
Random Representations Outperform Online Continually Learned Representations

Ameya Prabhu^{1*} Shiven Sinha^{2*} Ponnurangam Kumaraguru² Philip H.S. Torr¹
Ozan Sener³⁺ Puneet K. Dokania¹⁺
¹University of Oxford ²IIT Hyderabad ³Apple

Abstract

Continual learning has primarily focused on the issue of catastrophic forgetting and the associated stability-plasticity tradeoffs. However, little attention has been paid to the efficacy of continually learned representations, as representations are learned alongside classifiers throughout the learning process. Our primary contribution is empirically demonstrating that existing online continually trained deep networks produce inferior representations compared to a simple pre-defined random transforms. Our approach embeds raw pixels using a fixed random transform, approximating an RBF-Kernel initialized before any data is seen. We then train a simple linear classifier on top without storing any exemplars, processing one sample at a time in an online continual learning setting. This method, called RanDumb, significantly outperforms state-of-the-art continually learned representations across all standard online continual learning benchmarks. Our study reveals the significant limitations of representation learning, particularly in low-exemplar and online continual learning scenarios. Extending our investigation to popular exemplar-free scenarios with pretrained models, we find that training only a linear classifier on top of pretrained representations surpasses most continual fine-tuning and prompt-tuning strategies. Overall, our investigation challenges the prevailing assumptions about effective representation learning in online continual learning. Our code is available [here](#).

1 Introduction

Continual learning aims to develop models capable of learning from non-stationary data streams, inspired by the lifelong learning abilities exhibited by humans and the prevalence of such real-world applications (see Verwimp et al. [64] for a survey). It is characterized by sequentially arriving tasks, coupled with additional computational and memory constraints [33, 38, 54, 62, 49].

Building on the foundations of supervised deep learning, the prevalent approach in continual learning has been to jointly train representations alongside classifiers. This approach simply follows from the assumption that learned representations are expected to outperform fixed representation functions such as kernel classifiers, as demonstrated in supervised deep learning [34, 23, 57]. However, this assumption is never validated in continual learning, with scenarios having limited updates where networks might not be trained until convergence, such as online continual learning (OCL).

In this paper, we study the efficacy of representations derived from continual learning algorithms. Surprisingly, our findings suggest that these representations might not be as beneficial as presumed. To test this, we introduce a simple baseline method named RanDumb, which combines a random representation function with a straightforward linear classifier, illustrated in detail in Figure 1 (left). Our empirical evaluations, summarized in Table 1 (left, top), reveal that despite replacing

*authors contributed equally, + equal advising

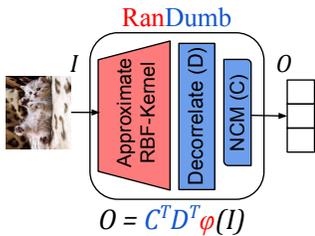


Figure 1: RanDumb projects raw pixels to a high dimensional space using random Fourier projections (φ), then decorrelates the features using Mahalanobis distance [43] and classifies with the nearest class mean. The online update only involves updating a single sample covariance matrix and class-means.

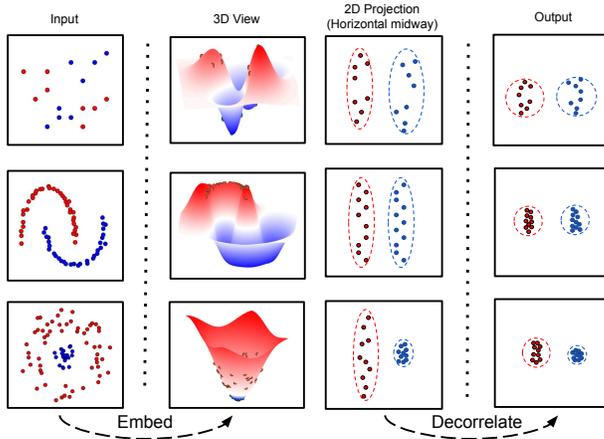


Figure 2: RanDumb projects the datapoints to a high-dimensional space to create a clearer separation between classes. Subsequently, it corrects the anisotropy across feature dimensions, scaling them to be unit variance each. This allows cosine similarity to accurately separate classes. The figure is adapted from [48].

Table 1: **(Left) Online Continual Learning.** Performance comparison of RanDumb on the PEC setup [75] and VAE-GC [63]. Setup and numbers borrowed from PEC [75]. RanDumb outperforms the best OCL method. **(Right) Offline Continual Learning.** Performance comparison with ImageNet21K ViT-B16 model using 2 initial classes and 1 new class per task. RanPAC-imp is an improved version of the RanPAC code which mitigates the instability issues in RanPAC. RanDumb nearly matches performance of joint for both online and offline, demonstrating the inefficacy of current benchmarks.

Method	MNIST	CIFAR10	CIFAR100	m-IMN	Method	CIFAR	IN-A	IN-R	CUB	OB	VTAB	Cars
Comparison with Best Method					Comparison with Best Method							
Best (PEC)	92.3	58.9	26.5	14.9	Best (RanPAC-imp)	89.4	33.8	69.4	89.6	75.3	91.9	57.3
RanDumb (Ours)	98.3	55.6	28.6	17.7	RanDumb (Ours)	86.8	42.2	64.9	88.5	75.3	92.4	67.1
Improvement	+6.0	-3.3	+2.1	+2.8	Improvement	-2.6	+8.4	-4.5	-1.1	+0.0	+0.5	+9.8
Random vs. Learned Representations					Random vs. Finetuned Representations							
VAE-GC	84.0	42.7	19.7	12.1	SLCA	86.8	-	54.2	82.1	-	-	18.2
RanDumb (Ours)	98.3	55.6	28.6	17.7	RanDumb (Ours)	86.8	42.2	64.9	88.5	75.3	92.4	67.1
Improvement	+14.3	+12.9	+8.9	+5.6	Improvement	+0.0	-	+10.7	+6.4	-	-	+48.9
Scope of Improvement					Scope of Improvement							
Joint (One Pass)	98.3	74.2	33.0	25.3	Joint	93.8	70.8	86.6	91.1	83.8	95.5	86.9
RanDumb (Ours)	98.3	55.6	28.6	17.7	RanDumb (Ours)	86.8	42.2	64.9	88.5	75.3	92.4	67.1
Gap Covered. (%)	100%	75%	87%	70%	Gap Covered. (%)	93%	60%	75%	97%	92%	97%	77%

the representation learning with a pre-defined random representation, RanDumb surpasses current state-of-the-art methods in latest online continual learning benchmarks [75].

We further expand our evaluations to scenarios incorporating methods that use pre-trained feature extractors [67]. By substituting our random projections with these feature extractors and retaining the linear classifier, RanDumb again outperforms leading methods as shown in Table 1 (right, top).

1.1 Technical Summary: Construction of RanDumb and Empirical Findings

Design. RanDumb first projects input pixels into a high-dimensional space using a fixed kernel based on random Fourier basis, which is a low-rank data-independent approximation of the RBF Kernel [52]. Then, we use a simple linear classifier which first normalizes distances across different feature dimensions (anisotropy) with Mahalanobis distance [43] and then uses nearest class means for classification [44]. In scenarios with pretrained feature extractors, we use the fixed pretrained model as embedder and learn a linear classifier as described above, similar to Hayes and Kanan [27].

Key Properties. RanDumb needs no storage of exemplars and requires only one pass over the data in a one-sample-per-timestep fashion. Furthermore, it only requires online estimation of the sample covariance matrix and nearest class mean.

Key Finding 1: Poor Representation Learning. We compare RanDumb with leading methods: VAE-GC [63] in Table 1 (left, middle) and SLCA [78] in Table 1 (right, middle). The primary distinction between them is their representation: RanDumb uses a fixed function (random/pretrained network), whereas VAE-GC and SLCA further continually trained deep networks. RanDumb consistently surpasses VAE-GC and SLCA by wide margins of 5-15%. This shows that state-of-the-art online continual learning algorithms fail to learn effective representations across standard exemplar-free continual learning benchmarks.

Finding 2: Over-Constrained Benchmarks. Given the demonstrated limitations of existing continual representation learning methods, an important question arises: Can better methods learn more effective representations? To explore this, we evaluated the performance of RanDumb against joint training, models trained without continual learning constraints, in both online and offline settings, as shown in Table 1 (left, bottom) and Table 1 (right, bottom). Our straightforward baseline, RanDumb, bridges 70-90% of the performance gap relative to the respective joint classifiers in both scenarios. This significant recovery of performance by such a simple method suggests that if our goal is to advance the study of representation learning, current benchmarks may be overly restrictive and not conducive to truly effective representation learning.

We highlight that the goal in our work is not to introduce a state-of-the-art continual learning method, but challenge prevailing assumptions and open a discussion on the efficacy of representation learning in continual learning algorithms, especially in online and low-exemplar scenarios.

2 RanDumb: Mechanism & Intuitions

RanDumb has two main elements: random projection and the dumb learner. We illustrate the mechanism of RanDumb using three toy examples in Figure 1 (right). To classify a test sample \mathbf{x}_{test} , we start with a simple classifier, the nearest class mean (NCM). It predicts the class among C classes by highest value of the similarity function f among class means μ_i :

$$y_{\text{pred}} = \arg \max_{i \in \{1, \dots, |C|\}} f(\mathbf{x}_{\text{test}}, \mu_i), \quad \text{where} \quad f(\mathbf{x}_{\text{test}}, \mu_i) := \mathbf{x}_{\text{test}}^\top \mu_i \quad (1)$$

and μ_i are the class-means in the pixel space: $\mu_i = \frac{1}{|C_i|} \sum_{\mathbf{x} \in C_i} \mathbf{x}$. RanDumb adds two additional components to this classifier: 1) Kernelization and 2) Decorrelation.

Kernelization: Classes are typically not linearly separable in the pixel space, unlike in the feature space of deep models. Hence, we apply the kernel trick to embed the pixels in a better representation space, computing all distances between the data and class-means in this embedding space. This phenomena is illustrated on three toy examples to build intuitions in Figure 1 (right, Embed). We use an RBF-Kernel, which for two points \mathbf{x} and \mathbf{y} is defined as: $K_{\text{RBF}}(\mathbf{x}, \mathbf{y}) = \exp(-\gamma \|\mathbf{x} - \mathbf{y}\|^2)$ where γ is a scaling parameter. However, calculating the RBF kernel is not possible due to the online continual learning constraints preventing computation of pairwise-distance between all points. Hence, we use a data-independent approximation, random Fourier projection $\phi(\mathbf{x})$, as given in [52]:

$$K_{\text{RBF}}(\mathbf{x}, \mathbf{y}) \approx \phi(\mathbf{x})^\top \phi(\mathbf{y})$$

where the random Fourier features $\phi(\mathbf{x})$ are defined by first sampling D vectors $\{\omega_1, \dots, \omega_D\}$ from a Gaussian distribution with mean zero and covariance matrix $2\gamma\mathbf{I}$, where \mathbf{I} is the identity matrix. Then $\phi(\mathbf{x})$ is a $2D$ -dimensional feature, defined as:

$$\phi(\mathbf{x}) = \frac{1}{\sqrt{D}} [\cos(\omega_1^\top \mathbf{x}), \sin(\omega_1^\top \mathbf{x}), \dots, \cos(\omega_D^\top \mathbf{x}), \sin(\omega_D^\top \mathbf{x})]$$

We keep these ω bases fixed throughout online learning. Thus, we obtain our modified similarity function from Equation 1 as:

$$f(\mathbf{x}_{\text{test}}, \mu_i) := \phi(\mathbf{x}_{\text{test}})^\top \bar{\mu}_i \quad (2)$$

where $\bar{\mu}_i$ are the class-means in the kernel space:

$$\bar{\mu}_i = \frac{1}{|C_i|} \sum_{\mathbf{x} \in C_i} \phi(\mathbf{x})$$

Decorrelation: Projected raw pixels have feature dimensions with different variances (anisotropic). Hence, instead of naively computing $\phi(\mathbf{x}_{\text{test}})^\top \bar{\mu}_i$, we further decorrelate the feature dimensions using a Mahalanobis distance with the shrunked covariance matrix \mathbf{S} using OAS shrinkage [15], inverse obtained by least squares minimization $(\mathbf{S} + \lambda \mathbf{I})$. We illustrate this phenomena as well on three toy examples in Figure 1 (right, Decorrelate) to build intuitions. Our similarity function finally is:

$$f(\mathbf{x}_{\text{test}}, \mu_i) := (\phi(\mathbf{x}_{\text{test}}) - \bar{\mu}_i)^\top \mathbf{S}^{-1} (\phi(\mathbf{x}_{\text{test}}) - \bar{\mu}_i) \quad (3)$$

Online Computation. Our random projection is fixed before seeing any data. During continual learning, we only perform online update on the running class mean and empirical covariance matrix².

3 Experiments

We compare RanDumb with algorithms across online continual learning benchmarks with an emphasis on exemplar-free and low-exemplar storage regime.

Benchmarks. The benchmarks which we used in our experiments are summarized in Table on the right. We aim for a comprehensive coverage and show results on four standard online continual learning benchmarks (A, B, D, E) which reflect the latest trends ('22-'24) across exemplar-free, contrastive-training³, meta-continual learning, and network-expansion based approaches respectively.

Setup	Num Passes	#Classes Per Task	#Samples Per Step	#Stored Exemplars	Contrastive Augment
Method: RanDumb	1	1	1	0	No
A (Zajac et al. [75])	1	1	10	0	No
B1 (Guo et al. [25])	1	2	10	100-2000	No
B2 (Guo et al. [25])	1	2	10	100-1000	Yes
C (Smith et al. [60])	Many	10	All	0	No
D (Wu et al. [70])	1	2-10	10	1000	No
E (Ye and Bors [74])	1	2-5	10	1000-5000	No
F (Wang et al. [67], modified)	Many	1	All	0	No

We also evaluate on a rehearsal-free offline continual learning benchmark C. These benchmarks are ordered by increasingly relaxed constraints, moving further away from the training scenario of RanDumb. Benchmark A closely matches RanDumb with one class per timestep and no stored exemplars. Benchmark B, D, E progressively relax the constraints on exemplars and classes per timestep. Benchmark C and E remove the online constraint by allowing unrestricted training and sample access within a task without exemplar-storage of past tasks. Benchmark F allows using large pretrained models, modified by us with one class per task, i.e. testing learning over longer timespans.

We further test on exemplar-free scenarios in offline continual learning using Benchmark F [67] with the challenging one-class per task constraint borrowed from [75]. This benchmark allows using pretrained models along with unrestricted training time and access to all class samples at each timestep. However, RanDumb is restricted to learning from a single pass seeing only one sample at a time. RanDumb only learns a linear classifier over a given pretrained model in Benchmark F.

We use LAMDA-PILOT [61] codebase for all methods, except RanPAC and SLDA for which use their codebases. We use the original hyperparameters. We only change initial classes to 2 and number of classes per task to 1 and test using both ImageNet21K and ImageNet1K ViT-B/16 models.

Implementation Details (RanDumb). We evaluate RanDumb using five datasets: MNIST, CIFAR10, CIFAR100, TinyImageNet200, and miniImageNet100. For the latter two, we downscale all images to 32x32. We augment each datapoint with flipped version, hence two images are seen by the classifier at each timestep (except for MNIST and Benchmark F). We normalize all images and flatten them into vectors, obtaining 784-dim input vectors for MNIST and 3072-dim input vectors for all the other. For Benchmark F, we compare RanDumb on seven datasets used in LAMDA-PILOT, replacing ObjectNet with Stanford Cars as ObjectNet license prohibits training models. We use the 768-dimensional features from the same pretrained ViT-B models used in this benchmark. We measure accuracy on the test set of all past seen classes after completing the full one-pass. We take the average accuracy after the last task on all past tasks [75, 25, 67]. In Benchmark A and F, since we have one class per task, the average accuracy across past tasks is the same regardless of the task ordering. In Benchmarks A-E, all datasets have the same number of samples, hence similarly the average accuracy across past tasks is the same regardless of the task ordering. We used the Scikit-Learn implementation of Random

²Online update for the inverse of the covariance matrix is possible using the Sherman–Morrison formula.

³Benchmark B is split into two sections: (B1) methods that do not rely on contrastive learning and heavy augmentation, and (B2) approaches that incorporate contrastive learning and extra augmentations.

Table 2: **Benchmark A** (Ref: Table 1 from PEC [75]). We compare RanDumb in a 1-class per task setting referred as ‘Dataset (num_tasks/1)’. We observe that RanDumb outperforms all approaches across all datasets by 2-6% margins, with an exception of latest work PEC [75] on CIFAR10.

	Method	Memory	MNIST (10/1)	CIFAR-10 (10/1)	CIFAR-100 (100/1)	miniImageNet (100/1)
	<i>Fine-tuning</i>	all	10.1± 0.0	10.0± 0.0	1.0± 0.0	1.0± 0.0
	<i>Joint, 1 epoch</i>	all	98.3± 0.0	74.2± 0.1	33.0± 0.2	25.3± 0.2
<i>Rehearsal Based Methods</i>	ER [13]	500	84.4± 0.3	40.6± 1.1	12.5± 0.3	5.7± 0.2
	A-GEM [12]	500	59.8± 0.8	10.2± 0.1	1.0± 0.0	1.1± 0.1
	iCaRL [54]	500	83.1± 0.3	37.8± 0.4	5.7± 0.1	7.5± 0.1
	BiC [71]	500	86.0± 0.4	35.9± 0.4	6.4± 0.3	1.5± 0.1
	ER-ACE [10]	500	87.8± 0.2	39.9± 0.5	8.2± 0.2	5.7± 0.2
	DER [9]	500	91.7± 0.1	40.0± 1.5	1.0± 0.1	1.0± 0.0
	DER++ [9]	500	91.9± 0.2	35.6± 2.4	6.2± 0.4	1.4± 0.1
	X-DER [8]	500	83.0± 0.1	43.2± 0.5	15.6± 0.1	8.2± 0.4
	GDumb [49]	500	91.0± 0.2	50.7± 0.7	8.2± 0.2	-
<i>Rehearsal Free Methods</i>	EWC [33]	0	10.1± 0.0	10.6± 0.4	1.0± 0.0	1.0± 0.0
	SI [76]	0	12.7± 1.0	10.1± 0.1	1.1± 0.0	1.0± 0.1
	LwF [37]	0	11.8 ± 0.6	10.1± 0.1	0.9± 0.0	1.0± 0.0
	LT [77]	0	10.9± 0.9	10.0± 0.2	1.1± 0.1	1.0± 0.0
	Gen-NCM [31]	0	82.0± 0.0	27.7± 0.0	10.0± 0.0	7.5± 0.0
	Gen-SLDA [27]	0	88.0± 0.0	41.4± 0.0	18.8± 0.0	12.9± 0.0
	VAE-GC [63]	0	84.0± 0.5	42.7± 1.3	19.7± 0.1	12.1± 0.1
	PEC [75]	0	92.3± 0.1	58.9± 0.1	26.5± 0.1	14.9± 0.1
	RanDumb (Ours)	0	98.3 (+5.9)	55.6 (-3.3)	28.6 (+2.1)	17.7 (+2.8)

Table 3: **Benchmark B.1** (Ref: Table adopted from OnPro [68], OCM[25]) We compare RanDumb in many-classes per task setting referred as ‘Dataset (num_tasks/num_classes_per_task)’. We categorize memory buffer sizes with ‘M’. RanDumb outperforms the competing approaches without heavy-augmentations by 3-20% margins despite being exemplar free. Only in one case, it is second best.

Method	MNIST (5/2)		CIFAR10 (5/2)		CIFAR100 (10/10)		CIFAR100 (50/2)		TinyImageNet (100/2)	
	M = 0.1k	M = 0.1k	M = 0.2k	M = 0.5k	M = 1k	M = 1k	M = 1k	M = 1k	M = 2k	M = 2k
AGEM [12]	56.9±5.2	17.7±0.3	22.7±1.8	5.8±0.2	5.9±0.1	1.8±0.2	0.8±0.1	0.9±0.1		
GSS [4]	70.4±1.5	18.4±0.2	26.9±1.2	8.1±0.2	11.1±0.2	4.3±0.2	1.1±0.1	3.3±0.5		
ER [13]	78.7±0.4	19.4±0.6	29.7±1.0	8.7±0.3	15.7±0.3	8.3±0.3	1.2±0.1	5.6±0.5		
ASER [58]	61.6±2.1	20.0±1.0	27.8±1.0	11.0±0.3	16.4±0.3	9.6±1.3	2.2±0.1	5.3±0.3		
MIR [3]	79.0±0.5	20.7±0.7	37.3±0.3	9.7±0.3	15.7±0.2	12.7±0.3	1.4±0.1	6.1±0.5		
ER-AML [10]	76.5±0.1	-	40.5±0.7	-	16.1±0.4	-	-	5.4±0.2		
iCaRL [54]	-	31.0±1.2	33.9±0.9	12.8±0.4	16.5±0.4	-	5.0±0.3	6.6±0.4		
DER++ [9]	74.4±1.1	31.5±2.9	44.2±1.1	16.0±0.6	21.4±0.9	9.3±0.3	3.7±0.4	5.1±0.8		
GDumb [49]	81.2±0.5	23.3±1.3	35.9±1.1	8.2±0.2	18.1±0.3	18.1±0.3	4.6±0.3	12.6±0.1		
CoPE [19]	-	33.5±3.2	37.3±2.2	11.6±0.4	14.6±1.3	-	2.1±0.3	2.3±0.4		
DVC [25]	-	35.2±1.7	41.6±2.7	15.4±0.3	20.3±1.0	-	4.9±0.6	7.5±0.5		
Co ² L [11]	83.1±0.1	-	42.1±1.2	-	17.1±0.4	-	-	10.1±0.2		
R-RT [6]	89.1±0.3	-	45.2±0.4	-	15.4±0.3	-	-	6.6±0.3		
CCIL [46]	86.4±0.1	-	50.5±0.2	-	18.5±0.3	-	-	5.6±0.9		
IL2A [81]	90.2±0.1	-	54.7±0.5	-	18.2±1.2	-	-	5.5±0.7		
BiC [71]	90.4±0.1	-	48.2±0.7	-	21.2±0.3	-	-	10.2±0.9		
SSIL [1]	88.2±0.1	-	49.5±0.2	-	26.0±0.1	-	-	9.6±0.7		
<i>Rehearsal-Free</i>										
PASS [81]	-	33.7±2.2	33.7±2.2	7.5±0.7	7.5±0.7	-	0.5±0.1	0.5±0.1		
RanDumb (Ours)	98.3 (+7.8)	55.6 (+20.4)	55.6 (+5.9)	28.6 (+12.6)	28.6 (+2.6)	28.6 (+10.5)	11.6 (+6.6)	11.6 (-1.0)		

Fourier Features [52] with 25K embedding size, $\gamma = 1.0$. We use progressively increasing ridge regression parameter (λ) with dataset complexity, $\lambda = 10^{-6}$ for MNIST, $\lambda = 10^{-5}$ for CIFAR10/100 and $\lambda = 10^{-4}$ for TinyImageNet200/miniImageNet100.

3.1 Results

Benchmark A. We assess continual learning models in the challenging setup of one class per timestep, closely mirroring our training assumptions, and present our results in Table 2. Comparing across rows, and see that RanDumb improves over prior state-of-the-art across all datasets with 2-6% margins.

Benchmark C. We compare against offline rehearsal-free continual learning approaches in Table 4 (Right) on CIFAR100. Despite online training, RanDumb outperforms PredKD by over 4% margins.

Benchmark D. We compare performance of RanDumb against meta-continual learning methods, which require large exemplars with buffer sizes of 1K in Table 5 (left). RanDumb achieves strong performance under these conditions, exceeding all prior work by a large margin of 9.1% on CIFAR100 and outperforms all but VR-MCL approach on the TinyImageNet dataset. GDumb performs the best on CIFAR10, indicating this is already in a large-exemplar regime uniquely unsuited for RanDumb.

Benchmark E. We compare RanDumb against network expansion-based online continual learning methods in Table 5 (right). These approaches grow model capacity to mitigate forgetting while dealing with shifts in the data distribution, and are allowed larger memory buffers. RanDumb matches the performance of the state-of-the-art method SEDEM [74] on MNIST, while exceeding it by 0.3% on CIFAR10 and 3.8% on CIFAR100.

Benchmark F. We compare performance of approaches which do not further train the deep network like RanDumb against popular continual finetuning and prompt-tuning approaches in Table 6. We discover that prompt-tuning approaches completely collapse under large timesteps and approaches which do not finetune their pretrained model achieve strong performance, even under challenging one class per timestep constraint. Note that RanPAC [41] adds a RP+ReLU and finetunes in a first-session adaptation fashion over RanDumb, yet fails to achieve higher accuracies.

Overall, despite RanDumb being exemplar-free, it outperforms nearly all online continual learning methods across various tasks when exemplar storage is limited. We specifically benchmark on lower exemplar sizes to complement settings in which GDumb does not perform well.

3.2 Analysis of RanDumb

Ablating Components of RanDumb. We ablate the contribution of only using Random Fourier features for embedding and decorrelation to the overall performance of RanDumb in Table 7 (left, top). Ablating the decorrelation and relying solely on random Fourier features, colloquially dubbed Kernel-NCM, has performance drops ranging from 6-25% across the datasets. Replacing random Fourier features with raw features, *ie.* the SLDA baseline, leads to pronounced drop in performance ranging from 3-14% across the datasets. Moreover, ablating both components results in the base nearest class mean classifier, and exhibits the poorest performance with an average reduction of 17%. Therefore, both decorrelation and random embedding are crucial for RanDumb.

Impact of Embedding Dimensions. We vary the dimensions of the random Fourier features ranging from compressing 3K input dimensions to 1K to projecting it to 25K dimensions and evaluate its impact on performance in Figure 3. Surprisingly, the random projection to a 3x compressed 1K dimensional space allows for significant performance improvement over not using embedding, given in Table 7 (left, top). Furthermore, increasing the dimension from 1K to 25K results in improvements of 3.6%, 10.4%, 7.0%, and 2.5% on MNIST, CIFAR10, CIFAR100, and TinyImageNet respectively.

Table 6: **Benchmark F** We compare RanDumb with prompt-tuning approaches using ViT-B/16 ImageNet-21K/1K pretrained models using 2 init classes and 1 class per task setting. Most prompt-tuning based methods collapse and RanDumb achieves either state-of-the-art or second-best performance. RanPAC-imp is an improved version of the RanPAC mitigating the instability issues.

Method	CIFAR	IN-A	IN-R	CUB	VTAB
ViT-B/16 (IN-1K Pretrained)					
Finetune	1.0	1.2	1.1	1.0	2.1
L2P [67]	2.4	0.3	0.8	1.4	1.3
DualPrompt [66]	2.3	0.3	0.8	0.9	4.2
CODA-Prompt [59]	2.6	0.3	0.8	1.9	6.3
Adam-Adapt [80])	76.7	49.3	62.0	85.2	83.6
Adam-SSF [80]	76.0	47.3	64.2	85.6	84.2
Adam-VPT [80]	79.3	35.8	61.2	83.8	86.9
Adam-FT [80]	72.6	49.3	61.0	85.2	83.8
Memo [79]	69.8	-	-	81.4	-
iCARL [54]	72.4	-	35.2	72.4	-
Foster [65]	52.2	-	76.8	86.6	-
NCM [31]	78.3	44.3	62.5	84.8	88.2
SLCA [78]	86.3	-	52.8	84.7	-
RanPAC [41]	88.2	39.0	72.8	77.7	93.0
RanPAC-imp [41]	87.8	43.5	72.6	89.6	93.0
RanDumb (Ours)	84.5	49.5	66.9	88.0	93.6
ViT-B/16 (IN-21K Pretrained)					
Finetune	2.8	0.5	1.2	1.2	0.5
Adam-Adapt [80]	82.4	48.8	55.4	86.7	84.4
Adam-SSF [80]	82.7	46.0	59.7	86.2	84.9
Adam-VPT [80]	70.8	34.8	53.9	84.0	81.1
Adam-FT [80]	65.7	48.5	56.1	86.5	84.4
Foster [65]	87.3	-	5.1	86.9	-
iCARL [54]	71.6	-	35.1	71.6	-
NCM [31]	83.5	41.4	54.8	86.5	88.5
SLCA [78]	86.8	-	54.2	82.1	-
RanPAC [41]	89.6	26.8	67.3	87.2	88.2
RanPAC-imp [41]	89.4	33.8	69.4	89.6	91.9
RanDumb (Ours)	86.8	42.2	64.9	88.5	92.4

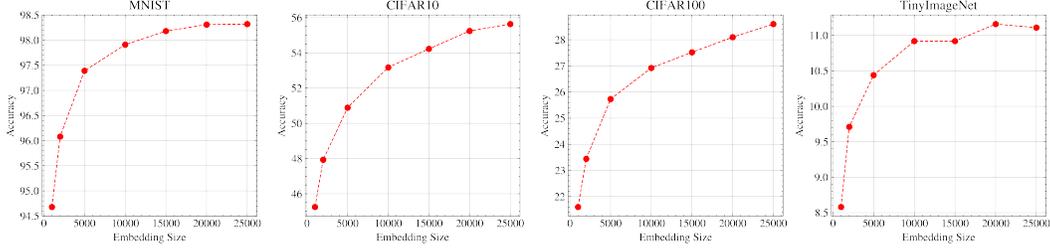


Figure 3: Accuracy of RanDumb with respect to embedding dimensionality across datasets.

Table 7: **(Left) Analysis of RanDumb:** We study contributions of decorrelation, random embedding, and data augmentation. We further vary the embedding sizes and regularisation parameter. Finally, we compare with alternate embeddings. **(Right) Architectures** (Ref: Table 1 from Mirzadeh et al. [45]) RanDumb surpasses continual representation learning across a wide range of architectures, achieving close to 94% of the joint performance.

Method	MNIST (10/1)	CIFAR10 (10/1)	CIFAR100 (10/1)	T-ImNet (200/1)	m-ImNet (100/1)	Model	CIFAR100
Ablating Components of RanDumb						Joint	79.58
RanDumb	98.3	55.6	28.6	11.1	17.7	CNN x1	62.2 ± 1.35
-Decorrelate	83.8 (-14.5)	30.0 (-25.6)	12.0 (-16.6)	4.7 (-6.4)	8.9 (-8.8)	CNN x2	66.3 ± 1.12
-Embed	88.0 (-10.3)	41.6 (-14.0)	19.0 (-9.6)	8.0 (-3.1)	12.9 (-4.8)	CNN x4	68.1 ± 0.5
-Both	82.1 (-16.2)	28.5 (-27.1)	10.4 (-18.2)	4.1 (-7.0)	7.28 (-10.4)	CNN x8	69.9 ± 0.62
Effect of Adding Flip Augmentation						CNN x16	76.8 ± 0.76
With	-	55.6	28.6	11.1	17.7	ResNet-18	45.0 ± 0.63
Without	98.3	52.5 (-3.1)	26.9 (-1.7)	10.7 (-0.4)	16.6 (-1.1)	ResNet-34	44.8 ± 2.34
Variation with Ridge Parameter λ						ResNet-50	56.2 ± 0.88
$\lambda = 10^{-6}$	98.3	53.9	27.8	10.3	15.8	ResNet-101	56.8 ± 1.62
$\lambda = 10^{-5}$	-	55.6	28.6	11.1	15.9	WRN-10-2	50.5 ± 2.65
$\lambda = 10^{-4}$	96.6	52.6	26.1	11.6	17.7	WRN-10-10	56.8 ± 2.03
Variation Across Embedding Projections						WRN-16-2	44.6 ± 2.81
No-Embed	88.0	41.6	19.0	8.0	12.9	WRN-16-10	51.3 ± 1.47
RP+ReLU (RanPAC)	95.2	48.8	23.1	9.7	15.7	WRN-28-2	46.6 ± 2.27
RanDumb (Ours)	98.3 (+3.1)	55.6 (+6.8)	28.6 (+5.5)	11.1 (+1.4)	17.7 (+2.0)	WRN-28-10	49.3 ± 2.02
						ViT-512/1024	51.7 ± 1.4
						ViT-1024/1546	60.4 ± 1.56
						RanDumb (Ours)	74.8 (-2.0)

Increasing the embedding sizes beyond 15K, however, only results in modest improvements of 0.1%, 1.4%, 1.1% and 0.2% on the same datasets, indicating 15K dimensions would be a good point for a performance-computational cost tradeoff.

Impact of Flip Augmentation. We evaluate the impact of adding the flip augmentation on the performance of RanDumb in Table 7 (left, middle). Note that MNIST was not augmented. Augmentation provided large gains of 3.1% on CIFAR10, 1.7% on CIFAR100, and 0.4% on TinyImageNet. We did not augment the data further with RandomCrop transform as done with standard augmentations.

Impact of Varying Ridge Parameter. All prior experiments use a ridge parameter (λ) that increases with dataset complexity: $\lambda = 10^{-6}$ for MNIST, 10^{-5} for CIFAR10 and CIFAR100, and 10^{-4} for TinyImageNet and miniImageNet. Table 7 (left, middle) shows the effect of varying λ on RanDumb’s performance. With a smaller $\lambda = 10^{-6}$, CIFAR10, CIFAR100, TinyImageNet and miniImageNet all exhibit minor drops of 0.1%-1.7%, 0.8%, 0.8%. Increasing shrinkage to a $\lambda = 10^{-4}$ reduces CIFAR10 and CIFAR100 performance more substantially by 3% and 2.5% versus their optimal $\lambda = 10^{-5}$. On the other hand, this larger λ leads to improvements of 0.5% and 1.8% on TinyImageNet and miniImageNet. This aligns with the trend that datasets with greater complexity benefit from more regularisation, with the optimal λ balancing under- and over-regularisation effects.

Comparison with Extreme Learning Machines. We compared our random Fourier features with random projections based extreme learning machines, as recently adapted to continual learning by RP+ReLU [41] in Table 7 (left, bottom) with their best embedding size. Our method performs significantly better on each dataset, averaging a gain of 3.4%.

Comparisons across Architectures. In table 7 (right), we compare whether using random Fourier features as embeddings outperforms models across various architectures for continual representation learning. We use experience replay (ER) baseline in the task-incremental CIFAR100 setup (for details, see Mirzadeh et al. [45] as it differs significantly from earlier setups). Our comparison spanned various architectures. The findings revealed that RanDumb surpassed the performance of nearly all considered architectures, and achieved close to 94% of the joint multi-task performance. This suggests that RanDumb outperforms continual representation learning across architectures.

Conclusion. Overall, both random embedding and decorrelation are critical components in the performance of RanDumb. Using random Fourier features is substantially better than RanPAC. Lastly, one can substantially reduce the embedding dimension without a large drop in performance for large gains in computational cost, additional augmentation may further significantly help performance and optimal shrinkage parameter increases with dataset complexity. RanDumb outperforms continual representation learning across a wide range of architectures.

4 Related Works and Equivalent Formulations

Random Representations. There have been extensive theoretical and empirical investigations into random representations in machine learning, compressed sensing, and other fields, often utilizing extreme learning machines [56, 14, 21] (see [30, 29] for a survey). Other investigations include efficient kernel methods using Fourier features and Nyström approximations [52, 69], and extensions to efficiently parameterize linear classifiers [2]. They are also embedded into deep networks [17, 35, 72, 16]. We tailored the already successful random fourier representations [52] to the problem at hand and applied to the online continual learning problem for the first time.

Continual Representation Learning. There are various works focusing on continual representation learning itself [53, 20, 39, 28], but they address the problem of alleviating the stability-plasticity dilemma in high-exemplar and offline continual learning scenarios where models are trained until convergence. In comparison, we focus on online and low-exemplar regime.

Representation Learning Free Methods in CL. Several works have developed the idea of using fixed pretrained networks after adapting on the first task across various settings [50, 41, 24]. Our work contributes to this growing evidence, however, we do not perform first-task adaptation [47], and propose OAS-shrunk SLDA as structurally simplest but highly accurate continual linear classifier without any extra bells-and-whistles. Moreover, we are the first work to introduce a representation learning free method with random features for continually learning from scratch.

Equivalent formulations to RanDumb. If the classes are equiprobable, which is the case for most datasets here, nearest class mean classifier with the Mahalanobis distance metric is equivalent to linear discriminant analysis (LDA) classifier [42]. Hence, one could say RanDumb is exactly equivalent to a Streaming LDA classifier with an approximate RBF Kernel. Alternatively, one could think of the decorrelation operation as explicitly decorrelating the features with ZCA whitening [7].

5 Discussion and Concluding Remarks

Our investigation reveals a surprising result — simply using random embedding (RanDumb) consistently outperforms learned representations from methods specifically designed for online continual training. Furthermore, using random/pretrained features also recovers 70-90% of the gap to joint learning, leaving limited room for improvement in representation learning techniques on standard benchmarks. Overall, our investigation questions our understanding of how to effectively design and train models that require efficient continual representation learning, and necessitates a re-investigation of the widely explored problem formulation itself. We believe adoption of computationally bounded scenarios without memory constraints and corresponding benchmarks [51, 50, 22] could be a promising way forward.

Limitations & Future Directions. We currently do not provide theory or justification for why training dynamics of continual learning algorithms fails to effectively learn good representations; doing so would provide deeper insights into continual learning algorithms. Moreover, our proposed method, RanDumb with random Fourier features is limited in scope towards low-exemplar scenarios and online-continual learning. Extending studies on representation learning to high-exemplar and offline continual learning scenarios might be exciting directions to investigate.

References

- [1] Hongjoon Ahn, Jihwan Kwak, Subin Lim, Hyeonsu Bang, Hyojun Kim, and Taesup Moon. Ss-il: Separated softmax for incremental learning. In *ICCV*, 2021.
- [2] Nir Ailon and Bernard Chazelle. The fast johnson–lindenstrauss transform and approximate nearest neighbors. *SIAM Journal on computing*, 2009.
- [3] Rahaf Aljundi, Lucas Caccia, Eugene Belilovsky, Massimo Caccia, Laurent Charlin, and Tinne Tuytelaars. Online continual learning with maximally interfered retrieval. In *NeurIPS*, 2019.
- [4] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. In *NeurIPS*, 2019.
- [5] Elahe Arani, Fahad Sarfraz, and Bahram Zonooz. Learning fast, learning slow: A general continual learning method based on complementary learning system. In *ICLR*, 2022.
- [6] Jihwan Bang, Heesu Kim, YoungJoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow memory: Continual learning with a memory of diverse samples. In *CVPR*, 2021.
- [7] Anthony Bell and Terrence J Sejnowski. Edges are the ‘independent components’ of natural scenes. In *NeurIPS*, 1996.
- [8] Matteo Boschini, Lorenzo Bonicelli, Pietro Buzzega, Angelo Porrello, and Simone Calderara. Class-incremental continual learning into the extended der-verse. *TPAMI*, 2022.
- [9] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. In *NeurIPS*, 2020.
- [10] Lucas Caccia, Rahaf Aljundi, Nader Asadi, Tinne Tuytelaars, Joelle Pineau, and Eugene Belilovsky. New insights on reducing abrupt representation change in online continual learning. In *ICLR*, 2022.
- [11] Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co2l: Contrastive continual learning. In *ICCV*, 2021.
- [12] Arslan Chaudhry, Marc’ Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. In *ICLR*, 2019.
- [13] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K Dokania, Philip HS Torr, and Marc’ Aurelio Ranzato. Continual learning with tiny episodic memories. In *ICML-W*, 2019.
- [14] CL Philip Chen. A rapid supervised learning neural network for function interpolation and approximation. *IEEE Transactions on Neural Networks*, 1996.
- [15] Yilun Chen, Ami Wiesel, Yonina C Eldar, and Alfred O Hero. Shrinkage algorithms for mmse covariance estimation. *IEEE transactions on signal processing*, 2010.
- [16] Yu Cheng, Felix X Yu, Rogerio S Feris, Sanjiv Kumar, Alok Choudhary, and Shi-Fu Chang. An exploration of parameter redundancy in deep networks with circulant projections. In *ICCV*, 2015.
- [17] Youngmin Cho and Lawrence Saul. Kernel methods for deep learning. *NeurIPS*, 2009.
- [18] Aristotelis Chrysakis and Marie-Francine Moens. Online bias correction for task-free continual learning. In *ICLR*, 2023.
- [19] Matthias De Lange and Tinne Tuytelaars. Continual prototype evolution: Learning online from non-stationary data streams. In *ICCV*, 2021.
- [20] Enrico Fini, Victor G Turrissi da Costa, Xavier Alameda-Pineda, Elisa Ricci, Karteek Alahari, and Julien Mairal. Self-supervised models are continual learners. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [21] Dmitriy Fradkin and David Madigan. Experiments with random projections for machine learning. In *KDD*, 2003.
- [22] Saurabh Garg, Mehrdad Farajtabar, Hadi Pouransari, Raviteja Vemulapalli, Sachin Mehta, Oncel Tuzel, Vaishaal Shankar, and Fartash Faghri. Tic-clip: Continual training of clip models. *ArXiv*, 2023.
- [23] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014.
- [24] Dipam Goswami, Yuyang Liu, Bartłomiej Twardowski, and Joost van de Weijer. Fecam: Exploiting the heterogeneity of class distributions in exemplar-free continual learning. *NeurIPS*, 2023.
- [25] Yiduo Guo, Bing Liu, and Dongyan Zhao. Online continual learning through mutual information maximization. In *ICML*, 2022.
- [26] Gunshi Gupta, Karmesh Yadav, and Liam Paull. Look-ahead meta learning for continual learning. *NeurIPS*, 2020.
- [27] Tyler L Hayes and Christopher Kanan. Lifelong machine learning with deep streaming linear discriminant analysis. In *CVPR-W*, 2020.

- [28] Timm Hess, Eli Verwimp, Gido M van de Ven, and Tinne Tuytelaars. Knowledge accumulation in continually learned representations and the issue of feature forgetting. *arXiv preprint arXiv:2304.00933*, 2023.
- [29] Guang-Bin Huang, Qin-Yu Zhu, and Chee-Kheong Siew. Extreme learning machine: theory and applications. *Neurocomputing*, 2006.
- [30] Guang-Bin Huang, Dian Hui Wang, and Yuan Lan. Extreme learning machines: a survey. *International journal of machine learning and cybernetics*, 2011.
- [31] Paul Janson, Wenxuan Zhang, Rahaf Aljundi, and Mohamed Elhoseiny. A simple baseline that questions the use of pretrained-models in continual learning. In *NeurIPS-W*, 2022.
- [32] Xisen Jin, Junyi Du, and Xiang Ren. Gradient based memory editing for task-free continual learning. In *NeurIPS*, 2021.
- [33] 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. *PNAS*, 2017.
- [34] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 2017.
- [35] Quoc Le, Tamás Sarlós, Alex Smola, et al. Fastfood-approximating kernel expansions in loglinear time. In *ICML*, 2013.
- [36] Soochan Lee, Junsoo Ha, Dongsu Zhang, and Gunhee Kim. A neural dirichlet process mixture model for task-free continual learning. In *ICLR*, 2020.
- [37] Zhizhong Li and Derek Hoiem. Learning without forgetting. *TPAMI*, 2017.
- [38] David Lopez-Paz and Marc’ Aurelio Ranzato. Gradient episodic memory for continual learning. In *NeurIPS*, 2017.
- [39] Divyam Madaan, Jaehong Yoon, Yuanchun Li, Yunxin Liu, and Sung Ju Hwang. Representational continuity for unsupervised continual learning. *arXiv preprint arXiv:2110.06976*, 2021.
- [40] 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 *CVPR*, 2021.
- [41] Mark D McDonnell, Dong Gong, Amin Parvaneh, Ehsan Abbasnejad, and Anton van den Hengel. Ranpac: Random projections and pre-trained models for continual learning. In *NeurIPS*, 2023.
- [42] Geoffrey J McLachlan. *Discriminant analysis and statistical pattern recognition*. John Wiley & Sons, 2005.
- [43] Geoffrey J McLachlan. Mahalanobis distance. *Resonance*, 4(6):20–26, 1999.
- [44] Thomas Mensink, Jakob Verbeek, Florent Perronnin, and Gabriela Csurka. Distance-based image classification: Generalizing to new classes at near-zero cost. *TPAMI*, 2013.
- [45] Seyed Iman Mirzadeh, Arslan Chaudhry, Dong Yin, Timothy Nguyen, Razvan Pascanu, Dilan Gorur, and Mehrdad Farajtabar. Architecture matters in continual learning. *arXiv preprint arXiv:2202.00275*, 2022.
- [46] Sudhanshu Mittal, Silvio Galesso, and Thomas Brox. Essentials for class incremental learning. In *CVPR*, 2021.
- [47] Aristeidis Panos, Yuriko Kobe, Daniel Olmeda Reino, Rahaf Aljundi, and Richard E Turner. First session adaptation: A strong replay-free baseline for class-incremental learning. *arXiv preprint arXiv:2303.13199*, 2023.
- [48] Karl Ezra Pilario, Mahmood Shafiee, Yi Cao, Liyun Lao, and Shuang-Hua Yang. A review of kernel methods for feature extraction in nonlinear process monitoring. *Processes*, 2020. doi: 10.3390/pr8010024.
- [49] Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *ECCV*, 2020.
- [50] Ameya Prabhu, Zhipeng Cai, Puneet Dokania, Philip Torr, Vladlen Koltun, and Ozan Sener. Online continual learning without the storage constraint. *arXiv preprint arXiv:2305.09253*, 2023.
- [51] Ameya Prabhu, Hasan Abed Al Kader Hammoud, Puneet Dokania, Philip HS Torr, Ser-Nam Lim, Bernard Ghanem, and Adel Bibi. Computationally budgeted continual learning: What does matter? In *CVPR*, 2023.
- [52] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *NeurIPS*, 2007.
- [53] Dushyant Rao, Francesco Visin, Andrei Rusu, Razvan Pascanu, Yee Whye Teh, and Raia Hadsell. Continual unsupervised representation learning. *NeurIPS*, 32, 2019.
- [54] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, 2017.

- [55] 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. In *ICLR*, 2019.
- [56] Wouter F Schmidt, Martin A Kraaijveld, Robert PW Duin, et al. Feed forward neural networks with random weights. In *ICPR*, 1992.
- [57] Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*, 2013.
- [58] Dongsub Shim, Zheda Mai, Jihwan Jeong, Scott Sanner, Hyunwoo Kim, and Jongseong Jang. Online class-incremental continual learning with adversarial shapley value. In *AAAI*, 2021.
- [59] James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim, Assaf Arbelle, Rameswar Panda, Rogerio Feris, and Zolt Kira. Coda-prompt: Continual decomposed attention-based prompting for rehearsal-free continual learning. In *CVPR*, 2023.
- [60] James Seale Smith, Junjiao Tian, Shaunak Halbe, Yen-Chang Hsu, and Zolt Kira. A closer look at rehearsal-free continual learning. In *CVPR-W*, 2023.
- [61] Hai-Long Sun, Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Pilot: A pre-trained model-based continual learning toolbox. *arXiv preprint arXiv:2309.07117*, 2023.
- [62] Guido M van de Ven and Andreas S Tolias. Three scenarios for continual learning. In *NeurIPS-W*, 2018.
- [63] Guido M Van De Ven, Zhe Li, and Andreas S Tolias. Class-incremental learning with generative classifiers. In *CVPR-W*, 2021.
- [64] Eli Verwimp, Shai Ben-David, Matthias Bethge, Andrea Cossu, Alexander Gepperth, Tyler L Hayes, Eyke Hüllermeier, Christopher Kanan, Dhireesha Kudithipudi, Christoph H Lampert, et al. Continual learning: Applications and the road forward. *arXiv preprint arXiv:2311.11908*, 2023.
- [65] Fu-Yun Wang, Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Foster: Feature boosting and compression for class-incremental learning. In *ECCV*, 2022.
- [66] Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for rehearsal-free continual learning. In *European Conference on Computer Vision (ECCV)*, 2022.
- [67] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [68] Yujie Wei, Jiaxin Ye, Zhizhong Huang, Junping Zhang, and Hongming Shan. Online prototype learning for online continual learning. In *ICCV*, 2023.
- [69] Christopher Williams and Matthias Seeger. Using the nyström method to speed up kernel machines. *NeurIPS*, 2000.
- [70] Yichen Wu, Long-Kai Huang, Renzhen Wang, Deyu Meng, and Ying Wei. Meta continual learning revisited: Implicitly enhancing online hessian approximation via variance reduction. In *ICLR*, 2024. URL <https://openreview.net/forum?id=TpD2aG1h0D>.
- [71] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *CVPR*, 2019.
- [72] Zichao Yang, Marcin Moczulski, Misha Denil, Nando De Freitas, Alex Smola, Le Song, and Ziyu Wang. Deep fried convnets. In *ICCV*, 2015.
- [73] Fei Ye and Adrian G Bors. Continual variational autoencoder learning via online cooperative memorization. In *ECCV*, 2022.
- [74] Fei Ye and Adrian G Bors. Self-evolved dynamic expansion model for task-free continual learning. In *ICCV*, 2023.
- [75] Michał Zając, Tinne Tuytelaars, and Guido M van de Ven. Prediction error-based classification for class-incremental learning. *ICLR*, 2024.
- [76] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *ICML*, 2017.
- [77] Chen Zeno, Itay Golan, Elad Hoffer, and Daniel Soudry. Task agnostic continual learning using online variational bayes. *arXiv preprint arXiv:1803.10123*, 2018.
- [78] Gengwei Zhang, Liyuan Wang, Guoliang Kang, Ling Chen, and Yunchao Wei. Slca: Slow learner with classifier alignment for continual learning on a pre-trained model. In *ICCV*, 2023.
- [79] Da-Wei Zhou, Qi-Wei Wang, Han-Jia Ye, and De-Chuan Zhan. A model or 603 exemplars: Towards memory-efficient class-incremental learning. *arXiv preprint arXiv:2205.13218*, 2022.

- [80] Da-Wei Zhou, Han-Jia Ye, De-Chuan Zhan, and Ziwei Liu. Revisiting class-incremental learning with pre-trained models: Generalizability and adaptivity are all you need. *arXiv preprint arXiv:2303.07338*, 2023.
- [81] Fei Zhu, Zhen Cheng, Xu-yao Zhang, and Cheng-lin Liu. Class-incremental learning via dual augmentation. *NeurIPS*, 2021.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper along with important assumptions.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Limitation Section is provided in the supplementary material.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: Rahimi and Recht [52] from our references details the theory for why random fourier representations perform so well quite beautifully. The random representations do not change (no continual aspect), hence the theory can be applied as-is in our case with no changes. We do not claim any novel theoretical contributions.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: RanDumb is fairly simple to implement. We dedicated half a page towards explaining hyperparameters and other information needed to reproduce all of our results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We provide code in the supplementary material to reproduce our results.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We use standard datasets and splits, we provide hyperparameters in experimental details along with ablations in experiment sections to understand the contribution of each component in our algorithm.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We state here that the only random component being the random fourier features kernel across, otherwise our method is simple and exactly reproducible. We conducted experiments with three different initialisations corresponding to seeds of the random kernel in sklearn to investigate this and different kernel initialisations lead to around ± 0.2 variation in the reported results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Detailed in Section 3 in Implementation Details. For easy access, we restate: All experiments were conducted on a CPU server with a153 48-core Intel Xeon Platinum 8268 CPU and 392GB of RAM, requiring less than 30 minutes per experiment.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: We have read the ethics guidelines and confirm that we do not use human subjects, use existing datasets, explicitly discuss social impacts.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The primary contribution in RanDumb was not to introduce a novel state-of-the-art continual learning method, but challenge prevailing assumptions and open a discussion on the efficacy of representation learning in continual learning algorithms. As

such, we do not recommend use of RanDumb for deployment in real-world production systems, hence no direct societal impact or explicit limitations on use in production systems is discussed.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We do not have any high-risk model or dataset introduced.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: MNIST, CIFAR-10, CIFAR-100, tinyImageNet and miniImagenet are cited appropriately. The licenses for these datasets is not explicitly released, hence we do not include that information.

Guidelines:

- The answer NA means that the paper does not use existing assets.

- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New Assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: Yes, we provide RanDumb with proper documentation under a GPL3 license.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: No crowdsourcing or research with human subjects was performed.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: No crowdsourcing or research with human subjects was performed.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.