# The Future of Continual Learning in the Era of Foundation Models: Three Key Directions

Jack Bell[1,*], Luigi Quarantiello[1], Eric Nuertey Coleman[1], Lanpei Li[1,2], Malio Li[1], Mauro Madeddu[1], Elia Piccoli[1] and Vincenzo Lomonaco[1]

[1]*Department of Computer Science, Università di Pisa, 56126 Pisa, Italy*

[2]*Institute of Information Science and Technologies, National Research Council, 56124 Pisa, Italy*

## Abstract

Continual learning—the ability to acquire, retain, and refine knowledge over time—has always been fundamental to intelligence, both human and artificial. Historically, different AI paradigms have acknowledged this need, albeit with varying priorities: early expert and production systems focused on incremental knowledge consolidation, while reinforcement learning emphasised dynamic adaptation. With the rise of deep learning, deep continual learning has primarily focused on learning robust and reusable representations over time to solve sequences of increasingly complex tasks. However, the emergence of Large Language Models (LLMs) and foundation models has raised the question: Do we still need continual learning when centralised, monolithic models can tackle diverse tasks with access to internet-scale knowledge? We argue that continual learning remains essential for three key reasons: (i) continual pre-training is still necessary to ensure foundation models remain up to date, mitigating knowledge staleness and distribution shifts while integrating new information; (ii) continual fine-tuning enables models to specialise and personalise, adapting to domain-specific tasks, user preferences, and real-world constraints without full retraining, avoiding the need for computationally expensive long context-windows; (iii) continual compositionality offers a scalable and modular approach to intelligence, enabling the orchestration of foundation models and agents to be dynamically composed, recombined, and adapted. While continual pre-training and fine-tuning are explored as niche research directions, we argue it is continual compositionality that will mark the rebirth of continual learning. The future of AI will not be defined by a single static model but by an ecosystem of continually evolving and interacting models, making continual learning more relevant than ever.

## Keywords

Continual Learning, Foundation Models, Continual Pre-training, Continual Fine-tuning, Continual Compositionality

## 1. Introduction

In recent years, artificial intelligence (AI) systems have begun to surpass human performance in many domains such as natural language processing (NLP) and computer vision. However, these models are typically static in nature and do not naturally update their understanding as new data emerges over time. In contrast, humans tend to approach problems as sequential learning tasks, building on past information without forgetting previously learned knowledge or requiring rehearsal to retain it [1]. Both human and AI systems require the ability to learn and adapt continuously, whilst avoiding so-called catastrophic forgetting, in which new learning erases previous knowledge. Addressing this challenge is one of the core aims of Continual Learning (CL).

Continual Learning research therefore revolves around two primary goals: **adaptation** and **memory consolidation**. Adaptation emphasises rapid responsiveness, enabling an agent to quickly adjust its behaviour or representations to maximise a utility function given the current task, situation or environment [2]. Memory consolidation, on the other hand, involves building durable, generalisable

knowledge and skills from past experiences deemed relevant to future tasks. This consolidation goes beyond mere retention; it focuses on developing abstract and hierarchical representations of knowledge, reusable across increasingly complex tasks over time.

Historically, these two goals have been addressed with varying degrees of importance in different research methodologies, contexts and communities. Early expert systems, for instance, focused on consolidating incremental domain knowledge, yet lacked flexibility in rapidly adapting to new information without significant manual effort. Reinforcement learning (RL) methods, such as CHILD, introduced by Ring [3], aimed instead at quick progressive learning, tackling easier tasks before addressing more complex ones, and adapting in a general approach towards continual reinforcement learning. Subsequent work on low-dimensional streaming data looked to address the issue of learning concept drifts [2], similarly focusing on rapid adaptation. With the steep rise and success of deep learning around 2012, the focus shifted prominently towards memory consolidation and generalisation. Deep continual learning leveraged neural networks to learn hierarchical, abstract representations directly from data, enabling the effective reuse of these representations across new tasks [4, 5]. Here, consolidation is not simply memory retention but involves the generalisation of latent knowledge and skills that facilitate adaptation to novel scenarios. We further explore the history of continual learning and its evolution to where it is today in section 2.

More recently, the emergence of LLMs has changed the focus of AI research towards transformer-based models with less focus on 'traditional' Machine Learning (ML). These models, pre-trained on vast datasets, have demonstrated remarkable capabilities to learn rich generalisations of the world [6], performing well across a range of tasks. Combining these impressive capabilities, a demonstrated reduced propensity for catastrophic forgetting [7] and access to the internet-scale knowledge, it is tempting to ask: *Is continual learning still necessary within the era of foundation models?*

Applying continual learning, with its two main aims of adaptability and memory consolidation, to Foundation Models (FMs) is a way to overcome some of their inherent shortcomings: since a FM's parameters are fixed at deployment, every model can be seen as a snapshot of the world at the point of training. Practical use of FMs demands *post-training* adaptation, through fine-tuning or personalisation for a downstream task. Therefore, their static nature poses a significant challenge — they lack the intrinsic adaptability required to stay current in rapidly changing environments.

So far, FM-based agents that actively interact with their environment have been proposed as a promising solution, leveraging continual adaptation to progressively improve capabilities [8]. However, the ability to adapt to new tasks is not enough; it is instead necessary to consolidate new knowledge over time to improve the overall understanding of the world. Challenges such as distributional shifts, long task sequences, task heterogeneity and inaccessible upstream data [9] necessitate a renewed focus on continual learning. Shi *et al.* [9] further observe that these hurdles have pushed recent work toward *task-incremental* and *domain-incremental* benchmarks, where the task identity is supplied or irrelevant. Consequently, while such settings simplify experimentation for academic settings, real-world deployments still confront *class-incremental* conditions in which entirely new tasks must be detected and learned on the fly [10]. We further detail the need for CL in foundation models in section 3.

Given this context, CL for foundational models is developing along three directions: **Continual Pre-Training (CPT)**, **Continual Fine-Tuning (CFT)** and **Continual Compositionality & Orchestration (CCO)**. Human perception provides an instructive analogy: during infancy, critical periods enable rapid specialisation and consolidation of broad sensory capabilities such as vision and language discrimination [11]. Similarly, in ML, extensive pre-training ideally establishes a general-purpose foundational model, but CPT still remains practically necessary. Continual fine-tuning, in turn, allows efficient specialisation and personalisation to specific downstream tasks or contexts. However, both CPT and CFT typically require relatively lower-frequency adaptation cycles. Furthermore, their reliance on large-scale datasets, substantial computational resources and incremental improvement constraints due to scaling laws limits their potential to drive substantial advances alone [9], particularly as these approaches predominantly extend established capabilities rather than enabling fundamentally new behaviours. However, recent work on test-time scaling laws demonstrates that, once model size passes a certain threshold, allocating additional inference-time compute delivers larger accuracy gains than

further parameter growth [12]. Parallel advances in multi-step chain-of-thought prompting and in multi-agent frameworks, where several specialised LLMs negotiate, critique, or divide labour, likewise point to performance improvements that arise from *coordination* rather than monolithic scale [13, 14]. Together, these trends expose the practical limits of an end-to-end foundation model and highlight the need for a modular, dynamically reconfigurable approach.

We therefore contend that **Continual Compositionality and Orchestration** represents the most promising and necessary direction for future continual learning research. Unlike CPT and CFT, CCO inherently supports high-frequency adaptation, allowing dynamic orchestration, recombination and collaborative interaction among multiple FMs or agents. Recent advances in FMs have primarily emerged not from additional computational resources, but rather through enhanced reasoning abilities [15] and longer context windows [16], both of which are forms of orchestration rather than scale. Consequently, the future of continual learning likely resides in decentralised ecosystems, where multiple adaptive agents continuously interact, evolve and collaboratively address increasingly complex problems. This paradigm exhibits parallels to society as a whole, where collections of different individuals can come together to solve difficult tasks. Taken from a large enough intelligent population, a random sample of people will outperform a sample of the best performing agents — with the intuition being that diversity is more important than individual ability [17]. In a similar vein, a diverse subset of agents may well be more adept at problem solving than a sample of the best performing agents at a given task.

In this paper, we first review Continual Pre-Training (4.1), which equips large-scale foundational models with adaptive, resource-efficient mechanisms to incorporate new knowledge without catastrophic forgetting. We then examine Continual Fine-Tuning (4.2), enabling precise specialisation to downstream tasks while retaining broad, generalisable representations. Finally, we argue that Continual Compositionality & Orchestration (4.3) - with its high-frequency, modular coordination of specialised agents - offers the most promising path forward. By moving from monolithic snapshots to dynamic, decentralised ecosystems of models, CCO can drive the next wave of resilient, scalable, and sustainable AI systems.

## 2. Background and Related Work

The rise of Deep Learning (DL) in 2012 marked a pivotal moment for the entire ML research community. Initial efforts in this domain concentrated on utilising deep neural networks for representation learning, allowing models to capture abstract and hierarchical features from data. Nevertheless, when trained sequentially on multiple tasks, these models remained susceptible to catastrophic forgetting. This highlights the need for CL, which has had a large impact on the broader field of ML.

A widely used definition, often taken to specify Deep Continual Learning, is offered by Lesort *et al.* [18]. They describe continual learning as a learning paradigm where a model learns from a continuous stream of data, adapting to new information while preserving previously acquired knowledge. This definition emphasises the importance of both stability and plasticity in the learning process, where there is an important trade-off between retaining past knowledge and being plastic enough to adapt to new data or domains. The objective of CL is to have a machine learning model that can be adapted quickly to shifts in data distribution or "tasks", enabling it to retain already acquired knowledge and concepts and reuse these representations to facilitate better learning across new tasks.

This is a key difference between CL and traditional machine learning approaches, which typically require retraining on a static dataset to incorporate new information.

Early studies focused on the problem of catastrophic forgetting in neural networks [19, 20], where authors discovered the degradation of model performance on previous tasks while learning a new one. To overcome the issue, different basic approaches were proposed [21, 22, 23, 24]. These early works laid the foundations for more sophisticated approaches that emerged in the following years.

The main approaches within CL can be categorised into three main schools of thought: regularisation-based methods, dynamic architectures, and memory-based techniques [25, 26].

**Regularisation-based** methods aim to mitigate catastrophic forgetting by adding constraints to

the learning process, ensuring that important weights from previous tasks or domains are preserved. EWC [24] adds a regularisation term to the loss function to preserve important weights from previous tasks; SI [27] uses a similar approach but focuses on the Fisher information matrix to identify important weights; LWF [23] introduces a distillation loss to retain knowledge from previous tasks.

**Dynamic architectures** involve modifying the model's structure to accommodate new tasks while retaining previously learned knowledge. PNNs [28] add new subnetworks for each task while keeping the previously learned parameters frozen; LWTA [21] divides the networks into different blocks and the forward propagation is done only by local winner weights; PackNet [29] prunes the networks and uses masks to filter weights for different tasks; Piggyback [30] drops the weights training [31] and focuses only on weights masking; a similar approach is used in HAT [32], but masking the units instead of the weights; SupSup [33] simplifies this even more using only a seed to generate random weights and a weighted sum of masks.

**Memory-based** methods utilise stored examples from past tasks to reinforce prior knowledge during training. GDumb [34] randomly selects a subset of exemplars from previous tasks and uses them to train the model; GEM [35] uses episodic memory to store exemplars from previous tasks and ensure that the model does not forget them during training; iCaRL [36] utilises stored examples from past tasks to replay and reinforce prior knowledge during training; potentially, one could also use latent generative replay to generate new samples from the previous tasks [37, 38].

For a clearer academic evaluation of different CL methods, three main scenarios have been defined and are widely used within the community [39]. In **Task-incremental learning**, the aim is to incrementally learn a set of distinct tasks based on a given task-id. Whereas in **Domain-incremental learning**, the context or input distribution varies over time, whilst the task remains constant (e.g. learning to drive in different weather conditions). Finally, in **Class-incremental learning**, the aim is to incrementally learn to discriminate between a growing number of objects or classes, where task identification is also required. This last scenario is naturally the most difficult to solve, as well as the most applicable within real-world scenarios.

Differently from CL models, Foundation Models (FM) are characterised by their large scale, pre-training on huge amounts of data, and ability to perform a wide range of tasks with a relatively small amount of fine-tuning. These models, such as BERT [40], GPT-3 [41], and CLIP [42], have demonstrated remarkable performance across various benchmarks and applications. However, one of the biggest problems faced by FM models is *homogenisation* [43], where a single unified model trained on diverse data results in generalised internal knowledge representations. While this approach enables model transfer across tasks, it also averages out critical domain-specific nuances, leading to inherent biases from the most dominant data sources.

More generally speaking, a monolithic AI system, which is a single large model that is trained on a wide range of tasks and domains [44], is typically pre-trained on massive datasets and then fine-tuned for specific tasks. Such models often require extensive computational resources for training and inference, making them less accessible for smaller organisations and researchers. Additionally, a monolithic model can suffer from knowledge staleness, where the model's performance degrades over time as new data becomes available. Furthermore, their centralised nature poses risks related to single points of failure making them less robust in critical applications.

On the other hand, continual learning offers a more sustainable and adaptive alternative by enabling models to learn incrementally, adapt to new tasks, and operate efficiently in dynamic settings [25]. This paradigm shift is essential for addressing the limitations of monolithic AI and fostering more equitable and resilient AI systems. Continual learning methods could benefit FMs in particular by enabling them to adapt to new tasks and domains without requiring extensive retraining [45], thereby improving their performance over time. By incorporating continual learning strategies, FMs can mitigate issues such as knowledge staleness and inefficiency in adaptation. Furthermore, these methods can help reduce the social and environmental impact of FMs by minimising the need for large-scale retraining, which often requires significant computational resources and energy consumption.

# 3. The Need for Continual Learning in the Foundation Model Era

Recently, it has been shown that foundational models such as LLMs have the ability to both "reason" and generalise through the use of techniques such as chain of thought (CoT) prompting, by breaking a complex problem into a series of intermediate steps [46]. This trend of utilising the zero-shot capabilities of LLMs, perhaps with the addition of prompt engineering and later downstream-task fine-tuning, is, however, subject to many different limitations such as the brittleness and inconsistency of the generalised reasoning steps.

Foundation models such as GPT-3, BERT and DALL-E can quickly become outdated as the real-world data they are trained on changes [41, 43]. This leads to model staleness over time, where often full model retraining is the only solution used to mitigate this. Recent estimates indicate that the cost of training models of a scale comparable to GPT-4.5 or similar architectures likely reached tens of millions of dollars due to the enormous compute and energy resources required [47]. Beyond the economic impact, the energy required for such massive compute workloads translates into significant environmental impacts. Studies by Schwartz *et al.* [48] and Strubell *et al.* [49] indicate that the high energy consumption involved in training these models contributes substantially to carbon emissions, emphasising the need for more sustainable approaches.

State of the art large scale models are typically trained on vast, diverse datasets to capture a wide range of linguistic patterns and knowledge. Whilst this training approach enables impressive generalisation across tasks, it results in a system that is too general to properly address the requirements of individual users or specialised domains [50]. In scenarios such as personalised recommendations or adaptive customer support, tailoring interactions based on a user's context or preferences is of paramount importance; however, static models tend to produce generic outputs even after fine-tuning.

Further, conventional personalisation methods, relying on post-training adjustments or Test Time Training (TTT) face significant challenges. In TTT or test time adaptation (TTA), a model adjusts its parameters during inference based on the current input or an auxiliary task, aiming to better align with the data distribution at test time [51]. Whilst this method allows for on-the-fly adaptation, it frequently requires additional compute post-deployment and may struggle to capture long-term user preferences. Additionally, TTT needs to be combined with continual adaptation with care as it *can* interfere with pre-trained representations if not done correctly, leading to instability and ultimately catastrophic forgetting.

The rising computational and financial demands for training state of the art models have led to a concentration of resources within a small number of organisations. This centralisation means that only a few giants of industry possess the capability to develop, maintain, and update foundation models at scale [52]. As a result, there is an inherent risk of monopolisation, where control over advanced AI technologies is restricted to those with the resources, potentially limiting progress and diversity within research, raising wider issues around transparency and accountability. For instance, when few entities dominate model training and deployment, issues such as biased data representation and the under-representation of marginalised groups can become more pronounced [53]. Additionally, the centralised model may raise the barrier to entry for smaller research teams, universities, and independent developers, inadvertently slowing innovation within the AI ecosystem. Furthermore, the monopolisation of AI resources restricts the development of more sustainable and decentralised approaches.

While the challenges of centralisation, high retraining costs, and limited personalisation have long constrained the evolution of large-scale AI systems, these can each be addressed in different ways by continual learning. There are two fundamentally different types of forgetting in the context of foundational models that necessitate continual learning as a solution: task-shift and time-shift forgetting [9]. The first, *task-shift forgetting*, arises when a broadly pre-trained model is adapted to a new downstream objective. Without careful safeguards, the updates that confer task skill can overwrite previously acquired general knowledge. This can be effectively mitigated using techniques such as continual pre-training (or domain adaptive pre-training) and continual fine-tuning. The second, *time-shift forgetting*, occurs even when the task remains unchanged: as the external world evolves, the data

distribution shifts and a static model's accuracy diminishes unless it is retrained. Additionally, the scalability and modularity of model architectures need to be considered in order to enable models to learn new tasks over time by dynamically composing task-specific modules to solve new tasks [54]. This notion of model compositionality is of the utmost importance to enable models to not only solve novel tasks, but also to be dynamically orchestrated, facilitating interaction with one another in a wider decentralised system of models. Whilst this ecosystem of interactive agents has already begun to come to the fore [55, 13], all of these models are fundamentally static, requiring centralised retraining over time, which is naturally prohibitive in encouraging more open, democratic and decentralised AI systems.

## 4. Three Key Research Directions for Continual Learning

Here, we outline the three main research directions that are crucial for the future of Continual Learning, namely Continual Pre-Training, Continual Fine-Tuning and Continual Compositionality & Orchestration, visualised in Figure 1. In the following sections, we will outline the requirements for each of these separate components, challenges they face, open problems that are yet to be solved and how these methods address the challenges within this context.
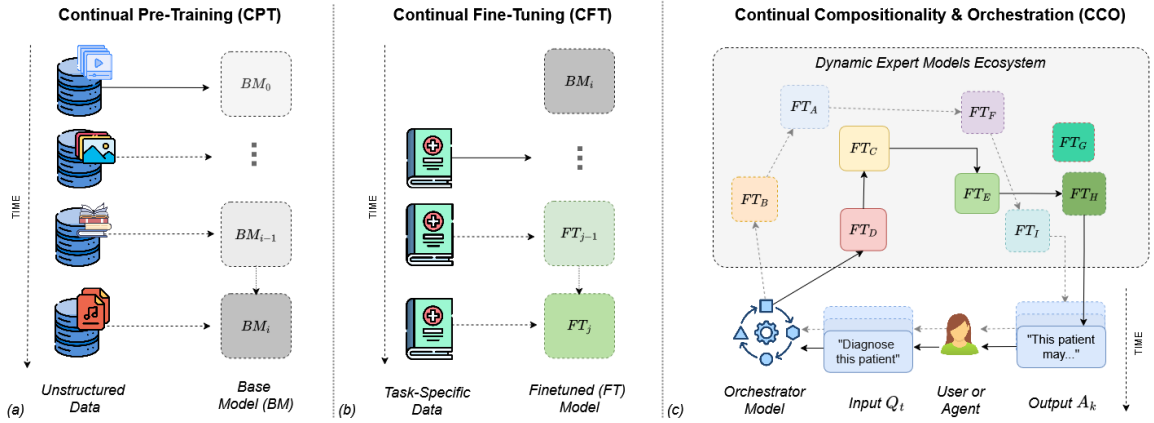


**Figure 1:** In (a), we see a base model pre-trained on video data is Continually Pre-trained on general corpora spanning different modalities (e.g. audio, images and text). Then in (b), this base model is Continually Fine-tuned over time, resulting in specialised fine-tuning (FT) modules trained on domain-specific datasets, such as medical texts. Finally in (c)—looking at model inference, an orchestrator routes a user's query through the appropriate FT modules and combines their outputs into a single response. User inputs may change over time, as may the configurations by which models are composed, the processes through which models evolve via Continual Fine-tuning, and the introduction of new models as they become available.

## 4.1. Continual Pre-Training

CPT refers to the process of incrementally updating the knowledge of FMs through exposure to new data after their initial pre-training phase [54]. This iterative updating allows FMs to maintain their foundational abilities established during the initial training while simultaneously adapting to assimilate emerging information, thereby extending their operational lifespan and enhancing their adaptability to the ever-changing landscape of data and knowledge [45].

### 4.1.1. Motivation

We now look to five main motivating factors to advance the field of CPT.

**Dynamic Knowledge Integration:** CPT is primarily driven by the necessity to keep FMs relevant and effective over time. Even large-scale foundation models can quickly become outdated as data distributions shift and new knowledge emerges [56]. Static pre-training, performed on a fixed dataset,

inevitably leads to models losing touch with rapidly developing fields, where trends, terminologies and societal norms continuously evolve, *e.g.* healthcare, law, technology. To address this challenge, CPT performs a **dynamic knowledge integration** by continuously learning on new data streams, enabling FMs to adapt and remain aligned with contemporary information [57]. Extensive studies have not only demonstrated the necessity of CPT for improved downstream performance, but also shown that when distributional shifts are gradual or somewhat correlated, CPT can effectively help models generalise to previously unseen data [9]. In particular, CPT enables FMs to handle distributional shifts such as temporal shifts (changes in data over time, leading to model drift), content shifts (changes in topic or domain of the data) and language shifts (introduction of new languages or significant vocabulary changes) [45].

**Methodological Evolution:** CPT is also motivated by the ongoing improvements in model architectures. Even when there is no new data or major distribution shift, updating foundation models can still be valuable for adopting architectural advances that improve efficiency or performance. For instance, transitioning from an encoder-decoder to a decoder-only architecture benefits from reusing existing pre-trained models, which avoids the need to retrain from scratch and significantly reduces computational cost. CPT enables this kind of update by initialising new models from previous checkpoints, allowing them to retain useful learned knowledge while gradually adapting to architectural changes [58].

**Resource Efficiency:** Retraining FMs entirely from scratch on increasingly vast datasets becomes computationally prohibitive over time. CPT significantly reduces these computational burdens by incrementally updating models with new or domain-specific data, circumventing the need for repeated, costly training cycles [9]. For instance, the LLaMA 4 Behemoth model [59], with its two trillion total parameters, makes full retraining prohibitively expensive, whereas CPT allows for efficient incremental adaptation of FMs.

**Mitigation of Catastrophic Forgetting:** Static models, once trained, are frozen at their initial knowledge cut-off, which creates a gap between initial training and real-world deployment needs. In contrast, CPT fosters a continual evolution of models, bridging the gap between initial pre-training and practical, lifelong learning scenarios. Notably, recent studies find that CPT can make models more robust to catastrophic forgetting of earlier knowledge, especially when using self-supervised objectives, highlighting CPT as a key enabler for foundation models to learn continuously like humans [7].

**Selective Forgetting:** In addition to mitigating catastrophic forgetting, CPT has the potential to support selective forgetting, where specific information embedded in a foundation model is intentionally removed over time [60]. This capability is particularly relevant as large-scale models may inadvertently memorise sensitive, outdated, or harmful content during pre-training [61]. Recent work has emphasised the importance of enabling continual forgetting to remove such undesirable knowledge while preserving overall model performance [62]. While earlier approaches explored this problem through fine-tuning, recent efforts have extended it into the CPT phase. For example, Zhu *et al.* introduce a regularised CPT method that enables the removal of backdoor behaviours in language models while maintaining their functionality on clean data [63]. These developments highlight selective forgetting within CPT as a promising direction for maintaining the safety, privacy, and reliability of foundation models over time.

### 4.1.2. Challenges & Open Problems

CPT is still in the early stages of development, and bridging the gap between research and production remains challenging: while CPT techniques show promise in controlled experiments, their long-term stability and effectiveness over months of deployment in real-world settings remain under-explored [9]. This highlights several key challenges that must be addressed to make CPT viable in practice, including:

**Handling catastrophic forgetting:** Catastrophic forgetting, the phenomenon where a continually updated model loses previously acquired knowledge, remains a critical challenge in CPT [64, 65]. Although scaling up pre-training tends to enhance knowledge transfer and resilience against forgetting during downstream CL, excessively extensive pre-training can significantly increase the risk of forgetting [66, 67]. Recent studies indicate that self-supervised CL exhibits significantly reduced catastrophic

forgetting compared to supervised approaches [68, 7]. Indeed, self-supervised pre-training has the advantage of reduced forgetting during upstream tasks, though effectively balancing upstream CL with downstream continual adaptation remains an open research question [69]. Additionally, model scale plays a pivotal role: larger models consistently demonstrate lower perplexity (an indicator of how unfamiliar or novel a document is to a language model) and reduced forgetting, whereas smaller models, despite achieving substantial learning gains, tend to exhibit the most pronounced forgetting effects [70].

**Balancing Efficiency vs. Model Drift:** FMs often have hundreds of billions of parameters, so retraining them on every new dataset is extremely computationally expensive [54]. Practical CPT must therefore be computationally efficient, for instance by updating only a subset of parameters or using limited data, but this can exacerbate the stability–plasticity dilemma. Insufficient or biased updates may lead to model drift, where performance on original domains degrades or the model's behaviour shifts unpredictably. Empirically, a "stability gap" has been observed: when an LLMs is first continually pre-trained on a new domain, its performance drops initially (due to distribution shift) before recovering [71]. Balancing efficient adaptation with stability (avoiding regressions on prior knowledge) is an open problem.

**Avoiding reinforcement of biases in pre-training:** Continuously ingesting new data can also reinforce biases or fairness issues if not carefully controlled. If the incoming data is skewed or uncurated, the model might amplify these biases over time, undermining responsible use. For example, biomedical FMs trained on federated data must address fairness across diverse populations while preserving privacy [72]. Ensuring that continual updates do not degrade the model's ethical alignment (e.g. with respect to bias and fairness) is crucial [54]. Methods to detect and mitigate bias drift during CPT (and to curate update data) are largely lacking and represent an important research frontier.

### 4.1.3. Potential Solutions & Future Research Directions

To address the challenges of CPT and unlock its full potential, several promising solutions and future research directions have been proposed, including:

**Incremental pre-training strategies:** Recent works propose reusing or initialising from previous model weights to maintain continuity, as in recyclable tuning methods that carry over knowledge from an old model to a new one [45]. Recent work has demonstrated the benefits of structured or multi-stage CPT. For instance, in [73], a two-step CPT was shown to enable a mixed-language neural machine translation system (first adapting to a language domain, then to a specific translation task) effectively. In [74], Dalla Noce *et al.* introduce a sequential CPT framework for neural machine translation, where a model is progressively exposed to new language pairs or domains in multiple stages. Their study finds that incrementally adding new languages during the pre-training phase does not substantially degrade the model's performance on previously seen language pairs during fine-tuning. Furthermore, incorporating CL strategies such as data rehearsal can further reduce performance loss on earlier language pairs compared to purely incremental pre-training but comes with increased computational cost during the training phase. This highlights a practical trade-off between training efficiency and performance robustness: reasonable downstream performance can be achieved through incremental pre-training, but further performance improvements can be attained when accepting the additional computational cost associated with CL strategies.

**Selective memory and rehearsal methods:** To combat forgetting, memory replay methods have shown promise in the context of foundation models. Rather than relying only on the latest data, the model can intermittently rehearse on representative samples of past data (or tasks). In practice, storing raw past data for a foundation model may be impractical or raise privacy concerns such that recent work leverages latent replay, where the model saves a cache of compact feature representations or embeddings of past examples instead of the raw inputs [54]. During CPT, these stored latent vectors can be replayed through the model to reinforce previously learned concepts. This memory-efficient replay has been shown to significantly mitigate forgetting in vision models and is especially valuable when sharing raw data is prohibited (e.g. user data privacy). Besides replay, selective sampling strategies can be used rather than naively mixing new data; the training scheduler might interleave the most

informative or relevant examples carefully chosen from older tasks or emphasise difficult examples that the model is starting to forget. There is evidence that the order and composition of training data in CPT can greatly affect retention [56]. For example, Xie *et al.* use perplexity and embedding similarity metrics to select a fraction of a domain corpus that achieves comparable adaptation with far less training cost [75]. Such sampling not only improves efficiency but can also prevent the model from drifting too far by ensuring the new training distribution is aligned with the model's original knowledge. Additionally, some approaches intermix new data with a small portion of the original pre-training data (or a similar distribution) during updates, explicitly to reduce distribution shift [71]. This kind of rehearsal or data mixing has been shown to narrow the stability gap and avoid performance deterioration on earlier tasks. Going forward, developing principled sampling and replay policies (potentially guided by uncertainty, importance weighting, or task identities) is a key research direction to make CPT robust and scalable.

**Self-supervised continual adaptation techniques:** CPT largely relies on self-supervised learning objectives such as predicting masked tokens, next sentence prediction, and image-text contrastive learning, because these allow the use of unlabelled streaming data. An intriguing finding from recent research is that self-supervised objectives can themselves be leveraged to improve CL. In [7], Cossu *et al.* provide strong empirical evidence that continuing pre-training models in a self-supervised manner yields better knowledge retention than supervised training in the continual setting. Intuitively, self-supervised learning updates may be softer or more diffuse in the parameter space (since they capture broad patterns in data) compared to task-specific fine-tuning which might overwrite more specialised parts of the model. In [76], Ostapenko *et al.* also observed that models pre-trained with broader or more diverse self-supervised signals tend to forget less and transfer better in downstream sequential tasks. This suggests that self-supervised CPT is a promising avenue: as new unlabelled data comes in, one can design auxiliary objectives that encourage the model to integrate new information while maintaining consistency with prior representations. Techniques like contrastive learning on a replay buffer, or predictive modelling that ties new concepts to old ones, fall in this category. Moreover, self-supervision can be combined with light supervision or prompts in an autonomous CL setup (where a model might generate pseudo-labels or questions on new data and learn from them). Overall, self-supervised learning-based CPT not only provides a means to utilise vast unlabelled streams, but also appears to inherently mitigate forgetting, making it a key research direction for lifelong FMs.

## 4.2. Continual Fine-tuning

Continual fine-tuning is the practice of applying a stream of lightweight, task-specific updates to a model after deployment, allowing it to evolve alongside newly arriving data rather than remaining fixed after a single adaptation.

### 4.2.1. Motivation

Fine-tuning is attractive because it it far cheaper than full retraining, requires only task-specific data, and can be executed on modest hardware [77, 78]. Yet, once the weights are written to disk, the model is frozen again. In realistic deployments, data arrives as a stream, such as customer queries, sensor logs, freshly published documents — therefore the ability to fine-tune **continually** is crucial [79, 80]. CFT turns a one-shot adaptation step into a standing capability that (i) personalises responses for each user or organisation [81], (ii) keeps proprietary data on-premise for privacy compliance [82], (iii) reacts quickly to domain drift without the latency of retrieval augmented generation (RAG) pipelines or very long context windows [16], (iv) does all of this with a fraction of the compute budget needed for repeated full-scale updates [83, 84]. The need for CFT, the process of incrementally fine-tuning a model to help it adapt to downstream tasks that involve shifting data distributions and temporal changes [79], cannot be overstated.

Although FMs have demonstrated impressive versatility across different tasks, with the ability to generalise effectively to various domains, their static nature limits the capacity to incorporate new knowledge, adapt to specialised fields, and personalise outputs over time. CFT presents an important

opportunity to make foundation models more flexible, efficient and responsive to real-world changes, making them more useful in more dynamic environments. While continual pre-training focuses on updating a model's general representations using broad, often unlabelled data, continual fine-tuning instead aims to incrementally adapt the model to specific downstream tasks using labelled or structured data, with an emphasis on retaining prior knowledge while learning new information.

### 4.2.2. Challenges & Open Problems

CFT in the context of foundation models like LLMs comes with several challenges such as:

**Balancing Specificity vs. Generalisation**: CFT must maintain a delicate equilibrium between tailoring a model to a specific downstream task and preserving its broad, generalisable knowledge. When a model undergoes CFT, its internal representation becomes optimised to capture patterns necessary for solving a specific downstream task. While this adaptation enhances the performance on domain specific tasks, it risks eroding the broad, general-purpose representations learned during pre-training. This comes back to the stability-plasticity dilemma where models must remain plastic enough to integrate task-specific knowledge while being stable enough to retain the broad representations acquired from prior experiences [85].

**Efficient Adaptation Without Catastrophic Forgetting**: In a similar vein to CPT, catastrophic forgetting is also experienced by CFT methods, however within this context *efficient adaptation* also needs to be considered. This concept refers to the process of updating a pre-trained foundation model to perform well on new tasks, domains, or data distributions while minimising computational resources, data requirements, and training time. In the context of LLMs and other foundation models, efficiency has become increasingly critical as these models grow to billions of parameters. Efficient finetuning techniques like LoRA [78] allow large pre-trained models to be adapted to downstream tasks by updating only a fraction of the model's original parameters, but these techniques are still prone to catastrophic forgetting.

**Data efficiency and privacy concerns in continual fine-tuning**: CFT deals with the dual challenges of data scarcity and privacy concerns, especially in specialised domains to effectively adapt models to new tasks. As foundation models are adapted to increasingly specialised domains, high-quality, domain-specific data becomes progressively scarcer [86].

### 4.2.3. Potential Solutions & Future Research Directions

Despite these limitations, numerous methods have been developed to tackle these challenges.

**Parameter Efficient Fine-Tuning** A key methodological tool in the context of CFT is represented by Parameter Efficient Fine-Tuning (PEFT) methods. These techniques aim to achieve performance comparable to or even surpass full model fine-tuning while updating only a small number of trainable parameters, either by selectively updating a subset of the model's parameters [78] or introducing new task-specific parameters [77]. PEFT methods are particularly advantageous in continual learning scenarios, where models must adapt to a sequence of tasks without forgetting previously learned information.

By updating only a limited number of parameters, PEFT approaches reduce computational overhead and, when applied properly, can help mitigate the risk of catastrophic forgetting. Among these PEFTs, Prompt Adapters and LoRA are the most widely used. LoRA works by introducing low-rank updates to the pre-trained model weights, expressed as:

$$W = W_0 + BA \tag{1}$$

where the pre-trained model $W_0$ is kept frozen, while the low-rank matrices $A$ and $B$ are updated.

Prompt-based techniques like L2P and CoDA Prompt [87, 88] incrementally learn from novel data by designing task-specific prompts that guide the model's attention toward relevant information for each new task, facilitating seamless integration of new knowledge without overwriting existing capabilities.

Similarly, LoRA-based CL approaches, such as C-LoRA and DualLoRA [81, 89], enhance LoRA's applicability in CL by introducing mechanisms like learnable routing matrices and orthogonal subspaces to manage parameter updates across tasks, thereby reducing computational overhead and mitigating catastrophic forgetting.

Adapter techniques like Continuous Adapter (C-ADA) and Adapter-based Continual Learning (ACL) [90, 91] instead offer more efficient solutions for CL. C-ADA introduces a Continual Adapter Layer that extends weights for new tasks while freezing old ones, preserving prior knowledge. It also employs a Scaling & Shifting module to align feature spaces between pre-training and downstream datasets. Similarly, ACL utilises lightweight, task-specific adapters within a fixed pretrained feature extractor and incorporates a task-specific head that groups previously learned classes into an "out-of-distribution" category, facilitating effective feature discrimination.

**Model Merging**: A particularly valuable approach when facing dynamic environments is given by Model Merging. The key point in this class of methods is to combine multiple specialised models learned over time, to create systems that preserve knowledge while adapting to new tasks. The central challenge in model merging is parameter interference, where integrating different models leads to performance degradation. Recent research has developed several innovative solutions to this problem. TIES-MERGING [92] addresses interference by strategically resetting minimally changed parameters and resolving sign conflicts between models. In contrast, DARE [93] employs a different strategy by randomly dropping redundant delta parameters and rescaling the remaining ones, effectively sparsifying merged models without significant performance loss.

While early model merging techniques focused on the static combination of pre-existing expert models, more recent approaches support dynamic integration as new tasks emerge over time. MagMax [94] introduces sequential fine-tuning with maximum magnitude weight selection to effectively incorporate new information while preserving earlier learning. Representation Surgery [95] tackles representation bias by inserting lightweight task-specific modules that realign internal representations between merged models. Adaptive LoRA Merging [96] moves beyond fixed-weight combinations by dynamically computing merging coefficients that balance contributions from new and old domains.

Recent trends in CFT have shifted towards the adaptive integration of lightweight modules, like adapters [77] and LoRA [78], in dynamic environments. This shift enables seamless integration of new tasks without extensive retraining of large models. By merging these modular components on demand, systems can efficiently handle real-world challenges while remaining practical for large-scale deployment.

**Meta Learning for Continual Adaptation**: An alternative perspective on the challenge of CFT is offered by meta learning approaches for continual adaptation. In fact, these methods integrate adaptability into the core learning objective, enabling models to rapidly adjust to new tasks with minimal data and computation. Traditional meta-learning approaches like Model-Agnostic Meta-Learning (MAML) [97] operate by finding parameter initialisations that enable rapid adaptation across a distribution of tasks. When applied to CL scenarios, these methods can be extended to discover parameter configurations that not only adapt quickly but also resist catastrophic forgetting. For instance, ANML [80] uses a neuromodulatory network that enables the model to focus on relevant tasks while minimising interference from previously learned tasks.

Recent works have combined meta-learning with parameter-efficient fine-tuning techniques to enhance CL. AutoLoRA [98] introduces a meta-learning framework that automatically identifies the optimal rank for each LoRA layer, improving adaptation efficiency to new tasks while maintaining performance on previous ones. Similarly, Meta-LoRA [99] presents a memory-efficient method for automatic sample re-weighting during fine-tuning, facilitating efficient continual adaptation across various domains. These approaches exemplify the potential of meta-learning to enhance the adaptability and efficiency of foundation models in dynamic environments.

**Federated Learning (FL) and Decentralised Fine-Tuning Strategies**: Moving towards the direction of a distributed and decentralised AI development, FL is essential for effectively adapting large FMs across organisations, while preserving data privay and optimising computational resources. When applied to FMs, however, FL faces unique challenges. The vast number of parameters in modern FMs

makes server-client communication prohibitively expensive. Different clients naturally generate data with varying distributions, creating potential conflicts in optimisation objectives. Frameworks like FATE-LLM [82] enable collaborative training of LLMs by employing parameter-efficient fine-tuning methods and incorporating privacy-preserving mechanisms. FibecFed [83] enhances this approach by utilising Fisher information for adaptive data sampling and dynamically selecting layers for global aggregation, thereby improving both performance and fine-tuning speed. Additionally, FedRewind [100] introduces a decentralised model exchange strategy inspired by continual learning principles, addressing data distribution shifts and enhancing generalisation performance in federated settings.

## 4.3. Continual Compositionality & Orchestration

Continual Compositionality & Orchestration refers to the dynamic integration of multiple AI agents over time, to solve higher-level tasks. It is the key component towards a distributed and decentralised AI framework.

### 4.3.1. Motivation

Large models solved most of the tasks addressed by AI methods for decades. In fact, using a Transformer-based architecture, pre-trained on some large dataset, we now know how to solve almost every task, provided enough data and computational power. Broadly speaking, modern FMs have achieved super-human performance on most of the traditional machine learning benchmarks, making them obsolete and less relevant for current AI research.

For this reason, we are now shifting towards higher level tasks, which require a higher level of intelligence, that current *state-of-the-art* models do not exhibit. This trend is demonstrated by the emergence of several more complex benchmarks, that could drive the research beyond current LLMs. One such benchmark is ARC-AGI [101], which consists of simple grid transformations. Despite being simple for a human solver, it poses great challenges, requiring strong abstraction skills and inductive reasoning. Another example is BIG-bench [102], a collection of more than 200 tasks, designed to test the limits of current large models. The benchmark covers a broad set of tasks, including linguistics, mathematics, common-sense reasoning, social bias detection and more. In both cases, human performance significantly surpasses the current best models, highlighting the need for alternative solutions.

Furthermore, latest large models, *e.g.* GPT 4.5, demonstrate that we are rapidly moving towards a diminishing returns regime, meaning that just increasing the model dimensions and the number of GPUs employed is no longer enough to obtain noticeable improvements [103]. We believe that a paradigmatic change is needed in the AI community in order to push research forward and to obtain more intelligent behaviours.

Our proposal is to address some of the shortcomings that current solutions exhibit, particularly regarding scalability and sustainability matters, with **Continual Compositionality and Orchestration** approaches. Nowadays, the *de facto* standard is to have single, monolithic models, trained once and deployed without any guarantees on their utility over time. CCO instead represents a framework built on the communication between multiple AI models, which share their knowledge and skills in terms of model parameters, deep representations or final predictions.

Within this paradigm, the orchestration among the *agents* becomes crucial: depending on the task, different modules are selected to be composed in various fashions. In other words, rather than adjusting a single network to a dynamic scenario, CCO employs a modular approach, where different module compositions can be used to adapt to non-stationary environments over time. Also, such a framework does not impose constraints on the scalability of the system, *i.e.* the number of modules involved, and it is more sustainable *by design*, since the modules are trained once and then re-used over time in multiple ways.

### 4.3.2. Challenges & Open Problems

With the advancements in capabilities of LLMs, there has been a growing research focus on building LLM-based agent architectures, in which multiple models are composed and coordinated to solve complex tasks [104, 105]. We can place agents on a spectrum based on their level of autonomy in order to differentiate their architectures. On one end, fully autonomous systems, in which agents interact with significant freedom; on the other, predictable and structured workflows, in which agents follow predefined steps and communication patterns [106]. More structured and predetermined workflows might be preferable for domains that require more precision and accountability, such as mathematical reasoning, scientific research, law, medicine, and software development.

**Task Decomposition and Specialisation**: An effective strategy to enable LLMs to solve complex problems is to break them down into simpler, more manageable sub-tasks [107, 108, 54]. However, fully automating the planning and decomposition of tasks into sub-tasks with LLMs is an area of open research. Some studies propose decomposing a problem with a single LLM request that generates a series of sub-steps [109, 110, 111]. Other studies propose more advanced search-based approaches, which iterate and further decompose each sub-task into smaller steps when necessary; the final execution plan can then be organised into a tree-like structure [112]. One major challenge of this area of research is generating plans for domain-specific problems; the use of external planners is one of the solutions that have been proposed to address this issue [105]. Additionally, the planning abilities of LLMs might still be limited by their lack of human-level comprehension of world dynamics and the ability to apply causal reasoning to them [106].

Within this context, Mixture-of-Experts (MoE) models have emerged as a prominent line of research, offering a natural implementation of the *divide and conquer* paradigm. MoE architectures aim to scale model capacity efficiently by activating only a subset of experts per input, leading to improved performance [113, 114]. However, a core challenge lies in achieving effective *expert specialisation* — ensuring that each expert acquires unique and distinct skills, with little overlap with the others. Preliminary works suggests that expert models specialise on superficial patterns, such as token IDs, rather than extracting high-level semantic information [115]. Although recent efforts have proposed solutions towards more meaningful expert specialisation [116], we still lack a clear understanding of these mechanisms, especially when considering distributed and decentralised AI frameworks. Further research is needed to ensure that MoE models can robustly and adaptively decompose complex tasks in dynamic, multi-agent environments.

**Role-based Collaboration and Interactions**: In a compositional framework, interactions among agents represents a key ingredient. One of the most common strategies to compose multiple agents to work together in solving a task is the so-called role-based collaboration. LLM-based agents assume clearly defined, specialised roles (such as domain experts, assistants, etc.) in order to solve a higher level goal, with each of them being assigned individual sub-tasks by other agents. Optimal role assignments and agent adaptability to dynamic tasks requirements are still areas of open research [104]. One challenge is that, while LLMs are able to simulate many common roles, there are still many roles that they might not be able to capture accurately, such as uncommon roles rarely seen in the training corpus, or roles corresponding to human characters with particular cognitive-psychological traits [117, 118, 105].

Additionally, agent interaction play an essential role for multi-round tasks, where multiple iterative feedbacks loops — both from the environment and from other agents — are required to achieve the objective. Such tasks demand dynamic coordination, contextual adaptation, and the ability to reason over partial progress. Recent studies have begun to explore this promising topic, introducing novel solutions for LLM collaboration [119, 120]. This underscores the relevance of multi-agent interactions as a key challenge for CCO.

**Propagation of errors**: Lastly, the propagation of errors is an additional open problem in ensuring robustness in model integration. Erroneous outputs, hallucinations and biases from one agent can have cascading effects, getting amplified and spread through model interactions and impacting the whole system [104].

### 4.3.3. Potential Solutions & Future Research Directions

The concept of multiple models collaborating within a shared environment is well established in the AI field, and is rooted in the foundational definition of **multi-agent systems (MAS)**. Such systems consist of multiple autonomous agents, each with its own goals and motivations, that are capable of interacting with one another. For this reason, a central focus of this paradigm is on **cooperation** and **coordination** among agents.

The CCO framework can be viewed as a concrete instantiation of MAS, where models such as LLMs must be orchestrated to achieve a common objective. A key research direction in this context is the dynamic selection and composition of the most appropriate agents for a given task. Indeed, in the CCO framework the goal is to enable automatic models composition, that can evolve dynamically over time in response to variations in the environment or in the task given by the user. This marks a significant advancement over existing compositionality frameworks, *e.g.* LangGraph[1], AutoGen [121], in which the orchestration is largely static and predefined by the programmer, thereby limiting both flexibility and the capacity to generalise across tasks and environments.

Another component enabling continual learning capabilities in the CCO framework could be the use of experience accumulation modules, in the form of memory modules and skill libraries, as proposed in works like GITM, Voyager, AppAgent and MemPrompt [122, 123, 124, 125]. These approaches allow models to dynamically acquire new knowledge and skills, as a result of interactions with other agents, humans and the environment [105, 108, 126]. Such knowledge and skills can be stored in natural language or code form, and later retrieved and added to the model input context as needed. These mechanisms rely on test-time inference and incorporate the new knowledge and skills in the input context window of the model. One limitation might be that, depending on the LLM used, the context window capacity might limit the amount of task information that can be incorporated; however, research advancements in this area are enabling ultra-long context windows of 1M tokens or more [16], albeit at the expense of added compute.

An additional crucial aspect lies in the communication between agents. In the case of LLM agents, communication can occur through natural language, which has the added benefit of being easily interpretable by human users — a property that enhances transparency and human-in-the-loop control [127]. However, the CCO framework is designed to be model-agnostic and general-purpose, extending beyond language models to integrate a diverse set of AI components — such as computer vision models, time series processors, rule-based systems, symbolic modules or even hard-coded functions.

To support such heterogeneity, the system requires a robust communication protocol, that accommodates decentralisation, asynchrony and different data formats, while enabling efficient knowledge exchange.

Furthermore, an important open question in the design of the CCO protocol is in what type of knowledge should be shared among agents, to maximise collaboration without unnecessary overheads. Depending on the use case, this could include (i) model parameters, either in entirety or specific subnetworks / modules; (ii) the model internal representation, *e.g.* latent vectors or output logits; (iii) training data, as raw samples or abstracted via data generators.

## 5. The Future of Continual Learning: From Niche Research to AI's Next Paradigm

Over the past decade, Continual Learning (CL) has established itself as a prominent research area exploring a variety of domains, problem settings and applications [25]. In Computer Vision, significant effort has been committed to *class incremental* and *domain incremental* scenarios, where models must progressively recognise new categories without forgetting previously learned ones, even as the input domains evolve [128]. In Reinforcement Learning, CL has focused on the ability of agents to adapt dynamically to evolving environments while retaining past knowledge and abilities, particularly in *task incremental* and *multi-task* settings [129]. More recently, CL research has extended its scope to

---

[1]https://www.langchain.com/langgraph

include LLMs and Foundational Models, where the challenge lies in enabling models to continuously acquire new linguistic capabilities or domain knowledge without catastrophic forgetting [54, 9]. Another critical aspect is memory management in lifelong learning AI systems. Several studies examine the inherent trade-off between limited computational and storage resources and the ever-growing volume of data that CL models are expected to handle. Researchers have proposed a wide range of methods to strike a balance between learning efficiency and memory constraints such as dynamic memory buffers, experience replay mechanisms, architectural approaches, and regularisation techniques [25, 26]. Leveraging this attention from the community, CL has become a well-established and recognised field, providing both theoretical foundations and practical methodologies for building adaptive, robust, and memory-efficient AI systems.

Given the increasingly consolidated position of CL within the research community, we believe the time is right for the field to take a decisive step forward. Rather than representing a separate area of study, CL should become a critical and fundamental component in the prototyping and evolution of modern AI systems. In particular, the rapid rise of foundational models has captured the attention of both academia and industry, asserting itself as one of the most prominent and transformative trends in contemporary AI research. These models showcase broad generalisation capabilities across tasks and modalities; however, they still exhibit a critical limitation: their knowledge is inherently static and fixed at training time. This immutability poses a major challenge in dynamic real-world environments, where new data and information continuously emerge.

To address this gap, foundational models must evolve towards true continuous adaptation, progressively updating, refining, and extending their knowledge over time. Continual learning principles can provide a concrete and resourceful solution in this context. As previously discussed, **Continual Pre-training** and **Continual Fine-tuning** represent emerging research directions that aim to integrate the principles of CL with large-scale models, enabling them to remain updated, accurate, and contextually relevant, thereby mitigating outdated or incorrect outputs [45]. However, CPT and CFT inherently require relatively low-frequency adaptation cycles and depend on substantial computational resources. Thus, while beneficial, these approaches alone might not address the dynamic adaptability required by real-world applications fully.

Real-world change, in contrast, is often expressed at the level of *orchestration*: new tools appear [55], regulations shift, a CoT must be revised, or a group of agents must re-organise to solve an emergent sub-problem [13]. This layer is inherently *high-frequency and pervasive*, with potential updates required minutes or even seconds after new data arrives, making repeated offline training cycles impractical [15]. Here **Continual Compositionality and Orchestration** is not merely advantageous - it is indispensable. CCO treats an AI system as a living assembly of modules composed of prompt routers, domain experts, external tools, episodic memories — all of which can be composed and adapted on the fly. Continual learning supplies the two capabilities such a system requires: *rapid adaptation* to integrate the next tool or component and *memory consolidation* to stabilise useful compositions so that they can be re-used rather than rediscovered.

Foundational models are generally effective across broad tasks; however, when maximising performance in specialised areas, such as mathematics or physics, monolithic models may underperform compared to models specifically distilled or trained for those individual tasks [130, 131]. We believe that creating and deploying CL models which can continually evolve and combine knowledge from different sources represents a promising and sustainable architectural solution for AI systems. Moreover, developing models that are a mixture of many components [54, 132] can allow for the distribution of computation and decentralisation of AI systems. Different institutions can collaborate and contribute to create the ad-hoc models that synergise to achieve the final and complete AI system.

Furthermore, aligning with recent advances, foundational models can also be improved through human feedback within a Reinforcement Learning from Human Feedback (RLHF) framework, where humans actively guide model evolution, shaping future AI outputs and capabilities in turn [53]. This human-AI feedback loop can help to ensure that AI development is in line with human preferences and can also be applied within the context of multi-agent systems. Human-in-the-loop approaches in agentic LLM systems can be used to provide guidance, supervision and feedback to individual agents,

facilitating alignment with human preferences [108, 133, 134, 14].

In summary, Continual Learning, specifically continual compositionality and orchestration, represents not only a promising research direction, but the cornerstone of AI's next paradigm shift. By transitioning from incremental improvements within individual models to dynamically composable and collaborative AI ecosystems, CL can drive a new generation of adaptive, scalable and human-aligned AI systems.

## 6. Conclusion

The remarkable capabilities of large foundational models highlight their potential in solving complex, diverse tasks across multiple domains. However, despite their robust generalisation abilities, these models are inherently static and struggle to adapt continually to evolving real-world data and tasks. To overcome these limitations, continual learning emerges as an indispensable tool, offering a multitude of methodologies to enhance adaptability, efficiency and sustainability within foundational models.

In this paper, we have highlighted three pivotal areas of continual learning critical to the evolution of FMs: continual pre-training, continual fine-tuning, and continual compositionality and orchestration. CPT equips very large, organisation-scale models with the mechanisms to incrementally incorporate new knowledge, capabilities or methodologies, maintaining their relevance and mitigating catastrophic forgetting through adaptive, resource-efficient updates. As such, it is mainly within the remit of industrial laboratories and cloud providers who possess the necessary data and compute. CFT remains valuable, although comparatively secondary, to enable precise adaptation to specialised tasks and domains, effectively balancing specific task performance with the retention of generalisable knowledge. Techniques such as PEFT, meta-learning and model merging were identified as promising approaches to achieving effective adaptation while managing computational resources and limiting data drift.

CCO, by contrast, is where academic research can and should place its primary emphasis. Moving from monolithic models to a decentralised ecosystems composed of specialised, modular agents facilitates adaptability, enhances scalability and reduces centralisation risks by enabling modular model replacement, upgrade and collaborative interaction. By studying and advancing CCO, the research community can catalyse an open, decentralised, and circular economy of AI components. Such a decentralised ecosystem not only encourages innovation and democratises access, but also mitigates the computational and environmental costs associated with continually retraining large, static models.

Ultimately, continual learning positions itself not merely as an optional enhancement but as a foundational requirement for future AI systems. As AI evolves from static to dynamic, from centralised to decentralised, and from monolithic to modular, the integration of continual learning methodologies will be crucial. Embracing continual learning will therefore be instrumental in building resilient, flexible and context-aware AI systems, capable of sustainably adapting to the ever-changing landscape of real-world challenges.

## Acknowledgements

## References

[1] R. M. French, Catastrophic forgetting in connectionist networks, Trends in Cognitive Sciences 3 (1999) 128–135. URL: https://www.cell.com/trends/cognitive-sciences/abstract/S1364-6613(99)01294-2. doi:10.1016/S1364-6613(99)01294-2, publisher: Elsevier.

[2] G. Ditzler, M. Roveri, C. Alippi, R. Polikar, Learning in nonstationary environments: A survey, IEEE Computational Intelligence Magazine 10 (2015) 12–25.

[3] M. Ring, Continual learning in reinforcement environments, 1994. URL: https://www.proquest.com/openview/2d2f13eb52fc09d3eadfd0c81fe5f181/1?cbl=18750&diss=y&pq-origsite=gscholar.

[4] A. Kumar, H. D. III, Learning Task Grouping and Overlap in Multi-task Learning, 2012. URL: http://arxiv.org/abs/1206.6417. doi:10.48550/arXiv.1206.6417, arXiv:1206.6417 [cs].

[5] F. Giannini, G. Ziffer, A. Cossu, V. Lomonaco, Streaming Continual Learning for Unified Adaptive Intelligence in Dynamic Environments, IEEE Intelligent Systems 39 (2024) 81–85. URL: https://ieeexplore.ieee.org/document/10779199/?arnumber=10779199. doi:10.1109/MIS.2024.3479469, conference Name: IEEE Intelligent Systems.

[6] W. Li, F. Qi, R. Yan, H. Zhang, W. Lei, J. Tang, J. Luo, Continual Learning meets Multimodal Foundation Models: Fundamentals and Advances, in: Proceedings of the 1st on Continual Learning meets Multimodal Foundation Models: Fundamentals and Advances, ACMMM CL'24, Association for Computing Machinery, New York, NY, USA, 2024, pp. 1–4. URL: https://dl.acm.org/doi/10.1145/3688859.3690083. doi:10.1145/3688859.3690083.

[7] A. Cossu, A. Carta, L. Passaro, V. Lomonaco, T. Tuytelaars, D. Bacciu, Continual pre-training mitigates forgetting in language and vision, Neural Networks 179 (2024) 106492.

[8] J. Zheng, C. Shi, X. Cai, Q. Li, D. Zhang, C. Li, D. Yu, Q. Ma, Lifelong Learning of Large Language Model based Agents: A Roadmap, 2025. URL: http://arxiv.org/abs/2501.07278. doi:10.48550/arXiv.2501.07278, arXiv:2501.07278 [cs].

[9] H. Shi, Z. Xu, H. Wang, W. Qin, W. Wang, Y. Wang, Z. Wang, S. Ebrahimi, H. Wang, Continual learning of large language models: A comprehensive survey, arXiv preprint arXiv:2404.16789 (2024).

[10] G. M. Van de Ven, T. Tuytelaars, A. S. Tolias, Three types of incremental learning, Nature Machine Intelligence 4 (2022) 1185–1197.

[11] T. K. Hensch, Critical period plasticity in local cortical circuits, Nature reviews neuroscience 6 (2005) 877–888.

[12] Z. Chen, W. Wang, Y. Cao, Y. Liu, Z. Gao, E. Cui, J. Zhu, S. Ye, H. Tian, Z. Liu, L. Gu, X. Wang, Q. Li, Y. Ren, Z. Chen, J. Luo, J. Wang, T. Jiang, B. Wang, C. He, B. Shi, X. Zhang, H. Lv, Y. Wang, W. Shao, P. Chu, Z. Tu, T. He, Z. Wu, H. Deng, J. Ge, K. Chen, K. Zhang, L. Wang, M. Dou, L. Lu, X. Zhu, T. Lu, D. Lin, Y. Qiao, J. Dai, W. Wang, Expanding Performance Boundaries of Open-Source Multimodal Models with Model, Data, and Test-Time Scaling, 2025. URL: http://arxiv.org/abs/2412.05271. doi:10.48550/arXiv.2412.05271, arXiv:2412.05271 [cs].

[13] S. Schmidgall, M. Moor, AgentRxiv: Towards Collaborative Autonomous Research, 2025. URL: http://arxiv.org/abs/2503.18102. doi:10.48550/arXiv.2503.18102, arXiv:2503.18102 [cs].

[14] Z. Cai, B. Chang, W. Han, Human-in-the-loop through chain-of-thought, arXiv preprint arXiv:2306.07932 (2023).

[15] K. Kumar, T. Ashraf, O. Thawakar, R. M. Anwer, H. Cholakkal, M. Shah, M.-H. Yang, P. H. S. Torr, F. S. Khan, S. Khan, LLM Post-Training: A Deep Dive into Reasoning Large Language Models, 2025. URL: http://arxiv.org/abs/2502.21321. doi:10.48550/arXiv.2502.21321, arXiv:2502.21321 [cs].

[16] G. Team, P. Georgiev, V. I. Lei, R. Burnell, L. Bai, A. Gulati, G. Tanzer, D. Vincent, Z. Pan, S. Wang, et al., Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024.

[17] L. Hong, S. E. Page, Groups of diverse problem solvers can outperform groups of high-ability problem solvers, Proceedings of the National Academy of Sciences 101 (2004) 16385–16389.

[18] T. Lesort, V. Lomonaco, A. Stoian, D. Maltoni, D. Filliat, N. Díaz-Rodríguez, Continual Learning for Robotics: Definition, Framework, Learning Strategies, Opportunities and Challenges, 2019. URL: http://arxiv.org/abs/1907.00182. doi:10.48550/arXiv.1907.00182, arXiv:1907.00182 [cs].

[19] M. McCloskey, N. J. Cohen, Catastrophic Interference in Connectionist Networks: The Sequential Learning Problem, in: G. H. Bower (Ed.), Psychology of Learning and Motivation, volume 24, Academic Press, 1989, pp. 109–165. URL: https://www.sciencedirect.com/science/article/pii/S0079742108605368. doi:10.1016/S0079-7421(08)60536-8.

[20] R. Ratcliff, Connectionist models of recognition memory: constraints imposed by learning and

forgetting functions., Psychological review 97 (1990) 285.

[21] R. K. Srivastava, J. Masci, S. Kazerounian, F. Gomez, J. Schmidhuber, Compete to Compute, in: C. J. Burges, L. Bottou, M. Welling, Z. Ghahramani, K. Q. Weinberger (Eds.), Advances in Neural Information Processing Systems, volume 26, Curran Associates, Inc., 2013. URL: https://proceedings.neurips.cc/paper_files/paper/2013/file/8f1d43620bc6bb580df6e80b0dc05c48-Paper.pdf.

[22] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, Y. Bengio, An Empirical Investigation of Catastrophic Forgetting in Gradient-Based Neural Networks, 2015. URL: http://arxiv.org/abs/1312.6211. doi:10.48550/arXiv.1312.6211, arXiv:1312.6211 [stat].

[23] Z. Li, D. Hoiem, Learning without Forgetting, 2017. URL: http://arxiv.org/abs/1606.09282. doi:10.48550/arXiv.1606.09282, arXiv:1606.09282 [cs].

[24] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, R. Hadsell, Overcoming catastrophic forgetting in neural networks, Proceedings of the National Academy of Sciences 114 (2017) 3521–3526. URL: http://arxiv.org/abs/1612.00796. doi:10.1073/pnas.1611835114, arXiv:1612.00796 [cs].

[25] L. Wang, X. Zhang, H. Su, J. Zhu, A comprehensive survey of continual learning: Theory, method and application, IEEE Transactions on Pattern Analysis and Machine Intelligence (2024).

[26] B. Wickramasinghe, G. Saha, K. Roy, Continual learning: A review of techniques, challenges, and future directions, IEEE Transactions on Artificial Intelligence 5 (2023) 2526–2546.

[27] F. Zenke, B. Poole, S. Ganguli, Continual Learning Through Synaptic Intelligence, 2017. URL: http://arxiv.org/abs/1703.04200. doi:10.48550/arXiv.1703.04200, arXiv:1703.04200 [cs].

[28] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, R. Hadsell, Progressive Neural Networks, 2022. URL: http://arxiv.org/abs/1606.04671. doi:10.48550/arXiv.1606.04671, arXiv:1606.04671 [cs].

[29] A. Mallya, S. Lazebnik, PackNet: Adding Multiple Tasks to a Single Network by Iterative Pruning, 2018. URL: http://arxiv.org/abs/1711.05769. doi:10.48550/arXiv.1711.05769, arXiv:1711.05769 [cs].

[30] A. Mallya, D. Davis, S. Lazebnik, Piggyback: Adapting a Single Network to Multiple Tasks by Learning to Mask Weights, 2018. URL: http://arxiv.org/abs/1801.06519. doi:10.48550/arXiv.1801.06519, arXiv:1801.06519 [cs].

[31] J. Frankle, M. Carbin, The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks, 2019. URL: http://arxiv.org/abs/1803.03635. doi:10.48550/arXiv.1803.03635, arXiv:1803.03635 [cs].

[32] J. Serrà, D. Surís, M. Miron, A. Karatzoglou, Overcoming catastrophic forgetting with hard attention to the task, 2018. URL: http://arxiv.org/abs/1801.01423. doi:10.48550/arXiv.1801.01423, arXiv:1801.01423 [cs].

[33] M. Wortsman, V. Ramanujan, R. Liu, A. Kembhavi, M. Rastegari, J. Yosinski, A. Farhadi, Supermasks in Superposition, in: Advances in Neural Information Processing Systems, volume 33, Curran Associates, Inc., 2020, pp. 15173–15184. URL: https://proceedings.neurips.cc/paper/2020/hash/ad1f8bb9b51f023cdc80cf94bb615aa9-Abstract.html.

[34] A. Prabhu, P. H. S. Torr, P. K. Dokania, GDumb: A Simple Approach that Questions Our Progress in Continual Learning, in: A. Vedaldi, H. Bischof, T. Brox, J.-M. Frahm (Eds.), Computer Vision – ECCV 2020, volume 12347, Springer International Publishing, Cham, 2020, pp. 524–540. URL: https://link.springer.com/10.1007/978-3-030-58536-5_31. doi:10.1007/978-3-030-58536-5_31, series Title: Lecture Notes in Computer Science.

[35] D. Lopez-Paz, M. Ranzato, Gradient Episodic Memory for Continual Learning, 2022. URL: http://arxiv.org/abs/1706.08840. doi:10.48550/arXiv.1706.08840, arXiv:1706.08840 [cs].

[36] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, C. H. Lampert, iCaRL: Incremental Classifier and Representation Learning, 2017. URL: http://arxiv.org/abs/1611.07725. doi:10.48550/arXiv.1611.07725, arXiv:1611.07725 [cs].

[37] L. Pellegrini, G. Graffieti, V. Lomonaco, D. Maltoni, Latent Replay for Real-Time Continual Learning, 2020. URL: http://arxiv.org/abs/1912.01100. doi:10.48550/arXiv.1912.01100,

arXiv:1912.01100 [cs].

[38] G. Graffieti, D. Maltoni, L. Pellegrini, V. Lomonaco, Generative negative replay for continual learning, Neural Networks 162 (2023) 369–383.

[39] G. M. Van de Ven, A. S. Tolias, Three scenarios for continual learning, arXiv preprint arXiv:1904.07734 (2019).

[40] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, 2019. URL: http://arxiv.org/abs/1810.04805. doi:10.48550/arXiv.1810.04805, arXiv:1810.04805 [cs].

[41] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, D. Amodei, Language Models are Few-Shot Learners, 2020. URL: http://arxiv.org/abs/2005.14165. doi:10.48550/arXiv.2005.14165, arXiv:2005.14165 [cs].

[42] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, I. Sutskever, Learning Transferable Visual Models From Natural Language Supervision, 2021. URL: http://arxiv.org/abs/2103.00020. doi:10.48550/arXiv.2103.00020, arXiv:2103.00020 [cs].

[43] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. v. Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, E. Brynjolfsson, S. Buch, D. Card, R. Castellon, N. Chatterji, A. Chen, K. Creel, J. Q. Davis, D. Demszky, C. Donahue, M. Doumbouya, E. Durmus, S. Ermon, J. Etchemendy, K. Ethayarajh, L. Fei-Fei, C. Finn, T. Gale, L. Gillespie, K. Goel, N. Goodman, S. Grossman, N. Guha, T. Hashimoto, P. Henderson, J. Hewitt, D. E. Ho, J. Hong, K. Hsu, J. Huang, T. Icard, S. Jain, D. Jurafsky, P. Kalluri, S. Karamcheti, G. Keeling, F. Khani, O. Khattab, P. W. Koh, M. Krass, R. Krishna, R. Kuditipudi, A. Kumar, F. Ladhak, M. Lee, T. Lee, J. Leskovec, I. Levent, X. L. Li, X. Li, T. Ma, A. Malik, C. D. Manning, S. Mirchandani, E. Mitchell, Z. Munyikwa, S. Nair, A. Narayan, D. Narayanan, B. Newman, A. Nie, J. C. Niebles, H. Nilforoshan, J. Nyarko, G. Ogut, L. Orr, I. Papadimitriou, J. S. Park, C. Piech, E. Portelance, C. Potts, A. Raghunathan, R. Reich, H. Ren, F. Rong, Y. Roohani, C. Ruiz, J. Ryan, C. Ré, D. Sadigh, S. Sagawa, K. Santhanam, A. Shih, K. Srinivasan, A. Tamkin, R. Taori, A. W. Thomas, F. Tramèr, R. E. Wang, W. Wang, B. Wu, J. Wu, Y. Wu, S. M. Xie, M. Yasunaga, J. You, M. Zaharia, M. Zhang, T. Zhang, X. Zhang, Y. Zhang, L. Zheng, K. Zhou, P. Liang, On the Opportunities and Risks of Foundation Models, 2022. URL: http://arxiv.org/abs/2108.07258. doi:10.48550/arXiv.2108.07258, arXiv:2108.07258 [cs].

[44] C. Zhou, Q. Li, C. Li, J. Yu, Y. Liu, G. Wang, K. Zhang, C. Ji, Q. Yan, L. He, et al., A comprehensive survey on pretrained foundation models: A history from bert to chatgpt, International Journal of Machine Learning and Cybernetics (2024) 1–65.

[45] Y. Yang, J. Zhou, X. Ding, T. Huai, S. Liu, Q. Chen, Y. Xie, L. He, Recent advances of foundation language models-based continual learning: A survey, ACM Computing Surveys 57 (2025) 1–38.

[46] A. Plaat, A. Wong, S. Verberne, J. Broekens, N. v. Stein, T. Back, Reasoning with Large Language Models, a Survey, 2024. URL: http://arxiv.org/abs/2407.11511. doi:10.48550/arXiv.2407.11511, arXiv:2407.11511 [cs] version: 1.

[47] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. d. L. Casas, L. A. Hendricks, J. Welbl, A. Clark, T. Hennigan, E. Noland, K. Millican, G. v. d. Driessche, B. Damoc, A. Guy, S. Osindero, K. Simonyan, E. Elsen, J. W. Rae, O. Vinyals, L. Sifre, Training Compute-Optimal Large Language Models, 2022. URL: http://arxiv.org/abs/2203.15556. doi:10.48550/arXiv.2203.15556, arXiv:2203.15556 [cs].

[48] R. Schwartz, J. Dodge, N. A. Smith, O. Etzioni, Green AI, Commun. ACM 63 (2020) 54–63. URL: https://dl.acm.org/doi/10.1145/3381831. doi:10.1145/3381831.

[49] E. Strubell, A. Ganesh, A. McCallum, Energy and Policy Considerations for Modern Deep Learning Research, Proceedings of the AAAI Conference on Artificial Intelligence 34 (2020) 13693–13696. URL: https://ojs.aaai.org/index.php/AAAI/article/view/7123. doi:10.1609/aaai.v34i09.7123, number: 09.

[50] H. R. Kirk, B. Vidgen, P. Röttger, S. A. Hale, Personalisation within bounds: A risk taxonomy and policy framework for the alignment of large language models with personalised feedback, 2023. URL: http://arxiv.org/abs/2303.05453. doi:10.48550/arXiv.2303.05453, arXiv:2303.05453 [cs].

[51] Y. Sun, X. Wang, Z. Liu, J. Miller, A. A. Efros, M. Hardt, Test-Time Training for Out-of-Distribution Generalization (2019). URL: https://openreview.net/forum?id=HyezmlBKwr.

[52] E. M. Bender, T. Gebru, A. McMillan-Major, S. Shmitchell, On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?, in: Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT '21, Association for Computing Machinery, New York, NY, USA, 2021, pp. 610–623. URL: https://dl.acm.org/doi/10.1145/3442188.3445922. doi:10.1145/3442188.3445922.

[53] D. Pedreschi, L. Pappalardo, E. Ferragina, R. Baeza-Yates, A.-L. Barabási, F. Dignum, V. Dignum, T. Eliassi-Rad, F. Giannotti, J. Kertész, et al., Human-ai coevolution, Artificial Intelligence (2024) 104244.

[54] T. Raheja, N. Pochhi, Foundation models meet continual learning: Recent advances, challenges, and future directions, in: NeurIPS 2024 Workshop on Scalable Continual Learning for Lifelong Foundation Models, 2024.

[55] Y. Shen, K. Song, X. Tan, D. Li, W. Lu, Y. Zhuang, HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in Hugging Face, 2023. URL: http://arxiv.org/abs/2303.17580. doi:10.48550/arXiv.2303.17580, arXiv:2303.17580 [cs].

[56] K. Roth, V. Udandarao, S. Dziadzio, A. Prabhu, M. Cherti, O. Vinyals, O. Hénaff, S. Albanie, M. Bethge, Z. Akata, A practitioner's guide to continual multimodal pretraining, arXiv preprint arXiv:2408.14471 (2024).

[57] S. Küchemann, K. E. Avila, Y. Dinc, C. Hortmann, N. Revenga, V. Ruf, N. Stausberg, S. Steinert, F. Fischer, M. Fischer, et al., On opportunities and challenges of large multimodal foundation models in education, npj Science of Learning 10 (2025) 11.

[58] J. Parmar, S. Satheesh, M. Patwary, M. Shoeybi, B. Catanzaro, Reuse, don't retrain: A recipe for continued pretraining of language models, arXiv preprint arXiv:2407.07263 (2024).

[59] The Llama 4 herd: The beginning of a new era of natively multimodal AI innovation, https://ai.meta.com/blog/llama-4-multimodal-intelligence/, 2025.

[60] A. Golatkar, A. Achille, S. Soatto, Eternal sunshine of the spotless net: Selective forgetting in deep networks, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020, pp. 9304–9312.

[61] L. Wang, X. Zeng, J. Guo, K.-F. Wong, G. Gottlob, Selective forgetting: Advancing machine unlearning techniques and evaluation in language models, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 39, 2025, pp. 843–851.

[62] H. Zhao, B. Ni, J. Fan, Y. Wang, Y. Chen, G. Meng, Z. Zhang, Continual forgetting for pre-trained vision models, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024, pp. 28631–28642.

[63] B. Zhu, G. Cui, Y. Chen, Y. Qin, L. Yuan, C. Fu, Y. Deng, Z. Liu, M. Sun, M. Gu, Removing backdoors in pre-trained models by regularized continual pre-training, Transactions of the Association for Computational Linguistics 11 (2023) 1608–1623.

[64] M. Brinner, T. A. Mustafa, S. Zarrieß, Enhancing domain-specific encoder models with llm-generated data: How to leverage ontologies, and how to do without them, arXiv preprint arXiv:2503.22006 (2025).

[65] C.-A. Li, H.-Y. Lee, Examining forgetting in continual pre-training of aligned large language models, arXiv preprint arXiv:2401.03129 (2024).

[66] S. V. Mehta, D. Patil, S. Chandar, E. Strubell, An empirical investigation of the role of pre-training in lifelong learning, Journal of Machine Learning Research 24 (2023) 1–50.

[67] T. Wu, L. Luo, Y.-F. Li, S. Pan, T.-T. Vu, G. Haffari, Continual learning for large language models: A survey, arXiv preprint arXiv:2402.01364 (2024).

[68] H. Zhu, Y. Wei, X. Liang, C. Zhang, Y. Zhao, Ctp: Towards vision-language continual pretraining

via compatible momentum contrast and topology preservation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 22257–22267.

[69] G. Zhang, L. Wang, G. Kang, L. Chen, Y. Wei, Slca: Slow learner with classifier alignment for continual learning on a pre-trained model, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 19148–19158.

[70] Ç. Yıldız, N. K. Ravichandran, N. Sharma, M. Bethge, B. Ermis, Investigating continual pretraining in large language models: Insights and implications, arXiv preprint arXiv:2402.17400 (2024).

[71] Y. Guo, J. Fu, H. Zhang, D. Zhao, Y. Shen, Efficient continual pre-training by mitigating the stability gap, arXiv preprint arXiv:2406.14833 (2024).

[72] X. Li, L. Peng, Y.-P. Wang, W. Zhang, Open challenges and opportunities in federated foundation models towards biomedical healthcare, BioData Mining 18 (2025) 2.

[73] M. Mendieta, B. Han, X. Shi, Y. Zhu, C. Chen, Towards geospatial foundation models via continual pretraining, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 16806–16816.

[74] N. D. Noce, M. Resta, D. Bacciu, Sequential continual pre-training for neural machine translation, in: 32nd European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning, ESANN, 2024. URL: https://doi.org/10.14428/esann/2024.ES2024-165.

[75] Y. Xie, K. Aggarwal, A. Ahmad, Efficient continual pre-training for building domain specific large language models, in: Findings of the Association for Computational Linguistics ACL 2024, 2024, pp. 10184–10201.

[76] O. Ostapenko, T. Lesort, P. Rodriguez, M. R. Arefin, A. Douillard, I. Rish, L. Charlin, Continual learning with foundation models: An empirical study of latent replay, in: Conference on lifelong learning agents, PMLR, 2022, pp. 60–91.

[77] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. de Laroussilhe, A. Gesmundo, M. Attariyan, S. Gelly, Parameter-efficient transfer learning for nlp, 2019. URL: https://arxiv.org/abs/1902.00751. arXiv:1902.00751.

[78] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen, Lora: Low-rank adaptation of large language models, 2021. URL: https://arxiv.org/abs/2106.09685. arXiv:2106.09685.

[79] D. Aggarwal, S. Damle, N. Goyal, S. Lokam, S. Sitaram, Exploring continual fine-tuning for enhancing language ability in large language model, 2024. URL: https://arxiv.org/abs/2410.16006. arXiv:2410.16006.

[80] S. Beaulieu, L. Frati, T. Miconi, J. Lehman, K. O. Stanley, J. Clune, N. Cheney, Learning to continually learn, 2020. URL: https://arxiv.org/abs/2002.09571. arXiv:2002.09571.

[81] X. Zhang, L. Bai, X. Yang, J. Liang, C-lora: Continual low-rank adaptation for pre-trained models, 2025. URL: https://arxiv.org/abs/2502.17920. arXiv:2502.17920.

[82] T. Fan, Y. Kang, G. Ma, W. Chen, W. Wei, L. Fan, Q. Yang, Fate-llm: A industrial grade federated learning framework for large language models, 2023. URL: https://arxiv.org/abs/2310.10049. arXiv:2310.10049.

[83] J. Liu, J. Ren, R. Jin, Z. Zhang, Y. Zhou, P. Valduriez, D. Dou, Fisher information-based efficient curriculum federated learning with large language models, 2024. URL: https://arxiv.org/abs/2410.00131. arXiv:2410.00131.

[84] A. Soutif-Cormerais, A. Carta, A. Cossu, J. Hurtado, V. Lomonaco, J. Van de Weijer, H. Hemati, A comprehensive empirical evaluation on online continual learning, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 3518–3528.

[85] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, S. Wermter, Continual lifelong learning with neural networks: A review, Neural Netw. 113 (2019) 54–71. URL: https://doi.org/10.1016/j.neunet.2019.01.012. doi:10.1016/j.neunet.2019.01.012.

[86] H. Du, S. Liu, L. Zheng, Y. Cao, A. Nakamura, L. Chen, Privacy in fine-tuning large language models: Attacks, defenses, and future directions, 2025. URL: https://arxiv.org/abs/2412.16504. arXiv:2412.16504.

[87] Z. Wang, Z. Zhang, C.-Y. Lee, H. Zhang, R. Sun, X. Ren, G. Su, V. Perot, J. Dy, T. Pfister, Learning to prompt for continual learning, 2022. URL: https://arxiv.org/abs/2112.08654. arXiv:2112.08654.

[88] J. S. Smith, L. Karlinsky, V. Gutta, P. Cascante-Bonilla, D. Kim, A. Arbelle, R. Panda, R. Feris, Z. Kira, Coda-prompt: Continual decomposed attention-based prompting for rehearsal-free continual learning, 2023. URL: https://arxiv.org/abs/2211.13218. arXiv:2211.13218.

[89] H. Chen, J. Li, N. Gazagnadou, W. Zhuang, C. Chen, L. Lyu, Dual low-rank adaptation for continual learning with pre-trained models, 2024. URL: https://arxiv.org/abs/2411.00623. arXiv:2411.00623.

[90] W. Zhang, Y. Huang, T. Zhang, Q. Zou, W.-S. Zheng, R. Wang, Adapter learning in pretrained feature extractor for continual learning of diseases, 2023. URL: https://arxiv.org/abs/2304.09042. arXiv:2304.09042.

[91] X. Gao, S. Dong, Y. He, Q. Wang, Y. Gong, Beyond prompt learning: Continual adapter for efficient rehearsal-free continual learning, 2024. URL: https://arxiv.org/abs/2407.10281. arXiv:2407.10281.

[92] P. Yadav, D. Tam, L. Choshen, C. Raffel, M. Bansal, Ties-merging: Resolving interference when merging models, 2023. URL: https://arxiv.org/abs/2306.01708. arXiv:2306.01708.

[93] L. Yu, B. Yu, H. Yu, F. Huang, Y. Li, Language models are super mario: Absorbing abilities from homologous models as a free lunch, 2024. URL: https://arxiv.org/abs/2311.03099. arXiv:2311.03099.

[94] D. Marczak, B. Twardowski, T. Trzciński, S. Cygert, Magmax: Leveraging model merging for seamless continual learning, 2024. URL: https://arxiv.org/abs/2407.06322. arXiv:2407.06322.

[95] E. Yang, L. Shen, Z. Wang, G. Guo, X. Chen, X. Wang, D. Tao, Representation surgery for multi-task model merging, 2024. URL: https://arxiv.org/abs/2402.02705. arXiv:2402.02705.

[96] E. N. Coleman, L. Quarantiello, J. Hurtado, V. Lomonaco, Adaptive LoRA merging for efficient domain incremental learning, in: Adaptive Foundation Models: Evolving AI for Personalized and Efficient Learning, 2024. URL: https://openreview.net/forum?id=tlB5eonGEk.

[97] C. Finn, P. Abbeel, S. Levine, Model-agnostic meta-learning for fast adaptation of deep networks, 2017. URL: https://arxiv.org/abs/1703.03400. arXiv:1703.03400.

[98] R. Zhang, R. Qiang, S. A. Somayajula, P. Xie, Autolora: Automatically tuning matrix ranks in low-rank adaptation based on meta learning, 2024. URL: https://arxiv.org/abs/2403.09113. arXiv:2403.09113.

[99] W. Li, L. Zou, M. Tang, Q. Yu, W. Li, C. Li, META-LORA: Memory-efficient sample reweighting for fine-tuning large language models, in: O. Rambow, L. Wanner, M. Apidianaki, H. Al-Khalifa, B. D. Eugenio, S. Schockaert (Eds.), Proceedings of the 31st International Conference on Computational Linguistics, Association for Computational Linguistics, Abu Dhabi, UAE, 2025, pp. 8504−8517. URL: https://aclanthology.org/2025.coling-main.568/.

[100] L. Palazzo, M. Pennisi, F. P. Salanitri, G. Bellitto, S. Palazzo, C. Spampinato, Fedrewind: Rewinding continual model exchange for decentralized federated learning, 2024. URL: https://arxiv.org/abs/2411.09842. arXiv:2411.09842.

[101] F. Chollet, Abstraction and Reasoning Corpus for Artificial General Intelligence (ARC-AGI), https://github.com/fchollet/ARC-AGI, 2019. URL: https://github.com/fchollet/ARC-AGI.

[102] A. Srivastava, A. Rastogi, A. Rao, A. A. M. Shoeb, A. Abid, A. Fisch, A. R. Brown, A. Santoro, A. Gupta, A. Garriga-Alonso, et al., Beyond the imitation game: Quantifying and extrapolating the capabilities of language models, arXiv preprint arXiv:2206.04615 (2022).

[103] C. Luo, Has llm reached the scaling ceiling yet? unified insights into llm regularities and constraints, arXiv preprint arXiv:2412.16443 (2024).

[104] K.-T. Tran, D. Dao, M.-D. Nguyen, Q.-V. Pham, B. O'Sullivan, H. D. Nguyen, Multi-agent collaboration mechanisms: A survey of llms, arXiv preprint arXiv:2501.06322 (2025).

[105] L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin, et al., A survey on large language model based autonomous agents, Frontiers of Computer Science 18 (2024) 186345.

[106] B. Liu, X. Li, J. Zhang, J. Wang, T. He, S. Hong, H. Liu, S. Zhang, K. Song, K. Zhu, Y. Cheng, S. Wang, X. Wang, Y. Luo, H. Jin, P. Zhang, O. Liu, J. Chen, H. Zhang, Z. Yu, H. Shi, B. Li, D. Wu, F. Teng, X. Jia, J. Xu, J. Xiang, Y. Lin, T. Liu, T. Liu, Y. Su, H. Sun, G. Berseth, J. Nie, I. Foster, L. Ward, Q. Wu,

Y. Gu, M. Zhuge, X. Tang, H. Wang, J. You, C. Wang, J. Pei, Q. Yang, X. Qi, C. Wu, Advances and challenges in foundation agents: From brain-inspired intelligence to evolutionary, collaborative, and safe systems, 2025. URL: https://arxiv.org/abs/2504.01990. arXiv:2504.01990.

[107] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al., Chain-of-thought prompting elicits reasoning in large language models, Advances in neural information processing systems 35 (2022) 24824–24837.

[108] Z. Xi, W. Chen, X. Guo, W. He, Y. Ding, B. Hong, M. Zhang, J. Wang, S. Jin, E. Zhou, et al., The rise and potential of large language model based agents: A survey, Science China Information Sciences 68 (2025) 121101.

[109] b. ichter, A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang, R. Julian, D. Kalashnikov, S. Levine, Y. Lu, C. Parada, K. Rao, P. Sermanet, A. T. Toshev, V. Vanhoucke, F. Xia, T. Xiao, P. Xu, M. Yan, N. Brown, M. Ahn, O. Cortes, N. Sievers, C. Tan, S. Xu, D. Reyes, J. Rettinghouse, J. Quiambao, P. Pastor, L. Luu, K.-H. Lee, Y. Kuang, S. Jesmonth, N. J. Joshi, K. Jeffrey, R. J. Ruano, J. Hsu, K. Gopalakrishnan, B. David, A. Zeng, C. K. Fu, Do as i can, not as i say: Grounding language in robotic affordances, in: K. Liu, D. Kulic, J. Ichnowski (Eds.), Proceedings of The 6th Conference on Robot Learning, volume 205 of *Proceedings of Machine Learning Research*, PMLR, 2023, pp. 287–318. URL: https://proceedings.mlr.press/v205/ichter23a.html.

[110] B. Xu, Z. Peng, B. Lei, S. Mukherjee, Y. Liu, D. Xu, Rewoo: Decoupling reasoning from observations for efficient augmented language models, CoRR abs/2305.18323 (2023). URL: https://doi.org/10.48550/arXiv.2305.18323.

[111] S. S. Raman, V. Cohen, E. Rosen, I. Idrees, D. Paulius, S. Tellex, Planning with large language models via corrective re-prompting, in: NeurIPS 2022 Foundation Models for Decision Making Workshop, 2022. URL: https://openreview.net/forum?id=cMDMRBe1TKs.

[112] S. Yao, D. Yu, J. Zhao, I. Shafran, T. Griffiths, Y. Cao, K. Narasimhan, Tree of thoughts: Deliberate problem solving with large language models, Advances in neural information processing systems 36 (2023) 11809–11822.

[113] N. Du, Y. Huang, A. M. Dai, S. Tong, D. Lepikhin, Y. Xu, M. Krikun, Y. Zhou, A. W. Yu, O. Firat, et al., Glam: Efficient scaling of language models with mixture-of-experts, in: International conference on machine learning, PMLR, 2022, pp. 5547–5569.

[114] B. Lin, Z. Tang, Y. Ye, J. Cui, B. Zhu, P. Jin, J. Huang, J. Zhang, Y. Pang, M. Ning, et al., Moe-llava: Mixture of experts for large vision-language models, arXiv preprint arXiv:2401.15947 (2024).

[115] F. Xue, Z. Zheng, Y. Fu, J. Ni, Z. Zheng, W. Zhou, Y. You, Openmoe: An early effort on open mixture-of-experts language models, arXiv preprint arXiv:2402.01739 (2024).

[116] D. Dai, C. Deng, C. Zhao, R. Xu, H. Gao, D. Chen, J. Li, W. Zeng, X. Yu, Y. Wu, et al., Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models, arXiv preprint arXiv:2401.06066 (2024).

[117] K. A. Fischer, Reflective linguistic programming (rlp): A stepping stone in socially-aware agi (socialagi), arXiv preprint arXiv:2305.12647 (2023).

[118] C. Li, J. Wang, K. Zhu, Y. Zhang, W. Hou, J. Lian, X. Xie, Emotionprompt: Leveraging psychology for large language models enhancement via emotional stimulus, arXiv preprint arXiv:2307.11760 (2023).

[119] X. Wang, S. Zhang, S. Li, K. Li, K. Kallidromitis, Y. Kato, K. Kozuka, T. Darrell, Segllm: Multi-round reasoning segmentation with large language models, in: The Thirteenth International Conference on Learning Representations, ????

[120] Z. Zhou, X. Zhang, S. Tan, L. Zhang, C. Li, Collaborative evolution: Multi-round learning between large and small language models for emergent fake news detection, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 39, 2025, pp. 1210–1218.

[121] Q. Wu, G. Bansal, J. Zhang, Y. Wu, B. Li, E. Zhu, L. Jiang, X. Zhang, S. Zhang, J. Liu, et al., Autogen: Enabling next-gen llm applications via multi-agent conversation, arXiv preprint arXiv:2308.08155 (2023).

[122] G. Wang, Y. Xie, Y. Jiang, A. Mandlekar, C. Xiao, Y. Zhu, L. Fan, A. Anandkumar, Voyager: An

open-ended embodied agent with large language models, Transactions on Machine Learning Research (2024). URL: https://openreview.net/forum?id=ehfRiF0R3a.

[123] X. Zhu, Y. Chen, H. Tian, C. Tao, W. Su, C. Yang, G. Huang, B. Li, L. Lu, X. Wang, Y. Qiao, Z. Zhang, J. Dai, Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory, 2023. URL: https://arxiv.org/abs/2305.17144. `arXiv:2305.17144`.

[124] C. Zhang, Z. Yang, J. Liu, Y. Han, X. Chen, Z. Huang, B. Fu, G. Yu, Appagent: Multimodal agents as smartphone users, arXiv preprint arXiv:2312.13771 (2023).

[125] A. Madaan, N. Tandon, P. Clark, Y. Yang, Memory-assisted prompt editing to improve gpt-3 after deployment, in: Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, 2022, pp. 2833–2861.

[126] T. Sumers, S. Yao, K. Narasimhan, T. Griffiths, Cognitive architectures for language agents, Transactions on Machine Learning Research (2024). URL: https://openreview.net/forum?id=1i6ZCvflQJ, survey Certification.

[127] G. Li, H. Hammoud, H. Itani, D. Khizbullin, B. Ghanem, Camel: Communicative agents for" mind" exploration of large language model society, Advances in Neural Information Processing Systems 36 (2023) 51991–52008.

[128] H. Liu, Y. Zhou, B. Liu, J. Zhao, R. Yao, Z. Shao, Incremental learning with neural networks for computer vision: a survey, Artificial intelligence review 56 (2023) 4557–4589.

[129] K. Khetarpal, M. Riemer, I. Rish, D. Precup, Towards continual reinforcement learning: A review and perspectives, Journal of Artificial Intelligence Research 75 (2022) 1401–1476.

[130] J. Gou, B. Yu, S. J. Maybank, D. Tao, Knowledge distillation: A survey, International Journal of Computer Vision 129 (2021) 1789–1819.

[131] C. Yang, Y. Zhu, W. Lu, Y. Wang, Q. Chen, C. Gao, B. Yan, Y. Chen, Survey on knowledge distillation for large language models: methods, evaluation, and application, ACM Transactions on Intelligent Systems and Technology (2024).

[132] A. Carta, A. Cossu, V. Lomonaco, D. Bacciu, Ex-model: Continual learning from a stream of trained models, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 3790–3799.

[133] Z. Kenton, T. Everitt, L. Weidinger, I. Gabriel, V. Mikulik, G. Irving, Alignment of language agents, arXiv preprint arXiv:2103.14659 (2021).

[134] W. Du, Z. M. Kim, V. Raheja, D. Kumar, D. Kang, Read, revise, repeat: A system demonstration for human-in-the-loop iterative text revision, in: Proceedings of the First Workshop on Intelligent and Interactive Writing Assistants (In2Writing 2022), 2022, pp. 96–108.