
Towards Robust Evaluations of Continual Learning

Sebastian Farquhar¹ Yarin Gal¹

Abstract

Experiments used in current continual learning research do not faithfully assess fundamental challenges of learning continually. Instead of assessing performance on challenging and representative experiment designs, recent research has focused on increased dataset difficulty, while still using flawed experiment set-ups. We examine standard evaluations and show why these evaluations make some continual learning approaches look better than they are. We introduce desiderata for continual learning evaluations and explain why their absence creates misleading comparisons. Based on our desiderata we then propose new experiment designs which we demonstrate with various continual learning approaches and datasets. Our analysis calls for a reprioritization of research effort by the community.

1. Introduction

Certain applications require *continual* learning—splitting training into a series of tasks and discarding the data after training each task. These settings vary tremendously. If data relate to individuals it may be unethical, illegal, or imprudent to retain old datasets: for example, a hospital might want to delete old patient data. In other applications, real-time systems could face distributional shifts with an underlying distribution which changes faster than the time it would take to retrain a new model with all the data. For example, a Mars rover might cross regions of different terrain that require adaptation without forgetting. Neural networks trained on such series of tasks tend to forget earlier tasks (often referred to as *catastrophic forgetting*).

Recent works have shown promising advances towards continual learning by adding regularization terms to the loss function which preserve important parameters (Kirkpatrick et al., 2017; Zenke et al., 2017; Nguyen et al., 2018; Chaudhry et al., 2018; Ritter et al., 2018). We call these

approaches ‘prior-focused’ because they can be interpreted as treating intermediate models as priors when learning updated weights. Although a useful starting point, experimental evaluations in these recent works have major blind-spots which mask weak-points of current approaches. The design of the evaluations does not fully reflect the core motivations for continual learning, regardless of the dataset used, including evaluations based on MNIST (LeCun et al., 1998), notMNIST (Bulatov, 2011), FashionMNIST (Xiao et al., 2017), CIFAR10 (Krizhevsky, 2009) and others. But evaluations which obscure the shortcomings of suggested continual learning solutions impede developments in the field, since researchers are not made aware of limitations of past research.

In this paper, we explore desiderata for continual learning evaluations based on real-world uses. Applications of continual learning are very diverse, and authors consider a wide range of settings. As a result, rather than propose a single benchmark which risks being overly narrow, we propose fundamental desiderata for empirical evaluations in continual learning to make them more representative of the challenges that motivate the field. We then demonstrate that experiments which neglect our core desiderata can be misleading—a continual learning system that performs well when a subset of the desiderata are in place might fail entirely when a core desideratum is imposed. No prior-focused continual learning system has so far been shown to succeed when all five desiderata are applied. We demonstrate that several leading prior-focused approaches stop working on a more robust evaluation. We therefore argue that existing evaluations are biased, because prior-focused approaches appear to succeed, but have major blind-spots. Lastly, we introduce new evaluations based on our list of desiderata which offer richer challenges for continual learning systems.

Specifically, our four main contributions are as follows:

1. We propose fundamental desiderata for future evaluations, which can be applied regardless of dataset.
2. We analyse the shortcomings of existing widely used evaluations in continual learning.
3. We show empirically that existing evaluations are biased towards prior-focused approaches.

¹OATML Research Group, Department of Computer Science, University of Oxford, United Kingdom. Correspondence to: Sebastian Farquhar <sebastian.farquhar@cs.ox.ac.uk>.

4. We propose new experimental designs which mitigate the issues of existing ones.

2. Continual Learning

We start by formalising the definition of *continual learning*. In a typical supervised learning setting, we aim to learn parameters \mathbf{w} using an independently and identically distributed (i.i.d.) labelled training dataset $\mathcal{D} \equiv \{(\mathbf{x}^{(i)}, y^{(i)})\}$ to accurately predict $p(y^*|\mathbf{w}, \mathbf{x}^*)$ for an unseen (\mathbf{x}^*, y^*) pair.

In the continual learning setting, members of \mathcal{D} are not i.i.d. Instead they may be split into disjoint subsets $\mathcal{D}_t \equiv \{(\mathbf{x}_t^{(i)}, y_t^{(i)})\}$. These sets are assumed to be drawn from T distinct i.i.d. distributions each of which represents a *task*. The challenge of continual learning is to learn a single model which is able to predict well on data from any task, despite training on each task in sequence without deliberately revisiting previous tasks. For some applications, stronger or weaker assumptions are made. For example, Nguyen et al. (2018) concern themselves with fast adaptation, so they allow themselves to retain a small coreset of data from old tasks. Similarly, Kirkpatrick et al. (2017); Schwarz et al. (2018b) consider a simulation in which they allow themselves to revisit each task multiple times. In this paper, we will frequently consider a model after training on some, but not necessarily all T tasks. The t 'th model is the state of the model after training on datasets $\{D_1 : D_t\}$.

3. Proposed Desiderata for Continual Learning Evaluations

The setting and motivation for continual learning can vary enormously, from image segmentation to medical diagnostics to control. Specialized continual learning datasets like Core50 (Lomonaco and Maltoni, 2017) or environments such as that in Schwarz et al. (2018a) are therefore useful for subfields but cannot be seen as an overall benchmark for continual learning. We therefore concentrate instead on developing set of sufficient desiderata for future evaluations of continual learning. We propose the following fundamental principles which ought either to be respected or to be explicitly mentioned as absent. In §5 we offer a critical analysis of existing experiments which suggests why core desiderata are important. In §6.3, we present a set of experiments showing how neglecting these core aspects can make a continual learning system look more successful than it is.

Use cases of continual learning vary tremendously. One example is a hospital developing a system for automated disease diagnosis using patient sub-population A which is then refined in a different hospital with patient sub-populations B (with sub-population A data not allowed to leave the first

hospital). A second example is a wind-turbine safety system which predicts wind velocity an hour ahead to decide whether to turn off the turbine in case of strong winds, with wind dynamics changing between the seasons of the year. A third example is an ad serving system that has to update with huge amounts of streaming data, without feasibly being able to keep all data around. Yet another use case is a Mars rover updating its behaviour as it encounters new types of terrain—perhaps much looser or coarser soil than it has seen before—but which must remain able to manoeuvre on old terrain. All use cases share common core desiderata, which we demonstrate next using the Mars rover example for concreteness. Note that this example is for illustrative purposes only—we are *not* proposing that this example should be used for evaluations. Our desiderata arise just as readily from other settings.

- A: Cross-task resemblances** Input data from later tasks must resemble old tasks enough that they at least sometimes result in confident predictions of old classes, early in training. The widely used Permuted MNIST (see §4.2.1) which violates this would correspond to every input sensor in our Mars rover being randomly rewired—unlikely to be representative of real cases.
- B: Shared output head** If each task is given a different output vector during training and testing, it must be explicitly and prominently mentioned. This has a large effect on the difficulty of the challenge. In our example, the rover is using the same outputs to control its motors in all terrain.
- C: No test-time assumed task labels** Related to desideratum B, if we knew in advance what the tasks were and had a way to distinguish them, one could have a different model for each task and switch between them. In our example, if the rover has a good way to distinguish rock from dust and can learn a separate policy for each it might not need continual learning.
- D: No unconstrained retraining on old tasks** Many motivations of continual learning preclude being able to retrain on task 1 after having trained on task T . In our example, the rover needs to actually move into terrain it has left behind before it can train there. In other examples, even retaining small amounts of data from old tasks might violate privacy laws
- E: More than two tasks** The more tasks a continual learning system can handle the better. Our rover, for example, may face a potentially open-ended and large number of tasks. Two-task transfer (§4.2.3) is therefore interesting, but succeeding at two-task transfer does not guarantee good performance on more tasks.

In certain applications we might deliberately ignore some of the desiderata. For example, Nguyen et al. (2018) neglect desiderata B and C because they consider fast adaptation to new tasks which have labels during both training and test-

ing. In a setting with easy access to a simulation of all tasks, retraining might be permissible, setting aside desideratum D (as in Schwarz et al. (2018b)). Desideratum E is mostly set aside for convenience—there might be interesting intermediate research which is nevertheless useful progress for the field. Desideratum A is mostly neglected for historical reasons discussed in §5.1. While the five listed desiderata will not always be applicable, it remains a surprising fact that there is **no prior-focused approach that has been shown to perform well in an evaluation with all of the five core desiderata**.

The core desiderata are a starting point which can lead to more difficult challenges, some of which are beginning to be addressed by recent work, including:

Unclear task demarcation Task boundaries are generally assumed knowledge during training.

Continuous tasks Current practice usually assumes that tasks are discretely different rather than varying continuously.

Overlapping tasks Current practice tends to have tasks where each task has its own classes which are disjoint from the other tasks. Lomonaco and Maltoni (2017) discuss this point.

Long task sequences Continual learning is most interesting over extremely long sequences of tasks.

Time/compute/memory constraints It is widely acknowledged that continual learning solutions that update quickly or use small/fixed memory are more useful than ones that have large computational overheads.

Strict privacy guarantees For settings motivated by privacy, tracking the differential privacy guarantees between tasks is important.

Each of these further desiderata offers a way to make a continual learning evaluation more representative of the needs of applications *without making the dataset itself harder*. In §6.4 we suggest two experiments that address the most interesting of these. In §6.4.1 we show how model uncertainty can be used to detect task changes, which is essential in many practical settings and allows streaming settings to be converted into continual learning ones. In §6.4.2 we show how training time and accuracy must be traded off against each other as design choices, rather than just reporting the time taken by the most accurate system. Memory can be treated like time, here, and the other desiderata are straightforward to construct from existing datasets by tweaking the dataloader.

We note that a very real risk for the field is to *embrace more complex datasets before establishing robust experimental design*. Until we have continual learning evaluations that test the core challenges of continual learning, moving from MNIST to more complex datasets does little to measure progress in continual learning itself.

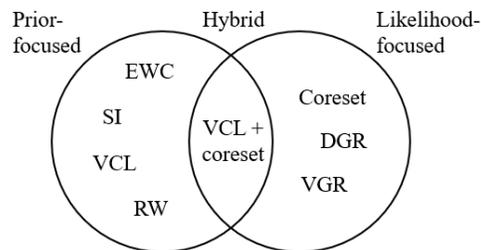


Figure 1. We contrast prior-focused with likelihood-focused continual learning. By comparing these, and hybrid forms, we can see which evaluations pose a bigger challenge to different approaches.

4. Existing Work

In this paper, we critically analyse common evaluations used by the majority of recent papers. In this section, we review existing work focusing first on the methods employed (§4.1) and second on the evaluations used (§4.2). Throughout, we provide a framing which anticipates our critical analysis.

4.1. Methods

Approaches to continual learning and catastrophic forgetting form three main families. First, prior-focused approaches use regularization to create ‘elastic’ parameters. These include: Elastic Weight Consolidation (EWC) (Kirkpatrick et al., 2017; Huszar, 2018), Synaptic Intelligence (SI) (Zenke et al., 2017), Variational Continual Learning (VCL) (Nguyen et al., 2018), Kronecker-factored approximated Laplace approximation (Ritter et al., 2018), and Riemannian Walk (RW) (Chaudhry et al., 2018). In fact, all of these can be seen as setting a distribution over the parameters of the $t - 1$ ’th model as the prior when training on the t ’th dataset.

VCL makes this most explicit. VCL is a variational inference (VI) method that uses Bayesian neural networks (MacKay, 1992; Neal, 1995; Blundell et al., 2015). VCL sets the posterior at the end of training a task to be the prior when beginning training for the next task. Schematically, as the authors of these papers note, all the approximately Bayesian approaches use a similar loss function:

$$\mathcal{L}_t = \sum_i \log(p(y_t^{(i)} | \mathbf{w}_t, \mathbf{x}_t^{(i)})) - \frac{c}{2} (\mathbf{w}_t - \mathbf{w}_{t-1})^T \Sigma_{t-1}^{-1} (\mathbf{w}_t - \mathbf{w}_{t-1}) \quad (1)$$

where Σ is the covariance of the prior, which is estimated differently under each approach. The first term reflects model log-likelihood on the t ’th task. The second reflects distance from the model prior.

Prior-focused approaches rely on the prior term to capture everything learned on previous tasks and only use the newest data to estimate the likelihood term. However, in practice the prior term is only approximated. Farquhar and Gal (2018b) contrast prior-focused approaches with a second family,

likelihood-focused approaches, that attempt to estimate the likelihood of the current model on past tasks (see figure 1). Examples of this include Deep Generative Replay (DGR) by Shin et al. (2017), which uses generative models to store summaries of old tasks, or Variational Generative Replay (VGR) a variational and simplified alternative presented in Farquhar and Gal (2018b). Simply storing a small random sample or coreset is an alternative approach in this family. Many members of this family are called rehearsal strategies after ‘pseudorehearsal’ introduced by Robins (1995), or called dual-memory strategies following a neuro-scientific motivation. However, some dual-memory strategies, e.g., GEM (Lopez-paz and Ranzato, 2017) or model adaptation (Hu et al., 2019), are not likelihood-focused.

Some authors have observed better results for prior-focused methods with the addition of some likelihood estimation, although this has not been rigorously explored. For example, VCL is also presented using coresets—fine-tuning on a small withheld sample just before testing. Similarly, RW presents results both with and without training on subsamples of previous datasets.

The last major family uses dynamic architectures. These change the structure of the networks in significant ways to incorporate learning from each task separately (e.g., Rusu et al. (2016); Li and Hoiem (2017)). These perform well, but in most cases their need for new models for each task creates growing compute/memory costs that do not seem to truly solve the continual learning problem. A more thorough survey of these as well as prior- and likelihood-focused approaches can be found in Appendix A. We focus our examination on prior-focused and their counterpart likelihood-focused methods, because they seem to represent the most exciting recent advances in continual learning, but the experiments used to support them may not be fairly evaluating their performance.

4.2. Experimental Evaluations

In this section we give an in-depth review of key experimental evaluations as they are used in *current research*. We illustrate our critique using the most commonly performed experiments using variations of the MNIST dataset: *Permuted MNIST* and *Split MNIST*, most commonly in *multi-headed* form. Although we emphasise MNIST for explanatory purposes, our critique in §5 applies to many other evaluations which share similar design choices.

4.2.1. PERMUTED MNIST

The Permuted MNIST experiment was introduced in Goodfellow et al. (2013a). In their experiment, a model is trained on MNIST as \mathcal{D}_1 . Each later \mathcal{D}_t for $1 < t \leq 10$ is constructed from the MNIST data but with the pixels of each digit randomly permuted. A fresh permutation is drawn

for each task and applied to all images in the same way for that task. After training on each dataset, one evaluates the model on each of the previous datasets as well as the current. Goodfellow et al. (2013a) used this experiment to investigate feature extraction, but it has since become a mainstay for continual learning evaluation (Zenke et al., 2017; Shin et al., 2017; Kirkpatrick et al., 2017; Lee et al., 2017; Lopez-paz and Ranzato, 2017; Nguyen et al., 2018; Ritter et al., 2018; Hu et al., 2019; Chaudhry et al., 2019).

4.2.2. SPLIT MNIST

The Split MNIST experiment was introduced by Zenke et al. (2017) in a *multi-headed* form and used by other authors including Shin et al. (2017); Nguyen et al. (2018); Chaudhry et al. (2018) (Ritter et al. (2018) use a two-task variant). The experiment constructs a series of five related tasks. The first task is to distinguish the digits (0, 1), then (2, 3) etc. Most papers use a multi-headed variant the model prediction is constrained to be only from the two classes represented in each task. For example, when evaluating the performance on the first task, the model only needs to predict probabilities for zero versus one. In some cases, multi-heading is taken even further and training is only done on the head governing the specific classes included in the task (Zenke et al., 2017; Nguyen et al., 2018; Ritter et al., 2018). A *single-headed* version does not limit predictions during either training or testing. As Chaudhry et al. (2018) note, the multi-headed variant is much easier to solve. Multi-heading is often used in similar non-MNIST evaluations, for example, in Zenke et al. (2017); Nguyen et al. (2018); Ritter et al. (2018). Chaudhry et al. (2018) use both a single- and multi-headed version of Split MNIST. Hu et al. (2019) use a single-headed set-up but their method effectively learns a different head for each task.

4.2.3. TWO-TASK TRANSFER

In some works, a two-task transfer learning evaluation is used for continual learning. For example, Shin et al. (2017) train first on MNIST and then on SVHN (Netzer et al., 2011) (or vice versa) in order to see whether their algorithm preserves performance on the first task after training on the second. Jung et al. (2016), Li and Hoiem (2017), and Pfulb and Gepperth (2018) perform similar experiments.

4.2.4. DESIDERATA, BENCHMARKING AND METRICS

Indicating the importance of robust evaluations, several authors have recently included some continual learning desiderata alongside their work. For example, Schwarz et al. (2018b) list desiderata for continual learning systems which includes being able to handle long sequences of tasks without knowing task labels and ideally no clear task demarcation, but do not introduce experimental set-ups which could

enforce all these desiderata, such as the ones we propose in §6.4. (Chaudhry et al., 2019) call for time and memory constraints as well as careful use of cross-validation. Pfulb and Gepperth (2018) call for testing on many datasets, care with cross-validation, and storage constraints. Farquhar and Gal (2018a) discuss the importance of tracking differential privacy guarantees in continual learning settings motivated by privacy.

Several authors have commented on the importance of systematic benchmarks for continual learning, including Lomonaco and Maltoni (2017), focused on object recognition, and Schwarz et al. (2018a), focused on control in a StarCraft environment. Others have worked towards richer metrics for measuring continual learning performance like forward/backwards transfer or learning curve area (Lopez-paz and Ranzato, 2017; Chaudhry et al., 2018; Schwarz et al., 2018b; Diaz-Rodriguez et al., 2018; Chaudhry et al., 2019). Both of these developments are orthogonal to the arguments of this paper — well-chosen metrics evaluated in a challenging environment are still unhelpful if the experiment is not designed to reflect the continual learning needs of diverse applications.

5. Critical Analysis of Existing Evaluations

Pixel permutation (e.g., Permuted MNIST), split classification (e.g., multi-headed Split MNIST), and two-task transfer experiments are commonly used to compare continual learning algorithms. Although they are a useful starting point, each of these makes continual learning easier for prior-focused approaches. All of our critiques of the way MNIST is used by Zenke et al. (2017) apply also to their use of the CIFAR100 dataset, the use of notMNIST, Fashion MNIST, SVHN and CIFAR10 in Ritter et al. (2018), and the use of notMNIST in (Nguyen et al., 2018). These critiques are about experimental design, not dataset choice.

5.1. Pixel Permutation

Permuted MNIST represents an unrealistic best case scenario for continual learning—although it satisfies the literal definition of continual learning. The positions of the pixels are fully randomized—which suited its original purpose. Goodfellow et al. (2013b) investigated whether neural networks have ‘high-level concepts’ that get re-mapped to pixel positions when the input space is permuted, and found evidence they did not. Now, however, the experiment has been repurposed for continual learning, where it is not suitable.

An image from each permuted dataset is practically unrecognizable given previous datasets. The actual world is almost never structured like this—new situations look confusingly similar to old ones—the sensor input in a Mars rover will never be permuted no matter what terrain one moves onto.

A model presented with a new task in Permuted MNIST will be uncertain, while in settings that are not deliberately randomized a model will tend to make confident but false predictions. This makes Permuted MNIST significantly different from real-world settings in a way that directly affects how new tasks are learned. In Appendix B.1, we offer an empirical investigation of this phenomenon and a hypothesis for the mechanism by which Permuted MNIST simplifies the continual learning challenge.

5.2. Split Classification

The multi-headed version of Split MNIST, which is most often used, is easier but less relevant to applications because it requires knowledge of the task and the classes represented in each task, as Chaudhry et al. (2018) point out. Usually, if that were possible, continual learning would be unnecessary, since one could use a separate classifier for each task, which generally would perform better.

Unfortunately, prior-focused continual learning systems tend to look much better in multi-headed evaluations than single-headed ones. Suppose some model has already been trained on the first four tasks of Split MNIST and is then tested on the digit ‘1’. In the single-headed variant, when shown a ‘1’ it may incorrectly predict the label is seven, which was seen more recently. In the multi-headed variant, we *knowingly assume* that the label comes from [0:1]. Because the model only needs to decide between 0 and 1 (and not even consider if the image is a 7), a multi-headed model could correctly predict the label is 1 even though the same approach will completely fail in a single-headed experiment setup. Most current prior-focused continual learning systems are good enough that models retain the ability to distinguish old classes, but only if they are able to know which task the example is from. A multi-headed evaluation can therefore make it seem as if an approach has solved a continual learning problem when it has not. *Single-headed Split MNIST*—a less-used evaluation—is the simplest experimental setting that meets all five of our core desiderata. We recommend it as a toy baseline, which is easily adapted to harder datasets. It satisfies our core desiderata: **A:** The tasks resemble each other—e.g., a seven can look like a one; **B:** Single-heading means all outputs are shared; **C:** A full prediction is made over all possible outputs each time; **D:** Each task is trained only once and only the model is carried forwards; and **E:** There are five tasks.

5.3. Other non-representative evaluations

Two-task transfer is not representative of realistic problems because continual learning use-cases might require a long series of tasks, not just two. An algorithm might have elements that perform well with just one previous task but fail with more. For example, the approximate Fisher information

matrix used in EWC is estimated using a Taylor expansion which is only locally accurate in one part of parameter-space. If, after many tasks, the model reaches a very different part of parameter-space, the estimated Fisher information for old tasks will be inaccurate. Moreover, in practice two-task transfer is just easier than true continual learning: Pfulb and Gepperth (2018) conclude that a simple fully connected network will show no catastrophic forgetting on a two-task Permuted MNIST but do not consider longer task sequences, while other authors have already shown that with more tasks there is indeed catastrophic forgetting (Kirkpatrick et al., 2017).

6. Empirical Analysis of Existing Evaluations

This section has three aims. First, we show that experiments which satisfy only a subset of the five main desiderata from §3 can have blind-spots that are important to applications that demand all five desiderata. We do this by showing that leading methods fail on an experiment that uses all five desiderata, while they look good on experiments that use just a subset. Second, we show that existing evaluations are biased towards prior-focused approaches. Third, we suggest two further experimental set-ups.

6.1. Experimental approach

We consider representatives of the prior-focused approach, likelihood-focused approach, and a hybrid approach. We use these to demonstrate the shortcomings of existing evaluations and to tease out which aspects of the architecture are challenged under different experimental design choices.

We use three variants of VCL because it has the clearest Bayesian interpretation. First, we use pure VCL without a coreset, exactly as described by Nguyen et al. (2018) (see Appendix A). Second, we use a small coreset of 40 datapoints extracted from each dataset (we use their k-center coreset approach, but this does not have a large effect). The second method is the same as pure VCL except that, at the end of training on each task, the model is trained on the coresets before testing, as described in their work. This reflects a hybrid approach. Third, we use a ‘coreset only’ approach as a likelihood-focused version. It is exactly like the second variant except that the prior used for variational inference is the initial prior each time—it is not updated after each task.¹

We use two further baselines to show that the effects we find are not an artefact of VCL specifically. EWC (Kirkpatrick et al., 2017) is used as a prior-focused alternative to VCL because it has been widely discussed in the literature.

¹This is not the coreset only algorithm used in Nguyen et al. (2018). Theirs is seeing *only* coresets—much less data—which is why it performs badly on even the first task.

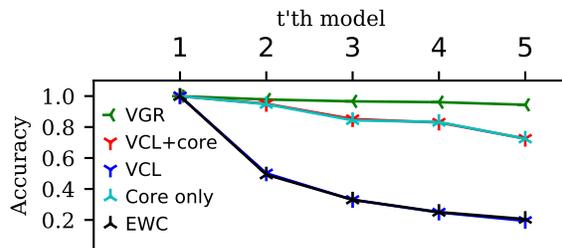


Figure 2. **Single-headed Split MNIST.** This experiment meets all core desiderata and shows big performance differences.

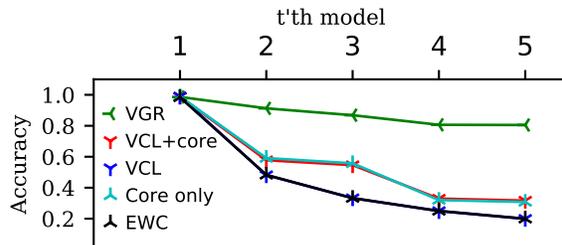


Figure 3. **Single-headed Split Fashion MNIST.** The harder dataset shows the prior-approximation starting to deteriorate, but does not change performance ranking.

VGR (Farquhar and Gal, 2018b) is used as a more powerful likelihood-focused alternative to coresets. VGR, VCL and its variants use a Bayesian neural network (BNN) and variational inference, while EWC does not, so performance comparisons between these architectures should be interpreted carefully. The fact that our observations for VCL and its variants are mirrored for EWC supports our hypothesis, based on our analysis, that the problems discussed in this section affect prior-focused approaches in general.

6.2. Experiments With All Five Desiderata Show Big Performance Differences

We use Single-headed Split MNIST, described in §5.2. This satisfies all main desiderata from §3. It reveals major differences in performance between the approaches (see figures 2 and 3). The performance of VCL with coreset appears to be entirely driven by the presence of the coresets. When coresets are removed, VCL alone completely forgets old tasks—its accuracy comes from correctly classifying the most recent task only. EWC performs exactly like VCL, suggesting that it is prior-focused approaches in general that are struggling. Meanwhile VGR, the more accurate likelihood-focused approach, is the only method that performs nearly perfectly. Using FashionMNIST rather than MNIST is harder and a worthwhile additional test, but does not reveal a radically different story. VCL performs much worse on FashionMNIST than VGR even though both use the same model. It seems the prior approximation struggles with data complexity more than the GAN.

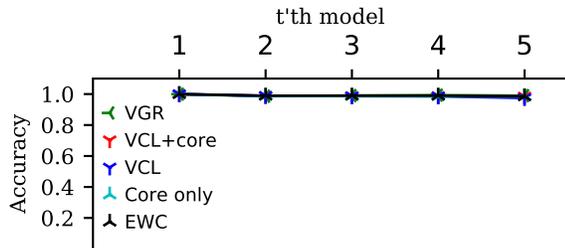


Figure 4. Multi-headed Split MNIST. All methods succeed.

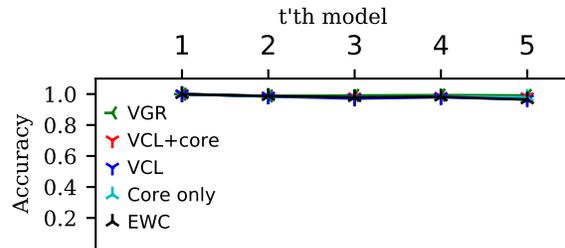


Figure 6. Test-time Knowledge Split MNIST. All succeed.

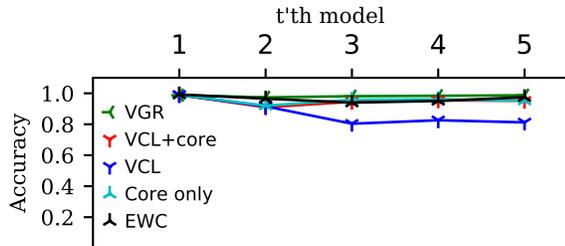


Figure 5. Multi-headed Split FashionMNIST. All perform similarly, so no clean differentiation. VCL performs slightly worse without coresets.

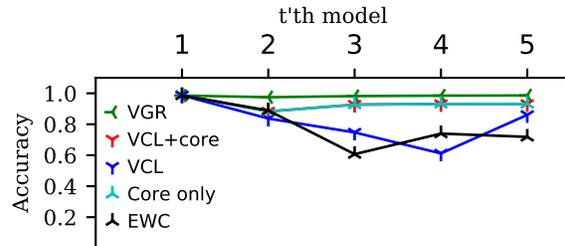


Figure 7. Test-time Knowledge Split FashionMNIST. The experiment is harder but still does not cleanly differentiate prior- and likelihood-focused methods.

6.3. Missing Desiderata Create Blindspots

Having shown that there are big differences in performance between our baselines on a well-constructed experiment that follows all desiderata, we now show that experiments that do not include all desiderata can fail to reveal these important performance differences. This means that these experimental evaluations can be misleading. The fact that they systematically hide shortcomings of prior-focused approaches suggests that the evaluations are inadvertently biasing research efforts in the field. We use MNIST and FashionMNIST to show that while using more complex datasets does improve evaluations, it does not make up for unrepresentative experimental design.

6.3.1. NO CROSS-TASK RESEMBLANCE

As we argue in §5.1, Permuted MNIST is unusual in that there is no resemblance between images in each task that might cause confident incorrect predictions. Previous authors have demonstrated the effectiveness of their methods on Permuted MNIST (Kirkpatrick et al., 2017; Nguyen et al., 2018). Virtually any continual learning method proposed has shown good performance on this task. This demonstrates the weakness of experiments that do not reflect desideratum A. Further experimental investigation of these phenomena is in Appendix B.1 with baseline results in Appendix B.2.

6.3.2. NO SHARED OUTPUT HEAD

Authors frequently use a multi-headed version of Split MNIST, following Zenke et al. (2017), which we describe in §4.2.2. On this evaluation, all approaches look similarly

good (see figures 4 and 5). This shows that multi-heading, and neglecting desideratum B, can create misleading experiments. Details of training and hyperparameters can be found in Appendix B.4.

6.3.3. TEST-TIME TASK KNOWLEDGE

Instead of entirely separate heads for each task, one might simply know which output classes to look at during testing. In this variant, which breaks desideratum C, everything is exactly the same as for the single-headed variant during training, but during testing all outputs except those for the correct class are masked with zeros before the most probable class is selected. Here, both on MNIST and FashionMNIST the evaluation does not differentiate the performance of prior- and likelihood-focused methods (see figures 6 and 7).

6.3.4. RETRAINING ON OLD TASKS

An experimental evaluation that is exactly like single-headed Split MNIST except that the tasks are looped through several times does still find that prior-focused methods perform badly (see Appendix B.6). Neglecting this desideratum, therefore, might not make evaluations biased, though it is certainly not in keeping with the core motivations of continual learning applications since it requires data to be kept indefinitely.

6.3.5. ONLY TWO TASKS

Many authors have tested their methods with only two tasks. This is an interesting challenge, but significantly simpler than a longer series. Showing that a method succeeds on a two-task transfer does not entail that it will work on a

longer series. For example, looking at the first two tasks in figure 2 would not differentiate VCL with coresets and VGR, where looking at the longer series shows a clear difference in performance.

6.4. Further Experimental Options

Having shown the importance of our five main desiderata in experimental evaluations, we further suggest experiments based on the extended desiderata listed in §3. The first tests the ability of the model to detect a change in task, which has not yet been demonstrated in continual learning. The second evaluates time against accuracy (or some other metric) as a trade-off that must be carefully selected for different settings.

6.4.1. MODEL UNCERTAINTY ON UNSEEN TASKS

Model uncertainty gives another tool for understanding forgetting. A good continual learning model should be more certain on tasks it has seen before than unseen tasks. To test the ability of model uncertainty to distinguish a task boundary, we assess the uncertainty of the model, for each batch of training, on data from the current task. We find, for a model trained using VGR, that model uncertainty spikes at the start of each new task (see figure 8). This means that VGR, or another method with sufficiently good model uncertainty estimation, does not need to be told task boundaries. Before each training epoch, we can compare the model uncertainty on the new data to the previous running average uncertainty. If the change in model uncertainty is above any sensible cut-off threshold we see a task boundary has occurred and trigger training of the generative model with a cache of the previous epoch’s data. One could either cache the last epoch’s data or train the generative model alongside the task training, but either way task detection comes with some additional memory and computational overheads. To measure uncertainty, we use the mutual information between predictions on each task and the model posterior, following Gal (2016).

6.4.2. TIME AND MEMORY

Timed Split MNIST evaluates the trade-off of training time against accuracy in continual learning, rather than reporting a single run-time. Much of the need for continual learning comes from situations where a system cannot feasibly retrain on all previous data due to constraints of cost and time. Architecture selection and hyperparameter selection, therefore involve a tradeoff of performance and speed. For example, in VGR, training on larger sampled datasets improves performance but takes more time. In figure 9 we show the average accuracy of various models plotted against the wall-clock time taken to finish training for a wide range of such hyper-parameters. While it is good to report the time

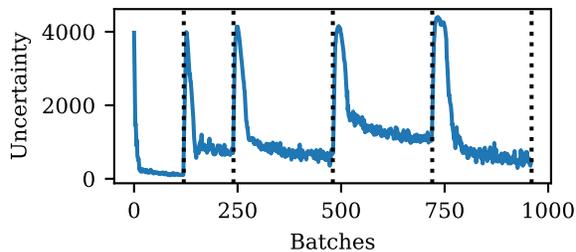


Figure 8. **Mutual Information on FashionMNIST with VGR.** Uncertainty spikes at the start of each new task (black dotted line), allowing task boundary detection. Uncertainty is assessed using the mutual information between predictions on each task and the model posterior.

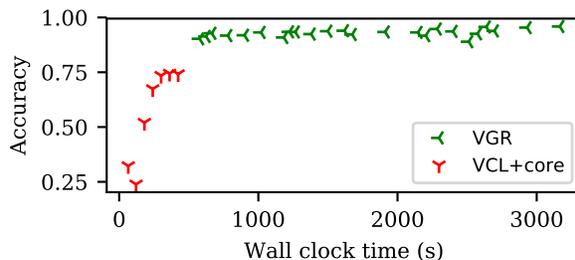


Figure 9. **Timed Single-headed Split MNIST.** Hyperparameter choice can lead to slower, more accurate, models for each architecture. Ideal performance is in the top left corner. VGR is accurate but slower than VCL.

taken by the most accurate model from one approach, representing the trade-offs within each approach gives a fairer reflection of the costs and benefits of each method. None of the configurations of VCL provides good accuracy on the single-headed task. At the same time, none of the configurations of VGR is able to provide good time performance (see Appendix B.7.1 for further details and hyperparameter choices).

7. Discussion

Experimental design shapes the field by providing reproducible comparisons between architectures, and testing whether an architecture satisfies fundamental objectives. Accidental blindspots in existing evaluations have created misleading comparisons in recent research. When experiments better reflect the continual learning problem, recent leading approaches fail even on simple datasets like MNIST. Using more realistic datasets is important, but until we improve the experimental design, the field will only show illusory progress on those datasets, without overcoming fundamental challenges inherent to continual learning itself. Authors of new papers currently feel compelled by the reviewing process to assess new methods on the same, incomplete, evaluations used by old papers. This paper reflects an attempt to reassess the experimental foundations of continual learning and provide a new direction for future work.

References

- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight Uncertainty in Neural Networks. *Proceedings of the 32nd International Conference on Machine Learning*, 37:1613–1622, 2015. 4.1
- Yaroslav Bulatov. notMNIST dataset, 2011. URL <http://yaroslavvb.com/upload/notMNIST/>. 1
- Arslan Chaudhry, Puneet K. Dokania, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Riemannian Walk for Incremental Learning: Understanding Forgetting and Intransigence. *Proceedings of the European Conference on Computer Vision*, 11215:556–572, 2018. 1, 4.1, 4.2.2, 4.2.4, 5.2, A.1
- Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient Lifelong Learning with A-GEM. *International Conference on Learning Representations*, 2019. 4.2.1, 4.2.4
- Natalia Diaz-Rodriguez, Vincenzo Lomonaco, David Filliat, and Davide Maltoni. Don’t forget, there is more than forgetting: new metrics for Continual Learning. *Continual Learning Workshop at NeurIPS*, October 2018. 4.2.4
- Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition. *Proceedings of the 31st International Conference on Machine Learning*, 31:647–655, 2014. A.5
- Sebastian Farquhar and Yarin Gal. Differentially Private Continual Learning. *Privacy in Machine Learning and AI workshop at ICML*, 2018a. 4.2.4
- Sebastian Farquhar and Yarin Gal. A Unifying Bayesian View of Continual Learning. *Bayesian Deep Learning Workshop at NeurIPS*, page 6, 2018b. 4.1, 6.1, A.4, B.2, B.3
- Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A. Rusu, Alexander Pritzel, and Daan Wierstra. PathNet: Evolution Channels Gradient Descent in Super Neural Networks. *arXiv*, 2017. A.5
- Yarin Gal. Uncertainty in Deep Learning. *PhD Thesis*, 2016. 6.4.1
- Ian J. Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An Empirical Investigation of Catastrophic Forgetting in Gradient-Based Neural Networks. *arXiv*, 2013a. 4.2.1, A.1, A.5
- Ian J. Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An Empirical Investigation of Catastrophic Forgetting in Gradient-Based Neural Networks. *arXiv*, 2013b. arXiv: 1312.6211 Citation Key: Goodfellow2013. 5.1
- Alex Graves. Practical Variational Inference for Neural Networks. *Neural Information Processing Systems*, 2011. B.2
- Wenpeng Hu, Zhou Lin, Bing Liu, Chongyang Tao, Zhengwei Tao, Jinwen Ma, Dongyan Zhao, and Rui Yan. Overcoming Catastrophic Forgetting via Model Adaptation. *International Conference on Learning Representations*, 2019. 4.1, 4.2.1, 4.2.2
- Ferenc Huszar. Note on the quadratic penalties in elastic weight consolidation. *PNAS*, 115(11):E2496–E2497, 2018. 4.1, B.3
- Heechul Jung, Jeongwoo Ju, Minju Jung, and Junmo Kim. Less-forgetting Learning in Deep Neural Networks. *arXiv*, 2016. 4.2.3, A.5
- Diederik P. Kingma, Tim Salimans, and Max Welling. Variational Dropout and the Local Reparameterization Trick. *Advances In Neural Information Processing Systems*, pages 2575–2583, 2015. B.2
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, March 2017. 1, 2, 4.1, 4.2.1, 5.3, 6.1, 6.3.1, A.1
- Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Technical report, 2009. 1
- Yann LeCun, Leon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *IEEE*, 86(11):2278–2324, 1998. 1
- Sang-woo Lee, Jin-hwa Kim, Jaehyun Jun, Jung-woo Ha, and Byoung-tak Zhang. Overcoming Catastrophic Forgetting by Incremental Moment Matching. *Advances in Neural Information Processing Systems*, 2017. 4.2.1, A.5
- Zhizhong Li and Derek Hoiem. Learning without Forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. 4.1, 4.2.3, A.5
- Vincenzo Lomonaco and Davide Maltoni. Core50: a new dataset and benchmark for continuous object recognition. *arxiv*, 2017. 3, 4.2.4
- David Lopez-paz and Marc’Aurelio Aurelio Ranzato. Gradient Episodic Memory for Continual Learning. *Proceedings of the 31st Conference on Neural Information Processing Systems*, 31, 2017. 4.1, 4.2.1, 4.2.4

- David J. C. MacKay. A Practical Bayesian Framework for Backpropagation Networks. *Neural Computation*, 4(3): 448–472, 1992. 4.1
- Radford M. Neal. *Bayesian learning for neural networks*. PhD thesis, University of Toronto, 1995. 4.1
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bis-sacco, Bo Wu, and Andrew Y. Ng. Reading Digits in Natural Images with Unsupervised Feature Learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011. 4.2.3
- Cuong V. Nguyen, Yingzhen Li, Thang D. Bui, and Richard E. Turner. Variational Continual Learning. *International Conference on Learning Representations*, 2018. 1, 2, 3, 4.1, 4.2.1, 4.2.2, 5, 6.1, 1, 6.3.1, A.3, A.4, B.2, B.3, B.4
- B. Pfulb and A. Gepperth. A comprehensive, application-oriented study of catastrophic forgetting in DNNs. *International Conference on Learning Representations*, September 2018. 4.2.3, 4.2.4, 5.3, B.5
- Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. CNN features off-the-shelf: An astounding baseline for recognition. *IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pages 512–519, 2014. A.5
- Hippolyt Ritter, Aleksandar Botev, and David Barber. On-line Structured Laplace Approximations For Overcoming Catastrophic Forgetting. *arXiv*, 2018. 1, 4.1, 4.2.1, 4.2.2, 5, A.1
- Anthony Robins. Catastrophic Forgetting, Rehearsal and Pseudorehearsal. *Connection Science*, 7(2):123–146, 1995. ISSN 13600494. 4.1
- Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive Neural Networks. *arXiv*, 2016. 4.1, A.5
- Jonathan Schwarz, Daniel Altman, Andrew Dudzik, Oriol Vinyals, Yee Whye Teh, and Razvan Pascanu. Towards a natural benchmark for continual learning. *Continual Learning Workshop at NeurIPS*, page 13, 2018a. 3, 4.2.4
- Jonathan Schwarz, Jelena Luketina, Wojciech M. Czarnecki, Agnieszka Grabska-Barwinska, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & Compress: A scalable framework for continual learning. *arXiv*, 2018b. 2, 3, 4.2.4
- Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual Learning with Deep Generative Replay. *Advances In Neural Information Processing Systems*, pages 2994–3003, 2017. 4.1, 4.2.1, 4.2.2, 4.2.3
- Rupesh Kumar Srivastava, Jonathan Masci, Sohrab Kazerou-nian, Faustino Gomez, and Jrgen Schmidhuber. Compete to Compute. *Advances in Neural Information Processing Systems*, pages 2310–2318, 2013. A.5
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a Novel Image Dataset for Benchmarking Ma-chine Learning Algorithms. *arXiv:1708.07747 [cs, stat]*, August 2017. 1
- Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lip-son. How transferable are features in deep neural net-works? *Advances In Neural Information Processing Systems*, pages 3320–3328, 2014. A.5
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Contin-ual Learning Through Synaptic Intelligence. *Proceedings of the 34th International Conference on Machine Learn-ing*, pages 3987–3995, 2017. 1, 4.1, 4.2.1, 4.2.2, 5, 6.3.2, A.2

Appendices

A. Further Prior Work

A.1. Elastic Weight Consolidation

Unlike VCL, EWC uses an ordinary neural networks with an L_2 regularization term added to the loss. This regularization reflects a Gaussian prior for each parameter. At the end of training on each dataset, EWC (Kirkpatrick et al., 2017) estimates the contribution each parameter makes to the gradient of the loss by approximating the Hessian of the likelihoods with:

$$F_t \propto \sum_{n=1}^{N_t} \left(\nabla_{\omega} \log p(y_t^{(n)} | \omega, \mathbf{x}_t^{(n)}) \right)^2 \Big|_{\omega=\omega_t}$$

This approximate Fisher information matrix is estimated at the end of each task. This implicitly sets the covariance of the parameter prior (which is assumed to be diagonal). The authors state that in the multi-task case the regularization term for task T is the sum of separate regularization terms for each past dataset \mathcal{D}_t .

Chaudhry et al. (2018) provide a generalization of both SI and EWC. Ritter et al. (2018) extend EWC by relaxing the assumption that the covariance matrix of the parameter distributions is diagonal.

Kirkpatrick et al. (2017) show EWC works well on the Permuted MNIST task introduced by Goodfellow et al. (2013a). They also show that EWC reduces forgetting on a succession of Atari games.

A.2. Synaptic Intelligence

Published soon after EWC, SI has a similar loss but computes the contribution of each parameter to the gradient of the loss over the entire course of training. The authors explicitly acknowledge that the derivation of their regularization term only extends to the two-task case, but point out that the model performs well even in multi-task settings.

Zenke et al. (2017) show that SI works in a *multi-headed* version of the Split MNIST task which they introduce, and matches EWC on the Permuted MNIST task. They also show some transfer learning on a multi-headed split CIFAR task of their own design.

A.3. Variational Continual Learning

Nguyen et al. (2018) motivate their regularization through VI. They argue that the posterior of the parameter distribution after learning a task should be the prior when learning the next task. They calculate the Kullback-Leibler (KL) divergence between the current distribution and the previous

posterior. This is in contrast with the more usual Variational Free Energy (VFE) loss which uses the KL divergence to a constant prior that is not updated after each task. The KL divergence can be represented approximately as a quadratic regularization with rotation, just like the two approaches above, as Nguyen et al. (2018) point out.

In addition to the Variational Continual Learning (VCL) regularization term, Nguyen et al. (2018) add a coreset which samples examples from old datasets, which are withheld from the main dataset. Rather than discarding all of an old dataset, as EWC and SI do, only the vast majority (about 99.7%) is discarded. Before being evaluated, the model is trained against the coreset. For multi-headed settings (see §5.2) coresets train each head separately.

Nguyen et al. (2018) show that VCL with or without the use of coresets outperforms both SI and EWC on the Permuted MNIST task. They show that VCL outperforms EWC and SI on the multi-headed Split MNIST task and on a very similar task with notMNIST. They take advantage of the VI setting to show reduced uncertainty in a generative task based on MNIST and notMNIST.

A.4. Variational Generative Replay

Each of VCL, EWC, and SI effectively set the parameters at the end of the last task as a prior, which is explicit in the Variational Inference (VI) derivation supporting VCL (Nguyen et al., 2018). Instead of changing priors between tasks, VGR adapts the log-likelihood component of the loss to depend on past datasets (Farquhar and Gal, 2018b). VGR estimates the log-likelihood component of the loss on past tasks using generative models. They train a DC-GAN at the end of each task on all of the classes that appeared in that task. Then, at the start of training a new task, they sample from all stored GANs training points, which are then mixed into the real training data.

A.5. Dynamic architecture approaches

Other work, which is not our focus, has also made progress with continual learning. Progressive neural networks (Rusu et al., 2016) conditions the newest model on the outputs of old, stored, models. Some approaches split the network into regions that have different roles. In Li and Hoiem (2017) one set of shared parameters governs a feature extraction component while each task is given parameters on top of that component. During training, the combined shared and old task networks are encouraged to classify similarly to a stored old version of the model, while the combined shared and new task networks are trained on the new task. Jung et al. (2016) stochastically freeze the final layer to encourage lower layers to extract features from all tasks that the final layer can use to classify. Others have semi- or fully-fixed a shared part of a network with task-specific layers on top

ENTROPY	
Split	0.003
Permuted	0.453

Figure 10. Average entropy of predictions on Task B, early in training; Note the 2 orders of magnitude difference between the two settings. Entropy is much higher in the Permuted setting.

(Razavian et al., 2014; Yosinski et al., 2014; Donahue et al., 2014). PathNet uses evolutionary selection to try to learn which patches of a larger network are helpful to each task rather than making a layer-by-layer assumption (Fernando et al., 2017). This broad family of approaches are especially useful in two-task transfer cases, but it can become impractical to introduce the many task-specific weights in the more general continual learning case that we consider.

Others find that adjustments to learning dynamics can reduce catastrophic forgetting, presumably by encouraging networks to use their capacity fully. This family of approaches includes selecting activations and using dropout to reduce forgetting (Srivastava et al., 2013; Goodfellow et al., 2013a). These are useful, but do not go far enough to solve forgetting fully. Lee et al. (2017) extend dropout by shifting the zero-point of each dropout parameter to the parameter value from training on the previous task.

B. Experimental Details and Further Figures

B.1. Shortcomings of Permuted MNIST

Although our work mostly aims to show that architectures which succeed on a Permuted MNIST task can fail in slightly different settings, we also carried out some investigations as to why this might be. We observe that in most settings, a model will confidently predict that an example from the new dataset \mathcal{D}_t is from a class that was in \mathcal{D}_{t-1} . This confident but false prediction creates large gradients in the output layer. The derivative of the likelihood term in the loss with respect to each output weight w_k in the output layer of a model is $\frac{\partial \mathcal{L}}{\partial w_k} = p_k - y_k$ where y_k is the k 'th entry of a one-hot vector of labels and p_k is the probability predicted for class k . This gradient is therefore biggest for confident but false predictions. In the permuted setting, however, no example from \mathcal{D}_t looks remotely like an example from \mathcal{D}_{t-1} . This means the model makes unconfident predictions and the gradients stemming from the likelihood term of the loss are unusually small. This may on its own be enough to explain part of the difference in the settings.

We may be able to probe more deeply. For prior-focused continual learning, the loss function is separated into a prior term and a likelihood term (c.f. eq. (1)). By examining this separation, we can see why the permuted setting, with unusually small likelihood-term gradients, is a best case

scenario for prior-focused continual learning. The prior term is *data-independent* and has the same magnitude no matter how incorrectly confident the model is when making predictions on \mathcal{D}_t . The likelihood term of the gradient, however, has a much larger magnitude when confident but false predictions are made. This means that in a typical setting the likelihood term of the loss dominates the prior early in training. In the permuted setting, however, the prior term is relevant throughout the training. This means that the prior-focused approaches are able to succeed in a permuted setting but fail in other settings. Two points are worth noting. First, if the prior term were to fully capture correct class of models, its size will change appropriately, and this effect would not exist. But in practice our priors are restrictive approximations. Second, the gradients in earlier layers are less straightforward to analyse. However, especially early in training, output layer gradients are several orders of magnitude larger than gradients in early layers. Moreover, our argument only depends on showing that the Permuted setting protects the prior-focused loss from systematic large mis-estimates for an important subset of its gradients.

We introduce two tests to show the shortcomings of the Permuted MNIST experiment discussed in §5.1. For each test we consider both the Permuted and Split settings. In each case, we train a model on Task A before training on Task B. In the split setting, Task A is the first five digits of MNIST and Task B the last five digits. In the permuted setting, Task A is MNIST and Task B is a random permutation of the pixels of MNIST. Here, we are assessing the evaluation framework itself, not continual learning performance.

B.1.1. VERIFYING PERMUTED SETTING UNDER-CONFIDENCE

We first validate our hypothesis that predictions are much more uniform in the Permuted setting. We evaluate the entropy of the output probability vector. Low entropy indicates a confident prediction of a single class, whereas a high entropy indicates a ‘uniform’ prediction, spreading the mass across different classes. These are summarised in Table 10, supporting our claim that models in the Permuted setting are much less confident than in a more typical setting.

B.1.2. VERIFYING UNUSUALLY LARGE GRADIENTS IN THE PERMUTED SETTING

Next, we test our hypothesis from §5.1 that more uniform predictions lead to slower ‘catastrophic forgetting’. To test this hypothesis, after training on Task A, we measured gradients for the final weight layer while training Task B in the two variants (figure 11). All results are averaged over 100 runs. Full experimental details are in Appendix B.1.3. We found that having digits resembling previously observed ones in the Split setting led to much larger likelihood-term

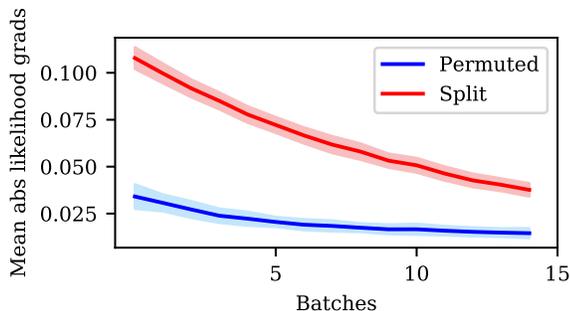


Figure 11. Early in training Task B, the likelihood term of gradients on the final layer is unusually low in the Permutated setting because permuted digits do not resemble any digits from Task A. This makes continual learning unrealistically easy in this evaluation. Averaged over 100 runs, shading is one standard deviation.

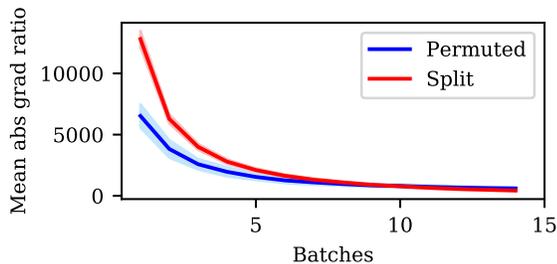


Figure 12. Early in training Task B, the ratio of the likelihood term of gradients the final layer against the prior term is much lower in the Permutated setting. This reduces forgetting because the prior has a larger impact on learning. Averaged over 100 runs, shading is one standard deviation.

gradients in the final layer than found in the Permutated setting for early batches. As a result, the ratio between the likelihood term and the prior term of the loss is much higher in the Split setting (figure 12). This means that the prior has a smaller influence over the gradients, leading to more forgetting.

B.1.3. EXPERIMENTAL SETTINGS FOR THE ABOVE GRADIENT ANALYSIS

To measure the gradients during training, we first trained on Task A for 120 epochs with batch size 256 and an Adam optimizer using Variational Continual Learning (VCL) with the same settings used in B.2. We then trained on Task B for 15 batches of 16 digits each, again using VCL. We measured the average absolute value of the gradients of the final layer of the model. We averaged over 100 training runs for each setting (with a different permutation each time in the Permutated setting), resetting the model to its initial position after each run. Graphs show the standard deviation of the average gradient of each batch. The ratio displayed in figure 12 is the log-likelihood term of the loss described in equation 1 divided by the prior term, using VCL. The ratio is calculated for each of the 100 runs on the averaged absolute gradients.

B.1.4. INVESTIGATING THE HYPOTHESIS THAT THE PRIOR TERM IS OVERLY SMALL

In order to investigate whether the prior-term of the gradient is overly small in the Split setting relative to the Permutated setting, we tried training a VCL model on Single-headed Split MNIST with an artificially upweighted prior-term in the cost. We experimented with factors of 1 (the default case), 10, 100, and 1000. We also tried upweighting the prior-term only during the first few epochs. In no case did we find that this improved performance significantly. This suggests that it is not only the magnitude, but also the direction, of the prior gradient which is inaccurate in the split setting.

B.2. Permutated MNIST as Baseline

As Nguyen et al. (2018) have shown, on Permutated MNIST, described in §5.1, VCL performs well with or without a coreset. We further show that the coreset on its own performs badly (figure 13, results averaged over 10 runs). This is consistent with our analysis of the Permutated setting, which suggested that the small and evenly spread gradients due to the likelihood term of the loss helped the prior term to work, despite a prior that does not fully support the function space. Note that this version of VCL was previously shown to outperform EWC and SI (Nguyen et al., 2018).

For permuted MNIST, we follow Nguyen et al. (2018) where possible. We use a Bayesian neural network (Graves, 2011)

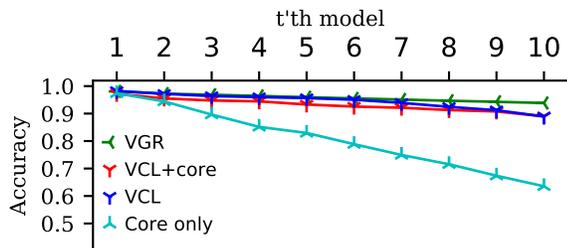


Figure 13. **Permuted MNIST.** VGR has marginally better performance than VCL with or without a coreset, which has been shown to outperform EWC and SI. Coresets without the VCL loss performs worse. Averaged over 10 runs.

with two hidden layers of 100 units with ReLU activations. The priors are initialized as a unit Gaussian and the parameters are initialized with the mean of a pre-trained maximum likelihood model and a small initial variance (10^{-6}). We use the Adam optimizer (Kingma et al., 2015) with learning rate 10^{-3} . We use a single head, again following Nguyen et al. (2018). For all results we present the average over 10 runs, with a different permutation each time. Standard errors are not shown as they are under a tenth of a percent.

We train for 100 epochs using a batch size of 256 on all the data except the with-held coresets. We train the whole single head on coresets for 100 epochs and use 200 digits of each permutation as a coreset chosen using the same k-center coresets algorithm used by Nguyen et al. (2018). VCL without coresets is exactly the same, but without a final training step on coresets. The coresets only algorithm is exactly the same as VCL, except that the prior is always initialized as though it were the first task.

Following Farquhar and Gal (2018b) where possible, we train VGR for 120 epochs using a GAN trained for 200 epochs on each MNIST digit. We use 6000 generated digits per class, sampled fresh for each task, and initialize network weights using the previous task. We use a batch size of 256 times the number of seen tasks, ensuring that the number of batches is held constant. The GAN is trained with an Adam optimizer with learning rate $2 * 10^{-4}$ and β_1 of 0.5. The network has four fully-connected hidden layers with 256, 512, 1024 and 784 weights respectively. It uses Leaky ReLU with α of 0.2.

B.3. Experimental Settings for Single-headed Split MNIST

The settings for Split MNIST follow Nguyen et al. (2018) where possible. For Bayesian neural network architectures, we use exactly the same settings as for Permuted MNIST, including the single head, except that each hidden layer has

256 weights, similarly to Nguyen et al. (2018) on their multi-headed Split MNIST. Results are shown averaged over 10 runs, with a different coreset selection each time. Standard errors are not shown as they are of the order of a tenth of a percent.

For all Bayesian neural network architectures we train for 120 epochs. We use batch sizes equal to the training set size. We use coresets of digits per task selected using the same k-center coreset algorithm as Nguyen et al. (2018), which are withheld from the training. We train for 120 epochs on the concatenated coreset across all the heads together.

Following Farquhar and Gal (2018b) where possible, for VGR we use 6000 digits per class generated by a convolutional GAN. Unlike VCL, we cap batch sizes at 30,000 rather than having the batch size equal the training set size. The GAN is trained for 50 epochs on each MNIST class using the same optimizer as the non-convolutional GAN used for Permuted MNIST. It has a fully connected layer followed by two convolutional layers with 64 and 1 channel(s) and 5x5 convolutions. Each convolutional layer is preceded by a 2x2 up-sampling layer. The activations are Leaky ReLU’s with α of 0.2.

For EWC, we use a neural network with two hidden layers with 256 weights. We train using the variant of EWC suggested in Huszar (2018). We estimate the Fisher information using 200 samples at the end of each task, train using SGD with learning rate 10^{-2} for 20 epochs per task. The coefficient of the regularization term in the loss was set to 10 but was found not to be important to performance.

B.4. Experimental Settings for Multi-headed Split MNIST

Multi-headed Split MNIST has almost precisely the same settings as the Single-headed Split MNIST above. Following Nguyen et al. (2018), for VCL we have five heads in the final layer, each of which is trained separately. Coresets are used to train each head in turn. Batch size equal is to training set size. Results are shown averaged over 10 runs with fresh coreset selection each time. Standard errors are not shown as they are of the order of a tenth of a percent.

EWC is configured in exactly the same way as in the single-headed setting.

B.5. Experimental Settings for Split Fashion MNIST

Single- and multi-headed Split Fashion MNIST are performed mostly the same as our Split MNIST experiments. We use a larger network, with four hidden layers with 200 units for the BNN and 256 units for EWC. Fashion MNIST has a much more diverse membership of its classes, which makes performance lower even with these larger networks. We deliberately did not optimize hyperparameters, given the

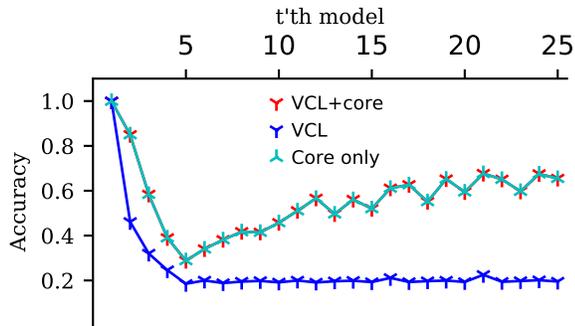


Figure 14. **Retraining MNIST.** With retraining, single-headed Split MNIST is still able to differentiate likelihood-focused and prior-focused approaches.

observation of Pfulb and Geppert (2018) concerning causality. However, once hyperparameters were set, we explored somewhat and found that using larger coresets and networks could improved performance slightly, but not dramatically. For VGR, the DC-GAN followed the implementation at <https://github.com/carpedm20/DCGAN-tensorflow>

B.6. Split MNIST with Retraining

To test the impact of retraining, we allowed ourselves to cycle through all the data 5 times, which is roughly what previous authors who have made use of retraining have done. Instead of finishing at the end of the fifth task, we begin fresh (keeping all coresets if applicable). We find that even with retraining, single-headed Split MNIST is still able to differentiate prior- and likelihood-focused approaches (see figure 14). This suggests that experiments with retraining may not be misleading in that they make models that fail look as though they work, though it certainly conflicts with some of the use cases that motivate continual learning. For these experiments, we reduced the number of epochs trained per task from 120 to 10, but other settings are the same.

B.7. Timed Single-headed Split MNIST

We perform the single-headed Split MNIST experiment with a range of different configurations. For VGR we allow between 10 and 50 epochs for training the convolutional GANs and between 1 and 120 epochs for training the main model. We use 2000, 4000, or 6000 generated images per class. We add GAN training time to the main figure, even though it is possible to do in parallel. For VCL, we use between 1 and 120 epochs for training and coresets. In all cases, we report elapsed wall-clock time from start to finish. We plot this against average accuracy over all tasks of the final model trained on all tasks. In all cases, the experiment was carried out on an NVIDIA Tesla K80.

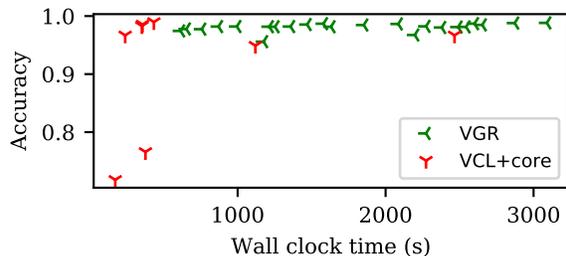


Figure 15. **Multi-headed Timed Split MNIST.** Ideal performance is in the top left corner. VGR can reach good accuracy, but slower than VCL. Now that VCL can use multi-headed training as well as coresets, the accuracy is good in the multi-headed experiment.

Here, we also show Timed Split MNIST performed as a multi-headed experiment in order to allow VCL to use multi-headed training (figure 15). Performance is better for all models than in the single-headed case, but for VCL it can get very good, and faster than VGR.

B.7.1. TIME AND MEMORY EFFICIENCY OF PRIOR- AND LIKELIHOOD-FOCUSED APPROACHES

Likelihood-focused approaches must train generative models and train on extra data relative to prior-focused approaches. But prior-focused approaches tend have slightly more complicated losses or add extra training phases. We found that VCL with coreset was faster than VGR. For the single-headed Split MNIST task it took roughly 7 minutes.² Training VGR with each generated set the same size as the true data took roughly 18 minutes plus 34 minutes to train the GANs. But a quicker training regime took only 2 minutes plus 8 for training GANs and was still more accurate than VCL.³ The prior-focused costs scale badly with model size, while VGR’s costs scale badly with data-space complexity and very long series of datasets. VGR is also potentially less memory intensive. Generated data can be sampled on demand, it need not be the same each epoch.

B.8. A Simpler Model Uncertainty Experiment

It is possible to measure the quality of model uncertainty for continual learning without tracking throughout training. After training on each task, we measure the model’s uncertainty when shown data from each task, both seen and unseen (see figures 18 and 19). We then set an uncertainty threshold, when the uncertainty is above this value

²Times are reported for a Nvidia Tesla K80.

³Average single-headed accuracy of the final model for this less-trained VGR was 90% rather than 96% for the fully trained model. Training used one third the generated data, only 20 epochs of training, and only a fifth the training for the GAN.

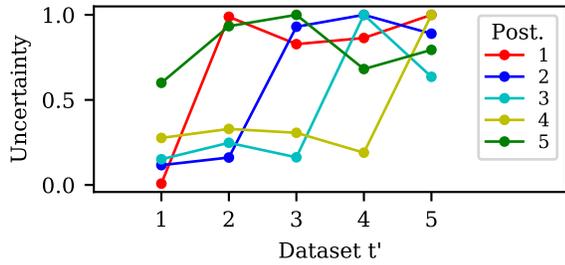


Figure 16. **Mutual Information MNIST VGR.** Uncertainty is high for unseen tasks, and low for all previously seen tasks, showing good calibration. Models trained on 1-5 tasks are each evaluated against all five tasks. Uncertainty is assessed using the mutual information between predictions on each task and the model posterior. It is normalized to 1 for the most uncertain task for each model.

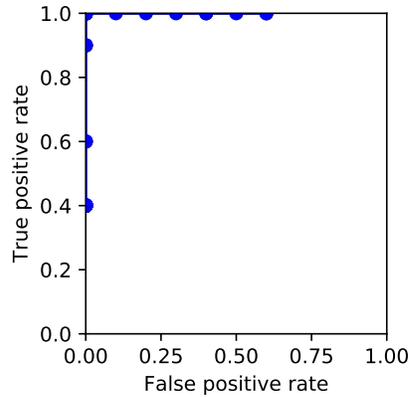


Figure 18. For VGR, sensible cut-off points cleanly differentiate previously-seen tasks from unseen ones with no false positives. This results in perfect recognition of previously seen tasks.

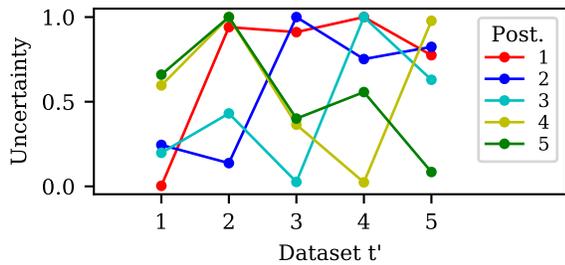


Figure 17. **Mutual Information MNIST VCL.** Uncertainty is high for unseen tasks, but rises for seen tasks that were seen some time ago. The fourth model, for example, is nearly as uncertain about the second task, which it has seen, as it is about the fifth, which it has not.

we predict that the data was previously unseen. By varying thresholds we generate an ROC plot. The area under the curve (AUC) of the ROC plot is a measure for the ability of the model to distinguish seen/unseen tasks. We compare the AUC of both these ROC plots and find that VGR’s AUC is 1 while VCL’s is 0.76. This means that VGR is in principle able to correctly detect all tasks it has not seen before, whereas VCL fails on a considerable number of those.

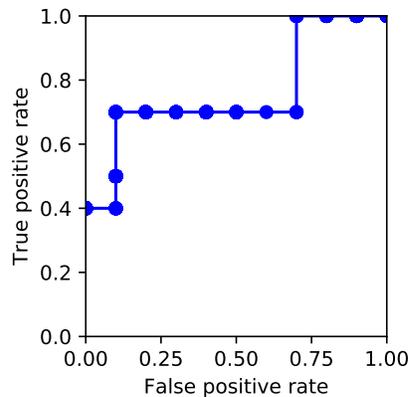


Figure 19. For VCL, false positives are a problem because the model is highly uncertain about many previously-seen tasks. This represents forgetting of old tasks.