# AUTOMATIC ATTACK DISCOVERY FOR FEW-SHOT CLASS-INCREMENTAL LEARNING VIA LARGE LANGUAGE MODELS

**Haidong Kang**[*]
School of Computer and Communication Engineering
Northeastern University
hdkang@stumail.neu.edu.cn

**Wei Wu**
School of Information and Communication
University of Electronic Science and Technology of China
202522010418@std.uestc.edu.cn

**Hanling Wang**
Department of Strategic and Advanced Interdisciplinary Research
Pengcheng Laboratory
wanghl03@pcl.ac.cn

December 4, 2025

## ABSTRACT

Few-shot class incremental learning (FSCIL) is a more realistic and challenging paradigm in continual learning to incrementally learn unseen classes and overcome catastrophic forgetting on base classes with only a few training examples. Previous efforts have primarily centered around studying more effective FSCIL approaches. By contrast, less attention was devoted to thinking the security issues in contributing to FSCIL. This paper aims to provide a holistic study of the impact of attacks on FSCIL. We first derive insights by systematically exploring how human expert-designed attack methods (i.e., PGD, FGSM) affect FSCIL. We find that those methods either fail to attack base classes, or suffer from huge labor costs due to relying on huge expert knowledge. This highlights the need to craft a specialized attack method for FSCIL. Grounded in these insights, in this paper, we propose a simple yet effective ACraft method to automatically steer and discover optimal attack methods targeted at FSCIL by leveraging Large Language Models (LLMs) without human experts. Moreover, to improve the reasoning between LLMs and FSCIL, we introduce a novel Proximal Policy Optimization (PPO) based reinforcement learning to optimize learning, making LLMs generate better attack methods in the next generation by establishing positive feedback. Experiments on mainstream benchmarks show that our ACraft significantly degrades the performance of state-of-the-art FSCIL methods and dramatically beyond human expert-designed attack methods while maintaining the lowest costs of attack.

## 1 Introduction

Deep neural networks have demonstrated remarkable success across various computer vision applications, showing superior performance in classification, detection, and segmentation tasks. However, this success typically relies on the assumption that the training and testing categories remain identical, and that models will not encounter new classes
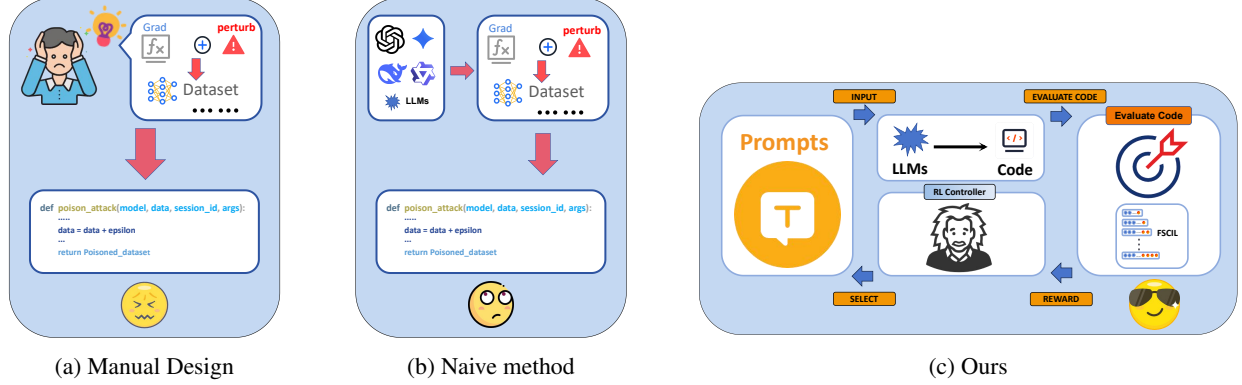
---

[*]Corresponding author

Figure 1: A comparison of the designed way of Attack. (a) Manual design relies on expert knowledge. (b) A naive method via LLMs. (c) Our method proposes an LLM-driven ACraft framework to automatically discover optimal Attack for FSCIL.

after deployment. This assumption often fails in real-world settings, where models must adapt to new categories over time, leading to catastrophic forgetting and degraded generalization. To mitigate this, Class-Incremental Learning (CIL) [8, 21] has emerged as a promising paradigm, enabling models to acquire new knowledge while retaining previously learned information.

Despite the progress of CIL, conventional paradigms still assume that each incremental task is supported by abundant labeled data, which is unrealistic in many applications such as robotic interaction, medical diagnosis updates, or recommendation systems. In such domains, new categories often emerge as few-shot instances, where only a handful of samples per class are available. This motivates the development of Few-Shot Class-Incremental Learning (FSCIL) [6, 12, 14, 19, 27], which aims to bridge the gap between continual learning and data scarcity. Although recent advances have led to promising performance, the security vulnerabilities of FSCIL remain largely unexplored, posing significant threats in safety-critical applications.

## 1.1 Challenges

Although recent works have made initial attempts to explore adversarial attacks in the FSCIL setting, these methods still suffer from several fatal drawbacks: (1) They require extensive expert knowledge through time-consuming manual design and tuning, which dramatically increases the cost and limits the scalability of attack algorithm development; (2) They are usually tailored for traditional class-incremental learning, where the model remains relatively stable. In contrast, FSCIL models are highly non-stationary, as their decision boundaries continuously evolve with each incremental session. Consequently, perturbations optimized for one state of the model quickly become obsolete in subsequent states, leading to inconsistent attack efficacy; (3) They fail to account for the inherent data sparsity in few-shot settings, resulting in biased or unstable attacks due to unreliable gradient estimation.

These limitations highlight the need to rethink the design paradigm of attack algorithms for FSCIL and explore new directions beyond the traditional handcrafted framework.

**Our New Observation.** Different from existing manually designed methods (as shown in Fig. 1b and Fig. 1c), we observe a new way to automatically design adversarial attack algorithms via Large Language Models (LLMs) in this work. This paradigm fundamentally disrupts the traditional manual design process by allowing LLMs to autonomously generate and refine attack strategies, significantly reducing human effort while improving adaptability to dynamic FSCIL scenarios. More details are illustrated in Section 3.

## 1.2 Contributions

In this work, we aim to analyze and overcome the aforementioned drawbacks by introducing a new paradigm of automated adversarial attack generation. To achieve this goal, we first explore a naive method that uses simple prompts as input to LLMs [17, 26, 29]. However, we find that this one-shot prompting paradigm performs poorly in FSCIL attack generation, as it degenerates into a black-box process with no interaction between the LLM and the target model. This observation raises a critical question: how can LLMs be effectively guided to generate more powerful and adaptive attack algorithms?

Table 1: **ACraft *v.s.* previous methods.** $T$ denotes the total number of attack iterations and $d$ represents the input dimensionality. $Cost(\cdot)$ indicates the asymptotic computational complexity of each algorithm with respect to $T$ and $d$.

| Method | Formula | Human Expert | Target | Attack Acc↓ | Cost (Algorithm) |
|---|---|---|---|---|---|
| FGSM [9] | $x_{adv} = \text{Clip}_{0,1}\left(x + \epsilon \cdot \text{sign}\left(\nabla_x \mathcal{L}(f(x), y)\right)\right)$ | ✓ | Image Classification | 58.01 | $O(d)$ |
| PGD [15] | $x_{adv}^{t+1} = \Pi_{x,\epsilon,[0,1]}\left(x_{adv}^t + \alpha \cdot \text{sign}\left(\nabla_x \mathcal{L}(f(x_{adv}^t), y)\right)\right)$ | ✓ | Image Classification | 58.00 | $O(T \cdot d)$ |
| C&W [5] | $\min_w \left\|\frac{1}{2}(\tanh(w)+1) - x\right\|_2^2 + c \cdot f\left(\frac{1}{2}(\tanh(w)+1)\right)$ | ✓ | Image Classification | 60.10 | $O(T \cdot d)$ |
| DeepFool [16] | $\delta^{(t)} = \frac{\left|f_k(x^{(t)}) - \max_{j \neq k} f_j(x^{(t)})\right|}{\left\|\nabla_x f_k(x^{(t)}) - \nabla_x f_l(x^{(t)})\right\|_2^2}\left(\nabla_x f_k(x^{(t)}) - \nabla_x f_l(x^{(t)})\right)$ | ✓ | Image Classification | 60.10 | $O(T \cdot d)$ |
| **ACraft** | $x_{adv} = \text{clip}\left(\text{clip}\left(x + \sum_{t=0}^{T-1}\left[-\alpha \cdot \text{sign}\left(\mu \cdot m^{(t)} + (1-\mu) \cdot \lambda_{\text{rev}} \cdot \nabla_x \mathcal{L}_{\text{comb}}^{(t)}\right)\right], x - \epsilon, x + \epsilon\right), 0, 1\right)$ | ✗ | FSCIL | 17.13 | $O(T \cdot d)$ |



Figure 2: ACraft *v.s.* existing attacks for FSCIL task in CUB200 dataset.

To address this challenge, inspired by actor-critic reinforcement learning and Chain-of-Thought (CoT) prompting [25], we propose a Proximal Policy Optimization (PPO)-based iterative framework that enables continuous refinement of LLM-generated attack algorithms. Specifically, the LLM first generates a population of diverse attack candidates, each of which is evaluated through a task-specific fitness function. A PPO-based controller then stabilizes policy updates through a clipping mechanism, producing improved prompts to guide the next generation of algorithm evolution. This closed-loop optimization fosters a positive feedback cycle between the LLM and the FSCIL task, leading to progressively stronger and more generalizable attack algorithms across iterations.

To ensure fairness and generality, we adopt a widely recognized FSCIL benchmark that supports training-time poisoning attacks, enabling comprehensive and standardized evaluation. Extensive experiments show that our proposed framework not only discovers novel and highly effective attack schemes but also outperforms expert-designed methods under the same computational constraints.

We summarize our main contributions as follows:

- **New Attack Generation Paradigm.** To the best of our knowledge, we are the first to propose a novel LLM-driven framework for automatically designing adversarial attacks in FSCIL, offering a fresh perspective for understanding and improving model robustness in continual few-shot scenarios.

- **PPO-driven Evolutionary Strategy.** Beyond the limitations of simple prompting, we reveal that the absence of feedback signals is the root cause of poor attack quality. To overcome this, we introduce a PPO-based reinforcement learning mechanism that builds a reasoning-driven optimization loop, significantly enhancing the coherence and effectiveness of LLM-generated attacks.

- **Numerical Verification.** Comprehensive experiments on standard FSCIL benchmarks demonstrate that our automated framework consistently outperforms manually designed attack methods, establishing a new state-of-the-art in automated adversarial algorithm discovery.

## 2 Rethinking the design of Attack Algorithms

The primary objective of traditional attack algorithms is to efficiently attack general models. As shown in Table 1, representative attack algorithms (such as FGSM, PGD, and C&W) can utilize heuristic or statistical methods to measure their expressiveness. However, these methods heavily rely on specialized expert knowledge, making them suboptimal when applied to specific models. Additionally, the design of such algorithms is particularly time-consuming. As illustrated in Table 1, the attack algorithms manually designed by experts exhibit extremely low effectiveness when used for FSCIL attacks. Moreover, to validate the statement that existing attack methods either fail to attack base
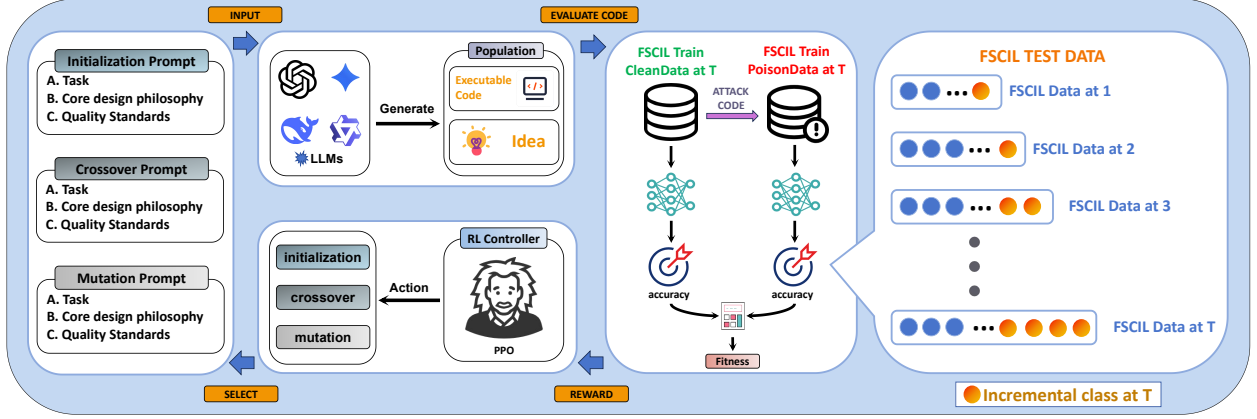
Figure 3: Overview of ACraft.

classes, we conduct an in-depth analysis in the CUB-200 dataset tailored for the FSCIL task. As depicted in Fig. 2, we can clearly observe that existing attack methods (e.g., PGD, DeepFool, etc.) indeed do not work on base classes. In addition, "C&W" and "DeepFool" fail to attack on new classes. By contrast, our Arcraft significantly attacks the base and new classes, leading to an approximate 70%, and 45% decrease in terms of accuracy, respectively.

Grounding in those findings, we analyze that those issues may be related to the task of FSCIL, where the small perturbations independently optimized for a few samples attacked by methods (e.g., C&W) are offset during the incremental update of the FSCIL model, resulting in a poisoning prototype that is almost identical to the clean prototype and unable to generate destructive parameter updates. The low performance of these manually crafted adversarial attack algorithms for FSCIL prompted us to seek a new paradigm for designing attack methods specifically tailored for FSCIL.

Inspired by the success of LLMs in generating new knowledge, we pose a natural question: **Can LLMs automatically generate adversarial attack algorithms for FSCIL?** To rigorously answer this question, we first designed a Naive method (as shown in **App. A**), asking the LLM to output a reasonable attack idea along with the corresponding implementation code. To validate the effectiveness proposed naive method, we conduct an in-depth analysis based on CLOSER [18] in the mainstream miniImageNet dataset. As shown in the Table. 4, the generated code reduced model accuracy by an average of 19.10% under four runs, which still remains lower than manually designed methods (i.e., PGD: 4.76% *v.s.* Naive: 19.10%). Although the results are promising, the Naive method is still far from the expectation of attacking for FSCIL. Further analysis demonstrates that the limitation of the Naive method arises from the absence of feedback from FSCIL: attack code generation conditioned only on prompts forms a blackbox process. To address this limitation, we propose an iterative optimization algorithm based on PPO reinforcement learning. Our framework incorporates a feedback mechanism that enables closed-loop optimization, allowing the LLM to iteratively refine the generated attack algorithms. This significantly improves the performance of ACraft (as shown in Fig. 2).

## 3 Automatic Attack Discovery for FSCIL

### 3.1 Overview of ACraft

To generate novel attack algorithms for FSCIL, ACraft leverages LLMs as a powerful tool to produce new attack strategies along with their corresponding code. However, LLMs can sometimes behave like black boxes. To address this issue, ACraft integrates the PPO reinforcement learning algorithm, which provides continuous feedback on the effectiveness of generated attacks and uses this feedback to guide the LLM's generation process intelligently. The ACraft system consists of three cooperating components: the Adversarial Algorithm Generator, the Attack Efficacy Evaluator, and a PPO-based Evolution Controller.

**Adversarial Algorithm Generator.** In the ACraft framework, the Large Language Model (LLM) is tasked with generating novel attack algorithms, guided by a series of carefully designed prompts. These prompts instruct the LLM to either refine an existing attack method or synthesize features from multiple methods to create a more effective one. Unlike typical code generators, each LLM output consists of two interconnected components for every attack algorithm. First, it produces a Strategic Rationale in natural language, explaining the core logic of the attack. Second, it generates the executable Algorithmic Code. This two-part output is crucial because it ensures that the generated method is interpretable and allows us to clearly understand how and why each attack works.

4

**Attack Efficacy Evaluator.** The primary task of this evaluator is to quantify the effectiveness of each generated attack algorithm. To this end, it performs attacks on a target FSCIL model and measures their impact on the model's performance. These measurements are then aggregated into a single fitness score, a comprehensive metric that primarily accounts for both the extent of damage caused to the model's performance and the computational cost of the attack. Ultimately, this fitness score serves as a crucial feedback signal to guide the entire optimization process.

**PPO-based Evolution Controller.** The PPO controller is the decision-maker for the whole evolution process. Its primary task is to learn the optimal policy to guide the generation process. In this way, the evolution can find the best attack algorithm. To this end, the controller uses the fitness score from the evaluator as a reward signal. It then employs the PPO algorithm to learn the best strategy. This strategy helps it decide the next action: whether to instruct the generator to create a new algorithm or to fine-tune an existing algorithm that shows promise.

## 3.2 Mathematical Formulation

We consider a standard FSCIL scenario. In this setting, a model $\mathcal{M}$ is trained on a series of incremental tasks $\mathcal{T} = \{\mathcal{T}_1, \mathcal{T}_2, ..., \mathcal{T}_N\}$ that arrive in order. Our goal is to design a general attack method. This method works by adding poisoned samples to the training data $\mathcal{D}_{train}^{\mathcal{T}_i}$ of each task $\mathcal{T}_i$. In this way, it continuously damages the model's performance. We structure the collection of all possible attack algorithms as a latent design space $\mathcal{G}$ encoded by LLM. Each attack algorithm $\in \mathcal{G}$ from this space is represented as a set of two-components. The first part of this set is a conceptual description $R$ of the attack algorithm that the LLM designs. It explains the logic of the attack, and then LLM translates this conceptual description into a fully executable implementation $C$. This implementation forms the second part of the set.

To measure the effectiveness of any algorithm $g$, we define an objective function $\mathcal{L}_{attack} : \mathcal{M} \times \mathcal{D} \to \mathbb{R}$. We use $P(M_\theta, \mathcal{D})$ to show how well a model with parameters $\theta$ performs on a dataset $\mathcal{D}$. For example, this performance can be the accuracy. Next, we define two models. We let $M_{clean}$ be the model that trains on the original, clean data. In contrast, $M_g$ is the model that trains on the poisoned data. Therefore, the objective function $\mathcal{L}_{attack}$ is the performance difference between these two models:

$$\mathcal{L}_{attack}(M_g, \mathcal{D}_{test}) = P(M_{clean}, \mathcal{D}_{test}) - P(M_g, \mathcal{D}_{test})$$

Formally, let $M_g$ represent the model after the original model $M$ undergoes poisoned training, as specified by $C$, the implementation of algorithm $g$. This implementation specifies the poisoned training procedure. The final goal of the ACraft framework is to find the best attack algorithm $g^*$. This algorithm aims to maximize the expected effectiveness of the attack on the distribution of the test data. This optimization problem can be represented as:

$$g^* = \arg \max_{g \in \mathcal{G}} \mathbb{E}_{\mathcal{D}_{test}}[\mathcal{L}_{attack}(M_g, \mathcal{D}_{test})]$$

Here, the expectation $\mathbb{E}$ shows that we are looking for an attack strategy that can generalize well. The entire ACraft framework is designed specifically to solve this optimization problem effectively.

## 3.3 Core Mechanisms of Attack Generation

**Algorithm Generator.** We model the Algorithm Generator as a set of discrete transformation functions $\mathcal{F}_{trans}$. These functions make up the executable operation space used by the generator when creating attack algorithms. Theoretically, these functions are specific, controllable mappings that the generator performs on its internally encoded latent design space $\mathcal{G}$.

First, we define $\Pi$ as the set of all possible instructions. Each function $F \in \mathcal{F}_{trans}$ has a fixed structure. Its main role is to turn the input context $\mathcal{C}$ into executable code $g'$:

$$F : \mathcal{C} \to \mathcal{G}$$

where $\mathcal{C}$ is a tuple made of all possible input information for the generator. It consists of an instruction subset $\Pi_{sub} \subseteq \Pi$ and a parent algorithm subset $\mathcal{G}_{sub} \subseteq \mathcal{G}$:

$$\mathcal{C} = (\Pi_{sub}, \mathcal{G}_{sub})$$

The set of transformation functions $\mathcal{F}_{trans}$ includes the following types:

- **Genesis Function** $F_{init}$**:** This function creates entirely new algorithms from scratch. In this case, the parent set $\mathcal{G}_{sub}$ is empty $\emptyset$. Therefore, the input context $\mathcal{C}$ relies only on the initial instruction $\Pi_{init} \in \Pi_{sub}$.

$$g' = F_{init}(\Pi_{init}, \emptyset)$$

  When the Algorithm Generator runs this function, its goal is to maximize the diversity measure $D(\mathcal{G}_{new})$ of the newly created subset $\mathcal{G}_{new}$. This ensures that the search is not limited to a narrow part of $\mathcal{G}$.

- **Refinement Function** $F_{refine}$**:** This function adjusts an existing successful algorithm $g \in \mathcal{G}$ using fine-tuning and local optimization (exploitation). This helps the Generator reach higher fitness peaks in the strategy space. Its input context $\mathcal{C}$ includes a single parent algorithm $\{g\} \subseteq \mathcal{G}$ and a specific optimization instruction $\Pi_{opt} \in \Pi_{sub}$.

$$g' = F_{refine}(\Pi_{opt}, \{g\})$$

When the Generator runs this function, it must ensure that the new executable code $g'$ maintains a high correlation with $g$ in its strategy: $Corr(g', g) \approx 1$ while simultaneously achieving a performance improvement: $\phi(g') > \phi(g)$.

- **Synthesis Function** $F_{synth}$**:** This function combines the best strategies from multiple parent algorithms. It encourages strategy mixing and breakthrough innovation (exploration). Its input context $\mathcal{C}$ includes a set of parent algorithms $\{g_1, g_2, ..., g_k\} \subseteq \mathcal{G}$ and a merging instruction $\Pi_{merge} \in \Pi_{sub}$.

$$g' = F_{synth}(\Pi_{merge}, \{g_1, g_2, ..., g_k\})$$

The main goal when the Generator runs this function is to maximize the expected effectiveness of the new executable code: $\mathbb{E}[\phi(g')]$. This particularly targets the nonlinear gain that results from combining $k$ different strategies.

By making the LLM's generation capability into this discrete and controllable function set $\mathcal{F}_{trans}$, we successfully change the traditional problem of black-box code generation into a decision problem that a reinforcement learning strategy can clearly select and manage. This function-based modeling not only makes the theoretical analysis and tuning of the algorithm generation process possible, but also ensures that the ACraft framework can perform a systematic, intelligent evolutionary search among diversity injection $F_{init}$, local optimization $F_{refine}$, and strategy breakthrough $F_{synth}$.

**Fitness Evaluator as a Multi-Attribute Utility Function.** The primary job of the Evaluator is to calculate a fitness score $\phi(g)$ for each algorithm $g$. This score is modeled as a Multi-Attribute Utility Function, balancing attack effectiveness and cost to reflect the algorithm's real-world value. First, we define an objective vector $\vec{J}(g)$. This vector captures all performance metrics that are important to us:

$$\vec{J}(g) = [J_{succ}(g), J_{cost}(g)]^T$$

where $J_{succ}(g)$ is further defined as a weighted sum of performance drops in new and old tasks. This shows how well the attack strategy targets catastrophic forgetting:

$$\begin{aligned} J_{succ}(g) = {} & \alpha \cdot (P(M_{clean}, \mathcal{D}_{old}) - P(M_g, \mathcal{D}_{old})) \\ & + (1 - \alpha) \cdot (P(M_{clean}, \mathcal{D}_{new}) - P(M_g, \mathcal{D}_{new})) \end{aligned}$$

where $J_{cost}(g)$ stands for the comprehensive cost, which includes the computational expense and the perturbation budget (constraining stealthiness). $\alpha \in [0, 1]$ is a hyperparameter used to adjust how much we prioritize the forgetting of old knowledge. To obtain a single fitness score $\phi(g)$, we perform a linear scalarization of the objective vector. We use a preset weight vector $\vec{w} = [w_{succ}, w_{cost}]^T$, where $w_{cost}$ is a negative value. The fitness is then calculated as:

$$\phi(g) = \vec{w} \cdot \vec{J}(g) = w_{succ}J_{succ}(g) + w_{cost}J_{cost}(g)$$

The weight vector $\vec{w}$ encodes our preference among objectives, ensuring that the search for the optimal algorithm $g^*$ balances high attack effectiveness with low execution cost.

**The Controller: Optimal Policy Learning via PPO.** The Controller is the decision core of the ACraft. Its task is to learn an optimal evolutionary strategy $\pi_\theta$. This strategy aims to maximize the expected fitness of the new algorithms generated by the set of transformation functions $\mathcal{F}_{trans}$. We model this strategy learning process as a policy optimization problem, and we use the PPO algorithm to solve it.

The parameters $\theta$ of the strategy $\pi_\theta$ define the probability distribution $\pi_\theta(a_t|s_t)$ that the Controller uses to select an evolutionary action $a_t$ given the population state $s_t$. A core benefit of PPO is its ability to update the policy stably, preventing performance collapse caused by excessively large policy steps.

- State ($s_t$): $s_t$ describes the statistical properties of the current algorithm population $P_t$ and the history of actions.

- Action ($a_t$): The action space $\mathcal{A}$ is the set of decisions to select a transformation function $F \in \mathcal{F}_{trans}$ and choose its required context $\mathcal{C}$. Thus, $a_t = (F, \mathcal{C})$.

---

**Algorithm 1** The ACraft Evolutionary Loop

---

**Require:** LLM Generator $\mathcal{L}$, Benchmark $\mathcal{B}$, Population Size $N$, Budget $T_{max}$, PPO Controller $(\pi_\theta, V_\psi)$

1: Define initial instruction set $\Pi_{init}$ based on $\mathcal{B}$
2: $\mathcal{P}_0 \leftarrow F_{init}(\Pi_{init}, \emptyset)$
3: $s_0 \leftarrow [\mu_{\Phi_0}, \sigma^2_{\Phi_0}, \emptyset]$
4: **for** $t = 1$ to $T_{max}$ **do**
5:     $a_t \leftarrow \pi_\theta(\cdot|s_t)$
6:     $G'_t \leftarrow F(\mathcal{C})$
7:     **for** each $g' \in G'_t$ **do**
8:         Compute fitness $\phi(g') \leftarrow \vec{w} \cdot \vec{J}(g')$
9:     **end for**
10:    $r_t \leftarrow \frac{1}{|G'_t|} \sum_{g' \in G'_t} \phi(g')$
11:    Update Controller $(\pi_\theta, V_\psi)$ using $L^{CLIP}(\theta)$ and $r_t$
12:    $s_{t+1} \leftarrow [\mu_{\Phi_{t+1}}, \sigma^2_{\Phi_{t+1}}, \text{History}(a_{t-k:t-1})]$
13:    $\mathcal{P}_{t+1} \leftarrow \text{SelectTopN}(\mathcal{P}_t \cup G'_t)$
14:    **if** average fitness has converged **then**
15:       **break**
16:    **end if**
17: **end for**
18: **return** Best algorithm found, $g^*$

---

To achieve the robust convergence of the strategy, PPO optimizes a clipped objective function $L^{CLIP}(\theta)$. This function limits how much the new policy $\pi_\theta$ can change relative to the old policy $\pi_{\theta_{old}}$.

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[ \min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \right]$$

The core components of the PPO objective are defined as follows:

- Policy Ratio $(r_t(\theta))$: This ratio measures the probability of the new policy, $\pi_\theta$, selecting the action $a_t$ relative to the old policy, $\pi_{\theta_{old}}$.

$$r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$$

- Reward $(r_t)$: This represents the immediate gain signal for the PPO optimization. It is defined as the average fitness score of the new generation $G'_t$ produced after executing action $a_t$:

$$\text{Reward } r_t = \frac{1}{|G'_t|} \sum_{g' \in G'_t} \phi(g')$$

where:

- $\phi(g')$ is the fitness score of the generated algorithm $g'$, which balances attack efficacy and cost (as defined by the Multi-Attribute Utility Function).

- $\hat{A}_t$ is the Advantage Estimate, which shows the expected extra return of the current action $a_t$ compared to the average policy.

- The Clipping Term limits the policy ratio $r_t(\theta)$ to the range $[1 - \epsilon, 1 + \epsilon]$, where $\epsilon$ is a clipping hyperparameter in PPO.

By iteratively optimizing $L^{CLIP}(\theta)$, the Controller learns a meta-level strategy $\pi_\theta$ efficiently. This strategy can intelligently orchestrate the algorithm generation process (that is, select the optimal transformation function $F \in \mathcal{F}_{trans}$) to systematically improve the overall fitness of the population. Consequently, it efficiently solves the optimization problem defined in the Mathematical Formulation section: $g^* = \arg\max_{g \in \mathcal{G}} \mathbb{E}[\mathcal{L}_{attack}(M_g)]$.

### 3.4 The Closed-Loop Algorithmic Workflow

This section integrates the Generator, Evaluator, and Controller into an iterative evolutionary cycle. This closed-loop workflow is essential for ACraft to efficiently solve for the optimal strategy $g^*$. The cycle systematically explores the latent design space $\mathcal{G}$ by continuously feeding algorithm performance back into the strategy optimizer.

Table 2: **Performance comparison in miniImageNet dataset.** Session-wise Acc (%) and the average of them (Avg) are reported. Attack Drop. represent the magnitude of the decrease in the accuracy of the last session after being attacked.

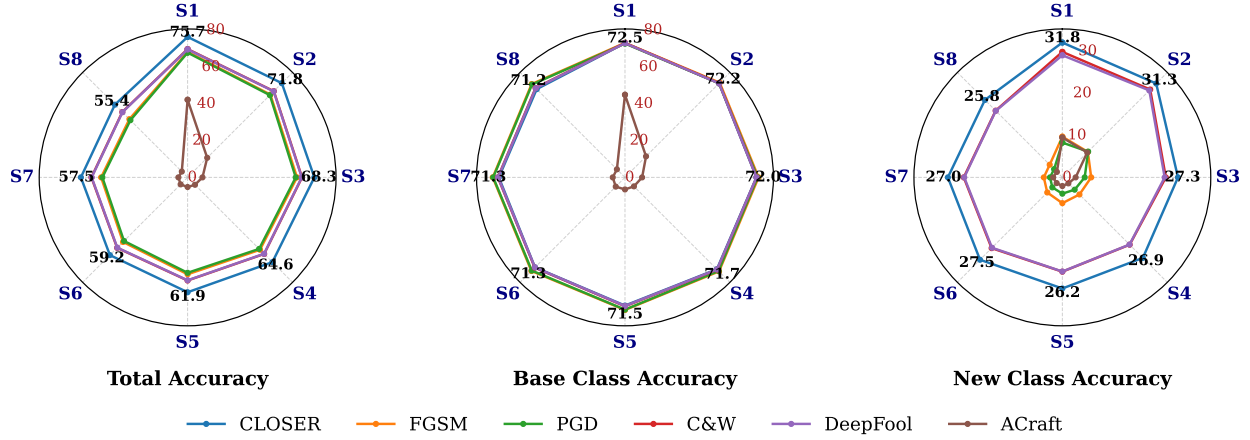| Methods | Year | Acc in each session (%)↓ | | | | | | | | | Avg↓ | Attack. Drop.↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | | |
| CLOSER [18] | ECCV 2024 | 76.02 | 71.61 | 67.99 | 64.69 | 61.70 | 58.94 | 56.23 | 54.52 | 53.33 | 62.78 | - |
| FGSM [9] | | 76.02 | 68.15 | 63.73 | 59.98 | 56.46 | 53.33 | 50.47 | 47.93 | 46.08 | 58.01 | 4.77 |
| PGD [15] | | 76.02 | 68.4 | 63.63 | 60.09 | 56.53 | 53.12 | 50.36 | 47.98 | 45.87 | 58.00 | 4.76 |
| C&W [5] | | 76.02 | 68.98 | 64.94 | 61.73 | 58.79 | 55.93 | 53.27 | 51.16 | 50.13 | 60.10 | 1.68 |
| DeepFool [16] | | 76.02 | 68.98 | 64.94 | 61.73 | 58.79 | 55.93 | 53.27 | 51.16 | 50.13 | 60.10 | 1.68 |
| **ACraft (our)** | | 76.02 | 12.20 | 11.34 | 10.52 | 9.88 | 9.31 | 8.69 | 8.28 | 7.96 | 17.13 | 58.89 |



Figure 4: ACraft *v.s.* existing attacks for FSCIL task in Cifar-100 dataset.

**Step 0: Initialization and Seed Generation.** The Controller selects the Genesis Function $F_{init}$, and the Generator executes $g' = F_{init}(\Pi_{init}, \emptyset)$ until $N$ initial algorithms form the population $P_0$.

**Step 1: Policy Sampling and Algorithm Synthesis.** The PPO Controller observes the current state $s_t$, samples action $a_t = (F, \mathcal{C})$ based on $\pi_\theta$, and the Generator executes $g' = F(\mathcal{C})$ to create the offspring $G'_t$.

**Step 2: Attack Deployment and Fitness Evaluation.** Each $g' \in G'_t$ is deployed against the target system. The Evaluator runs the executable code $C$, measures the impact, and calculates the fitness $\phi(g')$ using the Multi-Attribute Utility Function. Failed algorithms receive a penalty value.

**Step 3: Strategy Update via PPO.** The average fitness $r_t = \mathbb{E}[\phi(G'_t)]$ is fed back as the reward. PPO uses this reward and the policy ratio $r_t(\theta)$ to optimize the clipped objective function $L^{CLIP}(\theta)$, updating $\pi_\theta$ parameters $\theta$.

**Step 4: Population Management and Iteration.** The candidate set $P_t \cup G'_t$ is ranked by $\phi(\cdot)$. The top $N$ algorithms are kept to form $P_{t+1}$. The loop continues until the maximum iterations $T_{max}$ are reached or the average fitness $\mu_\Phi$ converges. The final output is the best discovered algorithm is $g^* = \max_{g \in \mathcal{G}} \phi(g)$.

## 4 Experiment and Discussions

### 4.1 Experimental Settings

We conduct extensive experimental comparison on benchmarks (i.e., miniImageNet [22], CIFAR100 [11], and CUB200 [23]) to validate the effectiveness of the proposed method. We take CLOSER [18] as our main baseline. The detailed experimental settings are presented in **App. B**.

### 4.2 Results on FSCIL Benchmarks

**miniImageNet Dataset.** As shown in Table 2, our proposed ACraft achieves the most destructive attack performance on the miniImageNet dataset, substantially outperforming all existing manually designed baselines. Specifically, com-

Table 3: **Generalizability of our ACraft for more FSCIL methods in miniImageNet dataset.**

| FSCIL methods | Year | Acc in each session (%)↓ | | | | | | | | | Avg↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | |
| Limit [28] | TPAMI2023 | 84.13 | 78.80 | 74.21 | 69.64 | 66.01 | 63.32 | 61.50 | 58.17 | 55.88 | 67.96 |
| Limit + ACraft | | 84.13 | 18.45 | 16.89 | 15.28 | 14.29 | 13.12 | 12.35 | 11.63 | 10.88 | 19.65 ↓ (48.31) |
| Approximation [24] | ECCV2024 | 84.13 | 75.25 | 70.25 | 67.54 | 64.85 | 61.34 | 59.74 | 56.38 | 53.45 | 65.88 |
| Approximation + ACraft | | 84.13 | 17.82 | 15.61 | 14.08 | 12.97 | 11.34 | 10.78 | 9.92 | 9.45 | 18.01 ↓ (47.87) |
| OrCo [4] | CVPR2024 | 84.13 | 77.84 | 73.42 | 69.22 | 65.85 | 62.75 | 60.56 | 57.34 | 55.33 | 67.38 |
| OrCo + ACraft | | 84.13 | 19.06 | 16.84 | 15.30 | 13.47 | 12.48 | 11.50 | 10.61 | 9.80 | 19.13 ↓ (48.07) |
| Tri-WE [28] | CVPR2025 | 84.13 | 81.41 | 76.65 | 73.59 | 70.1 | 65.13 | 63.42 | 61.02 | 60.13 | 70.62 |
| Tri-WE [28] + ACraft | | 84.13 | 20.35 | 17.99 | 15.76 | 14.09 | 12.62 | 11.08 | 10.92 | 10.12 | 19.56 ↓ (51.06) |

Table 4: Comparison of "Attack Drop" performance between Naive and ACraft methods with different LLMs on CLOSER [18] in the miniImageNet dataset.

| Methods | LLMs | Run 1 | Run 2 | Run 3 | Run 4 | Average |
|---|---|---|---|---|---|---|
| Naive | gpt-4o-mini [3] | 19.37 | 18.74 | 18.52 | 19.76 | 19.10 |
| ACraft | Claude 3.7 [1] | 53.76 | 55.28 | 52.72 | 48.93 | 52.67 |
| ACraft | Deepseek V3 [13] | 49.32 | 45.77 | 50.38 | 52.33 | 49.45 |
| ACraft | Gemini flash [7] | 57.28 | 50.83 | 49.17 | 52.65 | 52.48 |
| ACraft | Llama 3 [10] | 49.04 | 48.70 | 45.32 | 43.95 | 46.75 |
| ACraft | gpt-4o-mini [20] | 56.61 | 49.98 | 53.50 | 54.67 | 53.69 |
| ACraft | Grok 3 [2] | 52.98 | 53.89 | 49.67 | 55.36 | 52.98 |

pared with the original CLOSER model (62.78% average accuracy), ACraft reduces the accuracy to only 23.10%, resulting in an attack drop of 39.68 points, far exceeding traditional methods such as FGSM (4.77), PGD (4.76), and C&W (1.68). Remarkably, ACraft not only induces the largest performance degradation across all incremental sessions but also maintains extremely consistent and low accuracies (approximately 12–23%) throughout the entire incremental learning process. This demonstrates that ACraft can effectively destroy both base and novel class performance, overcoming the instability and inefficiency of expert-crafted attacks. These results clearly validate the superiority of our LLM-driven automatic attack generation framework, highlighting its capability to deliver both higher attack effectiveness and lower design cost in few-shot class-incremental learning scenarios. **In particular, ACraft is extremely efficient, only costing 47 minutes with gpt-4o-mini.**

**CIFAR-100 and CUB200 Datasets.** As shown in Fig. 4 and 2, ACraft exhibits outstanding attack performance on the CIFAR-100 and CUB200 datasets, further validating its robustness and generalization ability. Specifically, ACraft achieves the largest accuracy degradation across all incremental sessions, significantly outperforming classical hand-crafted methods such as FGSM, PGD, and C&W. Compared with the baseline CLOSER model, ACraft drastically reduces both the total and base accuracies while maintaining the lowest accuracy for new classes, indicating its ability to disrupt model learning consistently throughout the entire FSCIL process. Importantly, ACraft achieves these results without any manual design effort, highlighting the efficiency of our LLM-driven framework. These findings clearly demonstrate ACraft's superior attack effectiveness and strong generalization efficiency across different datasets, establishing a new benchmark for automated adversarial attack discovery in FSCIL.

### 4.3 Generalizability on more FSCIL Frameworks

To further validate the generalizability of ACraft across different mainstream FSCIL frameworks, we conduct extensive experiments on multiple representative methods, including Limit [28], Approximation [24], OrCo [4], and Tri-WE [28], on the miniImageNet dataset (see Table 3). From Table 3, we observe that ACraft consistently degrades the performance of all FSCIL methods by a large margin, confirming its strong generalization ability. For instance, when integrated with Limit, the average accuracy drops sharply from 67.96% to 19.65%. Similarly, Approximation, OrCo, and Tri-WE suffer significant performance declines from 65.88% to 18.01%, 67.38% to 19.13%, and 70.62% to 19.56%, respectively. These results demonstrate that ACraft retains strong generalizability and stability across different FSCIL methods, showing strong adaptability and transferability for diverse FSCIL frameworks.

### 4.4    Ablation study

**The impact of various LLMs:**    To scrutinize the influence of different large language models (LLMs), we conduct an ablation study on six representative LLMs, i.e., GPT4o, Claude 3.7, DeepSeek V3, Gemini Flash, Llama 3, and Grok 3, using the CLOSER benchmark on miniImageNet (see Table 4). Each result is averaged over four independent runs for fairness. As shown in Table 4, ACraft demonstrates strong robustness across all LLMs, with "Attack Drop" values remaining consistent (46.75–53.69), indicating stable adaptation to diverse model behaviors. Notably, GPT4o achieves the best performance (53.69%), reflecting its stronger reasoning and optimization capability in generating effective adversarial perturbations. Overall, ACraft maintains stable and reliable performance under different LLMs.

### 4.5    Additional Evaluation

Due to the page limit, this paper provides more experimental results in **App. A-G**. ❶ **Detailed Prompt Engineering of ACraft** are depicted in **App. C**. ❷ **More ablation studies of ACraft** are shown in **App. D**. ❸ **Limitations and Discussion** are presented in **App. E**. ❹ **Visualizations of searched attack method** are depicted in **App. F**.

**Related Work.** The related work is provided in **App. G**.

## 5    Conclusion and Future Works

In this work, we observe that existing attack methods in few-shot class-incremental learning (FSCIL) are manually designed, inefficient, and heavily dependent on expert knowledge. To address this limitation, we propose ACraft, an LLM-driven framework that automatically discovers effective adversarial attack algorithms without human intervention. By incorporating a Proximal Policy Optimization (PPO)-based reinforcement learning strategy, ACraft iteratively refines prompts to generate stronger and more adaptive attacks against FSCIL models. The proposed method significantly outperforms expert-designed attacks while maintaining the lowest computational cost on multiple FSCIL benchmarks. We believe ACraft introduces a novel paradigm for automated adversarial attack generation and will inspire future research on secure and robust continual learning systems. In future work, we plan to further enhance ACraft's reasoning and adaptation capabilities for broader applications.

## References

[1]   The claude 3 model family: Opus, sonnet, haiku. 9

[2]   grok-3. 9

[3]   Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 9

[4]   Noor Ahmed, Anna Kukleva, and Bernt Schiele. Orco: Towards better generalization via orthogonality and contrast for few-shot class-incremental learning. In *41st IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 2024. 9

[5]   Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 ieee symposium on security and privacy (sp)*, pages 39–57. Ieee, 2017. 3, 8

[6]   Yiyang Chen, Tianyu Ding, Lei Wang, Jing Huo, Yang Gao, and Wenbin Li. Enhancing few-shot class-incremental learning via training-free bi-level modality calibration. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 9881–9890, 2025. 2

[7]   Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 9

[8]   Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3366–3385, 2021. 2

[9]   Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014. 3, 8

[10]   Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David

Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang,

Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. 9

[11] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 8

[12] Dong Li, Aijia Zhang, Junqi Gao, and Biqing Qi. An efficient memory module for graph few-shot class-incremental learning. *Advances in Neural Information Processing Systems*, 37:130084–130108, 2024. 2

[13] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024. 9

[14] Ye Liu and Meng Yang. Sec-prompt: Semantic complementary prompting for few-shot class-incremental learning. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 25643–25656, 2025. 2

[15] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017. 3, 8

[16] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2574–2582, 2016. 3, 8

[17] Venkata Prabhakara Sarath Nookala, Gaurav Verma, Subhabrata Mukherjee, and Srijan Kumar. Adversarial robustness of prompt-based few-shot learning for natural language understanding. *arXiv preprint arXiv:2306.11066*, 2023. 2

[18] Junghun Oh, Sungyong Baik, and Kyoung Mu Lee. Closer: Towards better representation learning for few-shot class-incremental learning. In *ECCV*, 2024. 4, 8, 9

[19] Junghun Oh, Sungyong Baik, and Kyoung Mu Lee. Closer: Towards better representation learning for few-shot class-incremental learning. In *European Conference on Computer Vision*, pages 18–35. Springer, 2024. 2

[20] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. 9

[21] Eli Verwimp, Rahaf Aljundi, Shai Ben-David, Matthias Bethge, Andrea Cossu, Alexander Gepperth, Tyler L Hayes, Eyke Hüllermeier, Christopher Kanan, Dhireesha Kudithipudi, et al. Continual learning: Applications and the road forward. *arXiv preprint arXiv:2311.11908*, 2023. 2

[22] Oriol Vinyals, Charles Blundell, Tim Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. 2016. 8

[23] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. 8

[24] Xuan Wang, Zhong Ji, Xiyao Liu, Yanwei Pang, and Jungong Han. On the approximation risk of few-shot class-incremental learning. In *European Conference on Computer Vision*, pages 162–178. Springer, 2024. 9

[25] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35: 24824–24837, 2022. 3

[26] Cairong Zhao, Yubin Wang, Xinyang Jiang, Yifei Shen, Kaitao Song, Dongsheng Li, and Duoqian Miao. Learning domain invariant prompt for vision-language models. *IEEE Transactions on Image Processing*, 33:1348–1360, 2024. 2

[27] Li-Jun Zhao, Zhen-Duo Chen, Yongxin Wang, Xin Luo, and Xin-Shun Xu. Attraction diminishing and distributing for few-shot class-incremental learning. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 25657–25666, 2025. 2

[28] Da-Wei Zhou, Han-Jia Ye, Liang Ma, Di Xie, Shiliang Pu, and De-Chuan Zhan. Few-shot class-incremental learning by sampling multi-phase tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(11):12816–12831, 2023. 9

[29] Yiwei Zhou, Xiaobo Xia, Zhiwei Lin, Bo Han, and Tongliang Liu. Few-shot adversarial prompt learning on vision-language models. *Advances in Neural Information Processing Systems*, 37:3122–3156, 2024. 2