

1 We thank all reviewers for the insightful feedback. We are encouraged to note that our method, MERLIN, is novel
 2 (**R2,R3,R5**); our experimental setting captures the method’s effectiveness (**R3,R5**); we achieve ‘outstanding perfor-
 3 mance compared to strong baselines’(**R2**); our approach is ‘well-grounded’ (**R1**), ‘clearly different from prior works’
 4 (**R2**). **R2** notes that the task-free scenario used is a ‘realistic direction to the area’. Further, all reviewers unanimously
 5 agree that ‘the paper is well-written’. We address all comments below. We will use the additional page allowed in the
 6 final version to incorporate feedback and address any lack of clarity in presentation (**R2**).

7 **(R1) "Comparison with Bayesian Continual
 8 Learning (BCL) methods"**: Thanks for this
 9 question, it gives more completeness to our
 10 work. We compared MERLIN against the
 11 most recent BCL work, CN-DPM (Lee et al.,
 12 ICLR’20) as **R1** suggested. As shown in the
 13 table, CN-DPM performs better than MERLIN

Methods	Split MNIST	Permuted MNIST	Split CIFAR10	Split CIFAR100	Mini-ImageNet
Single	44.89 ± 0.30	73.13 ± 2.27	73.24 ± 3.08	30.81 ± 3.57	27.57 ± 2.64
EWC	45.01 ± 0.14	74.98 ± 2.04	74.28 ± 2.2	29.23 ± 3.38	28 ± 2.59
GEM	86.79 ± 1.56	82.05 ± 4.95	79.13 ± 1.68	40.65 ± 1.95	34.17 ± 1.23
iCaRL	89.91 ± 0.92	NA	72.65 ± 1.33	27.13 ± 2.99	38.86 ± 1.63
GSS	88.39 ± 0.81	81.44 ± 1.27	57.9 ± 2.65	19.19 ± 0.7	14.81 ± 0.98
MERLIN	90.67 ± 0.80	85.54 ± 0.5	82.93 ± 1.16	43.55 ± 0.61	40.05 ± 2.94
CN-DPM (ICLR’20)	92.12 ± 0.14	-	46.01 ± 1.23	14.29 ± 0.14	-
MERLIN - SN Prior	23.34 ± 0.24	32.51 ± 1.57	28.23 ± 2.21	12.32 ± 1.45	14.76 ± 0.23

14 on Split-MNIST, but drastically fails on harder datasets. We note that the baseline methods considered in this work
 15 also perform better than CN-DPM on non-MNIST datasets. We’d also like to add that not all VAE-based CL methods
 16 learn a posterior over model params, or operate in an online CL setting. For eg., VCL (Nguyen et al, ICLR’18) learns a
 17 posterior over the data distribution (also not an online CL method), and not model param distribution. This is a subtle
 18 difference to be noted. MERLIN performs variational CL in the model param space, can be adapted easily to class +
 19 domain incremental setting, and work in task-aware + task-agnostic settings.

20 **(R1) "Why task-specific learned priors, not std normal prior?"**: As correctly noted in **R1**’s ‘Summary’, the learned
 21 task-specific priors are necessary to generate task-specific weights and consolidate meta-model on previous task params,
 22 as well as to sample models for ensembling at inference. As suggested, we ran a study where we replaced the task-specific
 23 learned prior with a standard Normal prior and finetuned the corresponding generated model on task-specific exemplars.
 24 The last row of table above (MERLIN - SN Prior) shows the result (very poor),
 25 validating the usefulness of task-specific learned priors. We also visualized the
 26 learned-task specific prior in Fig 3 (Appendix), where we see good separability
 27 across tasks.

28 **(R2, R5)"Expts on heterogeneous datasets from HAT?; Results on more do-
 29 mains?"**: Thanks for the suggestion. We ran expts on the Heterogeneous dataset
 30 from HAT as well as the AudioMNIST dataset (Becker et al, 2018) (shown in adjoining table). We comfortably
 31 outperform baselines on these datasets and domains too.

Methods	HAT	AudioMNIST
Single	47.79 ± 0.94	76.37 ± 2.25
EWC	47.19 ± 2.58	79.89 ± 16.14
GEM	67.23 ± 0.97	89.45 ± 1.14
iCaRL	62.34 ± 2.45	84.73 ± 2.22
GSS	69.79 ± 1.51	92.81 ± 0.19
MERLIN	73.54 ± 1.71	96.47 ± 1.79

32 **(R1) "Use of bigger architectures in baselines may hurt their perfor-
 33 mance"**: We re-ran all baselines with the same smaller ResNet used in
 34 MERLIN (L248), and report the results in adjoining table. We see that
 35 MERLIN outperform all baselines here again; the performance of the
 36 baseline model drops significantly, possibly due to the smaller capacity.

Methods	Split CIFAR10	Split CIFAR100	Mini-ImageNet
Single	69.65 ± 0.79	18.8 ± 2.21	18.57 ± 2.31
EWC	67.98 ± 2.96	16.89 ± 3.95	19.29 ± 3.58
GEM	72.23 ± 1.56	26.71 ± 1.75	27.71 ± 2.56
iCaRL	69.23 ± 2.24	24.81 ± 2.88	23.84 ± 1.95
GSS	49.82 ± 2.01	13.99 ± 0.56	12.92 ± 0.17
MERLIN	82.93 ± 1.16	43.55 ± 0.61	40.05 ± 2.94

37 **(R1) "For CIFAR100/miniImageNet, 10 classes/task corresponds to 5000 samples/class and not 2500?
 38 Cited works used 5 classes/task"**: For these datasets, we randomly sampled 2500 samples from 5000,
 39 and used the same 2500 across all baselines, for fair comparison (results
 40 reported across 5 such trials).
 41 To further clarify, we ran experiments with 5 samples per task (20 tasks) and report
 42 the results in adjoining table. We note that baseline accuracy matches with values
 43 reported in GEM (Tab 2, Col 3). We perform better than baselines even in this setting.

Methods	Split CIFAR100	Mini-ImageNet
Single	36.44 ± 3.44	35.85 ± 2.08
EWC	37.03 ± 2.51	35.36 ± 2.07
GEM	57.02 ± 1.41	52.28 ± 1.53
iCaRL	50.23 ± 1.37	53.22 ± 1.56
GSS	18.74 ± 0.82	16.34 ± 0.12
MERLIN	64.83 ± 1.78	57.35 ± 1.92

44 **Other clarifications:** - **(R3) "Scale of subsets"**: Sec 4.1.1 has these details. The
 45 base model is trained with 1000 (MNIST) and 2500 (all datasets other than MNIST)
 46 samples per task; - **(R1,R3) "Amount of episodic memory used for baselines"**: All baselines had access to same amount
 47 of exemplar memory as in MERLIN: 100 for MNIST and 600 for all other datasets. In Appendix Sec C, we study effect of
 48 varying exemplar memory size of MERLIN and two of its best competitors (GEM and iCaRL); - **(R1) "Inference time"**:
 49 MERLIN take 745 ms while baseline methods take ~ 300ms for CIFAR datasets on a single 1080Ti GPU. It takes more
 50 time than baselines, but is still real-time; - **(R5) "Why forgetting measure is worse for MERLIN on two datasets."**: No
 51 method is perfect. As discussed in L308-312, iCaRL uses distillation loss to ensure that logits of previous task don’t
 52 alter much while learning a new task. This brings down the forgetting measure. We still outdo all baselines on 3 other
 53 datasets; - **(R1) "Baselines in task-aware setting"**: All reported results of baselines in the paper are task-aware. Task-
 54 agnostic MERLIN was put under disadvantage while being compared to task-aware counterparts, still outperforming
 55 them. We believe that this confusion arose because of L 331, which should have been "All results of MERLIN in Tab
 56 1 do not assume task information"; - **(R3) "vote of basic model...trick to improve performance"**: Our method design
 57 allows ensembling of models for CL, though each model may be weak by itself (i.e. number of models is 1) - each
 58 individual model is upto 8× smaller in param size than baseline models (L 353). Such an ensembling approach has not
 59 been tried before, and cannot be done easily with existing methods too.