









## NEWMIND AI JOURNAL WEEKLY CHRONICLES




7.4.2025 - 14.4.2025




- This week saw the launch of major new models, including Deep Cogito's hybrid Cogito v1 and Google's Gemini 2.5 Flash, designed for high efficiency.
- Google unveiled its Ironwood TPU, capable of delivering a remarkable 42.5 exaflops of computing power.
- LLM techniques are evolving, with a growing emphasis on smaller, specialized models and the Model Context Protocol gaining widespread use.
- Research continues to push boundaries in reasoning, evaluation methods, and multimodal systems.
- AI is being deployed in new domains such as nuclear power plants, legal drug discovery, and autonomous robotics.
- Applications of AI are growing across healthcare, manufacturing, and creative industries.
- Governments worldwide are accelerating the development of AI regulatory frameworks to promote responsible and ethical AI use.
- The Chronicle covers key developments across six core categories: Models, AI Chips, LLM Techniques & Metrics, AI Use Cases, AI Policies, Regulations & Strategies, and AI Events & People.
- Our goal is to provide a comprehensive snapshot of this fast-moving and transformative moment in artificial intelligence.




 Models					
#	Highlights	Summary	Author	Source	Date
1.1	<b>Cogito v1 Preview Introducing IDA as a path to general superintelligence</b>	Deep Cogito introduces Cogito v1, a new family of language models built to blend fast, general-purpose AI with advanced reasoning. These hybrid models shift between standard and reasoning modes, tackling both simple tasks efficiently and complex ones thoughtfully. With sizes from 3B to 70B parameters—and a 670B model on the way—Cogito v1 shows top-tier performance across reasoning and general benchmarks. Developed from LLaMA and Qwen foundations and refined through novel training methods, these models are now available via Fireworks AI and Together AI, opening up powerful capabilities for real-world AI applications.	By Deep Cogito Team		April 8, 2025



Models					
#	Highlights	Summary	Author	Source	Date
1.2	<b>DDT: Decoupled Diffusion Transformer</b>	Diffusion transformers achieve high-quality generation but suffer from slow training and conflicting optimization between semantic encoding and detail decoding. To address this, we propose Decoupled Diffusion Transformer (DDT), featuring a dedicated condition encoder for extracting semantic content and a velocity decoder for refining high-frequency details. This decoupling resolves the optimization conflict and improves training efficiency. DDT-XL/2 sets new state-of-the-art FID scores on ImageNet (1.31 at 256x256, 1.28 at 512x512) with nearly 4x faster training. Additionally, our architecture accelerates inference by reusing self-conditions across denoising steps, guided by a dynamic programming strategy for optimal performance.	By Shuai Wang, Zhi Tian, Weilin Huang, Limin Wang		April 8, 2025
1.3	<b>Amazon Unveils Nova and Sonic, Foundation Models for Multimodal and Speech AI</b>	Amazon has introduced <b>Nova</b> and <b>Sonic</b> , two new <b>foundation models</b> designed for <b>multimodal reasoning</b> and <b>speech processing</b> , respectively. <b>Nova</b> powers Alexa's next-gen capabilities with advanced text, image, and video understanding, while <b>Sonic</b> is a state-of-the-art speech model trained on 100,000+ hours of data for real-time transcription, understanding, and generation. These models support applications in customer service, accessibility, and smart home devices. Integrated into Alexa and AWS services, Nova and Sonic highlight Amazon's push toward <b>domain-specialized AI</b> , advancing interaction quality and enterprise use of generative AI.	By Amazon		April 8, 2025
1.4	<b>Google Unveils Gemini 2.5 Flash for Efficient AI Applications</b>	Google has introduced Gemini 2.5 Flash, a streamlined AI model designed for high-volume, cost-sensitive tasks. Set to launch on Vertex AI, it offers developers adjustable processing times to balance speed, accuracy, and cost. As a "reasoning" model, it takes slightly longer to respond, enhancing reliability through self-verification. Ideal for applications like customer	By Google Team		April 8, 2025




 Models					
#	Highlights	Summary	Author	Source	Date
		service and document parsing, Gemini 2.5 Flash emphasizes low latency and reduced costs. Additionally, Google plans to deploy this model on-premises via Google Distributed Cloud, collaborating with Nvidia to support clients with strict data governance needs.			
1.5	<b>ServiceNow Releases Apriel-5B-Instruct: A Lightweight, Open-Source Instruction-Tuned LLM</b>	ServiceNow has unveiled Apriel-5B-Instruct, a 4.8B parameter open-source language model designed for instruction following, reasoning, and safe dialogue. Built atop Apriel-5B-Base, it underwent continual pretraining (CPT), supervised fine-tuning (SFT), and post-training alignment using DPO and RLVR. The model merges domain-specific variants (e.g., code, math, instruction) into a general-purpose system. Evaluated via lm-eval-harness and evalchemy, Apriel-5B-Instruct demonstrates competitive performance across tasks like GSM8k and TruthfulQA, rivaling larger models such as LLaMA-3.1-8B-Instruct. It is optimized for efficiency and safety, and is available under the MIT license for research and enterprise use.	By ServiceNow		April 12, 2025
1.6	<b>Seaweed-7B: Cost-Effective Training of Video Generation Foundation Model</b>	This report introduces Seaweed-7B, a 7-billion-parameter video generation model trained from scratch using 665,000 H100 GPU hours. Despite limited computational resources, Seaweed-7B delivers performance on par with or better than much larger models. The paper emphasizes critical design decisions that optimize model performance in constrained environments. Seaweed-7B demonstrates strong generalization and can be efficiently adapted to various downstream tasks, including video generation, through lightweight fine-tuning or continued training. This work shows that with	By Team Seaweed		April 11, 2025



 Models					
#	Highlights	Summary	Author	Source	Date
		thoughtful design, mid-sized models can achieve competitive results, offering a cost-effective path for advancing video foundation models.			
1.7	<b>MINEWORLD: A REAL-TIME AND OPEN-SOURCE INTERACTIVE WORLD MODEL ON MINECRAFT</b>	MineWorld, a real-time interactive world model built on Minecraft. Using a visual-action autoregressive Transformer, MineWorld takes paired game scenes and actions, tokenizes them, and learns to predict future frames through next-token prediction. A novel parallel decoding method allows the model to generate 4–7 frames per second, enabling real-time gameplay interaction. The model excels at both visual fidelity and accurate action-following. Evaluated with new metrics tailored for world modeling, MineWorld significantly outperforms existing open-source diffusion-based models. Both the model and code are publicly available, supporting further research in embodied AI.	By Microsoft Research		April 11, 2025
1.8	<b>DeepCoder: A Fully Open-Source 14B Coder at O3-mini Level</b>	Together AI and Agentica released DeepCoder-14B-Preview, an open-source 14B parameter code reasoning model finetuned via reinforcement learning (RL) from DeepSeek-R1-Distill-Qwen-14B. It achieves 60.6% Pass@1 on LiveCodeBench, rivaling OpenAI’s o3-mini, and also scores 73.8% on AIME 2024 math tasks. Trained on 24,000 verifiable coding problems using 32 H100 GPUs over 2.5 weeks, its dataset was curated from TACO Verified, SYNTHETIC-1, and LiveCodeBench with rigorous filtering. The team introduced “verl-pipe” for a 2.5× RL training speedup. All models, datasets, training code, and optimizations have been open-sourced to foster community research and transparency.	By TogetherAI, Agentica		April 8, 2025




AI Chips					
#	Highlights	Summary	Author	Source	Date
2.1	<b>Broadcom Unveils \$10 Billion Share Buyback Amid AI Chip Surge</b>	Broadcom has announced a \$10 billion share buyback program running through the end of 2025, reflecting strong confidence in its semiconductor and infrastructure software businesses. CEO Hock Tan highlighted Broadcom's strategic positioning in the AI chip market, where demand from cloud providers seeking alternatives to Nvidia continues to rise. The announcement boosted shares nearly 3% in extended trading. Broadcom, a key Apple supplier, recently forecast robust Q2 revenue and hinted at acquiring new customers, underscoring its competitive edge in custom AI chip production.	By Reuters		April 8, 2025
2.2	<b>IBM Research Powers z17 with AI-Centric Telum II and Spyre Chips</b>	IBM Research played a pivotal role in developing the Telum II processor and Spyre Accelerator at the heart of the new IBM z17 mainframe, engineered to meet tomorrow's AI demands. Spyre, a 32-core AI chip available as a PCIe card, supports over 250 enterprise AI use cases like fraud detection and generative AI. Built through hardware-software co-design, it delivers over 3x efficiency per watt compared to GPUs. Spyre leverages low-precision computing, optimizing for on-premise inference at speed and scale while supporting IBM's broader AI hardware roadmap.	By IBM Research		April 8, 2025
2.3	<b>Ironwood: The first Google TPU for the age of inference</b>	Google has introduced Ironwood, its seventh-generation TPU built to power generative AI inference. It scales to 9,216 liquid-cooled chips, connected via advanced Inter-Chip Interconnect (ICI) networking, delivering 42.5 exaflops—over 24 times the compute of El Capitan, the top supercomputer. Each chip reaches 4,614 teraflops peak performance. Designed for large language models and complex reasoning, Ironwood includes improved SparseCore, greater high-bandwidth memory, and faster networking. Integrated with Google's Pathways stack, it enables efficient distributed computing. Ironwood marks a major leap in AI infrastructure, offering	By Google		April 9, 2025



AI Chips					
#	Highlights	Summary	Author	Source	Date
		developers immense power and scalability to meet today's toughest AI challenges.			
2.4	<b>TSMC Q1 Revenue Surges Past Forecasts on Strong AI Chip Demand</b>	<b>Taiwan Semiconductor Manufacturing Co (TSMC)</b> reported first-quarter 2025 revenue of <b>NT\$592.64 billion (\$18.4 billion)</b> , beating market expectations due to surging global demand for <b>AI-related semiconductors</b> . Revenue jumped <b>16.5% year-over-year</b> , driven by strong orders from key clients like <b>Apple and Nvidia</b> , who rely on TSMC's cutting-edge nodes for AI and data center chips. The performance signals continued momentum in AI infrastructure spending, despite broader economic uncertainties. TSMC's results reinforce its position as a leading global chipmaker and a key enabler of AI model deployment and innovation.	By <a href="#">Ben Blanchard</a> and <a href="#">Wen-Yee Lee</a>		April 10, 2025
2.5	<b>Google Unveils AI Hypercomputer Upgrades with Custom Chips and Liquid Cooling</b>	At Cloud Next 2025, Google introduced upgrades to its <b>AI Hypercomputer architecture</b> , combining <b>custom-built TPUs</b> , GPUs, and <b>liquid-cooled data center designs</b> optimized for large-scale AI workloads. New innovations include integration of <b>Nvidia's Blackwell GPU</b> , next-gen <b>TPU v5p</b> , and high-performance fabric interconnects, enabling faster training and inference. These updates support models with <b>trillions of parameters</b> and reduce energy usage per computation. Google also highlighted orchestration via <b>Vertex AI</b> and better utilization across compute clusters. The hypercomputer reflects Google's commitment to pushing the frontier in AI model scaling and infrastructure efficiency.	By Google		April 9, 2025
2.6	<b>Sarvam AI Unveils Tool to Run LLMs</b>	Indian startup <b>Sarvam AI</b> has launched a breakthrough tool that enables large language models (LLMs) to run efficiently on <b>non-GPU hardware</b> ,	By Bloomberg News		April 10, 2025




AI Chips					
#	Highlights	Summary	Author	Source	Date
	<b>Without Expensive GPUs</b>	significantly lowering the cost of AI deployment. By optimizing models to operate on standard CPUs or less powerful chips, the solution reduces dependency on costly GPUs, addressing accessibility and affordability for enterprises and governments. The innovation aligns with India's push for <b>self-reliant AI infrastructure</b> and could democratize AI adoption in resource-constrained environments. Sarvam AI plans to open-source the tool to foster broader community development.			
2.7	<b>AMD Launches 5th Gen EPYC Processors to Power Next-Gen AI and Cloud Workloads</b>	AMD has announced the 5th Gen EPYC processors, designed to deliver top-tier performance and efficiency for AI, cloud, and enterprise workloads. Built on the "Zen 4" architecture, the new chips feature up to 128 cores and industry-leading memory bandwidth, positioning them for data-intensive applications like large model training and real-time analytics. AMD claims significant improvements in performance-per-watt and total cost of ownership. Major cloud providers, including Microsoft and Oracle, plan to deploy the new chips, reinforcing AMD's role in supporting the expanding global AI infrastructure.	By AMD Newsroom		April 9, 2025
2.8	<b>Chhattisgarh Lays Foundation for State's First Semiconductor Manufacturing Unit</b>	Chhattisgarh Chief Minister Vishnu Deo Sai has laid the foundation stone for the state's first semiconductor manufacturing facility, marking a significant step in India's push for chip self-reliance. The unit, part of a ₹1,500 crore investment, will be set up in the Electronics Manufacturing Cluster in Nava Raipur and aims to support domestic AI and electronics industries. The project is expected to generate 1,500 jobs and reduce dependency on imported chips. This move aligns with India's broader semiconductor mission to boost local manufacturing for AI-driven innovation.	By Business Standard		April 13, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.1	<b>Meta Denies Claims of Inflated Llama 4 Benchmark Scores</b>	Meta's VP of Generative AI, Ahmad Al-Dahle, denied allegations that the company trained Llama 4 models (Maverick and Scout) on benchmark test sets to boost evaluation results. Rumors emerged on social media suggesting Meta used an unreleased version of Maverick on the LM Arena leaderboard to inflate scores. Al-Dahle called the claims "simply not true" and attributed mixed model performance to implementation differences across platforms. Meta released the models immediately after readiness, and continues refining public deployments. The controversy highlights growing scrutiny around benchmark transparency in LLM evaluation.	By Kyle Wiggers		April 7, 2025
3.2	<b>Gaussian Mixture Flow Matching Models</b>	GMFlow, a novel Gaussian mixture flow matching model designed to overcome limitations of diffusion and flow matching models in few-step sampling. Unlike previous models that predict a single Gaussian mean, GMFlow predicts dynamic Gaussian mixture parameters, capturing a multi-modal flow velocity distribution. This approach improves accuracy by minimizing KL divergence loss. GMFlow also introduces GM-SDE/ODE solvers for precise sampling, leveraging denoising distributions and velocity fields. Additionally, a new probabilistic guidance scheme mitigates over-saturation issues with classifier-free guidance (CFG). Extensive experiments show GMFlow outperforms baselines, achieving a Precision of 0.942 with just 6 steps on ImageNet 256x256.	By Hansheng Chen et al.		April 7, 2025
3.3	<b>Google Launches Deep Research for Gemini 2.5 Pro Experimental Users</b>	Google has launched a new feature called Deep Research for Gemini 2.5 Pro Experimental users, enabling more advanced, document-level information gathering from across the web. Integrated into the Gemini Advanced experience, Deep Research helps users explore topics in greater depth by surfacing relevant sources, quotes, and multi-page summaries. The feature is designed for students, analysts, and professionals seeking	By Google Keyword Blog		April 8, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		nuanced responses, building on Gemini's improved reasoning and synthesis capabilities. Google positions it as part of a broader push toward more intelligent, task-specific LLM applications.			
3.4	<b>Generative Evaluation of Complex Reasoning in Large Language Models</b>	Generative Evaluation of Complex Reasoning in Large Language Models introduces KUMO, a novel evaluation framework specifically designed to assess complex reasoning capabilities of large language models (LLMs). Unlike traditional benchmarks, KUMO dynamically generates challenging, multi-step reasoning tasks with adjustable complexity and partial observability. This method tests the models' true reasoning abilities by requiring iterative and adaptive problem-solving approaches. The authors demonstrate that existing evaluation methods inadequately measure advanced reasoning, whereas KUMO effectively differentiates between genuine reasoning performance and superficial problem-solving strategies, offering a more accurate, flexible, and insightful metric for evaluating sophisticated cognitive skills in modern LLMs.	By Haowei Lin et al.		April 3, 2025
3.5	<b>VAPO: Reliable and Efficient Reinforcement Learning for Long Reasoning Tasks</b>	This paper introduces VAPO (Value-based Augmented Proximal Policy Optimization), a reinforcement learning framework tailored for advanced reasoning tasks in large language models. Built on the Qwen 32B model, VAPO achieves state-of-the-art performance on the AIME 2024 dataset with a score of 60.4—outperforming DeepSeek-R1-Zero-Qwen-32B and DAPO by over 10 points. VAPO addresses key RL challenges: value model bias, sequence length heterogeneity, and reward sparsity. It reaches peak performance in just 5,000 steps with no crashes, showcasing stability, efficiency, and scalability for long chain-of-thought (long-CoT) reasoning tasks.	By ByteDance Seed		April 8, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.6	<b>Missing Premise exacerbates Overthinking: Are Reasoning Models losing Critical Thinking Skill?</b>	The paper investigates how large language models (LLMs) handle questions with missing premises—ill-defined prompts lacking necessary context. It reveals that such questions trigger "overthinking," where models generate unnecessarily complex or incorrect reasoning instead of recognizing the ambiguity. To analyze this, the authors create a benchmark called MAP (Missing-premise Analytical Probes) and test various LLMs across domains. Results show that even top-performing models fail to detect missing premises, often hallucinating justifications. The study emphasizes the need for better evaluation of critical thinking in LLMs and introduces techniques to measure and mitigate overthinking in reasoning tasks.	By Chenrui Fan, Ming Li, Lichao Sun, Tianyi Zhou		April 9, 2025
3.7	<b>OLMOTRACE: Tracing Language Model Outputs Back to Trillions of Training Tokens</b>	OLMoTrace, a system that enables tracing the outputs of large language models (LLMs) back to specific segments within trillions of training tokens. By analyzing token-level influence during training, OLMoTrace uncovers which data most significantly shaped a given output. This technique enhances interpretability, accountability, and debugging in LLMs. Using the open OLMo model suite, the authors demonstrate how influential training documents can be identified for various prompts. OLMoTrace offers a powerful tool for researchers to understand model behavior, investigate memorization, and improve transparency in LLM development and deployment at scale.	By Jiacheng Liu et al.		April 9, 2025
3.8	<b>Google Introduces A2A Protocol for Interoperable AI Agents Across Platforms</b>	Google has unveiled <b>A2A (Agents-to-Agents)</b> , an open protocol designed to enable seamless <b>interoperability between AI agents</b> across different ecosystems, devices, and platforms. Inspired by internet standards like HTTP and SMTP, A2A defines how agents discover, communicate, and collaborate while maintaining security and permission controls. The	By Google Cloud		April 9, 2025



✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		protocol allows for agent-to-agent task delegation and is built to support a decentralized, multi-agent future. Google invites developers and researchers to help shape the standard, promoting an open, composable ecosystem where <b>LLMs, tools, and services</b> can interact across boundaries.			
3.9	<b>A Sober Look at Progress in Language Model Reasoning: Pitfalls and Paths to Reproducibility</b>	Reasoning is a key challenge for language models, but current evaluation methods often lack rigor, transparency, and consistency. This study reveals that popular math reasoning benchmarks are highly sensitive to factors like decoding settings, prompt formatting, random seeds, and hardware/software differences. Reported performance gains frequently rely on unclear or unreported variables. To address this, the authors propose a standardized framework with best practices for reproducible evaluation. Reassessing recent methods, they find reinforcement learning (RL) offers limited improvements and risks overfitting, especially on small benchmarks. In contrast, supervised finetuning (SFT) demonstrates more robust and consistent generalization performance.	By Andreas Hochlehnert et al.		April 9, 2025
3.10	<b>FantasyTalking: Realistic Talking Portrait Generation via Coherent Motion Synthesis</b>	FantasyTalking, a novel framework for generating realistic talking portrait videos with coherent motion. Unlike prior methods that often produce jittery or unnatural facial movements, FantasyTalking introduces a two-stage motion synthesis pipeline: a diffusion-based motion generator for initial dynamics and a refinement module for enhancing realism and synchronization with speech. The model effectively decouples appearance and motion, ensuring identity preservation across various expressions and head poses. Extensive experiments demonstrate its superiority in visual quality, lip-sync accuracy, and temporal stability. FantasyTalking sets a new	By Mengchao Wang et al.		April 7, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		standard for high-fidelity talking-head generation in virtual avatars and media applications.			
3.11	<b>Google Launches Agent Development Kit for Building Multi-Agent AI Systems</b>	Google has introduced the <b>Agent Development Kit (ADK)</b> , an open-source framework designed to simplify the creation of <b>multi-agent AI applications</b> . Built with Google's <b>A2A (Agents-to-Agents)</b> protocol, ADK provides tools for defining agent roles, communication, and task coordination. Developers can build collaborative agents that interact through structured APIs and shared memory. ADK supports integration with <b>Gemini models</b> and encourages modular, interoperable AI systems. The toolkit aims to accelerate research and development of complex multi-agent workflows, enabling more flexible, scalable, and intelligent agent ecosystems.	By Google Developers		April 10, 2025
3.12	<b>Towards Visual Text Grounding of Multimodal Large Language Model</b>	Despite advances in Multimodal Large Language Models (MLLMs), they still struggle with visual text grounding, particularly in text-rich document images like scanned forms and infographics. Current benchmarks mainly target natural images, leaving this gap unaddressed. To tackle it, the authors introduce TRIG—a new task and instruction dataset focused on Text-Rich Image Grounding in document-based QA. They create 800 annotated QA pairs and a 90k synthetic dataset via an OCR-LLM-human pipeline. Evaluations reveal MLLM weaknesses in grounding. Two TRIG methods—instruction tuning and embedding-based—significantly boost spatial reasoning and grounding when models are fine-tuned on the new dataset.	By Ming Li, Ruiyi Zhang, Jian Chen, Jiuxiang Gu, Yufan Zhou, Franck Dernoncourt, Wanrong Zhu, Tianyi Zhou, Tong Sun		April 7, 2025
3.13	<b>SoTA with Less: MCTS-Guided</b>	ThinkLite-VL, a vision-language model that achieves strong visual reasoning performance through data-efficient training. Using Monte Carlo	By Xiyao Wang, Zhengyuan		April 10, 2025

✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
	<b>Sample Selection for Data-Efficient Visual Reasoning Self-Improvement</b>	Tree Search (MCTS), the authors select informative training samples to fine-tune Qwen2.5-VL-7B-Instruct, avoiding the need for full dataset usage. This method allows the model to achieve state-of-the-art results on the MathVista benchmark and seven other visual reasoning datasets, using only 25% of training data. The approach not only boosts efficiency but also enhances generalization to unseen tasks. This work demonstrates how intelligent sample selection can drive performance gains in large multimodal models without extensive data consumption.	Yang, Chao Feng, Hongjin Lu, Linjie Li, Chung-Ching Lin, Kevin Lin, Furong Huang, Lijuan Wang		
3.14	<b>C3PO: Critical-Layer, Core-Expert, Collaborative Pathway Optimization for Test-Time Expert Re-Mixing</b>	Mixture-of-Experts (MoE) LLMs often suffer from suboptimal expert routing, leaving a 10–20% performance gap. To address this, the authors propose C3PO (Critical-Layer, Core-Expert, Collaborative Pathway Optimization), a test-time optimization method that re-weights expert mixing in key layers using surrogate objectives based on similar samples. These include mode-finding, kernel regression, and average loss among neighbors. C3PO improves accuracy by 7–15% on six benchmarks and outperforms in-context learning and prompt tuning. Notably, it enables smaller MoE LLMs (1–3B active parameters) to surpass larger ones (7–9B), boosting both accuracy and efficiency through smarter expert selection.	By Zhongyang Li, Ziyue Li, Tianyi Zhou		April 10, 2025
3.15	<b>VCR-Bench: A Comprehensive Evaluation Framework for Video Chain-of-Thought Reasoning</b>	VCR-Bench, a benchmark for evaluating Video Chain-of-Thought (CoT) reasoning in large vision-language models (LVLMs). It includes 859 videos and 1,034 manually annotated question-answer pairs with stepwise rationales tagged for perception or reasoning. VCR-Bench defines seven task dimensions and proposes a CoT score to assess reasoning quality. Experiments reveal major weaknesses in current LVLMs—most score below 40%, with perception tasks posing greater challenges than reasoning. The best model achieves only 62.8% CoT score. Results	By Yukun Qi, Yiming Zhao, Yu Zeng, Xikun Bao, Wenxuan Huang, Lin Chen, Zehui Chen, Jie Zhao,		April 10, 2025


✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
		confirm the benchmark's validity and highlight the need for improved temporal-spatial reasoning in video-based AI systems.	Zhongang Qi, Feng Zhao		
3.16	<b>KIMI-VL TECHNICAL REPORT</b>	The Kimi-VL Technical Report introduces Kimi-VL, a powerful large vision-language model (LVLM) designed to handle complex multimodal tasks. By integrating visual perception and language understanding, Kimi-VL achieves strong performance in image captioning, visual question answering, and document understanding. The model leverages a two-stage training approach combining large-scale pretraining and task-specific instruction tuning. Kimi-VL demonstrates competitive results across multiple benchmarks and excels in processing both natural and document images. This report details the model architecture, training pipeline, and evaluation results, positioning Kimi-VL as a strong general-purpose LVLM capable of tackling diverse real-world vision-language tasks efficiently.	By Kimi Team		April 10, 2025
3.17	<b>Sakana AI's AI Scientist-v2 Achieves Autonomous Scientific Discovery</b>	Sakana AI has introduced AI Scientist-v2, an advanced autonomous system capable of conducting the full scientific research process without human intervention. Enhancements include agentic tree search, Vision-Language Model (VLM) feedback, and parallel experimentation. Notably, AI Scientist-v2 successfully authored and submitted three papers to a peer-reviewed ICLR workshop, with one surpassing the average human acceptance threshold. This marks a significant milestone in AI-driven research, showcasing the system's ability to generate hypotheses, design and execute experiments, analyze data, and write scientific manuscripts autonomously. The project is open-sourced to encourage further development in autonomous scientific discovery.	By Sakana AI		April 8, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.18	<b>SQL-R1: Training Natural Language to SQL Reasoning Model By Reinforcement Learning</b>	SQL-R1: Training Natural Language to SQL Reasoning Model By Reinforcement Learning introduces SQL-R1, a novel model that translates natural language queries into SQL statements using reinforcement learning (RL) techniques. Traditional NL2SQL models often rely on supervised fine-tuning, which can struggle with complex queries involving multi-table joins and nested structures. SQL-R1 addresses these challenges by employing a specialized RL-based reward function tailored for NL2SQL tasks. The model demonstrates strong performance, achieving execution accuracies of 88.6% on the Spider benchmark and 66.6% on BIRD, utilizing only a 7B base model and a limited amount of synthetic training data.	By Peixian Ma, Xialie Zhuang, Chengjin Xu, Xuhui Jiang, Ran Chen, Jian Guo		April 11, 2025
3.19	<b>Study Reveals Early-Fusion Multimodal Models Offer Efficiency and Performance Advantages</b>	A recent study titled "Scaling Laws for Native Multimodal Models" investigates the architectural design of native multimodal models (NMMs)—those trained from the ground up on all modalities. The research, encompassing 457 trained models with varying architectures and training mixtures, finds that early-fusion models, which process modalities jointly from the outset, outperform late-fusion models, especially at lower parameter counts. Early-fusion architectures are more efficient to train and easier to deploy. The study also demonstrates that incorporating Mixture of Experts (MoEs) allows models to learn modality-specific weights, significantly enhancing performance.	By Apple Research		April 11, 2025
3.20	<b>ModernBERT or DeBERTaV3? Examining Architecture and Data Influence on Transformer</b>	Transformer-based language models—ModernBERT and DeBERTaV3—to evaluate the impact of architecture and pretraining data. Both models are trained on the same French dataset to isolate the effects of architectural design. Results show DeBERTaV3 outperforms ModernBERT in accuracy and sample efficiency, while ModernBERT is faster in training and inference. The study also finds that high-quality pretraining data	By Wissam Antoun, Benoît Sagot, Djamé Seddah		April 11, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
	<b>Encoder Models Performance</b>	accelerates convergence but does not drastically improve final performance. This work highlights the trade-offs between model design and training efficiency, offering insights for building better multilingual language models.			
3.21	<b>Computer Agent Arena</b>	The Computer Agent Arena, developed by XLANG Lab, is a benchmark platform for evaluating AI agents in complex, open-ended tasks. It assesses agents' capabilities in reasoning, tool usage, and multi-agent collaboration through diverse scenarios like simulated environments and human-AI interactions. Key metrics include task completion accuracy, efficiency, and adaptability to novel challenges. The platform aims to standardize AI agent evaluation, addressing limitations of current benchmarks by emphasizing real-world applicability and generalization. It supports both rule-based and learning-based agents, fostering advancements in autonomous AI systems. The Arena's modular design allows customization for specific research needs, promoting innovation in agent development.	By Bowen Wang, Xinyuan Wang et al.		April 7, 2025
3.22	<b>Concise Reasoning via Reinforcement Learning</b>	This paper challenges the belief that longer responses improve reasoning in LLMs, showing verbosity results from RL optimization dynamics, not necessity. A mathematical analysis reveals RL inflates response length during loss minimization. The authors propose a two-phase training framework: (1) RL on complex tasks to build reasoning, then (2) conciseness tuning using solvable problems. Experiments with 1.5B/7B models show 54% shorter responses without accuracy loss. They also uncover a natural conciseness-accuracy correlation. Training requires $\leq 8$ examples, enabling resource-efficient deployment and improved robustness at low temperatures—offering a blueprint for leaner, cost-effective reasoning systems.	By Wand AI		April 7, 2025




✦ LLM Techniques & Metrics					
#	Highlights	Summary	Author	Source	Date
3.23	<b>One-Minute Video Generation with Test-Time Training</b>	This paper introduces Test-Time Training (TTT) layers to generate one-minute videos with complex narratives using LLMs, addressing the inefficiency of self-attention in handling ~300k-token video contexts. TTT replaces RNN layers with neural hidden states updated via gradient descent at inference, improving context compression over Mamba. Integrated into CogVideo-X 5B, the TTT-MLP architecture yields a 34 Elo gain over Mamba 2 and Gated DeltaNet in human evaluations. Trained on a 7-hour Tom & Jerry dataset with 63-second multi-scene clips, it enables coherent, artifact-limited storytelling across scenes with no post-processing and reduces inference latency to 2.5x local attention vs. 11x full attention.	By NVIDIA, Stanford University, UCSD, UC Berkeley and UT Austin		April 7, 2025
3.24	<b>LightPROF: A Lightweight Reasoning Framework for Large Language Model on Knowledge Graph</b>	LightPROF is a lightweight framework that boosts large language models' (LLMs) reasoning abilities by efficiently integrating knowledge graphs (KGs). Unlike traditional methods that rely on textual prompts and miss KG structures, LightPROF follows a three-step approach: Retrieve, Embed, and Reason. It first extracts relevant subgraphs, then uses a Transformer-based Knowledge Adapter to encode both factual and structural KG data into the LLM's embedding space. Only the adapter is trained, enabling easy integration with open-source LLMs. Tests on KGQA benchmarks show improved performance, fewer tokens, and faster reasoning, making it effective for complex reasoning tasks.	By Tu Ao, Yanhua Yu, Yuling Wang, Yang Deng, Zirui Guo, Liang Pang, Pinghui Wang, Tat-Seng Chua, Xiao Zhang, Zhen Cai		April 4, 2025
3.25	<b>Rethinking Reflection in Pre-Training</b>	Rethinking Reflection in Pre-Training investigates the emergence of self-reflection capabilities in large language models (LLMs) during the pre-training phase, prior to any reinforcement learning. The authors introduce deliberate errors into chains-of-thought (CoT) prompts to assess whether models can recognize and correct these mistakes. Their experiments with the OLMo-2 model family, trained on up to 4.8 trillion tokens, reveal that	By Essential AI		April 5, 2025





## ✦ LLM Techniques & Metrics




#	Highlights	Summary	Author	Source	Date
		even partially trained models exhibit both situational and self-reflection abilities. Notably, simple prompts like "Wait," can trigger models to identify and amend reasoning errors. This study suggests that reflective reasoning is an inherent capability developing during pre-training.			
3.26	<b>BrowseComp: A Benchmark for Persistent Web Browsing Agents</b>	This paper presents BrowseComp, a benchmark of 1,266 human-curated questions requiring persistent, multi-step web navigation to find entangled, obscure information (e.g., niche events). Questions are designed to be unsolvable by models like GPT-4o/4.5 within 10 minutes and have short, verifiable answers (e.g., "Ireland v Romania") for easy evaluation. The benchmark tests persistence (searching hundreds of pages) and creative reasoning (adaptive search). Humans solved 29.2% of questions (86.4% accuracy when solved), while OpenAI's Deep Research agent reached 51.5% accuracy, scaling with compute. However, model overconfidence in wrong answers reveals calibration issues, highlighting the need for more robust autonomous web agents.	By OpenAI		April 10, 2025



✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
4.1	<b>DeepRoute.ai and Qualcomm Partner to Develop Advanced Driver Assistance Systems</b>	Chinese autonomous driving company DeepRoute.ai has announced a strategic partnership with Qualcomm to co-develop advanced driver assistance systems (ADAS) using Qualcomm's Snapdragon Ride platforms. The collaboration will support both lidar-based and vision-only configurations, aiming to deliver features such as urban autopilot, highway navigation, and automated parking. The joint effort is focused on lowering costs and accelerating adoption of high-performance autonomous driving technologies in mainstream vehicles. By combining Qualcomm's hardware with DeepRoute's software, the two aim to deliver scalable and affordable ADAS solutions for next-generation automotive platforms.	By Reuters		April 8, 2025
4.2	<b>Rescale Raises \$115M to Advance AI-Driven Engineering Simulations</b>	Rescale has secured \$115 million in funding from investors including Nvidia and Applied Materials. The company specializes in cloud software that enables engineers to run advanced simulations—such as airflow over race cars or semiconductor designs. With the new funding, Rescale plans to scale its use of AI models trained on simulation data, delivering results in seconds with over 98% accuracy. While not as precise as full simulations, this approach dramatically reduces design time, supporting faster innovation in industries like aerospace, automotive, and chip manufacturing.	By Reuters		April 7, 2025
4.3	<b>Krea Raises \$83M to Build Creative GenAI Platform, Reaches \$500M Valuation</b>	AI startup Krea, founded by two Spanish engineers who rejected royal postgraduate fellowships, has raised \$83 million to build a generative AI platform for visual creators. Now valued at \$500 million, Krea offers a “one-stop shop” that integrates multiple GenAI models through an intuitive UI for image and video generation, with future plans for audio. Used by creators at Pixar, LEGO, and Samsung, Krea simplifies prompt engineering and	By Ingrid Lunden		April 7, 2025




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		emphasizes creator control. The Series B was led by Bain Capital Ventures with participation from Andreessen Horowitz and Abstract Ventures.			
4.4	<b>Galaxy S25's New Trick? Real-Time AI Chats Using Your Camera</b>	Samsung has introduced a new visual AI feature in the Galaxy S25, enhancing the smartphone's camera capabilities. The update allows users to capture more detailed and vibrant photos with improved automatic adjustments, including better lighting and color optimization. The AI can now intelligently recognize scenes and subjects, adjusting settings in real-time for optimal shots. Additionally, the update improves the device's performance, offering smoother multitasking and enhanced battery life. With this new update, Samsung continues to push forward in integrating AI into its devices, providing users with an upgraded and smarter experience.	By Jerri Ledford		April 7, 2025
4.5	<b>Quantization Hurts Reasoning? An Empirical Study on Quantized Reasoning Models</b>	This study explores the impact of quantization on reasoning language models, which are known for their high inference costs due to extended chain-of-thought processes. We evaluate the open-source DeepSeek-R1-Distilled Qwen, LLaMA models (ranging from 1.5B to 70B parameters), and QwQ-32B using state-of-the-art quantization techniques on weights, KV cache, and activations. Our evaluation includes benchmarks like AIME, MATH-500, GPQA, and LiveCodeBench. The results show that lossless quantization is possible with W8A8 or W4A16, but lower bit-widths risk significant accuracy loss. Model size, origin, and task difficulty are key factors in performance.	By Ruikang Liu et al.		April 7, 2025
4.6	<b>French AI clusters target collaborative, ethical use cases</b>	France is investing €360 million to launch nine AI clusters that merge research, innovation, and education, aiming to foster collaborative and ethical AI development. One notable project, PostGenAI@Paris at Sorbonne University, focuses on creating AI that is open, ethical, and	By Martin Greenacre		April 8, 2025



 AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		transparent. University President Nathalie Drach-Temam highlighted the importance of combating misinformation by improving data quality and educating people on the strengths and limitations of AI technologies. These clusters aim to build trust in AI, support responsible innovation, and ensure AI systems contribute positively to society by addressing real-world challenges with transparency and ethical integrity at their core.			
4.7	<b>Google hopes its experimental AI model can unearth new security use cases</b>	<p>Google has introduced Sec Gemini V1, an experimental AI reasoning model designed to assist cybersecurity professionals by handling data analysis and foundational tasks in vulnerability research. The model integrates data from sources like Mandiant threat intelligence and the open-source vulnerabilities database, aiming to enhance the efficiency of security workflows. Initial access is limited to select organizations for non-commercial research purposes, with the goal of identifying practical applications and refining the model based on user feedback. This initiative reflects Google's commitment to leveraging AI to bolster cybersecurity operations.</p>	<p>By Derek B. Johnson</p>		<p>April 7, 2025</p>
4.8	<b>Powering Personalized Learning with AI: Cengage Student Assistant Expansion Delivers New GenAI Capabilities to 1M+ Students</b>	<p>Cengage is expanding its AI-powered Student Assistant to over one million students across 100+ products by fall 2025. Integrated into the MindTap platform, this tool offers personalized, just-in-time feedback to enhance student learning. Recent updates include broader integration throughout the learning experience, support for complex question formats, and improved contextual linking to relevant resources. Additionally, instructors now have access to insights into student performance and usage statistics, aiding in addressing learning challenges. This expansion underscores Cengage's commitment to leveraging AI to personalize education and improve academic outcomes.</p>	<p>By Cengage Group</p>		<p>April 8, 2025</p>




 AI Use Cases					
#	Highlights	Summary	Author	Source	Date
4.9	<b>Diablo Canyon's the First U.S. Nuclear Plant to Use AI</b>	Diablo Canyon, California's only remaining nuclear power plant, has become the first in the U.S. to implement on-site generative AI. PG&E partnered with startup Atomic Canyon to deploy Neutron Enterprise, an AI tool designed to streamline the management of extensive technical documentation. This system aims to reduce the time plant personnel spend searching through millions of pages of regulatory and operational documents, enhancing efficiency and compliance. While currently focused on document retrieval, this deployment may pave the way for broader AI applications in nuclear facility operations.	By Alex Schultz, CalMatters		April 9, 2025
4.10	<b>The Hottest Pre-IPO Stock? An AI Robotics Startup With Bold Claims, Little Revenue</b>	Figure AI, founded by Brett Adcock, is a robotics startup planning to deploy over 200,000 humanoid robots by 2029, forecasting \$9 billion in revenue despite earning nothing last year. Only a few robots are currently in testing with BMW on assembly tasks. Pre-IPO shares are in high demand, reportedly outpacing SpaceX and OpenAI in popularity. Backed by Microsoft, Nvidia, Jeff Bezos' firm, and formerly OpenAI, the company claims a major breakthrough after ending its partnership with OpenAI. Still, experts express concern over its lack of audited financials and limited real-world testing of its ambitious robotic systems.	By Emily Glazer, Berber Jin, By Alexander Saeedy		April 9, 2025
4.11	<b>Fujitsu and Headwaters trial on-device generative AI solution to streamline JAL cabin crew workflows</b>	Fujitsu and Headwaters collaborated with Japan Airlines (JAL) to enhance cabin crew report creation through a field trial conducted from January 27 to March 26, 2025. They utilized Microsoft's Phi-4, a small language model optimized for offline use, to develop a chat-based system on tablets, facilitating efficient report generation during and after flights. The trial demonstrated significant time savings and improved report quality. Fujitsu's Kozuchi AI service fine-tuned Phi-4 with JAL's past reports, while Headwaters developed the application and provided technical support. This	By Fujitsu Limited, Headwaters Co., Ltd.		April 10, 2025



✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		initiative aims to streamline cabin operations and enhance passenger service.			
4.12	<b>Samsung and Google Cloud Expand Partnership, Bring Gemini to Ballie, a Home AI Companion Robot by Samsung</b>	Samsung Electronics and Google Cloud have expanded their partnership to integrate Google's generative AI technology, Gemini, into Samsung's home AI companion robot, Ballie. Set to be available in the United States and Korea this summer, Ballie will utilize Gemini's multimodal capabilities alongside Samsung's proprietary language models to process various inputs, including audio, visual, and environmental data. This integration enables Ballie to engage in natural conversations, assist with home management tasks like adjusting lighting and setting reminders, and provide personalized advice on health and well-being. The collaboration builds upon the successful integration of Gemini into Samsung's Galaxy S24 smartphone series.	By Samsung		April 9, 2025
4.13	<b>Nuro Hits \$6 Billion Valuation in New Fundan Round for Autonomous Delivery</b>	Autonomous delivery startup <b>Nuro</b> has secured fresh funding, pushing its valuation to <b>\$6 billion</b> , underscoring growing investor belief in <b>AI-powered last-mile logistics</b> . Nuro specializes in small, driverless vehicles designed for neighborhood deliveries and has partnered with major retailers and logistics firms. The latest round includes both existing and new investors, although specific funding amounts were not disclosed. The company plans to use the capital to scale operations, enhance its vehicle platform, and accelerate deployment in U.S. urban markets. Nuro aims to lead in autonomous delivery innovation.	By Reuters		April 9, 2025
4.14	<b>Google Expands Generative Media Capabilities for</b>	Google Cloud has expanded its <b>generative media offerings</b> on <b>Vertex AI</b> , empowering enterprises to create high-quality images, videos, audio, and 3D content using advanced models like <b>Imagen 2</b> , <b>MusicLM</b> , and	By Google Cloud		April 9, 2025


✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
	<b>Enterprises on Vertex AI</b>	DeepMind's video generation model. Businesses can now build <b>multimodal applications</b> , apply brand-safe filters, and customize content through <b>fine-tuning and parameter-efficient adapters</b> . Enhanced <b>governance and digital watermarking</b> via SynthID are included to ensure responsible use. Aimed at marketing, media, and retail industries, the update enables scalable, AI-driven content generation tailored to brand guidelines and creative workflows.			
4.15	<b>Google Unveils Unified Security Platform with AI at Its Core</b>	Google Cloud has introduced a <b>unified security platform</b> powered by AI to address emerging cybersecurity threats and protect AI innovation across enterprise environments. The new solution integrates threat detection, data protection, and zero trust capabilities across Google's ecosystem, including Mandiant and VirusTotal. Key to the launch is the use of <b>AI-driven threat intelligence and automated response tools</b> that help reduce manual effort and accelerate remediation. This move reflects Google's strategic vision for <b>security-by-design</b> in AI development, enabling organizations to innovate safely while complying with evolving regulations.	By Google Cloud		April 9, 2025
4.16	<b>Google Cloud Expands Cloud WAN to Support Global AI Infrastructure Demands</b>	Google Cloud has upgraded its <b>Cloud WAN</b> to meet the global connectivity demands of AI-era workloads. Designed for enterprises scaling AI applications, the enhanced Cloud WAN now offers <b>multi-region interconnectivity</b> , improved <b>SLA-backed latency guarantees</b> , and <b>application-aware routing</b> . These features support high-throughput, low-latency traffic essential for <b>training and deploying AI models</b> across distributed infrastructures. The update also integrates with Vertex AI and Google's AI-optimized compute services, ensuring secure and scalable	By Google Cloud		April 9, 2025




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		global data movement. It reflects Google's broader effort to optimize networking for AI-powered enterprise environments.			
4.17	<b>MIT Explores How LLMs Could Revolutionize Drug and Materials Design</b>	MIT researchers are exploring how <b>large language models (LLMs)</b> can accelerate the discovery of <b>new medicines and materials</b> by encoding chemical knowledge similarly to how they process natural language. These models can generate novel molecular structures, predict properties, and suggest experimental steps—reducing trial-and-error in <b>drug development</b> and <b>materials science</b> . While LLMs still face limitations in accuracy and real-world integration, early results show promise in enhancing scientific workflows. The initiative reflects a growing trend of applying foundation models beyond language to <b>physical sciences and engineering domains</b> .	By MIT News		April 9, 2025
4.18	<b>SC Capital Eyes Global Switch Acquisition Amid AI-Driven Data Center Boom</b>	Singapore-based SC Capital Partners is in talks to acquire Global Switch, a major British data center operator, signaling rising investor interest in infrastructure powering the AI revolution. Global Switch, valued at around \$6.5 billion, operates facilities across Europe and Asia, serving hyperscalers and cloud providers. As demand for AI accelerates, so does the need for high-performance, scalable data centers. The potential acquisition reflects how real estate and finance sectors are pivoting toward AI infrastructure investments, viewing data centers as essential to future digital economies.	By Business Times		April 10, 2025
4.19	<b>Google deploys AI to speed up connections at</b>	Google, in partnership with PJM Interconnection, is using artificial intelligence to streamline the process of connecting new power sources—like wind and solar—to the largest electricity grid in the U.S. The AI tools, developed with Alphabet-backed Tapestry, create a digital model of the grid	By Laila Kearney		April 11, 2025





✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
	<b>PJM, largest US power grid</b>	to automate and speed up interconnection reviews that traditionally take years. This innovation comes as rising energy demands from data centers and AI workloads place greater strain on the power grid. Regulators like FERC are monitoring developments to ensure energy costs and grid reliability remain balanced.			
4.20	<b>Emory-Led Team Uses AI to Discover New Family of Superconductors</b>	A research team led by <b>Emory University</b> has leveraged artificial intelligence to discover a <b>new family of superconducting materials</b> , potentially revolutionizing power transmission, computing, and medical technologies. Using AI to screen over 30,000 compounds, the team identified <b>five promising superconductors</b> , significantly accelerating what is traditionally a trial-and-error process. This marks one of the first successful AI-driven breakthroughs in superconductor research. The study showcases how <b>machine learning</b> can enhance materials science, opening doors to faster, more efficient scientific discovery in critical industrial and technological fields.	By Carol Clark		April 10, 2025
4.21	<b>MOSAIC: Modeling Social AI for Content Dissemination and Regulation in Multi-Agent Simulations</b>	MOSAIC, an open-source simulation framework combining generative language agents with a social graph to model user behaviors like liking, sharing, and flagging content. By assigning diverse personas to agents, MOSAIC simulates large-scale social content dissemination and user engagement. It analyzes emergent deceptive behaviors and explores how users assess online content veracity. The authors evaluate three content moderation strategies in simulated misinformation scenarios, finding that these not only reduce false content spread but also boost engagement. Additionally, they study whether agents' reasoning aligns with engagement outcomes. The framework supports interdisciplinary research in AI and social dynamics.	By Genglin Liu, Salman Rahman, Elisa Kreiss, Marzyeh Ghassemi, Saadia Gabriel		April 10, 2025




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
4.22	<b>Thumbtack Leverages AI to Revolutionize Home Services</b>	Home services platform Thumbtack reported \$400 million in revenue for 2024, a 27% increase from the previous year, attributing this growth to strategic investments in artificial intelligence. The company has integrated AI tools that analyze user-uploaded photos of home issues and respond to plain-language questions, streamlining the contractor hiring process. This approach transforms the experience from traditional keyword searches to intuitive, conversational interactions. Thumbtack's AI-driven model, coupled with integrations into platforms like Nextdoor and Alexa, exemplifies how vertical platforms can challenge traditional search engines by offering seamless, ambient booking experiences across digital ecosystems.	By SHRM		April 10, 2025
4.23	<b>Shopify CEO Adopts 'AI-First' Hiring Policy, Redefining Workforce Strategy</b>	Shopify CEO Tobi Lütke has declared an "AI-first" hiring policy, stating that new roles must prove AI cannot do the job before being filled by a human. This bold move signals a paradigm shift in workforce planning, prioritizing automation and cost-efficiency. While Lütke frames it as future-forward, critics argue it risks displacing skilled workers and accelerating job insecurity. The policy reflects a growing trend among tech leaders to restructure operations around generative AI tools, raising pressing questions about labor policy, productivity, and ethical adoption of workplace automation.	By <a href="#">Roger Dooley</a>		April 8, 2025
4.24	<b>EMO-X: Efficient Multi-Person Pose and Shape Estimation in One-Stage</b>	EMO-X introduces a real-time, single-stage framework for multi-person 3D pose and shape estimation, removing the need for complex multi-step pipelines. By combining detection and regression in one network, it reduces computational load while preserving accuracy. Using dense feature	By Haohang Jian et al.		April 11, 2025

✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		<p>correlations and spatial attention, it handles occlusions and crowded scenes effectively. Tested on CMU Panoptic and 3DPW, EMO-X achieves state-of-the-art results with significantly faster inference. Its lightweight design suits edge devices, enabling AR/VR and robotics applications. Ablation studies explore architectural trade-offs, highlighting the balance between speed and precision in monocular 3D reconstruction.</p>			
4.25	<p><b>PixelFlow: Pixel-Space Generative Models with Flow</b></p>	<p>PixelFlow presents a flow-based generative model that operates directly in pixel space, avoiding latent-space bottlenecks. Using invertible transformations, it models complex image distributions with efficient, tractable likelihood estimation. Compared to diffusion models, it offers better scalability, training stability, and high-fidelity synthesis with fewer resources. Experiments on benchmark datasets show strong performance in image generation and editing. Unlike autoregressive or GAN-based methods, its deterministic sampling enables faster inference. PixelFlow also supports super-resolution and inpainting, demonstrating versatility. This work bridges flow models and pixel-level generation, offering an efficient, high-quality alternative for image synthesis.</p>	<p>By Shoufa Chen et al.</p>		<p>April 10, 2025</p>
4.26	<p><b>Hypergraph Vision Transformers: Images are More than Nodes, More than Edges</b></p>	<p>This paper rethinks Vision Transformers (ViTs) by modeling images as hypergraphs, where higher-order relationships (beyond pairwise edges) capture complex spatial-semantic structures. The proposed Hypergraph ViT (HVT) dynamically learns hyperedges to group pixels or patches into semantically meaningful clusters, improving feature aggregation. Experiments on ImageNet and COCO show consistent gains over standard ViTs, particularly in fine-grained recognition and occlusion handling. The</p>	<p>By Joshua Fixelle</p>		<p>April 11, 2025</p>




✦ AI Use Cases					
#	Highlights	Summary	Author	Source	Date
		framework is modular, compatible with existing self-attention mechanisms, and scales linearly with input size. Ablations validate the importance of hypergraph sparsity and adaptive edge formation. HVT opens new avenues for integrating geometric priors into transformer-based vision models.			
4.27	<b>Seaweed-7B: Cost-Effective Training of Video Generation Foundation Model</b>	<p>Seaweed-7B addresses the prohibitive costs of training large-scale video generation models by introducing data-efficient strategies and architectural optimizations. Through curriculum learning and selective frame sampling, it reduces redundant computations while preserving temporal coherence. The model employs a sparse attention mechanism to scale to long video sequences, achieving competitive results on benchmarks like Kinetics and UCF101 with 30% fewer training resources. Notably, Seaweed-7B demonstrates zero-shot generalization to unseen domains, suggesting robust latent representations. The paper provides a detailed cost-performance analysis, comparing it to diffusion and transformer-based alternatives. This work enables broader accessibility to high-quality video synthesis without requiring massive compute infrastructure.</p>	By ByteDance Seed		April 11, 2025




AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
5.1	<b>Taiwan Accuses China of Using Generative AI for Disinformation Campaigns</b>	Taiwan's National Security Bureau reports that China is employing generative AI to intensify disinformation efforts aimed at dividing Taiwanese society. Over half a million controversial messages—mostly on platforms like Facebook and TikTok—have been detected in 2025 alone. These campaigns, often timed around sensitive political events or corporate announcements (e.g., TSMC's U.S. investment), are part of broader "cognitive warfare." The report also notes increased "grey-zone" tactics such as airspace incursions and balloon deployments. Taiwan accuses China of leveraging AI tools to automate and escalate information warfare.	By Reuters		April 8, 2025
5.2	<b>Anthropic Expands in Europe with 100 New Roles and New EMEA Head</b>	Anthropic, the U.S.-based AI firm behind the Claude chatbot, is significantly expanding its presence in Europe by creating over 100 new positions across engineering, research, sales, and operations, primarily in London and Dublin. As part of its strategic push, Anthropic has appointed Guillaume Princen as Head of EMEA to oversee its regional operations. The company, backed by tech giants Amazon and Google, was recently valued at \$61.5 billion. With enterprise clients such as BMW, WPP, and Novo Nordisk already using Claude, Anthropic sees Europe as vital to its global growth and long-term business development.	By Reuters		April 8, 2025
5.3	<b>IBM Acquires Hakkoda to Enhance Data Expertise for AI Transformations</b>	IBM has acquired Hakkoda Inc., a global data and AI consultancy, to bolster IBM Consulting's data transformation services. Hakkoda specializes in modern data platforms and cloud-native solutions, particularly within the Snowflake ecosystem. This acquisition aims to accelerate clients' AI-driven transformations by enhancing data management and analytics capabilities. The integration of Hakkoda's expertise is expected to provide clients with advanced tools to harness	By IBM Newsroom		April 7 ,2025




 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		data effectively, facilitating more informed decision-making and innovation in AI applications.			
5.4	<b>Expanding AI use, White House orders agencies to develop strategies and name leaders</b>	The White House has directed all U.S. federal agencies to appoint Chief Artificial Intelligence Officers as part of a broader push to adopt AI responsibly. This move aligns with President Biden’s executive order aimed at ensuring AI is used safely and ethically across government operations. Agencies must also implement AI governance boards and report AI use cases that could affect public rights or safety. The Office of Management and Budget (OMB) emphasized transparency, requiring public disclosures of AI systems and regular assessments to manage risks. The directive supports both innovation and accountability in federal AI deployment.	By David Shepardson		April 8, 2025
5.5	<b>Dr. Oz Pushed for AI Health Care in First Medicare Agency Town Hall</b>	In his first all-staff meeting as CMS administrator, Dr. Mehmet Oz promoted the use of AI in healthcare, suggesting that AI avatars could help diagnose conditions like diabetes at a fraction of the cost of human doctors. He emphasized how AI could make the Medicare system more efficient and affordable. This marks a significant policy shift, as CMS is one of the largest health agencies in the U.S. Oz’s push signals a political strategy to integrate emerging technologies like AI into national healthcare, sparking both interest and concern from healthcare professionals and government staff.	By Leah Feiger, Steven Levy		April 8, 2025
5.6	<b>How Trump’s Tariffs Could Make AI Development More Expensive</b>	Former President Donald Trump has proposed imposing high tariffs on Chinese imports if re-elected, including a 60% tariff on all Chinese goods. Experts warn this could significantly increase the cost of artificial intelligence (AI) development in the U.S., as many AI components—	By Billy Perrigo		April 8, 2025


AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		especially chips and hardware—are produced or assembled in China. The added expenses could hinder AI innovation, slow down progress, and make it harder for startups to compete. While intended to reduce reliance on Chinese manufacturing, the tariffs could unintentionally disrupt the U.S.'s ability to stay competitive in the global AI race.			
5.7	<b>Energy Secretary Links AI Power Demand to Coal Support, Warns Iran on Sanctions</b>	U.S. Energy Secretary Chris Wright warned Iran could face tighter sanctions if it fails to reach a nuclear agreement with President Trump. Simultaneously, he defended an executive order aimed at <b>reviving the coal industry</b> , emphasizing its importance in meeting rising energy demands from <b>AI data centers</b> and industrial operations. Wright argued that coal is essential for reliable base-load power as AI infrastructure grows. He also urged Europe to rely more on <b>U.S. energy exports</b> , predicting the region won't return to Russian supply post-Ukraine war.	By Reuters		April 9, 2025
5.8	<b>EU Plans to Ease AI Act Compliance for Startups</b>	The European Union is preparing to <b>lighten compliance requirements</b> under the upcoming AI Act for <b>startups and small businesses</b> . Proposed measures include reduced fees, simplified documentation, and access to <b>regulatory sandboxes</b> that allow real-world testing of AI systems. The goal is to <b>support innovation</b> without compromising safety, especially in high-risk applications. EU Internal Market Commissioner <b>Thierry Breton</b> stated that these adjustments will help smaller firms compete while maintaining core regulatory protections. The AI Act is expected to fully come into force by <b>2026</b> , marking a pivotal shift in EU tech policy.	By Foo Yun Chee		April 8, 2025
5.9	<b>Andreessen Horowitz Seeks \$20B</b>	Venture capital giant Andreessen Horowitz (a16z) is aiming to raise \$20 billion, its largest fund ever, to invest in growth-stage U.S. AI startups,	By Krystal Hu, Anna Tong and Kenrick Cai		April 8, 2025


 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
	<b>Megafund to Back U.S. AI Companies</b>	capitalizing on rising global interest in American AI innovation. The fund targets international investors eager to bypass geopolitical restrictions and back firms like xAI, Databricks, and OpenAI. A16z's ties to the Trump administration reportedly attract LPs seeking favorable U.S. alignment. The fund would support both new and follow-on AI investments, underscoring escalating capital demands in AI model development and the strategic role of U.S. tech leadership.			
5.10	<b>U.S. Senators Press Google, Microsoft on AI Cloud Deals Over Competition Fears</b>	Democratic U.S. Senators Elizabeth Warren and Ron Wyden are questioning Google's partnership with Anthropic and Microsoft's ties to OpenAI, citing concerns about reduced competition and potential antitrust violations. The senators asked for details on whether the deals give cloud providers exclusive AI model rights, or limit startups' independence. A recent FTC report flagged similar risks, noting that one agreement may prevent an AI firm from launching new models without the cloud provider. The inquiries signal rising scrutiny of Big Tech's dominance in AI infrastructure and partnership structures.	By Jody Godoy		April 9, 2025
5.11	<b>Trump order looks to tap coal in quest to power data centers</b>	President Donald Trump signed executive orders to revitalize the coal industry, aiming to meet growing energy demands from artificial intelligence (AI) data centers. The orders designate coal as a critical mineral, lifting barriers to mining on federal lands and prioritizing coal leasing. They also mandate agencies to rescind policies transitioning away from coal and preserve threatened coal plants. The administration seeks to promote coal exports and technology development while easing environmental reviews for coal projects. Critics argue these measures may increase greenhouse gas emissions and health risks.	By Bloomberg		April 8, 2025

AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
5.12	<b>China's quantum computer pioneers AI task with enhanced efficiency</b>	Origin Quantum (OQ), a Chinese quantum computing firm, has efficiently fine-tuned a billion-parameter AI model using its Origin Wukong superconducting quantum computer, marking a significant global milestone. This pioneering effort demonstrates quantum computing's practical application for refining large AI models. Reducing model parameters by 76% notably boosted training performance by 8.4%. The 72-qubit Origin Wukong quantum system serves over 23 million global users, completing 350,000 quantum computing tasks and showcasing quantum technology's potential to support specialized AI applications, addressing future challenges in computing power limitations.	By Xinhua		April 9, 2025
5.13	<b>Alphabet reaffirms \$75 billion spending plan in 2025 despite tariff turmoil</b>	Alphabet CEO Sundar Pichai reaffirmed the company's plan to spend \$75 billion in capital expenditures in 2025, with a strong focus on AI and data center infrastructure. Speaking at the Economic Club of Washington, Pichai emphasized that AI is central to Alphabet's long-term growth, and the investment will boost computing power and innovation. He also noted the company's efficiency efforts, including job reductions. This announcement reflects Alphabet's commitment to staying competitive in the rapidly evolving AI landscape, where major tech companies are racing to expand capabilities amid rising demand for generative AI tools and services.	By Kenrick Cai		April 10, 2025
5.14	<b>Amazon CEO sets out AI investment mission in annual shareholder letter</b>	In his annual shareholder letter, Amazon CEO Andy Jassy emphasized the company's significant investments in artificial intelligence (AI) as essential for maintaining competitiveness and enhancing customer experiences. He highlighted the necessity of substantial capital allocation for AI chips and data centers to support this initiative. Amazon has invested approximately \$8 billion in AI startup Anthropic, integrating its	By Greg Bensinger, Deborah Mary Sophia		April 10, 2025




 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		<p>Claude software into the newly introduced Alexa+. This move aligns with industry trends, as other tech leaders, including Alphabet's CEO Sundar Pichai, have also justified large AI-related expenditures.</p>			
5.15	<p><b>Stanford's AI Index 2025 Highlights Global AI Trends and Challenges</b></p>	<p>Stanford University's AI Index 2025 report reveals a rapidly evolving global AI landscape. The U.S. maintains leadership with 40 notable AI models in 2024, but China is closing the gap, producing 15 models and leading in AI publications and patents. The report notes a surge in open-weight models, with Meta's Llama and China's DeepSeek-R1 rivaling top U.S. models. AI hardware efficiency improved by 40%, reducing costs and enabling advanced models on personal devices. However, incidents of AI misuse have increased, prompting more safety research. The report underscores the need for responsible AI development amid rapid advancements.</p>	By Standford HAI		April 7, 2025
5.16	<p><b>National Academy of Medicine Calls for Collaborative Oversight of Generative AI in Healthcare</b></p>	<p>The National Academy of Medicine emphasizes that harnessing generative AI's potential in health and medicine necessitates robust collaboration and oversight. Key concerns include data privacy, security, and algorithmic bias. To address these, the Academy advocates for interdisciplinary cooperation among clinicians, technologists, policymakers, and patients. They recommend implementing governance frameworks, ethical guidelines, and continuous monitoring to ensure responsible AI integration. The report underscores that without coordinated efforts, the risks of AI misuse could overshadow its benefits in healthcare.</p>	By National Academies of Sciences, Engineering, and Medicine		April 10, 2025
5.17	<p><b>Xi Jinping's Southeast Asia Tour Aims to Strengthen</b></p>	<p>On April 14, 2025, Chinese President Xi Jinping commenced a diplomatic tour of Southeast Asia, starting with Vietnam, to reinforce China's commitment to global trade amid escalating tensions with the United States.</p>	By Phuong Nguyen and Khanh Vu		April 14, 2025

 AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
	<b>Regional Ties Amid U.S. Tariff Tensions</b>	States. Facing 145% U.S. tariffs, China seeks to deepen economic ties with neighboring countries. In Vietnam, Xi emphasized mutual benefits through collaboration in production, artificial intelligence, and green technologies, while warning against trade protectionism. Approximately 40 bilateral agreements are anticipated, covering areas such as defense, security, and infrastructure, including potential Chinese-funded railway projects. This tour underscores China's strategic positioning during a period of heightened global trade uncertainty.			
5.18	<b>Safe Superintelligence Inc. Achieves \$32B Valuation with Strategic Investments</b>	Safe Superintelligence Inc. (SSI), the AI startup founded by OpenAI co-founder Ilya Sutskever, has reached a valuation of \$32 billion following a \$2 billion funding round led by Greenoaks. Notably, Alphabet and Nvidia have invested in SSI, with Alphabet's cloud division providing Tensor Processing Units (TPUs) to support SSI's research efforts. SSI is dedicated exclusively to developing a safe superintelligence, deliberately avoiding interim products or commercialization. The company maintains a small team and operates with a singular focus on AI safety, distinguishing itself from other AI labs.	By Anthony Ha		April 12, 2025
5.19	<b>Germany to create 'super-high-tech ministry' for research, technology, and aerospace</b>	Germany is set to establish a new "super-high-tech ministry" dedicated to overseeing research, technology, and aerospace sectors. This initiative, outlined in the recent coalition agreement, aims to centralize and enhance the nation's innovation efforts. The ministry will coordinate scientific research, technological development, and aerospace initiatives, reflecting Germany's commitment to maintaining its competitive edge in these fields. By consolidating responsibilities previously spread across various departments, the government seeks to streamline decision-making processes and foster interdisciplinary collaboration. This strategic move	By Gretchen Vogel		April 11, 2025



🛡️ AI Policies Regulations & Strategies					
#	Highlights	Summary	Author	Source	Date
		underscores Germany's focus on advancing its position in global science and technology arenas.			
5.20	<b>NO FAKES Act re-introduced in US Senate</b>	Nearly a year after its initial introduction in July 2024, U.S. Senators Chris Coons, Marsha Blackburn, Amy Klobuchar, and Thom Tillis, along with Representatives María Elvira Salazar, Madeleine Dean, Nathaniel Moran, and Becca Balint, have reintroduced the NO FAKES Act (Nurture Originals, Foster Art, and Keep Entertainment Safe). The bill seeks to establish a federal intellectual property right over an individual's voice and likeness, protecting against unauthorized AI-generated use. Industry groups like SAG-AFTRA, RIAA, and MPA support the measure, along with tech leaders such as YouTube and OpenAI.	By The Senate Of The United States		April 11, 2025


★ AI Events & People					
#	Highlights	Summary	Author	Source	Date
6.1	<b>Welcome to Google Cloud Next '25</b>	At Google Cloud Next '25, Google introduced a wide range of AI-powered tools and infrastructure upgrades designed to boost enterprise adoption of generative AI. Key announcements include Gemini AI integrations into Google Workspace, upgraded agent-building tools in Vertex AI, and new custom hardware such as the TPU v5p and Axion CPUs. Google emphasized advancements in data privacy, cybersecurity, and responsible	By Google Cloud		April 9, 2025

☆ AI Events & People

#	Highlights	Summary	Author	Source	Date
		AI practices. The event also highlighted partnerships supporting open-source development and flexible deployment. These innovations aim to help businesses scale AI applications efficiently while maintaining governance and trust.			
6.2	<b>Optimizing Inference on Large Language Models With NVIDIA</b>	NVIDIA will host a free webinar titled “Optimizing Inference on Large Language Models” on April 17, 2025, at 2:00 PM IST. The session will guide developers on improving LLM performance using TensorRT-LLM and NVIDIA Triton. Topics include optimizing prompt processing, token generation, and real-world deployment strategies with a focus on latency, throughput, and cost. Attendees will see use cases like Tech Mahindra’s Hindi LLM. Participants also gain access to a \$90 course on LLM deployment. Ideal for AI engineers and students with basic LLM knowledge.	By NVIDIA		April 17, 2025
6.3	<b>Shaping the National Data Library: key considerations for the AI age</b>	The Open Data Institute (ODI) will host a key webinar, “Shaping the National Data Library: Key Considerations for the AI Age,” on April 10, 2025, from 11:00 to 12:00 BST. The event explores how a National Data Library (NDL) can support public interest and AI-driven innovation. Topics include data accessibility, centralized vs. federated systems, governance, technical design, ethical standards, and key datasets. Featuring experts like Professor Elena Simperl, the session invites input on how the UK can lead in building inclusive, effective data infrastructure. The event will be held on Zoom and recorded for later access.	By Open Data Institute		April 10, 2025
6.4	<b>Building Knowledge Graphs to Power Your AI Initiatives</b>	Ready to elevate your AI strategy? Join the Progress Semaphore webinar to explore how knowledge graphs can transform your approach to generative AI by enhancing accuracy, reducing hallucinations, and minimizing bias. This session covers the limitations of traditional AI, the	By Progress Semaphore		April 15, 2025

☆ AI Events & People

#	Highlights	Summary	Author	Source	Date
		benefits of knowledge graphs, and practical steps to build smarter, more reliable AI systems. You'll also see a live demo of knowledge graph construction and learn how metadata management and retrieval-augmented generation (RAG) enhance generative AI. If you're looking to overcome challenges in AI reliability, this webinar offers the insights and tools you need to advance.			
6.5	<b>Master Agentic AI with Managed MLflow</b>	Master Agentic AI with Managed MLflow to explore how Managed MLflow on Nebius AI Cloud can transform your LLM development process. This session will demonstrate how to enhance observability, benchmarking, and deployment of AI agents using tools like MLflow Tracing, Evaluation, Tracking, and Model Registry. Learn to analyze agent reasoning, run systematic experiments, and promote top-performing models to production. The webinar offers two sessions: 10 AM CEST for Europe and 10 AM PDT for the US. Ideal for developers seeking to build reliable, scalable AI applications.	By Nebius		April 17, 2025
6.6	<b>LLM Powered Browser Plugin</b>	LLM Powered Browser Plugin, to explore how large language models can enhance browser functionality. This session will delve into integrating LLMs into browser plugins, demonstrating how they can improve user experience, automate tasks, and provide intelligent assistance. Attendees will gain insights into the development process, best practices for implementation, and real-world applications of LLM-powered browser extensions. Whether you're a developer, product manager, or tech enthusiast, this webinar offers valuable knowledge on leveraging LLMs to create smarter browser tools. Don't miss this opportunity to learn from industry experts and advance your understanding of LLM integration.	By Intel Software		April 16, 2025

☆ AI Events & People					
#	Highlights	Summary	Author	Source	Date
6.7	<b>TechByte: 5 steps for laying the right data foundation for AI success</b>	Join Google Cloud's upcoming webinar, "5 Steps for Laying the Right Data Foundation for AI Success," to learn how to prepare your data infrastructure for effective AI implementation. This session will cover essential strategies for aligning data practices with AI goals, ensuring scalability, and maintaining data quality. Discover how to build a robust data foundation that supports AI-driven innovation and delivers reliable outcomes. Ideal for data professionals and business leaders aiming to harness AI's full potential. Register now to gain actionable insights and advance your organization's AI readiness.	By Google		April 15, 2025

## Conclusion

- The AI landscape is transitioning from experimental to operational, with infrastructure investments driving competitive advantage.
- Google's unified security platform, Amazon's custom chip investments, and growing industry-specific applications highlight this shift.
- The focus has moved from general-purpose AI to specialized, integration-ready solutions with measurable business impact.
- This evolution requires a dual investment strategy: building strong data foundations and developing domain-specific applications.
- Organizations that combine partnerships with established platforms and proprietary use case development will gain both speed and sustainable value.

- Success will favor those who embed AI into core operations, using their unique data, industry knowledge, and AI capabilities to create defensible advantages.