



中国移动
China Mobile

Model Distribution Network Considerations in China Mobile

March 2025

www.10086.cn

DeepSeek Opens the Era of Inclusive AI with Explosive Growth in Inference Services

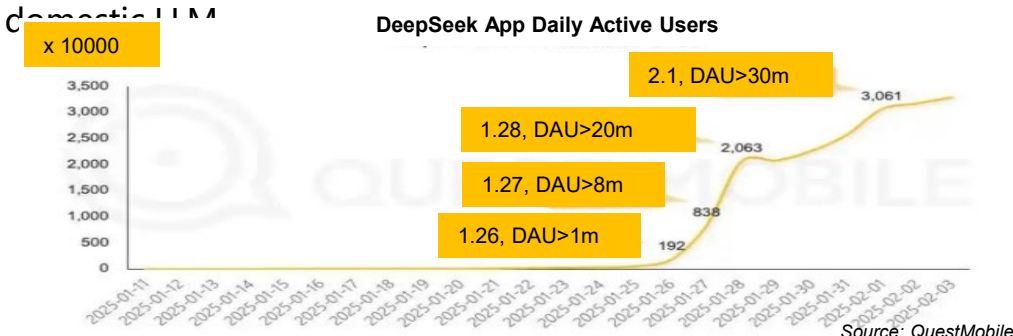
DeepSeek-R1 achieves accuracy comparable to OpenAI's o1 model with 70% lower computation requirements, becoming a global focus.

Market Success : Surpassed ChatGPT in downloads across **140** app markets; **fastest** application to exceed **30 million** daily active

Trend1: Paradigm shift from content to AI model delivery

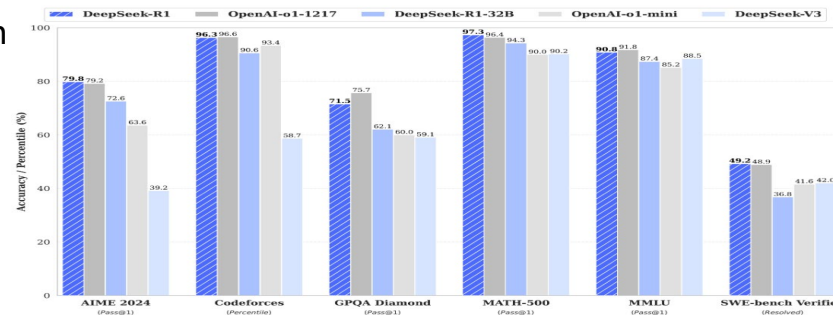
- **DeepSeek Integration**: Over 160 enterprises globally have announced integration with DeepSeek, covering cloud service providers, telecom operators, cybersecurity, automotive, smart hardware, finance, and semiconductor manufacturing.

- **DAU Record**: On Jan. 26, The number of daily active users exceeds 1 million. Two days later the number gets 20 million. Then after another four days it is 30 million, becoming the most popular domestic LLM



Trend 2: Focus on the "Inference Moment"

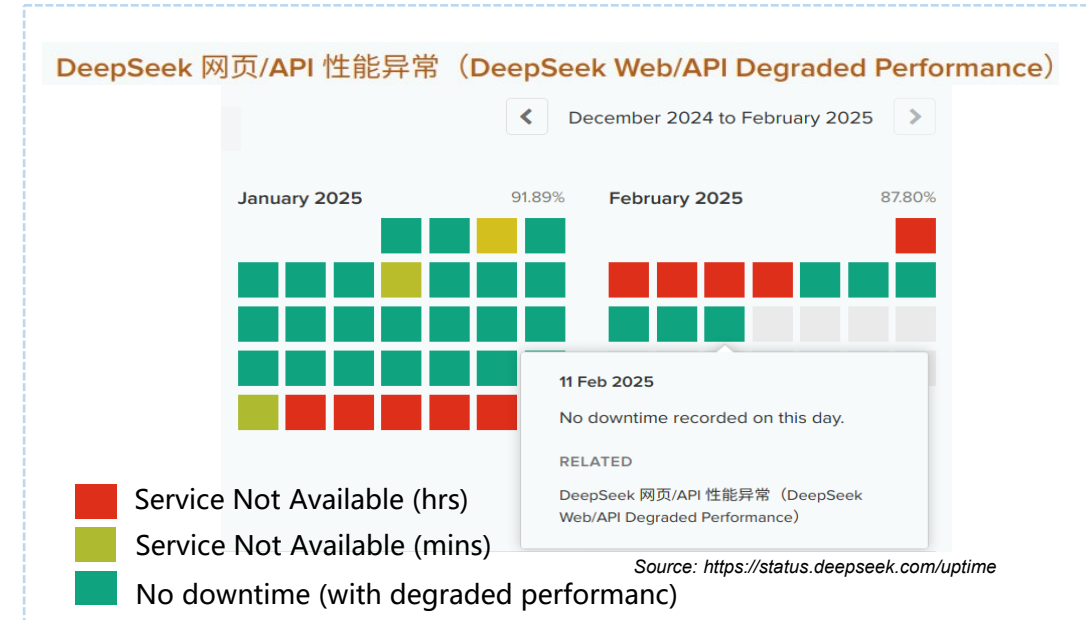
- **Performance**: **Leading industry performance** in mathematics(AIME, MATH-500), knowledge, reasoning(GPQA), and programming tasks(Codeforces)
- **Inference Speed**: Increased by **4** times, with API call costs at nearly **3%** of GPT-4-Turbo()
- **Model Distillation**: According to Ollama assessment data, **distilled smaller models (e.g., 32B)** perform very well in benchmark tests
- **Industry Requirement**: Vertical industries have a significant requirem



AI models are evolving from chat tools to production and daily life tools, forming irreversible new business scenarios, bringing huge opportunities for model distribution and edge inference

DeepSeek has experienced multiple outages due to surging access volumes, frequently becoming unavailable. This highlights the inability of the current deployment architecture to handle massive concurrent user access. Operators need to leverage network resources, technical capabilities, and market advantages to seize opportunities in the AI era

- **Phenomenon:** "The server is busy. Please try again later" issue remain unresolved
 - Access issues began as early as late January
 - On February 6, DeepSeek had to suspend API service recharges
- **Cause:** **centralized deployment** architecture cannot quickly adapt to the explosive growth of massive concurrent user requirements
 - **User Scale:** Over 30 million DAUs on February 1
 - **API Integration:** Major applications like Alibaba Dingding, Huawei Xiaoyi, Honor, and WPS have integrated DeepSeek's API
 - **Deployment Method:** Centralized deployment with a single inference pool in China (*single-point bottleneck*)



Comparison with Gemini: Gemini, with a similar DAU scale, uses a distributed architecture, allowing users to access the nearest nodes with much less network latency, without access issues.

- **Carrier Network Opportunity:** research a new **distributed inference service** architecture to handle the challenges of hundreds of millions of users and tens of millions of concurrent accesses for the AI inference age
 - **Internet era:** CDN distributed web/video content from centralized pools to edge nodes, solving congestion and large-scale user access
 - **Inclusive AI era:** MDN (**M**odel **D**istribution **N**etwork) uses a *distributed* inference architecture to provide ubiquitous low-latency inference services

Deepseek Inference Server: Centralized Deployment

chat.deepseek.com



60.204.2.9 (In China) Shanghai·Huawei Cloud
104.18.27.90 (Oversea) USA/CA·CloudFlare
104.18.26.90 (Oversea) USA/CA·CloudFlare

Only **3** IP addresses over the world for
chat.deepseek.com and *api.deepseek.com*

Location

Fastest

Slowest

Avg

①

上海电信 7ms

福建厦门移动 44ms

22ms

②

广东广州电信 30ms

海南海口联通 57ms

43ms

③

湖北武汉电信 21ms

湖南益阳联通 38ms

29ms

④

天津联通 25ms

内蒙古电信 38ms

29ms

⑤

贵州贵阳电信 35ms

西藏拉萨移动 74ms

49ms

⑥

陕西西安电信 25ms

新疆乌鲁木齐移动 71ms

46ms

⑦

吉林长春移动 37ms

吉林长春电信 51ms

42ms

**the fastest pings
are over 20ms**

Doubao LLM Inference Server: Distrubution Deployment

Doubao: A top popular LLM in China and made by Bytedance



Location

Fastest

Slowest

①

山东济南联通 <1ms

浙江宁波电信 16ms

6ms

②

海南海口电信 <1ms

广东潮州联通 29ms

7ms

③

河南郑州移动 <1ms

河南郑州电信 16ms

4ms

④

山西太原联通 <1ms

内蒙古移动 29ms

10ms

⑤

云南昆明电信 <1ms

西藏拉萨电信 36ms

9ms

⑥

甘肃兰州电信 <1ms

新疆昌吉电信 23ms

3ms

⑦

吉林长春电信 <1ms

黑龙江哈尔滨移动 3ms

1ms

the fastest pings
are all less 1ms

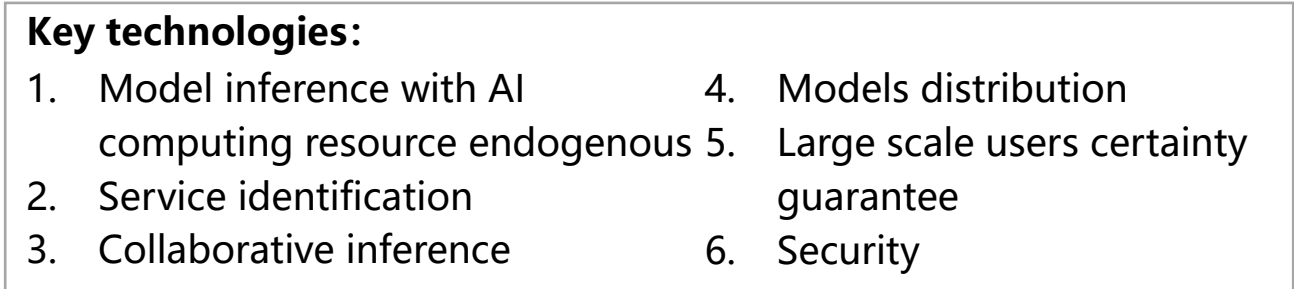
Average pings are
less than 10ms

More than 100+ IP addresses in China. So it can use near inference servers to distribute the inference requests

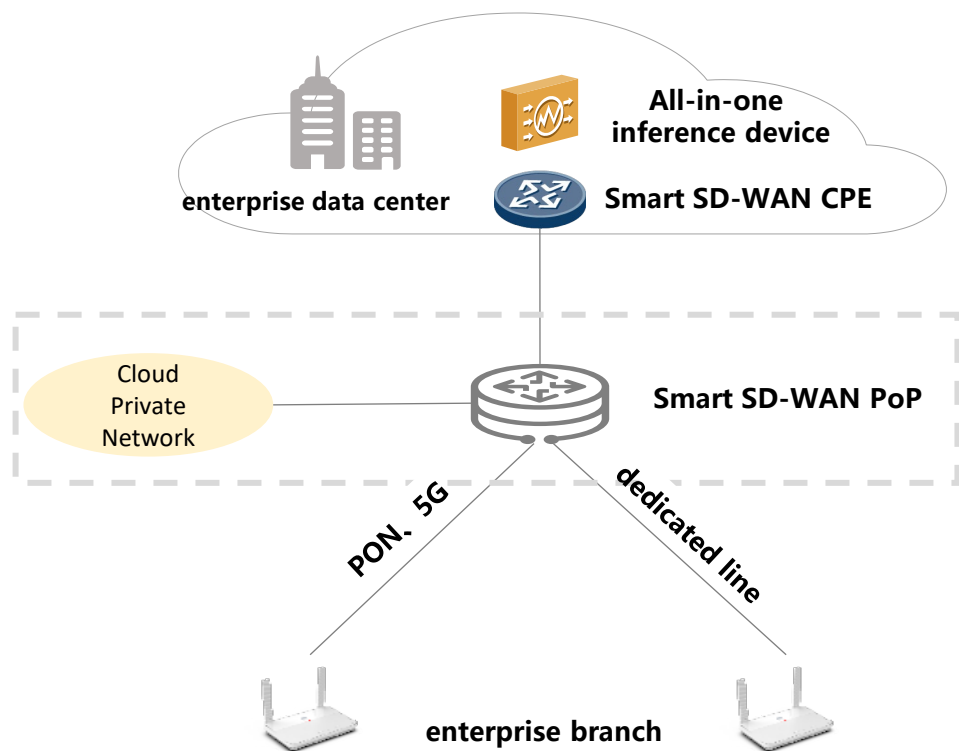
“Server busy” issue is never happened in
Doubao

Source: <https://www.itdog.cn/ping>

1. **Model Repository:** Stores and manages different LLM versions
2. **MDN Center Nodes:** Deployed centrally to obtain models from the repository and distribute them to edge nodes
3. **MDN Edge Nodes:** Deployed in various geographical locations to receive and deploy models from center nodes, providing local edge inference services
4. **User Side Client:** End-user devices/apps launch inference requests
5. **MDN Scheduler:** Dynamically adjusts inference task distribution based on network conditions and device load

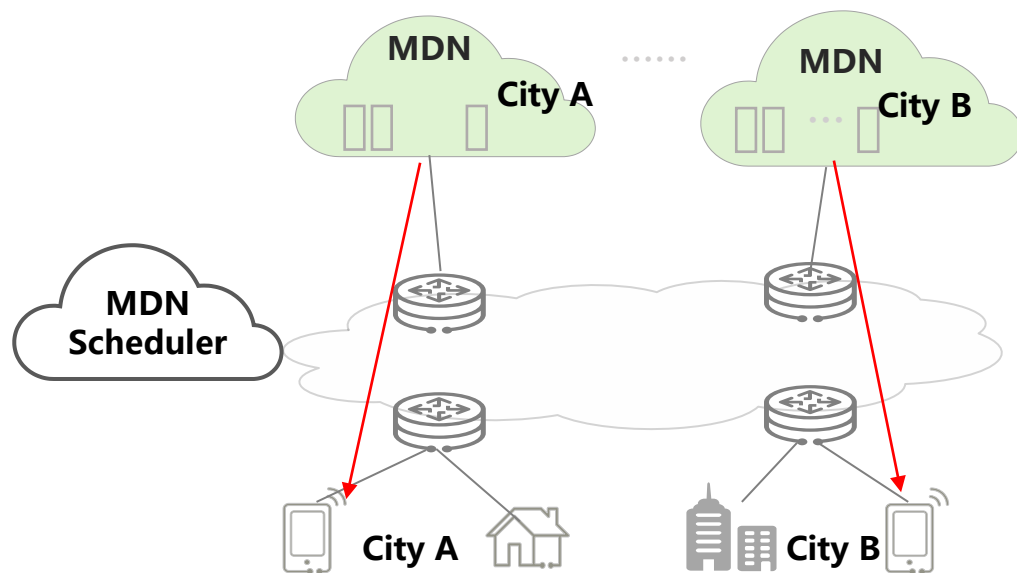


Scenario 1: 2B Standalone Deployment Integrated Secure Inference Service with Computing and Network



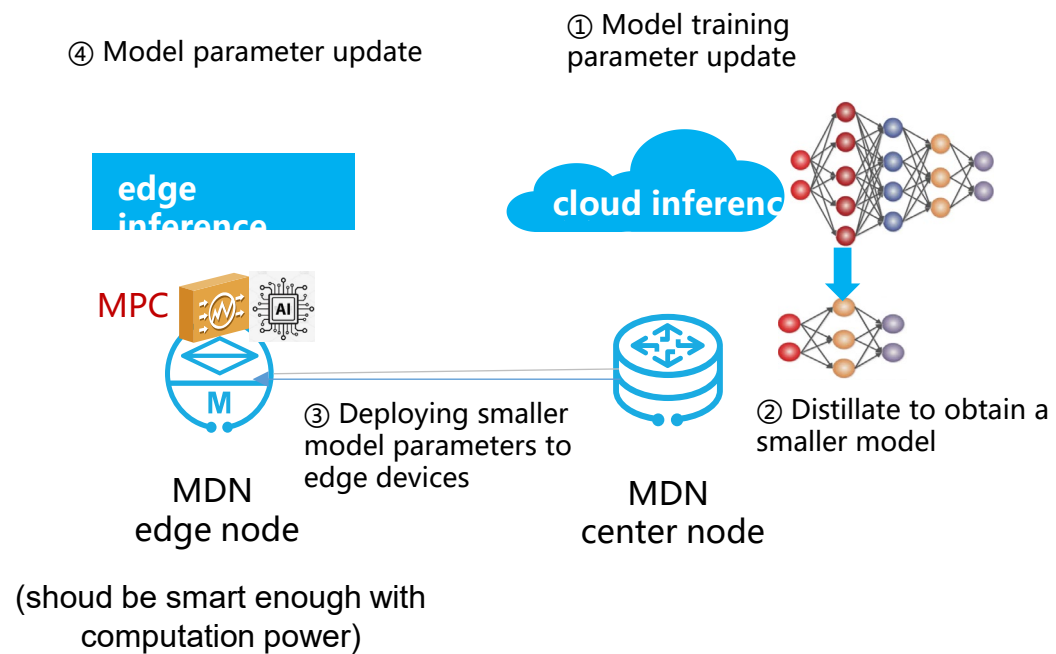
- For high-security 2B sectors, deploying an all-in-one inference device within the enterprise will protect the data.
- It is needed to protect the network accessing for enterprise branch accessing the private inference server from public internet or dedicated lines.
 - **"All-in-One Inference Server + Smart SD-WAN"** can provide **secure inference** service and **secure accessing**

Scenario 2: 2C&2H Dynamic Delivery Balanced Inference Scheduling Service



- Deploy the full size model (e.g., DeepSeek-R1 671B) in distributed MDN inference centers
- Use the **MDN scheduler** (DNS is a option) to **distribute inference tasks**
- **Dynamic Delivery**: the scheduler needs be aware inference resource load and user proximity for *multi-weighted factors* dynamic load balancing

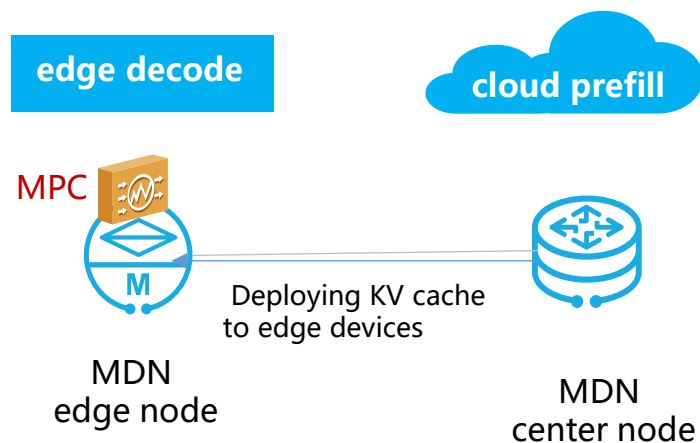
Scenario 3: Collaborative Inference (Full size and Distilled) Collaborative Inference Service for Large and Small Models



- Deploy the full size model (e.g., DeepSeek-R1 671B) in the inference service center
- Deploy distilled small models (e.g., 7B/70B) on edge nodes.
 - edge nodes can **schedule inference tasks to either the small models in edge or the full size model in the cloud** to reduce total computing resources
 - Edge node must support dynamic scheduling based on task type, complexity, and user priorities. So the edge node is not only a packet forwarding router (BNG) but also with computation and storage (*intelligent BNG: iBNG*)
 - Timely package the edge node inference tasks and results to center node for checking and updating.

Scenario 4: Collaborative Inference (Disaggregating Prefill and Decode)

Collaborative Inference Service for Prefill stage and Decode stage



- Deploy the full size model (e.g., DeepSeek-R1 671B) in the inference service center
 - Prefill phase (which is comput-intensive task) is executed on center node.
 - The center node focuses on generating the precomputed KV Cache, and then sends to MDN edge.
- Edge node leverages the precomputation results to perform decoding operations
 - Decode phase (which is memory-intensive task) is executed on edge node
- Disaggregating P/D **reduces** total computing resources, power consumption and network latency *without* loss any accuracy.
 - The network between center and edge should *zero packet loss, large bandwidth and latency guaranteed* (GP-C, Slicing, DotNet, etc.)

Thank You for Listening!

Corrections and Suggestions Welcome!