# Designing AI-Intensive Applications - swyx
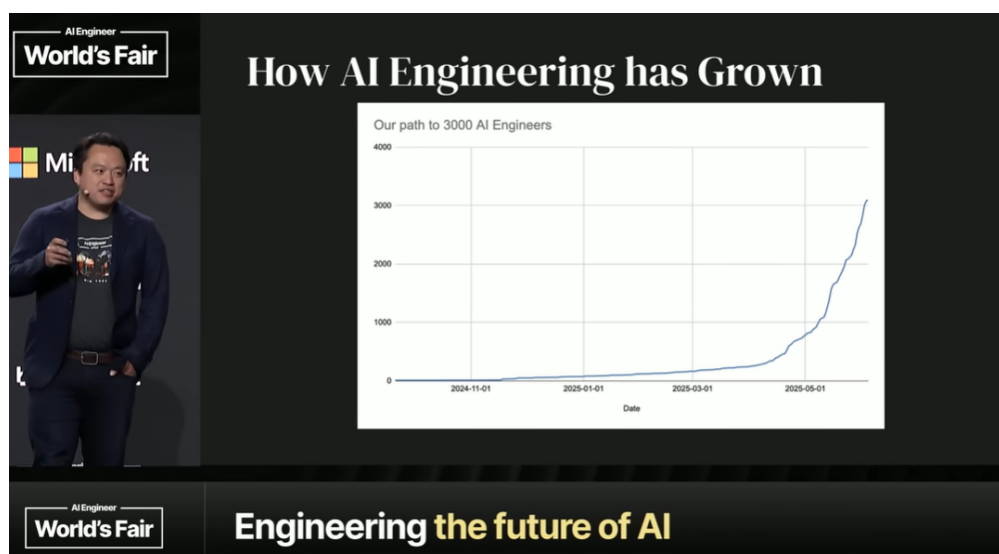
~ Prepared by Riya Kharade

---

## ➢ Introduction



The speaker welcomes everyone to the conference. He explains that this talk will cover the current state of AI engineering in 2025. He says the conference is growing and more topics are being added to help everyone learn more. The goal is to answer questions about AI, the conference, and how things are changing.
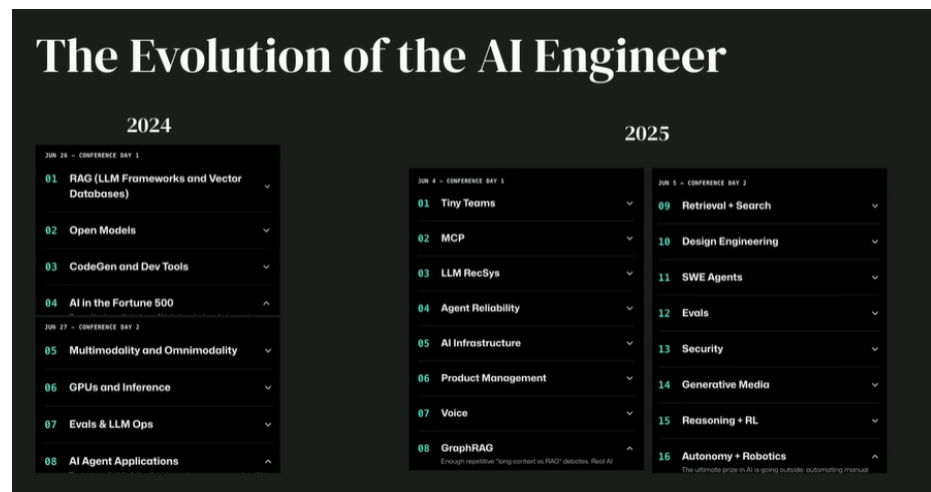
## ➢ Background

- Over 3,000 people registered for the conference, most at the last minute.

- The number of conference tracks has been doubled to give more value.

- The conference tries to cover all areas of AI engineering.

- Organizers use surveys and feedback to create content based on the community's needs.

- They are more technical than other events and respond quickly to attendees' requests.

➢ **All Points Discussed**
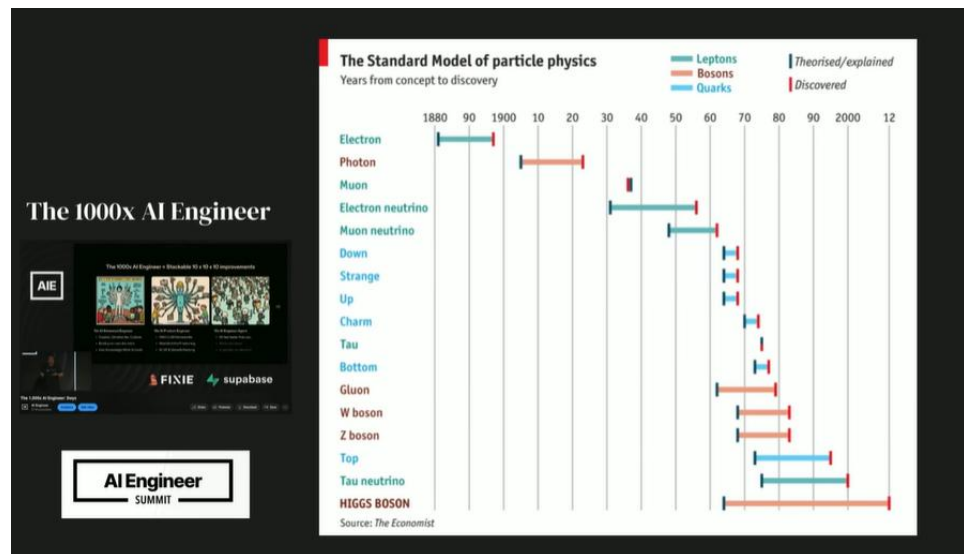
1. **AI Engineering is Growing**



- AI engineering was once mocked and considered a small field.

- Now it is respected and seen as a high-paying career.

- It's compared to the growth of physics in the early 1900s when big discoveries were made.
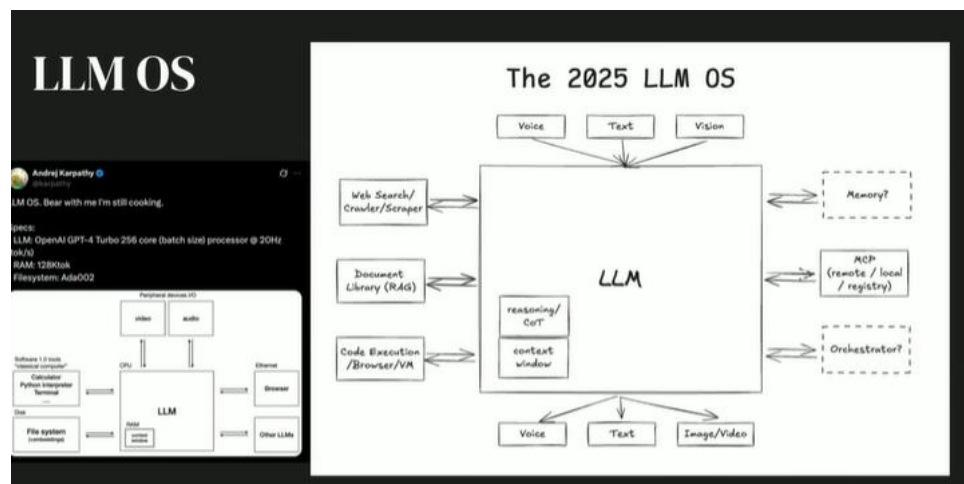
**2. The Need for Standard Models**

- Engineering fields use standard frameworks like MVC, ETL, MapReduce.

- The speaker asks: What will be the standard model for AI engineering?

- Having models helps organize ideas and make sense of complex tasks.
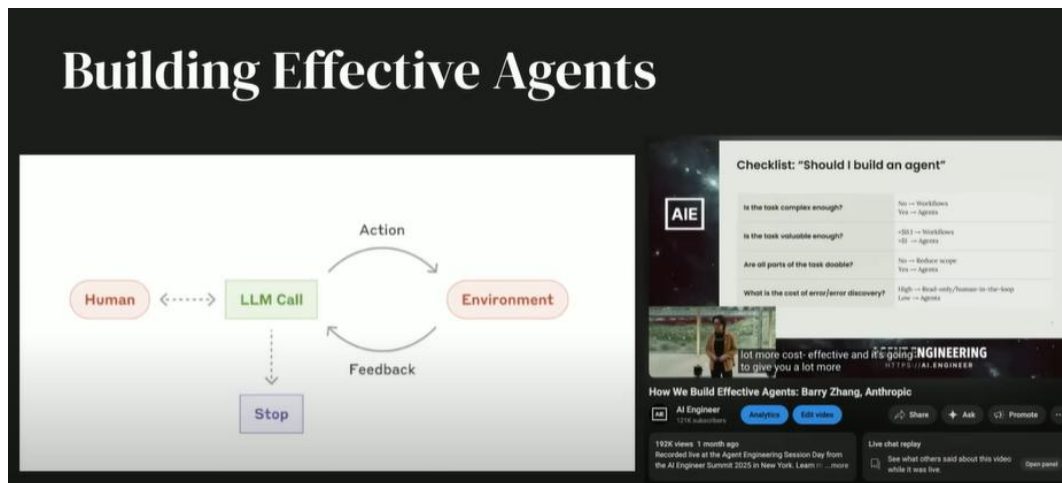
2. **Examples of AI Standard Models**



- **LM OS (Language Model Operating System):** Helps in building systems that use language tools.



- **MCP (Multipurpose Conversational Protocol):** Connects users to AI tools like chatbots and voice assistants.

- **AI SDLC (Software Development Life Cycle):** Guides how AI products are built, tested, and maintained.

### 3. Building Effective Agents



### 4. Standard Models in Engineering



### 5. Challenges in AI Engineering

- Early steps like training AI are becoming easier with free tools.

- Real challenges are in securing, evaluating, and scaling AI systems for businesses.

## 6. Agent Engineering



- An agent is a system that can perform tasks like planning and problem-solving.

- There is no single definition of what an agent is.

- Important parts include memory, planning, control flow, and tool use.

## 7. Human Input vs AI Output



- Instead of arguing about definitions, it's better to think about how much input humans give and how much output AI provides.

- This helps in designing systems that are efficient and useful.

**7. Mental Model of Input and Output**

- Simple tools like autocomplete need small human input and give small output.

- Advanced tools need more planning and give bigger, smarter output.

- Some systems can even operate with very little human involvement.

**8. AI News Example**

- The speaker built an AI news tool for himself and others.

- It is not technically an "agent" but still helps people by organizing and summarizing information.

- The process it follows is called **SPAD**:

    1. **Scrape** – Collect data from websites.

    2. **Plan** – Organize how to process the data.

    3. **Analyze/Summarize** – Understand and create summaries.

    4. **Deliver** – Present the information to users.

    5. **Evaluate** – Check if the information is useful and improve it.

**9. Scalable AI Systems**

- Systems can make thousands of AI calls every day to serve users.

- Structured processes like SPAD make this easier and more efficient.

**10. Beyond Text Outputs**

- AI tools can now create structured outputs like graphs and computer code.

- This makes AI even more useful for developers and businesses.

**11. Encouragement to Innovate**

- Attendees are encouraged to think about new models and improve their work.

- The field is still in its early phase with many opportunities to create something new.

## ➢ Key Takeaways

- AI engineering is growing fast and is now respected.

- Standard models help in organizing and building better systems.

- LM OS, MCP, and AI SDLC are important frameworks to guide work.

- The ratio of human input to AI output is a better way to think about systems than debating names.

- SPAD is a helpful method for building scalable AI products.

- AI tools can now produce structured outputs like code and graphs, not just text.

- Innovation is needed, and the community must work together to build better models.

## ➢ Conclusion

The speaker ends with a positive message. He encourages everyone to explore and define new models for AI engineering. The goal is to create systems that are helpful, efficient, and widely used. AI is still in its early days, and there are many chances to contribute, learn, and build something meaningful.