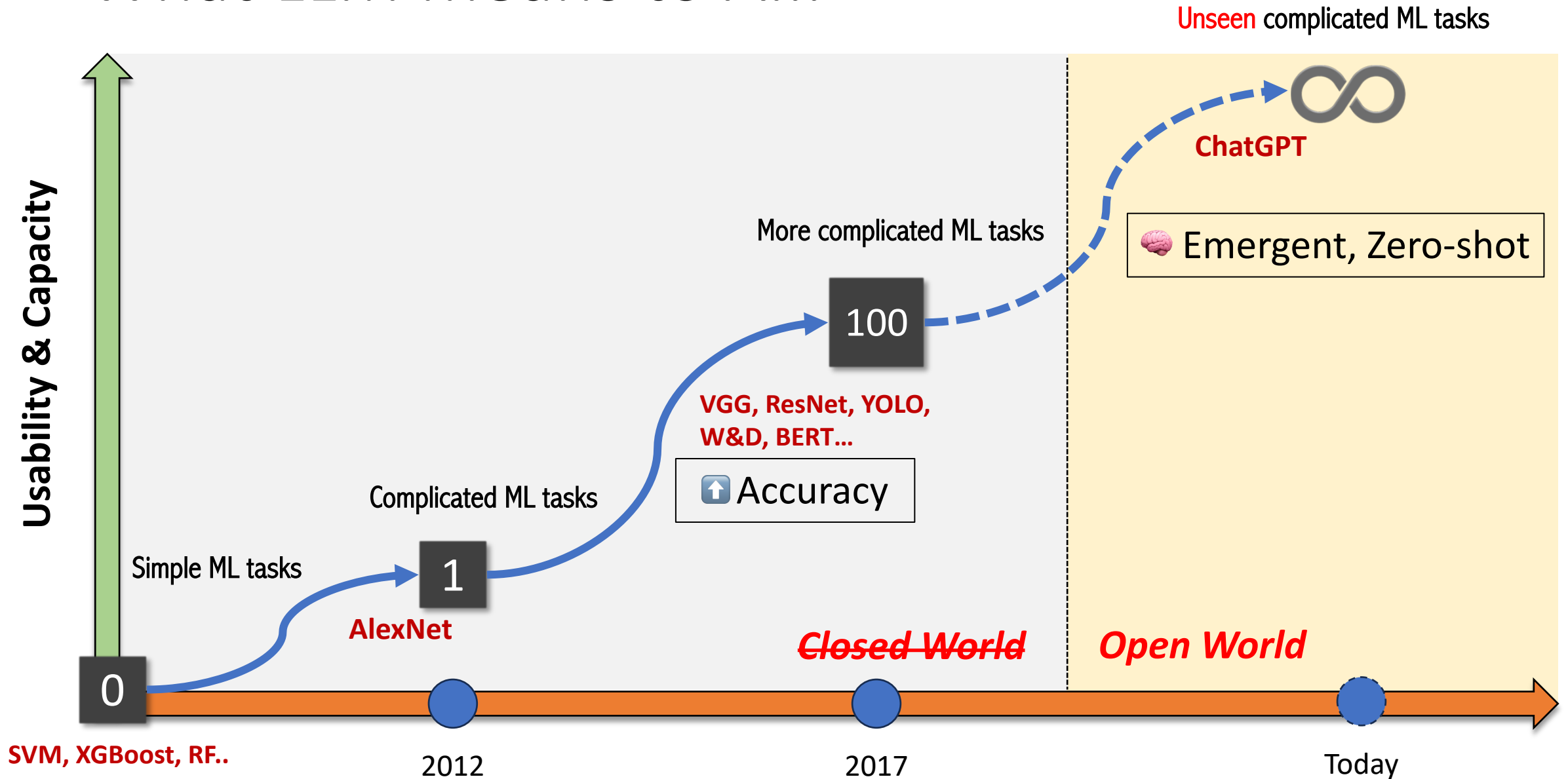


# What ChatGPT means to AI..

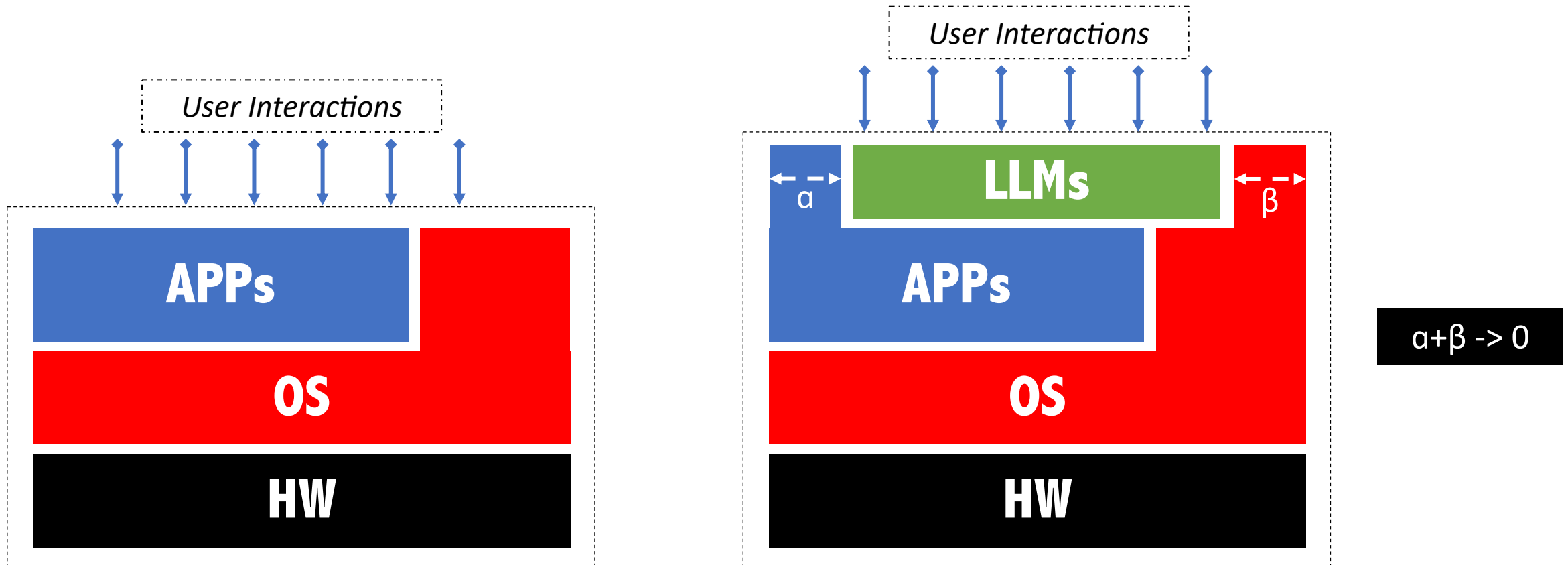
- “ChatGPT is just a smarter chatbot”
  - As a product, yes
    - But think about it: Moss is also a chatbot; robots/humans are chat bots with physical ability
  - As a research, hell no
    - It is a generative model that theoretically knows everything on Internet and can accomplish any NLP tasks
  - It's also
    - a series of papers cited by 10,000 times
    - a startup company worthy of 30,000,000,000 dollars.
  - It's also the one who opens the Pandora's box

# What LLM means to AI..



# LLM is the new Operating System

- Users interact with LLM, while LLM manages/utilizes old-time apps/OS and hardware



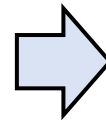
# The New Golden Era for Mobile Research

- Since iPhone 2007..
- The next long-term goal of mobile research: ChatGPT on smartphone
  - Takes 5 ~ 10 years
  - Takes collective efforts from hardware/architecture, mobile system, ML algorithm communities
  - LLM on edge vs. LLM-as-a-Cloud-Service
- **Old stories:** data privacy, low delay, low power consumption, etc..
- **New techniques:** memory-bounded LLMs, foundation model + adapters, generative and autoregressive, etc..

# Exploration Atop or Below LLM?

- Another way to go: build systems **for** LLM, or build systems **with** LLM
- When a software layer is finalized, most research/industry opportunities go above
  - Very very few system researchers rebuild OS now
  - Very very few network researchers rebuild network stacks now

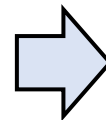
Auto-GPT, Agent-GPT, babyAGI, HuggingGPT,  
Web LLM, CAMEL, GPTRGB, PandaGPT..



*Easier to handle, potentially high impacts,  
but more crowded and competitive*

**LLMs**

GPTQ, Mixture-of-Experts, [EuroSys'23] Tabi,  
[MLSys'23] Flex, [OSDI'22] Orca, [ATC'22] PetS..



*More fundamental, potentially extremely-high  
impacts but technically/financially challenging*