



**ARE: scaling up agent environments and evaluations**

We introduce Meta Agents Research Environments (ARE), a research platform for scalable creation of environments, integration of synthetic or real applications, and execution of agentic orchestrations. ARE provides simple abstractions to build complex and diverse environments, each with their own rules, tools, content, and verifiers, helping to bridge the gap...

Submitted by  24 authors · Published on Sep 22, 2025

**ScaleCUA: Scaling Open-Source Computer Use Agents with Cross-Platform Data**

Vision-Language Models (VLMs) have enabled computer use agents (CUAs) that operate GUIs autonomously, showing great potential, yet progress is limited by the lack of large-scale, open-source computer use data and foundation models. In this work, we introduce ScaleCUA, a step toward scaling open-source CUAs. It offers a large-scale dataset spanning 6...

Submitted by  21 authors · Published on Sep 19, 2025

**Qwen3 Technical Report**

In this work, we present Qwen3, the latest version of the Qwen model family. Qwen3 comprises a series of large language models (LLMs) designed to advance performance, efficiency, and multilingual capabilities. The Qwen3 series includes models of both dense and Mixture-of-Expert (MoE) architectures, with parameter scales ranging from 0.6 to 235 billion. A key...

Submitted by  60 authors · Published on May 14, 2025

**Logics-Parsing Technical Report**

Recent advances in Large Vision-Language models (LVLM) have spurred significant progress in document parsing task. Compared to traditional pipeline-based methods, end-to-end paradigms have shown their excellence in converting PDF images into structured outputs through integrated Optical Character Recognition (OCR), table recognition, mathematical...

Submitted by  10 authors · Published on Sep 24, 2025

**Collaborating Action by Action: A Multi-agent LLM Framework for Embodied Reasoning**

Collaboration is ubiquitous and essential in day-to-day life -- from exchanging ideas, to delegating tasks, to generating plans together. This work studies how LLMs can adaptively collaborate to perform complex embodied reasoning tasks. To this end we introduce MINDcraft, an easily extensible platform built to enable LLM agents to control characters in the...

Submitted by deleted 8 authors · Published on Apr 25, 2025

**VoXtream: Full-Stream Text-to-Speech with Extremely Low Latency**

We present VoXtream, a fully autoregressive, zero-shot streaming text-to-speech (TTS) system for real-time use that begins speaking from the first word. VoXtream directly maps incoming phonemes to audio tokens using a monotonic alignment scheme and a dynamic look-ahead that does not delay onset. Built around an incremental phoneme transformer, a...

Submitted by deleted 3 authors · Published on Sep 19, 2025

**VolSplat: Rethinking Feed-Forward 3D Gaussian Splatting with Voxel-Aligned Prediction**

Feed-forward 3D Gaussian Splatting (3DGS) has emerged as a highly effective solution for novel view synthesis. Existing methods predominantly rely on a pixel-aligned Gaussian prediction paradigm, where each 2D pixel is mapped to a 3D Gaussian. We rethink this widely adopted formulation and identify several inherent limitations: it renders the reconstruct...

Submitted by deleted 10 authors · Published on Sep 24, 2025

**Youtu-GraphRAG: Vertically Unified Agents for Graph Retrieval-Augmented Complex Reasoning**

Graph retrieval-augmented generation (GraphRAG) has effectively enhanced large language models in complex reasoning by organizing fragmented knowledge into explicitly structured graphs. Prior efforts have been made to improve either graph construction or graph retrieval in isolation, yielding suboptimal performance, especially when domain shifts occur. In thi...

Submitted by deleted 9 authors · Published on Aug 27, 2025

**LIMI: Less is More for Agency**

We define Agency as the emergent capacity of AI systems to function as autonomous agents actively discovering problems, formulating hypotheses, and executing solutions through self-directed engagement with environments and tools. This fundamental capability marks the dawn of the Age of AI Agency, driven by a critical industry shift: the urgent need for AI...

Submitted by Yangxiao-nlp 21 authors · Published on Sep 22, 2025

**Easy Dataset: A Unified and Extensible Framework for Synthesizing LLM Fine-Tuning Data from Unstructured Documents**

Large language models (LLMs) have shown impressive performance on general-purpose tasks, yet adapting them to specific domains remains challenging due to the scarcity of high-quality domain data. Existing data synthesis tools often struggle to extract reliable fine-tuning data from heterogeneous documents effectively. To address this limitation, we propose...

Submitted by hiyouga 7 authors · Published on Jul 5, 2025

**MAPO: Mixed Advantage Policy Optimization**

Recent advances in reinforcement learning for foundation models, such as Group Relative Policy Optimization (GRPO), have significantly improved the performance of foundation models on reasoning tasks. Notably, the advantage function serves as a central mechanism in GRPO for ranking the trajectory importance. However, existing explorations encounter...

Submitted by WilliamHuang91 14 authors · Published on Sep 23, 2025

**3D and 4D World Modeling: A Survey**

World modeling has become a cornerstone in AI research, enabling agents to understand, represent, and predict the dynamic environments they inhabit. While prior work largely emphasizes generative methods for 2D image and video data, they overlook the rapidly growing body of work that leverages native 3D and 4D representations such as RGB-D imagery,...

Submitted by  23 authors · Published on Sep 5, 2025

**Is Diversity All You Need for Scalable Robotic Manipulation?**

Data scaling has driven remarkable success in foundation models for Natural Language Processing (NLP) and Computer Vision (CV), yet the principles of effective data scaling in robotic manipulation remain insufficiently understood. In this work, we investigate the nuanced role of data diversity in robot learning by examining three critical dimensions-task (what to do...

Submitted by yxluo 10 authors · Published on Jul 9, 2025

**VGGT: Visual Geometry Grounded Transformer**

We present VGGT, a feed-forward neural network that directly infers all key 3D attributes of a scene, including camera parameters, point maps, depth maps, and 3D point tracks, from one, a few, or hundreds of its views. This approach is a step forward in 3D computer vision, where models have typically been constrained to and specialized for single tasks. It is also...

Submitted by JianyangWang 6 authors · Published on Mar 15, 2025

**Paper2Code: Automating Code Generation from Scientific Papers in Machine Learning**

Despite the rapid growth of machine learning research, corresponding code implementations are often unavailable, making it slow and labor-intensive for researchers to reproduce results and build upon prior work. In the meantime, recent Large Language Models (LLMs) excel at understanding scientific documents and generating high-quality code. Inspired by...

Submitted by Seongyun 4 authors · Published on Apr 24, 2025

**SIM-CoT: Supervised Implicit Chain-of-Thought**

Implicit Chain-of-Thought (CoT) methods present a promising, token-efficient alternative to explicit CoT reasoning in Large Language Models (LLMs), but a persistent performance gap has limited the application of implicit CoT. We identify a core latent instability issue by scaling the computational budget of implicit CoT methods. Our proposed solution is to use a hybrid approach that combines implicit CoT with explicit CoT, which allows us to scale implicit CoT without losing the benefits of explicit CoT. This leads to a significant performance improvement over previous implicit CoT methods.

Submitted by  8 authors · Published on Sep 25, 2025

**MonkeyOCR: Document Parsing with a Structure-Recognition-Relation Triplet Paradigm**

We introduce MonkeyOCR, a vision-language model for document parsing that advances the state of the art by leveraging a Structure-Recognition-Relation (SRR) triplet paradigm. This design simplifies what would otherwise be a complex multi-tool pipeline (as in MinerU's modular approach) and avoids the inefficiencies of processing full pages with giant end-to-en...

Submitted by deleted 10 authors · Published on Jun 6, 2025

**CHARM: Control-point-based 3D Anime Hairstyle Auto-Regressive Modeling**

We present CHARM, a novel parametric representation and generative framework for anime hairstyle modeling. While traditional hair modeling methods focus on realistic hair using strand-based or volumetric representations, anime hairstyle exhibits highly stylized, piecewise-structured geometry that challenges existing techniques. Existing works often rely on...

Submitted by hyz317 9 authors · Published on Sep 25, 2025

**UniPixel: Unified Object Referring and Segmentation for Pixel-Level Visual Reasoning**

Recent advances in Large Multi-modal Models (LMMs) have demonstrated their remarkable success as general-purpose multi-modal assistants, with particular focuses on holistic image- and video-language understanding. Conversely, less attention has been given to scaling fine-grained pixel-level understanding capabilities, where the models are expected to...

Submitted by  7 authors · Published on Sep 23, 2025

**MixGROPO: Unlocking Flow-based GRPO Efficiency with Mixed ODE-SDE**

Although GRPO substantially enhances flow matching models in human preference alignment of image generation, methods such as FlowGRPO still exhibit inefficiency due to the necessity of sampling and optimizing over all denoising steps specified by the Markov Decision Process (MDP). In this paper, we propose MixGROPO, a novel framework that leverages t...

Submitted by  7 authors · Published on Jul 29, 2025

**FastVLM: Efficient Vision Encoding for Vision Language Models**

Scaling the input image resolution is essential for enhancing the performance of Vision Language Models (VLMs), particularly in text-rich image understanding tasks. However, popular visual encoders such as ViTs become inefficient at high resolutions due to the large number of tokens and high encoding latency caused by stacked self-attention layers. At different...

Submitted by hpouansari 11 authors · Published on Dec 18, 2024

**VideoFrom3D: 3D Scene Video Generation via Complementary Image and Video Diffusion Models**

In this paper, we propose VideoFrom3D, a novel framework for synthesizing high-quality 3D scene videos from coarse geometry, a camera trajectory, and a reference image. Our approach streamlines the 3D graphic design workflow, enabling flexible design exploration and rapid production of deliverables. A straightforward approach to synthesizing a video....

Submitted by  3 authors · Published on Sep 23, 2025