

Introduction and Motivation

Dataset Distillation

Synthesizing a small but informative dataset \mathcal{S} that has competitive performance to the original large training dataset \mathcal{T} .

Applications

- ❖ Replay exemplars in continual learning (CL)
- ❖ Accelerating neural architecture search (NAS)
- ❖ Privacy protection in federated learning
- ❖ Membership inference defense

Motivation

- ❖ Dataset distillation algorithms typically suffer from expensive computational costs.
- ❖ Large time consumption hinders scalability to high-quality and large-scale datasets.
- ❖ Overfitting to biased samples during dataset distillation procedures (biased data distribution).
- ❖ There still exists a substantial performance gap between models trained on condensed synthetic sets and those trained on the whole dataset.

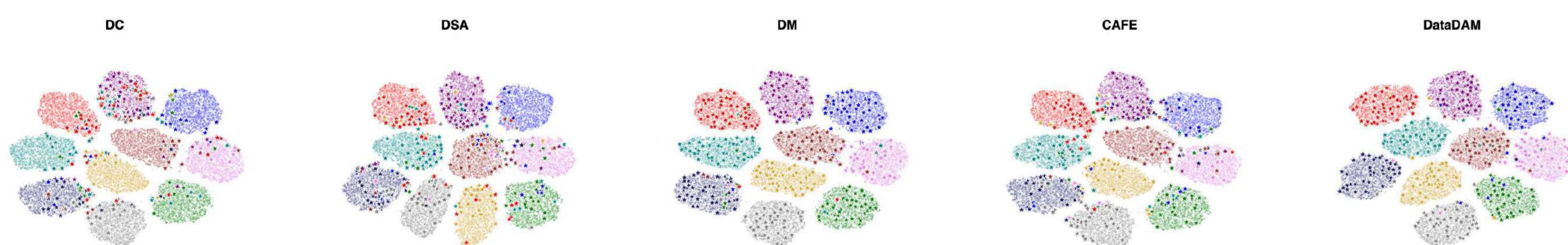


Figure 2. Data distributions of the synthetic images learned by prior methods on the CIFAR10 dataset with IPC 50. The stars represent the synthetic data dispersed amongst the original dataset.

Research Question

This study aims to answer the following research question: *Can we develop a simple yet effective data distillation algorithm capable of learning unbiased samples for any training datasets, regardless of their resolution and scale?*

Contributions

Keeping the research question in mind, we introduce a novel dataset distillation framework designed to overcome existing limitations, facilitating fast and data-efficient learning for visual classification tasks. Our contributions can be summarized as follows:

- ❖ We propose a simple method, **dataset distillation with attention matching (DataDAM)** to effectively approximate the distribution of the real dataset. This is achieved by matching the **spatial attention maps** of real and synthetic data generated by different layers within a family of randomly initialized neural networks.
- ❖ We evaluate **DataDAM** on computer vision datasets with **low, medium, and high resolutions**, where it achieves state-of-the-art results across multiple benchmark settings. Our approach also enables cross-architecture generalizations.
- ❖ We illustrate that **DataDAM** offers up to a **100x** reduction in run time costs while maintaining the lowest GPU memory consumption. Our approach also enables cross-architecture generalizations.
- ❖ We show that **DataDAM** can enhance downstream applications by improving memory efficiency for **continual learning** and accelerating **neural architecture search** through a more representative proxy dataset.

Dataset Distillation with Attention Matching (DataDAM)

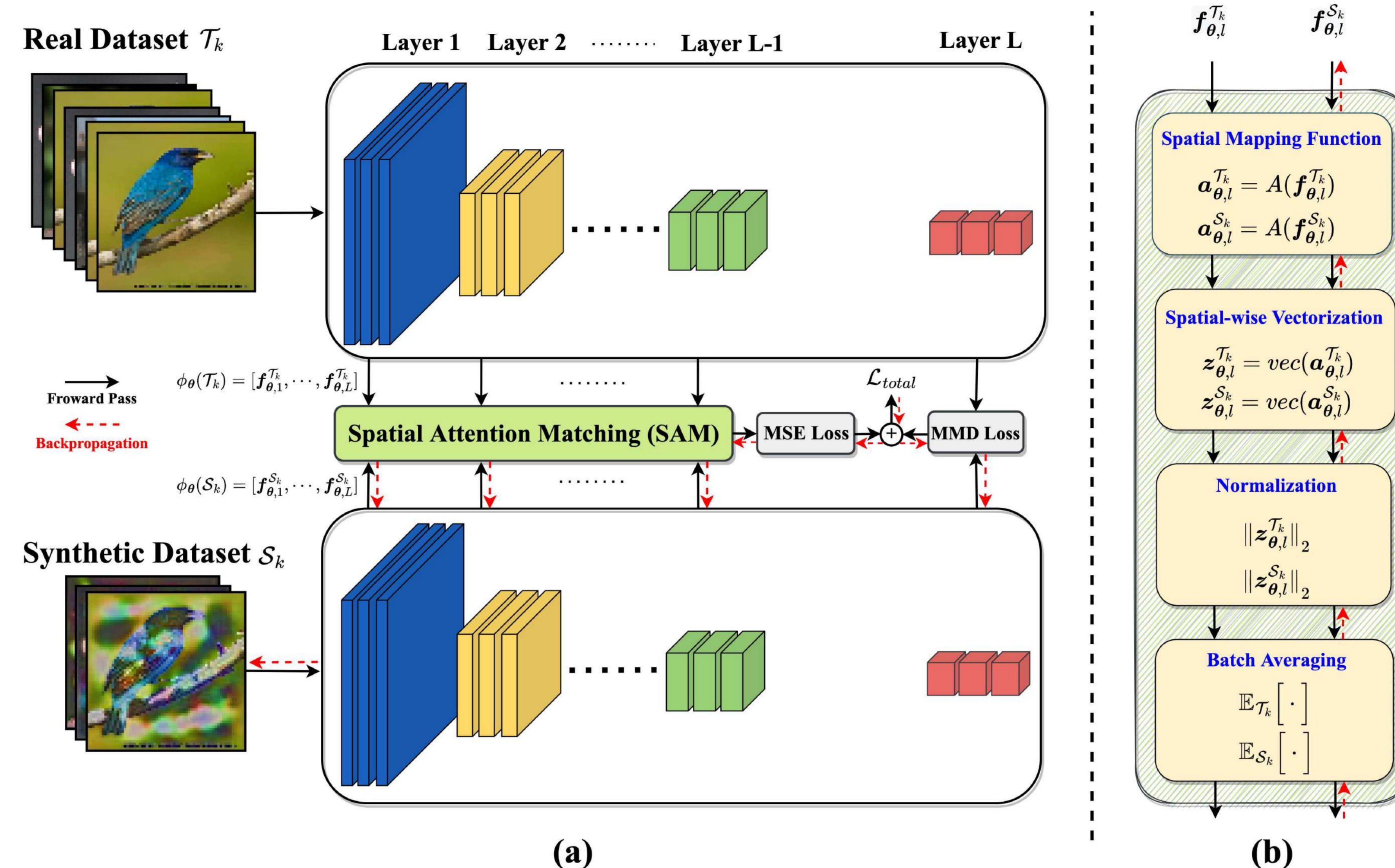


Figure 3. (a) Illustration of the DataDAM method. DataDAM includes a Spatial Attention Matching (SAM) module and a complementary MMD loss to capture the dataset's distribution. (b) The internal architecture of the SAM module.

Overall Performances on Benchmark Datasets

| Dataset | IPC | Coreset Selection | | Training Set Synthesis | | | | | | | | Whole Dataset |
|---------------|-----|-------------------|----------|------------------------|----------|----------|----------|----------|----------|-----------------|-----------------|---------------|
| | | Random | K-Center | DD | DC | DSA | DM | CAFE | KIP | MTT | DataDAM | |
| CIFAR-10 | 1 | 14.4±2.0 | 21.5±1.3 | - | 28.3±0.5 | 28.8±0.7 | 26.0±0.8 | 31.6±0.8 | 29.8±1.0 | 31.9±1.2 | 32.0±1.2 | 84.8±0.1 |
| | 10 | 26.0±1.2 | 14.7±0.9 | 36.8±1.2 | 44.9±0.5 | 52.1±0.5 | 48.9±0.6 | 50.9±0.5 | 46.1±0.7 | 56.4±0.7 | 54.2±0.8 | |
| | 50 | 43.4±1.0 | 27.0±1.4 | - | 53.9±0.5 | 60.6±0.5 | 63.0±0.4 | 62.3±0.4 | 53.2±0.7 | 65.9±0.6 | 67.0±0.4 | |
| CIFAR-100 | 1 | 4.2±0.3 | 8.4±0.3 | - | 12.8±0.3 | 13.9±0.3 | 11.4±0.3 | 14.0±0.3 | 12.0±0.2 | 13.8±0.6 | 14.5±0.5 | 56.2±0.3 |
| | 10 | 14.6±0.5 | 17.3±0.3 | - | 25.2±0.3 | 32.3±0.3 | 29.7±0.3 | 31.5±0.2 | 29.0±0.3 | 33.1±0.4 | 34.8±0.5 | |
| | 50 | 30.0±0.4 | 30.5±0.3 | - | 30.6±0.6 | 42.8±0.4 | 43.6±0.4 | 42.9±0.2 | - | 42.9±0.3 | 49.4±0.3 | |
| Tiny ImageNet | 1 | 1.4±0.1 | 1.6±0.1 | - | 5.3±0.1 | 5.7±0.1 | 3.9±0.2 | - | - | 6.2±0.4 | 8.3±0.4 | 37.6±0.4 |
| | 10 | 5.0±0.2 | 5.1±0.2 | - | 12.9±0.1 | 16.3±0.2 | 12.9±0.4 | - | - | 17.3±0.2 | 18.7±0.3 | |
| | 50 | 15.0±0.4 | 15.0±0.3 | - | 12.7±0.4 | 5.1±0.2 | 25.3±0.2 | - | - | 26.5±0.3 | 28.7±0.3 | |

Table 1. The testing accuracy % comparison to state-of-the-art methods for low- and medium-resolution datasets.

| Dataset | IPC | Random | | | Whole Dataset |
|-------------|-----|----------|----------|-----------------|---------------|
| | | DM | DataDAM | Whole Dataset | |
| ImageNet-1K | 1 | 0.5±0.1 | 1.3±0.1 | 2.0±0.1 | 33.8±0.3 |
| | 2 | 0.9±0.1 | 1.6±0.1 | 2.2±0.1 | |
| | 10 | 3.1±0.2 | 5.7±0.1 | 6.3±0.0 | |
| ImageNet | 1 | 23.5±4.8 | 32.8±0.5 | 34.7±0.9 | 87.4±1.0 |
| | 10 | 47.7±2.4 | 58.1±0.3 | 59.4±0.4 | |
| | 50 | 7.6±1.2 | 11.4±0.9 | 15.5±0.2 | |
| ImageWoof | 1 | 14.2±0.9 | 21.1±1.2 | 24.2±0.5 | 67.0±1.3 |
| | 10 | 27.0±1.9 | 31.4±0.5 | 34.4±0.4 | |
| | 50 | 21.8±0.5 | 31.2±0.7 | 36.4±0.8 | |
| ImageSquawk | 1 | 21.8±0.5 | 31.2±0.7 | 36.4±0.8 | 87.5±0.3 |
| | 10 | 40.2±0.4 | 50.4±1.2 | 55.4±0.9 | |
| | 50 | 40.2±0.4 | 50.4±1.2 | 55.4±0.9 | |

Table 2. The performance (testing accuracy %) comparison to state-of-the-art methods for large-scale and high-resolution computer vision datasets.

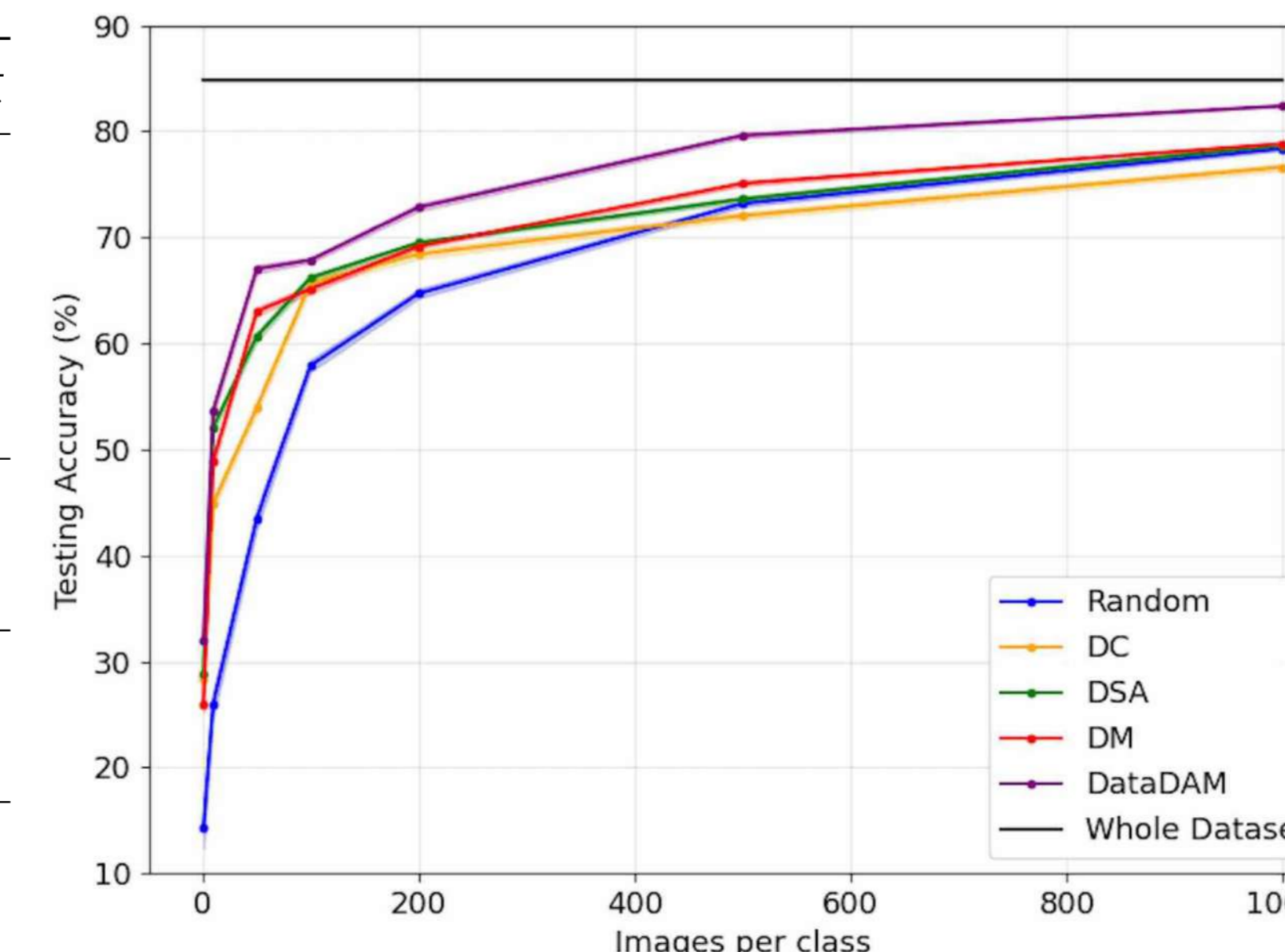


Figure 4. The testing accuracy % comparison with state-of-the-art methods on the CIFAR10 dataset for varying numbers of images per class (IPCs).

Cross-architecture Performances & Computational Cost

| Method | run time(sec) | | | GPU memory(MB) | | |
|--------|---------------|-----------------|-----------------|----------------|-------------|-------------|
| | IPC1 | IPC10 | IPC50 | IPC1 | IPC10 | IPC50 |
| DC | 0.16±0.0 | 3.31±0.0 | 15.74±0.1 | 3515 | 3621 | 4527 |
| CAFE | 0.22±0.0 | 4.47±0.1 | 20.13±0.6 | 3513 | 3639 | 4539 |
| DSA | 0.08±0.0 | 0.08±0.0 | 0.08±0.0 | 3323 | 3455 | 3605 |
| DM | 0.36±0.2 | 0.40±0.2 | OOM | 2711 | 8049 | OOM |
| MTT | 0.09±0.0 | 0.08±0.0 | 0.16±0.0 | 3452 | 3561 | 3724 |

Table 3. Cross-architecture testing performance (%) on CIFAR10 with 50 images per class. Table 4. Training time and GPU memory comparisons for state-of-the-art synthesis methods. Run time is expressed per step, averaged over 100 iterations.

Applications

DataDAM has the potential to significantly boost several downstream applications, such as enhancing memory efficiency for *continual learning* and expediting *neural architecture search* by utilizing a more representative proxy dataset.

Continual Learning:

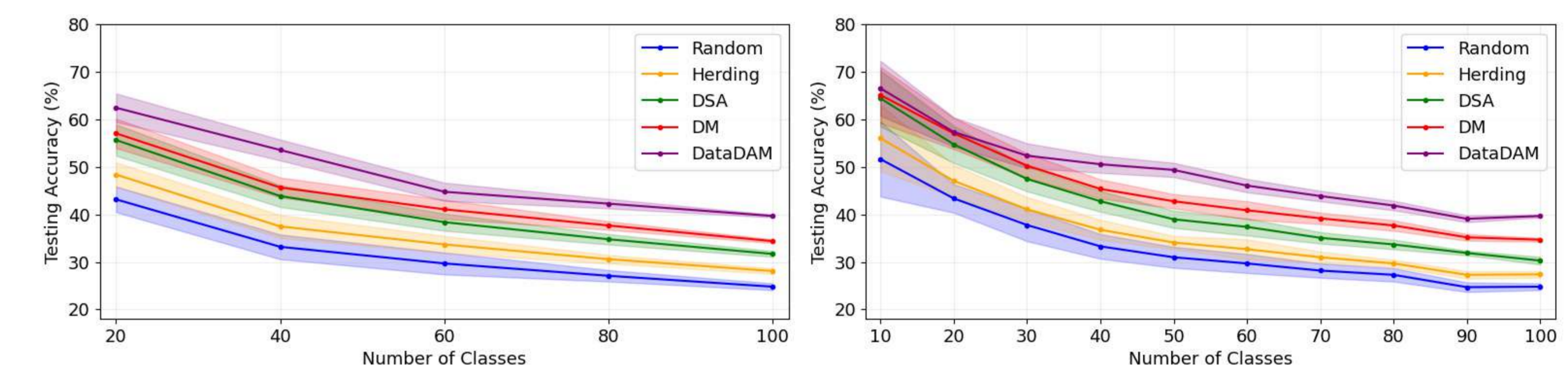


Figure 5. (Left): Showcases 5-step and (Right): Showcases 10-step continual learning with tolerance region.

Neural Architecture Search:

| Performance | Random | DM | CAFE | Ours | Early-stopping | Whole Dataset |
|-----------------|-------------|------------|------------|------------|---------------------|---------------------|
| | Correlation | 0.70 | 0.71 | 0.59 | 0.72 | 0.69 |
| Time cost (min) | 206.4 | 206.6 | 206.4 | 206.2 | 5168.9 | 5 × 10 ⁴ |
| Storage (imgs) | 500 | 500 | 500 | 500 | 5 × 10 ⁴ | 5 × 10 ⁴ |

Table 5. Neural architecture search on CIFAR10 dataset with a search space of the whole sample space.

| Performance | Random | DM | CAFE | Ours | Early-stopping | Whole Dataset |
|-----------------|-------------|------------|------------|------------|---------------------|---------------------|
| | Correlation | 0.44 | 0.51 | 0.36 | 0.69 | 0.64 |
| Time cost (min) | 33.0 | 32.2 | 30.7 | 34.8 | 37.1 | 5168.9 |
| Storage (imgs) | 500 | 500 | 500 | 500 | 5 × 10 ⁴ | 5 × 10 ⁴ |

Table 6. Neural architecture search on CIFAR10 with a search space of the top 20% of the sample space.

Distilled Image Visualization

The distilled images generated by DataDAM look real and are well-suited to be used with a variety of architectures that were not seen during training.

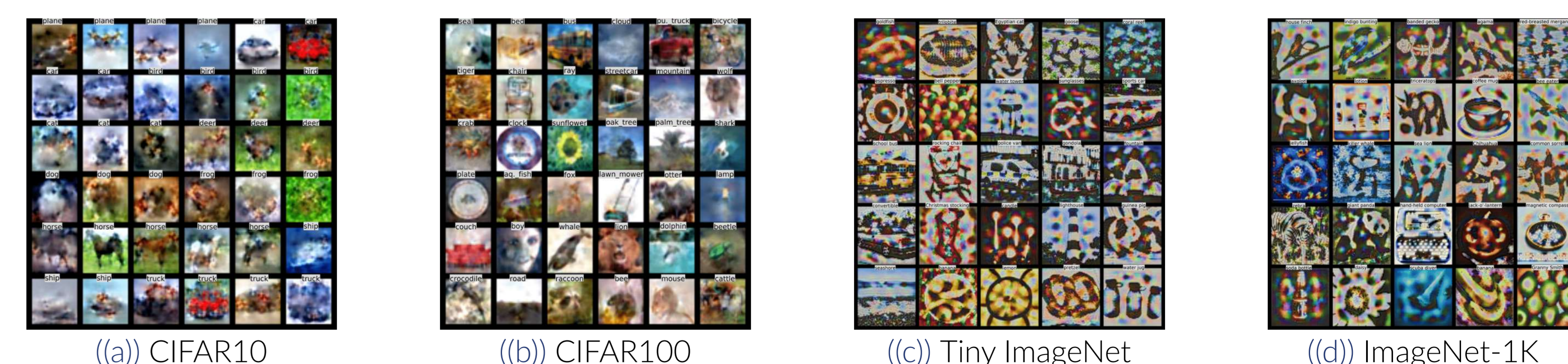


Figure 6. Example distilled images from 32x32 CIFAR10/100 (IPC10), 64x64 Tiny ImageNet (IPC1), and 64x64 ImageNet-1K (IPC1) datasets.