



# Cloud data centers in the 5G era



**By Dennis Gu**  
Chief Architect of  
Huawei Cloud

**T**he 5G ecosystem chain covers cloud, pipe, device, and everything in between. Innovative and cross-generation evolutionary wireless terminals, air interfaces in base stations, and network transmission technologies are accelerating the evolution of device and pipe architectures. However, cloud data centers are the core hub of the 5G digital ecosystem and will play a pivotal role in 5G's evolution. What kind of cloud data center will be able to meet the network and service requirements of 5G?

The answer is simple: A distributed full-stack cloud data center that is open, efficient, flexible, and intelligent.

**Open:** Open architecture means that cloud services at different layers aren't locked by a single vendor. Instead, cloud data centers in the 5G era adopt mainstream open-source northbound service APIs and applications in compliance with industry standards.

The infrastructure layer adopts mature OpenStack elastic computing, storage, and network service APIs. The data layer uses APIs for data operations and queries based on big data and database standards such as Hadoop, Spark, MySQL, and Redis. The platform layer uses Kubernetes container service APIs, which are now a mainstream solution for application deployment and the microservices framework.

Huawei cloud data centers employ standard service APIs on artificial intelligence (AI) platforms such as TensorFlow and MXNet, both of which represent benchmarks in machine learning and deep learning.

**Efficient:** 5G networks require 100 times higher transmission rates and bandwidth than 4G networks. 5G networks also demand much more in terms of reliability and latency for applications such as virtual reality (VR), ultra-HD video, intelligent manufacturing, and auto-pilot

functionality. However, delivering ultra-high throughput and ultra-low latency capabilities is challenging for 5G cloud data center platforms when faced with network-intensive workloads, such as vEPC, and storage-intensive workloads such as CDN and 4K/8K Video On Demand.

The Huawei solution delivers optimized energy and cost efficiency for compute-intensive workloads such as machine learning, deep learning, and 3D rendering. Obviously, general x86 architectures no longer meet these requirements. Data center vendors need to introduce heterogeneous computing architectures such as ARM, SoC-based intelligent network adapter, GPU/FPGA, and neural-network processing unit (NPU) chips, to ensure businesses can run at the highest energy efficiency ratio and cost-effectiveness.

Network elements (NEs) on the data plane of 5G networks, such as virtualized evolved packet core (vEPC) and CloudRAN base stations, normally require the throughput of a single NE or single server to jump from 10 Gbps to 100 Gbps. But, the conventional, software-only overlay cloud networks based on x86 CPUs face severe performance bottlenecks. Therefore, it's necessary to deploy SR-IOV direct pass-through or the DPDK user space mechanism to offload overlay networking functions, such as OVS, vRouter, and vFW, from

X86 to heterogeneous hardware, for example, Intelligent NIC or FPGAs with 100 Gbps port speed.

For 5G Internet of Things (IoT) scenarios, massive volumes of data collected from distributed IoT terminals and edge nodes in various industries need to be aggregated, stored, processed, analyzed, and managed in a centralized cloud data center based on the object storage service and Hadoop/Spark/Stream Big Data Pipeline services.

Huge amounts of machine learning and deep learning processing tasks require extensive resources for the following parallel computing operations: convolution, derivatives, logarithms, and matrix multiplication with floating numbers. These operations extract valuable business insights and train prediction models based on inputted raw data.

Offloading these computing capabilities to heterogeneous hardware, such as GPU/FPGA clusters or even the NPUs, improves the cost effectiveness of the high-density computing and energy efficiency ratio by 5 to 10 times. More customers require data interconnections between heterogeneous computing GPUs, FPGAs, and NPUs in massive parallel computing scenarios, which has prompted huge advancements in ultra-high-speed link connection technologies such as RDMA over converged Ethernet and NVLink.

This ensures that the performance advantages of heterogeneous computing clusters are fully brought into play.

For ultra-high performance IOPS in storage-intensive scenarios, the next-generation storage-class memory (SCM) is increasingly used as the default flash storage medium with an acceptable unit capacity price, while providing read and write speed and latency comparable to memory solutions. SCM-based storage nodes with RDMA/RoCE link connection to compute nodes are now more widely used in the distributed storage architecture, typically delivering sub-100 ms shared storage latency and storage bandwidth that's even higher than the local disk's PCI speed.

**Flexible:** Another challenge for 5G networks is to flexibly orchestrate and reassemble network slices. Network slicing requires the evolution of NEs for 5G services and protocols, as well as streamlining on the 5G IoT application data layer, core network layer, and wireless access layer of a cloud data center. These advancements will help implement end-to-end network functions on the management, control, and forwarding layers, as well as support dynamic on-demand isolation of capacity and QoS.

Based on the specific requirements of vertical service scenarios, flexible 5G network slicing requires that

*Dynamic orchestration helps simplify and shorten 5G network construction from several months to one-click, automated construction in minutes.*

deployment, including capacity, service configuration, network elements, service applications, and the in-between networking link, is completed as quickly as possible. Template-based orchestration services are introduced here to enable 5G networking elements and applications. Dependency can be automatically provisioned and configured based on the heterogeneous resource flavors of virtual machines, and physical machines with pre-defined topology dependency. In addition to static resource topology, PaaS requires dynamic orchestration capabilities such as sequential, conditional, and loop service logic control for orchestrated services with ensured transaction integrity. Dynamic orchestration helps simplify and shorten 5G network construction from several months to one-click, automated construction in minutes.

**Distributed:** In the 5G era, the majority of physical network access and routing network functions and applications, such as IoT core services, big data, and deep learning, will be deployed in centralized large-scale cloud data centers based on geo-redundant VMs. This will enable 5G IoT device access and corresponding application platforms, as well as diverse third-party applications located in distributed sites. The data centers in this layout can support tens of thousands of hosts. The access devices on the 5G data plane, such as vEPC gateways, are generally deployed near the metro aggregation access Points of Presence (POPs). This ensures that cloud services, like backup, disaster recovery (DR), video storage uploads, and other typical low-latency interactive services, can be accessed with ensured QoS/SLA over non-blocking dark-fiber/MPLS bandwidth. These services can run on hundreds of small-

scale satellite cloud sites configured as a one-stop cloud in box mode.

The ubiquitous access and coverage of 5G networks require ultra-low latency and ultra-high bandwidth. Functions that should be moved to the satellite cloud include radio air interface protocol processing; baseband control; wireless resource management and scheduling; network data tunnels; route forwarding, aggregation, and processing; and service processing.

In synergy with the centralized cloud region and in addition to the satellite cloud located around local city PoP, we still need large numbers of edge nodes beneath the satellite cloud, which is designed to better support IoT services. For example, predefined AI-enabled recognition of surveillance videos, stream processing filtering of raw IoT data, 3D content rendering of VR/AR games, and real-time user operations, can be migrated from the centralized data center to the access network edge. By doing so, centralized intelligent analysis, agile development iteration, and local and real-time access processing capabilities can complement each other.

Moreover, edge nodes can convert a large amount of high-bandwidth, unstructured, or multimedia data on the terminal side to high-value and low-bandwidth structured data that can be uploaded to the cloud data center for centralized analysis and processing. The edge nodes deliver control commands for the edge and terminals. This improves the overall throughput of E2E 5G networks and cloud service applications, which in turn will improve user experience. The seamless integration of edge computing services with full-stack cloud services, such as the cloud PaaS platform, big

data, and AI services, accelerates the development and rollout of innovative services such as IoT, video AI, and AR/VR games with ubiquitous access.

Typical reference architecture for edge clouds can be implemented by the reference architecture of centralized Kubernetes master nodes plus the remote minion node with Kublet agent. With northbound K8S APIs, the edge cloud is compatible with the edge computing services in the Kubernetes ecosystem. It supports access registration and security certificate management for tens of thousands of distributed edge nodes. Edge nodes also need to support parallel batch container instances, serverless instance deployment, and lifecycle management

**Intelligent:** Enabling IoT applications is one of the main objectives of 5G network construction. IoT applications generate massive amounts of data, so data centers in the 5G era need to provide ultra-large elastic storage capacity and computing capabilities. In addition, an intelligent, highly efficient engine that's easy to configure and use, combined with diverse domain knowledge and data models, is needed to quickly learn and extract targeted valuable information and strategies from this data. This facilitates closed-loop control on IoT terminals and edge devices.

Typical application scenarios include:

- Image and video recognition in wireless video surveillance scenarios
- Vehicle GPS location tracking and driving behavior preference analysis for the Internet of Vehicles
- Traffic congestion and violation detection in intelligent traffic scenarios
- Population density and mobility prediction in smart city scenarios
- Power use distribution and peak predictions in smart grid scenarios

The intelligent engine previously introduced in the IoT scenario requires that the cloud data center relies on IoT data lakes on the big data platform. This enables the cloud data center to provide rich platform services and APIs with pre-integrated machine learning, deep learning, a graphics engine, search capabilities, AI services, and APIs in common fields such as visual, voice, natural language, and optical character recognition (OCR). These platforms and general AI/machine learning service capabilities work closely with heterogeneous computing hardware, including GPUs, FPGAs, and NPU's and the scheduling system's 5G cloud platforms, to implement deep software optimization.

As 5G cloud data centers will be deployed over multiple distributed geo-locations with support for multi-

tenant network slices and millions of resource nodes, features such as maximum cloud region size must be supported and a truly intelligent and self-healing maintenance mechanism is urgently required.

Moreover, traditional local O&M and fault management needs to be replaced by proactive, predictive O&M and management based on powerful big data and AI services as part of a 5G full stack cloud. AI/ML algorithms deployed on the basic platform with supervised learning, semi-supervised learning, and unsupervised learning can help analyze the massive amounts of log information collected from software and hardware subsystems. These algorithms support the root cause analysis of failures, automatic identification of abnormal behavior patterns, and the prediction of network and hard disk faults, greatly improving hardware and software O&M efficiency, with a single O&M staff member able to maintain more than 1,000 servers.

Both equipment providers and carriers have great expectations for 5G. In some countries, 5G commercial deployment may be accelerated by government policy. So, carriers must plan ahead to build an open, efficient, flexible, and intelligent distributed full-stack cloud data center that ensures the openness and flatness of the 5G network. Then, 5G industrial applications can be quickly commercialized. 