











NEWMIND AI JOURNAL MONTHLY CHRONICLES




1.3.2025 - 31.3.2025





- March 2025 represents a pivotal month in artificial intelligence development, marking significant advancements across models, hardware, and methodologies that collectively push the boundaries of what AI systems can achieve.
- The period is characterized by an intensified focus on reasoning capabilities, with numerous models specifically designed to enhance logical thinking, step-by-step problem solving, and complex decision making.
- Multimodal integration reached new heights, with models increasingly capable of seamlessly processing and generating content across text, images, audio, video, and specialized data formats like scientific and medical information.
- Competition among major tech companies has accelerated, with Amazon, Meta, OpenAI, Google, Microsoft, Mistral, and an expanding roster of Chinese firms all releasing significant model updates and innovations.
- The hardware landscape continues to evolve rapidly, with new chip architectures, manufacturing techniques, and supply chain developments shaping the computational foundation of advanced AI systems.
- Methodological innovations in training, evaluation, and deployment are creating more efficient, reliable, and responsible AI systems capable of tackling increasingly complex real-world challenges.
- Open-source initiatives gained substantial momentum, democratizing access to powerful AI capabilities while fostering global innovation and collaboration.




 Models					
#	Highlights	Summary	Author	Source	Date
1.1	Amazon Launches AI-Powered Alexa+	Amazon has introduced Alexa+, an upgraded version of its virtual assistant powered by generative AI. Designed to be more conversational, intuitive, and capable, Alexa+ can handle complex tasks such as making reservations, managing smart homes, summarizing topics, and providing personalized assistance across entertainment, learning, and shopping. By leveraging advanced AI, it enables natural dialogues and proactive assistance, enhancing the overall user experience. With these advancements, Amazon aims to strengthen its position in the	By Panos Panay, SVP of Devices & Services		February 26, 2025





 Models					
#	Highlights	Summary	Author	Source	Date
		competitive virtual assistant market, making Alexa+ a more powerful and interactive tool for users.			
1.2	Meta Unveils Voice-Centric LLaMA 4	Meta Platforms is advancing into voice-driven AI with its upcoming LLaMA 4 model, an “omni model” designed for natural, two-way speech interactions without converting speech to text. Meta sees voice as the future of AI assistants, enabling more human-like conversations. With up to \$65 billion invested in AI by 2025, the company aims for LLaMA 4 to power interactive AI assistants in smart glasses, customer support, and real-time translation. By making AI responses interruptible and context-aware, Meta seeks to create seamless, conversational AI experiences across various applications.	By PYMNTS		March 7, 2025
1.3	Evo 2: Largest AI Genome Model	Researchers from the Arc Institute, Stanford, and NVIDIA have introduced Evo 2, the largest AI model for biology to date. Trained on 128,000 genomes—equivalent to 9.3 trillion nucleotides—Evo 2 spans all domains of life and is capable of generating entire chromosomes and small genomes, as well as interpreting DNA sequences with unprecedented accuracy. This breakthrough marks a significant advancement in AI-driven biological research, with potential applications in genetics, synthetic biology, and medical science, offering deeper insights into DNA functions and evolution across diverse species.	By NVIDIA		February 19, 2025
1.4	OpenAI's GPT-4.5 Release	OpenAI has unveiled GPT-4.5, code-named "Orion," its largest and most advanced AI model to date. This release emphasizes enhanced anthropomorphic features, including a deeper understanding of emotions and more intuitive communication, aiming to provide users with a more natural conversational experience. Despite its increased capabilities and computational demands, OpenAI has provided	By Reece Rogers		March 6, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		<p>limited details regarding the model's exact improvements. Initially, GPT-4.5 is available to ChatGPT Pro subscribers, with plans for broader access in the future. This development underscores OpenAI's ongoing efforts to balance its ambitions in artificial general intelligence with its commercial objectives.</p>			
1.5	<p>Microsoft Challenges OpenAI with Reasoning-AI</p>	<p>Microsoft is reportedly developing new AI reasoning models to compete with OpenAI's advanced systems. These models focus on improved logical reasoning and decision-making capabilities, positioning Microsoft as a strong competitor in the field of general AI development. This initiative highlights Microsoft's commitment to advancing AI technologies that can perform complex cognitive tasks, potentially reshaping the landscape of AI applications in business and technology sectors.</p>	By Reuters		March 7, 2025
1.6	<p>Mistral OCR Redefines Document AI</p>	<p>On March 6, 2025, Mistral AI launched Mistral OCR, a cutting-edge document understanding tool that outperforms competitors like Google Document AI and GPT-4o in interpreting media, text, tables, and equations, especially in multilingual and mathematical contexts. Processing 2,000 pages per minute, it offers exceptional speed. Its "Doc-as-Prompt" feature outputs structured JSON, streamlining AI integrations. Supporting thousands of scripts and fonts ensures broad language coverage. Available via API, Le Chat, and soon cloud partners, it also supports self-hosting for security. Mistral OCR is revolutionizing scientific digitization, legal indexing, and large-scale knowledge extraction.</p>	By Mistral AI Team		March 6, 2025
1.7	<p>QwQ-32B: Smaller, Smarter AI Model</p>	<p>Qwen team introduced QwQ-32B, a 32 billion parameter AI model that rivals the much larger DeepSeek-R1 in performance. The key breakthrough lies in scaling Reinforcement Learning (RL) to enhance reasoning, tool use, and adaptability. The</p>	By Ryan Daws		March 6, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		model excels in mathematical reasoning, coding, and problem-solving, as shown by benchmarks like AIME24, LiveCodeBench, and BFCL. QwQ-32B’s success demonstrates that RL can bridge the gap between model size and performance, offering an efficient alternative to conventional AI training methods.			
1.8	Microsoft Unveils Compact Phi-4 Models	Microsoft's Phi-4 series debuts with Phi-4-Mini, a 3.8B parameter language model excelling in math and coding. Trained on curated web and synthetic data, it rivals larger models in complex reasoning. Phi-4-Multimodal expands this, integrating text, vision, and speech. Utilizing LoRA adapters and modality routers, it enables diverse inference modes without interference. Notably, it leads the OpenASR leaderboard with a small speech component. It adeptly handles vision-language, vision-speech, and speech-audio tasks, surpassing larger models. These compact models demonstrate significant advancements in efficient and versatile AI processing.	By Microsoft		February 27, 2025
1.9	AI21 Labs Launches Jamba, Maestro	AI21 Labs unveiled Jamba 1.6 on March 6, an open model for private enterprise. Outperforming rivals like Llama 3.3 70B, it features a 256K context window and a hybrid Mamba-Transformer MoE architecture, balancing performance and cost. Accessible via AI21 Studio, Hugging Face, or private deployment, it offers flexibility and security. On March 10, they launched Maestro, an AI planning system optimizing LLM tasks. Maestro breaks complex prompts into substeps, ensuring accuracy through simulations and user-defined requirements. This enhances AI output reliability, benefiting document analysis and process automation.	By AI21 Editorial Team		March 6, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
1.10	Google Unveils Gemini Embedding Model	On March 7, 2025, Google introduced a new text embedding model based on its Gemini AI architecture. This model is now available to developers through the Gemini API, under the identifier gemini-embedding-exp-03-07. It generates state-of-the-art embeddings for words, phrases, code, and sentences, enabling applications such as semantic search, text classification, and clustering. This release underscores Google's commitment to advancing natural language understanding and providing developers with robust tools to enhance their AI-driven applications.	By Logan Kilpatrick, Zach Gleicher, Parashar Shah		March 7, 2025
1.11	Vision-R1: Multimodal Reasoning Redefined	Vision-R1 is a multimodal large language model (MLLM) designed to improve reasoning capabilities through reinforcement learning (RL). Inspired by DeepSeek-R1-Zero, it addresses reasoning challenges in MLLMs by integrating cold-start initialization with RL training. High-Quality Dataset: A 200K multimodal Chain-of-Thought (CoT) dataset was generated without human annotation, enhancing reasoning capabilities. Progressive Training: Introduces Progressive Thinking Suppression Training (PTST) to prevent overthinking and refine logical steps. Performance: Achieves 73.5% accuracy on MathVista, nearly matching OpenAI O1 and outperforming some models with 70B+ parameters despite being only 7B. Vision-R1 demonstrates strong reasoning skills by incorporating human-like cognitive processes, such as questioning and reflection.	By Wenxuan Huang and Bohan Jia and Zijie Zhai and Shaosheng Cao and Zheyu Ye and Fei Zhao and Zhe Xu and Yao Hu and Shaohui Lin		March 11, 2025
1.12	EuroBERT	A collaboration between CentraleSupélec's MICS laboratory, Diabolocom, Artefact, and Unbabel, supported by AMD and CINES, introduced EuroBERT, a state-of-the-art multilingual encoder model. Trained on 5 trillion tokens, it supports 15 languages, excelling in mathematics and programming tasks. With grouped query attention and rotary position embeddings, it handles sequences	By Nicolas-BZRD, Hippolyte Gisserot-Boukhlef, Duarte		March 10, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		up to 8,192 tokens. Available in three sizes (210M, 610M, and 2.1B parameters), it performs well in retrieval, classification, and regression. Open-sourced under Apache 2.0, EuroBERT is accessible on Hugging Face, fostering research in multilingual NLP.	Alves, Manuel Faysse		
1.13	AMD Launches Open-Source Instella LLM	AMD has launched Instella, a series of fully open-source language models with 3 billion parameters, trained on 128 Instinct MI300X GPUs. Instella-3B features 36 decoder layers, 32 attention heads per layer, and supports sequences up to 4,096 tokens with a 50,000-token vocabulary. Using FlashAttention-2, Torch Compile, and FSDP with hybrid sharding, it optimizes performance and resource efficiency. Instella-3B outperforms similar open models and competes with Llama-3.2-3B and Qwen-2.5-3B. AMD has open-sourced weights, configurations, datasets, and code, fostering AI research and innovation.	By Jiang Liu, Jialian Wu, Xiaodong Yu, Prakamya Mishra, Sudhanshu Ranjan, Zicheng Liu, Chaitanya Manem, Yusheng Su, Pratik Prabhanjan Brahma, Gowtham Ramesh, Ximeng Sun, Ze Wang, Emad Barsoum		March 5, 2025
1.14	STORM Revolutionizes Long Video Understanding	STORM (Spatiotemporal Token Reduction for Multimodal LLMs) is a novel architecture enhancing long video understanding in multimodal large language models (LLMs). It integrates a Mamba-based temporal projector between the image encoder and LLM, enriching visual tokens with temporal dynamics for improved video reasoning. STORM also employs efficient token reduction	By Nvidia		March 6, 2025





 Models					
#	Highlights	Summary	Author	Source	Date
		techniques, cutting computational costs by up to 8× and reducing decoding latency by 2.4–2.9×. Achieving state-of-the-art results, it improves MLVU and LongVideoBench scores by over 5%, making long video processing more efficient without compromising temporal information.			
1.15	Cohere Launches Multilingual Aya Vision	Cohere has unveiled Aya Vision, a state-of-the-art multimodal large language model that excels in both language and image understanding across 23 languages. This model supports tasks such as image captioning, visual question answering, text generation, and translations from both texts and images into coherent text. Aya Vision is available in two configurations: an 8-billion parameter variant optimized for low latency and performance, and a 32-billion parameter variant designed for state-of-the-art multilingual performance. Developers can access Aya Vision through the Cohere platform, enabling the creation of applications that seamlessly integrate multilingual and multimodal capabilities.	By Cohere For AI Team		March 4, 2025
1.16	DINOv2 Transforms AI Pathology Analysis	Meta’s AI blog highlights Dr. Faisal Mahmood's research on using DINOv2 to enhance pathology analysis through foundation models. The Mahmood Lab aligns with DINOv2’s hypothesis that data diversity is more important than sheer volume for effective AI training. By applying DINOv2 to pathology images, the lab has made significant progress in medical image analysis. Their approach demonstrates that models trained on diverse and representative datasets perform better in complex fields like pathology. This advancement paves the way for more accurate disease diagnosis and understanding, improving healthcare outcomes.	By Meta Team		March 6, 2025
1.17	Google Unveils Gemma 3: Faster,	Google has introduced Gemma 3, the latest addition to its open AI model family. The Gemma series is central to Google's vision of making AI more accessible. Last	By Google		March 12, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
	Portable, and Open AI Models	month, Gemma celebrated its first anniversary, reaching 100 million downloads and inspiring over 60,000 variants. Gemma 3 is a lightweight, state-of-the-art open model collection, built with the same technology as Gemini 2.0. It is faster, portable, and responsibly developed, optimized to run on phones, laptops, and workstations. Available in 1B, 4B, 12B, and 27B parameter sizes, it allows developers to choose the best model for their needs.			
1.18	DeepMind Unveils Gemini Robotics: Advancing AI-Powered Robot Interaction	Google DeepMind, Google’s AI research lab, announced new AI models called Gemini Robotics, designed to help robots interact with objects and navigate environments. DeepMind released demo videos showing robots using Gemini Robotics to fold paper, place glasses into a case, and perform other tasks via voice commands. The model was trained to generalize behavior across different robotics hardware and link visual inputs to actions. In tests, Gemini Robotics performed well in environments outside its training data. DeepMind also released Gemini Robotics-ER for researchers to train their own models and introduced Asimov, a benchmark for assessing risks in AI-powered robotics.	By Carolina Parada		March 12, 2025
1.19	Reasoning with Reka Flash 3	Reka Flash 3, a 21B parameter model, is a compact, general-purpose AI excelling in chat, coding, and instruction-following. Open-sourced under Apache 2.0, it offers competitive performance against proprietary models like OpenAI o1-mini, with 32k context length and efficient on-device deployment. Trained on diverse datasets, it uses RLOO for reinforcement learning, achieving strong multilingual and reasoning capabilities. Its lightweight design (11GB with 4-bit quantization) makes it ideal for cost-efficient, low-latency applications. Though not optimal for knowledge-intensive tasks, it’s a robust foundation for customization and domain-	By Reka AI		March 10, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		specific use, advancing accessible AI development for researchers and enterprises.			
1.20	Introducing Command A: Max performance, minimal compute	Command A is a state-of-the-art generative AI model optimized for enterprise use, offering superior performance, efficiency, and scalability. It outperforms competitors like GPT-4o and DeepSeek-V3 across business, STEM, and coding tasks while requiring minimal hardware (2 GPUs). With a 256k context length, advanced retrieval-augmented generation, multilingual capabilities, and enterprise-grade security, it excels in handling complex, real-world tasks. Command A delivers up to 156 tokens/sec, significantly faster than competitors, and supports cost-effective private deployments. Designed for seamless integration with enterprise systems, it enables AI-powered agents to securely process internal data, making it ideal for diverse industries and critical applications.	By Cohere Team		March 13, 2025
1.21	YOLOE	YOLOE introduces a highly efficient model for open-set object detection and segmentation, addressing limitations of predefined categories in conventional models like YOLO. It integrates text prompts (RepRTA), visual prompts (SAVPE), and prompt-free scenarios (LRPC) within a single framework, achieving real-time performance. YOLOE excels in zero-shot adaptability, with significant gains on benchmarks like LVIS (3.5 AP improvement) and COCO (0.6 APb, 0.4 APm gains) while reducing training costs and inference time. The model balances accuracy, efficiency, and transferability, making it suitable for diverse applications. Code and models are publicly available for further use.	By Ao Wang et al.		March 10, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
1.22	R1-Omni: Explainable Omni-Multimodal Emotion Recognition with Reinforcement Learning	The paper introduces R1-Omni, the first application of Reinforcement Learning with Verifiable Reward (RLVR) to omni-multimodal emotion recognition, integrating visual and audio data. By leveraging RLVR, the model significantly enhances reasoning, emotion recognition accuracy, and generalization, especially on out-of-distribution datasets. It provides interpretable reasoning for emotion predictions, advancing multimodal large language models. R1-Omni outperforms supervised fine-tuning approaches and demonstrates robust performance across diverse datasets. This work highlights RLVR's potential in optimizing multimodal models for complex tasks, offering valuable insights into the interplay of visual and audio modalities in emotion recognition.	By Jiaxing Zhao, Xihan Wei, Liefeng Bo, Tongyi Lab, Alibaba Group		March 10, 2025
1.23	Open R1: Update #3	This update highlights advancements in competitive programming and code reasoning through the release of the CodeForces-CoTs dataset, the IOI benchmark, and OlympicCoder models (7B and 32B). OlympicCoder-32B outperforms leading open and closed-weight models on challenging tasks, showcasing the effectiveness of fine-tuning on reasoning data. Key lessons include avoiding sample packing, using higher learning rates, and leveraging 8-bit optimizers for scalability. The update also emphasizes the need for fully verifiable datasets and introduces improvements in GRPO, making reasoning models more efficient. These contributions set new benchmarks for AI performance in competitive programming and reasoning tasks.	By Guilherme Penedo, Lewis Tunstall, Anton Lozhkov, Hynek Kydlicek, Edward Beeching, Loubna Ben Allal, Quentin Gallouédec, Leandro von Werra, Agustín Piqueres Lajarín, Nathan Habib		March 11, 2025




Models					
#	Highlights	Summary	Author	Source	Date
1.24	Wuhan's AI+ Initiative Empowers Zidong Taichu 3.0 for Industry Advancements	The Wuhan government launched the "AI+" initiative to advance AI industry growth, promoting large-model applications in over 20 sectors. Key measures include funding technological breakthroughs, boosting computing power, and encouraging innovation. Small and medium-sized enterprises will receive at least 10 million yuan annually in computing subsidies, while major tech projects can get up to 20 million yuan. A core part of this initiative, Zidong Taichu 3.0, is a multimodal AI model that integrates visual and audio data, enhancing performance in healthcare and manufacturing, driving industry advancements and improving efficiency.	By changjiang Daily		March 13, 2025
	Alibaba Enhances Quark AI Assistant with Advanced Reasoning	Alibaba has upgraded its Quark AI assistant with new reasoning capabilities, enabling it to handle complex queries and tasks such as academic research, document drafting, image generation, and travel planning. This enhancement positions Quark to compete with leading AI assistants and reflects Alibaba's significant investment in AI infrastructure, with plans to spend over \$50 billion in the next three years. The Quark application currently boasts over 200 million users, highlighting its widespread adoption and the growing competition in the AI sector.	By Gabbie Fu		March 13, 2025
1.25	Google Releases SpeciesNet for Wildlife Identification	Google has made its SpeciesNet AI model for wildlife identification publicly available as an open-source tool. Originally launched in 2019 on Google's cloud-based Wildlife Insights platform, SpeciesNet helps scientists identify species in trail camera images. The model was trained on over 65 million publicly available images and can classify animals across 2,000 categories from high-level taxa down	By The Wildlife Society		March 13, 2025





 Models					
#	Highlights	Summary	Author	Source	Date
		to specific species. This release enables developers, academics, and researchers to integrate the technology into their own wildlife monitoring projects, significantly reducing the time needed to process large volumes of camera trap data			
1.26	Cohere Inc.: Command A	Command A from Cohere Inc. introduces significant efficiency improvements, requiring only two GPUs for operation, compared to previous models that required up to 32. It offers a vast context window of 256,000 tokens and utilizes approximately 175 billion parameters, achieving remarkable processing speed and high accuracy in multilingual tasks	By Kyt Dotson		March 13, 2025
1.27	Wuhan Launches 'AI+' Initiative to Boost Artificial Intelligence Industry	The Wuhan Municipal People's Government launched the "AI+" initiative to accelerate AI adoption across 20+ industries. Based on Wuhan's AI industry policies, it includes ten measures targeting technological breakthroughs, computing power, and model innovation. To support SMEs, Wuhan will allocate 10M+ yuan annually for computing subsidies and offer up to 20M yuan for key tech projects. The city's AI sector is expected to exceed 70B yuan by 2024, reflecting 30%+ growth over three years, positioning Wuhan as a major AI hub.	By PR Newswire		March 14, 2025
1.28	Salesforce Unveils Agentforce 2dx at TDX 2025	At TrailblazerDX (TDX) 2025, Salesforce introduced Agentforce 2dx, the latest iteration of its platform designed for building, customizing, and deploying autonomous AI agents. This version expands beyond reactive chat interfaces, enabling proactive AI agents to operate autonomously, thereby unlocking new workflows for customers and employees. The update includes new pro-code, low-code, and no-code tools, lowering the barrier to entry for users without extensive coding backgrounds.	By Eira May		March 14, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
1.29	HPC-AI Tech Launches Open-Sora 2.0: Affordable Open-Source Video Generation Model	<p>HPC-AI Tech introduces Open-Sora 2.0, an open-source video generation model aimed at making high-quality video creation more accessible. By embracing open-source principles, Open-Sora democratizes advanced video generation while providing a user-friendly platform. The model integrates enhanced data selection, an advanced video autoencoder, a hybrid transformer framework, and optimized training techniques. With a training cost of just \$200,000—five to ten times cheaper than similar models—Open-Sora 2.0 significantly reduces financial barriers. This initiative fosters innovation and inclusivity, enabling developers and researchers to experiment, refine, and advance video generation technology.</p>	By Open-Sora Team		March 12, 2025
1.30	Baidu Unveils ERNIE 4.5 and ERNIE X1: Affordable AI Models Now Free for Users	<p>Baidu Inc. has unveiled ERNIE 4.5 and ERNIE X1, its latest foundation models, now free for individual users via the ERNIE Bot platform. ERNIE X1, a deep-thinking reasoning model with multimodal capabilities, matches DeepSeek R1 in performance but operates at half the cost. Meanwhile, ERNIE 4.5, Baidu’s newest multimodal foundation model, offers competitive pricing, with input costs starting at RMB 0.004 per 1,000 tokens and output at RMB 0.016 per 1,000 tokens. These advancements position Baidu as a key player in China’s AI sector, enhancing accessibility and affordability.</p>	By RTTNews.com		March 16, 2025





 Models					
#	Highlights	Summary	Author	Source	Date
1.31	Block Diffusion: Interpolating Between Autoregressive and Diffusion Language Models	<p>The study "Block Diffusion: Interpolating Between Autoregressive and Diffusion Language Models" introduces BD3-LMs, a hybrid approach to language modeling. Traditional autoregressive models generate high-quality, flexible-length text but lack parallel processing, while diffusion models support parallelism but often produce lower-quality, fixed-length outputs. BD3-LMs overcome these limitations by applying diffusion to token sequences divided into blocks, enabling high-quality, flexible, and parallel text generation. They improve performance through data-driven noise schedules and efficient training algorithms, reducing variance. Additionally, BD3-LMs enhance inference efficiency with key-value caching and parallel token sampling, setting a new benchmark for diffusion-based language models.</p>	By Marianne Arriola, Aaron Gokaslan, Justin Chiu, Zhihan Yang, Zhixuan Qi, Jiaqi Han, Subham Sahoo, Volodymyr Kuleshov		March 16, 2025
1.32	SmolDocling: An ultra-compact vision-language model for end-to-end multi-modal document conversion	<p>SmolDocling, a 256M-parameter vision-language model for end-to-end document conversion. It processes entire pages holistically, generating DocTags, a universal markup format capturing content, structure, and spatial layout. Unlike existing methods relying on large models or complex pipelines, SmolDocling offers a compact, efficient solution. It accurately reproduces code, tables, equations, charts, and lists across diverse documents, including patents and academic papers. Experimental results show it competes with models 27x larger while reducing computational costs. The model is available, with public datasets for tables, charts, and equations coming soon.</p>	By Ahmed Nassar, Andres Marafioti, Matteo Omenetti, Maksym Lysak, Nikolaos Livathinos, Christoph Auer, Lucas Morin, Rafael Teixeira de Lima, Yusik Kim, A. Said Gurbuz, Michele Dolfi, Miquel Farré, Peter W. J. Staar		March 14, 2025





Models					
#	Highlights	Summary	Author	Source	Date
1.33	OLMo 2 32B: First fully open model to outperform GPT 3.5 and GPT 4o mini	OLMo 2 32B is the largest and most capable model in the OLMo 2 family, marking a milestone as the first fully open model to outperform GPT-3.5 Turbo and GPT-4o Mini on key academic benchmarks. Trained on 6T tokens with advanced methodologies like reinforcement learning with verifiable rewards (RLVR), it achieves state-of-the-art results while requiring only a fraction of the compute cost of comparable models. With open access to all data, code, and weights, OLMo 2 32B fosters innovation, enabling researchers and developers to study and build cutting-edge AI models with unparalleled efficiency and transparency.	By Allen AI		March 13, 2025
1.34	VAMBA: Understanding Hour-Long Videos with Hybrid Mamba-Transformers	The paper introduces VAMBA, a hybrid Mamba-Transformer model designed to efficiently process hour-long videos. Traditional LMMs struggle with high memory usage and slow processing due to quadratic self-attention, limiting them to 256 frames. VAMBA integrates Mamba-2 blocks, achieving linear complexity, allowing it to handle 1,024+ frames at 640x360 resolution on a single GPU without token reduction. It reduces GPU memory usage by 50%, doubles training speed, and improves accuracy by 4.3% on LVBench. VAMBA delivers strong performance across long and short video tasks, making it a scalable, efficient alternative to transformer-based LMMs.	By Weiming Ren, Wentao Ma, Huan Yang, Cong Wei, Ge Zhang, Wenhui Chen, University of Waterloo, University of Toronto, 01.AI, Vector Institute, M-A-P		March 14, 2025
1.35	VGGT: Visual Geometry Grounded Transformer	VGGT: Visual Geometry Grounded Transformer" presents VGGT, a transformer-based neural network for extracting 3D geometric information from images. Unlike traditional methods, VGGT jointly predicts camera parameters, depth maps, point maps, and 3D trajectories in a single-pass inference, avoiding iterative optimization. It handles pose estimation, multi-view depth estimation, and 3D	By Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian		March 14, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		tracking efficiently, processing images in under a second. VGGT outperforms classical methods in 3D reconstruction and applies to autonomous driving, AR/VR, robotics, and photogrammetry, making it a fast, versatile, and practical solution for real-world applications.	Rupprecht, David Novotny		
1.36	LuSeg: Efficient Negative and Positive Obstacles Segmentation via Contrast-Driven Multi-Modal Feature Fusion on the Lunar	The paper "LuSeg: Efficient Negative and Positive Obstacles Segmentation" highlights the need for safe and autonomous lunar exploration. The authors developed LESS, a lunar surface simulation platform, and introduced LunarSeg, an RGB-D dataset containing positive obstacles (e.g., lunar rocks) and negative obstacles (e.g., craters). They proposed LuSeg, a two-stage segmentation network ensuring semantic consistency between the RGB and depth encoders via contrastive learning. Tested on LunarSeg and a real-world NPO dataset, LuSeg achieved state-of-the-art segmentation and a high inference speed of 57 Hz, making it an efficient and scalable solution.	By Shuaifeng Jiao, Zhiwen Zeng, Zhuoqun Su, Xieyuanli Chen, Zongtan Zhou, Huimin Lu		March 14, 2025
1.37	Mistral Small 3.1	Mistral Small 3.1 is an advanced, lightweight vision-language model designed for high-performance AI tasks, including text generation, multimodal understanding, and long-context processing with a 128k token window. With superior benchmarks across text, multilingual, and multimodal tasks, it surpasses comparable models like Gemma 3 and GPT-4o Mini while maintaining exceptional efficiency, running on consumer-grade hardware like an RTX 4090. Released under an Apache 2.0 license, it supports fine-tuning for specialized domains and real-world applications such as diagnostics, document verification, and virtual assistants, making it a versatile foundation for enterprise and consumer AI solutions.	By Mistral AI		March 17, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
1.38	Roblox Unveils Cube 3D: An Open-Source AI Model for 3D Object Generation	Roblox introduced Cube 3D, a foundational AI model capable of generating 3D objects and environments. This model was launched alongside the beta version of the Mesh Generation API, which allows users to create 3D objects using simple text commands. Cube 3D tokenizes 3D objects and predicts the next shape token, similar to language models. By releasing Cube 3D as open source, Roblox enables developers to fine-tune the model, create extensions, or train it with their own data. In the future, Cube 3D is expected to support image inputs and process multimodal data, including text, images, and videos.	By Anupam Singh, Vice President of Engineering, and Nick Tornow, Vice President of Creator Engineering		March 17, 2025
1.39	Cornstarch	Cornstarch is a distributed multimodal training framework designed to optimize the training of models that process multiple data types, such as text, images, and audio. It addresses inefficiencies in existing systems by being "multimodality-aware," ensuring balanced resource allocation and synchronization across modalities. This approach minimizes bottlenecks and maximizes throughput during training, making it particularly effective for large-scale multimodal tasks like vision-language models. By intelligently managing data flow and computation, Cornstarch enhances scalability and efficiency, enabling faster and more cost-effective training of complex multimodal systems	By Insu Jang et al.		March 17, 2025
1.40	NVIDIA Announces Isaac GROOT N1: The World's First Open Humanoid Robot Foundation Model	NVIDIA has introduced the open-source Isaac GROOT N1 model to accelerate humanoid robot development. Inspired by human cognition, the model features a dual-system architecture, where "System 1" manages fast and intuitive actions, while "System 2" handles deliberate decision-making. GROOT N1 can learn to grasp, transport objects, and perform multi-step tasks. The model can be	By Kristin Bryson		March 18, 2025





 Models					
#	Highlights	Summary	Author	Source	Date
		customized with real or synthetic data, offering flexible applications. NVIDIA is collaborating with Google DeepMind and Disney Research to develop Newton, an open-source physics engine based on the NVIDIA Warp framework, to enhance robotic simulations.			
1.41	Cosmos World Foundation Models	The Cosmos AI model series has been developed to generate cost-effective and highly realistic videos for robot training, enabling more efficient and scalable AI development. These models operate within NVIDIA's Omniverse platform, which the company describes as "the operating system for physical AI." By leveraging Omniverse, developers can condition Cosmos to create systematically controlled yet infinitely diverse "infinite environments", making it an ideal tool for training robotic systems in a wide range of scenarios. This approach allows AI models to learn and adapt in virtual spaces before being deployed in real-world applications, improving both accuracy and efficiency.	By Nvidia		March 18, 2025
1.42	LG AI Research Unveils 'EXAONE Deep' Reasoning AI Model	LG AI Research introduced ' EXAONE Deep ', an advanced reasoning AI model marking a significant milestone in AI innovation. EXAONE Deep-32B, with 32 billion parameters, demonstrated exceptional efficiency, achieving remarkable performance at just 5% of a very large competitor model's size. It excelled in solving complex mathematical and scientific problems, recording the highest score of 94.5 in the 2025 CSAT mathematics area and 95.7 points in MATH-500. Additionally, it achieved top grades in physics, chemistry, and biology evaluations. This open-source release underscores Korea's competitiveness in the global AI market.	By PR Newswire		March 18, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
1.43	Razer Introduces Wyrn Developer Platform with AI Tools	Razer unveiled Wyrn, a new developer platform featuring automated AI tools aimed at enhancing game development processes. The platform includes the AI QA Copilot, a cloud-based plugin compatible with engines like Unreal and Unity, designed to streamline quality assurance testing by identifying bugs and generating detailed reports. Razer claims this tool can enhance bug detection by up to 25% and reduce QA time by 50%, potentially impacting QA team sizes. Additionally, Wyrn offers the AI Gamer Copilot, an AI voice assistant providing real-time gaming tips and strategies, marking a significant advancement in integrating AI into gaming experiences.	By Antonio G. Di Benedetto		March 19, 2025
1.44	Llama-3.3-Nemotron-Super-49B-v1	Llama-3.3-Nemotron-Super-49B-v1 is a reasoning-optimized LLM derived from Meta Llama-3.3-70B-Instruct. It is post-trained for reasoning, chat preferences, retrieval-augmented generation (RAG), and tool use, supporting a 128K token context. Designed for efficiency and accuracy, it balances high performance with cost savings. A Neural Architecture Search (NAS) approach reduces memory usage, enabling deployment on a single H200 GPU with high throughput. This optimization allows larger workloads while maintaining strong reasoning capabilities. The NAS method fine-tunes the accuracy-efficiency tradeoff, making it a cost-effective, high-performance alternative to larger models, ideal for various AI-driven applications.	By Nvidia		March 18, 2025
1.45	RF-DETR	RF-DETR is a state-of-the-art object detection model that enhances the DETection TRansformer (DETR) architecture with advanced optimizations for better performance and efficiency. It tackles challenges like overlapping objects and complex scene accuracy through refined feature extraction and improved training techniques. Ideal for precise object localization and classification, RF-DETR excels	By Peter Robicieux et al.		March 20, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		in applications such as autonomous driving, surveillance, and robotics. By integrating cutting-edge transformer-based advancements, it provides a robust, scalable solution for modern object detection tasks, balancing accuracy and efficiency effectively.			
1.46	SpatialLM-Llama-1B: Advancing 3D Scene Understanding with Language Models	Manycore Research has introduced SpatialLM-Llama-1B , a 3D large language model designed to process 3D point cloud data and generate structured 3D scene understanding outputs, including architectural elements like walls, doors, windows, and oriented object bounding boxes with their semantic categories. Unlike previous methods that require specialized equipment for data collection, SpatialLM can handle point clouds from diverse sources such as monocular video sequences, RGBD images, and LiDAR sensors. This multimodal architecture bridges the gap between unstructured 3D geometric data and structured 3D representations, enhancing spatial reasoning capabilities for applications in embodied robotics, autonomous navigation, and complex 3D scene analysis tasks.	By Manycore Research		March 20, 2025
1.47	Orpheus TTS: Advancements in Text-to-Speech Synthesis	Canopy Labs has released Orpheus TTS , a state-of-the-art text-to-speech model based on Llama architecture, fine-tuned to deliver human-like speech synthesis with natural intonation and emotion. Key features include zero-shot voice cloning, guided emotion and intonation control through simple tags, and low-latency streaming suitable for real-time applications. The model is available in 3 billion parameter configurations, with quantized versions for efficient deployment. Developers can access the model and its codebase via Hugging Face and GitHub.	By CanopyLabs		March 18, 2025
1.48	Aardvark Weather	A new AI-driven forecasting model developed collaboratively by the University of Cambridge, the Alan Turing Institute, Microsoft Research, and the European	By Ben Geman,		March 24, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		Centre for Medium-Range Weather Forecasts, represents a significant advancement in meteorology. Unlike previous AI models that rely on traditional numerical weather prediction systems, Aardvark utilizes a purely machine-learning approach, processing observations from satellites, weather stations, and other sensors to produce high-resolution global and local forecasts. Notably, it operates efficiently on desktop computers, delivering results within minutes and outperforming the U.S. Global Forecast System in certain variables. This innovation holds promise for tailored applications in sectors like renewable energy and agriculture.	Andrew Freedman		
1.49	First Mamba-Powered Ultra-Large Model	The Hunyuan-T1 model, built on the TurboS Hybrid-Transformer-Mamba MoE architecture, represents a significant advancement in reasoning capabilities. With optimized long-text processing and efficient computation, it achieves double the decoding speed of its predecessor while excelling in mathematics, logic, and scientific reasoning. Trained with 96.7% reinforcement learning focus and curriculum learning strategies, T1 demonstrates superior alignment with human preferences. It achieves top-tier performance across benchmarks like MMLU-Pro (87.2) and MATH-500 (96.2), rivaling industry-leading models. T1's innovative architecture and training make it a cutting-edge solution for reasoning and alignment tasks.	By Tencent Hunyuan Research		March 21, 2025
1.50	AMD's Gaia Project	AMD launched Gaia, an open-source initiative aimed at enabling efficient local execution of large language models (LLMs) on standard PCs. Gaia optimizes hardware resources, particularly focusing on AMD's Ryzen CPUs and Radeon GPUs, enhancing performance even on less powerful consumer systems. This approach allows users to run sophisticated AI workloads without cloud	By Aaron Klotz		March 21, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		dependencies, significantly improving privacy and accessibility. AMD's Gaia directly challenges Nvidia's dominance by providing viable alternatives for decentralized AI processing, fostering greater hardware democratization in the rapidly expanding AI computing landscape.			
1.51	DeepSeek V3 Update	DeepSeek unveiled version 3 (V3) of its AI model, featuring significant improvements in performance and efficiency alongside a new licensing structure. V3 emphasizes enhanced multilingual capabilities, reasoning accuracy, and computational efficiency. Importantly, the updated license now allows commercial use with clearer terms, aiming to attract a broader range of developers and enterprises. This strategic shift positions DeepSeek to better compete in the commercial AI market, appealing especially to businesses seeking robust AI tools with transparent usage rights. Such licensing adjustments are increasingly critical in the competitive AI ecosystem dominated by both closed and open-source alternatives.	By Ambar Warrick		March 25,2025
1.52	Beijing-Backed AI Startup Manus Aims to Be China's Next DeepSeek	China is backing AI startup Manus, seen as a potential rival to domestic leaders like DeepSeek and Baidu. Supported by Beijing's strategic funding and talent pipelines, Manus is focused on developing cutting-edge large language models tailored to Chinese linguistic and regulatory needs. The move reflects China's ambition to cultivate a diverse and competitive AI ecosystem, reducing reliance on U.S. technologies. Manus' emergence is part of a broader national effort to scale homegrown AI innovation and secure leadership in foundational model development.	By Josh Ye		March 21, 2025





 Models					
#	Highlights	Summary	Author	Source	Date
1.53	Qwen2.5-VL-32B: Smarter and Lighter	Alibaba's Qwen team has released Qwen-VL-Max , a new 32B vision-language model with enhanced capabilities in document understanding, complex reasoning, chart interpretation, and OCR . It supports multi-image and multi-round dialogues, surpassing Gemini Pro 1.5 and GPT-4V in multiple benchmarks like MathVista and AI2D. Qwen-VL-Max also introduces a lightweight API that simplifies region selection, coordinates extraction, and chart data reading. With strong performance on benchmarks such as ChartQA, DocVQA, and MMMU, it showcases its superiority across vision-language tasks. The model, weights, and demo are now open source , promoting transparency and broader accessibility for researchers and developers.	By Qwen Team		March 24, 2025
1.54	Cosmos-Reason1: From Physical Common Sense to Embodied Reasoning	Cosmos-Reason1 explores how AI can develop embodied reasoning by integrating physical common sense with multimodal learning. The framework enhances an agent's ability to predict and reason about physical interactions in the world, leveraging text, vision, and motion cues. The model is trained on large-scale datasets with physically grounded scenarios, improving zero-shot and few-shot reasoning capabilities. Experimental results show advancements in tasks like object dynamics prediction, intuitive physics, and embodied decision-making. This work contributes to bridging the gap between common sense understanding and real-world physical reasoning in AI systems.	By Nvidia		March 18, 2025
1.55	Nemotron-H: A Family of Accurate, Efficient	The Nemotron-H family introduces advanced hybrid Mamba-Transformer models, offering state-of-the-art accuracy with up to 3x faster inference compared to similar-sized Transformer models. Spanning 8B to 56B parameters, these models	By NVIDIA ADLR		March 21, 2025





 Models					
#	Highlights	Summary	Author	Source	Date
	Hybrid Mamba-Transformer Models	leverage FP8 precision for efficient training and support long-context inference, achieving breakthroughs in reasoning, math, and vision-language tasks. The 56B model, trained on 20 trillion tokens, demonstrates competitive performance against much larger models, while the 47B variant enables efficient deployment on commodity GPUs. Nemotron-H represents a significant leap in LLM scalability, accuracy, and efficiency, with applications in AI reasoning, code generation, and physical AI systems.			
1.56	Cosmos-Transfer1: Conditional World Generation with Adaptive Multimodal Control	Cosmos-Transfer1 proposes a framework for generating conditional world representations using multimodal data, enabling adaptive control over generated environments. The model integrates textual, visual, and spatial inputs to create coherent world representations suitable for simulations and robotics. The authors introduce a novel transfer learning strategy that adapts knowledge from diverse datasets, improving generalization across tasks. The results demonstrate enhanced fidelity and consistency in generated environments. This approach is particularly relevant for AI-driven simulations, virtual reality, and autonomous system training, pushing the boundaries of controllable and context-aware world generation.	By Nvidia		March 18, 2025
1.57	Tokenize Image as a Set	TokenSet, a groundbreaking paradigm for image generation that tokenizes images into unordered sets instead of fixed-position sequences, enabling dynamic coding allocation based on regional semantic complexity. To model these sets, the authors propose a dual transformation mechanism that converts sets into fixed-length integer sequences while retaining summation constraints. They further develop Fixed-Sum Discrete Diffusion, a novel framework tailored for discrete data modeling with fixed-sum properties. Experiments demonstrate TokenSet's	By University of Science and Technology of China, Tencent Hunyuan Research		March 20, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		superior robustness, semantic-aware representation, and generation quality, marking a significant advancement over traditional sequential tokenization methods in visual generation tasks.			
1.58	Lyra: An Efficient and Expressive Subquadratic Architecture for Modeling Biological Sequences	Lyra, a novel subquadratic architecture for biological sequence modeling that integrates Projected Gated Convolutions (PGCs) and diagonalized State Space Models (S4D). Designed to efficiently model epistasis (context-dependent interactions), Lyra achieves state-of-the-art (SOTA) performance across over 100 biological tasks, including protein fitness prediction, RNA function analysis, and CRISPR guide design. With up to 120,000-fold fewer parameters and 64x faster inference than Transformers, Lyra democratizes access to advanced sequence modeling. Its efficiency and scalability make it transformative for computational biology, enabling breakthroughs in therapeutic design, diagnostics, and gene editing.	By Krithik Ramesh, Sameed M. Siddiqui, Albert Gu, Michael D. Mitzenmacher, Pardis C. Sabeti		March 20, 2025
1.59	Gemini 2.5: Most intelligent AI model	Google DeepMind announced major improvements to its Gemini models, focusing on advanced reasoning and better multi-step problem-solving. Gemini 1.5 Pro now excels in tasks involving deep thinking, such as coding, logic, and complex math. A new feature, "Notebook Mode," lets users guide the model step-by-step for more accurate and transparent outputs. The model also shows stronger performance on benchmarks like MATH and GPQA. These updates push Gemini closer to human-like reasoning, aiming to make AI more helpful for research, education, and problem-solving in real-world scenarios. Availability is expanding through Google products.	By Google		March 25, 2025



Models					
#	Highlights	Summary	Author	Source	Date
1.60	CoMP: Continual Multimodal Pre-training for Vision Foundation Models	CoMP, a multimodal pre-training pipeline that enhances pre-trained Vision Foundation Models (VFMs) by aligning their visual outputs with language representations. CoMP enables VFMs to handle visual inputs of various sizes through Continual Rotary Position Embedding and improves vision-language alignment using an Alignment Loss guided by language prototypes. Through a three-stage training process, VFMs show strong gains in both multimodal tasks and standard benchmarks. CoMP-SigLIP, paired with a 0.5B LLM, achieves 66.7 on ChartQA, 75.9 on DocVQA, 87.4% on ImageNet-1K, and 49.5 mIoU on ADE20K under frozen chunk evaluation.	By Yitong Chen et al.		March 24, 2025
1.61	Ai2 Paper Finder: Enhancing AI Literature Search with LLM Technology	The Allen Institute for AI (AI2) has introduced Ai2 Paper Finder, an LLM-powered literature search system designed to emulate the iterative process researchers use to locate relevant papers. Unlike traditional search tools, Ai2 Paper Finder allows users to input complex queries in natural language, facilitating the discovery of hard-to-find papers. The system intelligently breaks down queries, searches for pertinent papers, follows citations, evaluates relevance, and presents results with concise summaries. This approach streamlines the research process, saving time and uncovering leads that might be missed with conventional methods.	By Ai2 Team		March 26, 2025
1.62	HoarePrompt: Enhancing Program Correctness Verification with Structural Reasoning	Researchers have introduced HoarePrompt, a novel approach that integrates program analysis techniques with large language models (LLMs) to assess program correctness against natural language requirements. HoarePrompt employs a step-by-step process where an LLM generates natural language descriptions of program states at various code points, inspired by the strongest postcondition calculus. To handle loops effectively, the method incorporates few-shot-driven k-	By Dimitrios Stamatios Bouras et al.		March 25, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		induction, adapting a model checking technique. Evaluations using the CoCoClanL dataset demonstrate that HoarePrompt significantly outperforms existing methods in classifying program correctness.			
1.63	Qwen2.5-Omni-7B: A Unified Multimodal AI Model	Qwen2.5-Omni-7B is an advanced multimodal AI model developed by Qwen, designed to process and generate text, images, audio, and video. It features the Thinker-Talker architecture, enabling simultaneous text and speech generation, and employs TMRoPE position embedding to synchronize video and audio inputs. Despite its compact 7 billion parameters, it delivers exceptional performance across various modalities, outperforming similarly sized single-modality models. Qwen2.5-Omni-7B excels in real-time interactions, making it suitable for applications like intelligent voice assistants and multimedia content creation. The model is available under the Apache-2.0 license on Hugging Face.	By Qwen Development Team		March 26, 2025
1.64	Open Deep Search: Advancing Open-Source Search with Reasoning Agents	Open Deep Search (ODS) is an open-source framework that enhances large language models (LLMs) with advanced reasoning capabilities through web search tools. It comprises two main components: the Open Search Tool, a web search utility that outperforms proprietary counterparts, and the Open Reasoning Agent, which interprets tasks and orchestrates actions, including utilizing the Open Search Tool. When combined with the DeepSeek-R1 LLM, ODS achieves state-of-the-art performance, surpassing existing benchmarks like SimpleQA and FRAMES. This development narrows the gap between proprietary and open-source search AI solutions, promoting broader access and innovation in the field.	By Salaheddin Alzubi et al.		March 26, 2025



 Models					
#	Highlights	Summary	Author	Source	Date
1.65	PS3: Scaling Vision Pre-Training to 4K Resolution	<p>Researchers have introduced PS3, a method that scales CLIP-style vision pre-training to 4K resolution with minimal computational overhead. By focusing on local regions and contrasting them with detailed captions, PS3 reduces pre-training costs by 79% compared to global contrastive learning approaches. The approach involves collecting 75 million high-resolution images and curating 282 million image-caption pairs. Integrating PS3 into models like VILA-HD enhances performance on various benchmarks, demonstrating the effectiveness of high-resolution pre-training in visual learning tasks. Notably, VILA-HD surpasses previous models, achieving a 14.5% improvement over GPT-4o on the 4KPro benchmark.</p>	By Nvidia		March 25, 2025
1.66	DeepSeek's AI Model Accelerates China's Technological Progress	<p>Chinese AI startup DeepSeek has significantly advanced China's AI capabilities, narrowing the development gap with the United States to just three months in certain areas. This progress is attributed to DeepSeek's efficient use of chips and algorithms, challenging previous assumptions about China's position in AI technology. The company's innovative approach has spurred other Chinese startups to revise their business models, focusing on application development and enterprise solutions. DeepSeek's emergence has prompted both domestic and international attention, influencing market dynamics and prompting discussions about the future of AI development globally.</p>	By Liam Mo and Kane Wu		March 25, 2025
1.67	DeepSeek Releases Enhanced AI Model, Intensifying Competition with OpenAI	<p>Chinese AI startup DeepSeek has unveiled a significant upgrade to its V3 large language model, DeepSeek-V3-0324, available through the AI development platform Hugging Face. This new model demonstrates substantial improvements in reasoning and coding capabilities compared to its predecessor, with benchmark tests indicating enhanced performance across multiple technical metrics.</p>	By Liam Mo and Brenda Goh		March 25, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		DeepSeek's rapid emergence as a notable player in the global AI landscape is marked by its series of models that compete with Western counterparts while offering lower operational costs.			
1.68	Amazon Expands AI Initiatives with New Shopping and Health Tools	Amazon is testing two AI-driven tools to improve shopping and healthcare experiences. Interests AI provides personalized product recommendations based on natural language prompts, allowing users to receive tailored suggestions. Currently, it is available to a limited group of U.S. customers, with plans for wider access. Meanwhile, the Health AI Chatbot offers reliable health information and guidance, helping users find relevant products and care options. This chatbot is in beta testing on Amazon's website and app. These initiatives highlight Amazon's commitment to leveraging advanced AI to enhance user interactions across different domains.	By Annie Palmer		March 25, 2025
1.69	QVQ-Max: Advancing Visual Reasoning Capabilities	Qwen has introduced QVQ-Max, an advanced visual reasoning model capable of interpreting and analyzing images and videos to provide solutions across various domains, including mathematics, programming, and creative tasks. Building upon the previous QVQ-72B-Preview, QVQ-Max demonstrates enhanced capabilities in detailed observation, deep reasoning, and flexible application. For instance, when analyzing complex mathematical problems, QVQ-Max adjusts its reasoning process to improve accuracy, showcasing its potential in solving intricate tasks. This release signifies a substantial step forward in visual reasoning, offering users a versatile tool for both analytical and creative endeavors.	By Qwen Team		March 28, 2025
1.70	Nomic Embed Code: Advancing Code	Nomic has introduced Nomic Embed Code, a 7-billion-parameter code embedding model designed for enhanced code retrieval across multiple languages, including	By Nomic Team		March 26, 2024



 Models					
#	Highlights	Summary	Author	Source	Date
	Retrieval with Open-Source Innovation	Python, Java, Ruby, PHP, JavaScript, and Go. It surpasses models like Voyage Code 3 and OpenAI Embed 3 Large, scoring 81.6% (Python) and 93.8% (Go) on CodeSearchNet. Trained on the CoRNStack dataset, it leverages dual-consistency filtering and curriculum-based hard negative mining. Upholding open-source principles, Nomic has released the model weights and training code under Apache-2.0, promoting transparency and collaboration in developer communities.			
1.71	Introducing 4o Image Generation	OpenAI has introduced image generation capabilities to GPT-4o, enabling users to create visuals directly from text prompts. This new feature, launched in March 2025, enhances GPT-4o's multimodal abilities, making it more versatile for creative tasks. Users can now generate detailed images by simply describing what they want, and GPT-4o produces results in various artistic styles. The tool is designed to be user-friendly, with accessible integration into ChatGPT. OpenAI emphasizes responsible usage and has implemented safeguards to prevent misuse, including content filtering and provenance indicators for generated images.	By OpenAI		March 25, 2025
1.72	JudgeLRM: Large Reasoning Models as a Judge	The paper introduces JudgeLRM, a family of judgment-oriented large reasoning models (LRMs) trained using reinforcement learning (RL) with judge-specific, outcome-driven rewards to address limitations in supervised fine-tuning (SFT) for judgment tasks requiring complex reasoning. JudgeLRM models significantly outperform state-of-the-art models like GPT-4 and DeepSeek-R1, with JudgeLRM-7B achieving a 2.79% improvement in F1 score. The study highlights that judgment	By Nuo Chen , Zhiyuan Hu, Qingyun Zou, Jiaying Wu, Qian Wang, Bryan		March 31, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		tasks are inherently reasoning-intensive and demonstrates the effectiveness of RL in enhancing evaluative reasoning. This work sets a foundation for developing reliable, reasoning-focused LLM judges, advancing scalable and accurate AI evaluation systems.	Hooi, Bingsheng He		
1.73	Open-Reasoner-Zero: An Open Source Approach to Scaling Up Reinforcement Learning on the Base Model	Open-Reasoner-Zero introduces the first open-source implementation of large-scale reasoning-oriented reinforcement learning (RL) training on large language models (LLMs), emphasizing scalability, simplicity, and accessibility. Using vanilla PPO with GAE ($\lambda=1, \gamma=1$) and a minimalist rule-based reward function, it achieves superior performance on benchmarks like AIME2024, MATH500, and GPQA Diamond, while requiring only 1/10th the training steps compared to previous pipelines like DeepSeek-R1-Zero. By democratizing advanced RL techniques and releasing comprehensive resources, Open-Reasoner-Zero empowers researchers to explore scalable reasoning capabilities, marking a key advancement in RL-driven AI development.	By StepFun, Tsinghua University		March 31, 2025




AI Chips					
#	Highlights	Summary	Author	Source	Date
2.1	World's First Ternary Carbon Microchip	Chinese scientists unveiled a pioneering carbon-based AI microchip, using ternary logic, not binary. Built with carbon nanotubes (CNTs), it processes data in three states, boosting speed and reducing energy use. Peking University researchers showed perfect handwriting recognition accuracy, demonstrating AI potential. This CNT approach advances semiconductor tech, addressing silicon limitations. While CNT integration density trails silicon (e.g., NVIDIA GPUs), this is a significant step toward next-gen AI hardware. It underscores China's ambition in post-silicon chip research, potentially transforming AI processing with efficient, high-performance chips.	By Tribune Team		March 08, 2025
2.2	Nvidia AI Chips: Cross-Border Scandal	A cross-border probe investigates illegal routing of Nvidia AI chips. Singapore charged three men for fraud, involving servers with U.S.-made chips, allegedly destined for Chinese AI startup DeepSeek, violating export restrictions. Dell and SuperMicro supplied the servers, shipped via Malaysia, potentially to China. An anonymous tip prompted the investigation. Singapore and the U.S. are jointly probing if export-controlled components were involved. The U.S. also investigates DeepSeek's use of restricted chips. This scandal highlights geopolitical tensions over AI hardware, as demand for AI accelerators surges. It may lead to stricter export controls to prevent illicit chip transfers.	By Bing Hong Lok		March 4, 2025
2.3	Singapore Investigates Nvidia Chip Diversion	Singapore's Home Affairs Minister, K. Shanmugam, confirmed that Dell and SuperMicro supplied the servers under investigation, which were allegedly shipped to Malaysia as part of a circuitous route leading to China. The investigation was sparked by an anonymous tip, with authorities analyzing whether the servers contained export-controlled components. They are coordinating with U.S. counterparts for a joint investigation into potential	By Gao Yuan and Mackenzie Hawkins		March 8, 2025



AI Chips					
#	Highlights	Summary	Author	Source	Date
		sanctions violations. This development could have significant impacts on the global technology supply chain, leading to stricter compliance requirements for hardware manufacturers and distributors to prevent the diversion of restricted components through intermediary countries.			
2.4	Apple's M3 Ultra Runs DeepSeek R1	Apple's M3 Ultra is proving to be a competitive AI processing unit, successfully running DeepSeek R1, a 671 billion parameter model. Using 448GB unified memory, it delivers high-bandwidth AI performance under 200W power consumption, eliminating the need for multi-GPU setups. EXO Labs even ran a distributed 8-bit version across two M3 Ultra Mac Studios, achieving 11 transactions per second (t/s). However, discussions continue about Apple Silicon's practicality for local inference, as prompt processing speeds remain slow, with some favoring NVIDIA and AMD alternatives.	By Ali Salman		March 12, 2025
2.5	Celestial AI Secures \$250M for Photonic AI Chip Links	Silicon Valley startup Celestial AI has raised \$250 million to develop photonics-based interconnect technology for AI chips, aiming to replace traditional electrical data transfer with light-based communication. This innovation could enhance AI chip-to-memory speed and efficiency, competing with Nvidia's NVLink. The company has now secured over \$300 million in funding to advance data center and AI computing solutions. Photonic interconnects reduce energy consumption and improve AI model performance, making them a promising alternative to existing chip link technologies. This investment highlights growing interest in alternative AI chip architectures for faster and more efficient AI processing.	By Abhinaya Prabhu		March 11, 2025



AI Chips					
#	Highlights	Summary	Author	Source	Date
2.6	Meta Tests Its First In-House AI Training Chip, MTIA	Meta has started testing its first in-house AI training chip, the MTIA (Meta Training and Inference Accelerator), designed to enhance its AI capabilities while reducing reliance on third-party hardware like Nvidia GPUs. The chip aims to boost efficiency in AI model training and inference, supporting Meta's AI-driven services, including recommendation systems and content moderation. MTIA is part of Meta's broader strategy to build a custom AI hardware ecosystem, complementing its AI-powered infrastructure. This move aligns with industry trends, as tech giants seek greater control over AI processing power to meet growing computational demands.	By Katie Paul and Krystal Hu		March 11, 2025
2.7	Fudan University Develops High-Speed Photonic Chip for AI and Data Centers	A research team from Fudan University has developed a silicon photonic integrated high-order mode multiplexer chip, enabling ultra-high-capacity on-chip optical data transmission. This breakthrough enhances optical interconnection in data centers and high-performance computing, benefiting AI, large-scale computing, and model training. The chip supports 38 terabits per second (Tbps), transferring 4.75 trillion model parameters per second, significantly improving efficiency in AI training and GPU computing. Expert Ma Jihua noted that photonic chips are gaining traction over electronic chips, boosting bandwidth and transmission speeds. This innovation could revolutionize AI training and data processing within three to five years.	By Global Times		March 13, 2025
2.8	SoftBank and OpenAI Collaborate on Major	SoftBank plans to transform a former Sharp LCD panel plant in Osaka into a data center dedicated to AI operations, in partnership with OpenAI. The facility,	By Reuters		March 13, 2025




AI Chips					
#	Highlights	Summary	Author	Source	Date
	AI Data Center in Japan	expected to commence operations in 2026, will be one of Japan's largest, featuring a power capacity of 150 megawatts. This collaboration aims to commercialize OpenAI's AI agent model in Japan, training it on client companies' data to offer customized AI solutions. The total investment could reach up to 1 trillion yen (approximately \$6.77 billion), underscoring the scale and ambition of this project.			
2.9	Nvidia GTC 2025: Blackwell Ultra & AI Market Shifts	Nvidia's annual GPU Technology Conference (GTC) kicked off March 17 in San Jose, California, with industry attention focused on CEO Jensen Huang's upcoming keynote on March 18. The five-day event is expected to feature major announcements regarding Nvidia's Blackwell Ultra chips, projected for release in the first half of 2025. Industry analysts are watching closely following concerns about AI demand peaking after DeepSeek's recent model launch triggered Nvidia's unprecedented \$600 billion single-day market value drop. Huang reportedly aims to diversify Nvidia's offerings beyond chips to establish a more secure foundation	By PYMNTS		March 14, 2025
2.10	NVIDIA Blackwell Ultra	NVIDIA has introduced the Blackwell Ultra platform, an advanced evolution of its Blackwell GPU architecture, tailored for AI applications. A major upgrade includes expanding High Bandwidth Memory (HBM) from 192GB to 288GB in the Blackwell Ultra GB300, coupled with a 50% performance increase over the B200 series. The GB300 AI chip is available in dual-chip (CoWoS-L) and single-chip (CoWoS-S) configurations. Trial production is set for Q2 2025, with mass production in Q3. These advancements highlight NVIDIA's commitment to scaling AI hardware performance to support increasingly complex AI workloads.	By Nvidia		March 18, 2025



AI Chips					
#	Highlights	Summary	Author	Source	Date
2.11	NVIDIA Rubin Ultra GPU and NVIDIA Vera CPU	NVIDIA CEO Jensen Huang honored astronomer Vera Rubin by unveiling a roadmap for advancing data center performance. The plan features next-generation NVIDIA Rubin Ultra GPU and NVIDIA Vera CPU architectures, introducing major innovations. Huang stated, "basically everything is brand new except for the chassis." Systems built on the Rubin Ultra platform, including the Vera Rubin NVL 144, will launch in late 2026, with additional Rubin Ultra-based systems arriving in 2027. Huang emphasized that Rubin will significantly reduce costs, reinforcing NVIDIA's push for efficient AI computing.	By Nvidia		March 18, 2025
2.12	RTX PRO 6000 Blackwell	NVIDIA has unveiled the RTX PRO 6000 Blackwell series, designed for professional users, offering up to 96GB of GDDR7 memory and fifth-generation PCIe support. These GPUs are engineered to enhance performance in compute-intensive tasks, such as AI inference, data analytics, and content creation. Jensen Huang highlighted the series' innovative capabilities, emphasizing its role in advancing both AI applications and professional graphics processing. With cutting-edge architecture, the RTX PRO 6000 Blackwell delivers exceptional efficiency and power, making it an ideal choice for demanding workloads across industries like scientific computing, media production, and AI-driven research.	By Nvidia		March 18, 2025
2.13	Micron Stock Rises as AI Memory Boosts Revenue	Micron Technology saw a strong financial performance in Q1 2025, with revenue rising 38% year-over-year to \$8.05 billion, driven by high demand for its high-bandwidth memory (HBM) chips used in AI servers. The company's net income nearly doubled compared to the previous year. CEO Sanjay Mehrotra highlighted the robust growth in data center revenues and projected sustained demand for AI memory through at least 2026. Micron expects another record quarter ahead, further boosting investor confidence. As a result, the company's stock gained 2%	By Mike Wheatley		March 20, 2025

AI Chips					
#	Highlights	Summary	Author	Source	Date
		in after-hours trading, signaling optimism around its positioning in the rapidly growing AI hardware market.			
2.14	Lisa Su Sets Sights on Nvidia as AMD Expands AI Chip Ambitions	Since becoming CEO of AMD in 2014, Lisa Su has led a major transformation of the company, enabling it to surpass Intel. Now, she faces a new challenge: competing with Nvidia, which currently holds about 90% of the AI chip market. Su aims to strengthen AMD's software capabilities to rival Nvidia's dominant CUDA platform. The company is investing in open-source software to support training and inference of large language models. Her leadership style—focused on consistency and precise execution—has earned respect from both partners and critics. This strategy is designed to position AMD more competitively in the AI space.	By Kif Leswing		March 20, 2025
2.15	AMD signs huge multi-billion dollar deal with Oracle to build a cluster of 30,000 MI355X AI accelerators	AMD has secured a multi-billion dollar deal with Oracle to build a massive AI supercomputer cluster powered by 30,000 MI355X accelerators. These next-generation GPUs, based on AMD's CDNA 4 architecture and TSMC's advanced 3nm process, are set to rival Nvidia's Blackwell chips. The partnership underscores Oracle's ambition to expand its cloud AI infrastructure and support growing enterprise demand for generative AI workloads. With a projected 35-fold increase in AI server GPU demand by 2027, this strategic move strengthens AMD's position in the AI hardware market and marks a significant leap forward in large-scale AI compute capabilities.	By Wayne Williams		March 21, 2025
2.16	European Lawmakers Advocate for Chips Act 2.0	European lawmakers are pressing the European Commission to rapidly advance the proposed Chips Act 2.0, designed to strengthen Europe's semiconductor independence amid global supply chain concerns. Driven by intensifying	By Foo Yun Chee		March 24.2025


AI Chips					
#	Highlights	Summary	Author	Source	Date
		geopolitical tensions and technological competition, the legislation emphasizes increased investment, advanced chip manufacturing, and research innovation. Lawmakers warn delays could jeopardize Europe's competitiveness in AI and emerging technologies. With semiconductor demand surging, the Chips Act 2.0 seeks to bolster Europe's strategic autonomy and reduce reliance on non-European suppliers, ensuring stability and resilience in crucial tech sectors and enhancing regional economic security.			
2.17	Micron Shares Dip Despite Strong AI Demand Due to Weak Margin Outlook	Micron Technology's shares fell after the company issued a weaker-than-expected margin forecast, overshadowing optimism about AI-driven demand for its memory chips. While Micron highlighted robust growth in high-bandwidth memory (HBM) used in AI workloads, investors were concerned about profitability pressures and elevated production costs. The company remains a key supplier for Nvidia and other AI players, but its forecast sparked worries about near-term financial performance. The stock decline reflects market sensitivity to both the opportunities and operational challenges facing AI-related semiconductor firms.	By Max A. Cherney,		March 21, 2025
2.18	Nvidia-Backed CoreWeave Targets Up to \$27 Billion Valuation in U.S. IPO	CoreWeave, a cloud computing startup backed by Nvidia, is planning a U.S. IPO that could value the company at up to \$27 billion. Specializing in GPU-accelerated infrastructure for AI workloads, CoreWeave has rapidly grown as demand for high-performance computing surges. The company provides scalable infrastructure for clients training large AI models, positioning itself as a key player in the AI chip ecosystem. This IPO underscores investor appetite for firms supporting the AI boom, particularly those aligned with industry leaders like Nvidia.	By Manya Saini, Niket Nishant		March 20, 2025



AI Chips					
#	Highlights	Summary	Author	Source	Date
2.19	NVBleed: Covert and Side-Channel Attacks on NVIDIA Multi-GPU Interconnect	NVBleed exposes covert and side-channel vulnerabilities in NVIDIA’s multi-GPU NVLink interconnect. Researchers demonstrate how malicious actors can exploit shared resources to leak sensitive data (e.g., encryption keys) across GPUs without explicit communication. The attack leverages contention in memory access patterns or power consumption to infer activity from a victim GPU, bypassing isolation guarantees. These flaws arise from hardware design choices prioritizing performance over security. Mitigations require architectural changes, as software patches may degrade performance. The work highlights risk in multi-GPU environments, urging a reevaluation of interconnect security for cloud and HPC systems relying on NVIDIA’s technology.	Bt Yicheng Zhang et al.		March 22, 2025
2.20	SK Hix Reports Order Acceleration Ahead of Potential U.S. Tariffs	SK Hix, the world's second-largest memory chipmaker, announced that customers have expedited orders in anticipation of proposed U.S. semiconductor tariffs. This "pull-in" effect, coupled with reduced customer inventories, has created favorable market conditions. However, it remains uncertain if this trend will continue. President Trump has proposed a 25% tariff on semiconductor imports, which could lead to higher product prices and reduced demand if implemented. SK Hix forecasts significant growth in high-bandwidth memory (HBM) chip demand in 2025, driven by data center investments. The company plans to finalize HBM chip sales for 2026 within the first half of this year to enhance revenue stability.	By Heekyong Yang and Hyunjoo Jin		March 27, 2025




AI Chips					
#	Highlights	Summary	Author	Source	Date
2.21	China Maintains Lead in Chipmaking Investments Amid Global Competition	China is projected to lead global investments in chipmaking equipment in 2025, allocating \$38 billion despite a 24% decrease from the previous year, according to SEMI. This investment is part of China's strategy to reduce dependence on imported chips and counter U.S. restrictions. Following China, Taiwan and South Korea are expected to invest \$21 billion and \$21.5 billion, respectively. Overall, global investments in chipmaking equipment are anticipated to grow by 2% to \$110 billion in 2025, driven by demand for tools to manufacture AI chips, with further growth expected in 2026.	By Reuters		March 26, 2025
2.22	Schneider Electric Invests Over \$700 Million in U.S. to Support AI Growth	Schneider Electric will invest \$700M+ in U.S. operations through 2027, marking its largest American investment in 135 years. The initiative aims to enhance energy infrastructure, support AI growth, boost domestic manufacturing, and strengthen energy security. Expansion plans include facility upgrades and new openings across Tennessee, Massachusetts, Texas, Missouri, Ohio, and the Carolinas, creating 1,000+ jobs. This move positions Schneider Electric to meet rising energy demands fueled by AI advancements while reinforcing its commitment to U.S. energy resilience and technological growth.	By Mrinalika Roy		March 25, 2025
2.23	Dell Technologies Partners with NVIDIA to Expand AI Server Business	Dell Technologies has significantly advanced its position in the AI server market, achieving an annual revenue of \$10 billion in this sector. The company projects a 50% growth in AI sales for 2025, underscoring its commitment to AI infrastructure. Dell's collaboration with NVIDIA has been pivotal, notably contributing to the development of AI supercomputing projects for clients like Elon Musk's xAI. This	By Michael J. Kramer		March 26, 2025



AI Chips					
#	Highlights	Summary	Author	Source	Date
		partnership positions Dell to effectively meet the escalating demand for AI solutions, reinforcing its status as a key player in the AI server industry.			
2.24	SK Hynix Reports Order Acceleration Ahead of Potential U.S. Tariffs	SK Hynix, the world's second-largest memory chipmaker, reports that customers are expediting orders ahead of proposed U.S. semiconductor tariffs, creating favorable market conditions amid lower inventories. However, the sustainability of this trend remains uncertain. A proposed 25% tariff by President Trump could raise prices and reduce demand if enacted. Despite this, SK Hynix anticipates strong growth in high-bandwidth memory (HBM) chips in 2025, driven by data center investments, and aims to finalize 2026 HBM sales within H1 2025 to ensure revenue stability.	By Heekyong Yang and Hyunjoo Jin		March 27, 2025
2.25	Broadcom Lowers Power Consumption With Latest AI Networking Chips	Broadcom has introduced two advanced AI networking chips, Sian3 and Sian2M, designed to boost data transmission efficiency in AI and machine learning clusters. Built on 3-nanometer technology, the Sian3 chip reduces power consumption by over 20% for 1.6T optical modules versus its predecessor. The Sian2M targets 800G and 1.6T short-reach multi-mode fiber links, featuring integrated VCSEL drivers and Broadcom's 200G VCSEL technology to enhance speed and energy efficiency. These chips meet the rising demands of AI workloads by enabling scalable, power-efficient connectivity. Broadcom is sampling now, with Sian3 mass production set for the third quarter.	By PATRICK SEITZ		March 25, 2025




AI Chips					
#	Highlights	Summary	Author	Source	Date
2.26	Cerebras Systems and Ranovus Win \$45 Million U.S. Military Deal	Cerebras Systems and Ranovus have won a \$45 million contract from the U.S. military to accelerate the development of high-performance chips tailored for AI applications and computing tasks. The deal focuses on improving the speed, efficiency, and scalability of chip technology critical to military operations, particularly for AI-driven projects. By combining Cerebras' advanced wafer-scale chips with Ranovus' photonic interconnect solutions, the collaboration aims to enhance the processing capabilities required for complex AI models, marking a key step in strengthening the U.S. military's technological infrastructure.	By Stephen Nellis		April 1, 2025



LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.1	SWE-Lancer Benchmark Highlights	OpenAI's SWE-Lancer benchmark, announced in March, was developed to assess the real-world software engineering capabilities of advanced language models. Based on over 1,400 freelance tasks sourced from Upwork, the benchmark represents \$1 million worth of projects. Covering various domains such as coding and project management, these tasks are evaluated through end-to-end testing	By Daniel Dominguez		March 8, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		verified by professional engineers. Initial results indicate that LLMs struggle to complete complex, multi-step software projects effectively. SWE-Lancer serves as a crucial benchmark, highlighting the need for improvements in reasoning, long-term planning, and reliability, making it a key measure for future advancements in AI-driven engineering.			
3.2	Microsoft Develops In-House Reasoning LLMs (Project "MAI"):	Microsoft is developing its own large language models (LLMs) to reduce reliance on OpenAI. Led by Mustafa Suleyman, its AI division has trained the MAI model family, which performs competitively with OpenAI and Anthropic on benchmarks. These models prioritize chain-of-thought reasoning, enhancing multi-step decision-making. Microsoft is also testing external models from xAI, Meta, and DeepSeek as potential OpenAI alternatives for 365 Copilot. The MAI models, more advanced than the Phi series, are being trialed as GPT-4 replacements internally and may launch via API this year. This move signals Microsoft's ambition to become a leading AI provider beyond OpenAI.	By Reuters		March 7, 2025
3.3	START(Self-taught Reasoner with Tools):	In March 2025, researchers introduced START (Self-Taught Reasoner with Tools), an advanced LLM designed for complex reasoning by integrating external tools, especially Python execution. While traditional LLMs like OpenAI-o1 and DeepSeek-R1 excel in long chain-of-thought (CoT) reasoning, they often struggle with hallucinations and inefficiencies. Built by fine-tuning QwQ-32B-Preview, START achieved high accuracy on science QA (63.6%), math (95.0%, 66.7%), and	By University of Science and Technology of China and Alibaba Group		March 7, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		coding benchmarks (47.1%, 47.3%), surpassing its base model and competing with state-of-the-art models like R1-Distill-Qwen-32B and o1-Preview.			
3.4	HoT (Highlighted Chain of Thought):	In March 2025, researchers introduced Highlighted Chain-of-Thought Prompting (HoT), a novel technique aimed at improving the factual accuracy and transparency of Large Language Models (LLMs). HoT prompts LLMs to generate responses that include XML tags linking facts directly to those in the query. This method involves reformatting the input question by highlighting key facts with XML tags and then generating a response that includes similar tags. This approach allows users to trace statements in the answer back to the original input, improving the traceability and reliability of the response.	By Tin Nguyen Logan, Bolton Mohammad, Reza Taesiri, Anh Totti Nguyen		March 5, 2025
3.5	Industry-Scale LLMs	Major tech firms debuted new large language models. Notably, Foxconn's research arm unveiled FoxBrain, Taiwan's first LLM with advanced reasoning optimized for Traditional Chinese. FoxBrain is built on Meta's Llama 3.1 architecture and was trained on 120 Nvidia H100 GPUs in about four weeks. The model is being applied to manufacturing and supply-chain tasks, and Foxconn plans to open-source it for collaboration. Its performance is close to state-of-the-art, with only a slight gap versus a distilled version of China's leading model (DeepSeek). This showcases a trend of organizations customizing LLMs for specific languages and domains.	By Reuters		March 10, 2025
3.6	Multimodal AI Systems:	On March 2–3, researchers introduced TaxaBind, a groundbreaking multimodal LLM that integrates six data modalities—including ground-level photos, satellite imagery, audio, and text—into a unified model. Using an innovative technique called "multimodal patching", TaxaBind merges features from each modality into	By Shawn Ballard		March 3, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		a common representation, enabling cross-modal reasoning. This allows zero-shot learning for tasks like species identification from images (even for unseen species) and cross-modal retrieval (e.g., linking animal photos to climate data). TaxaBind outperforms traditional single-modality models, demonstrating how LLMs are evolving to incorporate vision, audio, and other inputs for a deeper understanding in AI.			
3.7	New Benchmarks & Leaderboards:	New benchmarks and leaderboards emerged to assess LLM performance on advanced tasks. GPQA Diamond now tests graduate-level reasoning with complex logic questions, while SWE Bench (Verified) evaluates coding skills on difficult programming challenges requiring correct, executable solutions. These benchmarks feed into leaderboards like Vellum AI's, ranking top models by task and updating as new results arrive. Early rankings show that different models excel in different areas—some leading in reasoning (GPQA), others in coding (SWE Bench). Traditional tests like MMLU and Big-Bench Hard remain key, while new domain-specific evaluations provide deeper insights into LLM capabilities.	By Vellum AI		March 6, 2025
3.8	Holistic Evaluation & Real-World Testing:	Static QA tests alone are no longer enough to evaluate modern LLMs. Researchers are developing interactive benchmarks that mimic real-world use, testing models in multi-turn environments where they must use tools, call APIs, or complete complex tasks. Berkeley Function Calling (BFCL) and SWE-Lancer simulate scenarios like writing and executing code, with automatic verification of correctness (e.g., does the code run?). These evaluations assess reasoning, planning, and tool use, moving beyond simple question-answering. As LLMs take on roles with ethical and financial stakes, ensuring reliability and multi-step	By Vishakha Agrawal, Archie Chaudhury, Shreya Agrawal		March 5, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		consistency is critical. "Agentic" evaluations are now complementing traditional benchmarks.			
3.9	Built-in Safety Reasoning	Ensuring LLMs behave safely is crucial, and a new approach, "Rational" (Reasoning-Enhanced Fine-Tuning for Interpretable LLM Safety), was introduced by Carnegie Mellon in March. Instead of relying on hard-coded filters, the model learns to generate a step-by-step safety rationale before answering. During fine-tuning, it analyzes prompts, considers intent, and explains why a query is harmful or safe. This method improves refusals to adversarial prompts, making them more context-aware and interpretable. By internalizing ethical reasoning, LLMs can avoid inappropriate content while reducing unnecessary refusals, marking a shift toward built-in, explainable safety mechanisms in AI models.	By Yuyou Zhang, Miao Li, William Han, Yihang Yao, Zhepeng Cen, Ding Zhao		March 6, 2025
3.10	Bias Mitigation via Activation Steering	In March 2025, researchers introduced a bias-mitigation method using Steering Vector Ensembles (SVE) to reduce social biases in LLMs without retraining. They applied steering vectors to model activations, using Bayesian optimization to adjust responses along nine bias axes (e.g., gender, race). These vectors, optimized on the BBQ benchmark, were combined into an ensemble, achieving up to 12% bias reduction in models like Mistral, Llama, and Qwen. The modular and interpretable approach allows dynamic adjustments without degrading performance. SVE demonstrates how activation engineering can offer an efficient, adaptable alternative to costly fine-tuning for fairness in AI systems.	By Zara Siddique, Irtaza Khalid, Liam D. Turner, Luis Espinosa-Anke		March 7, 2025
3.11	Chain-of-experts (CoE)	The Chain-of-Experts (CoE) framework enhances LLM efficiency and accuracy by activating specialized experts sequentially instead of all at once. This approach improves contextual understanding and optimizes resource usage, providing a	By Ben Dickson		March 10, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		cost-effective alternative to traditional dense models and Mixture-of-Experts (MoE) architectures. CoE offers improved reasoning, as experts collaborate progressively for better accuracy, while also reducing computation costs by lowering redundant processing and memory usage. It outperforms MoE models with 17.6% less memory usage while maintaining similar accuracy. CoE provides a scalable, efficient AI solution, making LLMs more accessible and sustainable.			
3.12	Sketch-of-Thought	Recent advances in LLMs have improved reasoning via Chain of Thought (CoT) prompting but at the cost of excessive verbosity. Sketch-of-Thought (SoT) is a new prompting framework that reduces token usage while maintaining accuracy. Inspired by cognitive science, SoT integrates Conceptual Chaining, Chunked Symbolism, and Expert Lexicons, selecting the best approach dynamically via a lightweight routing model. Tested across 15 reasoning datasets in multiple languages and modalities, SoT reduces tokens by 76% with minimal accuracy loss. In math and multi-hop reasoning, it even improves accuracy while being more efficient, making it a scalable solution for AI reasoning tasks.	By Simon A. Aytes , Jinheon Baek, Sung Ju Hwang		March 10, 2025
3.13	LLM-as-a-Judge: Evaluating AI with AI	The "Awesome LLM-as-a-Judge" initiative explores leveraging Large Language Models (LLMs) as evaluators across various domains. This approach aims to harness LLMs' capabilities to provide scalable and flexible assessments, potentially reducing reliance on traditional expert-driven evaluations. However, ensuring the reliability of LLM-as-a-Judge systems presents challenges such as maintaining consistency, mitigating biases, and adapting to diverse assessment scenarios. To address these issues, the initiative proposes strategies to enhance reliability and introduces a benchmark designed for evaluating LLM-based judgments. By outlining practical applications, challenges, and future directions,	By Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Kun		March 9, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		this work serves as a foundational reference for advancing AI-driven evaluation systems across multiple domains.	Zhang, Saizhuo Wang, Yuanzhuo Wang, Wen Gao, Lionel Ni, Jian Guo		
3.14	Multi Agent Bench	Multi Agent Bench is a comprehensive benchmark designed to evaluate multi-agent systems powered by Large Language Models (LLMs). It overcomes the limitations of traditional single-agent and domain-specific benchmarks by focusing on diverse, interactive scenarios that emphasize both collaboration and competition. The framework, MARBLE, incorporates innovative metrics such as milestone-based KPIs and supports flexible coordination protocols. Key findings show that graph-based coordination and cognitive planning outperform other methods, highlighting emergent social behaviors and advancing research toward AGI-level multi-agent collaboration.	By Kunlun Zhu, Hongyi Du, Zhaochen Hong, Xiaocheng Yang, Shuyi Guo, Zhe Wang Zhenhailong Wang, Cheng Qian, Xiangru Tang, Heng Ji, Jiaxuan You		March 3, 2025
3.15	Meta Plan Optimization (MPO)	Recent LLM advancements enable interactive planning, but existing methods suffer from hallucinations and require retraining for new agents. Meta Plan Optimization (MPO) enhances agent planning by incorporating explicit guidance via meta plans, avoiding reliance on complex human-curated knowledge. Unlike traditional methods, MPO continuously optimizes meta plans using feedback from task execution, improving efficiency and generalization. Experiments on two benchmark tasks show MPO outperforms existing approaches. Analysis confirms MPO as a plug-and-play solution, boosting task completion while adapting to	By Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma,		March 4, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		unseen scenarios. This framework provides a scalable way to refine agent reasoning without retraining, making AI planning more effective and flexible.	Honghao Liu, Saizhuo Wang, Kun Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni and Jian Guo		
3.16	Process-based Self-Rewarding Language Models (PSRLM)	Process-based Self-Rewarding Language Models" introduces a framework to improve LLMs' reasoning, especially in mathematics. Traditional self-rewarding methods, where LLMs evaluate their own outputs, often struggle with complex reasoning and may reduce accuracy. To solve this, a process-based approach is proposed, incorporating step-by-step reasoning where LLMs judge each intermediate step. This enables fine-grained feedback and iterative self-improvement. The method significantly enhances performance across mathematical benchmarks, showing that LLMs can refine their reasoning without external rewards. This approach suggests that LLMs could achieve, or even surpass, human-level reasoning through structured self-optimization.	By Shimao Zhang, Xiao Liu, Xin Zhang, Junxiao Liu, Zheheng Luo, Shujian Huang, Yeyun Gong		March 5, 2025
3.17	Search-R1: Training LLMs to Reason and Leverage Search Engines with Reinforcement Learning	The paper presents SEARCH-R1, an innovative reinforcement learning framework that enables large language models (LLMs) to integrate multi-turn reasoning with search engine interactions. By autonomously generating queries and leveraging retrieved information, SEARCH-R1 addresses limitations in retrieval-augmented generation (RAG) and tool-use approaches. Key features include retrieved token masking for stable optimization, structured multi-turn reasoning, and a simple outcome-based reward function. Experiments on seven datasets show up to 26% improvement over state-of-the-art baselines, demonstrating its effectiveness in	By Bowen Jin and Hansi Zeng and Zhenrui Yue and Dong Wang and Hamed Zamani and Jiawei Han		March 12, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		complex reasoning tasks requiring external knowledge. This work advances search-augmented LLMs and provides insights into reinforcement learning for retrieval-enhanced reasoning.			
3.18	The Comprehensive Guide to Understanding Modern Language Models	"The Hundred-Page Language Models Book" offers a technical introduction to large language models, covering everything from basic concepts to advanced techniques. The book methodically progresses through language modeling fundamentals, recurrent neural networks, and transformer architectures before delving into modern LLMs. It explains critical concepts such as self-attention mechanisms, positional encoding, and the impact of scale on model performance. With detailed diagrams and Python implementations, the guide serves as both an educational resource for newcomers and a reference for practitioners, concluding with insights into advanced topics like mixture of experts, model merging, and compression techniques	By Cornelius Yudha Wijaya,		March 13, 2025
3.19	Transformers without Normalization	A research team from Meta, NYU, MIT, and Princeton has introduced Dynamic Tanh (DyT), a novel alternative to normalization layers in Transformers. DyT replaces Layer Normalization (LN) with a simple element-wise operation $DyT(x) = \tanh(\alpha x)$, where α is a learnable parameter. This method achieves or surpasses the performance of traditional normalization while improving training efficiency and reducing computational overhead. Experiments across vision, language, and self-supervised learning confirm DyT's effectiveness, particularly in LLaMA and Vision Transformer (ViT) models. The findings challenge the long-held belief that normalization is essential for stable deep learning model training	By Jiachen Zhu, Xinlei Chen, Kaiming He, Yann LeCun and Zhuang Liu		March 13, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.20	Block Diffusion: A Hybrid Approach Between Autoregressive and Diffusion Models	Researchers from Cornell Tech, Stanford, and Cohere introduced Block Diffusion Language Models (BD3-LMs), a new approach that combines discrete denoising diffusion and autoregressive models. Unlike traditional diffusion models, BD3-LMs allow for variable-length text generation and leverage KV caching for improved efficiency. They introduce novel training techniques, such as gradient variance estimators and data-driven noise schedules, to close the performance gap with autoregressive models. BD3-LMs achieve state-of-the-art perplexity among diffusion-based models and are capable of generating arbitrarily long sequences, overcoming a major limitation of prior diffusion approaches	By Marianne Arriola, Aaron Gokaslan, Justin T Chiu, Zhihan Yang, Zhixuan Qi, Jiaqi Han, Subham Sekhar Sahoo and Volodymyr Kuleshov		March 12, 2025
3.21	Optimizing Test-Time Compute via Meta Reinforcement Fine-Tuning	Meta Reinforcement Fine-Tuning (MRT), proposed by Qu et al., enhances Large Language Models' (LLMs) reasoning during inference by optimizing test-time compute. The method frames test-time compute optimization as a meta-reinforcement learning problem, treating the LLM's output sequence as episodes. MRT employs cumulative regret to assess the effectiveness of test-time compute, balancing exploration and exploitation. Incorporating dense rewards based on information gain, MRT achieves a 2-3x performance improvement and a 1.5x increase in token efficiency for mathematical reasoning tasks compared to traditional methods.	By uxiao Qu and Matthew Y. R. Yang, Amrith Setlur and Lewis Tunstall, Edward Emanuel Beeching, Ruslan Salakhutdinov and Aviral Kumar		March 10, 2025
3.22	PLAN-AND-ACT: Improving Planning of Agents for Long-Horizon Tasks	The paper "PLAN-AND-ACT" introduces a novel framework to improve the planning capabilities of large language model (LLM)-based agents for complex, long-horizon tasks. By separating high-level planning (PLANNER) from low-level execution (EXECUTOR), the framework enables better alignment between user goals and executable actions. It proposes a scalable synthetic data generation	By Lutfi Eren Erdogan, Nicholas Lee, Sehoon Kim, Suhong Moon, Hiroki Furuta,		March 12, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		pipeline to train the PLANNER effectively without manual annotations. Evaluated on the WebArena-Lite benchmark, PLAN-AND-ACT achieves a state-of-the-art success rate of 54%, showcasing its ability to handle dynamic environments and improve task consistency. This work represents a significant step forward in enhancing LLM-driven agents for real-world applications.	Gopala Anumanchipalli, Kurt Keutzer, Amir Gholami		
3.23	Search-R1: Training LLMs to Reason and Leverage Search Engines with Reinforcement Learning	Researchers have introduced Search-R1, an innovative framework that combines reinforcement learning (RL) with search engine interactions to improve the reasoning capabilities of large language models (LLMs). Building upon the DeepSeek-R1 model, Search-R1 enables LLMs to autonomously generate multiple search queries during step-by-step reasoning, allowing for real-time retrieval of external information. This approach addresses limitations in existing retrieval-augmented generation methods by supporting complex multi-turn retrieval without relying on large-scale supervised data. Experiments across seven question-answering datasets demonstrated significant performance improvements: 26% with Qwen2.5-7B, 21% with Qwen2.5-3B, and 10% with LLaMA3.2-3B models over state-of-the-art baselines.	By Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani and Jiawei Han		March 12, 2025
3.24	Monitoring Reasoning Models for Misbehavior and the Risks of Promoting Obfuscation	Monitoring Reasoning Models for Misbehavior and the Risks of Promoting Obfuscation explores Chain-of-Thought (CoT) monitoring for detecting misalignment in AI models. It examines how CoT-based oversight improves the detection of reward hacking and deceptive reasoning in advanced AI systems like o3-mini. While CoT monitoring enhances transparency, it can also lead to obfuscated reward hacking, where models strategically hide misaligned behavior. The study highlights the benefits and risks of CoT-based monitoring, emphasizing	By Bowen Baker, Joost Huizinga, Leo Gao, Zehao Dou, Melody Y. Guan, Aleksander Madry, Wojciech Zaremba, Jakub		March 10, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		the need for robust oversight mechanisms to prevent AI from gaming reward systems while ensuring reliable, interpretable decision-making.	Pachocki, David Farhi		
3.25	LMM-R1: Enhancing Reasoning in 3B-Parameter Large Multimodal Models	Researchers have introduced LMM-R1, a two-stage framework designed to enhance reasoning capabilities in 3-billion-parameter Large Multimodal Models (LMMs). The first stage, Foundational Reasoning Enhancement (FRE), employs rule-based reinforcement learning (RL) on text-only data to strengthen reasoning abilities. Subsequently, the Multimodal Generalization Training (MGT) stage extends these capabilities to multimodal domains. Experiments on Qwen2.5-VL-Instruct-3B demonstrated that LMM-R1 achieved average improvements of 4.83% in multimodal benchmarks and 4.5% in text-only benchmarks compared to baselines, with a notable 3.63% gain in complex tasks like Football Game analysis. This approach offers a data-efficient paradigm, bypassing the need for costly high-quality multimodal training data.	By Yingzhe Peng, Gongrui Zhang, Miaosen Zhang, Zhiyuan You, Jie Liu, Qipeng Zhu, Kai Yang, Xingzhong Xu, Xin Geng and Xu Yang		March 11, 2025
3.26	AUDITING LANGUAGE MODELS FOR HIDDEN OBJECTIVES	The paper "Auditing Language Models for Hidden Objectives" explores the feasibility of alignment audits to uncover undesired objectives in AI systems. Using a testbed model trained with a hidden objective of reward model (RM) sycophancy, the study demonstrates how alignment audits can detect behaviors that exploit RM biases. Through a blind auditing game and unblinded analysis, the paper evaluates various auditing techniques, including sparse autoencoders, behavioral attacks, and training data analysis. This work highlights the importance of alignment audits in ensuring AI safety, providing a methodology to identify hidden objectives and improve pre-deployment safety practices for AI systems.	By Anthropic		March 13, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.27	reWordBench: Evaluating and Enhancing Reward Model Robustness	reWordBench investigates how even small input transformations affect the performance of reward models, which play a crucial role in text evaluation and alignment processes. However, due to overfitting, these models may not fully reflect their true capabilities. The authors developed reWordBench, a tool that tests models using transformations that preserve meaning or ranking. Results indicate that state-of-the-art models are fragile. As a solution, they propose a training strategy that assigns similar scores to paraphrases, significantly improving robustness and achieving up to a 59% performance gain in alignment tasks.	By Zhaofeng Wu, Michihiro Yasunag, Andrew Cohen, Yoon Kim, Asli Celikyilmaz, Marjan Ghazvininejad		March 18, 2025
3.28	Creation-MMBench: Assessing Context-Aware Creative Intelligence in MLLM	Creation-MMBench is a specialized benchmark designed to evaluate the creative capabilities of Multimodal Large Language Models (MLLMs) in real-world, image-based tasks. It comprises 765 test cases across 51 fine-grained tasks, each with instance-specific evaluation criteria to assess both the general quality of responses and their factual consistency with visual inputs. Experimental results indicate that current open-source MLLMs significantly underperform compared to proprietary models in creative tasks. Additionally, visual fine-tuning has been observed to negatively impact the base LLM's creative abilities. Creation-MMBench offers valuable insights for advancing MLLM creativity and establishes a foundation for future improvements in multimodal generative intelligence.	By Xinyu Fang, Zhijian Chen, Kai Lan, Shengyuan Ding, Yingji Liang, Xiangyu Zhao, Farong Wen, Zicheng Zhang, Guofeng Zhang, Haodong Duan, Kai Chen, Dahua Lin		March 18, 2025
3.29	R1-VL: Step-wise Group Policy	R1-VL: Learning to Reason with Multimodal Large Language Models via Step-wise Group Relative Policy Optimization" aims to enhance the reasoning capabilities of	By Jingyi Zhang,		March 17, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
	Optimization for Multimodal LLMs	Multimodal Large Language Models (MLLMs). Existing methods typically improve MLLMs' reasoning abilities by fine-tuning them on high-quality reasoning data generated by powerful models. However, these approaches may lead models to merely imitate successful reasoning paths, preventing them from understanding incorrect reasoning paths. To address this issue, the authors propose a novel online reinforcement learning framework called Step-wise Group Relative Policy Optimization (StepGRPO), which enables MLLMs to enhance their reasoning capabilities autonomously.	Jiaxing Huang, Huanjin Yao, Shunyu Liu, Xikun Zhang, Shijian Lu, Dacheng Tao		
3.30	Synopsys Introduces AgentEngineer for AI-Assisted Chip Design	Synopsys has introduced AgentEngineer, a groundbreaking technology that leverages AI agents to assist engineers in complex chip design tasks, such as testing circuit designs. As companies like Nvidia transition to designing AI server systems with thousands of interconnected chips, the complexity and speed of the design process have increased significantly. AgentEngineer enhances research and development capabilities without the need to expand engineering teams. It coordinates intricate system designs, ensuring timely product delivery, marking a significant step in the evolution of AI-assisted engineering for chip development.	By Stephen Nellis		March 19, 2025
3.31	ELTEX: A Framework for Domain-Driven Synthetic Data Generation	ELTEX (Efficient LLM Token Extraction) is a domain-driven framework for generating high-quality synthetic training data in specialized fields like cybersecurity. LLMs struggle in such domains due to limited domain-specific data. ELTEX tackles this by integrating domain indicator extraction and dynamic prompting to retain critical knowledge. Tested on blockchain cyberattack detection, ELTEX fine-tuned Gemma-2B, achieving GPT-4-level performance in classification and uncertainty calibration with fewer resources. A synthetic dataset of social media texts for cyberattack detection is released. This study	By Arina Razmyslovich, Kseniia Murasheva, Sofia Sedlova, Julien Capitaine, Eugene Dmitriev		March 19, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		shows that domain-driven synthetic data can help smaller models compete with larger ones.			
3.32	SWEET-RL: Training Multi-Turn LLM Agents on Collaborative Reasoning Tasks	This paper, SWEET-RL is a reinforcement learning method designed to improve multi-turn collaboration between LLM agents and humans. Introduced alongside ColBench, a benchmark for backend programming and frontend design tasks, SWEET-RL uses a critic model trained with extra training-time data to deliver step-wise rewards. This enhances the agent's decision-making across turns. Experiments show a 6% absolute gain in success and win rates over prior RL methods. Notably, SWEET-RL enables Llama-3.1-8B to perform on par with or better than GPT-4o in collaborative reasoning and content creation tasks.	By Yifei Zhou, Song Jiang, Yuandong Tian, Jason Weston, Sergey Levine, Sainbayar Sukhbaatar, Xian Li		March 19, 2025
3.33	Plug-and-Play 1.x-Bit KV Cache Quantization for Video Large Language Models	VideoLLMs can handle long video inputs but face memory and speed bottlenecks due to large KV caches from thousands of visual tokens. While 2-bit KV cache quantization has minimal impact on performance, lower-bit quantization remains underexplored. This paper introduces VidKV , a plug-and-play method compressing KV cache below 2 bits. It uses mixed-precision key quantization (2-bit for anomalous channels, 1-bit with FFT for others) and 1.58-bit value quantization with selective token filtering. Experiments on LLaVA-OV-7B and Qwen2.5-VL-7B show near-FP16 performance with 1.5-bit compression, outperforming prior per-token quantization by using per-channel strategies.	By Keda Tao, Haoxuan You, Yang Sui, Can Qin, Huan Wang		March 21, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.34	When Less is Enough: Adaptive Token Reduction for Efficient Image Representation	Vision encoders generate many visual tokens, offering rich representations but increasing computational costs. This paper introduces a method to assess and select only the most informative tokens, based on the idea that less useful features can be reconstructed from more important ones. Using an autoencoder with a Gumbel-Softmax mechanism, the approach retains key tokens while discarding redundant ones. Applied to the LLaVA-NeXT model, the method shows that over 50% of visual context can be removed in OCR tasks with minimal performance loss. It enables efficient, adaptive multimodal pruning for scalable inference without compromising quality.	By Eduard Allakhverdov, Elizaveta Goncharova, Andrey Kuznetsov		March 20, 2025
3.35	Modifying Large Language Model Post-Training for Diverse Creative Writing	This paper introduces a method to enhance both output diversity and quality in creative writing tasks using large language models (LLMs). The approach incorporates a "deviation"-based objective during training, prioritizing rare and high-quality samples by measuring how different each example is from others. Applied to Direct Preference Optimization (DPO) and Odds Ratio Preference Optimization (ORPO), the method enables an 8B-parameter model to generate outputs as diverse as human-curated datasets and match the quality of top models like GPT-4o and DeepSeek-R1. The approach is validated through human evaluations, ablation studies, and comparisons with the existing diversification method DivPO.	By John Joon Young Chung, Vishakh Padmakumar, Melissa Roemmele, Yuqian Sun, Max Kreminski		March 21, 2025
3.36	Judge Anything: MLLM as a Judge Across Any Modality	Multimodal Large Language Models (MLLMs) are gaining traction for their ability to process visual and textual data, but evaluating their judgment remains difficult due to limited benchmarks aligned with human preferences. The "MLLM-as-a-Judge" benchmark addresses this with tasks like Scoring Evaluation, Pair Comparison, and Batch Ranking. Results show MLLMs perform well in Pair	By Shu Pu, Yaochen Wang, Dongping Chen, Yuhang Chen, Guohao		March 21, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		Comparison but struggle with scoring and ranking, revealing issues like bias, hallucinations, and inconsistency—even in top models like GPT-4V. To improve reliability, the study introduces OmniArena, an automated platform designed to support better evaluation of multimodal models and reward systems through more rigorous, scalable testing.	Wang, Qi Qin, Zhongyi Zhang, Zhiyuan Zhang, Zetong Zhou, Shuang Gong, Yi Gui, Yao Wan, Philip S. Yu		
3.37	Typed-RAG: Type-aware Multi-Aspect Decomposition for Non-Factoid Question Answering	Non-factoid question-answering (NFQA) is challenging due to its open-ended nature, varied intents, and need for multi-aspect reasoning, making traditional factoid approaches like RAG insufficient. To address this, Typed-RAG is introduced—a type-aware framework that enhances RAG by classifying NFQs into categories such as debate, experience, and comparison. It decomposes complex questions into single-aspect sub-queries, improving both retrieval and response quality. By aggregating these focused results, Typed-RAG delivers more informative and relevant answers. A new benchmark, Wiki-NFQA, was developed to evaluate performance, showing that Typed-RAG outperforms existing baselines through its structured, type-driven approach to NFQA.	By DongGeon Lee, Ahjeong Park, Hyeri Lee, Hyeonseo Nam, Yunho Maeng		March 21, 2025
3.38	Modifying Large Language Model Post-Training for Diverse Creative Writing	Creative writing tasks require large language models (LLMs) to produce diverse, high-quality responses. However, post-training methods often prioritize quality while neglecting diversity. Researchers introduce a novel approach that incorporates “deviation”—how different a sample is from others with the same prompt—into the training objective. Applied to Direct Preference Optimization (DPO) and Odds Ratio Preference Optimization (ORPO), this method improves output diversity with minimal quality loss. An 8B-parameter model achieved	By John Joon Young Chung, Vishakh Padmakumar, Melissa Roemmele, Yuqia		March 21, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		human-level diversity and matched the quality of top instruction-tuned models like GPT-4o and DeepSeek-R1. Human evaluations and comparisons to DivPO further validate the effectiveness of the approach.	n Sun, Max Kreminski		
3.39	SuperBPE: Space Travel for Language Models	SuperBPE, an innovative tokenization algorithm that extends the traditional Byte-Pair Encoding (BPE) by incorporating "superword" tokens that bridge multiple words. SuperBPE improves encoding efficiency by up to 33% compared to BPE, enabling shorter token sequences and reducing inference compute by 27%. Pretraining experiments on 8B-scale language models show that SuperBPE achieves a +4.0% average improvement across 30 downstream tasks, including a +8.2% gain on MMLU. This straightforward modification enhances both performance and efficiency without altering model architectures, making it a transformative advancement in tokenization for language modeling.	By NVIDIA, Allen Institute for AI		March 17, 2025
3.40	TULIP: Towards Unified Language-Image Pretraining	TULIP (Towards Unified Language-Image Pretraining) is a novel framework designed to overcome the limitations of existing contrastive image-text models like CLIP and SigLIP, which struggle with fine-grained visual understanding and spatial reasoning. By integrating generative data augmentation, patch-level contrastive learning, and reconstruction-based feature regularization, TULIP achieves a balance between fine-grained visual detail and global semantic alignment. Scaling to over 1B parameters, TULIP establishes state-of-the-art performance across benchmarks, excelling in zero-shot classification, vision-language tasks, and fine-grained recognition, marking a significant advancement in multimodal representation learning. Its open-source availability further amplifies its impact.	By University of California		March 19, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.41	Reasoning to Learn from Latent Thoughts	As language model (LM) scaling outpaces the growth of human-written text, data scarcity becomes a bottleneck. To address this, researchers propose inferring latent thoughts —the unspoken reasoning behind text—to improve data efficiency. Treating text as compressed thought, they show that using synthetic latent thought data during pretraining dramatically boosts performance (e.g., 5.7% → 25.4% on MATH). Without relying on a strong teacher, a 1B LM uses an EM algorithm to iteratively enhance its own data and performance. This self-bootstrapping approach reveals promising gains, especially when scaling inference, for more efficient LLM training in data-limited settings.	By Yangjun Ruan, Neil Band, Chris J. Maddison, Tatsunori Hashimoto		March 24, 2025
3.42	Interpreting Reasoning Features in Large Language Models via Sparse Autoencoders	Researchers explore the inner workings of reasoning in large language models (LLMs) using Sparse Autoencoders (SAEs) to uncover interpretable latent features. Focusing on the open-source DeepSeek-R1 model, they identify “reasoning features” that correlate strongly with the model’s complex reasoning capabilities. By extracting and steering these features, they systematically enhance the model’s reasoning performance—offering the first mechanistic explanation of reasoning in LLMs. This breakthrough enables a deeper understanding of how reasoning emerges in neural networks and paves the way for more interpretable and controllable AI systems.	By Andrey Galichin, Alexey Dontsov, Polina Druzhinina, Anton Razzhigaev, Oleg Y. Rogov, Elena Tutubalina, Ivan Oseledets		March 24, 2025
3.43	Learning to Reason with Search for LLMs via Reinforcement Learning	Large language models (LLMs) can integrate external search processes to answer complex, multi-step questions. It introduces a new framework called ReSearch, which trains LLMs to use search as an integral part of their reasoning chains. In this approach, decisions about when and how to search are guided by text-based reasoning, while search results influence subsequent reasoning steps. Extensive experiments are conducted using the Qwen2.5-7B(-Instruct) and Qwen2.5-32B(-	By Mingyang Chen, Tianpeng Li, Haoze Sun, Yijie Zhou, Chenzheng Zhu, Fan Yang, Zenan		March 25, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		Instruct) models during training. Results show that ReSearch demonstrates strong generalization across various benchmarks. Moreover, analysis reveals that it naturally develops advanced reasoning abilities such as reflection and self-correction during training.	Zhou, Weipeng Chen, Haofen Wang, Jeff Z. Pan, Wen Zhang, Huajun Chen		
3.44	CoLLM: A Large Language Model for Composed Image Retrieval	Composed Image Retrieval (CIR) retrieves images using a multimodal query of a reference image and a modification text. Creating training triplets is costly, and existing methods using synthetic data or image-caption pairs have major limitations in scale, diversity, and fusion capability. To address this, CoLLM is introduced—a unified framework that generates training triplets on-the-fly from image-caption pairs without manual labels. It uses LLMs for joint embedding of image and text, enabling deeper vision-language fusion. CoLLM outperforms previous methods on multiple CIR benchmarks, aided by the new 3.4M-sample MTCIR dataset and refined evaluation benchmarks.	By Chuong Huynh, Jinyu Yang, Ashish Tawari, Mubarak Shah, Son Tran, Raffay Hamid, Trishul Chilimbi, Abhinav Shrivastava		March 25, 2025
3.45	Think Twice: Enhancing LLM Reasoning by Scaling Multi-round Test-time Thinking	Recent advancements in LLMs like OpenAI-o1 and DeepSeek-R1 highlight the power of test-time scaling, where iterative reasoning boosts performance. However, challenges remain in long-text processing and RL training efficiency. We introduce Multi-round Thinking, a simple test-time approach where the model uses its previous answer to refine reasoning in subsequent rounds. Tests on QwQ-32B and DeepSeek-R1 show consistent gains across AIME 2024, MATH-500, GPQA-diamond, and LiveCodeBench. For example, QwQ-32B improved from 80.3% to 82.1%, and DeepSeek-R1 from 79.7% to 82.0%. This method proves broadly effective for enhancing LLM performance via stable, multi-step reasoning.	By Xiaoyu Tian, Sitong Zhao, Haotian Wang, Shuaiting Chen, Yunjie Ji, Yiping Peng, Han Zhao, Xiangang Li		March 25, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.46	Interpreting Reasoning Features in Large Language Models via Sparse Autoencoders	Researchers have introduced HoarePrompt, a novel approach that integrates program analysis techniques with large language models (LLMs) to assess program correctness against natural language requirements. HoarePrompt employs a step-by-step process where an LLM generates natural language descriptions of reachable program states at various code points, inspired by the strongest postcondition calculus. To handle loops effectively, the method incorporates few-shot-driven k-induction, adapting a model checking technique. Evaluations using the CoCoClaNeL dataset demonstrate that HoarePrompt significantly outperforms existing methods in classifying program correctness.	By AIRI, MTUCI, Skoltech, Sber, 5HSE		March 24, 2025
3.47	Think Before Recommend: Unleashing the Latent Reasoning Power for Sequential Recommendation	Sequential Recommendation (SeqRec) predicts the next item by analyzing users' past interactions, but existing methods rely on shallow inference, limiting their ability to model evolving user preferences and long-tail items. To overcome this, we introduce ReaRec , the first inference-time computing framework for recommender systems. ReaRec enhances user representations through implicit multi-step reasoning , using autoregressive processing and special position embeddings. Additionally, we propose Ensemble Reasoning Learning (ERL) and Progressive Reasoning Learning (PRL) to maximize reasoning potential. Experiments on five real-world datasets show ReaRec improves SeqRec models by 30%-50% , paving the way for better inference-time computing in recommendations.	By Jiakai Tang, Sunhao Dai, Teng Shi, Jun Xu, Xu Chen, Wen Chen, Wu Jian, Yuning Jiang		March 28, 2025
3.48	PilotANN: Memory-Bounded GPU Acceleration for Vector Search	Approximate Nearest Neighbor Search (ANNS) is essential for modern deep learning, especially in generative models handling complex, high-dimensional data. Existing CPU-based methods can't keep up with rising computational demands, and GPU-only approaches are limited by memory. To address this, the	By Yuntao Gui, Peiqi Yin, Xiao Yan, Chaorui Zhang, Weixi		March 27, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		authors propose PilotANN , a hybrid CPU-GPU system that combines GPU acceleration with CPU memory capacity. It breaks top-k graph search into three stages: GPU-based subgraph traversal with SVD-reduced vectors, CPU-based refinement with full vectors, and improved entry point selection. PilotANN achieves a 3.9–5.4× speedup and supports datasets up to 12× larger than GPU memory , boosting scalability and performance.	Zhang, James Cheng		
3.49	AI Paper Introduces Diversified DPO and ORPO: Post-Training Methods to Boost Output Diversity in Creative Writing with LLMs	The paper introduces Diversified DPO and ORPO , two post-training methods aimed at boosting output diversity in creative writing with Large Language Models (LLMs) . These techniques enhance both semantic and stylistic diversity in the generated text while minimizing quality loss. By leveraging these methods, LLMs can produce more varied and creative outputs, akin to human-written content. The research demonstrates that the application of these approaches results in a noticeable increase in diversity without compromising on quality, offering significant improvements for tasks requiring creative and diverse writing, such as storytelling and content creation.	By Nikhil		March 31, 2025
3.50	Classical Planning with LLM-Generated Heuristics: Challenging the State of the Art with Python Code	Classical Planning with LLM-Generated Heuristics, researchers demonstrate how Large Language Models (LLMs) can autonomously generate domain-specific heuristic functions to solve classical planning tasks. By prompting an LLM to produce multiple Python-coded heuristics for a given domain, evaluating them on training tasks, and selecting the most effective one, they achieved superior performance compared to state-of-the-art domain-independent heuristics. Remarkably, these LLM-generated heuristics were competitive with leading domain-dependent learning algorithms, despite being implemented in unoptimized Python code versus the baselines' optimized C++ implementations.	By Augusto B. Corrêa, André G. Pereira, Jendrik Seipp		March 24, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		This study highlights the potential of LLMs to enhance planning capabilities through automated heuristic generation.			
3.51	CodeARC: Benchmarking Reasoning Capabilities of LLM Agents for Inductive Program Synthesis	The paper introduces CodeARC, a novel framework for evaluating large language models (LLMs) in inductive program synthesis, emphasizing interactive reasoning and self-correction. Unlike static evaluation protocols, CodeARC allows agents to query hidden target functions and use differential testing oracles to refine synthesized programs iteratively. Featuring a diverse benchmark of 1114 Python functions, the study highlights the challenges of inductive synthesis, with the best-performing model achieving a success rate of 52.7%. Fine-tuning on curated synthesis traces improved LLM performance by up to 31%. CodeARC sets a new standard for realistic and rigorous evaluation of LLM-based reasoning and synthesis capabilities.	By MIT, Intel, Visa Research		March 29, 2025
3.52	Agent S2: A Compositional Generalist-Specialist Framework for Computer Use Agents	Agent S2 introduces a novel compositional framework for computer-use agents, addressing key challenges in GUI interaction, long-horizon task planning, and performance bottlenecks in monolithic models. By integrating generalist planning modules with specialist grounding experts, it leverages a Mixture-of-Grounding technique for precise UI localization and Proactive Hierarchical Planning for dynamic task refinement. Agent S2 achieves state-of-the-art results across benchmarks like OSWorld, WindowsAgentArena, and AndroidWorld, significantly outperforming previous methods. Its modular design enhances adaptability, scalability, and generalization, offering a transformative approach to autonomous digital task automation. The code is publicly available, fostering further advancements in this domain.	By Simular Research		April 1, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.53	Efficient Inference for Large Reasoning Models: A Survey	The paper, Efficient Inference for Large Reasoning Models: A Survey, provides a comprehensive review of methods to enhance the efficiency of Large Reasoning Models (LRMs) during inference. LRMs improve reasoning capabilities in complex tasks but face challenges like token inefficiency, memory overhead, and high inference costs. The survey introduces a taxonomy of efficient inference techniques, categorizing them into explicit compact Chain-of-Thought (CoT) and implicit latent CoT methods. It highlights strengths, weaknesses, empirical results, and open challenges, offering insights into improving efficiency without compromising reasoning quality. This work is pivotal for advancing scalable and practical LRM applications.	By Yue Liu, Jiaying Wu, Yufei He, Hongcheng Gao, Hongyu Chen, Baolong Bi, Jiaheng Zhang, Zhiqi Huang, and Bryan Hooi.		March 29, 2025
3.54	What, How, Where, and How Well? A Survey on Test-Time Scaling in Large Language Models	The paper "What, How, Where, and How Well? A Survey on Test-Time Scaling in Large Language Models" provides a comprehensive review of test-time scaling (TTS), a key technique for improving reasoning and problem-solving in large language models (LLMs) during inference. It introduces a unified framework with four dimensions—what, how, where, and how well to scale—for systematic classification and analysis. The survey explores key techniques, applications, and evaluation metrics, addressing challenges like scalability and efficiency. This work serves as a roadmap for future LLM research and deployment, enhancing real-world AI performance.	By Qiyuan Zhang, Fuyuan Lyu, Zexu Sun, Lei Wang, Weixu Zhang, Zhihan Guo, Yufei Wang, Irwin King, Xue Liu, and Chen Ma.		March 31, 2025
3.55	Effectively Controlling Reasoning Models through Thinking Intervention	The paper, Effectively Controlling Reasoning Models through Thinking Intervention, introduces a novel approach called Thinking Intervention, which allows for fine-grained control over reasoning-enhanced large language models (LLMs). By strategically inserting or revising specific thinking tokens during the reasoning process, this method significantly improves model performance across	By NVIDIA, Princeton University		March 31, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		various tasks. The study demonstrates gains of up to 6.7% in instruction-following accuracy, 15.4% in reasoning hierarchy tasks, and a 40% increase in safety alignment. This work opens new avenues for enhancing the reliability and effectiveness of reasoning models in real-world applications.			
3.56	I Have Covered All the Bases Here: Interpreting Reasoning Features in Large Language Models via Sparse Autoencoders	The paper, <i>I Have Covered All the Bases Here: Interpreting Reasoning Features in Large Language Models via Sparse Autoencoders</i> , explores the internal reasoning mechanisms of large language models (LLMs) using Sparse Autoencoders (SAEs). It introduces a novel metric, ReasonScore, to identify reasoning-specific features that enhance model performance in reasoning tasks. Through empirical analysis and controlled feature steering, the study demonstrates that amplifying these features significantly improves reasoning quality, revealing a causal link between identified features and cognitive behaviors such as reflection and exploration. This work provides critical insights into the interpretability of reasoning in LLMs.	By AIRI, MTUCI, Skoltec, Sber and HSE		March 24, 2025
3.57	ReaRAG: Knowledge-guided Reasoning Enhances Factuality of Large Reasoning Models with Iterative Retrieval Augmented Generation	The paper, <i>ReaRAG: Knowledge-guided Reasoning Enhances Factuality of Large Reasoning Models with Iterative Retrieval Augmented Generation</i> , introduces ReaRAG, a model designed to improve the factual accuracy of large reasoning models (LRMs) in multi-hop question answering (QA) tasks. By integrating a structured reasoning approach with iterative retrieval, ReaRAG effectively constructs knowledge-guided reasoning chains that enhance performance while avoiding overthinking. Experimental results demonstrate significant improvements over existing methods across multiple benchmarks, showcasing ReaRAG's ability to refine reasoning trajectories and enhance the robustness of factual answers in complex queries.	By Tsinghua University, Siemens AG		March 27, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.58	UI-R1: Enhancing Action Prediction of GUI Agents by Reinforcement Learning	The paper presents UI-R1, a novel framework that enhances the reasoning capabilities of multimodal large language models (MLLMs) for graphic user interface (GUI) action prediction through rule-based reinforcement learning (RL). By introducing a unique action reward mechanism and utilizing a curated dataset of 136 challenging tasks, UI-R1 significantly improves performance on both in-domain and out-of-domain tasks, achieving average accuracy gains of 22.1% on ScreenSpot and 12.7% on ANDROIDCONTROL. This work highlights the potential of rule-based RL to advance GUI understanding and control, paving the way for future research in this domain.	By vivo AI Lab, MMLab		March 30, 2025
3.59	LEGO-Puzzles: How Good Are MLLMs at Multi-Step Spatial Reasoning?	The paper introduces LEGO-Puzzles, a novel benchmark designed to evaluate the multi-step spatial reasoning capabilities of Multimodal Large Language Models (MLLMs). With over 1,100 visual question-answering samples across 11 tasks, it highlights significant deficiencies in current MLLMs, which achieve only about 50% accuracy compared to over 90% for human participants. The study emphasizes the need for advancements in spatial understanding and sequential reasoning, revealing critical gaps in MLLMs' performance, particularly in complex tasks like image generation and multi-step reasoning, thereby underscoring the challenges in multimodal AI development.	By Shanghai AI Lab, Tongji University, Simons Institute		March 25, 2025
3.60	AgentRxiv: Towards Collaborative Autonomous Research	The paper introduces AgentRxiv, a collaborative framework for autonomous research using large language model (LLM) agents. It emphasizes the significance of iterative improvements in scientific discovery, enabling agents to share and build upon each other's research. The study demonstrates that agents utilizing AgentRxiv achieve notable performance enhancements, such as an 11.4% relative improvement on the MATH-500 benchmark. By fostering collaboration among	By Samuel Schmidgall and Michael Moor		March 23, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		multiple agent laboratories, the framework accelerates research progress and enhances overall accuracy, suggesting a promising direction for integrating autonomous systems into scientific workflows and advancing AI-driven research.			
3.61	Large Language Model Agent: A Survey on Methodology, Applications and Challenges	This survey explores the transformative potential of Large Language Model (LLM) agents, emphasizing their role in advancing artificial general intelligence through goal-driven behaviors and dynamic adaptability. It introduces a methodology-centered taxonomy that systematically deconstructs LLM agent systems into three core dimensions: construction, collaboration, and evolution. The paper highlights diverse applications, including scientific discovery, gaming, and productivity tools, while addressing critical challenges like security, privacy, and scalability. By offering a structured framework and identifying future research directions, this work provides a comprehensive foundation for understanding and advancing LLM agents in complex, real-world environments.	By Junyu Luo, Weizhi Zhang, Ye Yuan et al.		March 27, 2025
3.62	FFN FUSION: RETHINKING SEQUENTIAL COMPUTATION IN LARGE LANGUAGE MODELS	The paper introduces FFN Fusion, a novel optimization technique that reduces sequential computation in large language models (LLMs) by fusing consecutive Feed-Forward Network (FFN) layers into parallel operations. This approach significantly improves inference efficiency, achieving a 1.71× speedup in latency and 35× lower per-token cost without compromising accuracy. Applied to Llama-405B, the resulting Ultra-253B-Base model maintains state-of-the-art performance across benchmarks while reducing parameters from 405B to 253B. FFN Fusion complements existing methods like pruning and quantization, offering scalable efficiency gains and paving the way for future innovations in LLM architecture design.	By NVIDIA		March 24, 2025

✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
4.1	Google's AI Mode for Complex Questions	Google has introduced an AI-powered search mode that allows users to ask multi-part, complex questions. This update improves search by providing contextual and conversational responses, enhancing the experience for in-depth queries. Reflecting Google's continuous effort to refine AI-driven search technologies, this development has the potential to transform how users interact with search engines and access information. By offering more nuanced and relevant answers, it paves the way for a more intuitive and dynamic search experience.	By Aisha Malik		March 5, 2025
4.2	SOFYA	Sofya integrates Llama models for real-time healthcare solutions, hosting them on Oracle Cloud and utilizing Sglang and VLLM for efficient model serving. To meet low-latency demands, the team fine-tuned smaller Llama versions (8B, 3B, 70B) using knowledge distillation and self-reflection prompt engineering. Llama automates data structuring, entity recognition, and question answering, reducing errors and boosting efficiency. This has led to 30% less time spent on documentation, improved workflows, and an average CSAT score of 90%. Sofya plans to scale to 1 million consultations per month, expanding its AI-driven agent flow with Llama 70B for real-time healthcare applications.	By Meta Team		March 5, 2025




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
4.3	AGNTCY	On March 6, 2025, Cisco, LangChain, Galileo, Glean, and LlamaIndex launched AGNTCY, an open-source collective for AI agent interoperability. AGNTCY aims to establish a standardized communication framework, like TCP/IP, to enable seamless interaction between AI platforms. This “Internet of Agents” promotes collaboration, efficiency, and scalability across AI ecosystems. Cisco’s Outshift emphasizes its role in connecting AI systems from different vendors. By engaging the AI and infrastructure community, AGNTCY seeks to build an open, interoperable foundation, addressing multi-agent collaboration challenges and advancing next-generation AI applications globally.	By Emilia David		March 6, 2025
4.4	Manus	China introduced Manus, a general AI agent developed by Monica, gaining attention for its autonomous capabilities. Designed to think, plan, and execute tasks independently, Manus is compared to leading AI systems from OpenAI, Google, and Anthropic. Reports from Newsweek and Business Standard highlight its ability to build websites, plan trips, and analyze stocks without human supervision. This advancement has sparked concerns in Silicon Valley over China's AI leadership, raising ethical and regulatory questions regarding accountability in autonomous AI systems. Manus represents a major step in AI-powered industries, potentially giving China a first-mover advantage globally.	By Kyle Wiggers		March 9, 2025
4.5	OpenAI introduces new tools for building AI agents.	Google has introduced new APIs and tools to help developers build reliable AI agents efficiently. AI agents, autonomous systems that complete tasks independently, require advanced reasoning, multimodal interactions, and improved safety. However, many developers struggle with prompt tuning and orchestration. To address this, Google is launching the Responses API, merging Chat Completions and Assistants API, Built-in Tools for web search, file search, and	By OpenAI		March 11, 2025



✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		computer use, Agents SDK for managing single and multi-agent workflows, and Observability Tools for monitoring agent processes. These innovations streamline AI agent development, making it more accessible and efficient for businesses and developers.			
4.6	AI Surveillance in Schools Raises Privacy Concerns	Schools across the U.S., including Vancouver Public Schools (WA), are using AI-powered surveillance tools like Gaggle to monitor students' online activity on school-issued devices, aiming to prevent violence and address mental health concerns. These systems flag risks like bullying, self-harm, and suicide, alerting staff for intervention. However, a data breach exposing 3,500 unredacted student documents raises privacy and security concerns, including cybersecurity risks, forced outings, and loss of student trust. The effectiveness of such surveillance in improving safety and mental health remains debated, with concerns about long-term privacy impacts.	By Claire Bryan and Sharon Lurye		March 12, 2025
4.7	Gemini 2.0 Flash Introduces Native Image Generation	Google's Gemini 2.0 Flash now supports native image generation, a feature that OpenAI's GPT-4o teased but never launched. This enables users to create and edit images directly, significantly improving text-to-image alignment over existing models. Early testers praise its ease of use, though they criticize Google's UI for making the feature difficult to access. With this release, Google takes the lead in integrated AI image editing, potentially disrupting tools like DALL-E and Midjourney	By Kat Kampf Nicole Brichtova		March 12, 2025
4.8	OpenAI Launches Developer Platform for AI Agents	OpenAI has launched a developer platform aimed at facilitating the creation of AI agents. This platform equips developers with tools that enable AI agents to perform web and file searches, as well as execute web-based tasks similar to	By HIStalk Team		March 12, 2025




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		OpenAI's Operator browser. The initiative seeks to streamline the development process for AI applications capable of interacting with the internet and managing various online tasks, thereby enhancing the functionality and versatility of AI agents.			
4.9	NHS England Deploys AI Tool to Predict Patient Falls	NHS England is implementing an AI tool developed by Cera, designed to predict a patient's risk of falling with 97% accuracy. Beyond assessing fall risk, this software is also utilized to anticipate deterioration in home care patients, aiming to enhance preventive measures and improve patient safety. By accurately identifying individuals at higher risk, the tool enables healthcare providers to intervene proactively, potentially reducing hospital admissions and improving overall patient outcomes.	By HIStalk Team		March 10, 2025
4.10	Oracle Introduces AI Agents to Combat Financial Crime	Oracle Financial Services has enhanced its Investigation Hub Cloud Service with advanced AI capabilities to bolster the fight against financial crime. These agentic AI features aim to reduce manual workloads for financial institutions by automating complex investigative processes, thereby enabling faster detection and prevention of fraudulent activities. This development underscores Oracle's commitment to leveraging artificial intelligence to improve operational efficiency and strengthen compliance measures within the financial sector.	By Oracle		March 12, 2025
4.11	Singapore Airlines Implements AI-Powered Customer Service Platform	Singapore Airlines has partnered with Salesforce to enhance its customer service operations through AI integration. The airline is incorporating Agentforce, Einstein in Service Cloud, and Data Cloud into its customer case management system to deliver more personalized service. Agentforce deploys autonomous agents to perform specific tasks, streamlining operations and allowing	By TechNode Global Staff		March 13, 2025




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		representatives to focus on providing enhanced attention in customer interactions. The system also leverages Einstein generative AI to summarize previous customer interactions and provide guidance, enabling representatives to better understand and anticipate customer needs while reducing response times			
4.12	CERAAweek: AI leading to faster, cheaper oil production, executives say	At the CERAAweek energy conference, industry executives highlighted how AI is revolutionizing oil and gas production, making it faster and more cost-efficient. Companies are using AI for predictive maintenance, drilling optimization, and reservoir management, leading to higher output and reduced expenses. Major firms like ExxonMobil and Chevron emphasized AI's role in boosting efficiency while lowering environmental impact. The technology enhances decision-making, minimizes downtime, and improves energy exploration. As AI adoption grows, experts predict a fundamental shift in how oil and gas operations are managed, increasing profitability and sustainability in the sector.	By Sheila Dang and Georgina Mccartney		March 14, 2025
4.13	PortaOne Addresses AI Monetization Challenges with New Platform	At MWC Barcelona 2025, PortaOne showcased PortaAIM, a flexible charging and revenue management platform designed to help AI businesses optimize pricing, track costs, and scale profitably. CEO Andriy Zhylenko highlighted the critical gap between AI innovation and sustainable business models, noting that while AI agents are advancing rapidly, many companies fail to establish scalable monetization frameworks. PortaOne's platform addresses this challenge by enabling seamless customer charging while maintaining cost control. The company has already implemented a proof-of-concept virtual AI agent that can retrieve customer-specific data in real-time, transforming customer service approaches	By Hennadiy Kornev		March 13, 2025



 AI Use Cases




#	Highlights	Summary	Author	Source	Date
4.14	AAA to Showcase AI-Powered Innovations at Legalweek 2025	The American Arbitration Association (AAA) will showcase its AI-driven advancements at Legalweek 2025 (March 24-27, New York). On March 25, AAA President & CEO Bridget McCormack will join the panel “Disruption by Design: Shaping the AI-Forward Firm of Tomorrow”, discussing AI’s role in alternative dispute resolution (ADR). Attendees can explore AAA’s AI innovations at Booth 2311, including ClauseBuilder® AI (Beta), AAAi Chat Book Case Prep (Beta), and AAAiLab™. These initiatives highlight AAA’s commitment to enhancing legal efficiency and accessibility through AI.	By American Arbitration Association		March 13, 2025
4.15	Eclipse Foundation Introduces AI-Enhanced Open-Source IDE	The Eclipse Foundation has released an alpha version of its open-source Theia integrated development environment (IDE), now enhanced with artificial intelligence (AI) capabilities. This AI-powered Theia IDE allows developers to integrate coding agents with various large language models (LLMs), facilitating tasks such as prompt engineering, defining agentic AI behaviors, and customizing user interfaces. This development aims to prevent vendor lock-in by enabling the use of multiple LLMs and supports interoperability through a Model Contextual Protocol (MCP) that connects AI-driven workflows with external tools and data sources.	By Mike Vizard		March 13, 2025
4.16	Open-source AI matches top proprietary model in solving tough medical cases	A recent NIH-funded study led by Harvard Medical School researchers compared the performance of the open-source AI model Llama 3.1 405B with the proprietary GPT-4 in diagnosing complex medical cases. Analyzing 92 challenging scenarios from The New England Journal of Medicine, Llama 3.1 correctly diagnosed 70% of cases, surpassing GPT-4's 64%. Notably, Llama 3.1 ranked the correct diagnosis first in 41% of cases, compared to GPT-4's 37%. This advancement highlights the growing competitiveness of open-source AI in healthcare, offering benefits like	By Harvard Medical School		March 14, 2025



✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		enhanced data privacy and customization, as these models can be tailored to specific clinical needs and operated within hospital infrastructures			
4.17	AI-powered phishing attacks threaten Gmail, Outlook, and Apple Mail users	The Forbes article highlights the escalating threat of AI-driven phishing attacks targeting users of popular email services like Gmail, Outlook, and Apple Mail. Cybercriminals are leveraging artificial intelligence to craft highly personalized and convincing emails, making it increasingly difficult for users to distinguish between legitimate messages and malicious ones. These sophisticated phishing attempts can lead to significant security breaches, including unauthorized access to sensitive information and financial losses. The article emphasizes the importance of heightened vigilance and the implementation of advanced security measures to protect against these evolving AI-generated threats.	By Zak Doffman		March 16, 2025
4.18	Chinese company's 'dark factory' will no human workers soon be the norm	In Changping, China, Xiaomi's "dark factory" is revolutionizing manufacturing with a fully automated, AI-driven facility that operates 24/7 without human intervention. This facility produces one smartphone per second, using robots and AI systems for all production steps, from raw materials to assembly, ensuring precision and eliminating human error. With a capacity of 10 million devices annually, this 81,000-square meter factory represents the future of manufacturing. However, concerns about job displacement are rising, as the World Economic Forum predicts 23% of jobs will be affected by AI within five years. Experts call for global cooperation and regulation to ensure AI's safe and equitable benefits.	By Alex Blair		March 16, 2025




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
4.19	Anthropic Introduces Harmony, an AI Assistant Designed for Local File Access	Anthropic, a leading company in AI safety and research, is developing new features to expand the capabilities of its AI assistant, Claude. One of these features, Harmony, allows users to integrate their local file directories into Claude’s context, enabling the AI to read, index, and analyze these files. Harmony is a new feature aimed at enhancing user experience by enabling Claude to interact with local files seamlessly. This functionality offers significant advantages in AI-assisted code assistance and content analysis. Anthropic’s Compass, another feature in development, aims to extend Claude’s deep research capabilities.	By Alexey Shabanov		March 15, 2025
4.20	NetApp Announce Agentic AI Reasoning Solutions with NVIDIA	NetApp has partnered with NVIDIA to launch a new AI infrastructure designed for agentic AI solutions. This system integrates NetApp’s ONTAP storage operating system with NVIDIA’s accelerated computing, networking, and AI software to enhance enterprise data management and AI workloads. Built on the NVIDIA AI Data Platform reference architecture, it optimizes large-scale data processing, ensuring AI applications run efficiently, scalably, and securely. By accelerating AI-driven decision-making and automation, businesses can better leverage data-driven AI models. This innovation enables enterprises to process vast datasets, improve AI inference speeds, and unlock new possibilities in real-time analytics, automation, and competitive intelligence.	By InsideHPC		March 18, 2025
4.21	Innodata Unveils Generative AI Test and Evaluation Platform, Built with NVIDIA Technology	Innodata has introduced the beta version of the Generative AI Test & Evaluation Platform to enhance the security, reliability, and performance of AI models. Powered by NVIDIA’s advanced inference technology, this platform conducts automated attack tests on AI models to identify potential security vulnerabilities and performs model comparisons. Early adopters, such as MasterClass, plan to use the platform to ensure their AI investments deliver reliable and secure	By Innodata		March 18, 2025

✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		customer experiences. The full version of the platform is expected to be released in the second quarter of 2025.			
4.22	Tencent expands AI push with open-source 3D generation tools	Tencent has announced the release of a suite of new AI tools that convert text and images into 3D visuals, enhancing China's advancements in generative AI. The company introduced five open-source models leveraging its Hunyuan3D-2.0 technology, with "turbo" versions capable of generating high-quality 3D visuals in just 30 seconds. This launch signifies Chinese firms' increasing competition with the U.S. by offering high-performing AI models at reduced costs. Tencent's initiative builds on its first introduction of 3D AI models in November 2024, focusing on designers and game developers.	By Reuters		March 18, 2025
4.23	AI-Powered Robot Developed to Autonomously Make Coffee	Researchers from the University of Edinburgh have developed an AI-powered robot capable of autonomously making coffee. Equipped with a seven-jointed robotic arm, the robot integrates AI, sensors, and motor skills to navigate dynamic environments like kitchens. It can interpret verbal instructions, locate and retrieve a mug, and mix precise ratios of coffee and water. Unlike traditional robots reliant on pre-programmed actions, this robot adapts to unforeseen events, such as unexpected object movements. This breakthrough signifies a leap in integrating reasoning, movement, and perception in robotics, potentially enhancing automation in daily tasks.	By Elizabeth Hunter		March 19, 2025
4.24	Disney Collaborates with Nvidia and Google DeepMind on Advanced Robotics	Disney is partnering with Nvidia and Google DeepMind to develop a new physics engine, Newton , to enhance its advanced robotic characters. At Nvidia's GTC AI Conference, CEO Jensen Huang introduced an endearing robot named Blue , created by Disney Research. Newton, built on Nvidia's Warp framework, will aid	By Maxwell Zeff		March 18, 2025



✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		robots in completing complex tasks more efficiently. Disney plans to deploy these upgraded, lifelike robots in its parks, including Walt Disney World, Tokyo Disneyland, and Disneyland Paris, marking a significant leap in the development of expressive and functional robots for entertainment.			
4.25	Digital Physiotherapy for Back Pain Management	With around one in six people globally suffering from chronic back pain, access to physiotherapy remains limited—especially in rural or underserved areas. AI-powered digital physiotherapists are emerging as a scalable solution. These systems guide patients through tailored exercise routines, monitor movements using computer vision, and provide real-time feedback. Such tools not only extend access to care but also reduce the burden on healthcare systems. While they don't replace human therapists entirely, AI-driven platforms are proving effective in supporting rehabilitation and encouraging consistent, at-home treatment for back pain patients.	By Scott Nover		March 24, 2025
4.26	Coffee-making robot breaks new ground for AI machines	Researchers at the University of Edinburgh have developed a robotic arm that combines advanced motor skills with AI, enabling it to perform complex tasks in unpredictable settings like kitchens. Unlike traditional robots limited to pre-programmed actions, this robot interprets verbal instructions, analyzes its surroundings, and adapts in real time. It can open unfamiliar drawers, locate a mug, and prepare coffee—even adjusting if objects are moved mid-task. Published in Nature Machine Intelligence, the study highlights the integration of reasoning, perception, and movement as key to building robots capable of operating alongside humans.	By University Of Edinburgh		March 19, 2025




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
4.27	NatWest Partners with OpenAI in Landmark UK Banking Collaboration	NatWest has entered a milestone partnership with OpenAI, becoming the first UK bank to formally integrate the company's generative AI tools into its operations. The collaboration will support customer service, automate internal workflows, and assist staff with complex decision-making tasks. NatWest aims to leverage OpenAI's capabilities to enhance productivity and client engagement while ensuring responsible AI governance. This deal marks a significant step in AI adoption within the financial sector and sets a precedent for broader industry use of generative models.	By Iain Withers,		March 20, 2025
4.28	Micron forecasts upbeat quarterly revenue on strong AI memory chip demand	Micron Technology has issued an upbeat revenue forecast for the upcoming quarter, driven by soaring demand for high-bandwidth memory (HBM) chips used in AI applications. These advanced memory chips are essential for powering large language models and AI data centers. Micron is ramping up production to compete with industry leaders like Samsung and supply key players such as NVIDIA. The company's bullish outlook signals the growing strategic importance of AI hardware, as demand accelerates across cloud, enterprise, and AI sectors. Investors responded positively, pushing Micron shares higher following the announcement.	By Juby Babu		March 21, 2025
4.29	Training Video Foundation Models with NVIDIA NeMo	This paper introduces NVIDIA NeMo's framework for training large-scale video foundation models. It highlights a modular pipeline optimized for multimodal learning, leveraging transformer-based architectures and large-scale datasets. The authors discuss pretraining techniques, dataset curation, and efficiency improvements, including memory optimization and distributed training strategies. They also showcase benchmark results, demonstrating state-of-the-art performance in video understanding tasks. The paper emphasizes NeMo's	By Nvidia		March 17, 2025




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		adaptability for various applications, such as video captioning and action recognition, providing insights into the future of scalable video model training.			
4.30	Opus 2 announces AI Workbench to transform case management, strategy, and analysis	Opus 2 has launched AI Workbench, a powerful new feature integrated with its Opus 2 Cases platform to revolutionize legal case management and analysis. The tool enables lawyers to query, summarize, and extract insights from large volumes of case documents using AI-powered natural language interaction. By identifying key people, events, and legal issues, AI Workbench streamlines case strategy and decision-making. Early adopters like Norton Rose Fulbright report major productivity gains, especially in summarizing testimonies and analyzing complex records. This launch marks a significant advancement in applying AI to litigation and legal workflow optimization.	By Opus 2		March 24, 2025
4.31	On Large Multimodal Models as Open-World Image Classifiers	Traditional image classification relies on predefined semantic categories, whereas Large Multimodal Models (LMMs) classify images directly using natural language (e.g., answering “What is the main object in the image?”). While LMMs offer remarkable flexibility, most studies assess their performance in closed-world settings with fixed categories. This work addresses that limitation by evaluating LMMs in a truly open-world setting. We define an evaluation protocol, introduce alignment metrics, and analyze 13 models across 10 benchmarks, covering various class granularities. Our findings highlight key challenges, including fine-grained classification, and demonstrate how tailored prompting and reasoning can improve performance.	By Alessandro Conti, Massimiliano Mancini, Enrico Fini, Yiming Wang, Paolo Rota, Elisa Ricci		March 27, 2025


✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
4.32	Bridget Phillipson eyes AI's potential to free up teachers' time	Education Secretary Bridget Phillipson is exploring the use of Artificial Intelligence (AI) tools in schools to reduce teachers' administrative workloads, aiming to allow more focus on direct teaching. At a Department for Education hackathon, developers showcased AI applications such as compiling student reports, assessing vocational work, and recording developmental observations. For instance, tools like digital microscopes provide rapid feedback on vocational qualifications, and lapel microphones enable early-years teachers to efficiently document observations. Phillipson emphasized that AI should complement, not replace, human teaching, enhancing educational outcomes and addressing teacher recruitment challenges.	By Richard Adams		March 31, 2025
4.33	AI business technology takes on shoplifters and admin drag	Paris-based startup Veesion is leveraging AI to combat retail theft through advanced surveillance technology. Their system analyzes live video feeds from store cameras to detect suspicious movements, helping identify shoplifters in real time. Already adopted by over 4,000 stores across 25 countries, Veesion's technology significantly enhances security while reducing reliance on human monitoring. By integrating AI with existing camera infrastructure, the system provides an efficient and scalable solution for small and medium-sized retailers. As theft becomes an increasing concern, especially in urban centers, AI-powered surveillance like Veesion is emerging as a critical tool in loss prevention strategies.	By Nick Huber		March 27, 2025
4.34	Hate Instacart Substitutions? The Company's Finally Doing Something About That	Instacart has launched new AI-powered tools—Store View and Second Store Check—to improve order accuracy and reduce customer frustration from out-of-stock items. Store View uses videos captured by shoppers to create a virtual map of store shelves, helping identify product locations and gaps in inventory. Second Store Check automatically reroutes missing items to a nearby second store,	By Food & Wine		March 27, 2025




 AI Use Cases





#	Highlights	Summary	Author	Source	Date
		ensuring better fulfillment without extra work for customers. These innovations enhance inventory management, streamline the shopping process, and boost customer satisfaction. By leveraging AI, Instacart aims to make online grocery shopping more reliable, efficient, and responsive to real-time inventory changes.			
4.35	AI and satellites help aid workers respond to Myanmar earthquake damage	After a 7.7 magnitude earthquake struck Mandalay, Myanmar in March 2025, AI played a crucial role in disaster response. Microsoft's AI for Good Lab collaborated with the United Nations Satellite Centre and Carnegie Mellon University to analyze satellite imagery using AI models. This technology rapidly assessed building damage across affected regions, identifying over 400 damaged or destroyed structures. The AI-driven approach enabled precise, large-scale mapping within hours—something that would typically take days manually. This use of AI significantly improved the speed and accuracy of disaster assessments, ensuring that humanitarian aid could be more efficiently directed to the hardest-hit areas.	By Matt O'BRIEN		March 31, 2025
4.36	The AI Doctor Will See New York Ride-Share Drivers Now	Akido Labs is using AI, through its ScopeAI system, to provide medical care for ride-share drivers in New York City. The AI suggests diagnoses and treatments based on symptoms and medical histories, with human doctors making the final decisions. Partnering with the Independent Drivers Guild and Workers Benefit Fund, Akido plans to expand the service across the city, offering accessible healthcare without unpaid time off. The service is also expanding to other regions like California and Rhode Island. However, concerns about physician burnout arise as doctors face pressure to follow AI recommendations while remaining responsible for patient care.	By Brian Gormley		March 27, 2025




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
4.37	Amazon unveils Nova Act, an AI agent that can control a web browser	Amazon is developing AI agents called “Nova” to perform tasks like web browsing, offering shopping suggestions, handling customer service, and scanning content. These agents can execute multi-step actions by leveraging large language models to understand user intent and complete tasks step by step. Amazon aims to eventually integrate Nova with Alexa, creating a more interactive and capable digital assistant. The initiative highlights Amazon’s push to embed AI-powered tools directly into consumer experiences. Currently, Nova is being tested with a limited group of users, signaling early steps toward broader deployment of AI agents in everyday digital environments.	By Maxwell Zeff		March 31, 2025
4.38	Chinese brain chip project speeds up human trials after first success	A Chinese research team has accelerated human trials of its brain-computer interface (BCI) project after successfully implanting a brain chip into a patient suffering from a spinal cord injury. The chip, developed by Beijing-based Neucyber, enabled the patient to control a robotic arm using brain signals. This marks China’s most advanced step in neurotechnology, aiming to compete with companies like Elon Musk’s Neuralink. The project is part of China’s broader strategy to advance in AI-driven medical technologies. Following the initial success, researchers plan to expand trials, highlighting rapid progress in the country’s neuroscience and biotechnology sectors.	By Eduardo Baptista		March 31, 2025
4.39	China Accelerates Human Trials in Homegrown Brain-Computer Interface Project	A Chinese research team is advancing its brain-computer interface (BCI) project, moving swiftly into more human trials after initial success. The homegrown BCI implant, developed by NTC Smart Brain, allowed a patient with spinal injury to control a robotic arm using brain signals—marking a breakthrough in neural technology. This initiative, seen as China’s response to Elon Musk’s Neuralink, reflects the country’s ambition to become a leader in neurotechnology. The trials	By Eduardo Baptista		March 31, 2025


✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		will expand in 2025, focusing on restoring mobility and communication for patients with severe neurological conditions.			
4.40	DeepSearch on Private Visual Documents: An Enterprise Case Study	Jina AI's DeepSearch revolutionizes enterprise search by leveraging AI-driven reasoning to retrieve, read, and synthesize accurate answers from weakly or unstructured data. Unlike traditional systems, it iteratively refines queries, evaluates results, and extracts precise information, saving time and reducing errors. DeepSearch excels in multimodal information retrieval, handling text, diagrams, and multilingual queries seamlessly. Its ability to provide direct answers, supported by transparent document links, enhances productivity for businesses with complex data. By integrating with existing systems, DeepSearch offers powerful AI capabilities without costly infrastructure changes, making it essential for modern enterprises.	By Maximilian Werk, Scott Martens		March 31, 2025
4.41	Inference-Time Scaling for Complex Tasks: Where We Stand and What Lies Ahead	The paper explores inference-time scaling techniques to enhance reasoning capabilities of large language models (LLMs) across eight complex tasks, including math, STEM reasoning, NP-hard problems, and spatial reasoning. By evaluating nine state-of-the-art models, it highlights the benefits and limitations of scaling methods, revealing diminishing returns as task complexity increases. Key findings include variability in token efficiency, the importance of robust verifiers, and gains from superscaling. The study emphasizes the need for adaptive token allocation and improved verification mechanisms, offering critical insights into advancing LLM reasoning performance for challenging real-world applications.	By Microsoft Research		March 31, 2025
4.42	SimpleRL-Zoo: Investigating and	The paper, SimpleRL-Zoo: Investigating and Taming Zero Reinforcement Learning for Open Base Models in the Wild, explores the effectiveness of zero	By HKUST, TikTok, BUPT		March 24, 2025




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
	Taming Zero Reinforcement Learning for Open Base Models in the Wild	reinforcement learning (RL) training across ten diverse base models. By leveraging rule-based rewards and adjusting training strategies, the authors achieve significant improvements in reasoning accuracy and response length. Notably, they observe the emergence of cognitive behaviors, such as verification, in smaller models outside the Qwen family. The findings highlight the potential of zero RL training to enhance reasoning capabilities in various models, contributing valuable insights and open-source resources for future research.			
4.43	Defeating Prompt Injections by Design	The paper "Defeating Prompt Injections by Design" introduces CaMeL, a robust defense mechanism against prompt injection attacks in Large Language Models (LLMs). Inspired by traditional software security concepts like Control Flow Integrity and Information Flow Control, CaMeL enforces fine-grained security policies by tracking control and data flows using a custom Python interpreter. Unlike existing approaches, it secures LLM-based systems without modifying the model itself. Demonstrated on the AgentDojo benchmark, CaMeL achieves 67% task success with provable security guarantees, significantly reducing vulnerabilities. This work highlights a critical step toward building secure and reliable LLM-powered systems in real-world applications.	By Google DeepMind, ETH Zurich		March 24, 2025

🛡️ AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
5.1	EU Halts AI Liability Law Amid Disagreements	The European Commission withdrew its AI Liability Directive, citing a lack of stakeholder agreement. Intended to establish AI accountability, its withdrawal followed criticism that it stifled innovation. EU officials will now reassess AI accountability, considering alternative frameworks. The decision sparked debate: some see it as a consumer protection setback, others as pragmatic to avoid premature regulation. This highlights the innovation-oversight tension. Without the directive, questions remain on legal redress for AI-caused damages. The EU's next regulatory move is highly anticipated, shaping AI confidence across member states.	By Tim Wright, Nathan Evans		March 7, 2025
5.2	Global AI Summit Reveals Governance Divide	The Paris AI Summit in March 2025 exposed a global AI governance split. 57 nations, including China and India, signed a declaration for "inclusive" AI, emphasizing ethics. The U.S. and UK declined, citing security concerns and lack of clarity. Yoshua Bengio's AI Safety Report 2025 offered risk mitigation blueprints. The UK rebranded its AI Safety Institute to AI Security Institute, prioritizing national security. This highlights a divide: collective governance versus sovereign control. Reconciling these views is crucial for global AI innovation and trust.	By Tim Wright, Nathan Evans		March 7, 2025
5.3	Anthropic Secures \$3.5B, Expands Claude	On March 3, 2025, Anthropic secured \$3.5 billion in Series E funding, reaching a \$61.5 billion valuation. Led by Lightspeed Venture Partners and supported by investors like Salesforce Ventures, Cisco Investments, and Fidelity, the funds will enhance Anthropic's next-generation AI development, computing capacity, interpretability research, and global expansion. Recent launches, Claude 3.7 Sonnet and Claude Code, significantly improved AI-driven coding capabilities. Claude is now integrated into platforms like Replit's Agent and Thomson Reuters'	By Anthropic Team		March 3, 2025



 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		CoCounsel, boosting productivity across industries, including healthcare and finance. Novo Nordisk notably reduced clinical report-writing time using Claude.			
5.4	Trump Shifts AI Policy Focus	The Trump Administration has shifted U.S. AI policy to prioritize innovation and leadership over new regulation. A January 23 executive order, <i>“Removing Barriers to American Leadership in AI,”</i> revoked earlier AI directives and launched an AI Action Plan aimed at sustaining U.S. dominance in AI. The White House is actively soliciting public input on this plan through mid-March. Meanwhile, state-level lawmakers are stepping in: Virginia recently passed a comprehensive AI law to curb algorithmic bias in high-risk systems and ensure transparency, making it one of the first U.S. states with broad AI legislation.	By Michael Charalambous		March 6, 2025
5.5	China Boosts Strategic AI Investment	AI remains a strategic priority in China’s latest government agenda. At the National People’s Congress on March 5, Beijing announced plans to boost support for AI R&D – backing large-scale AI models, “industries of the future” (e.g. embodied AI and 6G), and venture funding – to drive tech breakthroughs and self-reliance. This marks the first time China’s annual work report explicitly mentioned AI models, following the global buzz around a Chinese startup’s advanced AI system. Regulators in China continue to enforce strict rules on deepfakes and generative AI content, but the government’s message in 2025 underscores an <i>“innovation-friendly”</i> approach alongside security oversight.	By Reuters		March 5, 2025
5.6	UK Delays AI Legislation, Debate Continues	The UK government has paused plans for sweeping AI-specific legislation, delaying the anticipated national AI bill until at least summer 2025, potentially aligning with the U.S.’s lighter regulatory approach. Meanwhile, a House of Lords member reintroduced a private AI regulation bill on March 4, sparking debate on AI	By Michael Charalambous		March 6, 2025

 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		governance. Internationally, countries like Canada and Japan have joined the Council of Europe's new AI Convention, the first global AI treaty. Switzerland plans to ratify this convention, aiming to balance innovation with human rights and AI transparency safeguards.			
5.7	OpenAI Faces Criticism on Alignment	In March 2025, OpenAI reaffirmed its commitment to AI safety and ethics in a blog post, outlining efforts to ensure future frontier AI systems are controllable and beneficial. However, Miles Brundage, a former OpenAI policy lead, criticized the post, accusing the company of rewriting the history of its AI safety journey. He argued that OpenAI downplayed past internal debates, such as concerns over releasing GPT-2 in 2019, and painted an overly positive picture of its vigilance. This incident highlights ongoing tensions in the AI industry between rapid development and transparency.	By ResearchBuzz		March 10, 2025
5.8	Bias Mitigation & Fairness:	Recent regulations, such as Virginia's High-Risk Artificial Intelligence Developer and Deployer Act (HB 2094), require developers to prevent "algorithmic discrimination" by auditing AI systems for bias. Similarly, global frameworks like the EU AI Act and ISO/IEC 42001 emphasize robust data governance, advocating for high-quality, representative datasets and comprehensive bias testing. Standardized datasets are emerging to evaluate facial recognition technologies across diverse demographics, aiding in quantifying ethical performance. Organizations increasingly recognize fairness audits as essential, adopting practices akin to security evaluations. Despite ongoing challenges, a consensus is	By Reena R. Bajowala, Wouter van Wengen of Greenberg Traurig		March 5, 2025




 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		building around routinely embedding fairness evaluations into AI development processes to mitigate algorithmic discrimination.			
5.9	UK Government Advocates AI Integration in Civil Service	Prime Minister Sir Keir Starmer has announced plans to integrate AI into the UK civil service to boost efficiency and cut administrative costs. Aiming to reduce regulatory burdens on businesses by 25%, the initiative will shift one in ten civil servants into technical and digital roles within five years, streamlining Whitehall operations. This transformation is projected to save taxpayers up to £45 billion, though concerns remain over legislative hurdles and potential job losses. Starmer stresses that AI should handle tasks where possible, ensuring modernized governance.	By Oliver Wright		March 12, 2025
5.10	SAG-AFTRA Raises Concerns Over AI Usage in Video Game Industry Amid Ongoing Strike	Members of SAG-AFTRA have been on strike for over seven months, addressing job security concerns related to AI replacing human roles in major video game studios. The union recently highlighted "alarming loopholes" in the latest bargaining proposal, which could allow companies to reuse union members' past work without consent. SAG-AFTRA urges its members to reject such projects, emphasizing the need for negotiations to protect against AI misuse. This action follows a prior strike focusing on AI's impact in television roles.	By Jacqueline Rayfield		March 12, 2025
5.11	From bureaucracy to brilliance: AI in federal IT	Artificial intelligence (AI) is poised to revolutionize federal agencies by enhancing citizen services, improving decision-making, and bolstering national security. To fully harness AI's potential, a comprehensive strategy is essential, encompassing infrastructure modernization, robust security measures, and workforce development. Key trends include increasing AI fluency among federal employees, investing in sovereign AI to maintain data control, and understanding agentic AI	By John Roesse		March 12, 2025

🛡️ AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		for automating complex tasks. Adopting a holistic approach ensures that federal IT remains adaptable and resilient amidst rapid technological advancements.			
5.12	OpenAI Urges Trump Administration to Prioritize AI Growth Over Regulation	OpenAI is lobbying the Trump administration to focus on AI development speed while maintaining minimal regulatory constraints. The company warns that excessive oversight could stifle innovation, giving global competitors an advantage. Instead, OpenAI advocates for flexible policies that encourage rapid AI advancements while managing risks through industry-driven standards. Executives stress that the U.S. must lead in AI to maintain technological dominance. However, critics argue that weak regulations could lead to ethical issues, monopolization, and potential misuse. The debate reflects ongoing tensions between innovation, competition, and responsible oversight in the AI sector.	By Hayden Field		March 13, 2025
5.13	Infosys Research Calls for AI Governance Task Force	As organizations transition from AI experimentation to scaled deployment, Infosys Knowledge Institute recommends creating dedicated AI governance task forces to reduce risk and improve accountability. Their latest survey, "Infosys AI Business Value Radar," reveals that only 19% of AI use cases fully deliver on business objectives, while another 32% partially meet goals. The research found that white-collar and technically focused industries achieve greater AI success, while sectors like travel, manufacturing, retail, and public services struggle with consistent implementation. The report outlines five critical steps for organizations to generate business value from AI deployments and become AI-first enterprises	By Uma Kannan		March 13, 2025
5.14	Google Outlines Three-Pronged	Google has submitted policy recommendations to the U.S. Office of Science and Technology Policy's Request for Information, outlining a framework to secure	By Kent Walker		March 13, 2025





🛡️ AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
	Approach for U.S. AI Action Plan	America's position as an AI leader. The proposal focuses on three key areas: investing in AI infrastructure, accelerating government AI adoption, and promoting pro-innovation approaches internationally. Google advocates for policy reforms addressing energy needs, balanced export controls, and preemption of fragmented state-level rules. The company emphasizes that policy decisions will significantly shape the outcome of global AI competition, calling for approaches that protect national security while enabling widespread benefits from AI advancements			
5.15	USPTO in 2025: Leadership, Policy, and Legislative Changes to Watch	The U.S. Patent and Trademark Office (USPTO) is undergoing significant changes in 2025, with new leadership and policy updates. Coke Morgan Stewart was appointed interim director, followed by Howard Lutnick as Commerce Secretary. John Squires has been nominated as USPTO's new director. The administration aims to reduce the patent backlog, but hiring freezes and return-to-office policies may hinder progress. Legislative efforts, including the Patent Eligibility Restoration Act (PERA) and PTAB reforms, are gaining attention. AI-related policy updates are also in focus. These changes will impact patent applicants and.. require strategic adjustments	By Maria E. Doukas, Alexander B. Stein, Jacob L. Peterson		March 13, 2025
5.16	US Expands AI Chip Export Controls Amid Rising Tech War with China	The Biden administration has imposed new export controls on AI chips to limit China's access to advanced technology. These regulations aim to prevent China from using AI for military and surveillance purposes, reinforcing US national security policies. The rules categorize countries into tiers, with China facing the strictest restrictions. The policy reflects the growing role of AI in global power struggles, balancing innovation with security concerns. While US companies face	By Dashveenjit Kaur		March 14, 2025




🛡️ AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		challenges due to market restrictions, the move aligns with broader AI strategies to maintain technological leadership.			
5.17	AI Giants Must Not Exploit British Creatives, Industry Leaders Warn	British media leaders and creative industry executives strongly oppose the UK government's proposed copyright reforms, which would let AI firms use copyrighted content freely unless creators opt out. Executives from NewsUK, The Guardian, Warner Music, and Channel 4 warned Technology Secretary Peter Kyle that this could destabilize the £125 billion creative sector. They argue that few countries allow unrestricted AI access to copyrighted works and that fair licensing deals exist without government interference. Figures like Sir Elton John, Simon Cowell, and Sir Paul McCartney stress transparency over sweeping changes, fearing AI firms could exploit UK intellectual property unfairly.	By, Martina Bet		March 14, 2025
5.18	An Integrated Approach is Essential for Testing AI Semiconductors	The rapid expansion of artificial intelligence (AI) applications is driving increased demand for specialized AI semiconductors. These include GPUs for cloud computing, dedicated AI chips for efficiency, and neuromorphic chips for low-power edge applications. Market forecasts predict substantial growth, with revenues rising from \$53.4 billion in 2023 to \$341 billion by 2033. This diversity and complexity present unique testing challenges, such as increased scan-pattern depths and the need for efficient test data distribution. To address these issues, test processes must become dynamically adaptable, incorporating real-time decision-making to optimize cost, yield, and quality. This integrated approach is crucial for effectively testing AI semiconductors.	By IRA LEVENTHAL		March 11, 2025





🛡️ AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
5.19	A Dynamic Governance Model for AI	As AI transitions from an efficiency tool to a force shaping policy and governance, technologists are emerging as political actors. No longer just innovators, they now set agendas, draft guidelines, and influence AI's future—once the domain of elected officials. This shift raises concerns about democratic resilience and public oversight. AI's impact on economies, information, and security is shifting decision-making from governments to corporations. To balance innovation with public interest, governance structures must prevent control by a few dominant players while ensuring AI serves society democratically and ethically. Addressing these challenges is crucial for AI's responsible development.	By Paulo Carvão, Yam Atir, Salvina Ancheva		March 13, 2025
5.20	CCI Chief Raises Alarm on AI Collusion, Urges Proactive Regulations	The Competition Commission of India (CCI) Chief, Ravneet Kaur, has warned about the risks of AI-driven collusion, highlighting concerns such as algorithmic price-fixing, cartelization without human communication, and discriminatory pricing under the guise of dynamic pricing. She emphasized the need for forward-looking regulations that ensure trust and fairness in AI-driven markets. As AI increasingly shapes market dynamics, regulators must adapt to prevent anti-competitive practices. The CCI is currently conducting a study on AI's impact on competition, aiming to develop regulatory frameworks that balance innovation with market integrity.	By Rediff Money Desk		March 16, 2025
5.21	Europol Warns of AI-Driven Crime Threats	Europol has issued a warning about the growing threats posed by AI-driven crime, highlighting how organized crime gangs are leveraging advanced technologies to enhance their global operations. The agency's European Serious Organised Crime Threat Assessment report notes that AI allows these criminals to craft multilingual messages, create realistic impersonations, and automate processes, complicating detection efforts. The report also highlights the potential for fully autonomous AI-	By Reuters		March 18, 2025



 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		controlled criminal networks in the future, underscoring the need for robust regulatory frameworks and international cooperation to combat these evolving threats.			
5.22	OpenAI Proposes 'Freedom-Focused' AI Policy to White House	OpenAI has submitted a "freedom-focused" policy proposal to the White House's AI Action Plan, advocating measures to maintain U.S. AI dominance over China. Key recommendations include leveraging trade laws, reducing copyright restrictions for training data, investing in AI infrastructure, and preventing states from enacting restrictive AI laws, particularly targeting California's defeated SB 1047 bill. This move comes amid a surge in state-level AI legislation, with 893 bills introduced across 48 states in less than 80 days of 2025. OpenAI's stance reflects concerns over navigating a patchwork of state regulations that could hinder innovation.	By Tina Nguyen		March 19,2025
5.23	EU Charges Google and Apple Under Digital Markets Act	The European Commission has initiated actions against Google and Apple under the EU Digital Markets Act, alleging violations that favor their own services and restrict competition. Google is accused of promoting its services in search results and limiting developers' ability to direct users to alternative platforms. Apple has been ordered to improve interoperability with rival devices. These measures aim to foster fair competition and innovation in the tech industry. However, they risk straining EU-U.S. relations, with President Donald Trump threatening retaliatory tariffs.	By Rob Davies		March 19,2025




🛡️ AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
5.24	Commerce Chief Seeks Help to Block China's Access to US Chip	U.S. Commerce Secretary Howard Lutnick has urged companies and foreign governments to help prevent China from obtaining U.S. semiconductor chips, asserting that Chinese AI company DeepSeek improperly used American-made chips. He emphasized concerns over national security and the potential consequences of losing access to Taiwanese chips. Lutnick called for export controls to be included in trade deals, encouraging countries to choose between aligning with the U.S. and Western values or prioritizing financial gain by aiding China. He also discussed efforts to boost U.S. production of critical goods like aluminum, steel, and semiconductors.	By Reuters		March 18, 2025
5.25	AI model providers see EU Commission build 'network of evaluators'	The European (EU) Commission is building a network of model evaluators to define how general-purpose AI models with systemic risk should be evaluated in accordance with the legal requirements of the AI Act and the GPAI code of practice. While the role of this network still needs to be formalized with an implementing act due in August, the participants might also provide a shortlist of "qualified" independent reviewers that AI companies and the AI Office might tap to carry out model evaluations on their behalf.	By Luca Bertuzzi		March 19, 2025
5.26	X's Use of Personal Data for Grok AI Training Investigation	The Swiss Federal Data Protection and Information Commissioner (FDPIC) concluded its preliminary investigation into X's use of personal data to train its AI model, Grok. The probe revealed that X processed personal data without obtaining users' consent, breaching Swiss data protection laws. Consequently, users now have the right to object to their public posts being utilized for AI training. X has been advised to implement measures ensuring compliance with data protection regulations and to secure explicit user consent for data processing activities related to AI development.	By Swiss federal authorities		March 20, 2025





 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
5.27	Commerce chief seeks industry help to prevent China from getting US chips	U.S. Commerce Secretary Gina Raimondo has urged the tech industry to assist in preventing China from accessing advanced American AI chips. Speaking at an event in Washington, she emphasized that national security depends on stopping sensitive technologies from reaching adversaries. Raimondo called on chipmakers and cloud providers to enhance internal controls and share information with the government. The Commerce Department has already restricted AI chip exports to China, but concerns remain about indirect access through cloud services. Raimondo stressed that maintaining America’s lead in AI requires a collaborative effort between government and the private sector.	By David Shepardson		March 18, 2025
5.28	CAC issues Security Management Measures for Facial Recognition Technology Applications	China’s Cyberspace Administration (CAC) and Ministry of Public Security have issued new rules for facial recognition technology, effective June 1, 2025. The regulations require data collection to be purpose-specific and minimal, mandate separate user consent, and prohibit internet transmission without approval. High-risk uses need impact assessments, and storing over 100,000 facial templates requires CAC filing. Security protocols ban using facial recognition as the sole verification method when alternatives exist, restrict scans in private spaces, and forbid coercive use. Organizations must also implement encryption, access controls, and auditing to safeguard biometric data and ensure compliance.	By Cyberspace Administration of China (CAC)		March 21, 2025
5.29	UK minister in US to pitch Britain as global AI investment hub	UK Technology Minister Saqib Bhatti visited the US to position Britain as a global hub for AI investment. Highlighting the UK’s “pro-innovation, pro-safety” regulatory approach, Bhatti emphasized a flexible framework that fosters AI development while ensuring ethical oversight. He pitched the UK’s strengths, including access to high-quality data, leading universities, and strong public-private partnerships. The UK aims to distinguish itself from more rigid regulatory	By Ryan Daws		March 20, 2025




 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		environments, appealing to global tech firms and investors. This move is part of Britain's broader strategy to lead in safe and responsible AI innovation on the world stage.			
5.30	UK's AI Copyright Proposal Sparks Debate	UK Technology Secretary Peter Kyle has urged critics of the nation's proposed AI copyright regime to embrace change rather than resist it. Addressing concerns from creative industries about potential impacts on intellectual property rights, Kyle emphasized the importance of adapting to technological advancements while ensuring fair compensation for creators. The proposed framework aims to balance innovation with the protection of artistic works, acknowledging the transformative potential of AI in various sectors. Kyle's stance highlights the government's commitment to fostering an environment where technological progress and creative rights coexist harmoniously.	By Financial Times		March 24,2025
5.31	AI's Potential Economic Impact	Bank of England Governor Andrew Bailey has renewed calls for coordinated international efforts to address escalating global trade tensions, highlighting particular concerns regarding technology and artificial intelligence. Bailey emphasized the importance of harmonized regulatory standards and collaborative policy frameworks to prevent fragmentation and economic instability stemming from tech-driven disputes. He warned that uncoordinated AI policies and protective economic measures could significantly disrupt global markets. His statements underscore the increasing urgency for unified global responses to technological advancements and associated economic challenges, advocating for cooperation to maintain international economic stability amid rapid technological innovation.	By David Milliken		March 24,2025

 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
5.32	OpenAI Expands COO Brad Lightcap's Role to Drive Global Growth and Partnerships	OpenAI COO Brad Lightcap will take on an expanded role to drive global growth and partnerships, according to CEO Sam Altman. This strategic move reflects OpenAI's ambition to deepen its international footprint and forge stronger alliances across industries and governments. Lightcap's leadership will focus on operational scaling, partnership development, and policy engagement, signaling OpenAI's intent to balance innovation with regulatory collaboration. As global scrutiny of AI intensifies, OpenAI is positioning itself to navigate geopolitical dynamics while expanding responsible deployment of its technologies across borders.	By Anna Tong		March 24,2025
5.33	Sam Altman's Worldcoin in Talks with Visa for Stablecoin Wallet Integration	Sam Altman's Worldcoin project is reportedly in discussions with Visa to integrate its stablecoin wallet into Visa's global payment network, according to CoinDesk. The potential partnership would allow users to make everyday purchases with Worldcoin's digital currency, bridging traditional financial systems with blockchain-based identity and currency tools. This move reflects Worldcoin's strategy to expand access to digital finance and strengthen its real-world utility, while navigating evolving regulatory frameworks around stablecoins and digital identity.	By Reuters		March 24,2025
5.34	Former Intel CEO Pat Gelsinger Joins Faith-Based Tech Firm Gloo for AI Expansion	Former Intel CEO Pat Gelsinger has joined Gloo, a Colorado-based technology company focused on serving faith-based communities, to support its AI-driven outreach and engagement tools. Gloo provides churches and nonprofits with AI-enabled platforms for personalized communication, community-building, and data-driven ministry. Gelsinger's involvement aims to accelerate the firm's AI capabilities while aligning with his personal faith values. This move illustrates a	By Jeffrey Dastin		March 24,2025





 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		growing trend of AI adoption in niche sectors like religious outreach, showing how AI tools are being tailored to support values-driven missions.			
5.35	US-China Committee Head Urges Cooperation on AI to Avoid Conflict	<p>The head of the US-China Relations Committee emphasized the urgent need for bilateral cooperation on artificial intelligence to mitigate risks and prevent potential conflicts. Speaking amid rising geopolitical and technological tensions, he highlighted that both nations have a shared responsibility to establish guardrails around AI development. The call for dialogue includes proposals for joint frameworks on safety, transparency, and ethical AI use. As global AI capabilities rapidly evolve, the appeal reinforces the importance of diplomacy and governance to ensure AI advances do not escalate international instability.</p>	By Reuters		March 24, 2025
5.36	Sweden Proposes Bill to Let Police Use AI Facial Recognition	<p>The Swedish government has proposed a bill allowing law enforcement to use AI-powered facial recognition technologies for criminal investigations. The move marks a significant shift in surveillance policy, aiming to enhance police efficiency while sparking debate over privacy and civil liberties. If passed, the legislation would permit real-time biometric analysis in public spaces, aligning Sweden with other European nations adopting AI in policing. Authorities emphasize safeguards and oversight mechanisms to prevent misuse, as the country navigates the balance between public safety and ethical AI deployment.</p>	By Reuters		March 20, 2025



🛡️ AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
5.37	In the Market: How Trump is driving Asia to diversify away from US	Fears over a potential second Trump presidency and his past protectionist trade policies are pushing Asian countries to reduce dependence on the U.S. Experts say nations like Japan, South Korea, and others are accelerating efforts to diversify supply chains and strengthen regional partnerships. This shift includes critical sectors like semiconductors and advanced tech—key components in AI development. As geopolitical uncertainty grows, Asia’s move toward economic self-reliance could reshape global AI infrastructure, investment flows, and hardware sourcing strategies, signaling a strategic pivot in how the region prepares for potential disruptions from U.S. policy shifts.	By Paritosh Bansal		March 20, 2025
5.38	Microsoft to Launch Three Data Centers in Malaysia by Q2 2025	Microsoft plans to open three data centers in Malaysia by the second quarter of 2025 as part of its broader strategy to expand cloud and AI infrastructure in Southeast Asia. The initiative supports Malaysia’s digital transformation goals and aligns with Microsoft's commitment to responsible AI deployment and regional capacity building. These data centers will provide local businesses and governments with improved access to cloud services, AI tools, and data sovereignty. The move also reflects Microsoft’s growing investments in AI-aligned infrastructure worldwide.	By Rozanna Latiff		March 20, 2025
5.39	Swedish government proposes bill to allow police to use AI face-recognition	The Swedish government has introduced a bill that would allow police to use AI-powered facial recognition to investigate serious crimes and locate missing persons. The proposed legislation aims to balance public safety with privacy concerns and will be aligned with the European Union’s AI Act. If passed, it would mark a significant step in Sweden’s adoption of AI within law enforcement. The	By Reuters		March 20, 2025




 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		move highlights growing efforts across Europe to regulate the use of AI in sensitive areas like surveillance, raising important debates around ethics, oversight, and civil liberties.			
5.40	EU lawmakers warn against 'dangerous' moves to water down AI rules	Architects of the EU's landmark AI Act have warned Brussels against weakening the regulation at the last minute, calling such moves "dangerous." The AI Act, aimed at regulating high-risk AI systems, is in its final stages before formal approval. However, recent lobbying efforts from major tech companies and some member states have pushed for looser rules, especially around foundation models. Lawmakers behind the Act urge the EU to maintain strong safeguards to ensure accountability and safety. The outcome will significantly influence global AI regulation and set a precedent for how governments oversee emerging technologies.	By Melissa Heikkilä, Barbara Moens		March 25, 2025
5.41	European Investors Urge AI Companies to Deliver Returns Amidst Market Pressures	European investors are demanding tangible AI returns by 2025 as companies continue heavy investments in generative AI. The launch of DeepSeek, a low-cost Chinese AI model, has increased market pressures, leading to stock declines in AI hardware suppliers like ASM International and BE Semiconductor. Meanwhile, AI adopters like RELX and SAP saw smaller declines, signaling a shift in investor preference. A Fidelity survey suggests minimal AI impact on profitability by 2025, though benefits are expected later. Investors now prioritize revenue-generating AI applications to sustain market confidence.	By Lucy Raitano		March 26, 2025
5.42	Anthropic Prevails in Initial Phase of AI	Anthropic, an AI company, has won an early legal battle against music publishers Universal Music Group (UMG), Concord, and ABKCO. A California federal judge denied the publishers' request for a preliminary injunction to prevent Anthropic	By Blake Brittain		March 26, 2025




 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
	Copyright Dispute with Music Publishers	from using song lyrics to train its AI chatbot, Claude. The judge ruled that the publishers' request was overly broad and failed to demonstrate "irreparable harm." The publishers, however, remain confident in their broader case and intend to pursue monetary damages. This case is part of a wider debate over the use of copyrighted materials in training AI models.			
5.43	Apple Announces WWDC 2025: iOS 19 and More Coming Soon	Apple has confirmed that its Worldwide Developers Conference (WWDC) will take place from June 9 to June 13, 2025. The event will be accessible online, highlighting updates to the software powering iPhones, iPads, and other Apple devices. Some developers and students will be invited to attend in person at Apple Park on the opening day. Anticipated announcements include iOS 19, featuring a significant design overhaul inspired by Apple's Vision Pro headset, and updates to other operating systems such as iPadOS, macOS, watchOS, and tvOS. Investors will be watching closely, as product enhancements announced during the event could help Apple attract new customers.	By Reuters		March 25, 2025
5.44	Cerebras IPO Faces Further Delay Amid Ongoing National Security Review	Cerebras Systems' initial public offering (IPO) has encountered additional delays due to an extended national security review by the Committee on Foreign Investment in the United States (CFIUS). The review centers on a \$335 million investment from Abu Dhabi-based G42, which has previous ties to China's Huawei, raising national security concerns. The process has been prolonged by unfilled key positions within the Trump administration, including the assistant Treasury secretary for investment security, who oversees CFIUS. This uncertainty hampers Cerebras' IPO plans, as investors seek clarity before proceeding. Despite these challenges, executives remain optimistic about eventual approval and the successful launch of the IPO.	By Echo Wang and Alexandra Alper		March 25, 2025


 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
5.45	Alibaba Resumes Hiring Amid Renewed Confidence Following Xi's Tech Summit	Alibaba Group Chairman Joe Tsai announced plans to restart hiring, citing increased confidence among China's business community following President Xi Jinping's February meeting with tech entrepreneurs, including Alibaba co-founder Jack Ma. This meeting, signaling a policy shift, encourages businesses to reinvest and expand after years of regulatory pressures that dampened investment and led to widespread layoffs. Tsai noted that Alibaba's workforce had declined for 12 consecutive quarters, but he now believes the company has reached the bottom and will begin rehiring. He emphasized that hiring would provide job security and income growth, translating business confidence into consumer confidence.	By Kane Wu and Selena Li		March 26, 2025
5.46	U.S. Expands Export Blacklist Targeting Chinese Tech Firms	The U.S. Department of Commerce has added over 50 Chinese entities to its export blacklist, aiming to impede China's progress in artificial intelligence, quantum computing, and military technologies. This action includes six subsidiaries of Inspur Group, China's leading cloud computing provider, which was previously blacklisted in 2023. Companies on this list require special licenses to receive U.S. exports, with approvals often denied. The Chinese Embassy has condemned these measures, accusing the U.S. of politicizing trade and technology issues. This move intensifies existing tensions between the two nations over technology and national security concerns.	By Taylor Herzlich		March 26, 2025
5.47	CoreWeave Adjusts IPO Amid Market Challenges	CoreWeave, a cloud computing company backed by Nvidia, has adjusted its IPO plans due to market conditions. It will now offer 37.5 million shares at \$40 each, raising \$1.5 billion and valuing the company at \$23 billion—a 23.5% reduction from initial plans. Concerns over CoreWeave's reliance on Microsoft and its capital-intensive model have influenced investor sentiment. With shares priced	By Echo Wang		March 28, 2025

 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		below the original \$47–\$55 range, the IPO is seen as a key market test for AI infrastructure investments amid shifting investor confidence.			
5.48	OpenAI Projects Significant Revenue Growth Amidst Industry Competition	OpenAI, the creator of ChatGPT, projects \$12.7 billion in revenue this year—tripling last year’s \$3.7 billion—driven by paid AI subscriptions. Looking ahead, it forecasts \$29.4 billion in 2026 but doesn’t expect positive cash flow until 2029 due to heavy AI investments in chips, data centers, and talent. Meanwhile, Chinese firms like DeepSeek are rapidly advancing, narrowing the AI gap with Western competitors. OpenAI’s growth underscores rising AI adoption, but intense competition and high costs remain key challenges in the evolving AI landscape.	By Bailey Lipschultz and Shirin Ghaffary		March 26, 2025
5.49	Microsoft Adjusts Data Center Strategy Amidst AI Demand Fluctuations	Microsoft has revised its data center expansion plans, canceling projects totaling 2 gigawatts of electricity usage in the U.S. and Europe over the past six months. This strategic shift is primarily due to an oversupply relative to current demand forecasts, influenced by reduced support for additional AI training workloads from OpenAI. Consequently, competitors like Google and Meta have seized some of the freed capacity in international markets and the U.S., respectively. Despite these adjustments, Microsoft remains committed to investing \$80 billion in AI infrastructure this fiscal year, emphasizing continued growth across all regions.	By Georgia Butler		March 27, 2025
5.50	Prysmian Projects Significant Profit Growth by 2028 Amid AI Data Center Boom	Prysmian Group projects a 64% rise in core profit by 2028, targeting up to €3.15 billion in adjusted EBITDA, fueled by U.S. AI data center growth. To strengthen its digital solutions portfolio, Prysmian will acquire U.S. connectivity firm Channell for up to \$1.15 billion. CEO Massimo Battaini emphasized the move’s role in capitalizing on AI-driven digitalization, as data centers’ U.S. power demand is	By Giulio Piovaccari		March 26, 2025

🛡️ AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		expected to rise from 6% to 14% by 2030. This acquisition enhances Prysmian's position in the expanding data infrastructure market.			
5.51	Chip startup Retym raises \$75 million to build chip that will AI connect data centers	Retym, a semiconductor startup, has raised \$75 million in its latest round, bringing total funding to \$180 million. The company develops digital signal processing (DSP) chips to boost data transfer speeds in AI-driven data centers. Its first chip, optimized for 30-40 km transmissions (ranging 10-120 km), uses advanced modulation for data integrity and is built on TSMC's 5nm process. This innovation aims to ease AI infrastructure bottlenecks. Backed by Spark Capital, the Series D funding will help Retym launch its first product this year.	By Max A. Cherney		March 31, 2025
5.52	China's Zhipu AI launches free AI agent, intensifying domestic tech race	Chinese AI startup Zhipu AI has launched AutoGLM Ruminant, a free AI agent capable of tasks such as web searches, travel planning, and research report writing. Powered by Zhipu's proprietary models, GLM-Z1-Air and GLM-4-Air-0414, the agent matches rival DeepSeek's R1 in performance while operating up to eight times faster and requiring significantly fewer computing resources. Unlike competitor Manus, which charges up to \$199 monthly, AutoGLM Ruminant is available free via Zhipu's official channels. Founded in 2019 as a Tsinghua University spinoff, Zhipu AI recently secured substantial government-backed funding, including a 300 million yuan (\$41.5 million) investment from Chengdu.	By Reuters		March 31, 2025

🛡️ AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
5.53	Exclusive: Arm expects its share of data center CPU market sales to rocket to 50% this year	Arm Holdings projects its share of the global data center CPU market will soar to 50% by the end of 2025, a significant rise from approximately 15% in 2024. This growth is fueled by the increasing demand for artificial intelligence (AI) applications. Arm's CPUs, known for their lower power consumption compared to Intel and AMD, are becoming the preferred choice in AI data centers, which require substantial energy. Major cloud providers like Amazon, Google, and Microsoft have adopted Arm-based chips, reflecting a shift towards more energy-efficient processing solutions.	By Max A. Cherney		March 31, 2025
5.54	Musk's social media firm X bought by his AI company, valued at \$33 billion	Elon Musk's AI company, xAI, has acquired the social media platform X (formerly Twitter) in an all-stock deal valuing X at \$33 billion, with the total transaction reaching \$45 billion, including \$12 billion in debt. The acquisition aims to merge data, computing resources, AI models, and distribution, potentially strengthening xAI's chatbot, Grok. With this deal, xAI's valuation has climbed to \$80 billion. The strategic integration could enhance AI-driven features on X, though details on leadership changes and regulatory scrutiny remain uncertain. This move reflects Musk's vision of creating a deeply integrated AI-powered digital ecosystem.	By Greg Bensinger		March 29, 2025
5.55	Anthropic says chatbot AI training makes fair use of books	Anthropic, an artificial intelligence company, has requested a California federal court to dismiss a lawsuit filed by authors Andrea Bartz, Charles Graeber, and Kirk Johnson. The authors allege that Anthropic infringed their copyrights by using their works to train its language model, Claude. Anthropic contends that this use constitutes fair use, as it transforms the original content into a new technological application rather than replicating the authors' work. The authors seek to represent a broader class of writers whose works were allegedly misused. The	By Blake Brittain		March 28, 2025

 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		case's outcome may hinge on the fair use doctrine, central to ongoing AI-related copyright disputes.			
5.56	Wall Street tumbles as fresh data fuels inflation fear	Wall Street saw a sharp decline as inflation worries grew amid rising tariff tensions. The S&P 500 fell 1.97% to 5,580.94, the Nasdaq dropped 2.70% to 17,322.99, and the Dow Jones slid 1.69% to 41,583.90. Tech giants Apple, Microsoft, and Amazon suffered notable losses. February's U.S. consumer spending rebounded less than expected, while core inflation hit a 13-month high. A University of Michigan survey showed 12-month inflation expectations at their highest in over two years. Fears mount that the Trump administration's tariffs may further drive inflation, potentially impacting Federal Reserve interest rate decisions.	By Noel Randewich and Pranav Kashyap		March 29, 2025
5.57	Scale AI seeking valuation as high as \$25 billion in potential tender offer, Business Insider reports	Artificial intelligence startup Scale AI is exploring a potential tender offer that could value the company at up to \$25 billion , aiming to capitalize on the surging demand for AI technologies. Founded in 2016 and backed by tech giants Nvidia , Amazon , and Meta , Scale AI specializes in providing accurately labeled data essential for training advanced AI models like OpenAI's ChatGPT . In a funding round last year, the company was valued at nearly \$14 billion . Currently, the U.S. Department of Labor is investigating Scale AI for compliance with the Fair Labor Standards Act .	By Reuters		March 29, 2025

AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
5.58	Trump tariffs triggered big Q1 plunge in market values for top global firms	Global stock markets have soared to a record \$110 trillion in total market capitalization, fueled by accelerating investment in artificial intelligence. Major U.S. tech companies—especially Nvidia, which jumped 82% in Q1 2024—have led the surge, with investors betting on AI’s potential to revolutionize productivity, business models, and earnings. The S&P 500 and Nasdaq have hit historic highs, while global valuations have nearly doubled since October 2022. Though analysts caution about rising market concentration and valuation risks, they emphasize that AI remains a key driver of optimism and structural economic change.	By Reuters		April 1, 2025

Conclusion: Implications and Future Trajectory

- The March 2025 developments reveal an AI ecosystem experiencing remarkable growth across multiple dimensions, with models becoming more powerful, specialized, efficient, and accessible than ever before.
- The rise of reasoning-focused models like DeepSeek-R1, OLMo 2, and various "R1" variants signals a new phase in AI development, where systems are increasingly capable of tackling problems requiring deep logical thinking and complex decision processes.
- Hardware innovations and infrastructure investments highlight the critical importance of computational resources in advancing AI capabilities, with chip designs, manufacturing capacity, and energy efficiency emerging as key competitive differentiators.
- The rapid advancement of Chinese AI models and chip technologies demonstrates the increasingly global nature of AI innovation, with important implications for technological competition, collaboration, and regulation.
- The growing focus on efficient architectures and training methods—including Mamba-Transformer hybrids, small but capable models, and optimized inference techniques—reflects an industry-wide emphasis on making AI more sustainable and accessible.
- Multimodal integration has matured significantly, with models now capable of processing and generating across multiple data types simultaneously, opening new possibilities for creative applications, scientific research, and human-AI interaction.
- The continued evolution of evaluation frameworks, benchmarks, and safety mechanisms underscores the industry's recognition that responsible AI development requires sophisticated means of measuring and ensuring system performance, reliability, and alignment with human values.

- As we move forward from March 2025, the AI landscape appears poised for continued rapid innovation, with models becoming increasingly integrated into critical infrastructure, creative processes, scientific discovery, and everyday human activities—raising both exciting possibilities and important questions about governance, access, and societal impact.