast-growing field in the domain of computer science. Instead of manually tweaking the features used for classification or regression, neural networks are able to extract them automatically. They have shown state-of-the-art performance on a variety of classification problems when trained on huge datasets and are well suited for parallelization.

A open question is how to deal with computational limitations and generic datasets. In many real-world applications (e.g. smartphones) computational resources, execution time and storage capacity is limited. Thus the current trend of using larger and larger networks for better classification results becomes unfeasible. In addition, it may be not known a priori what the data to be classified looks like, requiring the neural network to be applicable to a variety of different datasets. Multiple approaches emerged within the last years to tackle both challenges:

To reduce the model size (and thus implicitly also execution time, storage capacity and computational resources), one can remove the most “unimportant” parts of a neural network and/or quantize it, resulting in a compression of the neural network. The selection of the nodes/layers to be pruned can either be specified manually or derived automatically through reinforcement learning (Han et al., 2015; He et al., 2018). A second option is knowledge distillation where a smaller student “network” is trained to produce outputs similar to a previously trained larger “teacher” network (Hinton et al., 2014; Phuong & Lampert, 2019). By using the logits output instead of the softmax output, much more information is preserved in every training step. Consequently training the student network can be achieved with less samples. For both approaches there is only a negligible decrease in classification accuracy while the compression rate can be huge (up to a factor of 496x for certain layers of a neural network (Reagan et al., 2018)). Another more straightforward method is to search directly for smaller/faster models. This can be either done manually - see e.g. (Howard et al., 2017) - or through neural architecture search (Elsken et al., 2019). Note that the last approach can output networks that are state of the art with respect to model size, speed and classification accuracy (Tan & Le, 2019).

Other techniques are used if the time budget for training is tight or the dataset to be trained on is not known a priori. They all have in common that a network is pretrained on one or more datasets and information from pretraining is subsequently used for the final training step. One technique is transfer learning (Pan & Yang, 2010). Here the common use case consists of pretraining a network on a large dataset to initialize the network weights with reasonable values. If a new dataset is similar enough, only parts of the network (e.g. the last layers of a convolutional neural network) must be trained again. In addition, the new dataset can be rather small as less optimization steps are required (the network is not trained from scratch). Another option is meta learning or learning to learn (Hutter et al., 2018; Santoro et al., 2016): By learning multiple models for multiple tasks (or more specific: datasets) beforehand, it is quite likely that a new task will be similar to a previously learned task. In this case, knowledge gained from the previous task can be reused, thus accelerating the learn-

ing progress. The challenge with meta learning is a proper selection of the meta features to be extracted as their predictive capabilities vary heavily (Bilalli et al., 2017). One predominant meta learning approach at the moment is MAML (Finn et al., 2017), optimizing a model not for classification accuracy but specifically for good generalization to unseen datasets from the same dataset distribution.

The AutoDL challenge contains both aspects, tight computational limitations (both time- and resource-wise) and generic datasets. As such it provides an ideal testbed for advanced approaches in order to find time-efficient algorithms with a small memory footprint. Because the datasets are undisclosed, the options for manual tweaking are limited. Instead of highly specific solutions tied to a specific environment one is thus seeking for a fast yet robust approach. This makes the challenge different from “traditional” benchmarks like ImageNet classification where computing power is solely limited by the budget of the chair/institute and the dataset is known beforehand.

The remainder of the paper is organized as follows: Section 2 and section 3 provide general background information about the AutoDL challenge and an overview of the used datasets. A first method used during the AutoCV2 video challenge is presented in section 4. As focus shifted after the video challenge, section 5 illustrates two more elaborate approaches for video and image classification that can be used for the final AutoDL challenge. The documentation concludes with a short discussion in section 6 and a summary with possible future directions in section 7.

## 2. AutoDL challenge

This section gives a quick overview of the AutoDL 2019 – 2020 challenge: The AutoDL 2019-2020 challenges consist of a set of different, individual challenges depending on the underlying dataset modality: Image (AutoCV), video (AutoCV2), text (AutoNLP), speech (AutoSpeech), weakly supervised learning (AutoWSL) and all previously mentioned modalities combined (AutoDL). All challenges have in common that teams can upload code together with models during an online phase where it is tested immediately on a number of datasets. This gives every team immediate feedback about their relative performance. These online phases last between two weeks and two months. Some challenges (e.g. for video) have a second, full blind testing (on-site) phase at the end where every team’s last submission is evaluated on a new, previously unknown set of datasets.

Differing from other challenges, there are tight constraints: Rather than focusing on the final classification accuracy after a certain time limit, each team can choose freely how much time to spend on training and how much on testing within a specific time budget  $T = 1200s$ . Training

is performed by running a “train” function with the train dataset as input, test by running a “test” function with the test dataset as input and the classification results as output. Based on the output of every test run its score  $s(t)$  at time  $t$  is calculated as

$$s(t) = 2 * AUC - 1,$$

with  $AUC$  as the area under receiver operating characteristic curve. The final score  $ALC$  is then calculated as the time integral over all individual scores  $s(t)$ , weighted on a logarithmic time scale as

$$ALC = \frac{1}{\log\left(1 + \frac{T}{t_0}\right)} \int_0^T \frac{s(t)}{t + t_0} dt,$$

with  $t_0 = 20s$ . This way, the predictions of every test run at time  $t$  are weighted with  $\frac{1}{t+t_0}$ . Because early predictions are weighted more, focus shifts from a high classification accuracy at the end of the time budget to a high classification accuracy as early as possible. A example of  $ALC$  can be seen in Fig. 1. It is visible that the performance within the first 60s counts roughly as much as the performance within the last 600s. Additional limitations

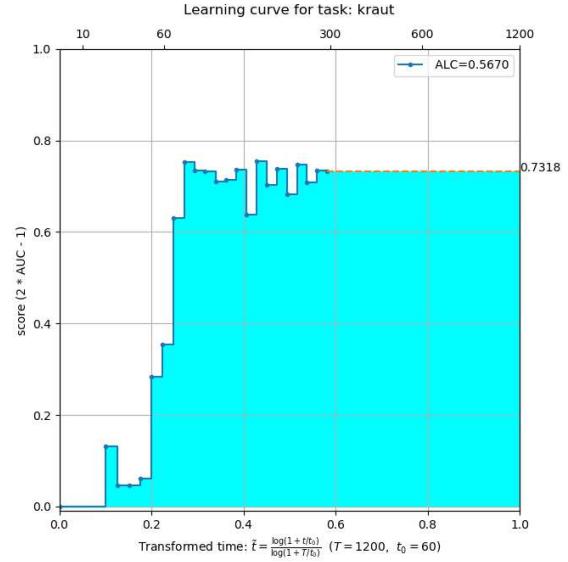


Figure 1. Visualization of the ALC plotted over time

e.g. for the AutoCV2 competition consider memory constraints (max. 300Mb per submission), computational constraints (NVIDIA Tesla P100 1GPU, 4vCPU, 26GB) and time constraints (max. 5 hours runtime per 24 hours, max. 5 submissions per 24 hours).

image dataset name	train samples	test samples	classes	source
Binary Alphadigits	1115	289	36	tfds
Caltech101	3059	6085	102	tfds
CUB-200	3000	3033	200	tfds
CUB-200-2011	5994	5794	200	tfds
Cats and Dogs	22263	999	2	tfds
CIFAR-10	50000	10000	10	tfds
CIFAR-100	50000	10000	100	tfds
COIL-100	6202	998	72	tfds
Colorectal Histology	4003	997	8	tfds
DeepWeeds	16520	989	5	tfds
EuroSAT	25976	1024	10	tfds
EMNIST	697932	116323	62	tfds
Fashion-MNIST	60000	10000	10	tfds
Food-101	75750	25250	101	tfds
Horses or Humans	1027	256	2	tfds
ImageNet	1281167	100000	1000	tfds
KMNIST	60000	10000	10	tfds
MNIST	60000	10000	10	tfds
Oxford-102 Flower	1020	6149	102	tfds
Oxford-IIIT Pet	3680	3669	37	tfds
PatchCamelyon	262144	32768	2	tfds
Rock-Paper-Scissors	2520	372	3	tfds
THE small NORB	24300	24300	5	tfds
SVHN cropped	73257	26032	10	tfds
tf flowers	2938	732	5	tfds
UC Merced	1678	422	21	tfds
Chucky	48061	11939	100	AutoDL
Decal	634	166	11	AutoDL
Hammer	8050	1965	7	AutoDL
Munster	60000	10000	10	AutoDL
Pedro	80095	19905	26	AutoDL

video dataset name	train samples	test samples	classes	source
Katze	1528	863	6	AutoDL
Kreatur	1528	863	4	AutoDL
Kraut	1528	863	4	AutoDL
EPIC-Kitchens	24699	3773	125*331	inet
HMDB51	3570	500	51	inet
JHMDB21	658	270	21	inet
Kinetics 400	238831	19675	400	inet
something-something v2	168913	1036	174	inet
UCF101	9537	3783	101	inet
Youtube8m	1459266	216409	3862	own
YFCC100m	521035	130258	1570	own

Table 1. Dataset overview

### 3. Datasets

Table 1 lists all available datasets used for training our models. Shown is the dataset name both for image and video datasets, the number of train samples, test samples and classes and the source from where the dataset can be obtained. Datasets marked with “tfds” can be downloaded using the Tensorflow Datasets module, the ones marked with “AutoDL” using the provided code from the AutoDL challenge and the ones with “inet” are publicly available. Both the Youtube8m and YFCC100m have been downloaded partially by using/extracting the individual video links in/from the dataset file. Note that these two datasets are also the only two multilabel datasets. Although looking promising in theory as they are both the largest video datasets available and the only multilabel datasets, Youtube8m and YFCC100m were impossible to use due to their sheer size. Even if the exact reason is unknown (possible inode overhead due to too many files), training any neural network on them was almost a magnitude slower compared to smaller datasets. Furthermore, results from the AutoCV2 challenge indicate that it is unnecessary to capture multiple frames from a video (one is sufficient), thus turning the video classification task into a image classification task.

### 4. Image/video classification using transfer learning

The approach used for the AutoCV2 challenge is described in this section. Focusing solely on video classification, the AutoCV2 challenge consists of 5 publicly available datasets that can be downloaded during the online phase, 5 disclosed datasets used for the preliminary ranking during the online phase and another 5 disclosed datasets for the final evaluation during the on-site phase.

The general structure of most video classification networks is bipartite. Given a video  $V$  consisting of individual frames  $v_i$  as  $V = [v_1, v_2, \dots, v_n]$ , a subset of frames  $V' = [v'_1, v'_2, \dots, v'_m] \subseteq V$  is selected. Typical selection choices are random, random within consecutive video segments or random with a bias towards the mid of the video. Every selected frame is then processed individually by a image classification network  $f(\cdot)$  (e.g. a ResNet or InceptionNet) as  $f(v'_i)$ . The results of all processed frames - either on logits or feature level - is fed to a second network  $g(\cdot)$  learning the temporal aspect between the frames. The classification result  $c$  is then calculated as

$$c = g(f(v'_1), f(v'_2), \dots, f(v'_m)). \quad (1)$$

Both the 300Mb limit for each submission and the logarithmic time scale impose tight constraints. As such focus is on fast training rather than a superb classification rate. Hence EfficientNet (Tan & Le, 2019) is chosen as primary network for the image classification  $f$ . Besides achieving state-of-the-art results on ImageNet, it can be almost arbitrarily scaled down for a better performance-accuracy trade-off and is almost a magnitude smaller and faster than comparable networks.

Different networks are evaluated for selecting the network  $g$ . In the end, a simple averaging layer shows sufficient performance while being fastest to train. Eq. (1) thus simplifies to

$$c = \frac{1}{m} \sum_{i=1}^m f(v'_i) \quad (2)$$

The averaged ALC scores based on three runs together with final ranks for the five datasets of the on-site phase are shown in Tab. 2. In total there are 20 participating teams. It is interesting to see how the ranks vary from first place for dataset 4 to last place for dataset. We assume that dataset 4 is an average to difficult video classification dataset - this is what our network has been optimized for. Conversely the used network is rather slow compared to other approaches like SVMs or very small networks. This causes problems if the test dataset is large or easy to classify.

dataset	avg. $ALC$ score	rank (of 20)
dataset 1	0.1836	20
dataset 2	0.9138	2
dataset 3	0.4009	10
dataset 4	0.5169	1
dataset 5	0.1031	14

Table 2. AutoCV2 results

## 5. Image/video classification using meta learning

The AutoCV2 challenge has shown severe limitations of the transfer learning approach using only one model, specifically low adaptability to changed conditions (e.g. very easy/difficult classification task or large test dataset). Thus different approaches based on meta learning are presented in this section. They all have in common that a classifier (from now on called meta classifier) is trained to measure the similarity between any new dataset (i.e. test dataset) and a set of known datasets (training datasets) using meta features. It shall be noted that all training and test datasets contain a train/test split used for training and evaluation of the models on every individual dataset. But as we are looking for good generalization across different datasets, the train/test set now refers primarily to datasets instead of samples within the dataset.

For every training dataset an optimized neural network has been selected a priori. Then the neural network of the most similar training dataset will be used for transfer learning on the test dataset. However, they differ on the type of meta features and the classifier that determines the similarity of the test dataset to the set of training datasets. The approaches belong to the larger class of non-parametric meta learning and are conceptually similar with (Snell et al., 2017; Sung et al., 2018; Vinyals et al., 2016) in the sense that they first learn an embedding of the data and then measure the similarity on the embedding space. The difference is that our principal goal is to distinguish datasets and not samples of different classes within one dataset. To describe our approach formally, we follow mostly the notation in (Hutter et al., 2018), chapter 2. As we are solely interested in dataset classification, every task in (Hutter et al., 2018) refers to a dataset classification problem:

For all presented approaches the dataset to be learned is named  $t_j \in T$  and contained in the set of all training datasets. Each dataset is described by a meta feature vector  $\mathbf{m}(t_j)$ . The neural network and its corresponding (hyper-)parameters are defined by a configuration  $\theta_i \in \Theta$  where  $\Theta$  represents the entire configuration space (discrete, continuous or mixed). In addition, there are offline evaluations  $P(\theta_i, t_j) \in \mathbf{P}$  from the set of all prior evaluations  $\mathbf{P}$ . Assuming that we are given a test dataset  $t_{test}$ , we

want to specify how similar it is to the training datasets  $t_j$  and derive an optimized configuration  $\theta_{test}^*$  to maximize  $P(\theta_{test}^*, t_{test})$ .

It is assumed that an optimized configuration  $\theta^*(t_j)$  is found for every training dataset  $t_j$  during the offline phase. As all datasets are represented by their meta feature vectors  $\mathbf{m}(t_j)$ , we want to find the most similar dataset  $t_{sim}$  for a test dataset  $t_{test}$  as

$$t_{sim} = \operatorname{argmax}_{t_j \in T} \|\mathbf{m}(t_{test}) - \mathbf{m}(t_j)\|, \quad (3)$$

and then use the configuration  $\theta^*(t_{sim})$  for  $t_{test}$ .

Two challenges arise from Eq. (3):

- Finding expressive meta features  $\mathbf{m}(t_j)$  for every dataset
- Finding a reasonable distance measure  $\|\cdot\|$

Expressive meta features are essential for distinguishing different datasets: If different datasets result in similar meta feature vectors they cannot be distinguished, regardless of the used distance measure. At the same time the variation of the meta features shall be roughly proportional to the variation of the datasets: An intuitive counterexample is the use of a generic hash function as meta feature, which leads to large variations even if two datasets are almost similar. The same applies for the distance measure: Having found good meta features, almost identical datasets shall have only a small distance whereas dissimilar datasets shall result in a large distance.

The next sections describe how both challenges are tackled by different approaches.

### 5.1. Approach 1

The first class of presented meta classifiers relies on the output of an image classification network as meta features. Being trained on a sufficiently large dataset with enough features/classes (corresponding to the output of the second last or last network layer), the distribution of the detected features/classes acts as a unique “fingerprint” for each dataset and constitutes the meta feature vector  $\mathbf{m}(t_j)$ .

Different classifiers are used to find a reasonable distance metric based on the output of the image classification network, see Fig. 2. By assigning an individual label to every dataset (note that every dataset represents a class), the distance measure is learned implicitly.

The first approach (which also serves as baseline for the experiments) relies on k-nearest centroids classification with a median as distance metric to be more robust against single outliers. It clusters the datasets according to their dataset



class label and calculates the centroid of each cluster. At test time it calculates the centroid of the test dataset and outputs class label of the dataset with the closest centroid. It also uses majority voting for more reliable predictions based on multiple samples (i.e. batch sizes  $> 1$ ).

The second approach is based upon a multi-layer perceptron (MLP) with two layers with a Swish activation function, dropout and batch normalization. Similar to the baseline approach, majority voting in combination with large batch sizes is used to make predictions less susceptible to outliers.

The third approach uses the same two-layer perceptron as the second approach. However, instead of relying on majority voting based on multiple samples per batch, the output of the image classification network is first fed to a precalculation layer. This layer calculates different p-quantile moments and related measures (mean, median, standard deviation, variance, skewness, ...) which are then fed to the two-layer perceptron. Note that the moments are calculated across the batch size, i.e. when calculating  $m$  moments for a batch size  $b$  and  $f$  meta features, the output of the precalculation layer is a vector with length  $m \cdot f$ .

A visual overview of the different approaches together with typical implementations is given in Fig. 2. Here, a ResNet18 is used as image classification network. It also shows the output dimensions of the different components for a batch size of 128, an input image size of  $96 \times 96$  pixels with three channels, a 512-dimensional ResNet18 output (the second-last layer), the use of five modes (e.g. mean, standard deviation, variance, skewness, kurtosis) resulting in the 2560-dimensional ( $512 \times 5$ ) output of the precalculation layer and 20 datasets to be distinguished.

## 5.2. Approach 2

The second approach uses the latent space of a variational autoencoder (VAE) as meta features. By pretraining the VAE on a large variety of images, similar images belonging to one dataset cluster together in latent space whereas dissimilar images from different datasets have a large (Euclidean) distance in latent space. When measuring both the mean and standard deviation of a batch of images in latent space for the test dataset and batches of images of the different training datasets, a two-sampled t-test is used to determine the “nearest” dataset.

The VAE consists of a DenseNet 121 (Huang et al., 2017) for the encoder and a PixelCNN++ (Salimans et al., 2017) for the decoder, see Fig. 3. A fully connected layer maps the 196-dimensional latent space (i.e. the logits output of the DenseNet) to the  $64 \times 64 \times 3$ -dimensional input of the PixelCNN++ network. Differing from the original DenseNet implementation, the first layer has been adjusted for the

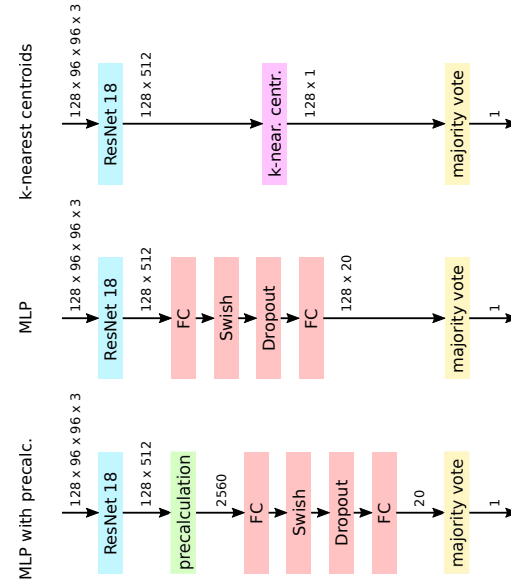


Figure 2. Overview of the different approaches based on an image classification network with typical implementation parameters

smaller image size to a  $3 \times 3$  convolution with stride 1, no padding and no pooling after it. The VAE trained on seven different datasets (Caltech101, Colorectal Histology, DeepWeeds, EuroSAT, Food-101, ImageNet) whereas the last one has been resized to  $64 \times 64$  pixel for faster calculation. These datasets have shown to yield a large variety of different images which are supposed to generalize well to other datasets.

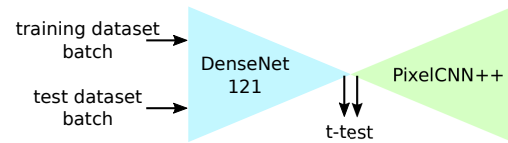


Figure 3. Overview of the VAE approach based on an image classification network for the encoder and a generative image network for the decoder.

## 5.3. Comparison of the two approaches

While the two approaches of section 5.1 and 5.2 are conceptually similar as they first learn a suitable embedding and then classify different datasets based on the information from the embedding, they have also some subtle differences:

**Number of parameters:** Both approaches are very small in terms of model parameters compared to current state-of-the-art image classification networks: Whereas a Resnet18 with the subsequent MLP with precalculation has around 12M parameters, the VAE has only around 7m parameters.

**Pretraining:** Both models can be pretrained well on a single dataset if it contains a large enough variety of images (e.g. ImageNet). Still, being a well-proven architecture a ResNet18 is probably even easier to train than a VAE. Moreover, there are multiple sources with well pretrained Resnet18 models, making manual pretraining obsolete.

**Adding more training datasets:** Unless the additional training datasets differ significantly from the ones used for pretraining, both the Resnet18 of the first approach and the VAE of the second approach can be kept as-is. However, the first approach requires pretraining of the MLP from scratch whenever a training dataset is added.

**Dependence on training datasets during test time:** Because the first approach encodes the relation between different training datasets in the MLP, it requires only test data during test time. On the contrary, the second approach compares a batch from the test dataset with batches from all training datasets during test time using a two-sample t-test. Thus it needs access to the entire training datasets or at least the latent representation (mean and standard deviation) of a few selected batches of every training dataset during test time.

**Dataset class classification:** Whereas the first approach is restricted to dataset classification, the second approach also allows for a more fine-grained classification: By comparing a batch of the test dataset to batches sampled from the same class of a training dataset, one can find not only the most similar dataset but also the most similar class of a dataset.

## 5.4. Experiments

BOHB (Falkner et al., 2018) is used for the a priori calculation of optimized configurations  $\theta^*(t_j)$  for different training datasets  $t_j$ . Every configuration consists of a pretrained network (either trained on ImageNet or CIFAR-100) and the following hyperparameters (if applicable):

- Optimizer: ADAM or SGD (with Nesterov momentum 0.9 and weight decay 1e-6)
- Train batch size: 16/32/64/128
- Learning rate: 1e-5 - 0.5
- Dropout 0.01 - 0.99

An overview of the used datasets and models is given in Fig. 4 and Fig. 5. The two figures show the best classification accuracy of the different models on the different datasets, both from a model-based perspective and a dataset-based perspective. The hyperparameters for every combination of dataset and model are optimized through 10 BOHB runs. Note that not all datasets listed in Tab. 1 appear in Fig. 4. Some are used as test datasets whereas

others (especially video datasets) are too large to be processed efficiently with the given resources.

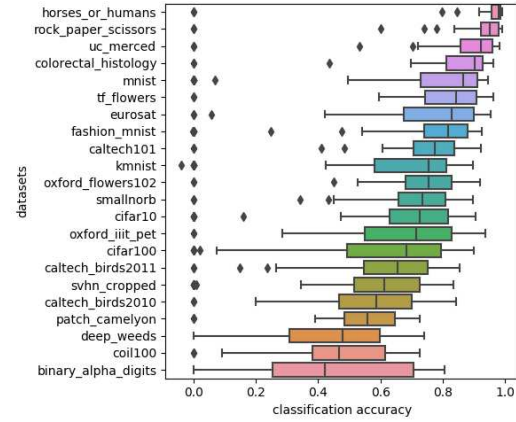


Figure 4. Classification accuracy of the different models across all datasets, sorted by dataset

In addition, multiple optimization runs find an optimal meta feature vector and a reasonable distance measure while keeping the configuration space  $\Theta$  small: Being reasonably small and fast, a ResNet18 network pretrained on ImageNet serves as basis to generate the meta feature vector  $\mathbf{m}(t_j)$ . However, using the output of the second last layer yields consistently better results than the last output and is therefore chosen.

The following parameters have been optimized for the different approaches using BOHB:

MLP:

- Train batch size: 1/2/4/8/16/32/64/128/256/512/1024
- Neurons per layer: 16/32/64/128/256/512/1024
- Learning rate: 1e-5 - 1e-2
- Dropout rate: 0 - 0.8

MLP with precalculation:

- Train batch size: 1/2/4/8/16/32/64/128/256/512/1024
- Neurons per layer: 16/32/64/128/256/512/1024
- Learning rate: 1e-5 - 1e-2
- Dropout rate: 0 - 0.8
- p-quantile: 0 - 0.4
- (Do not) use mean, median, standard deviation, variance, skewness, kurtosis

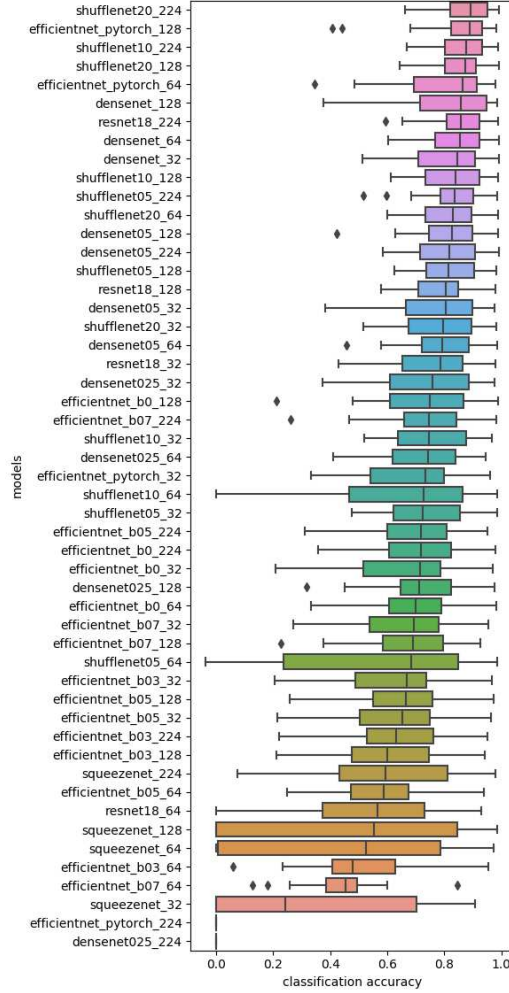


Figure 5. Classification accuracy of the different models across all datasets, sorted by model

Final results of the different approaches are shown in Fig. 6. The eight plots correspond to the eight test datasets (Chucky, Decal, Hammer, Munster, Pedro, Kraut, Katze, Kreatur). Every plot shows the test ALC according to Eq. (1) over the test batch size of the classification networks of section 5.1 and 5.2. The different thick solid lines correspond to:

- **kakaobrain**: ALC results from the winning AutoCV and AutoCV2 team kakaobrain (Baek et al., 2020).
- **same-config**: ALC of the optimized configuration  $\theta^*(t_s)$  for the same dataset  $t_s$ . Can be seen as ground truth.
- **precalc**: Third approach of section 5.1 with individual networks trained specifically for every test batch size.

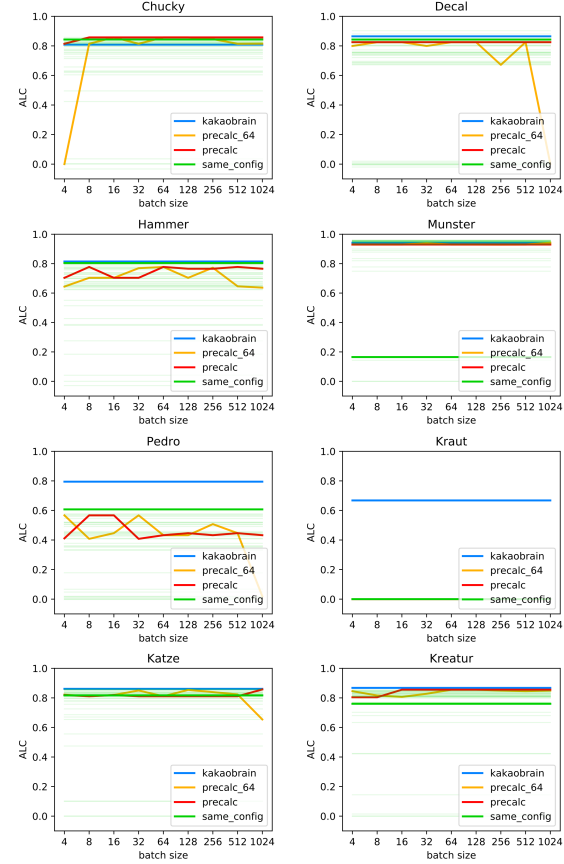


Figure 6. Final classification results

- **precalc-64**: Third approach of section 5.1 with one network trained for a test batch size of 64 samples.

There are also multiple horizontal thin green lines in every plot. Each one corresponds to the ALC of the optimized configuration  $\theta^*(t_p)$  of a prior dataset  $t_p$ , evaluated on the test dataset. Thus, a plot ideally has the following appearance: Kakaobrain has the highest ALC as we do not expect to perform better than a team of computer vision professionals. Next comes the ALC of the optimized configuration for the corresponding dataset (same-config). Note that this measure is usually not available and only added for comparative purpose. The individual ALCs of the different presented approaches are expected to score worse than the previous two as they rely on optimized configurations of different datasets (and not the test dataset).

There are several effects why the plots in Fig. 6 differ from the ideal plot described before:

Sometimes an optimized configuration has a better ALC than kakaobrain. A possible reason is the different ALC calculation: Due to heavily varying execution times on the cluster (up to a factor of 10) we did not measure the real

ALC. Rather, the execution time of an individual forward and backward pass for every model for images/videos of different input shapes has been measured on a reference system and been stored in a lookup table. The lookup table is then used on the cluster to estimate the ALC based on the number of batches, the number of samples per batch and the average shape of the images/videos of the dataset. This procedure allows one to be independent of the cluster execution time while having a good approximation of the true ALC but may overestimate the true ALC.

In addition, the same-config ALC is sometimes worse than the ALCs of optimized configurations of prior datasets. We are not sure about the exact cause. One reason can be the limited number of BOHB evaluations to find the optimized configuration  $\theta^*(t_s)$  for the same dataset. Only 90 different configurations (10 BOHB runs,  $\text{eta}=3$ ,  $\text{max budget} = 9 * \text{min budget}$ ) are evaluated for every dataset. Thus, with more than 20 other training datasets used, it is possible that one of these 20 optimized configurations scores better. Another option is the ALC variation during test time. Hence an optimized configuration which scored rather well during training may score considerably worse during testing (even on the same dataset).

For some datasets (e.g. Munster) the ALC of a presented approach is similar for multiple/all test batch sizes. Although this is the preferred case (the output of the meta classifier is ideally independent of the batch size), the effect is usually caused by the test dataset being similar to one of the training datasets. Then a well-trained meta classifier always correctly classifies the test dataset being most similar to the (identical) training dataset.

Choosing an optimal test batch size is crucial during test time. When being too small, the meta classifier may produce noisy (and thus possibly incorrect) predictions. When being too large, more samples than necessary are processed. This leads to computational overhead.

## 6. Discussion

Results of the AutoCV2 challenge have shown that it is not sufficient for any non-CV expert to rely on a single video classification network. Thus the meta learning approach with a portfolio of different classification networks has been developed. Although being effective, finding incumbent configurations for more than 20 datasets for more than 50 models is time-consuming and took around two weeks using 40 GPUs. In addition, Fig. 5 indicates that it is sufficient to use a much smaller portfolio of e.g. only the ten models with the average best performance. A small experiment leads to the result in Fig. 7 where the optimized averaged best performance across all datasets is plotted over the model portfolio size. With averaged best perfor-

mance we are referring to the best performance of all the models in the portfolio for every dataset, averaged over all datasets. The plot indicates that it is sufficient to use as few as five carefully selected models to achieve an almost similar performance as with all 52 models.

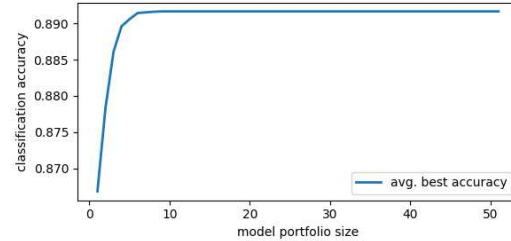


Figure 7. Averaged best performance of a reduced model portfolio across all datasets

A more substantial critic concerns the format of the AutoDL challenge: By using the ALC score (more specifically: the logarithmic time scale), submissions with very good results from the beginning on are favored. Albeit this is desired by the organizers, it became more an engineering problem than an interesting research challenge to optimize the data loading pipeline, strip state-of-the-art models to the bare minimum and evaluate how long they produce reasonable results. This is backed by the observation that the winners of both the AutoCV and AutoCV2 challenge used rather simple base models (Resnet9 and Resnet18) and their superior performance is the result of the combination of many small, existing tweaks instead of a completely new approach.

Using the wall clock time to calculate the ALC score was an major issue for reproducibility and portability across multiple systems. As such the final results for the AutoCV2 challenge were calculated as the average of three individual runs. The wall clock time issues was worst on the computing cluster in Freiburg. The combination of a high data throughput when classifying image/video datasets and the use of HDDs led to a huge variance in execution speed - sometimes up to a factor of 10. Thus the workaround with the execution time lookup table described in Sec. 5.4 has been chosen.

In the end we were suprised how conceptually simple our approach was. Distinguishing different datasets seems to be an easier problem than distinguishing different classes within a dataset. As such the simplest approach (Resnet18 + k-nearest centroids) provides surprisingly good results. However, this simplicity comes with the drawback of a large computational load when computing the incumbents for the different training datasets.



## 7. Summary and Outlook

We present two conceptually different methods for quick generalization of a pretrained network to new datasets in the context of the AutoCV2 and AutoDL challenge. Whereas the first method relies on transfer learning (i.e. pretraining a model and then fine tuning it on the new dataset), the second approach considers a portfolio of different models and selects the model that performed best on a dataset that is most similar to the new dataset. Thus focus shifts from finding a single state-of-the-art model working well on a variety of different datasets to measuring the similarity between two datasets and selecting the one that is most similar - the latter being much easier to implement for anyone who is not an expert in video classification.

There are several directions for future work: The first one considers redundancy in the datasets. It is shown in the discussion chapter how the used models are highly redundant and the model portfolio can be reduced by a factor of 10 without severe performance loss. It is unknown (but suspected) if the same applies for the datasets, i.e. if training the models on a small but carefully selected set of datasets leads to the same generalization performance.

Furthermore, the VAE approach of Sec. 5.2 is a promising idea for future expansions. It is an open question if more advanced classification methods using the latent space representation (e.g. neural networks) lead to better classification results than a two-sample t-test.

## References

- Baek, W., Kim, I., Kim, S., and Lim, S. Autoclint: The winning method in autocv challenge 2019. *arXiv*, 2020.
- Bilalli, B., Abelló, A., and Aluja-Banet, T. On the predictive power of meta-features in openml. *AMCS*, 27: 697–712, 2017.
- Elsken, T., Metzen, J. H., and Hutter, F. Neural architecture search: A survey. In *JMLR*, 2019.
- Falkner, S., Klein, A., and Hutter, F. BOHB: Robust and efficient hyperparameter optimization at scale. In *ICML*, 2018.
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017.
- Han, S., Mao, H., and Dally, W. J. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. In *ICLR*, 2015.
- He, Y., Lin, J., Liu, Z., Wang, H., Li, L., and Han, S. AMC: automl for model compression and acceleration on mobile devices. In *ECCV*, 2018.
- Hinton, G., Vinyals, O., and Dean, J. Distilling the knowledge in a neural network. In *NIPS*, 2014.
- Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., and Adam, H. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv*, 2017.
- Huang, G., Liu, Z., van der Maaten, L., and Weinberger, K. Q. Densely connected convolutional networks. In *CVPR*, 2017.
- Hutter, F., Kotthoff, L., and Vanschoren, J. (eds.). *Automatic Machine Learning: Methods, Systems, Challenges*. Springer, 2018.
- Pan, S. J. and Yang, Q. A survey on transfer learning. *IEEE Trans. on Knowl. and Data Eng.*, 22(10):1345–1359, October 2010. ISSN 1041-4347. doi: 10.1109/TKDE.2009.191.
- Phuong, M. and Lampert, C. Towards understanding knowledge distillation. In *ICML*, 2019.
- Reagan, B., Gupta, U., Adolf, B., Mitzenmacher, M., Rush, A., Wei, G.-Y., and Brooks, D. Weightless: Lossy weight encoding for deep neural network compression. In *ICML*, 2018.
- Salimans, T., Karpathy, A., Chen, X., and Kingma, D. P. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. In *ICLR*, 2017.
- Santoro, A., Bartunov, S., Botvinick, M., Wierstra, D., and Lillicrap, T. Meta-learning with memory-augmented neural networks. In *ICML*, 2016.
- Snell, J., Swersky, K., and Zemel, R. Prototypical networks for few-shot learning. In *NIPS*, 2017.
- Sung, F., Yang, Y., Zhang, L., Xiang, T., Torr, P. H. S., and Hospedales, T. M. Learning to compare: Relation network for few-shot learning. In *CVPR*, 2018.
- Tan, M. and Le, Q. EfficientNet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019.
- Vinyals, O., Blundell, C., Lillicrap, T., Kavukcuoglu, K., and Wierstra, D. Matching networks for one shot learning. In *NIPS*, 2016.