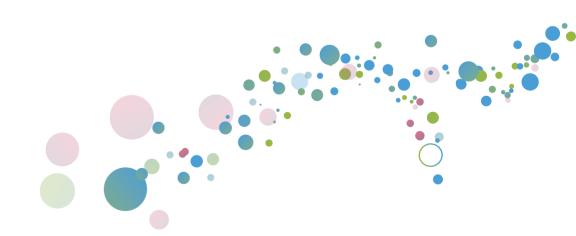


AI DC White Paper



A Guide for CIOs to Planning and Building AI DCs



Foreword

Computing power is becoming the new "black gold".

More than a decade ago, Time magazine noted that, after oil, network bandwidth would become the new black gold of the 21st century. At that time, no one could have foreseen how quickly artificial intelligence (AI) would evolve into the powerful force that it is today. With the impressive range of foundation models and the recent surge in generative AI, AI development is experiencing nothing short of a de Laval nozzle effect. As the global AI industry continues to flourish and mature, humanity will move towards an intelligent world at speeds far beyond our imagination. The revolutionary and far-reaching impact of computing power affirms its value to society and solidifies its status as the new "black gold" of the 21st century.

Al is a trend, not just hype.

Since the definition of artificial intelligence was first proposed in 1956, AI has gone through multiple ups and downs. However, even though AI dominates global technology headlines, a considerable number of people and organizations still remain skeptical, concerned, and uncertain about the future of AI. Meanwhile, AI is making strides in technical breakthroughs and industry scales. AI applications are expanding from standalone cases to diversified adoption, and evolving from general-purpose applications to industry-specific implementations. AI is set to reshape traditional industries and create many new industries.

The emergence of ChatGPT has made the pathways to artificial general intelligence (AGI) more clear-cut than ever before. Increasingly becoming a major driver of innovation and development, AI is shaping up to be a trend rather than a hype. Driven by AI, society will evolve from being data- and information-centric to being knowledge-and wisdom-centric. In the next few decades, we will undergo a cognitive revolution. Today's generative AI may just mark the start of this journey.

This is a white paper on data centers, not on Al.

When a vast range of models and applications are becoming an increasingly dominant presence, and industries are rushing to embrace intelligent transformation, the top priority for industry CIOs should not be AI applications, but rather AI infrastructure. China has this saying, "To get rich, build roads first." This proverb rings true worldwide. For any country or enterprise to capitalize on the opportunities represented by AI, they must first make sure they have a solid AI infrastructure in place. Data centers are at the very heart of this AI infrastructure.

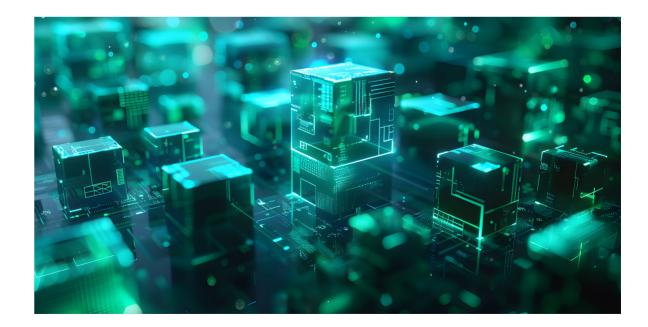
Data centers first began to emerge around 1940. Since then, and especially with the development of the Internet, big data, and cloud computing, data storage and processing have become increasingly important. As a result, data centers have become part of the core infrastructure for enterprises' informatization and digitalization. As the world moves towards an intelligent era, data centers are expected to provide computing power, carry AI training and inference workloads, and support key intelligent applications. We call these future-facing data centers AI data centers (AI DCs).

Tomorrow's data centers will be defined by AI.

Al DCs are not just a simple upgrade and makeover of traditional data centers. Rather, they are the refactoring of data centers in every aspect. This means a paradigm shift from cost centers to production centers, and from data storage and processing centers to value-generating centers.

The Internet and cloud computing have made software-defined infrastructure a mainstay. In the future, data center infrastructure will be defined by AI. Of course, AI will bring many different challenges, such as challenges related to computing power density, energy efficiency, AI-powered O&M, and sustainable development.

A powerful and solid AI computing backbone is the bedrock of intelligent transformation. As the functions of data centers evolve from storing data and supporting applications, to providing computing power and carrying AI training and inference workloads, and all the way through to generating and enabling intelligence, data centers are adding more and more value to industries and attracting attention from players along the entire industry chain. With these considerations in mind, we have written this white paper. We hope that it can serve as a helpful reference for planning and building AI DCs.





Yang Chaobin

—— Director of the Board, President of ICT Products & Solutions, Huawei

Recently, I have met with many customers, partners, and peers from the AI ecosystem and discussed the implementation of AI in enterprises and the importance of AI computing infrastructure. A common consensus is that AI DC construction should be prioritized as a strategic initiative for enterprises' intelligent transformation. However, AI DCs have some key differences that set them apart from traditional data centers. AI DCs are drastically different in terms of their positioning in enterprises' digital and intelligent infrastructure, the services they carry, and the requirements for their data processing and computing power provisioning capabilities. In addition, technologies are being continually innovated and upgraded. The entire industry should look for ways to build high-quality AI DCs efficiently. This white paper draws from practical experience and industry wisdom and can serve as a reference in this regard.



Li Peng

—— Corporate Senior Vice President, President of ICT Sales & Service, Huawei

In China, there is an old saying: "To get rich, build roads first." This concept holds particular weight in the world of ICT. High-quality ICT infrastructure is the cornerstone of digital and intelligent transformation and enterprise business success. AI DCs are the core of the newest generation of infrastructure powering digital and intelligent transformation. Huawei has years of experience and lessons learned from working with our customers on infrastructure construction and innovation exploration. Yet, there are still many new topics emerging every day that we must address together. We intend for this white paper to act as a starting point for these efforts. All industries will need to collaborate and innovate to jointly promote the development of AI DCs and help the world move towards a more intelligent era.



—— CIO, Aluminum Corporation of China Limited (Chinalco)

Industries are now actively embracing AI. By combining innovation, industry know-how, and foundation model capabilities, we are redefining production and organization in many traditional industries. The non-ferrous metal industry has been exploring how to use AI in our industry scenarios, with continuous practices in fields such as aluminum oxide, electrolytic aluminum, and high-end aluminum processing. This white paper provides many viewpoints that can be used as a reference for enterprises. It also provides suggestions and practical evaluation indicators for planning and building AI DCs – which lie at the core of all enterprises' digital and intelligent infrastructure. Properly planning and building AI DCs will be the most critical step for enterprises who want to implement AI.



Wang Lei

——Director of the Digital and Intelligence Research Institute, China Pacific Insurance (Group)

Generative AI is a new productive force for the insurance industry. Its scenario-specific implementation and commercial profitability are currently the core questions to answer. Enterprises have multiple pressing issues to consider when building and operating AI DCs, such as technical exploration, and the efficiency and cost of large-scale application deployment. This white paper draws from technology trends and industry practices, and systematically introduces construction strategies and implementation paths for AIGC industry applications. The white paper provides solutions for building AI DCs in different scenarios. It serves as a valuable, thought-provoking reference.



Zou Zhilei

—— Corporate Senior Vice President, Huawei

In the intelligent era, AI can unlock its immense value only when it is integrated into the core production scenarios of enterprises. This will drive enterprise service systems to shift from a traditional compositive style to a generative style. As the core of the digital and intelligent infrastructure, enterprise AI DCs will go from being cost centers to innovation centers, and its technical architecture will also undergo big changes. The construction mode, system architecture, and O&M of traditional data centers will likely change significantly. This white paper is a summary and reflection of current practices found within the AI industry. We need to take a more future-oriented approach to AI development and keep exploring more ways to plan and build AI DCs.



Ma Haixu

—— Corporate Vice President, President of ICT Product Portfolio Management & Solutions, Huawei

Computing power is the foundation of AI applications, and AI DCs are a key part of the digital and intelligent infrastructure that provides computing power. AI DCs will need to leverage the comprehensive advantages provided by computing, storage, network, cloud, and energy technologies, and turn to system architecture innovation to overcome the bottlenecks we face in large-scale computing power. Since the release of our first AI strategies and solutions in 2019, Huawei has worked with customers around the world to dive deep into AI computing infrastructure construction, and continuously pursued joint innovation with other industry stakeholders to build competitive products and solutions that will create even more value for customers. We have distilled our customers' valuable experience in infrastructure construction and the wisdom of multiple industries into this white paper. We hope these efforts will help more customers quickly and efficiently build AI DCs, and accelerate the intelligent transformation of industries.



He Baohong

—— Director, Cloud Computing and Big Data Research Institute, CAICT

As the world advances towards the intelligent era, AI DCs will serve as the cornerstones of our intelligent society. Governments at all levels in China have continuously taken measures in terms of distribution guidance, construction planning, technological innovation, and application enablement to promote the development of computing infrastructure. Enterprises are also accelerating research and development, and are promoting large-scale, high-quality, and application-facing AI DC development. This white paper offers systematic insights into the latest trends in the planning, construction, management, and use of AI DCs, and provides valuable guidance for industry development.



Lian Jye Su

—— Chief Al Analyst, Omdia

Al DCs carry Al application, training, and inference workloads, and this makes them very different from other types of data centers. Al is advancing at unprecedented speeds, and new technologies and applications are emerging one after another. Every enterprise needs to think about building a solid and reliable computing backbone to meet its long-term development requirements, and developing an infrastructure that can handle the iteration and evolution of Al applications.

Contents

Chapter 1

The General Vision for and Macroscopic Driving Forces of an AI World	1(
Al is a major new development direction	. 1
AI for All	. 1
The move towards AGI is driven by idealism and realism	. 1

Chapter 2

An "All in AI" Generative Service System	18
Certainty and uncertainty in enterprises' Al development	19
Architecture comes first: Turning uncertain challenges into definite opportunities	21
Application scenario centricity: Adopting a four-in-one framework to realize the value triangle	23
Data centers as the foundation	32

Chapter 3

Evolution and Changes of Data Centers in the Intelligent Era				
Evolving towards AI DCs	35			
Carrying the training, inference, and application of AI models	37			
Four AI DC construction scenarios and three AI DC types	39			
Five key features of AI DCs	43			
Data centers reshaped from layered decoupling to vertical integration	53			

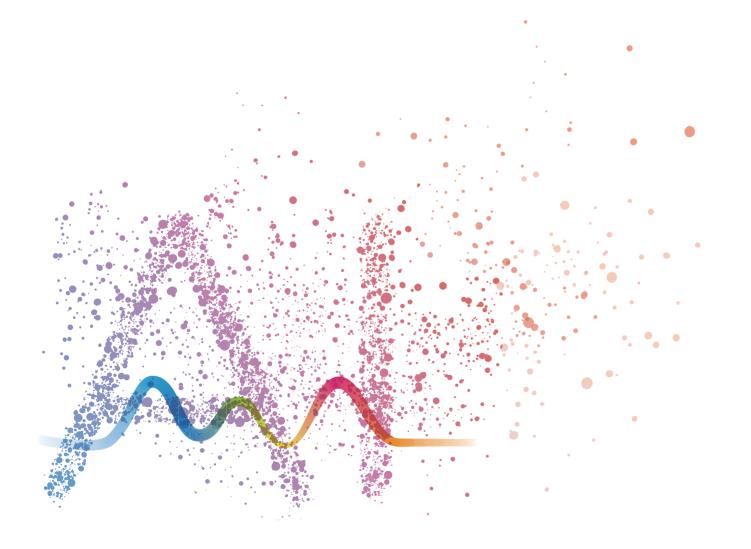
Chapter 4

Typical AI DC Planning and Construction	56
Ultra-large AI DCs	. 57
Large AI DCs	. 70
Small AI DCs	. 82

Chapter 5	
AI DC Construction and Development Initiatives	86
Initiative 1	
Taking a forward-thinking approach to AI DC construction	87
Initiative 2	
Building intensive and green AI DCs	90
Initiative 3	
Building an open and collaborative AI ecosystem	91
Initiative 4	
Ruilding three foundations to accelerate AL applications across industries	93

Chapter 1

The General Vision for and Macroscopic **Driving Forces of an** Al World



AI is a major new development direction

The rapid development of generative AI has catapulted AI into the spotlight.

According to reports from the China Academy of Information and Communications Technology (CAICT), as of July 2024, there were nearly 30,000 AI enterprises and 1,328 AI foundation models around the world. 478 of these foundation models were released by Chinese companies in less than two years.

Al is triggering a new round of industrial revolution along the entire industry chain and will bring enormous opportunities that will benefit humanity and drive social development. According to the Artificial Intelligence Index Report 2024 released by the Stanford Institute for Human-Centered AI (Stanford HAI), from 2023 to the first quarter of 2024, there were 234 AI unicorns worldwide. 37 of them

were classed as new unicorns, and they accounted for 40% of the total number of new unicorns. In 2023, although global investment in AI decreased to US\$189.2 billion, the investment in generative AI increased by a factor of almost 8 compared with investment levels in 2022, reaching US\$25.2 billion.

Sixty years of chip technology development, thirty years of Internet development, continuous breakthroughs in the transformer architecture, and abundant data combine to deepen AI technology advancements and bring AI applications to even more real-world scenarios. Following the launch of ChatGPT by OpenAI, Huawei launched Pangu Models 5.0 and Anthropic launched Claude 3.5 Sonnet in 2024, taking the application of foundation models one step beyond chatbots and into workflows.



AI is the culmination of the ICT industry's over 70 years of development

In 1956, John McCarthy, then an Assistant Professor of Mathematics at Dartmouth College, convened the now famous Dartmouth Summer Research Project on Artificial Intelligence, which is widely considered to be the founding event of artificial intelligence as a field. Since then, the field of AI has seen many ups and downs, including two "winters", but it has never stopped advancing.

Since Intel released the first microprocessor in 1971, Moore's Law has been witnessing and challenged by the booming development of the ICT industry. If we put together the development curves of the Al industry and ICT industry over the past 70 years, we can see that

there is significant correlation between the two, and that breakthroughs in academic research complement the development of engineering technologies. The two AI winters were due to the world's excessively high expectations for AI, which were way beyond the engineering capacity of the ICT industry at the time. As this metaphor indicates, winter is not an end. Winter eventually melted into spring, marking the beginning of a new round of AI development. Today, the world is ready to reap the rewards of this long-term innovation by and collaboration between leading academics and the global ICT industry over the past 70 years.

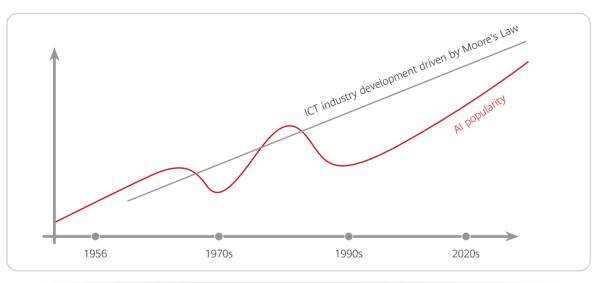


Figure 1-1 AI is the culmination of the ICT industry's over 70 years of development

The accurate positioning of technologies is the prerequisite for maximizing their value. This is essential for us to truly understand and adopt AI technologies.

At Huawei, we recognize AI as a combination of technologies that, together, form a new general-purpose technology (GPT). We have seen the rise of many GPTs before AI: The wheel and iron, which are both several thousand years old; railways and electricity in the 19th century; and automobiles, computers, and the Internet in the 20th century. In his book Economic Transformations: General Purpose Technologies and Long-Term Economic Growth, Canadian academic Richard G. Lipsey noted that new general-purpose technologies are the driving forces

behind sustainable socioeconomic growth. A GPT must have multiple uses, according to Lipsey, and have many technological complements (when two or more different technologies strengthen and reinforce each other) and spillover effects. Economists believe that to date, there have been 26 GPTs. Thanks to the cumulative development of the ICT industry over the past 70 years, AI has now become one of these GPTs.

Moving forward, we need to make full use of AI technology. We need to start reaping the benefits, and work hard to unleash its full value, so that it can thrive and we can continue to reap more benefits in the future.

AI will trigger transformations of an unprecedented magnitude this century

Throughout history, the large-scale application of general-purpose technologies has always been a catalyst for social change. Peter Diamandis defined AI as a leading technology in his book The Future Is Faster Than You Think. AI will lead to transformations of a magnitude unseen in a century. Since the emergence of steam engines in the 18th century, the history of technological innovation has been divided into the steam engine era, the industrial era, and the information era. Now, the intelligent era is coming – and fast. The driving force behind it is AI computing power. This force will not only offer more personalized and convenient experiences, but also increase efficiency and redefine empirical methods across many

different industries, opening up new pathways for scientific research. The popularization of AI will not only accelerate the intelligent transformations of traditional industries, optimize resource allocation, improve decision-making, and stimulate product and service innovation, but also further optimize our social and economic structures and lead the global economy into a new period of high-quality growth.

The transformations triggered by AI will revolutionize experience, efficiency, empirical methods, and scientific research, and usher in a new era characterized by intelligence.

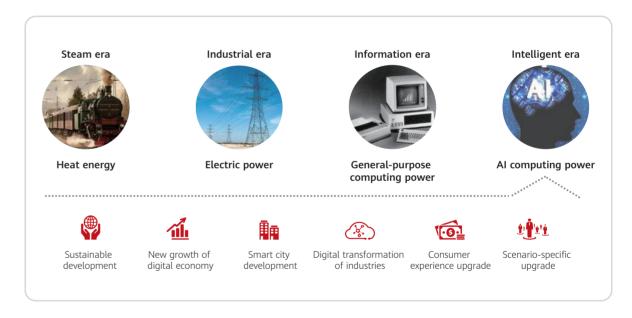


Figure 1-2 Human society is entering an intelligent era

The intelligent economy will be the next milestone in digital economy development

The global digital economy is continuing to develop rapidly. According to a CAICT report, the digital economies of the USA, China, Germany, Japan, and South Korea were worth more than US\$33 trillion in 2023, and had an average annual growth rate of over 8%. The digital economy accounted for around 60% of these countries' respective GDPs. This not only demonstrates the rapid advance of the digital economy, but also highlights its core role in the global economic landscape. AI has played a critical role in driving large-scale economic development.

The evolution of the digital economy began with the invention and popularization of personal computers, and then matured with the development of the Internet

of Things (IoT) and mobile Internet. Today, the digital economy is entering a new stage of intelligent economy with AI as the core. The intelligent economy is an economic structure and growth model that is focused on the goals of efficiency, harmony, and sustainability. The intelligent economy adopts a framework comprising physical devices, computer networks, and best practices, and key players include intelligent governments, advanced economies, and society at large. As a new engine of the global economy, the intelligent economy is committed to driving efficiency improvements, harmonious development, and sustainable growth in the global economy.

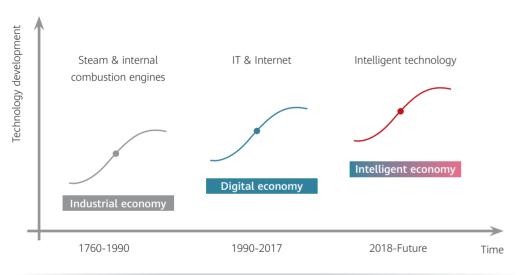


Figure 1-3 The intelligent economy will become a new engine for global economic growth

The intelligent economy driven by AI will bring significant changes in human-machine interactions, IT infrastructure, and new business forms. First, the intelligent economy optimizes human-machine interactions to make communication more natural and smoother. Second, it will reshape IT infrastructure and build a more efficient and intelligent system for information processing and transmission. Third, the intelligent economy will create a series of new business forms and stimulate cross-domain innovation. These three aspects do not exist in isolation

but rather evolve in a coordinated manner, synergizing to generate a compound effect.

Over the past four decades, informatization and digitalization have brought strategic opportunities worth trillions of dollars to the ICT industry. The full potential of the intelligent era is yet to unlock. Huawei predicts that the global intelligent economy will be worth more than US\$18.8 trillion by 2030, and that this will open the door to further strategic development in the ICT industry.

AI for All

The rapid development of AI and the emergence of foundation models signal that AI will reshape every organization and every aspect of people's lives. Experts and institutions predict that AI will impact the world profoundly. So, how is AI currently being perceived and applied by enterprises and individuals?

According to a McKinsey report in 2023, 55% of organizations have already adopted AI in at least one of their departments, and this represents a 275% increase compared to 2017. According to Gartner's Top 10 Strategic Technology Trends of 2024, over 80% of enterprises will use generative AI by 2026. In addition, 75% of enterprise software engineers will use AI coding assistants by 2028. which is an almost eightfold increase since the beginning of 2023.

Every industry will be reshaped by AI

As AI sets off a series of industry transformations, practically all industries will be reshaped. For example:

- Autonomous driving and electric vehicles will dramatically disrupt the automotive industry.
- Intelligent transportation will tremendously boost traffic efficiency.
- Personalized education will significantly improve teaching quality.
- Early prevention and precision treatment have the potential to increase life expectancy.
- With real-time translation between multiple languages, communication will be easier than ever
- Precision drug trials will cut the time and funding required to discover and bring new drugs to
- AI-enabled O&M for telecom networks will become more efficient.

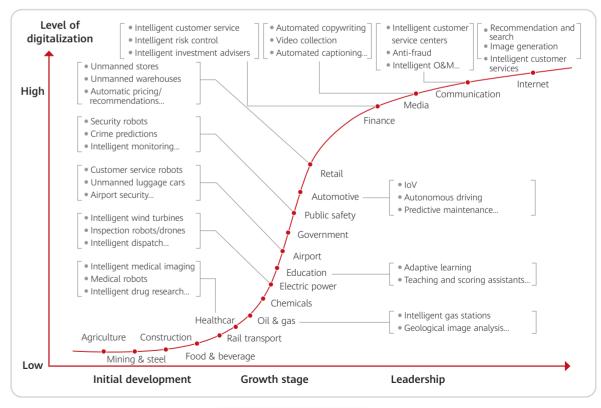


Figure 1-4 AI is reshaping industries

Al is reshaping industries far faster than expected. At the beginning of 2023. BYD Auto remarked that it would take some time to realize autonomous driving. Within a year of that comment, autonomous driving technologies have already been upgraded iteratively. In 2024, the penetration rate of autonomous driving functions in China's new energy vehicle market exceeded 51%. Several technical achievements are behind this leapfrog growth: advanced sensing systems, powerful computing platforms, and Aldriven decision-making and planning algorithms.

Al can also disrupt industries. India has long been a popular choice for outsourcing IT services because of its advantages in labor cost and English language skills. However, it's now facing serious challenges from AI technologies. Statistics show that as a result of AI, the five major IT service companies in India laid off 69,197 employees in the past year, the highest number of layoffs in 20 years. Al has taken over some of the tasks that used to be done by workers in India. Therefore, while AI can revolutionize and propel some industries forward, it can also disrupt certain industries that fail to keep pace. In the future, we have reason to believe that AI will be able to positively reshape every industry.

Every application and software is worth rewriting with AI

Generative AI represents a revolutionary leap. Some people call it AI 2.0, rather than an upgrade of AI 1.0. AI 2.0 can use massive amounts of data that do not require manual labeling to train a foundation model with crossdomain knowledge. It can be developed from scratch to truly generate intelligence. AI 2.0 enables everyone to be a creator, or even become a programmer. Al 2.0 also gives rise to products that have long been seen as

just imagination, such as digital avatars. The powerful generation capability of AI 2.0 can also reduce the costs of innovation implementation to almost zero and create more profitable business models.

The creative and commercial potential of AI 2.0 makes rewriting every application and software in the intelligent era worthwhile.

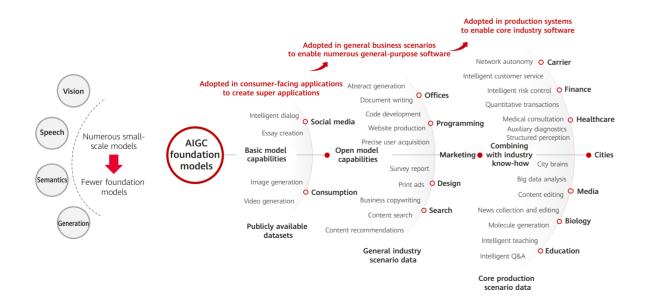


Figure 1-5 Every application and software is worth rewriting with Al

The move towards AGI is driven by idealism and realism

In 2015, the artificial general intelligence (AGI) experiment led by OpenAI became the starting point for AGI. In 2020, GPT-3 was launched, and scaling laws were established as the first principle of AGI, implying that we are making fast progress towards the goal of AGI. To support the continuous evolution of AI capabilities, Project Stargate was launched, with an initial investment of over US\$100 billion. This project aims to build a more powerful computing infrastructure. It is estimated that a data center with millions of xPUs will be launched in 2028.

Idealists believe that further AI development will bridge the technological chasm. They are committed to improving the computing capability of deep learning by a factor of at least 1 million over the next 10 years. Breakthroughs are being made and new papers and models are emerging in the AI field. From pre-training to supervised finetuning (SFT), data sources are being extended from public networks to synthetic data. The momentum that the rapid development of AI technologies has created as we move towards AGI has been felt in every industry.

However, there are still many challenges that AI needs to overcome in consumer-facing applications and businessfacing industries. Many AI applications and projects remain in the initial stages or even disappear after a short period of time. Therefore, the commercial profitability of AI applications has become a focus for the industry. The strategy makers in the AI industry need to think carefully about whether to push scaling laws to the extreme and pursue the infinite limits of AGI, or whether to commercialize and deploy AGI and prioritize generating

Most emerging technologies follow the same development pattern: They start from an idealistic vision but are restricted by natural limits. If we can balance and combine idealism and realism, the successful implementation of technologies will undoubtedly be accelerated. We believe that AI is an irreversible trend. In the vertical direction, the AI industry requires both idealists like scientists and realists like business leaders to strike a balance between technological possibilities and business implementation.

Engineers and scientists are idealists. They are inspired by the spirit of exploration and innovative thinking and

are committed to developing algorithms with higher levels of intelligence and autonomy. They tirelessly aim to maximize computing efficiency and minimize energy consumption. These efforts push the boundaries of what Al technologies are capable of and provide rich theoretical support and technical reserves for real-world applications. Rational participants in industry, on the other hand, are realists. They see AI technology as a key force for business transformation and social progress. They focus on the practicability and economic benefits of technologies, and aim to commercialize AI. They work to integrate AI into industries such as financial services, healthcare, and retail logistics. They hope to verify the market value of AI technologies, find new application scenarios, feedback data for sustainable development, and inspire new directions of research and innovation.

The evolution of AI technologies traces the fine line between idealism and realism and represents how they come together to shape the future of Al. Each technological leap will drive the innovation and expansion of business applications. Business success in turn will leverage more research funds and resources to support the scientific research and promote the maturing and advancement of technologies. Once this positive cycle is established, it can help enterprises create a new closedloop value chain when adopting new technologies. Success stories will accelerate the penetration of AI technologies in core production phases of diverse industries, promote the development of a series of efficient and intelligent solutions, and create considerable business value and social welfare.

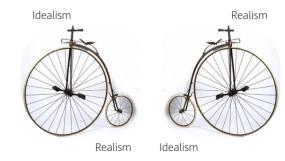


Figure 1-6 Idealism and realism act upon each other to drive AI development

Chapter 2

An "All in Al" Generative **Service System**



Certainty and uncertainty in enterprises' Al development

According to a survey conducted by McKinsey, more than 70% of enterprise leaders feel that AI will profoundly change their business landscape over the next five years. In addition, there are huge uncertainties in enterprises' Al development. Statistics from Deloitte show that 90% of large enterprises plan to invest in AI, but only 10% of them will be able to successfully deploy AI on a large

This is because while generative AI is a revolutionary technology, it has inherent limitations.

Case No. 1: ChatGPT can easily solve International Mathematical Olympiad questions, but struggles to compare the values of 13.11 and 13.8.

Case No. 2: While autonomous driving technology is disrupting the automobile industry and transforming mobility services for the public, the medical model that is designed to improve auxiliary image diagnosis is still in the innovation and research phase.

Case No. 3: 50 artists were able to create the first science fiction movie using AI. Meanwhile, many enterprises are debating whether investing heavily in AI is worth it when the benefit is just some writing assistants. There is no guarantee for when models will be able to answer questions consistently well.

For enterprises, the question is whether to swim out to chase the tide and lead industry innovation, or just wait till the tide ebbs and scoop up the fish that are washed up on the beach.

Large language models (LLMs) such as ChatGPT represent a revolutionary change in the world of AI thanks to their generalization capabilities, which were only made possible by aggregating the world's

knowledge. They've essentially condensed explicit knowledge and accumulated implicit experience to create a probability model that discovers internal rules based on structured data. Players in many industries, especially top enterprises, have massive amounts of data, service knowledge, and experience that is built into their service processes, making these resources extremely valuable. When these resources are used to train AI models, the models are able to naturally memorize the knowledge and experience. When enterprises begin using foundation models and industry-specific models, or start building their own private scenario-specific models, they are able to unlock a new source of energy to drive their business forward. They can develop models that best suit their needs while also leveraging the strengths of the larger models in a more cost-effective way. This enables them to efficiently pass on and use their internal experience, industry experience, as well as the world's collective public knowledge to achieve sustainable development.

The biggest type of waste that occurs in enterprises is the waste of experience and talent. And so, Huawei has evolved its enterprise AI from 1.0 to generative AI 2.0. Compared with the previous evolution, this AI is being applied to more of our own core business domains: from contract risk audit to pandemic-proof supply chain resilience management, from global network optimization to Internet information product experience management, and from intelligent customer service to high-volume, high-trust code generation. The goal of Huawei AI 2.0 is that a single top expert, with the support of a team of non-expert humans and AI-enhanced digital employees, will be as efficient as or even more efficient than a full team of top experts.

The biggest waste in enterprise: Wasted experience and talent

- 1 Massive amounts of data has yet to be mined
- 2 Inefficient use or loss of accumulated knowledge and experience



Sustainable enterprise growth: Unlock a new source of energy to drive their business forward

- 1 The ability to absorb social knowledge determines an enterprise's ability to achieve sustainable
- 2 An open mindset: Unlock a new source of energy to drive their business forward



Figure 2-1 The certainty in enterprises' AI development

The groundbreaking opportunities being created by AI foundation models come from the changes being made to our world's basic scientific paradigms and the unraveling of unknown rules from massive amounts of data. More and more enterprises want Al to become a core competitiveness of their products and services for core production scenarios, so that they can become pioneers with leading capabilities. For example, specialized steel manufacturers must observe strict production tolerances. One of the key parameters that must be controlled in steelmaking is mold fluctuation. But mold fluctuation levels are linked to more than 200 other parameters, such as the mold height, water volume, temperature, pressure, and batch sequence of raw materials. Expert experience and human calculations often fall short when trying to control all of these factors. Steel companies are now exploring how AI can optimize these production and manufacturing processes, and how to train scenario-specific models using their years of valuable historical data. They are also looking into how to gather continuous feedback to their models to further enhance real-time production



To develop AI, enterprises need to build enterpriselevel integrated intelligent twins. When a craftsman approaches a complex problem, they combine the explicit knowledge they've learned from study with the implicit experience they've gathered from practice. This creates a closed loop from perception and understanding, to predictive analysis and decision-making. We are pleased to see that AI foundation models are structuring massive, multi-source, and unstructured data throughout the entire process of sensing, prediction, and decision-making. When Al's vision extends from language and text prediction to scenarios that are more related to the real world. such as sounds, object images, time sequence sampling, molecular structure, and network load scheduling, infinite opportunities will emerge for enterprise AI.

The certainties of AI development can be addressed in enterprise business strategies, and the uncertainties of AI can be tackled during operational planning (we call them tactics). The best solution so far to develop AI, is to adopt an "All in AI" strategy, that is flexibly rolled out at a reasonable pace using an ecosystem-based model.



Architecture comes first: Turning uncertain challenges into definite opportunities

The core challenges to building an enterprise-level All-in-AI architecture can be boiled down to two simple geometric shapes: the architecture's unstable dumbbell-shaped structure and the impossible trilemma of industry-specific foundation models.

Challenge 1: The unstable dumbbell-shaped structure.

Traditional enterprise IT architectures are a stable triangle. Their infrastructure and technical platform are stable and change infrequently. Their data and application enablement platforms are also iterative, updating based on product and version releases, which means their changes are predictable. Applications update agilely and frequently based on user experience requirements. However, to support AI foundation models, IT architecture also includes an additional model layer. Models are developed and iterate rapidly, with their upgrades more frequent and more dramatic than those of traditional applications. This raises the question: How do we plan and design IT architectures capable of meeting this so that core parts of the architecture can be modularly replaced, like changing the engine of an airplane on the fly?

Challenge 2: The impossible trilemma of industryspecific models. Achieving generalization, professionalism, and cost-effectiveness at the same time during AI model development is very difficult. Generalization focuses on scenario-based learning capabilities with small samples; professionalism focuses on strong supervised learning capabilities; and cost-effectiveness focuses on moderate model scales. In addition, the focus points of foundation models vary depending on their types and scenarios. For example, language models have large numbers of parameters, and need high computing power and costeffectiveness. Models for product quality inspection only have access to a limited number of videos showing negative samples, and so demand high generalization capabilities. Risk identification models require high precision and professionalism. Due to the scarcity of industry data, the contradiction between generalization and professionalism in industry-specific models is particularly apparent.

A core concept in enterprise AI development is that the uncertainty of models can be addressed by leveraging the certainty of architectures. An unconventional vet stable architecture with a model layer is needed to support continuous development. If the application layer can use an iterative "All in AI" as the blueprint for longterm planning, and if the infrastructure and AI technology platform remain stable, then the model layer at the "shock center" can be decoupled from the application layer and base layer.

- Multi-source models: A computing backbone encapsulates software and hardware complexity, elastic resource scheduling improves computing efficiency, and service-oriented standard interfaces connect to the open model layer to support models from various sources.
- Triple evolution: Model capabilities are encapsulated using APIs, and applications are decoupled from models to form a replaceable "engine". This way, L0 foundation models evolve with industries, L1 industry-specific models evolve with the industry model market, industry ecosystem, or group center cloud, and L2 scenario-specific models can be fine-tuned and evolve within an enterprise.
- Application orchestration: Services range from edge and support applications to core production applications. Interaction understanding (NLP), perception (CV), simulation and prediction, decision-making optimization models, and search capabilities can be combined on demand, and APIs can be embedded into business processes in a lightweight manner or as assistants.

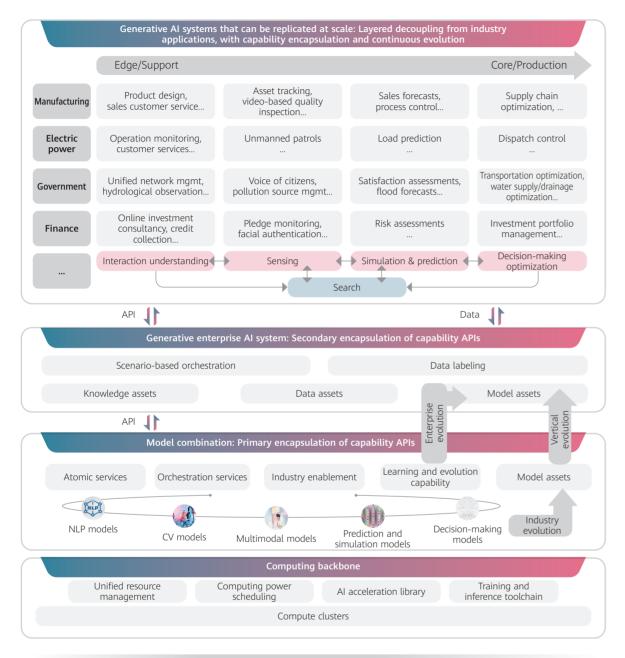


Figure 2-2 Open architecture to support a vast range of models and applications in evolution

A controllable open ecosystem can address the impossible trilemma of industry-specific models and help create an industry-specific model layer that supports on-demand combination. On the one hand, enterprises need to embrace standards and the industry ecosystem, and ensure on-demand integration and utilization of the industry ecosystem. On the other hand, enterprises can establish a pyramid structure for AI applications which allows applications to be managed by

their category, such as super applications, top applications, rigid-demand applications, and common applications. Enterprises can then flexibly pursue the AI development method most suitable for them, based on their own competitiveness strategy and capabilities. These different development methods include in-house development, joint development with strategic partners, and sourcing from ecosystem partners.

Application scenario centricity: Adopting a four-inone framework to realize the value triangle

In the early stages of AI development, enterprises tend to be model-centric. From a technological perspective, they replicate applications that are easy to implement based on the capabilities of their industry's foundation models. This can lead to siloed application, model, and computing infrastructure development.

(R&D, marketing, service, manufacturing, supply, finance, HR)

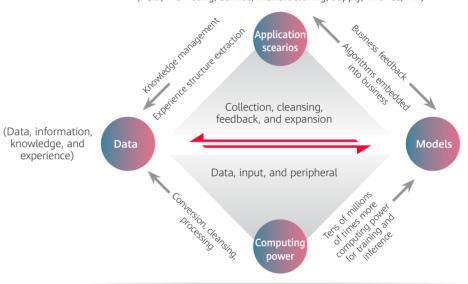


Figure 2-3 A four-in-one framework for enterprise AI development

The essence of application scenario centricity is to solve problems with "First Principle Thinking". Under this model, scenarios are both the start and the end of the development process, creating a closed value loop. Under this model, the focus of development should not be the foundation model itself or the number of model parameters. Instead, development focuses on: Whether the model can solve previously unsolvable or difficult-to-solve problems; whether the model can achieve a positive loop where the benefits are greater than the costs; and whether the model is widely applicable and replicable. In industry AI applications, Al models and mechanism models are often combined to improve their explanatory and deterministic performance. For example, in mine exploration, an AI model can be used to optimize the selection of drilling positions, while a mechanism model ensures the physical feasibility and safety of the mining solution.

The above four-in-one framework makes application scenarios, data, models, and computing power all indispensable in the process of implementing a closed value loop for application scenarios. Scenarios are the basis of this value loop. Focusing on scenarios that are of low business value and consuming a large amount of computing

power is like setting up a top-level expert team for a minor business domain. If models and data are not well matched. for example, if the model generalization is poor and there are insufficient samples, then the model will not be sufficiently professional and accurate. Using the previous metaphor, it is not possible to properly complete complex work by simply hiring a large number of workers, if all the workers lack relevant experience.

The four-in-one framework is divided into the technology triangle and business triangle to decouple technologies from business and facilitate the establishment of a platform-based technical architecture. The technology triangle relies on compute resources to implement data conversion, cleansing, and processing, accelerating model training and inference. The business triangle manages knowledge, extracts experience structures, continuously enriches enterprise datasets, and implements bidirectional interaction between data and models, all within the context of an application scenario, to provide business support and result feedback. The most typical example of "anything that is not normal is abnormal" illustrates the feedback and supplemental role of datasets in the use of models.

AI is gradually being applied to core and production scenarios along the enterprise value stream

Companies must systematically organize applications and establish a point-line-plane application map before developing AI. The AI value triangle shown below can be used as an experience-based paradigm that guides the determination of AI's potential value in different applications. Companies often use AI assistants to improve service efficiency and user experience. For example, AI assistants can be used to assist with day-today office work, HR tasks, and customer service. The use of AI in production scenarios, such as online consulting, process optimization, and demand and supply forecasting, can usually improve a company's productivity and competitiveness. In addition, AI can be applied to business continuity risk control, financial risk identification, and more to prevent black swan events.

Companies must apply AI step by step. When developing

an application map, companies should consider the big picture and their long-term plans. Model capabilities and ready-to-use data should not constrain the development of the map. Instead, companies should design and plan the map with their business development strategies, core principles related to AI technologies, and industry development trends in mind. By contrast, when formulating the implementation roadmap, companies should start with specific scenarios based on their short-term needs and avoid making big changes to service systems and processes. The roadmap should break down capabilities based on specific scenarios and combine model-based capabilities such as perception, understanding, prediction, and decision making. Breaking down tasks makes it easier to solve problems and make full use of the capabilities of the entire ecosystem so that companies can benefit from scenario-based flywheel effect.

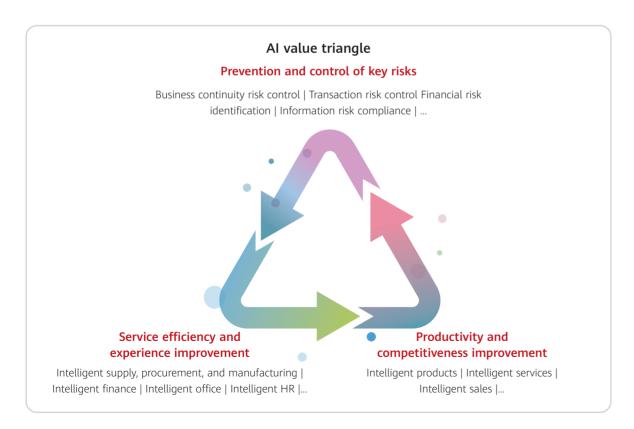


Figure 2-4 AI value triangle for application selection

Business, technology, and data readiness are determining factors when selecting the application scenarios for Al. Business readiness measures whether

the scenario-related business processes are clear, whether business rules are fixed, whether business organizations are willing and determined to invest, and whether experts who are familiar with the business rules are available. Technical readiness measures whether algorithms, models, equipment, services, and computing power that may be involved in the scenario are ready and whether they can create the expected value. Data readiness measures whether there is sufficient data of a high-enough quality and whether the data distribution and labeling required for the use of AI in a specific scenario are in place and ready.

The general principle of scenario selection is as follows: Companies should preferentially select scenarios that have a strong and frequent demand for Al support and where AI can be more easily applied in order to quickly identify the value of AI and cultivate AI talent to enable a continuous and virtuous cycle of AI development. Industryleading companies should select scenarios for which they have sufficient data and should focus on high-value "super scenarios", such as blast furnaces in the steel smelting industry and pilot production in the chemical industry. They should work with industry research institutes as well as AI and foundation model companies to make breakthroughs. Once a breakthrough is made, huge industry value will be unleashed.

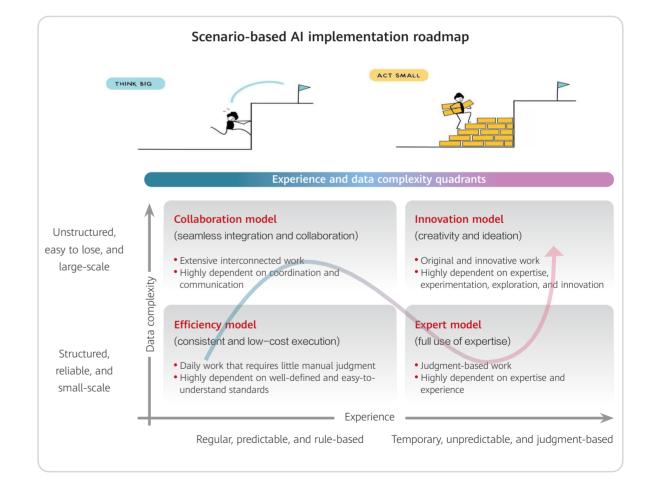
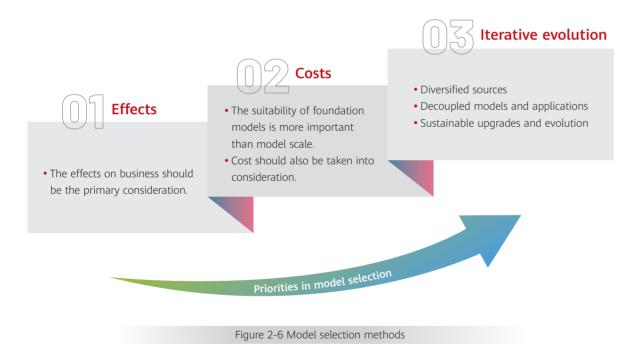


Figure 2-5 AI implementation roadmap

Models Considering alternatives to LLMs for AI application in industries

Large language models demonstrate powerful content generation capabilities, such as human-machine conversations and O&A. They can generate structured data based on unstructured information, and this capability enables image interpretation and emotion recognition. Also, they can generate unstructured text in accordance with strict pre-set rules, which is a useful function for engineering and code design. Large language models are being applied in many different knowledge-intensive digital assistant scenarios, such as customer service, code development, professional consultation, public opinion analysis, and auxiliary design.

However, real-world industry issues cannot be solved with large language models alone. Waterlogging warnings, electric vehicle charging fire risk predictions, water supply loss control, industrial process formulas and process parameter control, and financial credit assessments are all complex issues. Solving these issues requires combining mechanism analysis with AI models, integrating multiple models for perception, understanding, prediction, and decision making, and considering the real-time requirements and model costs.



As industries gradually come to understand these requirements, they are likely to prioritize suitable models over large and unified models. They will no longer blindly pursue model scale and a large number of parameters. Instead, they will make a trade-off based on their actual needs. Large and small models have their own advantages. Diverse and complex scenarios require the flexible use of different models, and this approach is becoming a trend that can be seen in many industries. The suitability and applicability of models are becoming more important than the number of parameters



Data remains a key company asset and AIGC transforms the data governance structure to maximize the value of data

In the long run, as foundation models converge and computing power becomes less scarce, the quantity and quality of personalized data will determine the specific capabilities of enterprise AI models. Data will remain a key company asset, but the characteristics of generative AI will transform its governance structure.

In the intelligent era, data will cement its position as a strategic corporate asset. Massive amounts of multidimensional raw data generated during processes, as well as the data generated from the service judgments made and the tasks executed by top industry experts in operations, are becoming the most valuable assets. The storage of massive amounts of historical and process data is no longer a pure cost. Instead, it is an asset for AI that is becoming ever more valuable.

Generative AI is also systematically transforming the data security governance structure. Data memorized and generated by models and models themselves are setting new data boundaries both inside and outside companies. Foundation models can accumulate data and knowledge within their parameters and generate data in the form of text, videos, and policies. As a result, there is no longer a clear difference between applications and data, and it is more difficult for companies to manage data boundaries. Integrating industry data and enterprise data is a top priority for companies. Traceable and manageable access control is implemented between domains along the dependency chain of raw data, training datasets, AIGC models, and model services based on the original security level of the data

Data directly affects the performance of models, and the use of data will facilitate its own flywheellike development. Data does not necessarily have to be large-scale but it has to be complete. It should be used in specific scenarios where it is easy to apply and then gradually extended to scenarios where it is more difficult to apply. The use of data will in turn facilitate its development. Companies can raise their data requirements based on a model's performance in a specific scenario, and obtain more data as needed to optimize model performance, thus forming a data flywheel.

Data governance ensures the quality of data. The best governance is source-based governance based on data collection. Smart cities, mines, oil fields, and factories and many other industry scenarios involve a large number of terminals, sensors, and other equipment. Unified standards and specifications for smart terminals and data collection can greatly reduce data governance costs, especially in scenarios involving multiple entities. Companies can obtain high-quality labeled data at a low cost through collaboration between inference at the edge and training at the center, irregularities in video sensing that are automatically labeled, or data labeling that is integrated into the operations of business personnel.

In the Data as a Service business model, maximizing the value of data is the overarching goal of data governance. Al is applied along the entire data value chain, from data reproduction and data labeling, to rule discovery. First, models are used to process and generate massive amounts of heterogeneous data. This aims to turn a lot of heterogeneous data, such as drawings, video surveillance footage, and public opinions on the Internet, into structured information, so as to provide a solid foundation for data analytics and risk assessment. Second, models help to implement trusted and accurate cross-department data sharing. By sharing high-level data, such as the security status of people or objects in videos, data can become available but invisible, and this ensures that the value of data is fully unleashed while also protecting privacy and ensuring data security. Finally, models can make predictions and decisions based on data from the entire domain. Each business unit can make more accurate predictions and discover more and increasingly complex rules based on the data of itself and associated entities.

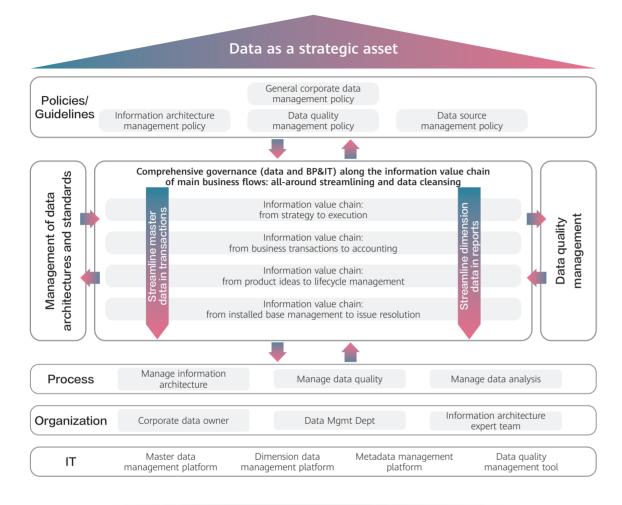


Figure 2-7 Data as a strategic asset



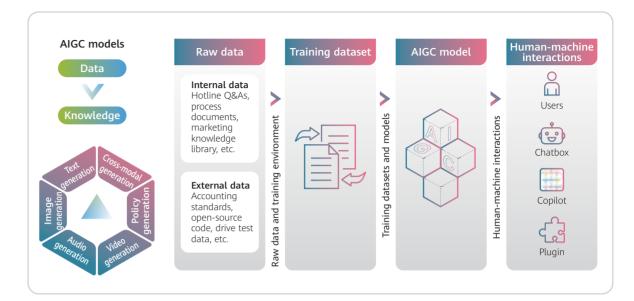


Figure 2-8 AIGC-driven data governance structural transformation

The processing and generation of massive amounts of heterogeneous data

Al can convert various heterogeneous data, such as drawings, video surveillance footage, and public opinions on the Internet, into structured information to provide a solid foundation for data analytics and risk assessment.

Trusted and accurate cross-department data sharing

High-level data (such as the security status of people or objects in videos) can be shared to ensure that data is available and invisible, and this ensures that data value is fully unleashed while also protecting privacy and ensuring data security.

Prediction and decision-making based on data from the entire domain

Each entity can make more accurate predictions and discover more and increasingly complex rules by using its own data and data from associated entities.

Figure 2-9 Maximized data value

Computing power

Making computing power a foundation and platform in collaboration with strategic like-minded partners

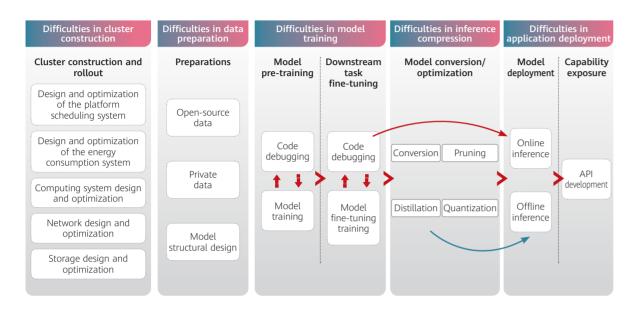
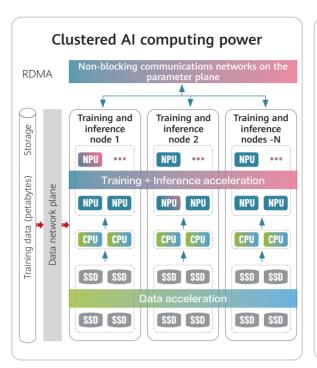


Figure 2-10 The engineering challenges faced during foundation model development

The development and application of a foundation model is a complex, systematic project. It requires a highly-integrated AI computing platform, with highly coupled internal software and hardware and standard external interfaces. The major challenges that need to be addressed include those stemming from cluster construction, model training, inference compression, and application implementation, and some examples are listed below.

- Cluster construction: Ensuring the high performance and long-term operational stability of ultra-large clusters; and building a lossless network on the parameter plane
- Model training: Selecting the most efficient parallel combination policy; implementing multi-task visualized optimization; achieving resumable training after a breakpoint; and predicting the scalability and performance of a foundation model
- Inference compression: Implementing distributed inference and inference acceleration; and performing lossless quantization of a foundation model
- Application implementation: Building a large-scale inference cluster scheduling system; designing attack defenses; and effectively rectifying and isolating faults



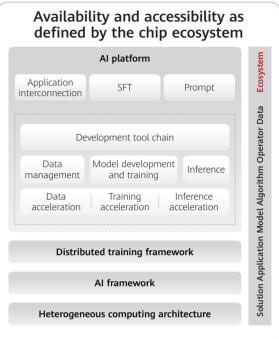


Figure 2-11 The highly coupled components in the computing backbone and their dependency on a thriving ecosystem

The selection of computing power is also the selection of technical tracks. A sustainable AI computing power supply chain is a prerequisite for technical track selection. Companies must consider both the long-term availability of hardware, and the sustainability of software stacks. During foundation model training and inference, the number of model parameters ranges from billions to trillions. The computing platform needs to have powerful parallel computing capabilities, and operators (the software modules that execute basic computing tasks) need to perform efficiently to unlock the full potential of hardware computing, memory access, and inter-card

communications. For example, Huawei's NPUs are specially designed for the matrix-based computing framework that features Al. They are more suitable for accelerating models such as convolutional neural networks. It is worth noting that Al computing chips need support not only from hardware, but also from the corresponding developer ecosystem, including the development tool chain, software library, framework, and developer community. Finally, the selection of computing power must take into account factors such as scheduling efficiency and development efficiency in addition to the training and inference requirements.



Data centers as the foundation

Networks play a key role in the information era because they connect corporate IT systems and many other things. Cloud plays an important role in the digital era by enabling agile application development. As we approach the intelligent era, computing power is poised to take center stage. As data centers provide computing power, their efficiency is a key determinant of enterprise Al efficiency. Data centers are no longer pure cost centers, but innovation centers.

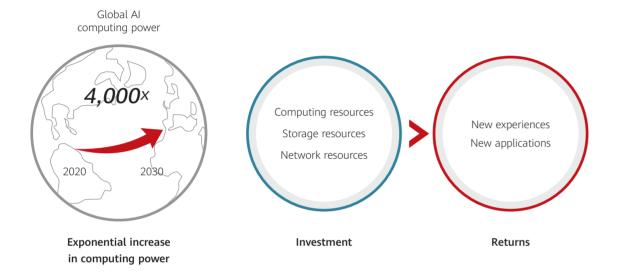


Figure 2-12 From cost centers to innovation centers

The scale, computing power efficiency, and development efficiency of data centers have become the core competitive edges of enterprise AI. As the parameter scale and data scale of models increase, the scale of data centers, effective computing power of clusters, and energy saving capabilities become key factors in enterprise model development and AI application implementation in the face of limited computing power supply and investment constraints. Enterprise AI's value is not showcased in low-frequency inference by some killer

models, but in the frequent use of hundreds of scenariospecific models in massive, repeated, and complex scenarios. Data center efficiency is obviously crucial when a common interaction requires tens of billions of operations. The training and inference of foundation models have become the most complex IT projects. Data centers are becoming the core of companies' intelligent infrastructure and an important factor in the "ROI inequality" of enterprise AI.

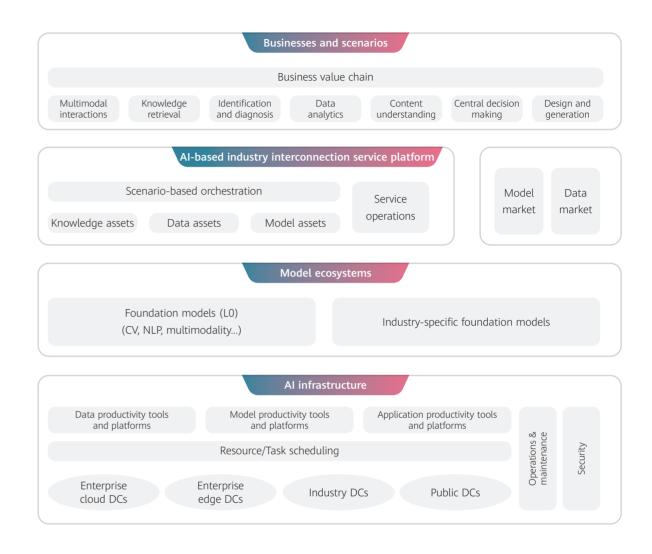


Figure 2-13 Enterprise-level AI architecture

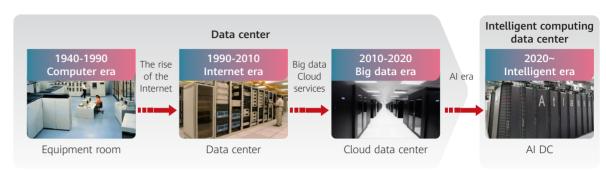
Data centers will be redefined by AI. We envision AI providing diversified computing power and enabling a vast range of models and innovate AI-native applications. Computing power will be no longer limited by equipment room infrastructure. The cluster scale will be no longer limited by communication networks. Tasks will be able to be flexibly scheduled, and computing power resources will be able to be shared between data centers. In this way, computing power will be able to keep pace

with the expansion of foundation models. Companies should support an open model ecosystem and provide flexible model selection and combination services for different service scenarios to ensure that each task can be matched with the most suitable algorithm and model combination. The agent-based task design mode can be integrated with knowledge, data, and model assets from companies and industries to implement scenario-based orchestration.

Chapter 3

Evolution and Changes of Data Centers in the **Intelligent Era**

Evolving towards AI DCs



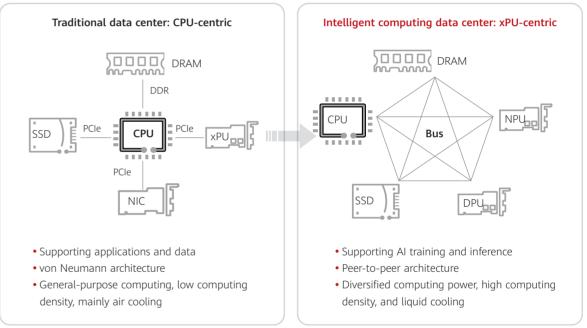


Figure 3-1 Evolving towards AI DCs

Over the past few decades, data centers have evolved towards AI DCs.

With the rise of the Internet, data centers have become the core carrier of IT infrastructure. Starting from 2010, with the rapid development of big data and cloud services, the architecture of data centers has undergone significant changes. The advent of cloud computing has made data centers more flexible and efficient and allows them to provide compute resources and services on demand.

By 2020, the rapid development of AI led to a surge in demand for computing power, marking a significant leap in the era of intelligence. Against this background, AI DCs have emerged, providing high-performance computing capabilities required for AI model training and inference. Examples of AI DCs include Google's data centers for machine learning, Meta's AI Research SuperCluster, and the Peng Cheng Cloud Brain II supercomputing platform designed for deep learning.

Traditional and AI DCs have the following differences:

Services involved

- Traditional data center: primarily supports enterprise applications and data storage, including routine information processing tasks like web services, database management, and file
- Al DC: primarily supports Al model training and inference, efficiently providing computing resources for processing large datasets.

Technical architecture

- Traditional data center: adopts the von Neumann primary-secondary architecture, where the CPU acts as the commander, assigning tasks to other components. However, for largescale parallel computing tasks, the primarysecondary architecture encounters limitations in computation, memory, and I/O, creating bottlenecks that hinder further performance improvements in the data center.
- AI DC: utilizes a fully interconnected peer-topeer architecture, enabling direct communication between processors, memory, and network adapters. This reduces latency from centralized control, overcomes the computational limitations of the primary-secondary architecture, and enables efficient distributed parallel computing.

Computing power

- Traditional data center: CPU-centric, suitable for general-purpose computing needs.
- Al DC: xPU-centric, processes parallel computing and matrix operations required for AI model

Cooling modes

- Traditional data center: typically has a power density of 3 to 8 kilowatts per cabinet, with limited server equipment and lower computing power density, and generally relies on traditional
- Al DC: typically features a power density of 20 to 100 kilowatts per cabinet, primarily utilizing liquid cooling or hybrid air-liquid cooling technology. Liquid cooling efficiently dissipates heat, ensuring the stable running of high-performance computing devices.



Carrying the training, inference, and application of AI models

The AI DC is designed and implemented for AI model training, inference, and application.

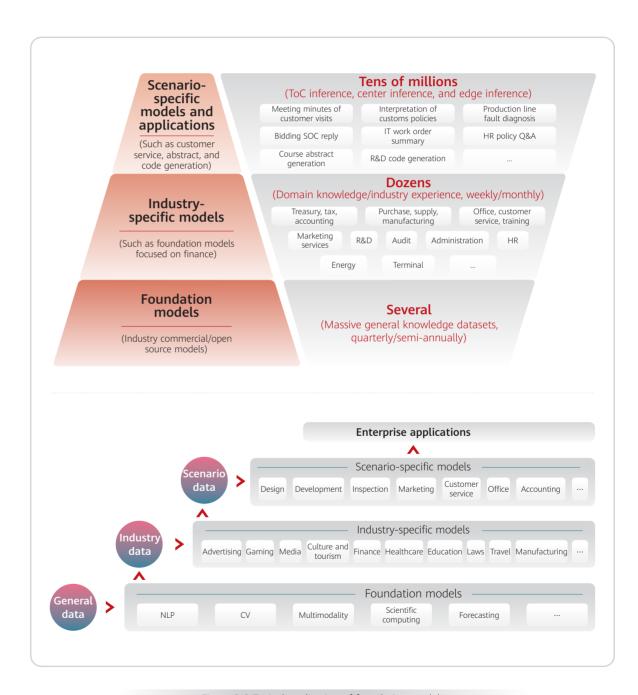


Figure 3-2 Typical application of foundation models

AI models are classified into foundation models. industry-specific models, and scenario-specific models. Foundation models are versatile and can perform well across a variety of tasks. Industry-specific models are tailored to specific industries, with a deep understanding of professional terminology and business processes within that domain. Scenario-specific models are customized for particular service scenarios or problems, ensuring they meet the precise requirements of specific

tasks and feature enhanced professionalism and service

The full application of AI models involves close collaboration from training to inference. This process includes pre-training for foundation models, secondary training for industry- or enterprise-specific models, and fine-tuning for scenario-specific models. Each step presents new challenges to the technical capabilities and resource management of the data center.

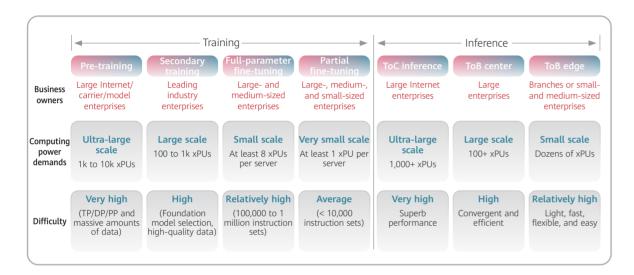


Figure 3-3 Computing power demands and engineering difficulties for training and inference across various scenarios

Large Internet enterprises and model companies focusing on foundation model training have clear AI DC planning and construction objectives to support foundation model pre-training. This is a big project, which requires the support of a computing platform with ultra-large clusters. In addition, trillions of tokens need to be collected and processed to ensure that the models can learn sufficient knowledge and skills. Such large-scale training is not only a technical challenge, but also poses a great challenge to resource allocation and system O&M management capabilities.

Leading industry enterprises prioritize the secondary training of industry-specific models during AI DC **planning.** These industry-specific models are developed through incremental training based on foundation models and extensive industry-specific data. Compared with foundation model training, secondary training of industry-specific models is less complex. However,

hundreds to thousands of NPUs/GPUs are required and hundreds of millions of tokens need to be processed.

For most enterprises, the AI DC primarily focuses on fine-tuning, inference, and deployment of AI models. Given the highly scenario-specific nature of Al applications, enterprises typically need to further fine-tune industry-specific or foundation models using their own scenario-specific data to enable the models to accurately understand and generate data in specific scenarios. This ensures that the models meet the standards required for deploying applications in real-world service environments. Key indicators of AI inference include latency, accuracy, concurrency, and efficiency. The planning and construction standards and technical requirements of AI DCs vary depending on the number of target users of inference services, such as consumer-facing services, business-facing services, and internal applications.

Four AI DC construction scenarios and three Al DC types

To address various needs, enterprises plan to build AI DCs tailored to four typical scenarios and functions.

Scenario 1: Full pre-training

Leading Internet companies, telecom carriers, and foundation model vendors are building ultra-large AI DCs to train foundation models and provide inference services for a large number of consumers.

Scenario 3: Secondary training + edge inference

In some group-based enterprises, large AI DCs are typically established at the headquarters for secondary training and center inference. Additionally, small AI DCs are set up in branches or near production centers for edge inference and fine-tuning. This creates a center-edge collaboration architecture that aligns with the enterprise's overall organizational structure. This architecture not only optimizes resource utilization but also enables real-time decision-making and enhances response speed.

Scenario 2: Secondary training + center inference

Top enterprises in finance, electric power, and other industries that have a significant impact on the national economy and people's livelihoods are actively promoting the construction of large AI DCs for secondary training of industry-specific models and center inference services.

Scenario 4: Lightweight inference

For enterprises in certain fields, even if the scale of AI applications is small, the importance of data security and privacy protection often leads them to build small AI DCs for lightweight inference jobs and model fine-tuning. For example, a top-tier hospital uses AI technologies for medical image analysis, helping doctors diagnose diseases more quickly and accurately while ensuring that patient data remains within the hospital's internal network, thereby enhancing data security.

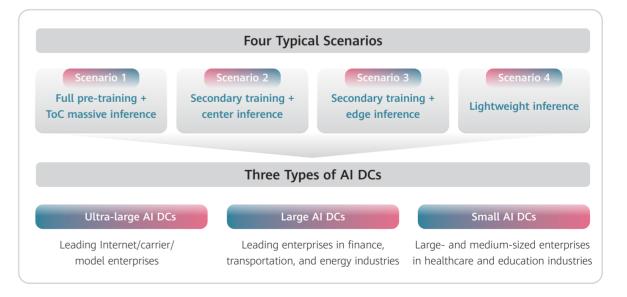


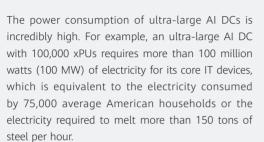
Figure 3-4 Typical scenarios and types of AI DCs

In conclusion, there are three types of typical AI DCs in the industry: ultra-large AI DCs, large AI DCs, and small AI DCs.

Ultra-large AI DCs

Ultra-large AI DCs are mainly responsible for foundation model pre-training and face the following challenges:

Power supply



Improvement of effective computing power

As the scale of AI DC computing, storage, and network devices expands, efficiently integrating these resources to maximize computing power has become a major research focus in the industry. To achieve effective interconnection of a large number of devices, challenges in network architecture, communication protocols, and data transmission efficiency must be addressed. This requires greater attention to scalability, flexibility, and reliability in network design to ensure efficient and stable data transmission and communication between devices. Moreover, simply stacking devices does not result in linear growth of computing power. More intelligent scheduling and management policies are needed to enable close collaboration between computing, storage, and network resources within a cluster. Industry statistics show that the computing power utilization of a leading intelligent computing cluster with 1,000 xPUs does not exceed 60%, with 10,000 xPUs does not exceed 55%, and with 100,000 xPUs does not exceed 40%. This highlights the importance and urgency of improving the effective computing power of ultra-large clusters.

Reliability and fault recovery

An ultra-large cluster consists of tens of millions of components. The training of a foundation model usually requires the cluster to run 24/7 at full load for hundreds of days. As a result, components such as optical modules, NPUs/GPUs, and high bandwidth memory (HBM) are prone to faults. However, the synchronous nature of training makes it less tolerant to faults. Any single point of failure may cause training job interruption, resulting in huge economic losses. Currently, a cluster with more than 10,000 xPUs is only able to run stably for several days. For example, while Meta was training the Llama3 405B model in its 16K cluster, there were 466 job interruptions over 54 days. Fault recovery typically takes hours or even days, severely impacting training efficiency.

To cope with the preceding challenges, the industry-leading ultra-large AI DC must have ultimate energy efficiency and computing efficiency.

Large Al DCs

Large AI DCs are usually planned and constructed by top enterprises in the industry. They are designed to train and finetune multiple models, as well as handle large-scale center inference and AI applications. They face the following challenges:

Optimized inference performance

How can an AI DC deliver optimal inference performance for a specific service scenario within a given latency?

Efficient training and fine-tuning

Enterprises aim to swiftly deploy intelligent applications into production environments to shorten development periods and stay ahead in a highly competitive market. At the same time, they seek to reduce costs and resource usage.

Improved computing power utilization

Building a large AI DC usually requires huge investment. Therefore, enterprises hope that the precious AI computing power resources can be reused as much as possible to avoid idle computing power resources.

Fast AI application innovation with multi-model orchestration

How can an AI DC flexibly combine and orchestrate multiple models to meet the requirements of rapid application innovation?

Simplified AI DC O&M

Enterprises typically handle the O&M of large AI DCs, and enterprise O&M personnel need ways to quickly identify and fix faults.

Generative AI's security guarantee

There are strict AI output requirements for certain scenarios in industries such as finance, government, and electric power. Generative AI outputs need to be accurate and comply with related laws and standards.

Intelligent computing consumes 10 times more rooms in advance.

Equipment room supplies and conditions

power and requires 10 times more cabling than general-purpose computing. In addition, liquid cooling has become a trend. Enterprises need to plan their requirements and prepare equipment

To sum up, desired large AI DCs are converged and efficient, as they can better adapt to future enterprise development.

Small AI DCs

Small AI DCs are intended for lightweight inference and AI service applications and usually built near production sites or users. Some can also provide model fine-tuning capabilities. The major challenges are as follows:

Higher computing power resource utilization

Small AI DCs can only provide limited computing power resources and these resources will be used to deploy as many service applications as possible.



There are few or even no dedicated O&M personnel for small AI DCs, which calls for minimized faults, simple routine O&M, and remote O&M when a fault occurs.

One-stop deployment

Some small AI DCs are far away from urban areas. In this situation, enterprises prefer one-stop deployment which allows the delivery personnel to complete AI DC deployment in one visit.

Security assurance

Small AI DCs are close to production sites and directly connect to sensors, smart cameras, and other sensing equipment that are prone to intrusion. In this case, security assurance of small AI DCs is imperative.

To sum up, desired small AI DCs are lightweight, simplified, flexible, and easy to maintain and use, which also support fast deployment and upgrades.





Five key features of AI DCs

From a technical perspective, major breakthroughs and innovations need to be made in five key technical fields to overcome the challenges faced by AI DCs and build leading AI DCs.

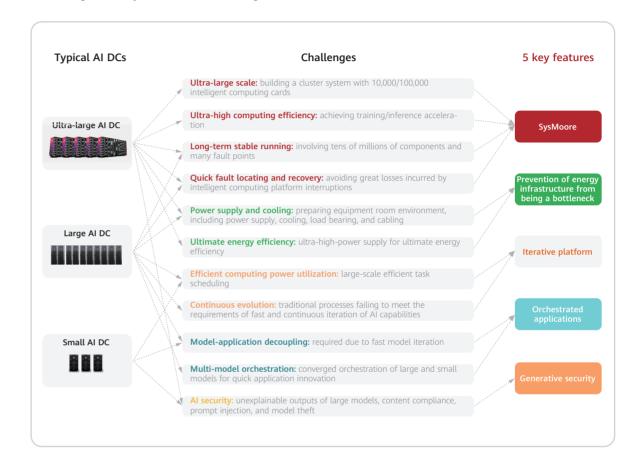


Figure 3-5 Key challenges and features of typical AI DCs

SysMoore

The computing power determines the upper limit of model capabilities. At present, the scaling law is still valid given that the upper limit of large models' capabilities has not been reached. The number of model parameters is expected to range from hundreds of trillions to thousands of trillions by 2028. Such large-scale model training requires breakthroughs in the scale and capabilities of computing power. However, Moore's Law, which dominates the development of computing power in traditional generalpurpose computing, is now reaching its physical and economic limits, and traditional silicon-based electronic technologies are also approaching their developmental limits. The growth rate of computing power is far slower than that of computing power demand, widening the computing power gap. As a result, the industry is in urgent need of a new computing power supply solution, which is where "SysMoore" comes in.

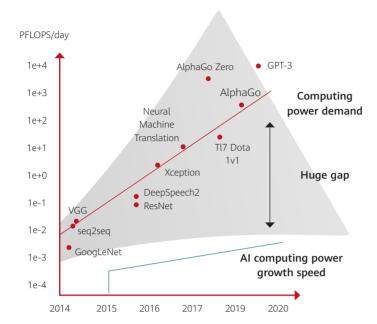


Figure 3-6 Widening computing power gap as intelligent development accelerates

SysMoore is defined as a new method for improving computing power in the Data Center 2030 report by Huawei. It mainly relies on system-level architecture innovation, collaboration between computing, storage, and network, and deep software-hardware synergy to improve computing power, meeting the requirements of exponential growth of computing power.

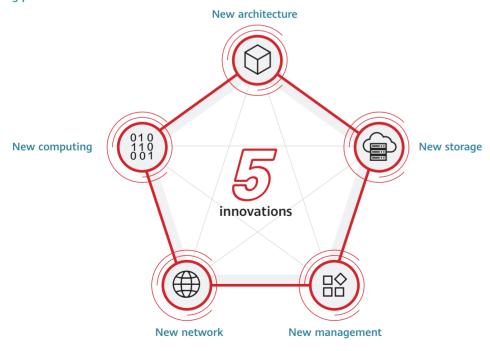


Figure 3-7 Five innovations centered on SysMoore

Al DC computing power supply solutions with SysMoore characteristics present five new features:

New architecture

Over the past 70 years, computers have been designed based on the von Neumann architecture, and a large amount of runtime data is transmitted between the processor and memory. However, this transmission mode involves an ultra-long latency in high-concurrency computing scenarios such as AI scenarios and thereby comes with a "communication wall". In addition, the performance of current memory systems improves at a much slower pace than the performance of processors. In this case, the limited memory bandwidth cannot ensure fast data transmission and thereby comes with a "memory wall". The fully interconnected peer-to-peer computing architecture has been designed to facilitate efficient data exchanges between NPUs, DPUs, CPUs, memories, and other heterogeneous chips, which breaks down the "communications and memory walls" in traditional computing architectures. It meets cross-Al-node, highbandwidth, and low-latency requirements in scenarios related to AI, implements DC as a Computer, facilitates deep collaboration between computing, network, and storage, as well as maximizes the efficiency of computing power through system-level architecture innovation.

New computing

New computing involves the following two aspects:

Evolving types of computing power. There is a shift from CPU-centric general-purpose computing to GPUand NPU-centric intelligent computing. This shift enables the large-scale parallel processing required by AI algorithms, and also greatly improves computing efficiency and flexibility. Parallel computing can process multiple computing tasks or data blocks at the same time, which substantially accelerates data processing and computing while also improving the compute resource utilization and overall computing efficiency. Furthermore, it also enables computing tasks that are more complex to be completed in a shorter amount of time, thereby further enhancing Al development.

Advancing chip technologies. The chiplet technology significantly improves die yield rates, reduces costs, and improves chip performance through flexible adaptation to desired product specifications. In addition, the singlebit energy consumption of solutions in 2.5D packaging is only half that of board-level interconnection solutions in traditional packaging, thereby further improving the energy efficiency ratio of the system.

New storage

In "new storage" fields, the wide application of large models comes with higher demand for high-performance storage. In particular, efficient data read/write has become a key factor for improving the overall efficiency of AI

During the training phase, sample data needs to be quickly loaded from the storage system to the NPU/GPU, and checkpoint data needs to be periodically written from the NPU/GPU back to the storage system. Therefore, improving storage I/O performance and shortening the data read/ write time are crucial ways to improve training efficiency. The NPU/GPU Direct Storage is introduced to provide a direct transmission path for memory access between the NPU/GPU and the storage system, cutting out the original CPU memory buffering and replication processes and greatly shortening the data read/write time.

During the inference phase, and especially in highconcurrency and long-sequence inference scenarios, the industry uses multi-level caching acceleration centered around the Key-Value Cache (KV Cache). This significantly improves the throughput performance of large-scale inference systems and ensures quick and efficient data processing by optimizing data access paths.

In conclusion, the use of both NPU/GPU Direct Storage during the training process and KV Cache multi-level caching acceleration during the inference process helps improve the overall performance and response speed of the system in high-concurrency scenarios with a large amount of data, and meet the requirements for large model applications.

New networks

As a key link between computing and storage systems, networks are rapidly evolving with the interconnection of 100,000 or even hundreds of thousands of xPUs (such as GPUs and NPUs) to meet the connection requirements of large-scale computing clusters. As network technologies develop, the network access rate has increased from 200GE to 400GE or even 800GE in the parameter plane.

The sharding methods of large models are also evolving, from tensor parallelism, data parallelism, and pipeline parallelism to Mixture of Experts (MoE). This poses higher requirements on load balancing at the network level. To overcome these challenges, vendors have launched their own load balancing solutions. For example, Huawei launched the dynamic network-scale load balancing (NSLB) technology for the Ascend platform. Test results show that this technology can improve the training efficiency of the Llama 2 13B model by 13% when no more than 512 intelligent computing cards are used.

To summarize, as network technologies and foundation models evolve, the network architecture and load balancing technologies are also being innovated to meet the computing power requirements for higher performance and a larger scale.

New management

New management modes require E2E system O&M capabilities for cross-domain collaborative management, covering the full-lifecycle O&M management, such as management, control, and the analysis of optical modules as well as computing, storage, and network equipment. This involves the following aspects:

- Full-link, visualized monitoring: By offering real-time monitoring of system performance, the platform ensures comprehensive oversight of computing, storage, and network resources, enabling timely detection of anomalies.
- Fast cross-domain fault locating: Advanced fault detection technologies are used to quickly and accurately locate fault points, reduce troubleshooting time, and avoid training job interruptions.
- Swift cross-domain fault recovery: A robust fault recovery mechanism ensures that, when an issue occurs, rapid response measures can restore system operations and minimize the downtime.

The preceding measures can greatly improve the training efficiency, reduce training costs, and ensure fast, stable, and high-quality large model trainings. Such allround system O&M management capability is the core competitiveness of large-scale or even ultra-large-scale AI DCs in the future.

Prevention of energy infrastructure from being a bottleneck

As the AI DC computing power becomes denser, it consumes much more power and this poses huge challenges for power supply, heat dissipation, and layout design, and necessitates a reshaping of the energy infrastructure of DCs.

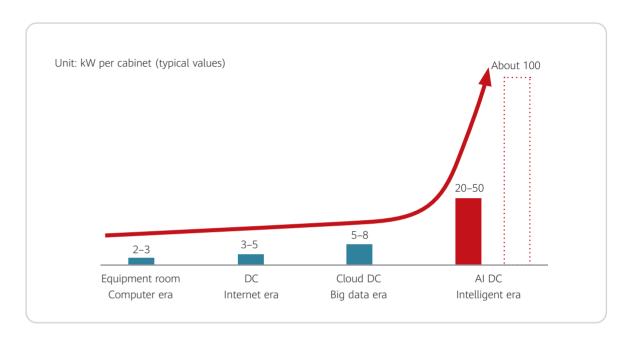


Figure 3-8 Typical powers of DC cabinets in different eras

Challenge 1

Acquiring and matching ultra-large-capacity power supply

As data center power consumption soars, and especially when the power consumption of a single data center reaches 200 MW or even over 500 MW, it is increasingly challenging for the existing urban power infrastructure to meet the demand for power. For example, OpenAI's Stargate project is estimated to use thousands of megawatts of power, requiring cross-region power supply solutions for its data centers. In this situation, the key to improving the computing power scale lies in efficiently and stably obtaining and matching a large number of power resources.



Innovatively dissipating from ultra-highdensity cabinets

High power density also brings about strict requirements on heat dissipation technologies. Although liquid cooling has become a consensus in the industry, it is important to identify new ways of increasing heat dissipation efficiency while ensuring reliability and easy maintenance so that we can overcome challenges related to high power density in the future.



Challenge 3

Providing prospective design of architectural space zoning

The design of AI DCs must consider the complex requirements of IT equipment rooms, cooling facilities, and power supply areas, which require a more forwardlooking layout to break the limits of traditional design. This includes reducing the coupling between IT facilities and electromechanical facilities, implementing modular and outdoor electromechanical facilities, and flexibly combining air cooling and liquid cooling technologies.

To prevent energy infrastructure from becoming the bottleneck of data center development and reduce costs and resource waste, the following measures must be taken:

- Optimizing data center layout: Scientific and reasonable planning and design ensure efficient collaboration between power supply, cooling systems, and computing power requirements to improve overall energy efficiency.
- Improving energy efficiency: Advanced energysaving technologies and management methods are used to reduce energy consumption and achieve the goal of developing green computing power.
- Developing renewable energy and energy storage technologies: Actively utilize renewable energy resources such as solar energy and wind energy, and build energy storage facilities to improve the power supply and risk resistance of data centers.
- Upgrading power supply and cooling devices: Keep pace with technology development and continuously introduce more efficient and reliable power supply and cooling devices to improve the operating efficiency and stability of data centers.

Facing challenges relating to the energy infrastructure of Al DCs, we need to explore and initiate effective responses from an innovative and forward-looking perspective to achieve sustainable development and green transformation while ensuring computing power supply.

Iterative platform

Compared with traditional data centers, AI DCs have a larger scale, more complex services, and faster technology updates. The AI platform faces the following challenges in terms of resource management and scheduling, model training and AI development, and O&M management:



• Efficient utilization of AI computing power resources: The purchase price of AI servers is several times that of traditional general-purpose computing servers. In addition, Al poses more stringent requirements on network and storage devices, resulting in high construction costs for AI DCs. Under such circumstances, how to properly manage and use Al computing power resources to increase the output per unit of computing power becomes a common concern of enterprise users.



 Sophisticated and costly AI development Traditional AI models have poor generalization capabilities, meaning that the performance may deteriorate for different users or data sources. Enterprises without algorithm experts face obstacles when it comes to debugging and optimizing models. Even if the generalization capability of foundation models is improved, the number of algorithm experts is still insufficient. High AI application development costs and long development period hinder AI technologies from serving all fields of enterprise services. Moreover, model maintenance is also a continuous challenge.



• Difficulty in AI DC O&M: As a new type of data center. AI DCs lack O&M personnel who have experience in managing large-scale AI servers and high-performance network and storage devices. They struggle with resource allocation, change management, and quick fault locating and restoration. To deal with these problems, O&M personnel need to enhance their skills and be furnished with comprehensive O&M tools.

To cope with the preceding challenges, an AI platform that supports continuous iteration is required to continuously integrate new technologies and architectures and provide them to users in a mature state. Such an AI platform needs to evolve towards enhanced performance, higher efficiency, streamlined O&M, and more comprehensive functions.



Enhanced performance

An excellent AI platform should continuously integrate cutting-edge technologies to help improve performance while reducing costs. Typically, technologies such as data parallelism and network optimization help improve training efficiency; quantization and compression improve inference efficiency; Prefill-Decode disaggregation enhances the performance of long sequence output; and prompt optimization improves inference accuracy at a low cost.





Streamlined O&M

Large-scale NPUs/GPUs and optical modules complicate AI cluster O&M. The next-generation O&M system is expected to provide functions such as comprehensive monitoring, fault prediction, and intelligent analysis to improve the mean time between failures (MTBF) of hardware and cluster efficiency. In the inference phase, the O&M system needs to monitor key metrics such as hardware usage, identify inefficient jobs, and assist in optimization to continuously improve cluster performance.



Higher efficiency

Improving the utilization of computing power clusters is critical due to the high cost of AI hardware. By optimizing the storage solution and communication algorithms, bottlenecks in parallel training can be overcome, data transmission efficiency can be improved, and training time can be shortened. For interactive inference applications, the platform is projected to support dynamic scheduling, such as API-based, scheduled, or load-based scaling, to release idle resources. Idle resources at night can be used for fine-tuning training. In addition, the platform needs to provide security isolation and flexible scheduling to ensure service continuity and effective resource utilization.



More comprehensive functions

There are multiple foundation model application development modes, such as retrieval-augmented generation (RAG) and agents. The AI platform provides support tools, such as the data engineering module that simplifies data preprocessing, the model development module that simplifies training, and the agent development module that simplifies service building, to improve development efficiency and streamline development efforts.

In a word, the future AI platform should provide enhanced performance, higher efficiency, streamlined O&M, and more comprehensive functions through continuous iteration and upgrades to better support the development of enterprises' Al services.

Orchestrated applications

As digitalization accelerates, many leading enterprises have dozens to hundreds of applications. Over the past year, the rapid development of AI technologies has promoted a concept that all industries, applications, and software are worth re-engineering with Al. In addition, foundation models greatly change software development and drive the advancement of a new development approach called application orchestration. As time goes by, enterprises will possess thousands of models during intelligent transformation. Such a huge model library means that enterprises must resort to application orchestration to fulfill intelligent transformation requirements and promote service innovation.

Application orchestration is fundamentally different from conventional application construction in terms of the entities involved, process decomposition, implementation, and processing. During application orchestration on a foundation model, service engineers and system engineers can instruct the foundation model to decompose and plan service processes according to the specific service logic by using natural language prompts. The process is implemented based on the planning result of the foundation model, and is changed from a static process to a dynamic process. In the future, application construction will rely more on service personnel than professional developers. Application orchestration makes it possible for service personnel and even end users to build agent applications independently.

Current: Conventional applications Future: Orchestrated applications Zero coding, making it possible for service Developers complete coding and testing according to the established process. personnel to build applications independently. Application Mainly developers Mainly service personnel development entity Complex process Service engineers & System engineers System engineers decomposition Manual decomposition, depending on Foundation model decomposition. Complex process coding implementation and extension implementation automated orchestration, zero coding Process handling Fixed process Dynamic process form

Figure 3-9 From conventional applications to orchestrated applications

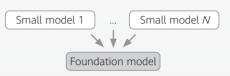
In application orchestration, it is important to make full use of the understanding and generation capabilities of foundation models and the expertise of small models in perception and execution. Through proper orchestration, the two kinds of models can complement each other

and jointly support application functions. Based on the analysis of AI application cases in multiple industries, we have summarized four major application orchestration



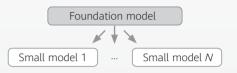


Perception by small models and understanding by the foundation model



A series of small models is used to collect and preliminarily analyze perception data, and then the data is inputted to the foundation model for indepth understanding and generation. For example, in smart city management, video analysis algorithms such as face recognition, vehicle recognition, and abnormal behavior detection can be used to obtain structured data. Such, data can then be inputted to the natural language processing (NLP) foundation model for comprehensive analysis, helping identify potential risks in cities.

Distribution by the foundation model and execution by small models



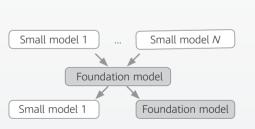
The NLP foundation model understands and distributes jobs, and then small models execute specific jobs. Take the diagnosis of eye diseases as an example. The NLP foundation model can understand a patient's medical records provided by doctors, plan the diagnosis process based on the records, allocate specific eye image analysis jobs to computer vision (CV) models for execution, and generate a diagnosis report. In this way, an efficient closed-loop workflow is formed to improve the diagnosis efficiency of doctors.

Collaboration between foundation and small models



The foundation model and small models work together to complete jobs. First, the foundation model understands the problem and generates a list of jobs. Then, multiple small models and foundation models are invoked to complete these jobs.

Combination of modes A and C



This is a comprehensive application of modes A and C. That is, small models perform perception, the foundation model understands and generates jobs, and small models and the foundation model jointly complete the entire job chain.

Figure 3-10 Four major application orchestration modes

Generative security

In addition to the security risks encountered by traditional data centers. AI DCs also face new security challenges. Firstly, the "black box" nature of AIpowered content production makes the output content uncertain and uninterpretable, bringing severe application risks, especially in scenarios with strict requirements on the output content. Next, Al systems are exposed to new security attacks. The prediction mechanism of foundation models based on statistics and language rules makes it difficult to distinguish between valid instructions and malicious input. Attackers can manipulate foundation models through well-designed prompts. For example, in mid-2023, the hackers utilized a ChatGPT trick known as the 'grandma exploit' to trick AI into performing prohibited operations. Finally, new data security risks may be

introduced. During the training of foundation models, a large amount of user data may be accessed, memorized, and stored. In the inference phase, customers' privacy information may be disclosed unintentionally. For example, employees of Samsung Semiconductor inadvertently disclosed sensitive information such as semiconductor device measurement data and product yield when using ChatGPT. Competitors can then obtain such sensitive information by simply asking ChatGPT, which greatly affects Samsung's market position and competitiveness. To address this issue, the world-leading Open Web Application Security Project (OWASP) community convened 500+ security experts from around the world and proposed the OWASP Top 10 for LLM Applications v1.1 in October 2023.



Figure 3-11 OWASP Top 10 for LLM

To address these security risks, we need to build a comprehensive and diversified security defense system to control risks from the source and ensure the security of foundation models. Firstly, we must ensure the security of training datasets, especially in terms of data copyright protection, privacy compliance, and data traceability. Secondly, in the model training phase, we must enhance

the intrinsic security capability of models and improve the robustness of foundation models by teaching them to learn various security knowledge. Finally, we must build a foundation model security "quardrail" to ensure that foundation models can effectively respond to diverse forms of security attacks and ensure input and output content

Data centers reshaped from layered decoupling to vertical integration

A traditional data center has a layered design, which consists of multiple layers including the energy infrastructure layer, IT hardware infrastructure layer, platform software layer, and application software layer. Components of different services such as compute, storage, network, cloud platform, and databases are purchased separately to construct the data center. This approach is commonly used in the general-purpose computing age, but it faces great challenges when it comes to constructing and planning AI DCs.

Shifted architecture of data centers

The architecture of data centers has undergone a fundamental change. Different from the layered architecture of traditional data centers, AI DCs gradually formed a new layered architecture, an AI DC architecture with the computing backbone layer, platform service layer, model enablement layer, and industry application layer as the core.

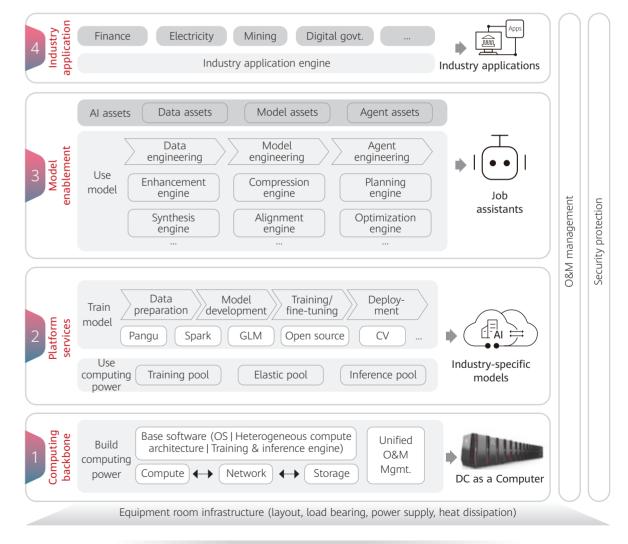


Figure 3-12 New AI DC architecture

Requirements for vertical integration

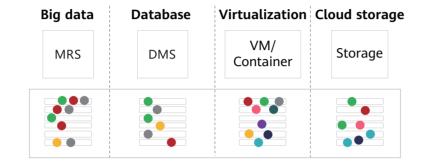
Vertical integration is also required in constructing AI DCs to meet demands of high-parallel intelligent computing. Take the computing backbone layer as an example. Conventional general-purpose computing services such as big data, databases, and virtualization can run different jobs on a loosely coupled system consisting of compute, storage, and network units. These jobs are running on different generalpurpose computing servers, and most of the jobs can be closed within a single server. Nodes are loosely coupled. Therefore, the component-based procurement and construction is still feasible.

However, intelligent computing tasks such as model training are different, especially for foundation model training in which trillions or even tens of trillions of parameters are involved, and the entire computing backbone is required to run a single training job. Running such a job requires the interworks in the whole chain from compute and storage, to network nodes of the computing backbone 24 hours a day for hundreds of days. Nodes must be closely coupled and synchronized. If any node is faulty, the entire job is interrupted and needs to be restarted, causing huge losses. To address this issue, it is urgent to build the computing backbone into a backbone system that works as a precise supercomputer, thereby implementing "DC as a Computer". This can ensure efficient and stable running of intelligent computing services.

In such case, the conventional componentbased procurement and construction can hardly realize "DC as a Computer". Therefore, the vertical integration of components for computing power, storage, and network is a necessity.

General-purpose computing: Jobs closed in a single server

Diversified loads, fewer relations, loose coupling between nodes



General-purpose computing: Multiple loads run on multiple nodes, most loads closed on one node due to loose coupling

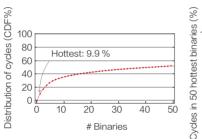


Figure 1: There is no "killer application" to optimize for. The top 50 hottest binaries only cover ≈60% of WSC cycles.

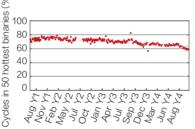
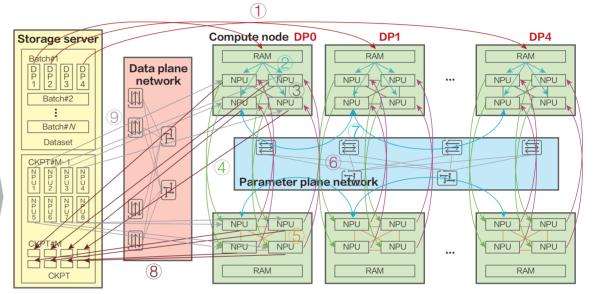


Figure 2: Workloads are getting more diverse. Fraction of cycles spent in top 50 hottest binaries is decreasing.

High concurrency

Intelligent computing: Jobs closed within the system

Single job, long-term full load, tight coupling between nodes (DC as a Computer)



- 1) Splits batch data based on DPs and transmits the data to target nodes.
- (2) Loads data to NPUs for computing.
- ③ Implements TP communication for forward propagation.
- 4 Implements PP communication for forward propagation.
- (5) Implements TP communication for backpropagation.
- ⑥ Implements PP communication for backpropagation.
- 7 Implements DP communication to synchronize gradient data.
- (8) [Periodic] Saves checkpoints.
- (9) [Faulty] Loads checkpoints for restoration.

Intelligent computing: One load runs synchronously on all nodes, with synchronous relationships and tight coupling

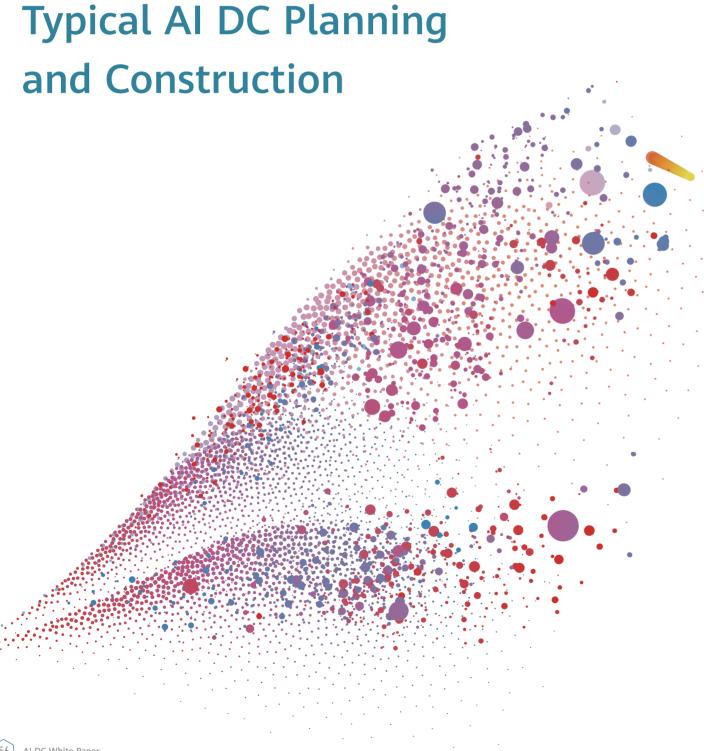
High parallelism

Figure 3-13 From layered decoupling of general-purpose computing to vertical integration of intelligent computing

In conclusion, to build AI DCs, decoupling between the computing backbone layer and model enablement layer are acceptable, as well as the decoupling between the model enablement layer and service application layer. However, intra-layer integration is required due to close relationship between the layers. Only with the integration, computing

power can be unleashed and comprehensive O&M management of computing, storage, and network can be implemented. This vertical integration ensures that Al DCs are able to run efficiently and stably when facing highly parallel computing tasks, thus meeting increasing demands in intelligent computing nowadays.

Chapter 4



Ultra-large AI DCs

This section describes the key requirements and challenges in the construction of ultra-large AI DCs, summarizes the key features of ultra-large AI DCs, and provides suggestions on constructing AI DCs.

Key requirements for construction

To be capable of undertaking foundation model pre-training and inference, there are three key requirements in the construction of ultra-large AI DCs:



Improving the foundation model pretraining efficiency and shortening the training duration

Top Internet and foundation model enterprises hope that the pre-training duration can be as short as possible to implement fast iteration of foundation models and secure a favorable position in the market. In addition, thousands of AI servers consume a large amount of electricity for long-term and heavy-load running. Therefore, improving training efficiency and shortening training duration can win market competition, while saving energy and reducing costs. Given a certain scale of computing power, to improve training efficiency, the key is to improve the effective computing power of a computing cluster.



Meeting the "LACE" requirements during inference

In the inference service oriented to a great number of users, what needs to be focused is the users' experience of "LACE", that is, "Latency" (response latency), "Accuracy" (response accuracy), "Concurrency" (throughput concurrency capability), and "Efficiency" (computing power utilization efficiency).

Latency

Performance measurement of a single inference

Accuracy

Number of variables that can be learned



Concurrency

Concurrent users /data volume

Efficiency

Computing power utilization

Figure 4-1 "LACE" inference indicators

- Latency: Latency directly impacts on user experience. Different application scenarios have different latency requirements. For example, a latency in the Internet application scenario is required not to exceed 30 milliseconds, while a latency in the text dialog scenario is required to range from 30 to 100 milliseconds, and a latency in the voice dialog scenario is required to range from 100 to 200 milliseconds. For a latency-insensitive service such as auxiliary programming or medical diagnosis, a latency of greater than 200 milliseconds is acceptable.
- Accuracy: The accuracy of the system output results needs to be ensured to meet user requirements and

expectations, especially in application scenarios that highly rely on accurate information.

- Concurrency: Hundreds of millions of concurrent requests are processed daily in Internet application scenarios. Therefore, the system must have a powerful throughput capability to cope with high concurrency requirements for
- Efficiency: The computing efficiency of the inference cluster has direct impacts on the final cost and cost control. To lower the costs, the effective computing power utilization of the inference cluster needs to be improved as much as possible.

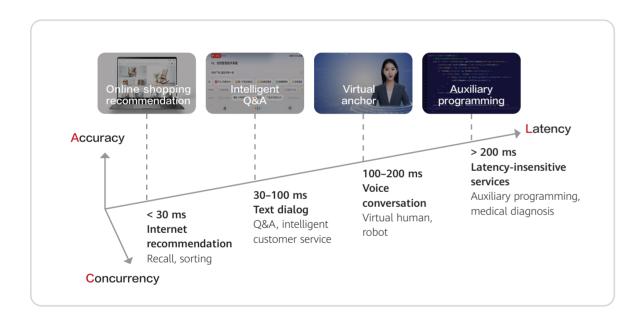


Figure 4-2 Requirements of typical applications for inference performance

Key requirement 3

Improving energy infrastructure to support the sustainable development of ultra-large AI DCs

To support for the sustainable development of ultra-large AI DCs, the energy infrastructure must retain the features of high efficiency, density, flexibility, and reliability.



Energy efficiency improved, PUE ≤1.15, zero carbon emission

Hiah density

Dozens-hundreds of kW/cabinet, power supply & heat dissipation enhanced

Hiah flexibility

type computing power deployment, adjustable power & air-liquid

Multi-gen & multi-

Hiah reliability

No energy infrastructurecaused breakdowns, precise temperature control for better reliability

Figure 4-3 Requirements of ultra-large AI DCs for energy infrastructure

High efficiency: The huge energy consumption of ultralarge AI DCs makes energy efficiency improvement essential. By reducing the power usage effectiveness (PUE) from 1.5 to 1.15, a data center with 100,000 xPUs and a total power capacity of 100 MW could save approximately 200 million kWh of electricity annually.

High density: As the energy consumption of AI chips keeps increasing, the power density of a single cabinet increases by 5 to 10 times. Therefore, the power supply and heat dissipation density must be improved to support more cabinets.

High flexibility: One generation of technology used to be outdated every three years, and now one generation can be updated every year. Hybrid deployment of multiple types and multiple generations of computing power technologies becomes a norm. The energy infrastructure

must be more flexible to support adjustable power and airliquid cooling proportioning, to ensure efficient utilization within the next 10 to 15 years.

High reliability: Intelligent computing devices are high in their costs. Any breakdown will cause great loss. A single point of failure (SPOF) can interrupt the entire cluster. Therefore, the power supply and heat dissipation systems must be highly reliable. In addition, the failure rate of components such as optical modules is closely related to the equipment room temperature. Therefore, more accurate temperature control is required to ensure stable running of the system.

In conclusion, for ultra-large AI DCs that feature ultralarge scale, ultra-high load, ultra-high cost, and ultra-high energy consumption, the ultimate computing efficiency and energy efficiency are the main advantages.



Foundation model pre-training

Distributed inference for massive number of users



Figure 4-4 Leading ultra-large AI DCs demand ultimate computing and energy efficiency

Planning and construction direction 1: Ultimate computing efficiency

To achieve ultimate computing efficiency, multiple strategies and technical means are required:



Foundation model pre-training

To accelerate foundation model pre-training, the effective computing power of ultra-large clusters needs to be improved. The effective computing power of a cluster is determined by three key indicators: cluster computing power scale, cluster model FLOPs utilization (MFU), and cluster availability.



Cluster computing power scale

The raw computing power of a cluster, which depends on the computing power of each node and the cluster scale. The computing power scale is limited by the computing performance of a single node and the number of nodes in the cluster.



Cluster MFU

The ratio of the time during which a computing device executes computing tasks to the theoretical computing time. High MFU means that computing resources are fully utilized, and idle time is reduced.



Cluster availability

The percentage of time during which a cluster is available. High availability (HA) means that the cluster can run properly for most of the time, with reduced downtime and fault time.

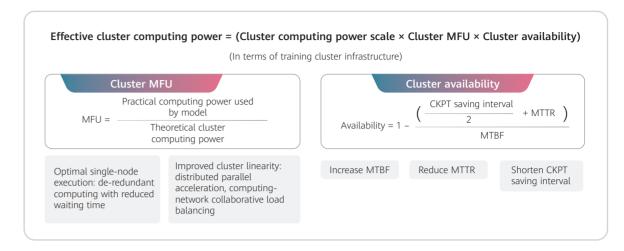


Figure 4-5 Effective cluster computing power

The entire industry is focusing on how to build ultra-large computing clusters, improve MUF and ensure cluster HA. The following key technologies are required:

Key technology

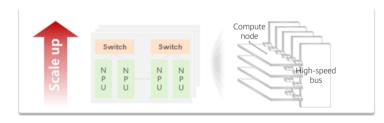
Improving cluster computing power based on supernodes and ultra-large scale networking architecture

In the industry, computing power can be enhanced from two dimensions: scale-up and scale-out.

In the scale-up dimension, the supernode technology is used to improve the unit computing power. The supernode technology uses an innovative peer-to-peer computing architecture to interconnect hundreds of AI chips through a high-speed bus, providing computing power far beyond that provided by the 8 (or 16) AI chips on a single node. This method significantly improves the computing density and performance of a single compute node.

In the scale-out dimension, the overall computing power is improved through the ultra-high-speed and ultra-large scale networking architecture. Ultra-high-speed network technology provides higher bandwidth to reduce network latency, ensuring more efficient data transmission in large-scale clusters. Ultra-large scale networking architecture,

Scale up: The **high-speed bus** increases supernode specifications and **improves the computing power per unit**.



Scale out: The **high-speed network** expands the cluster network scale and **improves the computing power at scale.**

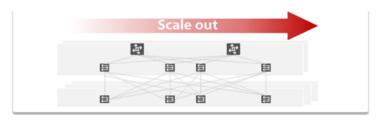


Figure 4-6 Innovation trends of ultra-large-scale networking

such as, Huawei Xinghe Network, uses a deterministic networking architecture consisting of two-layer modular and fixed switches and three-layer only fixed switches. This ensures stable and reliable data transmission in largescale clusters and implements ultra-high-speed network interconnection of hundreds of thousands of xPUs and beyond.

The combination of scale-up and scale-out not only significantly improves the computing power density of a single node, but also implements efficient network interconnection in a large-scale cluster, thereby improving the overall computing power of the cluster.

Key technology 2

Improving cluster MFU through single-node efficiency optimization and cluster parallel optimization

performance.

Cluster MFU = Single-node computing efficiency x Cluster linearity

There are two technical means to improve the cluster MFU: single-node efficiency optimization and cluster linearity improvement.

1 Single-node efficiency optimization

The concept of single-node efficiency optimization is software and hardware collaborative optimization, including the following measures:

• Small-to-large operator fusion: Utilizes technologies such as FlashAttention to fuse multiple small operators into a large operator to reduce scheduling times and HBM read and write overhead, improving computing efficiency.

2 Cluster linearity improvement

The concept of improving the cluster linearity is to optimize the computing power, network, and storage, including the following measures:

• FSPF policy: In an intelligent computing cluster, the interconnection bandwidth between cards, between nodes, and between supernodes converges layer by layer. The full-stack parallelism-friendly (FSPF) technology matches high-frequency and large-traffic parallel computing of Tensor Parallelism (TP) / Expert Parallelism (EP) with the layer-by-layer convergence feature, effectively hiding or reducing communication and increasing the computing proportion.

- Hardware-affinity operator optimization: Optimizes the operator adaptability to hardware, reducing unnecessary scheduling overhead and read and write times of the high bandwidth memory (HBM), improving computing efficiency. For example, Huawei uses algorithm optimization and operator fusion to fully utilize the advantages of Ascend hardware and improve computing
- High-speed bus technology: Uses the highspeed bus to reduce the communication duration. The high-speed bus uses the all-electrical interconnection architecture and supports ultralarge interconnection bandwidth and ultrahigh interconnection bandwidth-compute ratio, effectively reducing communication duration and increasing the computing proportion.
- Computing-network collaborative NSLB technology: Network-side load balancing (NSLB) supports the exchange of training job information between computing and network. Network routing affinity training load ensures that the network throughput reaches 95% or higher, effectively reducing the communication duration and increasing the computing proportion.

These technologies not only significantly improve the efficiency of a single compute node, but also implement efficient linear expansion in a large-scale cluster, thereby improving the overall computing power utilization of the cluster.

Key technology 5

Improving cluster availability using CKPT acceleration and fast fault recovery based on compute-storage collaboration

Cluster availability = 1 -
$$\frac{\frac{\text{CKPT backup interval}}{2} + \text{MTTR }}{\text{MTBF}}$$

Cluster availability is a key factor that affects the effective computing power of a cluster. The cluster availability can be improved by prolonging the mean time between failures (MTBF), shortening the mean time to repair (MTTR), and shortening the checkpoint (CKPT) backup interval. The specific measures are as follows:

1 System-level HA architecture design, prolonging MTBF

- Redundancy design: Redundant components, such as redundant power supplies, network connections, and storage systems, are added at the hardware layer to ensure that the entire system remains available in case of a single point of failure.
- Pre-training pressure test: Comprehensive pressure tests are carried out before training to ensure that the system runs stably under heavy
- Fault prediction: The smart fault prediction is used to monitor and warn about common components with a high failure rate, such as optical modules, NPUs/GPUs, motherboards, and DIMMs. For example, Huawei has developed the optical module channel loss resistance technology to implement uninterrupted training, improving reliability by 10 times. For NPU faults, innovative heat dissipation technologies and fan speed adjustment are used to reduce the NPU operating temperature by 7°C and the failure rate by 30%.

2 Computing-network-storage collaboration, shortening MTTR

- Fault detection: The intelligent monitoring system quickly detects faults and automatically reports alarms.
- Task scheduling: The task scheduling algorithm is used to ensure that standby resources can be activated when a fault occurs.
- CKPT acceleration: The CKPT saving and recovery processes are optimized to reduce the backup and recovery time.

- Computing acceleration: Computing acceleration is used to reduce the execution time of computing tasks.
- Collective communication acceleration: The collective communication mechanism is optimized to reduce communication latency.

For example, Huawei uses computing-network-storage collaborative optimization to accelerate the entire process from fault detection, task scheduling, CKPT acceleration, computing acceleration, to collective communication acceleration, significantly shortening the MTTR.

3 Compute-storage collaboration, accelerating CKPT and shortening the backup interval

- Asynchronous CKPT saving: Asynchronous CKPT saving is used to ensure that backup can be performed without affecting computing tasks.
- Local cache-based CKPT loading: Local cache is used to quickly load CKPT to reduce the recovery time.

For example, Huawei develops the local cache + Near Data Storage (NDS) solution allowing storage pass-through to computing memory, which significantly improves the CKPT read/write speed and shortens the backup interval.

The preceding measures not only prolong the MTBF and shorten the MTTR of the cluster, but also accelerate the CKPT backup, thereby improving the overall availability of the cluster. These methods not only improve cluster reliability, but also provide solid technical assurance for efficient running of ultra-large AI DCs.



Distributed inference for massive number of users

Both typical models with tens of thousands of calls per day and foundation models (with large-scale parameters, ultralong sequences, and multimodal features) evolved continuously face great challenges in improving inference. To address the challenges, the following key technologies can be employed.

Key technology

P/D disaggregation based on KV cache for high inference efficiency of massive number of users

The challenge for massive inference requests is how to quarantee the service quality for hundreds of millions of daily access requests at low costs while ensuring user experience. In other words, the time to first token (TTFT) should be less than 1 second while the time per output token (TPOT) be less than 50 ms. Currently, many optimization techniques for foundation model inference cannot achieve both objectives concurrently. Therefore, the inference process is usually broken down into two stages: Prefill and Decoding. During the Prefill stage, an

input prompt is processed to generate an initial key-value cache (KV cache), creating a context for decoding. Then, the Decoding stage outputs text gradually based on the initial output token and KV cache generated in the Prefill stage. With the process of the two stages and compute result reused through the KV cache, which is synchronized through the high-performance network connecting the two stages, the inference throughput is improved by two to five times while the TTFT is ensured.

The P/D disaggregated inference architecture includes the following components: task scheduling, a plurality of Prefill instances, a plurality of Decoding instances, and a high-performance network. This architecture depends on three key technical elements:

Prefill instances with high computing power: Prefill is a compute-intensive task that requires powerful NPUs/GPUs.

> Large-memory Decoding instances: Decoding is a memory-intensive task that requires large memory capacity

> > and high bandwidth.

High-performance RoCE network: The KV cache needs to be synchronized between the Prefill instances and Decoding instances through a high-speed network. Each NPU must be equipped with an RDMA over Converged Ethernet (RoCE) port of at least 200 Gbit/s to ensure low latency and high bandwidth for data transmission. In addition, the RoCE network adopts the 1:1 nonblocking Clos design to ensure network efficiency and reliability.

This P/D disaggregated inference architecture not only improves the quality of the inference service but also guarantees a good user experience in the case of massive number of concurrent requests, thereby achieving efficient and cost-effective inference services for numerous users.

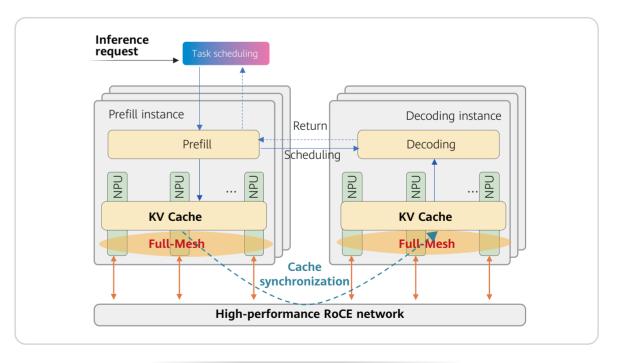


Figure 4-7 P/D disaggregation based on KV cache

Key technology 2

KV cache multi-level caching for high inference efficiency of ultra-long sequence models

Long-sequence models have garnered significant attention due to their robust capabilities in processing complex queries, long text, and even audio and video inputs within several hours. Currently, industry-leading models support long text inputs, such as sequences up to 32 KB, and are widely used for commercial purposes.

However, as sequence length increases, the size of the KV cache also grows during foundation model inference. This not only extends inference time but also significantly increases the memory space required. For example, a typical model with 70 billion parameters may need up to 800 GB of memory to cache the KV cache when processing 1 MB of text. Similarly, a model with 6B parameters requires 100 GB of memory to cache the KV cache for a

256-KB sequence length in single-channel concurrency, and 700 GB for 8-channel concurrency.

To address this issue, a hierarchical caching mechanism has been introduced. In this system, conventional HBM serves as the layer-1 KV cache, host DRAM as the layer-2 KV cache, and high-performance storage devices as the layer-3 KV cache. This multi-level cache management strategy allows for "trading storage for computation". By employing an efficient caching policy, the computing workload is reduced, which in turn decreases inference latency and costs. This approach not only enhances the models' ability to process long sequences but also offers a viable technical solution for large-scale applications.

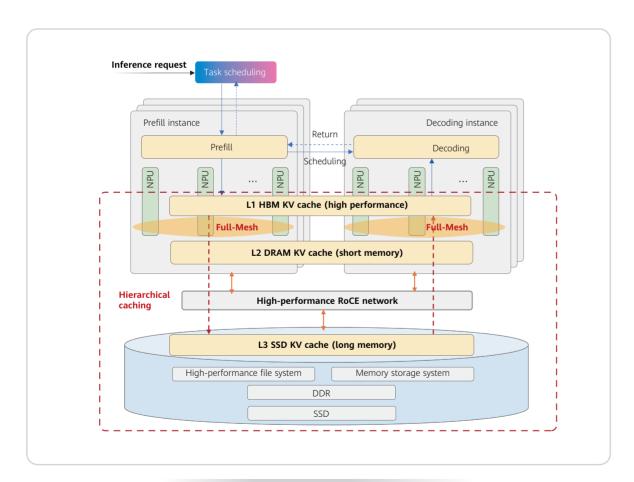


Figure 4-8 KV cache multi-level caching

Key technology 5

Multi-server parallel inference for high inference efficiency of ultra-large and multimodal foundation models

The number of model parameters directly affects the inference computing workload and memory usage. Consequently, ultra-large models suffer from low inference efficiency and high memory demands. For multimodal models, input sequences can extend to contain millions of data records, making long-sequence computing with the Attention mechanism a significant challenge for memory management and computing efficiency. Mainstream multimodal models use a decoder-transformer (DiT) architecture, which necessitates multiple iterations to produce a final result, leading to high resource consumption and prolonged inference duration.

To tackle these issues, a multi-server parallel inference architecture has been developed, significantly boosting inference efficiency and reducing latency. Specifically, NPUs within each compute node utilize a full-mesh architecture to ensure high-speed data exchange. Additionally, NPUs across nodes are connected via a high-performance RoCE network, enhancing the overall system's communication

Furthermore, a series of optimization measures are implemented to further improve inference efficiency.

Mixture of Experts Communication Automatic (MoE) dynamic hybrid parallelism and computing scheduling algorithm to minimize data **convergenc**e to reduce to ensure load transmission overhead the proportion of balancing between across nodes communication time expert modules in the entire inference process

These optimization methods can effectively alleviate the computing and memory pressures associated with ultra-large and multimodal models, thereby improving overall inference performance.

Planning and construction direction 2: Superb energy efficiency

Several key technologies are employed in the industry to achieve exceptional energy efficiency and promote the sustainable development of AI DCs. These technologies focus on creating infrastructure that is highly efficient, dense, flexible, and reliable.

Key technology

Elastic energy infrastructure modules

Elastic energy infrastructure modules must support the hybrid deployment of multi-generation and multiarchitecture computing power to address future service uncertainties. The modular and standardized design allows for rapid pipeline delivery of these modules. Specifically,

a data center is divided into several standardized energy infrastructure modules, with a pre-planned space layout to accommodate computing power devices using various cooling technologies (such as air cooling or liquid cooling) and different power density.

Key technology 2

Ultra-high power supply and distribution efficiency

Collaborative innovation between hardware and software can enhance the efficiency, density, and reliability of power supply and distribution. Advanced data centers in the IT industry currently use the ECO mode with 0 ms switching, which can increase power supply and distribution efficiency to 97.8%. For example, a 100 MW data center with 100,000 xPUs can reduce its annual power supply and distribution loss from 48 million kWh

to 18 million kWh. In addition, device vendors can replace multiple independent products such as transformers, low-voltage distribution products, UPSs, and output distribution products with power modules, and combine them with lithium batteries to double the power supply and distribution density. Finally, predictive maintenance of faults at distribution points and capacitors can significantly improve the reliability of power supply and distribution.

Key technology 3

Ultra-high cooling efficiency

As computing power consumes more power, liquid cooling has become essential. However, there are some concerns about the reliability and O&M complexity of liquid cooling. This is because the proximity of liquid cooling to the server has a significant impact on the operation of IT devices in the event of a leak or interruption. Additionally, new technologies, materials, and devices for liquid cooling require new O&M skills.

To address these challenges, liquid cooling needs comprehensive innovation at all levels, including chips, servers, cabinets, and cooling sources. Taking Huawei's Tiancheng liquid cooling system as an example, chips use cooling plates for cooling, achieving a heat flux of 180 W/cm² and thus meeting the chip cooling needs of over

1,000 W. Servers have a blind-mate design for cooling, electricity, and network, which simplifies deployment and maintenance. Servers also use drip-free quick connectors to improve reliability by integrating them into online leakage monitoring and leakage isolation processes. The cooling source system supports both liquid and air cooling. By upgrading the indirect evaporative cooling air handling unit (AHU), the cooling source system can directly supply liquid cooling servers with room-temperature water (18°C-35°C). The air-liquid convergence design for the cooling source system can reduce the required level of investment, make maintenance easier, and support a PUE as low as 1.10.

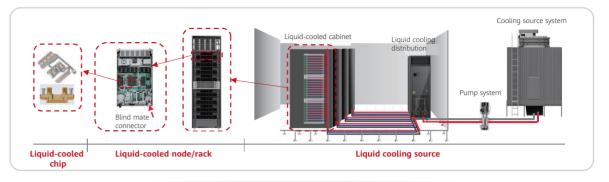


Figure 4-9 Liquid cooling system

Key technology

Collaborative optimization to reduce power consumption

AI-linked optimization reduces power consumption and enhances system reliability. The AI-linked optimization platform collects data on the environment, cooling devices, power supply devices, IT devices, and training and inference jobs. It uses a power consumption optimization model, a component fault warning model, and other models to predict the optimal working parameters in real time and delivers adjustment suggestions to optimize the

outputs of cooling towers, water pumps, and CDUs, as well as the power modes of IT devices. This reduces the total power consumption of data centers. For example, using AI to optimize the Huawei Cloud data center's power consumption based on cloud service perception achieves an accuracy rate of 99.5%, resulting in an 8% to 15% reduction in data center power consumption.

Comprehensive evaluation indicators for ultra-large AI DCs

Level-1 Indicator	Level-2 Indicator	Level-3 Indicator	Description
Training Computing efficiency Inference Computing efficiency	Computing power scale	Computing power scale	Computing power of a single node x Number of nodes (unit: PFLOPS in FP16)
	Computing power utilization	Computing power utilization	Ratio of the actual computing requirement of the model to the theoretical maximum computing capability, in percentage
	Availability	Fault recovery time	Recovery time required for a training job to change from the fault state to the working state, in minutes
		Mean time between failures (MTBF)	Mean time between two adjacent failures, in days
	g	Time to first token (TTFT)	Duration from the time when the model starts to process the input to the time when the first token is generated, in milliseconds or seconds
		Time between tokens (TBT)	Average interval between consecutive output tokens, in milliseconds or seconds
	Accuracy	Accuracy	Proportion of correct answers provided when the model executes inference jobs, in percentage
	Throughput	Throughput	Number of requests that can be processed or number of results that can be generated by the system within a unit time, in tokens/s
Power efficiency	Power efficiency	PUE	Total power consumption of the data center/Power consumption of IT devices (no unit)

Large Al DCs

Large AI DCs carry core services and assets during enterprises' intelligent transformation. Therefore, it is critical to plan and construct AI DCs scientifically and properly. This section describes the key requirements and challenges of large-AI-DC construction, summarizes the key features of large AI DCs, and provides five suggestions based on these features.

Key requirements for construction

Enterprises mainly face requirements and challenges from the following eight aspects.

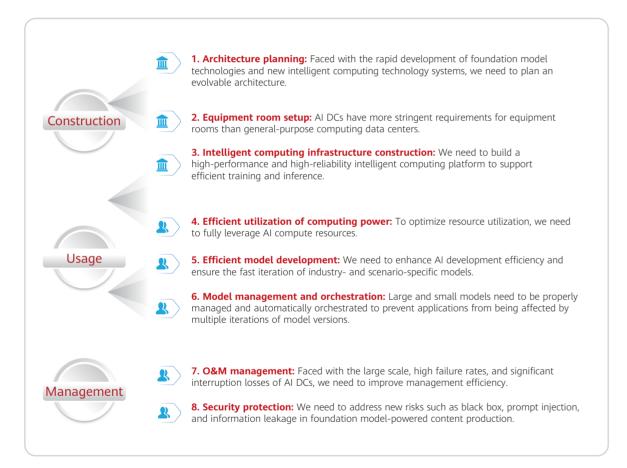


Figure 4-10 Key requirements and challenges faced by enterprises in large AI DC construction

Future-proof and evolvable AI DC architecture planning: Faced with the rapid development of Al technologies and new intelligent computing technology systems, we need to plan a futureproof and evolvable AI DC architecture, avoiding obsolescence upon completion and a significant waste of manpower, resources, and funds.

Intelligent computing infrastructure construction:

We need to make sure that the AI DC computing power infrastructure supports efficient and stable training and fine-tuning, and meets the "LACE" requirements during AI inference.

Full utilization of computing power resources:

To optimize resource utilization, we need to fully leverage AI computing power resources by implementing flexible scheduling (time division multiplexing) and multi-device management (space division multiplexing) during training and inference.

Efficient AI development: To support rapid service growth, we need to enhance AI development efficiency and ensure the fast iteration of largescale AI models and applications for enterprises.

Al DC power infrastructure setup: Al DCs have more stringent requirements than traditional data centers in terms of power supply, cooling, load bearing, and cabling. Therefore, we need to plan and prepare the power infrastructure of Al DCs in a forward-looking manner, based on the future trend of computing power infrastructure, to ensure it meets the evolving requirements of ALDCs.

Model iteration and application compatibility:

Foundation models are developing rapidly, with iteration speeds significantly outpacing those of Al service applications. Therefore, we need to avoid the impact of model updates on upperlayer applications. Meanwhile, we should port small models from the current network to new intelligent computing infrastructure and better combine the advantages of small and large models to achieve a synergistic effect.

Simplified O&M: Faced with the common issues of high failure rates and significant interruption losses in AI DCs, we need to simplify O&M procedures and enhance O&M efficiency.

Countermeasures to new security risks: The emergence of generative AI has brought new security risks, such as uncertainty in the output of foundation models and prompt injection attacks, which must be addressed to ensure the safety of data and systems.



To effectively address these critical needs and challenges, large AI DCs need to possess the characteristics of convergence and efficiency. Specifically, convergence is reflected in four aspects:

- Training and inference convergence: Large AI DCs support not only model training jobs but also AI inference jobs. Supporting both training and inference on the same platform enables resource sharing, simplifies management processes, and improves overall efficiency.
- Intelligent and general-purpose computing convergence: Large AI DCs support not only intelligent computing but also general-purpose computing. The convergence of both types of computing has become the norm in AI DCs.
- Air-liquid convergence: Different types of computing power require different cooling methods: air cooling is typically used for low-density general-purpose computing, networking, and storage devices; liquid cooling is becoming increasingly essential for highdensity intelligent computing. As large AI DCs require both intelligent and general-purpose computing, a coexistence of air and liquid cooling within AI DCs will become an inevitable trend in the future.

• Multi-model convergence: In practical enterprise applications, multiple models are often needed to support a complete AI application. For example, in an AIassisted ophthalmic diagnosis and treatment application, an NLP model is first called to interact with a user, while another NLP model is called to understand the user's queries and plan the diagnostic process. Based on this planning, several CV models are then called for imagebased eye disease screening. Through the collaboration of these models, a diagnostic report is generated. As indicated, large AI DCs often consist of mainstream models such as NLP and CV models, as well as other types of models such as multi-modal and predictive models. These models can range from large to small, allowing for a mixed deployment of multiple modalities.

Ophthalmological diagnosis and treatment assistant in a hospital LLM for interaction LLM for planning Small CV model Figure 4-11 Multi-model diagnosis and treatment assistant

Efficiency is reflected in five aspects: efficient architecture, efficient development, efficient computing power, efficient energy utilization, and efficient management. These are also the five major directions for planning and constructing leading large AI DCs.

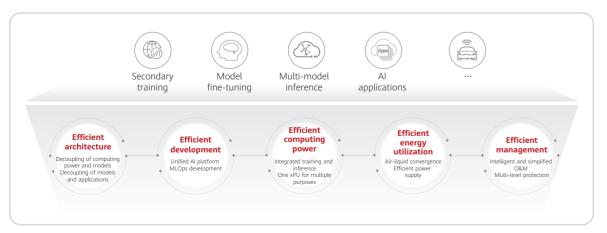


Figure 4-12 Five major directions for planning and constructing leading large Al DCs

Planning and construction direction 1: Efficient architecture

For large enterprises, AI DCs rely on large-scale and diverse computing power resources and models. The efficiency of an Al DC architecture primarily depends on the effectiveness of managing computing power resources and models. Here are specific requirements and solutions:

Efficient computing power management

Managers of enterprise intelligent computing power infrastructure face a significant challenge in efficiently managing and scheduling AI computing power resources. A large AI DC usually includes multiple AI computing power resource pools, such as a training resource pool, an inference resource pool, and an elastic training and inference resource pool. To meet service requirements and improve resource utilization, training and inference jobs need to be flexibly scheduled and switched between these resource pools.

In addition, some large enterprises may have multiple AI DCs distributed in different locations, requiring computing power scheduling across the data centers. Furthermore, in the event of unexpected service requirements, the computing power resources of the enterprises' self-built Al DCs may be insufficient, necessitating the rental of

This architecture needs to meet the following requirements:

- Unified management: A unified management platform is used to manage computing power resources on public clouds as well as in local and remote DCs, forming a logically unified resource pool.
- Flexible scheduling: Computing power resources from different resource pools can be flexibly scheduled based on service requirements, and training and inference jobs can be dynamically switched over.
- Standardized APIs: Standardized APIs are provided to enable upper-layer models and applications to seamlessly connect to underlying computing power resources to simplify the development and O&M processes.

computing power resources from public or industry clouds. In such cases, the capability of scheduling computing power resources across local and public clouds is required.

To address these issues, an efficient architecture for managing and scheduling computing power resources is needed to decouple models and computing power resources. The unified computing power resource management and scheduling platform provides a standard API encapsulated from computing power services for the training and inference jobs of upper-layer models. This approach hides the complexity of managing computing power resources across multiple pools, centers, and clouds, creating a logically unified AI computing power resource pool. This ensures that AI computing power resources are centrally visible, manageable, and controllable, which in return improves resource utilization.

- Cross-domain scheduling: Cross-DC and crosscloud computing power scheduling ensures that computing power resources can be quickly expanded in the event of burst service requirements.
- Resource visualization: Resource monitoring and visualization tools are provided to clearly display the status of computing power resources, which facilitates management and optimization.

With such an architecture, enterprises can improve the utilization of computing power resources, and maintain agility and flexibility in complex and changeable service environments to better cope with future challenges.

Efficient model management

When applying large models, enterprises also face some uncertainties. Due to fierce competition, large models are updated rapidly, and this brings both opportunities and challenges. Specifically, the capabilities of available large models become stronger, but enterprises also need to ensure that their upper-layer applications are not affected by the frequent changes to large models.

An efficient architecture is required to decouple models from applications. A unified model management and orchestration module provides a standardized API encapsulated from model capabilities to mitigate the impact of model changes or replacements on applications. This ensures the stability and reliability of upper-layer applications while also leveraging the latest model capabilities.

This architecture needs to meet the following requirements:

Model orchestration: A model management and orchestration module is used to separate models from applications, so that model changes do not directly affect upper-layer applications. This means that even if the underlying models are updated, the upper-layer applications can still run stably.

Automatic deployment: Models can be automatically deployed, simplifying the rollout process and ensuring that new models can be quickly and securely deployed in the production environment.

Standardized APIs: Standardized APIs are provided to enable upper-layer applications to seamlessly connect to the underlying models, simplifying the development and O&M processes. No matter how the underlying models change, upper-layer applications only need to call the unified APIs.

Monitoring and feedback: A mechanism for model running status monitoring and feedback is used to monitor real-time performance and effect of models in practice. Optimizations and adjustments are performed based on the feedback.

Model version management: Model version management ensures that models of different versions can coexist and be flexibly switched over as required. In this way, different versions of models can be better managed and utilized.

With such an architecture, enterprises can not only improve the stability and reliability of model applications, but also make full use of the latest model capabilities in the rapidly changing market environment to maintain a competitive edge. This architecture improves the flexibility of enterprises and enhances the overall stability and scalability of the system.

Through the preceding two aspects of efficient management, enterprises can not only improve the utilization of computing power resources, but also ensure the stability and flexibility of model applications, thereby enabling themselves to maintain leading positions in the market, even when facing fierce competition.

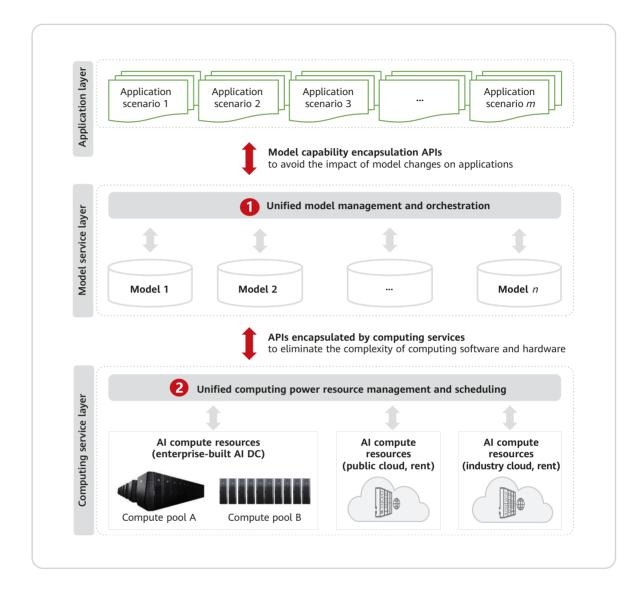


Figure 4-13 Technical architecture of layer-based decoupling, achieving efficient computing power and model management

Planning and construction direction 2: Efficient development

An efficient AI development platform needs to be built to cope with common problems such as low development efficiency, high costs, low model consistency, low reliability, and long model deployment time. The core capabilities of the platform are as follows:

Providing an integrated development environment which supports local development and remote training

During AI development, a user usually writes codes in the integrated development environment (IDE) on a PC. During model training, and especially in large-model development, NPU/GPU resources on the server side are required due to the limited local computing power. In this case, the user needs to upload codes to the server and prepare a complex training environment. This process usually involves multiple tools and interfaces. The AI development platform needs to provide the remote training IDE plug-in to allow the user to start model training jobs and monitor the training process in the local IDE. After installing and configuring the plug-in, the user can upload codes and complete model training using the plug-in without leaving the IDE environment. This simplifies the operations and improves the development efficiency.

Enabling wizard-based AI development workflows (MLOps)

An AI development workflow (MLOps) is an automated tool designed to simplify the AI model development process. Existing AI codes are valuable assets and reusing these assets effectively is challenging. Al development has high requirements and it is difficult for an IT engineer to independently complete model development tasks.

The AI development workflow helps solve this problem by institutionalizing the training and deployment processes for Al applications in specific service scenarios. During the reuse of a process, a workflow provides a wizard-based GUI which includes the entire process of model development and deployment, containing models, codes, and parameters, regardless of whether the process is being used in a similar scenario or for model data updating. In this way, even common IT engineers can complete the entire process from model training to deployment with the help of the wizard

Building an AI asset center to provide efficient Al asset management Data, model, and scenario-based workflows are valuable AI assets. By accumulating and reusing these AI assets, enterprises' development departments can build up their core competitiveness and improve the development efficiency of AI applications. An AI asset center can efficiently manage data, model, image, and workflow assets. A user can search for the assets they require in the AI asset center and learn more about the asset through figures and text. In addition, the user can easily import assets onto the AI development platform, with no need to spend a lot of time searching for or importing models or worry about complex asset management and asset losses. More importantly, the AI asset center promotes the sharing and reuse of assets within the team. One-click import further eliminates the need for manual importing and configuration, and this significantly improves the team's development efficiency.

Planning and construction direction 3: Efficient computing power

As AI computing power resources are limited, enterprises are doing their best to maximize computing power utilization by using the following two methods:

Dvnamic and flexible computing power scheduling (timebased reuse)

In actual AI applications, the secondary training of industry-specific models is usually performed every few months, and each training period lasts about one month. As a result, the computing power used for training is idle most of the time. However, the amount of inference compute resources required is closely related to the numbers of services and user visits and these both fluctuate a lot. For example, peak times would be during the day and on workdays, and off-peak times would be at night and on weekends.

Therefore, enterprises urgently need an AI computing power platform that supports flexible switchovers between training and inference jobs. For example, enterprises can choose to perform inference jobs on workdays and switch to model finetuning jobs on weekends. Alternatively, they can switch to inference jobs when the computing power used for training is idle. Such time-based reuse enables efficient utilization of computing power.

Computing power segmentation with one card for multiple purposes (space-based reuse)

In the actual AI construction of enterprises, there are both large language models (LLMs) and small models such as CV models. Small models require a small amount of computing power and memory resources and generally consume partial resources of an intelligent computing card.

In this case, computing power segmentation is required to enable multiple tasks to run concurrently on one intelligent computing card, thereby increasing the efficiency of resource usage. This is space-based reuse of computing power.

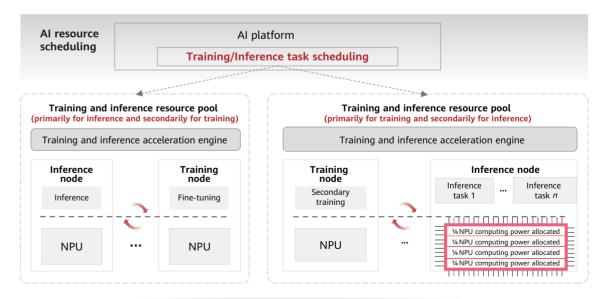


Figure 4-14 Time/Space-based reuse of computing power

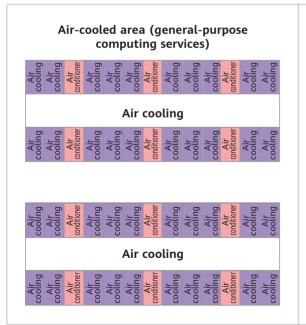
With these methods, enterprises can improve the utilization of computing power resources and flexibly allocate resources during peak hours to ensure that key tasks run smoothly. This not only helps reduce operational costs, but also improves the overall competitiveness of enterprises.

Planning and construction direction 4: Efficient energy utilization

The core part of efficient energy infrastructure is the hybrid deployment of liquid cooling (for intelligent computing) and air cooling (for general-purpose computing). The current challenges are as follows:

- Challenge 1: Dynamic adjustment of air-liquid hybrid cooling ratio. Both intelligent computing and general-purpose computing need to support independent expansion. The flexible air-liquid hybrid cooling ratio helps AI DCs adapt to service uncertainties.
- Challenge 2: Combination of high- and low-power-density cabinet deployment. The power of each cabinet ranges from dozens to hundreds of kilowatts for intelligent computing but remains below 10 kW for general-purpose computing. Enterprises need to adapt the power supply and cooling systems to meet the requirements for varying power densities.
- Challenge 3: Adaptation to complex equipment room conditions. In new deployment scenarios, areas need to be planned. In reconstruction scenarios, the interfaces need to be consistent with those in the existing system, while the running of existing services is

not affected.



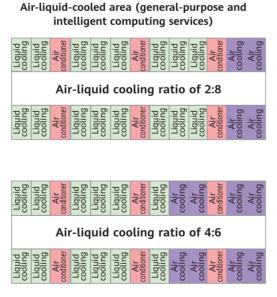


Figure 4-15 Area planning for large AI DC equipment rooms

It is recommended that the AI DC be divided into an air-cooled area and an air-liquid-cooled area. The former is oriented to deterministic general-purpose computing services, while the latter is used as elastic space for the on-demand expansion of both types of computing services. Assume that the total number of racks, the cooling system, and the maximum power supply capability remain unchanged. Power and cooling capacity requirements of liquid-cooled cabinets must be fully considered during space design. Power supply and cooling equipment need to be deployed based on the actual air-liquid cooling ratio.

In the air-liquid-cooled area, air-liquid hybrid cooling micro modules can be used to keep the airliquid cooling ratio flexible. These modules can be embedded with liquid cooling pipes, liquid cooling CDUs, power distribution units, and modular in-row air conditioners (or use a shared air conditioner pool in a room). Inside each module, liquid-cooled cabinets are deployed on demand with a proportion from 0% to 100%. Take Huawei micro modules as an example.

A liquid-cooled cabinet uses cold plate liquid cooling and air-cooled air-conditioning, which supports a maximum of 66 kW. An air-cooled cabinet is deployed with general-purpose computing, storage, network, and security equipment, and it supports a maximum of 35 kW. In this case, air-cooled air conditioning is used for heat dissipation.

The air-liquid hybrid micro-module provides standardized water, power, and management interfaces, enabling quick CI/CD. The water interface connects seamlessly to the data center's cooling source system, supporting water temperatures ranging from 18°C to 35°C. The electrical interface integrates with the low-voltage distribution system, utilizing intelligent busbar distribution with multiple hotswappable switches available in 16 A/32 A/40 A/63 A configurations to meet varying cabinet power density requirements. Management is unified through standard protocol interfaces, allowing centralized access to the DCIM O&M platform for streamlined administration.

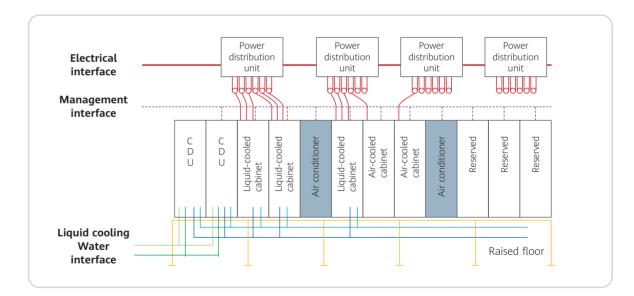


Figure 4-16 Air-liquid hybrid micro-module

Planning and construction direction 5: Efficient management

To address the increasingly complex O&M challenges, as well as emerging security risks in AI DCs, it is essential to develop a next-generation O&M platform and security protection system to enhance management efficiency.

The next-generation O&M platform should feature comprehensive monitoring, fault prediction, rapid fault locating, and recovery capabilities, simplifying routine operations and reducing O&M difficulties. Key functionalities include:

- Full-link, visualized monitoring: By offering real-time monitoring of system performance, the platform ensures comprehensive oversight of computing, storage, and network resources, enabling timely detection of anomalies.
- Automatic cross-domain fault detection: Leveraging a pre-built fault library and end-toend information flow analysis, the platform can automatically match and detect common faults, enabling proactive fault awareness.

- Rapid cross-domain fault locating: Utilizing advanced fault detection technologies such as log analysis, minute-level network traffic analysis, and storage fault and performance analysis, the platform allows for quick and accurate fault locating, minimizing troubleshooting time and preventing training job interruptions.
- Swift cross-domain fault recovery: A robust fault recovery mechanism ensures that, when an issue occurs, rapid response measures can restore system operations and minimize the downtime.

The next-generation security protection system **should focus** on model security on the basis of the basic environment security and secure operations of traditional data centers, while also reinforcing data and application security for intelligent computing services.

Comprehensive evaluation indicators for large AI DCs

Level-1 Indicator	Level-2 Indicator	Level-3 Indicator	Indicator Description
Architecture efficiency	Computing power management and scheduling capability	Resource pool type	Number and scale of AI resource pools that can be managed and scheduled by the platform, including self-built AI computing power resource pools as well as leased AI computing power resource pools within public or industry clouds
	Model management and orchestration capability	Model type	Model types supported by the platform, including traditional models (CV, OCR, and ASR), NLP models, and multi-modal models
Development efficiency	Development efficiency	Model development duration	Duration for developing a model. The unit is the number of person- days required for training or fine-tuning. The platform must have comprehensive model development engineering suites to support model design, model commissioning, training script writing, and model accuracy issue location.
	Deployment efficiency	Model deployment efficiency	Duration required for deploying an instance, in seconds. The platform must be capable of automatic model deployment and provide multiple deployment modes, such as single-node multi-xPU and multi-node multi-xPU, to implement efficient model deployment.
Computing power efficiency	Computing power utilization	Resource utilization	The unit is percentage. The platform monitors the resource utilization and supports dynamic job management and scheduling. This maximizes the utilization of AI computing power resources across centers, pools, and clusters.
	Training efficiency	Model training duration	Duration from the time when the model training starts to the time when the model accuracy converges to the target value, in days
	Inference performance	Inference TTFT	Duration from the time when the model starts to process an inference request to the time when the first token is generated, in milliseconds or seconds
		Inference TBT	Average interval between consecutive output tokens during model inference, in milliseconds or seconds
		Inference throughput	Number of tokens for the inference requests that can be processed plus the number of tokens for the results that can be generated by AI clusters within a unit time, in tokens/s
Power efficiency	Power efficiency	PUE	Total power consumption of the data center/Power consumption of IT devices (no unit)
O&M and security	O&M	Fault location time	Average duration from the time when an AI cluster fault occurs to the time when the fault is first detected during job running, in minutes. This duration is also called the mean time to detect (MTTD).
	Security	Qualification rate of generated content	Proportion of the generated content that meets the compliance requirements, in percentage. The generated content is randomly selected either manually or by keyword and classification model for security evaluation.

Small AI DCs

Constructing small AI DCs to carry local services is critical for smaller clients such as enterprise branches, healthcare institutions, and educational organizations as it facilitates intelligent service upgrades. This section describes the key requirements and challenges of small-AI-DC construction, summarizes the key features of small AI DCs, and provides four suggestions based on these features.

Key requirements for construction

The five core requirements for small AI DCs are as follows:

Fast application development Fast data processing, model fine-tuning, and RAG Fast deployment Simplified O&M and upgrades Unified full-stack One-stop quick intelligent O&M of deployment and Five core software and hardware remote deployment across multiple domains requirements **Efficient resources** Security assurance More applications with Security assurance for limited computing power, foundations, data, models, increasing computing and applications power resource utilization

Figure 4-19 Five core requirements for small AI DCs



Efficient resources: Due to environmental constraints, the available computing power resources in small AI DCs are relatively limited. Therefore, the design must ensure that more service applications can be carried by these limited resources to maximize the utilization of computing power resources.



Simplified O&M: Small AI DCs may lack dedicated O&M personnel. Therefore, the design should prioritize reducing the failure rate and supporting remote operability to simplify daily maintenance processes.



Fast deployment and upgrades: Some small AI DCs are located in remote areas, and enterprises want to deploy them with minimal manual intervention. Ideally, the delivery team would complete all installations with a single site visit, or even remotely.



Security assurance: Small AI DCs are typically directly connected to various sensing devices such as smart cameras and sensors, making it crucial to ensure the security of small AI DCs and their associated devices to prevent external intrusions.

Fast application development: Professional users who fine-tune models based on private datasets need a simple and easy-to-use toolchain to support the entire process, including data preparation, model training, inference deployment, and RAG, and make the entire process as simple and fast as drag-and-drop.

To meet these five requirements, we propose the concept of flexible, fast, lightweight, and easy construction of small AI DCs. This concept emphasizes that the data center design should be lightweight, easy to deploy and upgrade, flexible, and easy to manage and maintain.

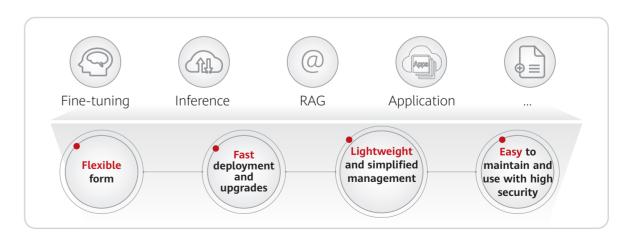


Figure 4-20 Leading small Al DCs that need to be flexible, fast, lightweight, and easy to use

Planning and construction direction 1: Flexibility

Flexibility refers to being flexible and having diversified forms. AI DCs can be classified into node-type and cabinet-type data centers based on their physical forms; they can be classified into those requiring independent deployment and those requiring cloud-edge deployment; they can be classified into NLP applications (such as intelligent customer service and code generation), CV applications (such as industrial quality inspections and medical imaging), and multi-modal applications (such as office assistants and text-to-image generation) based on their functions.

Setting up a lightweight AI platform with unified standards is crucial for diverse small AI DCs. First, a unified southbound access standard can accommodate diverse hardware bases and support access standards of various devices, including video, IoT, and smart devices. Second, a unified API openness standard allows quick calling for fine-tuning, inference, and application processes. Third, unified data- and management-plane standards enable multi-dimensional collaboration with the center, and this supports continuous learning and iterative upgrades for small AI DCs at the edge.

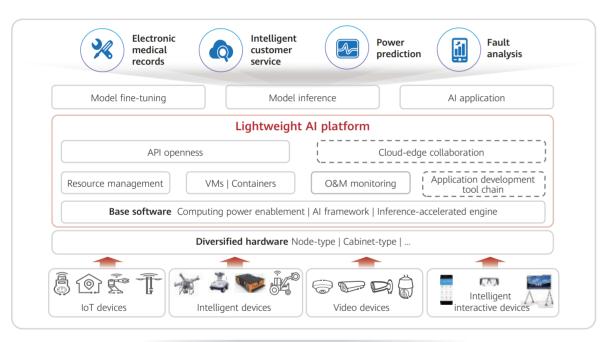


Figure 4-21 Lightweight AI platform with unified standards

Planning and construction direction 2: Speed

Here we are looking at the speed of deployment and upgrades. Multi-dimensional collaboration between applications, models, and data on the cloud and at the edge is supported. Small AI DCs at the edge can access the AppStore on the central cloud to download and deploy models and applications through online subscriptions. The central cloud can remotely upgrade the models and applications for scattered small AI DCs in a unified manner. Real-time data collected at the edge can be uploaded to the central cloud for continuous model iteration, and this enables "learning on the job".

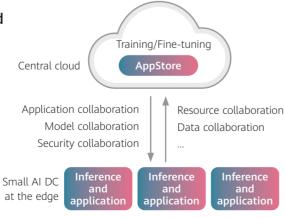


Figure 4-22 Cloud-edge collaboration architecture

Planning and construction direction 3: Weight

Weight refers to the load, resources, and management. **Light load:** Typical inference applications have less than 10 concurrent flows and small-scale model parameters. As a result, the computing power needs to be divided between multiple models concurrently, and this improves resource utilization. Light resources: Resources can be scaled

on demand by just one node or even one xPU, and this saves resources and reduces costs. Light management: Minimizing the resources required for deployment frees up resources for services. The goal is to enable deployment with a minimum of 1 CPU and 1 GB.

Planning and construction direction 4: Ease of use

This indicator considers how easy it is to maintain and use the AI DC and how easy it is to ensure its security.

Independent deployment: In terms of easy use, users need to fine-tune scenario-specific models at the edge. Small AI DCs integrate a complete toolchain to meet the requirements for data processing, model fine-tuning, inference deployment, and knowledge base generation, and thus enable AI application rollout in days. In terms of easy O&M, a unified approach is employed for the full-stack management of hardware, software, and applications related to computing, storage, and network. These provide visual O&M based on the tasks that are running and they reduce the manpower needed for O&M. In terms of security, basic capabilities in model content security, data security, and device access security are provided to meet the security requirements from data access to application deployment.

Edge cloud: In terms of easy use, edge AI DCs are an extension of the central cloud, and all operations can be completed in batches on the central cloud. In terms of easy O&M, the central cloud manages massive numbers of scattered small AI DCs in a unified manner, and this enables remote O&M at the edge. In terms of security, the edge supports capabilities such as OS security, device access authentication, and security detection to prevent it from becoming a weak point that is more susceptible to network intrusion.

Comprehensive evaluation indicators for small AI DCs

Level-1 Indicator	Level-2 Indicator	Level-3 Indicator	Indicator Description
Flexibility	Model adaptation	Number of adapted models	Number of models that can be directly downloaded and for which the model image environment is supported
Speed	Deployment and upgrade	Deployment time	Number of days required for deploying the hardware, OS, management platform, and O&M platform when the equipment room meets the deployment conditions
		Upgrade time	Number of minutes that it takes to upgrade the models and applications. If the upgrade is performed onsite, the round-trip travel time to the site must be included.
Weight	Resource requirements	Minimum number of nodes required	Optimal number of server nodes for building a small AI DC, including general-purpose and intelligent compute nodes
		Management resource requirements	CPU core quantity and memory capacity (in GB) required for establishing a management platform
Ease of use	Easy Number of DCs maintenance operated per person		Number of small AI DCs that can be maintained by each person, in sites/person

Chapter 5

AI DC Construction and Development **Initiatives**

Initiative 1: Taking a forward-thinking approach to AI DC construction

Enterprise adoption of AI should be planned and purposeful. When approaching intelligent transformation, enterprises must first consider how AI will benefit them before taking a step-by-step and efficient approach to implementation. This is how they can achieve sustainable business growth and innovation. We encourage all enterprises to focus on the following six pillars when adopting AI, as taking a forward-thinking approach to Al infrastructure construction is critical for enterprises to get a head start in the AI era.

Strategic resolve

Investment into innovation in intelligence should be considered strategic investment for enterprises. This kind of investment involves huge amounts of resources and returns are not guaranteed. In addition, it requires collaboration between different services and technological domains. Attention and support from top executives are the key to the success of such investment. Top executives need to consider their enterprises' actual needs and decisively establish a dedicated work team to ensure efficient collaboration between involved departments and business domains. They also need to organize comprehensive resource support for Al infrastructure construction and scenario-specific applications.

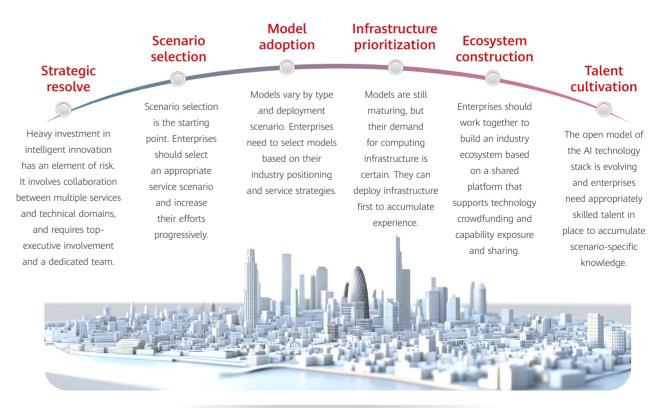


Figure 5-1 Six pillars enterprises should focus on when deploying AI

AI DC White Paper

Scenario selection

Selecting application scenarios is a key starting point for AI adoption, as it determines the scale, performance, reliability planning, and future-oriented evolution requirements of the AI infrastructure to be deployed. Based on actual business needs, enterprises should prioritize scenarios with clear challenges and high technological applicability, such as customer service, supply chain management, product design, and data wrangling. Enterprises should tackle smaller challenges before larger ones, to build upon their own momentum. This means first piloting AI in scenarios where AI applications can be most easily deployed and quickly create benefits. This will help them gather both experience and confidence, before expanding to broader and more complex fields, and ensure that AI applications are steadily advancing and constantly optimized.

Model adoption

When promoting enterprise-class AI applications, adopting a suitable combination of AI models is critical. Al models can be categorized multiple ways. By scale, there are large and small models. By application scenario, there are foundation models, industry-specific models, and domain-specific models. By function, there are natural language processing (NLP) models, computer vision (CV) models, predictive analysis models, multimodal models, and mechanistic models. By training and deployment scenario, AI models can be deployed for full training, incremental learning (secondary training), finetuning, and inference applications.

Current industry practices show that combinations of different types of models are often used to meet the complex and changing enterprise needs. Therefore, enterprises need to select the specific types of AI models that best meet their industry positioning and business strategies. In addition, they need to comprehensively evaluate technical difficulties, data security, and costeffectiveness, and develop an appropriate training policy to ensure that the models are accurate and practicable.

Infrastructure prioritization

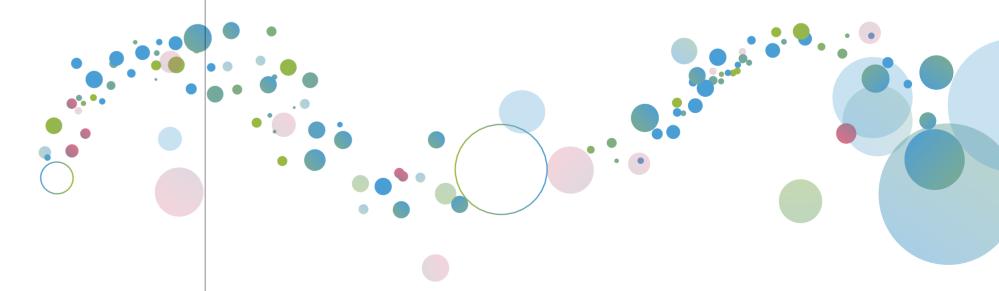
The stability of long-term training and iterative upgrade with AI models depends on the strength of an enterprise's computing power. Models are still maturing and we don't know everything about them yet, but their demand for computing power is certain. Enterprises need to prioritize the construction of high-performance, highly-elastic AI DCs to serve as their computing infrastructure platform for AI. This allows the platform to adopt advanced hardware equipment and software systems, and ensure efficient and stable supply of computing resources through software-hardware collaboration as well as traininginference collaboration. In addition, a comprehensive data governance system can be established based on AI DCs, spanning data collection, cleansing, labeling, storage, and analytics. This will allow enterprises to feed high-quality data into their AI models. In addition, enterprises will need to pay attention to the eco-friendliness and energy efficiency of AI DCs to ensure sustainability.

Industry ecosystems

Al development requires vertical integration, so enterprises looking to apply AI will need to actively participate in or lead the construction of an open, unified, and everevolving industry ecosystem centered around AI DCs that aggregates their own products alongside those of other industry players. This will allow them to benefit from industry-wide technology crowdfunding, capability exposure, and resource sharing. It also fosters unified industry-wide architectures, standards, and data specifications. These types of ecosystems accelerate AI innovation and application, lower their industry's threshold to AI development, improve the competitiveness of all industry players, and help the entire industry achieve shared success.

Talent cultivation

Al adoption in vertical industries and their core production scenarios will need large amounts of talent with comprehensive skills, including insights into both technical details and application scenarios. This means enterpriseled talent cultivation will be crucial in the field of AI application. Enterprises are uniquely situated to create a solid foundation for AI innovation, as well as for intelligent transformation, by leveraging effective recruitment and internal training programs.



Initiative 2: Building intensive and green Al DCs

To promote the intensive construction and green development of AI DCs, AI DC strategies should be developed with a focus on three areas: the distribution of intensive computing power, the adoption of innovative technologies, and the formulation of standards. More efforts are needed to guide the construction and application of AI computing power resources in order to avoid common issues such as redundant construction and low resource utilization, and ensure the high-quality development of the AI DC industry.

Promoting large-scale and intensive AI DC **development:** First, we should support enterprises as they shift from traditional data centers to AI DCs, gradually increase the scale of AI DCs, and promote the large-scale development of data centers. Second, we should increase the supply of intensive computing power and work with carriers, cloud service providers, and various computing platforms to coordinate the development of AI computing power and general-purpose computing power. This will satisfy the computing-power requirements for different types of services, such as balanced, computing-intensive, and storage-intensive services. Finally, we should improve computing-resource collaboration between AI DCs across different regions to balance loads and dynamically allocate

resources in order to improve system-wide elasticity and

Increasing the importance of green energy: It is challenging to make AI DCs run on green energy because there is a trade-off between ensuring a secure and stable power supply while also increasing the percentage of clean energy in the power mix. Technologies for microgrids, power source-grid-load-storage coordination, and new energy storage are attracting more attention from the industry. We should consolidate the foundation of green computing power and industry innovation. We call on relevant enterprises to increase their efforts to research and build microgrid systems, explore energy storage materials and develop related technologies, promote innovation in liquidcooling technology, and explore ways to recover dissipated heat in order to ensure a stable and uninterruptible power

Developing green computing standards: Standardization is the key to the development of the green computing industry. Although there already are data center standards and specifications on carbon usage effectiveness (CUE), carbon neutrality assessment, and IT equipment energy efficiency, a comprehensive framework of green computing standards based on industry consensus is needed. Currently, the most basic standards, such as terms and definitions, are not complete, and no consensus has been reached on the evaluation and classification of energy consumption. It is important to establish and then continually improve the green computing-power standards framework and promote its application in computing infrastructure.

Initiative 3: Building an open and collaborative AI ecosystem

When the development of a revolutionary technology starts to accelerate, there are usually two ways forward. One is to build a complete and closed ecosystem led by a single large enterprise and supported by numerous small- and medium-sized enterprises (SMEs). The other is to build a multi-core open ecosystem that is **initiated** and optimized by many different enterprises which are mutually compatible and synergistic.

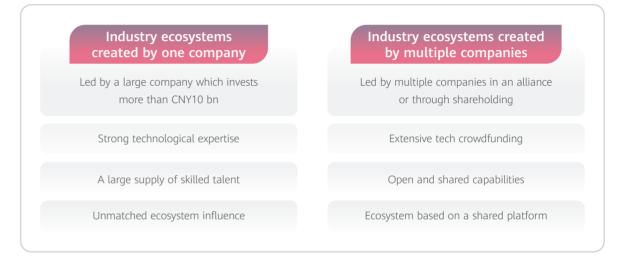
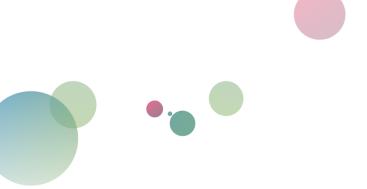


Table 5-2 Industry AI ecosystem construction models

Take the iPhone as an example. The launch of the iPhone took the world by storm and revolutionized the smartphone market. Apple and Google chose two drastically different development paths. Apple developed iOS, an operating system that was specifically designed for the company's own products, while Google developed the open-source Android operating system. These two operating systems have dominated the mobile phone industry. They can be seen as representing a singlecore closed ecosystem and a multi-core open ecosystem, respectively. Their ecosystem construction strategies provide a valuable reference for the developers of new revolutionary technologies like AI as they look to create their own ecosystems.

In view of the different characteristics of closed and open ecosystems, as well as the interdisciplinary nature and complexity of the AI domain, we call for enterprises, competent authorities, and academia to

work together to create an open and collaborative Al ecosystem within vertical industries. First, the combination of expertise from industry, academia, and the government may provide the necessary impetus for making technological breakthroughs and developing new applications. Second, unified industry standards and specifications can be established based on an open ecosystem platform, making cross-model, cross-company, and cross-industry interconnection and interoperability possible. Third, building an open communication and cooperation mechanism to jointly address ethical, legal, and social issues related to AI development will ensure that the development of new technologies and applications is both legitimate and sustainable. Fourth, promoting talent cultivation and mobility, as well as joint training, competitions, and cooperation programs will improve the technological skills and innovation capabilities of the industry as a whole and provide the talent needed for a thriving AI ecosystem.



Initiative 4: Building three foundations to accelerate AI applications across industries

Al is disrupting our expectations of future enterprise and market operations at an unprecedented speed. However, the value of AI does not come from the existence of technology in and of itself, but rather from the impact that this technology has when it is applied to different industries. Therefore, our ability to accelerate the application of AI and have it permeate more and more business processes in a wider range of industries is key to resolving key issues and creating substantial value for enterprises. Regardless of how AI came into being, it will ultimately be used to transform productivity, upgrade industries, and create economic benefits.

We recognize that the AI industry needs a stable and strategic framework as a solid foundation that will enable AI to thrive in the complex and everchanging market environment, overcome technological challenges, and achieve the goals that we have set. We propose the construction of three strategic foundations: a solution foundation, an ecosystem foundation, and a talent foundation. These three foundations can support the AI industry by ensuring that technology applications accurately meet enterprise needs so that AI applications can more effectively benefit industries.

The core is creating value and resolving key issues for enterprises.

Solution foundation

Driving cross-technology collaboration and systematic innovation to build solid computing infrastructure

Ecosystem foundation

Diving deep into industry scenarios and growing through joint innovation

Talent foundation

Developing industry AI application talent with technology expertise

Figure 5-3 The three foundations that will accelerate further AI applications across industries

The solution foundation: Building solid computing infrastructure based on AI DCs

Solutions are the foundation of technological practices. The goal of establishing a solution foundation is to consolidate AI-DC-centered AI infrastructure for enterprises based on unified technological architectures, industry standards, and data specifications. The core of the solution foundation is

to define a target reference AI DC architecture, increase computing efficiency, standardize the technological architecture, and standardize the ecosystem in order to align technologies and business objectives and accelerate technological innovation and applications.

The ecosystem foundation: Fostering an open and collaborative industry ecosystem

The ecosystem foundation is an open, shared, and closely connected system that facilitates ecosystem collaboration. Joint innovation across different industries and domains can be applied in core production scenarios and promote the application of AI. First, we should improve synergies with the industry ecosystem. We can encourage collaboration between hardware manufacturers, software developers, cloud service providers, and end users, and accelerate technology iteration and upgrades to improve

the performance and efficiency of the entire AI DC system. Second, we should foster an open and collaborative industry ecosystem. We can work with different enterprises within the industry to discuss and resolve key industry issues, promote technological innovation and standardization, and drive information sharing and technical exchange between industry players. This will enable the deployment of AI in different industries and create new value for society.

The talent foundation: Supporting the development of tech-savvy industry talent with the aim of driving new AI applications

Building an AI talent pool is essential for the growth of the AI industry. To build a talent foundation for AI, great importance needs to be attached to the cultivation of technical talent, and especially of professionals who are proficient in developing underlying operators and constructing industry acceleration libraries. They will become the driving force behind further exploration of cutting-edge technologies. AI deployment will lead to there being more AI application experts in vertical industries, such as industry-specific AI consultants and educators, who can combine AI with industry requirements to create practical value. We also need to advocate for

a new model that integrates industry with education.

We should provide refined and targeted training programs for developers, so that they can master heterogeneous computing power applications, to ensure a steady supply of innovative talent for different industries.

The three strategic foundations not only provide a solid support system for the adoption of AI applications, but also drive the long-term development of intelligent industries. By consolidating the solution, ecosystem, and talent foundations. Al can help industries become more competitive and mature faster.

Epilogue

As the world embraces digital and intelligent transformations, industries need to clear up any existing confusion and overcome the challenges that they face with regards to deploying data centers. In the meantime, the arrival of the intelligent era driven by AI foundation models is increasing the demand for AI DCs.

Currently, industries lack a comprehensive and systematic framework for planning and building AI DCs. Huawei presents this white paper, which is the result of collaboration between more than twenty industry experts. This white paper distills Huawei's in-depth research over the past few years, its recent innovation in key products and solutions, and its extensive industry exploration and practices. The white paper also features statistics and insights from multiple international organizations and research institutes, as well as the key findings from over a dozen workshops.

This white paper has five chapters. Chapter 1 lays out the general vision for and macroscopic driving forces of AI, and discusses how AI can trigger transformations of a magnitude unseen in a century and reshape the development paths of industries. Chapter 2 analyzes the certainties and uncertainties faced by enterprises during AI development, and outlines a four-pronged

approach for enterprise AI implementation, encompassing scenarios, models, data, and computing power. Chapter 3 describe development trends and looks at the shift from traditional data centers to AI DCs. It defines five features: SysMoore, prevention of energy infrastructure from being a bottleneck, iterative platforms, orchestrated applications, and generative security. Based on these five features, in Chapter 4, the white paper explores the key directions of planning and construction for ultra-large, large, and small Al DCs, and proposes a comprehensive set of evaluation indicators. This can serve as an important reference for efficient and high-quality AI DC construction. Finally, in Chapter 5, the white paper suggests that, decision makers looking to invest in AI DCs should prioritize a green, intensive, open, and collaborative approach, while also making sure that the data center is able to meet future needs. Decision makers should also focus on the construction of the three foundational layers: solutions, ecosystems, and talent.

We hope that this white paper can help industries assuredly continue on along the path of digital and intelligent transformation. By building powerful and solid AI DCs, we will be better equipped to enable industries and move towards a better, smarter future.





Trademark Notice

, HUAWEI, and Ware trademarks or registered trademarks of Huawei Technologies Co., Ltd.
Other trademarks, product, service and company names mentioned are the property of their respective owners.

General Disclaimer

The information in this document may contain predictive statements, including but not limited to, statements regarding future financial results, operating results, product portfolios, and new technologies. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purposes only, and constitutes neither an offer nor a commitment. Huawei may change the information at any time without notice, and is not responsible for any liabilities arising from your use of any of the information provided herein.

$Copyright @2024 \ Huawei \ Technologies \ Co., \ Ltd. \ All \ rights \ reserved.$

No part of this document may be reproduced or transmitted in any form or by any means without prior written consent of Huawei Technologies Co., Ltd.