# Artificial Intelligence
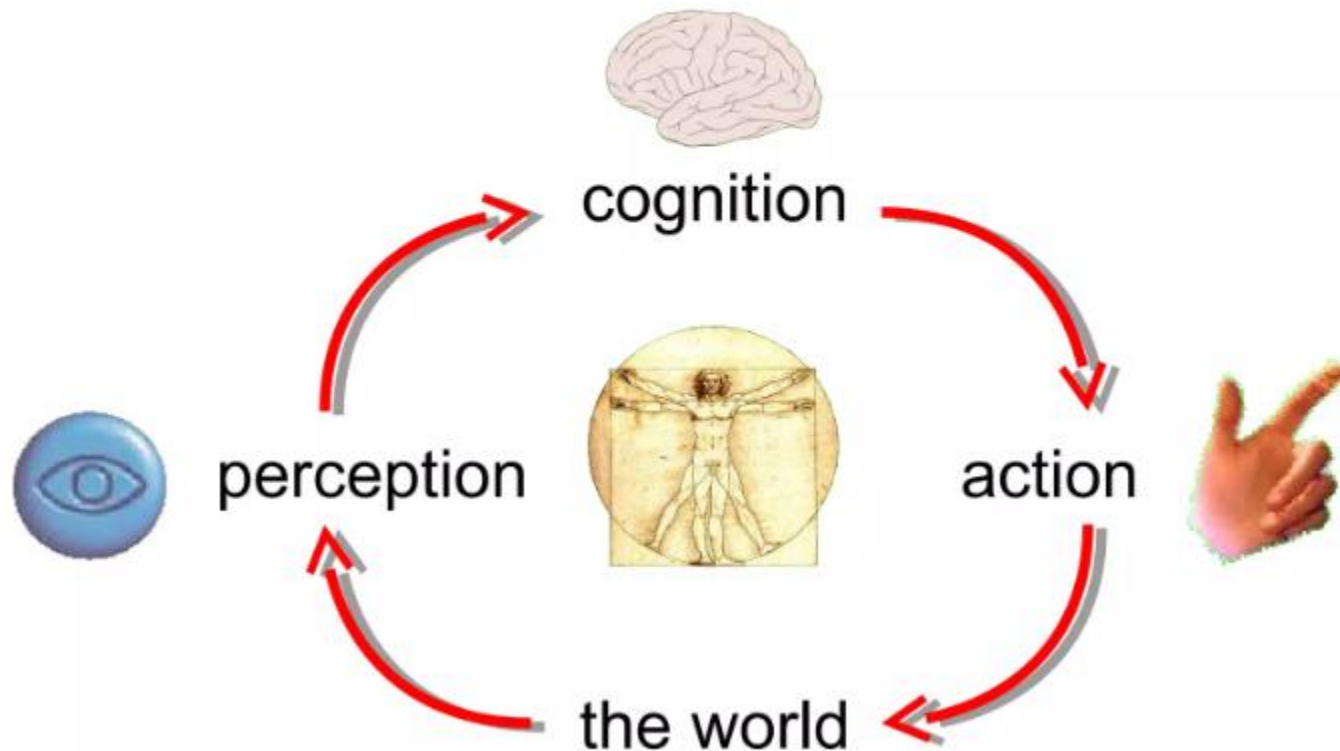# (Big-picture)

# Human Agent
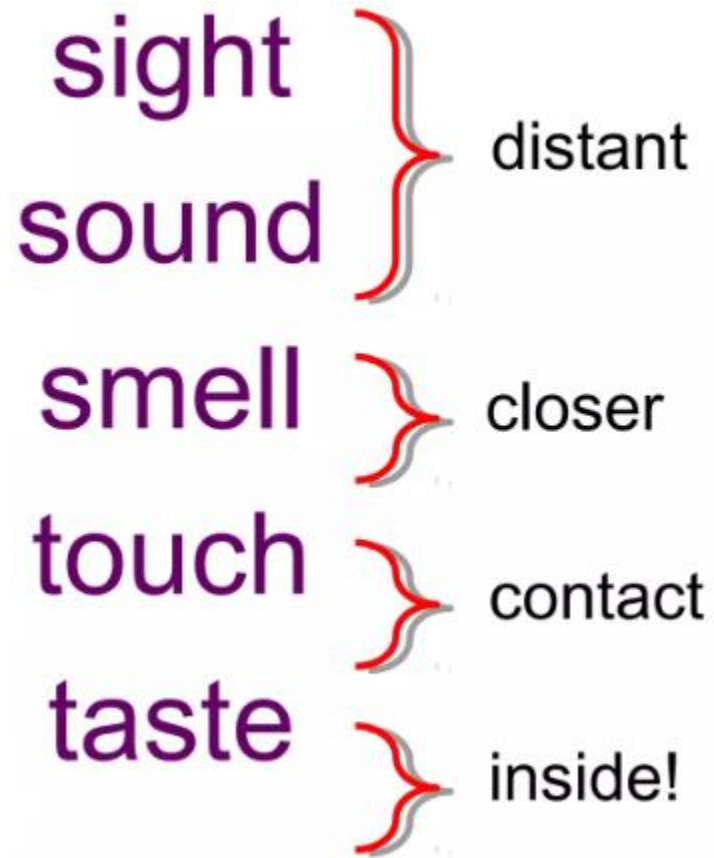
# Human Agent
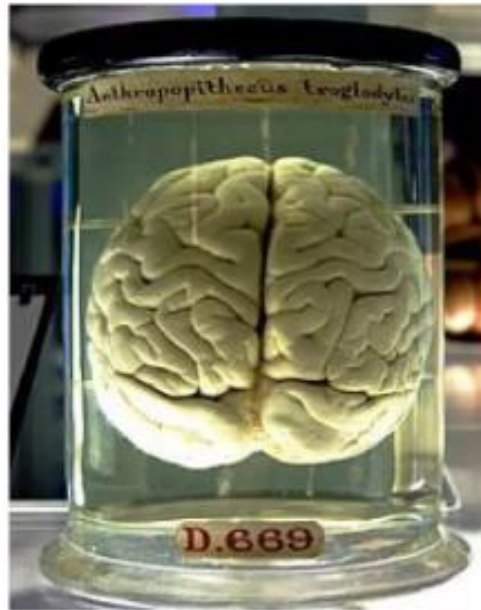
- Perception

- Cognition
  - Memory
  - Thinking
    - Learning -> Knowledge
    - Reasoning -> Decision Making
    - Imagination

- Action

- Autonomous

# Perception

sight
sound
} distant

smell
} closer

touch
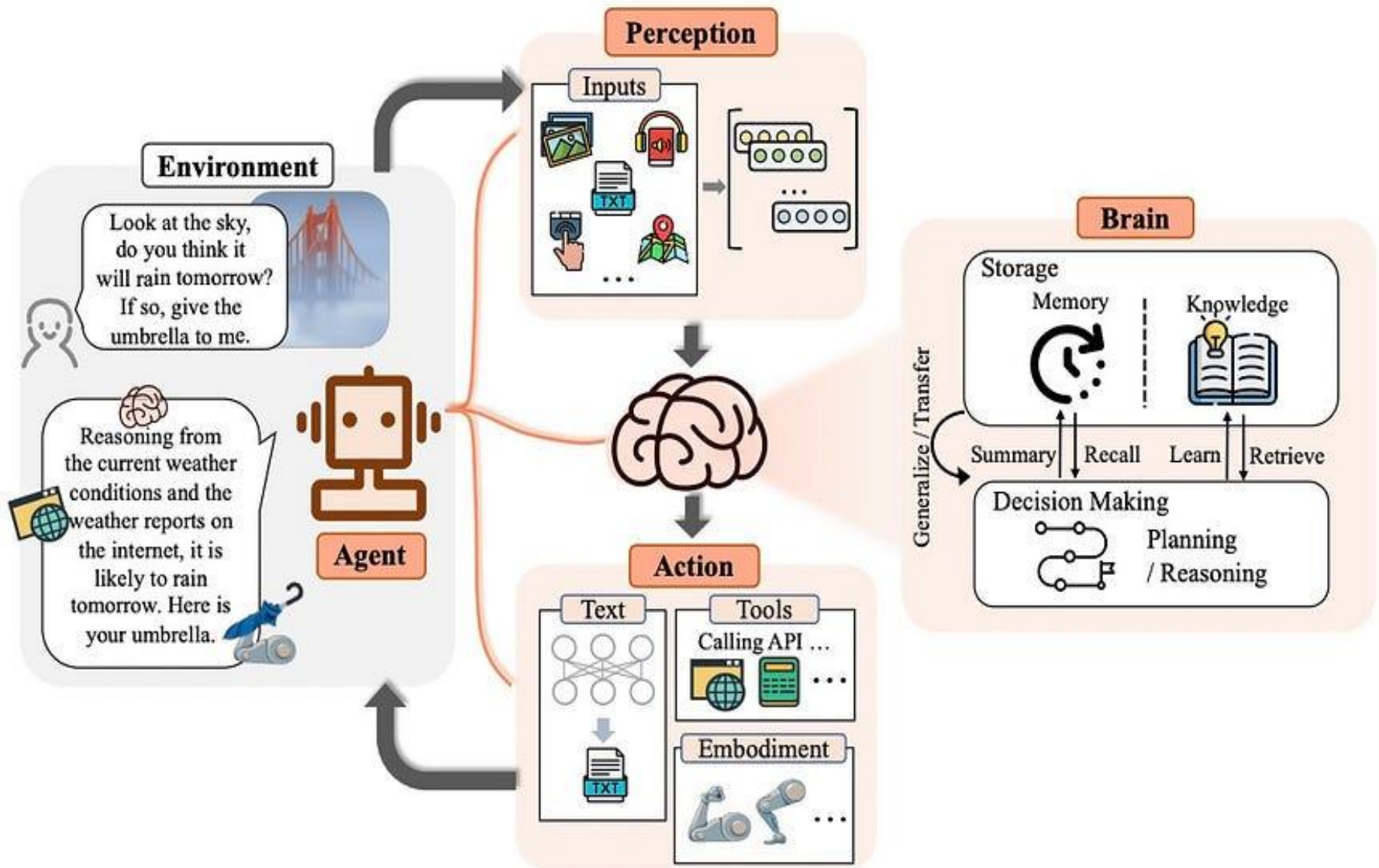} contact

taste
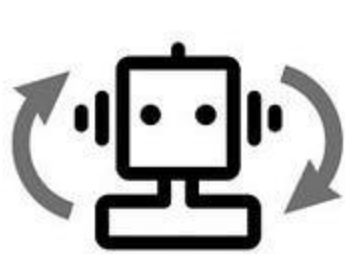} inside!

# Cognition

memory and thinking

# Artificial Intelligence

- A discipline of computer science used to create intelligent agents

- Intelligent Agents are systems that can learn, reason and act autonomously

# AI Agents

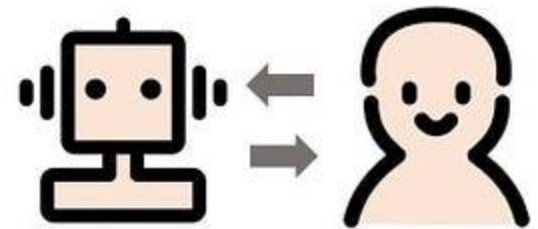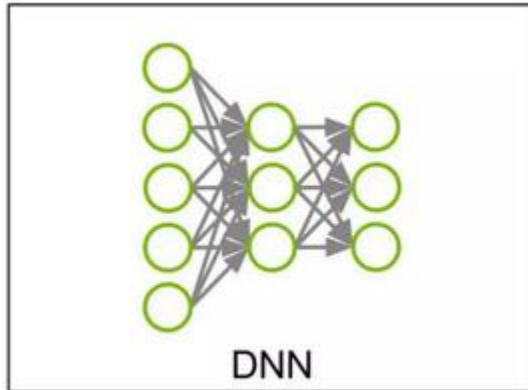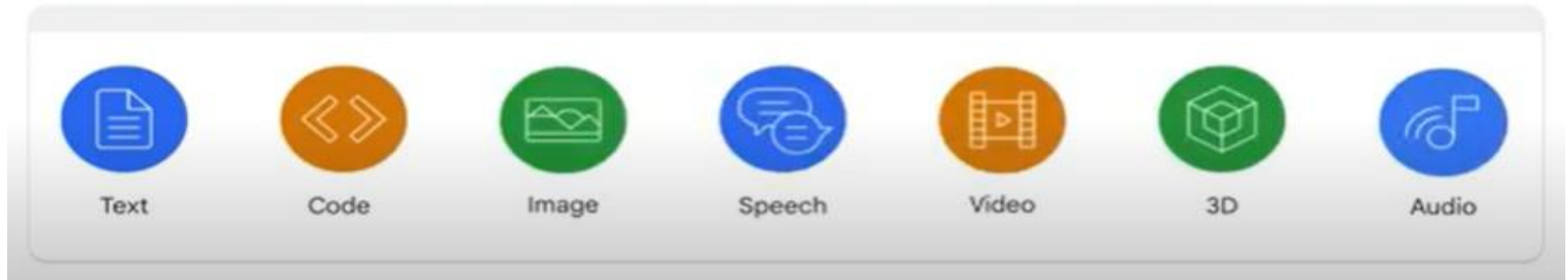# Agent Usage Scenarios



**Single Agent**    **Agent-Agent**    **Agent-Human**

# Fuel for AI Success



DNN

BIG DATA

GPU

- BIG DATA

- Deep Models

- Computing Hardware Innovations

# BIGDATA
# (Multimodal)

Text  Code  Image  Speech  Video  3D  Audio

## GPT-3 training data

| Dataset | # tokens | Proportion within training |
|---|---|---|
| Common Crawl | 410 billion | 60% |
| WebText2 | 19 billion | 22% |
| Books1 | 12 billion | 8% |
| Books2 | 55 billion | 8% |
| Wikipedia | 3 billion | 3% |

# GPU

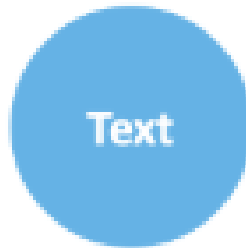- ## NVIDIA H100 Architecture

  - 80B transistors

  - >15K cores

  - Optimized for parallel multiply/add operations

  - >1 TeraFLOPS at 16bit precision

  - >900 GB/s bandwidth



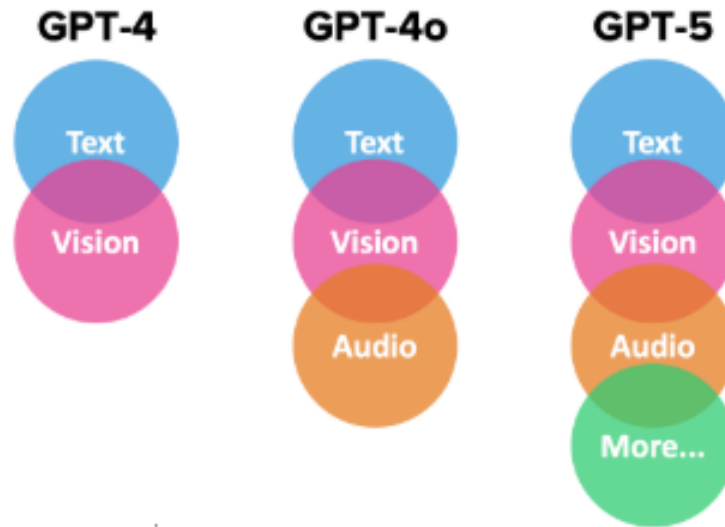**HOPPER H100 TENSOR CORE GPU**
80B Transistors, TSMC 4N

# Large Language Models (LLM)

**GPT-3**

Text

G BERT

# Large Multimodal Models (LMM)
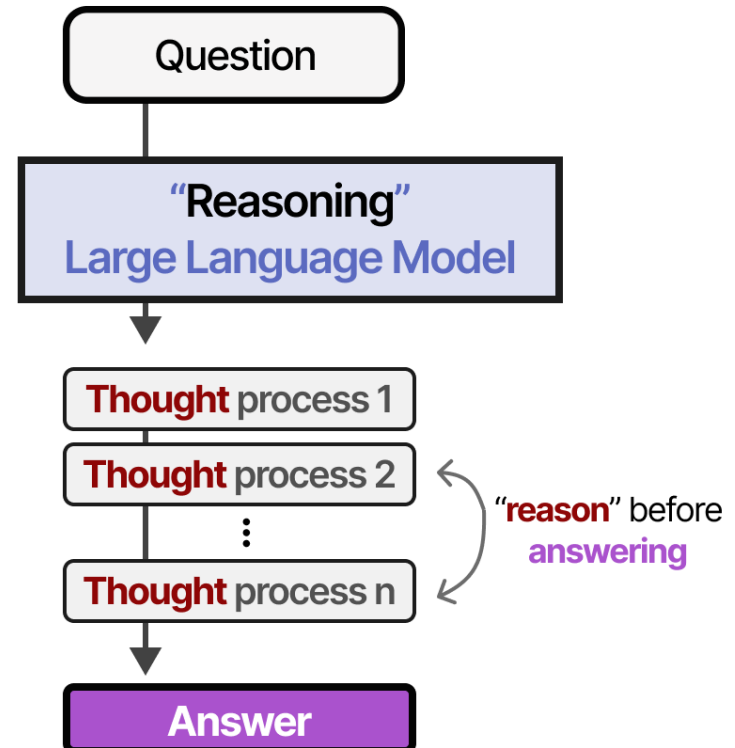
# Large Reasoning Models (LRM)

DeepSeek-R1

OpenAI o3-mini

Gemini 2.0 Flash Thinking

# Large Reasoning Models (LRM)

## "Regular" LLMs

Question

Large Language Model

**Answer**

## "Reasoning" LLMs

Question

"Reasoning"
Large Language Model

**Thought** process 1

**Thought** process 2

⋮
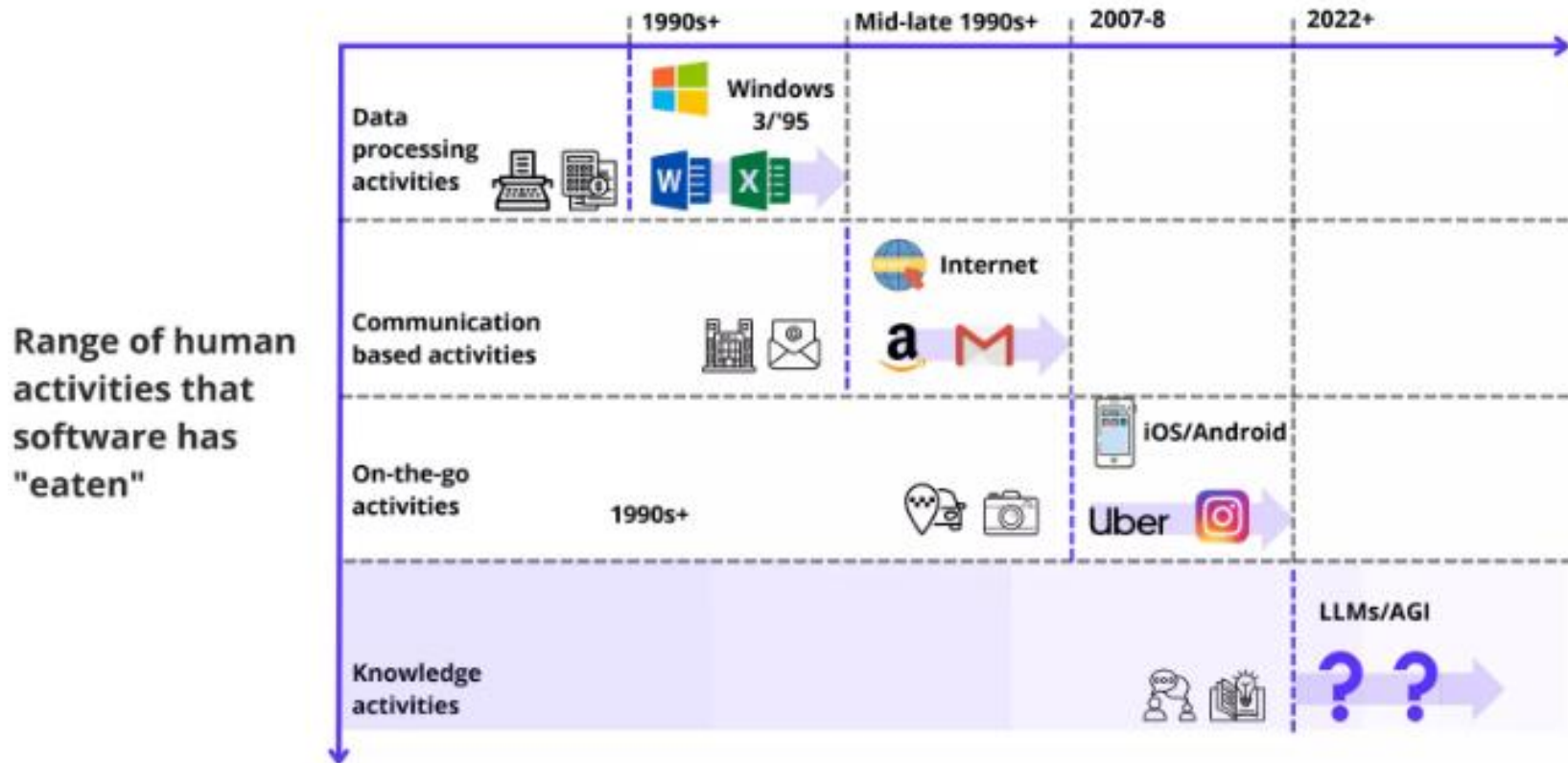
**Thought** process n

**Answer**

"**reason**" before
**answering**

# App development breakthroughs



Major platform launches that have enabled new types of applications, over time

LLM/LMMs are the Engine & AI Apps are the Product