# A Wholistic View of Continual Learning with Deep Neural Networks: Forgotten Lessons and the Bridge to Active and Open World Learning

Martin Mundt, Yong Won Hong, Iuliia Pliushch, and Visvanathan Ramesh

Abstract—Current deep learning research is dominated by benchmark evaluation. A method is regarded as favorable if it empirically performs well on the dedicated test set. This mentality is seamlessly reflected in the resurfacing area of continual learning, where consecutively arriving sets of benchmark data are investigated. The core challenge is framed as protecting previously acquired representations from being catastrophically forgotten due to the iterative parameter updates. However, comparison of individual methods is nevertheless treated in isolation from real world application and typically judged by monitoring accumulated test set performance. The closed world assumption remains predominant. It is assumed that during deployment a model is guaranteed to encounter data that stems from the same distribution as used for training. This poses a massive challenge as neural networks are well known to provide overconfident false predictions on unknown instances and break down in the face of corrupted data. In this work we argue that notable lessons from open set recognition, the identification of statistically deviating data outside of the observed dataset, and the adjacent field of active learning, where data is incrementally queried such that the expected performance gain is maximized, are frequently overlooked in the deep learning era. Based on these forgotten lessons, we propose a consolidated view to bridge continual learning, active learning and open set recognition in deep neural networks. Our results show that this not only benefits each individual paradigm, but highlights the natural synergies in a common framework. We empirically demonstrate improvements when alleviating catastrophic forgetting, querying data in active learning, selecting task orders, while exhibiting robust open world application where previously proposed methods fail.

Index Terms—Continual Deep Learning, Lifelong Machine Learning, Active Learning, Open Set Recognition, Open World Learning

#### 1 Introduction

W ITH the ongoing maturing of practical machine learning systems. ing systems, the community has found a resurfacing interest in continual learning [1], [2]. In contrast to the broadly practiced learning in isolation, where the algorithmic training phase of a system is constrained to a single stage based on a previously collected i.i.d. dataset, continuous learning entails a learning process that leverages data as it arrives over time. In spite of this paradigm having found various application in many machine learning systems, for a review see the recent book on lifelong machine learning [3], the advent of deep learning seems to have steered the focus of current research efforts towards a phenomenon known as "catastrophic inference" or alternatively "catastrophic forgetting" [4], [5], as suggested by recent reviews [6], [7], [8], [9] and empirical surveys of deep continual learning [8], [10], [11]. The latter is an effect particular to machine learning models that update their parameters greedily according to the presented data population, such as a neural network iteratively updating its weights with stochastic gradient estimates. When including continuously arriving data that leads to any shift in the data distribution, the set of learned representations is guided unidirectionally towards approximating any task's solution on the data instances the system is presently being exposed to. The natural consequence is

superseding former learned representations, resulting in an abrupt forgetting of previously acquired information.

Whereas current works predominantly concentrate on alleviating such forgetting in continual deep learning through the design of specialized mechanisms, we argue that there is a growing risk towards a very different form of catastrophic forgetting, namely the danger of forgetting the lessons learned from past literature. Notwithstanding the commendable efforts towards preserving neural network representations in continuous training, such a high focus is given on the practical requirements and trade-offs beyond metrics that only capture catastrophic forgetting [12], e.g. inclusion of memory footprint, computational cost, cost of data storage, task sequence length and amount of training iterations etc. [6], [13], that it could almost be seen as misleading when most current systems break immediately if unseen unknown data or minor corruptions are encountered during deployment [14], [15], [16]. The seemingly omnipresent assumption of a closed world, i.e. the belief that the model will always exclusively encounter data that stems from the same data distribution as encountered during training, is highly unrealistic in the real open world, where data can vary to extents that are impractical to capture into training sets or users have the ability to give almost arbitrary input to systems for prediction. In spite of the inevitable danger of neural networks generating entirely meaningless predictions when encountering unseen unknown data instances, a well known fact that has been exposed for multiple decades [14], current efforts towards benchmarking continual learning conveniently circumvent this challenge. Select exceptions

M. Mundt, I. Pliushch and V. Ramesh are with the Department of Computer Science, Goethe University, Frankfurt am Main, Germany. E-mail: {mmundt, vramesh}@em.uni-frankfurt.de

Yong Won Hong is with the Department of Computer Science, Yonsei University, Seoul, Republic of Korea.

attempt to solve the tasks of recognizing unseen and unknown examples, rejecting nonsensical predictions or setting them aside for later use, typically summarized under the umbrella of open set recognition. However, the majority of existing deep continual learning systems remain black boxes that unfortunately do not exhibit desirable robustness to respective miss-predictions on unknown data, dataset outliers or commonly present image corruptions [16].

Apart from current benchmarking practices still being constrained to the closed world, another unfortunate trend is a lack of understanding for the nature of created continual learning datasets. Both continual generative modelling (such as the works by the authors of [17], [18], [19], [20], [21], [22]), as well as the bulk of class incremental continuous learning works (such as the works presented in [12], [23], [24], [25], [26], [27], [28]) generally investigate sequentialized versions of time-tested visual classification benchmarks such as MNIST [29], CIFAR [30] or ImageNet [31], where individual classes are simply split into disjoint sets and are shown in sequence. In favor of retaining comparability on a benchmark, questions about the effect of task ordering or the impact of overlap between tasks are routinely overlooked. Notably, lessons learned from the adjacent field of active machine learning, a particular form of semi-supervised learning, do not seem to be integrated into modern continual learning practice. In active learning the objective is to learn to incrementally find the best approximation to a task's solution under the challenge of letting the system itself query what data to include next. As such, it can be seen as an antagonist to alleviating catastrophic forgetting. Whereas current continual learning is occupied with maintaining the information acquired in each step without endlessly accumulating all data, active learning has focused on the complementary question of identifying suitable data for the inclusion into an incrementally training system. Although early seminal works in active learning have rapidly identified the challenges of robust application and pitfalls faced through the use of heuristics [32], [33], [34], the latter are nonetheless once again dominant in the era of deep learning [35], [36], [37], [38] and the challenges are faced anew.

In this work we make a first effort towards a principled and consolidated view of deep continual learning, active learning and learning in the open world. We start by providing a review of each topic in isolation and then proceed to identify previously learned lessons that appear to receive less attention in modern deep learning. We will continue to argue that these seemingly separate topics do not only benefit from the viewpoint of the other, but should be regarded in conjunction. In this sense, we propose to extend current continual learning practices towards a broader view of continual learning as an umbrella term that naturally encompasses and builds upon prior active learning and open set recognition work. Whereas the main purpose of this paper is not to introduce novel techniques or advocate one specific method as a universal solution, we adapt and extend a recently proposed approach based on variational Bayesian inference in neural networks [39], [40] to illustrate one potential choice towards a comprehensive framework. Importantly, it serves as the basis of argumentation in an effort to illustrate the necessity of generative modelling as a key component in deep learning systems. We highlight

the importance of the viewpoints developed in this paper with empirical demonstrations and outline implications and promising directions for future research.

## 2 PREAMBLE: CONTINUAL MACHINE LEARNING

It is likely that the idea of continual machine learning dates back to a similar period of time to the surfacing of machine learning itself. There has been many attempts at defining concepts such as continuous, lifelong or continual machine learning. Often these terms feature negligible nuances and can generally be taken as synonyms. However it seems difficult, and perhaps is not constructive, to attempt to pinpoint the exact onset of when something should be referred to as continual or lifelong learning. Instead, in this section, we will present definitions and related paradigms that have come to enjoy great popularity in the machine learning community. Some of these paradigms are already, or if not yet, should be considered subsets of continual learning (CL) and as a standalone paradigm vary primarily in their current evaluation protocols. We will briefly introduce each of these paradigms and then proceed to summarise and identify characteristic differences with respect to the broader term of modern continual learning.

The first widely circulated definition of *lifelong machine learning* (LML) originated in the work proposed by Thrun [1], [2]. This definition is as follows:

**Definition 2.1.** Thrun - Lifelong Machine Learning [1], [2]: The system has performed N tasks. When faced with the (N+1)th task, it uses the knowledge gained from the N tasks to help the (N+1)th task.

Here, the unmentioned quintessence is that the data of the first N tasks is generally assumed to be no longer available at the time of learning about the N+1th task, i.e. observed data is not just endlessly accumulated and stored explicitly. While this definition captures the basic idea behind continued learning, it is also ambiguous with respect to the definition of task and knowledge. There has been many attempts to find a more concise definition across the literature over the years. One of the more succinct, yet still decently generic definitions followed in the work of Chen and Liu [3]:

**Definition 2.2.** Chen and Liu-Lifelong Machine Learning [3]: Lifelong Machine Learning is a continuous learning process. At any time point, the learner performed a sequence of N learning tasks,  $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_N$ . These tasks can be of the same type or different types and from the same domain or different domains. When faced with the (N+1)th task  $\mathcal{T}_{N+1}$  (which is called the new or current task) with its data  $\mathcal{D}_{N+1}$ , the learner can leverage past knowledge in the knowledge base (KB) to help learn  $\mathcal{T}_{N+1}$ . The objective of LML is usually to optimize the performance on the new task  $T_{N+1}$ , but it can optimize any task by treating the rest of the tasks as previous tasks. KB maintains the knowledge learned and accumulated from learning the previous task. After the completion of learning  $\mathcal{T}_{N+1}$ , KB is updated with the knowledge (e.g. intermediate as well as the final results) gained from learning  $\mathcal{T}_{N+1}$ . The updating can involve inconsistency checking, reasoning, and meta-mining of additional higher-level knowledge.

The authors of this latter definition argue that this definition can be summarized into three key characteristics:

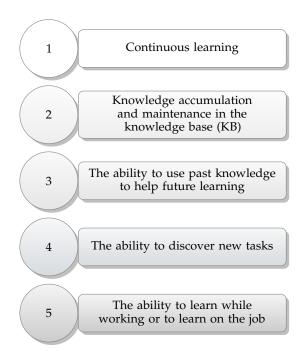


Fig. 1: The five main pillars of lifelong machine learning according to Chen and Liu [3]. Note that the first three pillars were originally proposed and the last two added recently in a second edition redefinition to emphasize new frontiers.

continuous learning; knowledge accumulation and maintenance in the knowledge base (KB); the ability to use past knowledge to help future learning. In contrast to the previous definition by Thrun, mainly the notion of a maintained knowledge base is introduced. Here LML is now defined such that at any given point in time performance can be optimized for any given task by treating all other tasks as previously presented, irrespective of their original order. Whereas the original definition unidirectionaly optimized towards benefiting  $\mathcal{T}_{N+1}$  and thus allowing for performance of previous tasks to degrade over time, Chen and Liu explicitly formulate the preservation of all accumulated information as a fundamental goal of LML. In a recent second iteration of this definition, the authors have added two additional desiderata: the ability to discover new tasks and the ability to learn while working. We have visualized these five essential pillars of LML in figure 1.

Although acknowledged by the authors themselves, this extended definition still lacks with respect to certain aspects:

- a coherent description of domain. This is currently not used unanimously in the literature and often applied interchangeably with task.
- a formalization of knowledge or respective representation thereof in the KB. Typically this is practically constrained to specific applications.
- the essential question of evaluation practice, i.e. choosing, ordering and evaluating the sequence of tasks. This generally requires a human in the loop and considered evaluation scenarios can vary immensely between individual works.

There's many more encountered open questions with

LML in practice, especially with respect to modern machine learning algorithms based on deep learning. As the latter is primarily based on the use of neural networks (NN), they will constitute the main focus of this paper. While the presented arguments will often be of generic nature, this has the advantage that the concept of a knowledge base and its maintenance collapses to the question of managing the model's learned representations. At the same time, this can make the question of how to leverage prior information quite involved as representations in NNs are densely entangled within layers, as well as distributed hierarchically across layers. Before delving into a review of contemporary works, their merits and current limitations, we will present various popular paradigms that are related to the former definitions. This will then be followed by a brief summary on evaluation practices to highlight the nuances.

### 2.1 Related paradigms: subsets of continual learning

Over the course of machine learning development, various different paradigms and evaluation practices have evolved. Throughout this paper, we will come to the already apparent conclusion that CL should ideally be defined as a superset. We will make an attempt towards such a definition at the end of this manuscript. For now, we start by introducing commonly considered machine learning paradigms. As a word of caution, the following definitions should be regarded as non-exhaustive. Even though we have made a considerable effort to provide a comprehensive amount of references, the practical use of certain terminology in particular may still vary largely from community to community. The following shall thus reflect the common use in modern deep learning.

We begin with transfer learning as it can intuitively be regarded as the most related concept. Originally, transfer learning has been proposed as converting a weak learner, one that performs marginally better than random guessing, to one that produces stronger hypotheses [41]. The corresponding formulation that is more specific to neural networks is how the representations obtained by learning through backpropagation can be "recycled" for new tasks [42], [43]. This challenge initially wasn't unanimously referred to as transfer learning, but often was referred to as boosting [44]. A pre-deep learning survey [45] has summarized efforts and formalized transfer learning in the way used today:

**Definition 2.3.** Transfer Learning [45]: Given a source domain  $D_S$  and learning task  $T_S$ , a target domain  $D_T$  and learning task  $T_T$ , transfer learning aims to help improve the learning of the target predictive function  $f_T()$  in  $D_T$  using the knowledge in  $D_S$  and  $T_S$ , where  $D_S \neq D_T$ ,  $orT_S \neq T_T$ .

Here, the authors formalize the use of the terms *domain* and *task* in the context of supervised transfer with datasets consisting of N data instances. They are defined by the following quote: "Given a specific domain,  $D = \{\mathcal{X}, p(\mathbf{x})\}$ , a task consists of two components: a label space Y and an objective predictive function f() (denoted by  $T = \{Y, f()\}$ ), which is not observed but can be learned from the training data, which consist of pairs  $\{\mathbf{x}^{(n)}, y^{(n)}\}$ , where  $\mathbf{x}^{(n)} \in X$  and  $y^{(n)} \in Y''$  [45]. The concept of a domain is therefore defined as the pair of marginal data distribution  $p(\mathbf{x})$  and a corresponding feature space  $\mathcal{X}$ . As it is generally implied that  $\mathcal{X}_S \neq \mathcal{X}_T$ 

or respectively  $p_S(x) \neq p_T(x)$ , an effortless translation of transfer learning to unsupervised or reinforcement learning settings is possible. Without further extensions, this definition of transfer learning is essentially a narrowed down version of the primitive lifelong learning definition 2.1, with the nuance that there typically only exist two tasks. It is similarly unidirectional in the sense that the source task is only used to improve learning the new target.

Since then an almost unending amount of works has sprouted, initiated by works that have started the investigation of transferability of deep neural network features beyond low-level patterns [46], [47], i.e. the higher abstractions and task-specific information believed to be encoded in deeper layers of the hierarchy. Weiss et al. [48] have provided a survey on recent advances. In this context of feature transferability, a variant named multi-task learning (MTL) has emerged. Caruana [49] summarizes the goal of MTL succinctly: "MTL improves generalization by leveraging the domain-specific information contained in the training signals of related tasks". Early works sometimes referred to this as including "hints" [50], [51] to improve learning. In contrast to transfer learning, generally multiple tasks are considered, with the requirement of the model performing well on all of them. However, in the MTL setting, tasks are all trained jointly and no sequence is assumed, corresponding to typical isolated learning practice. In modern day deep nets, MTL thus culminates in the question of how to exactly share the abundant amount of parameters in the architectural hierarchy, see e.g. the overview provided in [52] for variants of sharing architecture portions.

More recently, a very specific form of transfer or multitask learning has evolved. Few-shot Learning [53] developed due to the inability of deep learning techniques to cope with small datasets and empirical risk optimization being unrealiable in small sample regimes. Wang et al. [54] summarized few-shot learning as a type of machine learning problem, where the dataset only contains a limited number of examples with supervised information for the target domain (and generally no constraints on the source domain). This implies that few-shot learning also tackles the issue of rare cases, apart from computational cost and the issue of data collection and labelling. When there is only one example with a label, it is commonly referred to as one-shot learning [53], [55]. Respectively, if no supervised example is provided, the scenario is referred to as zero-shot learning [56]. These scenarios are typically regarded under the hood of transfer learning with additional constraints on data availability.

Apart from concerns about reasonably sized datasets, a different concern is as old as the quest for stochastic approximations itself, namely when to conduct updates. Already in Hebb 1949 [57], *online learning*, i.e. incorporating information immediately as data arrives as opposed to collecting batches before updating a model, was a natural requirement. This question has been elemental in later formalisation of frameworks for empirical risk optimization [58], [59]. Several works have elaborated on challenges in online learning in NNs [60], more generally online learning and stochastic approximations [61], [62] or specifically online gradient descent [63], the workhorse of modern optimization. Given the instance based update nature, online learning in neural networks is inherently tied to the question of how

to avoid catastrophic inference. It is thus not surprising that with the advent of DL immediate attempts have been made to consider online learning in DNNs [64], see a recent survey [65], but the quest for online learning nevertheless still revolves around the interaction between online desiderata and stochastic approximations, or the stochastic gradient descent with backpropagation procedure in particular.

While each paradigm arose for a reason and comes with its own value, namely that of providing better distinction to other works in concrete evaluation scenarios, it is important to remember that the emerging taxonomy is full of nuances that are at times indistinguishable in a more general framework. In consequence, evaluation protocols are central to any discussion. We therefore proceed with details of common evaluation methods in deep continual learning and then summarize the main differences to the paradigms introduced in this section for a compact overview.

### 2.2 Continual learning evaluation

In contrast to isolated machine learning, where the evaluation scenario can often be defined in a straightforward manner by employing performance or satisfying task metrics, continual learning does not directly allow for such an approach. Given that the interest lies in accumulation of information, there are many factors to consider in evaluation of corresponding algorithms. In general it is important to monitor the currently introduced task, yet also investigate semantic drift on previous tasks. One should consider the gain and the ability to leverage representations from task to task in progressive experimentation, yet take note of the task sequence that is crucial to the specific solution obtained. When introducing more tasks, the transfer behavior should be carefully examined, yet interpretation should be treated with caution as not all introduced tasks yield immediate benefits and thus a larger amount of tasks needs to be brought in to the system.

Before continuing with the discussion of evaluation difficulties and metrics, let us take a brief look at some currently employed evaluation methodology [3], summarized visually in figure 2. It seems that such an evaluation protocol is still largely inspired by the isolated machine learning practices. Whereas the notion of information transfer and the sequence of tasks is considered and benchmarked against isolated learning algorithms, such an approach to evaluating the value of continual learning algorithms disregards the relevance of the task sequence (or permutation thereof), choice of tasks or choice of data. Accordingly, recently developed experimental protocols in deep continual learning [6], [7], [8], [9], [10], [11], [12] seem to mainly occupy themselves with evaluation procedures that are heavily inspired by decades of benchmarking learning algorithms in isolation. As a reminder to the reader, we refer to isolated learning as the practice of end-to-end training on a static dataset and evaluation on its pre-defined test set, sans changes over time. As such, the majority of current empirical examination equates continual learning benchmarks with the monitoring of catastrophic forgetting in scenarios that are simple sequentialized versions of popular datasets, similarly to the steps shown in figure 2. With few exceptions, this means that existing datasets are simply split into  $t = 1, \dots, T$ 

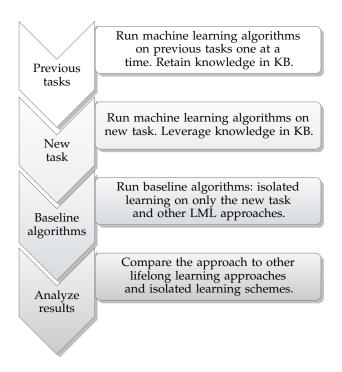


Fig. 2: A widely used approach to evaluation of lifelong machine learning algorithms in the literature [3].

sets, where each of these sets is referred to as one task. These task- or time-stamped sets are then presented one by one to a deep learning system. Typically, each step is assumed to consist of a disjoint set of classes or entire datasets, usually independently of whether the probed task is of supervised, unsupervised or semi-supervised nature, see figure 3 for an illustration. Respectively analyzed metrics [12] are based on this dataset sequentialization and routinely monitor e.g. the degradation of a first task's classification accuracy, the ability to encode new task increments, the overall development of a chosen metric as tasks accumulate or various similar measures to gain an intuition for generative models. It is obvious how this is inspired by isolated learning as these metrics can simply be extracted from a conventional confusion matrix. For this reason, multiple efforts have been made to emphasize the need for more diverse evaluation [6], [13]. Alas the persisting focus on catastrophic forgetting remains visible from the formulated criteria and questions that are deemed necessary to compare methods [6], [13]:

- Memory consumption: amount of required memory.
- Amount of stored data: how much past data does the method need to retain explicitly?
- Task boundaries: does the method require clear task divisions?
- Prediction oracle: does the method require knowing the task label for prediction?
- Amount of forgetting: how much information is retained as measured through proxy metrics.
- Forward transfer: do older tasks accelerate learning of new concepts?
- Backward transfer: do new tasks benefit old tasks?

At this stage the reader might already notice that some of

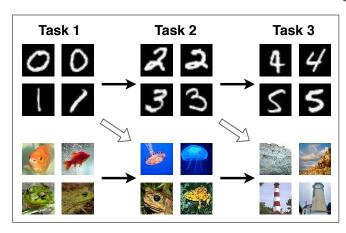


Fig. 3: A typical continual learning scenario dividing common benchmark datasets into a sequence of sub-tasks. Here, the digits one through six from the MNIST dataset [29] and the wordnet ids "n01443537": goldfish, "n01641577": bullfrog, "n01644900": tailed frog, "n01910747": jellyfish, "n09246464": cliff, "n02814860": beacon from the ImageNet dataset [31]. Common evaluation either follows the filled dark arrows to incrementally learn one dataset or alternatively also switches dataset, as denoted by the hollow light arrows.

these listed items are very particular to specific practices. For example, the idea that a prediction oracle would be required in the first place in order to give task labels is an artifact of several works that consider so called multi-head scenarios. The latter makes use of separate disconnected classifiers per task to circumvent explicitly dealing with task prediction interdependency. There exist recent reviews [8] that base their entire evaluation on such a scenario. Empirical surveys in the context of robotics [10], generative models [9] follow similar trends and conduct a "comprehensive application-oriented study of catastrophic forgetting" [11]. With catastrophic forgetting being the sole focus, these works at best cover the first three of the five earlier formulated continual learning pillars 1, if and only if they also conduct an analysis on how specific tasks benefit each other. The recent critiques that formulated above questions [6], [13] therefore present valid attempts to rid current evaluation from such practices that can be seen as inherently violating real continual learning scenarios. Nevertheless, we argue that there is even larger factors at play that transcend these arguments. Although transfer and the sequential nature is considered and benchmarked against isolated learning, crucial aspects such as the relevance of the task order (or permutation thereof), choice of tasks, choice of data and particularly any form of robustness in an open world and with respect to perturbations or attack scenarios are disregarded altogether. Open research areas such as curriculum learning [66], i.e. benefitting from a data ordering of increasing complexity, open world learning [67], i.e. equipping the model with awareness of unseen unknown data, and active learning, i.e. self-selecting data to query for the next step, try to address these crucial elements. We argue that it is imperative to take these perspectives into account in the evaluation of continual learning algorithms. Before proceeding to categorize individual works and consequently making an attempt at connecting the paradigms, we give a

brief summary of the present evaluation differences.

- Transfer Learning: Leverage a source task's representations to accelerate learning or improve a current target task.
  - Difference to CL: unidirectional knowledge transfer between two tasks.
- Multi-task Learning: Exploit tasks relatedness by forming a joint hypothesis space.
  - Difference to CL: isolated learning with multiple tasks
- Online Learning: Retaining and improving a task where data arrives sequentially and real-time constraints require online adaptation.
  - Difference to CL: typically continued learning of one task over time, however generally applicable to any paradigm.
- Few-shot Learning: Transfer or multi-task learning in a small data regime.
  - Difference to CL: unidirectional transfer or isolation similar to transfer or multi-task learning.
- Curriculum Learning: Finding a suitable curriculum that accelerates or improves training by means of introducing schedules of increasing data instance difficulty or data instance task specificity.
  - Difference to CL: isolated learning that prioritizes certain data instances
- Open World Learning: At any particular point in time the model needs to be able to identify and reject unseen data belonging to unknown tasks. These could be set aside and learned at a later stage.
  - Difference to CL: Current CL is typically evaluated in a closed world scenario.
- Active Learning: An iterative form of supervised learning, where the learner can query a user to provide labels for a subset of unlabelled examples that are deemed to yield the largest knowledge gain. Difference to CL: data and sampling efficiency is rarely taken into account in CL on pre-defined benchmarks.

## 3 AN OVERVIEW AND REVIEW OF THREE PERSPECTIVES

We provide a review of the plethora of practices and historically grown methods in the context of deep continual learning, active learning and open set recognition. What may at first seem like a tour de force review for the reader, is intended to first gain an overview of the vast landscape and the deluge of options. This will aid in delving into details of potential pitfalls and shortcomings, but also in highlighting synergies and the necessity for a consolidated view in consecutive sections. As the latter is the primary focus of this work we will limit our survey to concise summaries and will forgo lengthy elaborations on methodological details that are not essential to a generic understanding.

#### 3.1 Continual learning

As indicated in the introductory section, continual learning should ideally encompass a variety of research questions. Whereas our next section will continue to argue that currently considered scenarios are too reductive, resulting in potential difficulty to chose among existing algorithmic options, we will stick to the typical categorization of existing deep

continual works into the three categories of *regularization*, *rehearsal* and *architectural* approaches, in consistency with recent reviews [7], [8], [9]. We note that a strict organization into these groups is not always possible and hence also provide a forth category for works that combine multiple methods. In later sections we will argue that this is not only advantageous, but conceivably a necessity.

### 3.1.1 Regularization:

Continual learning approaches based on regularization aim to strike a balance between protecting already learned representations, while granting sufficient flexibility for new information to be encoded. Intuitively, a meaningful balance should be attainable for tasks with sufficient overlap in their high dimensional embeddings, i.e. if a considerable amount of the learned representations are shareable. Existing approaches can be further subdivided into regularization that explicitly protects parameters, which we refer to as *structural*, which constrains changes on every level of a model architecture, or *functional*, that is preserving a model's output for seen tasks while ensuring full adaptability with respect to each individual model stage that leads to the prediction.

Structural: Inspired by the neuroscientific stabilityplasticity dilemma [57], successful use of regularization of deep learning models for continual learning requires carefully balancing the trade-off between overwriting acquired representations in favor of sensitivity to new information and preservation of already existing formed patterns. Elastic Weight Consolidation (EWC) [24] aims to achieve this balance by estimating each parameter's importance through the use of Fisher information and respectively discouraging updates for parameters with greatest task specificity. Synaptic Intelligence (SI) [68] and Memory Aware Synapses (MAS) [69], where the biologically inspired term synapse is used synonymously with parameter, follow a similar approach by explicitly equipping each parameter with additional importance measures that keep track of past improvements to the objective. Assymetric Loss Approximation with Single-Side Overestimation (ALASSO) [70] can be seen as a direct extension to SI and aims to mitigate its limitations by introducing an assymetric loss approximation that is motivated from empirical observations. Riemannian Walk (RWalk) has generalized EWC and SI by taking into account both the Fisher information based importance, from a perspective of computing distances in the induced Riemann manifold, and the optimization trajectory based importance score. Incremental Moment Matching (IMM) [71] approaches structural regularization from a perspective of Bayesian approximations and matching the moments of tasks' posterior distributions. Uncertainty based Continual Learning (UCL) [72] makes use of Bayesian uncertainty estimates to adaptively regularize weights online. Similarly, Uncertaintyguided Continual Bayesian Neural Networks (UCB) [73] adapts the learning rate in dependence on the uncertainty defined in the probability distribution of the weights.

Functional: Functional regularization approaches are generally inspired by "knowledge distillation" [74], an approach originally proposed for model compression. A distillation loss is introduced by storing the prediction of a data sample for future use as a so called soft target. In learning without forgetting (LWF) [23] for class incremental

continual learning, the soft targets for existing classes are calculated using newly arriving data, even if these predictions might be nonsensical as the freshly added classes do not get correctly predicted yet, in hopes of regularizing towards preserving the output for old tasks. Encoder based lifelong learning (EBLL) [75] applies this concept to the unsupervised learning scenario by applying distillation to autoencoder reconstructions. Knowledge distillation seems to rarely be employed in isolation, but as will be apparent from the list of upcoming combined approaches is a popular technique in conjunction with other mechanisms.

### 3.1.2 Rehearsal:

As the name implies, rehearsal techniques for continual learning aim to preserve encoded information by replaying data from already seen tasks. Trivially, continual learning could be solved by simply storing and replaying all seen data, albeit at usually intolerable memory expense and growing computation time. Accordingly, a core aspect of rehearsal methods is to find a suitable subset of data that best approximates the entire observed data distribution, commonly referred to as selection of exemplars or construction of a core set. Alternatively, a generative modelling approach can be used to generates instances from a learned latent representation as an encoding of the observed data distribution. Most replay techniques indicate their inspiration to be drawn from the complex biological interplay between hippocampus and neocortex, wake + sleep cycles and dreaming in the brain.

Exemplar Rehearsal: GeppNet [76] explores the use of a dual-memory system that implements various short and long-term memory storages that serve to store newly arriving information or provide dedicated replay cycles of previously stored data. Selective experience replay (SER) [77] concentrates on exemplar selection techniques and investigates trade-offs between preferring surprising experiences over rewarding ones, or maximizing distribution coverage. Gradient Episodic Memory (GEM) [26] extends the use of a memory that gets replayed episodically with constraints on the gradients to be non-conflicting with updates for previous tasks. A respective extension called Averaged Gradient Episodic Memory (A-GEM) has introduced significant improvements on computational and memory cost for optimization under these constraints. CLEAR [78] uses experience replay together with off-policy learning to preserve old information and on-policy learning to learn new experiences in deep reinforcement learning. Bias Correction (BiC) [79] rehearses exemplars and additionally corrects for biases in the classification layer.

Generative: Generative replay is a specific version of rehearsal where the data to be rehearsed consists entirely of instances sampled from a generative model. Rather than making use of an episodic memory of previously seen data, generated samples of former tasks are typically interleaved with the current task's real data during training. The most elementary version of this procedure was coined Pseudorehearsal [80], where the generative model is of simple nature. Here, binary patterns are sampled at random, their target value or label computed given the current state of the classifier, and the classifier then needs to maintain the discrimination on these patterns and learn new classes. Such pseudo-rehearsal has then successfully been leveraged

in brain-inspired dual-memory architectures that use two distinct networks for acquisition and storage of information with generative rehearsal to consolidate the memory. Two early examples include pseudo recurrent networks [81] and coupling two reverberating neural networks [82]. Deep Generative Replay (DGR) [17] have introduced a deep learning variant of this practice, where the generative model is taken to be a separate generative adversarial network [83] that gets trained in alternation with a classification model. Replay through Feedback (RfF) [84] proposed generative replay using a single model that handles both classification and generation through the aid of feedback connections. Incremental learning using conditional adversarial networks (ILCAN) [28] follows a similar approach of using a single model, but additionally changes the generative replay component to rehearse feature embeddings instead of aiming at reconstructing original input data. Open-set Classifying Denoising Variational Auto-Encoder (OCDVAE) [39] further introduces the first approach to naturally integrate open set recognition with deep generative replay in a single architecture. This work will play a vital role for the remainder of this paper and we will demonstrate how suggested ideas can be extended to form one potential basis as means to broaden current continual learning practices.

#### 3.1.3 Architectural:

Architectural approaches attempt to alleviate catastrophic forgetting through modification of the underlying architecture. It might at this point be baffling to the reader why such modifications are listed distinctly from the works presented in previous subsections as they are almost by definition complementary to any method presented so far, and in fact most methods presented in this paper. For historical reasons, we will however stay consistent with former categorization of deep continual learning algorithms [7]. The importance of choice of architecture and the need for modifications over time will be another element of our upcoming proposition on an expanded view of continual learning. We will subcategorize architectural approaches further into implicit and explicit architecture modification, i.e. methods that use a fixed amount of maximum representational capacity and methods which dynamically increase capacity in the process of continued training.

Fixed maximum representational capacity: Approaches that use a static architecture rely on task specific information routing through the architecture. An early example is a technique coined activation sharpening towards semi-distributed representations [85], where the essence is to tune and limit the amount of high neural network activations to a maximum of k nodes, such that there is less activation overlap for different representations and consequently less potential for interference of new examples. While fixed architecture methods differ in the specifically employed technique to disambiguate the learned dense representations, the common denominator is the assumption of an over-parametrized architecture in order to warrant enough initial redundancy to permit overriding parameters without incurring catastrophic interference. PathNet [86] adopted this notion to deep neural networks and used a genetic algorithm to determine pathways through the network deemed particularly useful for a specific task in

order to freeze them. Instead of using a separate algorithmic layer to determine task specific network subsets, Piggyback [87] and hard attention to the task (HAT) [88] directly learn binary masks and use them to gate information propagation through the network. The UCB-P variant of the earlier introduced regularization approach Uncertainty-guided Continual Bayesian Neural Networks (UCB) [73] confronts this challenge from a Bayesian perspective. They use uncertainty to prune the model and identify binary masks per task to index into the weights' Gaussian mixture distributions.

Dynamic growth: Dynamic growth approaches administer representational capacity much more explicitly. The trivial solution would be to simply have one model per task and devise a mechanism to select the appropriate path for an input. Alas, such an arrangement doesn't fully leverage information from one task to positively transfer to another or respectively newly arriving information to aid already acquired tasks. First works in deep learning however nearly follow this naive but also intuitive approach to simply train on a task and consequently freeze all learned representations, such as demonstrated in Progressive Neural Networks (PNN) [89]. The amount of weights is then increased for a new task, with the twist that formerly learned representations laterally transmit their output to the new tasks' representations but not vice versa. Expert Gate [90] is comparable and differs mainly in the introduction of a gating mechanism that automates the choice of a suitable expert in an ensemble. Recent perhaps more practical approaches can be viewed as once again drawing their inspiration from decades of biological findings and discussion on neurogenesis. The latter refers to the process of creation and incorporation of new neurons into the existing system, see Aimone et al. or Vadodaria et al. for reviews [91], [92]. For the last two decades it has now been acknowledged that this process persist beyond early stage human development and continues its function in adults [93]. The seminal work of dynamic node creation in neural networks [94], where additional units are added whenever the loss plateaus, has thus found a renaissance in modern deep learning. Neurogenesis deep learning to accomodate new classes (NDL) [95] and lifelong learning with Dynamically Expandable Networks (DEN) [96] have adapted this heuristic approach for use in continual deep learning. The former by adding units whenever the reconstruction error of an autoencoder surpasses a predetermined threshold in the spirit of Zhou et al. [64], the latter based on an empirically found value of the classification loss in supervised learning. Reinforced Continual Learning (RCL) [97] or Learn-to-Grow [98] further attempt to overcome the challenge of finding suitable loss cut-offs and cast dynamic unit addition into a meta-learning framework in order to separate the learning of the network structure and estimation of its parameters.

### 3.1.4 Combined Approaches:

We list a number of, largely very recent works, that primarily advance the state of the art on a set of benchmark datasets by blending techniques from the previous categories. One of the most popularly cited works is iCarl [25], which couples a knowledge distillation based regularization approach with rehearsal of exemplars, assembled through a greedy herding

procedure [99]. Variational Continual Learning (VCL) [20] similarly fuses use of an episodic memory of exemplars with parameter regularization, but from a perspective of approximate Bayesian inference. FearNet [27] has later critiqued iCarl as a viable technique due to its heavy dependency on quantity of data in order to be successful. They have therefore additionally incorporated generative rehearsal to compensate the need to store large subsets of the original dataset. Variational Generative Replay (VGR) [19] can be seen as concurrent to VCL, where instead of exemplar rehearsal generative replay is made use of. Memory replay GAN (MRGAN) and Lifelong GAN (LLGAN) [22] are recent complements to these works and deviate in that they are based on GANs instead of variational inference in pure autoencoders. Whereas MRGAN uses a functional regularization approach to align the generator's output, LLGAN further applies such distillation loss based regularization across multiple places in the architecture to regularize encoders and discriminators. On the architectural front, Variational Autoencoder with Shared Embeddings (VASE) [18] adopts dynamic architecture growth in conjunction with generative replay. Their proposal is to allocate additional representational capacity for new concepts, determined through larger reconstruction loss in a variational autoencoder, however, is limited to expanding the latent space and leaving the rest of the architecture static. Lifelong Learning for Recurrent Neural Networks (LLRNN) [100] combines training of long short-term memory (LSTM) [101] with gradient episodic memory based exemplar rehearsal and a capacity expansion approach named Net2Net [102], which provides the means to transfer learned representations from an architecture to a larger untrained one before continuing to train the latter. While some of these works clearly exploit natural synergies, a generally desirable practice, we note that this can sometimes come at the expense of detailed analysis and comprehensive understanding of individual key ingredients and their necessity. While we agree that all approaches in this subsection pursue commendable directions, we argue that considerable future analysis is still required. We will discuss corresponding details and suggestions in later sections.

#### 3.2 Active learning

Rather than focusing on the question of how to preserve representations in incremental continual learning, the topic of active learning asks the reverse question of how to pick data increments for future inclusion. Generally, this is cast into the framework of semi-supervised learning. Here, it is assumed that the model is trained on labelled data  $oldsymbol{X}_L =$  $\{m{x}_L^1,\dots,m{x}_L^n\}$ , and a larger pool of unlabelled data  $m{X}_U$  exists. This is motivated from data acquisition being relatively cheap in the modern world, as opposed to human intensive data labelling that often requires highly skilled experts. The task of an active learner is thus to extract a set of M data instances  $\{x_{II}^1, \dots, x_{II}^m\}$  from the pool of unlabelled data, such that a maximum gain in performance on the inspected task is expected if a human in the loop provides the additional labels  $\{y^1,\ldots,y^m\}$  for further training. The underlying mechanism on which the query is based is referred to as the acquisition function and forms the main pillar of active learning research. We have visualized this active learning cycle in figure 4.

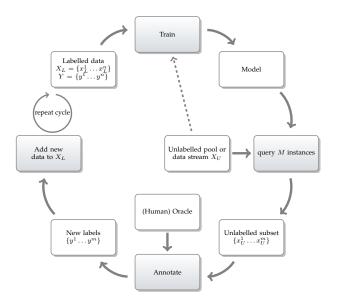


Fig. 4: Active learning cycle that repeatedly expands a labelled dataset by querying and then annotating data instances from a larger unlabelled pool. The dashed arrow from the latter to the training process indicates the common closed world active learning scenario, where the presence of all data at all times is assumed. Respective works typically include the entire unlabelled dataset into the training procedure by employing methods from semi-supervised learning. Shaded parts of the diagram correspond to processes, whereas light components represent objects.

There is multiple conceivable evaluation variants to gauge the usefulness of active learning acquisition function choices. They either explicitly assume the entirety of the unlabelled data to be accessible and usable upfront, or contrarily the query being informed solely by the available labelled data. Independently of the latter, the practical assessment of active learning strategies is generally conducted in a closed world scenario, i.e. the entire pool of unlabelled data is expected to stem from the same data distribution as the initially labelled set and the oracle is assumed to be infallible. In a crucial distinction to continual learning, evaluation of active learning however accumulates data and grows the labelled set, focusing primarily on the cost reduction of labour intensive annotation. In consequence, an active learner is deemed successful if each data query provides significant benefit over simply picking and labelling data at random.

"A probability analysis of the value of unlabelled data for classification problems" [103] provides an early analysis of the requirements for benefiting from semi-supervised or active learning approaches. The authors consider two types of models: parametric p(x,y|W) = p(x|W)p(y|x,W) and semi-parametric: p(x,y|W) = p(x)p(y|x,W). In the latter, the data probability p(x) is decoupled and can have an unknown (or non-parametric) form independent of the weights W, as is common in most discriminative models such as logistic regression or most neural networks. They argue that these models are particularly suited for active learning, as opposed to parametric models such as Gaussian mixtures being particularly suitable for semi-supervised

learning. This is because they do not need to rely on potentially inaccurate estimates of the entire data distribution when only a fraction of the data is observable. However, we will see in the subsequent review that both of these model types have been used to form different perspectives to address active learning and come with their respective advantages.

As with the majority of techniques, early active learning methods have rapidly cross-pollinated into applications with deep neural networks. However, due to the black-box nature of deep non-linear neural networks, many of these approaches are based on simple heuristics or approximations to uncertainty quantities that no longer have tractable closed-form solutions. We will start with these heuristic approaches, as they are often trivial to transfer to deep learning, and then continue to summarize more principled approaches, which can turn out to be genuinely challenging in the context of deep learning.

### 3.2.1 Uncertainty Heuristics

One theoretically sound approach to querying useful data is based on entropy [104] sampling and other information theoretic acquisition functions [105]. An early approach based on training two neural networks to estimate query areas in binary classification problems [106] remarks that this is difficult for neural networks as they are often overly confident in their outputs. This overconfidence is going to be one of the main subjects of our next major section on learning in an open world. Interestingly, while payed painstaking attention in early literature, this aspect seems to often be overlooked in the era of deep learning. Simply using neural network prediction confidence, predictive entropy or other derived heuristics [107] are still practically employed in comparisons today [36]. This is because many approaches have been shown to empirically work well in specific contexts, although there is no guarantee for them to succeed. Early works have shown uncertainty sampling based active learning for logistic regression [107] and neural networks [108], [109] based on "query by committee", an approach to estimate uncertainty by using an ensemble of neural networks. This idea has later found a one-to-one translation to deep ensembles for active learning [35]. Naturally, most black-box deep neural networks are not equipped with mechanisms to gauge uncertainty properly outside of using multiple parallel models. Bayesian active learning by disagreement (BALD) therefore provides an attempt at avoiding the necessity of ensembles and instead uses Monte Carlo dropout [37], [38] to calculate points of high variance in the output [110]. This has empirically been demonstrated to be effective and has been extended in Bayesian Generative Active Learning (BGAL). Here, BALD is used to query samples and then the labelled set is further augmented with generated examples [111]. Deep incremental learning with Neural Architecture Search (iNAS) [36] does not propose a new query mechanism and instead provides an evaluation of above acquisition functions in the context of architecture selection. They include the option of progressive architecture growth after each query, to illustrate that small models generally fare better in a small data regime, whereas large models are required when a certain degree of task complexity is reached. We will revisit this as an imperative insight in our later discussions.

### 3.2.2 Version Space and Expected Error Reduction:

A theoretically more substantiated approach to basing the acquisition function on heuristics is to query data that provably reduces the expected error. Clearly, this is beyond the current understanding of deep neural networks, but has been shown to be feasible in the context of parametric models such as Gaussian mixture models [112] or naive Bayes [32]. These works use the concept of a version space [113], i.e. the consistent set of hypotheses that separate the data in the induced feature space. An appropriate active learning strategy is to sequentially and monotonically reduce the size of this version space. In models such as SVMs for binary classification this is intuitively explained based on the margins [114], where new points are chosen according to hyperplanes that maximize the restriction with respect to the set of possible hyperplanes for correct classification. The latter was later extended to a multi-class SVM based approach [115], however still based on multiple binary classifiers. This allowed for theoretical guarantees on sample complexity and necessary amount of queries to be analyzed with respect to these binary classification problems with linear decision boundary in the context of greedy active learning strategies [116]. Whereas "learning active learning from data" [117] provides a recent effort to train a metalearning based regressor to predict expected error reduction for binary classification using random forests, the idea has not been adapted to deep neural networks yet.

### 3.2.3 Representation based approaches:

Although version space reduction can come with provable guarantees, respective application to deep neural networks is inconceivable before a mature theory of how their hypotheses are formed has evolved. At the same time, Roy et al. [32] have pointed out that the earlier summarized uncertainty sampling, or estimates thereof through ensembles, are generally insufficient. They argue that they are prone to querying outliers, as a result of sampled instances being viewed in isolation and without regarding the underlying density of the full data distribution. Similar conclusions were empirically observed in the large scale empirical evaluation of active learning for text applications [33]. As a solution, the authors suggest a representation based information density measure, and although heavy to compute, it implicitly takes into account the underlying data distribution. This can be seen as an approach that is orthogonal to minimizing the version space, where now typically the distribution coverage on the entire dataset according to the model representations is maximized instead of reducing the number of possible hypotheses. The often necessary core assumption is thus the presence of the entire unlabelled pool of data and its auxiliary use in optimization of the labelled set. We have attributed a third category of active learning to approaches that follow this objective.

Active learning using pre-clustering [118] uses a k-medoids algorithm in conjunction with a SVM or logistic regression to select data from the pre-clustered embedding of the unlabelled pool. Similarly, SVM based core vector machines [119] use a set of minimum enclosing balls to create a core set that best approximates the entire distribution. Li et al. estimate information density by using the unlabelled data

in a Gaussian process [34]. The idea in these works have since been abstracted to deep neural networks. Sener et al. [120] base their active learning procedure on construction of core sets based on a k-medians algorithm. Shui et al. [121] achieve distribution coverage by matching distributions through minimization of the Wasserstein distance in Autoencoders (WAAL). Variational adversarial active learning (VAAL) [122] approximates the data distribution by learning the latent space in a variational autoencoder [123] and simultaneously trains a latent based adversarial network to discriminate between unlabelled and labelled data.

In complement to these works, various querysynthesizing methods have been proposed [124], [125], [126]. Here, the challenge of active learning is tackled by using a deep generative model to generate informative queries. Instead of querying from an unlabelled pool directly, generative adversarial active learning (GAAL) [124] and "efficient active learning using conditional generative adversarial network" (Efficient cGAN AL) [125] both train GANs to synthesize and label queries. The core assumption is the ability to adequately capture the data distribution to generate meaningful instances. The usefulness of the generated samples with respect to a classifier can then either be assessed through uncertainty heuristics or by matching the synthesized data with samples from the pool and retrieving the most similar instance. The latter has been demonstrated in Adversarial Sampling for Active Learning (ASAL) [126].

In our later discussion, we will argue that the assumption of upfront presence of all data should, and in fact can be lifted when a natural bridge to the other paradigms is constructed. We proceed to conclude our review by delving into what will constitute the glue: learning in an open world and open set recognition.

## 3.3 Open set recognition

The term open set recognition was formally coined only recently [127], [128]. However, its foundation and associated challenge in neural networks dates back to at least several decades before, when discriminative neural networks were found to yield overconfident mispredictions on unseen unknown data [14]. To get an intuitive understanding, let us briefly consider the types of data we can expect our model to encounter. As soon as we move beyond the closed world benchmark scenario, we can no longer expect our trained models to be tested exclusively on some held-out data from the same distribution as observed during training. In the earlier introduced transfer learning parlance, for prediction, data can thus generally not be presumed to originate from the same domain. We can now distinguish three types of possible inputs to our model [127]:

- 1) **Knowns**: examples belonging to the distribution from which the training set was drawn. The model's prediction is accurate and confident.
- 2) Known unknowns: unknown instances that a model cannot predict confidently. Examples can optionally be labelled as not being affiliated with the set of known concepts for explicit training of negatives. Prediction uncertainty can indicate a model's awareness of its limitation.

 Unknown unknowns: unseen instances belonging to unexplored, unknown data distributions or classes for which the prediction is generally overconfident and false.

The broader inspiration for this categorization is commonly attributed to a notorious, machine learning unrelated, quote by Donald Rumsfeld [127], [129]: "We know that there are known knowns; these are things we think we know. We also know there are known unknowns; that is to say we know there are some things that we do not know. But there are also unknown unknowns; these are the ones we dont know, we dont know!" [130]. In the context of neural networks, known unknowns can be identified through gauging model uncertainty or relying on derived related heuristics, in correspondence to many of the methods employed in the active learning setting. However, as detailed in a recent survey [15], separating the known data from the essentially indistinguishable high-confidence mispredictions for unknown unknowns is far from trivial.

As any machine learning model is trained on a finite dataset, and the imaginable set of unknown unknowns is infinite, we refer to the challenge of recognizing the latter as open set recognition in analogy to prior works [15], [67], [127], [128], [131]. Formally, these works define the closed space as a union of balls  $S_K$  that enclose the entire training set  $X_K$ , whereas the open space  $\mathcal{O}$  constitutes the remainder of the input or feature space:  $\mathcal{O} \subset = \mathcal{X} - \mathcal{S}_K$ . Correspondingly, works that provide attempts at addressing open set recognition aim to find the respective boundaries between known and unknown spaces [39], [40], [127], [128], [131], [132], [133]. We will review these works last in favor of historically preceding approaches based on explicit inclusion of negative classes and rejection through anomalies in prediction patterns, even though the latter have been argued to be insufficient for open set recognition [14], [15], [127].

The above widespread categorization can technically be extended to encompass a fourth category, by splitting the knowns into known knowns and the set of unknown knowns [134]. We do not consider this further distinction as the existence of unknown knowns can be condensed to either a wilfully ignorant false prediction, because we in fact know the concept but choose to nevertheless treat it as unknown, or the more charitable alternative in which our chosen machine learning model has an inherent inability to represent the investigated concept and its structure altogether. We also note that there is other related concepts, such as novelty detection [135] or equipping classifiers with rejection options. These are different in such that they are typically still evaluated in the close world and data is generally still expected to reside in a similar domain. The aim is to recognise outliers of the distribution that are uninformative or represent a particularly interesting rare event. Although these works can have considerable merit in their respective closed world application context, we do not review them in favor of the more generic open set recognition, where considered inputs are allowed to be of almost arbitrary nature. We further note that we naturally cannot provide every example that has ever attempted open set recognition through simple heuristics like using the output values to distinguish examples.

### 3.3.1 Prior Knowledge

A conceivably simple effort to address unknown unknowns is by assuming that the human modeller has enough awareness about what forms of unknown inputs to expect during deployment to directly incorporate this prior knowledge into the model. As inclusion of prior knowledge into neural networks and other types of deep models turns out to be remarkably complex, the natural analogue is to steer efforts towards dataset design. "Inference with the universum" [136] has accordingly proposed to embrace prior knowledge by representing it through a collection of "non-examples", and hence letting the optimization algorithm decide how to include the presented information into the model. Unfortunately, this does not provide a general solution for open set recognition as upfront knowledge can only ever truly cover the family of known unknowns. At best, a mere workaround for major failure cases is therefore supplied, although without any associated guarantees for remaining unknown unknowns. This lack of guarantees is further enforced by the necessity to rely on machine learning algorithms extracting the information and composing abstractions from the supplied "non-example" data population.

Since then, the idea to include a "background" concept has been adopted so widely across applications, that singling out and thus giving preference to select works is difficult. Take as an example large-scale datasets surrounding the task of material classification and semantic segmentation. Because there is an abundance of material types, it has become the de-facto standard to collapse any available imagery that is connected to less important materials or where meager amounts of data are available into a single "other" material [137], [138]. Not only is it impractical to gather data for every material variation, but also unknown unknowns can feature other significant statistical deviations, due to e.g. previously unencountered illumination, acquisition and sensor differences, superposition of dirt and surface markings, or any type of perturbation and previously unencountered noise. Imaginably, in real applications beyond a closed world, inclusion of an endless universe is by definition infeasible. Nevertheless, multiple recent works follow this route and propose mechanism to calibrate output confidences in deep models [139], formulate a discrepancy loss between knowns and known unknowns [140], or modify the embedding to explicitly seperate them, e.g. in semantic categorical and contrastive mapping (SCM) or the objectosphere loss [141], [142]. Although these approaches are not tantamount to a comprehensive solution, we note that they can still in principle be sufficient for tasks in partially constrained environments that naturally limit the world's openness.

## 3.3.2 Predictive Anomalies

From an unsuspecting angle, a model will consistently yield accurate predictions only for observed data and produce highly uncertain output otherwise, yet still generalize correctly to data that is from the same domain but has not been included in training. In this view, determining a prediction threshold and obtaining an uncertainty estimate is sufficient to recognize any form of unknowns. This can work surprisingly well in models with thorough understanding of the decision boundary and its neighbourhood, such as

the Transduction Confidence Machine-k Nearest Neighbors (TCM-kNN) [143]. Even though it is well known that the entangled dense representations of neural networks result in overconfident predictions on any data [14], [15], a variety of practical approaches nevertheless proposed to simply rely on a hinge loss to reject during classification [144] or even to take the straightforward route and directly trust the softmax confidence [145]. As the quantitative outcome leaves room for improvement, multiple works have argued that uncertainty estimation is required to corroborate the decision to gain awareness of the unknown. In deep networks this could be achieved by assessing the variations of stochastic forward passes through a neural network with dropout [38], [146], [147], as a variational Bayesian approximation to a distribution on the weights [37], or by empirically estimating the output's variability with respect to introduced perturbations, such as done in ODIN (outlier detection in neural networks) [148], and by calibrating the prediction accordingly [139]. In similar spirit, an often employed argument is that generative modelling is required to obtain meaningful prediction values that allow to recognize out of distribution samples. For this purpose, Lis et al. [149] use image resynthesis and equate detection of unknown concepts with identification of discrepancies in poorly reconstructed image regions. Likewise, one-class novelty GAN (OCGAN) [150] generates examples from sparsely populated latent space regions in order to use them in explicit training of a binary out-of-distribution classifier. Although predictions and uncertainty from generative models have been shown to improve outlier and adversarial attack detection in contrast to purely discriminative models [39], [40], [151], there is strong empirical evidence that this is still insufficient to provide a generic solution [39], [40], [152], [153]. It is clear that former reported cases of success can be attributed to the specific constrained empirical studies and we illustrate some remarkably simple failure cases of prediction confidence and entropy in figure 5, even when uncertainty is assessed with Monte Carlo Dropout. This is to provide an intuitive picture of the challenge of open set recognition with neural networks and to summarize and repeat the findings of the much more detailed experiments presented in numerous prior works [39], [40], [152], [153].

## 3.3.3 Meta-recognition

Rather than assuming that predictions are somehow calibrated for any data, a more rigorous approach is to prevent overconfident misclassification by confining the model to the known closed space and averting any prediction from little-known open areas in the first place. Whereas it is evident how to achieve this when explicitly modelling the distribution, such as done in probabilistic mixture models, a straightforward approach is not typically applicable in the often complex feature hierarchies of modern discriminative machine learning approaches. A common technique is thus to resort to meta-recognition on top of the empirically emerged features obtained through black-box optimization procedures. Scheirer et al. [131] give an intuitive example based on support vector machines. Here, the menace of erratic predictions for unknown unknowns results from examples being projected close to the linear decision boundary, while at the same time being mapped arbitrarily far away from the train-

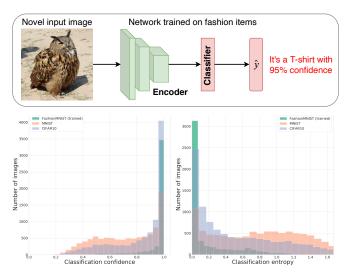


Fig. 5: Top panel: Qualtiative illustration of the challenge of open set recognition. A neural network that has been trained to discriminate fashion items misclassifies the unknown concept of an owl and assigns it to the t-shirt class with very high confidence. Bottom panel: A quantitative example of a deep wide residual neural network trained on the FashionMNIST dataset, asked to classify unrelated unencountered digits and objects from the MNIST and CIFAR10 datasets. Even though uncertainty is estimated using 50 Monte Carlo Dropout passes, misclassified unseen data still overlaps significantly with the known dataset in prediction confidence or entropy. Knowns and unknowns are largely indistinguishable. The shown quantiative results are a reproduced subset of our previous work investigating the limits of deep neural network unertainty for open set recognition [40].

ing data along a different dimension. The authors therefore define a compact abating probability (CAP) model, where the key idea is to make use of insights from extreme value theory (EVT). The essential notion is to take into account inherently present extreme statistical differences in the long tail of an extreme value distribution, here the Weibull distribution, and subsequently monotonously decrease a data point's probability of belonging to the observed closed set with increasing distance from the observed data population. In other words, a prediction is discarded in sparsely populated areas, independently of a sample's proximity to the decision boundary. Bendale et al. [67] have extended this approach to discriminative deep neural networks, where the above metarecognition idea is transferred to the network's penultimate layer. They propose the OpenMax algorithm that lowers softmax prediction probabilities with increasing distance from the average penultimate layer's activation values. A strongly related approach has been proposed in Lee et al. [132], where the affinity of a data point to the known set is measured based on a Mahalanobis distance in the feature space of the penultimate layer. More recent works have come to the conclusion that although the latter approaches have a strong theoretical foundation for open set recognition, they are still limited by activation values in discriminative neural networks being optimized exclusively towards predicting a

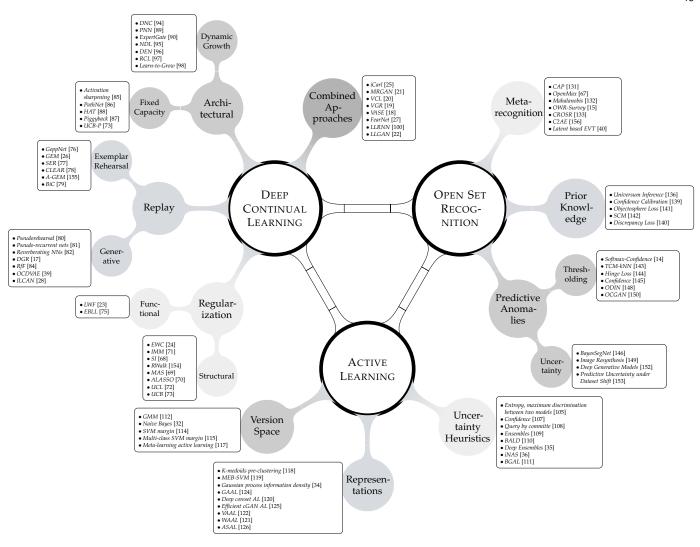


Fig. 6: A visual overview of the taxonomy of neural network based methods for continual learning, active learning and open set recognition. Distinctly categorized approaches are rarely coupled and synergies exploited only in select works, such as the combined continual learning approaches. More importantly, the intersection between the three machine learning paradigms remains largely unexplored. Highlighting the necessity for unification of the latter into a single viewpoint is the primary purpose of this work. A respective practical framework is conceptually and mathematically described in detail in sections four and five.

correct class [39], [40], [133]. In particular, the penultimate layer activation values do not generally encode all the information about the data x that might be required for open set recognition. "Classification Reconstruction learning for Open-Set Recognition" (CROSR) [133] has thus suggested to additionally append a generative model's latent variable z to the OpenMax classification procedure. Concurrently, open set classifiying denoising variational autoencoders (OCDVAE) [39], [40] translate the EVT based meta-recognition to a variational Bayesian setting. Here, the open set recognition is based directly on the approximate posterior in a deep generative model, which enables a natural interpretation based directly on the underlying generative factors of the data distribution p(x), instead of activation value heuristics. We believe that this approach offers one potential framework to consolidate research in active learning, open set recognition and continual learning. We will correspondingly revisit the

underlying approach, detail specific methods and introduce extensions in the next section.

## 4 Bridging perspectives: past insights and the challenge of evaluation

In the previous sections, we have kept up the tradition to treat continual machine learning, active learning and open set recognition as three distinct challenges. For convenience we provide a visual summary of the taxonomy in diagram 6. The remainder of the paper will now serve the purpose of revealing the natural interface. In fact, by identifying former lessons, stressing shortcomings of prevailing evaluation practices and bridging seemingly forgotten connections, we develop a wholistic view that simplifies the deluge of ongoing research questions into a single intuitive framework. To better understand why this is imperative for future progress,

let us briefly recall the earlier mentioned predominant evaluation routines and link insights from prior works to their current limitations.

If we look back at figure 2 and the corresponding section's discussion, we recall that deep continual learning typically collapses its practical evaluation to measuring catastrophic forgetting between task increments. These task increments belong to simple sequentialized versions of existing benchmark datasets and a continual learning technique is deemed successful if the model that is trained over time approaches the expected performance when trained in isolation. In almost complete analogy, active learning evaluation revolves around accuracy gains between query steps. In the majority of the aforementioned related works, the focus is exclusively on whether a specific query mechanism surpasses another in terms of quickly approaching the overall error achieved on a complete dataset. For empirical benchmarking purposes, the model is simply trained in isolation on multiple selected subsets of known data, where the difference between these subsets corresponds to the inclusion of one active query.

Before we continue with the limitations of such evaluation protocols, we emphasize that our intention at no point in this paper is to discredit and devalue the bulk of previously proposed methods. However, we would argue that claimed advances of individual methods are in grave danger from their constrained benchmark evaluation being non-indicative of the actual machine learning progress on a larger scale. We believe a major contributing factor is that key insights from past, often neural network unrelated, literature have surprisingly gone unnoticed or have been written off in the era of deep learning. To attach a slightly provocative connotation, we have termed these overlooked insights forgotten lessons. Although the term "forgotten" certainly is an exaggeration with regard to the ML field as a whole, the absence of derived practical implications is strongly manifested in deep learning evaluation schemes.

### 4.1 Forgotten lessons from past literature

**Forgotten lesson 1:** Machine learning models are by definition trained in a closed world, but real-world deployment is not similarly confined. Discriminative neural networks yield overconfident predictions on any sample.

Independently of whether additional metrics such as training speed-ups through representation transfer, computational cost or memory consumption are taken into account, currently considered experimentation features closed world train and test sets. This is occasionally amplified by continual learning works assuming the presence of a task oracle for testing or respectively the assumption of an infallible oracle to yield flawless data when labelling active learning queries. As such, open issues concerning continual training of a model or active learning queries in an open world are generally neglected. However, real-world deployment almost always inhabits an open world. In the extreme case, the model has to handle data from completely unknown type in previously unfamiliar conditions, think outdoor environments or uncontrolled arbitrary user inputs in web-based applications. Instead of the common overconfident misprediction that falsely attributes this data to any known concept, a multiple decade old seemingly

forgotten insight [14], any machine learning model should at least be equipped with the ability to identify unencountered scenarios and warn the practitioner. As a much milder, but heavily realistic form of an open world, even commonly occurring corruptions are disregarded, think blur or camera noise in images. The menace of the latter has recently been demonstrated in deep learning by Hendrycks and Dietterich [16], where the authors empirically demonstrate that current deep neural networks not only exhibit severe instability with respect to various simple perturbations, but advances in neural network architectures are reflected in only diminutive changes in robustness. Whereas certainly this hazard is universal to all machine learning research that is deployed in practice, continual and active learning are particularly prone to the threat of corrupted and unknown data as their goal is to accumulate knowledge from previously unseen sources already in the training process itself.

**Forgotten lesson 2:** Uncertainty is not predictive of the open set. Active learning resides in an open world and common heuristics based query mechanism are susceptible to meaningless or uninformative outliers.

Although early works have rapidly identified the fallacy that uncertainty sampling is a meaningful strategy to query [32], [33] in active learning or respectively detect unknown unknowns [14], [106], the belief that uncertainty provides a generic solution seems to have resurged with the advances of deep learning. This is apparent from the many approaches in our previous literature review basing querying strategies or detection of unseen examples on heuristics that rely on output variability or similar entropic quantities, see the branches labelled with uncertainty and predictive anomalies in our literature review diagram 6. Indeed, the challenge of accurate uncertainty quantification in deep learning is already genuinely difficult and does provide advantages in contrast to less principled empirical thresholding. However, paying homage to the detailed argumentation of the recent review by Boult et al [15], any machine learning model is still trained in a closed world scenario, independently of whether e.g. a Bayesian formalism is employed to obtain uncertainties. Predictions for *y* are known to be overconfident, uncertainty is not calibrated for points outside of  $p_{train}(\mathbf{x})$  and the posterior is often unusable, regardless of how well it is approximated.

In other words, given any parameters  $\phi$  and an unknown unseen input example  $\mathbf{x}^*$ , we don't know if evaluating  $q_{\phi}(\mathbf{z}|\mathbf{x}^*)$  will produce something meaningful. This issue is by no means exclusive to detecting unknown unknown examples, but comes with the same implications for realistic active learning scenarios. Take for example a more realistic set-up beyond a crafted benchmark where data is scarce and the investigated domain is demanding even for experts. The earlier reviewed VAAL has considered such a scenario with medical imaging, where correct oracle labelling and a noiseless image cannot always be expected. Sample selection based on uncertainty does not protect the query from such noise and there is a large chance that meaningless outliers are included into the system.

**Forgotten lesson 3:** Confidence or uncertainty calibration, as well as explicit optimization of negative examples can never be

sufficient to recognize the limitless amount of unknown unknowns.

At a first look, one might believe that impressive successes where demonstrated with approaches that extend the basic idea of "inference with the universum" [136]. Explicitly using prior knowledge in terms of expectations on what form of inputs can be anticipated, or respective inclusion of negative data that is believed to play a role in deployment, are popularly exhibited by works that have identified and attempt to address the first two lessons. The common presumption across all these works is the upfront presence of a larger, possibly unlabelled, dataset that can explicitly be included into the optimization process. Just as supposed out-of-distribution examples are made use of to modify loss functions and calibrate the output for detection of unknown unknowns [132], [138], [140], [141], [142], active learning techniques often resort to conditioning their procedure on the entire data pool [34], [118], [120], [121], [122], e.g. through clustering [118], [120] or fitting a generative model to the unseen data [34], [121]. Unfortunately, this impedes evaluation beyond a constrained closed set benchmark and more realistic continual and active learning scenarios where data becomes available at different times cannot be considered. In a sense the problem seems to be addressed from a reverse perspective. Instead of acquiring explicit knowledge about the nature of the trained data distribution, the challenge is sidestepped by reformulating it as an optimization problem that attempts to find the boundary between known and an existing set of unseen data, which by definition then does not consist of unknown unknowns. Thus, we receive no guarantees, as the pool of unlabelled data at any point in time is limited and can never truly approximate the unknown space.

Apart from this obvious argument that it is impossible to include all forms of variations and exceptions upfront, else we could have just modelled and hand-crafted the entire system from the start instead of falling back on purely data driven approaches, previous works have also asserted that the particular form of representations of discriminative deep neural networks can further confound predictions. The early 1992 work of French [85] has already pointed out that a major complication of continually training neural networks is their distributed representations and has subsequently investigated mechanism to obtain semi-distributed representations with sharp activations that are concept specific. We argue that with the onset of deep learning the challenge of distributed representations is further magnified due to distribution across the layer hierarchy. First, consider as an example a neural network that is trained to discriminate cars from airplanes, a scenario often assumed when incrementally training the popular CIFAR10 dataset [30]. As the neural network is not explicitly encouraged to encode information about the data distribution, the obstacle of predicting overconfidently on unseen data is further magnified by the ubiquitous option for any classifier to differentiate a concept based on a combination of noise patterns, the absence of a specific pattern, or background patterns altogether [157]. In the car vs. airplane scenario, depending on how well and diverse the dataset is constructed, this could be as trivial as distinguishing the two classes by identifying the presence of some feature that describes the sky. As neural networks

have been demonstrated to rely heavily on texture rather than object boundaries [158], this is not far fetched. In fact, a prominent recent work on "unmasking clever hans" predictors [159] has shown that the decision making of a discriminative deep neural network can be based on entirely trivial features, such as a certain object always occurring at a specific location in every image or almost imperceivable photography tags. "Adversarial examples are not bugs they are features" [160] takes this one step further and empirically showcases how classes can be distinguishable solely based on noise patterns. In a trivial case of our above car versus airplane example, presenting the trained model with images of ships that feature the similarly blue background of the sea is then not surprisingly resulting in overconfident misclassification. Using ships as a background class could initially solve this problem of attributing blue to airplanes. However, if a significant portion of our learned features were indeed to be composed of noise, background and adversarial patterns, then we would argue that overconfident mispredictions are impossible to overcome, as the extent of data on which these features activate is inconceivable to any human modeller. We believe this makes the approach to handle outlying and unknown unknown data through prior knowledge even less feasible.

Forgotten lesson 4: Data and task ordering are essential. Although this forms the quintessence of active learning it is yet untended to in continual learning.

It is well known that each dataset instance does not contribute equally to the overall objective. This forms the foundation and rationale behind active learning. In general, when conducting active learning queries, there is a trade-off between exploring the unknown space and exploiting more of the already known to avoid misclassification [115]. Alas, the implications of the latter statement are more nuanced and go beyond the simple question of whether a certain subset spans the entire data distribution. As an example, Joshi et al. [115] found certain active learning strategies to benefit primarily from creating a class imbalance, as more difficult classes might require a denser sampling than others. Bengio et al. [66] have similarly found that sorting data in a curriculum that introduces classes into the training process according to their difficulty improves the obtained accuracy. Recently, Hacohen et al. [161] have empirically observed that deep neural networks seem to build such a curriculum inherently during the training process. Consistently across multiple architectures, they always learn the same examples first when given access to the entire dataset, even though the mini-batch stochastic gradient descent shuffles the data differently every time. Intuitively, this notion of learning according to some measure of complexity seems only natural, as describing some inputs necessitates less complex and nuanced patterns than others.

Even though there is significant empirical evidence that data selection and task order plays a vital role for any learned algorithm, modern deep continual learning, to the authors' astonishment, seem to pay little attention to a careful experimental design.

Out of the numerous works of the previous review, less than a handful of works consider the question of task order at all. The rest remains in the comfort of benchmark datasets, where the classes are split and introduced in sequence for continual learning according to a class id that often just reflects an alphabetic ordering. However, there is no rigorous investigation of the effect of task order. Two out of the four works that examine task order [77], [88] only randomize the order across multiple experimental repetitions to obtain an average performance estimate. The other two [8], [162] follow this practice, but go even further and make the statement that task ordering has minimal influence towards continual learning methods. We will later demonstrate that this is obviously not the case, and can simply be attributed to the experimentation being a narrow trial of five randomly obtained orderings without any attached semantics. When selecting tasks from the overall pool of available data according to their similarity or dissimilarity with the already observed data distribution, we will observe a major divergence of obtained results.

Whether or not having access to all future tasks in order to select an ideal order is unrealistic in real-world continual learning scenarios, we believe task ordering to be an imperative factor that should be considered when designing our benchmarks to further our understanding. In particular, we note that a very common practice to reduce the computational cost of incrementally learning large scale datasets such as ImageNet [31] is to extract subsets [8], [25], [70], [79]. The main problematic here is that selecting e.g. 50 or a 100 from a larger pool of 1000 classes heavily influences the achievable result and using random selection mechanisms essentially renders works unreproducible.

**Forgotten lesson 5:** Parameter and architecture growth are not distinct methods to address any particular challenge such as catastrophic forgetting. They are at the core of the learning process.

We do not truly believe that the above lessons is forgotten, however, feel the need to call attention to it because an entire branch of continual learning seems to treat parameter addition and architecture growth as a separate solution. Our main goal for techniques that modify architectures on the fly is to point out that these should be analysed with particular caution. On the one hand, methods that use neural networks that are highly over-parametrized can implicitly expand their effective representational capacity due to the abundance of parameters when encountering new data. Investigated algorithms could thus always implicitly be accompanied with some form of representational expansion, depending purely on the initial choice of architecture. On the other hand, in active learning it has been shown that training in small sample scenarios is not only computationally more efficient with smaller neural networks but also yields more accurate estimates in these early stages if less representational capacity is available [36]. Whereas the latter statement might seem obvious to some reader, we note that this behaviour makes it tremendously difficult to attribute gains of active learning or continual learning experiments to a specific technique in contrast to innate advantages of the used architecture at any point in time.

## 4.2 Open set recognition: the natural interface between continual and active learning

As indicated in the previous sections, contemporary continual and active learning are prone to an alarming amount of threats due to their development and evaluation inhabiting a closed world. In this section we argue that awareness of an open world is not only required to overcome the threat of designing a non-robust system, but provide the natural means to merge techniques into a common perspective.

Recall that a majority of continual learning techniques alleviates the challenge of catastrophic inference by regularizing parameters for known tasks, rehearsing a subset of data from known tasks or respectively generating it with a generative model. Independent of the specific algorithm, a key concern is thus to identify exemplars, learn the generative factors of our known tasks or determine the parameters that are responsible for the majority of previously seen data. At the core, we need to thus find a good approximation of the known data distribution. In active learning, our task is very much alike, although the underlying question seems to be of reversed nature. Instead of protecting or sampling from the known data distribution, a query is conducted with respect to yet unobserved distributions. In a similar distinction to the continual learning mechanisms, query-acquiring active learning methods pick samples that are estimated to yield the best model improvement, whereas query-synthesizing methods attempt to tackle this challenge through generative modelling by generating these most informative examples.

Interestingly, in open set recognition, the task is to precisely gauge the boundary between the seen known data distribution and yet unseen unknown data. Although the original motivation stems from a perspective of outlier detection and thus model robustness in practical application in the presence of unknown unknowns, knowing this boundary also gives us the means to restrict a continual learning technique to protect the already seen knowns or respectively query active learning examples that are sufficiently statistically different without the fear of selecting uninformative noise. We argue that in general this forms the natural interface between active and continual learning.

We follow previously reviewed works that employ EVT based meta-recognition to identify unknown unknowns and schematically illustrate our proposed unified framework in figure 7. We will delve into the mathematical details of its realization in deep neural networks in the next section. For now, consider a generic embedding as a result of some deep neural network encoding. In the figure's leftmost panel, we have visualized an example embedding for three classes, with their mean indicated by a star and a potential decision boundary by dashed lines. In order to confine predictions to the known space, EVT based meta-recognition makes use of data instances with extreme distance values to the average embedding of a class. Typically, a Weibull distribution is used to model the distance distribution for the entire dataset and capture samples that feature stronger deviation in a heavy tail. In the original works that have proposed this model for open set recognition [127], [128], [131], the cumulative distribution function is then used to estimate whether a new unseen example should be regarded as an unknown unknown, outlying data point. In our own previous work [39], we have identified this technique to also be fundamental in judging whether a randomly sampled latent vector is proximate enough to the observed data such that it results in a clear output of a generated model.

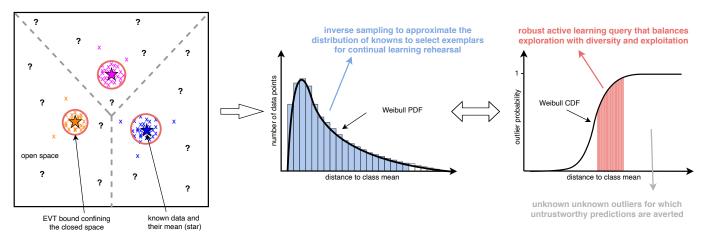


Fig. 7: Conceptual diagram to illustrate how extreme value theory based meta-recognition in neural networks can serve as a common denominator to protect knowledge in continual learning, conduct principlied queries in active data selection, while having the capability to reject or set aside unknown unknown data at any point in time. The leftmost figure of an embedding showcases the threat of the open space, where any examples that are very far away from known clusters always get falsely assigned to a known class and can be arbitrarily close to the decision boundary. The mid panel shows how a Weibull distribution, which models the extreme distance values to the mean of the correctly predicted trained data in a heavy tail, can enclose the known space (suggested by the red circles in the embedding). The corresponding cumulative distribution function in the right panel can be used to reject or set aside outliers and balance active learning queries to sample diverse, yet meaningful data (shaded red area). Alternatively either curves can be sampled inversely to select a subset of inlying data to approximate the entire known distribution in continual learning rehearsal (shaded blue area).

We now close the circle and tie this method to retention of a core set for continual learning, as well as a query mechanism for active learning, while retaining the method's innate ability to reject and set aside unknown unknowns. First, we postulate that the Weibull distribution for each data point's distance to the mean embedding equips us with a tool to approximate the known distribution with a subset. Specifically, we can employ inverse sampling from the Weibull probability density function to create a set of distance values with an arbitrary prior on how much of the distribution's tail should be disregarded, i.e. how many outliers are already assumed to be inherently present in the original dataset. Practically, we can then approximate the data distribution with a subset by selecting data instances whose embedded value lies closest to the drawn sample. Alternatively, as indicated in the diagram, we could discretize the distribution and sample a certain number of examples from each bin. Conversely, for active learning, we are less interested in sampling from the known distribution, but much more in the heavy tail. To our advantage, the long tail models data that is statistically deviating, but can still be attributed to the distribution of interest. We can thus balance exploitation with exploration. First and foremost, data instances for which the outlier probability is unity are avoided altogether in order to prevent sampling of uninformative noise or other corrupted data. Recall, that this is the primary pitfall of uncertainty sampling. At the same time, we want to avoid samples that have a minute probability of being an outlier, as these samples are too similar to previously observed data and are therefore also uninformative due to redundancy. As such, we can constrain our query to the center area of the cumulative distribution function (CDF), illustrated by the shaded area under the CDF in the diagram. The

rationale for this approach can intuitively be understood by looking back at the theoretically grounded works of version space maximization. We can implicitly reduce this space of possible hypothesis, even in complex models such as neural networks, as we incrementally expand the radius of the ball that encloses the closed space by sampling carefully along its boundary with each active learning query. This way, we avoid the vast open space and the redundant highly dense areas of known data, while making sure that previously unseen information is acquired.

Before we proceed with one imaginable realization of this unified framework in neural network and its mathematical formalism, we note that there is two works that have previously initiated a bridge between active learning and open set recognition, alas have not fully built it yet. The recently introduced open world learning [67] and the concurrently named cumulative learning [163] advance the pure open set identification step by proposing to set aside the unknown unknowns and including them into a later active learning cycle. Whereas these works made first steps towards formulating learning in an open world, they however assume the presence of labels for the entire dataset and the addition of classes itself is in the form of a fixed sequence that is injected by the human. The system is limited as it does not self-select which classes or instances should be learned next, nor does it protect its knowledge for continual learning, where the assumption of availability of all data at all times is lifted. As a result, the empirical evaluation is simply an investigation of the performance on the entire test set at each state of the growing known training set. Finally, the suggested open world learning [67] is based on nearest mean classifiers based on simple SIFT features and is yet to be extended to the context of modern deep neural networks.

## 5 Uniting perspectives with deep generative neural networks

How can we realize our proposed unified framework in a meaningful way in deep neural networks? As emphasized by prior work [40], [133], identification and correlation of unseen data with average activation patterns of known data is not necessarily sufficient in discriminative models, even when extreme values are modelled to obtain closed space boundaries (see prior works [39], [40] for empirical verification). This is because a neural network based classifier is generally not encouraged to aggregate the whole information describing the data, merely the features that allow for class distinction. These features themselves, come with a variety of further pitfalls, as summarized in the forgotten lessons. In our own previous work [39], [40], we have overcome this limitation by formulating the problem from a perspective of deep generative models trained with variational Bayesian inference, i.e. variational autoencoders (VAE) [164]. We will lean on this viewpoint, follow the notation of prior works and extend it towards one potential solution to consolidate continual and active learning through open set recognition.

The rationale to build upon VAEs is rather straightforward: the Bayesian formulation lets us learn about the distribution of seen data  $p(\mathbf{x})$  by capturing it through latent variables **z**. However, as  $p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) dz$  is untractable, we do this by optimizing a lower-bound to the marginal distribution  $p(\mathbf{x})$ , since the densities of the marginal and joint distribution are related through Bayes rule  $p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x},\mathbf{z})}{n(\mathbf{x})}$ . As we do not know our real posterior  $p(\mathbf{z}|\mathbf{x})$ , we typically resort to variational inference and introduce a variational approximation  $q(\mathbf{z}|\mathbf{x})$  to the posterior. In a neural network, this approximation  $q(\mathbf{z}|\mathbf{x})$  is learned through the parameters of a probabilistic encoder, whereas a probabilistic decoder is trained for the joint distribution  $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$  and thus forms the generative component. This generative model can effortlessly be augmented to additionally discriminate classes by including their label into the latent variable, e.g. by enforcing a linear class separation on z. The corresponding factorization and generative process is then simply  $p(\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{z}) = p(\boldsymbol{x}|\boldsymbol{z})p(\boldsymbol{y}|\boldsymbol{z})p(\boldsymbol{z})$  [39], [40]. Such formulation of a classifying variational autoencoder comes with the main advantage that using latent variables z allows us to base our decision regarding unknown unknowns on the underlying generative factors of variation and whether an example is close to the high density regions of our approximated data distribution.

### 5.1 The boundary between known and unknown

The first step towards open world aware active and continual learning is to train the above mentioned classifying variational autoencoder, followed by determining the boundary between the open and closed spaces for the observed distribution with the help of EVT. For ease of readability, we repeat the training and fitting procedure described in our previous work [39], [40]. The model's probabilistic encoder and decoder are trained jointly by minimizing the divergence between the variational approximation  $q_{\theta}(\boldsymbol{z}|\boldsymbol{x})$  and a chosen prior  $p(\boldsymbol{z})$ , typically  $\mathcal{N} \sim (0,I)$ , and the conjunction of reconstruction loss and the linear classification objective, parametrized through  $\phi$  and  $\boldsymbol{\xi}$  respectively. For a dataset

consisting of n = 1, ..., N elements, the following lower bound to the joint distribution p(x, y) is thus optimized:

$$\mathcal{L}\left(\boldsymbol{x}^{(n)}, \boldsymbol{y}^{(n)}; \boldsymbol{\theta}, \boldsymbol{\phi}, \boldsymbol{\xi}\right) = -\beta K L(q_{\boldsymbol{\theta}}(\boldsymbol{z}|\boldsymbol{x}^{(n)}) || p(\boldsymbol{z})) + \mathbb{E}_{q_{\boldsymbol{\theta}}(\boldsymbol{z}|\boldsymbol{x}^{(n)})} \left[\log p_{\boldsymbol{\phi}}(\boldsymbol{x}^{(n)}|\boldsymbol{z}) + \log p_{\boldsymbol{\xi}}(\boldsymbol{y}^{(n)}|\boldsymbol{z})\right]$$
(1)

At any point in time of training this model, there is a natural discrepancy between the prior and the approximate posterior. The added  $\beta$  factor in above equation serves the purpose of controlling this gap. Whereas one could belive this distributional mismatch to be an undesired property, we recall the arguments conjectured in multiple previous works [165], [166], [167]. In essence, they state that the overlap of the encoding needs to be reduced in order to avoid indistinguishability, but at the same time prevent latent variables to consist of individual uncorrelated data points that resemble a pure look-up table. In the intuitive picture of diagram 7, think of the former as multiple classes collapsing and thus being inseparable, and the latter as the dense clusters being scattered to allow differentiation of each and every single data point without a strong encoding of correlations. Therefore, the actually captured encoding of the data distribution should not simply be assumed to correspond to the prior, but rather corresponds to an empirically determinable distribution referred to as the aggregate posterior:

$$q_{\boldsymbol{\theta}}(\boldsymbol{z}) = \mathbb{E}_{p(\boldsymbol{x})} \left[ q_{\boldsymbol{\theta}}(\boldsymbol{z}|\boldsymbol{x}) \right] \approx \frac{1}{N} \sum_{n=1}^{N} q_{\boldsymbol{\theta}}(\boldsymbol{z}|\boldsymbol{x}^{(n)})$$
 (2)

Using EVT to find the boundaries of this distribution now corresponds to identification of our model's closed space. For emphasis, we repeat that this is necessary because VAEs generally assign non-zero density to any point in the latent space, the analogue of overconfident classifier predictions [152], [153], and that this boundary is not analogous to the extent of the prior because low density areas exist inside the prior as well. Practically, an EVT based fit can be obtained by empirically accumulating the mean latent variable for each class c for all correctly predicted known data points  $m=1,\ldots,M$ :

$$\bar{z}_c = \frac{1}{|M_c|} \sum_{m \in M} \mathbb{E}_{q_{\theta}(\boldsymbol{z}|\boldsymbol{x}^{(m)})}[\boldsymbol{z}]$$
 (3)

and defining a respective set of latent distances as:

$$\Delta_{c} \equiv \left\{ f_{d} \left( \bar{\boldsymbol{z}}_{c}, \mathbb{E}_{q_{\theta}(\boldsymbol{z} | \boldsymbol{x}_{t}^{(m)})} \left[ \boldsymbol{z} \right] \right) \right\}_{m \in M_{c}} \tag{4}$$

Here,  $f_d$  represents a chosen distance function, which prior works have typically chosen to be either euclidean or cosine distance [39], [127], [128], [131]. As this set represents the distances to the class conditional aggregate posterior, we can fit a Weibull distribution with parameters  $\rho_c = (\tau_c, \kappa_c, \lambda_c)$  on  $\Delta_c$  to model the trustworthy regions of high density that represent the observed data distribution, where the heavy-tail indicates a decaying reliability:

$$\omega_{\rho}(z) = \frac{\kappa}{\lambda} \left( \frac{|f_d(\bar{z}, z) - \tau|}{\lambda} \right)^{\kappa - 1} \exp\left( -\frac{|f_d(\bar{z}, z) - \tau|}{\lambda} \right)^{\kappa}$$
(5)

Here,  $\tau$  defines the location,  $\lambda$  the scale and  $\kappa$  the shape of the distribution. We can now make use of this distribution to pinpoint the observed data distribution, as a surrogate to the otherwise highly complex aggregate posterior. We proceed to highlight its various use cases in the following sections.

### 5.2 Approximate posterior based open set recognition

As described in previous works [39], [40], the most direct use of the aggregate posterior based Weibull parameters  $\rho$  is the identification, rejection or storage of unknown data. Using the corresponding cumulative distribution function (CDF) to the probability density function of equation 5, we can now estimate any data instance's statistical outlier probability for every known class:

$$\Omega_{\rho_c}(\boldsymbol{z}) = 1 - \exp\left(-\frac{|f_d(\bar{\boldsymbol{z}}_c, \boldsymbol{z}) - \tau_c|}{\lambda_c}\right)^{\kappa_c} \tag{6}$$

When we have observed multiple classes, we will typically take the minimum  $\min (\Omega_{\rho})$  of this equation across all known classes c and the respective mode's parameters  $\rho_c$ . This expresses the basic condition that a data point should be considered as a statistical anomaly only if its outlier probability is large for each known class. A respective decision should thus be based on the class where the smallest deviation to known data is observed. The more dissimilar a sample is with respect to the observed data distribution as approximated by the aggregate posterior, the more the outlier probability will approach unity. Irrespective of whether a machine learning algorithm is developed for active learning, continual learning or in fact any other paradigm, this robustness towards unknown unknown data is essential for any practically deployed system that operates outside of extremely narrow conditions.

#### 5.3 Outlier and redundancy aware active queries

Equation 6 gives us the direct means to estimate a sample's similarity with the already known data. For active learning this almost directly translates to the informativeness of a query. Small CDF values signify large similarity or overlap with already existing representations, larger values indicate previously unobserved data. Naively, one would follow the earlier strategies developed in uncertainty based active learning and simply query batches that consist of the most outlying data points. However, this would neither grant protection from exploring noisy, perturbed and uninformative data, nor balance it with exploitation to fester partially known concepts. Our proposition is thus to query a variety of data that is well distributed across the center part of the CDF, i.e. data that surpasses an outlier probability of e.g. 0.5 and at the same time is limited on the upper end by e.g. a value of 0.95. As explained in the earlier introduction of the framework, this is tantamount to sampling on the outer edge of the sphere that encloses the currently known closed space. Naturally, as a repetition of the ultimate statement of the last subsection, if the employed active learner is simultaneously deployed or used in application once it has finished learning, avoiding predictions for unknown unknown data is imperative.

## 5.4 Core set selection for continual learning rehearsal

In contrast to active queries that need to select meaningful unknown data, in the currently formulated continual learning paradigm the main goal is to protect the known knowledge while learning a predetermined new task. We will question the role of the order prearrangement in the next subsection. Here, we focus on open world aware techniques to preserve previously acquired representations. Depending on available memory, the most successful approaches either store and rehearse a small subset of exemplars or alternatively generate data for former tasks with a generative model. In our previous work [39] we have shown how we can use equation 6 to reject samples from the prior  $z \sim p(z)$  that do not fall into the obtained bounds of the aggregate posterior for generative rehearsal. The choice for this sampling with rejection originated from the decision to employ the cosine distance, which collapses the distance to a scalar. A different distance function, such as a euclidean distance per dimension would allow to directly inversely sample a highly multimodal Weibull distribution, i.e. with one mode per dimension per class. Independently of the selected distance metric, we can leverage inverse sampling for the construction of a small data subset. Specifically, drawing at uniform from the inverse of the CDF in equation 6 is guaranteed to yield samples that approximate the aggregate posterior:

$$f_d(\bar{z}, z) = \Omega^{-1}(p|\tau, \lambda, \kappa) = \lambda \left(-\log(1-p)^{\frac{1}{\kappa}}\right) - \tau$$
 (7)

The core set can now simply be obtained by picking the data points that are closest to the obtained distance values, if the chosen distance metric collapses the distance to a scalar, or directly to the latent vector, if the chosen distance metric preserves the dimensionality. Note that we have chosen to inversely sample the CDF of equation 6 in favor of a more compact equation. It should however be clear that eq 5 can alternately be sampled equivalently. The advantage of such a core set selection procedure is that we always attempt to approximate the underlying distribution, with the quality being defined by the desired amount of exemplars, while excluding statistical anomalies by limiting outlier probability values to e.g. p < 0.95. As anticipated, the latter plays the additional crucial role of robust application when the system has finished learning and is deployed.

#### 5.5 Class incremental curricula and task order

Continual learning methods are mostly evaluated in the context of class incremental learning. The classes of a benchmark dataset are typically split into disjoint sets and introduced to the learner in alphabetical or class index sequence. Due to the large computational effort of training neural networks to convergence on long task sequences, several works choose to evaluate on subsets of classes [8], [25], [70], [79]. An important remaining question is thus how such evaluation affects comparability and reproducibility, or more generally the role of task order. As mentioned earlier, selecting a meaningful ordering is in most cases non-trivial. Large-scale dataset such as ImageNet are often composed by scraping data from the internet, social media or through uncontrolled acquisition that prioritizes as large as possible datasets. We as humans thus lack the knowledge to build an intuitive learning curriculum when paired with our lack

of understanding of deep neural network representations. Consequently, scarcely any works have attempted to address this challenge beyond a simple randomization of the class order. Fortunately, we can provide at least a partial remedy to the seemingly arbitrary class incremental evaluation setting. Although we do not have access to explicit data distributions for any task, equation 6 allows us to assess the similarity of new tasks with the aggregate posterior for known tasks. In the spirit of our earlier formulated active learning query, we can start with any task t and proceed to select future tasks  $t \in T$  that feature the least overlap with already encountered tasks (or most overlap, depending on what is desired):

$$t_{\text{next}} = \underset{t \in T}{\arg \max} \left\{ \mathbb{E}_{p_t(\boldsymbol{x})} \Omega_{\boldsymbol{\rho}} \left( \mathbb{E}_{q_{\boldsymbol{\theta}}(\boldsymbol{z}|\boldsymbol{x})} \left[ z \right] \right) \right\}$$
(8)

To provide an example, if our objective was to incrementally expand a system to recognize individual animal species, one assumption could be to accelerate training by always including the species that is most similar to what has already been learned, as this could be hypothesized to require only small representational updates. An alternative objective could be to design a system that expands its knowledge in an attempt to cover and generalize to an as large as possible variety of concepts. In this scenario, one could choose to always include the next task with the smallest amount of overlap with existing tasks to maximize learning of diverse representations.

We could now delve into a philosophical debate on when it is reasonable to assume access to future tasks in continual learning to undergo above selection, and when the task sequence is unavoidably dictated by other external factors. We refrain from this discussion at this point and will instead focus on highlighting the large effect on performance when the task order is chosen by above mechanism in the following empirical investigation. At the very least, we hope that this will invoke a more careful and consistent evaluation on existing benchmarks, instead of picking arbitrary data subsets, selecting different random class orders and nevertheless attempting to compare results across methods.

#### 6 EXPERIMENTAL VERIFICATION AND ANALYSIS

In this section we provide the empirical verification for the earlier introduced framework and its specific realization in deep neural networks. For this purpose, we start with a quantitative comparison of exemplar selection mechanisms to prevent catastrophic forgetting in continual learning and querying strategies in active learning. Here, we will first show that the proposed common EVT based foundation surpasses several conventionally employed techniques. We then proceed to further highlight the method's superiority in the open world. In contrast to most methods that are developed with a unidirectional focus on improving a specific active learning or continual learning benchmark, our framework has the critical advantage of not breaking down in the presence of corruptions that commonly occur in practical application in the wild. To conclude the experimental section, we investigate the role of task order for evaluation. We show that a task curriculum constructed through our framework consistently results in considerable improvements.

We base our experiments on the MNIST [29], CIFAR10 and 100 datasets [30]. Although these datasets could be regarded as fairly simple, they are advocated as the predominant benchmarks in all of the presented continual learning works and still present a significant challenge in this context. They are further sufficient to point out major differences between methods, particularly with respect to robustness, showcasing a disconnect with real application and realistic evaluation. We use a 14 layer wide residual network (WRN) [168], [169] encoder and decoder with a widening factor of 10, rectified linear unit activations, weight initialization according to He et. al [170] and batch normalization [171] with  $\epsilon = 10^{-5}$  at every layer, to reflect popular state-ofthe-art practice. To avoid finding elaborate learning rate schedules or resorting to other excessive hyperparameter tuning, we use the Adam optimizer [123] with a learning rate of 0.001 and a sufficiently high-dimensional latent space of size 60 for all training. We use this common setting to corroborate our wholistic view and describe further details for specific experiments in consecutive subsections.

### 6.1 Exemplar selection and core set extraction

Before we dive into a quantitative comparison of methods that aim to alleviate catastrophic forgetting through the selection and maintenance of a core set, we need to address a potential evaluation obstacle. In continual learning works, the typical evaluation relies on monitoring the decay of a metric over time when training is conducted on new tasks and old tasks are retained by continued training on a few select exemplars. However, there seemingly is no common protocol of how these exemplars are interleaved. Apart from obvious factors such as the amount of chosen exemplars, works such as variational continual learning [20] use the exemplars only at the end of each task's training cycle to finetune and recover old tasks, whereas most other works [25], [77], [79] simply concatenate exemplars with newly arriving data. Ultimately, the different works make use of different methods for exemplar selection and attempt to compare their effectiveness through the final metric, even though they are generally not trivially comparable due to their distinct choices of the training procedure.

To highlight this argument we have trained the typical split MNIST and CIFAR10 scenarios, where classes are introduced sequentially in pairs of two and only the new task's data is available to an incrementally growing single head classifier. The old task is approximated through a core set of size 2400 and 3000 respectively, i.e. we pick 240 and 300 exemplars per class that correspond to retention of 4% and 6% of the original data. We train the model for 150 epochs per task to assure convergence and interleave exemplars selected by our proposed EVT approach in three different manners: 1.) We conduct the predominant naive concatenation of the core set with the new task's data and continue training with mini-batch gradient descent that samples data uniformly (unbalanced mini-batch sampling). 2.) We recognize that the former combination and sampling leads to a heavy imbalance as the core set size is generally much smaller than the new task's available data. We naively correct this through weighted sampling that samples a mini-batch such that it consists in equal portions of former tasks' exemplars

and new task's data, generally oversampling the exemplars (balanced mini-batch sampling). 3.) We identify that the latter weighted balanced sampling always results in an equal amount of exemplars and new data in a mini-batch, independently of the number of classes that the core set or the new task increment are comprised of. To correct for the number of classes, we further investigate class balanced sampling, where each mini-batch is sampled such that each class is equally represented. To give an example, if we have seen two tasks of two classes and proceed to learn the next task, the core set with its four classes will be oversampled to constitute two thirds of a mini-batch and the remaining third is made up of the two classes of the third task.

We show the obtained empirical continual learning accuracies in figure 8. With gaps of over 5% it is evident that balancing mini-batches is essential. More so, it is clear that a comparison of different core set works, just because they have used a similar core set size, can result in an apples to pears comparison if other aspects such as the detailed training procedure and mini-batch sampling are not taken into account. As our main focus is to analyze the core set selection strategies and their limitations, we proceed to compare different core set selection strategies in isolation from the precise continual learning setting. In analogy to Bachem et al. [172] and the "reverse accuracy" evaluated in LLGAN [22], we first train the model on the entire dataset, then select core sets of different sizes, and finally retrain the model exclusively on the core set to assess the approximation quality of our strategy. We repeat this entire procedure five times to gauge statistical consistency and estimate deviations. Without a doubt, methods that select a core set that yields a better approximation of the overall population and results in larger accuracies when trained in isolation, also provide better means to alleviate catastrophic forgetting in continual learning. We compare six different methods:

- 1) **Random:** select exemplars uniformly at random.
- 2) **Greedy k-center:** greedy k-center approximation [173] for coreset selection as used in Variational Continual Learning [20]. In essence, exemplars get picked one by one to obtain a cover of the distribution by maximizing their distance in latent space to all existing data points in the core set.
- 3) **Input k-means:** k-means clustering with k being equal to the number of exemplars. Raw data points get selected that are closest to each obtained mean. Suggested as an alternative to greedy k-center in variational continual learning [20].
- 4) Latent k-means: analogous to above input based k-means, but with the difference that the clustering is conducted on the lower dimensional latent embedding.
- 5) Latent herding: an adaptation of the herding procedure, used in Rebuffi et al. and Wu et al. [25], [79], to operate on the latent space instead of an arbitrary neural network feature space. Herding greedily selects exemplars one by one such that each exemplar addition best approximates the overall data's mean embedding.
- 6) **Latent EVT:** our proposed EVT based inverse Weibull sampling introduced in sections 4 and 5.

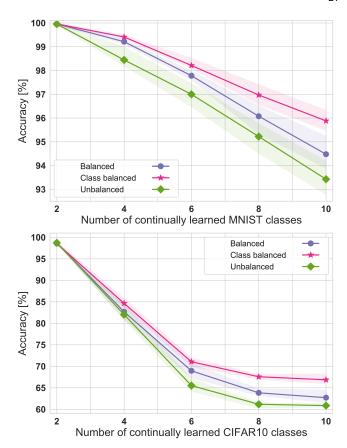


Fig. 8: Influence of mini-batch sampling in continual learning with core sets on MNIST and CIFAR10. The green squared line represents unbalanced sampling, the naive practice of sampling mini-batches uniformly from the concatenated pool of the new task's data and the retained core set. The purple dotted line weights the sampling to oversample the much smaller core set to balance the mini-batch equally. The latter is further corrected with respect to classes in the pink starred line, where the sampling is adjusted to draw mini-batches that are comprised of the same amount of instances per class independently of their origin. We have repeated the experiments five times, illustrated by the shaded regions ranging from the minimum to the maximum obtained values. We can observe that such training details result in very significant performance differences beyond the statistical deviations of a specific core set selection strategy. This imposes an additional challenge in the evaluation of core sets for continual learning. Core sets have been selected with the proposed EVT based method and consist of 240 and 300 exemplars per class for MNIST and CIFAR10 respectively.

We show the obtained accuracies by training on differently sized core sets selected by the above mechanisms in figure 9. As expected, random sampling features large variations, with the best attempts rivalling the other methods and in the worst case yielding substantially worse results. The k-means methods both perform similarly, with the latent space version operating on a lower-dimensional embedding showing minor improvements over the clustering obtained on the original image data. The smaller the core set size, the worse these methods seem to perform. This is not

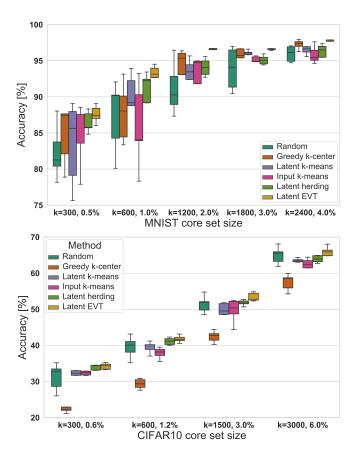


Fig. 9: Training accuracy on core sets constructed by different popular strategies. Results for different core set sizes, characterized through their size k and the respective percentage of the dataset, are illustrated in a box plot to show the median, first and third quartile and minimum and maximum values obtained from five experimental repetitions. If viewed without color, methods are displayed from left to right in order of the legend from top to bottom.

surprising and Bachem et al. [172] have already argued that kmeans with well separated clusters with sufficiently different amount of data points per cluster can be prone to inaccurately estimating multiple cluster centers in highly populated areas versus none in more sparsely populated clusters. This is further amplified by k-means generally necessitating a subsampled initialization to operate in high dimensions and at large scale. As such, we also observe larger variations for these methods. Latent herding is subject to much less overall variation and seems to initially do very well. However, in contrast to the proposed latent based EVT procedure, we notice an increasing gap in accuracy with larger core set sizes. Intuitively, we attribute this to herding picking increasingly redundant samples due to the objective relying exclusively on the best mean approximation, which does not simultaneously tend to diversity. Our latent based EVT approach that aims to approximate the underlying distribution features by far the least deviation and consistently outperforms all other methods.

To provide a better intuition, we have re-trained the model with a two-dimensional latent space to visualize the aggregate posterior and compare it with the selected

core sets. Figure 10 shows the latent embedding with the first four CIFAR10 classes. The colored points correspond to the embedding of the entire set of data points and the respective curves correspond to kernel density estimates of the aggregate posterior. The black crosses indicate the points selected for a small core set of size 200, i.e. 50 per class. The left panel illustrates the greedy k-center approach, whereas the right panel shows the EVT aggregate posterior based approximation. Evidently, the approximation of the distribution is almost impeccable for our proposed approach, with the greedy alternative leaving much to be desired. We argue that this is due to the greedy k-center procedure optimizing for a cover based on maximal distances, alas without explicitly replicating the density or taking into account inherently present outliers and unrepresentative examples. While this might not be much of an issue for the highly redundant clean MNIST dataset, the arbitrarily collected real world data of the CIFAR10 dataset entails complete failure for the greedy k-center approach. In fact, by introducing a few naturally occurring image corruptions, we will show that such lack of robustness can be observed for all but our proposed method in a later section. Before we dive into this aspect of robust application in the open world, we first proceed with a quantitative analysis of the active learning perspective.

### 6.2 Active queries

In addition to the last section showing the advantages of our proposed framework for the construction of core sets that approximate the aggregate posterior, we empirically demonstrate the benefits when conducting EVT based queries for active learning. Recall that active learning is challenging because we generally desire to query batches of informative data at a time instead of querying, re-training and re-evaluating one by one. This is particularly imperative for computationally expensive deep learning and adds a further constraint of not only querying meaningful samples, but also making sure to query diversely without too much redundancy between the queried examples. We consider this typical deep active learning scenario for MNIST and CIFAR10, where we start with a random subset of 50 and 100 data points respectively, train for 100 epochs to assure convergence and then make a query to include 100 further data points. We then proceed to train the network with the additional instances before repeatedly querying and training again. In a crucial distinction to the majority of active learning works that only investigate the quality of the query by re-training the entire model from scratch, we do not reset our weights in continued incremental training. This implicitly introduces a stronger impact of ordering and further acknowledges that not only labelling, but also training itself is expensive. Each experiment is repeated five times, alas always with the same initial random subset to preserve comparability between individual repetitions and across methods.

We investigate popular metrics and mechanisms on which current deep active learning is based. The majority of these are techniques that attempt to take optimal action without explicitly approximating the entire set of unknown data. To estimate and account for uncertainty we make use

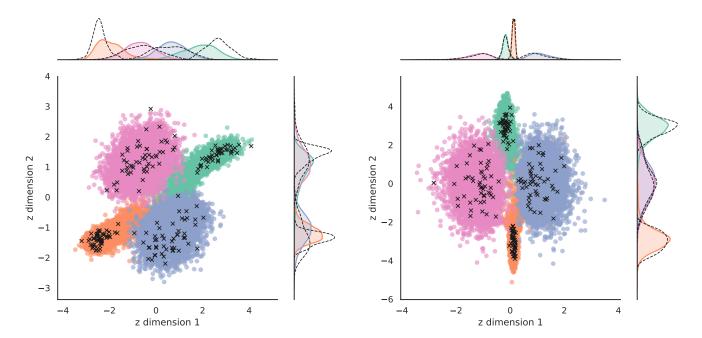


Fig. 10: Visualization of the aggregate posterior for a model with two-dimensional latent space trained on the first four classes of the CIFAR10 dataset and 200 selected core set exemplars. The left panel shows the greedy k-center approach, whereas the right panel shows our proposed EVT based core set construction. Classes are color coded points and the core set elements are illustrated through black crosses. A kernel density estimate of the per class aggregate posterior (in color) and the corresponding distributional approximation of the selected core set elements (dashed black) are added on each dimension. In contrast to the greedy k-center approach that features large discrepancies, insignificant differences are observable for our proposed method, painting an intuitive picture for our methods quantitative success of figure 9.

of Monte Carlo Dropout (MCD) [37] where appropriate. Although we believe that there is an inherent limitation in earlier introduced approaches that explicitly use the entire unlabelled pool for optimization, we also investigate the proposed technique to query based on a k-means core set extracted from the unknown data [118], [120]. Whereas we certainly regard such methods as valuable in a closed world context, we note that these methods are infeasible without prior knowledge outside of a constrained pool or for sequentially arriving data subsets. As we will see in the next section, they feature little robustness to nonsensical data that might be present in the pool, as the entire unlabelled pool is included and assumed to be useful. The metrics and methods that we investigate are:

- 1) **Random:** sampling uniformly at random from the unlabelled pool.
- 2) Reconstruction loss: in our particular scenario, because our proposed framework includes a generative model, we can query examples based on largest reconstruction loss. This is typically unavailable in a purely discriminative neural network classifier.
- 3) **K-means core set:** use the entire unlabelled pool to base the query on an extracted core set that is equivalent in size to the query amount. Nguyen et al. had suggested such pre-clustering [118] and it was later used in deep active learning with k-means as the core set algorithm [120].
- 4) MCD classification confidence: query based on lowest softmax confidence [107]. As neural network

- classifiers are known to be overconfident, we additionally gauge uncertainty with MCD as a suggested remedy by Gal et al. [110].
- 5) MCD classification entropy: query based on largest predictive entropy [105]. Similar to lowest confidence, we use uncertainty from MCD to obtain better entropy estimates [110].
- 6) Latent EVT: our proposed EVT based approach that balances exploration with exploitation by querying instances that distribute across outlier probabilities, but limited by an upper rejection prior to avoid uninformative outliers.

We first note that we have included classification confidence and entropy with MCD because omitting uncertainty estimates resulted in no improvement of the active learning query upon simple random selection. This has previously been argued and corresponds to the empirical observations made by Sinha et al. [122]. For our proposed EVT approach we empirically distribute the query uniformly across examples that fall into the range of 0.5 to 0.95 outlier probability, as estimated by equation 6. Although it never occurred in practice, we note that it would likely be preferential to extend this range to the lower end if not enough samples in the pool were available in the mentioned range, rather than including complete outliers. We will provide empirical evidence for this in the next section.

Figure 11 shows the quantitative results of our active learning experiments. On both datasets, the k-means based core set is either similar or slightly worse than simply

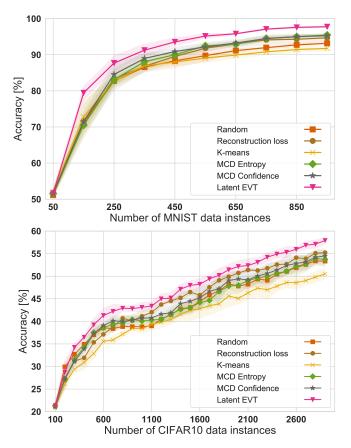


Fig. 11: Active learning accuracy for different methods on the MNIST and CIFAR10 datasets. All experiments start with the same randomly sampled 50 and 100 dataset examples. In each step, an additional 100 data instances are queried from the remaining unlabelled pool and included for further continued training. Results show the average over five experiments, with the shaded areas ranging from the minimum to the maximum obtained values.

sampling at random. This reflects our previous observations in the core set continual learning section. On the contrary, the uncertainty based methods surpass random sampling. Using largest reconstruction loss similar results can be accomplished, although at the additional computational expense of calculating the decoding. However, all methods are significantly outperformed by our proposed latent EVT method at all times. The respective rationale behind this improvement is quite intuitive. Our strategy balances completely novel examples with less novel examples that are still required to strengthen the existing learned features. More importantly, it rejects uninformative outliers that are inherently present in the pool, a threat that uncertainty based methods can be particularly prone to. This threat is magnified with even less knowledge about the acquired dataset and even more unconstrained data acquisition. The past two subsections have focused on showing our methods advantage in the typical continual and active learning benchmark perspective in the closed world scenario, devoid of any analysis with respect to robustness. In the next section we extend this evaluation to analyze each individual methods' behavior in the presence of corruptions.

## 6.3 Robustness to open world corruptions

Prior works that address open set recognition or in general application of machine learning algorithms in an open world have argued that prediction on previously unseen unknown classes results in inevitable misprediction [15], [67], [127], [128], [131]. For example, if a user is given the freedom to provide any image input to a neural network based classifier, an arbitrarily chosen image's prediction will be indistinguishable from the typical training set output. We have previously empirically demonstrated that the proposed EVT based approach overcomes this challenge, much in contrast to relying on uncertainty based measures that fail to even distinguish the most trivially disparate datasets such as visual and audio data [39], [40]. Although this poses a serious threat to building a user's trust, just imagine your own faith in a classifier that assigns an image of a car the label of a t-shirt (recall the earlier figure 5), we can naturally question if this scenario could simply be circumvented by including guidelines with respect to the expected model input, i.e. "this model has been trained on fashion-items, it is not designed for other types of data". The more sensible solution would be to have the model reject unknown unknown data. Whether or not we consider the latter scenario as meaningful, unknown unknown data is not necessarily always composed of completely dissimilar classes. A perhaps at least equivalently large threat is data that is statistically deviating for other reasons: corruption and perturbation. In any real-world scenario, we can no longer assume that our machine learning model is faced exclusively with the carefully curated data that benchmarks are comprised off. Often a simple change in camera can dramatically skew the statistics of the acquired image. In an almost endless list, low lighting conditions can introduce various forms of noise, small jitter can cause blur, weather conditions change, the condition of the object of interest correspondingly changes, etc.. The gullible solution would again be to attempt to model all forms of corruptions and perturbations, but this simply connects back to the infeasibility of the earlier introduced "inference with the universum" approach.

In a recent effort to benchmark the performance against 15 types of various corruptions, Hendrycks and Dietterich [16] have shown that none of the developed neural network models feature any intrinsic robustness, even if they converge to more accurate solutions on the initial benchmark. This was concluded from experiments where neural networks are trained on the uncorrupted benchmark dataset and evaluated on the corrupted data. We extend this evaluation by investigating the presence of a minor portion of corrupted data in the training process, as can realistically be assumed for active or continual learning. We examine whether common query strategies in active learning and core set construction in continual learning are robust, or whether querying and including this unrepresentative corrupted data into core sets leads to performance degradation in comparison with the clean benchmark. We believe that this is critical for two reasons: 1.) The necessity to carefully curate every single example in the unknown data pool can outweigh the active learning human labelling effort and thus renders active learning ineffective in the first place. 2.) Data cleaning itself is extremely challenging and it is often not immediately clear

TABLE 1: Active learning with and without partial dataset corruption. Uncorrupted values correspond to those visualized in figure 11.

|                               | Accuracy [%]: meandifference to minimum |                         |                         |                         |                         |                         |  |  |  |
|-------------------------------|---|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|--|--|--|
| CIFAR10 queries, dataset size | 8,900                                   |                         | 18, 1900                |                         | 28, 2900                |                         |  |  |  |
| Dataset                       | regular                                 | corrupted               | regular                 | corrupted               | regular                 | corrupted               |  |  |  |
| Random                        | $38.80^{+0.69}_{-1.75}$                 | $38.97^{+1.03}_{-1.87}$ | $47.81^{+2.02}_{-3.93}$ | $47.91^{+2.13}_{-3.58}$ | $53.36^{+1.17}_{-2.34}$ | $53.53^{+1.13}_{-2.42}$ |  |  |  |
| Reconstruction loss           | $41.14^{+2.06}_{-3.89}$                 | $38.26^{+0.64}_{-1.89}$ | $50.70^{+0.69}_{-1.50}$ | $46.49^{+0.82}_{-2.13}$ | $55.22^{+1.37}_{-1.92}$ | $50.85^{+1.03}_{-1.57}$ |  |  |  |
| K-means                       | $38.34_{-2.63}^{+1.46}$                 | $36.05^{+1.65}_{-2.53}$ | $45.08^{+1.50}_{-3.23}$ | $42.93^{+1.59}_{-3.65}$ | $50.52^{+0.94}_{-3.15}$ | $47.58^{+1.93}_{-3.39}$ |  |  |  |
| MCD Entropy                   | $40.05^{+1.15}_{-2.99}$                 | $38.83^{+0.68}_{-1.03}$ | $47.96^{+2.91}_{-5.28}$ | $44.73^{+0.61}_{-1.02}$ | $53.72^{+2.35}_{-4.76}$ | $50.06^{+0.37}_{-0.75}$ |  |  |  |
| MCD Confidence                | $40.67^{+0.87}_{-1.89}$                 | $37.93^{+0.35}_{-0.81}$ | $49.40^{+2.86}_{-4.44}$ | $47.16^{+1.29}_{-3.22}$ | $54.51^{+1.15}_{-3.13}$ | $51.91^{+1.78}_{-2.67}$ |  |  |  |
| Latent EVT                    | $44.67^{+0.32}_{-0.63}$                 | $43.79^{+0.74}_{-1.72}$ | $51.66^{+1.05}_{-1.69}$ | $51.12^{+0.38}_{-0.91}$ | $57.43^{+0.51}_{-1.09}$ | $56.83^{+0.41}_{-0.78}$ |  |  |  |

TABLE 2: Coreset selection and training with and without dataset corruption. Uncorrupted values correspond to those visualized in figure 9.

|                                 | Accuracy [%]: meandifference to maximum         |  |   |   |   |  |  |  |  |
|---------------------------------|---|--|---|---|---|--|--|--|--|
| CIFAR10 coreset size            | 300   |  | 60  | 00  | 1500  |  |  |  |  |
| Dataset                         | regular   | corrupted  | regular   | corrupted                                       | regular   | corrupted  |  |  |  |
| Random<br>Greedy k-center       | $31.23_{-9.14}^{+3.94} \\22.82_{-1.65}^{+3.05}$ | $30.35_{-5.92}^{+1.88} \\ 22.19_{-3.37}^{+1.76}$ | $\begin{array}{ c c c c c c }\hline 39.52^{+3.61}_{-7.95} \\ 29.33^{+1.50}_{-3.23} \\ \hline \end{array}$ | $39.05^{+1.99}_{-5.89}$ $29.48^{+1.91}_{-5.11}$ | $\begin{array}{ c c c c c }\hline 51.43^{+3.33}_{-6.12} \\ 42.41^{+1.97}_{-4.13} \\ \hline \end{array}$ | $51.01_{-4.49}^{+2.30}$ $42.37_{-2.44}^{+1.49}$  |  |  |  |
| Latent k-means                  | $32.76^{+2.29}_{-3.35}$                         | $29.00^{+2.12}_{-4.05}$                          | $39.49^{+1.71}_{-4.17}$   | $35.71^{+1.69}_{-4.08}$                         | $50.01^{+1.80}_{-3.28}$   | $48.52^{+2.59}_{-3.86}$                          |  |  |  |
| Image k-means<br>Latent herding | $32.85_{-3.76}^{+2.57}$ $33.92_{-1.45}^{+0.61}$ | $30.74_{-3.16}^{+1.43}$ $33.81_{-1.39}^{+0.82}$  | $\begin{array}{c c} 37.86^{+1.66}_{-3.98} \\ 41.13^{+1.18}_{-2.29} \end{array}$                           | $36.38^{+0.90}_{-2.75}$ $40.77^{+1.34}_{-1.57}$ | $\begin{array}{c c} 49.62^{+2.83}_{-8.09} \\ 51.87^{+1.12}_{-1.85} \end{array}$                         | $48.23_{-2.50}^{+1.78} \\ 51.06_{-2.30}^{+2.43}$ |  |  |  |
| Latent EVT                      | $34.16^{+1.10}_{-2.27}$                         | $34.18^{+1.07}_{-2.55}$                          | $\begin{array}{ c c c c c c }\hline 41.13_{-2.29} \\ 41.78_{-2.57}^{+1.34} \\ \end{array}$                | $41.67^{+1.57}_{-2.53}$                         | $\begin{array}{ c c c c c c }\hline 51.87_{-1.85} \\ 53.35_{-2.53}^{+1.48} \\ \hline \end{array}$       | $53.28^{+1.06}_{-2.17}$                          |  |  |  |

whether the inclusion of a data instance is beneficial or is accompanied by side effects.

We make use of corruptions across four categories: noise, blur, weather and digital corruptions, as introduced by Hendrycks and Dietterich [16]. These can further be distinguished into 15 types: low-lighting Gaussian noise, electronic shot noise, bit error impulse noise, speckle noise, Gaussian blur, defocus blur, glass blur, zoom blur, motion blur, snow, fog, brightness, contrast, saturation and elastic deformations. Each corruption is algorithmically generated with five discretized levels of severity, of which the first two are at times barely discernible from a typical image by a human. We accordingly corrupt 7.5% of the data across these 75 corruptions. We add the additional constraint that each image can only be corrupted once. Note that in principle some corruptions, such as noise resulting from low lighting conditions and out of focus blurring, could occur simultaneously. We have deliberately chosen this amount of corruption to, on the one hand be small enough to not affect overall performance if trained on the entire dataset, on the other hand be larger than the core set size or active learning query amounts used in previous sections. Hypothetically, in the absolute worst case this could result in only corrupted images being selected and the entire chosen set being much less representative of the complete dataset than a selection of clean examples would be. We repeat the previous CIFAR10 experiments under these conditions. For better visualization and quantification we do not show plots, but have instead picked three evenly spaced points of figures 9 and 11.

We show the originally obtained results in direct com-

parison with the results obtained under inclusion of the corrupted data in tables 1 and 2. From these quantitative results it is evident that only two techniques are robust in active learning: random sampling and our proposed EVT based approach. The logical explanation is that random sampling on average will pick roughly 7.5% corrupted data, of which another 40% feature only minor low severities. The small amount thus only has minor effect on the optimization. The EVT based algorithm is similarly unaffected as it does not query statistical outliers in the first place, or if it includes corrupted examples then only those with minor severity that are statistically still largely similar to the uncorrupted data. All other methods are prone to the corrupted outliers in one way or another. Classifier uncertainty and reconstruction loss tend to pick very corrupted examples by definition, the k-means approach will have shifted centers or falsely query from new clusters that are centered around corruptions of the unknown pool. Looking at the quantitative accuracy values, we can in fact even conclude that all these methods perform worse than a simple random query. The continual learning core set construction picture is quite similar. Here, we can observe corruption robustness for random sampling, latent herding and our proposed approach. Latent herding is robust to outliers because it picks samples greedily one by one to best approximate the mean, which intuitively involves picking the next best example that is close to the class mean and does not involve outliers (potentially only in a minor fashion through a drifted mean if the outliers are not embedded symmetrically around the class mean). However, the issue of including redundant samples into the





Fig. 12: Typically selected dataset examples in the core set construction using a greedy k-center algorithm. Qualitative illustration is intended to provide intuition for a method's failure. The left panel shows how picked exemplars from an uncorrupted dataset are unrepresentative of the average image, with unusual backgrounds, occlusion and scaling issues. The right panel shows how the core set is comprised of many corrupted examples if a small portion of the dataset is corrupted, a lack of robustness that many methods in tables 1 and 2 suffer from.

core set remains unaddressed, and our EVT based method nevertheless outperforms all other approaches.

Interestingly, the greedy k-center approach also seems to be robust to the corruptions, although it performs equally miserably to the uncorrupted scenario. Recall that this algorithm greedily chooses the next data point for inclusion in a farthest-first traversal, by maximizing the distance to all presently existing core set elements. In other words, outliers are always queried as they by definition are farthest away. Only after a sufficiently large cover is obtained will representative data be queried. Because such unrepresentative outliers are already present in the uncorrupted data, the performance is consequently always low for small core set sizes. To visually illustrate this statement we show a uniform sub-sample of the acquired core set for the first four classes with and without corruption in figure 12. In the left panel we can observe the core set being comprised of atypical airplanes with deep green or black background, a captured overexposed sunset, partially occluded cars and birds by bushes and fences or images where the animal is almost not discernible and comprises only a fraction of the image. Arguably these do not represent good exemplars. In the right panel, we can see that in the presence of corruption, the core set is comprised of noisy, blurry and otherwise distorted images. Ultimately neither of these core sets are a particularly good approximation of the dataset, intuitively explaining the abysmal performance of this technique.

## 6.4 Choosing the curriculum - the importance of task order

As detailed in the earlier introduction of our framework, we can apply our proposed EVT based active learning strategy to the construction of a continual, class incremental learning curriculum. In this context, a task's outlier probability is synonymous with its dissimilarity to already accumulated tasks. Conversely, a task that is deemed to be largely inlying has a large representational overlap with existing knowledge, even though it might have been assigned a distinct label. In

the best case scenario, this implies that only fine-tuning is necessary to sufficiently include a proximate task. In the worst case scenario, the representational entanglement severely limits the discriminability. Unless a major addition or overhaul of the learned representations ensues, this leads to confusion with existing concepts. In contrast, most outlying tasks are hypothesized to be distinct enough to not interfere with previous tasks, assuming the old task's data is still available or a continual learning mechanism prevents its catastrophic forgetting.

We investigate the importance of task order and whether the construction of a curriculum beyond alphabetical class order provides substantial learning benefits. For this purpose we consider four conceivable scenarios:

- Class sequential ordering: learn the classes in order of their integer class label. For many datasets this is in alphabetical order.
- 2) Random order: randomized class order.
- 3) **Most outlying, dissimilar tasks first:** determine the next class to add by evaluating equation 8, i.e. pick the next class that is most outlying and dissimilar with respect to the already seen classes.
- 4) Most inlying, similar tasks first: determine the next class to add by evaluating equation 8, but with a minimum over task outlier probabilities to include the most similar task in each increment.

Note that for all strategies we always start with the same first task for comparability. To make sure that obtained results and found curricula are not just a result of sheer luck, we repeat each experiment five times, report the average and the minimum and maximum obtained accuracies at each step to gauge deviations. We conduct experiments on two datasets: the CIFAR100 and the AudioMNIST [174] dataset. We follow the typical continual incremental learning procedure of adding classes in pairs of two. We chose the first dataset because it allows for the construction of a long task sequence. We chose the latter because it represents a nonimage dataset and previous work has observed that some classes can provide strong retrospective improvement [39], an early indicator that the class ordering should be investigated further. In order to show the impact of task ordering, we provide an analysis, both, when independently evaluated from, or coupled to specific techniques that alleviate continual learning catastrophic forgetting. As such, we evaluate CIFAR100 in what is typically referred to as a continual learning upper-bound, i.e. the maximum obtainable accuracy given a specific model choice and training procedure in which the data of each task is simply accumulated with each subsequent task. For the AudioMNIST we use generative replay to prevent catastrophic forgetting, where old tasks' data is rehearsed based on the trained generative model. We do not make use of any data augmentation.

The achieved accuracies at each task increment are shown in figure 13. We can observe that for the CIFAR100 dataset, random sampling seems to yield a very similar accuracy trajectory in comparison to sequentially learning the classes in order of their alphabetical class id, resembling earlier observations [8], [162]. However, in contrast to the conclusion that task order is negligible, we can observe that our proposed framework's selection schemes, that rank order the data

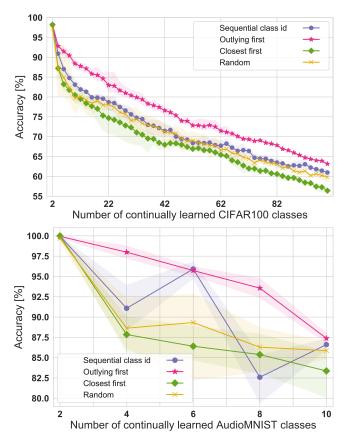


Fig. 13: Continual learning accuracy of learning classes in increments of two in dependence on the choice of task order. Top panel shows the incremental upper-bound, i.e. a simple accumulation of the real data, for the CIFAR100 dataset. The bottom panel shows obtained performance on the AudioMNIST dataset with alleviated catastrophic forgetting through generative replay. For each of the order selection mechanisms the experiment has been repeated five times. The corresponding average together with the maximum and minimum deviation are reported respectively.

according to their similarity with the existing encoding, paint a dramatically different picture. Selecting the most dissimilar task for inclusion consistently improves the accuracy by several percent, even at the end of training. Conversely, including tasks that are very proximate to existing concepts results in an all-time performance decrease. We hypothesize that this is due to the classifier experiencing immediate confusion. Our initial classes consist of "apples" and "aquarium fish" and the query consensus across repeated experiments is to continue with selecting the classes "pears" and "whale" or "shark". The opposite strategy that prioritizes dissimilarity in the curriculum instead includes unrelated classes such as "lawnmower", "mountain" or "oak". We believe that this allows the model to more rapidly acquire a diverse set of representations.

We can draw almost analogous conclusions for continually learning the AudioMNIST dataset with generative replay. Here, we additionally see that the conventional order of learning the sounds from "zero" to "nine" is accompanied by a pattern of repeated retrospective improvement. The first task increment results

in a larger accuracy drop, that is rectified through backwards improvement of the next task increment. This pattern repeats for the next two classes and its consistent strong emergence is only visible when learning sequentially in order of class id. The accuracy at any time is again best for our proposed measure of dissimilarity and worst when selecting according to task proximity. For the latter, in analogy to the earlier hypothesized confusion of the classifier, the generative model is faced with difficulty to disambiguate the resembling classes and produce unambiguous output.

Our results indicate that using active learning techniques in continual learning can have critical impact on the achieved performance. More so, the results provide an important signal for reproducibility and significance of various conjured continual learning benchmarks. In a world of benchmarking methods and regularly claiming advances when a method surpasses another by 1-2 %, the observed absolute discrepancy between the different task orders for CIFAR100 is as large as 10%. This is a substantial gap. Whereas we obviously believe that there is value in analyzing and contrasting different techniques to alleviate catastrophic forgetting on a common dataset, it is clear that there is still much we need to learn about neural network training and evaluation that can only be discovered by moving away from our current rigid benchmarks.

## 7 CONCLUSION: TOWARDS A WHOLISTIC DEFINITION OF DEEP CONTINUAL LEARNING

We have presented a common viewpoint to naturally unite robust continual and active learning in the presence of the unknown. For each aspect, we have conducted an empirical investigation that demonstrated the benefits of the viewpoint's realization in a variational Bayesian deep neural network framework. Needless to say, each of our individually presented experiments can be extended with multiple facets and several nuanced applications can be derived and thoroughly investigated. At this point, we remark that we do not wish to claim that our proposed method provides the generally best solution or selects optimal task sequences. Although our framework clearly shows quantitative promise, our main goal is to highlight the importance of the introduced consolidated viewpoint. In the ideal case, we would encourage future works to adopt our framework or take a similarly wholistic approach. At the very minimum, we would expect future works to rethink current practices and question whether current benchmarks are a realistic reflection of our desiderata for continual machine learning systems. As illustrated throughout the paper, this necessitates stepping out of our closed world benchmark routines. In hopes of providing some guidelines for the latter, we make an attempt at a revised continual learning definition and suggestions towards more systematic assessment.

**Definition 7.1.** Continual Machine Learning - this work :The learner performs a sequence of N continual learning tasks,  $\mathcal{T}_1, \mathcal{T}_2, \ldots, \mathcal{T}_N$ , that are distinct from each other in terms of shifts in the underlying data distribution. The latter can imply a change in objective, transitions between different domains or inclusion of new modalities. At any point in time, the learner must be able to robustly identify unseen unknown data instances and

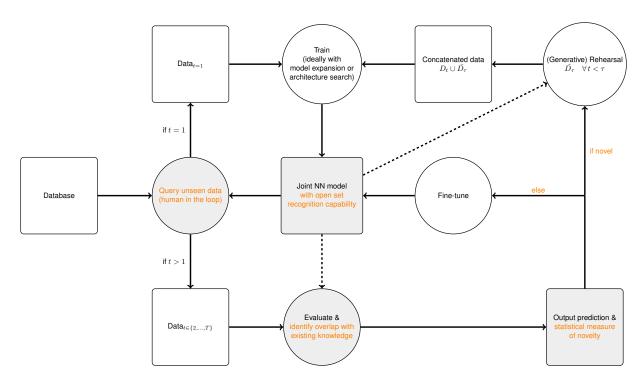


Fig. 14: A suggestion for a more comprehensive, system oriented evaluation. In contrast to the conventional continual learning pipeline, the system is extended with a (optionally human in-the-loop) data querying mechanism and a measure of novelty that is used for robust application in the open world and to select adequate ensuing optimization techniques. These suggested additions to the conventional continual learning process are emphasized through orange text and shading in the diagram. Rectangles represent objects and circles correspond to processes. Dashed arrows indicate a process' dependency on the model.

rank order them according to similarity with existing tasks, in order to actively build a learning curriculum. If the system is desired to be supervised, a human in the loop may group and label the set of identified unseen unknowns to explicitly guide future learning. When faced with a selected (N+1)th task  $\mathcal{T}_{N+1}$  (which is called the new or current task) with its data  $\mathcal{D}_{N+1}$ , the learner should leverage its dictionary of representations to accelerate learning of  $\mathcal{T}_{N+1}$  (forward transfer), extend the dictionary with unique representations obtained from the new task's data (this can be completely new types of dictionary elements), while simultaneously maintaining and improving the existing representational dictionary with respect to former tasks (backward transfer).

In comparison with former continual learning definitions, reiterated at the beginning of this paper, the definition is now extended to include active data queries, the corresponding importance of data choice and task order, in coherence with awareness of the open world.

## 7.1 Outlook: a suggestion for a more comprehensive, system oriented evaluation

We show one example of how a revised outline of a continual learning system that satisfies the above definition could look like in figure 14. Again, this example can be realized with our proposed specific EVT based framework, although several other implementations are conceivable. The main idea of the system can be summarized as follows: After initial training on some seeding data, ideally by finding a baseline architecture through architecture search or through progressive architecture growth, a new task is queried

through an inherent model mechanism to optimize the effect of order and the queried data is consecutively labelled. Alternatively, specific data can be introduced by a human in the loop, if it is desired that the system is constrained to very specific tasks. The new data is then evaluated with respect to existing tasks and associated with a measure of novelty. This measure of novelty serves the dual purpose of introducing robustness into the system when applied in the wild, and at the same time is used as the foundation to decide on how to proceed with further optimization. If the overlap with existing knowledge is very large, it is sufficient to conduct minor fine-tuning steps. If there is a large amount of expected novelty, the optimization needs to proceed with a mechanism to protect previously acquired tasks, typically through means of core set or generative rehearsal. Because the amount of expected novelty is large, it is recommended to then continue training with model expansion in order to ensure sufficient representational capacity is available to accommodate entirely new concepts. The cycle is then repeated.

In comparison with the classical continual learning evaluation pipeline, presented in the beginning of this work in figure 2, we thus suggest to extend the system with essential robust evaluation and active queries to address questions concerning the importance of input data selection. As demonstrated, integration of these aspects can be achieved through prediction of a statistical measure of novelty based on overlap with existing knowledge, e.g. with

our suggested posterior based EVT open set recognition approach. This measure of novelty serves a natural triple purpose: 1.) Rejection or setting aside of unknown unknown data in robust application. 2.) Querying data from an unlabelled pool in a suitable order that provides large expected benefit to the model. 3.) If the data order is pre-imposed, e.g by a human or a stream, the novelty metric can be used to dynamically switch the training procedure to incorporate dissimilar novel data, while preserving prior representations through extensive continual learning mechanisms that alleviate catastrophic forgetting, or to simply fine-tune in the presence of sufficient overlap with previously seen data.

Even though the advantages of expanding the effective representational capacity during training are clear, we have put the use of model expansion and progressive architecture search in brackets. Although its use is theoretically and empirically desirable, we understand that this ideal evaluation involves several challenges that can limit its practicality. It is clear from previously discussed works that continuous model growing is advantageous, but we note that heavily over-parametrized models have shown satisfactory results. We thus encourage future research to first and foremost focus on the questions about benchmark construction, data point selection, and the voiced concerns regarding robust application in the open world. We would then expect future work to additionally include model expansion techniques.

We anticipate that this work leads to increased awareness of the dangers of our current closed world practices and the necessity of expanding our views towards more realistic real-world relevant evaluation. In doing so, we believe that further synergies between presently separately treated machine learning paradigms will be exposed and can be exploited. This should ultimately lead to improved, more robust and simpler machine learning systems.

## REFERENCES

- [1] S. Thrun, "Is Learning The n-th Thing Any Easier Than Learning The First?" Advances in Neural Information Processing Systems, 1996.
- [2] S. Thrun, Explanation-Based Neural Network Learning A Lifelong Learning Approach. Springer US, 1996.
- [3] Z. Chen and B. Liu, *Lifelong Machine Learning*. Morgan and Claypool, 2017, vol. 33.
- [4] M. McCloskey and N. J. Cohen, "Catastrophic Interference in Connectionist Networks: The Sequential Learning Problem," Psychology of Learning and Motivation - Advances in Research and Theory, vol. 24, no. C, pp. 109–165, 1989.
- [5] R. Ratcliff, "Connectionist Models of Recognition Memory: Constraints Imposed by Learning and Forgetting Functions," Psychological Review, vol. 97, no. 2, pp. 285–308, 1990.
- [6] S. Farquhar and Y. Gal, "Towards Robust Evaluations of Continual Learning," International Conference on Machine Learning (ICML), Lifelong Learning: A Reinforcement Learning Approach Workshop, 2018.
- [7] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual Lifelong Learning with Neural Networks: A Review," Neural Networks, vol. 113, pp. 54–71, 2019.
- [8] M. De Lange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and T. Tuytelaars, "A continual learning survey: Defying forgetting in classification tasks," arXiv preprint arXiv: 1909.08383, 2019.
- [9] T. Lesort, V. Lomonaco, A. Stoian, D. Maltoni, D. Filliat, and N. Díaz-Rodríguez, "Continual learning for robotics: Definition, framework, learning strategies, opportunities and challenges," *Information Fusion*, vol. 58, no. September 2019, pp. 52–68, 2020.

- [10] T. Lesort, H. Caselles-Dupré, M. Garcia-Ortiz, A. Stoian, and D. Filliat, "Generative Models from the perspective of Continual Learning," *International Joint Conference on Neural Networks* (IJCNN), 2019.
- [11] B. Pfülb and A. Gepperth, "A Comprehensive, Application-Oriented Study of Catastrophic Forgetting in DNNs," *International Conference on Learning Representations (ICLR)*, 2019.
- [12] R. Kemker, M. McClure, A. Abitino, T. Hayes, and C. Kanan, "Measuring Catastrophic Forgetting in Neural Networks," AAAI Conference on Artificial Intelligence (AAAI), 2018.
- [13] N. Díaz-Rodríguez, V. Lomonaco, D. Filliat, and D. Maltoni, "Don't forget, there is more than forgetting: new metrics for Continual Learning," Neural Information Processing Systems (NeurIPS), Continual Learning Workshop, 2018.
- [14] O. Matan, R. Kiang, C. E. Stenard, and B. E. Boser, "Handwritten Character Recognition Using Neural Network Architectures," 4th USPS Advanced Technology Conference, vol. 2, no. 5, pp. 1003–1011, 1990.
- [15] T. E. Boult, S. Cruz, A. Dhamija, M. Gunther, J. Henrydoss, and W. Scheirer, "Learning and the Unknown: Surveying Steps Toward Open World Recognition," AAAI Conference on Artificial Intelligence (AAAI), 2019.
- [16] D. Hendrycks and T. Dietterich, "Benchmarking neural network robustness to common corruptions and perturbations," *Interna*tional Conference on Learning Representations (ICLR), 2019.
- [17] H. Shin, J. K. Lee, and J. J. Kim, "Continual Learning with Deep Generative Replay," Neural Information Processing Systems (NeurIPS), 2017.
- [18] A. Achille, T. Eccles, L. Matthey, C. P. Burgess, N. Watters, A. Lerchner, and I. Higgins, "Life-Long Disentangled Representation Learning with Cross-Domain Latent Homologies," Neural Information Processing Systems (NeurIPS), 2018.
- [19] S. Farquhar and Y. Gal, "A Unifying Bayesian View of Continual Learning," Neural Information Processing Systems (NeurIPS) Bayesian Deep Learning Workshop, 2018.
- [20] C. V. Nguyen, Y. Li, T. D. Bui, and R. E. Turner, "Variational Continual Learning," International Conference on Learning Representations (ICLR), 2018.
- [21] C. Wu, L. Herranz, X. Liu, Y. Wang, J. van de Weijer, and B. Raducanu, "Memory Replay GANs: learning to generate images from new categories without forgetting," *Neural Information Processing Systems (NeurIPS)*, 2018.
- [22] M. Zhai, L. Chen, F. Tung, J. He, M. Nawhal, and G. Mori, "Lifelong GAN: Continual Learning for Conditional Image Generation," International Conference on Computer Vision (ICCV), 2019.
- [23] Z. Li and D. Hoiem, "Learning without forgetting," European Conference on Computer Vision (ECCV), 2016.
- [24] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell, "Overcoming catastrophic forgetting in neural networks," Proceedings of the National Academy of Sciences (PNAS), vol. 114, no. 13, pp. 3521–3526, 2017
- [25] S. A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "iCaRL: Incremental classifier and representation learning," Computer Vision and Pattern Recognition (CVPR), 2017.
- [26] D. Lopez-Paz and M. A. Ranzato, "Gradient Episodic Memory for Continual Learning," Neural Information Processing Systems (NeurIPS), 2017.
- [27] R. Kemker and C. Kanan, "FearNet: Brain-inspired model for incremental learning," *International Conference on Learning Repre*sentations (ICLR), 2018.
- [28] Y. Xiang, Y. Fu, P. Ji, and H. Huang, "Incremental Learning Using Conditional Adversarial Networks," *International Conference on Computer Vision (ICCV)*, 2019.
- [29] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2323, 1998.
- [30] A. Krizhevsky, "Learning Multiple Layers of Features from Tiny Images," Toronto, Tech. Rep., 2009.
- [31] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
- 32] N. Roy, A. Mccallum, and M. W. Com, "Toward optimal active learning through monte carlo estimation of error reduction."

- Proceedings of the International Conference on Machine Learning (ICML), pp. 441–448, 2001.
- [33] B. Settles and M. Craven, "An analysis of active learning strategies for sequence labeling tasks," Empirical Methods in Natural Language Processing (EMNLP), pp. 1070–1079, 2008.
- [34] X. Li and Y. Guo, "Adaptive active learning for image classification," Computer Vision and Pattern Recognition (CVPR), pp. 859–866, 2013.
- [35] W. H. Beluch, T. Genewein, A. Nürnberger, and J. M. Köhler, "The Power of Ensembles for Active Learning in Image Classification," Computer Vision and Pattern Recognition (CVPR), 2018.
- [36] Y. Geifman and R. El-Yaniv, "Deep Active Learning with a Neural Architecture Search," Neural Information Processing Systems (NeurIPS), 2019.
- [37] Y. Gal and Z. Ghahramani, "Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning," *International Conference on Machine Learning (ICML)*, vol. 48, 2015.
- [38] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting," *Journal of Machine Learning Research* (*JMRL*), vol. 15, pp. 1929–1958, 2014.
- [39] M. Mundt, S. Majumder, I. Pliushch, and V. Ramesh, "Unified Probabilistic Deep Continual Learning through Generative Replay and Open Set Recognition," arXiv preprint arXiv:1905.12019, 2019.
- [40] M. Mundt, I. Pliushch, S. Majumder, and V. Ramesh, "Open Set Recognition Through Deep Neural Network Uncertainty: Does Out-of-Distribution Detection Require Generative Classifiers?" International Conference on Computer Vision (ICCV), First Workshop on Statistical Deep Learning for Computer Vision (SDL-CV), 2019.
- [41] R. E. Schapire, "The Strength of Weak Learnability," *Machine Learning*, vol. 5, no. 2, pp. 197–227, 1990.
- [42] L. Pratt, J. Mostow, and C. Kamm, "Direct Transfer of Learned Information Among Neural Networks," AAAI Conference on Artificial Intelligence (AAAI), 1991.
- [43] L. Y. Pratt, "Discriminability-Based Transfer between Neural Networks," Neural Information Processing Systems (NeurIPS), 1993.
- [44] Y. Freund and R. Schapire, "A decision theoretic generalisation of online learning and an application to boosting," *Computer and System Sciences*, vol. 55, no. 1, pp. 119–139, 1997.
- [45] S. J. Pan and Q. Yang, "A Survey on Transfer Learning," *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, vol. 22, no. 10, 2010.
- [46] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, "Learning and transferring mid-level image representations using convolutional neural networks," Computer Vision and Pattern Recognition (CVPR), 2014
- [47] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?" Neural Information Processing Systems (NeurIPS), 2014.
- [48] K. Weiss, T. M. Khoshgoftaar, and D. D. Wang, "A survey of transfer learning," *Journal of Big Data*, vol. 3, no. 1, 2016.
- [49] R. Caruana, "Multitask Learning," Machine Learning, vol. 28, pp. 41–75, 1997.
- [50] S. C. Suddarth and Y. L. Kergosien, "Rule-injection hints as a means of improving network performance and learning time," Neural Networks. EURASIP 1990. Lecture Notes in Computer Science, vol. 412, 1990.
- [51] Y. S. Abu-Mostafa, "Learning from hints in neural networks," Journal of Complexity, vol. 6, no. 2, pp. 192–198, 1990.
- [52] S. Ruder, "An Overview of Multi-Task Learning in Deep Neural Networks," arXiv preprint arXiv: 1706.05098, 2017.
- [53] L. Fei-Fei, R. Fergus, and P. Perona, "One-shot learning of object categories," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 28, no. 4, pp. 594–611, 2006.
- [54] Y. Wang, Q. Yao, J. Kwok, and L. M. Ni, "Generalizing from a Few Examples: A Survey on Few-Shot Learning," ACM Computing Surveys, 2020.
- [55] M. Fink, "Object classification from a single example utilizing class relevance metrics," Neural Information Processing Systems (NeurIPS), 2005.
- [56] C. H. Lampert, H. Nickisch, and S. Harmeling, "Learning to detect unseen object classes by between-class attribute transfer," Computer Vision and Pattern Recognition (CVPR), 2009.
- [57] D. O. Hebb, The Organization of Behavior; A Neuropsychological Theory. John Wiley & Sons, Chapman & Hall, 1949.
- [58] Y. Tsypkin, Adaptation and Learning in Automatic Systems. New York: Academic Press, 1971.

- [59] V. Vapnik, Estimation of Dependences Based on Empirical Data: Springer Series in Statistics. Berlin, Heidelberg: Springer-Verlag, 1982.
- [60] T. M. Heskes and B. Kappen, "On-line learning processes in artificial neural networks," *Mathematical Foundations of Neural Networks*, vol. 51, no. C, pp. 199–233, 1993.
- [61] L. Bottou, "Online Learning and Stochastic Approximations," in *Online Learning in Neural Networks*, 1999, pp. 9–42.
- [62] D. Saad, Ed., On-line learning in neural networks. New York, NY, USA: Cambridge University Press, 1999.
- [63] M. Zinkevich, "Online Convex Programming and Generalized Infinitesimal Gradient Ascent," International Conference On Machine Learning (ICML), 2003.
- [64] G. Zhou, S. Kihyuk, and H. Lee, "Online Incremental Feature Learning with Denoising Autoencoders," *International Conference* on Artificial Intelligence and Statistics (AISTATS), vol. 22, pp. 1453– 1461, 2012.
- [65] D. Sahoo, Q. Pham, J. Lu, and S. C. Hoi, "Online deep learning: Learning deep neural networks on the fly," *International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 2660–2666, 2018.
- [66] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," International Conference on Machine Learning (ICML), 2009
- [67] A. Bendale and T. E. Boult, "Towards Open Set Deep Networks," Computer Vision and Pattern Recognition (CVPR), 2016.
- [68] F. Zenke, B. Poole, and S. Ganguli, "Continual Learning Through Synaptic Intelligence," *International Conference on Machine Learning* (ICML), vol. 70, pp. 3987–3995, 2017.
- [69] R. Aljundi, F. Babiloni, and M. Elhoseiny, "Memory Aware Synapses: Learning what (not) to forget," European Conference on Computer Vision (ECCV), 2018.
- [70] D. Park, S. Hong, B. Han, and K. M. Lee, "Continual Learning by Asymmetric Loss Approximation with Single-Side Overestimation," *International Conference on Computer Vision (ICCV)*, 2019.
- [71] S. W. Lee, J. H. Kim, J. Jun, J. W. Ha, and B. T. Zhang, "Overcoming catastrophic forgetting by incremental moment matching," *Neural Information Processing Systems (NeurIPS)*, pp. 4653–4663, 2017.
- [72] H. Ahn, S. Cha, D. Lee, and T. Moon, "Uncertainty-based Continual Learning with Adaptive Regularization," Neural Information Processing Systems (NeurIPS), 2019.
- [73] S. Ebrahimi, M. Elhoseiny, T. Darrell, and M. Rohrbach, "Uncertainty-guided Continual Learning with Bayesian Neural Networks," *International Conference on Learning Representations* (ICLR), 2020.
- [74] G. E. Hinton, O. Vinyals, and J. Dean, "Distilling the Knowledge in a Neural Network," NeurIPS Deep Learning Workshop, 2014.
- [75] A. Rannen, R. Aljundi, M. B. Blaschko, and T. Tuytelaars, "Encoder Based Lifelong Learning," *International Conference on Computer Vision (ICCV)*, 2017.
- [76] A. Gepperth and C. Karaoguz, "A Bio-Inspired Incremental Learning Architecture for Applied Perceptual Problems," Cognitive Computation, vol. 8, no. 5, pp. 924–934, 2016.
   [77] D. Isele and A. Cosgun, "Selective experience replay for lifelong
- [77] D. Isele and A. Cosgun, "Selective experience replay for lifelong learning," AAAI Conference on Artificial Intelligence (AAAI), 2018.
- [78] D. Rolnick, A. Ahuja, J. Schwarz, T. P. Lillicrap, and G. Wayne, "Experience Replay for Continual Learning," Neural Information Processing Systems (NeurIPS), 2018.
- [79] Y. Wu, Y. Chen, L. Wang, Y. Ye, Z. Liu, Y. Guo, and Y. Fu, "Large Scale Incremental Learning," Computer Vision and Pattern Recognition (CVPR), 2019.
- [80] A. Robins, "Catastrophic Forgetting, Rehearsal and Pseudorehearsal," Connection Science, vol. 7, no. 2, pp. 123–146, 1995.
- [81] R. M. French, "Pseudo-recurrent Connectionist Networks: An Approach to the 'Sensitivity-Stability' Dilemma," Connection Science, vol. 9, no. 4, pp. 353–380, 1997.
- [82] B. Ans, P. Rousset, C. E. P. Gi, and F. Les, "Avoiding catastrophic forgetting by coupling two reverberating neural networks," *Life Sciences*, pp. 989–997, 1997.
- [83] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative Adversarial Nets," Neural Information Processing Systems (NeurIPS), 2014.
- [84] G. M. van de Ven and A. S. Tolias, "Generative replay with feedback connections as a general strategy for continual learning," arXiv preprint arXiv:1809.10635, 2018.
- [85] R. M. French, "Semi-distributed Representations and Catastrophic Forgetting in Connectionist Networks," Connection Science, vol. 4, no. 3-4, pp. 365–377, 1992.

- [86] C. Fernando, D. Banarse, C. Blundell, Y. Zwols, D. Ha, A. A. Rusu, A. Pritzel, and D. Wierstra, "PathNet: Evolution Channels Gradient Descent in Super Neural Networks," arXiv preprint arXiv:1701.08734, 2017.
- [87] A. Mallya, D. Davis, and S. Lazebnik, "Piggyback: Adapting a Single Network to Multiple Tasks by Learning to Mask Weights," in European Conference on Computer Vision (ECCV), 2018.
- [88] J. Serra, D. Suris, M. Mirón, and A. Karatzoglou, "Overcoming Catastrophic forgetting with hard attention to the task," *Interna*tional Conference on Machine Learning (ICML), 2018.
- [89] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirk-patrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell, "Progressive Neural Networks," arXiv preprint arXiv:1606.04671, 2016.
- [90] R. Aljundi, P. Chakravarty, and T. Tuytelaars, "Expert gate: Lifelong learning with a network of experts," Computer Vision and Pattern Recognition (CVPR), 2017.
- [91] J. B. Aimone, Y. Li, S. W. Lee, G. D. Clemenson, W. Deng, and F. H. Gage, "Regulation and function of adult neurogenesis: from genes to cognition," *Physiological reviews*, vol. 94, no. 4, pp. 991–1026, 2014.
- [92] K. C. Vadodaria and S. Jessberger, "Functional neurogenesis in the adult hippocampus: Then and now," Frontiers in Neuroscience, vol. 8, no. 8 MAR, pp. 1–3, 2014.
- [93] C. G. Gross, "Neurogenesis in the adult brain: Death of a dogma," *Nature Reviews Neuroscience*, vol. 1, no. 1, pp. 67–73, 2000.
- [94] T. Ash, "Dynamic Node Creation in Backpropagation Networks," Connection Science, vol. 1, no. 4, pp. 365–375, 1989.
- [95] T. J. Draelos, N. E. Miner, C. C. Lamb, J. A. Cox, C. M. Vineyard, K. D. Carlson, W. M. Severa, C. D. James, and J. B. Aimone, "Neurogenesis deep learning: Extending deep networks to accommodate new classes," *International Joint Conference on Neural Networks (IJCNN)*, pp. 526–533, 2017.
- [96] J. Yoon, E. Yang, J. Lee, and S. J. Hwang, "Lifelong Learning with Dynamically Expandable Networks," *International Conference on Learning Representations (ICLR)*, 2018.
- [97] J. Xu and Z. Zhu, "Reinforced continual learning," Neural Information Processing Systems (NeurIPS), 2018.
- [98] X. Li, Y. Zhou, T. Wu, R. Socher, and C. Xiong, "Learn to Grow: A Continual Structure Learning Framework for Overcoming Catastrophic Forgetting," *International Conference on Machine Learning* (ICML), 2019.
- [99] M. Welling, "Herding dynamical weights to learn," *International Conference On Machine Learning (CML)*, pp. 1121–1128, 2009.
- [100] S. Sodhani, S. Chandar, and Y. Bengio, "Towards Training Recurrent Neural Networks for Lifelong Learning," Neural Computation, 2019
- [101] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," Neural Computation, vol. 9, no. 8, pp. 1735–1780, 1997.
- [102] T. Chen, I. Goodfellow, and J. Shlens, "Net2Net: Accelerating learning via knowledge transfer," *International Conference on Learning Representations (ICLR)*, 2016.
- [103] T. Zhang and F. J. Oles, "A Probability Analysis on the Value of Unlabelled Data for Classification Problems," *International Conference on Machine Learning (ICML)*, 2000.
- [104] C. É. Shannon, "A Mathematical Theory of Communication," Bell System Technical Journal, vol. 27, no. 4, pp. 623–656, 1948.
- [105] D. J. C. MacKay, "Information-Based Objective Functions for Active Data Selection," *Neural Computation*, vol. 4, no. 4, pp. 590–604, 1992.
- [106] L. E. Atlas, D. Cohn, and R. Ladner, "Training connectionist networks with queries and selective sampling," Neural Information Processing Systems (NeurIPS), 1990.
- [107] D. D. Lewis and W. A. Gale, "A sequential algorithm for training text classifiers," International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 1994, 1994.
- [108] H. S. Seung, M. Opper, and H. Sompolinsky, "Query by committee," Proceedings of the Fifth Annual ACM Workshop on Computational Learning Theory, pp. 287–294, 1992.
- [109] A. K. McCallum and K. Nigam, "Employing EM and Pool-Based Active Learning for Text Classification," *International Conference* on Machine Learning (ICML), 1998.
- [110] Y. Gal, R. Islam, and Z. Ghahramani, "Deep Bayesian Active Learning with Image Data," *International Conference on Machine Learning (ICML)*, 2017.
- [111] T. Tran, T. T. Do, I. Reid, and G. Carneiro, "Bayesian generative active deep learning," *International Conference on Machine Learning* (ICML), 2019.

- [112] D. A. Cohn, Z. Ghahramani, and M. I. Jordan, "Active learning with statistical models," *Journal of Artificial Intelligence Research*, vol. 4, pp. 129–145, 1996.
- [113] T. M. Mitchell, "Generalization as search," Artificial Intelligence, vol. 18, no. 2, pp. 203–226, 1982.
- [114] S. Tong and D. Koller, "Support Vector Machine Active Learning with Applications to Text Classification," Journal of Machine Learning Research (JMRL), 2001.
- [115] A. J. Joshi, F. Porikli, and N. Papanikolopoulos, "Multi-class active learning for image classification," Computer Vision and Pattern Recognition (CVPR), 2009.
- [116] S. Dasgupta, "Analysis of a greedy active learning strategy," Neural Information Processing Systems (NeurIPS), 2005.
- [117] K. Konyushkova, S. Raphael, and P. Fua, "Learning active learning from data," Neural Information Processing Systems (NeurIPS), 2017.
- [118] H. T. Nguyen and A. Smeulders, "Active learning using preclustering," *International Conference on Machine Learning (ICML)*, 2004.
- [119] I. W. Tsang, J. T. Kwok, and P. M. Cheung, "Core vector machines: Fast SVM training on very large data sets," *Journal of Machine Learning Research*, vol. 6, pp. 363–392, 2005.
- [120] O. Sener and S. Savarese, "Active learning for convolutional neural networks: A core-set approach," *International Conference on Learning Representations (ICLR)*, 2018.
- [121] C. Shui, F. Zhou, C. Gagné, and B. Wang, "Deep Active Learning: Unified and Principled Method for Query and Training," AISTATS, 2020.
- [122] S. Sinha, S. Ebrahimi, and T. Darrell, "Variational adversarial active learning," *International Conference on Computer Vision (ICCV)*, 2019.
- [123] D. P. Kingma and J. L. Ba, "Adam: a Method for Stochastic Optimization," *International Conference on Learning Representations* (ICLR), 2015.
- [124] J.-J. Zhu and J. Bento, "Generative Adversarial Active Learning," arXiv preprint arXiv: 1702.07956, 2017.
- [125] D. Mahapatra, B. Bozorgtabar, J. P. Thiran, and M. Reyes, "Efficient active learning for image classification and segmentation using a sample selection and conditional generative adversarial network," Medical Image Computing and Computer Assisted Intervention (MICCAI), 2018.
- [126] C. Mayer and R. Timofte, "Adversarial sampling for active learning," Winter Conference on Applications of Computer Vision (WACV), 2020.
- [127] W. J. Scheirer, A. Rocha, A. Sapkota, and T. E. Boult, "Towards Open Set Recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 35, no. 7, pp. 1757–1772, 2013.
- [128] A. Bendale and T. Boult, "Towards Open World Recognition," Computer Vision and Pattern Recognition (CVPR), 2015.
- [129] R. Naylor, "Known knowns, known unknowns and unknown unknowns: A 2010 update on carotid artery disease," Surgeon, vol. 8, no. 2, pp. 79–86, 2010.
- [130] D. Rumsfeld, "U.S. Department of Defense news briefing addressing unknown unknowns," 2002. [Online]. Available: https://archive.defense.gov/Transcripts/Transcript. aspx?TranscriptID=2636
- [131] W. J. Scheirer, L. P. Jain, and T. E. Boult, "Probability Models For Open Set Recognition," IEEE Transactions on Pattern Analysis and Machine Intelligence, 2014.
- [132] K. Lee, K. Lee, H. Lee, and J. Shin, "A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks," Neural Information Processing Systems (NeurIPS), 2018.
- [133] R. Yoshihashi, W. Shao, R. Kawakami, S. You, M. Iida, and T. Naemura, "Classification-Reconstruction Learning for Open-Set Recognition," Computer Vision and Pattern Recognition (CVPR), 2019.
- [134] R. Munro, Human-in-the-Loop Machine Learning. Manning Publications, Manning Early Access Program, 2020.
- [135] C. M. Bishop, "Novelty detection and neural network validation," *IEE Proceedings: Vision, Image and Signal Processing*, vol. 141, no. 4, pp. 217–222, 1994.
- [136] J. Weston, R. Collobert, F. Sinz, L. Bottou, and V. Vapnik, "Inference with the Universum," *International Conference on Machine Learning* (ICML), 2006.
- [137] M. Cimpoi, S. Maji, and A. Vedaldi, "Deep convolutional filter banks for texture recognition and segmentation," in *Computer Vision and Pattern Recognition (CVPR)*, 2015.

- [138] S. Bell, P. Upchurch, N. Snavely, and K. Bala, "Material recognition in the wild with the Materials in Context Database," in Computer Vision and Pattern Recognition (CVPR), 2015.
- [139] K. Lee, H. Lee, K. Lee, and J. Shin, "Training Confidence-Calibrated Classifiers for Detecting Out-of-Distribution Samples," International Conference on Learning Representations (ICLR), 2018.
- [140] Q. Yu and K. Aizawa, "Unsupervised Out-of-Distribution Detection by Maximum Classifier Discrepancy," *International Conference on Computer Vision (ICCV)*, 2019.
- [141] A. R. Dhamija, M. Günther, and T. E. Boult, "Reducing Network Agnostophobia," Neural Information Processing Systems (NeurIPS), 2018.
- [142] Q. Feng, G. Kang, H. Fan, and Y. Yang, "Attract or Distract: Exploit the Margin of Open Set," *International Conference on Computer Vision (ICCV)*, 2019.
- [143] F. Li and H. Wechsler, "Open set face recognition using transduction," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 11, pp. 1686–1697, 2005.
- [144] P. L. Bartlett and M. H. Wegkamp, "Classification with a reject option using a hinge loss," *Journal of Machine Learning Research*, vol. 9, pp. 1823–1840, 2008.
- [145] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," *International Conference on Learning Representations (ICLR)*, 2017.
- [146] A. Kendall, V. Badrinarayanan, and R. Cipolla, "Bayesian segnet: Model uncertainty in deep convolutional encoder-decoder architectures for scene understanding," *British Machine Vision Conference* (BMVC), 2017.
- [147] D. Miller, L. Nicholson, F. Dayoub, and N. Sunderhauf, "Dropout Sampling for Robust Object Detection in Open-Set Conditions," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 3243–3249.
- [148] S. Liang, Y. Li, and R. Srikant, "Enhancing the Reliability of Out-ofdistribution Image Detection in Neural Networks," *International Conference on Learning Representations (ICLR)*, 2018.
- [149] K. Lis, K. Nakka, P. Fua, and M. Salzmann, "Detecting the Unexpected via Image Resynthesis," *International Conference on Computer Vision (ICCV)*, 2019.
- [150] P. Perera, R. Nallapati, and B. Xiang, "OCGAN: One-class Novelty Detection Using GANs with Constrained Latent Representations," Computer Vision and Pattern Recognition (CVPR), 2019.
- [151] Y. Li, J. Bradshaw, and Y. Sharma, "Are Generative Classifiers More Robust to Adversarial Attacks?" International Conference on Machine Learning (ICML), 2019.
- [152] E. Nalisnick, A. Matsukawa, Y. W. Teh, D. Gorur, and B. Lakshminarayanan, "Do Deep Generative Models Know What They Don't Know?" *International Conference on Learning Representations (ICLR)*, 2019.
- [153] Y. Ovadia, E. Fertig, J. Ren, Z. Nado, D. Sculley, S. Nowozin, J. V. Dillon, B. Lakshminarayanan, and J. Snoek, "Can You Trust Your Model's Uncertainty? Evaluating Predictive Uncertainty Under Dataset Shift," Neural Information Processing Systems (NeurIPS), 2019.
- [154] A. Chaudhry, P. K. Dokania, T. Ajanthan, and P. H. Torr, "Riemannian Walk for Incremental Learning: Understanding Forgetting and Intransigence," European Conference on Computer Vision (ECCV), 2018.
- [155] A. Chaudhry, R. Marc'Aurelio, M. Rohrbach, and M. Elhoseiny, "Efficient lifelong learning with A-GEM," International Conference on Learning Representations (ICLR), 2019.
- [156] P. Oza and V. M. Patel, "C2AE: Class Conditioned Auto-Encoder for Open-set Recognition," Computer Vision and Pattern Recognition (CVPR), 2019.
- [157] K. Xiao, L. Engstrom, A. Ilyas, and A. Madry, "Noise or Signal: The Role of Image Backgrounds in Object Recognition," ArXiv preprint arXiv: 2006.09994, 2020.
- [158] R. Geirhos, C. Michaelis, F. A. Wichmann, P. Rubisch, M. Bethge, and W. Brendel, "Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness," *International Conference on Learning Representations (ICLR)*, 2019.
- [159] S. Lapuschkin, S. Wäldchen, A. Binder, G. Montavon, W. Samek, and K. R. Müller, "Unmasking Clever Hans predictors and assessing what machines really learn," *Nature Communications*, vol. 10, no. 1, 2019.
- [160] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, "Adversarial Examples are not Bugs, they are Features," *Neural Information Processing Systems (NeurIPS)*, 2019.

- [161] G. Hacohen, L. Choshen, and D. Weinshall, "Let's Agree to Agree: Neural Networks Share Classification Order on Real Datasets," International Conference on Learning Representations (ICLR), 2020.
- [162] K. Javed and F. Shafait, "Revisiting Distillation and Incremental Classifier Learning," Asian Conference on Computer Vision (ACCV), 2018.
- [163] G. Fei, S. Wang, and B. Liu, "Learning cumulatively to become more knowledgeable," Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 1565– 1574, 2016.
- [164] D. P. Kingma and M. Welling, "Auto-Encoding Variational Bayes," International Conference on Learning Representations (ICLR), 2013.
- [165] M. D. Hoffman and M. J. Johnson, "ELBO surgery: yet another way to carve up the variational evidence lower bound," Neural Information Processing Systems (NeurIPS), Advances in Approximate Bayesian Inference Workshop, 2016.
- [166] C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins, and A. Lerchner, "Understanding disentangling in beta-VAE," Neural Information Processing Systems (NeurIPS), Workshop on Learning Disentangled Representations, 2017.
- [167] E. Mathieu, T. Rainforth, N. Siddharth, and Y. W. Teh, "Disentangling disentanglement in variational autoencoders," *International Conference on Machine Learning (ICML)*, pp. 7744–7754, 2019.
- [168] S. Zagoruyko and N. Komodakis, "Wide Residual Networks," British Machine Vision Conference (BMVC), 2016.
- [169] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," Computer Vision and Pattern Recognition (CVPR), 2016.
- [170] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," International Conference on Computer Vision (ICCV), 2015.
- [171] S. Ioffe and C. Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift," International Conference on Machine Learning (ICML), 2015.
- [172] O. Bachem, M. Lucic, and A. Krause, "Coresets for Nonparametric Estimation - the Case of DP-Means," *International Conference on Machine Learning (ICML)*, vol. 37, pp. 209–217, 2015.
- [173] T. F. Gonzalez, "Clustering to minimize the maximum intercluster distance," Theoretical Computer Science, vol. 38, no. C, pp. 293–306, 1985
- [174] S. Becker, M. Ackermann, S. Lapuschkin, K.-R. Müller, and W. Samek, "Interpreting and Explaining Deep Neural Networks for Classification of Audio Signals," arXiv preprint arXiv: 1807.03418, 2018.