

Introduction

Continual learning aims to learn incrementally when data arrive in a possibly non i.i.d. way whereby tasks may change over time, without revisiting all previous data.

Methods

Bayes Rule gives a well-defined way to perform Continual Learning:

$$p(\theta|\mathcal{D}_{1:T}) \propto p(\theta|\mathcal{D}_{1:T-1})p(\mathcal{D}_T|\theta)$$

new posterior
old posterior/new prior
likelihood

The intractability of posteriors is tackled by Variational Inference by $q(\theta) \approx p(\theta|\mathcal{D})$:

$$\log p(\mathcal{D}) = \underbrace{\mathbb{E}_{q(\theta)} [\log p(\mathcal{D}|\theta)] - \text{KL} [q(\theta)||p(\theta)] + \text{KL} [q(\theta)||p(\theta|\mathcal{D})]}_{\text{ELBO}}$$

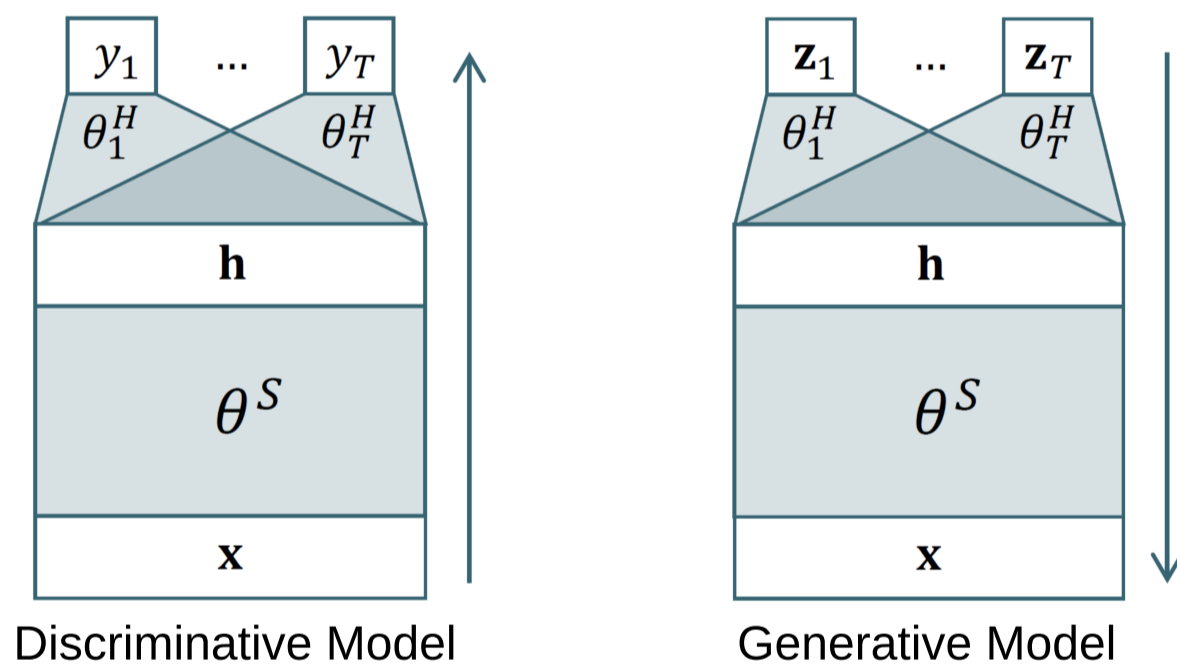
The error accumulated by sequential approximation is corrected by keeping a small “coreset” to avoid catastrophic forgetting:

$$p(\theta) \xrightarrow[\mathcal{D}_1 \setminus \mathcal{C}_1]{\text{propagation}} q(\theta|\mathcal{D}_1 \setminus \mathcal{C}_1) \xrightarrow[\mathcal{D}_2 \cup \mathcal{C}_1 \setminus \mathcal{C}_2]{\text{propagation}} q(\theta|\mathcal{D}_1 \cup \mathcal{D}_2 \setminus \mathcal{C}_2) \xrightarrow{\text{propagation}} \dots$$

$\downarrow c_1$
 $\downarrow c_2$

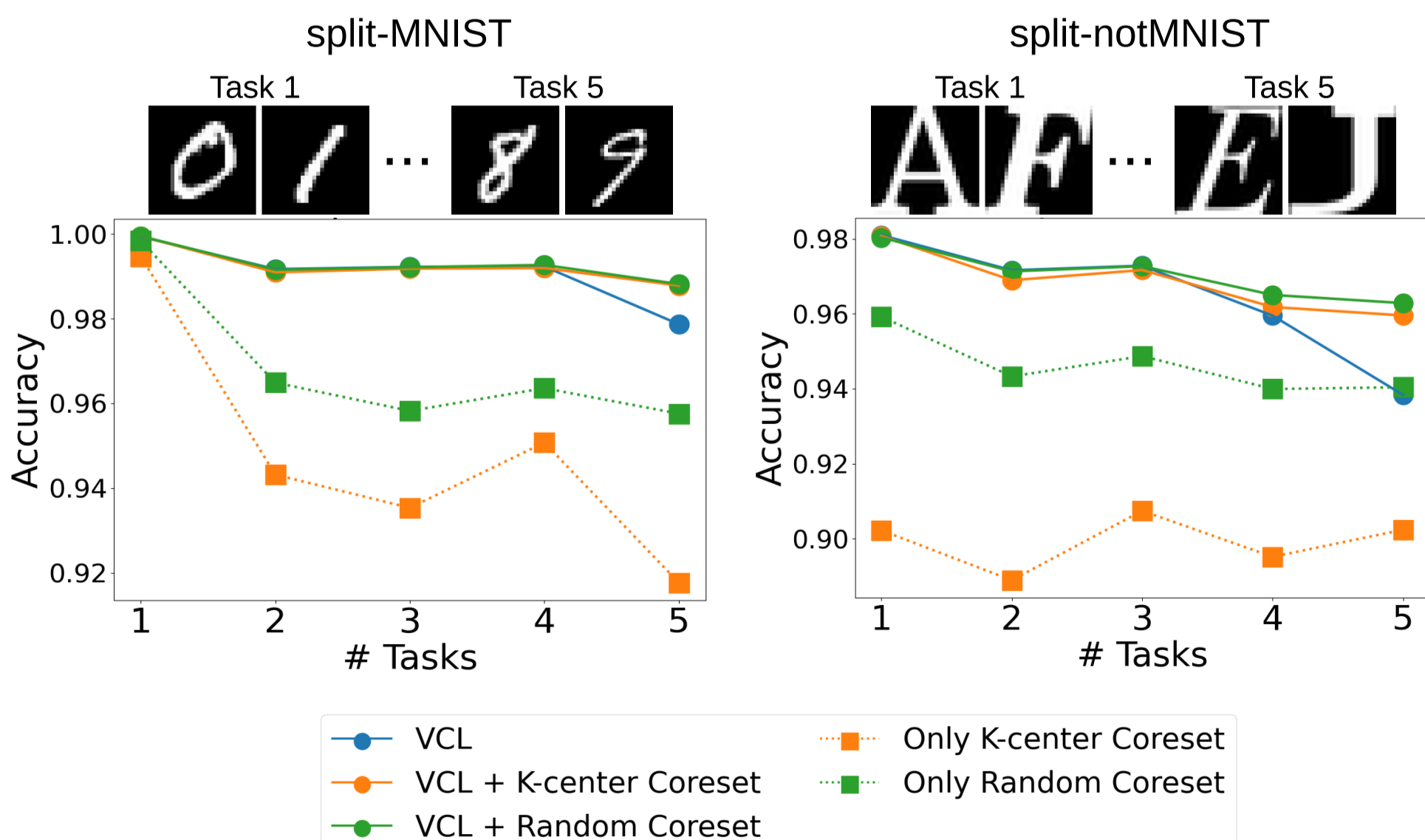
$q(\theta|\mathcal{D}_1)$
 $q(\theta|\mathcal{D}_1 \cup \mathcal{D}_2)$

Bayesian Neural Networks with the following architectures are used:



Split-MNIST/notMNIST

Datasets for each task are generated by splitting MNIST/notMNIST into subsets of two classes each.



Permuted MNIST

Datasets for each task are generated by permuting pixels of MNIST images.

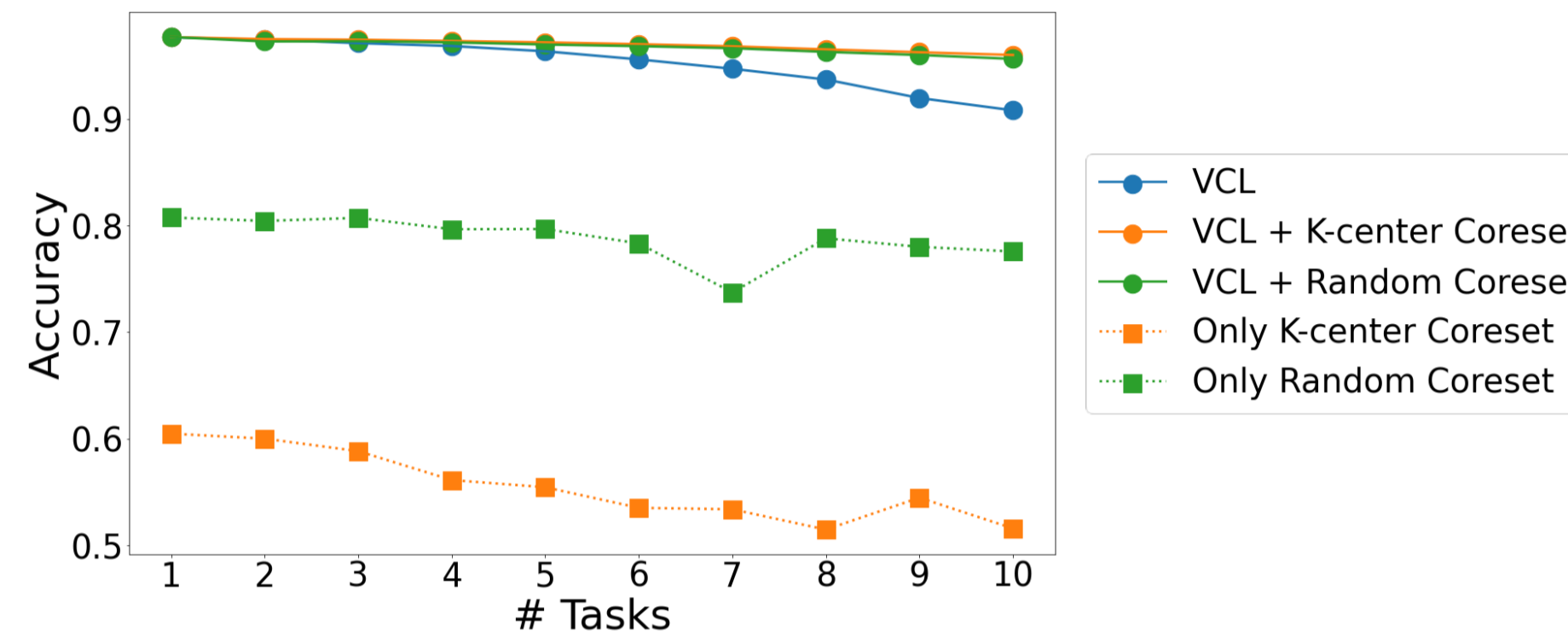
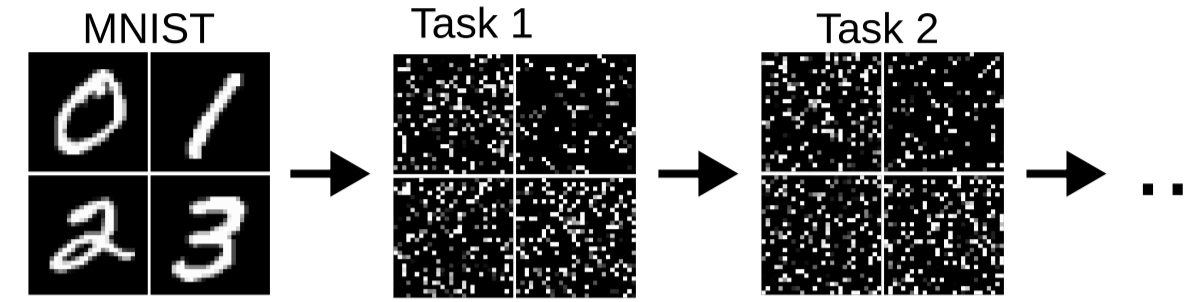
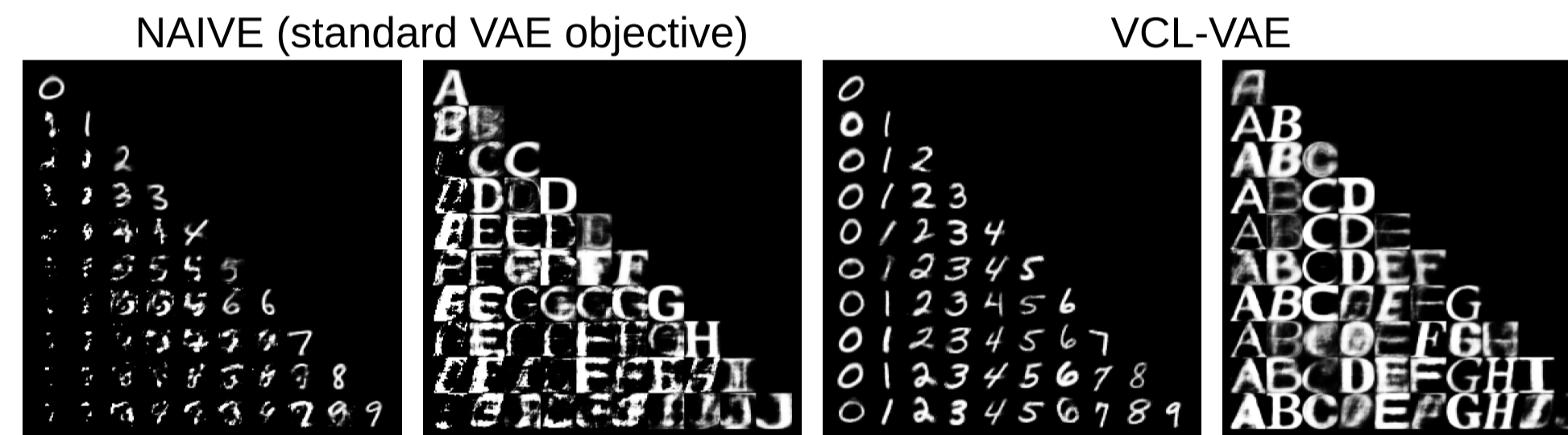
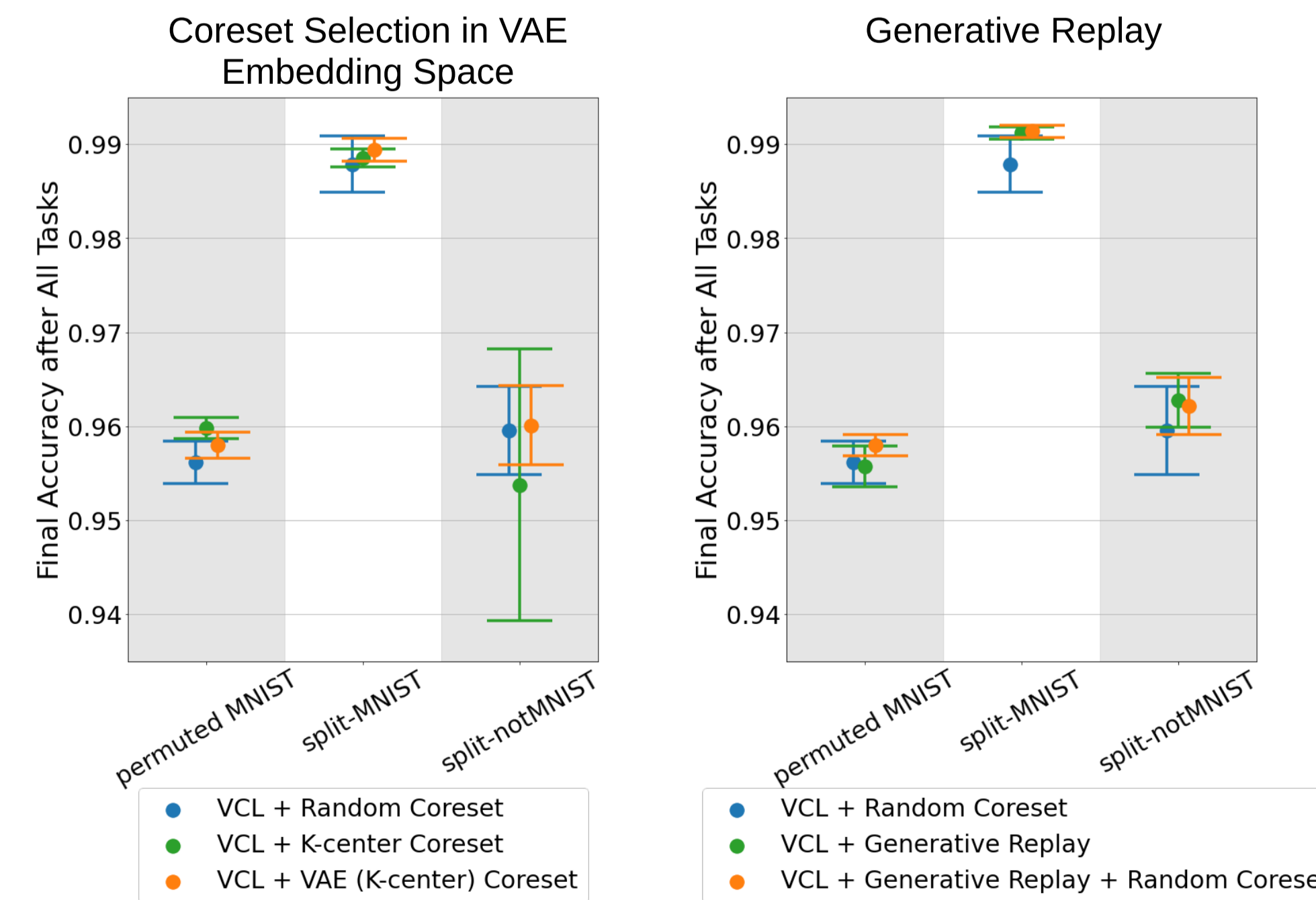


Image Generation (VCL-VAE)

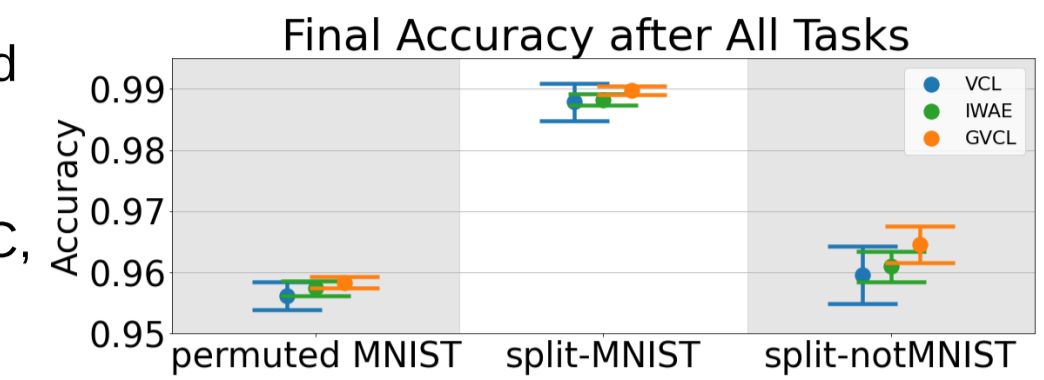


Extension: Memory Mechanism



Extension: Alternative Loss

- IWAE Bound provides a tighter bound of marginal likelihood than ELBO;
- GVCL Bound unifies VCL and EWC, another continual learning algorithm.



Extension: VCL-GAN

VCL on GAN is tricky:

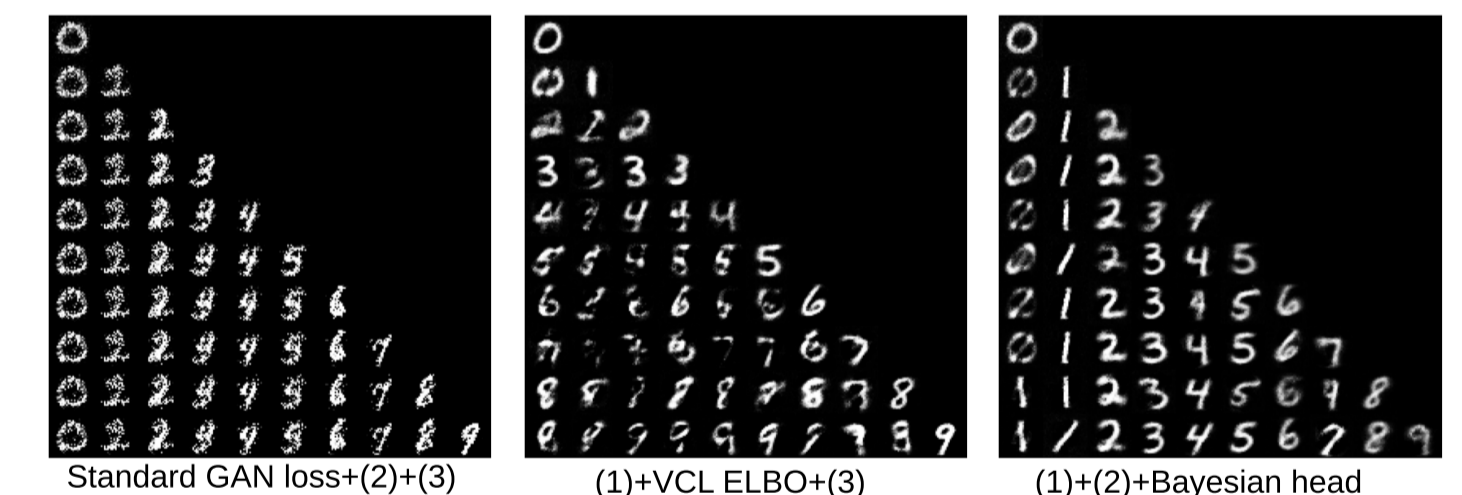
- It is difficult to balance the Bayesian-generator and discriminator;
- GAN loss is not a well-defined likelihood.



A successful VCL-GAN involves 3 key components:

- Wasserstein GAN Loss as the “negative log-likelihood” term
- Generalized-VCL Bound
- Bayesian body with task-specific non-Bayesian heads

Ablation study:



Conclusions

- VCL is a universal continual learning framework for both discriminative models and generative models;
- A more representative coreset tends to improve knowledge retention – Generative Replay combined with coreset provides a consistently better memory mechanism;
- Both IWAE Bound and GVCL Bound present superior results to the original;
- VCL was successfully applied to GAN. WGAN loss, GVCL Bound and non-Bayesian heads are key to fully functional VCL-GAN.

References

[1] Burda, Y., Grosse, R., & Salakhutdinov, R. (2015). Importance weighted autoencoders. arXiv preprint arXiv:1509.00519.
 [2] Nguyen, C. V., Li, Y., Bui, T. D., & Turner, R. E. (2017). Variational continual learning. arXiv preprint arXiv:1710.10628.
 [3] Shin, H., Lee, J. K., Kim, J., & Kim, J. (2017). Continual learning with deep generative replay. Advances in neural information processing systems, 30.
 [4] Swaroop, S., Nguyen, C. V., Bui, T. D., & Turner, R. E. (2019). Improving and understanding variational continual learning. arXiv preprint arXiv:1905.02099.
 [5] Loo, N., Swaroop, S., & Turner, R. E. (2020). Generalized variational continual learning. arXiv preprint arXiv:2011.12328.