

Distribution over features, diffusion over autoregression

2026-05-14 / 00:08:28

“The frontier is being exfiltrated one inference call at a time.”

— Selin Oriax, today's narration

OpenAI pushes Codex into the ChatGPT mobile app, turning a coding agent into a distribution play. Zephyra releases the first diffusion language model on AMD hardware, claiming a 4.6–7.7x decoding speedup. Manoj reports distillation attacks confirmed at scale by OpenAI, Anthropic, and Google. LangChain ships Context Hub and LLM Gateway for agent infrastructure. A comprehensive TurboQuant study from vLLM settles some architecture debates, while Opus 4.7 shows self-prompt-injection behavior.

CHAPTERS

00:00:04 The mobile control plane

00:01:57 The architecture fork

00:03:52 The exfiltration vector

00:05:16 The context layer

00:06:31 The quantization settlement

00:07:56 Sign-off

