

# Continual Conceptual Entity Learning for Text-to-Image Generative Models

Yabin Wang, Xiaopeng Hong *Senior Member, IEEE*, Zhiheng Ma *Member, IEEE*, Zhou Su *Senior Member, IEEE*, Jinpeng Zhang, Zhiwu Huang

**Abstract**—Current Text-to-Image generative models struggle to continuously learn multiple distinct entities or concepts, limiting their scalability and hindering practical deployment in dynamic environments. We formulate this task as Continual Conceptual Entity Learning (CEL) and propose a novel framework called Continual Entity Adapter Learning (CEAL). CEAL leverages a compact set of tunable parameters, termed SuperLoRA, to efficient and scalable learning of new entities. We propose a dynamic rank-increasing strategy to train the SuperLoRA, balancing computational efficiency with performance. To evaluate our method, we create three benchmarks encompassing generic objects, human faces, and artistic styles. Experimental results demonstrate that CEAL effectively learns new entities while preserving prior knowledge, outperforming existing methods in both entity fidelity and parameter efficiency.

**Index Terms**—Continual learning, text-to-image synthesis, diffusion models.

## I. INTRODUCTION

Recently, Text-to-Image (T2I) generative models, such as GLIDE [1], DALL·E [2], and Stable Diffusion (SD) [3], have achieved remarkable success in generating diverse and complex images based on textual descriptions with high fidelity. However, they often struggle to generate styles, objects, or characters that fall outside the scope of their pre-trained datasets. For instance, while current T2I models are experts at generating images such as “a dog on the moon,” they cannot create images such as “our pet dog Buddy on the moon” without being trained on Buddy’s images during the pre-training phase. Thus, personalizing pre-trained T2I models for new user-provided concepts (e.g., styles [4], [5], objects [6], [7], or characters [8]) is increasingly in demand and has emerged as a research hotspot.

To bridge this gap, we introduce the task of Continual Conceptual Entity Learning (CEL). CEL involves sequentially integrating new conceptual entities (e.g., specific objects, persons, or styles) into existing T2I models in an ongoing manner over time, crucially, without extensive retraining on previous data. Unlike previous continual learning approaches for generative models [9]–[14] focused on class-level generation or single-session multi-concept customization methods (e.g.,

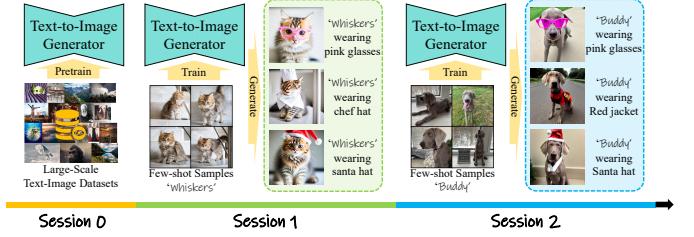


Fig. 1: In CEL, the model sequentially learns new entities (e.g., pets named ‘Whiskers’ or ‘Buddy’) from limited samples in each learning session. After learning, the T2I generator can generate images of both newly learned and previously known entities based on text prompts.

[33]) that learn multiple concepts simultaneously, CEL emphasizes entity-level sequential learning. The model must capture and reproduce individual entities’ unique characteristics and details (see Fig. 1). The primary challenges of CEL stem from three core requirements: (1) sequentially learning new, distinct entities effectively from limited data, (2) mitigating catastrophic forgetting of previously acquired entities, and (3) achieving this in a parameter-efficient manner to ensure scalability.

Existing approaches [6], [7] for T2I personalization are mainly designed for single-stage incremental learning and are not well-suited for the continuous, sequential learning of multiple new entities over time. When multiple new entities are continually introduced, these existing methods struggle to maintain both newly learned and previously acquired knowledge, leading to catastrophic forgetting. Recently, Low-Rank Adaptation (LoRA) [15], [16] has emerged as an efficient method for fine-tuning T2I models using low-rank adapters, allowing for scalable learning of numerous entities with small parameter expansion. However, the storage requirements for new entities remain a significant challenge, particularly in resource-constrained environments. Additionally, our findings indicate that directly merging multiple LoRA weights can result in conflicts, preventing the generation of high-quality outputs of new concepts.

To overcome these challenges, we propose a novel framework, termed Continual Entity Adapter Learning (CEAL), for CEL. Grounded in the theoretical understanding of parameter adaptation and network function, CEAL proactively searches for a layer-wise LoRA allocation for each entity that can be integrated with minimal conflict, thereby ensuring entity fidelity, knowledge preservation, and generalization ability.

CEAL’s design is grounded in the following theoretical

Y. Wang and X. Hong are with Harbin Institute of Technology, Harbin, China (e-mail: wang-yabin@outlook.com and hongxiaopeng@ieee.org). Z. Su is with Xi'an Jiaotong University, Xi'an, China (e-mail: zhousu@ieee.org). Z. Ma is with Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China (e-mail: zh.ma@siat.ac.cn). J. Zhang is with Intelligent Science & Technology Academy Limited of CASIC, Beijing 100043, China. (e-mail: 2076707060@qq.com). Z. Huang is with the University of Southampton, Southampton, UK (e-mail: Zhiwu.Huang@soton.ac.uk).

principles and empirical observations: First, following the principle of Functional Localisation observed in deep learning models [17], diffusion models feature functionally distinct layers, each specializing in specific image generation aspects such as content, style, or color attributes [18], [19]. This Layer-wise Functional Specialization thus requires optimized, layer-specific LoRA configurations and provides the theoretical grounding for strategically allocating these parameters with varying ranks to mitigate interference. Second, the foundation of LoRA relies on the Intrinsic Dimension Hypothesis [20], suggesting adaptation lies in a low-dimensional subspace. Furthermore, consistent with the Lottery Ticket Hypothesis [21], which implies sparse weight updates can be sufficient, LoRA weights typically exhibit high sparsity. This means inter-entity conflicts are largely confined to a small parameter subset, thus leaving ample room to integrate new knowledge with minimal interference. Third, LoRA modules typically lack orthogonality, and their naive merging often leads to destructive interference from conflicting entity information.

Thus, CEAL allocates LoRAs with varying ranks to each layer instead of using a uniform rank. To implement this, we introduce a new structure called *Superlora*, a rank-degradable LoRA architecture viewed as a combination of multiple low-rank LoRAs, which enables weight-sharing. We train *Superlora* for each entity and search for the optimal configuration based on multiple criteria. The optimized LoRA weights are then merged into the base model, allowing for the generation of new entities without increasing the overall parameter count.

The training quality of *Superlora* directly affects the search results. Since exhaustively training all possible LoRA configurations is computationally infeasible, an efficient strategy is crucial to explore the vast configuration space effectively. We introduce a dynamic rank-increasing strategy that gradually expands the search space by introducing new rank parameters alongside the training steps. This approach allows *Superlora* to initially focus on learning essential low-rank configurations before progressively exploring more complex, higher-rank configurations. After training *Superlora* for a new entity, an optimizer explores *Superlora* using multiple criteria to determine the optimal LoRA rank configuration for that entity.

To facilitate standardized evaluation of methods addressing the unique challenges of CEL, we construct three benchmarks that encapsulate a range of different concepts, including generic objects, human faces, and distinctive artistic styles. In addition, we design evaluation metrics for CEL that evaluate the model's performance on new tasks and its susceptibility to forgetting in continuous learning. Extensive experiments on three benchmarks demonstrate the effectiveness of CEAL. Code and data will be publicly available.

The technical contributions are manifold:

- We propose a framework, CEAL, to continuously learn entities using a compact set of tunable parameters. To optimize CEAL, we define three key CEL criteria that encompass concept fidelity, parameter efficiency, and aesthetic quality.
- We propose a novel rank-degradable LoRA architecture called *Superlora*. This architecture allows for weight-sharing and dynamic allocation of LoRA ranks across different layers, enabling efficient and scalable searching.

Our dynamic rank-increasing strategy further stabilizes the training process, ensuring effective learning of new entities.

- We establish three benchmarks for CEL. Each benchmark focuses on different aspects of the task, allowing for comprehensive evaluation.

## II. RELATED WORKS

### A. Text-to-image Personalization

Text-to-image (T2I) generation models [1], [22]–[31] have received unprecedented interest from the community in recent years. To achieve text-to-image generation, these models typically first employ a language encoder, such as CLIP [32], to encode user text input into a latent representation. This latent representation then serves as a conditional input, with the model subsequently trained on large-scale paired image-text datasets to generate corresponding images.

With the rapid advance of T2I models, their personalization (*a.k.a.* customization) is becoming increasingly crucial. Personalization methods [33]–[35] tailor a model to the specific needs of an individual or group by utilizing data unique to the intended users. DreamBooth [6] fine-tunes all parameters of the diffusion network by just giving a few images of the target entity. DreamBooth uses the pre-trained model to generate regularization data about similar concepts to relieve forgetting. However, tuning all parameters on few-shot images would increase the risk of the model forgetting previously learned knowledge and overfitting. Textual Inversion [7] proposes a concise way for customization by solely optimizing a set of word embeddings to portray a novel concept while leaving the denoising model fixed. Although it effectively preserves the original model's knowledge, it may have restricted learning capability and may encounter difficulty in comprehending intricate entities. Cones [19] identifies and manipulates concept neurons within diffusion models to enable efficient and customizable multi-subject image generation. InstantBooth [36] captures subject identity by representing the general concept as a textual embedding and feeding fine-grained visual details through lightweight adapter layers into the frozen model. HyperDreamBooth [37] uses a HyperNetwork to generate compact, personalized weights from a single image, enabling extremely fast personalization of T2I models. Custom Diffusion (CD) [38] combines these two strategies for multi-concept learning. CD learns a set of textual embeddings for a given concept and fixes the parameters of the text encoder, only fine-tuning a small subset of parameters in the cross-attention layers of the denoising model, specifically the key and value matrices. Concept Weaver [39], building upon CD, then addresses how to effectively fuse the visual features of these multiple distinct concepts into a single image at inference time. It first generates a template image aligned with the text prompt, then utilizes spatial region masks to inject the visual appearances of multiple customized concepts into their respective regions. Similarly, ED-LoRA [40] merges multiple LoRAs into a single model using Gradient Fusion and employs Regionally Controllable Sampling to accurately place the customized concepts and their attributes in images. LoRA-Composer [41] is a training-free method that

leverages the cross-attention mechanism to inject concept-specific LoRA features into image regions designated by user-provided layout conditions. MC2 [42] enables training-free composition of heterogeneous single-concept models (LoRA or TI) using inference-time optimization. Its Multi-concept Guidance (MCG) refines attention weights to spatially disentangle concepts and ensure their faithful representation while minimizing interference. However, these works primarily focus on composing multiple pre-existing customized concepts. The CEL task, in stark contrast, centers on the continual, sequential learning of new entities and the critical mitigation of catastrophic forgetting.

Recent works [16] adopt parameter-efficient fine-tuning methods, such as LoRA [15], to learn new entities without tuning the entire denoising network. Building upon the parameter efficiency of LoRA, some works [43]–[45] focus on optimizing the allocation of trainable parameters under the constraint of limited storage resources. DyLoRA [43] trains LoRA with a range of ranks, enabling dynamic adjustment of model capacity without retraining. AdaLoRA [44] adaptively allocates the parameter budget among different weight matrices based on their importance scores. IncreLoRA [45] incrementally adds trainable parameters during training. However, while these methods enhance LoRA’s parameter efficiency, they are not tailored for the CEL task. They typically lack a proactive mechanism to automatically determine optimal, layer-wise LoRA configurations that minimize inter-entity conflict during sequential learning, thus often requiring users to define these specific CEL-focused setups. Moreover, *stand-alone* addition of each new concept brings substantial modifications to the model parameters, leading to an inevitable rise in storage burden and knowledge forgetting.

### 220 B. Continual Learning

Continual learning is a machine learning paradigm that learns new data over time without forgetting previously learned knowledge. Three primary strategies are used to address the challenge of *catastrophic forgetting*. *Rehearsal* methods allow a model to partially access and utilize data from previous sessions when learning new tasks. Previous data can be obtained by storing samples from past tasks [46]–[56] or synthesizing through generative models [9], [57]–[60]. *Regularization* methods [48], [61]–[68] use the knowledge from the old model to guide and constrain the learning process for new tasks, by techniques such as knowledge distillation. *Network Expansion* methods [69]–[73] augment the model’s parameters to incorporate new knowledge without interfering with previous parameters, using methods like network expanding-pruning [74] and Neural Architecture Search (NAS) [75], [76]. Recently, parameter efficient tuning, such as prompt-tuning [77]–[80] and adapter [81], shows a promising way for continual learning on pre-trained models. However, most studies focus only on a class-level continual learning scenario, where the *class* serves as the learning unit. Consequently, they cannot be directly applied to the CEL scenario, where the learning unit is the *entity* rather than the *class*. Recently, C-LoRA [81], [82] introduces the use of LoRA to learn each

entity for T2I diffusion models. However, C-LoRA stores all LoRA parameters in memory, which becomes problematic when dealing with a large number of entities.

## III. METHODS

### A. Preliminary and Problem Formulation

**Preliminary for T2I Diffusion Models:** To achieve high-quality and creative image generation, T2I diffusion models [3], [83], [84] have become a hot topic of current research. T2I diffusion models are multimodal generators that learn precise correspondences between textual descriptions and images. Based on these learned correspondences, trained models generate images aligned with given textual prompts through a progressive denoising process. Given a diffusion model  $\theta$ , the training objective can be derived as Eq. 1.

$$\mathbb{E}_{(\mathbf{x}_0, \mathcal{C}, \epsilon, t)} \left[ \|\epsilon - \hat{\epsilon}_\theta(\mathbf{x}_t, \mathcal{C}, t)\|^2 \right], \quad (1)$$

where  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$  is a random noise,  $\mathbf{x}_0$  is the original clean image,  $\mathbf{x}_t$  corresponds to the noised image at the  $t$ -th timestep, and  $\mathcal{C}$  is the text embeddings as condition.

**Problem Formulation:** In CEL, we continually insert knowledge of new entities into a pre-trained T2I diffusion model by reducing the reconstruction loss 1. An *entity* refers to any conceptual or physical item, such as an object, character, or artistic style.

Formally, the model learns sequentially over  $S$  sessions, where each session introduces a new entity. In the session  $s$ , the T2I model  $\theta$  encounters new data  $D_s = \{(x_i^s, c_i^s) | i \in N^s\}$  of a specific entity  $v$ . Here,  $N^s$  represents the number of training samples.  $c_i^s$  denotes the textual description corresponding to the image  $x_i^s$ . It follows a structured template like “a photo of  $\text{id}$ ,” where  $\text{id}$  is the identifier for the newly introduced entity. Once trained, we can use  $\text{id}$  in text prompts and the model  $\theta$  to create the desired images. This process is repeated until all  $S$  entities are learned.  $\theta$  is expected to have the capability to generate images containing any combination of the learned  $S$  entities, as shown in Fig. 5.

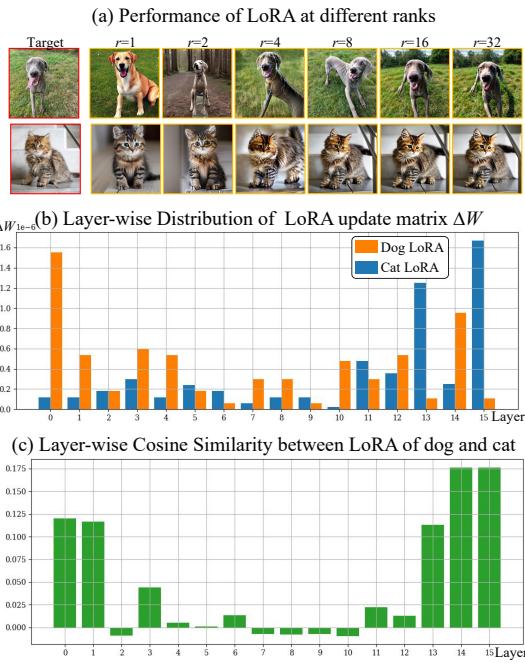
### B. Empirical Analysis of LoRA

LoRA [15] is an effective approach for fine-tuning the denoising model to incorporate new entities. Unlike traditional fine-tuning, LoRA introduces compact, low-rank matrices to adapt existing pre-trained weights. Only these new matrices are updated during training, keeping the original model parameters frozen. Formally, given a pre-trained T2I model  $\theta$ , the weights of the  $l^{th}$  layer is  $W^l \in \mathbb{R}^{n_l \times m_l}$ , where  $n_l$  and  $m_l$  represent the input and output dimensions of the  $l^{th}$  layer, respectively. The calculation of low-rank adaptation applied to  $W^l$  can be expressed as Eq.2.

$$\hat{W}^l = W^l + A_r^l (B_r^l)^T, \quad (2)$$

where  $A_r^l \in \mathbb{R}^{n_l \times r}$  and  $B_r^l \in \mathbb{R}^{m_l \times r}$  are the low-rank matrices for the  $l^{th}$  layer and  $r$  represents the rank of LoRA.

To analyze the behavior of LoRA in personalizing T2I models, we conduct a series of experiments focusing on the entities ‘dog’ and ‘cat’ within our CEL benchmarks.



**Fig. 2: Analysis of LoRA in T2I model personalization.**  
(a) Performance of pruned LoRA models for 'dog' and 'cat' concepts at different ranks. (b) Layer-wise distribution of LoRA update matrix  $\Delta W$ . (c) Layer-wise cosine similarity between 'dog' and 'cat' LoRA weights.

We train two independent LoRAs with a rank of  $r = 32$  for each entity. After training, we prune the trained LoRA to certain ranks using SVD decomposition. Interestingly, as illustrated in Fig. 2 (a), we can find that despite removing half of the parameters, the LoRA models maintain their capability to effectively generate the intended concepts. This suggests the sparsity, where many components of the LoRA weights have minimal contributions. This observation strongly supports the Lottery Ticket Hypothesis [21] in the context of LoRA, implying that a sparse weight update is often sufficient. Consequently, sacrificing some of these weights may not significantly impair LoRA's overall functionality. Similar observations have been reported in related studies [85], [86].

Then, we conduct a layer-wise analysis of the parameter distribution in LoRA. As shown in Fig. 2 (b), we find that a significant portion of the parameters  $\Delta W$  in LoRA have magnitudes very close to zero, indicating they have little impact on the output of the pre-trained models. Moreover, the weight distributions of LoRA vary across different entities, suggesting that various entities need distinct parts of the model. This observation implies that different layers within diffusion models are not uniform in function. This provides direct empirical evidence for the principle of Functional Localisation [17] within diffusion models, aligning with findings in related research [18], [19], and theoretically justifying our strategy of allocating layer-specific ranks rather than a uniform rank.

Finally, we computed the layer-wise cosine similarity distribution between the LoRA weights for 'dog' and 'cat', as

depicted in Fig. 2 (c). Despite a few layers exhibiting non-zero cosine similarity, the majority of the layers demonstrate values exceedingly close to zero, implying a high degree of orthogonality. This observation suggests that directly merging the two LoRA weights together might cause the knowledge to overlap, potentially leading to less distinct representations.

### C. Continual Entity Adapter Learning Framework

Based on the aforementioned theoretical principles and empirical analysis, we identify several limitations of traditional LoRA methods, including inefficient parameter usage and challenges in merging knowledge. These drawbacks make traditional LoRA unsuitable for CEL scenarios, where LoRA must sequentially incorporate new entities while maintaining parameter efficiency and ensuring seamless integration of multiple adaptations.

To address these challenges, we propose the framework, CEAL, to continuously and efficiently integrate a sequence of new entities into a pre-trained T2I diffusion model, as illustrated in Fig. 3.

The CEAL framework operates in two main stages. First, we introduce *Superlora*, a novel structure designed to effectively manage a wide spectrum of LoRA configurations within a single adapter. Through a dynamic rank-increasing strategy, *Superlora* is a flexible foundational adapter with adjustable rank complexity. This core characteristic enables the direct derivation of diverse sub-LoRA configurations, without the need for retraining. This significantly simplifies the search for optimal LoRA configurations across the extensive parameter space.

Second, leveraging the trained *Superlora*, we search for the optimal layer-wise rank allocation. This involves sampling LoRA configurations, generating images, and evaluating them against multi-objective criteria, specifically entity fidelity, parameter efficiency, and aesthetic quality, to identify compact and high-performing LoRAs.

The resulting LoRA can then be seamlessly merged into the base T2I model, providing a parameter-efficient and effective solution for the continual integration of new entities.

### D. Superlora Structure and Efficient Training Strategy

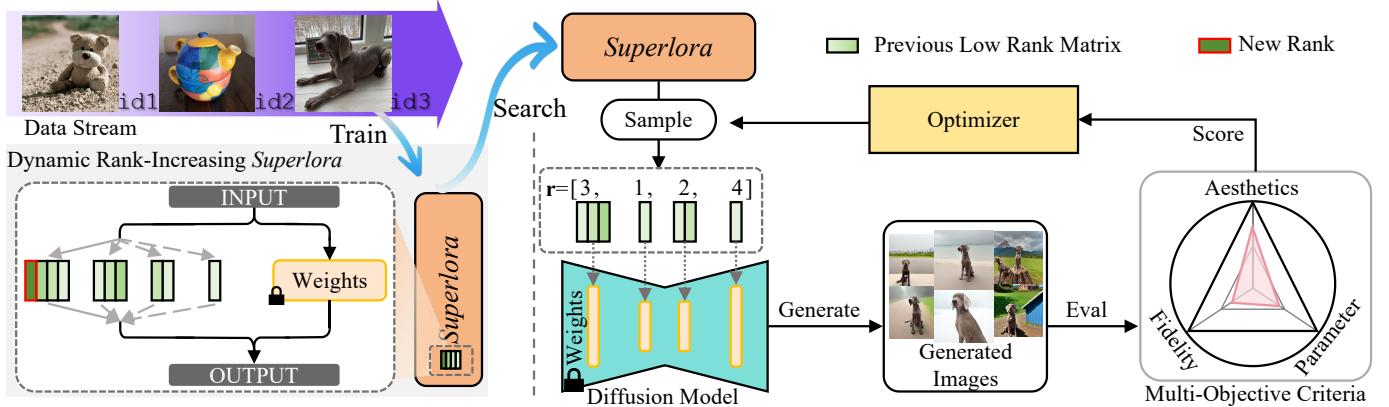
This section details the *Superlora* structure and the efficient training strategies, specifically weight-sharing and dynamic rank-increasing, designed to manage its vast configuration space and enable its rank-degradable properties. To detail the *Superlora* architecture, we first establish the following definitions:

**Definition 1:** For a given layer  $l$ , a LoRA module attached to the original weight matrix  $W^l$  can be expressed as Eq.3.

$$\hat{W}^l = W^l + A_{r_l}^l (B_{r_l}^l)^T, \quad (3)$$

where  $A_{r_l}^l \in \mathbb{R}^{n_l \times r_l}$  and  $B_{r_l}^l \in \mathbb{R}^{m_l \times r_l}$  are the low-rank matrices for the  $l^{th}$  layer with rank  $r_l$ .

**Definition 2:** The *Superlora* structure, mathematically defined as a set  $\mathcal{R}$ , encompasses all possible combinations of LoRA configurations across the network's layers. It is defined as:



**Fig. 3: Overview of CEAL Framework.** In a training stage, CEAL continually integrates new entities into a pre-trained T2I diffusion model through a two-step optimization process, which consists of a prerequisite training process (**left**) with a rank-degradable LoRA structure, *Superlora*, and a search and optimization process (**right**) to identify the optimal LoRA configuration and optimize the parameter settings according to multiple objectives.

376  $\mathcal{R} = \{(A_{r_l}^l, B_{r_l}^l) \mid l = 1, 2, \dots, N\}$  where  $N$  is the total  
377 number of layers in the model.

378 **Definition 3:** A specific LoRA configuration within *Superlora*  
379 is represented by a rank vector  $\mathbf{r} = [r_1, r_2, \dots, r_N]$ , where  
380 each  $r_l$  satisfies  $0 \leq r_l \leq \hat{r}$  and  $r_l \in \mathbb{Z}$ . This vector  $\mathbf{r}$  specifies  
381 the rank of the LoRA adaptation at each layer of the network.  
382  $\hat{r}$  is the maximum rank of *Superlora*. When  $r_l = 0$ , the LoRA  
383 module is not applied to the  $l^{\text{th}}$  layer.

384 Consequently, exhaustively training all these possible  
385 LoRA configurations is computationally intensive and time-  
386 consuming. For instance, a diffusion model with  $N = 16$   
387 layers and a maximum rank  $\hat{r} = 8$  per layer would yield  
388 approximately  $1 \times 10^{14}$  possible rank configurations. Training  
389 such a vast number of configurations is impractical.

390 To address this impracticality and train *Superlora* efficiently,  
391 we employ a weight-sharing strategy. The training process  
392 does not involve training each sub-network independently.  
393 Instead, at each iteration, a specific configuration  $\mathbf{r}$ , is ran-  
394 domly sampled. Subsequently, only the parameters within the  
395 sampled rank are updated, with all other parameters kept  
396 frozen during that step. This strategy significantly reduces  
397 the computational burden of training the SuperLoRA struc-  
398 ture itself, as it avoids the need to train all its numerous  
399 potential configurations separately, thereby enabling efficient  
400 exploration of its vast configuration space specifically during  
401 this training phase.

402 While weight-sharing makes *Superlora* training feasible, its  
403 vast search space can still lead to inefficient learning from  
404 uniform random sampling. To effectively train the *Superlora*  
405 in such a large search space, we introduce the dynamic  
406 rank-increasing strategy, which is similar to the curriculum  
407 learning paradigm. Initially, we set the maximum rank  $\hat{r}$   
408 to 1, focusing the training on configurations with minimal  
409 complexity. We gradually allocate more rank parameters to  
410 *Superlora* as training progresses. This rank increase operates  
411 at regular intervals, denoted by  $\nu$ , throughout the training  
412 process. At each interval, we increment  $\hat{r}$  by a predefined step  
413 size until it reaches a maximum rank hyperparameter,  $r_{\text{final}}$ .

414 By gradually increasing the rank, we allow the model  
415 to first learn low-rank adaptations, which capture the most  
416 significant features with fewer parameters. This facilitates a  
417 more stable and efficient training process. As higher ranks are  
418 introduced, the model progressively captures more complex  
419 representations without overwhelming the training dynamics.  
420 This strategy also ensures that all ranks up to  $r_{\text{final}}$  are  
421 adequately trained.

#### E. Searching with Multi-objective Criteria

422 To proactively search for the optimal LoRA for any entity  
423 and ensure high-quality image generation, we define a  
424 multi-objective function  $L$  that incorporates three essential  
425 components: aesthetic quality, parameter efficiency, and entity  
426 fidelity.

427 Firstly, to assess aesthetic quality, we employ an aesthetic  
428 score predictor [87], [88] trained on human aesthetic prefer-  
429 ence data. This predictor evaluates the generated images based  
430 on prompts and assigns a score, referred to as the *Human*  
431 *Preference Metric* ( $L_{\text{hpm}}$ ), indicating the aesthetic appeal of  
432 the images. Higher scores reflect images that are more likely  
433 to be preferred by humans.

434 Secondly, the *Memory Cost Metric* ( $L_{\text{mcm}}$ ) quantifies the  
435 number of parameters in the LoRA modules. An unconstrained  
436 number of parameters may lead to overfitting and catastrophic  
437 forgetting [67]. We calculate  $L_{\text{mcm}}$  by summing the number  
438 of parameters of the selected LoRA modules, ensuring that  
439 the LoRA has a compact storage size. This ensures that only  
440 the necessary parameter changes are performed.

441 Thirdly, the *Concept Similarity Metric* ( $L_{\text{csm}}$ ) assesses the  
442 content fidelity between generated images and target entities.  
443 This is achieved by S-Prompts [79], a continual AI-generated  
444 detector. Specifically, S-Prompts is a binary classifier built  
445 upon the pre-trained CLIP. For each new entity, this classifier  
446 is trained to distinguish the entity's authentic training images  
447 from those generated by the model. The resulting classification  
448 accuracy of S-Prompts directly constitutes the  $L_{\text{csm}}$  score,

450 where higher accuracy indicates greater fidelity of the generated images to the characteristics of the target entity.  
 451

452 The overall function of  $L$  is a weighted sum of these three  
 453 metrics as Eq. 4.

$$L = -\alpha L_{csm} + \beta L_{mcm} - \gamma L_{hpm}. \quad (4)$$

454 where  $\alpha$ ,  $\beta$ , and  $\gamma$  are hyperparameters that balance the  
 455 magnitude of differences metrics.

---

**Algorithm 1** Optimizer for CEAL
 

---

```

1: Input: Search space  $R$ , Initial temperature  $T_0$ , Cooling
   rate  $\alpha$ , Total iterations  $N_I$ , Objective Function  $L$ 
2: Output: Optimal solution
3: Initialize a solution  $r$  of length 16, with randomly selected
   values from  $R$ 
4: Set current temperature  $T = T_0$ 
5: for  $i$  in 1 to  $N_I$  do
6:   Generate a new solution  $r'$  by random modifying  $r$ 
7:   Generate images by  $r$  and  $r'$ 
8:   Evaluate  $r$  and  $r'$  using  $L$ , get fitness scores  $f$  and  $f'$ 
9:   if  $f' < f$  then
10:    Accept  $r'$  as the new current solution
11:   else
12:    Accept  $r'$  with probability  $\exp(-\frac{f'-f}{T})$ 
13:   end if
14:   Update temperature  $T = \alpha \times T$ 
15: end for
16: Return the final solution  $r$ 
```

---

456 Optimizing this multi-objective function poses significant
 457 challenges due to the discrete and high-dimensional nature
 458 of the search space. Traditional optimization methods, such
 459 as grid search or random search, are inefficient and may
 460 not effectively explore the vast combinations of rank values
 461 across different layers. Gradient-based methods are unsuitable
 462 because the rank values are discrete and non-differentiable.

463 To address these challenges, we use Simulated Annealing
 464 (SA) as our optimization algorithm, detailed in Algorithm 1.
 465 SA is a probabilistic technique well-suited for discrete
 466 optimization problems and can effectively explore large search
 467 spaces. It mimics the annealing process in metallurgy, where a
 468 material is heated and then slowly cooled to decrease defects
 469 and find a low-energy state.

470 In the context of *Superlora*, each potential combination
 471 of rank values across all block layers of the denoising U-
 472 Net in Stable Diffusion constitutes a solution, represented
 473 as a sequence  $[r_1, r_2, \dots, r_n]$ . SA effectively explores this
 474 discrete space by initiating with a high temperature ( $T_0 = 100$ ), allowing broad exploration of the solution landscape.
 475 At each iteration, it generates a new solution by making
 476 slight modifications to the current one, promoting diversity
 477 and avoiding premature convergence. We evaluate the fitness
 478 of both solutions using the objective function  $L$ , ensuring
 479 progress toward optimality. As the algorithm progresses, we
 480 decrease the temperature according to a cooling schedule with
 481 a rate of  $\alpha$ . This strategy allows SA to extensively explore
 482 the search space initially and progressively focus on refining

483 the solution, thus avoiding local minima and increasing the likelihood of finding the global optimum.  
 484

485 Upon obtaining a solution, we extract the corresponding LoRA modules from *Superlora* and integrate them into the base model 2. We then generate images with fixed parameters to evaluate performance  $L$ .

486 In CEL, we sequentially obtain a set of optimal LoRA modules for each entity and merge them into the base model. This seamless integration results in no additional computational overhead during inference and eliminates memory requirements for continual learning.

## IV. BENCHMARK DESIGN AND CONSTRUCTION

### A. Datasets Collection

495 We collect three sets as benchmarks for CEL. The benchmarks cover general objects, human faces, and artistic styles. **CEL-Objects** comprises 15 common objects, each represented by 4-6 images, sourced from [7], [38]. The dataset includes various categories such as buildings, animals, and household items. The benchmark sequence consists of 'barn', 'cat', 'cat statue', 'clock', 'colorful teapot', 'dog', 'elephant', 'mug skulls', 'physics mug', 'red teapot', 'round bird', 'teddybear', 'thin bird', 'tortoise plushy', and 'wooden pot'.  
**CEL-Faces** is derived from the high-resolution CelebA-HQ dataset [89]. We select seven identities, each with a minimum of 20 samples, based on the highest number of available images. These celebrity identities serve as unique conceptual entities for our incremental learning task. The continual learning sequence comprises '#1 woman', '#2 man', '#3 man', '#4 man', '#5 man', '#6 woman', and '#7 woman'.

500 **CEL-Styles** features 1,207 images representing 6 distinct artistic styles, collected from Wikiart [90]. The dataset showcases works by renowned artists representing various art movements, including Neo-impressionism, Impressionism, Post-Impressionism, Mathematical art, Russian Revivalism, and Modernism. Each image is accompanied by a descriptive caption. The dataset is split into training and validation sets using an 8:2 ratio for each artist. The sequence of continual learning comprises 'Georges Seurat', 'Konstantin Korovin', 'Maurice Prendergast', 'M.C. Escher', 'Boris Kustodiev', and 'Marc Chagall'.

### B. Metrics for Evaluation

524 The CEL task encompasses the requirements of both continuous learning and T2I generation. Existing metrics for 525 both of these two tasks are insufficient for comprehensive 526 evaluation. To comprehensively assess the performance of 527 CEL, we integrate the T2I metrics into continuous learning 528 scenarios to form a (single-stage) performance indicator and 529 then derive two CEL performance metrics.  
 530

531 The Entity Generation Performance Indicator (EGI, denoted 532 as  $\rho$ ) measures the performance during a single learning 533 session. We define  $\rho$  as the *average* of three commonly used 534 T2I generation metrics [3], [6], [38]: the DINO Score, the 535 CLIP-I Score, and the CLIP-T Score.  
 536

537 The DINO and CLIP-I scores assess visual similarity, 538 measuring the similarity between features of the generated

539 images and the target images using DINO [91] and CLIP [32],  
 540 respectively. The CLIP-T Score evaluates textual similarity,  
 541 defined as the average cosine similarity between the CLIP  
 542 features of the generated images and their corresponding text  
 543 prompts. Formally,  $\rho_i^j$  denotes the performance of the model  
 544 on the  $i$ -th task after learning  $j$  sessions. Based on  $\rho$ , we define  
 545 two metrics:

546 **Incremental Concept Accuracy (ICA)** measures the genera-  
 547 tion performance of new entities after the model has completed  
 548 all continual learning sessions. It is defined as:

$$549 \quad ICA = \frac{1}{S} \sum_{i=1}^S \rho_i^S, \quad (5)$$

549 where  $S$  is the total number of continual learning sessions.

550 **Incremental Concept Forgetting (ICF)** quantifies the de-  
 551 crease in performance on previously learned concepts after  
 552 learning new ones. It is defined as:

$$553 \quad ICF = \frac{1}{S-1} \sum_{i=1}^{S-1} |\rho_i^i - \rho_i^S|, \quad (6)$$

553 where  $\rho_i^i$  is the performance on the  $i$ -th task immediately after  
 554 learning it, and  $\rho_i^S$  is the performance on the  $i$ -th task after  
 555 all sessions have been completed.

## 556 V. EXPERIMENTS

557 In this section, we present a comprehensive evaluation of  
 558 our proposed CEAL approach for continual concept learning,  
 559 encompassing quantitative and qualitative analysis.

### 560 A. Evaluation Setup

561 We implemented baselines of continual learning methods  
 562 for CEL task: generative replay (Replay), EWC [61] and C-  
 563 LoRA [82]. Additionally, we report the performance of state-  
 564 of-the-art personalization methods, including Textual Inversion  
 565 (TI) [7], DreamBooth (DB) [6], ED-LoRA [40] and Custom  
 566 Diffusion (CD) [38]. We also compare our method with  
 567 LoRA [15], using a rank of  $r = 8$  for each concept. We train  
 568 independent LoRA adapters for each entity and keep them in  
 569 memory, which serves as a potential upper bound for CEL  
 570 and offers additional context for our results. For CEAL, we  
 571 set  $r_{final} = 8$ , which means that the maximum rank obtained  
 572 by CEAL is consistent with that of LoRA with  $r = 8$ .

### 573 B. Implementation Details

574 We conducted our experiments using Stable Diffusion  
 575 v1.5 [3] as the base model. The training process utilized the  
 576 Mean Squared Error (MSE) loss function, and the LoRA was  
 577 integrated into all attention layers 16 of the denoising U-Net  
 578 architecture. *Superlora* begins with an initial rank of 1, which  
 579 is progressively incremented by 1 after every 1,000 training  
 580 iteration, resulting in a final rank ( $r_{final}$ ) of 8 after 8,000  
 581 iterations. Training was carried out with a batch size of 4. For  
 582 optimization, we used the AdamW optimizer with a learning  
 583 rate of  $1 \times 10^{-4}$  and a weight decay of 0.01. All input images  
 584 were preprocessed to a resolution of  $512 \times 512$  pixels.

585 For the search optimizer, we employ Simulated Annealing  
 586 (SA), as presented in Algorithm 1. In SA, a solution is a  
 587 sequence of rank values, where each value is selected for a  
 588 corresponding layer in *Superlora*. The length of each solution  
 589 is set to 16, corresponding to the 16 block layers in the  
 590 denoising U-Net of SD. We initialize the SA parameters with  
 591 an initial temperature ( $T_0$ ) of 100, a cooling rate ( $\alpha_{SA}$ ) of  
 592 0.95, and a total number of iterations ( $N_I$ ) of 1000. Our  
 593 search space is defined as the set  $\{0, 1, 2, 4, 6, 8\}$ , meaning that  
 594 each element in a solution can take on one of these values.  
 595 The SA algorithm explores this designated space to find the  
 596 optimal solution. The fitness function of SA is determined  
 597 by the objective function  $L$ . Given a specific solution, we  
 598 first extract the LoRA from *Superlora* with the corresponding  
 599 ranks. Then we merge the extracted LoRA into the SD model  
 600 and generate several images (specifically, 4 images per prompt,  
 601 using prompts listed in the Supplementary Material). These  
 602 images have dimensions of  $512 \times 512$  pixels, a guidance scale  
 603 of 7.5, and 50 sample steps. The seeds for sampling are also  
 604 fixed. In our main results, we simply set the hyperparameters  
 605  $\alpha$ ,  $\beta$ , and  $\gamma$  in the objective function to 100, 0.00001, and 0.1,  
 606 respectively, across all benchmarks. The purpose of this setting  
 607 is to normalize the magnitudes of these objective values to the  
 608 same scale, ensuring a balanced contribution from each term  
 609 during the optimization process. This serves as a standardized  
 610 baseline, and with minimal tuning, the search results can be  
 611 adapted to meet specific requirements.

### 612 C. Main Qualitative Results

613 To evaluate the effectiveness of our continual learning  
 614 approach, we conducted qualitative assessments for single-  
 615 concept fidelity and multi-concept compositionality.

616 Figure 4 presents a visual comparison of different methods  
 617 on the CEL-Faces dataset. The images shown are for the final  
 618 three identities learned after all methods completed the full  
 619 sequence of seven continual learning sessions. The results  
 620 clearly highlight the limitations of existing approaches. For  
 621 instance, the DB baseline suffers from complete mode col-  
 622 lapsed, failing to generate any coherent images. Other methods,  
 623 while avoiding outright collapse, exhibit significant forgetting,  
 624 producing faces that bear little resemblance to the target  
 625 identities. In stark contrast, our method (CEAL) successfully  
 626 generates high-fidelity images that are visually consistent  
 627 with the target subjects, demonstrating its superior ability to  
 628 mitigate catastrophic forgetting.

629 We further challenged the models with a multi-subject gen-  
 630 eration task on the CEL-Objects dataset, with results displayed  
 631 in Figure 5. This evaluation was performed after the models  
 632 had continually learned all 15 distinct entities. Although the  
 633 CEL-Objects dataset primarily consists of simple, distinct  
 634 items, and methods like TI and LoRA can achieve high-  
 635 fidelity representations for single entities, they face significant  
 636 challenges in multi-concept generation. This limitation likely  
 637 stems from their training strategies, which do not explicitly  
 638 optimize for the simultaneous generation of multiple entities.

639 In contrast, our proposed method, CEAL, excels in this  
 640 complex task. It not only generates high-quality images that



Fig. 4: **Visualizations of continual generation results on CEL-Faces.** We present visualizations of generated images from various methods for newly learned entities upon completion of all 7 continual learning sessions.

641 accurately depict newly learned entities but also seamlessly  
 642 integrates them with previously learned concepts. As shown,  
 643 our approach successfully captures the complex interactions  
 644 between concepts in a single frame, demonstrating remarkable  
 645 performance in multi-concept compositional generation. Ad-  
 646 dditional qualitative results are provided in the supplementary  
 647 material.

#### 648 D. Main Quantitative Results

649 Table I presents a comparative analysis of *ICA* and *ICF*  
 650 across three benchmarks: CEL-Objects, CEL-Faces, and CEL-  
 651 Styles. The primary focus of *ICA* is to evaluate the fidelity of  
 652 new entities after continual learning, whereas *ICF* measures  
 653 the degree of forgetting after learning new entities.

654 Methods such as Replay, EWC, and DB sequentially fine-  
 655 tune the denoising network without increasing the number of  
 656 parameters ( $\#P$ ). TI only requires adding a single token, so its  
 657  $\#P$  is almost negligible. Although these methods incorporate  
 658 knowledge replay and regularization to prevent forgetting,  
 659 they still exhibit inferior performance compared to parameter  
 660 expansion methods, as indicated by the increase in *ICF*.  
 661 While TI generally performs well on CEL-Faces and CEL-  
 662 Styles benchmarks with no forgetting problem, its performance  
 663 is suboptimal on the more challenging CEL-Objects dataset.  
 664 This may indicate that learning human faces and art styles  
 665 is noticeably easier for these methods than CEL-Objects,  
 666 possibly due to the T2I model having a more intuitive grasp  
 667 of these concepts. Our method significantly outperforms these  
 668 approaches without requiring additional parameters.

669 For methods that expand parameters for each new entity,  
 670 we report the average increase in parameters ( $\#P$ ). LoRA

671 increases parameters by  $2.59M$  per session, CD requires an  
 672 increase of  $19.17M$  parameters, and ED-LoRA needs  $0.81M$   
 673 per session. Although ED-LoRA and CD both employ LoRA  
 674 fusion strategies, their fusion is centralized, meaning that when  
 675 generating multiple concepts together, they merge LoRAs for  
 676 better results but must keep all independent LoRAs in memory.  
 677 C-LoRA merges LoRAs together but keeps separate from the  
 678 base model, resulting in a small parameter increase of  $0.78M$   
 679 per entity, as reported in their paper. CEAL significantly  
 680 outperforms these existing approaches without any additional  
 681 parameters. This success can be attributed to our innovative  
 682 multi-objective searching strategy, which strikes a remarkable  
 683 balance between performance enhancement and parameter  
 684 efficiency.

#### 685 E. Ablation Studies

686 Table II presents an ablation study on the impact of the main  
 687 components of CEAL in terms of *ICA*. We individually  
 688 remove each main component from CEAL to assess its influence.  
 689 Our full CEAL significantly outperforms the Naive LoRA  
 690 Merging baseline (48.12), which highlights the effectiveness  
 691 of CEAL’s overall design.

692 We first evaluate CEAL’s core architectural and training  
 693 strategies. Replacing the entire *Superlora* structure with a  
 694 standard fixed max-rank LoRA, upon which our search process  
 695 is then applied (CEAL w/o Superlora), reduces performance  
 696 to 58.95. This underscores Superlora’s inherent advantage in  
 697 facilitating the discovery of optimal layer-wise configurations.  
 698 Similarly, omitting the Dynamic Rank-Increasing Strategy dur-  
 699 ing Superlora training (CEAL w/o Dynamic Rank-Increasing)  
 700 results in an ICA score of 59.12, confirming this strategy’s

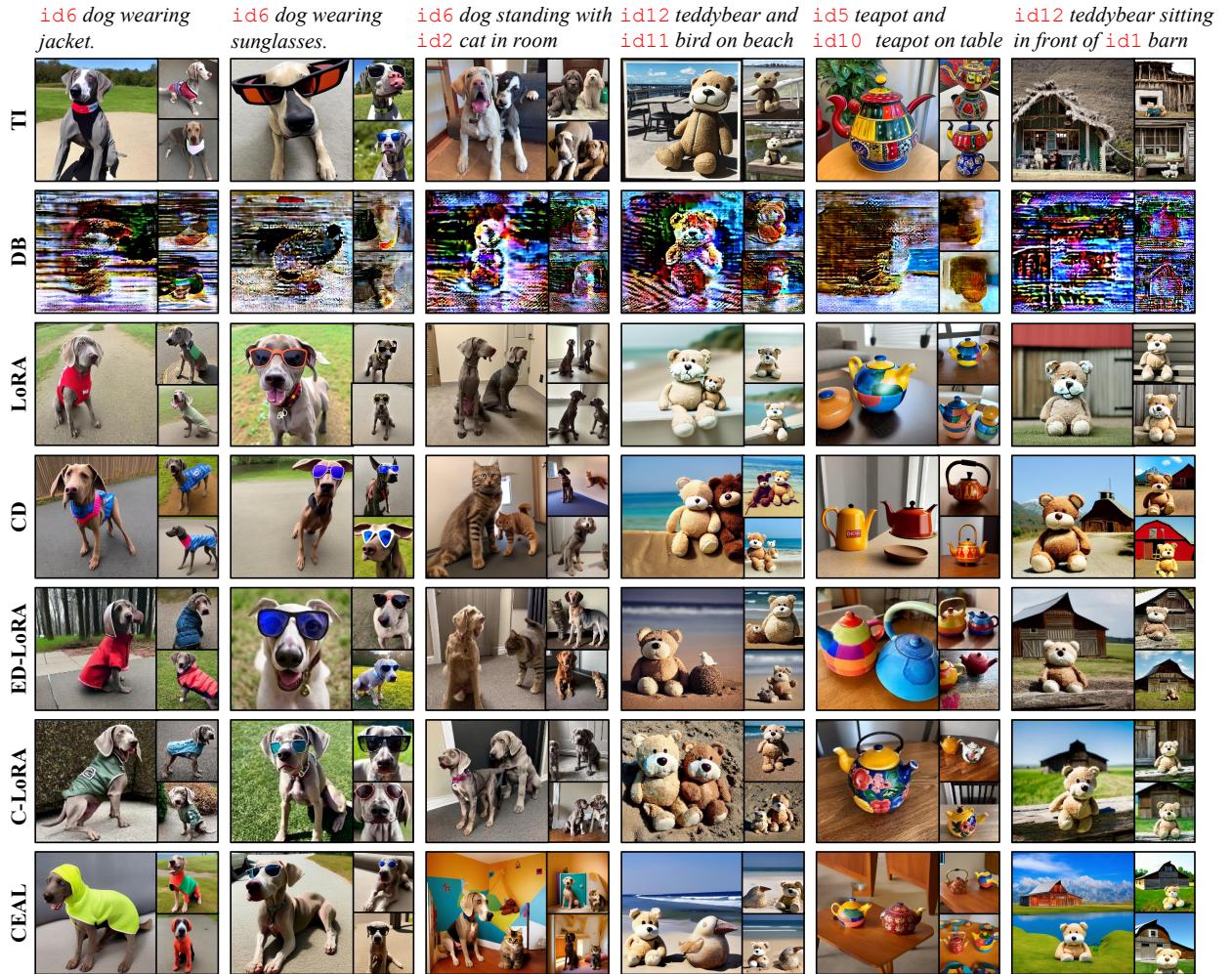


Fig. 5: **Visualization of multi-subject generation on CEL-Objects.** We show the results of multi-subject generation on CEL-Objects after all competing methods have completed 15 continual sessions. Entities are marked by their identifiers [ $\text{id}$ ], with their target samples provided at the bottom and prompts at the top.

contribution to effective LoRA learning. Constraining the search to a Uniform Rank across all layers, as opposed to distinct layer-wise ranks (CEAL w/o Layer-wise Ranks), not only leads to suboptimal performance with an ICA of 60.12 but also potentially increases parameter count due to less efficient allocation. Finally, substituting the SA Optimizer with a random search mechanism (CEAL w/o SA Optimizer) causes a substantial performance drop to 53.89, accompanied by high variance. This result affirms the critical role of Simulated Annealing in effectively and consistently navigating the complex search space of LoRA configurations.

Next, we individually remove each of the three search criteria—Human Preference Metric, Memory Cost Metric, and Concept Similarity Metric—and perform the search each time. The Concept Similarity Metric proves to be the most crucial component in the CEAL method, as removing it results in a substantial accuracy drop of 12.30. This is mainly because when only two metrics remain, the search process tends to select the LoRA with the smallest number of parameters, which may not align with the objectives of the CEL.

Finally, investigating the merging process, we found that

applying the optimized LoRAs independently without merging them into the base model (CEAL w/o Merging) achieves an  $\text{ICA}_{\text{new}}$  of  $65.25 \pm 0.76$ . This is slightly higher than the full CEAL with merging, suggesting a marginal trade-off between peak concept fidelity and the parameter efficiency gained by merging. However, the full CEAL's merging provides the significant benefit of zero additional parameters post-learning.

Then we study the impact of the **main** hyperparameters  $r_{\text{final}}$ ,  $\nu$ ,  $\alpha$ ,  $\beta$ , and  $\gamma$  in Fig. 6.  $r_{\text{final}}$  constrains the maximum LoRA rank that can be searched for each entity. The results show that the performance is poor when  $r_{\text{final}} = 1$ , indicating that the rank of LoRA may be too limited to learn LoRA sufficiently. Interestingly, the performance is suboptimal when  $r_{\text{final}} = 32$ . This result supports our claim that a large search space is difficult to train, and it is essential to design mechanisms that ensure *Superlora* is adequately trained and to develop search algorithms with better evaluation metrics.  $\nu$  determines the training steps required before adding a new rank to the *Superlora*. The results demonstrate that using 1000 steps as the increment is sufficient for achieving good performance, and the model's performance is not highly sensitive to the

TABLE I: Results for all benchmarks.

Datasets		CEL-Objects		CEL-Faces		CEL-Style	
Method	# $P$ (↓)	$ICA$ (↑)	$ICF$ (↓)	$ICA$ (↑)	$ICF$ (↓)	$ICA$ (↑)	$ICF$ (↓)
LoRA	2.59M	62.38±0.42	-	57.32±0.83	-	41.55±0.10	-
CD	19.17M	58.86±1.22	-	55.91±1.22	-	41.47±0.13	-
ED-LoRA	0.81M	55.35±0.83	-	58.77±0.44	-	40.33±0.20	-
C-LoRA	0.78M	61.94±0.30	6.12±0.49	57.48±0.29	2.83±0.23	41.45±0.48	1.03±0.09
TI	0.0	48.61±1.38	-	59.33±1.25	-	38.30±0.42	-
Replay	0.0	27.27±6.97	25.36±3.10	14.22±0.29	8.84±0.36	20.88±1.79	5.98±0.72
EWC	0.0	49.15±0.22	12.39±1.28	56.63±0.73	4.84±0.89	37.11±0.83	1.83±0.22
DB	0.0	32.58±5.08	38.07±2.72	49.02±0.30	4.61±0.62	41.99±0.40	1.01±0.29
CEAL	0.0	<b>63.74±0.31</b>	<b>9.12±0.88</b>	<b>59.01±0.73</b>	<b>2.71±0.72</b>	<b>45.21±0.53</b>	1.04±0.04

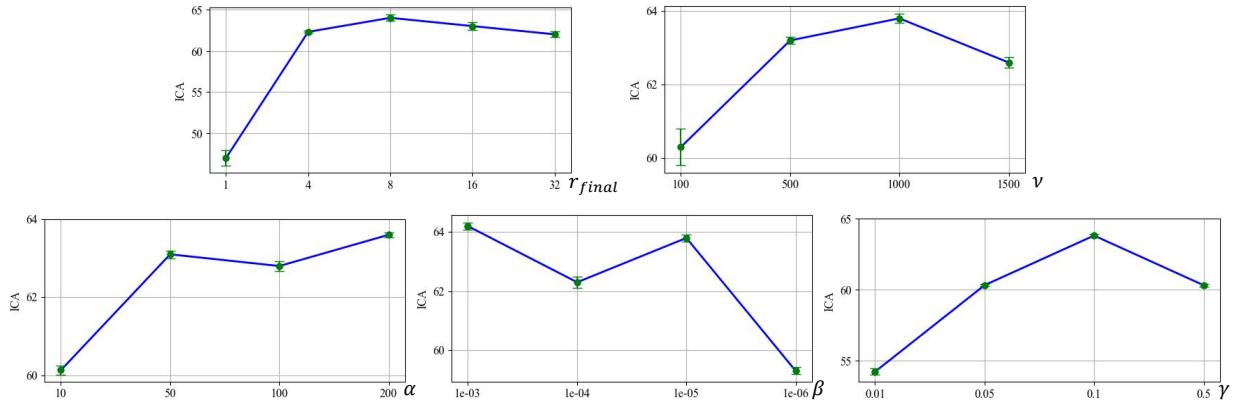
Fig. 6: Impact of  $r_{final}$ ,  $\nu$ ,  $\alpha$ ,  $\beta$ , and  $\gamma$  on CEL-Objects.

TABLE II: Ablation study on the main components of CEAL on CEL-Objects.

Component Variant	ICA <sub>new</sub>
Naive LoRA Merging	48.12±1.31
CEAL	63.74±0.31
CEAL w/o Superlora (Fixed Max-Rank LoRA + Search)	58.95±0.69
CEAL w/o Dynamic Rank-Increasing	59.12±0.69
CEAL w/o Layer-wise Ranks (Uniform Rank Search)	60.12±0.40
CEAL w/o SA Optimizer (e.g., Random Search)	53.89±10.34
CEAL w/o $L_{hpmp}$ (Human Preference Metric)	60.48±0.34
CEAL w/o $L_{mcm}$ (Memory Cost Metric)	62.89±0.75
CEAL w/o $L_{csm}$ (Concept Similarity Metric)	51.44±0.26
CEAL w/o Merging (Independent Optimized LoRAs)	65.25±0.76

choice of  $\nu$ . Even with smaller or larger step sizes, the model still maintains relatively good accuracy. For hyperparameters  $\alpha$ ,  $\beta$ , and  $\gamma$  that balance the searching objective, we find that although there are fluctuations and impacts, the variations are not significant.

#### F. Computational Costs

Training and optimizing *Superlora* for each concept requires roughly 2 gpu hours on a single Nvidia RTX 4090 GPU.

The subsequent search phase for optimal rank configurations, however, is considerably more time-consuming. The dominant factor in this cost is the image generation required to evaluate each candidate solution presented by the SA optimizer. Other operations, such as extracting specific LoRA modules from the *Superlora* backbone or calculating aesthetic and CLIP scores, have a negligible computational overhead in comparison. For each prompt used in the evaluation, we generate 4 images, a process that takes approximately 4.5

seconds on an NVIDIA RTX 4090. To illustrate, consider the CEL-Objects benchmark, where each concept is typically evaluated using 20 distinct prompts. With the SA optimizer set to a maximum of  $N_I = 1000$  iterations, evaluating one candidate solution (i.e., 20 prompts) takes  $20 \times 4.5\text{s} = 90\text{s}$ . This would lead to a maximum theoretical search duration of  $90\text{s}/\text{iteration} \times 1000 \text{ iterations} = 90,000 \text{ seconds}$ , or approximately 25 GPU hours per concept.

However, in practice, the SA algorithm often converges to a satisfactory solution well before reaching the maximum iteration limit. Furthermore, if the same rank configuration is proposed multiple times by the SA optimizer, its previously computed fitness score can be directly reused, avoiding redundant image generation. Due to these factors, the average search time per concept is significantly lower, typically around 4 GPU hours on the same hardware.

## VI. CONCLUSION

This paper introduces the concept of Continual Conceptual Entity Learning (CEL) for pre-trained T2I models. We introduce *Superlora*, a search space to organize potential LoRA configurations into a comprehensive set. To provide optimal customized LoRA rank configurations reasonably, we propose a novel dynamic rank-increasing strategy for training *Superlora*. We then use an efficient optimizer to continuously search for the optimal LoRA in line with multiple competing criteria. In addition, three comprehensive evaluation metrics and benchmarks for CEL are introduced. Extensive experiments on three benchmarks demonstrate the effectiveness of CEAL. However, it is worth noting that the proposed method may introduce additional computational overhead due to the search

process. On a positive note, our approach enhances parameter efficiency, potentially making AI models more accessible by reducing storage resource requirements.

Furthermore, evaluating CEAL on more diverse and challenging large-scale real-world datasets constitutes an important direction for future work to more broadly assess its generalization capabilities and practical applicability.

## ACKNOWLEDGMENTS

This work was funded in part by the National Natural Science Foundation of China (62376070, 62076195) and in part by the Fundamental Research Funds for the Central Universities (AUGA5710011522).

## REFERENCES

- [1] A. Q. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen, “Glide: Towards photorealistic image generation and editing with text-guided diffusion models,” in *International Conference on Machine Learning*. PMLR, 2022, pp. 16 784–16 804.
- [2] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, “Hierarchical text-conditional image generation with clip latents,” *arXiv preprint arXiv:2204.06125*, vol. 1, no. 2, p. 3, 2022.
- [3] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 684–10 695.
- [4] M. N. Everaert, M. Bocchio, S. Arpa, S. Süsstrunk, and R. Achanta, “Diffusion in style,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 2251–2261.
- [5] Y. Zhang, N. Huang, F. Tang, H. Huang, C. Ma, W. Dong, and C. Xu, “Inversion-based style transfer with diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 10 146–10 156.
- [6] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman, “Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 22 500–22 510.
- [7] R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-or, “An image is worth one word: Personalizing text-to-image generation using textual inversion,” in *The Eleventh International Conference on Learning Representations*, 2023. [Online]. Available: <https://openreview.net/forum?id=NAQvF08TcyG>
- [8] Q. Wang, X. Bai, H. Wang, Z. Qin, and A. Chen, “Instantid: Zero-shot identity-preserving generation in seconds,” *arXiv preprint arXiv:2401.07519*, 2024.
- [9] H. Shin, J. K. Lee, J. Kim, and J. Kim, “Continual learning with deep generative replay,” *Advances in neural information processing systems*, vol. 30, 2017.
- [10] J. Pomponi, S. Scardapane, and A. Uncini, “Continual learning with invertible generative models,” *Neural Networks*, vol. 164, pp. 606–616, 2023.
- [11] B. Yang, X. Deng, H. Shi, C. Li, G. Zhang, H. Xu, S. Zhao, L. Lin, and X. Liang, “Continual object detection via prototypical task correlation guided gating mechanism,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 9255–9264.
- [12] M. Zhai, L. Chen, F. Tung, J. He, M. Nawhal, and G. Mori, “Lifelong gan: Continual learning for conditional image generation,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 2759–2768.
- [13] Y. Zhang, X. Wang, and D. Yang, “Continual sequence generation with adaptive compositional modules,” in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022, pp. 3653–3667.
- [14] F. Mi, L. Chen, M. Zhao, M. Huang, and B. Faltings, “Continual learning for natural language generation in task-oriented dialog systems,” in *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2020, pp. 3461–3474.
- [15] E. J. Hu, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen *et al.*, “Lora: Low-rank adaptation of large language models,” in *International Conference on Learning Representations*, 2021.
- [16] Cloneofsimo, “LoRA: Using Low-rank Adaptation to Quickly Fine-tune Diffusion Models,” <https://github.com/cloneofsimo/lora>, 2023, accessed: 2024-02-19.
- [17] M. D. Zeiler and R. Fergus, “Visualizing and understanding convolutional networks,” in *Computer Vision-ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I 13*. Springer, 2014, pp. 818–833.
- [18] A. Voynov, Q. Chu, D. Cohen-Or, and K. Aberman, “p+: Extended textual conditioning in text-to-image generation,” *arXiv e-prints*, pp. arXiv-2303, 2023.
- [19] Z. Liu, R. Feng, K. Zhu, Y. Zhang, K. Zheng, Y. Liu, D. Zhao, J. Zhou, and Y. Cao, “Cones: Concept neurons in diffusion models for customized generation,” in *Proceedings of the 40th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, Eds., vol. 202. PMLR, 23–29 Jul 2023, pp. 21 548–21 566. [Online]. Available: <https://proceedings.mlr.press/v202/liu23j.html>
- [20] A. Aghajanyan, S. Gupta, and L. Zettlemoyer, “Intrinsic dimensionality explains the effectiveness of language model fine-tuning,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 2021, pp. 7319–7328.
- [21] J. Frankle and M. Carbin, “The lottery ticket hypothesis: Finding sparse, trainable neural networks,” in *International Conference on Learning Representations*, 2018.
- [22] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. N. Metaxas, “Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5907–5915.
- [23] M. Tao, B.-K. Bao, H. Tang, and C. Xu, “Galip: Generative adversarial clips for text-to-image synthesis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14 214–14 223.
- [24] A. Sauer, T. Karras, S. Laine, A. Geiger, and T. Aila, “Stylegan-t: Unlocking the power of gans for fast large-scale text-to-image synthesis,” in *International conference on machine learning*. PMLR, 2023, pp. 30 105–30 118.
- [25] J. Yu, Y. Xu, J. Y. Koh, T. Luong, G. Baid, Z. Wang, V. Vasudevan, A. Ku, Y. Yang, B. K. Ayan *et al.*, “Scaling autoregressive models for content-rich text-to-image generation,” *Transactions on Machine Learning Research*, 2022.
- [26] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, “Zero-shot text-to-image generation,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 8821–8831.
- [27] M. Ding, Z. Yang, W. Hong, W. Zheng, C. Zhou, D. Yin, J. Lin, X. Zou, Z. Shao, H. Yang *et al.*, “Cogview: Mastering text-to-image generation via transformers,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 19 822–19 835, 2021.
- [28] C. Zhang, W. Yang, X. Li, and H. Han, “Mmginpainting: Multi-modality guided image inpainting based on diffusion models,” *IEEE Transactions on Multimedia*, vol. 26, pp. 8811–8823, 2024.
- [29] B. Yuan, Y. Sheng, B.-K. Bao, Y.-P. P. Chen, and C. Xu, “Semantic distance adversarial learning for text-to-image synthesis,” *IEEE Transactions on Multimedia*, vol. 26, pp. 1255–1266, 2024.
- [30] Y. Huang, J. Huang, J. Liu, M. Yan, Y. Dong, J. Lv, C. Chen, and S. Chen, “Wavedm: Wavelet-based diffusion models for image restoration,” *IEEE Transactions on Multimedia*, vol. 26, pp. 7058–7073, 2024.
- [31] J. Pan, Z. Hu, Z. Su, and M.-H. Yang, “ $l_0$ -regularized intensity and gradient prior for deblurring text images and beyond,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 2, pp. 342–355, 2016.
- [32] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [33] G. Yuan, X. Cun, Y. Zhang, M. Li, C. Qi, X. Wang, Y. Shan, and H. Zheng, “Inserting anybody in diffusion models via celeb basis,” in *Advances in Neural Information Processing Systems*, A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, Eds., vol. 36. Curran Associates, Inc., 2023, pp. 72 958–72 982. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/e6d37cc5723e810b793c834cb6647cf-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/e6d37cc5723e810b793c834cb6647cf-Paper-Conference.pdf)
- [34] H. Xue, Z. Huang, Q. Sun, L. Song, and W. Zhang, “Freestyle layout-to-image synthesis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14 256–14 266.

- [35] L. Zhang, A. Rao, and M. Agrawala, "Adding conditional control to text-to-image diffusion models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 3836–3847.
- [36] J. Shi, W. Xiong, Z. Lin, and H. J. Jung, "Instantbooth: Personalized text-to-image generation without test-time finetuning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 8543–8552.
- [37] N. Ruiz, Y. Li, V. Jampani, W. Wei, T. Hou, Y. Pritch, N. Wadhwa, M. Rubinstein, and K. Aberman, "Hyperdreambooth: Hypernetworks for fast personalization of text-to-image models," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 6527–6536.
- [38] N. Kumari, B. Zhang, R. Zhang, E. Shechtman, and J.-Y. Zhu, "Multi-concept customization of text-to-image diffusion," 2023.
- [39] G. Kwon, S. Jenni, D. Li, J.-Y. Lee, J. C. Ye, and F. C. Heilbron, "Concept weaver: Enabling multi-concept fusion in text-to-image models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 8880–8889.
- [40] Y. Gu, X. Wang, J. Z. Wu, Y. Shi, Y. Chen, Z. Fan, W. Xiao, R. Zhao, S. Chang, W. Wu *et al.*, "Mix-of-show: Decentralized low-rank adaptation for multi-concept customization of diffusion models," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [41] Y. Yang, W. Wang, L. Peng, C. Song, Y. Chen, H. Li, X. Yang, Q. Lu, D. Cai, B. Wu *et al.*, "Lora-composer: Leveraging low-rank adaptation for multi-concept customization in training-free diffusion models," *arXiv preprint arXiv:2403.11627*, 2024.
- [42] J. Jiang, Y. Zhang, K. Feng, X. Wu, W. Li, R. Pei, F. Li, and W. Zuo, "Mc<sup>2</sup>: Multi-concept guidance for customized multi-concept generation," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 2802–2812.
- [43] M. Valipour, M. Rezagholizadeh, I. Kobyzhev, and A. Ghodsi, "Dylora: Parameter-efficient tuning of pre-trained models using dynamic search-free low-rank adaptation," in *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, 2023, pp. 3266–3279.
- [44] Q. Zhang, M. Chen, A. Bukharin, P. He, Y. Cheng, W. Chen, and T. Zhao, "Adaptive budget allocation for parameter-efficient fine-tuning," in *The Eleventh International Conference on Learning Representations*, 2022.
- [45] F. Zhang, L. Li, J. Chen, Z. Jiang, B. Wang, and Y. Qian, "Increlora: Incremental parameter allocation method for parameter-efficient fine-tuning," *arXiv preprint arXiv:2308.12043*, 2023.
- [46] Y. Ma, Z. Xie, J. Wang, K. Chen, and L. Shou, "Continual federated learning based on knowledge distillation," in *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence*, vol. 3, 2022.
- [47] A. Chaudhry, N. Khan, P. Dokania, and P. Torr, "Continual learning in low-rank orthogonal subspaces," *Advances in Neural Information Processing Systems*, vol. 33, pp. 9900–9911, 2020.
- [48] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "icarl: Incremental classifier and representation learning," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 2001–2010.
- [49] N. A. Cayco-Gajic and R. A. Silver, "Re-evaluating circuit mechanisms underlying pattern separation," *Neuron*, vol. 101, no. 4, pp. 584–602, 2019.
- [50] S. Hou, X. Pan, C. C. Loy, Z. Wang, and D. Lin, "Learning a unified classifier incrementally via rebalancing," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 831–839.
- [51] Y. Wu, Y. Chen, L. Wang, Y. Ye, Z. Liu, Y. Guo, and Y. Fu, "Large scale incremental learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 374–382.
- [52] B. Zhao, X. Xiao, G. Gan, B. Zhang, and S.-T. Xia, "Maintaining discrimination and fairness in class incremental learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 13208–13217.
- [53] H. Ahn, J. Kwak, S. Lim, H. Bang, H. Kim, and T. Moon, "Ssil: Separated softmax for incremental learning," in *Proceedings of the IEEE/CVF International conference on computer vision*, 2021, pp. 844–853.
- [54] A. Douillard, M. Cord, C. Ollion, T. Robert, and E. Valle, "Podnet: Pooled outputs distillation for small-tasks incremental learning," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI 16*. Springer, 2020, pp. 86–102.
- [55] H. Cha, J. Lee, and J. Shin, "Co2l: Contrastive continual learning," in *Proceedings of the IEEE/CVF International conference on computer vision*, 2021, pp. 9516–9525.
- [56] F.-Y. Wang, D.-W. Zhou, H.-J. Ye, and D.-C. Zhan, "Foster: Feature boosting and compression for class-incremental learning," in *European conference on computer vision*. Springer, 2022, pp. 398–414.
- [57] C. Wu, L. Herranz, X. Liu, J. Van De Weijer, B. Raducanu *et al.*, "Memory replay gans: Learning to generate new categories without forgetting," *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [58] Y. Cong, M. Zhao, J. Li, S. Wang, and L. Carin, "Gan memory with no forgetting," *Advances in Neural Information Processing Systems*, vol. 33, pp. 16481–16494, 2020.
- [59] X. Liu, C. Wu, M. Menta, L. Herranz, B. Raducanu, A. D. Bagdanov, S. Jui, and J. v. de Weijer, "Generative feature replay for class-incremental learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 226–227.
- [60] G. Petit, A. Popescu, H. Schindler, D. Picard, and B. Delezoide, "Fetril: Feature translation for exemplar-free class-incremental learning," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 3911–3920.
- [61] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska *et al.*, "Overcoming catastrophic forgetting in neural networks," *Proceedings of the national academy of sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.
- [62] F. Zenke, B. Poole, and S. Ganguli, "Continual learning through synaptic intelligence," in *International conference on machine learning*. PMLR, 2017, pp. 3987–3995.
- [63] L. Wang, M. Zhang, Z. Jia, Q. Li, C. Bao, K. Ma, J. Zhu, and Y. Zhong, "Afec: Active forgetting of negative transfer in continual learning," *Advances in Neural Information Processing Systems*, vol. 34, pp. 22379–22391, 2021.
- [64] M. M. Derakhshani, X. Zhen, L. Shao, and C. Snoek, "Kernel continual learning," in *International Conference on Machine Learning*. PMLR, 2021, pp. 2621–2631.
- [65] H. Ahn, S. Cha, D. Lee, and T. Moon, "Uncertainty-based continual learning with adaptive regularization," *Advances in neural information processing systems*, vol. 32, 2019.
- [66] S. Jung, H. Ahn, S. Cha, and T. Moon, "Continual learning with node-importance based adaptive group sparse regularization," *Advances in neural information processing systems*, vol. 33, pp. 3647–3658, 2020.
- [67] Z. Li and D. Hoiem, "Learning without forgetting," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 12, pp. 2935–2947, 2017.
- [68] P. Buzzega, M. Boschini, A. Porrello, D. Abati, and S. Calderara, "Dark experience for general continual learning: a strong, simple baseline," *Advances in neural information processing systems*, vol. 33, pp. 15920–15930, 2020.
- [69] J. Serra, D. Suris, M. Miron, and A. Karatzoglou, "Overcoming catastrophic forgetting with hard attention to the task," in *International conference on machine learning*. PMLR, 2018, pp. 4548–4557.
- [70] A. Mallya and S. Lazebnik, "Packnet: Adding multiple tasks to a single network by iterative pruning," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 7765–7773.
- [71] S. Ebrahimi, F. Meier, R. Calandra, T. Darrell, and M. Rohrbach, "Adversarial continual learning," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*. Springer, 2020, pp. 386–402.
- [72] A. Douillard, A. Ramé, G. Couairon, and M. Cord, "Dytox: Transformers for continual learning with dynamic token expansion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 9285–9295.
- [73] Y. Wang, Z. Ma, Z. Huang, Y. Wang, Z. Su, and X. Hong, "Isolation and impartial aggregation: A paradigm of incremental learning without interference," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 8, 2023, pp. 10209–10217.
- [74] S. Yan, J. Xie, and X. He, "Der: Dynamically expandable representation for class incremental learning," in *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Los Alamitos, CA, USA: IEEE Computer Society, jun 2021, pp. 3013–3022. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1109/CVPR46437.2021.00303>
- [75] H. Liu, K. Simonyan, and Y. Yang, "Darts: Differentiable architecture search," in *International Conference on Learning Representations*, 2018.

- 1086 [76] H. Cai, L. Zhu, and S. Han, "Proxylessnas: Direct neural architecture  
1087 search on target task and hardware," in *International Conference on*  
1088 *Learning Representations*, 2018.
- 1089 [77] Z. Wang, Z. Zhang, C.-Y. Lee, H. Zhang, R. Sun, X. Ren, G. Su, V. Perot,  
1090 J. Dy, and T. Pfister, "Learning to prompt for continual learning," in  
1091 *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
1092 *Pattern Recognition*, 2022, pp. 139–149.
- 1093 [78] Z. Wang, Z. Zhang, S. Ebrahimi, R. Sun, H. Zhang, C.-Y. Lee, X. Ren,  
1094 G. Su, V. Perot, J. Dy *et al.*, "Dualprompt: Complementary prompting for  
1095 rehearsal-free continual learning," in *European Conference on Computer*  
1096 *Vision*. Springer, 2022, pp. 631–648.
- 1097 [79] Y. Wang, Z. Huang, and X. Hong, "S-prompts learning with pre-trained  
1098 transformers: An ocam's razor for domain incremental learning," *Advances*  
1099 *in Neural Information Processing Systems*, vol. 35, pp. 5682–  
1100 5695, 2022.
- 1101 [80] J. S. Smith, L. Karlinsky, V. Gutta, P. Cascante-Bonilla, D. Kim,  
1102 A. Arbelle, R. Panda, R. Feris, and Z. Kira, "Coda-prompt: Continual  
1103 decomposed attention-based prompting for rehearsal-free continual  
1104 learning," in *Proceedings of the IEEE/CVF Conference on Computer*  
1105 *Vision and Pattern Recognition*, 2023, pp. 11909–11919.
- 1106 [81] J. S. Smith, Y.-C. Hsu, Z. Kira, Y. Shen, and H. Jin, "Continual diffusion  
1107 with stamina: Stack-and-mask incremental adapters," in *Proceedings of*  
1108 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
1109 2024, pp. 1744–1754.
- 1110 [82] J. S. Smith, Y.-C. Hsu, L. Zhang, T. Hua, Z. Kira, Y. Shen, and H. Jin,  
1111 "Continual diffusion: Continual customization of text-to-image diffusion  
1112 with c-lora," *Transactions on Machine Learning Research*.
- 1113 [83] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models,"  
1114 *Advances in neural information processing systems*, vol. 33, pp. 6840–  
1115 6851, 2020.
- 1116 [84] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models,"  
1117 in *International Conference on Learning Representations*, 2020.
- 1118 [85] Y. Frenkel, Y. Vinker, A. Shamir, and D. Cohen-Or, "Implicit style-  
1119 content separation using b-lora," in *European Conference on Computer*  
1120 *Vision*. Springer, 2024, pp. 181–198.
- 1121 [86] V. Shah, N. Ruiz, F. Cole, E. Lu, S. Lazebnik, Y. Li, and V. Jampani,  
1122 "Ziplora: Any subject in any style by effectively merging loras," in  
1123 *European Conference on Computer Vision*. Springer, 2024, pp. 422–  
1124 438.
- 1125 [87] shunk031, "simple-aesthetics-predictor: Clip-based aesthetics predictor  
1126 inspired by the interface of huggingface transformers." 2023. [Online].  
1127 Available: <https://github.com/shunk031/simple-aesthetics-predictor>
- 1128 [88] Y. Kirstain, A. Polyak, U. Singer, S. Matiana, J. Penna, and O. Levy,  
1129 "Pick-a-pic: An open dataset of user preferences for text-to-image  
1130 generation," *Advances in neural information processing systems*, vol. 36,  
1131 pp. 36 652–36 663, 2023.
- 1132 [89] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing  
1133 of gans for improved quality, stability, and variation," in *International*  
1134 *Conference on Learning Representations*, 2018.
- 1135 [90] B. Saleh and A. Elgammal, "Large-scale classification of fine-art paintings:  
1136 Learning the right metric on the right feature," 05 2015.
- 1137 [91] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and  
1138 A. Joulin, "Emerging properties in self-supervised vision transformers,"  
1139 in *Proceedings of the IEEE/CVF international conference on computer*  
1140 *vision*, 2021, pp. 9650–9660.

1141

## VII. BIOGRAPHY SECTION

1142  
1143  
1144  
1145  
1146  
1147

**Yabin Wang** received the Ph.D. degree from Xi'an Jiaotong University, P. R. China, in 2025. He is an Associate Researcher with Harbin Institute of Technology ( HIT). His research interests include continual learning, multi-modal learning and deepfake detection.

1148



**Xiaopeng Hong** (Senior Member, IEEE) received his Ph.D. degree from Harbin Institute of Technology ( HIT), P. R. China, in 2010. He is currently a professor at HIT. He has authored over 80 publications in leading journals and conferences, including IEEE T-PAMI, TIP, PIEEE, CVPR, ICCV, NeurIPS, AAAI, and ACM MM. His research on subtle facial movement analysis has been featured by international media such as MIT Technology Review. He is a co-author of the Top Paper Award at ACM Multimedia 2022 and the IEEE Finland Section Best Student Conference Paper in 2020. His current research interests include multi-modal learning, continual learning, and model interpretability.

1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161

**Zhiheng Ma** (Member, IEEE) received the PhD degree from Xi'an Jiaotong University, in 2021. He is a research assistant professor with the Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences (SIAT). He has authored articles in journals and conferences, such as IEEE Transactions on Image Processing, CVPR, ICCV, and AAAI. His current research interests include incremental learning, crowd counting, and novelty detection.

1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171

**Zhou Su** (Senior Member, IEEE) has authored or coauthored technical papers, including top journals and top conferences, such as the IEEE Journal on Selected Areas in Communications, IEEE Transactions on Information Forensics and Security, IEEE Transactions on Dependable and Secure Computing, IEEE Transactions on Mobile Computing, IEEE/ACM Transactions on Networking, and INFOCOM. His research interests include multimedia communication, wireless communication, and network traffic. He is an associate editor for the IEEE Internet of Things Journal, IEEE Open Journal of the Computer Society, and IET Communications.

1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184

**Jinpeng Zhang** received the Ph.D. degree from the Institute of Automation, Chinese Academy of Sciences, Beijing, China. He works for the Intelligent Science & Technology Academy Limited of CASIC. His research interests include deep learning and multi-modal learning.

1185  
1186  
1187  
1188  
1189  
1190

**Zhiwu Huang** received his Ph.D. in Computer Science and Technology from the University of Chinese Academy of Sciences in 2015. He then worked as a Postdoctoral and Guest Researcher in the Computer Vision Lab (CVL) at ETH Zurich from September 2015 to July 2021. From September 2021 to December 2022, he served as an Assistant Professor of Computer Science at Singapore Management University. He is currently a Lecture (Assistant Professor) in the Vision, Learning, and Control (VLC) research group at the University of Southampton. His research interests include Computer Vision, Machine Learning, Generative Artificial Intelligence, and Geometric Deep Learning.

1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204