

Recent Advances and Emerging Frontiers in Machine Learning: Extending Foundational Research

Introduction

This report explores recent, high-impact research that builds upon seminal works in machine learning, including foundational models such as "DeepFace," "ArcFace," "Transformers," "GPT," and "FaceXformer." The primary focus is on developments published after 2023, aiming to uncover extensions to original ideas, improvements in architectures and benchmarks, new theoretical understandings, and emerging subfields. The investigation uses these landmark articles as seeds to discover and analyze new research across several key domains within machine learning. The subsequent sections are organized by these domains: Face Recognition, Natural Language Processing, Computer Vision, Reinforcement Learning, and General Machine Learning Theory or Optimization, providing a structured overview of the evolving landscape.

I. Advancements in Face Recognition

Face Recognition (FR) technology has undergone substantial evolution since the introduction of foundational deep learning models like DeepFace. Contemporary research endeavors focus on enhancing accuracy, improving robustness against challenging variations such as aging and pose, ensuring fairness across diverse demographic groups, and harnessing novel architectural paradigms, notably Transformers. Furthermore, the increasing sophistication of synthetic data generation techniques presents both new opportunities and challenges for the field.

A. Novel Loss Functions and Metric Learning (Beyond ArcFace)

The development of effective loss functions is central to training discriminative face recognition models. While ArcFace¹ introduced a significant improvement with its Additive Angular Margin Loss, pushing for better class separability, recent work explores augmenting such established losses or developing new ones to tackle specific, persistent challenges in FR. One notable direction is the integration of principles from other successful architectures, like Transformers, into the loss formulation itself. An example of this is presented in:

- **Full Title:** Transformer-Based Auxiliary Loss for Age-Invariant Face Recognition
- **Author(s):** Mahdi Abbasi, Sobhan Shafiei, Pooya Salehi, Akram Gholamzadeh, Mohammad H. Rohban, Mahdi Eftekhari
- **Year:** 2025 (based on arXiv:2412.02198 publication date; PDF indicates submission to a 2024 conference)

- **Key Contribution:** This paper proposes the addition of a transformer-based auxiliary loss to existing metric learning loss functions, such as ArcFace, to specifically enhance age-invariant face recognition. The transformer loss component, which processes the output of the final convolution layer of the main CNN backbone, is designed to capture global and long-range dependencies within feature maps. These dependencies are considered crucial for effectively handling the significant facial changes that occur due to aging. This approach aims to improve performance on age-challenging datasets without altering the structure of the primary metric loss function.²
- **Direct Link:** <https://arxiv.org/pdf/2412.02198>

This work directly extends the conceptual lineage of improving upon loss functions like ArcFace.¹ It exemplifies a trend where the strengths of Convolutional Neural Networks (CNNs) are hybridized with the capabilities of Transformers, not only in the main architectural backbone but also innovatively within the loss formulation itself.² The motivation stems from the understanding that while ArcFace and similar losses provide strong general discriminative power, certain FR sub-problems, like age-invariant recognition, present unique difficulties due to extreme intra-class variations over time.² Transformers, renowned for their proficiency in capturing long-range dependencies and global context, offer a complementary strength. This points towards a discernible trend: rather than solely focusing on inventing entirely new primary loss functions, researchers are finding value in augmenting well-established ones with *targeted auxiliary losses*. These auxiliary components are designed to bring in complementary architectural strengths to handle specific, difficult variations. Such a "plug-and-play" auxiliary loss strategy could offer a flexible mechanism for adapting robust, general-purpose FR models to a variety of specialized FR tasks—such as pose-invariant or occlusion-robust recognition—without necessitating complete retraining or a fundamental redesign of the core discriminative loss.

B. Mitigating Demographic Bias and Ensuring Fairness

A critical challenge in the deployment of FR systems is ensuring fairness and mitigating demographic bias, where performance varies significantly across groups defined by race, gender, or age. This issue has garnered substantial attention due to the ethical implications and real-world consequences of biased systems.³

- **Full Title:** Review of Demographic Fairness in Face Recognition
- **Author(s):** Ketan Kotwal, Sebastien Marcel
- **Year:** 2025 (arXiv:2502.02309v2, revised Apr 2025)
- **Key Contribution:** This comprehensive review consolidates the extensive body of research focused on demographic fairness in FR. It systematically examines the primary causes of demographic bias, including imbalanced training datasets, variability in skin-tone affecting sensor performance, algorithmic factors, and image quality issues. The paper also covers datasets available for fairness research, assessment metrics for quantifying disparities, and various mitigation approaches categorized into pre-processing, in-processing, and post-processing techniques. The authors emphasize that this is the first comprehensive review singularly dedicated to

demographic fairness in FR, highlighting the critical need for equitable and trustworthy systems, particularly as FR technologies are increasingly deployed in sensitive domains like law enforcement and security.⁴

- **Direct Link:** <https://arxiv.org/abs/2502.02309>

This survey underscores that despite significant research progress, demographic bias remains a persistent and critical hurdle in FR.⁴ The formal incorporation of demographic effects into the evaluation frameworks of prominent initiatives like the National Institute of Standards and Technology's (NIST) Face Recognition Vendor Tests (FRVT) since 2019 further highlights the recognized importance of addressing this issue.⁴

Addressing these biases often involves improving the data used for training. Synthetic data generation has emerged as a promising avenue:

- **Full Title:** VariFace: Fair and Diverse Synthetic Dataset Generation for Face Recognition
- **Author(s):** Michael Yeung, Toya Teramoto, Songtao Wu, et al.
- **Year:** 2024 (arXiv:2412.06235)
- **Key Contribution:** VariFace introduces a novel two-stage diffusion-based pipeline designed to create synthetic face datasets that are both demographically fair and exhibit high diversity in terms of intraclass and interclass variations. The methodology incorporates techniques such as Face Recognition Consistency for refining demographic labels, Face Vendi Score Guidance to enhance interclass diversity, and Divergence Score Conditioning to manage the trade-off between identity preservation and intraclass diversity. Significantly, the authors report that in unconstrained settings, VariFace can train FR models that outperform those trained on established real-world datasets like CASIA-WebFace.⁷
- **Direct Link:** <https://arxiv.org/abs/2412.06235>

The development of VariFace directly tackles the bias issues highlighted in the Kotwal and Marcel review⁴ by leveraging advanced generative models—specifically diffusion models, which are also noted for their capabilities in broader computer vision contexts.¹¹ The claim that models trained on VariFace can outperform those trained on real datasets marks a potentially significant milestone for the utility of synthetic data in FR.

These developments reveal a strong confluence of research in synthetic data generation and the imperative for demographic fairness in FR. Real-world datasets are often plagued by inherent biases and privacy issues.⁷ Synthetic data offers a pathway to create more balanced and diverse datasets, potentially mitigating these problems.⁷ The success of approaches like VariFace suggests that synthetic data is becoming a viable, and perhaps in some cases superior, alternative for training fair FR models. This could mean that the future of FR model development may lean heavily on highly controlled, ethically generated synthetic datasets. However, this also necessitates ongoing research into potential new biases that could be introduced by the synthetic generation processes themselves. The ability to finely *control* the generation process, as demonstrated in VariFace, is a key aspect of this approach.

Beyond algorithmic fairness achieved through data or model improvements, there's an emerging focus on human-in-the-loop systems and a deeper understanding of human-algorithm interaction in FR. This is particularly driven by documented real-world

failures and biases of fully automated systems, such as wrongful arrests based on erroneous FR matches, which have disproportionately affected certain racial groups.³ Research indicates that humans rarely produce false positive errors in face matching tasks and can be effective at correcting machine errors, especially when ML-derived similarity scores can predict potential algorithmic mistakes.³ While algorithms might outperform human annotators in straightforward or moderately difficult tasks, humans, particularly experts, can exhibit superior performance in complex, real-life scenarios or in instances where algorithms fail entirely.³ This body of evidence points to the necessity of hybrid human-machine strategies, not only for enhancing fairness but also for improving overall reliability and trustworthiness in critical FR applications. Future FR systems deployed in high-stakes environments will likely need to integrate mechanisms for human oversight and error correction. This, in turn, will drive research into optimizing the presentation of information to human reviewers to maximize their efficacy in identifying and rectifying machine errors, connecting closely with the broader field of Explainable AI (XAI) in FR.

C. Transformer-based Architectures in Face Recognition (Beyond FaceXformer)

The success of Transformer architectures in other domains has spurred their adoption and adaptation for face recognition and broader facial analysis tasks.

- **Full Title:** FaceXFormer: A Unified Transformer for Facial Analysis
- **Author(s):** Kartik Narayan, Vibashan VS, Rama Chellappa, Vishal M. Patel
- **Year:** 2025 (based on project page¹²; arXiv preprint 2024)
- **Key Contribution:** FaceXFormer stands out as an end-to-end unified transformer model engineered to perform ten distinct facial analysis tasks within a single, cohesive framework. These tasks include face recognition, face parsing, landmark detection, head pose estimation, attribute prediction, age, gender, and race estimation, facial expression recognition, and face visibility assessment. The model employs a transformer-based encoder-decoder architecture where each specific task is represented as a learnable token. This design enables seamless multi-task processing. To enhance efficiency, FaceXFormer introduces FaceX, a lightweight decoder featuring a novel bi-directional cross-attention mechanism that jointly processes face and task tokens to learn robust and generalized facial representations. The model has demonstrated state-of-the-art or competitive performance across multiple benchmarks, showcasing the versatility of Transformers for comprehensive facial analysis.¹²
- **Direct Link:** <https://kartik-3004.github.io/facexformer/> (Project Page), <https://arxiv.org/abs/2403.12960> (arXiv)

FaceXformer itself represents a significant development, illustrating the capacity of Transformer architectures to unify a wide array of face analysis tasks, including recognition, under one umbrella. Building on this theme of comprehensive facial understanding using advanced neural architectures, another notable recent development is:

- **Full Title:** Face-LLaVA: Facial Expression and Attribute Understanding through

Instruction Tuning

- **Author(s):** Ashutosh Chaubey, Xulang Guan, Mohammad Soleymani
- **Year:** 2025 (arXiv:2504.07198)
- **Key Contribution:** Face-LLaVA is a multimodal large language model (MLLM) specifically designed for face-centered in-context learning. It is capable of handling a variety of facial analysis tasks, such as facial expression recognition, attribute recognition, action unit (AU) detection, age estimation, and deepfake detection. A key feature of Face-LLaVA is its ability to generate natural language descriptions that can be used for reasoning about its predictions. The system leverages a novel face-specific visual encoder powered by Face-Region Guided Cross-Attention. While not primarily a face *recognition* method in the traditional sense of identity verification, Face-LLaVA represents a significant advancement in comprehensive face *understanding* by applying MLLM principles. This is thematically aligned with the evolution of face analysis beyond simple identification and is a notable development post-2023.¹⁴
- **Direct Link:** <https://arxiv.org/abs/2504.07198> (Project page link available within the paper)

The progression from specialized Transformer applications in vision (like the original Vision Transformer for image classification) to models like FaceXformer, and subsequently to systems like Face-LLaVA, indicates a significant evolution. FaceXformer demonstrated that a single Transformer could adeptly manage multiple, distinct facial analysis tasks by conceptualizing tasks as tokens.¹² This was a pivotal step in unifying the field of facial analysis. Face-LLaVA extends this trajectory by constructing a multimodal large language model specifically for faces. It not only performs tasks like expression and attribute recognition but, crucially, *generates natural language descriptions to facilitate reasoning*.¹⁴ This evolution signifies a shift: Transformers are moving from being primarily feature extractors for a singular task (like FR) to becoming the foundational backbones of more generalized "face understanding" systems. These systems can perceive, analyze, and reason about a multitude of facial aspects.

This trend suggests that the future of face-related AI may involve fewer siloed "face recognition" systems and more integrated "facial intelligence" platforms powered by MLLMs. Such platforms could enable richer human-computer interactions and more nuanced applications. However, this deeper level of facial information processing and interpretation also brings forth more complex ethical considerations that will need careful navigation.

D. 3D and Video Face Recognition

While 2D image-based FR has seen widespread adoption, research into 3D and video-based FR persists due to their inherent potential to offer more information and robustness, especially in unconstrained real-world scenarios.

- **Full Title:** A Comprehensive Review of Face Recognition Techniques, Trends and Challenges
- **Author(s):** (Authors not explicitly listed in the provided text ¹⁷, but the publication is IEEE Access, January 2024)

- **Year:** 2024
- **Key Contribution:** This review article published in IEEE Access provides an in-depth examination of a wide array of FR methodologies, covering aspects from feature extraction and pre-processing to face detection and classification algorithms. Significantly, its scope transcends still-image recognition to specifically include video-based FR. The review also analyzes datasets pertinent to both 2D and 3D FR techniques, highlighting ongoing applications and challenges in these areas.¹⁷
- **Direct Link:** (A specific link is not available in the provided text, but it would be accessible via the IEEE Xplore database using its DOI).

Foundational models like DeepFace were primarily designed for 2D images. However, the limitations of static 2D images in handling real-world complexities—such as motion, significant pose variations, and occlusions—have long motivated research into leveraging the richer data available from video sequences and 3D sensors. The 2024 survey¹⁷ confirms the continued relevance and active development in these multi-modal FR approaches. Challenges related to pose, illumination, and expression, which were identified even in earlier surveys of deep FR methods¹⁸, are inherently better addressed by the more comprehensive data captured by video or 3D sensors.

The sustained interest in 3D and video FR suggests a persistent need for robust solutions that can operate effectively in dynamic and unconstrained environments. As 3D capture technologies become increasingly ubiquitous (e.g., integrated into smartphones) and as computational capabilities continue to advance, a resurgence or continued strong focus on 3D and video FR methods can be anticipated. These approaches hold the promise of leveraging richer data streams to achieve more reliable and secure recognition, potentially overcoming some of the inherent limitations faced by purely 2D systems.

The advancements in Face Recognition are not occurring in isolated silos. Instead, progress appears to be driven by an interconnected "triad" of themes. Firstly, architectural innovations, particularly the adaptation and refinement of Transformers and their variants², are providing new backbones and capabilities. Secondly, the pursuit of fairness and the mitigation of demographic bias is a major thrust, often intertwined with the development and use of sophisticated synthetic data generation techniques.⁴ Thirdly, the continual refinement of advanced loss functions aims to improve discriminability and address specific challenges like age-related changes.² These three areas are synergistic: new architectures may necessitate novel loss considerations, and the generation of high-quality synthetic data often benefits from advanced generative architectures, which themselves are frequently Transformer-based (e.g., diffusion models). Thus, future breakthroughs in FR are likely to emerge from combined progress across these fronts, potentially culminating in state-of-the-art models that integrate novel Transformer variants, trained on highly realistic and fair synthetic data, and optimized using sophisticated, multi-component loss functions.

II. Breakthroughs in Natural Language Processing

The advent of the Transformer architecture and subsequent models like GPT has fundamentally reshaped the landscape of Natural Language Processing (NLP). Current

research trajectories are focused on constructing even larger and more proficient Large Language Models (LLMs), enhancing their reasoning abilities, factuality, and safety, improving their computational efficiency, and extending their functionalities to novel areas such as knowledge editing and multimodal understanding. A recent data-driven survey of LLM limitations (LLMs) between 2022 and 2024 highlights the explosive growth in LLM-related research, with investigations into their limitations growing even faster. Key areas of concern and active research include reasoning failures, generalization issues, hallucinations, bias, and security. Notably, research published on arXiv shows a trend towards safety, controllability, and multimodality, particularly from 2022 to 2024.¹⁹

A. Novel LLM Architectures and Efficiency Improvements

While Transformers form the backbone of most current LLMs, their quadratic computational complexity with respect to sequence length poses a significant bottleneck, especially as the demand for processing longer contexts increases.²² This challenge is spurring research into alternative architectures and a deeper understanding of Transformer variants to optimize for efficiency and other desirable properties like bias robustness.

- **Full Title:** Birdie: Efficient State Space Models through Bidirectional Input Processing and Specialized Pre-Training Objective Mixtures
- **Author(s):** Posu Chen, Hitu Seth, Parijat Dube, et al.
- **Year:** 2024 (arXiv:2411.01030)
- **Key Contribution:** Birdie is a novel training procedure designed for State Space Models (SSMs), aiming to significantly bolster their in-context retrieval capabilities without necessitating architectural modifications. This enhancement makes SSMs more competitive with Transformer models on tasks that are recall-intensive. The Birdie approach integrates bidirectional input processing with dynamic mixtures of specialized pre-training objectives, which are optimized using reinforcement learning. The work also introduces a new bidirectional SSM architecture that facilitates a seamless transition from bidirectional context processing to causal generation.²²
- **Direct Link:** <https://arxiv.org/abs/2411.01030>

This work directly addresses the efficiency imperative in LLMs. SSMs, including variants like Mamba and Hawk²², represent a class of alternative architectures explored for their potential to overcome the Transformer's scaling limitations. Birdie's contribution lies in enhancing the performance of these efficient alternatives, particularly on tasks where Transformers have traditionally held an advantage, by innovating at the training procedure level. This signifies a notable post-2023 direction in exploring architectures beyond the standard Transformer. Simultaneously, research continues to refine the Transformer paradigm itself, focusing on understanding and mitigating inherent issues.

- **Full Title:** From N-grams to Transformers and Back: A Comparative Behavioral Framework for Bias in Language Models
- **Author(s):** Debjanee Barua, Gaoge WANG, Elissa M. Redmiles, Michelle Mazurek, Wojtek Palubicki
- **Year:** 2025 (arXiv:2505.12381, May 2025)

- **Key Contribution:** This study investigates the propagation of bias across different language model architectures, specifically comparing n-gram models with Transformers. By systematically varying parameters such as context window size, the temporal provenance of training data, and architectural details within Transformers (e.g., depth, number of attention heads, attention types), the research finds that Transformer architectures exhibit greater robustness to contextual bias compared to n-gram models. The framework developed offers an interpretable tool for analyzing and understanding bias propagation mechanisms in language models.²⁴
- **Direct Link:** <https://arxiv.org/pdf/2505.12381>

While not proposing a new LLM architecture, this research provides crucial insights into how specific architectural choices within the dominant Transformer family influence critical issues like social bias. Such understanding is vital for designing future LLMs that are not only capable but also fairer and more equitable.

The tension between improving LLM capabilities and managing their computational cost is a primary driver for architectural diversity. The need for longer context processing in many applications further amplifies the efficiency concerns associated with Transformers.²² This has catalyzed the exploration of more efficient sequential models like SSMs. However, these alternatives often do not match Transformers in certain capabilities, particularly long-range in-context retrieval.²² Innovations like "Birdie"²² seek to narrow this gap by enhancing the capabilities of efficient architectures through novel training methodologies, rather than solely through architectural redesign. Concurrently, research such as that by Barua et al.²⁴ delves into the nuances of Transformer architectural choices to better understand and mitigate inherent problems like bias, indicating ongoing optimization within the Transformer paradigm itself. This suggests that the future LLM architectural landscape may not be monolithic but rather a diverse ecosystem. Highly efficient SSMs might be favored for tasks aligning with their strengths, while Transformers continue to be refined for applications demanding their unique capabilities. The overarching trend is a shift towards selecting the optimal architectural tool—and associated training procedure—for a given task, carefully balancing performance with computational resources and ethical considerations.

B. Enhancing Reasoning in LLMs

Despite their fluency, LLMs often fall short of human-level complex reasoning, including logical deduction, mathematical problem-solving, and multi-step inference.²⁵ Addressing this limitation is a major focus of current NLP research.

- **Full Title:** Advancing Reasoning in Large Language Models: Promising Methods and Approaches
- **Author(s):** Avinash Patil, Aryan Jadon
- **Year:** 2025 (arXiv:2502.03671v2, May 2025)
- **Key Contribution:** This survey provides a comprehensive review of emerging techniques designed to enhance the reasoning capabilities of LLMs. It categorizes existing methods into three main approaches: (1) prompting strategies (e.g., Chain-of-Thought (CoT), Self-Consistency, Tree-of-Thought (ToT)), which guide LLMs

through step-by-step reasoning; (2) architectural innovations (e.g., retrieval-augmented models, modular reasoning networks, neuro-symbolic integration); and (3) learning paradigms (e.g., fine-tuning with reasoning-specific datasets, reinforcement learning). The paper also explores evaluation frameworks and discusses significant open challenges, such as hallucinations and the fundamental differences between the statistical learning of LLMs and formal symbolic logic.²⁵

- **Direct Link:** <https://arxiv.org/abs/2502.03671>
- **Full Title:** Empowering LLMs with Logical Reasoning: A Comprehensive Survey
- **Author(s):** Fengxiang Cheng, Haoxuan Li, Fenrong Liu, Robert van Rooij, Kun Zhang, Zhouchen Lin
- **Year:** 2025 (arXiv:2502.15652v2, Feb 2025)
- **Key Contribution:** This survey focuses specifically on *logical* reasoning within LLMs, addressing challenges in both logical question answering and maintaining logical consistency. It offers a detailed taxonomy of methods, including solver-aided approaches, prompting techniques, and pretraining/fine-tuning strategies. Furthermore, it reviews commonly used benchmark datasets, evaluation metrics, and discusses promising future research directions, such as the integration of modal logic to handle uncertainty.²⁷
- **Direct Link:** <https://arxiv.org/abs/2502.15652>

These surveys reveal that enhancing LLM reasoning is being pursued through a multi-pronged strategy. LLMs, primarily based on architectures like Transformers (e.g., GPT), excel at learning statistical patterns from data but inherently lack explicit symbolic logic capabilities.²⁵ This deficiency leads to difficulties in tasks requiring complex, multi-step, or formal logical reasoning.²⁵ The identified solutions span various approaches: prompting strategies like Chain-of-Thought aim to elicit more structured, step-by-step reasoning from the existing model. Solver-aided methods, as detailed in ²⁷, involve translating natural language problems into a symbolic representation that can be processed by external logical solvers, with the LLM then translating the solver's output back into natural language; this effectively uses the LLM as an intelligent interface to traditional symbolic reasoning systems. Architectural innovations, such as neuro-symbolic models ²⁵, endeavor to integrate logical components or rule-based systems more directly into the LLM's architecture or operational workflow. Finally, specialized fine-tuning and pretraining regimes are being developed to imbue the models with improved intrinsic reasoning skills.

The pursuit of robust, human-like reasoning in LLMs will likely necessitate hybrid systems that effectively combine the pattern-recognition strengths of LLMs with the precise manipulative capabilities of symbolic AI. The key challenge lies in achieving seamless and efficient integration of these distinct paradigms. This points towards an important emerging subfield focused on neuro-symbolic LLMs and their applications.

C. Mitigating Hallucination and Bias

Hallucinations—the generation of plausible but factually incorrect or nonsensical information—and inherent biases remain critical limitations of LLMs.¹⁹ Significant research

effort is directed towards understanding their causes and developing mitigation strategies, particularly for Multimodal LLMs (MLLMs) where inconsistencies can arise between textual outputs and visual inputs.

- **Full Title:** Cross-Images Contrastive Decoding: Precise, Lossless Suppression of Language Priors in Large Vision-Language Models
- **Author(s):** Jianfei Zhao, Feng Zhang, Xin Sun, Chong Feng
- **Year:** 2025 (arXiv:2505.10634v2, May 2025)
- **Key Contribution:** This paper introduces Cross-Images Contrastive Decoding (CICD), a training-free method designed to reduce hallucinations in Large Vision-Language Models (LVLMs). CICD operates on the observation that problematic language priors, which often cause hallucinations, stem from the LLM backbone and remain consistent across different images. The method uses different images to construct negative visual contexts, allowing for the selective suppression of "detrimental" language priors (those causing hallucinations) while preserving "essential" priors (those necessary for fluency and coherence).²⁹
- **Direct Link:** <https://arxiv.org/abs/2505.10634>
- **Full Title:** Mitigating Hallucinations in Large Vision-Language Models via Summary-Guided Decoding
- **Author(s):** Kyungmin Min, Minbeom Kim, Kang-il Lee, Dongryeol Lee, Kyomin Jung
- **Year:** 2025 (Findings of NAACL 2025)
- **Key Contribution:** This work proposes Summary-Guided Decoding (SumGD), a novel decoding strategy that encourages LVLMs to focus more on the provided image information, thereby mitigating hallucinations. It achieves this by reducing the textual context available to the model through summarization and by selectively controlling the generation of tokens associated with image-related Parts-of-Speech (POS). The method aims to counteract the tendency of LVLMs to over-rely on language priors, especially when generating longer sequences, while striving to maintain overall text quality.³³
- **Direct Link:** <https://aclanthology.org/2025.findings-naacl.235.pdf>

A significant trend in combating LLM (and particularly LVLM) hallucinations involves the development of sophisticated decoding-time strategies, as opposed to relying solely on model retraining or architectural modifications. Hallucinations frequently arise from an over-reliance on the language priors learned by the LLM backbone.²⁹ Given that retraining these massive models is exceptionally costly and may not entirely resolve the underlying issue, training-free decoding methods offer a more agile approach. Both CICD²⁹ and SumGD³³ exemplify this trend. CICD employs contrastive decoding with cross-image negative contexts to selectively filter out detrimental priors. SumGD utilizes summaries to curtail the textual context and guides the decoding process based on image-related POS tags. This emphasis on modifying the generation process at inference time provides a more flexible and resource-efficient means of enhancing the faithfulness of already trained models. This could lead to the development of a diverse toolkit of decoding algorithms that can be dynamically applied based on the specific task or desired output characteristics (e.g., prioritizing factual

accuracy versus creative generation), potentially resulting in more adaptable and reliable LLM systems without the need for constant, expensive retraining cycles.

D. Knowledge Editing in LLMs

As LLMs are deployed, maintaining the accuracy and currency of their embedded knowledge is crucial. Knowledge Editing (KE) has emerged as a computationally efficient strategy to correct inaccuracies or update LLMs by modifying their internal parameters without full-scale retraining.³⁴

- **Full Title:** ConKE: Conceptualization-Augmented Knowledge Editing in Large Language Models for Commonsense Reasoning
- **Author(s):** Liyu Zhang, Weiqi Wang, Tianqing Fang, Yangqiu Song
- **Year:** 2025 (arXiv:2412.11418v2, May 2025)
- **Key Contribution:** ConKE is a novel knowledge editing framework specifically tailored for updating commonsense knowledge within LLMs. The system employs VERA, an automated commonsense plausibility verifier, to identify erroneous knowledge within the LLM that requires editing. For the identified inaccuracies, ConKE integrates conceptualization and instantiation techniques to enrich the semantic coverage of the edits. This approach aims to make the edits more generalizable, affecting not only the targeted piece of knowledge but also other potentially relevant yet implausible information stored by the LLM. ConKE is designed for end-to-end scalability and can handle open-ended knowledge structures, moving beyond traditional (head, relation, tail) triplets.³⁴
- **Direct Link:** <https://github.com/HKUST-KnowComp/ConKE> (Code), <https://arxiv.org/abs/2412.11418> (Paper)

Knowledge editing is rapidly evolving from methods focused on targeted factual corrections to more automated and sophisticated systems. These advanced KE frameworks aim to handle complex types of knowledge, such as commonsense, and strive for edits that generalize beyond the specific fact being altered. While LLMs store vast quantities of information, this knowledge can be inaccurate or become outdated, and retraining such large models is prohibitively expensive. KE offers a more efficient alternative.³⁴ Early KE techniques often concentrated on simple factual triplets. In contrast, ConKE³⁴ specifically addresses commonsense knowledge, which is typically more unstructured and demands more nuanced editing than straightforward facts. A key innovation in ConKE is the incorporation of automated detection of erroneous knowledge using an external verifier (VERA), coupled with conceptualization and instantiation processes. This allows edits to be more generalizable, influencing related concepts rather than just isolated pieces of information. The ambition for an "end-to-end scalable" pipeline³⁴ signals a drive to transform KE into a practical and robust tool for the ongoing maintenance and improvement of LLMs.

The development of effective and scalable knowledge editing techniques is critical for the long-term viability and trustworthiness of LLMs. If successful, such methods could enable LLMs to be continuously updated and corrected without undergoing full retraining cycles. This would make them more adaptable to new information and more resilient to propagating

errors. The particular focus on editing commonsense knowledge is vital for developing AI systems that exhibit more human-like understanding and reasoning.

E. Multimodal Large Language Models (MLLMs)

The capabilities of LLMs are being extended to process and integrate information from multiple modalities, such as text, images, audio, and video, leading to the rise of Multimodal Large Language Models (MLLMs).

- **Full Title:** Multimodal Large Language Models Can Significantly Advance Scientific Reasoning
- **Author(s):** Jiahua Dong, Zhibin Gou, Yifu Geng, et al.
- **Year:** 2025 (arXiv:2502.02871, Feb 2025)
- **Key Contribution:** This position paper argues that MLLMs, by their ability to integrate diverse data types like text and images, possess the potential to significantly advance scientific reasoning. It reviews the current state of MLLM applications in scientific domains, identifies challenges (such as open-source MLLMs like LLaVA, Qwen-VL, and LLaMA-3.2-Vision lagging behind closed-source counterparts like GPT-4o in complex reasoning tasks), and outlines future steps needed to harness their full potential in science.³⁷
- **Direct Link:** <https://arxiv.org/pdf/2502.02871>
- **Full Title:** Large Multimodal Models for Low-Resource Languages: A Survey
- **Author(s):** Mohamed Abdalla, Muhammad Abdul-Mageed, et al.
- **Year:** 2025 (Implied by survey content reviewing works up to 2024/early 2025)
- **Key Contribution:** This survey provides an analysis of LMMs tailored for low-resource (LR) languages, examining 106 studies that cover 75 different languages. It highlights that the majority of LMM development efforts are concentrated on high-resource languages, predominantly English. The paper discusses the unique challenges and specific techniques for developing LMMs for LR languages, where text-image combinations are the most common multimodal pairing, while more complex integrations involving video are less prevalent.³⁸
- **Direct Link:** <https://arxiv.org/pdf/2502.05568>

The MLLM field is experiencing rapid expansion, yet it faces a set of bifurcated challenges. On one hand, there is a push towards enhancing high-end reasoning capabilities, particularly in complex domains like scientific research which inherently require the synthesis of multimodal data.³⁷ On the other hand, there are significant efforts to broaden the accessibility and applicability of MLLMs, for instance, by adapting them for low-resource languages.³⁸ Despite this progress, critical gaps persist. Open-source MLLMs often trail proprietary models in sophisticated reasoning tasks.³⁷ Furthermore, there is a general scarcity of resources, including data and benchmarks, for non-English languages, which hinders the global reach of MLLM technology.³⁸ The focus in low-resource MLLM development is frequently on foundational tasks, such as text-image understanding, while more complex modalities like video remain less explored.³⁸

The MLLM landscape is thus evolving along two primary axes: increasing the *depth of*

capability (i.e., enabling sophisticated reasoning for high-resource scenarios) and broadening the *breadth of applicability* (i.e., making these models effective for diverse languages and varied contexts). Progress on both these fronts is essential if MLLMs are to become truly transformative and equitably beneficial technologies. This situation also underscores a pressing need for more robust open-source initiatives to bridge the reasoning capability gap with proprietary models and to foster development for a wider range of languages and modalities.

F. Efficient LLMs (Quantization, Pruning, Distillation)

The immense size of modern LLMs presents significant challenges for deployment, particularly in terms of memory footprint and computational requirements. Research into model compression techniques like quantization, pruning, and distillation is crucial for making these powerful models more accessible and practical.

- **Full Title:** Quantizing Large Language Models for Code Generation: A Differentiated Replication
- **Author(s):** Alessandro Giagnorio, Antonio Mastropaolo, Saima Afrin, Massimiliano Di Penta, Gabriele Bavota
- **Year:** 2025 (arXiv:2503.07103)
- **Key Contribution:** This study replicates and extends prior work on quantizing LLMs, focusing on recent and larger models (up to 34 billion parameters) specialized for code generation. The research investigates quantization down to 2-4 bits per parameter. The findings indicate that 4-bit precision can achieve a substantial average memory footprint reduction of 70% with only a limited impact on code generation performance. For more aggressive quantization levels (3 and 2 bits), the study shows that using a code-specific calibration dataset can help mitigate performance degradation.³⁹
- **Direct Link:** <https://doi.org/10.5281/zenodo.13752774> (Replication package), <https://arxiv.org/abs/2503.07103> (Paper)
- **Full Title:** When Reasoning Meets Compression: Benchmarking Compressed Large Reasoning Models on Complex Reasoning Tasks
- **Author(s):** Nan Zhang, Yusen Zhang, Prasenjit Mitra, Rui Zhang
- **Year:** 2025 (arXiv:2504.02010)
- **Key Contribution:** This paper benchmarks various compressed versions of DeepSeek-R1, a large reasoning model, using techniques such as quantization, distillation, and pruning. The evaluation is conducted across several complex reasoning datasets. A key finding is that the total parameter count of a model has a more significant impact on its knowledge memorization capabilities than on its reasoning ability. Additionally, the study observes that shorter model outputs generally achieve better performance across multiple benchmarks for both the original R1 model and its compressed variants.⁴¹
- **Direct Link:** <https://arxiv.org/abs/2504.02010>

These studies demonstrate that LLM compression techniques are maturing and yielding practical outcomes. They allow for significant reductions in model size with manageable

performance trade-offs, even for large, specialized models designed for demanding tasks like code generation or complex reasoning. The work on quantizing LLMs for code generation shows the feasibility of 4-bit quantization for models up to 34B parameters, resulting in a ~70% memory reduction without a substantial drop in performance.³⁹ This is a concrete and practical advancement. The research on compressing DeepSeek-R1 specifically investigates the impact of compression on large reasoning models, finding that reasoning capabilities can often be preserved more effectively than general knowledge memorization during the compression process.⁴¹ This suggests that compression is not limited to general-purpose LLMs but can be effectively applied to specialized, high-capability models, thereby making them more accessible.

The maturation of these compression techniques is poised to accelerate the democratization and wider adoption of powerful LLMs. It is increasingly plausible to envision more capable LLMs running efficiently on edge devices or within resource-constrained environments. The ability to apply compression selectively—for instance, prioritizing the preservation of reasoning abilities over broad, general knowledge—will be a key factor in optimizing these models for specific applications.

A dominant meta-theme in current NLP research is the intricate balancing act between three crucial objectives: enhancing LLM *capabilities* (such as advanced reasoning and multimodal understanding), improving their *efficiency* (through architectural innovation and compression techniques), and ensuring they are *trustworthy* (by mitigating hallucinations and biases, and enabling reliable knowledge editing). The initial success of models like GPT demonstrated immense capability, fueling a drive for even more powerful systems, evident in research on sophisticated reasoning²⁵ and MLLMs.³⁷ However, the sheer scale of these models introduces significant efficiency challenges, which are being addressed by investigations into novel architectures like SSMS (e.g., Birdie²²) and advanced compression methods.³⁹ Simultaneously, the widespread deployment and potential societal impact of LLMs bring critical trustworthiness issues to the forefront: the propensity for hallucinations²⁹, the presence of biases¹⁹, and the need for factual accuracy and updatability, addressed by knowledge editing techniques.³⁴ These three dimensions—capability, efficiency, and trustworthiness—are often in tension. For example, increasing model parameters might enhance capability but negatively impact efficiency, while certain reasoning techniques could inadvertently increase the risk of hallucination if not carefully controlled. Future breakthroughs in NLP will likely involve co-design across these three dimensions. The ideal LLM would be highly capable, resource-efficient, and verifiably trustworthy. Research efforts that holistically address this "LLM Trilemma," rather than focusing on one aspect in isolation, are expected to be the most impactful.

III. Innovations in Computer Vision

Computer Vision (CV) continues to be a field of rapid transformation, largely propelled by advancements in deep learning. While Convolutional Neural Networks (CNNs) laid much of the groundwork, Vision Transformers (ViTs) have emerged as a prominent architectural paradigm. Generative models are producing visual content of increasing realism and controllability.

Self-supervised learning techniques are reducing the dependency on vast labeled datasets. A newer frontier, physics-aware AI, aims to imbue vision models with an intrinsic understanding of real-world physical laws and dynamics.

A. Vision Transformers (ViTs) and Architectural Hybrids

Vision Transformers (ViTs) have demonstrated remarkable performance but also come with challenges related to computational cost and data requirements. Recent research focuses on refining ViTs through hybridization with CNNs and enhancing their adaptability.

- **Full Title:** ECViT: Efficient Convolutional Vision Transformer with Local-Attention and Multi-scale Stages
- **Author(s):** Zhoujie Qian
- **Year:** 2025 (arXiv:2504.14825)
- **Key Contribution:** ECViT is a hybrid architecture that synergistically combines the strengths of CNNs and Transformers. It aims to mitigate common ViT limitations by introducing CNN-derived inductive biases, such as locality and translation invariance, into the Transformer framework. This is achieved by extracting image patches from low-level features generated by convolutional layers and enhancing the transformer encoder with convolutional operations. Furthermore, ECViT incorporates local attention mechanisms and a pyramid structure to facilitate efficient multi-scale feature extraction and representation. The design prioritizes an optimal balance between model performance and computational efficiency, reportedly outperforming state-of-the-art models on various image classification tasks while maintaining low computational and storage demands.⁴³
- **Direct Link:** <https://arxiv.org/abs/2504.14825>
- **Full Title:** UniViTAR: Unified Vision Transformer with Native Resolution
- **Author(s):** Limeng Qiao, et al.
- **Year:** 2025 (arXiv:2504.01792v2, Apr 2025)
- **Key Contribution:** UniViTAR introduces a family of Vision Transformer-based foundation models specifically designed to uniformly process visual modalities (images or videos) at their native resolutions and dynamic aspect ratios, a departure from the fixed-resolution inputs typically required by ViTs. The framework incorporates several architectural upgrades inspired by recent advancements in LLMs, such as 2D Rotary Position Embedding (RoPE), SwiGLU activation functions, RMSNorm, and Query-Key Normalization. Complementing these architectural changes is a progressive training paradigm that includes resolution curriculum learning and a hybrid loss function combining sigmoid-based contrastive loss with feature distillation from a frozen teacher model. This approach aims to enhance adaptability and scalability, with models ranging from 0.3B to 1B parameters.⁴⁶
- **Direct Link:** <https://arxiv.org/abs/2504.01792>

The initial ViT concept, while groundbreaking, faced practical limitations such as data inefficiency due to a lack of inherent inductive biases (compared to CNNs) and high computational costs stemming from the quadratic complexity of self-attention.⁴³ Current

research trends show a maturation of ViTs through two primary strategies: hybridization and adaptability enhancements. ECViT⁴³ exemplifies hybridization by re-integrating convolutional operations and principles like locality and pyramid structures directly into the Transformer architecture. This creates a hybrid model that is more efficient and often performs better, especially when training data is limited. UniViTAR⁴⁶, on the other hand, addresses another practical constraint: the typical requirement for ViTs to process fixed-size, often square, input patches, which is incongruent with the variability of real-world image and video data. UniViTAR focuses on architectural and training modifications (e.g., 2D RoPE, resolution curriculum learning) to enable ViTs to handle native resolutions effectively, thereby making them more versatile. Both these approaches signify a move away from "pure" ViT architectures towards more practical, robust, and efficient variants that draw lessons from both the successes of CNNs and recent innovations in the LLM space (such as the adoption of SwiGLU activation in UniViTAR).

This evolutionary path mirrors that of many successful technological paradigms: an initial breakthrough is followed by periods of refinement, hybridization with established technologies, and adaptation to overcome initial limitations and broaden practical applicability. Future research in ViTs is likely to continue along these lines, with the goal of developing models that are not only powerful in their representational capacity but also highly practical and efficient for diverse real-world vision tasks.

B. Advanced Generative Models (Images, Video, 3D)

Generative AI is rapidly advancing, enabling the creation of highly realistic images, videos, and 3D/4D content. A significant emerging trend is the shift from focusing solely on visual fidelity to incorporating physical plausibility and controllability, aiming to create "world simulators".¹¹

- **Full Title:** Generative Physical AI in Vision: A Survey
- **Author(s):** Daochang Liu, Junyu Zhang, Anh-Dung Dinh, Eunbyung Park, Shichao Zhang, Ajmal Mian, Mubarak Shah, Chang Xu
- **Year:** 2025 (arXiv:2501.10928v2, Apr 2025)
- **Key Contribution:** This survey provides a systematic review of physics-aware generative models in computer vision. These models aim to generate content that is not only visually realistic but also adheres to real-world physical laws, a critical limitation of conventional generative models that primarily optimize for visual fidelity. The paper categorizes methods based on how they incorporate physical knowledge—either through explicit physical simulation or via implicit learning from data. It also analyzes key paradigms in the field, discusses evaluation protocols, and identifies future research directions toward developing generative models that can function as comprehensive "world simulators".¹¹
- **Direct Link:** <https://arxiv.org/abs/2501.10928> (A summary of reviewed papers is available at: <https://tinyurl.com/Physics-Aware-Generation>)

The development of sophisticated video generation models like Sora, Kling, and Veo 2⁵⁰ necessitates more nuanced evaluation benchmarks that go beyond simple visual quality assessments.

- **Full Title:** VBench-2.0: Advancing Video Generation Benchmark Suite for Intrinsic Faithfulness
- **Author(s):** (A large collaborative effort, full author list in the paper)
- **Year:** 2025 (arXiv:2503.21755v1, Mar 2025)
- **Key Contribution:** VBench-2.0 is a new benchmark suite designed to evaluate advanced video generation models by focusing on "intrinsic faithfulness" rather than just superficial visual quality. It introduces five key dimensions for evaluation: Human Fidelity (anatomical correctness, temporal consistency), Controllability (adherence to complex prompts, dynamic changes), Creativity (diversity, novel compositions), Physics (adherence to physical laws), and Commonsense (motion rationality, instance preservation). The benchmark aims to drive progress towards models that can simulate and reason about the real world, highlighting current limitations of even state-of-the-art models in areas like complex plot generation and robust commonsense reasoning.⁵⁰
- **Direct Link:** <https://arxiv.org/abs/2503.21755>

In the realm of 3D content generation, controllability remains a key challenge.

- **Full Title:** Controllable 3D Outdoor Scene Generation via Scene Graphs
- **Author(s):** Yiran Xing, Zhaoxiang Cai, et al.
- **Year:** 2025 (arXiv:2503.07152v1, Mar 2025)
- **Key Contribution:** This work proposes a novel method for controllable generation of outdoor 3D scenes by leveraging scene graphs as an intuitive and user-friendly control format. The system features an interactive component that allows users to transform a sparse scene graph into a dense Bird's Eye View (BEV) Embedding Map. This map then guides a conditional diffusion model to generate 3D scenes that accurately match the structural and semantic descriptions provided by the scene graph.⁵²
- **Direct Link:** <https://arxiv.org/abs/2503.07152>

The GitHub repository "Awesome-Indoor-Scene-Synthesis"⁵⁵ curates numerous recent papers (many from 2024 and slated for 2025 conferences) on 3D scene generation, such as "SceneFactor: Factored Latent 3D Diffusion for Controllable 3D Scene Generation" and "MIDI: Multi-Instance Diffusion for Single Image to 3D Scene Generation." These works often utilize diffusion models and focus on enhancing controllability and compositional capabilities in 3D synthesis. Academic blogs frequently serve to synthesize and explain such cutting-edge research to a broader audience.⁵⁶ For example, a blog post discussing these papers might be titled: "*SceneFactor and MIDI: New Frontiers in Controllable and Compositional 3D Scene Generation with Diffusion Models*", authored by researchers in the field or expert commentators, and published in 2024. Such a post would elucidate how these methods leverage diffusion models for more nuanced and controllable 3D scene creation, potentially by decomposing complex generation tasks into more manageable components.

A significant shift is occurring in the field of generative AI: the focus is expanding from achieving mere "visual fidelity" to striving for "world fidelity." Early generative models, like Generative Adversarial Networks (GANs), concentrated heavily on producing photorealistic images. Current leading models, including advanced diffusion models and large-scale video

generation systems like Sora, have achieved remarkable success in visual realism.¹¹ However, a crucial limitation remains their often-superficial understanding of underlying physical principles or commonsense logic.¹¹ The survey on "Generative Physical AI" ¹¹ explicitly identifies this gap and reviews emerging methods that aim to create "world simulators" capable of generating content consistent with physical laws. Similarly, the VBench-2.0 benchmark ⁵⁰ introduces evaluation dimensions such as "Physics" and "Commonsense," thereby pushing the assessment of generative models beyond surface-level appearance. Research in controllable 3D scene generation ⁵² further underscores this trend by emphasizing not just the generation of 3D environments but also the provision of meaningful, intuitive control over the generated world's content and structure. This evolution towards "world fidelity" is paramount for applications that require interaction with the physical world, such as robotics and autonomous systems, or for creating complex, high-fidelity simulations. "World models" that possess a deeper understanding of physics, causality, and object interactions could unlock entirely new AI capabilities. However, they also necessitate far more sophisticated evaluation methodologies and present new safety challenges if their understanding of the world is flawed or incomplete.

C. Self-Supervised Learning (SSL) in Vision

Self-Supervised Learning (SSL) has emerged as a powerful paradigm for learning rich visual representations from unlabeled data, thereby reducing the reliance on costly and time-consuming manual annotation. Current research is focused on improving the efficiency and adaptability of SSL methods.

- **Full Title:** Escaping the Big Data Paradigm in Self-Supervised Representation Learning from Images
- **Author(s):** Carlos Velez-García, Miguel Cazorla, Jorge Pomares
- **Year:** 2025 (arXiv:2502.18056)
- **Key Contribution:** This paper introduces SCOTT (Sparse Convolutional Tokenizer for Transformers), a shallow tokenization architecture, and proposes its use within a Masked Image Modeling framework based on a Joint-Embedding Predictive Architecture (MIM-JEPA). This combined approach is designed to enable Vision Transformers (ViTs) to be trained effectively via SSL on significantly smaller datasets than traditionally required, and notably, without the need for pretraining on massive external datasets. The work aims to make SSL more accessible and practical, particularly for fine-grained visual tasks.⁵⁸
- **Direct Link:** <https://github.com/inescopresearch/scott> (Code), <https://arxiv.org/abs/2502.18056> (Paper)
- **Full Title:** UpStep: Unsupervised Parameter-efficient Source-free Post-pretraining
- **Author(s):** Ulas Gul, Oguz Kaan Yüksel, et al.
- **Year:** 2025 (Implied by context, paper references work from 2024/2025)
- **Key Contribution:** UpStep is an unsupervised, parameter-efficient, and source-free post-pretraining method designed to adapt a pre-existing pretrained visual model to an unlabeled target domain. It employs an SSL training scheme along with a novel "center

vector regularization" (CVR) technique. CVR aims to minimize catastrophic forgetting of the original pretraining and also reduces computational costs by, for instance, skipping backpropagation in a percentage of training iterations.⁶⁰

- **Direct Link:** <https://www.arxiv.org/pdf/2502.21313> (Paper PDF)

SSL is undergoing a pivot towards greater efficiency and adaptability. While the initial successes of SSL often depended on the availability of large models (like ViTs) and vast quantities of unlabeled data⁵⁸, this "resource-hungry" paradigm inherently limits its accessibility.⁵⁸ The SCOTT/MIM-JEPA framework⁵⁸ directly confronts this issue by developing methods (such as convolutional tokenizers and the JEPA learning strategy) that empower ViTs to learn meaningful representations from substantially smaller datasets, even when trained from scratch. UpStep⁶⁰ addresses efficiency from a different angle: the adaptation of *existing* large pretrained models to new target domains using SSL. Crucially, UpStep operates in a parameter-efficient and *source-free* manner, which is vital in scenarios where the original pretraining data or even the specifics of the source model are unavailable due to privacy or proprietary reasons. Both these research directions aim to make the potent capabilities of SSL more practical and broadly applicable under a wider range of real-world constraints. These advancements are poised to democratize SSL, enabling a larger community of researchers and organizations, including those with limited computational resources or access to massive datasets, to leverage its benefits. Furthermore, source-free adaptation techniques like UpStep offer significant advantages in terms of privacy and data governance. The field is clearly maturing from the fundamental question of "Can we perform SSL?" to the more nuanced and practical question of "How can we perform SSL efficiently, adaptively, and responsibly?"

The landscape of contemporary computer vision research is significantly shaped by the advancements and synergies between three key pillars: Vision Transformers (and their hybrid derivatives), sophisticated Generative Models (particularly those based on diffusion processes), and Self-Supervised Learning techniques. Vision Transformers provide powerful and flexible architectural backbones, which are continually being refined for improved efficiency and adaptability to diverse visual inputs.⁴³ Generative Models are pushing the frontiers of visual content creation across images, video, and 3D domains, with an increasing emphasis on physical plausibility, commonsense consistency, and user controllability.¹¹ Many of these advanced generative models themselves leverage Transformer architectures in their core design. Self-Supervised Learning offers a compelling pathway to train these large and complex models with reduced reliance on extensively labeled datasets, and current SSL research is increasingly focused on enhancing efficiency and adaptability for various downstream tasks and data regimes.⁵⁸ These three areas are not evolving in isolation but are becoming increasingly interconnected. For instance, SSL can be employed to pretrain powerful ViT backbones, which can subsequently serve as the foundation for high-capacity generative models. Conversely, generative models might also be used to synthesize diverse data augmentations for SSL. Future progress in computer vision will likely emerge from the continued co-evolution and deeper integration of these three technological pillars. Innovations that effectively bridge these areas—such as SSL techniques specifically optimized

for training generative models, or novel ViT architectures tailored for generative tasks—are expected to be particularly impactful. The overarching objective appears to be the creation of models that can understand, generate, and learn about the visual world with ever-increasing fidelity, efficiency, autonomy, and alignment with real-world principles.

IV. Developments in Reinforcement Learning

Reinforcement Learning (RL) continues to advance its capacity to solve complex sequential decision-making problems across a diverse array of domains. Key areas of active research and significant progress include improving the sample efficiency of learning algorithms, developing more robust and intelligent exploration strategies, enhancing the explainability and interpretability of RL agents (Explainable RL or XRL), making offline RL (learning from fixed datasets) more practical and effective, scaling multi-agent RL (MARL) to handle more complex interactions, and refining Reinforcement Learning from Human Feedback (RLHF) for better alignment of AI systems with human preferences and values.

A. Explainable Deep Reinforcement Learning (XRL)

As Deep Reinforcement Learning (DRL) models, often relying on opaque neural network architectures, are deployed in increasingly complex and high-stakes applications, the need for transparency and interpretability becomes paramount. XRL aims to address these challenges.

- **Full Title:** A Survey on Explainable Deep Reinforcement Learning
- **Author(s):** Zelei Cheng, Jiahao Yu, Xinyu Xing
- **Year:** 2025 (arXiv:2502.06869v1, Feb 2025)
- **Key Contribution:** This survey provides a comprehensive review of the field of XRL. It categorizes existing explanation techniques into feature-level, state-level, dataset-level, and model-level approaches. The paper also evaluates their qualitative and quantitative assessment frameworks and explores the role of XRL in policy refinement, enhancing adversarial robustness, and improving security. A notable aspect of this survey is its examination of the integration of RL with Large Language Models (LLMs), particularly through Reinforcement Learning from Human Feedback (RLHF), and the implications this synergy has for the explainability of both RL agents and the LLMs they help align.⁶¹
- **Direct Link:** <https://arxiv.org/abs/2502.06869>

The imperative for explainability in DRL is intensifying as these systems tackle more intricate tasks and are deployed in sensitive domains such as healthcare, autonomous driving, and finance, where understanding the rationale behind an agent's decisions is crucial for trust, compliance, and debugging.⁶¹ The opacity of DRL agents, typically deep neural networks, can lead to unreliable decision-making if not properly understood and validated.⁶¹ The survey⁶¹ explicitly highlights that the "increasing integration of RL with LLMs" via methods like RLHF "further amplifies the explainability challenge." This is because RLHF uses RL to fine-tune already complex LLMs, making it difficult to discern how RL updates interact with the LLM's vast parameter space and neural representations to produce specific behaviors or align with human preferences. Thus, understanding *why* an RL-tuned LLM behaves in a particular way

becomes even more critical and challenging.

XRL is therefore not merely an academic pursuit but a practical necessity for the responsible development and deployment of advanced AI systems. Future DRL frameworks, especially those involving interactions with or modifications of LLMs, may need to incorporate explainability features by design, rather than treating interpretability as an optional add-on. This could drive further research into inherently more interpretable DRL algorithms or the development of more powerful post-hoc explanation techniques specifically tailored for the complex dynamics of RL-LLM interactions.

B. Advancements in Offline RL and Multi-Agent RL (MARL)

Offline RL, which focuses on learning policies from pre-existing static datasets without further environment interaction, and MARL, which deals with multiple interacting agents, are both highly active research areas. Recent work aims to make offline RL more robust and to tackle the complexities of MARL.

- **Full Title:** Video-Enhanced Offline RL (VeoRL): A Model-Based Approach
- **Author(s):** Yecheng Moon, et al.
- **Year:** 2025 (arXiv:2505.06482v2, May 2025)
- **Key Contribution:** VeoRL is a model-based offline RL method that addresses the limitations of learning from static datasets by constructing an interactive world model. Uniquely, this world model is enhanced using diverse, unlabeled video data that is readily available online. The system employs a hierarchical world model with two state transition branches: one predicting future state evolution based on the agent's real actions (from the offline dataset) and another predicting longer-term environmental feedback derived from latent behaviors extracted from the auxiliary video data. VeoRL has demonstrated significant improvements over existing offline RL methods across various visual control benchmarks, including robotic manipulation, autonomous driving, and open-world video games.⁶⁴
- **Direct Link:** <https://arxiv.org/abs/2505.06482>
- **Full Title:** Variational OOD State Correction for Offline Reinforcement Learning
- **Author(s):** Zicheng Zhang, et al.
- **Year:** 2025 (arXiv:2505.00503v1, May 2025)
- **Key Contribution:** This paper introduces Density-Aware Safety Perception (DASP), a novel method for out-of-distribution (OOD) state correction in offline RL. The core idea of DASP is to encourage the learning agent to prioritize actions that are likely to lead to outcomes with higher data density, thereby promoting its operation within, or its return to, in-distribution (and thus safer) regions of the state space. This is achieved by optimizing an objective within a variational framework that concurrently considers both the potential outcomes of decision-making and their associated density in the offline dataset.⁶⁶
- **Direct Link:** <https://arxiv.org/abs/2505.00503>
- **Full Title:** Addressing Rotational Learning Dynamics in Multi-Agent Reinforcement Learning

- **Author(s):** Di-An Jan, et al.
- **Year:** 2025 (arXiv:2410.07976v2, Feb 2025)
- **Key Contribution:** This work tackles the issue of instability and poor convergence in MARL, which often arises from "rotational optimization dynamics" caused by the competing objectives of interacting agents. The authors reframe MARL problems using the mathematical framework of Variational Inequalities (VIs), which is better suited for equilibrium-finding problems than standard gradient-based optimization. They propose general approaches, LA-MARL (Lookahead-MARL) and EG-MARL (Extragradient-MARL), for integrating gradient-based VI methods into existing MARL algorithms, demonstrating significant performance improvements and enhanced convergence across benchmarks.⁶⁷
- **Direct Link:** <https://arxiv.org/abs/2410.07976>
- **Full Title:** Offline Multi-agent Reinforcement Learning via Score Decomposition
- **Author(s):** Lingheng Meng, et al.
- **Year:** 2025 (arXiv:2505.05968v1, May 2025)
- **Key Contribution:** This paper addresses the compounded challenges of offline MARL, including distributional shifts, the high dimensionality of joint action spaces, and the diversity in agent coordination strategies present in offline data. The proposed novel two-stage framework first employs a diffusion-based generative model to explicitly capture the complex joint behavior policy from the static dataset, enabling accurate modeling of diverse multi-agent coordination patterns. Second, it introduces a sequential score function decomposition mechanism to regularize individual agent policies and facilitate decentralized execution, achieving state-of-the-art performance on several offline MARL benchmarks.⁶⁹
- **Direct Link:** <https://arxiv.org/abs/2505.05968>

The open-source library CORL (Clean Offline Reinforcement Learning)⁷³, presented at NeurIPS 2023, also contributes significantly by providing thoroughly benchmarked single-file implementations of deep offline and offline-to-online RL algorithms. Such tools are vital for improving reproducibility and accessibility in this complex research area. An academic blog post discussing CORL might be titled: "*CORL: Simplifying and Standardizing Offline RL Research with Benchmarked Single-File Implementations*", authored by researchers involved or RL bloggers, and published around 2023/2024. It would emphasize CORL's role in making performance-relevant details easier to recognize and its integration of experiment tracking. A key trend in making offline RL more practical and powerful is the innovative use of generative models and auxiliary unlabeled data sources. The primary challenge in offline RL is learning effectively from a fixed, potentially suboptimal or incomplete dataset, which often leads to problems with distributional shift.⁶⁴ VeoRL⁶⁴ directly confronts this by incorporating unlabeled video data to construct a richer world model. This effectively expands the "experience" available to the agent beyond the confines of the static dataset. Similarly, the work on Offline MARL via Score Decomposition⁶⁹ utilizes a diffusion-based generative model to explicitly capture the complex behavior policy evidenced in the offline data. This approach aids in modeling diverse multi-agent coordination patterns. These methods signify a move

beyond purely algorithmic solutions for handling fixed datasets, towards actively enriching the learning process with external data or through powerful generative modeling of the available data. This suggests that the distinction between offline and online RL might become increasingly blurred. Future "offline" methodologies could routinely integrate vast amounts of unlabeled environmental data (e.g., internet videos for robotics training) or employ generative models to create diverse "simulated" experiences, thereby making the learned policies more robust and sample-efficient when eventually deployed or fine-tuned in an online setting. Concurrently, MARL research is maturing by addressing fundamental issues of stability and the complexity of coordination. MARL training is notoriously unstable and highly sensitive to initialization and hyperparameters, partly due to what is described as "rotational dynamics" arising from the interactions of competing or collaborating agents.⁶⁷ The paper "Addressing Rotational Learning Dynamics..."⁶⁷ proposes a foundational shift by using Variational Inequalities (VIs) as a more suitable mathematical framework than standard optimization for these scenarios, leading to more stable algorithms like LA-MARL. In the *offline* MARL context, these challenges are amplified by high-dimensional joint action spaces and the often unobserved diversity of coordination strategies within the static dataset.⁶⁹ "Offline MARL via Score Decomposition"⁶⁹ leverages diffusion models to capture these intricate multi-agent behavior policies and employs score function decomposition to enable effective decentralized execution. This focus on tackling both the modeling and execution challenges in offline MARL, using sophisticated mathematical tools like VIs and powerful generative techniques like diffusion models, indicates a significant maturation of the field. As MARL systems are envisioned for increasingly complex real-world applications, such as the coordination of autonomous vehicles or robotic teams, resolving these core issues of stability, reproducibility, and effective coordination from potentially limited data is paramount.

C. Improvements in RL from Human Feedback (RLHF) and Alignment

RLHF has become a cornerstone technique for aligning LLMs and other AI systems with human preferences. Current research is focused on making RLHF more personalized, efficient, and robust.

- **Full Title:** A Shared Low-Rank Adaptation Approach to Personalized RLHF
- **Author(s):** Renpu Liu, Peng Wang, Donghao Li, Cong Shen, Jing Yang
- **Year:** 2025 (AISTATS 2025; arXiv:2503.19201, Mar 2025)
- **Key Contribution:** This paper introduces P-ShareLoRA (Personalized LoRA with Shared Component) for personalized RLHF. It addresses the limitation of standard RLHF approaches that typically assume homogeneous human preferences by developing a method to efficiently learn personalized reward functions. P-ShareLoRA leverages the shared components of Low-Rank Adaptation (LoRA) modules to achieve this personalization, aiming to improve the adaptability of RLHF-tuned models to individual user needs and preferences.⁷⁴
- **Direct Link:** <https://arxiv.org/abs/2503.19201>
- **Full Title:** MA-RLHF: Reinforcement Learning from Human Feedback with Macro Actions
- **Author(s):** Yekun Chai, Haoran Sun, Huang Fang, Shuohuan Wang, Yu Sun, Hua Wu

- **Year:** 2025 (ICLR 2025; arXiv:2410.02743, Oct 2024/Feb 2025)
- **Key Contribution:** MA-RLHF (Macro-Action RLHF) aims to improve the efficiency and effectiveness of token-level RLHF, particularly for aligning LLMs that generate long sequences of text. It achieves this by incorporating "macro-actions"—which can be sequences of tokens or higher-level linguistic constructs—into the learning process. This approach helps to address the credit assignment problem, where it is difficult to determine which specific early actions contributed to a delayed reward, thereby enhancing learning efficiency and speeding up convergence.⁷⁷
- **Direct Link:** <https://github.com/ernie-research/MA-RLHF> (Code), <https://arxiv.org/abs/2410.02743> (Paper)

RLHF is evolving towards greater personalization and improved learning efficiency. Standard RLHF methodologies, widely used for fine-tuning models like ChatGPT, typically train a single reward model based on aggregated human preferences.⁷⁴ However, human preferences are inherently diverse and heterogeneous⁷⁴; a single, monolithic reward model may not align optimally with all individual users. P-ShareLoRA⁷⁴ directly confronts this "homogeneous preference" limitation by proposing an efficient method for personalized RLHF, enabling the learning of individual-specific reward functions through the clever use of shared LoRA components. This represents a significant step towards more nuanced alignment. Separately, a common challenge in applying RLHF at the token level, especially for LLMs generating lengthy text, is the credit assignment problem: delayed rewards make it difficult for the model to discern which specific actions (tokens generated) early in a sequence were responsible for a preferred outcome.⁷⁷ MA-RLHF⁷⁷ tackles this by introducing the concept of macro-actions. Operating at a higher level of temporal abstraction simplifies the decision-making process and improves credit attribution, leading to faster and more stable convergence.

For AI systems to be truly aligned and genuinely helpful, they must possess the capacity to adapt to individual user needs, values, and contexts. Personalized RLHF is a key enabler in this direction. Concurrently, enhancing the core learning mechanisms of RLHF, such as improving credit assignment, will make the alignment process more scalable and robust. This will allow for more complex and nuanced behaviors to be learned and aligned effectively.

The major advancements in RL appear to be converging at the intersection of three critical frontiers: enhancing agent *explainability* (XRL), improving the ability to *leverage diverse data sources* more effectively (as seen in offline RL innovations using video and generative models), and ensuring that agent behaviors robustly *align with human expectations and preferences* (evidenced by the evolution of RLHF towards personalization and greater efficiency). As RL models grow in power and are deployed in increasingly impactful real-world scenarios, the demand for transparency (XRL⁶¹) becomes essential for building trust, enabling effective debugging, and ensuring accountability. Simultaneously, the inherent limitations of online data collection—such as cost, safety concerns, and sample inefficiency—are driving offline RL to innovate. This includes incorporating vast amounts of unlabeled data like internet videos (e.g., VeoRL⁶⁴) or employing sophisticated generative models to better understand behavioral policies from static datasets (e.g., Offline MARL with diffusion models⁶⁹). This signifies a broader shift towards more data-centric RL paradigms. RLHF stands as a prime

example of direct human-AI interaction for alignment, and it is currently evolving to accommodate the diversity of human preferences (e.g., P-ShareLoRA ⁷⁴) and to refine its underlying learning mechanics for better efficiency (e.g., MA-RLHF ⁷⁷). These trends are not isolated; for instance, effectively explaining the decisions of an LLM fine-tuned via RLHF (a challenge highlighted in the XRL survey ⁶¹) requires a deep understanding of both the RL optimization process and the LLM's internal behavior. Similarly, offline RL agents trained from diverse and potentially noisy data sources might also require XRL techniques to understand what they have learned and why they make certain decisions. The future of RL seems to be moving towards more holistic systems that are not only highly capable but also interpretable, data-aware, and human-compatible. This necessitates interdisciplinary approaches that draw insights and methodologies from XAI, data science, and human-computer interaction, pushing RL beyond purely algorithmic development into a more integrated and application-focused discipline.

V. Progress in General ML Theory and Optimization

This section delves into broader theoretical underpinnings of machine learning, including novel optimization algorithms essential for training large and complex models, the evolving landscape of federated learning with its focus on privacy, robustness, and efficiency, and deeper investigations into AI alignment and the fundamental principles of generalization in deep learning.

A. Novel Optimization Algorithms for Deep Learning

The development of efficient and effective optimization algorithms is crucial for training increasingly large and complex deep learning models. Recent research explores gradient-free methods and techniques to enhance generalization.

- **Full Title:** ZeroFlow: Overcoming Catastrophic Forgetting is Easier than You Think
- **Author(s):** Tao Feng, Wei Li, Didi Zhu, Hangjie Yuan, Wendi Zheng, Dan Zhang, Jie Tang
- **Year:** 2025 (arXiv:2501.01045v3, Jan 2025, Submitted to ICML)
- **Key Contribution:** ZeroFlow introduces the first benchmark specifically designed to evaluate gradient-free (zeroth-order) optimization algorithms in the context of continual learning, with a focus on their ability to overcome catastrophic forgetting. The study examines a suite of forward pass-based methods across various forgetting scenarios, model types, and datasets. A key finding is that forward passes alone can be sufficient to mitigate forgetting. The research reveals new optimization principles highlighting the potential of forward-pass methods in managing task conflicts and reducing memory demands, even with just a single forward pass in some cases.⁷⁹
- **Direct Link:** <https://arxiv.org/abs/2501.01045>
- **Full Title:** LORENZA: Enhancing Generalization in Low-Rank Gradient LLM Training and Fine-Tuning via Efficient Zeroth-Order Adaptive SAM Optimization
- **Author(s):** Yehonathan Refael, Iftach Arbel, Ofir Lindenbaum, Tom Tirer
- **Year:** 2025 (arXiv:2502.19571v1, Feb 2025, Submitted to ICML)
- **Key Contribution:** LORENZA proposes AdaZo-SAM, a novel framework that combines

the Adam optimizer with Sharpness-Aware Minimization (SAM) but requires only a single gradient computation per iteration. This efficiency is achieved by using a stochastic zeroth-order estimation to find SAM's ascent perturbation. Additionally, the paper introduces LORENZA itself as a memory-efficient version of AdaZo-SAM that employs adaptive low-rank gradient updates. This approach is designed to improve generalization ability for Large Language Model (LLM) parameter-efficient fine-tuning (PEFT) and pre-training.⁸²

- **Direct Link:** <https://arxiv.org/abs/2502.19571>

There is a growing interest and demonstrated success in employing zeroth-order (ZO) optimization methods for training and fine-tuning deep learning models, particularly in complex scenarios such as continual learning and the optimization of LLMs. ZO methods, which rely solely on function evaluations (forward passes) rather than explicit gradient computations, offer compelling advantages. Traditional deep learning heavily depends on first-order (gradient-based) optimizers like SGD and Adam. However, gradients can be inaccessible (e.g., when interacting with black-box model APIs), computationally prohibitive for extremely large models, or present unique challenges in specific settings like continual learning, where they can contribute to catastrophic forgetting.⁷⁹ The ZeroFlow benchmark and associated research⁷⁹ demonstrate that ZO methods can effectively mitigate catastrophic forgetting in continual learning scenarios, sometimes achieving this with remarkably few forward passes, thereby challenging the conventional wisdom that backpropagation is indispensable for this problem. Furthermore, the LORENZA framework⁸² ingeniously incorporates ZO estimation into Sharpness-Aware Minimization (SAM)—a technique known for finding flatter minima that generalize better—to render SAM more computationally efficient (requiring only a single gradient computation per iteration) for LLM fine-tuning, especially when combined with low-rank adaptation strategies. This body of work indicates that ZO methods are transitioning from theoretical concepts to engineered, practical solutions that can enhance efficiency and enable optimization in situations where traditional gradient-based approaches are problematic or infeasible.

The maturation of ZO optimization techniques could significantly lower the barriers to training and fine-tuning large-scale models by reducing computational demands (e.g., fewer or no backward passes) and by opening up optimization possibilities in new contexts, such as when interacting with proprietary black-box APIs. This could lead to more efficient continual learning systems and more accessible pathways for customizing and deploying LLMs.

B. Generalization Theory in Deep Learning

Understanding why deep learning models generalize well from training data to unseen data, especially when they are highly overparameterized, remains a central question in machine learning theory.

- **Full Title:** Deep Multi-Task Learning Has Low Amortized Intrinsic Dimensionality
- **Author(s):** Hossein Zakerinia, et al.
- **Year:** 2025 (arXiv:2501.19067v1, Jan 2025, Submitted to ICML)
- **Key Contribution:** This research confirms the phenomenon that deep multi-task

learning (MTL) models, despite their high ambient dimensionality (number of parameters), effectively learn within a low intrinsic dimensional subspace. The authors introduce a method to parameterize MTL networks directly in this low-dimensional space using random expansion techniques. They demonstrate that this low-dimensional representation, when combined with weight compression and PAC-Bayesian reasoning, leads to the first non-vacuous generalization bounds for deep MTL networks, providing concrete theoretical guarantees.⁸⁴

- **Direct Link:** <https://arxiv.org/abs/2501.19067>
- **Full Title:** Survey on Generalization Theory for Graph Neural Networks
- **Author(s):** Kajetan Schweighofer, et al.
- **Year:** 2025 (arXiv:2503.15650v1, Mar 2025)
- **Key Contribution:** This survey systematically reviews the existing literature on the generalization abilities of Message-Passing Neural Networks (MPNNs), a prominent class of Graph Neural Networks (GNNs). It analyzes the strengths and limitations of various theoretical studies that employ tools such as VC dimension, Rademacher complexity, stability analysis, PAC-Bayesian bounds, covering numbers, and other frameworks to understand and quantify GNN generalization.⁸⁷
- **Direct Link:** <https://arxiv.org/pdf/2503.15650>

The pursuit of meaningful (i.e., non-vacuous) generalization bounds for deep learning models is driving researchers to move beyond classical complexity measures. It necessitates incorporating deeper insights about the learned structure of these models (such as low intrinsic dimensionality and compressibility) and the dynamics of the learning process itself (often analyzed through frameworks like PAC-Bayesian theory). Deep learning models often generalize surprisingly well even when they are massively overparameterized, a behavior that seems to defy classical learning theory, which typically links generalization to a favorable trade-off between model complexity and the amount of training data.⁸⁴ Traditional generalization bounds based on measures like VC dimension or Rademacher complexity frequently prove to be vacuous (i.e., too loose to be informative) for modern deep networks.⁸⁴ The work on MTL by Zakerinia et al. ⁸⁴ explicitly connects the concepts of low intrinsic dimensionality and network compressibility with PAC-Bayesian reasoning to derive *non-vacuous* generalization bounds for deep MTL systems. This is a significant achievement as it means the derived bounds can provide actual numerical guarantees on generalization performance, rather than just conceptual insights. Similarly, the survey on GNN generalization ⁸⁷ reviews a wide array of theoretical tools, including PAC-Bayes and stability analysis, indicating a broad and active search for effective measures of generalization in these structured data domains. This collective body of work signals a trend towards developing theories that account for *what is actually learned* by the model and *how that learned knowledge is represented internally*, rather than relying solely on worst-case capacity measures of the model architecture.

A more robust theoretical understanding of generalization can provide invaluable guidance for designing more reliable, data-efficient, and robust deep learning models. The achievement of non-vacuous bounds, even if initially for specific settings like MTL or GNNs, represents crucial

progress towards this overarching goal. Such theoretical advances can help explain why certain architectural choices, training strategies, or data characteristics lead to better generalization, ultimately fostering more principled and effective model development.

C. Federated Learning: Privacy, Robustness, Communication Efficiency

Federated Learning (FL) enables collaborative model training across decentralized data sources while preserving data privacy. Research in FL is rapidly evolving to address challenges related to the scale of modern models (like LLMs), data heterogeneity (non-IID data), robustness against adversarial attacks, and communication efficiency.

- **Full Title:** Ferret: Federated Full-Parameter Tuning at Scale for Large Language Models
- **Author(s):** Yao Shu, Wenyang Hu, See-Kiong Ng, Bryan Kian Hsiang Low, Fei Richard Yu
- **Year:** 2024 (arXiv:2409.06277v2, Sep 2024)
- **Key Contribution:** Ferret is presented as the first first-order FL method that utilizes shared randomness to facilitate scalable full-parameter tuning of LLMs. It addresses the challenges of high communication overhead and computational cost by employing efficient first-order methods for local updates on client devices, projecting these updates into a low-dimensional space to significantly reduce communication burdens, and then reconstructing the local updates from this compressed space using shared randomness for effective global aggregation. This approach aims to strike a balance between maintaining model accuracy, achieving high computational efficiency, reducing communication overhead, and ensuring fast convergence.⁸⁸
- **Direct Link:** <https://github.com/allen4747/Ferret> (Code), <https://arxiv.org/abs/2409.06277> (Paper)
- **Full Title:** FedDDL: Federated Deconfounding and Debiasing Learning for Out-of-Distribution Generalization
- **Author(s):** Qi Zhuang, Ming-Chang Lee, Han Yu, et al.
- **Year:** 2025 (Implied by recent publication activity and references within the paper from Han Yu's group ⁹²)
- **Key Contribution:** FedDDL addresses the problem of attribute bias in FL, which can cause local models to learn spurious correlations and optimize inconsistently, leading to degraded out-of-distribution (OOD) generalization. The proposed method includes a Disentangled Effect Calibration (DEC) module designed to decouple background features from object features. By generating counterfactual samples locally, FedDDL aims to establish a more robust association between features and classes, thereby improving the model's ability to focus on causal features and generalize better to unseen data distributions. The authors report that FedDDL significantly enhances model capability to focus on main objects in unseen data, outperforming several state-of-the-art existing methods.⁹³
- **Direct Link:** <https://arxiv.org/abs/2505.04979> (Paper)
- **Full Title:** FLTG: Byzantine-Robust Federated Learning via Angle-Based Defense and Non-IID-Aware Weighting

- **Author(s):** Jiahao Wang, et al.
- **Year:** 2025 (arXiv:2505.12851v1, May 2025)
- **Key Contribution:** FLTG is a novel aggregation algorithm for FL designed to be robust against Byzantine attacks, where malicious clients attempt to corrupt the global model by sending manipulated updates. The method integrates an angle-based defense mechanism, which uses cosine similarity calculations against a small, clean dataset held by the server to filter out misaligned or malicious updates. It also incorporates dynamic reference selection and non-IID-aware weighting to improve performance, particularly when client data is not independent and identically distributed (non-IID), a common scenario where traditional Byzantine-robust metrics can struggle.⁹⁵
- **Direct Link:** <https://arxiv.org/abs/2505.12851>
- **Full Title:** Communication-efficient Vertical Federated Learning via Compressed Error Feedback
- **Author(s):** Pedro Valdeira, João Xavier, Cláudia Soares, Yuejie Chi
- **Year:** 2024 (arXiv:2406.14420v3, accepted to IEEE Transactions on Signal Processing)
- **Key Contribution:** This paper introduces EF-VFL, an error feedback compressed vertical federated learning method specifically for training split neural networks. Vertical FL (VFL) is a setting where clients hold different feature subsets of the same samples. EF-VFL addresses the communication bottleneck in VFL by using lossy compression on communicated information and leveraging error feedback to maintain convergence. The method is shown to achieve improved convergence rates compared to prior art in VFL with compression and also supports the use of private labels.⁹⁷
- **Direct Link:** <https://github.com/pedromvaldeira/EF-VFL> (Code), <https://arxiv.org/abs/2406.14420> (Paper)

The field of Federated Learning is rapidly specializing to address the complex, real-world challenges that arise when applying collaborative learning principles to diverse applications and data types. Initial FL research, such as the development of FedAvg, laid the foundational framework. However, practical deployments quickly revealed significant hurdles, including substantial communication overhead when dealing with large models, performance degradation due to heterogeneous (non-IID) client data, vulnerabilities to malicious attacks, and the need to accommodate various data distribution setups (e.g., horizontal versus vertical FL). Recent research, as exemplified by the cited works, shows a clear trend towards developing tailored solutions for these specific issues. Ferret⁸⁹ directly targets the problem of conducting full-parameter tuning of massive LLMs within an FL paradigm, with a keen focus on optimizing communication efficiency and overall scalability. FedDDL⁹³, emerging from research groups like Han Yu's lab known for extensive FL contributions⁹², tackles the critical challenge of out-of-distribution generalization and attribute bias in FL, which is crucial for ensuring model robustness when client data distributions vary significantly. FLTG⁹⁵ concentrates on Byzantine robustness, a key security concern in FL, particularly in non-IID settings where distinguishing malicious updates from benign variations due to data heterogeneity is difficult. Lastly, EF-VFL⁹⁷ addresses communication efficiency specifically within the Vertical Federated Learning setting, a distinct configuration from the more

commonly studied Horizontal FL.

This increasing specialization indicates that FL is maturing and moving closer to widespread, practical adoption. However, it also underscores that "one-size-fits-all" FL algorithms are generally insufficient for the diverse range of potential applications. Future FL systems will likely need to be highly customized, taking into account the specific model architecture being trained, the characteristics of the distributed data, the prevailing security requirements, and the particular federated setup (e.g., cross-silo vs. cross-device, horizontal vs. vertical).

Academic blogs, such as those potentially stemming from active research labs like Han Yu's ⁹², often play a role in synthesizing these specialized advancements, for example, by discussing overarching challenges in "Federated LLMs" and pointing towards emerging solutions.

D. AI Alignment Research

Ensuring that advanced AI systems behave in accordance with human intentions and values—a field broadly known as AI alignment—is a critical area of research, particularly as AI models become more autonomous and capable.

- **Full Title:** A Statistical Case Against Empirical Human-AI Alignment
- **Author(s):** Julian Rodemann, Esteban Garces Arias, Christoph Luther, Christoph Jansen, Thomas Augustin
- **Year:** 2025 (arXiv:2502.14581v2, May 2025)
- **Key Contribution:** This position paper presents a critique of naive empirical approaches to human-AI alignment, particularly those involving forward (a priori) alignment where alignment is attempted during the training phase before deployment. The authors argue that such methods, which often rely on learning from observed human behavior or preferences (e.g., RLHF), can inadvertently introduce significant statistical biases and may fail to capture the true, underlying human purpose. As alternatives, they advocate for prescriptive alignment (based on explicitly defined axioms or rules) and a posteriori empirical alignment (involving monitoring, evaluation, and adjustment of AI systems after deployment).⁹⁹
- **Direct Link:** <https://arxiv.org/abs/2502.14581>

The field of AI alignment is currently undergoing a period of critical re-evaluation, with a discernible movement towards developing more robust and statistically grounded approaches. There is a growing apprehension that prevailing empirical methods for AI alignment, such as Reinforcement Learning from Human Feedback (RLHF), might possess inherent flaws. These concerns stem from the potential for statistical biases to be introduced when learning from observed human behavior and the fundamental difficulty of accurately capturing true human intent solely from such observations. This critical perspective is leading to proposals for more principled, axiom-based, or post-hoc alignment strategies. The position paper by Rodemann et al. ⁹⁹ cogently argues that empirical alignment—especially when attempted *a priori* during the training phase—can inadvertently introduce statistical biases that skew the AI's behavior away from the "purpose which we really desire" towards what they term a "colorful imitation of it." They specifically critique this forward empirical alignment and instead champion *prescriptive alignment*, which involves guiding AI behavior through explicitly

defined axioms or rules, and *backward (or a posteriori) empirical alignment*, which focuses on monitoring, evaluating, and adjusting AI systems once they are deployed. This suggests a potential shift in the field away from a naive trust in the ability of AI to learn human values purely from observational data, towards a more cautious and principled approach that incorporates explicit human-defined guidelines and mandates ongoing oversight and intervention.

This critical discourse implies that the future of AI alignment is likely to involve more debate and a diversification of methodologies. If purely empirical methods are indeed susceptible to significant and potentially harmful biases, then a greater emphasis on formal methods, the development of ethical frameworks encoded as operational rules (prescriptive alignment), and the establishment of robust post-deployment auditing and correction mechanisms will become necessary for creating AI systems that are both safe and genuinely aligned with human values. This re-evaluation could also influence how existing techniques like RLHF are approached, perhaps by integrating more explicit constraints, principles, or oversight mechanisms into the feedback and learning loop.

Cross-cutting themes in general ML theory and optimization revolve around achieving a "practicality triad": robust *generalization*, enhanced computational and data *efficiency*, and the enablement of *trustworthy collaboration*. The quest for a deeper understanding of generalization remains a core theoretical pursuit, with research striving for non-vacuous bounds by meticulously considering model architecture and learning dynamics.⁸⁴ Such understanding is fundamental to building reliable and predictable models. Efficiency is a major driver of innovation, particularly given the escalating scale of modern models. This is evident in the exploration of novel optimization algorithms, including zeroth-order methods⁷⁹, and in the continuous refinement of federated learning techniques to minimize communication overhead and computational demands.⁸⁹ Trustworthy collaboration encompasses several facets: developing federated learning systems that are private by design and robust against various failures or attacks⁹³, and advancing AI alignment research to ensure that AI systems operate in a manner consistent with human values and intentions.⁹⁹ These three pillars—generalization, efficiency, and trustworthy collaboration—are not independent but are often interconnected and sometimes in tension, requiring holistic approaches for continued progress in the field.

VI. Conclusion

The landscape of machine learning is characterized by rapid advancements across multiple interconnected domains. In **Face Recognition**, the pursuit of higher accuracy and robustness is now intricately linked with the critical need for demographic fairness, driving innovations in loss functions, architectural designs like Transformers, and the strategic use of synthetic data. Human-machine collaboration is also emerging as a key factor for enhancing reliability in high-stakes FR applications.

Natural Language Processing is dominated by the evolution of Large Language Models. Research is intensely focused on a "trilemma": enhancing their diverse capabilities (e.g., reasoning, multimodality), improving their operational efficiency (through novel architectures

and compression), and ensuring their trustworthiness (by mitigating hallucinations, biases, and enabling knowledge editing). The field is exploring alternatives to standard Transformers, refining methods to imbue LLMs with better reasoning, and developing sophisticated techniques to manage their knowledge and behavior post-training.

Computer Vision is witnessing the maturation of Vision Transformers through hybridization and adaptability enhancements, making them more practical and efficient. Generative models are moving beyond mere visual fidelity towards "world fidelity," aiming to create content that is physically plausible and controllable. Concurrently, Self-Supervised Learning is pivoting towards greater resource efficiency and adaptability, democratizing access to powerful representation learning. These three pillars—Transformers, generative models, and SSL—are increasingly synergistic.

Reinforcement Learning is advancing on several fronts. Explainable RL is becoming crucial for deploying agents in complex, high-stakes environments, especially when integrated with LLMs. Offline RL is leveraging external data sources like unlabeled videos and generative models to overcome the limitations of static datasets. Multi-Agent RL is tackling fundamental issues of stability and coordination complexity using more sophisticated mathematical frameworks and modeling techniques. Furthermore, RL from Human Feedback is evolving towards greater personalization and improved learning efficiency to better align AI with nuanced human preferences.

Finally, **General ML Theory and Optimization** are seeing progress in foundational areas. Novel optimization techniques, including gradient-free (zeroth-order) methods, are emerging as viable solutions for complex deep learning scenarios like continual learning and efficient LLM fine-tuning. The theoretical understanding of generalization in deep learning is advancing through efforts to derive non-vacuous bounds that consider model structure and learning dynamics. Federated Learning is rapidly specializing to address real-world complexities such as scaling for LLMs, ensuring robustness against out-of-distribution data and Byzantine attacks, and improving communication efficiency across various FL configurations. AI alignment research is undergoing critical re-evaluation, with a push towards more statistically grounded and principled approaches beyond purely empirical methods. Across all these domains, a common thread is the drive towards creating AI systems that are not only more capable and intelligent but also more efficient, reliable, fair, interpretable, and aligned with human values. The post-2023 research landscape indicates a clear trend towards tackling more nuanced and practical challenges, often through interdisciplinary approaches and by building upon the foundational breakthroughs of previous years. The continued exploration of these frontiers will be pivotal in shaping the future impact of machine learning.

Works cited

1. ArcFace - Additive Angular Margin Loss For Deep Face Recognition - Scribd, accessed June 1, 2025, <https://www.scribd.com/document/706535195/ArcFace-Additive-Angular-Margin-Loss-for-Deep-Face-Recognition>
2. Transformer-Based Auxiliary Loss for Face Recognition Across Age Variations -

- arXiv, accessed June 1, 2025, <https://arxiv.org/pdf/2412.02198?>
3. A Comparison of Human and Machine Learning Errors in Face Recognition - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2502.11337v1>
 4. Review of Demographic Fairness in Face Recognition - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2502.02309v2>
 5. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2502.02309>
 6. Review of Demographic Fairness in Face Recognition - arXiv, accessed June 1, 2025, <https://arxiv.org/pdf/2502.02309>
 7. VariFace: Fair and Diverse Synthetic Dataset Generation for Face Recognition - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2412.06235v2>
 8. [Literature Review] VariFace: Fair and Diverse Synthetic Dataset Generation for Face Recognition - Moonlight | AI Colleague for Research Papers, accessed June 1, 2025, <https://www.themoonlight.io/en/review/variface-fair-and-diverse-synthetic-dataset-generation-for-face-recognition>
 9. VariFace: Fair and Diverse Synthetic Dataset Generation for Face Recognition, accessed June 1, 2025, https://www.researchgate.net/publication/386577486_VariFace_Fair_and_Diverse_Synthetic_Dataset_Generation_for_Face_Recognition
 10. arxiv.org, accessed June 1, 2025, <https://arxiv.org/pdf/2412.06235.pdf>
 11. Generative Physical AI in Vision: A Survey - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2501.10928v2>
 12. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2403.12960>
 13. FaceXFormer: A Unified Transformer for Facial Analysis - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2403.12960v3>
 14. Face-LLaVA: Facial Expression and Attribute Understanding through Instruction Tuning, accessed June 1, 2025, <https://arxiv.org/html/2504.07198>
 15. [2504.07198] Face-LLaVA: Facial Expression and Attribute Understanding through Instruction Tuning - arXiv, accessed June 1, 2025, <https://arxiv.org/abs/2504.07198>
 16. [Literature Review] Face-LLaVA: Facial Expression and Attribute Understanding through Instruction Tuning - Moonlight, accessed June 1, 2025, <https://www.themoonlight.io/en/review/face-llava-facial-expression-and-attribute-understanding-through-instruction-tuning>
 17. A Comprehensive Review of Face Recognition Techniques, Trends and Challenges, accessed June 1, 2025, https://www.researchgate.net/publication/382098748_A_Comprehensive_Review_of_Face_Recognition_Techniques_Trends_and_Challenges
 18. (PDF) Deep Face Recognition: A Survey - ResearchGate, accessed June 1, 2025, https://www.researchgate.net/publication/324600003_Deep_Face_Recognition_A_Survey
 19. LLLMs: A Data-Driven Survey of Evolving Research on Limitations of Large Language Models - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2505.19240v1>
 20. [2505.19240] LLLMs: A Data-Driven Survey of Evolving Research on Limitations of Large Language Models - arXiv, accessed June 1, 2025,

- <https://arxiv.org/abs/2505.19240>
21. Artificial Intelligence May 2025 - arXiv, accessed June 1, 2025, <http://arxiv.org/list/cs.AI/2025-05?skip=3400&show=25>
 22. arXiv:2411.01030v5 [cs.CL] 21 Feb 2025, accessed June 1, 2025, <https://arxiv.org/pdf/2411.01030?>
 23. arxiv.org, accessed June 1, 2025, <https://arxiv.org/pdf/2411.01030.pdf>
 24. arXiv:2505.12381v1 [cs.CL] 18 May 2025, accessed June 1, 2025, <https://arxiv.org/pdf/2505.12381>
 25. Advancing Reasoning in Large Language Models: Promising Methods and Approaches, accessed June 1, 2025, <https://arxiv.org/html/2502.03671v2>
 26. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2502.03671>
 27. arXiv:2502.15652v2 [cs.AI] 24 Feb 2025, accessed June 1, 2025, <https://arxiv.org/pdf/2502.15652?>
 28. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2502.15652>
 29. Mitigate Language Priors in Large Vision-Language Models by Cross-Images Contrastive Decoding - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2505.10634v2>
 30. [2505.10634] Cross-Image Contrastive Decoding: Precise, Lossless Suppression of Language Priors in Large Vision-Language Models - arXiv, accessed June 1, 2025, <https://arxiv.org/abs/2505.10634>
 31. Mitigate Language Priors in Large Vision-Language Models by Cross-Images Contrastive Decoding - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2505.10634v1>
 32. arxiv.org, accessed June 1, 2025, <https://arxiv.org/pdf/2505.10634.pdf>
 33. aclanthology.org, accessed June 1, 2025, <https://aclanthology.org/2025.findings-naacl.235.pdf>
 34. ConKE: Conceptualization-Augmented Knowledge Editing in Large Language Models for Commonsense Reasoning - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2412.11418v2>
 35. arXiv:2412.11418v2 [cs.CL] 28 May 2025, accessed June 1, 2025, <http://arxiv.org/pdf/2412.11418>
 36. arXiv:2412.11418v2 [cs.CL] 28 May 2025, accessed June 1, 2025, <https://arxiv.org/abs/2412.11418>
 37. Position: Multimodal Large Language Models Can Significantly Advance Scientific Reasoning - arXiv, accessed June 1, 2025, <https://arxiv.org/pdf/2502.02871?>
 38. arXiv:2502.05568v1 [cs.CL] 8 Feb 2025, accessed June 1, 2025, <https://arxiv.org/pdf/2502.05568?>
 39. Quantizing Large Language Models for Code Generation: A Differentiated Replication, accessed June 1, 2025, <https://arxiv.org/html/2503.07103v1>
 40. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2503.07103>
 41. arXiv:2504.02010v1 [cs.LG] 2 Apr 2025, accessed June 1, 2025, <https://arxiv.org/pdf/2504.02010>
 42. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2504.02010>
 43. ECViT: Efficient Convolutional Vision Transformer with Local-Attention and Multi-scale Stages - arXiv, accessed June 1, 2025,

- <https://arxiv.org/html/2504.14825v1>
44. [2504.14825] ECViT: Efficient Convolutional Vision Transformer with Local-Attention and Multi-scale Stages - arXiv, accessed June 1, 2025, <https://arxiv.org/abs/2504.14825>
 45. Artificial Intelligence Apr 2025 - arXiv, accessed June 1, 2025, <http://arxiv.org/list/cs.AI/2025-04?skip=1275&show=1000>
 46. UniViTAR: Unified Vision Transformer with Native Resolution - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2504.01792v1>
 47. UniViTAR: Unified Vision Transformer with Native Resolution - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2504.01792v2>
 48. [2504.01792] UniViTAR: Unified Vision Transformer with Native Resolution - arXiv, accessed June 1, 2025, <https://arxiv.org/abs/2504.01792>
 49. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2501.10928>
 50. VBench-2.0: Advancing Video Generation Benchmark Suite for Intrinsic Faithfulness - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2503.21755v1>
 51. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2503.21755>
 52. Controllable 3D Outdoor Scene Generation via Scene Graphs - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2503.07152v1>
 53. Controllable 3D Outdoor Scene Generation via Scene Graphs | Request PDF - ResearchGate, accessed June 1, 2025, https://www.researchgate.net/publication/389715074_Controllable_3D_Outdoor_Scene_Generation_via_Scene_Graphs
 54. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2503.07152>
 55. A curated list of awesome Indoor Scene Synthesis papers - GitHub, accessed June 1, 2025, <https://github.com/YandanYang/Awesome-Indoor-Scene-Synthesis>
 56. The 12 Best AI Blogs You Should Be Following - University of San Diego Online Degrees, accessed June 1, 2025, <https://onlinedegrees.sandiego.edu/ai-blogs/>
 57. 20 best ai blogs to read in 2024 - SMIAT, accessed June 1, 2025, <https://smiatblogs.com/20-best-ai-blogs-to-read-in-2024/>
 58. arXiv:2502.18056v1 [cs.CV] 25 Feb 2025, accessed June 1, 2025, <https://arxiv.org/pdf/2502.18056>
 59. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2502.18056>
 60. arXiv:2502.21313v1 [cs.CV] 28 Feb 2025, accessed June 1, 2025, <https://www.arxiv.org/pdf/2502.21313>
 61. A Survey on Explainable Deep Reinforcement Learning - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2502.06869v1>
 62. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2502.06869>
 63. arXiv:2502.06869v1 [cs.LG] 8 Feb 2025, accessed June 1, 2025, <https://arxiv.org/pdf/2502.06869?>
 64. Video-Enhanced Offline Reinforcement Learning: A Model-Based Approach - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2505.06482v2>
 65. accessed January 1, 1970, <https://arxiv.org/pdf/2505.06482.pdf>
 66. Variational OOD State Correction for Offline Reinforcement Learning - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2505.00503v1>
 67. Addressing Rotational Learning Dynamics in Multi-Agent Reinforcement Learning

- arXiv, accessed June 1, 2025, <https://arxiv.org/html/2410.07976v2>
- 68. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2410.07976>
- 69. Offline Multi-agent Reinforcement Learning via Score Decomposition - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2505.05968v1>
- 70. [2505.05968] Offline Multi-agent Reinforcement Learning via Score Decomposition - arXiv, accessed June 1, 2025, <https://arxiv.org/abs/2505.05968>
- 71. Machine Learning May 2025 - arXiv, accessed June 1, 2025, <https://www.arxiv.org/list/cs.LG/2025-05?skip=275&show=250>
- 72. accessed January 1, 1970, <https://arxiv.org/pdf/2505.05968.pdf>
- 73. Research-oriented Deep Offline Reinforcement Learning Library - NeurIPS Poster CORL, accessed June 1, 2025, <https://neurips.cc/virtual/2023/poster/73613>
- 74. arXiv:2503.19201v1 [cs.LG] 24 Mar 2025, accessed June 1, 2025, <https://arxiv.org/pdf/2503.19201>
- 75. Federated Learning: A Survey on Privacy-Preserving Collaborative Intelligence - arXiv, accessed June 1, 2025, <https://arxiv.org/pdf/2504.17703>
- 76. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2503.19201>
- 77. ma-rlhf: reinforcement learning from hu - arXiv, accessed June 1, 2025, <https://arxiv.org/pdf/2410.02743?>
- 78. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2410.02743>
- 79. ZeroFlow: Overcoming Catastrophic Forgetting is Easier than You Think - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2501.01045v3>
- 80. [2501.01045] ZeroFlow: Overcoming Catastrophic Forgetting is Easier than You Think, accessed June 1, 2025, <https://arxiv.org/abs/2501.01045>
- 81. accessed January 1, 1970, <https://arxiv.org/pdf/2501.01045.pdf>
- 82. LORENZA: Enhancing Generalization in Low-Rank Gradient LLM Training and Fine-Tuning via Efficient Zeroth-Order Adaptive SAM Optimization - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2502.19571v1>
- 83. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2502.19571>
- 84. Deep Multi-Task Learning Has Low Amortized Intrinsic Dimensionality - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2501.19067v1>
- 85. [2501.19067] Deep Multi-Task Learning Has Low Amortized Intrinsic Dimensionality - arXiv, accessed June 1, 2025, <https://arxiv.org/abs/2501.19067>
- 86. accessed January 1, 1970, <https://arxiv.org/pdf/2501.19067.pdf>
- 87. Survey on Generalization Theory for Graph Neural Networks arXiv:2503.15650v1 [cs.LG] 19 Mar 2025, accessed June 1, 2025, <https://arxiv.org/pdf/2503.15650>
- 88. Publications | GLOW.AI - NUS Computing, accessed June 1, 2025, <https://www.comp.nus.edu.sg/~lowkh/publications.html>
- 89. Ferret: Federated Full-Parameter Tuning at Scale for Large Language Models - arXiv, accessed June 1, 2025, <https://arxiv.org/pdf/2409.06277>
- 90. Ferret: Federated Full-Parameter Tuning at Scale for Large Language Models - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2409.06277v2>
- 91. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2409.06277>
- 92. Han Yu's Homepage - The Federated Learning Portal, accessed June 1, 2025, <https://federated-learning.org/han.yu/>
- 93. Federated Deconfounding and Debiasing Learning for Out-of-Distribution

- Generalization - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2505.04979v2>
94. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2505.04979>
95. FLTG: Byzantine-Robust Federated Learning via Angle-Based Defense and Non-IID-Aware Weighting - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2505.12851v1>
96. accessed January 1, 1970, <https://arxiv.org/abs/2505.12851>
97. [2406.14420] Communication-efficient Vertical Federated Learning via Compressed Error Feedback - arXiv, accessed June 1, 2025, <https://arxiv.org/abs/2406.14420>
98. Free-Rider and Conflict Aware Collaboration Formation for Cross-Silo Federated Learning - arXiv, accessed June 1, 2025, <https://arxiv.org/pdf/2410.19321?>
99. A Statistical Case Against Empirical Human-AI Alignment - arXiv, accessed June 1, 2025, <https://arxiv.org/html/2502.14581v1>
100. A Statistical Case Against Empirical Human-AI Alignment - arXiv, accessed June 1, 2025, <https://arxiv.org/pdf/2502.14581>
101. arxiv.org, accessed June 1, 2025, <https://arxiv.org/abs/2502.14581>