# Executive Summary

This week reveals critical gaps in AI agent capabilities, with frontier LLMs achieving less than 60% success on real-world multi-tool tasks. Notable advances include new GUI automation frameworks reaching 73.3% on AndroidWorld, single-pass 3D scene generation, and concerning evidence that "over-reasoning" impairs confidence calibration. DeepSeek-V3.1 shows significant benchmark improvements while environmental impact measurements suggest 33× energy reduction per prompt year-over-year.

# Breakthrough Papers

**1** **LiveMCP-101: Stress Testing and Diagnosing MCP-enabled Agents**

**Authors:** Ming Yin, Dinghan Shen, Silei Xu, et al.

Even frontier LLMs score less than 60% success on 101 real-world multi-tool tasks, exposing critical agent tool-use failure modes. The study reveals systematic issues in planning, execution, and error recovery across different agent architectures.

**Key Impact:** Challenges assumptions about current agent capabilities and provides diagnostic framework for identifying failure patterns in multi-tool scenarios.

[ar5iv] LiveMCP-101 Paper

**2** **Mobile-Agent-v3: Foundational Agents for GUI Automation**

**Authors:** Jiabo Ye, Xi Zhang, Haiyang Xu, et al.

New GUI-Owl framework sets open-source SOTA: AndroidWorld 73.3, OSWorld 37.7. GUI-Owl-7B achieves 66.4 and 29.4 respectively,

demonstrating significant progress in cross-platform GUI understanding and automation.

**Significance:** First open-source framework to achieve competitive performance with proprietary systems on mobile and desktop GUI automation tasks.

[ar5iv] Mobile-Agent-v3 Paper

---

**3** ## Don't Think Twice! Over-Reasoning Impairs Confidence Calibration

**Authors:** Romain Lacombe, Kerrie Wu, Eddie Dilworth

More "thinking" worsens calibration; search-augmented generation reaches 89.3% confidence-assessment accuracy vs 48.7% for pure reasoning. Challenges the assumption that more computation always improves performance.

**Implications:** Suggests optimal stopping points for reasoning and highlights the importance of external knowledge over extended internal computation.

[ar5iv] Over-Reasoning Paper

## Recent Model Releases

### DeepSeek-V3.1 (Aug 21)

- SWE-bench Verified: 66.0
- SWE-bench Multilingual: 54.5
- MMLU-Pro: 81.2 (vs 75.9 in V3)
- GPQA: 68.4 (vs 59.1 in V3)
- AIME: 59.4 (vs 39.6 in V3)
- LiveCodeBench: 49.2 (vs 39.2 in V3)

DeepSeek API Docs

### Other Notable Releases

**Genie 3 (DeepMind)**
World model for playable video at 720p, 24 fps, with minutes-long temporal consistency.

**Surya (IBM+NASA)**
Open-source heliophysics foundation model; 16% improvement in solar-flare classification.

**Gemini 2.5 Deep Think**
SOTA claims on LiveCodeBench V6; bronze-level 2025 IMO performance variant.

# Fresh Metrics You Can Cite

## Performance Benchmarks

- **Agent tool-use:** Frontier LLMs <60% success on LiveMCP-101
- **GUI automation:** AndroidWorld 73.3, OSWorld 37.7 (Mobile-Agent-v3)
- **Confidence calibration:** 89.3% with search vs 48.7% pure reasoning
- **DeepSeek improvements:** MMLU-Pro +5.3, GPQA +9.3, AIME +19.8

## Environmental Impact

- **Energy per prompt:** 0.24 Wh (Gemini Apps median)
- **Carbon footprint:** 0.03 g $CO_2$e per prompt
- **Water usage:** 0.26 mL per prompt
- **YoY improvement:** 33× energy, 44× carbon reduction

# What's Trending in Research

- Agentic AI and tool use at scale (MCP, planning, diagnostics)
- GUI automation agents across desktop/mobile
- World models for interactive simulation and video
- Energy, water, and carbon accounting for inference
- Calibration under test-time scaling
- Single-pass 3D scene generation from images

# Notable Industry Moves

- OpenAI will open its first India office in New Delhi this year.

- Sweden's Wallenberg groups with AstraZeneca, Ericsson, Saab, and SEB launch national AI company.

- Amazon signals "agents" as its near-term differentiator in consumer AI.

- IBM+NASA open-source "Surya" solar-weather model and dataset on Hugging Face.

---

- Sweden's Wallenberg groups with AstraZeneca, Ericsson, Saab, and SEB launch national AI company.

- Amazon signals "agents" as its near-term differentiator in consumer AI.

- IBM+NASA open-source "Surya" solar-weather model and dataset on Hugging Face.