



Fine-Grained and Continual Visual Recognition for Assisting Visually Impaired People

MARCUS KLASSON

Doctoral Thesis
Stockholm, Sweden, 2022

KTH Royal Institute of Technology
School of Electrical Engineering and Computer Science
TRITA-EECS- AVL-2020:4 Division of Robotics, Perception, and Learning
ISBN 100- SE-10044 Stockholm, Sweden

Akademisk avhandling som med tillstånd av Kungl Tekniska högskolan framlägges till offentlig granskning för avläggande av Teknologie doktorexamen i elektroteknik fredagen den 18 januari 2022 klockan 14.00 i Sal F3, Lindstedtsvägen 26, Kungliga Tekniska Högskolan, Stockholm.

© Marcus Klasson, date

Tryck: Universitetsservice US AB

Abstract

In recent years, computer vision-based assistive systems have enabled visually impaired people to use automatic visual recognition on their mobile phones. These systems should be capable of recognizing objects on fine-grained levels to provide the user with accurate predictions. Additionally, the user should have the option to update the system continuously to recognize new objects of interest. However, there are several challenges that need to be tackled to enable such features with assistive vision systems in real and highly-varying environments. For instance, fine-grained image recognition usually requires large amounts of labeled data to be robust, and image classifiers struggle with retaining performance of previously learned abilities when adapting to new tasks. This thesis is divided into two parts where we address these challenges. First, we focus on the application of using assistive vision devices for grocery shopping where items are naturally structured based on fine-grained details. We demonstrate how image classifiers can be trained with a combination of natural images and web-scraped information about the groceries to obtain more accurate classification performance compared to only using natural images for training. Secondly, we bring forward a new approach for continual learning called replay scheduling where we select which tasks to replay at different times to improve memory retention. Furthermore, we propose a novel framework for learning replay scheduling policies that can generalize to new continual learning scenarios for mitigating the catastrophic forgetting effect in image classifiers. This thesis provides insights on practical challenges that need to be addressed to enhance the usefulness of computer vision for assisting the visually impaired in real-world scenarios.

Sammanfattning

hej

List of Papers

- A *A Hierarchical Grocery Store Image Dataset with Visual and Semantic Labels*
Marcus Klasson, Cheng Zhang, Hedvig Kjellström
In *IEEE Winter Conference on Applications of Computer Vision (2019)*
- B *Using Variational Multi-View Learning for Classification of Grocery Items*
Marcus Klasson, Cheng Zhang, Hedvig Kjellström
In *Patterns, Volume 1(8) (2020)*
- C *Learn the Time to Learn: Replay Scheduling for Continual Learning*
Marcus Klasson, Hedvig Kjellström, Cheng Zhang
Under submission
- D *Policy Learning for Replay Scheduling in Continual Learning*
Marcus Klasson, Hedvig Kjellström, Cheng Zhang
Under preparation

Acronyms

List of commonly used acronyms:

A2C	Advantage Actor-Critic
DQN	Deep Q-Network
CL	Continual Learning
CNN	Convolutional Neural Network
FGIR	Fine-Grained Image Recognition
LSTM	Long Short-Term Memory
MCTS	Monte Carlo Tree Search
MLP	Multilayer Perceptron
PCA	Principal Component Analysis
RL	Reinforcement Learning
RNN	Recurrent Neural Network
SGD	Stochastic Gradient Descent
VAE	Variational Autoencoder
VCCA	Variational Canonical Correlation Analysis
VI	Visually Impaired

Contents

List of Papers	iii
Acronyms	v
Contents	1
I Overview	3
1 Introduction	5
1.1 Visual Impairments	6
1.2 Assistive Vision Technologies	6
1.3 Scope of Thesis	7
1.4 Thesis Contributions	9
1.5 Thesis Outline	12
2 Background	13
2.1 Problem Settings in Machine Learning	13
2.2 Deep Learning	14
3 Fine-Grained Image Recognition	21
3.1 Related Work	21
3.2 Dataset Collection	23
3.3 Fine-Grained Image Recognition of Grocery Items	24
3.4 Experiments	26
3.5 Discussion	29
4 Continual Learning	31
4.1 Related Work	31
4.2 Replay Scheduling in Continual Learning	33
4.3 Policy Learning for Replay Scheduling	36
4.4 Experiments	37
4.5 Discussion	41
5 Conclusions and Future Directions	43
5.1 Conclusions	43
5.2 Future Directions	45

Bibliography	49
II Included Papers	75
A A Hierarchical Grocery Store Image Dataset with Visual and Semantic Labels	77
B Using Variational Multi-View Learning for Classification of Grocery Items	93
C Learn the Time to Learn: Replay Scheduling in Continual Learning	139
D Meta Policy Learning for Replay Scheduling in Continual Learning	167

Part I

Overview

Chapter 1

Introduction

Vision is probably the most important of all senses that humans possess. Our society is built on having the ability to see. For example, if we want to cross a street, thick colored stripes on the road and signs above our heads indicate where the cross walk is located in order for us to cross the street appropriately. Another example is how we use text to communicate with each other, where messages that we read and write consist of sequences of symbols that are structured into words of some specific language. Furthermore, visuals are used frequently in education to help students comprehend new information and concepts. Basically, many everyday tasks such as traveling, communication, and learning become more convenient with normal vision.

Millions of people have visual impairments that cannot be treated with correctable eyewear [1]. To enhance the mobility of visually impaired (VI) people, there exist various kinds of assistive tools, such as screen readers and Braille typewriting machines, for supporting them with receiving information and communicating through text. More recently, several assistive vision technologies have emerged in the form of mobile phone applications [2–7] and wearable devices [8–10]. Furthermore, several works in computer vision research has been specifically focused developing methods for assisting VI people with e.g. object recognition [11–18] and wayfinding [19–24].

Despite the recent advances in deep learning for computer vision [25–29], several challenges have been demonstrated when deploying computer vision models in real-world scenarios. For example, the models may surprisingly fail to recognize known objects appearing in new poses [30], lighting conditions [31], and geographical locations [32], models trained on images from one set of hospitals may perform poorly on image conditions from another [33], and models can suffer from biases that disadvantage certain users based on their gender and ethnicity [34]. Part of the reason for these challenges is that specifying a model that can be injected with the rich complexity of the real world is difficult [35]. To enhance their applicability in the real-world, researchers are striving for developing computer vision models that can recognize objects in detail [36], learn new abilities continuously [37, 38], and adapt fast to new scenarios from few examples [39–41]. Furthermore, such tasks should be possible to execute in a time-efficient and robust manner in assistive vision devices to improve the user experience.

In this thesis, we address challenges on computer vision models for recognizing objects and their fine-grained details, as well as continually adapting the models to

new concepts while retaining its previously learned knowledge. We will continue this chapter by briefly describing visual impairments in Section 1.1, followed by a summary of assistive vision technologies in Section 1.2. Section 1.3 describes the scope of the thesis, and we summarize the contributions of the included papers in Section 1.4. In Section 1.5, we present the overall outline of this thesis.

1.1 Visual Impairments

A visual impairment is defined as the decrease of an individual's ability to see from various distances [42]. The conditions varies from having issues with seeing from far or near distances to blindness, i.e., complete or nearly-complete vision loss. Visual capabilities are in general assessed by measuring the visual acuity (sharpness) of seeing from some fixed distance, e.g., letters of different size on an eye chart. The visual acuity is measured differently based on whether near- or far-sighted VI people are being assessed. For far-sighted VI people, visual acuity is calculated by the ratio between the distance that the subject and a normal-sighted person can recognize the target object. For near-sighted VI people, visual acuity is measured using a standardized point system defining the font size of letters that the subject can recognize [43].

In 2020, it was estimated that 338 million people possess a moderate to severe visual impairment globally, including 43 million people that are blind [1]. Furthermore, the World Health Organization (WHO) have estimated that at least 2.2 billion people live with a near or distance visual impairment, where at least 1 billion cases could have been prevented or yet has to be addressed [43]. The number of untreated cases is projected to be 1.7 billion people by 2050, mainly due to global population growth as well as increased aging among the populations [1]. The main causes for having vision loss are uncorrected refractive errors, untreated cataracts, age-related macular degeneration, glaucoma, diabetic retinopathy, where 90% of such cases are preventable and treatable [44]. These causes also varies between different countries and regions with different socioeconomic status.

There exists many kinds of tools for assisting VI people with everyday tasks. The white cane and guiding dogs are probably the most common tools for enhancing their mobility by preparing the user for what is present in their near surroundings. Furthermore, there are different tools for various recognition tasks, such as currency markers for identifying different bills in wallets, color indicators for notifying the user of the color of clothes, and labeling apparatus with Braille writing for distinguishing between similar items [45]. There also exist Braille keyboards and screen readers compatible with both computers and mobile phones for helping VI people with e.g. sending text messages and emails, surfing the web, and using calendar apps [46]. In the next section, we introduce a selection of assistive vision technologies aimed for assisting VI people with recognition of objects in their environment.

1.2 Assistive Vision Technologies

Cameras are used by blind people to record events and memories similarly as people with normal vision [47]. Moreover, there is an urge among VI people to use their

cameras for daily tasks other than recording events, such as object recognition [48], face recognition [49], and currency identification [50]. In this thesis, we focus on object recognition which can be challenging in many ways without visual information. For instance, VI people may need assistance from non-visual cues indicating how to capture the whole item in the camera view [51], and also may need to point the camera towards some specific side of the object to identify details of interest to the user. These challenges have encouraged the development of technical aids that use computer vision for assisting VI people which we will cover next.

In the last decade, we have seen several variants of assistive vision technologies emerging on the market. The advances of such technologies have made it possible for VI people to reduce the need for human assistance for performing daily tasks where vision is a necessity. Currently, these devices have taken the form of both mobile phone applications [2–5] and wearable devices [8–10] where computer vision methods are used for visual tasks, such as object and face recognition, barcode scanning, color and currency identification, and text recognition. An alternative to the computer vision-based technologies are other mobile phone apps where VI users have video calls with sighted volunteers that help them with any kind of task requiring visual capabilities [6, 7].

Despite the increased availability of assistive vision technologies, there are still several challenges remaining to be tackled in order to enhance the user experience. Although machine learning techniques have been successfully applied in computer vision applications such as object recognition [25, 26, 52], generating scene descriptions [28, 29, 53], and visual question answering [54–56], these methods require immense amounts of human-annotated data for providing accurate predictions. The annotation becomes even more labor-intensive in scenarios where the assistive vision system must distinguish between objects with subtle differences, making it challenging to provide users with more detailed information than the general object class. Another challenge is how to continuously adapt the system to recognize new objects while ensuring that the system is still capable of classifying previous known items correctly. Furthermore, assistive vision technologies need system requirements to provide better usability, where the system should process images in real-time, be robust when applied in new environments, as well as ensuring privacy for the user.

1.3 Scope of Thesis

This thesis focuses on two topics for machine learning and computer vision-based assistive technologies, namely fine-grained image recognition (FGIR) [36] and continual learning (CL) [37, 38]. The FGIR task involves identifying visually similar subcategories of some general object classes. An example is when one has to distinguish between two kinds of visually similar juice packages with different ingredients. In Paper A and B, we study FGIR from the real-world application of recognizing groceries from mobile phone with an assistive vision app. The CL setting targets the scenarios where a classifier receives data from new object classes to learn from a continuous stream, while retaining its capability of recognizing previously learned objects. In Paper C and D, we introduce a novel approach for improving the retention of previously learned abilities of the classifier. The common denominator of these topics is image classification, However, both have challenges of their own that have to be

addressed for enabling such adding such features into assistive vision technologies. Next, we describe the challenges that we have focused on in this thesis.

Challenges in Fine-Grained Image Recognition

One main challenge for FGIR is the required data collection procedure. Firstly, annotating the specific details about the objects in the collected data becomes more time-consuming for the annotators, and may also require expert knowledge of the domain. Secondly, as some fine-grained classes could be rare, the dataset may be unbalanced making it difficult for the classifier classify the rare classes accurately from only a few training samples. A real-world application for helping VI people using assistive vision technologies that contains both of these challenges is grocery shopping [12, 57–62]. Grocery items usually require visual information to distinguish between them, e.g., when distinguishing between two juice packages with the same shape but different ingredients. Similarly for raw groceries, it can be difficult to tell the difference between two different kinds of apples without the ability to see, unless the customer knows how the apples smell or how the texture of their peel differs when touching them. Furthermore, the conditions of the grocery store environment can disturb the recognition performance for computer vision-based assistive devices, e.g., when multiple and misplaced items appear in the camera view as well as poor lighting conditions in particular areas. Collecting training data that covers all possible scenarios that can occur in the store would be a nearly impossible task. Our goal is to reduce the need for large amounts of training data in the grocery stores by collecting web-scraped information about the items to utilize when training the image classifier.

Challenges in Continual Learning

The main challenge in CL is catastrophic forgetting [63] where the classifier overwrites previously learned knowledge with information about recently learned objects. Replay-based CL [37, 38] is a simple yet efficient approach for mitigating catastrophic forgetting where samples from old classes are stored in a small memory which can be re-trained on when learning new classes. In contrast to the small memory requirement, many machine learning systems in real-world applications are limited by constraints on processing times rather than data storage capacity [64–66]. In such settings, retraining on all seen data is often prohibited, the system would need a method for selecting what data to replay from a huge amount of stored data to mitigate catastrophic forgetting. Motivated by human learning techniques for improving memory retention [67–71], we propose a novel approach called replay scheduling where we have policy for selecting which tasks to replay at different times. Our goal is to demonstrate that scheduling over which tasks to replay can be crucial for CL systems in settings where immense amounts of historical data is accessible. To enable our approach in real-world CL scenarios, we need to develop efficient methods that can propose replay schedules to mitigate catastrophic forgetting in new CL scenarios without additional computations.

Table 1.1: Two examples of grocery item classes in the Grocery Store dataset presented in Paper A. We display two class labels (coarse-grained class in parenthesis), two natural images taken with a mobile phone inside grocery stores, and the web-scraped iconic images and text description of the items.

Class Labels	Natural Images	Iconic Images	Text Descriptions
Granny Smith (Apple)	 		<i>“...green apple with white, firm pulp and a clear acidity in the flavor.”</i>
Tropicana Mandarin (Juice)	 		<i>“...is a ready to drink juice without pulp pressed on orange, mandarin and grapes. Not from concentrate. Mildly pasteurized.”</i>

1.4 Thesis Contributions

In this section, we provide summaries of the included papers as well as stating the contributions of each author to the manuscripts.

Paper A: A Hierarchical Grocery Store Image Dataset with Visual and Semantic Labels

Marcus Klasson, Cheng Zhang, Hedvig Kjellström. In *IEEE Winter Conference on Applications of Computer Vision (WACV) 2019*.

Summary: We collect a dataset with natural images of raw and refrigerated grocery items taken in grocery stores in Stockholm, Sweden, for evaluating image classifiers on a challenging real-world scenario. The data collection was performed by taking photos of groceries with a mobile phone to simulate a scenario of grocery shopping using an assistive vision app. Furthermore, we downloaded iconic images and text descriptions of each grocery item by web-scraping a grocery store website to enhance the dataset with information describing the semantics of each individual item. The items are grouped based on their type, e.g., apple, juice, etc., to provide the dataset with a hierarchical labeling structure. We show two examples of grocery item classes and their corresponding web-scraped information in Table 1.1. We provide benchmark results evaluated using pre-trained and fine-tuned convolutional networks for image classification.

Author Contributions: CZ and HK presented the idea and the data collection procedure for the natural images and web-scraped information. MK performed the data collection including visiting the grocery stores for taking the natural images and the web-scraping of the grocery store website for iconic images and text descriptions. MK performed all the experiments and wrote most of the text. All authors took part in discussing the results and contributed to writing the manuscript.

Paper B: Using Variational Multi-View Learning for Classification of Grocery Items

Marcus Klasson, Cheng Zhang, Hedvig Kjellström. In *Patterns, Volume 1(8)* (2020).

Summary: We investigate whether training image classifiers benefits from combining the natural images and web-scraped information in the Grocery Store dataset from Paper A. We employ a deep generative model called Variational Canonical Correlation Analysis (VCCA) [72] for learning joint representations of the different data views, i.e., natural images, iconic images, and text descriptions, that are used for training the classifiers. We performed a thorough ablation study over all data views to demonstrate how they contribute individually to enhancing the classification performance. Our results show that the iconic images help to group the items based on their color and shape while text descriptions separate the items based on differences in ingredients and flavor. Figure 1.1 shows visualizations of the latent representations learned by a variational autoencoder (VAE) [73] and VCCA. Finally, we concluded that utilizing the iconic images and text descriptions yielded better classification results than only using natural images.

Author Contributions: CZ and HK presented the idea and all authors contributed to formalizing the methodology. MK performed all the experiments, created the visualizations, and wrote most of the text. All authors took part in discussing the results and contributed to writing the manuscript.

Paper C: Learn the Time to Learn: Replay Scheduling for Continual Learning

Marcus Klasson, Hedvig Kjellström, Cheng Zhang. *Under submission.*

Summary: We present a slightly altered CL setting where the historical data is accessible at any time to mitigate catastrophic forgetting. However, retraining the CL system from all historical data is prohibitive due to constraints on processing times. To this end, we propose to learn the time to learn for a CL system, in which we learn schedules over which tasks to replay at different times. Figure 1.2 illustrates the differences between our replay scheduling approach and the traditional replay approach in CL. We study the benefits of replay scheduling by allowing multiple episodes in the CL environment to enable searches for efficient replay schedules using Monte Carlo tree search. We show that the replay schedules from MCTS can

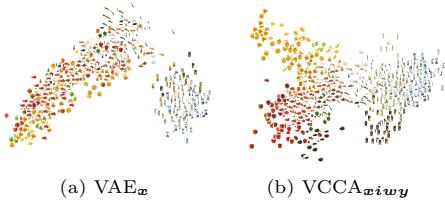


Figure 1.1: Visualizations of the latent representations projected in 2D space with PCA from models VAE_x and VCCA_{xiwy} , where we plot the corresponding iconic image for each latent representation. We observe that VCCA_{xiwy} structures the items based on visual similarities by incorporating the web-scraped information into the learning.

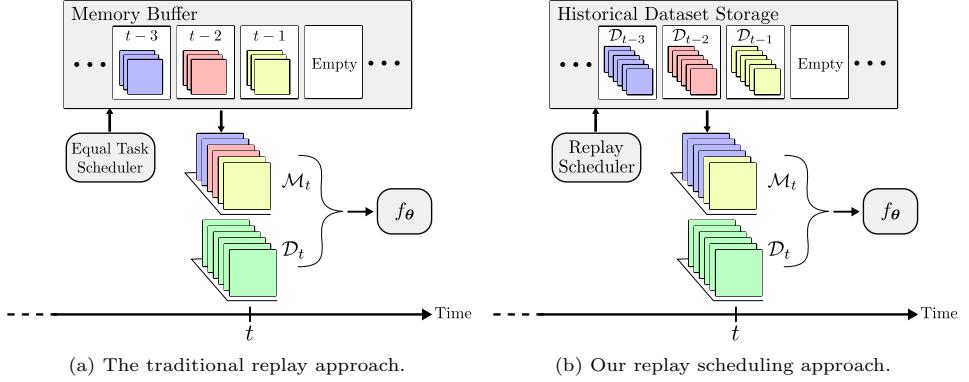


Figure 1.2: Illustrations of the (a) traditional replay approach where an equal amount of old samples in a small memory buffer are replayed at every time, and (b) our replay scheduling approach from Paper C where the scheduler selects which tasks to replay by selecting samples from the stored previously seen datasets.

significantly improve the performance compared to baselines with scheduling policies. We conclude that learning the time to replay different tasks can be critical for the performance of replay-based methods in our new proposed CL setting.

Author Contributions: CZ presented the idea, and MK and CZ contributed to formalizing the methodology. MK performed all the experiments, created the visualizations, and wrote the text. All authors took part in discussing the results and contributed to writing the manuscript.

Paper D: Policy Learning for Replay Scheduling in Continual Learning

Marcus Klasson, Hedvig Kjellström, Cheng Zhang. *Under preparation.*

Summary: To enable replay scheduling in real-world CL scenarios, we propose a reinforcement learning-based framework for learning replay scheduling policies. Figure 1.3 illustrates the pipeline of our framework. Our intuition is that there may exist general patterns regarding the replay scheduling, e.g., tasks that are harder or have been forgotten should be replayed more often. Therefore, we define the state as the validation performance of the classifier, which is passed to the policy that selects actions for how to compose the replay memory to efficiently mitigate catastrophic forgetting. Our framework enables learning policies that can generalize to new CL scenarios without additional computational cost or training in the test environment. In the experiments, we show that the learned policies can generalize to CL scenarios with new task orders and datasets unseen during training.

Author Contributions: CZ presented the idea, and MK and CZ contributed to formalizing the methodology. MK performed all the experiments, created the visualizations, and wrote most of the text. All authors took part in discussing the results and contributed to writing the manuscript.

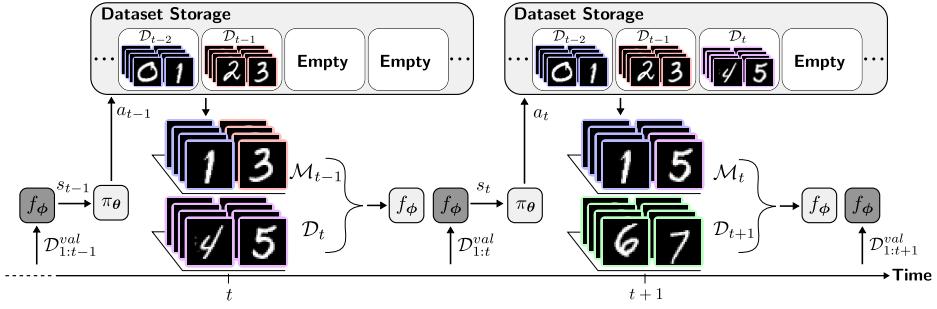


Figure 1.3: Illustration over our policy learning framework for replay scheduling proposed in Paper D. The policy π_θ receives state s_t of the classifier f_ϕ to select action a_t determining how the replay memory \mathcal{M}_t should be constructed. Dark gray-shaded box means that classifier is in evaluation mode.

1.5 Thesis Outline

The rest of the thesis is organized as follows: Chapter 2 provides some preliminaries on deep learning methods in computer vision that we used in this thesis. Chapter 3 is focused on our contributions in FGIR where we summarize Paper A and B on grocery item classification with an assistive vision app. Similarly, in Chapter 4, we focus on our contributions in CL and summarize the content in Paper C and D on our proposed replay scheduling approach in CL. Finally, in Chapter 5, we provide our conclusions of the presented works and present some potential research directions that we believe could be interesting to look deeper into in future research on computer vision-based assistive technologies.

Chapter 2

Background

In this chapter, we provide an overview of the preliminaries for comprehending the methodology in the included papers in this thesis. We assume that the reader has some basic knowledge in calculus, linear algebra, and probability theory. Section 2.1 gives an introduction to machine learning, including different problem settings, some terminology, as well as the notation that will be used throughout this thesis. In Section 2.2, we give an overview of deep learning [74], where we cover the different types of network architectures that have been used in the thesis, as well as a brief introduction to deep reinforcement learning.

2.1 Problem Settings in Machine Learning

Machine learning tasks can be divided into three main fields, namely, *supervised learning*, *unsupervised learning*, and *reinforcement learning* (RL). In this section, we will give a brief introduction to these topics to provide some context on the challenges that we are addressing in this thesis.

We will begin by introducing some algebraic notation that will be used throughout this thesis. For representing data, we use the vector \mathbf{x} which is a column vector

$$\mathbf{x} = [x_1, \dots, x_D]^T,$$

where x_i for $1 \leq i \leq D$ is the i -th scalar, real-valued feature of the data vector, and the superscript T transposes the column vector into a row vector. In some cases, each data vector \mathbf{x} have an assigned class label y which is an integer number. Scalars can hence both represent real values and integers. The data vectors \mathbf{x} will often represent a 2-dimensional image, where each feature x_i is a pixel, in this thesis, but we still represent images as column vectors for simplicity.

We assume that the data follows some underlying distribution, $p_{data}(\mathbf{x})$, usually referred to as the data generating distribution. In machine learning, we are often interested in approximating p_{data} with a distribution $p_{\theta}(\mathbf{x})$ parameterized by θ given a finite dataset $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ of N samples where $\mathbf{x}^{(i)} \sim p_{data}(\mathbf{x})$. Estimating the parameters θ from the dataset \mathcal{D} is called the *training phase*. A central goal for most applications in machine learning is to make predictions on new data which is called *generalization*. We measure the generalization capability by evaluating the task of interest on a held-out dataset during the *test phase*.

Supervised Learning. In this setting, we are given a dataset $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ where each data point $\mathbf{x}^{(i)}$ is accompanied with a target $y^{(i)}$. In classification problems, the target $y \in \{1, \dots, C\}$ belongs to one of C discrete categories, while in regression problems, the targets $y \in \mathbb{R}$ are continuous and real-valued. The goal is to estimate $p_{\theta}(y|\mathbf{x})$ which is the probability of assigning the target y given data \mathbf{x} . We represent $p_{\theta}(y|\mathbf{x}) = f_{\theta}(x)$ with the parameterized function $f_{\theta}(x)$ which maps data \mathbf{x} to target variables y .

Unsupervised Learning. The goal in the unsupervised setting is to reveal hidden structures in the given dataset $\{\mathbf{x}^{(i)}\}_{i=1}^N$ without access to target variables. For example, we might be interested in estimating the data generating distribution p_{data} with p_{θ} (as mentioned above), discovering groups of similar data points using clustering techniques, or project high-dimensional data into two or three dimensions for visualization purposes.

Reinforcement Learning. For these problems, we have a *learning agent* that performs actions based on perceived states of the environment to reach some specific goal represented by a reward signal r . An example of such a task Mountain Car problem [75] where the agent tries to drive a car over the top of a hill. The action selection is modeled by the policy $\pi_{\theta}(a|s)$ which is a mapping from states s to actions a . The goal for the agent is to learn a policy that maximizes the reward within a specified time limit.

2.2 Deep Learning

In this section, we give an introduction to deep learning [74] which is a family of machine learning methods that we have used in this thesis. Deep learning methods are capable of learning tasks from large and high-dimensional datasets due to more efficient machine learning algorithms and advances in computer hardware in the last decades. A deep learning architecture is constructed by stacking layers of parameters that extract intermediate representations of the input data, where the last layer usually outputs a predicted target answer from the queried input. This approach has been successfully applied in various number of fields such as computer vision [25, 26], natural language processing [76], and reinforcement learning [77, 78].

Deep networks are most commonly optimized using the Stochastic Gradient Descent (SGD) algorithm. This procedure uses gradients of a loss function $\mathcal{L}(f_{\theta}(\mathbf{x}), \mathbf{y})$ with respect to every parameter to perform the update step

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}), \quad (2.1)$$

where η is the learning rate which determines how much the weights should be updated with the computed gradient. Common loss functions \mathcal{L} are the cross-entropy loss for classification and the mean-squared error for regression tasks. The gradients are computed using the backpropagation algorithm [79] which computes the gradients iteratively backwards one layer at a time. To improve the computational efficiency, we compute the average of the gradients from a randomly sampled mini-batch of data from the dataset during the parameter updates.

Many innovations have been designed in recent years to stabilize the training of deep neural network architectures. For instance, dropout [80] is used for miti-

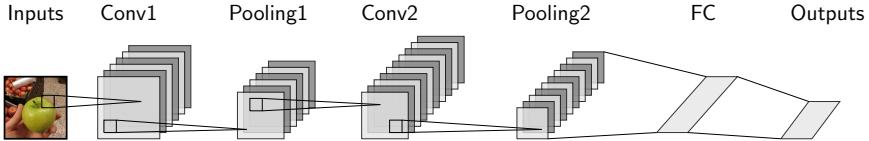


Figure 2.2: Simple CNN architecture with two convolutional layers (Conv) for extracting feature maps, two subsampling layers performing a pooling operation (Pooling), and one fully-connected layer (FC) to produce the network output.

gating overfitting to the training data by randomly keeping layer outputs active to encourage the parameters to take more (or less) responsibility for the inputs. Batch normalization [81, 82] normalizes the layer inputs to enable faster and more stable learning. Residual connections [26] eases training networks with hundreds of layers by mitigating the degradation problem, where adding more layers resulted in higher training error. Furthermore, various adaptive learning rate optimizers [83–85] for the SGD updates, learning rate schedulers [86], and network initializations [87, 88] have been proposed to enable more efficient training. We refer the reader to the provided references for the details on these techniques. Next, we will describe three popular types of neural networks that have been applied in the experiments of this thesis, namely, multilayer perceptrons (MLPs), convolutional neural networks (CNNs), and recurrent neural networks (RNNs).

Multilayer Perceptrons. The simplest form of feedforward neural networks is the MLP based on the Perceptron [89]. Basically, MLPs constitutes of multiple hidden layers performing matrix multiplication to extract intermediate representations of the input. Each matrix multiplication is followed by an activation function, such as the Rectified Linear Unit (ReLU) activation, which is an essential step for enabling neural networks to learn non-linear functions. Figure 2.1 shows an illustration of an MLP with two hidden layers. The edges between two columns of units represent the matrix multiplication where each edge corresponds to a parameter in the network. The hidden layers are often called fully-connected layers in MLPs.

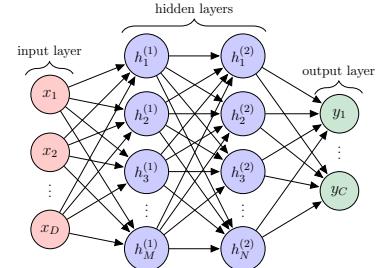


Figure 2.1: Illustration of MLP with two hidden layers.

Convolutional Neural Networks. Convolutional layers constitutes of a set of filters with adaptable parameters to capture relationships between neighboring features in the data [90]. The parameters in every filter are shared across every spatial location in the data when extracting the succeeding feature representation. This design choice enables each filter to extract important features at any location in the data. Furthermore, parameter sharing also reduces the amount of parameters in the network and lowers the number of operations for computing the outputs. Figure 2.2 shows an illustration of a CNN architecture. Each filter slides across the width and height of the input to compute dot products between the filter weights and the local input region to produce a 2-dimensional feature map. The feature maps from all filters are then stacked depth-wise to obtain the output volume, which are often

subsampled along their spatial dimensions with a pooling operation after a non-linear activation function. The feature maps in the last subsampling layer are flattened and fed into an arbitrary number of fully-connected layers (an MLP) to produce the final output.

Recurrent Neural Networks. RNNs are a family of neural network models specialized for processing sequences of data. These networks predicts the next output given the past sequence and uses an extension of backpropagation to compute gradients across different time steps. Similar to CNNs, RNNs apply parameter sharing across time steps to enable extracting relevant information that can appear at different positions in the sequence. RNNs makes an assumption that all relevant information about the past sequence $\mathbf{x}_{1:t-1}$ can be summarized in a hidden state \mathbf{h}_t . The state evolves over time through the updated equation $\mathbf{h}_t = f_{\theta}(\mathbf{h}_{t-1}, \mathbf{x}_{t-1})$ where f_{θ} is a neural network for capturing the long-term dependencies. Common choices for this network are the Long Short-Term Memory (LSTM) [91] and Gated Recurrent Unit (GRU) [92] which use gating functions for determining what information is relevant for the future time steps. Figure 2.3 shows a graphical representation of an RNN. The initial hidden state \mathbf{h}_0 is usually set to zero or learned.

Attention-based network architectures, such as Transformers [52, 93] and Graph Neural Networks [94, 95], are not covered in this thesis, so we refer the reader to the provided references for information on these networks.

Deep Autoencoders

A common network architecture for learning latent representations of high-dimensional data in unsupervised fashions with deep learning are autoencoders. These architectures are typically built using two networks called decoder f_{θ} and encoder g_{ϕ} with a bottleneck layer between the networks that extracts a low-dimensional representation \mathbf{z} of the data \mathbf{x} . The idea is that the learned representation \mathbf{z} should contain enough relevant information for reconstructing the input data. Figure 2.4 shows an illustration of an autoencoder network. The encoder extracts the latent representation with $\mathbf{z} = g_{\phi}(\mathbf{x})$, while the decoder aims to reconstruct the original input, such that $\hat{\mathbf{x}} = f_{\theta}(\mathbf{z})$ where $\hat{\mathbf{x}} \approx \mathbf{x}$. There exist various autoencoder approaches for improving the quality of the learned representations. For example, denoising autoencoders [96] adds noise to the input data to avoid overfitting, and sparse autoencoders [97] induces adds sparsity constraints on the activation functions to improve robustness.

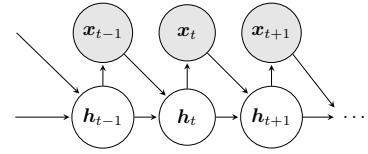


Figure 2.3: Graphical representation of RNN.

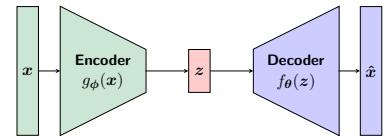


Figure 2.4: Illustration of an autoencoder network.

Multimodal Autoencoders. Autoencoders can easily be extended to handling multiple types of data. The goal here is usually to learn one joint latent repre-

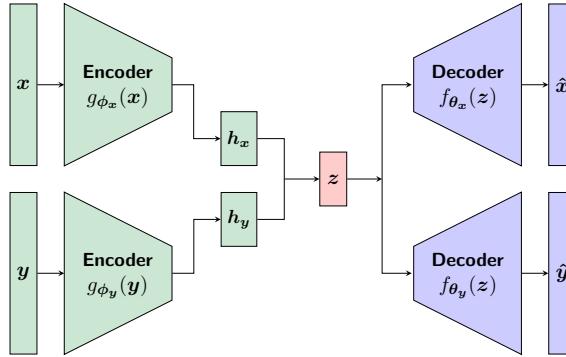


Figure 2.5: Illustration of multimodal autoencoder network for learning joint representation \mathbf{z} from the different but related data types \mathbf{x} and \mathbf{y} .

sentation for capturing correspondences between the data types [98]. Consider the multimodal dataset $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ where \mathbf{x} and \mathbf{y} originate from two different data views (real-world images, web images) or modalities (visual signals, text descriptions) that share some high-level information. For example, \mathbf{x} could be an image of a living room and \mathbf{y} is a text sentence describing the appearance of the room, where objects are located etc. Figure 2.5 shows an illustration of the network architecture. The input \mathbf{x} and \mathbf{y} have two separate encoders, $g_{\phi_x}(\mathbf{x})$ and $g_{\phi_y}(\mathbf{y})$, that extracts the intermediate latent representations \mathbf{h}_x and \mathbf{h}_y respectively. The intermediate representations are then combined into the joint latent representation \mathbf{z} , often through some element-wise operation (addition, multiplication) or concatenation. The joint multimodal representation \mathbf{z} is then passed through two separate decoder networks, $f_{\theta_x}(\mathbf{z})$ and $f_{\theta_y}(\mathbf{z})$, for predicting the original input data separately. Combining different modalities, such as image and video, natural language, and audio signals, with autoencoders for learning more rich representations of data has been studied frequently the last decade [99–102] [Add REFs]. One advantage is that the network architecture can be trained in an end-to-end fashion to be used for both learning representations and making predictions of the used data types. However, a major challenge is how to handle scenarios where some data types could be missing. An option here is to learn a joint representation by using a single encoder for the data type that is available at both training and test phases that is used for predicting each data type [100].

Variational Autoencoders

In this section, we give an overview of the variational autoencoder (VAE) [73, 103] that we apply for representation learning in Paper A and B. The VAE is a deep generative model that enable learns representations of data by combining ideas from deep learning and probabilistic graphical models [104]. As commonly done in unsupervised learning, VAEs aim to learn a parametric distribution $p_{\theta}(\mathbf{x})$ that approximates the data distribution p_{data} given a dataset $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$ where $\mathbf{x}_i \sim p_{data}$. We frame the problem of learning $p_{\theta}(\mathbf{x})$ by expressing the distribution using latent variables

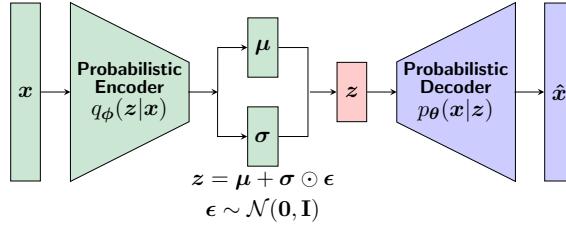


Figure 2.6: Illustration of a VAE.

z , such that

$$p_{\theta}(x) = \int p_{\theta}(x, z) dz = \int p_{\theta}(x|z)p(z) dz, \quad (2.2)$$

where the joint distribution $p_{\theta}(x, z)$ is the model which factorizes into the likelihood $p_{\theta}(x|z)$ and the prior distribution $p(z)$ for the latent variables. The latent variables z are introduced for representing underlying structures and relationships in x that can be difficult to observe from the data itself. We can view z as containing meaningful information for the data generation process by assuming that the data was generated with the following steps:

1. sample the latent vector $z \sim p(z)$ from the prior.
2. generate data point $x \sim p_{\theta}(x|z)$ given the sampled z .

The Evidence Lower Bound. The problem with learning the parameters θ that maximizes $p_{\theta}(x)$ is that the integral in Equation 2.2 is almost always intractable due to the immense amount of values for z that needs to be evaluated. Instead, we can compute $p_{\theta}(x)$ using Bayes' rule

$$p_{\theta}(x) = \frac{p_{\theta}(x|z)p(z)}{p_{\theta}(z|x)}, \quad (2.3)$$

where the posterior distribution $p_{\theta}(z|x)$ can facilitate faster learning by inferring the latents z from data x . However, we need to resort to approximating the posterior using variational inference [105, 106] since the posterior needs access to $p_{\theta}(x)$ to be computed. We then define an approximate posterior distribution $q_{\phi}(z|x)$ parameterized by ϕ to infer the latents z . This allows us to a form lower bound on the marginal log-likelihood $\log p_{\theta}(x)$ given by

$$\log p_{\theta}(x) \geq \mathbb{E}_{z \sim q_{\phi}(z|x)}[\log p_{\theta}(x|z)] - D_{\text{KL}}[q_{\phi}(z|x) \parallel p(z)], \quad (2.4)$$

which is called the evidence lower bound (ELBO). The expected value over $\log p_{\theta}(x|z)$ can be estimated with Monte Carlo sampling using $q_{\phi}(z|x)$, and the Kullback-Leibler (KL) divergence between $q_{\phi}(z|x)$ and $p(z)$ can be computed analytically depending on how we choose these distributions. Commonly, the approximate posterior is selected to be Gaussian distribution $q_{\phi}(z|x) = \mathcal{N}(z; \mu_{\phi}(x), \text{diag}(\sigma_{\phi}(x)))$ where the parameters ϕ are used for computing the mean $\mu_{\phi}(x)$ and standard deviation $\sigma_{\phi}(x)$ of the approximate posterior given x , and $\text{diag}(\cdot)$ is for constraining the covariance matrix to be a diagonal matrix. The prior distribution is also selected to be a zero-mean, standard Gaussian distribution $p(z) = \mathcal{N}(z; \mathbf{0}, \mathbf{I})$ where \mathbf{I} is the identity matrix.

Network Architecture. The VAE architecture resembles a standard autoencoder where the encoder and decoder networks are probabilistic in the sense that they output parameters of probability distributions. In contrast to autoencoders that learn latent representations \mathbf{z} that minimizes the reconstruction error on the data, VAEs learns an approximate posterior $q_\phi(\mathbf{z}|\mathbf{x})$ that maximizes the likelihood of $p_\theta(\mathbf{x})$. Similar to autoencoders, we can extract latent representations of the data with the probabilistic encoder. However, VAEs can also estimate the likelihood of unseen data points under the learned model, as well as generate new data by decoding latents sampled from the prior. Figure 2.6 shows an illustration of a VAE architecture. The encoder represents the approximate posterior $q_\phi(\mathbf{z}|\mathbf{x})$ which outputs the mean $\mu_\phi(\mathbf{x})$ and standard deviation $\sigma_\phi(\mathbf{x})$ of a Gaussian distribution over \mathbf{z} , which is the standard choice for q_ϕ [73]. The latent vector is sampled using the "reparametrization trick" [73, 107] by computing

$$\mathbf{z} = \mu_\phi(\mathbf{x}) + \sigma_\phi(\mathbf{x}) \odot \epsilon, \quad \text{where } \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (2.5)$$

where \odot denotes element-wise multiplication. The reparametrization enables computing gradients of the expectation $\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]$ in Equation 2.4 with the backpropagation algorithm. The decoder represents the likelihood $p_\theta(\mathbf{x}|\mathbf{z})$ that can be used to generate data \mathbf{x} from sampled latents \mathbf{z} , where the distributional form depends on the data type of \mathbf{x} .

VAEs have been used in various applications for modeling images, text, and audio data [Add REFS]. Furthermore, they have been extended to modeling data from different modalities due to their capability of handling missing modalities to enable cross-modal data generation [72, 108–111]. In Paper A and B, we employ Variational Canonical Correlation Analysis (VCCA) [72] for learning joint representations of natural images and web-scraped information of grocery items to facilitate learning image classifiers.

Deep Reinforcement Learning

In this section, we give a brief overview RL to provide some preliminaries for Paper D. The RL setup considers a learning agent that interacts with an environment E over a number of discrete time steps [75]. The environment is modeled with a Markov Decision Process (MDP) [112] represented as a tuple $E = (\mathcal{S}, \mathcal{A}, P, R, \mu, \gamma)$ consisting of the state space \mathcal{S} , action space \mathcal{A} , state transition probability $P(s'|s, a)$, reward function $R(s, a)$, initial state distribution $\mu(s_1)$, and discount factor γ . At each time step t , the agent receives a state s_t from the environment, selects an action $a_t \in \mathcal{A}$ using a policy $\pi(a|s)$, and enters the next state s_{t+1} with transition probability $P(s_{t+1}|s_t, a_t)$ and receives a numerical reward following r_t from the environment. This procedure is repeated until the agent reaches a terminal state in which the procedure can be restarted. The return $G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$ is the discounted accumulated reward from time step t . The goal for the agent is to learn a policy that maximizes the expected return.

The value function $V^\pi(s) = \mathbb{E}[G_t|s_t = s]$ defines the expected return following policy π from state s . Value functions are essential for learning as they can be used for comparing policies, such that $\pi \geq \pi'$ if and only if $V^\pi(s) \geq V^{\pi'}(s)$ [75]. An optimal policy π^* is defined as a policy that is better than or equal to all other

policies. There may exist several optimal policies π^* that share the same optimal value function given by $V^*(s) = \max_{\pi} V^{\pi}(s)$ for all states $s \in \mathcal{S}$. Similar to the value function, we also have the action value function $Q^{\pi}(s, a) = \mathbb{E}[G_t | s_t = s, a_t = a]$ which is the expected return for selecting action a in state s and following policy π . Optimal policies also share the same action value function given by $Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a)$ for all states $s \in \mathcal{S}$ and actions $a \in \mathcal{A}$. The optimal action value function Q^* can be written in terms of the optimal value function V^* as

$$Q^*(s, a) = \mathbb{E}[r_{t+1} + \gamma V^*(s') | s_t = s, a_t = a]. \quad (2.6)$$

This means that if we have access to Q^* , we can directly obtain the optimal action a^* in state s by using $a^* = \arg \max_a Q^*(s, a)$ to maximize the expected return.

Deep Q-Networks. The Deep Q-Network (DQN) [77, 113] is a value-based RL algorithm which aims to approximate the optimal action value function as $Q^*(s, a) \approx Q_{\theta}(s, a)$ with a neural network Q_{θ} parameterized by θ . For learning the function Q_{θ} , we collect data from the environment with by using an epsilon-greedy policy, i.e., $a_t = \arg \max_a Q_{\theta}(s_t, a)$, in state s_t . For estimating Q_{θ} , we minimize the following loss based on the Bellman equation:

$$\mathcal{L}_{\text{DQN}}(\theta) = (r_t + \gamma \max_{a'} Q_{\theta^-}(s_{t+1}, a') - Q_{\theta}(s_t, a_t))^2, \quad (2.7)$$

where θ^- is a previous copy of the network θ referred to as the target network. In addition to the introduction of target networks, several methods have been proposed for stabilizing the learning process, such as using experience replay [114] to sample training data stored in a replay buffer, applying L1-smoothing to the loss, correcting the action value estimates [115], and improving generalization across actions [116].

Policy Gradient Methods. An alternative to value-based RL is using policy gradient methods for learning the policy. Here, the optimal policy is estimated directly with a parameterized form $\pi_{\theta}(a|s)$ where θ represents the parameters of a neural network. The policy network takes the state s_t as input and outputs a distribution over the possible actions a_t . Policy gradient methods often use an actor-critic approach [75, 117], where the policy is known as the actor that selects actions, and the critic evaluates the actions made by the actor by estimating the value function $V_{\lambda}(s_t)$ parameterized by λ . For updating the policy, the agent collects experiences from the environment with the current policy and aims to minimize the loss

$$\mathcal{L}_{\text{PG}}(\theta) = \mathbb{E}_{s_t \sim E, a_t \sim \pi_{\theta}} [\log \pi_{\theta}(a_t | s_t) \hat{A}(s_t, a_t)], \quad (2.8)$$

where $\hat{A}(s, a)$ is an estimate of the advantage function, which describes how much better or worse a taken action is on average, and is given by

$$\hat{A}(s_t, a_t) = \sum_{i=0}^{k-1} \gamma^i r_i + \gamma^k V_{\lambda}(s_{t+k}) - V_{\lambda}(s_t), \quad (2.9)$$

where k can vary depending on if we reached a terminal state from s_t and is upper-bounded by the number of steps between updates t_{\max} [118]. A popular actor-critic method is Proximal Policy Optimization [119] that incorporates constraints on the policy updates, as well as mini-batches and epochs of learning between each interaction with the environment.

Chapter 3

Fine-Grained Image Recognition

In this chapter, we provide an overview for Paper A and B that focuses on fine-grained image recognition (FGIR) [36] of grocery items. Object recognition in grocery shopping scenarios is a challenging task for visually impaired (VI) people [22, 48, 120] due to the necessity of visual capabilities to navigate and identify the products. Regarding identifying products, groceries are often structured in fine-grained manners where details on appearance and shape help us to distinguish between different items. However, since distinguishing between packaged items often requires reading text and raw food items can have similar shapes, recognizing fine-grained details can be difficult without the ability to see. In recent years, there has been an emerging interest of development of computer vision-based assistive technologies for assisting VIs with grocery shopping [12, 57–62]. Nevertheless, fine-grained recognition of groceries still remain challenging due to the need for real-world training data for such systems to be robust in grocery stores. We aimed to enhance the robustness and predictive performance of image classifiers by combining natural images with external web information with the items. In Paper A, we have collected a dataset with natural images of groceries with a mobile phone camera, and web-scraped images and text descriptions about the items, for evaluating image classifiers in real-world supermarket scenarios. We perform an ablation study over the different data types in the data set in Paper B, where we show that utilizing all data types improves the fine-grained classification performance over image classifiers trained with only the natural images. We cover some related work (Section 3.1), describe the dataset collection process (Section 3.2) and our used approach for utilizing multi-view data (Section 3.3), as well as summarize the experimental results from the ablation study (Section 3.4).

3.1 Related Work

In this section, we will briefly discuss the related work on FGIR, datasets for FGIR, and multi-view learning.

Fine-Grained Image Recognition. The goal with FGIR [36] is to distinguish between a set of visually similar object classes that belong to some super-class. The main challenge is to recognize the fine-grained visual details that are required for

discriminating between visually similar objects. In recent years, the successes with deep learning in computer vision have encouraged researchers to explore various approaches for FGIR in different domains with the goal of discriminating between sub-classes of different animals [121], cars [122], fruits [123], and retail products [124]. These approaches can broadly be divided into three directions where the deep networks utilize (i) localization-classification subnetworks, (ii) end-to-end feature encodings, and (iii) external information. In (i), the goal is to find parts that are shared across the sub-classes of some general object and then discover how these localized parts differ in appearance. This can be achieved by utilizing feature maps from convolutional layers as local descriptors [125–127], employing detection and segmentation techniques for localizing object parts [128–130], or leveraging attention mechanisms in scenarios where common object parts are difficult to represent [131–133]. With (ii), the goal has been to learn features that are better at capturing local differences by computing second-order feature interactions [134, 135], as well as designing novel loss functions [136–138]. In (iii), data from external domains, such as the web or different modalities (e.g., text or audio), is utilized to influence the learning of fine-grained details in the objects. These additional data types can be relatively cheap to collect and serves as extra guidance for learning useful representations to the FGIR task. The goal with learning from noisy web data is to reduce the need for human-annotated data [139–141], while learning from multimodal data aims to obtain a joint representation that captures the correspondences between the images and the other modalities [142–144]. In addition to FGIR, both web-scraped and multimodal data has been used for zero-shot learning to transfer knowledge from densely annotated categories to learn new fine-grained classes [145–147]. Our work in Paper A and B is mostly related to (iii) FGIR with external information, where we use web-scraped images and text descriptions in combination with natural images to improve the classification performance of fine-grained grocery recognition.

Datasets for Fine-Grained Image Recognition There exist many image datasets for benchmarking computer vision models [15, 121, 148–153]. These datasets ranges from large-scale datasets annotated by groups of classes, such as [148, 154], to smaller but densely-annotated datasets [152, 155–157]. There also exist domain-specific datasets for FGIR tasks of animals and species [121, 158–162], fashion [163, 164], airplanes and cars [165–169], faces [170–173], and foods [123, 174]. Several of the mentioned datasets have also been extended to benchmarking in zero/few-shot image classification [145, 175–177]. In computer vision for visually impaired people, we have recently seen an emergence of datasets collected by people who are blind/low-vision [13–15, 59, 178]. Datasets for FGIR of grocery items in their natural environments, such as grocery stores, shelves, and kitchens, have been addressed in plenty of previous works [124, 179–183]. Similarly, there exist other kinds of food datasets with images of various food dishes [174, 184–186], cooking videos [187, 188], recipes [189–191], and also restaurant-oriented information [192, 193]. Our grocery item dataset (see Section 3.2) shares several similarities with the mentioned works above. For instance, all images of raw and packaged groceries are taken in their natural environment, images are taken with a mobile phone camera from an egocentric view-point, grocery classes are hierarchically labeled, and each class have an additional web-scraped iconic image and text description.

Multi-view Learning. Machine learning with multiple types of data is an active research field [98, 194]. In particular, the intersection of computer vision and natural language processing has been of great interest and has led to applications such as image captioning [28, 195, 196], visual question answering [54, 56, 197]. In this thesis, we have focused on learning joint representations across different data types by applying the subspace learning approach from multi-view learning [194]. A view can be defined as data that has been recorded from a specific sensor [198]. For example, we consider natural and web images to be from different views as the images have been recorded using different cameras. Subspace learning approaches aims to learn joint representations from multiple views by assuming that each view have been generated from a latent space that is shared among the views. Most methods originate from Canonical Correlation Analysis [199] (CCA) where the goal is to linearly project pairs of different views into a lower-dimensional space where the correlation between the projections is maximized. There exist various extensions of CCA which uses deep neural networks for learning nonlinear mappings to extract more rich features of the views, such as Deep CCA [200] and Deep Canonically Correlated Autoencoders [201]. Inspired from the Deep CCA variants, Variational CCA [72] (VCCA) is a multi-view deep generative model that can both learn joint representations and generate new data by sampling from the shared latent space. An advantage of VCCA is the straight-forward extension to modeling both a shared latent space and private latent spaces which capture the view-specific variations to focus the learning of shared variations into the shared latent space [72, 198, 202–205]. In Paper B, we employ VCCA, and its extention to private latent spaces, for learning joint representations of the different data views in our grocery item dataset to obtain more robust image classifiers.

3.2 Dataset Collection

In this section, we describe our procedure for collecting the image dataset of grocery items in Paper A. We visited several supermarkets in Stockholm, Sweden, to collect natural images of groceries with a mobile phone camera to imitate a shopping scenario using an assistive vision device. The collected images captures situations that can be challenging for assistive vision devices, such as misplaced items, hand occlusions, multiple instances and classes present, and various lighting conditions (see Figure 3.1). All images were taken with a single targeted item in mind, such that each image is paired with a single label. For items which belong to a clear super-class, e.g., apples and milk packages, we also provided the general class of the items to establish a hierarchical labeling structure of the data. Furthermore, we complemented the image dataset with web-scraped information of each grocery item. More specifically, we visited the online shopping website of a Swedish supermarket chain and downloaded 1) an iconic image of the item on a white background, 2) a text description that describes the flavor and ingredients of the item,



Figure 3.1: Example images of challenging scenarios.

Table 3.1: Examples of grocery item classes in the Grocery Store dataset. We display four different items (coarse-grained class in parenthesis), followed by two natural images taken with a mobile phone inside grocery stores. Next comes the web-scraped information of the items consisting of an iconic image and a text description. We have highlighted ingredients and flavors in the text description that are characteristic for the specific item.

Class Labels	Natural Images	Iconic Images	Text Descriptions
Granny Smith (Apple)	 		<i>“...green apple with white, firm pulp and a clear acidity in the flavor.”</i>
Royal Gala (Apple)	 		<i>“...crispy and very juicy apple, with yellow-white pulp. The peel is thin with a red yellow speckled color.”</i>
Tropicana Mandarin (Juice)	 		<i>“...is a ready to drink juice without pulp pressed on orange, mandarin and grapes. Not from concentrate. Mildly pasteurized.”</i>
Yoggi Vanilla (Yoghurt)	 		<i>“...creamy vanilla yoghurt original... added sugar than regular flavored yoghurt. Great for both breakfast and snacks.”</i>

and 3) the nutrition values for items where it was applicable. We show four examples of grocery items and their web-scraped information in Table 3.1. Since these data types are on a class-based level, we can use the web-scraped information as weak supervision to guide the classifier to learn fine-grained details that helps discriminating between visually similar items. We have established training, validation, and test sets of the natural images where the split was made according to the grocery store location to avoid mixing images from different stores. Finally, we made the dataset publicly available under the MIT license¹.

3.3 Fine-Grained Image Recognition of Grocery Items

In this section, we describe the problem setting of the grocery item classification task from the dataset collected in Paper A. Moreover, we describe our approach from Paper B for learning joint representations of the available data views to improve the classification performance.

Problem Setting

We focus on the application of training an image recognition app in mobile phones for assisting VIs with grocery shopping. Training such image classifiers typically requires immense amounts of annotated images of the groceries. Additionally, the natural images should be collected in scenarios where the app would be used, such as grocery stores or home environments, which further complicates the collection procedure. To

¹Grocery Store Dataset URL: <https://github.com/marcusklasson/GroceryStoreDataset>

reduce the need for collecting and annotating natural images from mobile phones, we aim to use web-scraped information about the groceries as additional supervision when training the image classifiers. The web-scraped information should help the classifier to focus on learning the learn fine-grained details of the items to improve the classification performance and robustness in real environments.

The following data views are available in the dataset for training the image classifiers:

- \mathbf{I} : Natural images of the grocery items taken inside grocery stores.
- \mathbf{x} : Feature vectors of the natural images \mathbf{I} extracted from some CNN f_φ parameterized by φ .
- \mathbf{y} : Class labels of the grocery items in the corresponding natural images.
- \mathbf{i} : Iconic images of the grocery items scraped from a supermarket website.
- \mathbf{w} : Text description of the grocery items scraped from the same supermarket website as \mathbf{i} .

See Table 3.1 for examples of the data views. The straightforward approach is to train a CNN with pairs of natural images \mathbf{I} and class labels \mathbf{y} in a supervised setting. However, the number of samples may be insufficient for training a CNN from scratch, which can lead to overfitting and low generalization capability to new images. One alternative is to employ transfer learning [206] where parameters from a CNN pre-trained on some large dataset, such as Imagenet [148], are used for initialization and we adapt the network to the grocery item recognition task by either 1) fine-tune some of the final layers, or 2) train a linear classifier from extracted features [207]. We take the latter approach and use extracted features for training to mitigate overfitting. Furthermore, we use the web-scraped data views for learning more rich representations of the grocery items that can be used for training more accurate and robust image classifiers, which we present in the next section.

Multi-View Representation Learning of Grocery Items

We describe the multi-view learning approach we applied for learning joint representations of grocery items to use for training the image classifiers. More specifically, we employed a deep latent variable model called Variational Canonical Correlation Analysis [72] (VCCA) for learning the joint representations. In the VCCA approach, we assume that the different data views have been generated from the same, shared latent space before they are observed in their original data type. The goal is to obtain joint representations by learning the shared latent space that captures the correspondences across the available views. The joint representations can then be used for training more accurate and robust image classifiers.

Figure 3.2 shows the network architecture of VCCA that we used in Paper B for learning the joint representations from all available data views. In this example, we have an input image \mathbf{I} of a juice package with the fine-grained class label \mathbf{y} *God Morgan Orange and Red Grapefruit Juice*. We obtain the latent representation \mathbf{z} by encoding the feature vector $\mathbf{x} = f_\varphi(\mathbf{I})$ of the input \mathbf{I} extracted from a CNN f_φ pretrained on Imagenet. In Paper B, we use a DenseNet [208] architecture as f_φ . We assume that the latent representation \mathbf{z} is shared by all available views to establish a more rich representation of the grocery item images by using the each available view as additional supervision to the natural images. We capture the correspondences

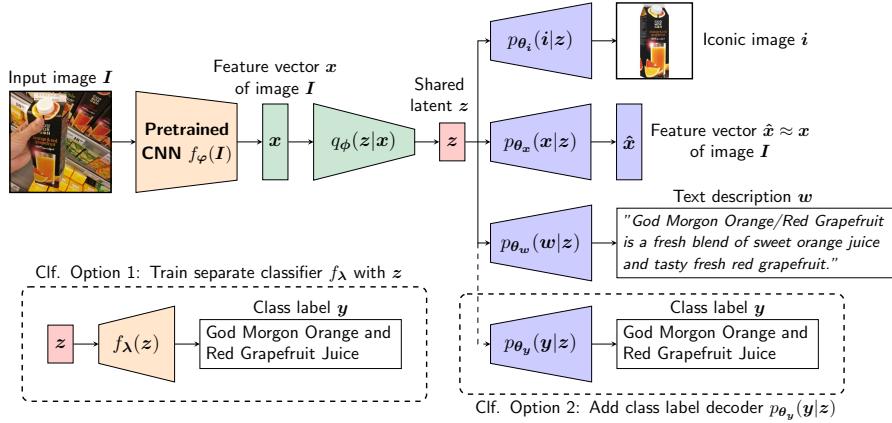


Figure 3.2: Illustration of the VCCA architecture used in Paper B for fine-grained classification.

between each view by predicting each view separately using the latent z as input to decoder networks that are specific for each view. The decoder $p_{\theta_x}(x|z)$ is an MLP similar to the encoder and outputs a reconstruction \hat{x} of the feature vector x . The decoder $p_{\theta_i}(i|z)$ is a CNN that predicts the iconic image i , while the decoder $p_{\theta_w}(w|z)$ is an RNN predicting the next word in the text description w . We have two options for training the image classifier. The first option is a two-step procedure where we 1) train the VCCA model from all available views except the class labels y , then 2) train a separate classifier $f_\lambda(z)$ on the learned latent representations z and the class labels y . The second option is to learn the classifier simultaneously as the latent space by adding a decoder $p_{\theta_y}(y|z)$ that predicts the class label y from z . Note that we only use x in the encoder part since the natural images is the only view that is available at test time when predicting the item class of new images taken in the grocery store. In the next section, we perform an ablation study over the available views to analyze how each view affects the fine-grained classification performance.

3.4 Experiments

In this section, we summarize the main experimental results from Paper B. The experiments involved an ablation study using VCCA with the possible combinations of available data views to investigate how each data view contributes to the classification performance. We present results on fine-grained classification performance, visualize the learned latent space to provide insights in how the web-scraped information affect the performance, and demonstrate how the iconic image decoder can be used for explaining the misclassifications by generating iconic images from natural images. We use subscripts to denote which data views that are utilized by VCCA. Adding y to the subscript means that the VCCA model uses a class label decoder $p_{\theta_y}(y|z)$ for classification, otherwise we have trained a separate linear classifier $f_\lambda(z)$ with softmax activation on the latent representation after training VCCA.

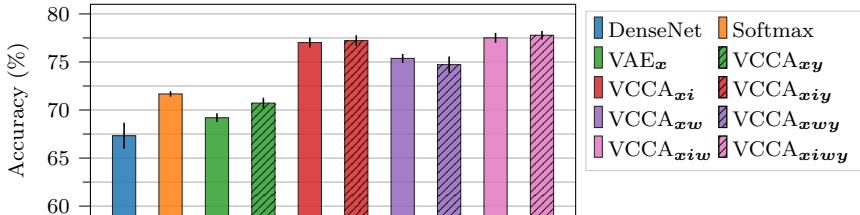


Figure 3.3: Fine-grained accuracies on the Grocery Store dataset for all classification methods. We show the means and standard deviations averaged over 5 seeds. Adding the iconic image i and text description w for learning joint representations with VCCA improves the classification performance over approaches that only utilize the natural images and class labels.

Fine-Grained Classification Results. Adding the web-scraped views in VCCA improves the classification performance over approaches that only utilize the natural images and class labels. We performed an ablation study over the available views that can be used by VCCA. The fine-grained classification performance was measured with a separate linear classifier trained on the learned latent representations and a class label decoder. We also compared the performance against training a DenseNet from scratch on the dataset (DenseNet) and a linear classifier with softmax activation trained on extracted features from a pretrained DenseNet (Softmax). Figure 3.3 shows a bar plot over the fine-grained accuracies achieved by all classification methods. For the methods only using natural images and class labels, i.e., DenseNet, Softmax, VAE_x , and VCCA_{xy} , we see that Softmax performs best which could be due to loss of information in VAE_x and VCCA_{xy} when the feature vectors are compressed into lower-dimensional latent representations. However, the performance of all VCCA models utilizing the iconic image i and/or the text description w outperforms Softmax significantly, which shows that both web-scraped views are useful for enhancing the fine-grained classification performance. Only utilizing i in VCCA performs on par with combining both i and w which could potentially be improved by processing the text descriptions with some word filtering or pretrained word embeddings [76, 209]. Finally, we observe that both classification options for VCCA performs similar, which shows that both options can be practical for classification using the learned latent representation.

Latent Space Visualizations. The web-scraped iconic images and text descriptions structures the grocery items based on view-specific similarities in the latent space that are beneficial for fine-grained classification. Figure 3.4 demonstrates how adding either the iconic image i or the text description w can change the structure of the latent space. We used PCA for projecting the learned latent representations of the test set into a 2-dimensional space. We have plotted the corresponding iconic image for the latent representations of VAE_x , VCCA_{xi} , VCCA_{xw} for visualization purposes in Figure 3.4(a-c). The VAE_x model in Figure 3.4(a) separates raw and packaged grocery items into two separate clusters in the latent space. When adding the iconic image in in Figure 3.4(b), we observe that the latent space becomes structured according to the color of the items. For VCCA_{xw} in in Figure 3.4(c), the latent space becomes structured according to the ingredients of the items, where we see in the raw grocery item cluster that bell peppers are placed in the upper region while apples are in the lower region.

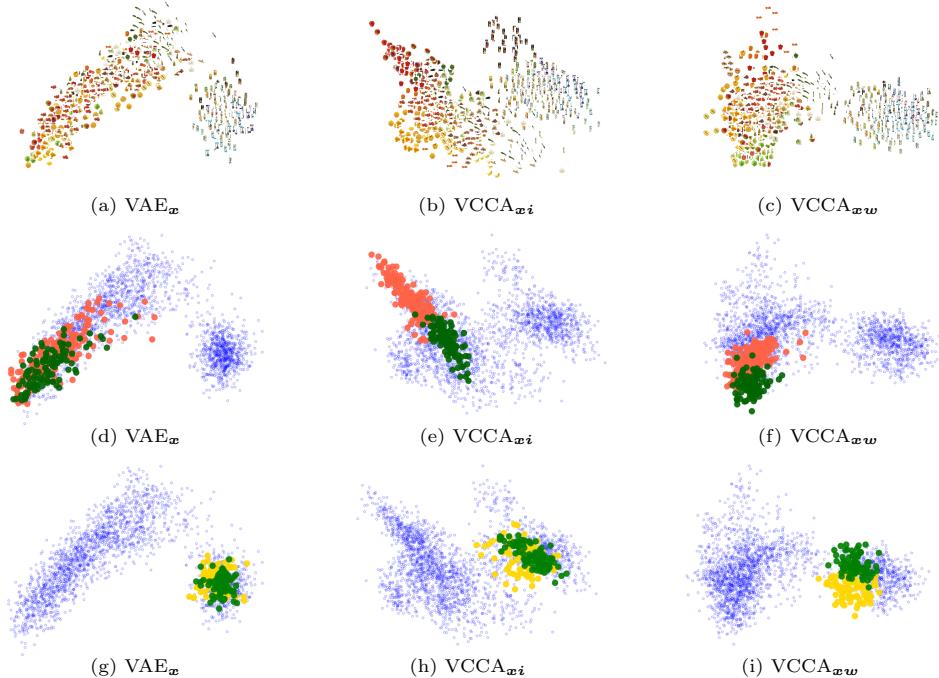


Figure 3.4: Visualizations of the latent representations from VAE_x , VCCA_{xi} , and VCCA_{xw} projected in 2-dimensional space with PCA. In (a-c), we show the latent representations plotted using the iconic images of the corresponding object class. In (d-f), we illustrate how the iconic images structures the items based on visual similarities by focusing on the green and red apple classes in the dataset plotted in their corresponding colors. Similarly, in (g-i), we show how the text descriptions structure items based on ingredients and flavor by focusing on visually similar yoghurt (green) and juice (yellow) packages. The blue dots correspond to all other grocery item classes.

Next, we focus on a certain set of classes to gain more insights in how the web-scraped views affect the latent space. In Figure 3.4(d-f), we inspect how items with different colors are separated in the latent space by utilizing i or w in VCCA by plotting the green and red apple classes as green and red points respectively. All other classes are plotted as smaller points in blue. We observe that both VCCA_{xi} and VCCA_{xw} manages to separate the apple classes better than VAE_x where VCCA_{xi} obtains the most clear separation. We also want to inspect how the views separates visually similar items with different ingredients and flavors. Hence, in Figure 3.4(g-i), we have plotted white-colored yoghurt and juice package classes in the colors green and yellow respectively. We see that VCCA_{xw} manages to separate these items better than VCCA_{xi} due to different ingredient information about yoghurt and juice that is mentioned in the text descriptions.

Iconic Image Generation. The iconic image decoder can be used for interpreting the outcomes of predicted classes. Figure 3.5 shows two examples of iconic images decoded from VCCA_{xiwy} using two natural images with class labels *Orange Bell Pepper* and *Anjou Pear*

Natural Image	Iconic Image	Decoded Image

Figure 3.5: Examples of decoded

from the test set. On the first row, we see that the model has generated an iconic image of an orange-green colored bell pepper with the shape of the green bell pepper in its iconic image. On the second row, we see that the decoded iconic image is a *Granny Smith* apple instead of the true class *Anjou Pear*. Hence, the iconic image decoder can be used as a tool for providing insights in why the classifier makes prediction errors.

3.5 Discussion

In this chapter, we summarized the contributions in Paper A and B on FGIR for assisting VI people with grocery shopping using an assistive vision app on their mobile phone. In Paper A, we collected a publicly available dataset with mobile phone images of groceries taken in grocery stores as well as web-scraped information about the items from a supermarket website. In Paper B, we show how the web-scraped information can be multi-view learning can be used for use the web-scraped information to train more accurate and robust image classifiers. More specifically, the iconic images structures the items after visual similarities such as colors and shapes in the learned latent space, while the text descriptions pushes items with similar ingredients and flavors closer to each other. This shows that the cheap-to-collect web-scraped views can serve as a good alternative for improving the fine-grained classification performance instead of collecting more natural images in the grocery stores.

Regarding the dataset collection in Paper A, we have some belated suggestions on how to make the collection procedure easier and also how to enrich the dataset. The images should be collected by recording videos of the grocery items rather than taking still images in the grocery stores. This would potentially enhance the dataset size and could also mitigate the class imbalance in the training and test splits. Furthermore, we could evaluate the classifiers on video sequences which would be a more realistic and user-friendly setting for how an assistive vision app would be used by VI people. Another suggestion to enrich the dataset would be to download more instances of iconic images and text descriptions from supermarket websites, which could potentially allow the VCCA model to capture the view-specific variations into the learned latent representations. Finally, it would have been valuable to have natural images collected by VI people to benchmark the classifiers on scenarios with the actual target users [15, 47].

In Paper B and Section 3.4, we observed that utilizing the iconic images in VCCA affects the classification performance significantly better than the text descriptions. We believe that this could be due to the noise in the text descriptions where often there are only a few words that are relevant to the image recognition task. In the future, we could pre-process the text descriptions by removing stop words, e.g., 'the', 'it', 'and', etc., and other words that are irrelevant for recognizing the fine-grained details of the items. Moreover, we could use pretrained word embeddings [209–211] for considering word similarities, add attention mechanisms [93, 212] in the model for

focusing on relevant words, and use data augmentation techniques for text [213] to learn more robust joint latent representations.

Finally, we believe the Grocery Store dataset could be used for benchmarking image classifiers on zero/few-shot learning [39, 40] and continual learning [37, 38] settings. This would be an important extension of the dataset as these applications are crucial to enhance the usability of assistive vision systems in real-world scenarios.

Chapter 4

Continual Learning

In this chapter, we give an overview of Paper C and D where we focus on continual learning (CL) [37, 38]. A CL system receives samples of classes to learn continuously over time while maintaining its overall performance on every seen class. Such capability is critical for assistive vision systems as the system may need to adapt to unseen items in new environments. Furthermore, CL on mobile devices would enable personalizing the devices to specific items that the VI user wants to locate and recognize. The main challenge with the CL setup is catastrophic forgetting [63] where previously learned abilities are overwritten by recent acquired knowledge (see Figure 4.1). However, retraining on all previously seen data may be prohibited due to lack of computational budget or accessibility to the seen datasets. Replay-based CL methods efficiently mitigate catastrophic forgetting by revisiting old samples stored in a small memory buffer [214–216].

Most replay-based methods in CL use simple strategies for retrieving replay samples and ignore the time to learn which is important for memory retention in human learning [67, 69, 70]. Furthermore, most works rely on tiny memory buffers although historical data is almost always available in many real-world machine learning applications [64–66]. Nevertheless, the requirement of small replay memories remains as due to limitations on computational budget and data transmission times. To this end, we propose in Paper C a new CL setting where historical data is available while the processing time is limited. For this new setting, we propose *replay scheduling* where we select which tasks to replay at different times, and demonstrate the advantages of considering time to replay in CL. In Paper D, we present a framework based on reinforcement learning (RL) [75] for learning replay scheduling policies that can be transferred to new CL scenarios for mitigating catastrophic forgetting without added computational cost. Our proposed replay scheduling approach opens up for new research directions that can bring current CL research closer to real-world needs.

4.1 Related Work

In this section, we give an overview of different approaches in CL, especially replay-based methods which our proposed replay scheduling approach in Paper C and D is based on.

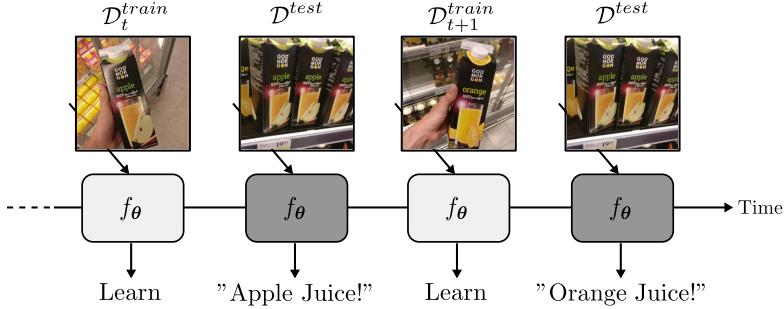


Figure 4.1: Illustration of catastrophic forgetting in a continual learning scenario. The network f_θ receives a task dataset of an Apple Juice class in D_t^{train} to learn, and correctly predicts the Apple Juice class in the unseen test image in D^{test} . The next task dataset involves learning an Orange Juice class from D_{t+1}^{train} which f_θ adapts to without retaining knowledge about the Apple Juice class, such that f_θ makes the wrong prediction during evaluation on D^{test} again.

Continual Learning. The existing CL approaches for mitigating catastrophic forgetting in neural networks can broadly be divided into three main areas, namely, regularization-based, architecture-based, and replay-based methods. Regularization-based methods mainly focus on preventing the parameters important for previous tasks from wide changes and fit the remaining parameters to new tasks [217–219]. Knowledge distillation methods [220] also belong to these methods where stored classification logits are used for regularizing the output units for previous tasks in the network [221, 222]. More recently, projection-based approaches have been proposed for constraining the parameter updates to subspaces which avoid interference with previous tasks [223, 224]. Architecture-based methods focus on adding task-specific network modules for every seen task [225–229], or isolating parameters for predicting specific task in fixed-size networks [230–232]. Replay-based methods revisits samples from old tasks when learning new tasks. The old samples are either stored in an external memory [214, 216, 233–235], or synthesized with a generative model [109, 236–238]. Both regularization- and architecture-based methods can be combined with replay for improving the models capability of remembering tasks [219, 228, 239–244]. Our replay scheduling idea is originated from replay-based methods which we will cover more in detail next.

Replay-based Continual Learning. Much research effort in replay-based CL has focused on selecting higher quality samples to store in memory [214, 215, 219, 245–251]. In [214], several selection strategies are reviewed in scenarios with tiny memory capacity, such as reservoir sampling [252], first-in first-out buffer [215], k-Means, and Mean-of-Features [250]. However, elaborate selection strategies have been shown to give little benefit over random selection for image classification problems [216, 253]. More recent work has been focused on enhancing the memory capacity by storing compressed features rather than raw data [216, 254, 255], evolving the memory samples using data augmentation to avoid overfitting [234, 256], and using contrastive learning to improve memory retention [257, 258]. Selecting the time to replay old tasks has mostly been ignored in the literature, with an exception in [259] which replays memory samples that would most interfere with a foreseen parameter update. Our replay scheduling approach differs from the above mentioned works since we focus on

learning to select which tasks to replay. Nevertheless, our scheduling can be combined with any selection strategy and replay-based method.

4.2 Replay Scheduling in Continual Learning

In this section, we present our proposed replay scheduling approach studied in Paper C and D for mitigating catastrophic forgetting in CL. More specifically, we focus on the contributions in Paper C, where we introduce a slightly new CL setting that considers real-world needs where all historical data can be available since data storage is cheap. Additionally, we demonstrate the advantages of replay scheduling in such setting where we select which tasks to replay, where we use Monte Carlo tree search (MCTS) [260] to search for efficient replay schedules by allowing multiple episodes in single CL environments.

New Problem Setting for Continual Learning

In this section, we describe the proposed CL setting presented in Paper C. We assume that historical data is accessible at any time for mitigating catastrophic forgetting in the model. However, re-training the model on all available historical data is prohibited due to limitations on compute and data transmission times. Therefore, the goal is to determine which historical tasks to replay and sample a small replay memory from the selected tasks to mitigate catastrophic forgetting as efficiently as possible.

The notation of this problem setting resembles the traditional CL setting for image classification. We let a neural network f_ϕ , parameterized by ϕ , learn T tasks from the datasets $\mathcal{D}_1, \dots, \mathcal{D}_T$ arriving sequentially one at a time. The t -th dataset $\mathcal{D}_t = \{(\mathbf{x}_t^{(i)}, y_t^{(i)})\}_{i=1}^{N_t}$ consists of N_t samples where $\mathbf{x}_t^{(i)}$ and $y_t^{(i)}$ are the i -th data point and class label respectively. Additionally, the dataset is split into training, validation, and test sets such that $\mathcal{D}_{1:T} = \{(\mathcal{D}_t^{train}, \mathcal{D}_t^{val}, \mathcal{D}_t^{test})\}_{t=1}^T$. The training objective at task t is given by

$$\min_{\phi} \sum_{i=1}^{N_t} \mathcal{L}(f_\phi(\mathbf{x}_t^{(i)}), y_t^{(i)}), \quad (4.1)$$

where $\mathcal{L}(\cdot)$ is the loss function, which in our case is the cross-entropy loss. The challenge is for the network f_ϕ to retain its classification performance on the previous $t - 1$ tasks.

We assume that historical data from old tasks are accessible at any time step. However, we can only fill a small replay memory \mathcal{M} with M historical samples for mitigating catastrophic forgetting of old tasks due to processing time constraints. The challenge then becomes how to select the M samples from the historical data to include in the replay memory that efficiently retain the previous knowledge. We focus on selecting the samples on task-level by deciding on the task proportions (p_1, \dots, p_{t-1}) on samples to fetch from each task, where $\sum_{i=1}^{t-1} p_i = 1$ and $p_i \geq 0$. The number of samples from task i to place in \mathcal{M} is given by $p_i \cdot M$. To simplify the selection of which tasks to replay, we construct a finite set of possible task proportions that can be used for constructing \mathcal{M} .

This problem setting shares several assumptions with traditional CL, such as:

- The model receives the datasets at different time steps from a continuous stream.
- Re-training on all historical data is prohibited.
- The model should perform well overall across all seen tasks and, hence, mitigate catastrophic forgetting.
- Replay memory size M is small.

The only difference with this setting and traditional CL is that all historical has been stored rather than deleted, and we can access this data to fill the replay memory \mathcal{M} to mitigate catastrophic forgetting. We argue that this setting aligns well with real-world needs as data storage is cheap and easy to maintain but retraining large machine learning models is computationally expensive. Therefore, we use small replay memory sizes to limit the processing times of replay when the model is learning new tasks.

Monte Carlo Tree Search for Replay Scheduling

In this section, we describe how we used MCTS in Paper C to show the advantages of replay scheduling in the new CL setting. We focus on an ideal CL environment where executing multiple episodes is allowed to enable searching for replay schedules. For demonstration purposes, we use MCTS for finding replay schedules that efficiently mitigate catastrophic forgetting by selecting compositions of replay memories at different time steps.

We define a replay schedule as a sequence $S = (\mathbf{p}_1, \dots, \mathbf{p}_{T-1})$, where $\mathbf{p}_i = (p_1, \dots, p_{T-1})$ for $1 \leq i \leq T-1$ is the sequence of task proportions for determining how many samples per task to fill the replay memory with at time step i . To make the selection of task proportions tractable, we construct a discrete action space with a finite number of choices for how to construct the replay memory at different time steps. Algorithm 1 shows the procedure for creating this action space. At task i , we have $i-1$ historical tasks that we can choose from. We then generate all possible bin vectors $\mathbf{b}_i = [b_1, \dots, b_i] \in \mathcal{B}_i$ of size i where each element are a task index $1, \dots, i$.

The elements in each bin vector are sorted by task index, and we only keep the unique bin vectors after sorting. For example, at $i=2$, the unique choices of vectors are $[1, 1], [1, 2], [2, 2]$, where $[1, 1]$ indicates that all samples in the replay memory should be from task 1, $[1, 2]$ indicates that half memory is from task 1 and the other half are from task etc. The task proportions are then computed by counting the number of occurrences of each task index in \mathbf{b}_i and dividing by i , such that $\mathbf{p}_i = \text{bincount}(\mathbf{b}_i)/(i)$. From this specification, we have built a tree \mathcal{T} with different task proportions that can be selected at different time steps. A replay schedule S is constructed by traversing through \mathcal{T} and appending a single task proportion from every task level to S .

Algorithm 1 Discretization of action space with task proportions

```

Require: Number of tasks  $T$ 
1:  $\mathcal{T} = ()$ 
2: for  $i = 1, \dots, T-1$  do
3:    $\mathcal{P}_i = \{\}$ 
4:    $\mathcal{B} = \text{combinations}([1 : i], i)$ 
5:    $\mathcal{B}^* = \text{unique}(\text{sort}(\mathcal{B}))$ 
6:   for  $\mathbf{b}_i \in \mathcal{B}^*$  do
7:      $\mathbf{p}_i = \text{bincount}(\mathbf{b}_i)/i$ 
8:      $\mathcal{P}_i = \mathcal{P}_i \cup \{\mathbf{p}_i\}$ 
9:   end for
10:   $\mathcal{T}[i] = \mathcal{P}_i$ 
11: end for
12: return  $\mathcal{T}$ 
```

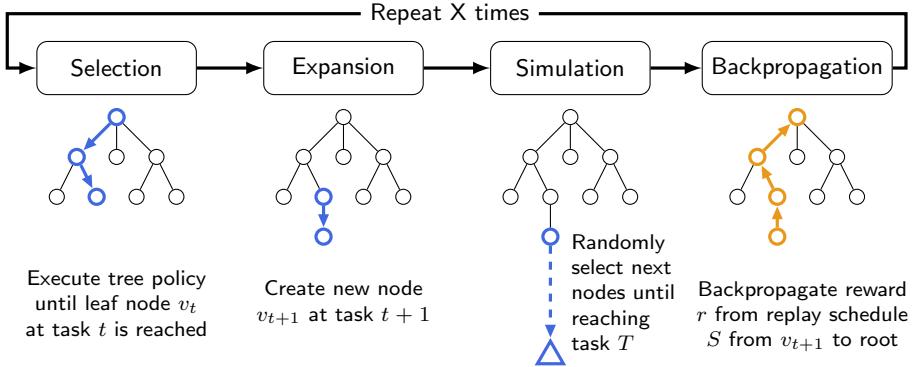


Figure 4.2: Outline of MCTS when searching for replay schedules.

The tree-shaped action space of task proportions grows fast with the number of tasks, which complicates studying replay scheduling in datasets with longer task-horizons. Therefore, we proposed using MCTS since it has been successful in applications with large action spaces [78, 261, 262]. In our case, MCTS concentrates the search for replay schedules in directions with promising CL performance in the environment. Each memory composition in the action space corresponds to a node that can be visited by MCTS. For instance, at task t , the node v_t is related to a task proportions \mathbf{p}_t that can be used for retrieving a replay memory from the historical data. We store the related task proportion \mathbf{p}_t from every visited node v_t in the replay schedule S . The final replay schedule is then used for constructing the replay memories at each task during the CL training. Figure 4.2 shows an illustration of the procedure for executing MCTS to search for replay schedules, where each step involves the following:

Selection. During a rollout, the tree policy either selects an unvisited child randomly from the current node v_t , or selects the next node by evaluating the Upper Confidence Tree (UCT) [263] if all children have been visited earlier. The child v_{t+1} with the highest UCT score is selected using the function from [262]:

$$UCT(v_t, v_{t+1}) = \max(q(v_{t+1})) + C \sqrt{\frac{2 \log(n(v_t))}{n(v_{t+1})}}, \quad (4.2)$$

where $q(\cdot)$ is the reward function, C the exploration constant, and $n(\cdot)$ the number of node visits. The tree policy is executed until a leaf node v_t is reached.

Expansion. If the leaf node v_t has unvisited children, the search tree is expanded with one of the unvisited child nodes v_{t+1} selected with uniform sampling.

Simulation and Reward. After the expansion step, the succeeding nodes are selected randomly until reaching a terminal node v_T . The task proportions $\mathbf{p}_{1:T-1}$ from the visited nodes in the rollout constitutes the replay schedule S . After training the network using S for replay, we calculate the reward for the rollout given by $r = \frac{1}{T} \sum_{i=1}^T A_{T,i}^{(val)}$, where $A_{T,i}^{(val)}$ is the validation accuracy of task i at task T .

Backpropagation. The reward r is backpropagated from the expanded node v_{t+1} to the root node, where the reward function $q(\cdot)$ and number of visits $n(\cdot)$ are

updated at each node in the path.

As mentioned earlier, we use MCTS to search for replay schedules to illustrate the importance of learning the time to learn. To this end, we study replay scheduling in an ideal CL environment where multiple episodes is allowed to find replay schedules. However, in the standard CL setting, multiple episodes are prohibited which currently prevents replay scheduling to be applied in real-world CL scenarios. In the next section, we present a policy learning framework based on RL to learn replay scheduling policies that generalize.

4.3 Policy Learning for Replay Scheduling

In this section, we introduce the RL-based framework for learning replay scheduling policies that was proposed in Paper D. In real-world CL settings, learning the policy with multiple episodes is infeasible, which is why we need a general policy for replay scheduling. Our intuition is that there may exist general patterns regarding the replay scheduling, e.g., tasks that are harder or have been forgotten should be replayed more often. Moreover, the policy may non-trivially take task properties into consideration. Therefore, we aim to learn policies that select which tasks to replay from states representing the current task performance in the CL environments. The policy can then be applied for mitigating catastrophic forgetting in new CL scenarios.

We model the CL environments as Markov Decision Processes [112] (MDPs) where each MDP is represented as a tuple $E_i = (\mathcal{S}_i, \mathcal{A}, P_i, R_i, \mu_i, \gamma)$ consisting of the state space \mathcal{S}_i , action space \mathcal{A} , state transition probability $P_i(s'|s, a)$, reward function $R_i(s, a)$, initial state distribution $\mu_i(s_1)$, and discount factor γ . We assume that we have access to a fixed set of training environments $\mathcal{E}^{(train)} = \{E_1, \dots, E_K\}$ sampled from a distribution of CL environments, i.e., $E_i \sim p(E)$ for $i = 1, \dots, K$. Each environment E_i contains of a network f_ϕ and T datasets $\mathcal{D}_{1:T}$ where the t -th dataset is learned by f_ϕ at time step t . Note that we will use task and time step interchangeably. To generate diverse CL environments, we obtain environments with different initializations of the network f_ϕ and shuffled task orders in the dataset when we sample environments from $p(E)$.

We define the state s_t of the environment as the validation accuracies $A_{t,1:t}^{(val)}$ on each seen task $1, \dots, t$ from f_ϕ at task t , i.e., $s_t = [A_{t,1}, \dots, A_{t,t}, 0, \dots, 0]$, where we use zero-padding on future tasks. The action space \mathcal{A} is constructed as described in Algorithm 1 (Section 4.2), such that the $a_t \in \mathcal{A}$ corresponds to a task proportion \mathbf{p}_t used for sampling the replay memory \mathcal{M}_t . We use a dense reward based on the average validation accuracies at task t , i.e., $r_t = \frac{1}{t} \sum_{i=1}^t A_{t,i}^{(val)}$. The state transition distribution $P_i(s'|s, a)$ represents the dynamics of the environment, which depend on the initialization of f_ϕ and also the task order in the dataset. More specifically, the dynamics represent how the validation accuracies varies after f_ϕ replays according to the action a_t when learning the current task.

The procedure for training the policy goes as follows: The state s_t is obtained by evaluating the network f_ϕ on the validation sets $\mathcal{D}_{1:t}^{(val)}$ after learning the t -th task from $\mathcal{D}_t^{(train)}$. Action a_t is selected under the policy $\pi_\theta(a|s_t)$ parameterized by θ . The action is converted into task proportion \mathbf{p}_t that is used for sampling the replay memory \mathcal{M}_t from the historical datasets. We then train classifier f_ϕ with $\mathcal{D}_{t+1}^{(train)}$

and \mathcal{M}_t , and obtain the reward r_{t+1} and the next state s_{t+1} by evaluating f_ϕ on the validation sets $\mathcal{D}_{1:t+1}^{(val)}$. The collected transitions $(s_t, a_t, r_{t+1}, s_{t+1})$ are used for updating the policy. A new episode starts after f_ϕ has learned the final task T .

We evaluate the learned policy by applying it to mitigate catastrophic forgetting in new CL environments at test time. To foster generalization across environments, we train the policy on multiple environments with different dynamics, such as task orders and datasets, to learn from a diverse set of training data. The goal for the agent is to maximize the sum of rewards in each training environment. At test time, the policy is applied on new CL classifiers and datasets in the test environments without added computational cost nor experience collection. In Section 4.4, we test the policies generalization capability to new CL environments where the task orders and datasets are unseen during training.

4.4 Experiments

In this section, we give an overview of the experimental results from Paper C and D. In Paper C, we evaluate the benefits of replay scheduling in single CL environments using MCTS for finding replay schedules. We perform extensive evaluation to show that the replay schedules from MCTS can outperform the baselines across several CL benchmarks and different backbones. In Paper D, we evaluate our RL-based framework using DQN [113] and A2C [118] for learning policies that generalize to new CL scenarios. We show that the learned policies can efficiently mitigate catastrophic forgetting in CL environments with new task orders and datasets that are unseen during training.

Results on Replay Scheduling with Monte Carlo Tree Search

Here, we provide a selection of the experimental results from Paper C. Our approach has been named Replay Scheduling MCTS (RS-MCTS). We evaluate the CL performance of RS-MCTS for varying memory sizes to illustrate the importance of replay scheduling. Moreover, we show that replay scheduling can benefit any replay-based method by combining scheduling with three recent CL methods, namely, HAL [264], MER [265], and DER [239]. Finally, we demonstrate the efficiency benefits of replay scheduling in settings when the memory size is smaller than the number of classes. We give an overview of the experimental settings, see Paper C for the full experimental results and hyperparameter settings.

Datasets and Architectures. We perform experiments on common CL benchmark datasets such as Split MNIST [218], Split Fashion-MNIST [266], Split notMNIST [267], Permuted MNIST [268], and Split CIFAR-100 [151, 215, 250] and Split miniImagenet [269]. We use a 2-layer MLP with 256 hidden units and ReLU activations for Split MNIST, Split FashionMNIST, Split notMNIST, and Permuted MNIST. For Split CIFAR-10 and CIFAR-100, we use the CNN architecture used in [222, 269, 270], and a reduced ResNet-18 from [215] for Split mini-Imagenet. We use a multi-head output layer for each dataset, except Permuted MNIST where the network uses single-head output layer, such that we assume that task labels are available at test time [271].

Table 4.1: Performance comparison with ACC between scheduling methods RS-MCTS (Ours), Random, ETS, and Heuristic combined with replay-based methods HAL, MER, and DER. Memory sizes are $M = 10$ and $M = 100$ for the 5-task and 10/20-task datasets respectively. We report the mean and std. averaged over 5 seeds. * denotes where some seed in Heuristic did not converge.

Method	Schedule	5-task Datasets			10- and 20-task Datasets		
		S-MNIST	S-FashionMNIST	S-notMNIST	P-MNIST	S-CIFAR-100	S-miniImagenet
HAL	Random	97.24 (± 0.70)	86.74 (± 6.05)	93.61 (± 1.31)	88.49 (± 0.99)	36.09 (± 1.77)	38.51 (± 2.22)
	ETS	94.02 (± 4.25)	95.81 (± 3.53)	91.01 (± 1.39)	88.46 (± 0.86)	34.90 (± 2.02)	38.13 (± 1.18)
	Heuristic	97.69 (± 0.19)	*74.16 (± 11.19)	93.64 (± 0.93)	*66.63 (± 28.50)	35.07 (± 1.29)	39.51 (± 1.49)
	Ours	97.93 (± 0.56)	98.27 (± 0.17)	94.64 (± 0.39)	89.14 (± 0.74)	40.22 (± 1.57)	41.39 (± 1.15)
MER	Random	93.07 (± 0.81)	85.53 (± 3.30)	91.13 (± 0.86)	75.90 (± 1.34)	42.96 (± 1.70)	31.48 (± 1.65)
	ETS	92.89 (± 3.53)	96.47 (± 0.85)	93.80 (± 0.82)	73.01 (± 0.96)	43.38 (± 1.81)	33.58 (± 1.53)
	Heuristic	94.30 (± 2.79)	96.91 (± 0.62)	90.90 (± 1.30)	83.86 (± 3.19)	40.90 (± 1.70)	34.22 (± 1.93)
	Ours	98.20 (± 0.16)	98.48 (± 0.26)	93.61 (± 0.71)	79.72 (± 0.71)	44.29 (± 0.69)	32.74 (± 1.29)
DER	Random	98.23 (± 0.53)	96.56 (± 1.79)	92.89 (± 0.86)	87.51 (± 1.10)	56.83 (± 0.76)	42.19 (± 0.67)
	ETS	98.17 (± 0.35)	97.69 (± 0.58)	94.74 (± 1.05)	85.71 (± 0.75)	52.58 (± 1.49)	35.50 (± 2.84)
	Heuristic	94.57 (± 1.71)	*72.49 (± 19.32)	*77.88 (± 12.58)	81.56 (± 2.28)	55.75 (± 1.08)	43.62 (± 0.88)
	Ours	99.02 (± 0.10)	98.33 (± 0.51)	95.02 (± 0.33)	90.11 (± 0.18)	58.99 (± 0.98)	43.46 (± 0.95)

Baselines. We compare our method against the following scheduling baselines:

- **Random Schedule:** Randomly select which tasks to replay.
- **Equal Task Schedule (ETS):** Replay all seen tasks equally.
- **Heuristic Schedule:** Replay tasks which validation performance is below a tuned threshold found through hyperparameter searches.

Evaluation. We use the following evaluation metric for measuring the CL performance

$$\text{ACC} = \frac{1}{T} \sum_{j=1}^T A_{T,j}^{test}, \quad (4.3)$$

where $A_{T,j}^{test}$ is the accuracy of task j after learning the final task T .

Varying Memory Size. We show that the replay schedules found by MCTS improves the CL performance across different memory sizes. In Figure 4.3, we observe that RS-MCTS obtains better task accuracies than ETS, especially for small memory sizes. Both RS-MCTS and ETS perform better than Heuristic as M increases showing that Heuristic requires careful tuning of the validation accuracy threshold. These results show that replay scheduling can outperform the baselines, especially for small M , on both small and large datasets across different backbone choices.

Applying Scheduling to Recent Replay Methods. We show that replay scheduling can be combined with any memory-based CL method to enhance the performance. Table 4.1 shows the performance comparison between our RS-MCTS against using Random, ETS, and Heuristic schedules for each method. Especially for HAL and DER, the schedule found by RS-MCTS mostly outperforms the other schedules

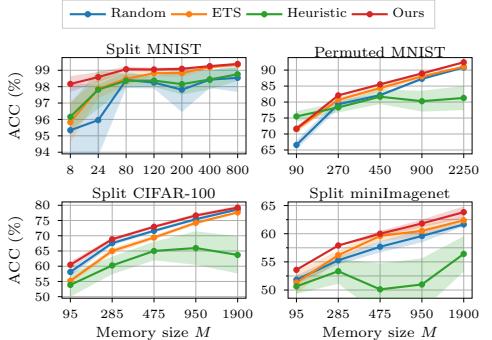


Figure 4.3: Performance comparisons with ACC over different replay memory sizes M for RS-MCTS (Ours) and the baselines. All results have been averaged over 5 seeds.

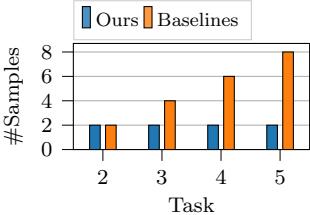


Figure 4.4: Number of replayed samples per task for the 5-task datasets for RS-MCTS (Ours) and the baselines in the tiny memory setting.

across the different datasets. For MER, the Heuristic baseline performs better than RS-MCTS and Permuted-MNIST and Split miniImagenet which could potentially be adjusted with more hyperparameter searches. Furthermore, the Heuristic baseline occasionally did not converge for some seeds which could indicate that this scheduling method is more sensitive to the settings than the other scheduling methods. These results confirm that replay scheduling is important for the final performance given the same memory constraints and it can benefit any existing CL framework.

Efficiency of Replay Scheduling. We illustrate the efficiency of replay scheduling with comparisons to several common replay-based CL baselines in a memory setting where we can only store 1 sample/class. Our goal is to investigate if scheduling over which tasks to replay can be more efficient in situations where the memory size is even smaller than the number of classes. To this end, we set the replay memory size for our method to $M = 2$ for the 5-task datasets. We then compare against the most memory efficient CL baselines, namely A-GEM [272], ER-Ring [214] which have shown promising results with 1 sample/class for replay, as well as with uniform memory selection as reference. Additionally, we compare to using random replay schedules (Random) with the same memory setting as for RS-MCTS. We visualize the memory usage for our method and the baselines when training on a 5-task dataset in Figure 4.4. Table 4.2 shows the ACC for each method across the 5-task datasets. Despite using significantly fewer samples for replay, RS-MCTS performs better or on par with the best baselines on each dataset. These results show that replay scheduling can be even more efficient than the common memory efficient replay-based methods, which indicate that learning the time to learn is an important research direction in CL.

Policy Generalization to New Continual Learning Scenarios

In this section, we present the experimental results from Paper D showing that the policies learned with our RL-based framework using DQN and A2C are capable of generalizing across CL environments with new task orders and datasets unseen during training. We run experiments on 5-task datasets Split MNIST [218], Split FashionMNIST [266], Split notMNIST [267], and Split CIFAR-10 [151] where there are 2 classes/task. We resort to CL benchmark datasets with $T = 5$ tasks to have feasible size of the discrete action space. The replay memory size is set to $M = 10$ in all experiments. See Paper D for full details on the experimental setup.

Method	5-task Datasets		
	S-MNIST	S-FashionMNIST	S-notMNIST
Random	92.56 (± 2.90)	92.70 (± 3.78)	89.53 (± 3.96)
A-GEM	94.97 (± 1.50)	94.81 (± 0.86)	92.27 (± 1.16)
ER-Ring	94.94 (± 1.56)	95.83 (± 2.15)	91.10 (± 1.89)
Uniform	95.77 (± 1.12)	97.12 (± 1.57)	92.14 (± 1.45)
Ours	96.07 (± 1.60)	97.17 (± 0.78)	93.41 (± 1.11)

Table 4.2: Performance comparison with ACC averaged over 5 seeds in the 1 ex/class memory setting. RS-MCTS (Ours) has replay memory size $M = 2$, while the baselines replay all available memory samples.

DQN and A2C Architectures. The input layer has size $T - 1$ where each unit is inputting the task performances since the states are represented by the validation accuracies $s_t = [A_{t,1}^{(val)}, \dots, A_{t,t}^{(val)}, 0, \dots, 0]$. The current task can therefore be determined by the number of non-zero state inputs. The output layer has 35 units representing the possible actions at $T = 5$ in the discrete action space. We use action masking on the output units to prevent the network from selecting invalid actions for constructing the replay memory at the current task. The DQN is a 2-layer MLP with 512 hidden units and ReLU activations. For A2C, we use separate networks for parameterizing the policy and the value function, where both networks are 2-layer MLPs with 64 hidden units of Tanh activations.

Baselines. We add two more heuristic scheduling baselines:

- **Heuristic Local Drop (Heur-LD).** Heuristic policy that replays tasks with validation accuracy below a threshold proportional to the previous achieved validation accuracy on the task.
- **Heuristic Accuracy Threshold (Heur-AT).** Heuristic policy that replays tasks with validation accuracy below a fixed threshold.

Here, we name the Heuristic from the experiments in Paper C as Heuristic Global Drop (Heur-GD). The Random and ETS baselines are the same as earlier.

Evaluation Protocol. We use ACC from Equation 4.3 to evaluate the performance in the CL environments. The reported results have been evaluated on the test set where the replay schedules are selected from evaluating the validation sets, where 15% of the training data is used for creating the validation sets. To assess generalization capability, we use a ranking method based on the ACC between the methods in every test environment for comparison (see Paper D for details).

Generalization to New Task Orders.

We show that the learned replay scheduling policies can generalize to CL environments with previously unseen task orders. We generate training and test environments with unique task orders for three datasets, namely, Split MNIST, FashionMNIST, and CIFAR-10. Table 4.3 shows the average ranking for the DQN, A2C, and the baselines when being applied in 10 test environments. Our learned policies obtain the best average ranking across the datasets, where the DQN performs best for Split MNIST and FashionMNIST while A2C outperforms all methods on Split CIFAR-10. We provide further insights in the benefits of learning the replay scheduling policy by visualizing the replay schedule and task accuracy progress from A2C in one Split CIFAR-10 test environment in Figure 4.5.

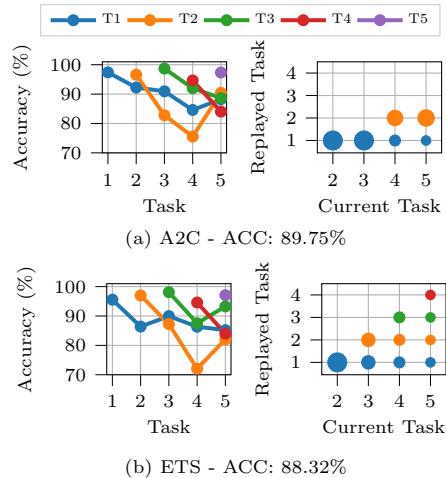


Figure 4.5: Task accuracies and replay schedules for A2C and ETS for a Split CIFAR-10 environment.

Table 4.3: Average ranking (lower is better) for experiments on generalizing policies to environments with new task orders or a new dataset. We average the results over 10 test environments.

Method	New Task Order			New Dataset	
	S-MNIST	S-FashionMNIST	S-CIFAR-10	S-notMNIST	S-FashionMNIST
Random	4.23	3.60	5.03	3.87	3.95
ETS	3.80	4.57	5.37	4.27	3.63
Heur-GD	4.48	4.25	3.98	4.58	2.77
Heur-LD	4.65	3.65	3.75	4.97	5.10
Heur-AT	4.33	3.87	3.43	4.28	3.72
DQN (Ours)	3.23	3.55	3.68	3.20	4.37
A2C (Ours)	3.27	4.52	2.75	2.83	4.47

Additionally, we show the replay schedule and task accuracies from the ETS baseline in the same environment. The replay schedules are visualized with bubble plots showing the selected task proportion to use for composing the replay memories at each task. In Figure 4.5a, we observe that A2C decides to replay task 2 more than task 1 as the performance on task 2 decreases, which results in a slightly better ACC metric achieved by A2C than ETS. These results show that the learned policy can flexibly consider replaying forgotten tasks to enhance the CL performance.

Generalization to New Datasets. We show that the learned replay scheduling policy is capable of generalizing to CL environments with new datasets unseen in the training environments. We perform two sets of experiments, 1) train with environments generated with Split MNIST and FashionMNIST and test on environments generated with Split notMNIST, and 2) train with environments generated with Split MNIST and notMNIST and test on environments generated with Split FashionMNIST. Table 4.3 shows the average ranking for DQN, A2C, and the baselines when generalization to test environments with new datasets. We observe that both A2C and DQN successfully generalize to Split notMNIST environments outperforming all baselines. However, the learned policies have difficulties to generalize to Split FashionMNIST environments, which could be due to high variations in the dynamics between training and test environments. The policy may exhibit state transition dynamics which has not been experienced during training, which makes generalization difficult for both DQN and A2C. Potentially, the performance could be improved by generating more training environments for the agent to exhibit more variations in the CL scenarios or by using other advanced RL methods which may generalize better [273].

4.5 Discussion

In this chapter, we summarized our contributions in CL presented in Paper C and D. In Paper C, we focused on a new CL setting where historical data is accessible at any time for mitigating catastrophic forgetting, but replaying all historical data is prohibited due to processing time constraints. We then proposed learning the time to learn for a CL system where we learn replay schedules over which tasks to replay at different time steps. To demonstrate the importance of learning the time to learn, we use MCTS to find the proper replay schedule and show that it can significantly outperform fixed scheduling policies in terms of CL performance. In Paper D, we address the efficiency of replay scheduling by proposing to use RL to learn the replay

scheduling policies that can generalize to new continual learning scenarios without added computational cost. In the experiments, we show that the learned policies can propose replay schedules that efficiently mitigate catastrophic forgetting in environments with previously unseen task orders and datasets. The proposed approaches opens up for new research directions in replay-based CL that stems well with real-world needs where training CL systems can be limited by processing time constraints rather than storing historical data for maintaining overall performance.

We argue that replay scheduling brings current CL research closer to real-world CL scenarios. Take for example image classification apps on mobile phones that should assist VI people with recognizing objects in their surroundings. The users may want to customize their apps to recognize new and personal items continuously. The app should adapt fast from few samples while maintaining its performance on the already known classes which is where CL comes in. However, as we saw in Paper C and D, fixed replay policies show difficulties performing well given different task orders, which potentially complicates generalizing to different users and new environments. To this end, we could transfer a learned replay scheduling policy that generalizes to new users and their unseen environments to mitigate catastrophic forgetting efficiently for any user’s personal image recognition app. Moreover, a more general case for the importance of replay scheduling in CL is scenarios where the allowed replay memory size is smaller than the number of seen classes to remember. Scheduling over which task to replay at different times then becomes crucial to balance efficient learning of new classes while retaining previously learned abilities.

Generalization in RL is a challenging research topic by itself [274–276]. The main limitation with the RL-based framework proposed in Paper D is the requirements of large amounts of diverse data and training time to enable the learned policy to generalize well. This can be expensive due to the computational cost to generate the CL environments since each state transition involves training the classifier on a CL task. Moreover, we have currently only considered a discrete action space which is hard to construct especially when the number of tasks is large. Thus, in future work, we would explore more advanced RL methods which can handle continuous action spaces and generalize well [273].

Finally, we want to highlight the privacy issues that can arise when applying replay scheduling. Any replay-based CL method inherently can violate privacy if raw samples are stored in the memory, but these risks can be mitigated by instead storing compressed feature representations. Furthermore, as mentioned above, we needed substantial amounts of data and training time for learning policies that generalize which increases the need for secure data storage solutions. As the states were defined using validation accuracies rather than images, we have circumvented the potential privacy issues as only the state transitions are stored from the generated CL environments to use for training the policy.

Chapter 5

Conclusions and Future Directions

In this chapter, we provide our conclusions from the works described in this thesis (Section 5.1). Thereafter, we discuss potential future directions for work towards better computer vision methods for assisting visually impaired (VI) people (Section 5.2).

5.1 Conclusions

We have divided our conclusions into two parts, where we first provide our insights on the work on fine-grained image recognition (FGIR) of grocery items presented in Chapter 3, and thereafter our insights on our proposed replay scheduling method for continual learning (CL) presented in Chapter 4.

Fine-Grained Image Recognition of Groceries

Here, we provide our insights around the work presented in Chapter 3.

Dataset Collection. In Paper A, we presented a realistic dataset of images of groceries with a mobile phone to simulate a grocery shopping scenario using an assistive vision app. Collecting these images in the grocery stores was indeed time consuming and required effort. However, as there is a gap between benchmark datasets and reality, we argue that evaluating computer vision methods on real-world data and scenarios is critical for assistive vision. One major learned lesson is that we should have collected the images by recording videos rather than taking snapshot images. First of all, recording videos to classify the groceries would be more user-friendly and probably more accurate than processing static images one at a time. The dataset size would probably be enhanced as well, and the dataset could be used for evaluating both image and video models on supervised classification. In addition to data quality, the selection of network architecture can also have high influence on the final performance. Therefore, we should have performed a broader study over different network types, especially those specialized for mobile computing [277, 278], while measuring memory consumption and compute requirements as computation should be on-device to preserve user privacy [279]. It is also important to collect meta-data, such as hardware information, image resolutions, and geographical loca-

tion, about the mobile phone images to develop image recognition models that are device-agnostic and robust in new environments. Finally, it would have been valuable to include images from blind/low-vision collectors to measure the FGIR performance with data from the actual users of assistive vision devices [280].

Utilizing Web-Scraped Views for FGIR. In Paper B, we showed that the web-scraped iconic images and text descriptions in the dataset from Paper A are useful for enhancing the FGIR performance. Our ablation study showed that the learned latent representations improves the classification performance as the iconic images structures the items based on visual similarities (e.g. color and shape) in the latent space, while the text descriptions structures the items based on ingredients and flavors. To provide more insights on how the web-scraped views boost the performance, we could analyze the natural images with visual explanation methods, e.g., with class activation maps [281, 282], to investigate whether the classifier recognizes the fine-grained details in the items. However, performing such evaluation automatically probably requires human annotations on the location of the details. Adding more examples of iconic images and text descriptions may enable better modeling of the view-specific variations to more accurately capture the shared information about the views into the learned latent representations. In combination with more advanced text processing [211, 213], this could potentially enhance the separation between visually similar packaged items with different ingredients to further improve the FGIR performance. Finally, it would be interesting to study how the web-scraped views can be utilized to enhance the data-efficiency in few-shot learning and CL scenarios, which are important applications for assistive vision systems.

Replay Scheduling in Continual Learning

In this section, we provide our insights around the work presented in Chapter 4.

Benefits of Replay Scheduling. In Paper C, we proposed a slightly new CL setting where 1) the historical is stored rather than deleted and 2) the replay memory size is limited by processing time constraints instead of storage capacity. This setting stems well with real-world needs as many major companies stores all their recorded data [66, 283, 284]. We then showed in an ideal CL environment that scheduling over which tasks to replay can significantly improve the final CL performance compared to using fixed policies, e.g., replaying all tasks equally, especially with small replay memory sizes. Consider a CL system where the number of seen classes is much larger than the number of samples that the system has time to replay. In such scenario, it becomes crucial for the system to have a policy scheduling over which tasks to replay for reducing catastrophic forgetting efficiently. Furthermore, advances in CL are critical for assistive vision devices to maintain recognition of new objects in ever-changing environments without the need for re-training the systems from scratch. As these systems could require on-device learning to maintain privacy, replay scheduling could potentially retain the previous knowledge better than fixed scheduling policies and reduce idle times when learning new abilities.

Learning Replay Scheduling Policies. In Paper D, we proposed a framework based on reinforcement learning (RL) for learning policies scheduling which tasks to replay based on the current CL performance. Furthermore, we showed that the

replay scheduling policies could generalize to new CL scenarios without additional computational cost nor interactions with the test environments. The proposed approach enables replay scheduling to be applied in standard CL settings. However, as we only considered a discrete action space in the experiments, we need to extend the framework to handle continuous action spaces that would scale better to larger number of tasks. Furthermore, learning a policy that generalizes well to new scenarios required large amounts of diverse data and training time, so we need to explore more sample-efficient methods. We believe that our policy learning framework could be beneficial for user personalization of assistive vision systems using CL. Since fixed scheduling policies show difficulties generalizing to new CL scenarios in our experiments, we could instead transfer a learned replay scheduling policy that potentially can generalize to new users and their surroundings to reduce catastrophic forgetting in their personal devices.

5.2 Future Directions

In this section, we discuss two future directions that we believe are important to advance assistive vision technologies. First, we suggest future researchers in machine learning and computer vision to work with federated learning [285] applied to assistive vision devices to enhance the robustness in their image recognition performance. Secondly, we highlight the importance of involving visually impaired (VI) people throughout research processes to enhance the usefulness of the next generation of assistive vision technologies.

Federated Learning

In recent years, there has been a growing interest for federated learning [285–287] in the machine learning community. Federated learning involves leveraging the storage and computational capabilities of remote devices, such as mobile phones, for training local models while keeping the data on the device to preserve privacy.

Example Scenario. We give an example scenario in the context of assistive vision where Alice and Bob are grocery shopping in a supermarket using an assistive vision app installed on their personal smartphones. Alice needs to buy milk and wants to recognize a milk package held in her hand. She capture several frames of the milk package by recording a video with the camera that is processed in the app. However, the app responds that this milk package is an unknown item since this type of milk brand is new in the supermarket. Alice asks Charles the clerk what item she is holding, gets help with updating her app with the new milk class from Charles, and purchases the milk. The next day, Bob does his grocery shopping in the same supermarket as Alice. Bob is curious about a new milk package he has never seen before, which is of the same class as the milk Alice bought the day before. Charles app can successfully recognize the correct class of the milk package although this is the first time Bob takes photos of this milk with the app. This is enabled through a central server communicating with all mobile phones with the assistive app installed for learning a global model from local updates from the connected devices in the network. Thus, after the local update of the new milk package in Alice’s app has

been incorporated, the server transmits the new global model to Bob’s app. This is the standard procedure for federated learning which has the potential to utilize the computational resources for predictions on smartphones while maintaining the user experience and preserving privacy.

There are many interesting challenges that remains to tackle in federated learning. Li et al. [285] mentions challenges with statistical heterogeneity in the data from different users [285, 288, 289], variability of the systems on each device [290], privacy concerns [291, 292], and expensive communication rounds between the server and the devices [293]. From the perspective of this thesis, it would be valuable to study the statistical heterogeneity as the generated data can vary immensely across different users in terms of sample frequency, data quality, and context of where the data has been recorded. Here, it may be important to consider issues regarding fairness in addition to accuracy, since the learned model may become biased towards devices with large amounts of data or to commonly occurring groups of devices. Another relevant extension to this thesis would be on systems heterogeneity among different devices, since it will be important to push research and development towards image recognition methods that are accurate, robust, and performs in real-time regardless of the hardware, network connectivity, and power of the mobile devices. This could enable assistive vision apps to be less reliant on the mobile phone type, which in return may yield assistive vision apps that work equally well across different geographical locations and income levels.

Disability-First Approach in Development of Assistive Vision Technologies

This section is for encouraging future researchers in computer vision and machine learning to have their target user in mind for projects motivated by real-world applications such as assistive technologies. There exist many studies on computer vision-based assistive technologies describing the benefits and challenges to ease everyday life for VI people [23, 47, 294–297]. The past years, there has been an increasing interest of deploying computer vision and machine learning models on mobile devices [298, 299]. Additionally, it is well-known that VI people use mobile phones, including the camera to record events, and want help with object recognition and navigation [14, 22, 47, 51]. Even if the combination is evidently a good fit, many of the mentioned works stress that researchers should keep VI people involved throughout the process to ensure that the developed technology is requested and, in fact, useful for blind/low-vision people. An example of close collaboration between the two communities is the collection of the recent ORBIT dataset [15, 280] where the data collection was performed by blind/low-vision participants. They take on a *disability-first approach* in the process to produce good quality dataset that should primarily serve the VI community by continuously engaging with and supporting the participants to balance the usefulness versus the effort of data to be collected [280]. Involving VI users in the research process should provide benefits on for both the VI users as well as the vision community by bringing insights on how to build more robust image recognizers.

Aspects on usefulness, fairness, and privacy needs to be considered when developing assistive vision technologies. Regarding usefulness, it will be important to engage

with VI people and people from technical disciplines, such as human-computer interaction and user experience designers, to understand what VI people want their app to assist them with and how to build functionality to achieve those goals. For computer vision researchers, we have seen that the models suffer from biases that can disadvantage users based on their gender [34,300,301], geographical location [32], and co-occurrences between objects and surrounding contexts [302]. Such challenges on fairness [303–306] and various types of distribution shifts [307,308] will be crucial to tackle for establishing trust in image recognition systems for practitioners. Furthermore, we should strive for developing privacy-preserving systems to avoid unintended leakages of user-sensitive information [178,309–311]. However, we also need to take into account the point-of-view of bystanders and whether they are comfortable with being photographed by mobile and wearable devices. This aspect has been studied recently [24,312,313], however, it is still non-trivial how to approach users and ask them to stop using their device. This matter of privacy and social acceptance between assistive vision users and surrounding people needs to be further studied in realistic environments for the purpose of understanding the needs of those who will be involved in the usage of the technology.

Bibliography

- [1] Rupert Bourne, Jaimie D Steinmetz, Seth Flaxman, Paul Svitil Briant, Hugh R Taylor, Serge Resnikoff, Robert James Casson, Amir Abdoli, Eman Abu-Gharbieh, Ashkan Afshin, et al. Trends in prevalence of blindness and distance and near vision impairment over 30 years: an analysis for the global burden of disease study. *The Lancet global health*, 9(2):e130–e143, 2021.
- [2] Microsoft Corporation. Seeing AI, 2017. URL <https://www.microsoft.com/en-us/ai/seeing-ai>. Accessed 2022-03-14.
- [3] Patrick Clary. Lookout: an app to help blind and visually impaired people learn about their surroundings, 2018. URL <https://blog.google/outreach-initiatives/accessibility/lookout-app-help-blind-and-visually-impaired-people-learn-about-their-surroundings/>. Accessed 2022-04-05.
- [4] Inc Cloudsight. TapTapSee, 2013. URL <https://taptapseeapp.com/>. Accessed 2022-03-22.
- [5] Envision. Envision App, 2018. URL <https://www.letsenvision.com/envision-app>. Accessed 2022-03-22.
- [6] Be My Eyes. Be My Eyes, 2017. URL <https://www.bemyeyes.com/>. Accessed 2022-03-22.
- [7] Aira Tech Corp. Aira, 2017. URL <https://aira.io/>. Accessed 2022-03-14.
- [8] OrCam. OrCam MyEye 2, 2019. URL <https://www.orcam.com/sv/myeye2/>. Accessed 2022-03-22.
- [9] Envision. Envision Glasses, 2020. URL <https://www.letsenvision.com/envision-glasses>. Accessed 2022-03-22.
- [10] S. Caraiman, A. Morar, M. Owczarek, A. Burlacu, D. Rzeszotarski, N. Botezatu, P. Hergheliegiu, F. Moldoveanu, P. Strumillo, and A. Moldoveanu. Computer vision for the visually impaired: the sound of vision system. In *IEEE International Conference on Computer Vision Workshops*, 2017.
- [11] Dragan Ahmetovic, Daisuke Sato, Uran Oh, Tatsuya Ishihara, Kris Kitani, and Chieko Asakawa. Recog: Supporting blind people in recognizing personal objects. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–12, 2020.

- [12] Rabia Jafri, Syed Abid Ali, Hamid R Arabnia, and Shameem Fatima. Computer vision-based object recognition for the visually impaired in an indoors environment: a survey. *The Visual Computer*, 30(11):1197–1222, 2014.
- [13] Hernisa Kacorri. Teachable machines for accessibility. *ACM SIGACCESS Accessibility and Computing*, (119):10–18, 2017.
- [14] Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3608–3617, 2018.
- [15] Daniela Massiceti, Luisa Zintgraf, John Bronskill, Lida Theodorou, Matthew Tobias Harris, Edward Cutrell, Cecily Morrison, Katja Hofmann, and Simone Stumpf. Orbit: A real-world few-shot dataset for teachable object recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10818–10828, 2021.
- [16] Hernisa Kacorri, Kris M Kitani, Jeffrey P Bigham, and Chieko Asakawa. People with visual impairment training personal object recognizers: Feasibility and challenges. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, pages 5839–5849, 2017.
- [17] Kyungjun Lee and Hernisa Kacorri. Hands holding clues for object recognition in teachable machines. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–12, 2019.
- [18] Kyungjun Lee, Jonggi Hong, Simone Pimento, Ebrima Jarjue, and Hernisa Kacorri. Revisiting blind photography in the context of teachable object recognizers. In *The 21st International ACM SIGACCESS Conference on Computers and Accessibility*, pages 83–95, 2019.
- [19] James Coughlan and Roberto Manduchi. Functional assessment of a camera phone-based wayfinding system operated by blind and visually impaired users. *International Journal on Artificial Intelligence Tools*, 18(03):379–397, 2009.
- [20] Hernisa Kacorri, Eshed Ohn-Bar, Kris M Kitani, and Chieko Asakawa. Environmental factors in indoor navigation based on real-world trajectories of blind users. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pages 1–12, 2018.
- [21] Jack M Loomis, Reginald G Golledge, Roberta L Klatzky, and James R Marston. Assisting wayfinding in visually impaired travelers. In *Applied Spatial Cognition*, pages 179–202. Psychology Press, 2020.
- [22] Sarit Szpiro, Yuhang Zhao, and Shiri Azenkot. Finding a store, searching for a product: a study of daily challenges of low vision people. In *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, pages 61–72, 2016.

- [23] YingLi Tian, Xiaodong Yang, Chucai Yi, and Aries Arditi. Toward a computer vision-based wayfinding aid for blind persons to access unfamiliar indoor environments. *Machine vision and applications*, 24(3):521–535, 2013.
- [24] Kyungjun Lee, Daisuke Sato, Saki Asakawa, Hernisa Kacorri, and Chieko Asakawa. Pedestrian detection with wearable cameras for the blind: A two-way perspective. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–12, 2020.
- [25] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [27] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7263–7271, 2017.
- [28] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057. PMLR, 2015.
- [29] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018.
- [30] Michael A Alcorn, Qi Li, Zhitao Gong, Chengfei Wang, Long Mai, Wei-Shinn Ku, and Anh Nguyen. Strike (with) a pose: Neural networks are easily fooled by strange poses of familiar objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4845–4854, 2019.
- [31] Dengxin Dai and Luc Van Gool. Dark model adaptation: Semantic image segmentation from daytime to nighttime. In *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pages 3819–3824. IEEE, 2018.
- [32] Terrance De Vries, Ishan Misra, Changhan Wang, and Laurens Van der Maaten. Does object recognition work for everyone? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 52–59, 2019.
- [33] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*, pages 5637–5664. PMLR, 2021.

- [34] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR, 2018.
- [35] Richard Szeliski. *Computer vision: algorithms and applications*. Springer Science & Business Media, 2010.
- [36] Xiu-Shen Wei, Yi-Zhe Song, Oisin Mac Aodha, Jianxin Wu, Yuxin Peng, Jinhui Tang, Jian Yang, and Serge Belongie. Fine-grained image analysis with deep learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [37] Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [38] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- [39] Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. Generalizing from a few examples: A survey on few-shot learning. *ACM computing surveys (csur)*, 53(3):1–34, 2020.
- [40] Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. *IEEE transactions on pattern analysis and machine intelligence*, 41(9):2251–2265, 2018.
- [41] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *arXiv preprint arXiv:2004.05439*, 2020.
- [42] World Health Organization. International Classification of Diseases 11th Revision (ICD-11), 2022. URL <https://icd.who.int/en>. Accessed 2022-03-22.
- [43] World Health Organization. *World report on vision*. World Health Organization, 2019.
- [44] Jaimie D Steinmetz, Rupert RA Bourne, Paul Svitil Briant, Seth R Flaxman, Hugh RB Taylor, Jost B Jonas, Amir Aberhe Abdoli, Woldu Aberhe Abrha, Ahmed Abualhasan, Eman Girum Abu-Gharbieh, et al. Causes of blindness and vision impairment in 2020 and trends over 30 years, and prevalence of avoidable blindness in relation to vision 2020: the right to sight: an analysis for the global burden of disease study. *The Lancet Global Health*, 9(2):e144–e160, 2021.
- [45] Synskadades Riksförbund (SRF). Vardagstips för personer med synnedsättning, 2017. URL https://www.srf.nu/globalassets/informationsmaterial/for-medlemmar-och-andra-synskadade/vardagstips_2017_tillganglig.pdf. The Swedish Association of the Visually Impaired (SRF).

- [46] Henrik Götesson. Challenges and possibilities when using screen readers, 2019. URL https://www.srf.nu/contentassets/d26ce8ef340141768331b0fbc6e20f59/screenreader_survey_english_summary.pdf. The Swedish Association of the Visually Impaired (SRF).
- [47] Chandrika Jayant, Hanjie Ji, Samuel White, and Jeffrey P Bigham. Supporting blind photography. In *The proceedings of the 13th international ACM SIGACCESS conference on Computers and accessibility*, pages 203–210, 2011.
- [48] Erin Brady, Meredith Ringel Morris, Yu Zhong, Samuel White, and Jeffrey P Bigham. Visual challenges in the everyday lives of blind people. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 2117–2126, 2013.
- [49] Yuhang Zhao, Shaomei Wu, Lindsay Reynolds, and Shiri Azenkot. A face recognition application for people with visual impairments: Understanding use beyond the lab. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2018.
- [50] Xu Liu. A camera phone based currency reader for the visually impaired. In *Proceedings of the 10th international ACM SIGACCESS conference on Computers and accessibility*, pages 305–306, 2008.
- [51] Marynel Vázquez and Aaron Steinfeld. Helping visually impaired users properly aim a camera. In *Proceedings of the 14th international ACM SIGACCESS conference on Computers and accessibility*, pages 95–102, 2012.
- [52] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [53] Justin Johnson, Andrej Karpathy, and Li Fei-Fei. Densecap: Fully convolutional localization networks for dense captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4565–4574, 2016.
- [54] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015.
- [55] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019.
- [56] Ronghang Hu, Anna Rohrbach, Trevor Darrell, and Kate Saenko. Language-conditioned graph networks for relational reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10294–10303, 2019.

- [57] Patrick E Lanigan, Aaron M Paulos, Andrew W Williams, Dan Rossi, and Priya Narasimhan. Trinetra: Assistive technologies for grocery shopping for the blind. In *2006 10th IEEE International Symposium on Wearable Computers*, pages 147–148. IEEE, 2006.
- [58] Tess Winlock, Eric Christiansen, and Serge Belongie. Toward real-time grocery detection for the visually impaired. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition-Workshops*, pages 49–56. IEEE, 2010.
- [59] Joan Sosa-García and Francesca Odone. “hands on” visual recognition for visually impaired users. *ACM Transactions on Accessible Computing (TACCESS)*, 10(3):1–30, 2017.
- [60] Roger Boldu, Denys JC Matthies, Haimo Zhang, and Suranga Nanayakkara. Aisee: an assistive wearable device to support visually impaired grocery shoppers. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 4(4):1–25, 2020.
- [61] Peter A Zientara, Sooyeon Lee, Gus H Smith, Rorry Brenner, Laurent Itti, Mary B Rosson, John M Carroll, Kevin M Irick, and Vijaykrishnan Narayanan. Third eye: A shopping assistant for the visually impaired. *Computer*, 50(2):16–24, 2017.
- [62] Marian George, Dejan Mircic, Gabor Soros, Christian Floerkemeier, and Friedemann Mattern. Fine-grained product class recognition for assisted shopping. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 154–162, 2015.
- [63] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989.
- [64] Tom M. Mitchell. Machine learning and data mining. *Communications of the ACM*, 42:30–36, 1999.
- [65] Kim Hazelwood, Sarah Bird, David Brooks, Soumith Chintala, Utku Diril, Dmytro Dzhulgakov, Mohamed Fawzy, Bill Jia, Yangqing Jia, Aditya Kalro, et al. Applied machine learning at facebook: A datacenter infrastructure perspective. In *2018 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pages 620–629. IEEE, 2018.
- [66] Peter Bailis, Edward Gan, Samuel Madden, Deepak Narayanan, Kexin Rong, and Sahaana Suri. Macrobase: Prioritizing attention in fast data. In *Proceedings of the 2017 ACM International Conference on Management of Data*, pages 541–556, 2017.
- [67] Frank N Dempster. Spacing effects and their implications for theory and practice. *Educational Psychology Review*, 1(4):309–330, 1989.

- [68] Hermann Ebbinghaus. Memory: A contribution to experimental psychology. *Annals of neurosciences*, 20(4):155, 2013.
- [69] Karri S Hawley, Katie E Cherry, Emily O Boudreux, and Erin M Jackson. A comparison of adjusted spaced retrieval versus a uniform expanded retrieval schedule for learning a name–face association in older adults with probable alzheimer’s disease. *Journal of Clinical and Experimental Neuropsychology*, 30(6):639–649, 2008.
- [70] T. Landauer and Robert Bjork. Optimum rehearsal patterns and name learning. *Practical aspects of memory*, 1, 11 1977.
- [71] Paul Smolen, Yili Zhang, and John H Byrne. The right time to learn: mechanisms and optimization of spaced learning. *Nature Reviews Neuroscience*, 17(2):77, 2016.
- [72] Weiran Wang, Xincheng Yan, Honglak Lee, and Karen Livescu. Deep variational canonical correlation analysis. *arXiv preprint arXiv:1610.03454*, 2016.
- [73] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [74] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [75] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [76] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [77] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [78] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- [79] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- [80] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- [81] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.

- [82] Shibani Santurkar, Dimitris Tsipras, Andrew Ilyas, and Aleksander Madry. How does batch normalization help optimization? *Advances in neural information processing systems*, 31, 2018.
- [83] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [84] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond. *arXiv preprint arXiv:1908.03265*, 2019.
- [85] Tijmen Tieleman, Geoffrey Hinton, et al. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.
- [86] Leslie N Smith. Cyclical learning rates for training neural networks. In *2017 IEEE winter conference on applications of computer vision (WACV)*, pages 464–472. IEEE, 2017.
- [87] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.
- [88] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [89] Frank Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.
- [90] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [91] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [92] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [93] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [94] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.

- [95] Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications. *AI Open*, 1:57–81, 2020.
- [96] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103, 2008.
- [97] Andrew Ng et al. Sparse autoencoder.
- [98] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):423–443, 2018.
- [99] Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 631–648, 2018.
- [100] Jiquan Ngiam, Aditya Khosla, Mingyu Kim, Juhan Nam, Honglak Lee, and Andrew Y Ng. Multimodal deep learning. In *ICML*, 2011.
- [101] Carina Silberer and Mirella Lapata. Learning grounded meaning representations with autoencoders. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 721–732, 2014.
- [102] Michelle A Lee, Yuke Zhu, Peter Zachares, Matthew Tan, Krishnan Srinivasan, Silvio Savarese, Li Fei-Fei, Animesh Garg, and Jeannette Bohg. Making sense of vision and touch: Learning multimodal representations for contact-rich tasks. *IEEE Transactions on Robotics*, 36(3):582–596, 2020.
- [103] Diederik P Kingma and Max Welling. An introduction to variational autoencoders. *arXiv preprint arXiv:1906.02691*, 2019.
- [104] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. 2009.
- [105] Cheng Zhang, Judith Bütepage, Hedvig Kjellström, and Stephan Mandt. Advances in variational inference. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):2008–2026, 2018.
- [106] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- [107] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pages 1278–1286. PMLR, 2014.

- [108] Mike Wu and Noah Goodman. Multimodal generative models for scalable weakly-supervised learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- [109] Yuge Shi, Brooks Paige, Philip Torr, et al. Variational mixture-of-experts autoencoders for multi-modal deep generative models. *Advances in Neural Information Processing Systems*, 32, 2019.
- [110] Ramakrishna Vedantam, Ian Fischer, Jonathan Huang, and Kevin Murphy. Generative models of visually grounded imagination. *arXiv preprint arXiv:1705.10762*, 2017.
- [111] Masahiro Suzuki, Kotaro Nakayama, and Yutaka Matsuo. Joint multimodal learning with deep generative models. *arXiv preprint arXiv:1611.01891*, 2016.
- [112] Richard Bellman. A markovian decision process. *Journal of Mathematics and Mechanics*, 6(5):679–684, 1957.
- [113] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [114] Long-Ji Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine learning*, 8(3-4):293–321, 1992.
- [115] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
- [116] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In *International conference on machine learning*, pages 1995–2003. PMLR, 2016.
- [117] Vijay Konda and John Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.
- [118] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.
- [119] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [120] Verena R Cimarolli, Kathrin Boerner, Mark Brennan-Ing, Joann P Reinhardt, and Amy Horowitz. Challenges faced by older adults with vision loss: a qualitative study with implications for rehabilitation. *Clinical rehabilitation*, 26(8):748–757, 2012.

- [121] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8769–8778, 2018.
- [122] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *4th International IEEE Workshop on 3D Representation and Recognition (3dRR-13)*, Sydney, Australia, 2013.
- [123] Yushan Feng, Saihui Hou, and Zilei Wang. Vegfru: A domain-specific dataset for fine-grained visual categorization. In *IEEE International Conference on Computer Vision*, 2017.
- [124] Xiu-Shen Wei, Quan Cui, Lei Yang, Peng Wang, and Lingqiao Liu. Rpc: A large-scale retail product checkout dataset. *arXiv preprint arXiv:1901.07249*, 2019.
- [125] Xiaopeng Zhang, Hongkai Xiong, Wengang Zhou, Weiyao Lin, and Qi Tian. Picking deep filter responses for fine-grained image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1134–1142, 2016.
- [126] Yaming Wang, Vlad I Morariu, and Larry S Davis. Learning a discriminative filter bank within a cnn for fine-grained recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4148–4157, 2018.
- [127] Yao Ding, Yanzhao Zhou, Yi Zhu, Qixiang Ye, and Jianbin Jiao. Selective sparse sampling for fine-grained image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6599–6608, 2019.
- [128] Steve Branson, Grant Van Horn, Serge Belongie, and Pietro Perona. Bird species categorization using pose normalized deep convolutional nets. *arXiv preprint arXiv:1406.2952*, 2014.
- [129] Ning Zhang, Jeff Donahue, Ross Girshick, and Trevor Darrell. Part-based rcnn for fine-grained category detection. In *European conference on computer vision*, pages 834–849. Springer, 2014.
- [130] Yu Zhang, Xiu-Shen Wei, Jianxin Wu, Jianfei Cai, Jiangbo Lu, Viet-Anh Nguyen, and Minh N Do. Weakly supervised fine-grained categorization with part-based image representation. *IEEE Transactions on Image Processing*, 25(4):1713–1725, 2016.
- [131] Jianlong Fu, Heliang Zheng, and Tao Mei. Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4438–4446, 2017.
- [132] Heliang Zheng, Jianlong Fu, Tao Mei, and Jiebo Luo. Learning multi-attention convolutional neural network for fine-grained image recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 5209–5217, 2017.

- [133] Ming Sun, Yuchen Yuan, Feng Zhou, and Errui Ding. Multi-attention multi-class constraint for fine-grained image recognition. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 805–821, 2018.
- [134] Tsung-Yu Lin, Aruni RoyChowdhury, and Subhransu Maji. Bilinear cnn models for fine-grained visual recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 1449–1457, 2015.
- [135] Yin Cui, Feng Zhou, Jiang Wang, Xiao Liu, Yuanqing Lin, and Serge Belongie. Kernel pooling for convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2930, 2017.
- [136] Abhimanyu Dubey, Otkrist Gupta, Ramesh Raskar, and Nikhil Naik. Maximum-entropy fine grained classification. *Advances in neural information processing systems*, 31, 2018.
- [137] Abhimanyu Dubey, Otkrist Gupta, Pei Guo, Ramesh Raskar, Ryan Farrell, and Nikhil Naik. Pairwise confusion for fine-grained visual classification. In *Proceedings of the European conference on computer vision (ECCV)*, pages 70–86, 2018.
- [138] Dongliang Chang, Yifeng Ding, Jiyang Xie, Ayan Kumar Bhunia, Xiaoxu Li, Zhanyu Ma, Ming Wu, Jun Guo, and Yi-Zhe Song. The devil is in the channels: Mutual-channel loss for fine-grained image classification. *IEEE Transactions on Image Processing*, 29:4683–4695, 2020.
- [139] Jonathan Krause, Benjamin Sapp, Andrew Howard, Howard Zhou, Alexander Toshev, Tom Duerig, James Philbin, and Li Fei-Fei. The unreasonable effectiveness of noisy data for fine-grained recognition. In *European Conference on Computer Vision*, pages 301–320. Springer, 2016.
- [140] Zhe Xu, Shaoli Huang, Ya Zhang, and Dacheng Tao. Webly-supervised fine-grained visual categorization via deep domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 40(5):1100–1113, 2016.
- [141] Xiaoxiao Sun, Liyi Chen, and Jufeng Yang. Learning from web data using adversarial discriminative neural networks for fine-grained classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 273–280, 2019.
- [142] Xiangteng He and Yuxin Peng. Fine-grained image classification via combining vision and language. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5994–6002, 2017.
- [143] Hua Zhang, Xiaochun Cao, and Rui Wang. Audio visual attribute discovery for fine-grained object recognition. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [144] Tianshui Chen, Liang Lin, Riquan Chen, Yang Wu, and Xiaonan Luo. Knowledge-embedded representation learning for fine-grained image recognition. *arXiv preprint arXiv:1807.00505*, 2018.

- [145] Scott Reed, Zeynep Akata, Honglak Lee, and Bernt Schiele. Learning deep representations of fine-grained visual descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 49–58, 2016.
- [146] Li Niu, Ashok Veeraraghavan, and Ashutosh Sabharwal. Webly supervised learning meets zero-shot learning: A hybrid approach for fine-grained classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7171–7180, 2018.
- [147] Zeynep Akata, Scott Reed, Daniel Walter, Honglak Lee, and Bernt Schiele. Evaluation of output embeddings for fine-grained image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2927–2936, 2015.
- [148] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [149] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pages 722–729. IEEE, 2008.
- [150] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. Technical report, 2011.
- [151] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [152] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [153] Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [154] Ivan Krasin, Tom Duerig, Neil Alldrin, Vittorio Ferrari, Sami Abu-El-Haija, Alina Kuznetsova, Hassan Rom, Jasper Uijlings, Stefan Popov, Andreas Veit, Serge Belongie, Victor Gomes, Abhinav Gupta, Chen Sun, Gal Chechik, David Cai, Zheyun Feng, Dhyanesh Narayanan, and Kevin Murphy. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from <https://github.com/openimages>*, 2017.
- [155] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.

- [156] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123(1):32–73, 2017.
- [157] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5356–5364, 2019.
- [158] Grant Van Horn, Elijah Cole, Sara Beery, Kimberly Wilber, Serge Belongie, and Oisin Mac Aodha. Benchmarking representation learning for natural world image collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12884–12893, 2021.
- [159] Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Li Fei-Fei. Novel dataset for fine-grained image categorization. In *First Workshop on Fine-Grained Visual Categorization, IEEE Conference on Computer Vision and Pattern Recognition*, Colorado Springs, CO, June 2011.
- [160] Grant Van Horn, Steve Branson, Ryan Farrell, Scott Haber, Jessie Barry, Panos Ipeirotis, Pietro Perona, and Serge Belongie. Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 595–604, 2015.
- [161] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2012.
- [162] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *2009 IEEE conference on computer vision and pattern recognition*, pages 951–958. IEEE, 2009.
- [163] Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaou Tang. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [164] Menglin Jia, Mengyun Shi, Mikhail Sirotenko, Yin Cui, Claire Cardie, Bharath Hariharan, Hartwig Adam, and Serge Belongie. Fashionpedia: Ontology, segmentation, and an attribute localization dataset. In *European conference on computer vision*, pages 316–332. Springer, 2020.
- [165] Andrea Vedaldi, Siddharth Mahendran, Stavros Tsogkas, Subhransu Maji, Ross Girshick, Juho Kannala, Esa Rahtu, Iasonas Kokkinos, Matthew B Blaschko, David Weiss, et al. Understanding objects in detail with fine-grained attributes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3622–3629, 2014.

- [166] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- [167] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 554–561, 2013.
- [168] Linjie Yang, Ping Luo, Chen Change Loy, and Xiaoou Tang. A large-scale car dataset for fine-grained categorization and verification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3973–3981, 2015.
- [169] Timnit Gebru, Jonathan Krause, Yilun Wang, Duyun Chen, Jia Deng, and Li Fei-Fei. Fine-grained car detection for visual census estimation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- [170] Omkar M Parkhi, Andrea Vedaldi, and Andrew Zisserman. Deep face recognition. 2015.
- [171] Gary B Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. In *Workshop on faces in'Real-Life'Images: detection, alignment, and recognition*, 2008.
- [172] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. Vgface2: A dataset for recognising faces across pose and age. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*, pages 67–74. IEEE, 2018.
- [173] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *European conference on computer vision*, pages 87–102. Springer, 2016.
- [174] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101-mining discriminative components with random forests. In *European conference on computer vision*, pages 446–461. Springer, 2014.
- [175] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Attribute-based classification for zero-shot learning of object categories. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (01):1–1, 2013.
- [176] Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, et al. Meta-dataset: A dataset of datasets for learning to learn from few examples. *arXiv preprint arXiv:1903.03096*, 2019.
- [177] Sebastian Bujwid and Josephine Sullivan. Large-scale zero-shot image classification from rich and diverse textual descriptions. *arXiv preprint arXiv:2103.09669*, 2021.

- [178] Danna Gurari, Qing Li, Chi Lin, Yinan Zhao, Anhong Guo, Abigale Stangl, and Jeffrey P Bigham. Vizwiz-priv: A dataset for recognizing the presence and purpose of private visual information in images taken by blind people. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 939–948, 2019.
- [179] Philipp Jund, Nichola Abdo, Andreas Eitel, and Wolfram Burgard. The freiburg groceries dataset. *arXiv preprint arXiv:1611.05799*, 2016.
- [180] Georg Waltner, Michael Schwarz, Stefan Ladstätter, Anna Weber, Patrick Luley, Horst Bischof, Meinrad Lindschinger, Irene Schmid, and Lucas Paletta. Mango - mobile augmented reality with functional eating guidance and food awareness. In *International Workshop on Multimedia Assisted Dietary Management*, 2015.
- [181] Marian George and Christian Floerkemeier. Recognizing products: A per-exemplar multi-label image classification approach. In *European Conference on Computer Vision*, pages 440–455. Springer, 2014.
- [182] Michele Merler, Carolina Galleguillos, and Serge Belongie. Recognizing groceries in situ using in vitro training data. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2007.
- [183] Weidong Geng, Feilin Han, Jiangke Lin, Liuyi Zhu, Jieming Bai, Suzhen Wang, Lin He, Qiang Xiao, and Zhangjiong Lai. Fine-grained grocery product recognition by one-shot learning. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 1706–1714, 2018.
- [184] Y. Kawano and K. Yanai. Automatic expansion of a food image dataset leveraging existing categories with domain adaptation. In *Proc. of ECCV Workshop on Transferring and Adapting Source Knowledge in Computer Vision (TASK-CV)*, 2014.
- [185] Weiqing Min, Linhu Liu, Zhengdong Luo, and Shuqiang Jiang. Ingredient-guided cascaded multi-attention network for food recognition. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 1331–1339, 2019.
- [186] Jaclyn Rich, Hamed Haddadi, and Timothy M Hospedales. Towards bottom-up analysis of social food. In *Proceedings of the 6th International Conference on Digital Health Conference*, pages 111–120, 2016.
- [187] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. Scaling egocentric vision: The epic-kitchens dataset. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 720–736, 2018.
- [188] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, , Antonino Furnari, Jian Ma, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100. *International Journal of Computer Vision (IJCV)*, 2021.

- [189] Javier Marin, Aritro Biswas, Ferda Ofli, Nicholas Hynes, Amaia Salvador, Yusuf Aytar, Ingmar Weber, and Antonio Torralba. Recipe1m+: A dataset for learning cross-modal embeddings for cooking recipes and food images. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2019.
- [190] Amaia Salvador, Nicholas Hynes, Yusuf Aytar, Javier Marin, Ferda Ofli, Ingmar Weber, and Antonio Torralba. Learning cross-modal embeddings for cooking recipes and food images. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [191] Semih Yagcioglu, Aykut Erdem, Erkut Erdem, and Nazli Ikizler-Cinbis. Recipeqa: A challenge dataset for multimodal comprehension of cooking recipes. *arXiv preprint arXiv:1809.00812*, 2018.
- [192] Oscar Beijbom, Neel Joshi, Dan Morris, Scott Saponas, and Siddharth Khullar. Menu-match: Restaurant-specific food logging from images. In *2015 IEEE Winter Conference on Applications of Computer Vision*, pages 844–851. IEEE, 2015.
- [193] Ruihan Xu, Luis Herranz, Shuqiang Jiang, Shuang Wang, Xinhang Song, and Ramesh Jain. Geolocalized modeling for dish recognition. *IEEE transactions on multimedia*, 17(8):1187–1199, 2015.
- [194] Chang Xu, Dacheng Tao, and Chao Xu. A survey on multi-view learning. *arXiv preprint arXiv:1304.5634*, 2013.
- [195] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Neural baby talk. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7219–7228, 2018.
- [196] Jyoti Aneja, Harsh Agrawal, Dhruv Batra, and Alexander Schwing. Sequential latent spaces for modeling the intention during diverse image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4261–4270, 2019.
- [197] Drew A Hudson and Christopher D Manning. Compositional attention networks for machine reasoning. *arXiv preprint arXiv:1803.03067*, 2018.
- [198] Mathieu Salzmann, Carl Henrik Ek, Raquel Urtasun, and Trevor Darrell. Factorized orthogonal latent spaces. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 701–708. JMLR Workshop and Conference Proceedings, 2010.
- [199] Harold Hotelling. Relations between two sets of variates. *Biometrika*, 28(3/4):321–377, 1936.
- [200] Galen Andrew, Raman Arora, Jeff Bilmes, and Karen Livescu. Deep canonical correlation analysis. In *International conference on machine learning*, pages 1247–1255, 2013.
- [201] Weiran Wang, Raman Arora, Karen Livescu, and Jeff Bilmes. On deep multi-view representation learning. In *International conference on machine learning*, pages 1083–1092, 2015.

- [202] Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. *Neural networks*, 13(4-5):411–430, 2000.
- [203] Yao-Hung Hubert Tsai, Paul Pu Liang, Amir Zadeh, Louis-Philippe Morency, and Ruslan Salakhutdinov. Learning factorized multimodal representations. In *International Conference on Learning Representations*, 2019.
- [204] Cheng Zhang, Hedvig Kjellström, and Carl Henrik Ek. Inter-battery topic representation learning. In *European Conference on Computer Vision*, pages 210–226. Springer, 2016.
- [205] Andreas Damianou, Neil D Lawrence, and Carl Henrik Ek. Multi-view learning as a nonparametric nonlinear inter-battery factor analysis. *Journal of Machine Learning Research*, 22(86):1–51, 2021.
- [206] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2020.
- [207] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 806–813, 2014.
- [208] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- [209] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013.
- [210] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [211] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- [212] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*, 2015.
- [213] Jason Wei and Kai Zou. Eda: Easy data augmentation techniques for boosting performance on text classification tasks. *arXiv preprint arXiv:1901.11196*, 2019.
- [214] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K Dokania, Philip HS Torr, and Marc’Aurelio Ranzato. On tiny episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*, 2019.

- [215] David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. *arXiv preprint arXiv:1706.08840*, 2017.
- [216] Tyler L Hayes, Kushal Kafle, Robik Shrestha, Manoj Acharya, and Christopher Kanan. Remind your neural network to prevent catastrophic forgetting. In *European Conference on Computer Vision*, pages 466–483. Springer, 2020.
- [217] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [218] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*, pages 3987–3995. PMLR, 2017.
- [219] Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational continual learning. *arXiv preprint arXiv:1710.10628*, 2017.
- [220] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2(7), 2015.
- [221] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- [222] Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual learning. In *International Conference on Machine Learning*, pages 4528–4537. PMLR, 2018.
- [223] Gobinda Saha, Isha Garg, and Kaushik Roy. Gradient projection memory for continual learning. *arXiv preprint arXiv:2103.09762*, 2021.
- [224] Ta-Chu Kao, Kristopher Jensen, Gido van de Ven, Alberto Bernacchia, and Guillaume Hennequin. Natural continual learning: success is a journey, not (just) a destination. *Advances in Neural Information Processing Systems*, 34, 2021.
- [225] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- [226] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically expandable networks. *arXiv preprint arXiv:1708.01547*, 2017.
- [227] Jaehong Yoon, Saehoon Kim, Eunho Yang, and Sung Ju Hwang. Scalable and order-robust continual learning with additive parameter decomposition. *arXiv preprint arXiv:1902.09432*, 2019.
- [228] Sayna Ebrahimi, Franziska Meier, Roberto Calandra, Trevor Darrell, and Marcus Rohrbach. Adversarial continual learning. *arXiv preprint arXiv:2003.09553*, 2020.

- [229] Ju Xu and Zhanxing Zhu. Reinforced continual learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- [230] Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7765–7773, 2018.
- [231] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *International Conference on Machine Learning*, pages 4548–4557. PMLR, 2018.
- [232] Jonathan Schwarz, Siddhant Jayakumar, Razvan Pascanu, Peter E Latham, and Yee Teh. Powerpropagation: A sparsity inducing weight reparameterisation. *Advances in Neural Information Processing Systems*, 34:28889–28903, 2021.
- [233] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy P Lillicrap, and Greg Wayne. Experience replay for continual learning. *arXiv preprint arXiv:1811.11682*, 2018.
- [234] Xisen Jin, Arka Sadhu, Junyi Du, and Xiang Ren. Gradient based memory editing for task-free continual learning. *arXiv preprint arXiv:2006.15294*, 2020.
- [235] Eli Verwimp, Matthias De Lange, and Tinne Tuytelaars. Rehearsal revealed: The limits and merits of revisiting samples in continual learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9385–9394, 2021.
- [236] Gido M van de Ven and Andreas S Tolias. Generative replay with feedback connections as a general strategy for continual learning. *arXiv preprint arXiv:1809.10635*, 2018.
- [237] Gido M van de Ven, Hava T Siegelmann, and Andreas S Tolias. Brain-inspired replay for continual learning with artificial neural networks. *Nature communications*, 11(1):1–14, 2020.
- [238] Chenshen Wu, Luis Herranz, Xialei Liu, Joost van de Weijer, Bogdan Raducanu, et al. Memory replay gans: Learning to generate new categories without forgetting. *Advances in Neural Information Processing Systems*, 31, 2018.
- [239] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *arXiv preprint arXiv:2004.07211*, 2020.
- [240] Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In *Computer Vision – ECCV 2020*, pages 86–102. Springer International Publishing, 2020.
- [241] KJ Joseph and Vineeth N Balasubramanian. Meta-consolidation for continual learning. *arXiv preprint arXiv:2010.00352*, 2020.

- [242] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Dilan Gorur, Razvan Pascanu, and Hassan Ghasemzadeh. Linear mode connectivity in multitask and continual learning. *arXiv preprint arXiv:2010.04495*, 2020.
- [243] Pingbo Pan, Siddharth Swaroop, Alexander Immer, Runa Eschenhagen, Richard E Turner, and Mohammad Emtiyaz Khan. Continual deep learning by functional regularisation of memorable past. *arXiv preprint arXiv:2004.14070*, 2020.
- [244] Johannes von Oswald, Christian Henning, João Sacramento, and Benjamin F Grewe. Continual learning with hypernetworks. *arXiv preprint arXiv:1906.00695*, 2019.
- [245] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. *arXiv preprint arXiv:1903.08671*, 2019.
- [246] Zalán Borsos, Mojmir Mutný, and Andreas Krause. Coresets via bilevel optimization for continual learning and streaming. *arXiv preprint arXiv:2006.03875*, 2020.
- [247] Aristotelis Chrysakis and Marie-Francine Moens. Online continual learning from imbalanced data. In *ICML*, 2020.
- [248] Tyler L Hayes, Nathan D Cahill, and Christopher Kanan. Memory efficient experience replay for streaming learning. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 9769–9776. IEEE, 2019.
- [249] David Isele and Akansel Cosgun. Selective experience replay for lifelong learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [250] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010, 2017.
- [251] Jaehong Yoon, Divyam Madaan, Eunho Yang, and Sung Ju Hwang. Online coresnet selection for rehearsal-based continual learning. *arXiv preprint arXiv:2106.01085*, 2021.
- [252] Jeffrey S Vitter. Random sampling with a reservoir. *ACM Transactions on Mathematical Software (TOMS)*, 11(1):37–57, 1985.
- [253] Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 532–547, 2018.
- [254] Ahmet Iscen, Jeffrey Zhang, Svetlana Lazebnik, and Cordelia Schmid. Memory-efficient incremental learning through feature adaptation. In *European Conference on Computer Vision*, pages 699–715. Springer, 2020.

- [255] Lorenzo Pellegrini, Gabriele Graffetti, Vincenzo Lomonaco, and Davide Maltoni. Latent replay for real-time continual learning. *arXiv preprint arXiv:1912.01100*, 2019.
- [256] Jihwan Bang, Heesu Kim, YoungJoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow memory: Continual learning with a memory of diverse samples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8218–8227, 2021.
- [257] Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co2l: Contrastive continual learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9516–9525, 2021.
- [258] Zheda Mai, Ruiwen Li, Hyunwoo Kim, and Scott Sanner. Supervised contrastive replay: Revisiting the nearest class mean classifier in online class-incremental continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3589–3599, 2021.
- [259] Rahaf Aljundi, Eugene Belilovsky, Tinne Tuytelaars, Laurent Charlin, Massimo Caccia, Min Lin, and Lucas Page-Caccia. Online continual learning with maximal interfered retrieval. *Advances in neural information processing systems*, 32, 2019.
- [260] Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *International conference on computers and games*, pages 72–83. Springer, 2006.
- [261] Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in games*, 4(1):1–43, 2012.
- [262] Muhammad Umar Chaudhry and Jee-Hyong Lee. Feature selection for high dimensional data using monte carlo tree search. *IEEE Access*, 6:76036–76048, 2018.
- [263] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *European conference on machine learning*, pages 282–293. Springer, 2006.
- [264] Arslan Chaudhry, Albert Gordo, Puneet Dokania, Philip Torr, and David Lopez-Paz. Using hindsight to anchor past knowledge in continual learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 6993–7001, 2021.
- [265] Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauro. Learning to learn without forgetting by maximizing transfer and minimizing interference. *arXiv preprint arXiv:1810.11910*, 2018.
- [266] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

- [267] Yaroslav Bulatov. The notMNIST dataset. <http://yaroslavvb.com/upload/notMNIST/>, 2011.
- [268] Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*, 2013.
- [269] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. *arXiv preprint arXiv:1606.04080*, 2016.
- [270] Tameem Adel, Han Zhao, and Richard E Turner. Continual learning with adaptive weights (claw). *arXiv preprint arXiv:1911.09514*, 2019.
- [271] Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. *arXiv preprint arXiv:1904.07734*, 2019.
- [272] Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. *arXiv preprint arXiv:1812.00420*, 2018.
- [273] Maximilian Igl, Kamil Ciosek, Yingzhen Li, Sebastian Tschiatschek, Cheng Zhang, Sam Devlin, and Katja Hofmann. Generalization in reinforcement learning with selective noise injection and information bottleneck. *Advances in neural information processing systems*, 32, 2019.
- [274] Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A survey of generalisation in deep reinforcement learning. *arXiv preprint arXiv:2111.09794*, 2021.
- [275] Amy Zhang, Nicolas Ballas, and Joelle Pineau. A dissection of overfitting and generalization in continuous reinforcement learning. *arXiv preprint arXiv:1806.07937*, 2018.
- [276] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [277] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [278] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- [279] Tyler L Hayes and Christopher Kanan. Online continual learning for embedded devices. *arXiv preprint arXiv:2203.10681*, 2022.

- [280] Lida Theodorou, Daniela Massiceti, Luisa Zintgraf, Simone Stumpf, Cecily Morrison, Edward Cutrell, Matthew Tobias Harris, and Katja Hofmann. Disability-first dataset creation: Lessons from constructing a dataset for teachable object recognition with blind and low vision data collectors. In *The 23rd International ACM SIGACCESS Conference on Computers and Accessibility*, pages 1–12, 2021.
- [281] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016.
- [282] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- [283] Nathan Bronson, Thomas Lento, and Janet L Wiener. Open data challenges at facebook. In *2015 IEEE 31st international conference on data engineering*, pages 1516–1519. IEEE, 2015.
- [284] Anthony Asta. Observability at twitter: technical overview, part i, 2016. URL https://blog.twitter.com/engineering/en_us/a/2016/observability-at-twitter-technical-overview-part-i. Accessed 2022-05-04.
- [285] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020.
- [286] Jakub Konečný, Brendan McMahan, and Daniel Ramage. Federated optimization: Distributed optimization beyond the datacenter. *arXiv preprint arXiv:1511.03575*, 2015.
- [287] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- [288] Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S Talwalkar. Federated multi-task learning. *Advances in neural information processing systems*, 30, 2017.
- [289] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.
- [290] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingeman, Vladimir Ivanov, Chloe Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMahan, et al. Towards federated learning at scale: System design. *Proceedings of Machine Learning and Systems*, 1:374–388, 2019.

- [291] Robin C Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A client level perspective. *arXiv preprint arXiv:1712.07557*, 2017.
- [292] Abhishek Bhowmick, John Duchi, Julien Freudiger, Gaurav Kapoor, and Ryan Rogers. Protection against reconstruction and its applications in private federated learning. *arXiv preprint arXiv:1812.00984*, 2018.
- [293] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- [294] Juan R Terven, Joaquín Salas, and Bogdan Raducanu. New opportunities for computer vision-based assistive technology systems for the visually impaired. *Computer*, 47(4):52–58, 2013.
- [295] Marco Leo, G Medioni, M Trivedi, Takeo Kanade, and Giovanni Maria Farinella. Computer vision for assistive technologies. *Computer Vision and Image Understanding*, 154:1–15, 2017.
- [296] Ruxandra Tapu, Bogdan Mocanu, and Titus Zaharia. Wearable assistive devices for visually impaired: A state of the art survey. *Pattern Recognition Letters*, 137:37–52, 2020.
- [297] Linda Wang and Alexander Wong. Implications of computer vision driven assistive technologies towards individuals with visual impairment. *arXiv preprint arXiv:1905.07844*, 2019.
- [298] Chaoyun Zhang, Paul Patras, and Hamed Haddadi. Deep learning in mobile and wireless networking: A survey. *IEEE Communications surveys & tutorials*, 21(3):2224–2287, 2019.
- [299] Wei Yang Bryan Lim, Nguyen Cong Luong, Dinh Thai Hoang, Yutao Jiao, Ying-Chang Liang, Qiang Yang, Dusit Niyato, and Chunyan Miao. Federated learning in mobile edge networks: A comprehensive survey. *IEEE Communications Surveys & Tutorials*, 22(3):2031–2063, 2020.
- [300] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29, 2016.
- [301] Kaylee Burns, Lisa Anne Hendricks, Kate Saenko, Trevor Darrell, and Anna Rohrbach. Women also snowboard: Overcoming bias in captioning models. *arXiv preprint arXiv:1803.09797*, 2018.
- [302] Krishna Kumar Singh, Dhruv Mahajan, Kristen Grauman, Yong Jae Lee, Matt Feiszli, and Deepti Ghadiyaram. Don’t judge an object by its context: learning to overcome contextual bias. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11070–11078, 2020.
- [303] Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning*. fairmlbook.org, 2019. <http://www.fairmlbook.org>.

- [304] Kenneth Holstein, Jennifer Wortman Vaughan, Hal Daumé III, Miro Dudik, and Hanna Wallach. Improving fairness in machine learning systems: What do industry practitioners need? In *Proceedings of the 2019 CHI conference on human factors in computing systems*, pages 1–16, 2019.
- [305] Kaiyu Yang, Klint Qinami, Li Fei-Fei, Jia Deng, and Olga Russakovsky. Towards fairer datasets: Filtering and balancing the distribution of the people subtree in the imagenet hierarchy. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 547–558, 2020.
- [306] Priya Goyal, Adriana Romero Soriano, Caner Hazirbas, Levent Sagun, and Nicolas Usunier. Fairness indicators for systematic assessments of visual feature extractors. *arXiv preprint arXiv:2202.07603*, 2022.
- [307] Joaquin Quionero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. Dataset shift in machine learning, 2009.
- [308] Olivia Wiles, Sven Gowal, Florian Stimberg, Sylvestre Alvise-Rebuffi, Ira Ktena, Taylan Cemgil, et al. A fine-grained analysis on distribution shift. *arXiv preprint arXiv:2110.11328*, 2021.
- [309] Tousif Ahmed, Roberto Hoyle, Kay Connelly, David Crandall, and Apu Kapadia. Privacy concerns and behaviors of people with visual impairments. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, pages 3523–3532, 2015.
- [310] Tousif Ahmed, Roberto Hoyle, Patrick Shaffer, Kay Connelly, David Crandall, and Apu Kapadia. Understanding the physical safety, security, and privacy concerns of people with visual impairments. *IEEE Internet Computing*, 21(3):56–63, 2017.
- [311] Taslima Akter, Bryan Dosono, Tousif Ahmed, Apu Kapadia, and Bryan Se-maan. ” i am uncomfortable sharing what i can’t see”: Privacy concerns of the visually impaired with camera based assistive applications. In *29th USENIX Security Symposium (USENIX Security 20)*, pages 1929–1948, 2020.
- [312] Tousif Ahmed, Apu Kapadia, Venkatesh Potluri, and Manohar Swaminathan. Up to a limit? privacy concerns of bystanders and their willingness to share additional information with visually impaired users of assistive technologies. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 2(3):1–27, 2018.
- [313] Taslima Akter, Tousif Ahmed, Apu Kapadia, and Manohar Swaminathan. Shared privacy concerns of the visually impaired and sighted bystanders with camera based assistive technologies. *ACM Trans. Access. Comput.*, 2021.

Part II

Included Papers

Paper A

A Hierarchical Grocery Store Image Dataset with Visual and Semantic Labels

Marcus Klasson*, Cheng Zhang†, Hedvig Kjellström*

*KTH Royal Institute of Technology, Stockholm, Sweden

†Microsoft Research, Cambridge, United Kingdom

Abstract

Image classification models built into visual support systems and other assistive devices need to provide accurate predictions about their environment. We focus on an application of assistive technology for people with visual impairments, for daily activities such as shopping or cooking. In this paper, we provide a new benchmark dataset for a challenging task in this application – classification of fruits, vegetables, and refrigerated products, e.g. milk packages and juice cartons, in grocery stores. To enable the learning process to utilize multiple sources of structured information, this dataset not only contains a large volume of natural images but also includes the corresponding information of the product from an online shopping website. Such information encompasses the hierarchical structure of the object classes, as well as an iconic image of each type of object. This dataset can be used to train and evaluate image classification models for helping visually impaired people in natural environments. Additionally, we provide benchmark results evaluated on pretrained convolutional neural networks often used for image understanding purposes, and also a multi-view variational autoencoder, which is capable of utilizing the rich product information in the dataset.

1 Introduction

In this paper, we focus on the application of image recognition models implemented into assistive technologies for people with visual impairments. Such technologies already exist in the form of mobile applications, e.g. Microsoft’s Seeing AI [1] and Aipoly Vision [2], and as wearable artificial vision devices, e.g. Orcam MyEye [3] and the Sound of Vision system introduced in [4]. These products have the ability to

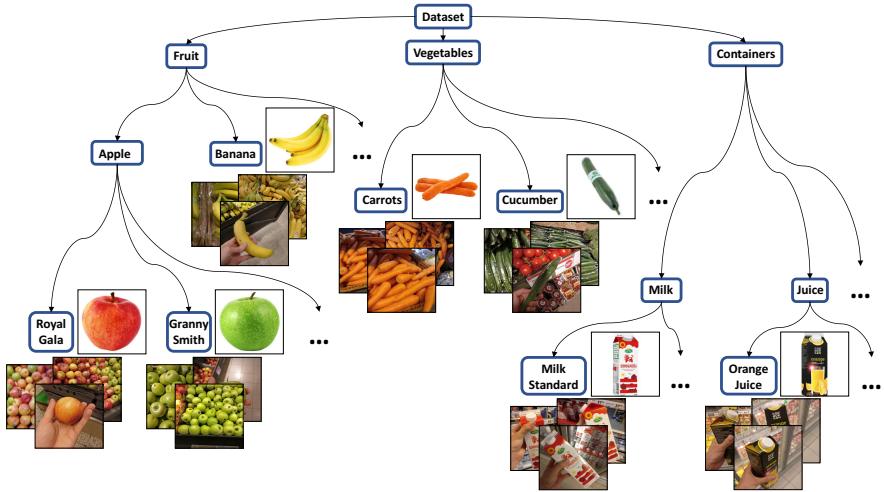


Figure 1: The primary contribution of this paper is a dataset of grocery items, for the purpose of training a visual recognition system to aid visually impaired people. The dataset is organized according to a hierarchical class structure, as illustrated above. A novel aspect of the dataset is that each class, apart from the semantic label, also has a visual label in the form of an iconic image.

support people with visual impairments in many different situations, such as reading text documents, describing the user’s environment and recognizing people the user may know.

We here address a complementary scenario not handled by current systems on the market: visual support when shopping for grocery items considering a large range of eatable objects, including fruits, vegetables, milk, and juices. In the case of fruits and vegetables, these are usually stacked in large bins in grocery stores as shown in Figure 2(a-f). A common problem in grocery stores is that similar items are often stacked next to each other; therefore, items are often misplaced into neighboring bins. Figure 2a shows a mix of red and green apples, where it might be difficult for the system to determine which kind of apple is the actual target. Humans can distinguish between groceries without vision to some degree, e.g. by touching and smelling them, but it requires prior knowledge about texture and fragrance of food items.

Moreover, in addition to raw grocery items, there are also items that can only be differentiated with the help of visual information, e.g. milk, juice, and yogurt cartons, see Figure 2(g-i). Such items usually have barcodes, that are readable using the existing assistive devices described above. However, the barcodes are not easily located by visually impaired persons. Thus, an assistive vision device that fully relies on natural image understanding would be of significant added value for a visually impaired person shopping in a grocery store.

Image recognition models used for this task typically require training images collected in similar environments. However, current benchmark datasets, such as ImageNet [5] and CIFAR-100 [6], do contain images of fruits and vegetables, but are not suitable for this type of assistive application, since the target objects are commonly not presented in this type of natural environments, with occlusion and cluttered backgrounds. To address this issue, we present a novel dataset containing natural images of various raw grocery items and refrigerated products, e.g. milk,

juice, and yogurt, taken in grocery stores. As part of our dataset, we collect images taken with single and multiple target objects, from various perspectives, and with noisy backgrounds.

In computer vision, previous studies have shown that model performance can be improved by extending the model to utilize other data sources, e.g. text, audio, in various machine learning tasks [7–10]. Descriptions of images are rather common to computer vision datasets, e.g. Flickr30k [11], whereas the datasets in [8, 12] includes both descriptions and a reference image with clean background to some objects. Therefore, in addition to the natural images, we have collected iconic images with a single object centered in the image (see Figure 3) and a corresponding product description to each grocery item. In this work, we also demonstrate how we can benefit from using additional information about the natural images by applying the multi-view generative model.

To summarize, the contribution of this paper is a dataset of natural images of raw and refrigerated grocery items, which could be used for evaluating and training image recognition systems to assist visually impaired people in a grocery store. The dataset labels have a hierarchical structure with both coarse- and fine-grained classes (see Figure 3). Moreover, each class also has an iconic image and a product description, which makes the dataset applicable to multimodal learning models. The dataset is described in Section 3.

We provide multiple benchmark results using various deep neural networks, such as Alexnet [13], VGG [14], DenseNet [15], as well as deep generative models, such as VAE [16]. Furthermore, we adapt a multi-view VAE model to make use of the iconic images for each class (Section 4), and show that it improves the classification accuracy given the same model setting (Section 5). Last, we discuss possible future directions for fully using the additional information provided with the dataset and adopt more advanced machine learning methods, such as visual-semantic embeddings, to learn efficient representations of the images.

2 Related Work

Many popular image datasets have been collected by downloading images from the web [5, 6, 8, 12, 17–21]. If the dataset contains a large amount of images, it is convenient to make use of crowdsourcing to get annotations for recognition tasks [5, 6, 22]. For some datasets, the crowdsourcers are also asked to put bounding boxes around the object to be labeled for object detection tasks [8, 17, 20]. In [18] and [6], the target objects are usually centered and takes up most content of the image itself. Another significant characteristic is that web images usually are biased in the sense that they have been taken with the object focus in mind; they have good lighting settings and are typically clean from occlusions, since the collectors have used general search words for the object classes, e.g. *car*, *horse*, or *apple*.

Some datasets include additional information about the images beyond the single class label, e.g. text descriptions of what is present in the image and bounding boxes around objects. These datasets can be used in several different computer vision tasks, such as image classification, object detection, and image segmentation. Structured labeling is another important property of a dataset, which provides flexibility when classifying images. In [8, 12], all of these features exist and moreover they



(a) Royal Gala (b) Golden Delicious (c) Orange



(d) Aubergine (e) Onion (f) Zucchini



(g) Apple Juice (h) Milk Medium Fat (i) Yoghurt Natural



(a) Royal Gala (b) Golden Delicious (c) Orange



(d) Aubergine (e) Onion (f) Zucchini



(g) Apple Juice (h) Milk Medium Fat (i) Yoghurt Natural

Figure 2: Examples of natural images in our dataset, where each image have been taken inside a grocery store. Image examples of fruits, vegetables and refrigerated products are presented in each row respectively.

Figure 3: Examples of iconic images downloaded from a grocery shopping website, which corresponds to the target items in the images in Figure 2.

include reference images to each object class, which in [12] is used for labeling multiple categories present in images, while in [8] these images are used for fine-grained recognition. Our dataset includes a reference image, i.e. the iconic image, and a product description for every class, and we have also labeled the grocery items in a structured manner.

Other image datasets of fruits and vegetables for classification purposes are the FIDS30 database [23] and the dataset in [24]. The images in FIDS30 were downloaded from the web and contain background noise as well as single or multiple instances of the object. In [24], all pixels belonging to the object are extracted from the original image, such that all images have white backgrounds with the same brightness condition. There also exist datasets for detecting fruits in orchards for robotic harvesting purposes, which are very challenging since the images contain plenty of background and various lighting conditions, and the targeted fruits are often occluded or of the same color as the background [25, 26].

Another dataset that is highly relevant to our application need is presented in [27]. They collected a dataset for training and evaluating the image classifier by extracting images from video recordings of 23 main classes, which are subdivided into 98 classes, of raw grocery items (fruits and vegetables) in different grocery stores. Using this dataset, a mobile application was developed to recognize food products in grocery store environments, which provides the user with details and health recommendations about the item along with other proposals of similar food items. For each class, there exists a product description with nutrition values to assist the user in shopping

scenarios. The main difference between this work and our dataset is firstly the clean iconic images (visual labels) for each class in our dataset, and secondly that we have also collected images of refrigerated items, such as dairy and juice containers, where visual information is required to distinguish between the products.

3 Our Dataset

We have collected images from fruit and vegetable sections and refrigerated sections with dairy and juice products in 18 different grocery stores. The dataset consists of 5125 images from 81 fine-grained classes, where the number of images in each class range from 30 to 138. Figure 4 displays a histogram over the number of images per class. As illustrated in Figure 3, the class structure is hierarchical, and there are 46 coarse-grained classes. Figure 2 shows examples of the collected natural images. For each fine-grained class, we have downloaded an iconic image of the item and also a product description including origin country, an appreciated weight and nutrient values of the item from a grocery store website. Some examples of downloaded iconic images can be seen in Figure 3.

Our aim has been to collect the natural images under the same condition as they would be as part of an assistive application on a mobile phone. All images have been taken with a 16-megapixel Android smartphone camera from different distances and angles. Occasionally, the images include other items in the background or even items that have been misplaced in the wrong shelf along with the targeted item. It is important that image classifiers that are used for assisting devices are capable of performing well with such noise since these are typical settings in a grocery store environments. The lighting conditions in the images can also vary depending on where the items are located in the store. Sometimes the images are taken while the photographer is holding the item in the hand. This is often the case for refrigerated products since these containers are usually stacked compactly in the refrigerators. For these images, we have consciously varied the position of the object, such that the item is not always centered in the image or present in its entirety.

We also split the data into a training set and test set based on the application need. Since the images have been taken in several different stores at specific days and time stamps, parts of the data will have similar lighting conditions and backgrounds for each photo occasion. To remove any such biasing correlations, all images of a certain class taken at a certain store are assigned to either the test set or training set. Moreover, we balance the class sizes to as large extent as possible in both the training and test set. After the partitioning, the training and test set contains 2640 and 2485 images respectively. Predefining a training and test set also makes it easier for other users to compare their results to the evaluations in this paper.

The task is to classify natural images using mobile devices to aid visually impaired people. The additional information such as the hierarchical structure of the class labels, iconic images, and product descriptions can be used to improve the

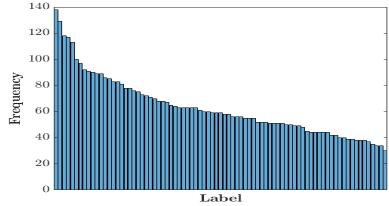


Figure 4: Histogram over the number of images in each class in the dataset.

performance of the computer vision system. Every class label is associated with a product description. Thus, the product description itself can be part of the output for visually impaired persons as they may not be able to read what is printed on a carton box or a label tag on a fruit bin in the store.

The dataset is intended for research purposes and we are open to contributions with more images and new suitable classes. Our dataset is available at <https://github.com/marcusklasson/GroceryStoreDataset>. Detailed instructions on how to contribute to the dataset can be found on our dataset webpage.

4 Classification Methods

We here describe the classification methods and approaches that we have used to provide benchmark results to the dataset. We apply both deterministic deep neural networks as well as a deep generative model used for representation learning to the natural images that we have collected. Furthermore, we utilize the additional information – iconic images – from our dataset with a multi-view deep generative model. This model can utilize different data sources and obtain superior representation quality as well as high interpretability. For a fair evaluation, we use a linear classifier with the learned representation from the different methods.

Deep Neural Networks. CNNs have been the state-of-the-art models in image classification ever since AlexNet [13] achieved the best classification accuracy in ILSVRC in 2012. However, in general, computer vision models require lots of labeled data to achieve satisfactory performance, which has resulted in interest for adapting CNNs that have already been trained on a large amount of training data to other image datasets. When adapting pretrained CNNs to new datasets, we can either use it directly as a feature extractor, a.k.a use the off-the-shelf features, [28, 29], or fine-tune it [30–34]. Using off-the-shelf features, we need to specify which feature representation we should extract from the network and use these for training a new classifier. Fine-tuning a CNN involves adjusting the pretrained model parameters, such that the network can e.g. classify images from a dataset different from what the CNN was trained on before. We can either choose to fine-tune the whole network or select some layer parameters to adjust while keeping the others fixed. One important factor on deciding which approach to choose is the size of the new dataset and how similar the new dataset is to the dataset which the CNN was previously trained on. A rule of thumb here is that the closer the features are to the classification layer, the features become more specific to the training data and task [33].

Using off-the-shelf CNN features and fine-tuned CNNs have been successfully applied in [28, 29] and [30, 31, 34] respectively. In [28, 29], it is shown that the pretrained features have sufficient representational power to generalize well to other visual recognition tasks with simple linear classifiers, such as Support Vector Machines (SVMs), without fine-tuning the parameters of the CNN to the new task. In [30, 34], all CNN parameters are fine-tuned, whereas in [31] the pretrained CNN layer parameters are kept fixed and only an adaptation layer of two fully connected layers are trained on the new task. The results from these works motivate why we should evaluate our dataset on fine-tuned CNNs or linear classifiers trained on off-the-shelf feature representations instead of training an image recognition model from scratch.

Variational Autoencoders with only natural images. Deep generative models, e.g. the variational autoencoder (VAE) [16, 35, 36], have become widely used in the machine learning community thanks to their generative nature. We thus use VAEs for representation learning as the second benchmarking method. For efficiency, we use low-level pretrained features from a CNN as inputs to the VAE.

The latent representations from VAEs are encodings of the underlying factors for how the data are generated. VAEs belongs to the family of latent variable models, which commonly has the form $p_{\theta}(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z})$, where $p(\mathbf{z})$ is a prior distribution over the latent variables \mathbf{z} and $p_{\theta}(\mathbf{x}|\mathbf{z})$ is the likelihood over the data \mathbf{x} given \mathbf{z} . The prior distribution is often assumed to be Gaussian, $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$, whereas the likelihood distribution depends on the values of \mathbf{x} . The likelihood $p_{\theta}(\mathbf{x}|\mathbf{z})$ is referred to as a decoder represented as a neural network parameterized by θ . An encoder network $q_{\phi}(\mathbf{z}|\mathbf{x})$ parameterized by ϕ is introduced as an approximation of the true posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$, which is intractable since it requires computing the integral $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$. When the prior distribution is a Gaussian, the approximate posterior is also modeled as a Gaussian, $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\sigma}^2(\mathbf{x}) \odot \mathbf{I})$, with some mean $\boldsymbol{\mu}(\mathbf{x})$ and variance $\boldsymbol{\sigma}^2(\mathbf{x})$ computed by the encoder network. The goal is to maximize the marginal log-likelihood by defining a lower bound using $q_{\phi}(\mathbf{z}|\mathbf{x})$:

$$\begin{aligned} \log p_{\theta}(\mathbf{x}) &\geq \mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] \\ &\quad - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})). \end{aligned} \quad (1)$$

The last term is the Kullback-Leibler (KL) divergence of the approximate posterior from the true posterior. The lower bound \mathcal{L} is called the evidence lower bound (ELBO) and can be optimized with stochastic gradient descent via backpropagation [16, 37]. VAE is a probabilistic framework. Many extensions such as utilizing structured priors [38] or using continual learning [39] have been explored. In the following method, we describe how to make use of the iconic images while retaining the unsupervised learning setting in VAEs.

Utilizing iconic images with multi-view VAEs. Utilizing extra information has shown to be useful in many applications with various model designs [38, 40–43]. For computer vision tasks, natural language is the most commonly used modality to aid the visual representation learning. However, the consistency of the language and visual embeddings has no guarantee. As an example with our dataset, the product description of a Royal Gala apple explains the appearance of a red apple. But if the description is represented with word embeddings, e.g. word2vec [44], the word ‘royal’ will probably be more similar to the words ‘king’ and ‘queen’ than ‘apple’. Therefore, if available, additional visual information about objects might be more beneficial for learning meaningful representations instead of text. In this work, with our collected dataset, we propose to utilize the iconic images for the representation learning of natural images using a multi-view VAE. Since the natural images can include background noise and grocery items different from the targeted one, the role of the iconic image will be to guide the model to which features that are of interest in the natural image.

The VAE can be extended to modeling multiple views of data, where a latent variable \mathbf{z} is assumed to have generated the views [40, 42]. Considering two views \mathbf{x} and \mathbf{y} , the joint distribution over the paired random variables (\mathbf{x}, \mathbf{y}) and latent variable \mathbf{z} can be written as $p_{\theta}(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(\mathbf{z})p_{\theta(1)}(\mathbf{x}|\mathbf{z})p_{\theta(2)}(\mathbf{y}|\mathbf{z})$, where both $p_{\theta(1)}(\mathbf{x}|\mathbf{z})$

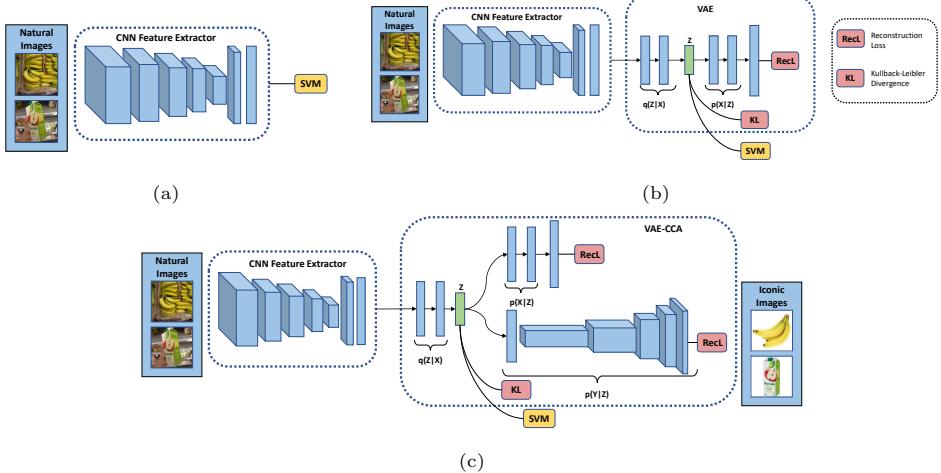


Figure 5: The architectures for the classification methods described in Section 4. In this paper, we use either a pretrained AlexNet, VGG16 or DenseNet-169 as the CNN feature extractor, but it may be replaced with any CNN architecture. Note that the pretrained CNN can be fine-tuned. The encoder and decoder of the VAE in 5b consist of two fully-connected layers. VAE-CCA in 5c uses the DCGAN architecture as an iconic image decoder and the same encoder and feature vector decoder as the VAE.

and $p_{\theta(2)}(\mathbf{y}|\mathbf{z})$ are represented as neural networks with parameters $\theta^{(1)}$ and $\theta^{(2)}$. Assuming that the latent variable \mathbf{z} can reconstruct both \mathbf{x} and \mathbf{y} when only \mathbf{x} is encoded into \mathbf{z} by the encoder $q_\phi(\mathbf{z}|\mathbf{x})$, then the ELBO is written as

$$\begin{aligned} \log p_{\theta}(\mathbf{x}, \mathbf{y}) &\geq \mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{y}) \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_{\theta(1)}(\mathbf{x}|\mathbf{z}) + \log p_{\theta(2)}(\mathbf{y}|\mathbf{z})] \\ &\quad - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})). \end{aligned} \quad (2)$$

This model is referred to as variational autoencoder canonical correlation analysis (VAE-CCA) and was introduced in [42]. The main motivation for using VAE-CCA is that the latent representations need to contain information about reconstructing both natural and iconic images. The main motivation for using VAE-CCA is that the latent representation needs to preserve information about how both the natural and iconic images are reconstructed. This also allows us to produce iconic images from new natural images to enhance the interpretability of the latent representation of VAE-CCA (see Section 5) [40].

5 Experimental Results

We apply the three different types of models described in Section 4 to our dataset and evaluate their performance. The natural images are propagated through a CNN pretrained on ImageNet to extract feature vectors. We experiment with both the off-the-shelf features as well as fine-tuning the CNN. When using off-the-shelf features, we simply extract feature vectors and train an SVM on those. For the fine-tuned CNN, we report both results from the softmax classifier used in the actual fine-tuning

Table 1: Fine-grained classification (81 classes) accuracies with the methods described in Section 5.1. Each row displays from which network architecture and layer that we extracted the feature vectors of the natural images. The columns show the result from the classifiers that we used (see Section 5.1).

	SVM	SVM-ft	VAE+SVM	VAE+SVM-ft	VAE-CCA+SVM	VAE-CCA+SVM-ft
AlexNet ₆	69.2	72.6	65.6	70.7	67.8	71.5
AlexNet ₇	65.0	70.7	63.0	68.7	65.0	70.9
VGG16 ₆	62.1	73.3	57.5	71.9	60.7	73.0
VGG16 ₇	57.3	71.7	56.8	67.8	56.8	71.3
DenseNet-169	72.5	85.0	65.4	79.1	72.6	80.4

procedure and training an SVM with extracted fine-tuned feature vectors.

These extracted feature vectors are also used for VAE and VAE-CCA which makes further compression. We perform classification for those VAE based models by training a classifier, e.g. an SVM, on the data encoded into the latent representation. We use this classification approach for both VAE and VAE-CCA. In all classification experiments, except when we fine-tune the CNN, we use a linear SVM trained with the one-vs-one approach as in [29].

We experiment with three different pretrained CNN architectures, namely AlexNet [13], VGG16 [14] and DenseNet-169 [15]. For AlexNet and VGG16, we extract feature vectors of size 4096 from the two last fully connected (FC) layers before the classification layer. The features from the n^{th} hidden layer are denoted as AlexNet _{n} and VGG16 _{n} . As an example, the last hidden FC layer in AlexNet is denoted as AlexNet₇, the input of which is output from AlexNet₆. For DenseNet-169, we extract the features of size 1664 from the average pooling layer before its classification layer.

5.1 Experimental Setups

The following setups were used in the experiments:

Setup 1. Train an SVM on extracted off-the-shelf features from a pretrained CNN, which is denoted as SVM in the results. We also fine-tune the CNN by replacing the final layer with a new softmax layer and denote these results as Fine-tune. We denote training an SVM on extracted finetuned feature vectors as SVM-ft.

Setup 2. Extract feature vectors with a pretrained CNN of the natural images and learn a latent representation \mathbf{z} with a VAE. Then the data is encoded into the latent space and we train an SVM with these latent representations, which used for classification. We denote the results as VAE+SVM when using off-the-shelf feature vectors, whereas using the fine-tuned feature vectors are denoted as VAE+SVM-ft. In all experiments with the VAE, we used the architecture from [45], i.e. the latent layer having 200 hidden units and both encoder and decoder consisting of two FC layers with 1,000 hidden units each.

Setup 3. Each natural image is paired with its corresponding iconic image. We train VAE-CCA similarly as the VAE, but instead, we learn a joint latent representation that is used to reconstruct the extracted feature vectors \mathbf{x} and the iconic images \mathbf{y} . The classification is performed with the same steps as in Setup 2 and

Table 2: Coarse-grained classification (46 classes) accuracies with an SVM for the methods described in Section 5.1 that uses off-the-shelf feature representations. Each row displays from network architecture and layer that we extracted the feature vectors of the natural images and the columns show the result for the classification methods.

	SVM	VAE+SVM	VAE-CCA+SVM
AlexNet ₆	78.0	74.2	76.4
AlexNet ₇	75.4	73.2	74.4
VGG16 ₆	76.6	74.2	74.9
VGG16 ₇	72.8	71.7	72.3
DenseNet-169	85.2	79.5	82.0

Table 3: Fine-grained classification accuracies from fine-tuned CNNs pretrained on ImageNet, where the column shows which architecture that has been fine-tuned. A standard softmax layer is used as the last layer.

	Fine-tune
AlexNet	69.3
VGG16	73.8
DenseNet-169	84.0

denotes the results similarly with VAE-CCA+SVM and VAE-CCA+SVM-ft. Our VAE-CCA model takes the feature vectors \mathbf{x} as input and encodes them into a latent layer with 200 hidden units. The encoder and the feature vector decoder uses the same architecture, i.e. two FC layers with 512 hidden units, whereas the iconic image decoder uses the DCGAN [46] architecture.

Figure 5 displays the three experimental setups described above. We report both fine-grained and coarse-grained classification results with an SVM in Table 1 and 2 respectively. In Table 3, we report the fine-grained classification results from fine-tuned CNNs.

When fine-tuning the CNNs, we replace the final layer with a softmax layer applicable to our dataset with randomly initialized weights drawn from a Gaussian with zero mean and standard deviation 0.01 [34]. For AlexNet and VGG16, we fine-tune the networks for 30 epochs with two different learning rates, 0.01 for the new classification layer and 0.001 for the pretrained layers. Both learning rates are reduced by half after every fifth epoch. The DenseNet-169 is fine-tuned for 30 epochs with momentum of 0.9 and an initial learning rate of 0.001, which decays with 10^{-6} after each epoch. We report the classification results from the softmax activation after the fine-tuned classification layer. We also report classification results from an SVM trained with feature representations from a fine-tuned CNN, which are extracted from FC6 and FC7 of the AlexNet and VGG16 and from the last average pooling layer in DenseNet-169.

The VAE and VAE-CCA models are trained for 50 epochs with Adam [47] for optimizing the ELBOs in Equation 3 and 2 respectively. We use a constant learning rate of 0.0001 and set the minibatch size to 64. The extracted feature vectors are rescaled with standardization before training the VAE and VAE-CCA models to stabilize the learning.

5.2 Results

The fine-grained classification results for all methods using an SVM as classifier are shown in Table 1. We also provide coarse-grained classification results for some of the methods in Table 2 to demonstrate the possibility of hierarchical evaluation that our labeling of the data provides (see Figure 3). The accuracies in the coarse-grained classification are naturally higher than the accuracies in the corresponding columns in Table 1. Table 3 shows fine-grained classification accuracies from a softmax classifier

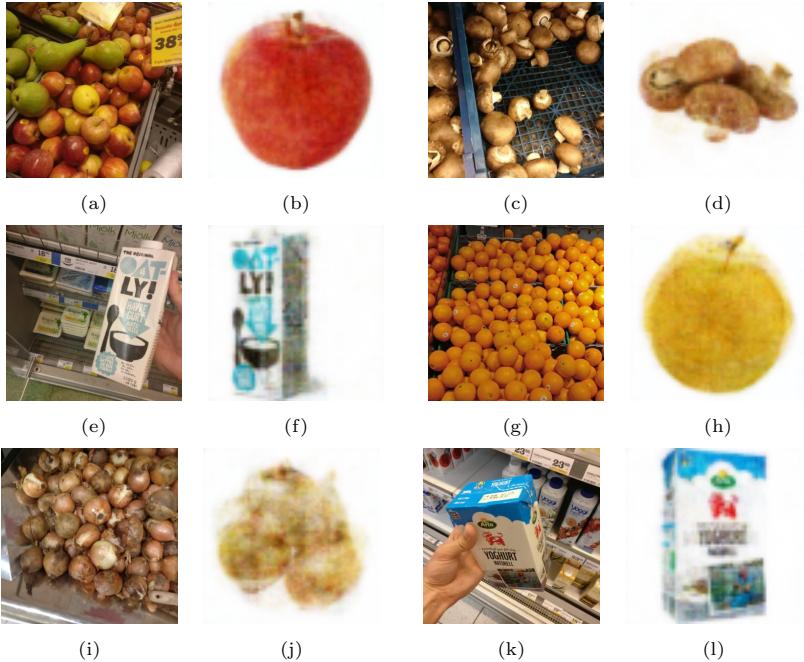


Figure 6: Examples of six natural images in the test set that have been decoded into product iconic images by the iconic image decoder $p_{\theta(2)}(\mathbf{y} | \mathbf{z})$ using VAE-CCA model as in Figure 5c. This result is obtained with the fine-tuned DenseNet-169 features, which corresponds to VAE-CCA+SVM-ft in Table 1. We show the natural image and corresponding decoded iconic image next to each other. The classes for all images are (a, b) Royal Gala Apple, (c, d) Brown Cap Mushroom, (e, f) Oatly Oatgurt, (g, h) Orange, (i, j) Onion, and (k, l) Arla Natural Yogurt.

in the fine-tuned CNNs. We note that fine-tuning the networks gives consistently better results than training an SVM on off-the-shelf features (see Table 1).

Fine-tuning the entire network results improves the classification performance consistently for each method in Table 1. The performance is clearly enhanced for features extracted from fine-tuned VGG16 and DenseNet-169, which improves the classification accuracy by 10% in most cases for SVM-ft, VAE+SVM-ft, and VAE-CCA+SVM-ft. For AlexNet and VGG16, we see that the performance drops when extracting the features from layer FC7 instead of FC6. The reason might be that the off-the-shelf features in FC7 are more difficult to transfer to other datasets since the weights are biased towards classifying objects in the ImageNet database. The performance drops also when we use fine-tuned features, which could be due to the small learning rate we use for the pretrained layers, such that the later layers are still ImageNet-specific. We might circumvent this drop by increasing the learning rate for the later pretrained layers and keeping the learning rate for earlier layers small.

The VAE-CCA model achieves mostly higher classification accuracies than the VAE model in both Table 1 and 2. This indicates that the latent representation separates the classes more distinctly than the VAE by jointly learning to reconstruct the extracted feature vectors and iconic images. However, further compressing the feature vectors with VAE and VAE-CCA will lower the classification accuracy com-

pared to applying the feature vectors to a classifier directly. Since both VAE and VAE-CCA compresses the feature vectors into the latent representation, there is a risk of losing information about the natural images. We might receive better performance by increasing the dimension of the latent representation at the expense of speed in both training and classification.

In Figure 6, we show results from the iconic image decoder $p_{\theta(2)}(\mathbf{y} | \mathbf{z})$ when translating natural images from the test set into iconic images with VAE-CCA and a fine-tuned DenseNet-169 as feature extractor. Such visualization can demonstrate the quality of the representation using the model, as well as enhancing the interpretability of the method. Using VAE-CCA in the proposed manner, we see that with challenging natural images, the model is still able to learn an effective representation which can be decoded to the correct iconic image. For example, some pears have been misplaced in the bin for Royal Gala apples in Figure 6a, but still the image decoder manages to decode a blurry red apple seen in Figure 6b. In Figure 6h, a mix of an orange and an apple are decoded from a bin of oranges in Figure 6g, which indicates these fruits are encoded close to each other in the learned latent space. Even if Figure 6e includes much of the background, the iconic image decoder is still able to reconstruct the iconic images accurately in Figure 6f, which illustrates that the latent representation is able to explain away irrelevant information in the natural image and preserved the features of the oatgurt package. Thus, using VAE-CCA with iconic images as the second view not only advances the classification accuracy but also provides us with the means to understand the model.

6 Conclusions

This paper presents a dataset of images of various raw and packaged grocery items, such as fruits, vegetables, and dairy and juice products. We have used a structured labeling of the items, such that grocery items can be grouped into more general (coarse-grained) classes and also divided into fine-grained classes. For each class, we have a clean iconic image and a text description of the item, which can be used for adding visual and semantic information about the items in the modeling. The intended use of this dataset is to train and benchmark assistive systems for visually impaired people when they shop in a grocery store. Such a system would complement existing visual assistive technology, which is confined to grocery items with barcodes. We also present preliminary benchmark results for the dataset on the task of image classification.

We make the dataset publicly available for research purposes at <https://github.com/marcusklasson/GroceryStoreDataset>. Additionally, we will both continue collecting natural images, as well as ask for public contributions of natural images in shopping scenarios to enlarge our dataset.

For future research, we will advance our model design to utilize the structured nature of our labels. Additionally, we will design a model that use the product description of the objects in addition to the iconic images. One immediate next step is to extend the current VAE-CCA model to three views, where the third view is the description of the product.

References

- [1] Microsoft Seeing AI app. <https://www.microsoft.com/en-us/seeing-ai/>. Accessed on 2018-02-22.
- [2] Aipoly Vision app. <https://www.aipoly.com/>. Accessed on 2018-02-28.
- [3] Orcam. <https://www.orcam.com/en/>. Accessed on 2018-02-28.
- [4] S. Caraiman, A. Morar, M. Owczarek, A. Burlacu, D. Rzeszotarski, N. Botezatu, P. Herghelegiu, F. Moldoveanu, P. Strumillo, and A. Moldoveanu. Computer vision for the visually impaired: the sound of vision system. In *IEEE International Conference on Computer Vision Workshops*, 2017.
- [5] J. Deng, W. Dong, R. Socher, L. J. Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2009.
- [6] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- [7] Andrea Frome, Greg Corrado, Jonathon Shlens, Samy Bengio, Jeffrey Dean, Marc'Aurelio Ranzato, and Tomas Mikolov. Devise: A deep visual-semantic embedding model. In *Advances in Neural Information Processing Systems*, 2013.
- [8] Timnit Gebru, Jonathan Krause, Yilun Wang, Duyun Chen, Jia Deng, and Li Fei-Fei. Fine-grained car detection for visual census estimation. In *AAAI Conference on Artificial Intelligence*, 2017.
- [9] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [10] Jiquan Ngiam, Aditya Khosla, Mingyu Kim, Juhan Nam, Honglak Lee, and Andrew Y. Ng. Multimodal deep learning. In *International Conference on Machine Learning*, 2011.
- [11] Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *IEEE International Conference on Computer Vision*, 2015.
- [12] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, 2014.
- [13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, 2012.
- [14] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [15] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [16] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.

- [17] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338, Jun 2010.
- [18] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. Technical Report 7694, California Institute of Technology, 2007.
- [19] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [20] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-UCSD Birds 200. Technical Report CNS-TR-2010-001, California Institute of Technology, 2010.
- [21] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. SUN database: Large-scale scene recognition from abbey to zoo. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2010.
- [22] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *International Conference on Computer Vision*, 2015.
- [23] Škrjanec Marko. Automatic fruit recognition using computer vision. (Mentor: Matej Kristan), Fakulteta za računalništvo in informatiko, Univerza v Ljubljani, 2013. FIDS30 dataset was accessed 2018-02-24 at <http://www.vicos.si/Downloads/FIDS30>.
- [24] Horea Muresan and Mihai Oltean. Fruit recognition from images using deep learning. Technical report, Babes-Bolyai University, 2017.
- [25] Suchet Bargoti and James Patrick Underwood. Deep fruit detection in orchards. In *IEEE International Conference on Robotics and Automation*, 2017.
- [26] Inkyu Sa, Zongyuan Ge, Feras Dayoub, Ben Upcroft, Tristan Perez, and Chris McCool. Deepfruits: A fruit detection system using deep neural networks. *Sensors*, 16(8):1222, 2016.
- [27] Georg Waltner, Michael Schwarz, Stefan Ladstätter, Anna Weber, Patrick Luley, Horst Bischof, Meinrad Lindschinger, Irene Schmid, and Lucas Paletta. Mango - mobile augmented reality with functional eating guidance and food awareness. In *International Workshop on Multimedia Assisted Dietary Management*, 2015.
- [28] 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. In *International Conference on Machine Learning*, 2014.
- [29] Ali Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. CNN features off-the-shelf: An astounding baseline for recognition. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2014.
- [30] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, June 2014.
- [31] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.

- [32] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10):1345–1359, Oct 2010.
- [33] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Advances in Neural Information Processing Systems*, 2014.
- [34] Ning Zhang, Jeff Donahue, Ross B. Girshick, and Trevor Darrell. Part-based r-cnns for fine-grained category detection. In *European Conference on Computer Vision*, 2014.
- [35] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic back-propagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, 2014.
- [36] Cheng Zhang, Judith Buntage, Hedvig Kjellstrom, and Stephan Mandt. Advances in variational inference. *arXiv preprint arXiv:1711.05597*, 2017.
- [37] Carl Doersch. Tutorial on variational autoencoders. *CoRR*, abs/1606.05908, 2016.
- [38] Judith Buntage, Jiawei He, Cheng Zhang, Leonid Sigal, and Stephan Mandt. Informed priors for deep representation learning. In *Symposium on Advances in Approximate Bayesian Inference*, 2018.
- [39] Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational continual learning. In *International Conference on Learning Representations*, 2018.
- [40] Ramakrishna Vedantam, Ian Fischer, Jonathan Huang, and Kevin Murphy. Generative models of visually grounded imagination. In *International Conference on Learning Representations*, 2018.
- [41] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3156–3164, 2015.
- [42] Weiran Wang, Honglak Lee, and Karen Livescu. Deep variational canonical correlation analysis. *CoRR*, abs/1610.03454, 2016.
- [43] Cheng Zhang, Hedvig Kjellström, and Carl Henrik Ek. Inter-battery topic representation learning. In *European Conference on Computer Vision*, pages 210–226. Springer, 2016.
- [44] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, 2013.
- [45] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. In *Advances in Neural Information Processing Systems*. 2015.
- [46] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [47] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.

Paper B

Using Variational Multi-View Learning for Classification of Grocery Items

Marcus Klasson*, Cheng Zhang†, Hedvig Kjellström*

*KTH Royal Institute of Technology, Stockholm, Sweden

†Microsoft Research, Cambridge, United Kingdom

Abstract

An essential task for computer vision-based assistive technologies is to help visually impaired people to recognize objects in constrained environments, for instance, recognizing food items in grocery stores. In this paper, we introduce a novel dataset with natural images of groceries – fruits, vegetables, and packaged products – where all images have been taken inside grocery stores to resemble a shopping scenario. Additionally, we download iconic images and text descriptions for each item that can be utilized for better representation learning of groceries. We select a multi-view generative model, which can combine the different item information into lower-dimensional representations. The experiments show that utilizing the additional information yields higher accuracies on classifying grocery items over only using the natural images. We observe that iconic images help to construct representations separated by visual differences of the items, while text descriptions enable the model to distinguish between visually similar items by different ingredients.

1 Introduction

In recent years, computer vision-based assistive technologies have been developed for supporting people with visual impairments. Such technologies exist in the form of mobile applications, e.g., Microsoft’s Seeing AI [1] and Aipoly Vision [2], and as wearable artificial vision devices, e.g., Orcam MyEye [3], Transsense [4], and the Sound of Vision system [5]. These products can support people with visual impairments in many different situations, such as reading text documents, describing the user’s environment, and recognizing people the user may know. In this paper, we focus on an application that is essential for assistive vision, namely visual support when



(a) Royal Gala

(b) Golden Delicious

(c) Orange

(d) Aubergine

(e) Onion

(f) Zucchini

(g) Apple Juice

(h) Milk Medium Fat

(i) Yoghurt Natural



(a) Royal Gala

(b) Golden Delicious

(c) Orange



(d) Aubergine



(e) Onion



(f) Zucchini



(g) Apple Juice



(h) Milk Medium Fat



(i) Yoghurt Natural

Figure 1: Examples of natural images in our dataset, where each image has been taken inside a grocery store. Image examples of fruits, vegetables, and refrigerated products are presented in each row respectively.

Figure 2: Examples of iconic images downloaded from a grocery shopping website, which corresponds to the target items in the images in Figure 1.

shopping for grocery items considering a large range of eatable objects, including fruits, vegetables, and refrigerated products, e.g., milk and juice packages.

Grocery shopping with low vision capabilities can be difficult for various reasons. For example, in grocery store sections for raw groceries, the items are often stacked in large bins as shown in Figure 1a-f. Additionally, similar items are usually stacked next to each other, and therefore, items can be misplaced into neighboring bins. Humans can distinguish between groceries without vision to some degree, e.g., by touch and smell, but it requires prior knowledge about texture and fragrance of food items. Furthermore, packaged items, e.g., milk, juice, and yoghurt cartons, can only be differentiated with the help of visual information (see Figure 1g-i). Such items usually have barcodes, that are readable using the existing assistive vision devices described above. Even if using a barcode detector is a clever solution, it can be inconvenient and exhausting always having to detect barcodes of packaged items. Therefore, an assistive vision device that relies on natural images would be of significant value for a visually impaired person in a grocery store.

For an image classification model to be capable of recognizing grocery items in the setting of grocery shopping, we need image data from grocery store environments. In our previous work [6], we addressed this by collecting a novel dataset containing natural images of various raw grocery items and refrigerated products, where all images have been taken with a smartphone camera inside grocery stores to simulate a realistic shopping scenario. In addition to the natural images, we also looked upon alternative types of grocery item data that could facilitate the classification task. Grocery store chains commonly maintain websites where they store information about

the grocery items they sell and are currently available in the store for online grocery shopping. For every item, there is usually an iconic image of the item with white background, a text description about the item, and also information about nutrition values, origin country, etc. We downloaded such information from an online grocery shopping website to complement the natural images in the dataset.

To utilize the available information from the dataset for grocery item classification, we use a multi-view generative model Variational Canonical Correlation Analysis (VCCA) [7] that learns a shared representation of the natural images and the downloaded information. A view can be defined as any signal or data measured by some appropriate sensor and combining information from multiple views has previously been shown to be helpful for various image classification tasks [8–15]. However, naively adding more additional information may not lead to improved results and even harm the performance of the model [16, 17]. We, therefore, perform an ablation study over the available views in the dataset with VCCA to gain insights into how each view can affect the classification performance. Moreover, VCCA allows separating the latent space into shared and private components, where the private latent spaces should hold information about a single view. This might prove useful for reducing view-specific noise in the shared latent space, which can ease the learning process of the representations we want to use for classification. We investigate the effects each view has on the learned representations of VCCA by measuring classification performances as well as visualizing the latent space for the different model settings. The contributions of this paper are the following:

- We present a novel dataset with natural images of grocery items as well as iconic images and text descriptions for every item class (see subsection The Grocery Store Dataset in Results) [6]. The natural images are taken in grocery stores in different lighting conditions and backgrounds that can be challenging settings for a computer vision-based assistive device. The additional iconic images and text descriptions make the dataset suitable for multi-view learning by combining the natural images with the extra information to obtain better classification performance.
- We investigate how to use information from different views for the task of grocery item classification with the deep generative model VCCA (see subsection Methods in Experimental Procedures). This model combines information from multiple views into a low-dimensional latent representation that can be used for classification. We also select a variant of VCCA denoted VCCA-private that separates shared and private information about each view through factorization of the latent representation (see subsection Extracting Private Information of Views with VCCA-private in Experimental Procedures). Furthermore, we use a standard multi-view autoencoder model called Split Autoencoder (SplitAE) [17, 18] for benchmarking against VCCA and VCCA-private on classification.
- We conduct experiments with SplitAE, VCCA, and VCCA-private on the task of grocery item classification with our dataset (Results). We perform a thorough ablation study over all views in the dataset to demonstrate how each view contributes to enhancing the classification performance and conclude that

utilizing the web-scraped views yields better classification results than only using the natural images (see subsection Classification Results in Results). To gain further insights into the results, we visualize the learned latent representations of the VCCA models and discuss how the iconic images and textual descriptions impose different structures on the latent space that are beneficial for classification (see subsection Investigation of the Learned Representations in Results).

This work is an extended version of Klasson et al. [6], where we first presented this dataset. In this paper, we have added a validation set of natural images from two stores that were not present in the training and test set splits from [6] to avoid overfitting effects. We also demonstrate how the text descriptions can be utilized alone and along with the iconic images in a multi-view setting, while [6] only experimented with the combination of natural and iconic images to build better representations of grocery items. Finally, we decode iconic images from unseen natural images as an alternative to evaluate the usefulness of the latent representations (see subsection Decoding Iconic Images from Unseen Natural Images in Results). As we only evaluated the decoded iconic images qualitatively in Klasson et al. [6], we have extended the assessment by comparing the quality of the decoded images from different VCCA models with multiple image similarity metrics.

Next, we discuss image datasets, food datasets, and multi-view models related to our work:

Image Datasets Many image datasets used in computer vision have been collected by downloading images from the web [19–30]. Some datasets [19, 22, 24, 26, 28] use search words with the object category in isolation, which typically returns high-quality images where the searched object is large and centered. To collect images from more real-world scenarios, searching for combinations of object categories usually returns images of two searched categories but also numerous other categories [25, 30]. The simplest annotation of these images is to provide a class label for the present objects. Occasionally, the dataset can use a hierarchical labeling structure and provide a fine- and coarse-grained label to objects where it is applicable. The annotators can also be asked to increase the possible level of supervision for the objects by, for instance, providing bounding boxes, segmentation masks, keypoints, text captions that describe the scene, and reference images of the objects [21, 23, 25, 26, 28, 30]. Our dataset includes reference (iconic) images of the objects that were web-scraped from a grocery store website. We also downloaded text descriptions that describe general attributes of the grocery items, such as flavor and texture, rather than the whole visual scene. The grocery items have also been labeled hierarchically in a fine- and coarse-grained manner if there exist multiple variants of specific items. For example, fine-grained classes of apples such as *Golden Delicious* or *Royal Gala* belongs to the coarse-grained class *apple*.

Food Datasets Recognizing grocery items in their natural environments, such as grocery stores, shelves, and kitchens, have been addressed in plenty of previous works [6, 31–41]. The addressed tasks range from hierarchical classification, object detection, segmentation, and 3D model generation. Most of these works collect a

dataset that resembles shopping or cooking scenarios, where the datasets vary in the degree of labeling, different camera views, and the data domain difference between the training and test set. The training sets in GroZi-120 [36], Grocery Products [32], and CAPG-GP [31] datasets were obtained by web-scraping product images of single instances on grocery web stores, while the test sets were collected in grocery stores where there can be single and multiple instances of the same item and other different items. The RPC [39] and TGFS [41] datasets are used for object detection and classification of grocery products, where RPC is targeted for automatic checkout systems and TGFS is for the task of recognizing items purchased from self-service vending machines. The BigBIRD [37] dataset and datasets from [?, 33] contain images of grocery items from multiple camera views, segmentation masks, and depth maps for 3D reconstruction of various items. The Freiburg Groceries [34] dataset contains images taken with smartphone cameras of items inside grocery stores, while its test set consists of smartphone photos in home environments with single or multiple instances from different kinds of items. The dataset presented in [38] also contains images taken with smartphone cameras inside grocery stores to develop a mobile application for recognizing raw food items and provide details about the item, such as nutrition values and recommendations of similar items. Other works that collected datasets of raw food items, such as fruits and vegetables, focused on the standard image classification task [42, 43] and on detecting fruits in orchards for robotic harvesting [44, 45]. Our dataset – the Grocery Store dataset – shares many similarities with the mentioned works above, for instance, that all images of groceries are taken in their natural environment, the hierarchical labeling of the classes, and the iconic product images for each item in the dataset. Additionally, we have provided a text description for each item that was web-scraped along with the iconic image. As most grocery item datasets only include packaged products, we have also collected images of different fruit and vegetable classes along with packages in our dataset.

Other examples of food datasets are those with images of food dishes, where [46] provides a summary of existing benchmark food datasets. The contents of these datasets range from images of food dishes [47–50], cooking videos [51], recipes [52–54], and also restaurant-oriented information [55, 56]. One major difference between recognizing groceries and food dishes is the appearance of the object categories in the images. For instance, a fruit or vegetable is usually intact and present in the image, while ingredients that are significant for recognizing a food dish may not be visible at all depending on the recipe. However, recognizing raw food items and dishes share similarities since they can appear with many different visual features in the images compared to packaged groceries, e.g., carton boxes, cans, and bottles, where the object classes have identical shape and texture. Another key difference is the natural environments and scenes where the images of grocery items and food dishes have been taken. Food dish images usually show the food on a plate placed on a table and, occasionally, with cutlery and a glass next to the plate. Images taken in grocery stores can cover many instances of the same item stacked close to each other in shelves, baskets, and refrigerators, while there can be multiple different kinds of items next to each other in a kitchen environment. To summarize, recognizing grocery items and food dishes are both challenging tasks because examples from the same category can look very different and also appear in various realistic settings in images.

Multi-View Learning Models There exist lots of multi-view learning approaches for data fusion of multiple features [7, 9, 10, 17, 57–70]. A common approach is to obtain a shared latent space for all views with the assumption that each view has been generated from this shared space [68]. A popular example of this is approach is Canonical Correlation Analysis (CCA) [71], which aims to project two sets of variables (views) into a lower-dimensional space so that the correlation between the projections is maximized. Similar methods propose maximizing other alignment objectives for embedding the views, such as ranking losses [9, 10, 59, 67]. There exist nonlinear extensions of CCA, e.g., Kernel CCA [72] and Deep CCA [73], that optimize their nonlinear feature mappings based on the CCA objective. Deep Canonically Correlated Autoencoders (DCCAE) [18] is a Deep CCA model with an autoencoding part, which aims to maximize the canonical correlation between the extracted features as well as reconstructing the input data. Removing the CCA objective reduces DCCAE to a standard multi-view autoencoder, e.g., Bimodal Autoencoders and Split Autoencoders (SplitAEs) [17, 18], that only aim to learn a representation that best reconstructs the input data. SplitAEs aim to reconstruct two views from a representation encoded from one of the views. This approach was empirically shown to work better than Bimodal Autoencoders in [17] in situations where only a single view is present at both training and test time.

Variational CCA (VCCA) [7] can be seen as an extension of CCA to deep generative models, but can also be described as a probabilistic version of SplitAEs. VCCA uses the amortized inference procedure from Variational Autoencoders (VAEs) [74, 75] to learn the shared latent space by maximizing a lower bound on the data log-likelihood of the views. Succeeding works have proposed new learning objectives and inference methods for enabling conditional generation of views, e.g., generating an image conditioned on some text and vice versa [58, 62, 63, 65, 66]. These approaches rely on fusing the views into the shared latent space as the inference procedure, which often requires tailored training and testing paradigms when views are missing. However, adding information from multiple views may not lead to improved results and can even make the model perform worse on the targeted task [16, 17]. This is especially the case when views have noisy observations, which complicates learning a shared latent space that combines the commonalities between the views. To avoid disturbing the shared latent space with noise from single views, some works design models which extend the shared latent space with private latent spaces for each view that should contain the view-specific variations to make learning the shared variations easier [7, 61, 64, 69, 76]. VCCA can be extended to extract shared and private information between different views through factorization of the latent space into shared and private parts. In this paper, we investigate if the classification performance of grocery items in natural images can be improved by extracting the view-specific variations in the additional views (iconic images and product descriptions) from the shared latent space with this variant of VCCA called VCCA-private. We will experiment with treating each data point as pairs of natural images and either iconic images or text descriptions as well as triplets of natural images, iconic images, and text descriptions. A difference between how we apply VCCA to our dataset compared to the works above is that the iconic image and text description are the same for every natural image of a specific class.

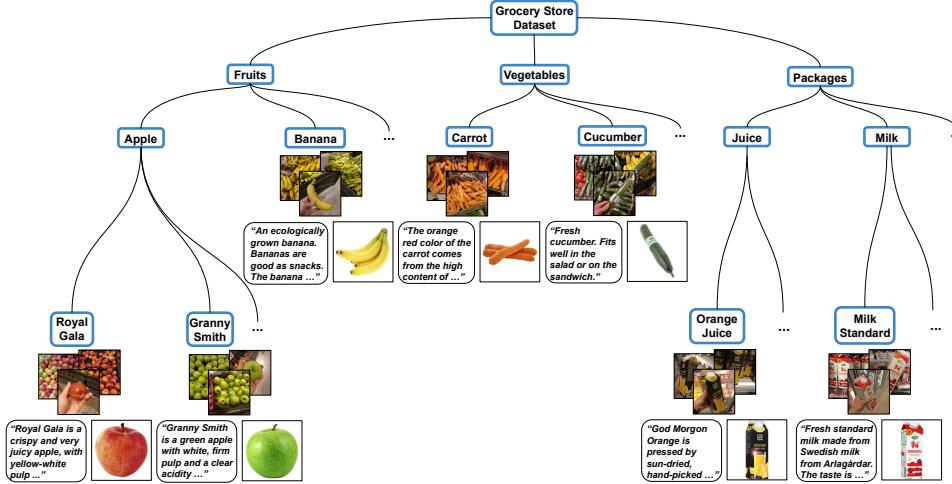


Figure 3: Illustration of the hierarchical class structure of the dataset. First, the classes are divided by their grocery item type, i.e., *Fruits*, *Vegetables*, and *Packages*, followed by a separating the items into coarse-grained classes, e.g., *Apple*, *Carrot*, and *Milk*, and then into fine-grained classes. There are 81 fine-grained classes in total and 46 coarse-grained classes. The figure also shows the iconic image and text description of the items next to the class label.

2 Results

In this section, we begin by providing the details about the collected dataset, which we have named the Grocery Store dataset. Furthermore, we illustrate the utility of the additional information in the Grocery Store dataset to classify grocery items in the experiments. We compare SplitAE, VCCA, and VCCA-private with different combinations of views against two standard image classifiers. Additionally, we experiment with a vanilla autoencoder (denoted as AE) and a VAE that post-processes the natural image features to train a linear classifier cheaper and compare the performance against the multi-view models. We measure the classification performance on the test set for every model and also compare the classification accuracies when the number of words in the text description varies for the models concerned (see subsection Classification Results in Results). To gain insights on how the additional views affect the learned representations, we visualize the latent spaces of VCCA and VCCA-private with PCA and discuss how different views changes the structure of the latent space (see subsection Investigation of the Learned Representations in Results). Finally, we show how iconic images can be used for enhancing the interpretability of the classification (see subsection Decoding Iconic Images from Unseen Natural Images in Results), which was also illustrated by Klasson et al. [6].

2.1 The Grocery Store Dataset

In Klasson et al. [6], we collected images from fruit, vegetable, and refrigerated sections with dairy and juice products in 20 different grocery stores. The dataset consists of 5421 images from 81 different classes. For each class, we have downloaded

an iconic image of the item, a text description, and information including origin country, appreciated weight and nutrient values of the item from a grocery store website. Some examples of natural images and downloaded iconic images can be seen in Figure 1 and 2 respectively. Furthermore, Table S1 displays a selection of text descriptions with their corresponding iconic image. The text description length varies between 6 and 91 words with an average of 36 words. We also structure that classes hierarchically with three levels as illustrated in Figure 3. The first level divides the items into three categories; *Fruits*, *Vegetables*, and *Packages*. The next level consists of 46 coarse-grained classes which divides the items into groups contain groups of items types, e.g., *Apple*, *Carrot*, *Milk*. The bottom level consists of the 81 fine-grained classes, where the items are completely separated. Note that a coarse-grained class without any fine-grained classes in the dataset, e.g., *Banana* and *Carrot*, is considered as a fine-grained class in the classification experiments (see subsection Classification Results in Results), where we report both fine-grained and coarse-grained classification performance of the models.

We aimed to collect the natural images under similar conditions as if they would be taken with an assistive mobile application. All images have been taken with a 16-megapixel Android smartphone camera from different distances and angles. Occasionally, the images include other items in the background or even items that have been misplaced in incorrect shelves along with the targeted item. The lighting conditions in the images can also vary depending on where the items are located in the store. Sometimes the images are taken while the photographer is holding the item in the hand. This is often the case for refrigerated products since these containers are usually stacked compactly in the refrigerators. For these images, we have consciously varied the position of the object, such that the item is not always centered in the image or present in its entirety. We split the natural images into a training set and test set based on the application need. Since the images have been taken in several different stores at specific days and time stamps, parts of the data will have similar lighting conditions and backgrounds for each photo occasion. To remove any such biasing correlations, all images of a certain class taken at a certain store are assigned to either the training set, validation set, or test set. In the first version of the dataset [6], we balanced the class sizes to a large extent as possible in both the training and test set, which resulted in a training and test set containing 2640 and 2485 images respectively. In this paper, we have extended the dataset with a validation set containing 296 images taken with the same smartphone camera as before. The validation set was collected from two different grocery stores than the ones in the first version to avoid the biasing correlations described above. Initially, we experimented with grabbing a validation set from the current training set and noticed that the trained classifiers performed exceptionally well on the validation set. However, the classifiers generalized poorly to images from the test set that were taken in other stores and therefore decided to collect the validation set in two unseen stores to avoid the biases from the training set. We show histograms representing the class distributions for the training, validation, and test splits in Figure S1.

The scenario that we wanted to depict with the dataset was using a mobile device to classify natural images visually impaired people with grocery shopping. The additional information such as the iconic images, text descriptions, and hierarchical structure of the class labels can be used to enhance the performance of the computer

vision system. Since every class label is associated with a text description, the description itself can be part of the output for visually impaired persons as they may not be able to read what is printed on a carton box or a label tag on a fruit bin in the store.

2.2 Models

In this section, we briefly describe the models that we apply to grocery classification. The notation of the views that are available in the Grocery Store dataset are the following:

- x : Natural image encoded into image feature space with an off-the-shelf convolutional neural network.
- i : Iconic image of the object class in the natural image.
- w : Text description of the object class in the natural image.
- y : Class label of the natural image.

We mainly focus on analyzing VCCA [7] for utilizing different combinations of the views and investigate how each view contributes to the classification performance of grocery items. Our primary baseline model is the SplitAE which extracts shared representations by reconstructing all views, while VCCA aims to maximize a lower bound on the data log-likelihood for all views. We also study a variant of VCCA called VCCA-private [7], which is used for extracting private information about each view in addition to shared information across all views by factorizing the latent space.

We compare the multi-view models against single-view methods only using the natural images for classification. As our first single-view baseline, we customize the output layer of DenseNet169 [77] to our the Grocery Store dataset and train it from scratch to classify the natural images, which we refer to as DenseNet-scratch in the experiments. The second baseline called Softmax is a Softmax classifier trained on the off-the-shelf features from DenseNet169 pre-trained on the ImageNet dataset [19], where we extract 1664-dimensional from the average pooling layer before the classification layer in the architecture. We also experiment with AEs and VAEs for post-processing the off-the-shelf features and use a linear classifier to evaluate the models on classification. See subsection Methods in Experimental Procedures for a thorough description of the single- and multi-view autoencoders used in this paper. We name the single- and multi-view autoencoders using subscripts to denote the views utilized for learning the shared latent representations. For example, VCCA _{xi} utilizes natural image features x and iconic images i , while VCCA _{$xiwy$} uses natural image features x and iconic images i , text descriptions w , and class labels y .

2.3 Classification Results

We evaluated the classification accuracy on the test set for each model. We also calculate the accuracy for the coarse-grained classes with the following method: Let the input $x^{(i)}$ have a fine-grained class $y_{fg}^{(i)}$ and a coarse-grained class $y_{cg}^{(i)}$. Each fine-grained class can be mapped to its corresponding coarse-grained class using

$Pa(y_{fg}^{(i)}) = y_{cg}^{(i)}$, where $Pa(\cdot)$ stands for "parent". Then we compute the coarse-grained accuracy using

$$\text{coarse accuracy} = \frac{1}{N} \sum_{i=1}^N \left[Pa(\hat{y}_{fg}^{(i)}) = y_{cg}^{(i)} \right] \quad (1)$$

$$\hat{y}_{fg}^{(i)} = \arg \max_y p(y|x^{(i)})$$

where $\left[Pa(\hat{y}_{fg}^{(i)}) = y_{cg}^{(i)} \right] = 1$ when the condition is true and $\hat{y}_{fg}^{(i)}$ is the predicted fine-grained class from the selected classifier. The classification results for all models are shown in Table 1. We group the results in the table according to the utilized views and classifier. We see that Softmax trained on off-the-shelf features outperforms DenseNet-scratch by 4%. This result is common when applying deep learning to small image datasets, where pre-trained deep networks are transferred to a new dataset usually performs well compared to training neural networks on the dataset from scratch. Therefore, we present results using the off-the-shelf features for all other models.

The SplitAE and VCCA models surpass the Softmax baseline in classification performance when incorporating either the iconic image or text description view. We believe that the models using the iconic images achieve better classification accuracy over models using the text description because the iconic images contain visual features, e.g., color and shape of items, that are more useful for the image classification task. The text descriptions include more often information about the flavor, origin, and cooking details rather than describing visual features of the item, which can be less informative when classifying items from images. In most cases, the corresponding SplitAE and VCCA models perform on par for classification performance. However, VCCA-private with iconic images results in a significant drop in accuracy compared to its counterpart. We observed that the private latent variable simply models noise since there is only a single iconic image (and text description) for each class. We provide a further explanation of this phenomenon in subsection Investigation of the Learned Representations in Results.

We observe that VCCA models compete with their corresponding SplitAE models in the classification task. The main difference between these models is the Kullback-Leibler (KL) divergence [78] term in the ELBO that encourages a smooth latent space for VCCA (see Experimental Procedures). In contrast, SplitAE learns a latent space that best reconstructs the input data, which can result in parts of the space that does not represent the observed data. We showcase these differences by plotting the latent representations of SplitAE_{xiwy} and VCCA_{xiwy} using PCA in Figure 4. In Figure 4a and 4b, we have plotted the corresponding iconic image for the latent representations. We observe that VCCA_{xiwy} tries to establish a smooth latent space by pushing visually similar items closer to each other, but at the same time prevent spreading out the representations too far from the origin. Figure 4c and 4d shows the positions of the bell peppers items in the latent spaces, where the color of the point corresponds to the specific bell pepper class. In Figure 4c, we observe that SplitAE_{xiwy} has spread out the bell peppers across the space, while VCCA_{xiwy} establishes shorter distances between them in Figure 4d due to the regularization.

We evaluated the classification performance achieved by each SplitAE, VCCA,

Table 1: Classification accuracies on the test set for all models in percentage (%) along for each model. The subscript letters in the model names indicate the data views used in the model. The column Accuracy corresponds to the fine-grained classification accuracy. The column Coarse Accuracy corresponds to the classifying a class within the correct parent class. Results are averaged using 10 different random seeds and we report both means and standard deviations. Abbreviations: AE, Autoencoder; VAE, Variational Autoencoder; SplitAE, Split Autoencoder; VCCA, Variational Canonical Correlation Analysis.

Model	Accuracy (%)	Coarse Accuracy (%)
DenseNet-scratch	67.33 ± 1.35	75.67 ± 1.15
Softmax	71.67 ± 0.28	83.34 ± 0.32
AE _x +Softmax	70.69 ± 0.82	82.42 ± 0.58
VAE _x +Softmax	69.20 ± 0.46	81.24 ± 0.63
SplitAE _{xy}	70.34 ± 0.56	82.11 ± 0.38
VCCA _{xy}	70.72 ± 0.56	82.12 ± 0.61
SplitAE _{xi} +Softmax	77.68 ± 0.69	87.09 ± 0.53
VCCA _{xi} +Softmax	77.02 ± 0.51	86.46 ± 0.42
VCCA-private _{xi} +Softmax	73.04 ± 0.56	84.16 ± 0.51
SplitAE _{xiy}	77.43 ± 0.80	87.14 ± 0.57
VCCA _{xiy}	77.22 ± 0.55	86.54 ± 0.51
VCCA-private _{xiy}	74.04 ± 0.83	84.59 ± 0.83
SplitAE _{xw} +Softmax	76.27 ± 0.66	86.45 ± 0.56
VCCA _{xw} +Softmax	75.37 ± 0.46	86.00 ± 0.32
VCCA-private _{xw} +Softmax	75.11 ± 0.81	85.91 ± 0.55
SplitAE _{xwy}	75.78 ± 0.84	86.13 ± 0.63
VCCA _{xwy}	74.72 ± 0.85	85.59 ± 0.78
VCCA-private _{xwy}	74.92 ± 0.74	85.59 ± 0.67
SplitAE _{xiw} +Softmax	77.79 ± 0.48	87.12 ± 0.62
VCCA _{xiw} +Softmax	77.51 ± 0.51	86.69 ± 0.41
SplitAE _{xiwy}	78.18 ± 0.53	87.26 ± 0.46
VCCA _{xiwy}	77.78 ± 0.45	86.88 ± 0.47

and VCCA-private model using the text descriptions with different description lengths T . Figure 5 shows the fine-grained classification accuracies for the concerned models. For models using only the text descriptions, the classification accuracies increase as T increases in most cases. Setting $T \geq 32$ results in good classification performance, potentially since the models have learned to separate the grocery items based on that the text descriptions have become more dissimilar and unique. The classification accuracies are mostly stable as T varies for the models with the additional iconic images. Since including iconic images significantly increases the classification performance over models only using text descriptions, we conclude that the iconic images are more helpful when we want to classify the grocery items from natural images.

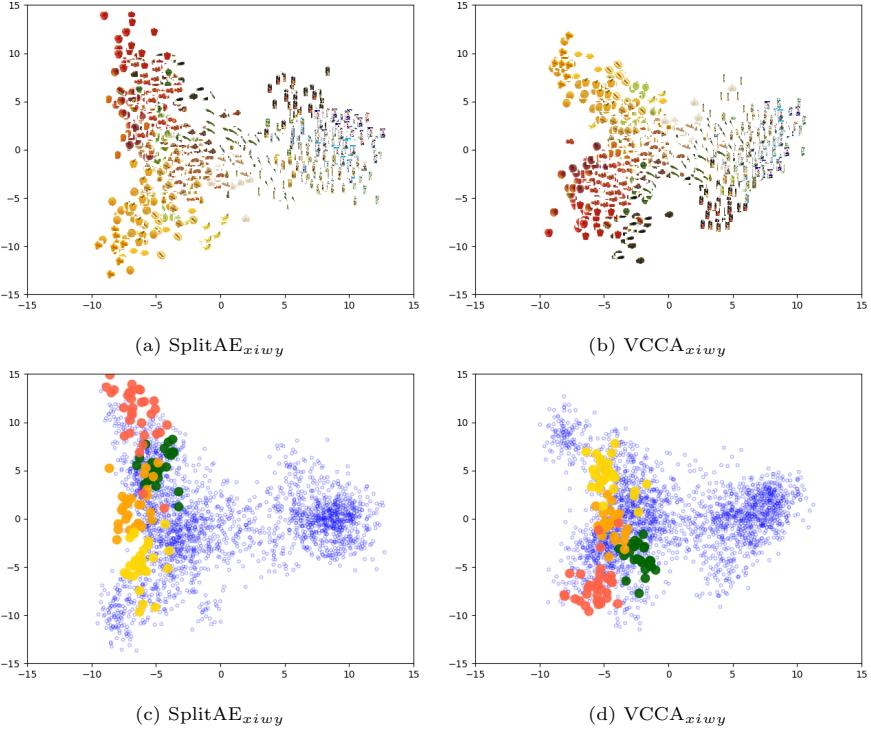


Figure 4: Visualizations of the latent representations of the test set from SplitAE_{xiwy} and VCCA_{xiwy} . We plot the corresponding iconic image to each latent representation in Figure 4a and 4b. In Figure 4c and 4d, we plot the bell pepper representations according to the color of the class, while the blue points correspond to the other grocery items. Abbreviations: SplitAE, Split Autoencoder; VCCA, Variational Canonical Correlation Analysis.

2.4 Investigation of the Learned Representations

To gain insights about the effects that each view has on the classification performance, we visualize the latent space by plotting the latent representations using PCA. Utilizing the additional views showed to have similar effects on the structure of the latent spaces from SplitAE and VCCA. Since our main interest lies in representation learning with variational methods, we focus on studying the latent representations of VCCA and VCCA-private. Firstly, we use PCA to visualize the latent representations in 2D and plot the corresponding iconic images of the representations (see Figure 6). Secondly, to illustrate the effects that the iconic images and text descriptions have on the learned latent space, we focus on two cases of grocery items where one of the views helps to separate two different types of classes and the other one does not (see Figure 7 and 8). Finally, we look into the shared and private latent spaces learned by VCCA-private_{xw} and observe that variations in image backgrounds and structures of text sentences have been separated from the shared representations into the private ones.

In Figure 6, we show the latent representations for the VCCA models that were used in Table 1 (see subsection Classification Results in Results). We also plot the

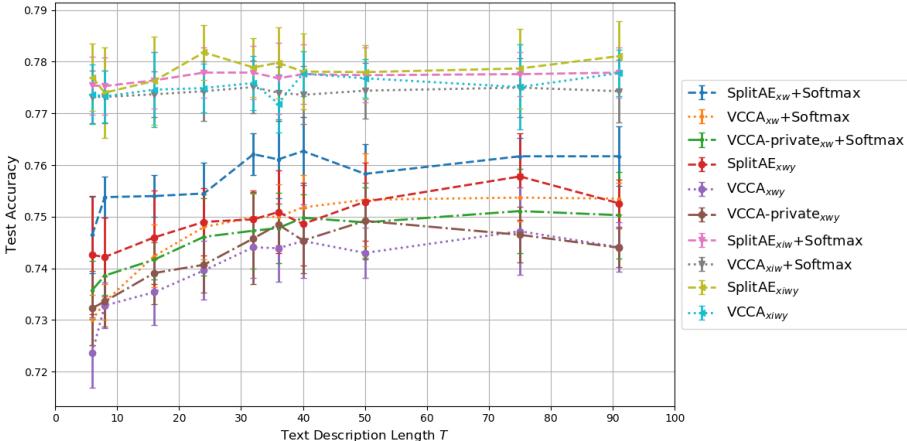


Figure 5: Test accuracy over the text description length T for all SplitAE, VCCA, and VCCA-private models using the text description. We show the accuracy for the models trained with $T = 6, 8, 16, 24, 32, 36, 40, 50, 75$, and 91 words. The results have been averaged for 10 different seeds and the error bars show the standard deviations for every setting of T . Abbreviations: SplitAE, Split Autoencoder; VCCA, Variational Canonical Correlation Analysis.

PCA projections of the natural image features from the off-the-shelf DenseNet169 in Figure 6a as a baseline. Figure 6b and 6c shows the latent space learned by n VAE_x and VCCA_{xy}, which are similar to the DenseNet169 feature space since these models are performing compression of the natural image features into the learned latent space. We observe that these models have divided packages and raw food items into two separate clusters. However, the fruits and vegetables are scattered across their cluster and the packages have been grouped close to each other despite having different colors, e.g., black and white, on the cartons.

The structure of the latent spaces becomes distinctly different for the VCCA models that use either iconic images or text descriptions as an additional view and we can observe the different structures that the views bring to the learned latent space. In Figure 6d and 6g, we see that visually similar objects, in terms of color and shape, have moved closer together by utilizing iconic images in VCCA_{xi} and VCCA_{xiwy}. When using text descriptions in VCCA_{xw} and VCCA_{xwy}, we also observe in the fruit and vegetable cluster that the items are more grouped based on their color in Figure 6e and 6h. Figure 6f and 6i shows the latent spaces in VCCA_{xiw} and VCCA_{xiwy} respectively. These latent spaces are similar to the ones learned by VCCA_{xi} and VCCA_{xiwy} in the sense that these latent spaces also group items based on their color and shape. We believe that this structure imposed by the iconic images could be softened by reducing the scaling weight λ_i , which potentially could reduce the classification accuracy as a consequence. The difference between the latent spaces is not evident comparing the models using the class label.

Red and Green Apples To showcase how the iconic images help to learn good representations, we consider all of the apple classes in the dataset, namely the red apples *Pink Lady*, *Red Delicious* and *Royal Gala*, and also the green apples *Golden*

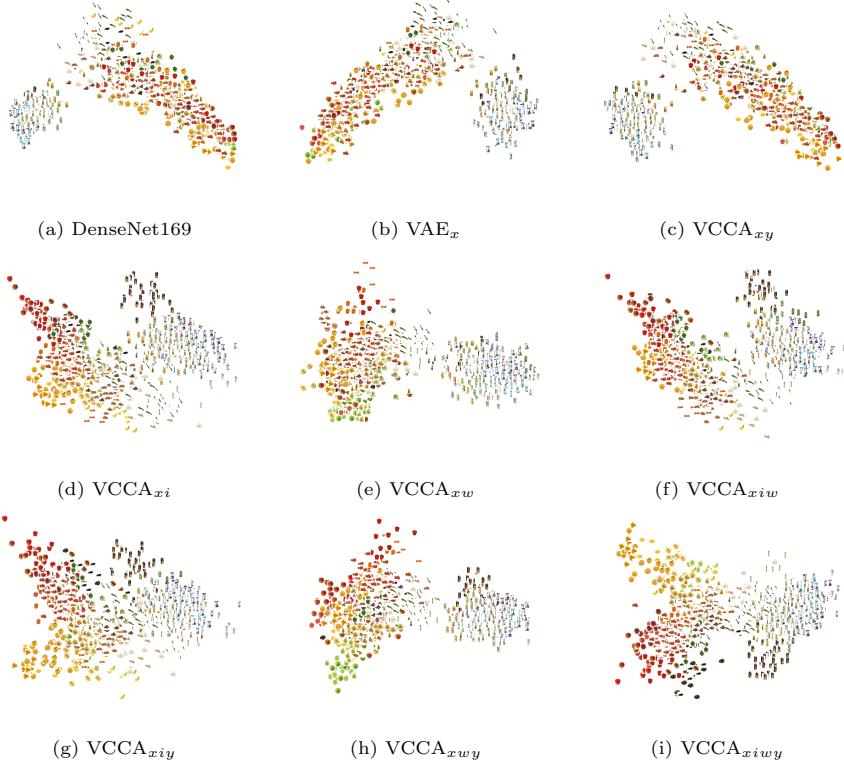


Figure 6: Visualizations of the latent representations from the test set, where we plot the iconic image of the corresponding object classes. We also plot the PCA projection of the natural image features from the off-the-shelf DenseNet169 in Figure 6a. All models have been initialized with the same random seed before training.

Delicious and *Granny Smith*. In Figure 7, we group the red apple classes and visualize their latent representations by red points. The green apples are grouped similarly and we visualize their latent representations with green points. Latent representations of all other grocery items are visualized as blue points. The models using iconic images as one view in Figure 7a, 7d, 7c, and 7f have managed to separate the red and green apples based on their color differences. The models using text description have instead moved the apples closer together in one part of the latent space, possibly because of their similarities mentioned in the description.

Juice and Yoghurt Packages To illustrate how the text descriptions can establish more useful latent representations, we consider a selection of juice and yoghurt classes. These packaged items have similar shapes and colors, which makes it difficult for a classifier to distinguish their content differences using only visual input. In Figure 8, we visualize the latent representations of the juice and yoghurt packages using yellow and green points respectively. We observe that only VCCA_{xw} and VCCA_{xwy} manages to separate the packages in Figure 8b and 8e due to their different text descriptions. Since the iconic images of the packages are visually similar, adding this

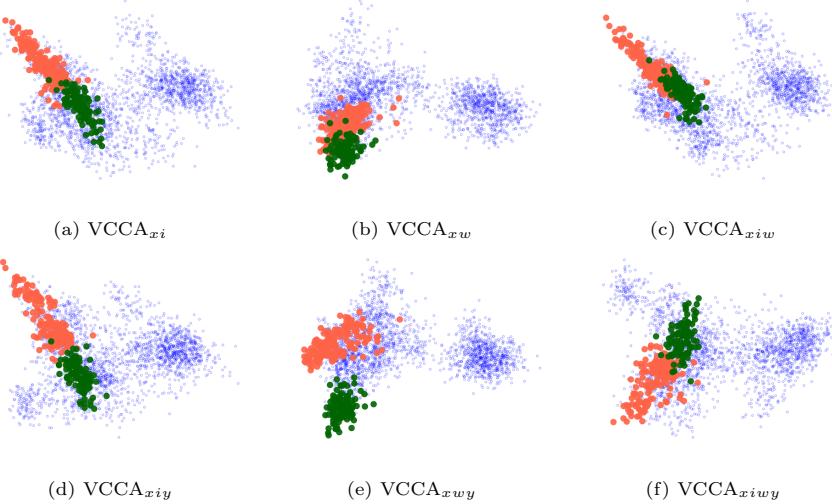


Figure 7: Visualizations of the latent representations μ_z of the red and green apples in the Grocery Store dataset. The red points correspond to the red apple classes, while the green points correspond to the green apple. The blue points correspond to the other grocery items.

information is insufficient for separating these packages in the latent space. This indicates that we gain different benefits from the iconic images and text descriptions when it comes to classifying grocery items.

Latent Spaces of VCCA-private We show the shared and private latent spaces of VCCA-private_{xw} in Figure 9b-d, as well as the single latent space of VCCA_{xw} for comparison in Figure 9a. Comparing with the latent space of VCCA_{xw}, the shared latent space of VCCA-private_{xw} in Figure 9b, has structured the raw food items based on their class, color, and shape better than standard VCCA_{xw}. In Figure 9c, we plot the natural images corresponding to the latent representation for the private latent variable u_x . We zoom in on some natural images and found that the images are structured based on their similarities in background and camera view. On the left and bottom sides of the cluster, we found images of grocery items closely packed together in bins. Single items that are held in the hand of the photographer are placed on the right side, whereas images of items and the store floor are placed on the top and middle of the latent space. The model has therefore managed to separate the variations within the natural image view into the private latent variable u_x from the shared latent variable z , which probably is the main reason why similar raw food items are closer to each other in Figure 9b than in Figure 9a. We also plot the corresponding iconic image on the position of the text description representation for the private latent variable u_w in Figure 9d. Note that every text description is projected at the same location in the latent space since the text descriptions are the same for every class item. We highlighted some specific words in the descriptions and observed that descriptions with the same words are usually close to each other. Visually dissimilar items can be grouped close to each other in this latent space, which indicates that the private latent variable u_w contains information about the structure

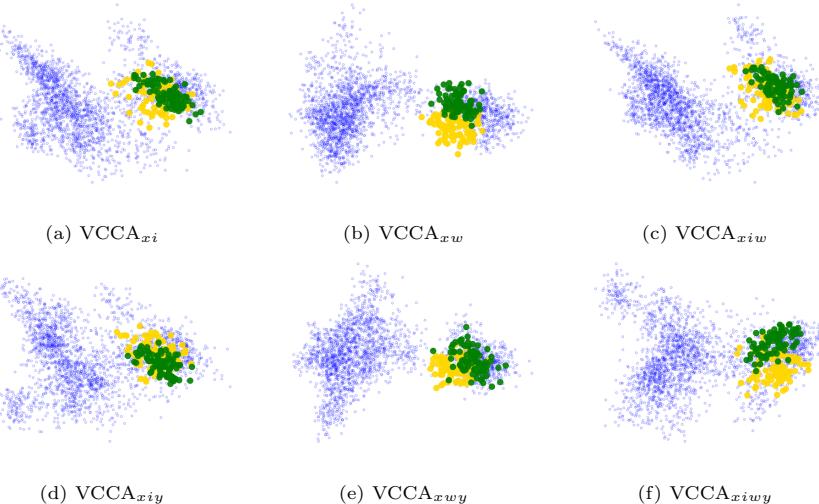


Figure 8: Visualizations of the latent representations μ_z of a selection of juice packages and yoghurt packages in the Grocery Store dataset. The yellow and green points correspond to the juice and yoghurt packages respectively. The blue points correspond to the other grocery items.

of the text sentences, i.e., word occurrences and how they are ordered in the text description. In Figure S4, we show the shared and private latent spaces of VCCA-private_{*xi*} and provide a conclusion to the results in the Supplemental Experimental Procedures.

2.5 Decoding Iconic Images from Unseen Natural Images

In this section, we show how the iconic image decoder can decode plausible iconic images of grocery items from unseen natural images in the test set. We apply the same approach as in [?], where we encode unseen natural images and then decode the retrieved latent representation back into an iconic image. We also report several image similarity metrics to investigate if the decoded image quality is correlated with classification performance. We report peak signal-to-noise ratio (PSNR), structural similarity (SSIM) [79], and the KL divergence by comparing the decoded iconic images to the true ones. For computing the KL divergence, we model the decoded and true iconic image using two Gaussian Mixture Model (GMM) [80, 81]. The images are then represented as density functions, such that we can measure the similarity between the two densities with the KL divergence and use it as an image similarity metric. Since the KL divergence between two GMMs is not analytically tractable, we apply Monte Carlo simulation using n i.i.d. samples drawn from the decoded image density for approximating the KL divergence [82]. Due to the simple structure of the iconic images, we fit the GMMs with $K = 2$ Gaussian components using the RGB color values and draw $n = 100$ Monte Carlo samples to estimate the KL divergence in all experiments. Table 2 shows the image similarity metrics between the VCCA models using the iconic images. We also show the model classification accuracy for each model, which have been taken from Table 1. The models perform on par on the

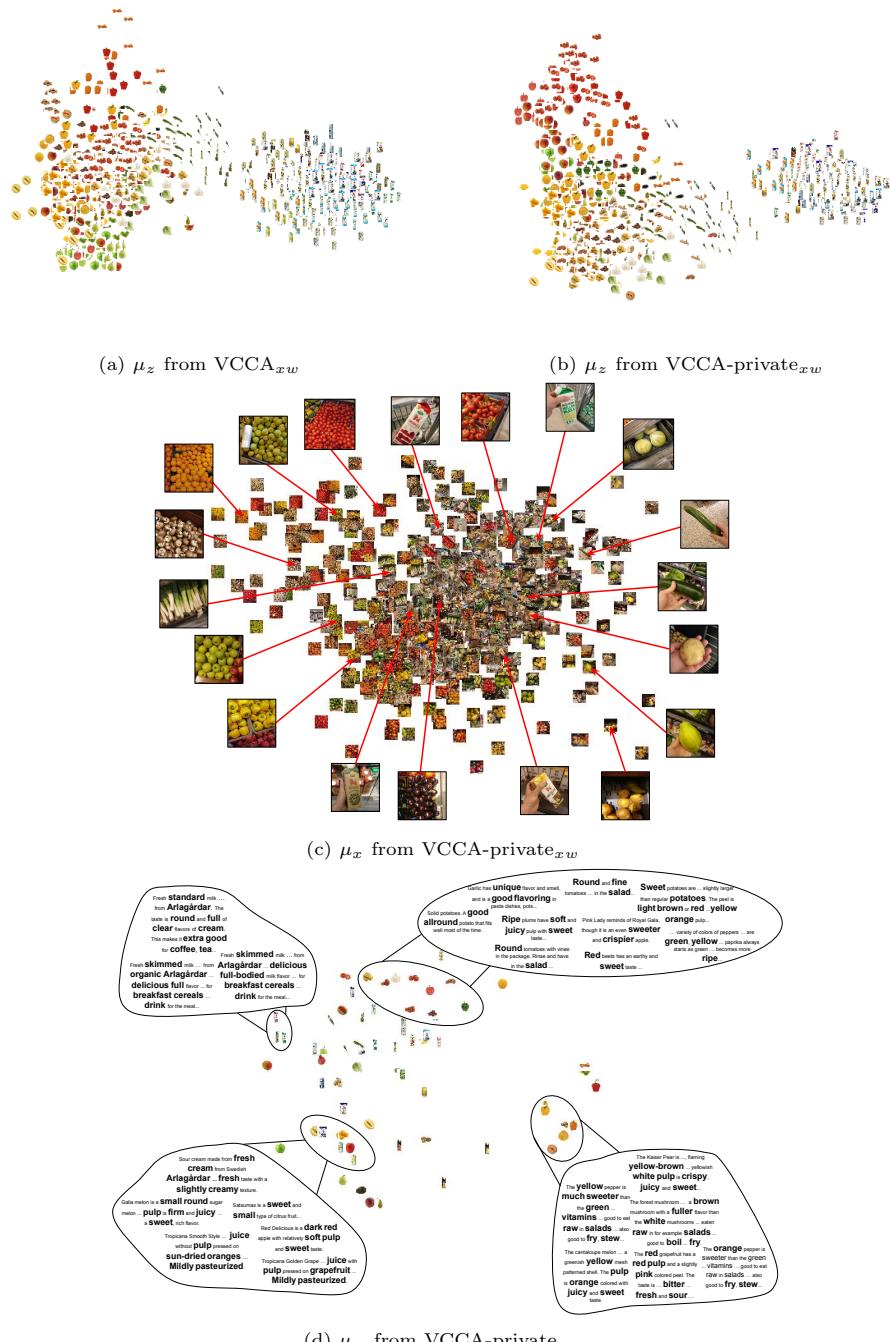


Figure 9: Visualizations of the latent representations μ_z of a selection of juice packages and yoghurt packages in the Grocery Store dataset. The yellow and green points correspond to the juice and yoghurt packages respectively. The blue points correspond to the other grocery items.

image similarity metrics, which indicates that the quality of the decoded images is intact if the model extends to utilizing text descriptions and class labels in addition to the iconic images.

In Figure 10, we display five different natural images from the test set, their true corresponding iconic image, the decoded iconic image from VCCA_{xiwy} . We also show the true and predicted labels from the class label decoder (see Pred. Label). Additionally, we report the image similarity metrics PSNR, SSIM, and KL divergence between the decoded and true iconic images. For the Mango and Royal Gala images, we observe that the decoded images are visually plausible, in terms of recognized item, color, and shape, in both cases, which coheres with the high PSNR and SSIM values and low KL values. The third row shows a shelf with orange and green bell peppers where the decoded image has indeed been decoded into a mix of a green and orange bell pepper. In the two succeeding rows, we display failure cases where the model confuses the true class label with other grocery items. We observe that each metric drops according to the mismatch between decoded and true iconic images. The fourth row shows a basket of Anjou pears, where the model confuses the pear with a Granny Smith apple which can be seen in the decoded image. In the fifth row, there are red milk packages stacked behind a handheld sourcream package, where the decoded image becomes a blurry mix of the milk and sourcream package. Although the predicted class is incorrect, we observe that the prediction is reasonable based on the decoded iconic image.

Table 2: Results on image quality of decoded iconic images for the Variational Canonical Correlation Analysis (VCCA) models using the iconic images. The subscript letters in the model names indicate the data views used in the model. \uparrow denotes higher is better, \downarrow lower is better. Peak Signal-to-Noise Ratio (PSNR), Structural Similarity (SSIM), and Kullback-Leibler (KL) divergence are measured by comparing the true iconic image against the decoded one. Accuracy shows the classification performance for each model and has been taken from Table 1. We report means and standard deviations averaged over 10 random seeds for all metrics.

Model	PSNR \uparrow	SSIM \uparrow	KL \downarrow	Accuracy (%) \uparrow
VCCA_{xi}	20.13 ± 0.05	0.72 ± 0.00	4.43 ± 0.21	77.02 ± 0.51
VCCA_{xiy}	20.12 ± 0.09	0.73 ± 0.00	4.35 ± 0.22	77.22 ± 0.55
VCCA_{xiw}	20.11 ± 0.09	0.73 ± 0.00	4.29 ± 0.24	77.51 ± 0.51
VCCA_{xiwy}	20.16 ± 0.08	0.73 ± 0.00	4.32 ± 0.22	77.78 ± 0.45

3 Discussion

In this section, we summarize the experimental results and discuss our findings.

Classification Results In the first experiments (see subsection Classification Results in Results), we showed that utilizing all four views with SplitAE_{xiwy} and VCCA_{xiwy} resulted in the best classification performance on the test set. This indicates that these models take advantage of each view to learn representations that enhance their generalization ability compared to the single-view models. Moreover, using either the iconic images, the text descriptions or both views yields better representations for classification compared to using the natural image features alone.

Natural Image	Iconic Image	Decoded Image	Classification	Metrics
			True Label: Mango Pred. Label: Mango	PSNR: 25.31 SSIM: 0.88 KL: 0.36
			True Label: Royal Gala Pred. Label: Royal Gala	PSNR: 25.31 SSIM: 0.88 KL: 0.36
			True Label: Orange Bell Pepper Pred. Label: Orange Bell Pepper	PSNR: 25.31 SSIM: 0.88 KL: 0.36
			True Label: Anjou Pred. Label: Granny Smith	PSNR: 25.31 SSIM: 0.88 KL: 0.36
			True Label: Arla Eco. Sourcream Pred. Label: Arla Std. Milk	PSNR: 25.31 SSIM: 0.88 KL: 0.36

Figure 10: Examples of decoded iconic images from $VCCA_{xiwy}$ with their corresponding natural image and true iconic image as well as predicted labels and image similarity metrics. The column Classification shows the true label for the natural image (True Label) and the label predicted by the model (Pred. Label).

Note that it was necessary to properly set the scaling weights λ on the reconstruction losses of the additional views to achieve good classification performance (see Table S2). Whenever the weight values are increased, the model tries to structure the representations according to variations between the items in the upweighted view rather than structuring the items based on visual features in the natural images, e.g., shape, color, and pose of items, image backgrounds, etc. Thus, the latent representations are structured based on semantic information that describes the item itself, which is important for constructing robust representations that generalize well in new environments. Furthermore, the class label decoders performed on par with the separate Softmax classifiers in most cases. The upside with training a separate classifier is that we only would have re-train the classifier if we receive a new class in the dataset, while we would have to train the whole model from scratch when the model uses a class label decoder. Note that the encoder for any of the multi-view models can be used for extracting latent representations for new tasks, whether the model utilizes the label or not, since the encoder only uses natural images as input.

Iconic Images vs. Text Descriptions In Table 1, the iconic images yielded higher classification accuracies compared to using the text descriptions. This was also evident in Figure 5 where the classification performance remains more or less

the same regardless of the text description length T when the models utilize iconic images. We believe that the main reasons for the advantages with iconic images lie in the clear visual features of the items in these images, e.g., their color and shape, which carry lots of information that is important for image classification tasks. However, we also observed that iconic images and text descriptions can yield different benefits for constructing good representations. In Figure 6, we see that iconic images and text descriptions make the model construct different latent representations of the grocery items. Iconic images structures the representations with respect to color and shape of the items (see Figure 7), while the descriptions groups items based on their ingredients and flavor (see Figure 8). Therefore, the latent representations benefit differently from utilizing the additional views and a combination of all of them yields the best classification performance as shown in Table 1. We want to highlight the results in Figure 8, where the model manages to separate juice and yoghurt packages based on their text description. Refrigerated items, e.g., milk and juice packages, have in general very similar shapes and the same color if they come from the same brand. There are minor visual differences between items of the same brand that makes it possible to differentiate between them, e.g., the picture of the main ingredient on the package and ingredient description. Additionally, these items can be almost identical depending on which side of the package that is present on the natural image. When utilizing the text descriptions, we add useful information on how to distinguish between visually similar items that have different ingredients and contents. This is highly important for using computer vision models to distinguish between packaged items without having to use other kinds of information, e.g., barcodes.

Text Description Length We showed that the text descriptions are useful for the classification task, and that careful selection of the description length T is important for achieving the best possible performance (see subsection Classification Results in Results). In Figure 5, we observed that most models achieve significantly better classification performance when the text description length T increases up until $T = 32$. The reason for this increase is due to the similarities between the descriptions of items from the same kind or brand, such as milk and juice packages. For instance, in Table S1, the first sentence in the descriptions for the milk packages only differ by the ninth word, which is *organic* for the ecological milk package. This means that their descriptions will be identical when $T = 8$. Therefore, the descriptions will become more different from each other as we increase T , which helps the model to distinguish between items with similar descriptions. However, the classification accuracies have more or less saturated when setting $T > 32$, which is also due to the similarity between the descriptions. For example, the bell pepper descriptions in Table S1 only differ by the second word that describes the color of the bell pepper, i.e., the words *yellow* and *orange*. We also see that the third and fourth sentences in the descriptions of the milk packages are identical. The text descriptions typically have words that separate the items in the first or second sentence, whereas the following sentences provide general information on ingredients and how the item can be used in cooking. For items of the same kind but of different colors or brand, e.g., bell peppers or milk packages respectively, the useful textual information for constructing good representations of grocery items typically comes from a few words in the description that describes features of the item. Therefore, the models yield

better classification performance when T is set to include at least the whole first sentence of the description. We could select T more cleverly, e.g., by using different T for different descriptions to make sure that we utilize the words that describe the item or filter out non-informative words for the classification task.

VCCA vs. VCCA-private The main motivation for using VCCA-private is to use private latent variables for modeling view-specific variations, e.g., image backgrounds and writing styles of text descriptions. This could allow the model to build shared representations that more efficiently combine salient information shared between the views for training better classifiers. This would then remove noise from the shared representation since the private latent variables are responsible for modeling the view-specific variations. For VCCA-private_{xw}, we observed that the private latent spaces managed to group different image backgrounds and grocery items with similar text descriptions respectively in Figure 9c and 9d respectively. This model also performed on par with VCCA_{xw} regarding classification performance in Table 1. However, we also saw in the same table that the VCCA-private models using the iconic image perform poorly on the classification task compared to their VCCA counterpart. The reason for why this model fails is because of a form of *posterior collapse* [83] in the encoder for the iconic image, where the encoder starts outputting random noise. We noticed this as the KL divergence term for the private latent variable converged to zero when we trained the models for 500 epochs (see Figure S2 and S3), which means that the encoder outputs a distribution which equals a Gaussian prior. The same phenomenon occurs for VCCA-private with the text description as well. We have also experimented with other models with encoders, such as Deep CCA and DCCAE, which also suffered from the collapsing encoder problem for the additional views. Therefore, we believe that the collapsing effect is a consequence of only having access to a single iconic image and text description for every grocery item. Therefore, a potential solution would be to extend the dataset with multiple web-scraped iconic images and text descriptions for every grocery item, which would then establish some variability within the view for each item class. Another possible solution would be to use data augmentation techniques to create some variability in the web-scraped views. For example, we could take a denoising approach and add noise to the iconic images which would force the decoder to reconstruct the real iconic images [84]. For the text descriptions, we could mask words at random in the encoder and let the decoder predict the whole description, which would work as a form of *word dropout* [83, 85]. We leave this for future work if such augmentation techniques can create the needed variability for learning more robust representations as well as discovering the structures of the private latent spaces that this approach would bring.

Decoded Iconic Images We observed with image similarity metrics that the quality of decoded iconic images coheres to some extent with the classification performance for VCCA models using the iconic images (see subsection Decoding Iconic Images from Unseen Natural Images in Results). We also showed that the decoded images are visually plausible decoded images with respect to colors, shapes, and identities of the grocery items in the dataset. The values for the metrics PSNR, SSIM, and KL divergence are similar across the different VCCA models. Since we

used RGB values for estimating KL, we believe that including spatial information of pixels or using other color spaces, e.g., Lab, could provide more information about the dissimilarities between the decoded and true iconic images. To thoroughly assess the relationship between good classification performance and accurately decoding the iconic images, we also suggest evaluating the image quality on other image similarity metrics, e.g., perceptual similarity [86]. Finally, we see the decoding of iconic images as a promising method to evaluate the quality of the latent representations as well as enhancing the interpretability of the classification. For example, we could inspect decoded iconic images qualitatively or by using image similarity metrics to determine how certain the model was about the present items in the natural images, which could then be used as a tool for explaining misclassifications.

3.1 Conclusions

In this paper, we introduce a dataset with natural images of grocery items taken in real grocery store environments. Each item class is accompanied by web-scraped information in the form of an iconic image and a text description of the item. The main application for this dataset is for training image classifiers that can assist visually impaired people when shopping for groceries but is not limited to this use case only.

We selected the multi-view generative model VCCA that can utilize all of the available data views for image classification. To evaluate the contribution to the classification performance for each view, we conducted an ablation study comparing classification accuracies between VCCA models with different combinations of the available data types. We showed that utilizing the additional views with VCCA yields higher accuracies on classifying grocery items over models only using the natural images. The iconic images and text descriptions impose different structures of the shared latent space, where we observed that iconic images help to group the items based on their color and shape while text descriptions separate the items based on differences in ingredients and flavor. These types of semantics that VCCA has learned can be useful for generalizing to new grocery items and other object recognition tasks. We also investigated VCCA-private, which introduces private latent variables for view-specific variations, that separates the latent space into shared and private spaces for each view to provide high-quality representations. However, we observed that the private latent variables for the web-scraped views became uninformative by modeling noise due to the lack of variations in the additional web-scraped views. This encourages to explore new methods for extracting salient information from such data views that can be beneficial for downstream tasks.

An evident direction of future work would be to investigate other methods for utilizing the web-scraped views more efficiently. For instance, we could apply pre-trained word representations for the text description, e.g., BERT [85] or GloVe [87], to see if they enable the construction of representations that can more easily distinguish between visually similar items. Another interesting direction would be to experiment with various data augmentation techniques in the web-scraped views to create view-specific variations without the need for collecting and annotating more data. It is also important to investigate how the model can be extended to recognize multiple items. Finally, we see zero- and few-shot learning [67] of new grocery items and

transfer learning [88] as potential applications where our dataset can be used for benchmarking of multi-view learning models on classification tasks.

4 Experimental Procedures

4.1 Resource Availability

4.1.1 Lead Contact

Marcus Klasson is the lead contact for this study and can be contacted by email at mklas@kth.se.

4.1.2 Materials Availability

There are no physical materials associated with this study.

4.1.3 Data and Code Availability

1. The Grocery Store dataset along with documentation is available at the following Github repository: <https://github.com/marcusklasson/GroceryStoreDataset>
2. The source code for the multi-view models along with documentation is available at the following Github repository: https://github.com/marcusklasson/vcca_grocerystore

4.2 Methods

In this section, we outline the details of the models we use for grocery classification. We begin by introducing autoencoders and SplitAEs [18]. Then we describe VAEs [74] and how it is applied to single-view data, followed by the introduction of VCCA [7] and how we adapt it to our dataset. We also discuss a variant of VCCA called VCCA-private [7], which is used for extracting private information about each view in addition to shared information across all views by factorizing the latent space. The graphical model representations of the VAE, VCCA, and VCCA-private models that have been used in this paper are shown in Figure S5. The model names use subscripts to denote the views utilized for learning the shared latent representations. For example, VCCA_{*xi*} utilizes natural image features *x* and iconic images *i*, while VCCA_{*xiwy*} uses natural image features *x* and iconic images *i*, text descriptions *w*, and class labels *y*.

4.2.1 Autoencoders and Split Autoencoders

The autoencoding framework can be used for feature extraction and learning latent representations of data in unsupervised manners [89]. It begins with defining a parameterized function called the encoder for extracting features. We denote the encoder as f_ϕ where ϕ includes its parameters, which commonly are the weights and bias vectors of a neural network. The encoder is used for computing a feature vector $h = f_\phi(x)$ from the input data *x*. Another parameterized function g_θ called the decoder is also defined, that maps the feature *h* back into input space, i.e., $\hat{x} = g_\theta(h)$.

The encoder and decoder are learned simultaneously to minimize the reconstruction loss between the input and its reconstruction of all training samples. By setting the dimension of the feature vector smaller than the input dimension, i.e., $d_h < d_x$, the autoencoder can be used for dimensionality reduction which makes the feature vectors suitable for training linear classifiers in a cheap way.

As in the case for the Grocery Store dataset, we have multiple views available during training, while only the natural image view is present at test time. In this setting, we can use a Split Autoencoder (SplitAE) to extract shared representations by reconstructing all views during training from the one view that is available during the test phase [17, 18]. As an example, we have the two-view case with x present at both training and test while y is only available during training. We therefore define an encoder f_ϕ and two decoders g_{θ_x} and g_{θ_y} where both decoders inputs the same representation $h = f_\phi(x)$. The objective of the SplitAE is to minimize the sum of the reconstruction losses, which will encourage representations h that best reconstructs both views. The total loss is then

$$\mathcal{L}_{\text{SplitAE}}(\theta, \phi; x, y) = \lambda_x \mathcal{L}_x(x, g_{\theta_x}(h)) + \lambda_y \mathcal{L}_y(y, g_{\theta_y}(h)), \quad (2)$$

where $\theta_x, \theta_y \in \theta$ and λ_x, λ_y are scaling weights for the reconstruction losses. For images, the reconstruction loss can be the mean squared error, while the cross-entropy loss is commonly used for class labels and text. This architecture can be extended to more than two views by simply using view-specific decoders that input the shared representation extracted from natural images. Note that in the case when the class labels are available, we can use the class label decoder g_{θ_y} as a classifier during test time. Alternatively, we can train a separate classifier with the learned shared representations after the SplitAE has been trained.

4.2.2 Variational Autoencoders with only Natural Images

The Variational Autoencoder (VAE) is a generative model that can be used for generating data from single views. Here, we describe how the VAE learns latent representations of the data and how the model can be used for classification. VAEs define a joint probability distribution $p_\theta(x, z) = p(z)p_\theta(x|z)$, where $p(z)$ is a prior distribution over the latent variables z and $p_\theta(x|z)$ is the likelihood over the natural images x given z . The prior distribution is often assumed to be an isotropic Gaussian distribution, $p(z) = \mathcal{N}(z | \mathbf{0}, \mathbf{I})$, with the zero vector $\mathbf{0}$ as mean and the identity matrix \mathbf{I} as the covariance. The likelihood $p_\theta(x|z)$ takes the latent variable z as input and outputs a distribution parameterized by a neural network with parameters θ which is referred to as the decoder network. A common distribution for natural images is a multivariate Gaussian, $p_\theta(x|z) = \mathcal{N}(x | \mu_x(z), \sigma_x^2(z) \odot \mathbf{I})$, where $\mu_x(z)$ and $\sigma_x^2(z)$ are the means and standard deviations of the pixels respectively outputted from the decoder and \odot denotes element-wise multiplication. We wish to find latent variables z that are likely to have generated x , which is done by approximating the intractable posterior distribution $p_\theta(z|x)$ with a simpler distribution $q_\phi(z|x)$ [75]. This approximate posterior $q_\phi(z|x)$ is referred to as the encoder network since it is parameterized by a neural network ϕ which inputs x and outputs a latent variable z . Commonly, we let the approximate posterior to be Gaussian $q_\phi(z|x) = \mathcal{N}(z | \mu_z(x), \sigma_z^2(x) \odot \mathbf{I})$, where the mean $\mu_z(x)$ and variance $\sigma_z^2(x)$ are the outputs of the encoder. The latent variable z is then sampled using the mean and variance from the encoder with the

reparameterization trick [74, 90]. The goal is to maximize a tractable lower bound on the marginal log-likelihood of x using $q_\phi(z|x)$:

$$\log p_\theta(x) \geq \mathcal{L}(\theta, \phi; x) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] - D_{\text{KL}}(q_\phi(z|x) || p(z)). \quad (3)$$

The last term is the Kullback-Leibler (KL) divergence [78] of the approximate posterior from the prior distribution, which can be computed analytically when $q_\phi(z|x)$ and $p(z)$ are Gaussians. The expectation can be viewed as a reconstruction loss term, which can be approximated using Monte Carlo sampling from $q_\phi(z|x)$. The lower bound \mathcal{L} is called the evidence lower bound (ELBO) and can be optimized with stochastic gradient descent via backpropagation. The mean outputs $\mu_z(x)$ from the encoder $q_\phi(z|x)$ are commonly used as the latent representation of the natural images x . We can use the representations $\mu_z(x)$ for training any classifier $p(y|\mu_z(x))$, e.g., a Softmax classifier, where y is the class label of x . We can therefore see training the VAE as a pre-processing step, where we extract a lower-dimensional representation of the data x which hopefully makes the classification problem easier to solve. We can also extend the VAE with a generative classifier by incorporating the class label y in the model [91, 92]. Hence, the VAE defines a joint distribution $p_\theta(x, y, z) = p(z)p_{\theta_x}(x|z)p_{\theta_y}(y|z)$ where the class label decoder $p_{\theta_y}(y|z)$ is used as the final classifier. We therefore aim to maximize the ELBO on the marginal log-likelihood over x and y :

$$\log p_\theta(x, y) \geq \mathcal{L}(\theta, \phi; x, y)$$

$$= \lambda_x \mathbb{E}_{q_\phi(z|x)} [\log p_{\theta_x}(x|z)] + \lambda_y \mathbb{E}_{q_\phi(z|x)} [\log p_{\theta_y}(y|z)] - D_{\text{KL}}(q_\phi(z|x) || p(z)). \quad (4)$$

The parameters λ_x and λ_y are used for scaling the magnitudes of the expected values. When predicting the class label for an unseen natural image x^* , we can consider multiple output predictions of the class label by sampling K different latent variables for x^* from the encoder to determine the final predicted class. For example, we could either average the predicted class scores over K or use the maximum class score from K samples as the final predictions. In this paper, we compute the average of the predicted class scores using

$$\hat{y} = \arg \max_y \frac{1}{K} \sum_{k=1}^K p_{\theta_y}(y|z^{(k)}), \quad z^{(k)} \sim q_\phi(z|x^*), \quad (5)$$

where \hat{y} is the predicted class for the natural image x^* . We denote this model as VCCA _{xy} due to its closer resemblance to VCCA than VAE in this paper.

4.2.3 Variational Canonical Correlation Analysis for Utilizing Multi-View Data

In this section, we describe the details of Variational Canonical Correlation Analysis (VCCA) [7] for our application. In the Grocery Store dataset, the views can be the natural images, iconic images, text descriptions, or class labels and we can use any combination of those three in VCCA. To illustrate how we can employ this model to the Grocery Store dataset, we let the natural images x and the iconic images i be the two views. We assume that both views x and i have been generated from a single latent variable z . Similarly as with VAEs, VCCA defines a joint probability distribution $p_\theta(x, i, z) = p(z)p_{\theta_x}(x|z)p_{\theta_i}(i|z)$. There are now two likelihoods for

each view modeled by the decoders $p_{\theta_x}(x|z)$ and $p_{\theta_i}(i|z)$ are represented as neural networks with parameters θ_x and θ_i . Since we want to classify natural images, the other available views in the dataset will be missing when we have received a new natural image. Therefore, the encoder $q_\phi(z|x)$ only uses x as input to infer the latent variable z shared across all views, such that we do not have to use inference techniques that handle missing views. With this choice of approximate posterior, we receive the following ELBO on the marginal log-likelihood over x and i that we aim to maximize:

$$\begin{aligned} \log p_\theta(x, i) &\geq \mathcal{L}(\theta, \phi; x, i) \\ &= \lambda_x \mathbb{E}_{q_\phi(z|x)} [\log p_{\theta_x}(x|z)] + \lambda_i \mathbb{E}_{q_\phi(z|x)} [\log p_{\theta_i}(i|z)] - D_{\text{KL}}(q_\phi(z|x) || p(z)). \end{aligned} \quad (6)$$

The parameters λ_x and λ_i are used for scaling the magnitude of the expected values for each view. We provide a derivation of the ELBO for three or more views in the Supplemental Experimental Procedures. The representations $\mu_z(x)$ from the encoder $q_\phi(z|x)$ can be used for training a separate classifier. We can also add a class label decoder $p_{\theta_y}(y|z)$ to the model and use Equation 5 to predict the class of unseen natural images.

4.2.4 Extracting Private Information of Views with VCCA-private

In the following section, we show how the VCCA model can be altered to extract shared information between the views as well as view-specific private information to enable more efficient posterior inference. Assuming that a shared latent variable z is sufficient for generating all different views may have its disadvantages. Since the information in the views is rarely fully independent nor fully correlated, information only relevant to one of the views will be mixed with the shared information. This may complicate the inference of the latent variables, which potentially can harm the classification performance. To tackle this problem, previous works have proposed learning separate latent spaces for modeling shared and private information of the different views [7, 61, 69]. The shared information should represent the correlations between the views, while the private information represents the independent variations within each view. As an example, the shared information between natural and iconic images are the visual features of the grocery item, while their private information is considered as the various backgrounds that can appear in the natural images and the different locations of non-white pixels in the iconic images. For the text descriptions, the shared information would be words that describe visual features in the natural images, whereas the private information would be different writing styles for describing grocery items with text. We adapt the approach from [7] called VCCA-private and introduce private latent variables for each view along with the shared latent variable. To illustrate how we employ this model to the Grocery Store dataset, we let the natural images x and the text descriptions w be the two views. The joint distribution of this model is written as

$$p_\theta(x, w, z, u_x, u_w) = p_{\theta_x}(x|z, u_x)p_{\theta_w}(w|z, u_w)p(z)p(u_x)p(u_w), \quad (7)$$

where u_x and u_w are the private latent variables for the x and w respectively. To enable tractable inference of this model, we employ a factorized approximate posterior

distribution of the form

$$q_\phi(z, u_x, u_w | x, w) = q_{\phi_z}(z|x)q_{\phi_x}(u_x|x)q_{\phi_w}(u_w|w), \quad (8)$$

where each factor is represented as an encoder network inferring its associated latent variable. With this approximate posterior, the ELBO for VCCA-private is given by $\log p_\theta(x, w) \geq \mathcal{L}_{\text{private}}(\theta, \phi; x, w)$

$$\begin{aligned} &= \lambda_x \mathbb{E}_{q_{\phi_z}(z|x), q_{\phi_x}(u_x|x)} [\log p_{\theta_x}(x|z, u_x)] + \lambda_w \mathbb{E}_{q_{\phi_z}(z|x), q_{\phi_w}(u_w|w)} [\log p_{\theta_w}(w|z, u_w)] \\ &\quad - D_{\text{KL}}(q_{\phi_z}(z|x) || p(z)) - D_{\text{KL}}(q_{\phi_x}(u_x|x) || p(u_x)) - D_{\text{KL}}(q_{\phi_w}(u_w|w) || p(u_w)). \end{aligned} \quad (9)$$

The expectations in $\mathcal{L}_{\text{private}}(\theta, \phi; x, w)$ in Equation 9 can be approximated using Monte Carlo sampling from $q_\phi(z|x)$. The sampled latent variables are concatenated and then used as input to their corresponding decoder. We let the approximated posteriors over both shared and private latent variables to be multivariate Gaussian distributions and their priors distributions to be standard isotropic Gaussians $\mathcal{N}(\mathbf{0}, \mathbf{I})$. The KL divergences in Equation 8 can then be computed analytically. Since only natural images are present during test time and because the shared latent variable z since it should contain information about similarities between the views, e.g., the object class, we use the encoder $q_\phi(z|x)$ to extract latent representations $\mu_z(x)$ for training a separate classifier. As for the VAE and standard VCCA, we can also add a class label decoder $p_{\theta_y}(y|z)$ only conditioned on z to the model and use Equation 5 to predict the class of unseen natural images. We evaluated the classification performance of VCCA-private and compare it to the standard VCCA model only using a single shared latent variable (see subsection Classification Results in Result).

4.3 Experimental Setup

This section briefly describes the network architecture designs and the selection of hyperparameters for the models. See the Supplemental Experimental Procedures for full details on the network architectures and hyperparameters that we use.

Processing of Natural Images We use a DenseNet169 [77] as the backbone for processing the natural images since this architecture showed good classification performance in [6]. As our first baseline, we customize the output layer of DenseNet169 to our the Grocery Store dataset and train it from scratch to classify the natural images. For the second baseline, we train a Softmax classifier on off-the-shelf features from DenseNet169 pre-trained on the ImageNet dataset [19], where we extract 1664-dimensional from the average pooling layer before the classification layer in the architecture. Using pre-trained networks as feature extractors for smaller datasets has previously been proven to be a successful approach for classification tasks [93], which makes it a suitable baseline for the Grocery Store dataset. We denoted the DenseNet169 trained from scratch and the Softmax classifier trained on off-the-shelf features as DenseNet-scratch and Softmax respectively in the Results section.

Network Architectures We use the same architectures for SplitAE and VCCA for a fair comparison. We train the models using off-the-shelf features extracted from a pre-trained DenseNet169 for the natural images. No fine-tuning of the DenseNet

backbone was used in the experiments, which we leave for future research. The image feature encoder and decoder consist of a single hidden layer, where the encoder outputs the latent representation and the decoder reconstructs the image feature. We use a DCGAN [94] generator architecture for the iconic image decoder reconstructing the iconic images. The text description is predicted sequentially using an LSTM network [95]. The label decoder uses a single hidden layer with 512 hidden units and Leaky ReLU activation. The latent space dimension to $d_z = 200$ for all SplitAE and VCCA models in the experiments. In the VCCA-private models, the encoder and decoders have the same architectures as in the VCCA models. We use a reversed DCGAN generator as the iconic image encoder. The text description encoder is an LSTM. We obtain an embedding for the description by averaging all of the hidden states h_t generated from the LSTM, i.e., $\frac{1}{T} \sum_{t=1}^T h_t$, and input it to a linear layer that outputs the latent representation. The dimensions of the private latent spaces are the same as the shared latent space dimension $d_z = 200$.

Training Details In all experiments, we train all models for 200 epochs using the Adam optimizer [96] with initial learning rate 10^{-4} and hyperparameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$. We apply the sum-of-squares loss for the natural and iconic images and the categorical cross-entropy loss for text descriptions and class labels. The latent representations for the AE and SplitAE are the output from the encoder, while for the VAE, VCCA, and VCCA-private we instead use the mean outputs $\mu_z(x)$ from the encoder. In VCCA-private, we use the the mean outputs $\mu_{u_x}(x)$ and $\mu_{u_w}(w)$ for visualizing the private latent representations (see Figure 9). The Softmax classifiers are trained for 100 epochs using with initial learning rate 10^{-4} and hyperparameters $\beta_1 = 0.5$ and $\beta_2 = 0.999$. We run a hyperparameter search for the scaling weights λ for the reconstruction losses of each view (see the Supplemental Experimental Procedures for more details).

Choice of Text Description Length T We wanted to investigate how the classification performance is affected by the text description length T for the SplitAE and VCCA models using the text description w . Since the LSTM predicts the text description sequentially, the computational time increases as the text description length increases. We began by setting $T = 16$ because of how the sentence length was set for the image captioning model in [12]. Then we created an interval by taking steps with 8 words up to $T = 40$ and included the minimum, mean, and maximum text description lengths which are $T = 6, 36$, and 91 respectively. We added two additional points at $T = 50$ and 75 since for most models there was a slight increase in classification performance when T was increased from 40 to 91.

Computational Time The computational time for the SplitAE and VCCA models differs depending on which views that are being used. We measured the number of seconds per epoch (s/epoch) on an Nvidia GeForce RTX 2080 Ti, where an epoch consists of the 2640 training samples. Running experiments with the text description view took around 0.4–2.7 s/epoch with text description length $T \in [6, 91]$. Adding the iconic image view and training the models took around 4.5–6.2 s/epoch depending on the text description length.

Visualization Method We use PCA to visualize the latent space by plotting the latent representations from the trained encoders. PCA is used for projecting the representations from dimension d_z to a 2D space. We obtain the principal components for the representations with the natural images from the training set. In all figures, we plot the representations of test set images given the principal components obtained from the training images (see subsection Investigation of the Learned Representations in Results). To distinguish more easily between the class labels of the images, we occasionally use the iconic images from the dataset and plot the corresponding iconic image of the representation on its location in the 2D space.

1 Supplemental Figures

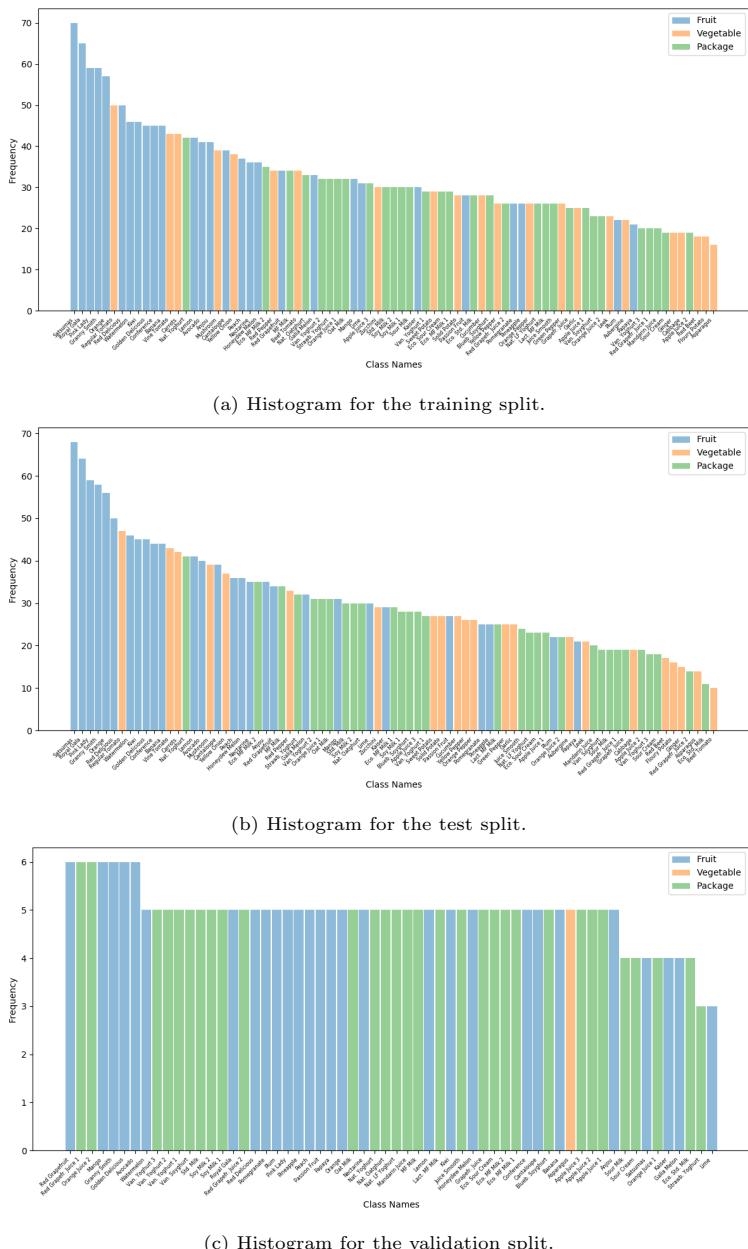


Figure 11: Histograms of natural images for every class in the training (a), test (b), and validation (c) splits of the Grocery Store dataset. We also show with different colors on the bins if the class is either a Fruit, Vegetable, or Package item.

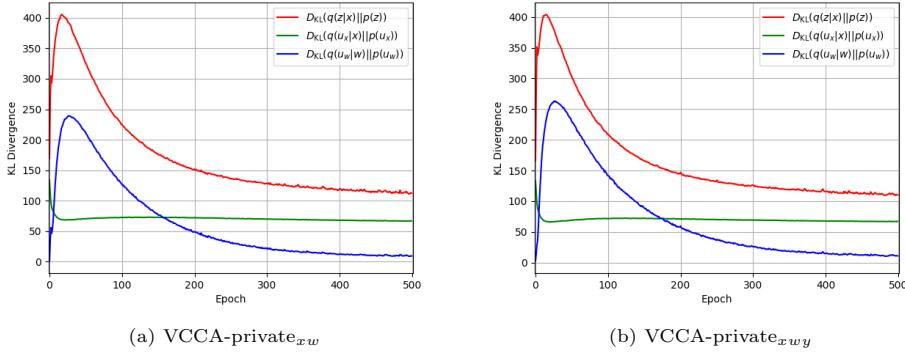


Figure 12: The measured KL divergences $D_{KL}(q_{\phi_z}(z|x) || p(z))$ (red), $D_{KL}(q_{\phi_x}(u_x|x) || p(u_x))$ (green), and $D_{KL}(q_{\phi_w}(u_w|w) || p(u_w))$ (blue) over epochs. We increased the number of epochs from 200 to 500 to demonstrate that the KL divergence for the private latent variable u_w goes to zero. The number of words $T = 24$ is the same for both models. The likelihood weights are set to $\lambda_w = 1000$ and $\lambda_y = 100$. Abbreviations: VCCA, Variational Canonical Correlation Analysis; KL, Kullback-Leibler.

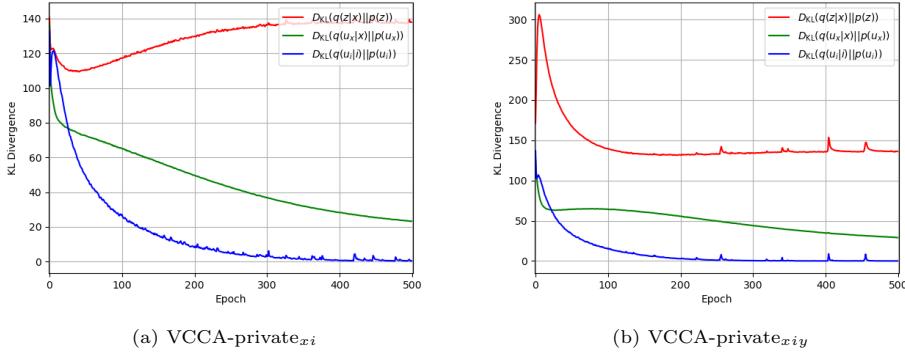


Figure 13: The measured KL divergences $D_{KL}(q_{\phi_z}(z|x) || p(z))$ (red), $D_{KL}(q_{\phi_x}(u_x|x) || p(u_x))$ (green), and $D_{KL}(q_{\phi_i}(u_i|i) || p(u_i))$ (blue) over epochs. We increased the number of epochs from 200 to 500 to demonstrate that the KL divergence for the private latent variable u_i goes to zero. The likelihood weights are set to $\lambda_i = 10$ and $\lambda_y = 1000$. Abbreviations: VCCA, Variational Canonical Correlation Analysis; KL, Kullback-Leibler.

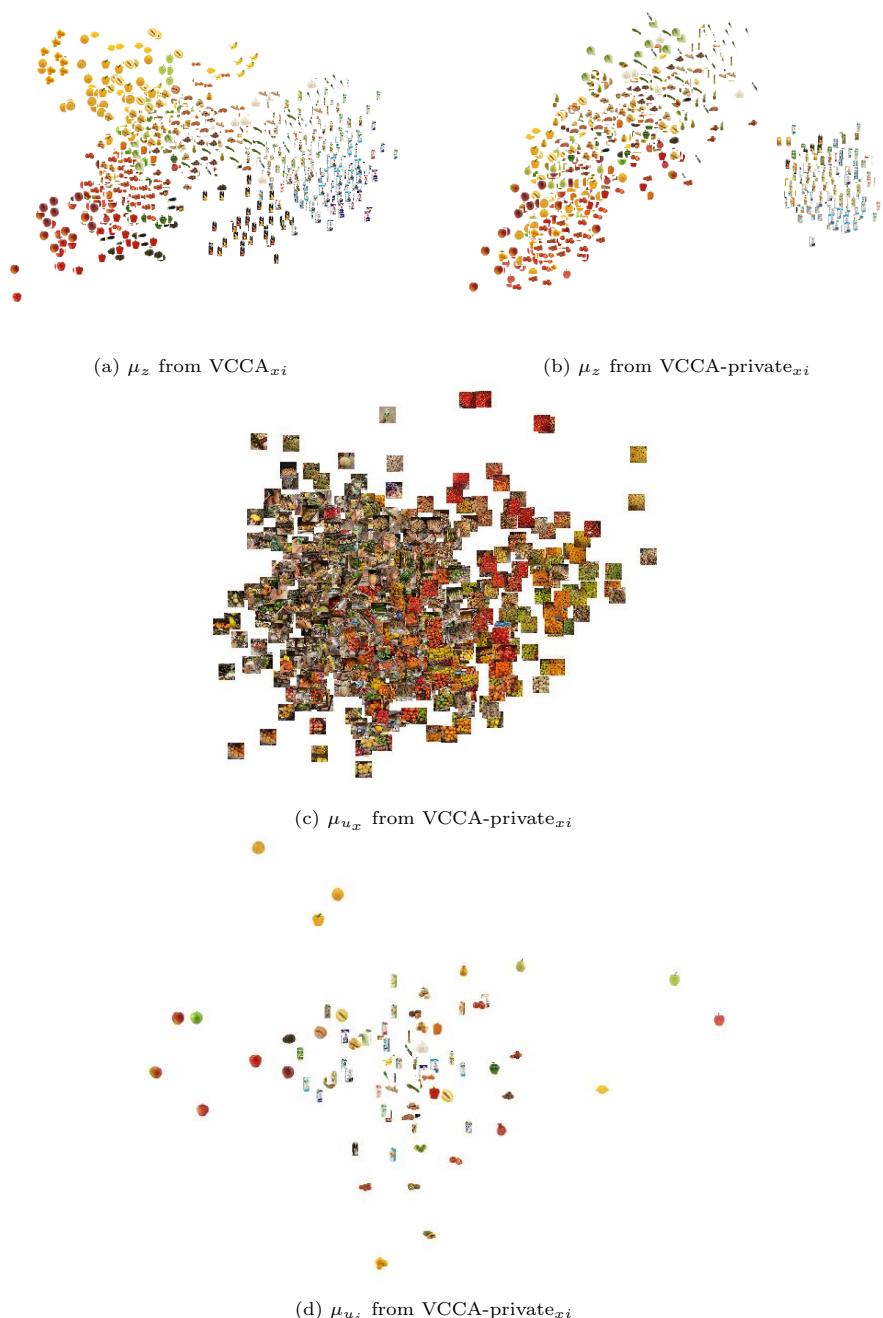


Figure 14: Visualizations of the latent representations $\mu_z(x)$ from VCCA_{xi} and $\text{VCCA-}\text{private}_{xi}$ on the first row followed by $\mu_{u_x}(x)$ and $\mu_{u_i}(i)$ for $\text{VCCA-}\text{private}_{xi}$. Abbreviations: VCCA, Variational Canonical Correlation Analysis.

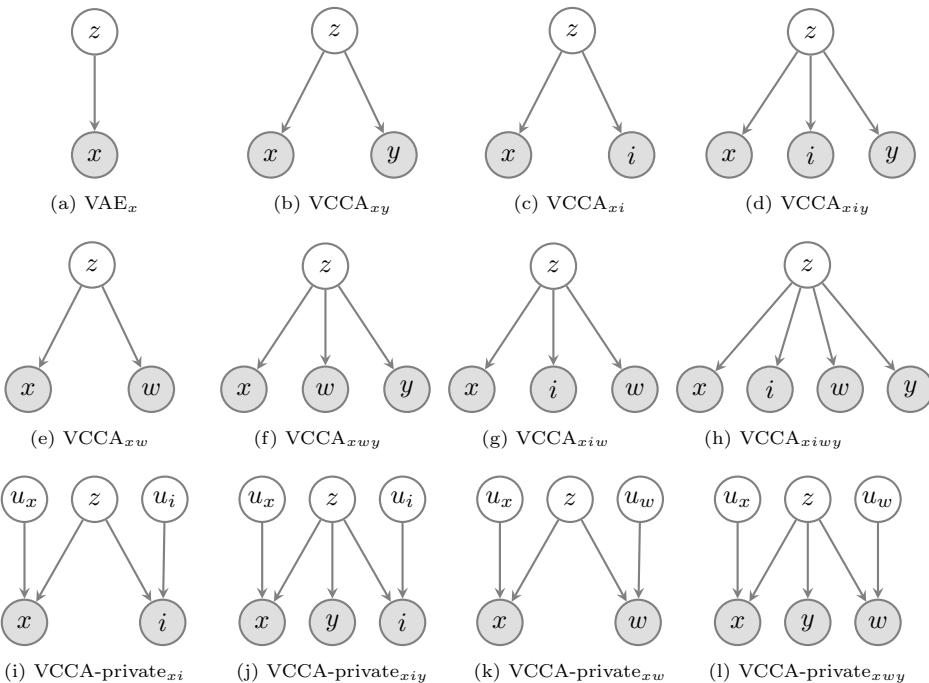


Figure 15: The graphical models of the VAE, VCCA, and VCCA-private, where nodes represent random variables and edges indicate possible dependence. Grey nodes are observed random variables, while white nodes are latent random variables. The joint distribution is given by the product of the conditional distributions of nodes given their parents. Abbreviations: VAE, Variational Autoencoder; VCCA, Variational Canonical Correlation Analysis.

Table 3: Examples of text descriptions and their corresponding class label and iconic image from the Grocery Store dataset.

Class Label	Iconic Image	Text Description
Golden Delicious Apple		Golden Delicious has a white juicy pulp and a greenish yellow shell. The taste is mellow and sweet, making Golden Delicious suitable for desserts.
Red Delicious Apple		Red Delicious is a dark red apple with relatively soft pulp and sweet taste.
Orange		There are many different types of oranges that ripen and are sold during different parts of the year. The orange is a very important vitamin C source and the vitamins are best kept if the fruit is eaten naturally.
Yellow Bell Pepper		The yellow pepper is much sweeter than the green. It also contains more vitamins and antioxidants than the green. Peppers are good to eat raw in salads and as garnish, but are also good to fry, stew or gratinate, for example with filling. Paprika also fits well in pots, gratins and pies.
Orange Bell Pepper		The orange pepper is sweeter than the green. It also contains more vitamins and antioxidants than the green. Peppers are good to eat raw in salads and as garnish, but are also good to fry, stew or gratinate, for example with filling. Paprika also fits well in pots, gratins and pies.
Arla Medium Fat Milk		Fresh skimmed milk made from Swedish milk from Arlagårdar. Skimmed milk has a delicious full-bodied milk flavor and is a popular choice for breakfast cereals, porridge or as a drink for the meal. Milk is a natural source of, for example, protein, calcium and vitamin B12. Protein contributes to muscle building and calcium is needed to maintain a normal bone structure. The brand Arla Ko guarantees that the product is made of 100% Swedish milk.
Arla Ecological Medium Fat Milk		Fresh skimmed milk made from Swedish milk from organic Arlagårdar. Skimmed milk has a delicious full flavor and is a popular choice for breakfast cereals, porridge or as a drink for the meal. Milk is a natural source of, for example, protein, calcium and vitamin B12. Protein contributes to muscle building and calcium is needed to maintain a normal bone structure. The brand Arla Ko guarantees that the product is made of 100% Swedish milk. The new brown carton has 24 percent lower climate impact compared to the previous white carton.
God Morgan Orange Juice		God Morgan Orange is pressed by sun-dried, hand-picked oranges. The package contains juice from 2 kilograms of oranges!
Tropicana Golden Grapefruit Juice		Tropicana Golden Grape is a ready to drink juice with pulp pressed on grapefruit. Not from concentrate. Mildly pasteurized.

Table 4: Classification accuracies on the test set for all models in percentage (%) along with the best hyperparameter setting, i.e., the scaling weights λ and text description length T , for each model. The column Accuracy corresponds to the fine-grained classification accuracy. The column Coarse Accuracy corresponds to the classifying a class within the correct parent class. Results are averaged using 10 different random seeds and we report both means and standard deviations. Abbreviations: AE, Autoencoder; VAE, Variational Autoencoder; SplitAE, Split Autoencoder; VCCA, Variational Canonical Correlation Analysis.

Model	λ_x	λ_i	λ_w	λ_y	T	Accuracy (%)	Coarse Accuracy (%)
DenseNet-scratch	-	-	-	-	-	67.33 ± 1.35	75.67 ± 1.15
Softmax	-	-	-	-	-	71.67 ± 0.28	83.34 ± 0.32
AE _x +Softmax	1	-	-	-	-	70.69 ± 0.82	82.42 ± 0.58
VAE _x +Softmax	1	-	-	-	-	69.20 ± 0.46	81.24 ± 0.63
SplitAE _{xy}	1	-	-	1000	-	70.34 ± 0.56	82.11 ± 0.38
VCCA _{xy}	1	-	-	1000	-	70.72 ± 0.56	82.12 ± 0.61
SplitAE _{xi} +Softmax	1	1000	-	-	-	77.68 ± 0.69	87.09 ± 0.53
VCCA _{xi} +Softmax	1	1000	-	-	-	77.02 ± 0.51	86.46 ± 0.42
VCCA-private _{xi} +Softmax	1	10	-	-	-	73.04 ± 0.56	84.16 ± 0.51
SplitAE _{xiy}	1	1000	-	1000	-	77.43 ± 0.80	87.14 ± 0.57
VCCA _{xiy}	1	1000	-	1000	-	77.22 ± 0.55	86.54 ± 0.51
VCCA-private _{xiy}	1	10	-	1000	-	74.04 ± 0.83	84.59 ± 0.83
SplitAE _{xw} +Softmax	1	-	1000	-	40	76.27 ± 0.66	86.45 ± 0.56
VCCA _{xw} +Softmax	1	-	1000	-	75	75.37 ± 0.46	86.00 ± 0.32
VCCA-private _{xw} +Softmax	1	-	1000	-	75	75.11 ± 0.81	85.91 ± 0.55
SplitAE _{xwy}	1	-	1000	10	75	75.78 ± 0.84	86.13 ± 0.63
VCCA _{xwy}	1	-	1000	10	75	74.72 ± 0.85	85.59 ± 0.78
VCCA-private _{xwy}	1	-	1000	1000	50	74.92 ± 0.74	85.59 ± 0.67
SplitAE _{xiw} +Softmax	1	1000	1000	-	24	77.79 ± 0.48	87.12 ± 0.62
VCCA _{xiw} +Softmax	1	1000	1000	-	32	77.51 ± 0.51	86.69 ± 0.41
SplitAE _{xiwy}	1	1000	1000	1000	24	78.18 ± 0.53	87.26 ± 0.46
VCCA _{xiwy}	1	1000	1000	1000	91	77.78 ± 0.45	86.88 ± 0.47

Supplemental Experimental Procedures

Details on Experimental Setup

In this section, we provide the full details on the experimental setups.

Training DenseNet169 from Scratch We train a DenseNet169 on the dataset from scratch as a baseline. We train the network using stochastic gradient descent for 300 epochs and follow the learning rate schedule of Huang et al. [77], i.e., using an initial learning rate of 0.1 and dividing it by 10 after 150 and 225 epochs. We use a weight decay of 10^{-4} and Nesterov momentum of 0.9 without dampening. We denote this network as DenseNet-scratch in the experiments.

Training the Softmax Classifier We use a Softmax classifier trained on off-the-shelf features as another baseline. We use a DenseNet169 [77] pre-trained on ImageNet 1K as the feature extractor, where we extract 1664-dimensional from the average pooling layer before the classification layer in the architecture. The Softmax classifier is trained for 100 epochs with batch size 64 to minimize the cross-entropy

loss. We use the Adam optimizer [96] with initial learning rate 10^{-4} and hyperparameters $\beta_1 = 0.5$ and $\beta_2 = 0.999$. Note that we used no regularization when training the Softmax classifier. We denote this classifier as Softmax in the experiments. This training setup is also used when training the Softmax classifiers for SplitAE, VCCA, and VCCA-private.

Architectures for Single-View Autoencoders We use a vanilla autoencoder and a VAE as baselines that learn latent representations of the natural image features only. These models are denoted as AE_x and VAE_x respectively in the experiments. Their latent representations have dimension $d_z = 200$ throughout all experiments. The encoder and decoder networks consist of one hidden layer of 512 hidden units with Leaky ReLU activation. The decoder aims to reconstruct the natural image features and we use the sum-of-squares as the reconstruction loss during training. For the VAE, we use the mean outputs $\mu_z(x)$ from the encoder as the latent representations of the image features to train a Softmax classifier, which we denote as $\text{VAE}_x + \text{Softmax}$.

Architectures for SplitAE and VCCA We use the same encoder and decoder architecture for the natural image features as in the single-view autoencoders for all SplitAE, VCCA, and VCCA-private models, i.e., one hidden layer of 512 hidden units with Leaky ReLU activation. We set the latent dimension to $d_z = 200$ in all experiments. The class label decoders use the same architecture as the image feature decoder and predict the class label by optimizing the cross-entropy loss. The iconic image decoders use the generator architecture of the DCGAN [94]. This decoder reconstructs the iconic images $i \in \mathbb{R}^{64 \times 64 \times 3}$ by minimizing the sum of squares loss. The text descriptions are assumed to be a sequence of words $w = (w_1, \dots, w_T)$, where T is the length of the description. We create a vocabulary $V \in \mathbb{R}^{658}$ of the total number of unique words from all text descriptions in the dataset. The text description decoder is an LSTM [95] followed by a linear layer that predicts the next word in the description. More specifically, at time step t , the decoder yields a multinomial probability distribution $p_{\theta_w}(w_t | w_{t-1}, z)$ defined over $w_t \in V$. The LSTM is trained using teacher forcing, meaning that we use the previous ground truth word as input at every time step during the training phase. We project each word into an embedding space of dimensionality $d_{emb} = 200$ by using a lookup table before the word is input to the LSTM. The hidden state h and memory state c of the LSTM is initialized with a linear projection of the latent representation z to provide the LSTM with some context about the natural image. We use the same dimension for h and c as for z , i.e., $d_z = d_h = d_c = 200$. We minimize the cross-entropy loss at every time step by comparing the predicted word with the true word in the description. We apply dropout [97] with a keep rate of 0.5 on the output hidden state h_t before we input it through the linear layer that predicts the word at time step t .

Architectures for VCCA-private In the VCCA-private models, the decoders has the same architectures as the VCCA models. The encoder for the shared latent variable z is also the same as the encoder for the previous described models. The encoder for the private latent variable u_x uses an identical architecture as encoder for

z . We use a convolutional encoder for the iconic image private latent variable, which is a reversed DCGAN generator outputting the mean $\mu_{u_i}(i)$ and variance $\sigma_{u_i}^2(i)$. The text description encoder is an LSTM with the same architectural details as the LSTM decoder. We obtain an embedding for the description by averaging all of the hidden states h_t generated from the LSTM, i.e., $\frac{1}{T} \sum_{t=1}^T h_t$, and input it to a linear layer outputting the mean $\mu_{u_w}(w)$ and the variance $\sigma_{u_w}^2(w)$ for the private latent variable for the text description. We use the same latent dimensions for the shared and private latent variables, i.e., $d_z = d_{u_x} = d_{u_i} = d_{u_w} = 200$.

Training the Single- and Multi-View Autoencoders The single- and multi-view autoencoding models are trained for 200 epochs with batch size 64 and aims to minimize either their reconstruction losses or their corresponding ELBOs. We use the Adam optimizer with initial learning rate 10^{-4} and hyperparameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$ in all experiments. The mean outputs $\mu_z(x)$ from the encoder are used as the latent representations of the image features to train a Softmax classifier. For the class label decoder, we use $K = 1$ posterior samples for predicting the class label during training, while we set $K = 5$ in the validation and test stages.

Grid Search We run a hyperparameter search for the scaling weights λ for the reconstruction losses of the views using a grid search with grid points $\{0.1, 1, 10, 100, 1000\}$. The grid search is performed for all VCCA and VCCA-private models. The weight for the natural image feature loss λ_x is used as a reference by setting it to $\lambda_x = 1$ throughout all experiments. This means that we only vary the scaling weights λ_i , λ_w , and λ_y . We run the grid searches using three different random seeds and average the resulting validation accuracies to select the best hyperparameter setting. Table 4 shows the best hyperparameter settings for the VCCA and VCCA-private models with their validation accuracies. Note that we use the same scaling weights for SplitAE as the ones we found for its corresponding VCCA model. To summarize the grid search results, it is always beneficial to use scaling weights $\lambda > 1$ for the models to enhance the classification accuracy. This will make the models add some semantically meaningful information to the latent space from the additional views, which makes the models more suitable for downstream tasks such as classification.

Posterior Collapse in VCCA-private

We noticed in Table 1 that VCCA-private achieves similar classification performance as standard VCCA when utilizing the text description w and even worse results when utilizing the iconic image i . The private latent variables cannot capture any view-specific variations when each class uses the same iconic image and text description for every natural image. A consequence of this is that the iconic image and text description can be identified by only using the shared latent variable z in the decoding phase. The private latent variables are thus not necessary for determining which iconic image or text description to generate, the generated view will be the same anyway for every natural image of a specific class. The model then finds out that the ELBO can be maximized by letting the approximate posteriors be equal to their prior distributions, which minimize the KL divergences of the private latent variables. This phenomenon is referred to as *posterior collapse* [83].

In Figure 12 and 13, we illustrate that the approximate posterior deduces to its prior distribution – a zero-mean Gaussian with unit variance – during training. Figure 12(a) and 12(b) shows the KL divergences over epochs for VCCA-private_{xw} and VCCA-private_{xwy} respectively. The number of words $T = 24$ is the same for both models and we train the models for 500 epochs to emphasize that $D_{\text{KL}}(q_{\phi_w}(u_w|w) \parallel p(u_w))$ goes towards zero. We believe that this is due to that there are no variations in the text descriptions within one class (see subsection Investigations of the Learned Representations in Results). We perform a similar experiment with VCCA-private_{xi} and VCCA-private_{xiy} and plot their KL divergences in Figure 13. The KL divergence $D_{\text{KL}}(q_{\phi_i}(u_i|w) \parallel p(u_i))$ decreases to zero faster for these models than when we use the text description. We also observe that the KL divergence $D_{\text{KL}}(q_{\phi_x}(u_x|x) \parallel p(u_x))$ decreases to zero as well for both models. Why the KL divergence for the private latent variables decreases faster in the models with the iconic image i is mainly because their likelihood weight λ_i is smaller than λ_w . We noticed that the KL divergences of the private latent variables decreases slower towards zero when λ_i is increased, probably because the model foremost focuses on minimizing the reconstruction loss.

Investigating Latent Representations in VCCA-private_{xi}

Figure 14 shows the shared and private latent spaces of VCCA-private_{xi}, as well as the latent space of the standard VCCA_{xi} for comparison. Note that the shared latent spaces in Figure 14(a) and 14(b) have different structures mainly due to the different settings of their likelihood weight λ_i , which is $\lambda_i = 1000$ for VCCA_{xi} and $\lambda_i = 10$ for VCCA-private_{xi}. We observe in Figure 14(c) that the natural images are structured based on their similarities in background and camera setup. Across the upper right and the lower left parts of the cluster, we find images of grocery items closely packed together in bins. The middle and upper left part includes images with the hand of the photographer and grey backgrounds, e.g., the floor and shelves in the grocery store. Note that this private latent space is rather densely packed mainly because the KL divergence $D_{\text{KL}}(q_{\phi_x}(u_x|x) \parallel p(u_x))$ is steadily decreasing towards zero (see Figure 13(a)). This means that the approximate posterior $q_{\phi_x}(u_x|x)$ is collapsing to its prior distribution, which is why the mean values μ_{u_x} are close to the origin of the space. The mean values μ_{u_i} for the private latent variable u_i are shown in Figure 14(d) using the iconic image. Note that every iconic image i is projected at the same location in the latent space. We observe that similar iconic images are projected close to each other. For instance, an orange, grapefruit and yellow bell pepper have been projected in the upper part, packaged items are in the center parts and the left region we find round objects with a dark red color. To the far right, we see a green and a red apple that has been pushed far away from the similar iconic images on the left. If we look closer into these iconic images, we observe that the apples on the right have a higher stalk on top of the apple than the apples on the left side has, which could be the reason why the model has separated them in the latent space. Another interesting observation is that iconic images with multiple items, e.g., tomatoes, kiwis and satsumas, have been projected into the lower region of the space. We conclude that similar iconic images, based on color and their appearance in the image, are grouped closely in the private latent space.

Derivation of the ELBO for VCCA

Let $x_{1:M}$ denote the all observed data, i.e., $x_{1:M} = x_1, \dots, x_M$, for M different views. We derive the ELBO for VCCA by introducing the approximate posterior $q_\phi(z|x_m)$, where x_m is the only view that we use to infer the latent variable z . We now derive the ELBO from the marginal log-likelihood $\log p_\theta(x_{1:M})$ as

$$\begin{aligned}\log p_\theta(x_{1:M}) &= \log p_\theta(x_{1:M}) \int q_\phi(z|x_m) dz \\ &= \int q_\phi(z|x_m) \log p_\theta(x_{1:M}) dz \\ &= \int q_\phi(z|x_m) \log \frac{p_\theta(x_{1:M}, z)}{p_\theta(z|x_{1:M})} dz \\ &= \int q_\phi(z|x_m) \log \frac{p_\theta(x_{1:M}, z)}{p_\theta(z|x_{1:M})} \frac{q_\phi(z|x_m)}{q_\phi(z|x_m)} dz \\ &= \int q_\phi(z|x_m) \left(\log \frac{q_\phi(z|x_m)}{p_\theta(z|x_{1:M})} + \log \frac{p_\theta(x_{1:M}, z)}{q_\phi(z|x_m)} \right) dz \\ &= D_{\text{KL}}(q_\phi(z|x_m) || p_\theta(z|x_{1:M})) + \mathbb{E}_{q_\phi(z|x_m)} \left[\log \frac{p_\theta(x_{1:M}, z)}{q_\phi(z|x_m)} \right] \\ &\geq \mathbb{E}_{q_\phi(z|x_m)} \left[\log \frac{p_\theta(x_{1:M}, z)}{q_\phi(z|x_m)} \right] = \mathcal{L}(x_{1:M}; \theta, \phi)\end{aligned}$$

We use the factorization property of the joint distribution $p_\theta(x_{1:M})$ to further derive the ELBO $\mathcal{L}(x_{1:M}; \theta, \phi)$:

$$\begin{aligned}\mathcal{L}(\theta, \phi; x_{1:M}) &= \mathbb{E}_{q_\phi(z|x_m)} \left[\log \frac{p_{\theta_1}(x_1|z) \cdots p_{\theta_M}(x_M|z)p(z)}{q_\phi(z|x_m)} \right] \\ &= \mathbb{E}_{q_\phi(z|x_m)} \left[\log p_{\theta_1}(x_1|z) + \dots + \log p_{\theta_M}(x_M|z) + \log \frac{p(z)}{q_\phi(z|x_m)} \right] \\ &= \mathbb{E}_{q_\phi(z|x_m)} [\log p_{\theta_1}(x_1|z)] + \dots + \mathbb{E}_{q_\phi(z|x_m)} [\log p_{\theta_M}(x_M|z)] - D_{\text{KL}}(q_\phi(z|x_m) || p(z))\end{aligned}$$

It may be necessary to balance the terms in the ELBO with some constant for every term, especially when the dimensions and magnitudes differ between the modalities. Thus, we introduce the likelihood weights $\lambda_1, \dots, \lambda_M$, such that the ELBO will be written as

$$\begin{aligned}\mathcal{L}(\theta, \phi; x_{1:M}) &= \lambda_1 \mathbb{E}_{q_\phi(z|x_m)} [\log p_{\theta_1}(x_1|z)] + \dots + \lambda_M \mathbb{E}_{q_\phi(z|x_m)} [\log p_{\theta_M}(x_M|z)] \\ &\quad - D_{\text{KL}}(q_\phi(z|x_m) || p(z))\end{aligned}$$

The likelihood weights $\lambda_1, \dots, \lambda_M$ that are optimal for the task at hand can be found with some hyperparameter search.

Acknowledgments

This research is funded by the Promobilia foundation and the Swedish e-Science Research Centre.

References

- [1] Microsoft Seeing AI app. <https://www.microsoft.com/en-us/seeing-ai/>. Accessed on 2020-04-13.
- [2] Aipoly Vision app. <https://www.aipoly.com/>. Accessed on 2020-04-13.
- [3] Orcam. <https://www.orcam.com/en/>. Accessed on 2020-04-13.
- [4] Transsense. <https://www.transsense.ai/>. Accessed on 2020-04-13.
- [5] S. Caraiman, A. Morar, M. Owczarek, A. Burlacu, D. Rzeszotarski, N. Botezatu, P. Herghelegiu, F. Moldoveanu, P. Strumillo, and A. Moldoveanu. Computer vision for the visually impaired: the sound of vision system. In *IEEE International Conference on Computer Vision Workshops*, 2017.
- [6] Marcus Klasson, Cheng Zhang, and Hedvig Kjellström. A hierarchical grocery store image dataset with visual and semantic labels. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 491–500. IEEE, 2019.
- [7] Weiran Wang, Xinchen Yan, Honglak Lee, and Karen Livescu. Deep variational canonical correlation analysis. *arXiv preprint arXiv:1610.03454*, 2016.
- [8] Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [9] Andrea Frome, Greg Corrado, Jonathon Shlens, Samy Bengio, Jeffrey Dean, Marc’Aurelio Ranzato, and Tomas Mikolov. Devise: A deep visual-semantic embedding model. In *Advances in Neural Information Processing Systems*, 2013.
- [10] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137, 2015.
- [11] Girish Kulkarni, Visruth Premraj, Vicente Ordonez, Sagnik Dhar, Siming Li, Yejin Choi, Alexander C Berg, and Tamara L Berg. Babytalk: Understanding and generating simple image descriptions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2891–2903, 2013.
- [12] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Neural baby talk. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7219–7228, 2018.
- [13] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Hierarchical question-image co-attention for visual question answering. In *Advances In Neural Information Processing Systems*, pages 289–297, 2016.
- [14] Devi Parikh and Kristen Grauman. Relative attributes. In *2011 International Conference on Computer Vision*, pages 503–510. IEEE, 2011.
- [15] Nitish Srivastava and Ruslan Salakhutdinov. Multimodal learning with deep boltzmann machines. *Journal of Machine Learning Research*, 15:2949–2980, 2014.
- [16] Püren Güler, Yasemin Bekiroglu, Xavi Gratal, Karl Pauwels, and Danica Kragic. What’s in the container? classifying object contents from vision and touch. In *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3961–3968. IEEE, 2014.

- [17] Jiquan Ngiam, Aditya Khosla, Mingyu Kim, Juhan Nam, Honglak Lee, and Andrew Y Ng. Multimodal deep learning. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 689–696, 2011.
- [18] Weiran Wang, Raman Arora, Karen Livescu, and Jeff Bilmes. On deep multi-view representation learning. In *International conference on machine learning*, pages 1083–1092, 2015.
- [19] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [20] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338, Jun 2010.
- [21] Timnit Gebru, Jonathan Krause, Yilun Wang, Duyun Chen, Jia Deng, and Li Fei-Fei. Fine-grained car detection for visual census estimation. In *AAAI Conference on Artificial Intelligence*, 2017.
- [22] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. Technical Report 7694, California Institute of Technology, 2007.
- [23] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. 2016.
- [24] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- [25] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, 2014.
- [26] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *Indian Conference on Computer Vision, Graphics and Image Processing*, Dec 2008.
- [27] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [28] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.
- [29] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. SUN database: Large-scale scene recognition from abbey to zoo. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2010.
- [30] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *TACL*, 2:67–78, 2014.
- [31] Weidong Geng, Feilin Han, Jiangke Lin, Liuyi Zhu, Jieming Bai, Suzhen Wang, Lin He, Qiang Xiao, and Zhangjiong Lai. Fine-grained grocery product recognition by one-shot learning. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 1706–1714, 2018.

- [32] Marian George and Christian Floerkemeier. Recognizing products: A per-exemplar multi-label image classification approach. In *European Conference on Computer Vision*, pages 440–455. Springer, 2014.
- [33] Edward Hsiao, Alvaro Collet, and Martial Hebert. Making specific features less discriminative to improve point-based 3d object recognition. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 2653–2660. IEEE, 2010.
- [34] Philipp Jund, Nichola Abdo, Andreas Eitel, and Wolfram Burgard. The freiburg groceries dataset. *arXiv preprint arXiv:1611.05799*, 2016.
- [35] Kevin Lai, Liefeng Bo, Xiaofeng Ren, and Dieter Fox. A large-scale hierarchical multi-view rgb-d object dataset. In *2011 IEEE international conference on robotics and automation*, pages 1817–1824. IEEE, 2011.
- [36] Michele Merler, Carolina Galleguillos, and Serge Belongie. Recognizing groceries in situ using in vitro training data. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2007.
- [37] Arjun Singh, James Sha, Karthik S Narayan, Tudor Achim, and Pieter Abbeel. Bigbird: A large-scale 3d database of object instances. In *2014 IEEE international conference on robotics and automation (ICRA)*, pages 509–516. IEEE, 2014.
- [38] Georg Waltner, Michael Schwarz, Stefan Ladstätter, Anna Weber, Patrick Luley, Horst Bischof, Meinrad Lindschinger, Irene Schmid, and Lucas Paletta. Mango - mobile augmented reality with functional eating guidance and food awareness. In *International Workshop on Multimedia Assisted Dietary Management*, 2015.
- [39] Xiu-Shen Wei, Quan Cui, Lei Yang, Peng Wang, and Lingqiao Liu. Rpc: A large-scale retail product checkout dataset. *arXiv preprint arXiv:1901.07249*, 2019.
- [40] Tess Winlock, Eric Christiansen, and Serge Belongie. Toward real-time grocery detection for the visually impaired. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition-Workshops*, pages 49–56. IEEE, 2010.
- [41] Hao Yu, Fu Yanwei, and Jiang Yu-Gang. Take goods from shelves: A dataset for class-incremental object detection. In *ACM International Conference on Multimedia Retrieval (ICMR '19)*, 2019.
- [42] Horea Muresan and Mihai Oltean. Fruit recognition from images using deep learning. Technical report, Babes-Bolyai University, 2017.
- [43] Škrjanec Marko. Automatic fruit recognition using computer vision. (Mentor: Matej Kristan), Fakulteta za računalništvo in informatiko, Univerza v Ljubljani, 2013.
- [44] Suchet Bargoti and James Patrick Underwood. Deep fruit detection in orchards. In *IEEE International Conference on Robotics and Automation*, 2017.
- [45] Inkyu Sa, Zongyuan Ge, Feras Dayoub, Ben Upcroft, Tristan Perez, and Chris McCool. Deepfruits: A fruit detection system using deep neural networks. *Sensors*, 16(8):1222, 2016.
- [46] Weiqing Min, Shuqiang Jiang, Linhu Liu, Yong Rui, and Ramesh Jain. A survey on food computing. *ACM Computing Surveys (CSUR)*, 52(5):1–36, 2019.
- [47] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In *European Conference on Computer Vision*, 2014.

- [48] Y. Kawano and K. Yanai. Automatic expansion of a food image dataset leveraging existing categories with domain adaptation. In *Proc. of ECCV Workshop on Transferring and Adapting Source Knowledge in Computer Vision (TASK-CV)*, 2014.
- [49] Weiqing Min, Linhu Liu, Zhengdong Luo, and Shuqiang Jiang. Ingredient-guided cascaded multi-attention network for food recognition. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 1331–1339, 2019.
- [50] Jaclyn Rich, Hamed Haddadi, and Timothy M Hospedales. Towards bottom-up analysis of social food. In *Proceedings of the 6th International Conference on Digital Health Conference*, pages 111–120, 2016.
- [51] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. Scaling egocentric vision: The epic-kitchens dataset. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 720–736, 2018.
- [52] Javier Marin, Aritro Biswas, Ferda Ofli, Nicholas Hynes, Amaia Salvador, Yusuf Aytar, Ingmar Weber, and Antonio Torralba. Recipe1m+: A dataset for learning cross-modal embeddings for cooking recipes and food images. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2019.
- [53] Amaia Salvador, Nicholas Hynes, Yusuf Aytar, Javier Marin, Ferda Ofli, Ingmar Weber, and Antonio Torralba. Learning cross-modal embeddings for cooking recipes and food images. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [54] Semih Yagcioglu, Aykut Erdem, Erkut Erdem, and Nazli Ikizler-Cinbis. Recipeqa: A challenge dataset for multimodal comprehension of cooking recipes. *arXiv preprint arXiv:1809.00812*, 2018.
- [55] Oscar Bejbom, Neel Joshi, Dan Morris, Scott Saponas, and Siddharth Khullar. Menu-match: Restaurant-specific food logging from images. In *2015 IEEE Winter Conference on Applications of Computer Vision*, pages 844–851. IEEE, 2015.
- [56] Ruihan Xu, Luis Herranz, Shuqiang Jiang, Shuang Wang, Xinhang Song, and Ramesh Jain. Geolocalized modeling for dish recognition. *IEEE transactions on multimedia*, 17(8):1187–1199, 2015.
- [57] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(2):423–443, 2018.
- [58] Chris Cremer and Nate Kushman. On the importance of learning aggregate posteriors in multimodal variational autoencoders. *1st Symposium on Advances in Approximate Bayesian Inference*, 2018.
- [59] Yanwei Fu and Leonid Sigal. Semi-supervised vocabulary-informed learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5337–5346, 2016.
- [60] Alessandro Pieropan, Giampiero Salvi, Karl Pauwels, and Hedvig Kjellström. Audio-visual classification and detection of human manipulation actions. In *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3045–3052. IEEE, 2014.

- [61] Mathieu Salzmann, Carl Henrik Ek, Raquel Urtasun, and Trevor Darrell. Factorized orthogonal latent spaces. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 701–708, 2010.
- [62] Yuge Shi, N Siddharth, Brooks Paige, and Philip Torr. Variational mixture-of-experts autoencoders for multi-modal deep generative models. In *Advances in Neural Information Processing Systems*, pages 15692–15703, 2019.
- [63] Masahiro Suzuki, Kotaro Nakayama, and Yutaka Matsuo. Joint multimodal learning with deep generative models. *arXiv preprint arXiv:1611.01891*, 2016.
- [64] Yao-Hung Hubert Tsai, Paul Pu Liang, Amir Zadeh, Louis-Philippe Morency, and Ruslan Salakhutdinov. Learning factorized multimodal representations. In *International Conference on Learning Representations*, 2019.
- [65] Ramakrishna Vedantam, Ian Fischer, Jonathan Huang, and Kevin Murphy. Generative models of visually grounded imagination. *arXiv preprint arXiv:1705.10762*, 2017.
- [66] Mike Wu and Noah Goodman. Multimodal generative models for scalable weakly-supervised learning. In *Advances in Neural Information Processing Systems*, pages 5575–5585, 2018.
- [67] Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. *IEEE transactions on pattern analysis and machine intelligence*, 41(9):2251–2265, 2018.
- [68] Chang Xu, Dacheng Tao, and Chao Xu. A survey on multi-view learning. *arXiv preprint arXiv:1304.5634*, 2013.
- [69] Cheng Zhang, Hedvig Kjellström, and Carl Henrik Ek. Inter-battery topic representation learning. In *European Conference on Computer Vision*, pages 210–226. Springer, 2016.
- [70] Jing Zhao, Xijiong Xie, Xin Xu, and Shiliang Sun. Multi-view learning overview: Recent progress and new challenges. *Information Fusion*, 38:43–54, 2017.
- [71] Harold Hotelling. Relations between two sets of variates. *Biometrika*, 28(3/4):321–377, 1936.
- [72] Shotaro Akaho. A kernel method for canonical correlation analysis. *arXiv preprint cs/0609071*, 2006.
- [73] Galen Andrew, Raman Arora, Jeff Bilmes, and Karen Livescu. Deep canonical correlation analysis. In *International conference on machine learning*, pages 1247–1255, 2013.
- [74] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [75] Cheng Zhang, Judith Butepage, Hedvig Kjellström, and Stephan Mandt. Advances in variational inference. *IEEE transactions on pattern analysis and machine intelligence*, 2018.
- [76] Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. *Neural networks*, 13(4-5):411–430, 2000.
- [77] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.

- [78] Solomon Kullback and Richard A Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
- [79] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [80] Shiyong Cui and Mihai Datcu. Comparison of kullback-leibler divergence approximation methods between gaussian mixture models for satellite image retrieval. In *2015 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, pages 3719–3722. IEEE, 2015.
- [81] Jacob Goldberger, Shiri Gordon, and Hayit Greenspan. An efficient image similarity measure based on approximations of kl-divergence between two gaussian mixtures. In *IEEE International Conference on Computer Vision*, 2003.
- [82] John R Hershey and Peder A Olsen. Approximating the kullback leibler divergence between gaussian mixture models. In *2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP’07*, volume 4, pages IV–317. IEEE, 2007.
- [83] Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*, 2015.
- [84] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol, and Léon Bottou. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of machine learning research*, 11(12), 2010.
- [85] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [86] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.
- [87] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.
- [88] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10):1345–1359, Oct 2010.
- [89] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- [90] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, 2014.
- [91] Yingzhen Li, John Bradshaw, and Yash Sharma. Are generative classifiers more robust to adversarial attacks? In *International Conference on Machine Learning*, pages 3804–3814, 2019.
- [92] Chao Ma, Sebastian Tschiatschek, Konstantina Palla, José Miguel Hernández-Lobato, Sebastian Nowozin, and Cheng Zhang. Eddi: Efficient dynamic discovery of high-value information with partial vae. *arXiv preprint arXiv:1809.11142*, 2018.

- [93] Ali Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. CNN features off-the-shelf: An astounding baseline for recognition. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2014.
- [94] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [95] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [96] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [97] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.

Paper C

Learn the Time to Learn: Replay Scheduling in Continual Learning

Marcus Klasson*, Hedvig Kjellström*, Cheng Zhang[†]

*KTH Royal Institute of Technology, Stockholm, Sweden

[†]Microsoft Research, Cambridge, United Kingdom

Abstract

Replay-based continual learning has shown to be successful in mitigating catastrophic forgetting, where most works focus on improving the sample quality in the commonly small replay memory. However, in many real-world applications, replay memories would be limited by constraints on processing time rather than storage capacity as most organizations store all historical data in the cloud. Inspired by human learning, we demonstrate that scheduling over the time to replay is critical to the final performance with finite memory resources. To this end, we propose to learn the time to learn for a continual learning system, in which we learn schedules over which tasks to replay at different times. We use Monte Carlo tree search to illustrate this idea and show that our method can be combined with any replay-based method and memory selection technique. We perform extensive evaluation showing that learning replay schedules can significantly improve the performance compared to baselines without learned scheduling. Our results indicate that the learned schedules are also consistent with human learning insights.

1 Introduction

Many organizations deploying machine learning systems receive large volumes of data daily where these new data are often associated with new tasks. Although all historical data are stored in the cloud in practice, retraining machine learning systems on a daily basis is prohibitive both in time and cost. In this setting, the systems must continuously adapt to new tasks without forgetting the previously learned abilities. Continual learning (CL) methods [1, 2] address this challenge where, in particular, replay-based methods [3, 4] have shown to be very effective in achieving great prediction performance and retaining knowledge of old tasks. Replay-based methods mitigate catastrophic forgetting by revisiting a small set of samples, which

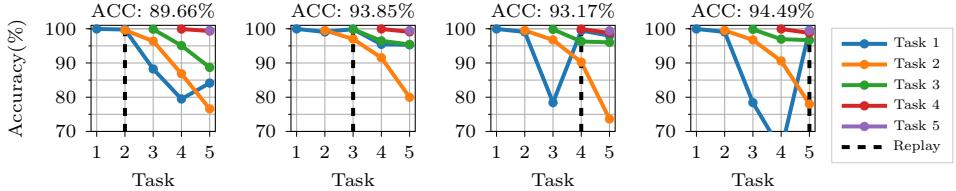


Figure 1: Task accuracies on Split MNIST [16] when replaying only 10 samples of classes 0/1 at a single time step. The black vertical line indicates when replay is used. ACC denotes the average accuracy over all tasks after learning Task 5. Results are averaged over 5 seeds. These results show that the time to replay the previous task is critical for the final performance.

is feasible to process compared to the size of the historical data. In the traditional CL literature, the replay memory is limited due to the assumption that historical data are not available. In the real-world setting where historical data are in fact always available, the requirement of small memory remains due to processing time and cost issues.

Most research on replay-based CL has been focused on the sample quality in the memory [3, 5–10] or data compression to increase the memory capacity [4, 11, 12]. Common for these methods is that the memory allocates an equal amount of space for storing samples from old tasks. When learning new tasks, the whole memory is replayed to mitigate catastrophic forgetting. However, in life-long learning settings, this simple strategy would be inefficient as the memory must store a large number of tasks. Furthermore, these methods ignore the time to learn old tasks again which is important in human learning. Humans are CL systems, and different methods have been developed to enhance memory retention, such as spaced repetition [13–15] which is often used in education. These education methods focus on the scheduling of learning and rehearsal of previous learned knowledge.

We argue that finding the proper schedule of which tasks to replay in the fixed memory setting is critical for CL. To demonstrate our claim, we perform a simple experiment on the Split MNIST [16] dataset where each task consists of learning the digits 0/1, 2/3, etc. arriving in sequence. The replay memory contains data from task 1 and can only be replayed at one point in time. Figure 1 shows how the task performances progress over time when the memory is replayed at different time steps. In this example, the best final performance is achieved when the memory is used when learning task 5. Note that choosing different time points to replay the same memory leads to noticeably different results in the final performance. These results indicate that scheduling the time when to apply replay can influence the final performance significantly of a CL system.

To this end, we propose learning the time to learn, in which we learn replay schedules of which tasks to replay at different times inspired from human learning [13]. We illustrate the advantages with replay scheduling by using Monte Carlo tree search (MCTS) [17] to learn policies for replay. More specifically, we train a neural network on the current task dataset mixed with the scheduled replay samples and measure the final performance of the network to evaluate the replay schedules selected by MCTS. In summary, our contributions are:

- We propose a new CL setting where historical data is available while the processing time is limited, in order to adjust current CL research closer to real-

world needs (Section 3.1). In this new setting, we introduce replay scheduling in CL where we learn the time to learn which tasks to replay (Section 3.2).

- We use MCTS as an example method to illustrate how to learn replay scheduling policies by creating finite sets of memory compositions that can be replayed at every task (Section 3.3).
- We demonstrate with six benchmark datasets that learned scheduling can improve the CL performance significantly in the fixed size memory setting (Section 4.1 and 4.5). Moreover, we show that our method can be combined with any memory selection technique and replay-based method (Section 4.3 and 4.4), as well as being efficient in situations where the memory size is even smaller than the number of classes (Section 4.6).

2 Related Work

In this section, we give a brief overview of CL methods, essentially replay-based methods, as well as spaced repetition techniques for human CL.

Continual Learning. Traditional CL can be divided into three main areas, namely regularization-based, architecture-based, and replay-based approaches. Regularization-based methods aim to mitigate catastrophic forgetting by protecting parameters influencing the predictive performance on known tasks from wide changes and use the rest of the parameters for learning new tasks [8, 18–23, 218]. Architecture-based methods isolate task-specific parameters by either increasing network capacity [24–26] or freezing parts of the network [27, 28] to maintain good performance on previous tasks. Replay-based methods mix samples from old tasks with the current dataset to mitigate catastrophic forgetting, where the replay samples are either stored in an external memory [3, 4, 29, 30] or generated using a generative model [31, 32]. Regularization-based approaches and dynamic architectures have been combined with replay-based approaches to methods to overcome their limitations [8, 12, 19, 33–42]. Our work relates most to replay-based methods with external memory which we spend more time on describing in the next paragraph.

Replay-based Continual Learning. A commonly used memory selection strategy of replay samples is random selection. Much research effort has focused on selecting higher quality samples to store in memory [3, 5–10, 29, 30, 43]. Chaudhry et al. [3] reviews several selection strategies in scenarios with tiny memory capacity, e.g., reservoir sampling [44], first-in first-out buffer [30], k-Means, and Mean-of-Features [9]. However, more elaborate selection strategies have been shown to give little benefit over random selection for image classification problems [4, 19]. More recently, there has been work on compressing raw images to feature representations to increase the number of memory examples for replay [4, 11, 12]. Our approach differs from the above mentioned works since we focus on learning to select which tasks to replay at the current task rather than improving memory selection or compression quality of the memory samples. Replay scheduling can however be combined with any selection strategy and feature compression method.

Human Continual Learning. Humans are CL systems in the sense of learning tasks and concepts sequentially. Furthermore, humans have the ability to memo-

rize experiences but forgets learned knowledge gradually rather than catastrophically [45]. Different learning techniques have been suggested for humans to memorize better [46, 47]. An example is spaced repetition where time intervals between rehearsal are gradually increased to improve long-term memory retention [13], where the earliest documented works are from Ebbinghaus [14]. Further studies have shown that memory training schedules with adjusted spaced repetition are better at preserving memory than using uniformly spaced rehearsal times [15, 48]. Several works in CL with neural networks are inspired by human learning techniques, including spaced repetition [49–51], sleep mechanisms [23, 27, 52], and memory reactivation [4, 53]. Replay scheduling is also inspired by spaced repetition, where we learn schedules of which tasks to replay at different times.

3 Method

In this section, we describe our new problem setting in CL where historical data are available while the processing time is limited when learning new tasks. To this end, we present our method for learning replay schedules in CL where the idea is to learn schedules of which tasks the network should replay at different times. We use Monte Carlo tree search (MCTS) [17] to learn a scheduling policy by encouraging searches for promising replay schedules based on the classification accuracy.

3.1 Problem Setting

We focus on a slightly new setting, considering the needs for CL in the real-world where all historical data can be available since data storage is cheap. However, as this data volume is typically huge, we are constrained to use a small amount of historical data for replay when learning new tasks due to limitations on the amount of compute rather than data storage capability. Thus, the goal is to learn how to select subsets of historical data to efficiently mitigate catastrophic forgetting. We refer to these subsets of historical data as the *replay memory* throughout the paper, where the size of the replay memory affects the processing time when learning a new task. Moreover, we focus on composing the replay memory based on the seen tasks in the historical data rather than single stored instances.

Here, we introduce the notation of our problem setting which resembles the traditional CL setting for image classification. We let a neural network f_{θ} , parameterized by θ , learn T tasks sequentially from the datasets $\mathcal{D}_1, \dots, \mathcal{D}_T$ arriving one at a time. The t -th dataset $\mathcal{D}_t = \{(\mathbf{x}_t^{(i)}, y_t^{(i)})\}_{i=1}^{N_t}$ consists of N_t samples where $\mathbf{x}_t^{(i)}$ and $y_t^{(i)}$ are the i -th data point and class label respectively. The training objective at task t is given by

$$\min_{\theta} \sum_{i=1}^{N_t} \ell(f_{\theta}(\mathbf{x}_t^{(i)}), y_t^{(i)}), \quad (1)$$

where $\ell(\cdot)$ is the loss function, e.g., cross-entropy loss in our case. When learning task t , the network f_{θ} is at risk of catastrophically forgetting the previous $t - 1$ tasks. Replay-based methods mitigate forgetting historical tasks by storing old examples in an external replay memory, that is mixed with the current task dataset during training. Next, we describe our method for constructing this replay memory.

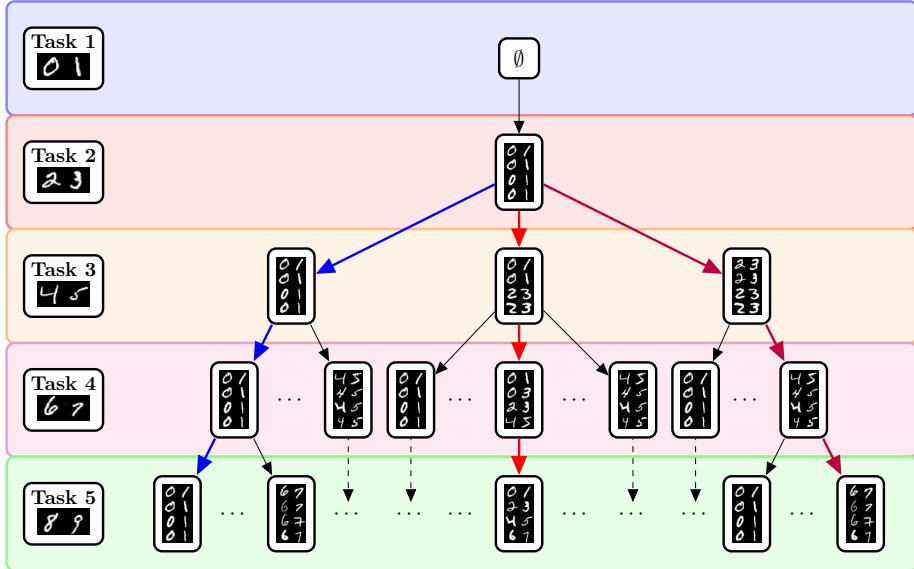


Figure 2: An exemplar tree of replay memory compositions from the proposed discretization method described in Section 3.2 for Split MNIST. The replay memories from one replay schedule are found by traversing from task 1-5 through the tree on the right hand side. The replay memory compositions have been structured according to the task where they can be used for replay. Note that the replay memory at task 1 is the empty set, i.e., $\mathcal{M} = \emptyset$. Example images for each task are shown on the left.

We assume that historical data from old tasks are accessible at any time step t . However, since the historical data is prohibitively large, we can only fill a small replay memory \mathcal{M} with M historical samples for replay due to processing time constraints. The challenge is how to fill the replay memory with M samples that efficiently retain the knowledge of old tasks when learning new tasks. We focus on selecting the samples on task-level by deciding on the task proportion (a_1, \dots, a_{t-1}) of samples to fetch from each task, where $a_i \geq 0$ is the proportion of M examples from task i to place in \mathcal{M} and $\sum_{i=1}^{t-1} a_i = 1$. Consequently, we need a method for choosing these task proportions of which old tasks to replay. To simplify this selection, we construct a discrete set of choices for possible task proportions that can be used for constructing the replay memory \mathcal{M} .

3.2 Replay Scheduling in Continual Learning

In this section, we describe our replay scheduling method for selecting the replay memory at different time steps. We define a replay schedule as a sequence $S = (\mathbf{a}_1, \dots, \mathbf{a}_{T-1})$, where $\mathbf{a}_i = (a_1, \dots, a_{t-1})$ for $1 \leq i \leq T-1$ is the sequence of task proportions used for determining how many samples per task to fill the replay memory with at task i . To make the selection of task proportions tractable, we construct an action space with a discrete number of choices for task proportions from old tasks. We use the following method to construct this action space: At the time to learn task t , we have $t-1$ historical tasks that we can choose from. We create $t-1$ bins $\mathbf{b}_t = [b_1, \dots, b_{t-1}]$ and choose a task index to sample for each bin

$b_i \in \{1, \dots, t - 1\}$. We treat the bins as interchangeable and only keep the unique choices. For example, at task 3, we have seen task 1 and 2, so the unique choices of vectors are $[1, 1]$, $[1, 2]$, $[2, 2]$, where $[1, 1]$ indicates that all memory samples are from task 1, $[1, 2]$ indicates that half memory is from task 1 and the other half are from task etc. The task proportions are then computed by counting the number of occurrences of each task index in \mathbf{b}_t and dividing by $t - 1$, such that $\mathbf{a}_t = \text{bincount}(\mathbf{b}_t)/(t - 1)$. If the memory size M is not divisible by the task proportion value for task i , we round the number of replay samples from task i , i.e., $a_i \cdot M$, up or down accordingly while keeping the memory size fixed. From this specification, we can build a tree of different replay schedules to evaluate with the network.

Figure 2 shows an example of a replay schedule tree with Split MNIST [16] where the memory size is $M = 8$. The figure shows the current tasks with example images of the task classes on the left, and the right side shows examples of possible replay memories that can be evaluated. The memory starts as the empty set, i.e. $\mathcal{M}_1 = \emptyset$, at task 1. Before learning task 2, \mathcal{M}_2 is filled with M task 1 examples since this is the only task seen so far. At task 3, the memory compositions we can choose from are M examples from either task 1 or 2, as well as equally filling \mathcal{M}_3 with four examples each from both tasks. A replay schedule is represented as a path traversal of different replay memory compositions from task 1 to task 5. We have color-coded three examples of possible schedules in Figure 2 to use for illustration: the blue path represents a replay schedule where only task 1 examples are replayed. The red path represents using an equally distributed amount of memory samples per task in the memory, and the purple path represents a schedule where the memory is filled M examples from the previously visited task. Note that all other possible paths in the tree are also valid replay schedules.

3.3 Monte Carlo Tree Search for Replay Schedules

Our proposed setup to learn replay schedules is to discretize the number of possible task proportions per task. With our discretization method, the action space of task proportions becomes shaped as a tree where each node represents a task proportion that corresponds to a specific composition of the replay memory \mathcal{M} , as can be seen in Figure 2. To find the most efficient replay schedule S in mitigating catastrophic forgetting, we can perform an exhaustive search (e.g., breadth-first search) where all possible replay schedules are evaluated. However, since the tree grows fast with the number of tasks, we need a scalable method that enables tree searches in large action spaces. To this end, we propose to use Monte Carlo tree search (MCTS) [17] since MCTS has been successful in applications with similar conditions [54–57]. In our setting, MCTS can help us to concentrate the search in directions towards promising replay schedules in terms of classification performance and thus learn the time to learn.

Before outlining the steps for performing MCTS, we describe the notation for the tree search. The tree has T levels where level t corresponds to task t . Each tree level t contains a set of nodes $V_t = \{v_t^i\}_{i=1}^{K_t}$ where K_t is the number of nodes at level t . Every node v_t^i is related to a sequence of task proportions \mathbf{a}_t^i , which can be used to retrieve the replay memory from the historical data. Referring back to Figure 2, the replay memory composition on the root node v_1^1 is the empty set

$\mathcal{M} = \emptyset$, since the historical data is empty. At task 2, the only task proportion is $\mathbf{a}_2^1 = (a_1, a_2, \dots, a_{T-1}) = (1.0, 0.0, \dots, 0.0)$, so the replay memory at node v_2^1 is only filled with samples from task 1. In the rest of the paper, we denote a node at task t as v_t by ignoring the superscript i to avoid clutter in the notation.

At every visited node v_t , we store the corresponding task proportions \mathbf{a}_t in the replay schedule S . The final replay schedule is then used for constructing the replay memory at each task during the CL training. Next, we briefly outline the MCTS steps for performing the search of replay schedules:

Selection. An MCTS iteration begins by selecting the next node to visit from the root node. If the current node v_t has visited all its children during the previous iterations, the next node is selected by evaluating the Upper Confidence Tree (UCT) [58] function for all children. We use the following UCT function proposed by [55] to evaluate the score for moving from node v_t to its child v_{t+1} :

$$UCT(v_t, v_{t+1}) = \max(Q(v_{t+1})) + C \sqrt{\frac{2 \log(N(v_t))}{N(v_{t+1})}}. \quad (2)$$

The reward function $Q(\cdot)$ contains the rewards for previous searches that has passed through the child v_{t+1} . The exploration constant $C \geq 0$ determines the degree of exploration to consider based on the number of visits $N(v_t)$ and $N(v_{t+1})$ to the corresponding nodes in the tree. The child node to visit next is the one with the highest UCT score.

Expansion. Whenever the current node v_t has unvisited child nodes, the search tree is expanded with one of the unvisited child nodes v_{t+1} selected with uniform sampling.

Simulation and Reward. After the expansion step, the search proceeds by selecting the succeeding nodes uniformly at random until reaching a terminal node v_T . The replay schedule S is then collected by appending the task proportions from the visited nodes during the iteration. We train the network with the replay memories retrieved from S and evaluate the average accuracy over all tasks after learning task T as the reward r , i.e., $r = \frac{1}{T} \sum_{i=1}^T A_{T,i}^{(val)}$, where $A_{T,i}^{(val)}$ is the validation accuracy of task i at task T .

Backpropagation. Reward r is backpropagated from the expanded node v_t to the root v_1 , where the reward function $Q(\cdot)$ and number of visits $N(\cdot)$ are updated at each node.

We provide pseudocode in Algorithm 2 in Appendix 2. Finally, we emphasize that we use MCTS as an example method for illustrating the benefits of learning replay schedules in our new proposed CL setting.

4 Experiments

We evaluated our replay scheduling method empirically on six common benchmark datasets for CL. We denote our method as Replay Scheduling MCTS (RS-MCTS) and select the result on the held-out test sets from the replay schedule that yielded the best reward on the validation set during the search. Full details on experimental settings are in Appendix ?? and additional results are in Appendix 4. Our code is available in the supplementary material.

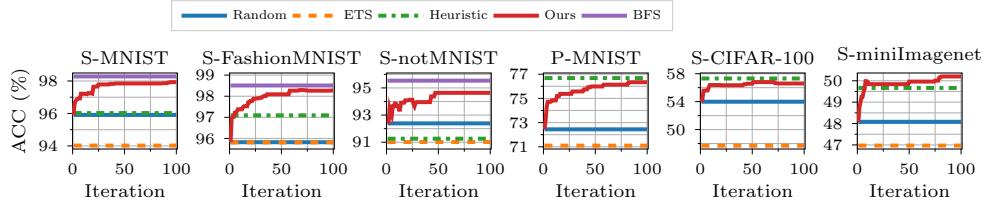


Figure 3: Average test accuracies over tasks after learning the final task (ACC) over the MCTS simulations for all datasets, where ‘S’ and ‘P’ are used as short for ‘Split’ and ‘Permuted’. We compare performance for RS-MCTS (Ours) against random replay schedules (Random), Equal Task Schedule (ETS), and Heuristic Scheduling (Heuristic) baselines. For the first three datasets, we show the best ACC found from a breadth-first search (BFS) as an upper bound. All results have been averaged over 5 seeds. These results show that replay scheduling can improve over ETS and outperform or perform on par with Heuristic across different datasets and network architectures.

Datasets. We conduct experiments on six datasets commonly used as benchmarks in the CL literature: Split MNIST [16, 59], Fashion-MNIST [60], Split notMNIST [61], Permuted MNIST [62], Split CIFAR-100 [63], and Split miniImagenet [64]. We randomly sample 15% of the training data from each task to use for validation when computing the reward for the MCTS simulations.

Baselines. We compare RS-MCTS to using 1) random replay schedules (Random), 2) equal task schedules (ETS), and 3) a heuristic scheduling method (Heuristic). The ETS baseline uses equal task proportions, such that $M/(t - 1)$ samples per task are replayed during learning of task t , and use both training and validation sets for training such that ETS use the same amount of data as RS-MCTS. The Heuristic baseline replays the tasks which accuracy on the validation set is below a certain threshold proportional to the best achieved validation accuracy on the task. Here, the replay memory is filled with M/k samples per task where k is the number of selected tasks. If $k = 0$, then we skip applying replay at the current task. See Appendix ?? for more details on the Heuristic baseline.

Architectures. We use a multi-head output layer for all datasets except for Permuted MNIST where the network uses single-head output. We use a 2-layer MLP with 256 hidden units for Split MNIST, Split FashionMNIST, Split notMNIST, and Permuted MNIST. For Split CIFAR-100, we use the CNN architecture used in Vinyals et al. [64] and Schwartz et al. [23]. For Split miniImagenet, we apply the reduced ResNet-18 from [30].

Evaluation Metric. We use the average test accuracy over all tasks after learning the final task, i.e., $\text{ACC} = \frac{1}{T} \sum_{i=1}^T A_{T,i}$ where $A_{T,i}$ is the test accuracy of task i after learning task T . We report means and standard deviations of ACC using 5 different seeds on all datasets.

4.1 Performance Progress of RS-MCTS

In the first experiments, we show that the replay schedules from RS-MCTS yield better performance than replaying an equal amount of samples per task. The replay memory size is fixed to $M = 10$ for Split MNIST, FashionMNIST, and notMNIST, and $M = 100$ for Permuted MNIST, Split CIFAR-100, and Split miniImagenet.

Table 1: Performance comparison with ACC between RS-MCTS (Ours), Random scheduling (Random), Equal Task Schedule (ETS), and Heuristic Scheduling (Heuristic) with various memory selection methods evaluated across all datasets. Replay memory sizes are $M = 10$ and $M = 100$ for the 5-task and 10/20-task datasets respectively. We report the mean and standard deviation averaged over 5 seeds. Ours performs better or on par with the baselines on most datasets and selection methods, where MoF yields the best results in general.

Selection	Method	5-task Datasets			10- and 20-task Datasets		
		S-MNIST	S-FashionMNIST	S-notMNIST	P-MNIST	S-CIFAR-100	S-miniImagenet
Uniform	Random	95.91 (± 1.56)	95.82 (± 1.45)	92.39 (± 1.29)	72.44 (± 1.15)	53.99 (± 0.51)	48.08 (± 1.36)
	ETS	94.02 (± 4.25)	95.81 (± 3.53)	91.01 (± 1.39)	71.09 (± 2.31)	47.70 (± 2.16)	46.97 (± 1.24)
	Heuristic	96.02 (± 2.32)	97.09 (± 0.62)	91.26 (± 3.99)	76.68 (± 2.13)	57.31 (± 1.21)	49.66 (± 1.10)
	Ours	97.93 (± 0.56)	98.27 (± 0.17)	94.64 (± 0.39)	76.34 (± 0.98)	56.60 (± 1.13)	50.20 (± 0.72)
<i>k</i> -means	Random	94.24 (± 3.20)	96.30 (± 1.62)	91.64 (± 1.39)	74.30 (± 1.43)	53.18 (± 1.66)	49.47 (± 2.70)
	ETS	92.89 (± 3.53)	96.47 (± 0.85)	93.80 (± 0.82)	69.40 (± 1.32)	47.51 (± 1.14)	45.82 (± 0.92)
	Heuristic	96.28 (± 1.68)	95.78 (± 1.50)	91.75 (± 0.94)	75.57 (± 1.18)	54.31 (± 3.94)	49.25 (± 1.00)
	Ours	98.20 (± 0.16)	98.48 (± 0.26)	93.61 (± 0.71)	77.74 (± 0.80)	56.95 (± 0.92)	50.47 (± 0.85)
<i>k</i> -center	Random	96.40 (± 0.68)	95.57 (± 3.16)	92.61 (± 1.70)	71.41 (± 2.75)	48.46 (± 0.31)	44.76 (± 0.96)
	ETS	94.84 (± 1.40)	97.28 (± 0.50)	91.08 (± 2.48)	69.11 (± 1.69)	44.13 (± 1.06)	41.35 (± 1.23)
	Heuristic	94.55 (± 2.79)	94.08 (± 3.72)	92.06 (± 1.20)	74.33 (± 2.00)	50.32 (± 1.97)	44.13 (± 0.95)
	Ours	98.24 (± 0.36)	98.06 (± 0.35)	94.26 (± 0.37)	76.55 (± 1.16)	51.37 (± 1.63)	46.76 (± 0.96)
MoF	Random	95.18 (± 3.18)	95.76 (± 1.41)	91.33 (± 1.75)	77.96 (± 1.84)	61.93 (± 1.05)	54.50 (± 1.33)
	ETS	97.04 (± 1.23)	96.48 (± 1.33)	92.64 (± 0.87)	77.62 (± 1.12)	60.43 (± 1.17)	56.12 (± 1.12)
	Heuristic	96.46 (± 2.41)	95.84 (± 0.89)	93.24 (± 0.77)	77.27 (± 1.45)	55.60 (± 2.70)	52.30 (± 0.59)
	Ours	98.37 (± 0.24)	97.84 (± 0.32)	94.62 (± 0.42)	81.58 (± 0.75)	64.22 (± 0.65)	57.70 (± 0.51)

Uniform sampling is used as the memory selection method for all methods in this experiment. For the 5-task datasets, we provide the optimal replay schedule found from a breadth-first search (BFS) over all 1050 possible replay schedules in our action space (which corresponds to a tree with depth of 4) as an upper bound for RS-MCTS. As the search space grows fast with the number of tasks, BFS becomes computationally infeasible when we have 10 or more tasks.

Figure 3 shows the progress of ACC over iterations by RS-MCTS for all datasets. We also show the best ACC metrics for Random, ETS, Heuristic, and BFS (where appropriate) as straight lines. We observe that RS-MCTS outperforms Random and ETS successively with more iterations. Furthermore, RS-MCTS approaches the upper limit of BFS on the 5-task datasets. For Permuted MNIST and Split CIFAR-100, the Heuristic baseline and RS-MCTS perform on par after 50 iterations. This shows that Heuristic with careful tuning of the validation accuracy threshold can be a strong baseline when comparing replay scheduling methods. The top row of Table 1 shows the ACC for each method for this experiment. We note that RS-MCTS outperforms ETS significantly on most datasets and performs on par with Heuristic.

4.2 Replay Schedule Visualizations

We visualize a learned replay schedule from Split CIFAR-100 with memory size $M = 100$ to gain insights into the behavior of the scheduling policy from RS-MCTS. Figure 4 shows a bubble plot of the task proportions that are used for filling the replay memory at every task. Each circle color corresponds to a historical task and the circle size represents its proportion of replay samples at the current task. The sum of all points from all circles at each column is fixed at the different time steps since the memory size M is fixed. The task proportions vary dynamically over time in a sophisticated nonlinear way which would be hard to replace by a heuristic method. Moreover, we can observe space repetition style scheduling on many tasks, e.g., task 1-3 are replayed with similar proportion at the initial tasks but eventually starts

Table 2: Performance comparison with ACC between scheduling methods RS-MCTS (Ours), Random, ETS, and Heuristic combined with replay-based methods HAL, MER, DER, and DER++. Replay memory sizes are $M = 10$ and $M = 100$ for the 5-task and 10/20-task datasets respectively. We report the mean and standard deviation averaged over 5 seeds. Results on Heuristic where some seed did not converge is denoted by *. Applying RS-MCTS to each method can enhance the performance compared to using the baseline schedules.

Method	Schedule	5-task Datasets			10- and 20-task Datasets		
		S-MNIST	S-FashionMNIST	S-notMNIST	P-MNIST	S-CIFAR-100	S-miniImagenet
HAL	Random	97.24 (± 0.70)	86.74 (± 6.05)	93.61 (± 1.31)	88.49 (± 0.99)	36.09 (± 1.77)	38.51 (± 2.22)
	ETS	94.02 (± 4.25)	95.81 (± 3.53)	91.01 (± 1.39)	88.46 (± 0.86)	34.90 (± 2.02)	38.13 (± 1.18)
	Heuristic	97.69 (± 0.19)	*74.16 (± 11.19)	93.64 (± 0.93)	*66.63 (± 28.50)	35.07 (± 1.29)	39.51 (± 1.49)
	Ours	97.93 (± 0.56)	98.27 (± 0.17)	94.64 (± 0.39)	89.14 (± 0.74)	40.22 (± 1.57)	41.39 (± 1.15)
MER	Random	93.07 (± 0.81)	85.53 (± 3.30)	91.13 (± 0.86)	75.90 (± 1.34)	42.96 (± 1.70)	31.48 (± 1.65)
	ETS	92.89 (± 3.53)	96.47 (± 0.85)	93.80 (± 0.82)	73.01 (± 0.96)	43.38 (± 1.81)	33.58 (± 1.53)
	Heuristic	94.30 (± 2.79)	96.91 (± 0.62)	90.90 (± 1.30)	83.86 (± 3.19)	40.90 (± 1.70)	34.22 (± 1.93)
	Ours	98.20 (± 0.16)	98.48 (± 0.26)	93.61 (± 0.71)	79.72 (± 0.71)	44.29 (± 0.69)	32.74 (± 1.29)
DER	Random	98.23 (± 0.53)	96.56 (± 1.79)	92.89 (± 0.86)	87.51 (± 1.10)	56.83 (± 0.76)	42.19 (± 0.67)
	ETS	98.17 (± 0.35)	97.69 (± 0.58)	94.74 (± 1.05)	85.71 (± 0.75)	52.58 (± 1.49)	35.50 (± 2.84)
	Heuristic	94.57 (± 1.71)	*72.49 (± 19.32)	*77.88 (± 12.58)	81.56 (± 2.28)	55.75 (± 1.08)	43.62 (± 0.88)
	Ours	99.02 (± 0.10)	98.33 (± 0.51)	95.02 (± 0.33)	90.11 (± 0.18)	58.99 (± 0.98)	43.46 (± 0.95)
DER++	Random	97.90 (± 0.52)	97.10 (± 1.03)	93.29 (± 1.43)	87.89 (± 1.10)	58.49 (± 1.44)	48.40 (± 0.69)
	ETS	97.98 (± 0.52)	98.12 (± 0.40)	94.53 (± 1.02)	85.25 (± 0.88)	52.54 (± 1.06)	41.36 (± 2.90)
	Heuristic	92.35 (± 2.42)	*67.31 (± 21.20)	93.88 (± 1.33)	79.17 (± 2.44)	56.70 (± 1.27)	45.73 (± 0.84)
	Ours	98.84 (± 0.21)	98.38 (± 0.43)	94.73 (± 0.20)	89.84 (± 0.22)	59.23 (± 0.83)	49.45 (± 0.68)

varying the time interval between replay. Also, task 4 and 6 need less replay in their early stages, which could potentially be that they are simpler or correlated with other tasks. We provide a similar visualization for Split MNIST in Figure 7 in Appendix ?? to bring more insights to the benefits of replay scheduling.

4.3 Alternative Memory Selection Methods

We show that our method can be combined with any memory selection method for storing replay samples. In addition to uniform sampling, we apply various memory selection methods commonly used in the CL literature, namely k -means clustering, k -center clustering [8], and Mean-of-Features (MoF) [9]. We compare our method and ETS combined with these different selection methods. The replay memory sizes are $M = 10$ for the 5-task datasets and $M = 100$ for the 10- and 20-task datasets. Table 1 shows the results across all datasets. We note that using the replay schedule from RS-MCTS outperforms the baselines when using the alternative selection methods, where MoF performs the best on most datasets.

4.4 Applying Scheduling to Recent Replay Methods

In this experiment, we show that replay scheduling can be combined with any replay method to enhance the CL performance. We combine RS-MCTS with Hindsight Anchor Learning (HAL) [35], Meta-Experience Replay (MER) [65], Dark Experience Replay (DER) [33], and DER++. We provide the hyperparameter settings in Appendix 4.3. Table 2 shows the performance comparison between our proposed replay scheduling against using Random, ETS, and Heuristic schedules for each method. The results confirm that replay scheduling is important for the final performance given the same memory constraints and it can benefit any existing CL framework.

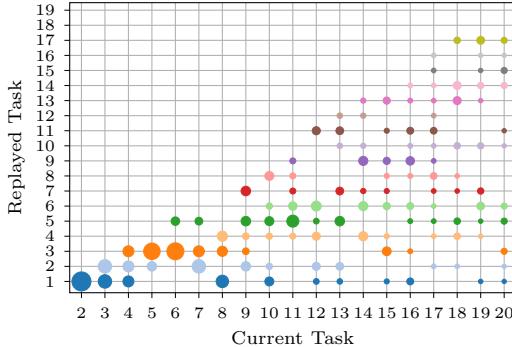


Figure 4: Replay schedule learned from Split CIFAR-100 visualized as a bubble plot. The task proportions vary dynamically over time which would be hard to replace by a heuristic method.

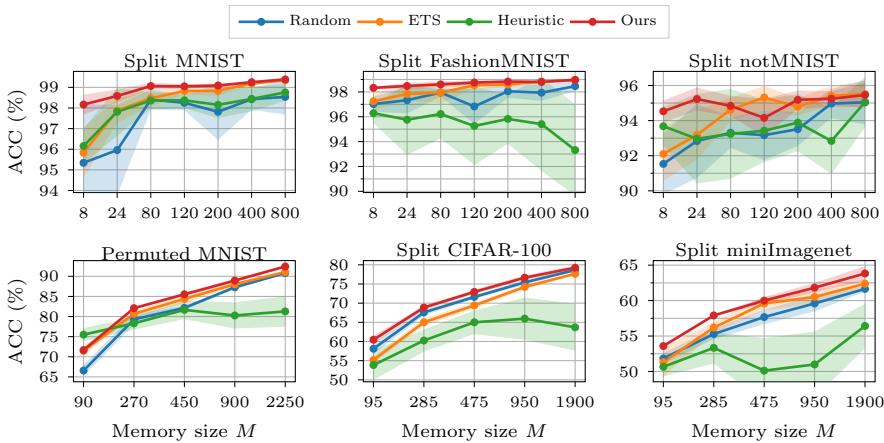


Figure 5: Average test accuracies over tasks after learning the final task (ACC) over different replay memory sizes M for the RS-MCTS (Ours) and the Random, ETS, and Heuristic baselines on all datasets. All results have been averaged over 5 seeds. The results show that replay scheduling can outperform replaying with random and ETS schedules, especially for small M , on both small and large datasets across different backbone choices. Furthermore, our method requires less careful tuning than the Heuristic baseline as M increases.

4.5 Results with Varying Replay Memory Size

We show that our method can improve the performance across different choices of memory size. In this experiment, we set the replay memory size to ensure the ETS baseline replays an equal number of samples per class at the final task. The memory size is set to $M = n_{cpt} \cdot n_{spc} \cdot (T - 1)$, where n_{cpt} is the number of classes per task in the dataset and n_{spc} are the number of samples per class we wish to replay at task T for the ETS baseline. In Figure 5, we observe that RS-MCTS obtains better task accuracies than ETS, especially for small memory sizes. Both RS-MCTS and ETS perform better than Heuristic as M increases showing that Heuristic requires careful tuning of the validation accuracy threshold.

4.6 Efficiency of Replay Scheduling

Table 3: Performance comparison with ACC averaged over 5 seeds in the 1 ex/class memory setting evaluated across all datasets. RS-MCTS (Ours) has replay memory size $M = 2$ and $M = 50$ for the 5-task and 10/20-task datasets respectively. The baselines replay all available memory samples. RS-MCTS performs on par with the best baselines on all datasets except S-CIFAR-100.

Method	5-task Datasets			10- and 20-task Datasets		
	S-MNIST	S-FashionMNIST	S-notMNIST	P-MNIST	S-CIFAR-100	S-miniImagenet
Random	92.56 \pm 2.90	92.70 \pm 3.78	89.53 \pm 3.96	70.02 \pm 1.76	48.62 \pm 1.02	48.85 \pm 1.38
A-GEM	94.97 \pm 1.50	94.81 \pm 0.86	92.27 \pm 1.16	64.71 \pm 1.78	42.22 \pm 2.13	32.06 \pm 1.83
ER-Ring	94.94 \pm 1.56	95.83 \pm 2.15	91.10 \pm 1.89	69.73 \pm 1.13	53.93 \pm 1.13	49.82 \pm 1.69
Uniform	95.77 \pm 1.12	97.12 \pm 1.57	92.14 \pm 1.45	69.85 \pm 1.01	52.63 \pm 1.62	50.56 \pm 1.07
RS-MCTS (Ours)	96.07 \pm 1.60	97.17 \pm 0.78	93.41 \pm 1.11	72.52 \pm 0.54	51.50 \pm 1.19	50.70 \pm 0.54

We illustrate the efficiency of replay scheduling with comparisons to several common replay-based CL baselines in an even more extreme memory setting. Our goal is to investigate if scheduling over which tasks to replay can be more efficient in situations where the memory size is even smaller than the number of classes. To this end, we set the replay memory size for our method to $M = 2$ for the 5-task datasets, such that only 2 samples can be selected for replay at all times. For the 10- and 20-task datasets which have 100 classes, we set $M = 50$. We then compare against the most memory efficient CL baselines, namely A-GEM [34], ER-Ring [3] which show promising results with 1 sample per class for replay, and with uniform memory selection as reference. Additionally, we compare to using random replay schedules (Random) with the same memory setting as for RS-MCTS. We visualize the memory usage for our method and the baselines when training on a 5-task dataset in Figure 6. Figure 8 in Appendix ?? shows the memory usage for the other datasets.

Table 3 shows the ACC for each method across all datasets. Despite using significantly fewer samples for replay, RS-MCTS performs better or on par with the best baselines on all datasets except Split CIFAR-100. These results show that replay scheduling can be even more efficient than the state-of-the-art memory efficient replay-based methods, which indicate that learning the time to learn is an important research direction in CL.

5 Conclusions

We propose learning the time to learn, i.e., in a real-world CL context, learning schedules of which tasks to replay at different times. To the best of our knowledge, we are the first to consider the time to learn in machine learning inspired by human learning techniques. We illustrate with MCTS how replay schedules can be learned and show that they can produce significantly improved results when comparing against methods without learned scheduling under the same memory budgets. Furthermore, the dynamic behavior of the learned schedules showed similarities to human learning techniques, such as spaced repetition, by replaying previous tasks with varying time intervals. Finally, we showed that replay scheduling can be combined with any

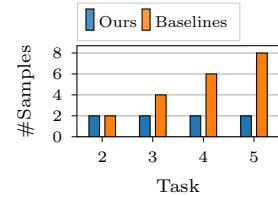


Figure 6: Number of replayed samples per task for the 5-task datasets in the tiny memory setting. Ours use $M = 2$ samples for replay, while the baselines increment their memory per task.

replay-based method as well as utilize the memory efficiently even when the memory size is smaller than the number of classes.

In future work, we would like to explore choosing memory samples on an instance level as the current work selects samples on task level. This would need a policy learning method that scales to large action spaces which is a research challenge by itself. Also, using MCTS for learning replay scheduling policies can be inefficient, especially in CL settings, in terms of training time and since it needs to learn for each dataset and application separately. In future work, we will incorporate reinforcement learning methods that generalize and extend to learning general policies that can be directly applied to any new application and domain.

Appendix

This supplementary material is structured as follows:

- Appendix ??: Full details of the experimental settings.
- Appendix ??: Pseudocode of RS-MCTS in Algorithm 2 to provide more details about our method.
- Appendix ?? describes the Heuristic scheduling baseline.
- Appendix 4 shows additional experimental results.
- Appendix 5 includes additional figures.
- Our code will be made publicly available upon acceptance.

1 Experimental Settings

In this section, we describe the full details of the experimental settings used in this paper.

Datasets. We conduct experiments on six datasets commonly used in the continual learning literature. Split MNIST [16] is a variant of the MNIST [59] dataset where the classes have been divided into 5 tasks incoming in the order 0/1, 2/3, 4/5, 6/7, and 8/9. Split Fashion-MNIST [60] is of similar size to MNIST and consists of grayscale images of different clothes, where the classes have been divided into the 5 tasks T-shirt/Trouser, Pullover/Dress, Coat/Sandals, Shirt/Sneaker, and Bag/Ankle boots. Similar to MNIST, Split notMNIST [61] consists of 10 classes of the letters A-J with various fonts, where the classes are divided into the 5 tasks A/B, C/D, E/F, G/H, and I/J. We use training/test split provided by [37] for Split notMNIST. Permuted MNIST [62] dataset consists of applying a unique random permutation of the pixels of the images in original MNIST to create each task, except for the first task that is to learn the original MNIST dataset. We reduce the original MNIST dataset to 10k samples and create 9 unique random permutations to get a 10-task version of Permuted MNIST. In Split CIFAR-100 [63], the 100 classes are divided into 20 tasks with 5 classes for each task [9, 30]. Similarly, Split miniImagenet [64] consists of 100 classes randomly chosen from the original Imagenet dataset where the 100 classes are divided into 20 tasks with 5 classes per task.

Network Architectures. We use a 2-layer MLP with 256 hidden units and ReLU activation for Split MNIST, Split FashionMNIST, Split notMNIST, and Permuted MNIST. We use a multi-head output layer for each dataset except Permuted MNIST where the network uses single-head output layer. For Split CIFAR-100, we use a multi-head CNN architecture built according to the CNN in [18, 23, 64], which consists of four 3x3 convolutional blocks, i.e. convolutional layer followed by batch normalization [66], with 64 filters, ReLU activations, and 2x2 Max-pooling. For Split miniImagenet, we use the reduced ResNet-18 from [30] with multi-head output layer.

Hyperparameters. We train all networks with the Adam optimizer [67] with learning rate $\eta = 0.001$ and hyperparameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$. Note that the learning rate for Adam is not reset before training on a new task. Next, we give details on number of training epochs and batch sizes specific for each dataset:

- Split MNIST: 10 epochs/task, batch size 128.
- Split FashionMNIST: 30 epochs/task, batch size 128.
- Split notMNIST: 50 epochs/task, batch size 128.
- Permuted MNIST: 20 epochs/task, batch size 128.
- Split CIFAR-100: 25 epochs/task, batch size 256.
- Split miniImagenet: 1 epoch/task (task 1 trained for 5 epochs as warm up), batch size 32.

Monte Carlo Tree Search. We run RS-MCTS for 100 iterations in all experiments. The replay schedules used in the reported results on the held-out test sets are from the replay schedule that gave the highest reward on the validation sets. The exploration constant for UCT in Equation 2 is set to $C = 0.1$ in all experiments [55].

Computational Cost. All experiments were performed on one NVIDIA GeForce RTX 2080Ti. The wall clock time for ETS on Split MNIST was around 1.5 minutes, and RS-MCTS and BFS takes 40 seconds on average to run one iteration, where BFS runs 1050 iterations in total for Split MNIST.

Implementations. We adapted the implementation released by Borsos et al. [6] for the memory selection strategies Uniform sampling, k -means clustering, k -center clustering [8], and Mean-of-Features [9]. For HAL [35], MER [65], DER [33], and DER++, we follow the implementations released by [33] for each method to apply them to our replay scheduling methods. Furthermore, we follow the implementations released by C:chaudhry2019tiny and C:mirzadeh2021cl-gym for A-GEM [34] and ER-Ring [3]. For MCTS, we adapted the implementation from <https://github.com/int8/monte-carlo-tree-search> to search for replay schedules.

Experimental Settings for Single Task Replay Memory Experiment. We motivated the need for replay scheduling in continual learning with Figure 1 in Section 1. This simple experiment was performed on Split MNIST where the replay memory only contains samples from the first task, i.e., learning the classes 0/1. Furthermore, the memory can only be replayed at one point in time and we show the performance

on each task when the memory is replayed at different time steps. We set the memory size to $M = 10$ samples such that the memory holds 5 samples from both classes. We use the same network architecture and hyperparameters as described above for Split MNIST. The ACC metric above each subfigure corresponds to the ACC for training a network with the single task memory replay at different tasks. We observe that choosing different time points to replay the same memory leads to noticeably different results in the final performance, and in this example, the best final performance is achieved when the memory is used when learning task 5. Therefore, we argue that finding the proper schedule of what tasks to replay at what time in the fixed memory situation can be critical for continual learning.

2 Replay Scheduling Monte Carlo Tree Search Algorithm

We provide pseudo-code in Algorithm 2 outlining the steps for our method Replay Scheduling Monte Carlo tree search (RS-MCTS) described in the main paper (Section 3.3). The MCTS procedure selects actions over which task proportions to fill the replay memory with at every task, where the selected task proportions are stored in the replay schedule S . The schedule is then passed to EVALUATEREPLAYSCHEDULE(\cdot) where the continual learning part executes the training with replay memories filled according to the schedule. The reward for the schedule S is the average validation accuracy over all tasks after learning task T , i.e., ACC, which is backpropagated through the tree to update the statistics of the selected nodes. The schedule S_{best} yielding the best ACC score is returned to be used for evaluation on the held-out test sets.

The function GETREPLAYMEMORY(\cdot) is the policy for retrieving the replay memory \mathcal{M} from the historical data given the task proportion a . The number of samples per task determined by the task proportions are rounded up or down accordingly to fill \mathcal{M} with M replay samples in total. The function GETTASKPROPORTION(\cdot) simply returns the task proportion that is related to given node.

3 Heuristic Scheduling Baseline

We implemented a heuristic scheduling baseline to compare against RS-MCTS. The baseline keeps a validation set for the old tasks and replays the tasks which validation accuracy is below a certain threshold. We set the threshold in the following way: Let $A_{t,i}$ be the validation accuracy for task t evaluated at time step i . The best evaluated validation accuracy for task t at time i is given by $A_{t,i}^{(best)} = \max(\{A_{t,1}, \dots, A_{t,i}\})$. The condition for replaying task t on the next time step is then $A_{t,i} < \tau A_{t,i}^{(best)}$, where $\tau \in [0, 1]$ is a ratio controlling how much the current accuracy on task t is allowed to decrease w.r.t. the best accuracy. The replay memory is filled with M/k , where k is the number of tasks that need to be replayed according to their decrease in validation accuracy. This heuristic scheduling corresponds to the intuition of re-learning when a task has been forgotten. Training on the current task is performed without replay if the accuracy on all old tasks is above their corresponding threshold.

Algorithm 2 Replay Scheduling Monte Carlo tree search

Require: Tree nodes $v_{1:T}$, Datasets $\mathcal{D}_{1:T}$, Learning rate η

Require: Replay memory size M

```

1:  $ACC_{best} \leftarrow 0, S_{best} \leftarrow ()$ 
2: while within computational budget do
3:    $S \leftarrow ()$ 
4:    $v_t, S \leftarrow \text{TREEPOLICY}(v_1, S)$ 
5:    $v_T, S \leftarrow \text{DEFAULTPOLICY}(v_t, S)$ 
6:    $ACC \leftarrow \text{EVALUATEREPLAYSCHEDULE}(\mathcal{D}_{1:T}, S, M)$ 
7:    $\text{BACKPROPAGATE}(v_t, ACC)$ 
8:   if  $ACC > ACC_{best}$  then
9:      $ACC_{best} \leftarrow ACC$ 
10:     $S_{best} \leftarrow S$ 
11:   end if
12: end while
13: return  $ACC_{best}, S_{best}$ 

14: function TREEPOLICY( $v_t, S$ )
15:   while  $v_t$  is non-terminal do
16:     if  $v_t$  not fully expanded then
17:       return EXPANSION( $v_t, S$ )
18:     else
19:        $v_t \leftarrow \text{BESTCHILD}(v_t)$ 
20:        $S.append(\mathbf{a}_t)$ , where  $\mathbf{a}_t \leftarrow \text{GETTASKPROPORTION}(v_t)$ 
21:     end if
22:   end while
23:   return  $v_t, S$ 
24: end function

25: function EXPANSION( $v_t, S$ )
26:   Sample  $v_{t+1}$  uniformly among unvisited children of  $v_t$ 
27:    $S.append(\mathbf{a}_{t+1})$ , where  $\mathbf{a}_{t+1} \leftarrow \text{GETTASKPROPORTION}(v_{t+1})$ 
28:   Add new child  $v_{t+1}$  to node  $v_t$ 
29:   return  $v_{t+1}, S$ 
30: end function

31: function BESTCHILD( $v_t$ )
32:    $v_{t+1} = \arg \max_{v_{t+1} \in \text{children of } v_t} \max(Q(v_{t+1})) + C \sqrt{\frac{2 \log(N(v_t))}{N(v_{t+1})}}$ 
33:   return  $v_{t+1}$ 
34: end function

35: function DEFAULTPOLICY( $v_t, S$ )
36:   while  $v_t$  is non-terminal do
37:     Sample  $v_{t+1}$  uniformly among children of  $v_t$ 
38:      $S.append(\mathbf{a}_{t+1})$ , where  $\mathbf{a}_{t+1} \leftarrow \text{GETTASKPROPORTION}(v_{t+1})$ 
39:     Update  $v_t \leftarrow v_{t+1}$ 
40:   end while
41:   return  $v_t, S$ 
42: end function

43: function EVALUATEREPLAYSCHEDULE( $\mathcal{D}_{1:T}, S, M$ )
44:   Initialize neural network  $f_\theta$ 
45:   for  $t = 1, \dots, T$  do
46:      $\mathbf{a} \leftarrow S[t - 1]$ 
47:      $\mathcal{M} \leftarrow \text{GETREPLAYMEMORY}(M, \mathbf{a})$ 
48:     for  $\mathcal{B} \sim \mathcal{D}_t^{(train)}$  do
49:        $\theta \leftarrow SGD(\mathcal{B} \cup \mathcal{M}, \theta, \eta)$ 
50:     end for
51:   end for
52:    $A_{1:T}^{(val)} \leftarrow \text{EVALUATEACCURACY}(f_\theta, \mathcal{D}_{1:T}^{(val)})$ 
53:    $ACC \leftarrow \frac{1}{T} \sum_{i=1}^T A_{T,i}^{(val)}$ 
54:   return  $ACC$ 
55: end function

56: function BACKPROPAGATE( $v_t, R$ )
57:   while  $v_t$  is not root do
58:      $N(v_t) \leftarrow N(v_t) + 1$ 
59:      $Q(v_t) \leftarrow R$ 
60:      $v_t \leftarrow \text{parent of } v_t$ 
61:   end while
62: end function

```

Grid search for τ . We performed a coarse-to-fine grid search for the ratio τ on each dataset. The best value for τ is selected according to the highest mean accuracy

on the validation set averaged over 5 seeds. The validation set consists of 15% of the training data and is the same for RS-MCTS. We use the same experimental settings as described in Appendix ???. The memory sizes are set to $M = 10$ and $M = 100$ for the 5-task datasets and the 10/20-task datasets respectively, and we apply uniform sampling as the memory selection method. We provide the ranges for τ that was used on each dataset and put the best value in **bold**:

- Split MNIST: $\tau = [0.9, 0.93, 0.95, \mathbf{0.96}, 0.97, 0.98, 0.99]$
- Split FashionMNIST: $\tau = [0.9, 0.93, 0.95, 0.96, \mathbf{0.97}, 0.98, 0.99]$
- Split notMNIST: $\tau = [0.9, 0.93, 0.95, 0.96, 0.97, \mathbf{0.98}, 0.99]$
- Permuted MNIST: $\tau = [0.5, 0.55, 0.6, 0.65, 0.7, \mathbf{0.75}, 0.8, 0.9, 0.95, 0.97, 0.99]$
- Split CIFAR-100: $\tau = [0.3, 0.4, 0.45, \mathbf{0.5}, 0.55, 0.6, 0.65, 0.7, 0.8, 0.9, 0.95, 0.97, 0.99]$
- Split miniImagenet: $\tau = [0.5, 0.6, 0.65, 0.7, \mathbf{0.75}, 0.8, 0.85, 0.9, 0.95, 0.97, 0.99]$

Note that we use these values for τ on all experiments with Heuristic for the corresponding datasets. The performance for this heuristic highly depends on careful tuning for the ratio τ when the memory size or memory selection method changes, as can be seen in in Figure 5 and Table 1.

4 Additional Experimental Results

In this section, we bring more insights to the benefits of replay scheduling in Section 4.1 as well as provide metrics for catastrophic forgetting in Section ??.

4.1 Replay Schedule Visualization for Split MNIST

In Figure 7, we show the progress in test classification performance for each task when using ETS and RS-MCTS with memory size $M = 10$ on Split MNIST. For comparison, we also show the performance from a network that is fine-tuning on the current task without using replay. Both ETS and RS-MCTS overcome catastrophic forgetting to a large degree compared to the fine-tuning network. Our method RS-MCTS further improves the performance compared to ETS with the same memory, which indicates that learning the time to learn can be more efficient against catastrophic forgetting. Especially, Task 1 and 2 seems to be the most difficult task to remember since it has the lowest final performance using the fine-tuning network. Both ETS and RS-MCTS manage to retain their performance on Task 1 using replay, however, RS-MCTS remembers Task 2 better than ETS by around 5%.

To bring more insights to this behavior, we have visualized the task proportions of the replay examples using a bubble plot showing the corresponding replay schedule from RS-MCTS in Figure 7(right). At Task 3 and 4, we see that the schedule fills the memory with data from Task 2 and discards replaying Task 1. This helps the network to retain knowledge about Task 2 better than ETS at the cost of forgetting Task 3 slightly when learning Task 4. This shows that the learned policy has considered

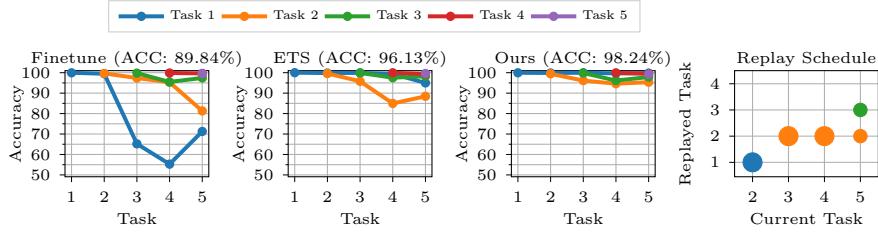


Figure 7: Comparison of test classification accuracies for Task 1-5 on Split MNIST from a network trained without replay (Finetune), ETS, and RS-MCTS (Ours). The ACC metric for each method is shown on top of each figure. We also visualize the replay schedule found by RS-MCTS as a bubble plot to the right. The memory size is set to $M = 10$ with uniform memory selection for ETS and RS-MCTS. Results are shown for 1 seed.

the difficulty level of different tasks. At the next task, the RS-MCTS schedule has decided to rehearse Task 3 and reduces replaying Task 2 when learning Task 5. This behavior is similar to spaced repetition, where increasing the time interval between rehearsals helps memory retention. We emphasize that even on datasets with few tasks, using learned replay schedules can overcome catastrophic forgetting better than standard ETS approaches.

4.2 Analysis of Catastrophic Forgetting

We have compared the degree of catastrophic forgetting for our method against the baselines by measuring the backward transfer (BWT) metric from lopez2017gradient, which is given by

$$\text{BWT} = \frac{1}{T-1} \sum_{i=1}^{T-1} A_{T,i} - A_{i,i}, \quad (3)$$

where $A_{t,i}$ is the test accuracy for task t after learning task i . Table 4 shows the ACC and BWT metrics for the experiments in Section 4.3. In general, the BWT metric is consistently better when the corresponding ACC is better. We find an exception in Table 4 on Split CIFAR-100 and Split miniImagenet between Ours and Heuristic with uniform selection method, where Heuristic has better BWT while its mean of ACC is slightly lower than ACC for Ours. Table 5 shows the ACC and BWT metrics for the experiments in Section 4.6, where we see a similar pattern that better ACC yields better BWT. The BWT of RS-MCTS is on par with the other baselines except on Split CIFAR-100 where the ACC on our method was a bit lower than the best baselines.

4.3 Applying Scheduling to Recent Replay Methods

In Section 4.4, we showed that RS-MCTS can be applied to any replay method. We combined RS-MCTS together with four recent replay methods, namely Hind-sight Anchor Learning (HAL) [35], Meta Experience Replay (MER) [65], and Dark Experience Replay (DER) [33]. Table 6 shows the ACC and BWT for all methods combined with the scheduling from Random, ETS, Heuristic, and RS-MCTS. We observe that RS-MCTS can further improve the performance for each of the replay methods across the different datasets.

Hyperparameters. We present the hyperparameters used for each method below. We used the same architectures and hyperparameters as described in Appendix ?? for all datasets if not mentioned otherwise. We used the Adam optimizer with learning rate $\eta = 0.001$ for MER, DER, and DER++. For HAL, we used the SGD optimizer since using Adam made the model diverge in our experiments.

- Hindsight Anchor Learning (HAL):
 - Split MNIST: learning rate $\eta = 0.1$, regularization $\lambda = 0.1$, mean embedding strength $\gamma = 0.5$, decay rate $\beta = 0.7$, gradient steps on anchors $k = 100$
 - Split FashionMNIST: learning rate $\eta = 0.1$, regularization $\lambda = 0.1$, mean embedding strength $\gamma = 0.1$, decay rate $\beta = 0.5$, gradient steps on anchors $k = 100$
 - Split notMNIST: learning rate $\eta = 0.1$, regularization $\lambda = 0.1$, mean embedding strength $\gamma = 0.1$, decay rate $\beta = 0.5$, gradient steps on anchors $k = 100$
 - Permuted MNIST: learning rate $\eta = 0.1$, regularization $\lambda = 0.1$, mean embedding strength $\gamma = 0.1$, decay rate $\beta = 0.5$, gradient steps on anchors $k = 100$
 - Split CIFAR-100: learning rate $\eta = 0.03$, regularization $\lambda = 1.0$, mean embedding strength $\gamma = 0.1$, decay rate $\beta = 0.5$, gradient steps on anchors $k = 100$
 - Split miniImagenet: learning rate $\eta = 0.03$, regularization $\lambda = 0.3$, mean embedding strength $\gamma = 0.1$, decay rate $\beta = 0.5$, gradient steps on anchors $k = 100$
- Meta Experience Replay (MER):
 - Split MNIST: across batch meta-learning rate $\gamma = 1.0$, within batch meta-learning rate $\beta = 1.0$
 - Split FashionMNIST: across batch meta-learning rate $\gamma = 1.0$, within batch meta-learning rate $\beta = 0.01$
 - Split notMNIST: across batch meta-learning rate $\gamma = 1.0$, within batch meta-learning rate $\beta = 1.0$
 - Permuted MNIST: across batch meta-learning rate $\gamma = 1.0$, within batch meta-learning rate $\beta = 1.0$
 - Split CIFAR-100: across batch meta-learning rate $\gamma = 1.0$, within batch meta-learning rate $\beta = 0.1$
 - Split miniImagenet: across batch meta-learning rate $\gamma = 1.0$, within batch meta-learning rate $\beta = 0.1$
- Dark Experience Replay (DER):
 - Split MNIST: loss coefficient memory logits $\alpha = 0.2$
 - Split FashionMNIST: loss coefficient memory logits $\alpha = 0.2$

- Split notMNIST: loss coefficient memory logits $\alpha = 0.1$
- Permuted MNIST: loss coefficient memory logits $\alpha = 1.0$
- Split CIFAR-100: loss coefficient memory logits $\alpha = 1.0$
- Split miniImagenet: loss coefficient memory logits $\alpha = 0.1$
- Dark Experience Replay++ (DER++):
 - Split MNIST: loss coefficient memory logits $\alpha = 0.2$, loss coefficient memory labels $\beta = 1.0$
 - Split FashionMNIST: loss coefficient memory logits $\alpha = 0.2$, loss coefficient memory labels $\beta = 1.0$
 - Split notMNIST: loss coefficient memory logits $\alpha = 0.1$, loss coefficient memory labels $\beta = 1.0$
 - Permuted MNIST: loss coefficient memory logits $\alpha = 1.0$, loss coefficient memory labels $\beta = 1.0$
 - Split CIFAR-100: loss coefficient memory logits $\alpha = 1.0$, loss coefficient memory labels $\beta = 1.0$
 - Split miniImagenet: loss coefficient memory logits $\alpha = 0.1$, loss coefficient memory labels $\beta = 1.0$

5 Additional Figures

Memory usage. We visualize the memory usage for Permuted MNIST and the 20-task datasets Split CIFAR-100 and Split miniImagenet in Figure 8 for our method and the baselines used in the experiment in Section 4.6. Our method uses a fixed memory size of $M = 50$ samples for replay on all three datasets. The memory size capacity for our method is reached after learning task 6 and task 11 on the Permuted MNIST and the 20-task datasets respectively, while the baselines continue incrementing their replay memory size.

Table 4: Performance comparison with ACC and BWT metrics for all datasets between RS-MCTS and the baselines with various memory selection methods. The memory size is set to $M = 10$ and $M = 100$ for the 5-task and 10/20-task datasets respectively. We report the mean and standard deviation of ACC and BWT, where all results have been averaged over 5 seeds. RS-MCTS performs better or on par than the baselines on most datasets and selection methods, where MoF yields the best performance in general.

Selection	Method	Split MNIST		Split FashionMNIST		Split notMNIST	
		ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑
Uniform	Random	95.91 (± 1.56)	-4.79 (± 1.95)	95.82 (± 1.45)	-4.35 (± 1.79)	92.39 (± 1.29)	-4.56 (± 1.29)
	ETS	94.02 (± 4.25)	-7.22 (± 5.33)	95.81 (± 3.53)	-4.45 (± 4.34)	91.01 (± 1.39)	-6.16 (± 1.82)
	Heuristic	96.02 (± 2.32)	-4.64 (± 2.90)	97.09 (± 0.62)	-2.82 (± 0.84)	91.26 (± 3.99)	-6.06 (± 4.70)
	Ours	97.93 (± 0.56)	-2.27 (± 0.71)	98.27 (± 0.17)	-1.29 (± 0.20)	94.64 (± 0.39)	-1.47 (± 0.79)
k -means	Random	94.24 (± 3.20)	-6.88 (± 4.00)	96.30 (± 1.62)	-3.77 (± 2.05)	91.64 (± 1.39)	-5.64 (± 1.77)
	ETS	92.89 (± 3.53)	-8.66 (± 4.42)	96.47 (± 0.85)	-3.55 (± 1.07)	93.80 (± 0.82)	-2.84 (± 0.81)
	Heuristic	96.28 (± 1.68)	-4.32 (± 2.11)	95.78 (± 1.50)	-4.46 (± 1.87)	91.75 (± 0.94)	-5.60 (± 2.07)
	Ours	98.20 (± 0.16)	-1.94 (± 0.22)	98.48 (± 0.26)	-1.04 (± 0.31)	93.61 (± 0.71)	-3.11 (± 0.55)
k -center	Random	96.40 (± 0.68)	-4.21 (± 0.84)	95.57 (± 3.16)	-7.20 (± 3.93)	92.61 (± 1.70)	-4.14 (± 2.37)
	ETS	94.84 (± 1.40)	-6.20 (± 1.77)	97.28 (± 0.50)	-2.58 (± 0.66)	91.08 (± 2.48)	-6.39 (± 3.46)
	Heuristic	94.55 (± 2.79)	-6.47 (± 3.50)	94.08 (± 3.72)	-6.59 (± 4.57)	92.06 (± 1.20)	-4.70 (± 2.09)
	Ours	98.24 (± 0.36)	-1.93 (± 0.44)	98.06 (± 0.35)	-1.59 (± 0.45)	94.26 (± 0.37)	-1.97 (± 1.02)
MoF	Random	95.18 (± 3.18)	-5.73 (± 3.95)	95.76 (± 1.41)	-4.41 (± 1.75)	91.33 (± 1.75)	-5.94 (± 1.92)
	ETS	97.04 (± 1.23)	-3.46 (± 1.50)	96.48 (± 1.33)	-3.55 (± 1.73)	92.64 (± 0.87)	-4.57 (± 1.59)
	Heuristic	96.46 (± 2.41)	-4.09 (± 3.01)	95.84 (± 0.89)	-4.39 (± 1.15)	93.24 (± 0.77)	-3.48 (± 1.37)
	Ours	98.37 (± 0.24)	-1.70 (± 0.28)	97.84 (± 0.32)	-1.81 (± 0.39)	94.62 (± 0.42)	-1.80 (± 0.56)
Permutated MNIST		Split CIFAR-100		Split miniImagenet			
Selection	Method	ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑
Uniform	Random	72.44 (± 1.15)	-25.56 (± 1.39)	53.99 (± 0.51)	-34.19 (± 0.66)	48.08 (± 1.36)	-15.98 (± 1.08)
	ETS	71.09 (± 2.31)	-27.39 (± 2.59)	47.70 (± 2.16)	-41.68 (± 2.37)	46.97 (± 1.24)	-18.32 (± 1.34)
	Heuristic	76.68 (± 2.13)	-20.82 (± 2.41)	57.31 (± 1.21)	-30.76 (± 1.45)	49.66 (± 1.10)	-12.04 (± 0.59)
	Ours	76.34 (± 0.98)	-21.21 (± 1.16)	56.60 (± 1.13)	-31.39 (± 1.11)	50.20 (± 0.72)	-13.46 (± 1.22)
k -means	Random	74.30 (± 1.43)	-23.50 (± 1.64)	53.18 (± 1.66)	-35.15 (± 1.61)	49.47 (± 2.70)	-14.10 (± 2.71)
	ETS	69.40 (± 1.32)	-29.23 (± 1.47)	47.51 (± 1.14)	-41.77 (± 1.30)	45.82 (± 0.92)	-19.53 (± 1.10)
	Heuristic	75.57 (± 1.18)	-22.11 (± 1.22)	54.31 (± 3.94)	-33.80 (± 4.24)	49.25 (± 1.00)	-12.92 (± 1.22)
	Ours	77.74 (± 0.80)	-19.66 (± 0.95)	56.95 (± 0.92)	-30.92 (± 0.83)	50.47 (± 0.85)	-13.31 (± 1.24)
k -center	Random	71.41 (± 2.75)	-26.73 (± 3.11)	48.46 (± 0.31)	-39.89 (± 0.27)	44.76 (± 0.96)	-18.72 (± 1.17)
	ETS	69.11 (± 1.69)	-29.58 (± 1.81)	44.13 (± 1.06)	-45.28 (± 1.04)	41.35 (± 0.96)	-23.71 (± 1.45)
	Heuristic	74.33 (± 2.00)	-23.45 (± 2.27)	50.32 (± 1.97)	-37.99 (± 2.14)	44.13 (± 0.95)	-18.26 (± 1.05)
	Ours	76.55 (± 1.16)	-21.06 (± 1.32)	51.37 (± 1.63)	-37.01 (± 1.62)	46.76 (± 0.96)	-16.56 (± 0.90)
MoF	Random	77.96 (± 1.84)	-19.44 (± 2.13)	61.93 (± 1.05)	-25.89 (± 1.07)	54.50 (± 1.33)	-8.64 (± 1.26)
	ETS	77.62 (± 1.12)	-20.10 (± 1.26)	60.43 (± 1.17)	-28.22 (± 1.26)	56.12 (± 1.12)	-8.93 (± 0.83)
	Heuristic	77.27 (± 1.45)	-20.15 (± 1.63)	55.60 (± 2.70)	-32.57 (± 2.77)	52.30 (± 0.59)	-9.61 (± 0.67)
	Ours	81.58 (± 0.75)	-15.41 (± 0.86)	64.22 (± 0.65)	-23.48 (± 1.02)	57.70 (± 0.51)	-5.31 (± 0.55)

Table 5: Performance comparison with ACC and BWT metrics for all datasets between RS-MCTS and the baselines in the setting where only 1 sample per class can be replayed. The memory sizes are set to $M = 10$ and $M = 100$ for the 5-task and 10/20-task datasets respectively. We report the mean and standard deviation of ACC and BWT, where all results have been averaged over 5 seeds. RS-MCTS performs on par with the best baselines for both metrics on all datasets except S-CIFAR-100.

Method	Split MNIST		Split FashionMNIST		Split notMNIST	
	ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑
Random	92.56 (± 2.90)	-8.97 (± 3.62)	92.70 (± 3.78)	-8.24 (± 4.75)	89.53 (± 3.96)	-8.13 (± 5.02)
A-GEM	94.97 (± 1.50)	-6.03 (± 1.87)	94.81 (± 0.86)	-5.65 (± 1.06)	92.27 (± 1.16)	-4.17 (± 1.39)
ER-Ring	94.94 (± 1.56)	-6.07 (± 1.92)	95.83 (± 2.15)	-4.38 (± 2.59)	91.10 (± 1.89)	-6.27 (± 2.35)
Uniform	95.77 (± 1.12)	-5.02 (± 1.39)	97.12 (± 1.57)	-2.79 (± 1.98)	92.14 (± 1.45)	-4.90 (± 1.41)
RS-MCTS (Ours)	96.07 (± 1.60)	-4.59 (± 2.01)	97.17 (± 0.78)	-2.64 (± 0.99)	93.41 (± 1.11)	-3.36 (± 1.56)

Method	Permuted MNIST		Split CIFAR-100		Split miniImagenet	
	ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑
Random	70.02 (± 1.76)	-28.22 (± 1.92)	48.62 (± 1.02)	-39.95 (± 1.10)	48.85 (± 1.38)	-14.55 (± 1.86)
A-GEM	64.71 (± 1.78)	-34.41 (± 2.05)	42.22 (± 2.13)	-46.90 (± 2.21)	32.06 (± 1.83)	-30.81 (± 1.79)
ER-Ring	69.73 (± 1.13)	-28.87 (± 1.29)	53.93 (± 1.13)	-34.91 (± 1.18)	49.82 (± 1.69)	-14.38 (± 1.57)
Uniform	69.85 (± 1.01)	-28.74 (± 1.17)	52.63 (± 1.62)	-36.43 (± 1.81)	50.56 (± 1.07)	-13.52 (± 1.34)
RS-MCTS (Ours)	72.52 (± 0.54)	-25.43 (± 0.65)	51.50 (± 1.19)	-37.01 (± 1.08)	50.70 (± 0.54)	-12.60 (± 1.13)

Table 6: Performance comparison with ACC and BWT metrics between scheduling methods RS-MCTS (Ours), Random, ETS, and Heuristic when combining them with replay-based methods Hindsight Anchor Learning (HAL), Meta Experience Replay (MER), Dark Experience Replay (DER), and DER++. Replay memory sizes are $M = 10$ and $M = 100$ for the 5-task and 10/20-task datasets respectively. We report the mean and standard deviation averaged over 5 seeds. Results on Heuristic where some seed did not converge is denoted by *. Applying RS-MCTS to each method can enhance the performance compared to using the baseline schedules.

Method	Schedule	Split MNIST		Split FashionMNIST		Split notMNIST	
		ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑
HAL	Random	97.24 (± 0.70)	-2.77 (± 0.90)	86.74 (± 6.05)	-15.54 (± 7.58)	93.61 (± 1.31)	-2.73 (± 1.33)
	ETS	97.21 (± 1.25)	-2.80 (± 1.59)	96.75 (± 0.50)	-2.84 (± 0.75)	92.16 (± 1.82)	-5.04 (± 2.24)
	Heuristic	97.69 (± 0.19)	-2.22 (± 0.24)	*74.16 (± 11.19)	-31.26 (± 14.00)	93.64 (± 0.93)	-2.80 (± 1.20)
	Ours	97.96 (± 0.15)	-1.85 (± 0.18)	97.56 (± 0.51)	-2.02 (± 0.63)	94.47 (± 0.82)	-1.67 (± 0.64)
MER	Random	93.07 (± 0.81)	-8.36 (± 0.99)	85.53 (± 3.30)	-0.56 (± 3.36)	91.13 (± 0.86)	-5.32 (± 0.95)
	ETS	92.97 (± 1.73)	-8.52 (± 2.15)	84.88 (± 3.85)	-3.34 (± 5.59)	90.56 (± 0.83)	-6.11 (± 1.06)
	Heuristic	94.30 (± 2.79)	-6.46 (± 3.50)	96.91 (± 0.62)	-1.34 (± 0.76)	90.90 (± 1.30)	-6.24 (± 1.96)
	Ours	96.44 (± 0.72)	-4.14 (± 0.94)	86.67 (± 4.09)	0.85 (± 3.85)	92.44 (± 0.77)	-3.63 (± 1.06)
DER	Random	98.23 (± 0.53)	-1.89 (± 0.65)	96.56 (± 1.79)	-3.48 (± 2.22)	92.89 (± 0.86)	-3.75 (± 1.16)
	ETS	98.17 (± 0.35)	-2.00 (± 0.42)	97.69 (± 0.58)	-2.05 (± 0.71)	94.74 (± 1.05)	-1.94 (± 1.17)
	Heuristic	94.57 (± 1.71)	-6.08 (± 2.09)	*72.49 (± 19.32)	-20.88 (± 11.46)	*77.88 (± 12.58)	-12.66 (± 4.17)
	Ours	99.02 (± 0.10)	-0.91 (± 0.13)	98.33 (± 0.51)	-1.26 (± 0.63)	95.02 (± 0.33)	-0.97 (± 0.81)
DER++	Random	97.90 (± 0.52)	-2.32 (± 0.67)	97.10 (± 1.03)	-2.77 (± 1.29)	93.29 (± 1.43)	-3.11 (± 1.60)
	ETS	97.98 (± 0.52)	-2.24 (± 0.66)	98.12 (± 0.40)	-1.59 (± 0.52)	94.53 (± 1.02)	-1.82 (± 1.02)
	Heuristic	92.35 (± 2.42)	-8.83 (± 2.99)	*67.31 (± 21.20)	-24.86 (± 16.34)	93.88 (± 1.33)	-2.86 (± 1.49)
	Ours	98.84 (± 0.21)	-1.14 (± 0.26)	98.38 (± 0.43)	-1.17 (± 0.51)	94.73 (± 0.20)	-1.21 (± 1.12)

Method	Schedule	Permuted MNIST		Split CIFAR-100		Split miniImagenet	
		ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑	ACC(%)↑	BWT(%)↑
HAL	Random	88.49 (± 0.99)	-7.03 (± 1.05)	36.09 (± 1.77)	-17.49 (± 1.78)	38.51 (± 2.22)	-6.65 (± 1.43)
	ETS	88.46 (± 0.86)	-7.26 (± 0.90)	34.90 (± 2.02)	-18.92 (± 0.91)	38.13 (± 1.18)	-8.19 (± 1.73)
	Heuristic	*66.63 (± 28.50)	-29.68 (± 27.90)	35.07 (± 1.29)	-24.76 (± 2.41)	39.51 (± 1.49)	-5.65 (± 0.77)
	Ours	89.14 (± 0.74)	-6.29 (± 0.74)	40.22 (± 1.57)	-12.77 (± 1.30)	41.39 (± 1.15)	-3.69 (± 1.86)
MER	Random	75.90 (± 1.34)	-21.69 (± 1.47)	42.96 (± 1.70)	-34.01 (± 2.07)	31.48 (± 1.65)	-6.99 (± 1.27)
	ETS	73.01 (± 0.96)	-25.19 (± 1.10)	43.38 (± 1.81)	-34.84 (± 1.98)	33.58 (± 1.53)	-6.80 (± 1.46)
	Heuristic	83.88 (± 3.19)	-12.48 (± 3.60)	40.90 (± 1.70)	-44.10 (± 2.03)	34.22 (± 1.93)	-7.57 (± 1.63)
	Ours	79.72 (± 0.71)	-17.42 (± 0.78)	44.29 (± 0.69)	-32.73 (± 0.88)	32.74 (± 1.29)	-5.77 (± 1.04)
DER	Random	87.51 (± 1.10)	-8.81 (± 1.28)	56.83 (± 0.76)	-27.34 (± 0.63)	42.19 (± 0.67)	-10.60 (± 1.28)
	ETS	85.71 (± 0.75)	-11.15 (± 0.87)	52.58 (± 1.49)	-32.93 (± 2.04)	35.50 (± 2.84)	-10.94 (± 2.21)
	Heuristic	81.56 (± 2.28)	-15.06 (± 2.51)	55.75 (± 1.08)	-31.27 (± 1.02)	43.62 (± 0.88)	-8.18 (± 1.16)
	Ours	90.11 (± 0.18)	-5.89 (± 0.23)	58.99 (± 0.98)	-24.95 (± 0.64)	43.46 (± 0.95)	-9.32 (± 1.37)
DER++	Random	87.89 (± 1.10)	-8.35 (± 1.33)	58.49 (± 1.44)	-25.48 (± 1.55)	48.40 (± 0.69)	-4.20 (± 0.86)
	ETS	85.25 (± 0.88)	-11.60 (± 1.03)	52.54 (± 1.06)	-33.22 (± 1.51)	41.36 (± 2.90)	-4.07 (± 2.28)
	Heuristic	79.17 (± 2.44)	-17.68 (± 2.68)	56.70 (± 1.27)	-30.33 (± 1.41)	45.73 (± 0.84)	-6.09 (± 1.24)
	Ours	89.84 (± 0.22)	-6.13 (± 0.29)	59.23 (± 0.83)	-24.61 (± 0.91)	49.45 (± 0.68)	-3.12 (± 0.89)

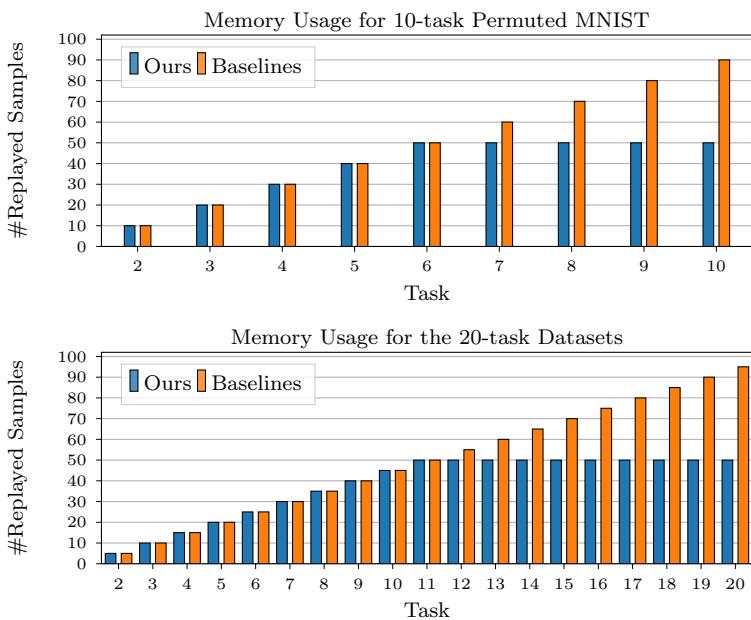


Figure 8: Number of replayed samples per task for 10-task Permutated MNIST (top) and the 20-task datasets in the experiment in Section 4.6. The fixed memory size $M = 50$ for our method is reached after learning task 6 and task 11 on the Permutated MNIST and the 20-task datasets respectively, while the baselines continue incrementing their number of replay samples per task.

References

- [1] Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. Continual learning: A comparative study on how to defy forgetting in classification tasks. *arXiv preprint arXiv:1909.08383*, 2(6), 2019.
- [2] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- [3] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K Dokania, Philip HS Torr, and Marc’Aurelio Ranzato. On tiny episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*, 2019.
- [4] Tyler L Hayes, Kushal Kafle, Robik Shrestha, Manoj Acharya, and Christopher Kanan. Remind your neural network to prevent catastrophic forgetting. In *European Conference on Computer Vision*, pages 466–483. Springer, 2020.
- [5] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. *arXiv preprint arXiv:1903.08671*, 2019.
- [6] Zalán Borsos, Mojmir Mutný, and Andreas Krause. Coresets via bilevel optimization for continual learning and streaming. *arXiv preprint arXiv:2006.03875*, 2020.
- [7] Aristotelis Chrysakis and Marie-Francine Moens. Online continual learning from imbalanced data. In *ICML*, 2020.
- [8] Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational continual learning. *arXiv preprint arXiv:1710.10628*, 2017.
- [9] Sylvestre-Alvise Rebiffé, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010, 2017.
- [10] Jaehong Yoon, Divyam Madaan, Eunho Yang, and Sung Ju Hwang. Online coresnet selection for rehearsal-based continual learning. *arXiv preprint arXiv:2106.01085*, 2021.
- [11] Ahmet Iscen, Jeffrey Zhang, Svetlana Lazebnik, and Cordelia Schmid. Memory-efficient incremental learning through feature adaptation. In *European Conference on Computer Vision*, pages 699–715. Springer, 2020.
- [12] Lorenzo Pellegrini, Gabriele Graffieti, Vincenzo Lomonaco, and Davide Maltoni. Latent replay for real-time continual learning. *arXiv preprint arXiv:1912.01100*, 2019.
- [13] Frank N Dempster. Spacing effects and their implications for theory and practice. *Educational Psychology Review*, 1(4):309–330, 1989.
- [14] Hermann Ebbinghaus. Memory: A contribution to experimental psychology. *Annals of neurosciences*, 20(4):155, 2013.
- [15] T. Landauer and Robert Bjork. Optimum rehearsal patterns and name learning. *Practical aspects of memory*, 1, 11 1977.
- [16] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*, pages 3987–3995. PMLR, 2017.

- [17] Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *International conference on computers and games*, pages 72–83. Springer, 2006.
- [18] Tameem Adel, Han Zhao, and Richard E Turner. Continual learning with adaptive weights (claw). *arXiv preprint arXiv:1911.09514*, 2019.
- [19] Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 532–547, 2018.
- [20] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [21] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- [22] Amal Rannen, Rahaf Aljundi, Matthew B Blaschko, and Tinne Tuytelaars. Encoder based lifelong learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1320–1328, 2017.
- [23] Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual learning. In *International Conference on Machine Learning*, pages 4528–4537. PMLR, 2018.
- [24] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- [25] Jaehong Yoon, Saehoon Kim, Eunho Yang, and Sung Ju Hwang. Scalable and order-robust continual learning with additive parameter decomposition. *arXiv preprint arXiv:1902.09432*, 2019.
- [26] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically expandable networks. *arXiv preprint arXiv:1708.01547*, 2017.
- [27] Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7765–7773, 2018.
- [28] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *International Conference on Machine Learning*, pages 4548–4557. PMLR, 2018.
- [29] David Isele and Akansel Cosgun. Selective experience replay for lifelong learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [30] David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. *arXiv preprint arXiv:1706.08840*, 2017.
- [31] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. *arXiv preprint arXiv:1705.08690*, 2017.
- [32] Gido M van de Ven and Andreas S Tolias. Generative replay with feedback connections as a general strategy for continual learning. *arXiv preprint arXiv:1809.10635*, 2018.

- [33] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *arXiv preprint arXiv:2004.07211*, 2020.
- [34] Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. *arXiv preprint arXiv:1812.00420*, 2018.
- [35] Arslan Chaudhry, Albert Gordo, Puneet Dokania, Philip Torr, and David Lopez-Paz. Using hindsight to anchor past knowledge in continual learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 6993–7001, 2021.
- [36] Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In *Computer Vision – ECCV 2020*, pages 86–102. Springer International Publishing, 2020.
- [37] Sayna Ebrahimi, Franziska Meier, Roberto Calandra, Trevor Darrell, and Marcus Rohrbach. Adversarial continual learning. *arXiv preprint arXiv:2003.09553*, 2020.
- [38] KJ Joseph and Vineeth N Balasubramanian. Meta-consolidation for continual learning. *arXiv preprint arXiv:2010.00352*, 2020.
- [39] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Dilan Gorur, Razvan Pascanu, and Hassan Ghasemzadeh. Linear mode connectivity in multitask and continual learning. *arXiv preprint arXiv:2010.04495*, 2020.
- [40] Pingbo Pan, Siddharth Swaroop, Alexander Immer, Runa Eschenhagen, Richard E Turner, and Mohammad Emtiyaz Khan. Continual deep learning by functional regularisation of memorable past. *arXiv preprint arXiv:2004.14070*, 2020.
- [41] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy P Lillicrap, and Greg Wayne. Experience replay for continual learning. *arXiv preprint arXiv:1811.11682*, 2018.
- [42] Johannes von Oswald, Christian Henning, João Sacramento, and Benjamin F Grewe. Continual learning with hypernetworks. *arXiv preprint arXiv:1906.00695*, 2019.
- [43] Tyler L Hayes, Nathan D Cahill, and Christopher Kanan. Memory efficient experience replay for streaming learning. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 9769–9776. IEEE, 2019.
- [44] Jeffrey S Vitter. Random sampling with a reservoir. *ACM Transactions on Mathematical Software (TOMS)*, 11(1):37–57, 1985.
- [45] Robert M. French. Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, 3(4):128–135, 1999.
- [46] John Dunlosky, Katherine A. Rawson, Elizabeth J. Marsh, Mitchell J. Nathan, and Daniel T. Willingham. Improving students’ learning with effective learning techniques: Promising directions from cognitive and educational psychology. *Psychological Science in the Public Interest*, 14(1):4–58, 2013.
- [47] Judy Willis. Review of research: Brain-based teaching strategies for improving students’ memory, learning, and test-taking success. *Childhood Education*, 83(5):310–315, 2007.
- [48] Karri S Hawley, Katie E Cherry, Emily O Boudreux, and Erin M Jackson. A comparison of adjusted spaced retrieval versus a uniform expanded retrieval schedule for learning a name–face association in older adults with probable alzheimer’s disease. *Journal of Clinical and Experimental Neuropsychology*, 30(6):639–649, 2008.

- [49] Hadi Amiri, Timothy Miller, and Guergana Savova. Repeat before forgetting: Spaced repetition for efficient and effective training of neural networks. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2401–2410, 2017.
- [50] Kanyin Feng, Xiao Zhao, Jing Liu, Ying Cai, Zhifang Ye, Chuansheng Chen, and Gui Xue. Spaced learning enhances episodic memory by increasing neural pattern similarity across repetitions. *Journal of Neuroscience*, 39(27):5351–5360, 2019.
- [51] Paul Smolen, Yili Zhang, and John H Byrne. The right time to learn: mechanisms and optimization of spaced learning. *Nature Reviews Neuroscience*, 17(2):77, 2016.
- [52] Philip J Ball, Yingzhen Li, Angus Lamb, and Cheng Zhang. A study on efficiency in continual learning inspired by human learning. *arXiv preprint arXiv:2010.15187*, 2020.
- [53] Gido M van de Ven, Hava T Siegelmann, and Andreas S Tolias. Brain-inspired replay for continual learning with artificial neural networks. *Nature communications*, 11(1):1–14, 2020.
- [54] Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in games*, 4(1):1–43, 2012.
- [55] Muhammad Umar Chaudhry and Jee-Hyong Lee. Feature selection for high dimensional data using monte carlo tree search. *IEEE Access*, 6:76036–76048, 2018.
- [56] Sylvain Gelly, Yizao Wang, Rémi Munos, and Olivier Teytaud. Modification of uct with patterns in monte-carlo go. Technical Report RR-6062, INRIA, 2006. inria-00117266v3f.
- [57] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- [58] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *European conference on machine learning*, pages 282–293. Springer, 2006.
- [59] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [60] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [61] Yaroslav Bulatov. The notMNIST dataset. <http://yaroslavvb.com/upload/notMNIST/>, 2011.
- [62] Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*, 2013.
- [63] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [64] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. *arXiv preprint arXiv:1606.04080*, 2016.

- [65] Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauro. Learning to learn without forgetting by maximizing transfer and minimizing interference. *arXiv preprint arXiv:1810.11910*, 2018.
- [66] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- [67] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Paper D

Meta Policy Learning for Replay Scheduling in Continual Learning

Marcus Klasson*, Hedvig Kjellström*, Cheng Zhang†

*KTH Royal Institute of Technology, Stockholm, Sweden

†Microsoft Research, Cambridge, United Kingdom

Abstract

Learning the time to replay different tasks have been demonstrated to be important in continual learning. However, a replay scheduling policy that can be applied in any continual learning scenario is currently absent, which makes replay scheduling infeasible in real-world scenarios. To this end, we propose using reinforcement learning to enable learning general policies that can generalize across different data domains. Thus, an efficient replay schedule can be used in any new real-world continual learning scenario to improve the continual learning performance without any additional computational cost. We compare the learned policies to several replay scheduling baselines and show that the learned policies can improve the continual learning performance given unseen tasks and datasets.

1 Introduction

The timing when to rehearse on previously learned tasks is important for memory retention in human continual learning (CL) [1–3], however, there currently exist no such method for CL in neural networks. Memory-based CL [4, 5] have been shown to efficiently retain knowledge of previously learned tasks by rehearsing the network on samples from historical tasks stored in a limited replay memory. Most methods simply apply random sampling for retrieving mini-batches from the replay memory since more elaborate retrieval strategies have been shown to yield small performance improvements [6, 7]. Recently, Klasson et al. [8] showed that replay scheduling over which tasks to replay at different time steps can be critical for the final CL performance. However, a general replay scheduling policy that can be applied in any CL scenario is currently missing, which makes it infeasible to apply in real-world scenarios.

Assuming we have a policy for replay scheduling, we would like the policy to have the property of being transferable to any CL scenario. Ideally, the policy should be capable of selecting the tasks to replay in new CL environments without additional training with experiences from the environment. This would be useful in real-world applications where a machine learning system is continuously updated from streaming data and there are vast amounts of historical data accessible, e.g., from the cloud, that can be used for mitigating catastrophic forgetting. In such situations, the system would likely have constraints on the processing time due to operational costs prohibiting training on all the available historical data. Therefore, replay scheduling is important for selecting what historical data to replay for obtaining the best overall performance when new tasks are learned. In this paper, we present a novel method for learning a general replay scheduling policy that satisfies the above mentioned properties to move CL research closer to real-world needs.

We propose using reinforcement learning (RL) [9] for learning a general policy for replay scheduling that can be applied to any new CL scenario. We focus on the setting for memory-based CL proposed in Klasson et al. [8] where limitations are on the processing time of memory samples rather than the memory storage capacity. We take inspiration from meta-learning [10] and learn the policy from experiences from multiple episodes of CL procedures. The RL agent takes actions within the CL procedures where each action represents which tasks to replay for the CL model, where the actions are taken based on the current performance on the seen tasks. Thus, the goal for the agent is to learn the time when the CL model should replay different tasks and by what proportion to use of memory samples to fill the replay memory with. To foster generalization between new CL scenarios, we let the agent interact with several different CL models with various variations on the task datasets they should learn from. Finally, we evaluate the obtained policy by applying it to new CL scenarios without adding any computational overhead by training on experiences from the test environments. Our contributions are the following:

- We propose a novel method based on RL for learning a policy for replay scheduling in CL settings. The learned policy can then be transferred to any new CL scenario without any additional training.
- We evaluate the generalization capabilities of the learned policy on unseen CL scenarios where we have varied the task splits in the datasets as well as introducing unseen datasets for testing the policy.

2 Background

In this section, we recall the CL setting from Klasson et al. [8] and give a brief overview of RL and Deep Q-Networks (DQNs) [11, 12].

Continual Learning Setting. We target the CL setting presented in Klasson et al. [8] where historical data is assumed to be accessible at any time to connect the current CL research with real-world scenarios where historical data is accessible at any time. The CL model parameterized by ϕ should learn T tasks arriving sequentially in the form of T datasets $\mathcal{D}_{1:T} = \{\mathcal{D}_1, \dots, \mathcal{D}_T\}$. The dataset for the t -th task $\mathcal{D}_t = \{(\mathbf{x}_t^{(i)}, y_t^{(i)})\}_{i=1}^{N_t}$ where $\mathbf{x}_t^{(i)}$ and $y_t^{(i)}$ are the i -th data point and corresponding

class label among a total of N_t examples in the dataset. In contrary to traditional CL, all seen data are available to use for replay here. However, we are only allowed to fill a tiny replay memory \mathcal{M} of size M with historical data due to limitations on processing time. Hence, we need a method for selecting which tasks to replay and the amount of samples to add to the replay memory from each selected task.

Reinforcement Learning. We treat the CL setting described above as an RL problem where an agent interacts with an environment with a discrete sequence of states, actions, and rewards of length T . The environment is formalized as a Markov Decision Process (MDP) represented as a tuple $E = (\mathcal{S}, \mathcal{A}, P, r, \mathcal{S}_1, \gamma)$ consisting of state space \mathcal{S} , action space \mathcal{A} , transition distribution $P(s_{t+1}|s_t, a_t)$, reward function $r(s_t, a_t)$, initial state distribution \mathcal{S}_1 , and discount factor $\gamma \in [0, 1]$ [13]. The action selection strategy is represented as the policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$. The goal for the agent is to learn an optimal policy π^* for selecting actions that maximize the expected sum of rewards. In this paper, the learned policy should select actions to maximize the CL performance for the model measured using classification accuracy. Therefore, we represent the states $s_t \in \mathcal{S}$ as the classification performance of model ϕ and the actions $a_t \in \mathcal{A}$ represent task proportions over the tasks to fill the replay memory \mathcal{M} with.

Deep Q-Networks. We employ DQNs for learning a policy for replay scheduling. For DQNs, the goal is to estimate the state-action value function recursively using the Bellman equation:

$$Q^\pi(s_t, a_t) = \mathbb{E}_{s_{t+1} \sim P}[r_t + \gamma \max_{a'} Q^\pi(s_{t+1}, a')]. \quad (1)$$

In many applications, the transition distribution P is unknown, so we collect data from the environment to learn a function $Q_\theta(s, a)$ parameterized by θ that estimates the Q-function by minimizing the loss $L(\theta) = (y_t - Q_\theta(s_t, a_t))^2$ where $y_t = r_t + \gamma \max_{a'} Q_\theta(s_{t+1}, a')$. The learned Q-function is parameterized by a neural network to ease the estimation in problems with high-dimensional state spaces. To stabilize the learning, DQNs use target networks from previous iterations to compute the target values y_t , incorporates experience replay [14], and can apply Double Q-learning [15] to mitigate overestimation of the state-action values.

3 Related Work

Continual Learning. One of the main goals of CL methods is to mitigate catastrophically forgetting previous tasks when learning new tasks. Approaches for retaining knowledge of previously learned tasks can be broadly divided into regularization-based methods [16–19], dynamic architectures [20–22], and memory-based methods [6, 7, 23], where the latter is the approach we focus on in this paper. Regularization-based methods put constraints on parameters important for previous tasks to mitigate forgetting and use less salient parameters for learning new tasks. Dynamic architecture approaches increase the capacity of the network by adding more parameters as new tasks are encountered. Memory-based methods mitigate forgetting by replaying samples from old tasks that are stored in a small memory buffer. This work is based on memory-based approaches, but in contrast to the traditional CL setting, we assume that we have a large amount of storage for historical data and can

therefore fill the replay memory with historical data from the storage before learning new tasks [8]. Due to processing time constraints, we need a method for selecting which historical tasks to replay at every time step.

Large and Changing Action Spaces in RL. Several works focus on bringing RL closer to real-world scenarios where the action spaces are large as well as can change size over time or given certain states. [24] proposes to learn action embeddings where k -nearest neighbors are used for selecting valid actions which the policy can select from to enable RL in large action spaces. [25] proposes to let the agent reason about which actions to take in a lower-dimensional representation space and transforms decisions in this space into actual actions which improve generalization in large action spaces when the agent infers outcomes of actions that are similar to actions already taken. [26] takes this a step further and learns action representations that should generalize to completely new actions in a zero-shot setting. [?, 27] present how RL can be applied to stochastic action sets, while [28] proposes algorithms that adapt to action spaces where the set of available actions changes over time. Several attempts have been made to scale up Q-learning to high-dimensional or continuous action spaces. [29] applies Q-learning to action spaces for natural language where the Q-function is approximated by the dot-product between state and action embeddings for measuring the relevance of performing the action. [30] comes around the maximization over all actions by instead taking the maximization over a subset of actions sampled from a learned proposal distribution, and [31] discretizes each action dimension into bins where the maximization is performed and predicts Q-value for each dimension sequentially. We assume that the structure of the action space is known as the number of tasks to replay grows linearly with every seen task. Since we also assume the task-horizon is known, we know the number of available actions at the end of the episodes which allows us to use a static architecture and use masking for selecting valid actions.

Generalization in RL. Generalization in RL is an active research field where much focus has been put on developing proper benchmark datasets for evaluating generalization capabilities of RL agents [?, ?, ?, ?]. Regularization techniques from supervised learning have been used to investigate whether these can enhance generalization in RL [?, 32, 273], and also how variations in the environments help to obtain agents that generalize better [33, 275]. The survey by [34] focuses on zero-shot policy transfer where a policy is evaluated on environments different from it was trained on. Furthermore, access to rewards or additional training is disallowed in the test environments with the aim to enable RL in real-world scenarios [35]. We evaluate our learned policies on the zero-shot policy transfer setting where the goal is to mitigate catastrophic forgetting in new CL scenarios. The rewards are accessible through accuracies calculated on validation datasets, however, fine-tuning the policy during CL training is prohibited.

4 Methodology

In this section, we describe our problem setting for learning replay scheduling policies and also how we employ DQNs to learn such policies.

4.1 Problem Setting

Our goal is to learn a replay scheduling policy for selecting which tasks to replay at different times that can generalize to new CL scenarios without additional training in the new environment. We apply RL for learning the policy where a Deep Q-Network (DQN) acts as a scheduler of the replay tasks for a CL model. We train the policy by collecting experience from multiple learning episodes in different CL environments to improve policy transfer. At test time, we let the policy be used to select replay schedules for new CL scenarios.

Each environment is modeled by its own MDP $E = (\mathcal{S}, \mathcal{A}, P, r, \mathcal{S}_1, \gamma) \in \mathcal{E}$ with an accompanied classifier ϕ^E and CL datasets $\mathcal{D}_{1:T}^E$ of T tasks. We assume the task-horizon T is finite and known. The agent passes action a_t^E at task t to environment E for building the replay memory \mathcal{M}_t before the classifier ϕ^E learns task $t+1$. After learning the task, the agent receives the reward r_t^E from applying action a_t^E and the state of the environment s_{t+1}^E . We calculate the rewards using the average validation accuracy over all seen tasks, i.e.

$$r_t^E = \frac{1}{t} \sum_{i=1}^t A_{t,i}^{(val)}, \quad (2)$$

where the accuracy $A_{t,i}^{(val)}$ is calculated using the validation dataset \mathcal{V}_i^E . The goal for the agent is to maximize the sum of rewards after the classifier has learnt the final task T . We have summarized the CL procedure in Algorithm 3. Next, we describe how the states and action spaces are defined.

We are dealing with state and action spaces that are growing per seen task. We wish for the DQN to select which tasks to replay based on the current performance of the CL model ϕ . Therefore, we represent the states with a vector of the currently seen accuracies, such that $s_t \in [0, 1]^t = \mathcal{S}_t$ where the state space domain \mathcal{S}_t grows linearly with t . We keep a validation set for each task to calculate the task accuracies to use for the states. The action space is also growing per task, such that $\mathcal{A} = (\mathcal{A}_1, \dots, \mathcal{A}_{T-1})$ where $|\mathcal{A}_t| < |\mathcal{A}_{t+1}|$ for $t = 1, \dots, T-1$. Each action $a_t \in \mathcal{A}_t$ is represented as a vector of task proportions to use for filling the replay memory. More specifically, the action $a_t = [p_1, \dots, p_t]$ where p_i is the task proportion for task i which satisfies $\sum_{i=1}^t p_i = 1$ for $t \in \{1, \dots, T-1\}$. We construct the action space and the number of actions to use in the same way as in Klasson et al. [8] which we describe in more detail in Appendix 2.

4.2 Deep Q-Networks for Learning Replay Scheduling Policy

We employ Deep Q-Networks [11] (DQNs) as the RL agent for learning the replay scheduling policy. We assume that the number of tasks T is known which makes it easier to set up the DQN architecture. The output layer size to be the same as the number of actions in the last action space, i.e., $|\mathcal{A}_{T-1}|$. As earlier action spaces has fewer actions, i.e. $|\mathcal{A}_{t-1}| < |\mathcal{A}_t|$, we use masking for preventing the DQN to select

Algorithm 3 Continual Learning Step

Require: ϕ : Continual learning model
Require: $\mathcal{D}_{1:T}$: Task datasets, a_t : action from policy

- 1: Sample replay memory \mathcal{M}_t of historical data using the task proportions in action a_t
- 2: Train ϕ on task $t + 1$ with data $\mathcal{D}_{t+1}^{(train)} \cup \mathcal{M}_t$
- 3: Calculate reward r_t using Eq. 2
- 4: Get next state s_{t+1} with validation performance of ϕ **return** reward r_t and next state s_{t+1}

Algorithm 4 Learning Replay Scheduling Policy with DQN

Require: θ : DQN parameters,
Require: \mathcal{E} : Set of training environments,

- 1: Initialize replay buffer \mathcal{B}
- 2: **for** episode = 1, ..., n_{episodes} **do**
- 3: Initialize s_1^E for every $E \in \mathcal{E}$
- 4: **for** $t = 1, T - 1$ **do**
- 5: **for** $E \in \mathcal{E}$ **do**
- 6: Take action $a_t^E = \arg \max_{a \in \mathcal{A}_t} Q_\theta(s_t^E, a)$ or random action with probability ϵ
- 7: Execute a_t^E in E with Alg. 3 and observe reward r_t^E and next state s_{t+1}^E
- 8: Store transition $(s_t^E, a_t^E, r_t^E, s_{t+1}^E)$ in \mathcal{B}
- 9: Sample mini-batch $(s_j, a_j, r_j, s_{j+1}) \sim \mathcal{B}$
- 10: Perform gradient descent step on $(y_j - Q_\theta(s_j, a_j))^2$ with target values y_j
- 11: **end for**
- 12: **end for**
- 13: **end for**

invalid actions. The input layer is of fixed size $T - 1$, such that we can input the task performances of the seen tasks as states. We apply zero-padding on the input state vectors for future tasks that have not been evaluated yet.

We let the DQN agent collects experience from multiple CL environments for training the policy to enhance the generalization capability. We hypothesize that the policy will get better at learning the tasks to replay based on the task performance through learning experiences from multiple environments with different variations on the datasets in each environment, e.g., different task splits or datasets (Split MNIST, Split FashionMNIST, etc.). We store the experiences from all environments in a single replay buffer used for updating the DQN, where we denote the experience from environment E as $(s_t^E, a_t^E, r_t^E, s_{t+1}^E)$. The replay buffer is shared among the environments to obtain a diverse set of outcomes for learning a general policy that can be used for replay scheduling by only selecting actions with the current task performances of a classifier. See Algorithm 4 for the steps on how to train the DQN.

At test time, we evaluate the policy on CL environments with different datasets

than seen during training. For example, a different dataset can mean that we have changed the task split of a dataset we used for training, or that we evaluate on a new dataset. We assume that the task-horizons are the same on the datasets used for both training and testing such that we can keep the DQN architecture the same. The DQN acts greedily during testing selecting the next action using

$$a_t^{E'} = \arg \max_a Q_\theta(s_t^{E'}, a) \quad (3)$$

in test environment E' at task t . The hope is that we should have obtained a general policy for replay scheduling that generalizes to new CL scenarios.

5 Experiments

We perform two experiments to evaluate how well the learned replay scheduling policy can generalize to new CL scenarios. First, we investigate whether the policy can be applied to environments with new dataset task splits which are different from the splits that the policy was trained on. Secondly, we study if the policy can be applied to environments with new datasets than it was originally trained on. We run experiments on CL datasets of $T = 5$ tasks to have a discrete action space of reasonable size. The replay memory size is set to $M = 10$ in all experiments. We use 15% of the training set for validation for the methods where it applies. See Appendix 1 for more details on the experimental settings.

For measuring the generalization performance, we use a ranking method where we calculate of average ranking of each method based on the CL performance the scheduling method achieved in each environment. We use this ranking method is used because the accuracies for measuring the CL performance differs between the environments, which can make it difficult to draw conclusions from averaging the accuracies instead of the ranks. To measure the CL performance of the classifiers, we use the average test accuracy over all tasks after learning the final task, i.e.,

$$\text{ACC} = \frac{1}{T} \sum_{i=1}^T A_{T,i}^{(test)}. \quad (4)$$

For comparing each method across the environments, we compare the ACC achieved by each method in each environment and rank the methods in descending order of the ACC score. The average ranking is then computed by taking the average over all ranking positions in every test environment for each method separately. Moreover, we perform qualitative assessments of our method by visualizing the learned policies.

5.1 Baselines

We compare our proposed method to five baselines for scheduling. Three of the baselines use heuristic rules for scheduling the replay tasks which are based on the intuition that forgotten tasks should be replayed. Hence, we use a validation set for selecting which task to replay for the heuristic baselines. Note that heuristic baselines skip replaying any tasks if no tasks satisfies the conditions on the validation accuracy.

Random. Random policy that randomly selects task proportions from the action space on how to structure the replay memory at every task.

Table 1: Average ranking (Avg. Rank) and median of rankings for different label split experiment with Split MNIST and Split CIFAR-10. We train the DQN on 10 training envs and evaluate on 10 test envs. We average the results over 5 DQNs and the 10 test envs.

Method	Split MNIST			Split CIFAR-10		
	Avg. Rank (\downarrow)	Median (\downarrow)	IQR	Avg. Rank (\downarrow)	Median (\downarrow)	IQR
Random	3.70 ± 1.88	4.00	3.00	3.98 ± 1.88	5.00	3.75
ETS	4.38 ± 1.70	5.00	2.00	3.74 ± 2.08	4.00	4.00
Heur-GD	4.16 ± 0.95	4.00	2.00	3.84 ± 1.17	4.00	2.00
Heur-LD	2.80 ± 1.46	2.00	1.00	3.28 ± 1.30	3.50	1.00
Heur-AT	2.86 ± 1.25	3.00	2.00	3.26 ± 1.52	3.00	2.00
DQN (Ours)	3.10 ± 2.05	3.00	4.00	2.90 ± 1.85	3.00	3.00

Equal Task Schedule (ETS). Policy that selects equal task proportion such that the replay memory aims to fill the memory with an equal number of samples from every seen task.

Heuristic Global Drop (Heur-GD). This heuristic uses ratio constant $\tau \in [0, 1]$ representing the degree of how much the validation accuracy of a task is allowed to drop compared to a reference accuracy. The reference is the best measured validation accuracy $A_{t,i}^{(best)} = \max\{(A_{1,i}^{(val)}, \dots, A_{t,i}^{(val)})\}$ for task i , such that we check the task performance for task i globally over time. Therefore, task i is replayed if $A_{t,i} < \tau A_{t,i}^{(best)}$ after learning task t .

Heuristic Local Drop (Heur-LD). Similar to Heur-GD, this heuristic also uses a ratio constant $\tau \in [0, 1]$ representing the degree of allowed drop in validation performance. The reference accuracy that the current accuracy is compared locally against is from the previous time step, such that we replay task i when $A_{t,i} < \tau A_{t-1,i}$ after learning task t .

Heuristic Accuracy Threshold (Heur-AT). The third heuristic has a fixed threshold $\tau \in [0, 1]$ on the accuracy such that task i is replayed after learning task t if $A_{t,i} < \tau$.

5.2 Generalization to New Continual Learning Environments

In this experiment, we evaluate the generalization capability of the DQN to new CL environments unseen during training. First, we compare the performance of the DQN and the scheduling baselines when they are tested on environments with task splits different from the training environments, where we evaluate each method Split MNIST [36] and Split CIFAR-10 [?]. Secondly, we evaluate the capability of the DQN to select replay schedules to new CL datasets by training on environments with Split MNIST and Split notMNIST [37] datasets and testing on environments with Split FashionMNIST [38]. We use 10 test environments for evaluating the generalization performance in all experiments where the task split is different in each environment.

Table 1 shows the average ranking across 10 testing environments for our method and the baselines when experimenting with Split MNIST and Split CIFAR-10. We

Table 2: Average ranking (Avg. Rank) and median of rankings for training with mixed datasets and test on split FashionMNIST experiment. We train the DQN on 10 training environments (5 MNIST + 5 notMNIST) and evaluate on 10 test environments of FashionMNIST. We average the results over 5 DQNs and the 10 test environments.

Method	Split FashionMNIST		
	Avg. Rank (\downarrow)	Median (\downarrow)	IQR
Random	3.80 ± 1.71	4.00	3.00
ETS	4.60 ± 1.55	5.00	2.00
Heur-GD	2.56 ± 1.00	3.00	1.00
Heur-LD	4.12 ± 1.61	5.00	3.75
Heur-AT	2.96 ± 1.44	3.00	2.00
DQN (Ours)	2.96 ± 1.83	3.00	3.00

also show the median and interquartile range (IQR) of the ranks for each method. Our scheduling method performs better than the Random, ETS, and Heur-GD policies and is almost on par with Heur-LD and Heur-AT on both datasets. In the CIFAR-10 experiment, we observe that the DQN has a slightly better average ranking than the Heur-LD and Heur-AT baselines. This could mean that learning the replay scheduling policy can be more beneficial than creating heuristic scheduling rules for challenging datasets where the classifiers are more sensitive to the selection of replay tasks.

Table 2 shows the average ranking across 10 testing environments with Split FashionMNIST datasets for our method when 10 (5+5) environments of both Split MNIST and notMNIST have been used for training. We observe that the DQN is the second-best method and performs on par with Heur-AT while Heur-GD achieves the best result. Furthermore, the DQN performs significantly better than the ETS policy which shows that scheduling can outperform replay strategies that are common in memory-based CL. These results indicate that learning the replay scheduling policy is a promising approach when applying memory-based CL methods in new scenarios.

5.3 Visualization of Learned Policies

We visualize the learned policies for the DQN to bring further insights into the advantages of learning the replay scheduling policy. We analyze the task accuracy progress and the replay schedule used for mitigating catastrophic forgetting by visualizing the schedule as a bubble plot. In the bubble plots shown in Figure ??, the x- and y-axis shows the current and replay task respectively, and the bubbles represent the proportion of the replayed task in the replay memory at the current task. Each color of the circles corresponds to a seen task that is being replayed at the current task. The sum of all points in the circles at a current task column is fixed since the replay memory size is constant.

Figure ?? shows the task accuracy progression and replay schedules from the DQN and Heur-LD in one test environment with Split MNIST. Using the replay schedule from the DQN gives higher overall CL performance since the DQN manages

to mitigate forgetting tasks 1 and 2 better than Heur-LD. The DQN decides to replay task 1 and task 2 as soon as these tasks have been learned. At task 4, the DQN decides to replay both tasks 1 and 2 which helps the classifier to retain its knowledge of task 2. However, as the classifier forgets task 1, the DQN selects to only replay task 1 at task 5 which makes the classifier improve its performance on task 1. The heuristic baseline Heur-LD enables replay on tasks 3 and 5 according to its rule since the performance of task 1 and 2 has decreased below its threshold at those time steps, but the classifier forgets these tasks more than when using the learned replay schedule from the DQN. This example shows creating a good heuristic for replay scheduling can be difficult which is why learning such a policy is preferable.

In Figure ??, we further illustrate the flexibility of the learned policy when the DQN is applied to Split FashionMNIST which was unseen during training. Here, we observe that the DQN schedule exhibits a replay pattern for task 1 that is similar to spaced repetition used in human learning, which helps to achieve better performance on task 2 as the DQN only replays this task when learning task 3. This shows the benefits of learning policies for replay scheduling as this enables scheduling behaviors that would be difficult to create with heuristics.

6 Discussion

In this section, we provide a discussion of the experimental results in Section 5. Our goal was to answer whether it was possible to learn a policy for replay scheduling that can be applied in new CL scenarios without additional training cost. We evaluate the policy on two CL scenarios where (i) the training and test environments share the same dataset but the task splits are different, and (ii) the datasets in the training and test environments are completely different. We observed that the DQN receives a similar average ranking of the performance in the test environments as the best heuristic scheduling baselines in both CL scenarios. Furthermore, we analyzed the learned policies by visualizing them together with the test accuracy progression for each task to illustrate the benefits of the DQN. We then showed an example in Figure ?? where learning the replay scheduling policy indeed opens up for scheduling behaviors similar to human learning insights, such as spaced repetition, that would be difficult to create a heuristic scheduling rule for. We conclude that our results indicate that using RL for learning policies for replay scheduling is a promising direction for making replay scheduling in CL practical for real-world applications where training at test time is prohibited.

Several challenges need to be tackled to make the best out of our approach to learn replays scheduling policies. As the policy is learned with RL, the training requires careful hyperparameter tuning to succeed to learn a policy that generalizes well to the studied CL scenarios. We will therefore perform more hyperparameter searches for the DQN as well as increase the number of training environments since adding more experience for training should intuitively improve the generalization capabilities. Moreover, each step in the environment requires training the classifier on the current CL task which makes each environment step computationally expensive. Hence, we will apply RL algorithms that can be easier to tune, e.g., PPO [39] and A2C [40], and investigate whether these can be more sample-efficient than DQN for our application.

Another challenge is how to learn policies for replay scheduling in CL scenarios with long task-horizons. Continuous action spaces are a scalable option that would be appropriate for representing the task proportions since the action space would simply add one dimension per seen task. However, as learning policies for continuous actions are considered to be more difficult than learning how to take discrete actions, we could also look into various ways for discretizing continuous action spaces into finite sets of actions.

MK: Add a paragraph for limitations.

7 Conclusions

We proposed a novel method for learning policies for replay scheduling in CL which previously have been unavailable. The goal of our method is to learn a general policy that can be applied to any CL scenario for replay scheduling without requiring additional training costs at test time. We used RL for learning the policy where we employed DQNs for selecting which tasks to replay at different times. In the experiments, we showed that the replay schedules from the learned policy perform similarly to carefully tuned heuristic scheduling baselines. Furthermore, we gained more insights by visualizing the learned replay schedules and observed that the policy is capable of learning scheduling behaviors similar to human education techniques such as spaced repetition. In future work, we will investigate how alternative RL algorithms than DQN are capable of learning policies in this setting as well as experimenting with continuous action spaces to enable scaling to longer task-horizons. We hope that our work can encourage more research within this CL setting where the limitations are on compute rather than memory storage which stems well with the needs in real-world CL scenarios.

Appendix

This supplementary material is structured as follows:

- Appendix 1: Details of the experimental settings on the continual learning (CL) part, hyperparameter settings for the heuristic baselines, and hyperparameters for the Deep Q-Networks (DQNs).
- Appendix 2: Details on the construction of the action space with task proportions.
- Appendix 3: Tables over the task splits for Split FashionMNIST and Split CIFAR-10.

1 Experimental Settings

In this section, we provide more details on the experimental settings.

1.1 Hyperparameter Settings

Continual Learning Training. We use a 2-layer MLP with 256 hidden units and ReLU activation for the MNIST, FashionMNIST, and notMNIST datasets. For CIFAR10, we use the ConvNet from [41] consisting of 4 blocks with 3×3 convolution and 64 filters followed by 2×2 max pooling before the classification layer. For all architectures, we use the Adam optimizer [42] with learning rate $\eta = 0.0001$ with default hyperparameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The training details for each dataset are given below:

- MNIST: Batch size 128, 10 epochs/task
- FashionMNIST: Batch size 128, 10 epochs/task
- notMNIST: Batch size 128, 20 epochs/task
- CIFAR10: Batch size 256, 20 epochs/task

Heuristic Baselines. We performed a grid search for the threshold ratio constant τ used in the heuristic scheduling baselines Heuristic Global Drop (Heur-GD), Heuristic Local Drop (Heur-LD), and Heuristic Accuracy Threshold (Heur-AT). We used the search range $\tau \in [0.85, 0.999]$ for each baseline and selected the best parameter value based on the mean test set performance in the environments used for training for each experiment.

- Test on MNIST with new task splits, training with 10 MNIST environments
 - **Heur-GD:** Selected $\tau^{(best)} = 0.90$.
 - **Heur-LD:** Selected $\tau^{(best)} = 0.95$.
 - **Heur-AT:** Selected $\tau^{(best)} = 0.90$.
- Test on CIFAR-10 with new task splits, training with 10 CIFAR-10 environments
 - **Heur-GD:** Selected $\tau^{(best)} = 0.96$.
 - **Heur-LD:** Selected $\tau^{(best)} = 0.999$.
 - **Heur-AT:** Selected $\tau^{(best)} = 0.91$.
- Test on FashionMNIST, train on 5 MNIST + 5 notMNIST environments
 - **Heur-GD:** Selected $\tau^{(best)} = 0.99$.
 - **Heur-LD:** Selected $\tau^{(best)} = 0.95$.
 - **Heur-AT:** Selected $\tau^{(best)} = 0.93$.

Table 3: DQN hyperparameter choices used for generating the results in the experiments (Section 5). The hyperparameters in the column for Split FashionMNIST correspond to the values used when generalizing to this dataset when training on Split MNIST and Split notMNIST environments. The **bold** numbers correspond to the ones used in the experimental results from a hyperparameter search.

Hyperparameters	Split MNIST	Split CIFAR-10	Split FashionMNIST
Double DQN	True	True	True
Learning Rate	0.0001	0.0001	0.0001
Optimizer	Adam	Adam	Adam
Buffer Size	[1000 , 5000]	[1000, 5000]	[1000 , 5000]
Target Update per step	500	500	500
Batch Size	32	32	32
Discount Factor γ	1.0	1.0	1.0
Exploration Start ϵ_{start}	1	1	1
Exploration Final ϵ_{final}	[0.02, 0.05]	[0.02, 0.05]	[0.02 , 0.05]
Exploration Annealing (episodes)	1250	1250	1250
Training Episodes	5000	5000	5000

1.2 DQN Hyperparameter Choices

2 Construction of the Action Space

In this section, we describe the action space construction used in [8] that we use in our experiments. Each action $a_t = (p_1, \dots, p_{T-1}) \in \mathcal{A}_t$ consists of a sequence of task proportions p where the valid proportions sum to one, i.e., $\sum_{i=1}^t p_i = 1$, and invalid proportions are set to zero as $p_i = 0$ for $i > t$. To simplify the selection of task proportions, we create a discrete number of choices for the proportions of valid tasks. We use a binning method to construct the discrete action space. At task t , there are $t - 1$ historical tasks that can be replayed. Thus, we create $t - 1$ bins $\mathbf{b}_t = [b_1, \dots, b_{t-1}]$ and choose a task index to sample for each bin $b_i \in \{1, \dots, t - 1\}$. The bins are treated as interchangeable and only the unique choices are kept. As an example, the unique choices of vectors at task 3 are [1, 1], [1, 2], [2, 2], where [1, 1] indicates that all memory samples are from task 1, [1, 2] indicates that half memory is from task 1 and the other half are from task etc. The task proportions for each choice are then computed by counting the number of occurrences of each task index in \mathbf{b}_t and dividing by $t - 1$, such that $a_t = \text{bincount}(\mathbf{b}_t)/(t - 1)$. The task proportions are used to determine the number of samples to place in the replay memory from each task, such that $\text{round}(p_t \cdot M)$ are the number of samples from task t where we round up or down accordingly while keeping the memory size M fixed. From this specification, the whole action space over the task-horizon becomes a tree of different replay schedules that can be evaluated by the network.

3 Task Split Tables

Here, we provide tables of the task splits for Split FashionMNIST and split CIFAR-10. We provide the random seeds that we used for shuffling the label splits as well as the class names for the datasets concerned.

Table 4: Split FashionMNIST task splits with their corresponding seed.

Seed	Task 1	Task 2	Task 3	Task 4	Task 5
0	T-shirt/top, Trouser	Pullover, Dress	Coat, Sandal	Shirt, Sneaker	Bag, Ankle boot
1	Pullover, Ankle boot	Shirt, Coat	T-shirt/top, Dress	Trouser, Sneaker	Bag, Sandal
2	Coat, Trouser	Sandal, T-shirt/top	Sneaker, Pullover	Dress, Shirt	Ankle boot, Bag
3	Sandal, Coat	Trouser, Pullover	Ankle boot, Shirt	Sneaker, T-shirt/top	Dress, Bag
4	Dress, Bag	Coat, Ankle boot	Pullover, Shirt	T-shirt/top, Trouser	Sandal, Sneaker
5	Ankle boot, Sandal	Pullover, Coat	Sneaker, Trouser	T-shirt/top, Bag	Shirt, Dress
6	Bag, Trouser	Sneaker, T-shirt/top	Shirt, Sandal	Pullover, Coat	Dress, Ankle boot
7	Bag, Sandal	T-shirt/top, Pullover	Trouser, Ankle boot	Sneaker, Dress	Shirt, Coat
8	Bag, Shirt	Ankle boot, T-shirt/top	Pullover, Sandal	Sneaker, Trouser	Coat, Dress
9	Bag, Coat	Sneaker, Pullover	Trouser, Ankle boot	Dress, T-shirt/top	Shirt, Sandal

Table 5: Split CIFAR-10 task splits with their corresponding seed.

Seed	Task 1	Task 2	Task 3	Task 4	Task 5
0	Airplane, Automobile	Bird, Cat	Deer, Dog	Frog, Horse	Ship, Truck
1	Bird, Truck	Frog, Deer	Airplane, Cat	Automobile, Horse	Ship, Dog
2	Deer, Automobile	Dog, Airplane	Horse, Bird	Cat, Frog	Truck, Ship
3	Dog, Deer	Automobile, Bird	Truck, Frog	Horse, Airplane	Cat, Ship
4	Cat, Ship	Deer, Truck	Bird, Frog	Airplane, Automobile	Dog, Horse
5	Truck, Dog	Bird, Deer	Horse, Automobile	Airplane, Ship	Frog, Cat
6	Ship, Automobile	Horse, Airplane	Frog, Dog	Bird, Deer	Cat, Truck
7	Ship, Dog	Airplane, Bird	Automobile, Truck	Horse, Cat	Frog, Deer
8	Ship, Frog	Truck, Airplane	Bird, Dog	Horse, Automobile	Deer, Cat
9	Ship, Deer	Horse, Bird	Automobile, Truck	Cat, Airplane	Frog, Dog
10	Ship, Bird	Dog, Frog	Cat, Automobile	Airplane, Horse	Deer, Truck
11	Horse, Ship	Bird, Frog	Deer, Dog	Automobile, Cat	Airplane, Truck
12	Dog, Ship	Horse, Airplane	Deer, Truck	Cat, Bird	Automobile, Frog
13	Cat, Dog	Frog, Automobile	Deer, Horse	Ship, Truck	Airplane, Bird
14	Cat, Truck	Airplane, Dog	Deer, Bird	Automobile, Horse	Frog, Ship
15	Bird, Frog	Automobile, Cat	Horse, Airplane	Truck, Deer	Dog, Ship
16	Frog, Bird	Airplane, Horse	Ship, Deer	Cat, Automobile	Dog, Truck
17	Horse, Bird	Dog, Cat	Deer, Airplane	Truck, Ship	Frog, Automobile
18	Horse, Truck	Airplane, Deer	Bird, Automobile	Frog, Dog	Ship, Cat
19	Automobile, Horse	Truck, Frog	Ship, Deer	Cat, Airplane	Bird, Dog

Method	New Task Order		New Dataset
	Split MNIST	Split CIFAR-10	Split FashionMNIST
Random	3.70	3.98	3.80
ETS	4.38	3.74	4.60
Heur-Global Drop	4.16	3.84	2.56
Heur-Local Drop	2.80	3.28	4.12
Heur-Fixed	2.86	3.26	2.96
DQN	3.10	2.90	2.96

References

- [1] Frank N Dempster. Spacing effects and their implications for theory and practice. *Educational Psychology Review*, 1(4):309–330, 1989.
- [2] Hermann Ebbinghaus. Memory: A contribution to experimental psychology. *Annals of neurosciences*, 20(4):155, 2013.
- [3] T. Landauer and Robert Bjork. Optimum rehearsal patterns and name learning. *Practical aspects of memory*, 1, 11 1977.
- [4] Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [5] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- [6] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K Dokania, Philip HS Torr, and Marc'Aurelio Ranzato. On tiny episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*, 2019.
- [7] Tyler L Hayes, Kushal Kafle, Robik Shrestha, Manoj Acharya, and Christopher Kanan. Remind your neural network to prevent catastrophic forgetting. In *European Conference on Computer Vision*, pages 466–483. Springer, 2020.
- [8] Marcus Klasson, Hedvig Kjellström, and Cheng Zhang. Learn the time to learn: Replay scheduling for continual learning. *International Conference on Machine Learning (ICML) 2021 Workshop on Theory and Foundations of Continual Learning*, 2021.
- [9] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [10] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *arXiv preprint arXiv:2004.05439*, 2020.
- [11] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [12] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg

- Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [13] Richard Bellman. A markovian decision process. *Journal of Mathematics and Mechanics*, 6(5):679–684, 1957.
 - [14] Long-Ji Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine learning*, 8(3-4):293–321, 1992.
 - [15] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
 - [16] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
 - [17] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
 - [18] Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational continual learning. *arXiv preprint arXiv:1710.10628*, 2017.
 - [19] Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual learning. In *International Conference on Machine Learning*, pages 4528–4537. PMLR, 2018.
 - [20] Sayna Ebrahimi, Franziska Meier, Roberto Calandra, Trevor Darrell, and Marcus Rohrbach. Adversarial continual learning. *arXiv preprint arXiv:2003.09553*, 2020.
 - [21] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
 - [22] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically expandable networks. *arXiv preprint arXiv:1708.01547*, 2017.
 - [23] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *arXiv preprint arXiv:1706.08840*, 2017.
 - [24] Gabriel Dulac-Arnold, Richard Evans, Hado van Hasselt, Peter Sunehag, Timothy Lillicrap, Jonathan Hunt, Timothy Mann, Theophane Weber, Thomas Degris, and Ben Coppin. Deep reinforcement learning in large discrete action spaces. *arXiv preprint arXiv:1512.07679*, 2015.
 - [25] Yash Chandak, Georgios Theocharous, James Kostas, Scott Jordan, and Philip Thomas. Learning action representations for reinforcement learning. In *International conference on machine learning*, pages 941–950. PMLR, 2019.
 - [26] Ayush Jain, Andrew Szot, and Joseph J Lim. Generalization to new actions in reinforcement learning. *arXiv preprint arXiv:2011.01928*, 2020.
 - [27] Craig Boutilier, Alon Cohen, Amit Daniely, Avinatan Hassidim, Yishay Mansour, Ofer Meshi, Martin Mladenov, and Dale Schuurmans. Planning and learning with stochastic action sets. *arXiv preprint arXiv:1805.02363*, 2018.

- [28] Yash Chandak, Georgios Theocharous, Chris Nota, and Philip Thomas. Lifelong learning with a changing action set. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3373–3380, 2020.
- [29] Ji He, Jianshu Chen, Xiaodong He, Jianfeng Gao, Lihong Li, Li Deng, and Mari Ostendorf. Deep reinforcement learning with a natural language action space. *arXiv preprint arXiv:1511.04636*, 2015.
- [30] Tom Van de Wiele, David Warde-Farley, Andriy Mnih, and Volodymyr Mnih. Q-learning in enormous action spaces via amortized approximate maximization. *arXiv preprint arXiv:2001.08116*, 2020.
- [31] Luke Metz, Julian Ibarz, Navdeep Jaitly, and James Davidson. Discrete sequential prediction of continuous actions for deep rl. *arXiv preprint arXiv:1705.05035*, 2017.
- [32] Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization in reinforcement learning. In *International Conference on Machine Learning*, pages 1282–1289. PMLR, 2019.
- [33] Charles Packer, Katelyn Gao, Jernej Kos, Philipp Krähenbühl, Vladlen Koltun, and Dawn Song. Assessing generalization in deep reinforcement learning. *arXiv preprint arXiv:1810.12282*, 2018.
- [34] Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A survey of generalisation in deep reinforcement learning. *arXiv preprint arXiv:2111.09794*, 2021.
- [35] Jiachen Yang, Brenden Petersen, Hongyuan Zha, and Daniel Faissol. Single episode policy transfer in reinforcement learning. *arXiv preprint arXiv:1910.07719*, 2019.
- [36] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [37] Yaroslav Bulatov. The notMNIST dataset. <http://yaroslavvb.com/upload/notMNIST/>, 2011.
- [38] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [39] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [40] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.
- [41] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. *arXiv preprint arXiv:1606.04080*, 2016.
- [42] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

