

# **AI Engineer Summit Introduction and Welcome**

# Summary & Key Points

**Summary:** The hosts welcome attendees to the AI Engineer Summit, emphasizing the focus on builders and real-world applications of AI agents. They highlight the shift from theory to practice and the importance of addressing open questions in scaling, accuracy, and memory. The summit will feature use cases, deployment stories, and insights from various companies.

**Counter-Intuitive Points:** The emphasis on real-world applications and open questions, rather than just future possibilities, is a refreshing and practical approach.

**Topics to Explore:** AI Agents, Agent Engineering, AI Summit, Real-world AI, AI Challenges

# Q&A

Q: What is the main focus of the AI Engineer Summit?

A: The main focus is on AI agents, agent engineering, and the builders shaping the future of agents.

Q: What are some of the key challenges in AI agent development?

A: Key challenges include scaling, accuracy, memory, evaluation, reliability, cost management, and transitioning from theory to real-world practice.

Q: Which companies will be sharing their insights at the summit?

A: Experts from Google, OpenAI, Anthropic, Jane Street, BlackRock, Weights & Biases, Datadog, Morgan Stanley, Bloomberg, Brightwave, and Gaule will share insights.

# Swix's Context Setting: The Year of Agents

# Summary & Key Points

Summary: Swix discusses the evolution of AI engineering, the resistance from MLE and software engineering perspectives, and the pivot to agent engineering. He defines agents, explores why they are working now, and highlights key use cases, emphasizing the growth of ChatGPC and the importance of agency in AI products.

Counter-Intuitive Points: The idea that closing doors to certain AI areas (like RAG) can open up others in agent engineering is a surprising but insightful point.

Topics to Explore: AI Engineering, Agent Engineering, LLMs, ChatGPC, AI Trends

# Q&A

Q: What are the two main perspectives resisting the emergence of AI engineering as a distinct discipline?

A: MLE (Machine Learning Engineer) sees AI engineering as mostly MLE plus prompts, while software engineering sees it as mostly software engineering calling AI APIs.

Q: Why is 2025 considered the 'Year of Agents'?

A: There is a desire and prediction from figures like Satin Adela, Roman, Greg Brockman, and Sam Altman that 2025 will be the year when AI agents become prominent and widely adopted.

Q: What are the key factors driving the current success of AI agents?

A: Improved reasoning, better tool use, model diversity, the decreasing cost of intelligence, and the availability of fine-tuning options.

# **Building and Evaluating AI Agents Effectively**

# Summary & Key Points

Summary: Sayash Kapoor discusses the challenges in building effective AI agents, focusing on the difficulty of evaluation, the misleading nature of static benchmarks, and the confusion between capability and reliability. He emphasizes the need for rigorous evaluation, cost-controlled benchmarks, and a reliability-focused mindset in AI engineering.

Counter-Intuitive Points: The idea that static benchmarks can be misleading and that verifiers can be imperfect in practice challenges conventional evaluation methods.

Topics to Explore: AI Agents, AI Evaluation, Reliability Engineering, Static Benchmarks, AI Engineering Challenges



# Q&A

Q: What are the three main reasons why AI agents don't yet work effectively?

A: Evaluating agents is hard, static benchmarks are misleading, and there's confusion between capability and reliability.

Q: Why are static benchmarks misleading for evaluating AI agents?

A: Agents take actions in the real world, evaluations have no cost ceiling, and agents are often purpose-built, requiring meaningful multidimensional metrics.

Q: What is the key difference between capability and reliability in AI agent deployment?

A: Capability is what a model could do, while reliability is consistently getting the answers right each time, which is crucial for consequential decisions.

# **Gemini Deep Research: Product and Technical Challenges**

# Summary & Key Points

Summary: Arish and Mukhun from Google discuss the motivation, product challenges, and technical challenges in building Gemini Deep Research, a web research agent. They cover asynchronous experiences, user expectations, long outputs, failure robustness, iterative planning, noisy web environments, and context management.

Counter-Intuitive Points: The idea of showing users what Gemini is doing under the hood in real-time to build trust is a unique and effective approach.

Topics to Explore: AI Agents, Web Research, Gemini, Product Challenges, Technical Challenges

# Q&A

Q: What were the main product challenges in building Gemini Deep Research?

A: The product challenges included building asynchronous experiences in a synchronous product, setting user expectations for a specific use case, and making long outputs easy to engage with in a chat experience.

Q: How does Gemini Deep Research handle intermediate failures?

A: It uses a state management solution to recover from errors effectively, preventing the entire research task from failing due to one issue.

Q: What are some of the challenges in planning and context management for a web research agent?

A: Challenges include planning iteratively, spending compute effectively, interacting with a noisy web environment, and managing growing context size.

# **Anthropic's Approach: Building Effective AI Agents**

# Summary & Key Points

Summary: Barry Zhang from Anthropic shares insights on building effective AI agents, emphasizing the importance of not building agents for everything, keeping it simple, and thinking like your agents. He discusses the agent architecture, the role of tools, and the need for budget awareness and self-evolving tools.

Counter-Intuitive Points: The recommendation to keep agent architecture as simple as possible and to focus on the core components is a valuable reminder to avoid unnecessary complexity.

Topics to Explore: AI Agents, Agent Architecture, Tool Use, System Prompt, AI Engineering Best Practices

# Q&A

Q: What are the key criteria for determining when to build an AI agent?

A: Consider the complexity and value of the task, de-risk critical capabilities, and assess the cost of error and error discovery.

Q: What are the three basic components of an AI agent?

A: The environment, the set of tools, and the system prompt.

Q: Why is it important to 'think like your agent' during development?

A: To understand the agent's limited context and ensure it's sufficient and coherent for decision-making.

# **Sierra's Agent Development Lifecycle for Brands**



# Summary & Key Points

Summary: Zach Renaud Wadine from Sierra discusses how they build and improve AI agents for consumer brands, emphasizing the importance of treating each agent as a product. He introduces the Agent Development Lifecycle, which includes quality assurance, issue reporting, and continuous improvement, and highlights the benefits of building for voice.

Counter-Intuitive Points: The idea that large language models remind us of ourselves and allow us to be great designers by having empathy is a thought-provoking perspective.

Topics to Explore: AI Agents, Conversational AI, Customer Experience, Agent Development Lifecycle, Voice AI

# Q&A

Q: What is Sierra's approach to building and improving AI agents?

A: Sierra believes every agent is a product and uses the Agent Development Lifecycle to iteratively improve agents, involving a fully featured developer platform and customer experience operations platform.

Q: What are the key components of the Agent Development Lifecycle?

A: The Agent Development Lifecycle includes quality assurance, issue reporting, test creation, and new releases, with iterative refinement and customer feedback incorporated at each stage.

Q: How does Sierra think about building for voice?

A: Sierra builds AI agents that are responsive to whatever channel someone reaches out in, and whatever modality you're operating in, similar to responsive web design.

# Reinforcement Learning's Role in AI Agents

# Summary & Key Points

Summary: Will Brown from Morgan Stanley discusses the potential of reinforcement learning (RL) for enhancing AI agents. He explains how RL can enable agents to learn and improve their skills through interaction with an environment, and explores the challenges and opportunities in integrating RL into agent engineering.

Counter-Intuitive Points: The idea that long chains of thought can emerge as a byproduct of reinforcement learning, rather than being manually programmed, is a significant insight.

Topics to Explore: Reinforcement Learning, AI Agents, Model Training, GRPO Algorithm, AI Engineering Future

# Q&A

Q: What is the key idea behind reinforcement learning?

A: The key idea is to explore and exploit, trying new things and doing more of what works while doing less of what doesn't.

Q: How does the GRPO algorithm work?

A: For a given prompt, sample end completions, score them all, and tell the model to be more like the ones with higher scores.

Q: What are the challenges in extending reinforcement learning to more agentic systems?

A: The challenges involve figuring out how to give models skills and have them learn to get better, particularly in conjunction with environments, tools, and verification.

# **Windsorv: AI Agent-Powered Editor**

## **Principles**

# Summary & Key Points

Summary: Kevin Howe from Windsorv discusses the principles behind building their AI agent-powered editor, emphasizing trajectories, meta-learning, and scaling with intelligence. He explains how Windsorv uses these principles to reduce human input, predict user actions, and adapt to user preferences.

Counter-Intuitive Points: The decision to delete chat in favor of a more integrated agent experience is a bold move that challenges conventional UI design.

Topics to Explore: AI Agents, Code Editor, Trajectories, Meta-learning, Scale with Intelligence

# Q&A

Q: What are the three main principles behind Windsorv's AI agent-powered editor?

A: The three main principles are trajectories, meta-learning, and scaling with intelligence.

Q: How does Windsorv use trajectories to understand user actions?

A: Windsorv uses trajectories by building a unified timeline of user actions, including viewing files, navigating code, editing, searching, grepping, and making commits, allowing the agent to understand the user's workflow and predict their next steps.

Q: What is the concept of meta-learning in Windsorv?

A: Meta-learning in Windsorv refers to the agent's ability to adapt and remember user preferences and organizational guidelines over time, using auto-generated memories and custom MCP servers to personalize the coding experience.



# Scaling AI Agents: Method's Approach

# Summary & Key Points

Summary: Mustafa Ali from Method and Kyle Corbett from OpenPipe discuss how Method scaled to 500 million AI agents in production with just two engineers. They highlight the challenges of cost, quality, and latency, and explain how fine-tuning open-source models helped them overcome these challenges.

Counter-Intuitive Points: The idea that a smaller, fine-tuned model can outperform a larger, more general model in a specific use case is a key takeaway.

Topics to Explore: AI Agents, Scaling AI, Fine-tuning, Open-source Models, Cost Optimization

# Q&A

Q: What were the main challenges Method faced in scaling their AI agents?

A: The main challenges were high API costs, prompt engineering limitations, slow baseline latency, and AI errors (hallucinations).

Q: How did fine-tuning help Method overcome these challenges?

A: Fine-tuning allowed Method to use a smaller, cheaper model that met their accuracy requirements, significantly reduced latency, and lowered costs.

Q: What is the key takeaway for productionizing AI agents?

A: Productionizing AI agents requires openness and patience from the engineering and leadership teams, as it takes time to achieve production readiness.

# **Voice AI Agents: Reliability and Scalability**

# Summary & Key Points

Summary: Nick Carriotakis from SuperDial discusses how to make reliable and scalable voice AI agents. He emphasizes the importance of conversation design, choosing the right tools, and addressing last-mile problems like pronunciation and spelling. He also highlights the ethical considerations in building voice AI.

Counter-Intuitive Points: The emphasis on reliability over realism in voice AI challenges the common focus on creating human-like interactions.

Topics to Explore: Voice AI, AI Agents, Conversation Design, Speech-to-Text, Text-to-Speech

# Q&A

Q: What are the key challenges in building reliable and scalable voice AI agents?

A: Key challenges include audio hallucinations, pronunciation/spelling errors, maintaining low latency, and ensuring realistic yet reliable conversations.

Q: Why is conversation design important in voice AI?

A: Conversation design is important because it helps to create a more natural and engaging user experience, and it can also help to improve the accuracy and reliability of the agent.

Q: What are some of the ethical considerations in building voice AI?

A: Ethical considerations include potential biases against certain accents or dialects, the risk of creating spooky or misleadingly realistic interactions, and the need for accessibility and collaboration.

# **Scaffolding AI Agents Wisely for Scalability**

# Summary & Key Points

Summary: Rahul Sengotuvalu from Ramp discusses how to scaffold AI agents wisely for scalability, emphasizing the importance of building systems that scale with compute. He explores different architectures for AI agents and highlights the benefits of using fuzzy compute and allowing LLMs to decide when to break into classical compute.

Counter-Intuitive Points: The idea that the backend can be the LLM itself, rather than just using LLMs for code generation, is a radical but potentially transformative concept.

Topics to Explore: AI Agents, Scalability, Fuzzy Compute, LLMs, Backend Architecture



# Q&A

Q: What is the core idea behind scaffolding AI agents wisely?

A: Build systems that improve with more intelligence and compute, leveraging exponential trends.

Q: What are the different architectures for AI agents?

A: Classical compute only, fuzzy compute called from classical compute, and classical compute called from fuzzy compute.

Q: Why is it beneficial to allow LLMs to decide when to break into classical compute?

A: It leverages the strengths of both approaches and allows the system to scale with improvements in LLMs.

# **Creating Agents that Co-Create and Innovate**

# Summary & Key Points

Summary: Karina Nguyen from OpenAI discusses the scaling paradigms in AI research, including next token prediction and reinforcement learning on chain of thought. She explores the design challenges in creating new interaction paradigms with humans and highlights the potential of AI agents to become co-innovators through human-AI collaboration.

Counter-Intuitive Points: The idea that the interface to AI will become a blank canvas that self-morphs into your intent challenges conventional UI design.

Topics to Explore: AI Agents, Co-creation, Human-AI Collaboration, Reinforcement Learning, Chain of Thought

# Q&A

Q: What are the two main scaling paradigms in AI research?

A: The two main scaling paradigms are next token prediction (pre-training) and scaling reinforcement learning on chain of thoughts.

Q: What are the design challenges in creating new interaction paradigms with AI agents?

A: Design challenges include bringing unfamiliar capabilities into familiar forms, bridging real-time interaction with asynchronous task completion, and enabling human verification and editing of model outputs.

Q: How can AI agents become co-innovators through human-AI collaboration?

A: AI agents can become co-innovators by combining reasoning, tool use, and long context with creativity, enabled through human-AI collaboration, leading to new knowledge creation.

# **Educating the Next Generation of AI Engineers**

# Summary & Key Points

Summary: Tafania Druga from Google discusses the importance of opening up AI knowledge to the next generation of AI engineers, starting at a young age. She highlights the potential of multimodal AI to transform education and shares her work on Cognomates and other projects that enable kids to learn about AI by building games and training their own models.

Counter-Intuitive Points: The idea that kids are actually little scientists and can engage in the scientific process if given the right tools is a valuable insight.

Topics to Explore: AI Education, AI Literacy, Cognomates, Multimodal AI, Scratch Programming

# Q&A

Q: Why is it important to start AI education at a young age?

A: To ensure AI literacy, as it's now part of the EU AI Act, and to demystify AI, allowing young people to shape its future.

Q: What are some of the key features of Cognomates?

A: It expands scratch to allow children to learn about AI by building games, training their own AI models, and programming hardware, using a visual programming language.

Q: How can AI be used to support creativity in teachers?

A: By designing agents that support creativity in teachers, such as allowing them to see simulations of agents programming or generating assets directly on the stage.

# **Building Personal Local Private AI Agents**



# Summary & Key Points

Summary: Sumith Chintala from Meta discusses the challenges and opportunities in building personal, local, and private AI agents. He emphasizes the importance of context, the limitations of cloud-based services, and the need for catastrophic action classifiers. He also highlights the potential of open-source models and the importance of privacy.

Counter-Intuitive Points: The idea that you should be able to run your own AI agent locally and privately, rather than relying on cloud-based services, challenges conventional wisdom.

Topics to Explore: AI Agents, Privacy, Local AI, Open-source Models, Catastrophic Action Classifiers

# Q&A

Q: Why is it important to build personal AI agents that are local and private?

A: To maintain control over personal data, avoid vendor lock-in, and prevent potential punishment for private thoughts.

Q: What are the technical challenges in building personal AI agents?

A: Slow and limited local model inference, and the need for better open multi-modal models for computer use.

Q: What are catastrophic action classifiers and why are they important?

A: They identify potentially irreversible or harmful actions an agent might take, notifying the user to prevent unintended consequences.