# AI Trends

Andrew Ng

# The AI Stack: Where are the biggest opportunities?



**APPLICATIONS**

Workhelix · BEARING AI · meeno · Woebot Health · kira.Learning · Workera · VALIDMIND · SpeechLab · credo|ai · SKYFIRE|AI · profitmind · common sense privacy · esteam

**Foundation Models**

OpenAI · ANTHROP\C · Meta

**Cloud**

aws · Google Cloud · Azure

**Semiconductors**

NVIDIA · AMD · intel.

Even though a lot of attention is on AI technology (esp. foundation models) most of the opportunities will be in building AI applications.

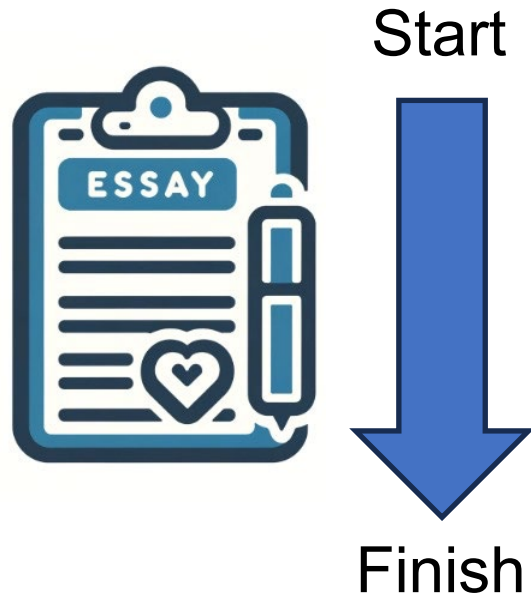This is driving significant demand for inference workloads (more, faster tokens).

Andrew Ng

Agentic AI workflows

# Agentic AI

## Non-agentic workflow (zero-shot):

Please type out an essay on topic X from start to finish in one go, without using backspace.

Start
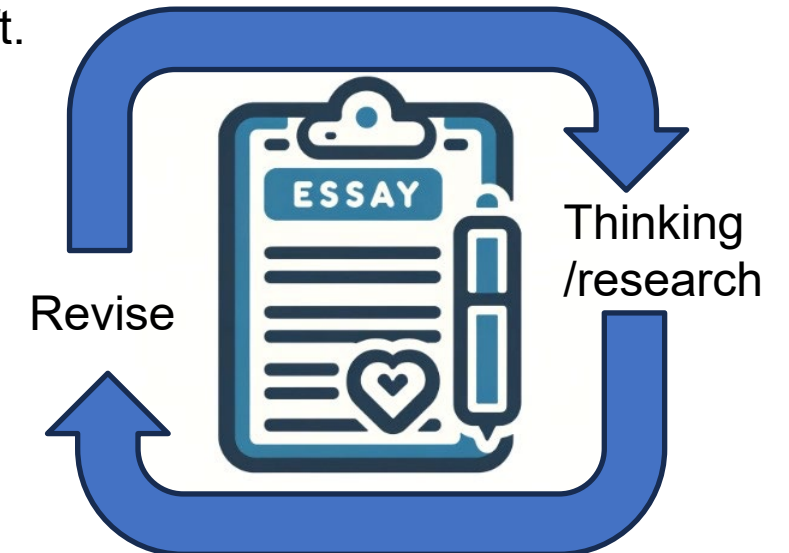
Finish

## Agentic workflow:

Write an essay outline on topic X

Do you need any web research?

Write a first draft.

Consider what parts need revision or more research.

Revise your draft.

....

Revise

Thinking /research

# The AI Stack: Where are the biggest opportunities?

**APPLICATIONS**

Workhelix · BEARING AI · meeno · Woebot Health · kira.Learning · Workera · VALIDMIND · SpeechLab · credo|ai · SKYFIRE|AI · profitmind · common sense privacy · esteam

**Agentic Orchestration Layer**

LangChain · crewAI · AG

**Foundational Models**

OpenAI · ANTHROPIC · Meta

**Cloud**

aws · Google Cloud · Azure

**Semiconductors**

NVIDIA · AMD · intel

Even though a lot of attention is on AI technology (esp. foundation models) most of the opportunities will be in building AI applications.

← New agentic orchestration layer

Andrew Ng

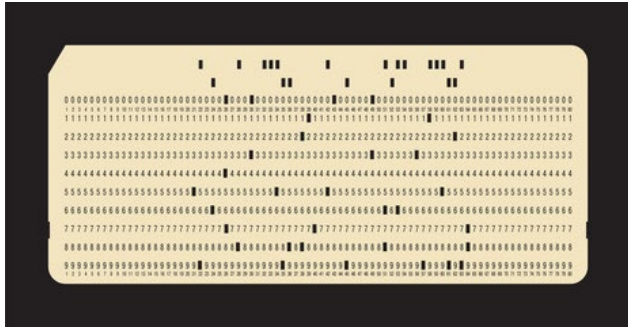# Important AI trends

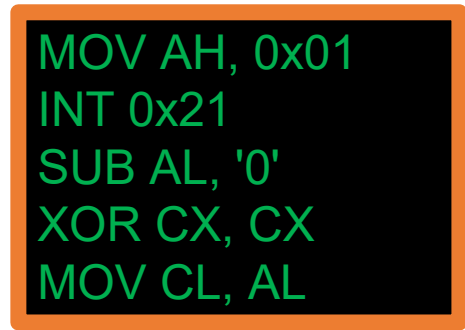# Five important technical AI trends

1. AI coding assistance

2. Fast prototyping

3. Visual AI

4. Data engineering

5. Talent gap

Andrew Ng

# 1. AI coding assistance − making coding pervasive

Some are advising people not to learn coding on the grounds AI will automate it. This is bad career advice.
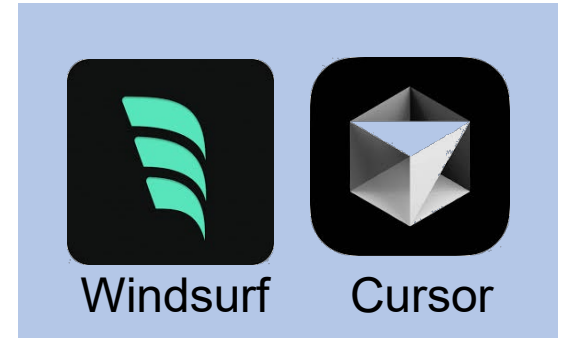


Punch cards
→ Keyboard

Assembly
→ Modern languages

Text editor
→ IDE

AI coding assistance

As coding becomes easier, people should code a lot more!

For professional engineers, this is already significantly boosting efficiency.

Andrew Ng

# Getting computers to do what you want



From: Generative AI for Everyone

Andrew Ng

# 2. Fast prototyping

| Build Prototypes | Write/maintain production software |
|:---:|:---:|
| **10x faster!** | **30-50% faster?** |

- Standalone prototypes/applications require less integration with legacy data sources and infrastructure
- Lower requirements on reliability, scalability, or even security, if prototyping to test basic functionality

- To pursue innovation, you can build 20 prototypes to see what works    Andrew Ng

# 2. Fast prototyping

Build Prototypes

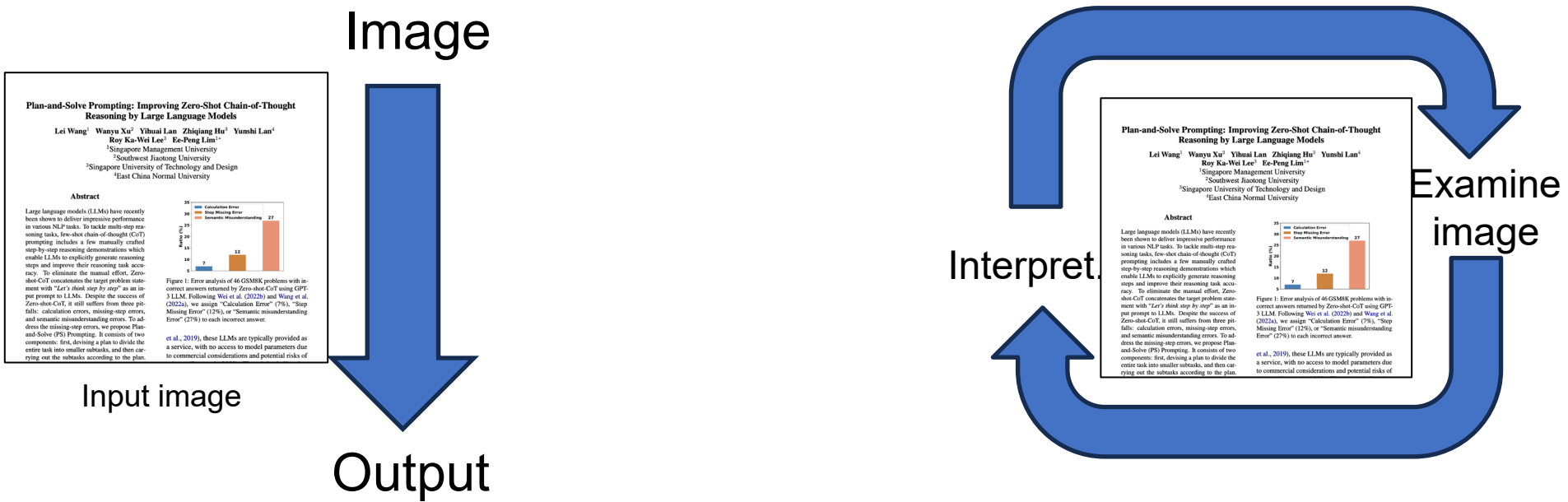Write/maintain production software

10x faster!

30-50% faster?

~~Move fast and break things.~~

Move fast and be responsible.

Andrew Ng

# 3. Visual AI

- The text processing revolution has arrived. The image processing revolution is coming, and will enable many new visual applications in manufacturing, self-driving, etc.
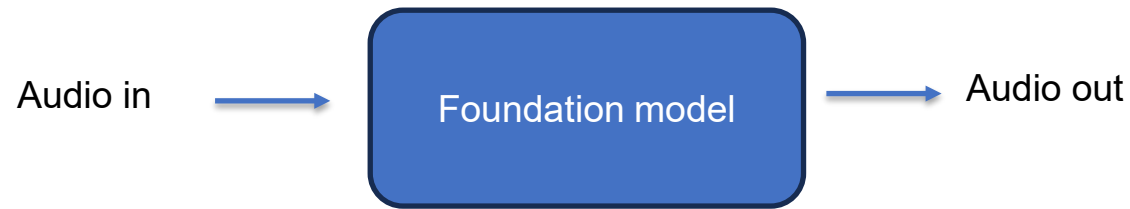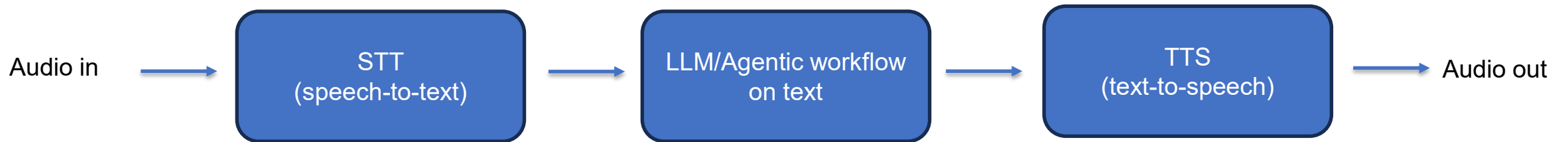
**Agentic Document Extraction**

Image

Input image

Output

Examine image

Interpret...

# 4. Voice stack

- Building voice applications is much easier than a year ago.

- Voice stack options:

## Direct generation

Audio in → **Foundation model** → Audio out

## Voice pipeline based generation

Audio in → **STT (speech-to-text)** → **LLM/Agentic workflow on text** → **TTS (text-to-speech)** → Audio out

Andrew Ng

# 5. Data engineering

- Data engineering's importance is rising, particularly on management of text, images, video, audio (also called unstructured data)

- Getting the data to AI for processing has high value

Data gravity is decreasing.

- Data Gravity is the idea that data tends to attract other data and compute. You would not transmit 1TB of data across clouds for processing.

- But for GenAI workloads, transmission costs are dwarfed by processing costs:
    - 1GB of data
    - Cost of transmission: ~$0.10
    - Cost of processing: ~$30-40 (gpt-4o-mini prices)

- This is leading to highly distributed software architectures, where we bring together many best-of-breed AI services.

Andrew Ng

# Talent gap

- Shortage of skilled AI engineers continues

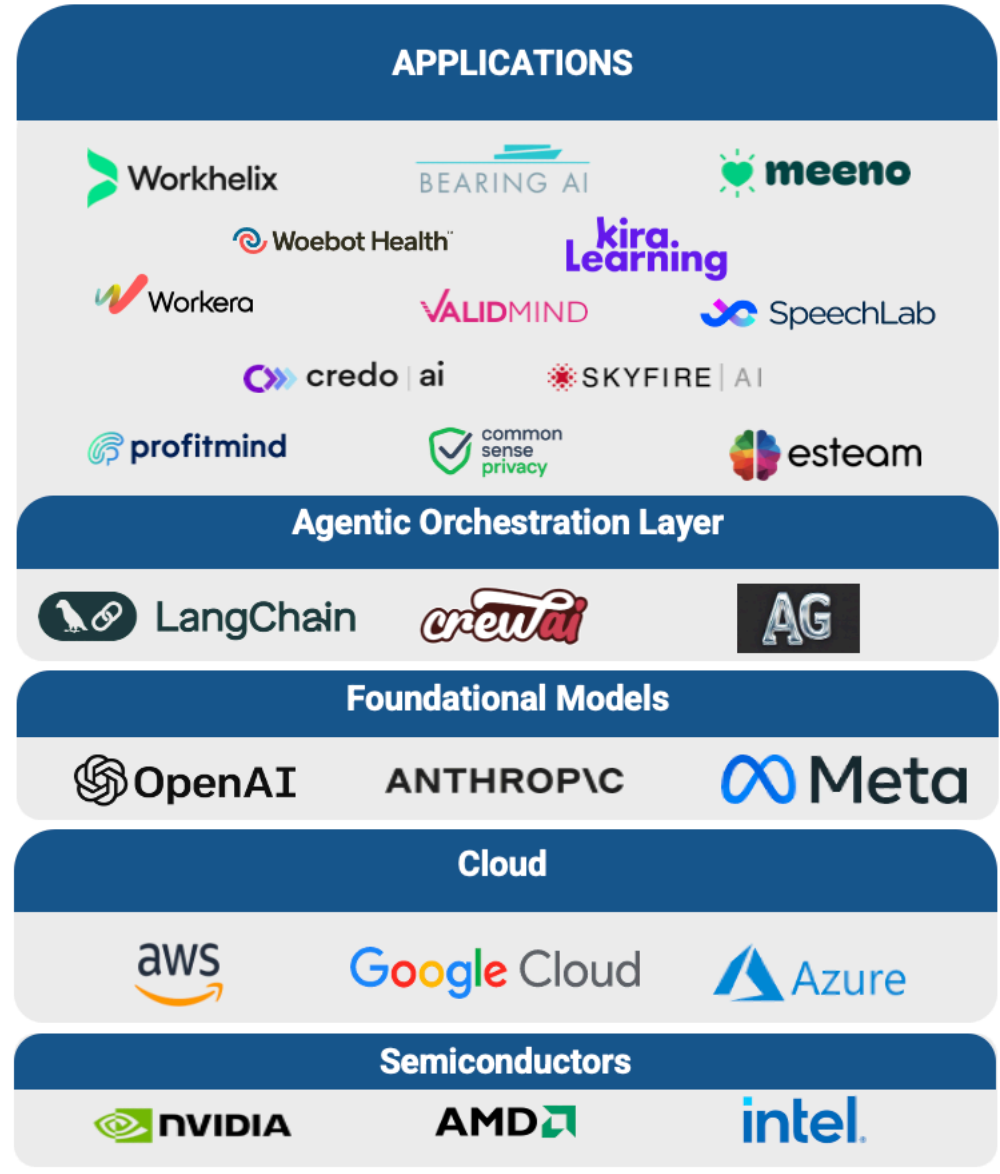- Easily 10x difference in productivity between skilled and less skilled AI engineers.



Kirsty Tan
AI Aspire

Andrew Ng

# Five important technical AI trends

Technology trends are making it easier to identify and build AI applications.

1. **AI coding assistance.** This is software development much more pervasive and efficient .

2. **Fast prototyping.** Generative AI is making it possible to build AI prototypes very efficiently. This is changing how we innovate and invent new things.

3. **Visual AI**. The text processing revolution has arrived. The image processing revolution is coming, and will enable many new visual applications in manufacturing, self-driving, etc.

4. **Voice stack**. There will be many compelling voice-based applications.

5. **Data engineering**'s importance is rising, particularly on management of unstructured data (text, images, video, audio).

Andrew Ng

# Conclusion



This is a wonderful
time to build!

Andrew Ng

# END