アプサラ・カンファレンス2025 最新情報

Jack Wang

Solution Architect
Alibaba Cloud Intelligence Group



王志正 (Jack Wang) Wang Zhizheng

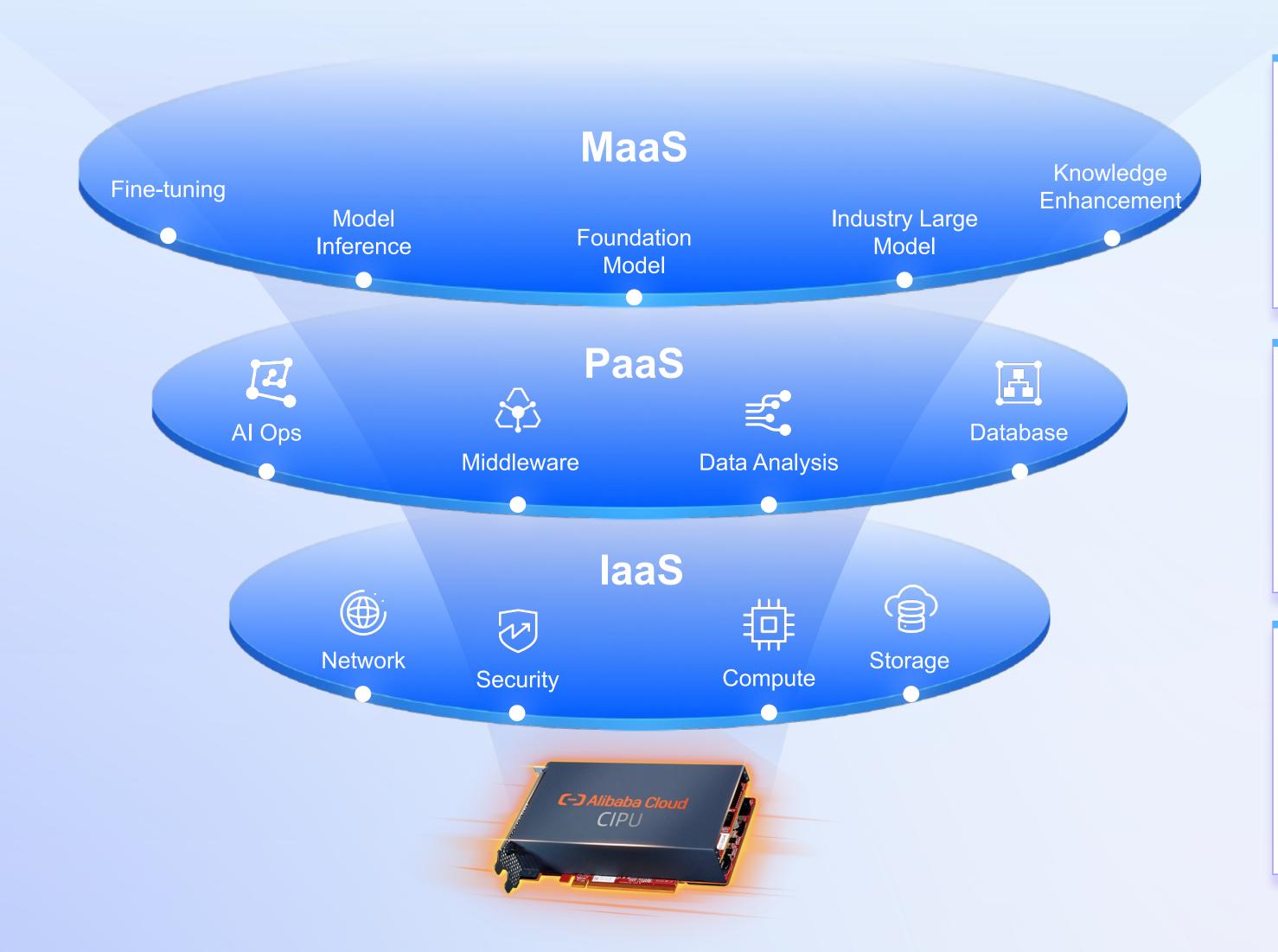
キャリア経歴:

- 2024年4月アリババクラウドジャパンに入社
- 前職 Datadog Japan、SBCloud、EMC
- 現職 アリババクラウド ソリューションアーキテクト





Full-stack AI としてのクラウドサービスはグローバル リーダーの位置付けになっています



費用対効果が高く、信頼性 の高いエンドツーエンドの ソリューション

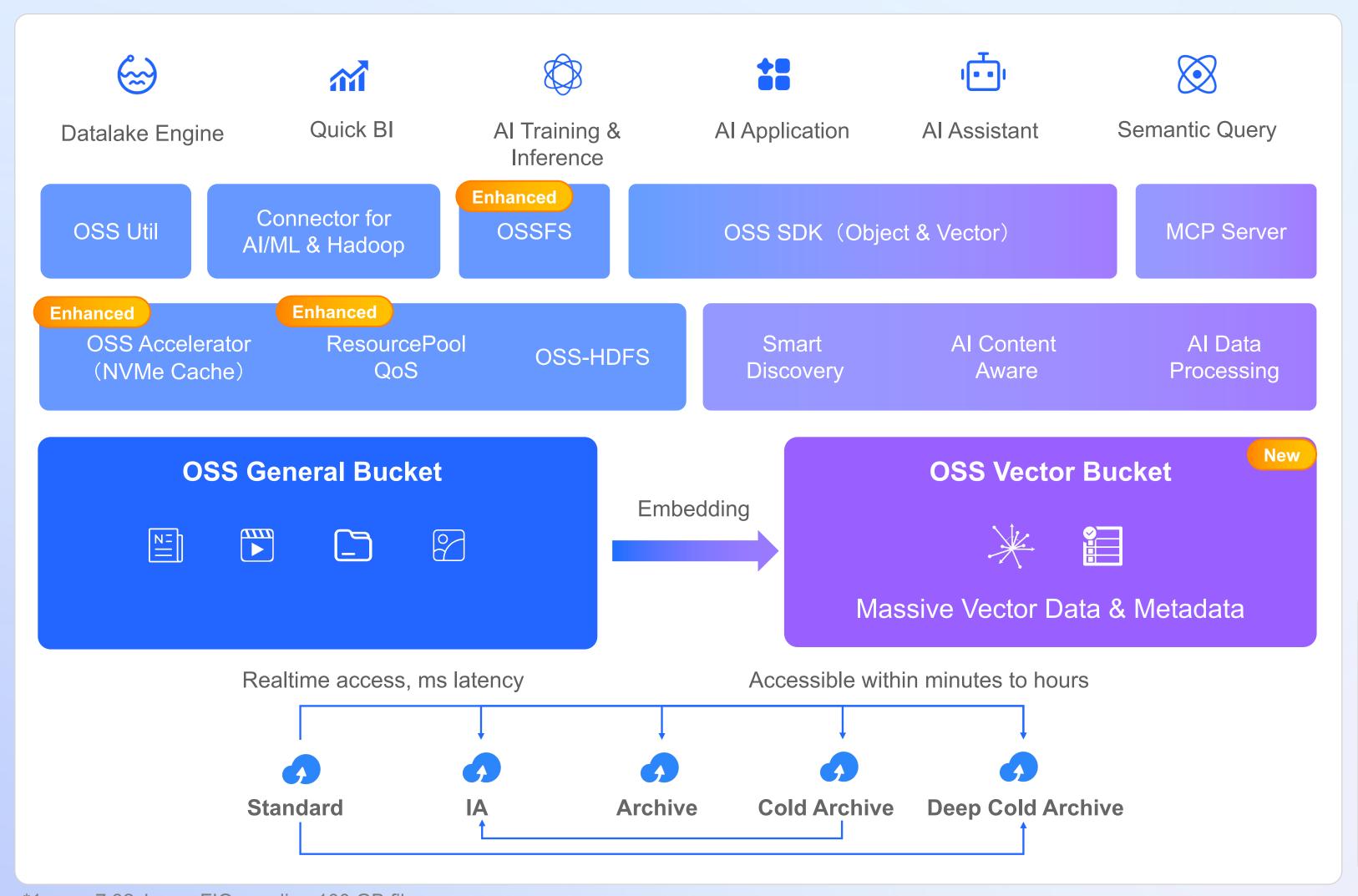
水平および垂直 製品ポートフォリオの統合

自社開発のソフトウェアとハー ドウェアによる技術リーダー シップ





AIワークロード向けの大規模オブジェクトストレージ



組み込みのVector Bucketにより 従来のVector DBのコストを大幅に 削減します。

パフォーマンスQoS – 優先制御と最 小スループットの保証

OSSアクセラレータ – 最大100GB/秒 OSSFS 2.0 – 新しいシンプルな FUSEクライアント

7.65 x

Single Thread Reading(*1)

Model Loading Speed (*2)

^{*1} ecs.g7.32xlarge, FIO, reading 100 GB file

^{*2} ecs.8gi.48xlarge, llama, loading DeepSeek-R1-GGUF model



コンテンツ認識ストレージ

OSS:データのコンテキストを理解するAIストレージ

RAGアプリ

セマンティック検索

AIエージェント

AIデータ管理





複雑な操作なしでワン クリックで機能を起動



グローバル展開

迅速なグローバル展開 を可能にする包括的な SDK ツール



マルチディスカバ リーモード

既存のコンテンツを検 出してクエリし、新し いデータに自動的にタ グを付ける

既存コンテンツの検出

セマンティッククエリ



キーワードを生成する



"cat", "lion", "feline"

新規データに自動タグ を付ける

新しいコンテンツの発見



OSS組み込みメタデータとベクターインデックス

New "kitten" uploaded



信頼できる結果

85%の精度を誇る成熟したアルゴリズム技術



コスト効率が高い

初期設定コストが低く、 従量課金制

OSS Raw Data



One-click activation "Al Content Aware"



- UserMeta, Vector Data,
- Content Description...



AIワークロード向けALB: フルシナリオロードバランサーの構築



GA サポート エニーキャスト IP

Enhanced

• 高可用性ノード

- 最適なルーティング
- 最寄りのPoP経由の低遅延アクセス

近接アクセスのためのエニーキャストを備えた GA

セキュリティが組み込まれた統合エントリポイント

Enhanced

- 統合アクセスアーキテクチャ
- 集中管理
- WAFは悪意のあるトラフィックをブロックします

GA WAF

LLMインテリジェントプロキシとスケジューリング

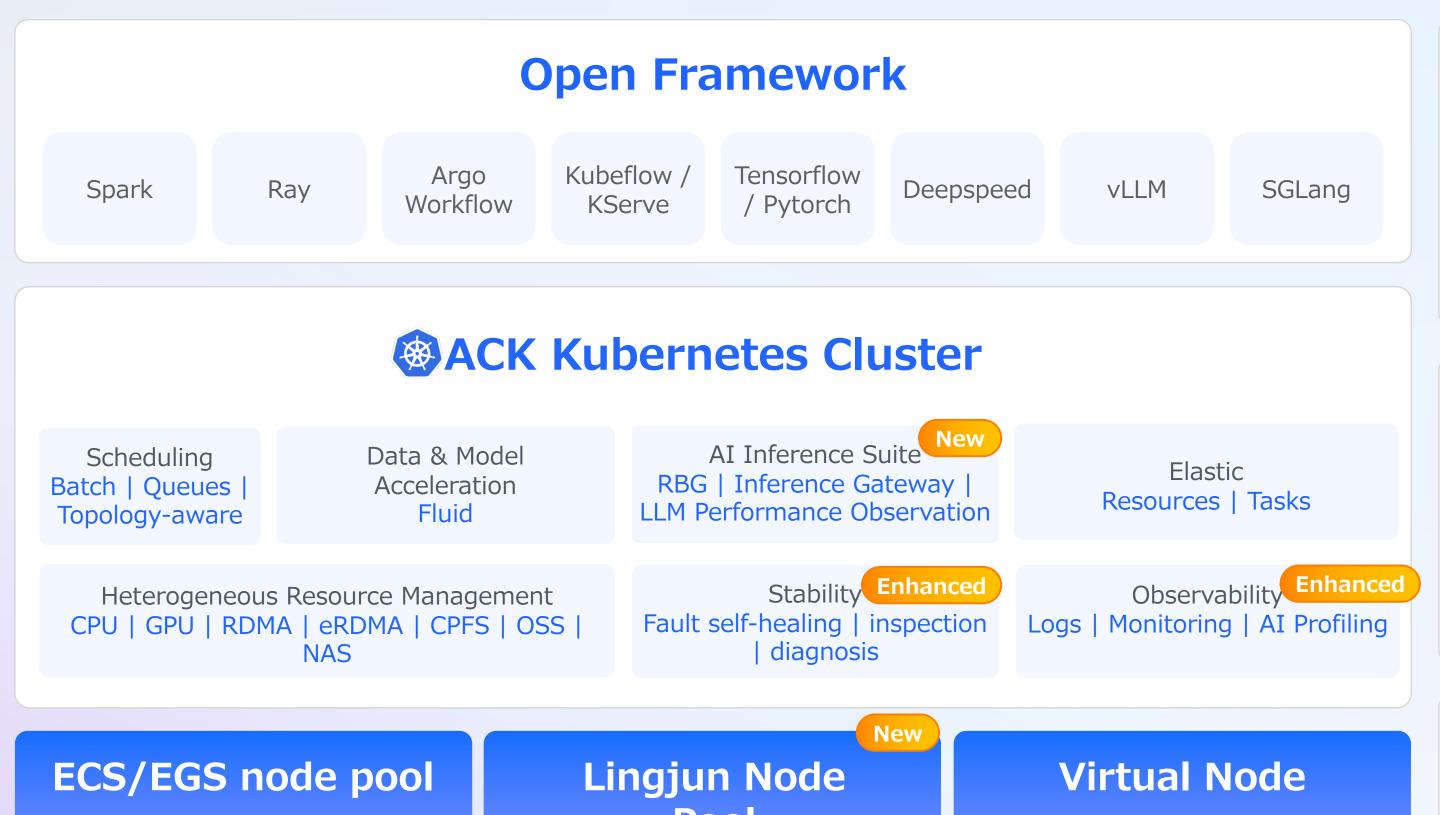
- 異なるAIバックエンド統合
- 最適なモデル選択のためのスケジューリング
- 認証+レート制限によるコスト管理
- 完全な可観測性

ALB for AI workloads

New



ACK Kubernetes: AIコンピューティングの効率的で安定した基盤



AIフレームワークのオープンエコシステム

最適化されたスケジューリング | 柔軟性 可観測性 | セキュリティ

AI Suite: 効率的なトレーニングと推論

- LLM推論スループットが2倍に向上
- モデルのコールドスタート時間が90%短縮

General Computing Instance

Accelerated Computing Instance

Pool

High-performance **GPU Bare Metal cluster** ACS Pod

ACS Pod

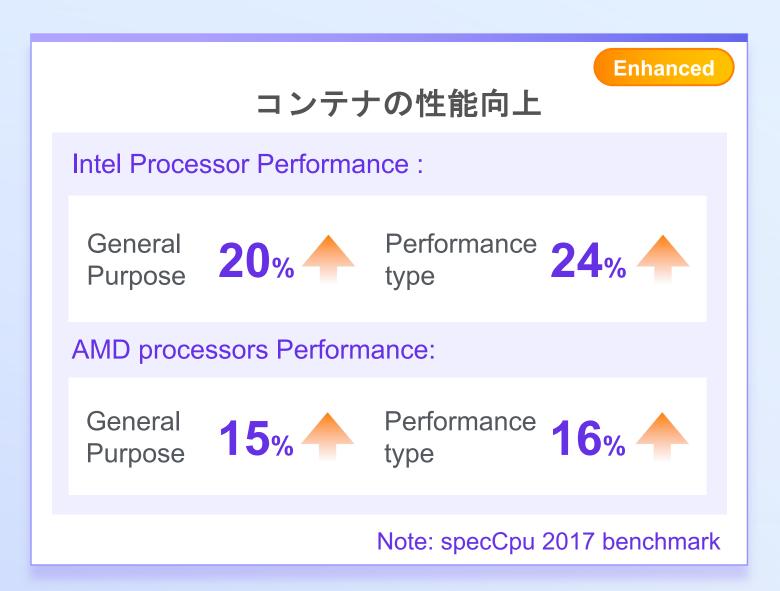
ACS Pod

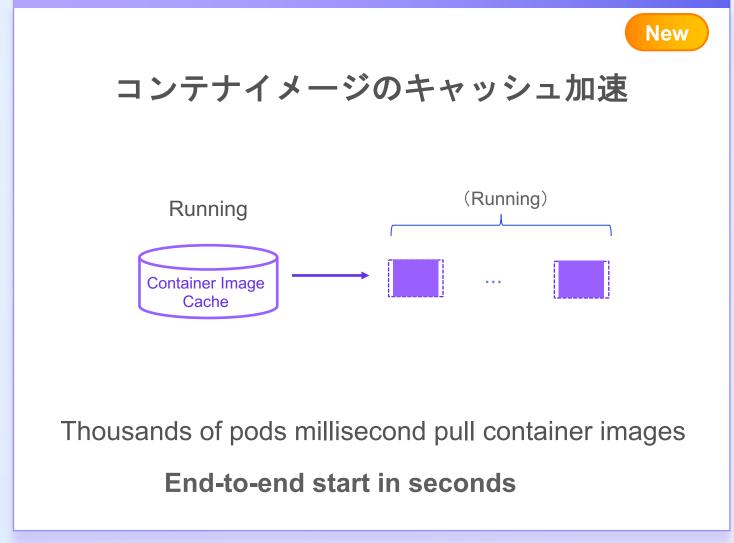
インテリジェントなO&M

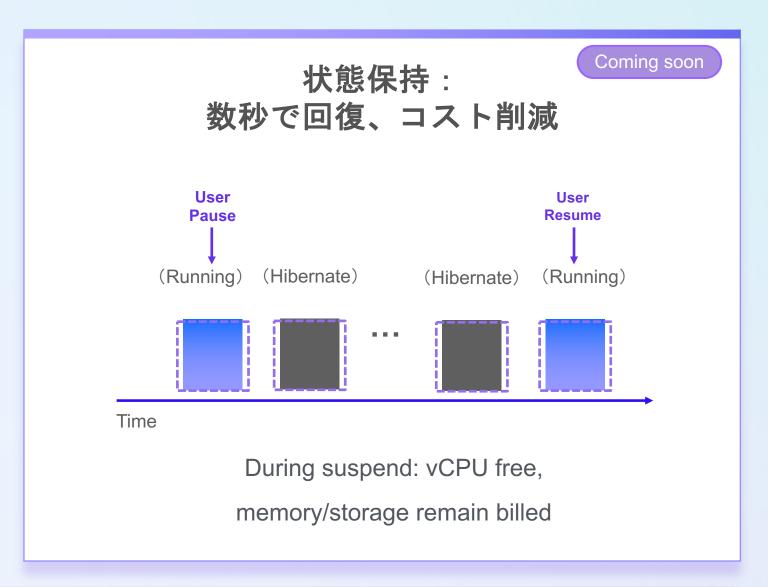
自己修復: GPU効率が85%向上

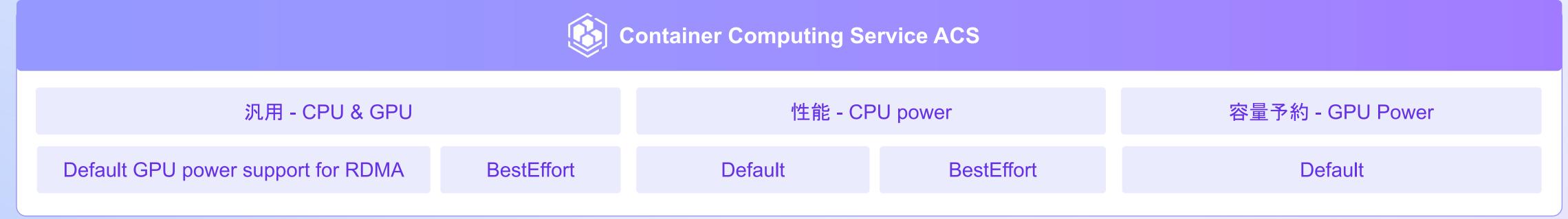


ACS コンテナ: Smarter, Simpler, Stronger









高性能 Container performance ahead of competition 極めて高い弾力性

15,000 pods/min startup per region

利便性

12-month GA: 10M+ Pods/day on average



PAI-EAS:推論のためのコスト効率の高いサービス



より迅速な展開

89.8% Cold Start 97.6% Scale out

低レイテンシー

20.3% TTFT 70.6% TPOT

より高いスループット

71.0% TPS

モデル展開レイヤー

負荷分散

スケジュールの最適化

モデル展開の加速

推論エンジンレイヤー

大規模EP

コンピューティング通信 の重複

演算子の量子化

PD分離から AF分離へ

KVスマートキャッシュ

MTP



Al アプリケーション向けのフルスタック Observable Cloud Monitor 2.0

モデルアプリ ケーション

アプリケーション

RAG, MCP call monitoring Al application monitoring

Model call monitoring

Safety diagnosis

モデル

Model monitoring

Model Log model link

フルスタック監視

Intelligent laaS, PaaS, Application Full Linインテリ ジェントなlaaS、PaaS、アプリケーションのフルリ ンクk

AIプラット フォーム (PaaS)

モデル推論

Service and Component Monitoring

Inference Engine Monitoring

サービスコンポーネント

Vector database monitoring

AI Gateway Monitoring

モデルのトレーニング

Training task monitoring

Data Flow Monitoring

Container

Kubernetes

コンテナのスケジュール

ワークロードの監視、コントロールプレーンの監視、GPU モニタリング、 イングレス モニタリング、イベントおよびログ監査、AI スイートの監視

AIインフラストラクチャ

インテリジェ ントコン ピューティン グサーバー

(laaS)

Lingjun コンピューティング

ノードの監視

ネットワーク RDMA 監視

高速ストレージ CPFS 監視

GPU クラウドサーバー 監視

エンドツーエンドの診断

10 種類以上の主流エージェント フレーム ワークをサポート

重要なアクセスデータ収集

コンテナ、PAI、その他のプラットフォーム とデフォルトで統合

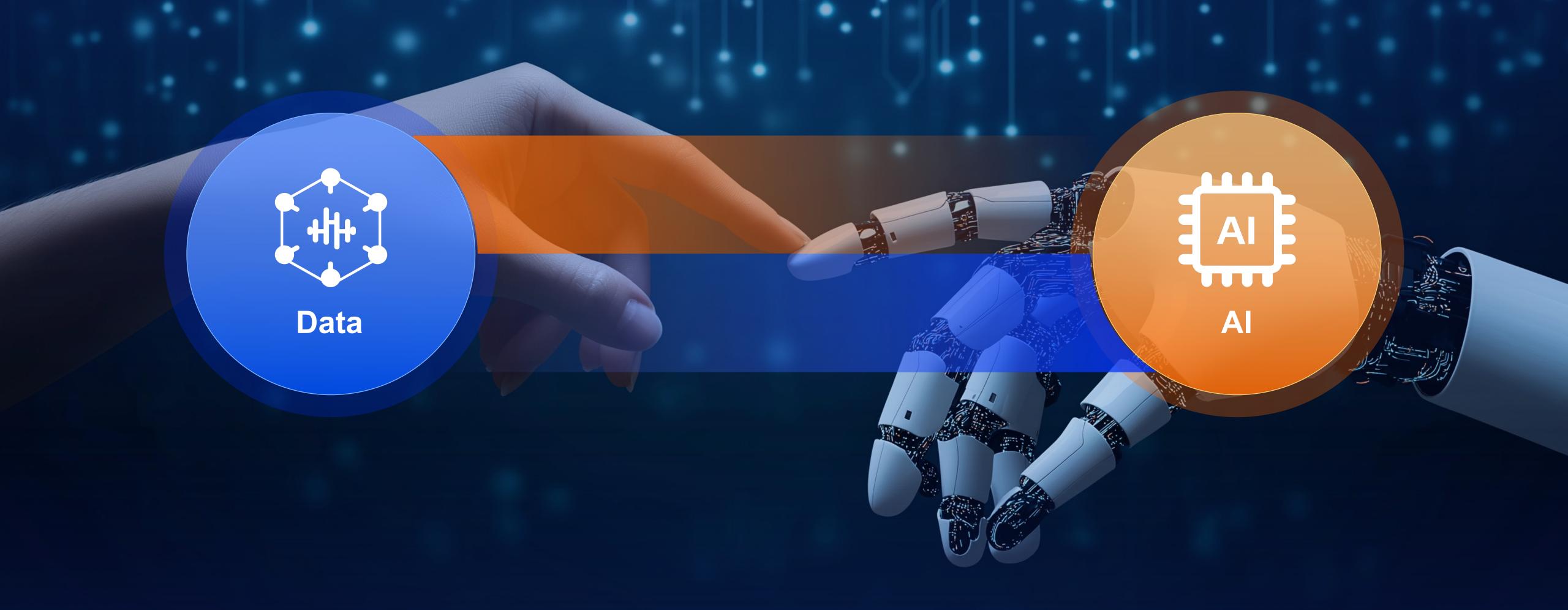
> コスト管理とモデル 品質評価

トークン、GPU使用コスト監視

モデルの入力と出力の品質評価 セマンティクス、出力品質、セキュリティリ スクをカバーする



Revolutionize Data Intelligence with Al





データ管理とデータ準備のための単一のプラットフォーム



すぐに使用できるデータ エージェントと開発プラットフォーム、並列クエリ + SIMD と自動インデックスによる高速化によるパフォーマンスの向上

統合された メタデータ

統合データガバナンス

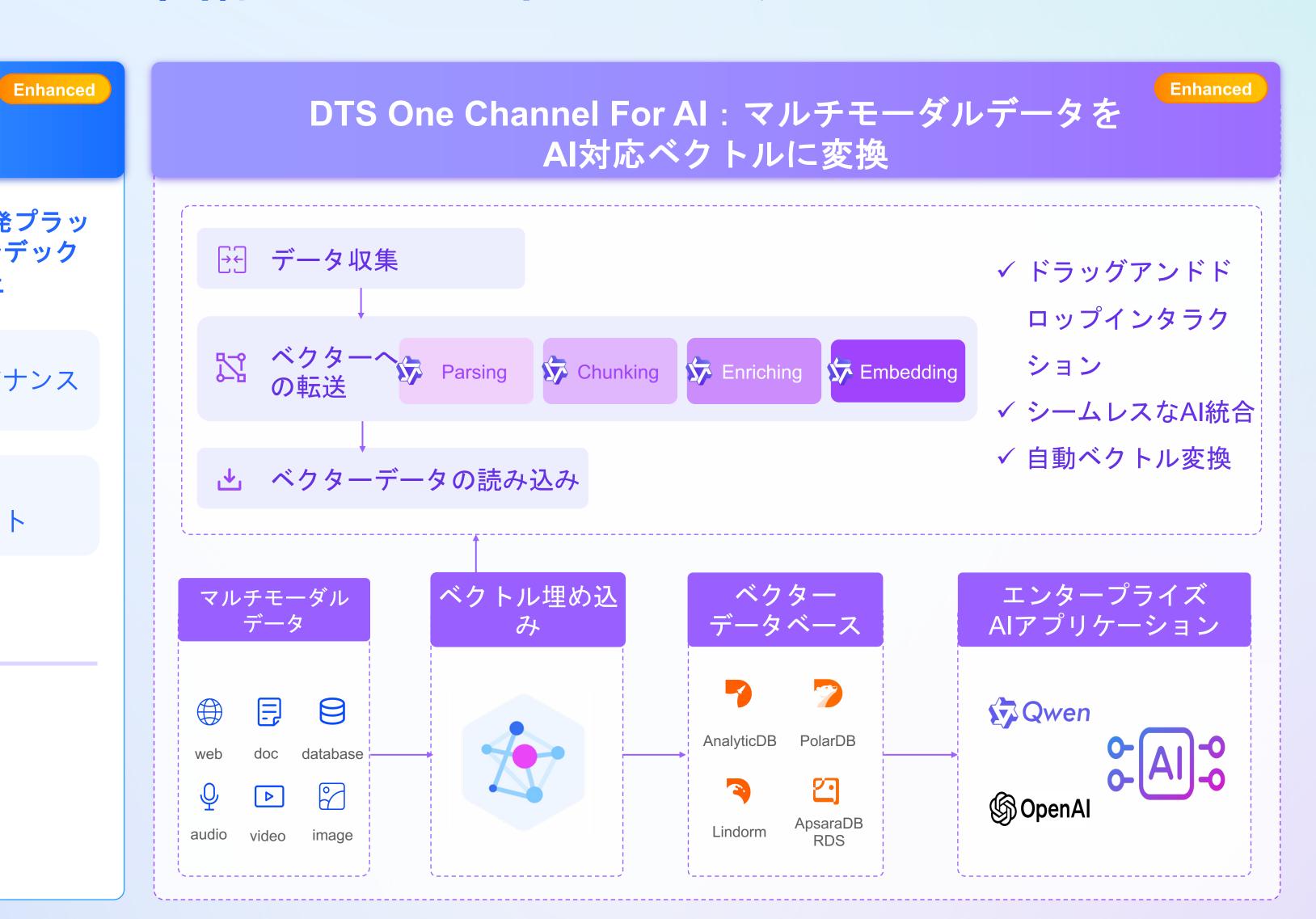
統合データ運用

データ エージェント

データ管理

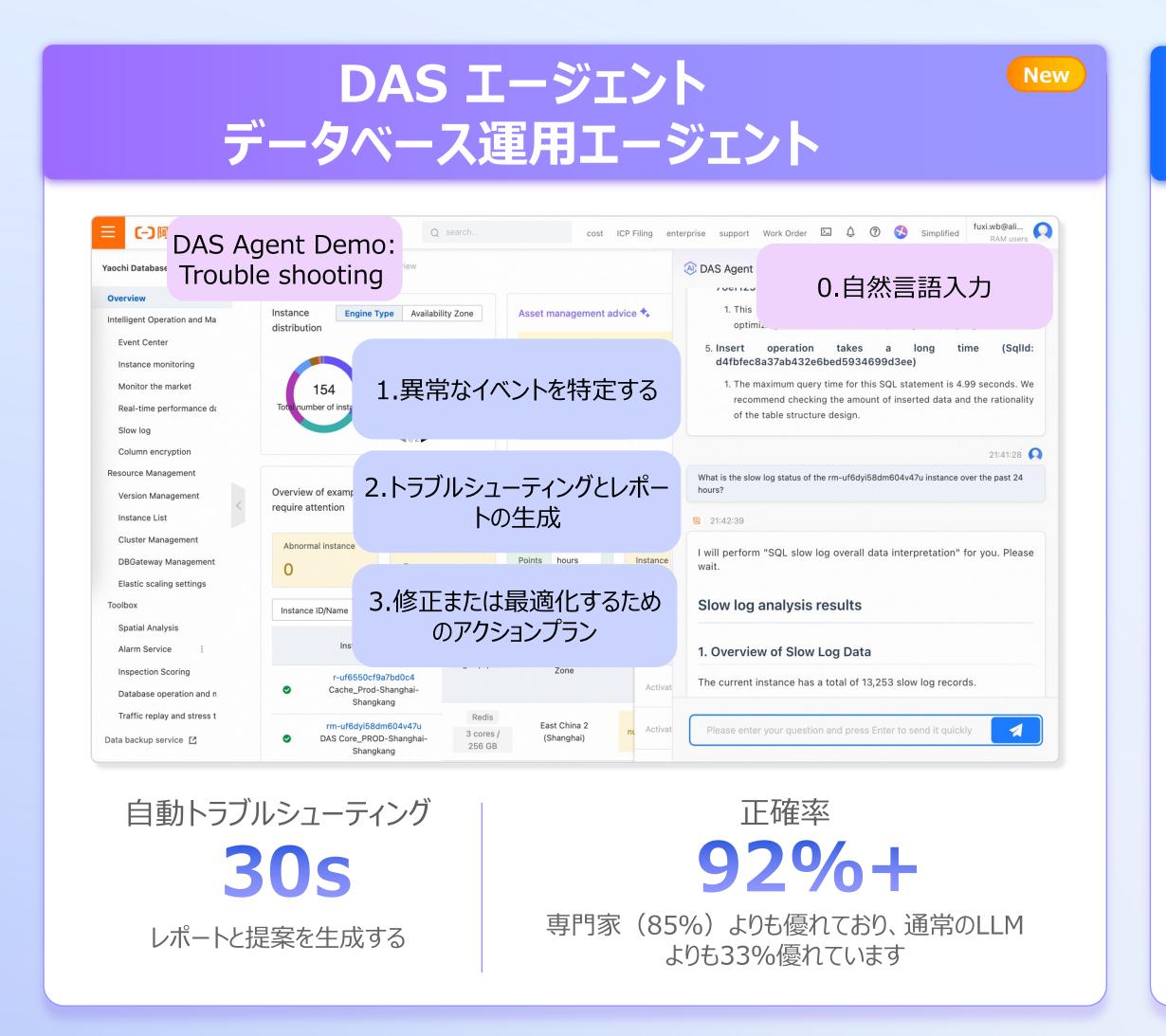


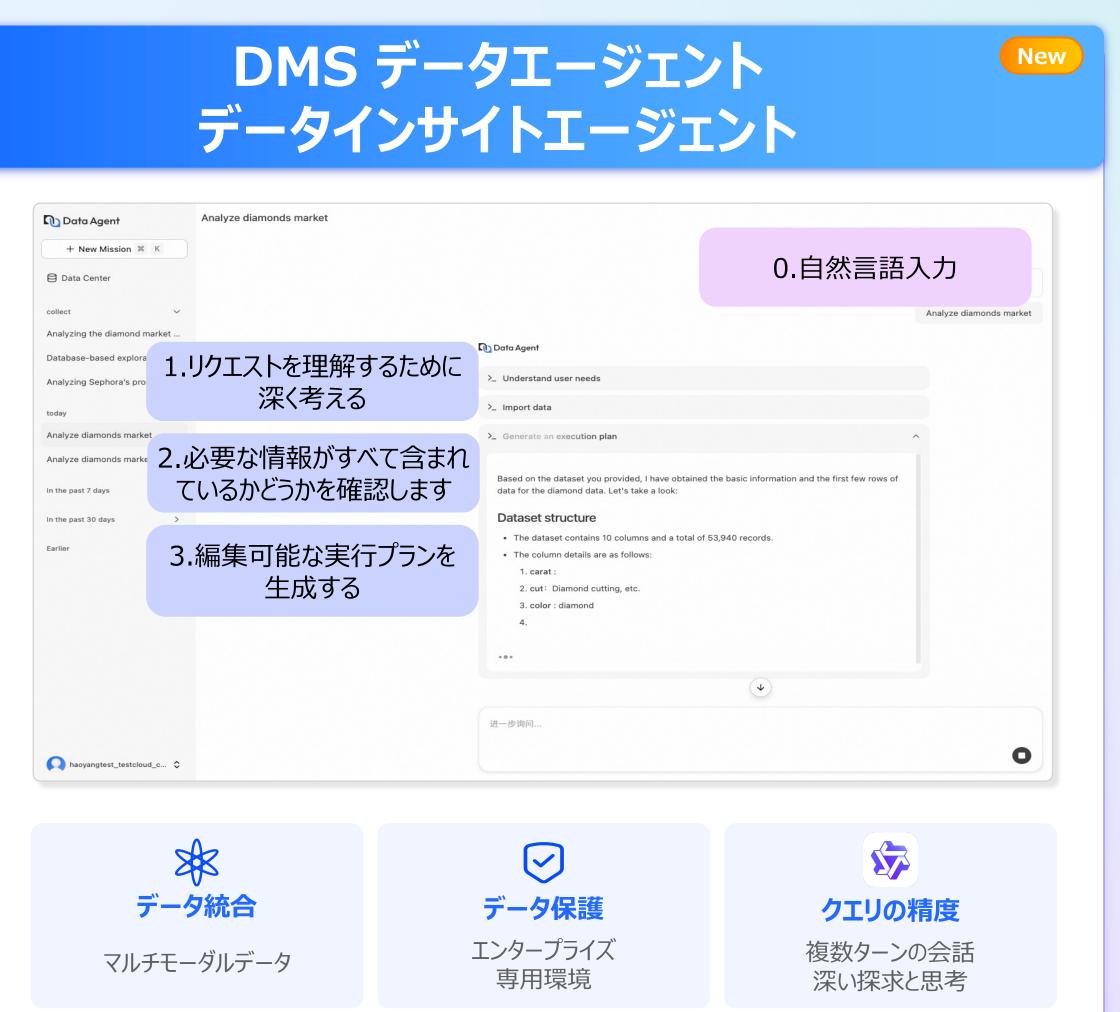
Alibaba Cloud DMS (Data Management)





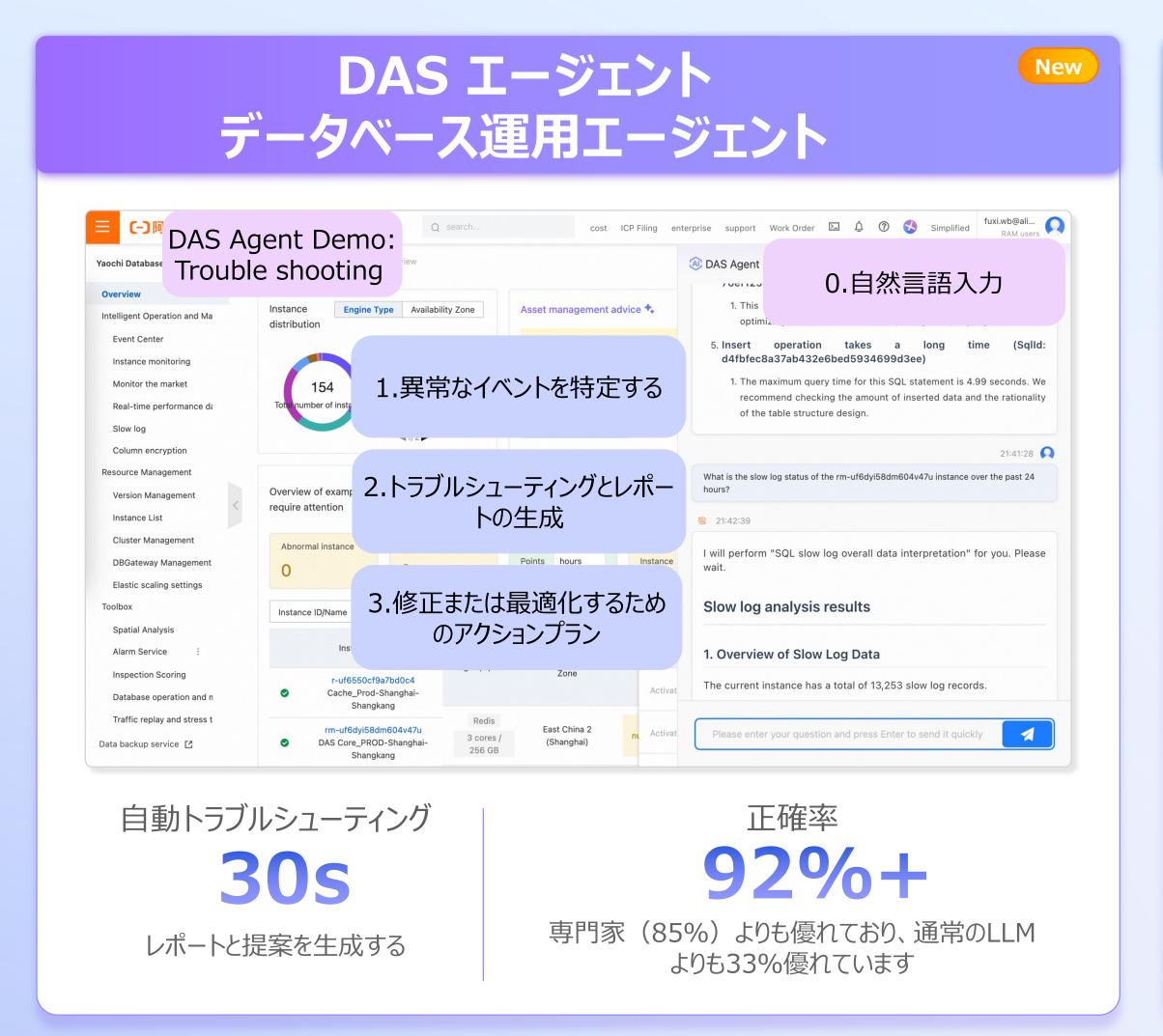
ApsaraDB エージェント - NEW

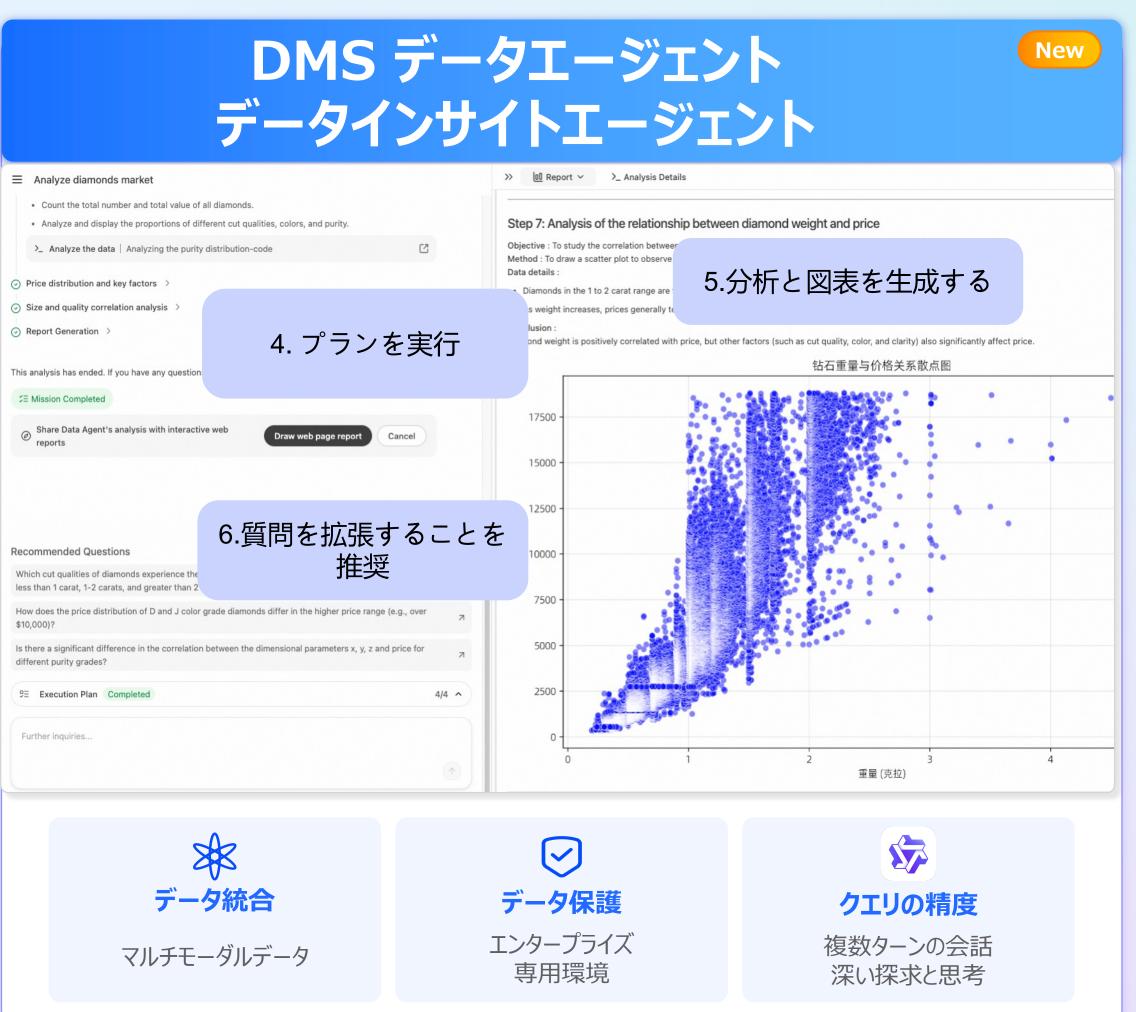






ApsaraDB エージェント - NEW









クラウド製品 AI アシスタント: オープンで 柔軟なクラウド エクスペリエンス









クロスプロダクトクラウドAIOpsエージェント



AIOpsエージェント New

自然言語リクエスト: 「なぜネットワークが遅いのですか?」

フルスタックへの高速アクセス

200以上のクラウド製品、 オープンソースコンポーネントをカバー 多言語サポート

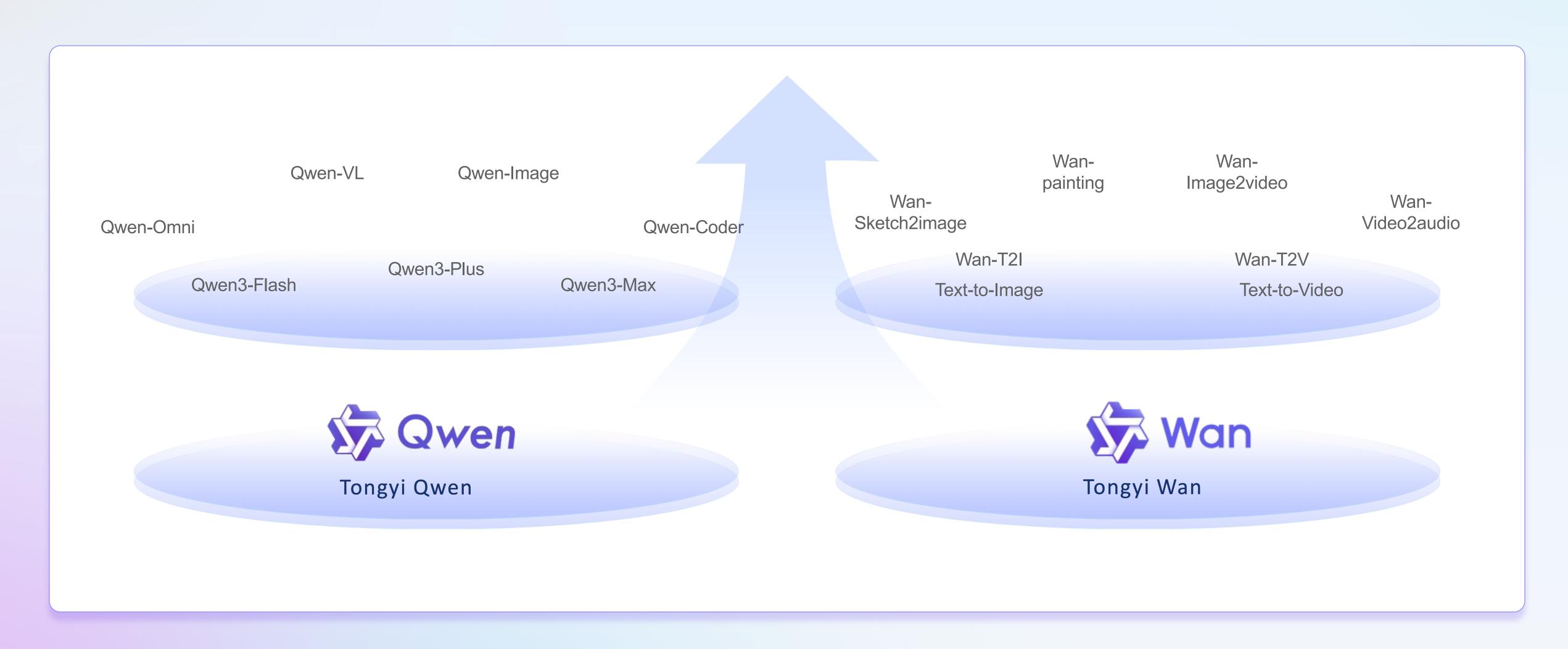
> 統合データストレージおよび 処理プラットフォーム

データ統合ストレージがデータサイロを解消 | 効率的なデータ処理がデータ価値を向上



MaaS: TONGYI Model Family

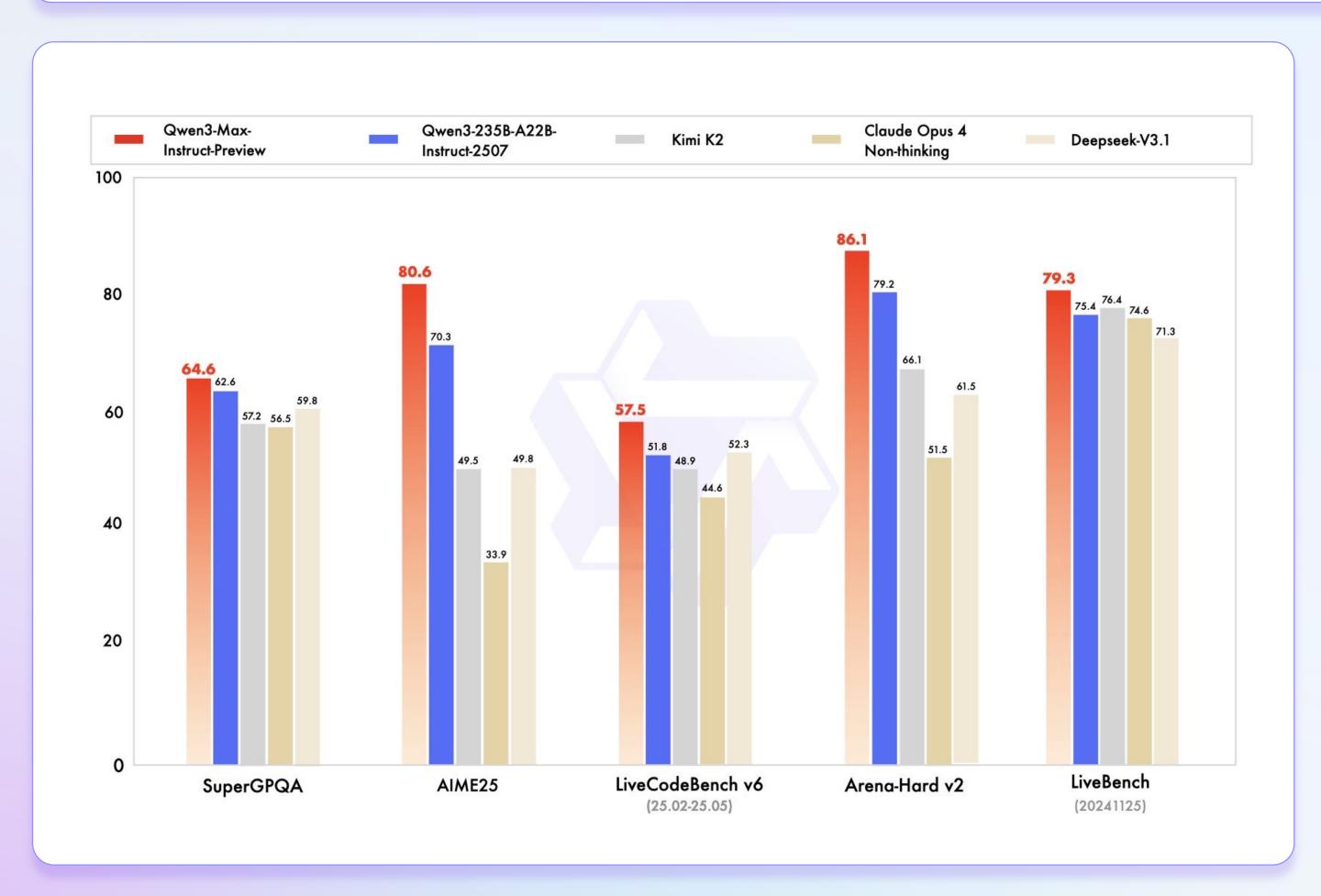
Varied Sizes, Multi-Modalities, Open Source





Qwen3 - Max

Qwen3は、1兆を超えるパラメータを持ち、包括的なパフォーマンス改善を備えた最初のモデルです。 Qwen3-Max, the first model with over one trillion parameters, achieves comprehensive performance improvement.



Highlight

- Comprehensive leap in general capabilities: Based on the Qwen3 series, the overall general capabilities of the Qwen3-Max model have been significantly enhanced compared to Qwen2.5-Max, with notably superior performance in Chinese and English general text understanding and other related abilities.
- 総合的な能力の向上: Qwen3 シリーズをベースにしたこのモデルの全体的な一般機能は Qwen2.5-Max と比較して大幅に向上しており、中国語と英語の一般テキスト理解機能が大幅に向上しています。
- **Significantly improved reasoning ability:** Achieving higher accuracy in tasks requiring reasoning, such as mathematics, coding, logical reasoning, and science subjects.
- **推論能力が大幅に向上:**数学、コーディング、論理的推論、科学など、推論を必要とするタスクにおける高い精度
- Complex instruction following: Better meeting users' specific task requirements and response format needs.
- **複雑な指示に従う:**ユーザーの特定のタスク要件、応答形式などをより適切に満たすことができます。
- Broader Knowledge Coverage: Better Mastery of Long-tail Knowledge,
 Reduced Hallucinations, More Reliable and Credible Responses.
- ・ **より広い知識範囲:**ロングテール知識の理解が深まり、錯覚が減り、より信頼性の高い回答が得られます



Qwen3 - Next

基本モデルの将来における2つの大きなトレンド Two Major Trends in the Future of Foundation Models

コンテキスト拡張 Context-Length Scaling 総パラメータ拡張 Total Parameter Scaling

Qwen3-Next

超高スパース比 Ultra-high sparsity ratio

活性化された専門家の割合1:16 (Qwen)

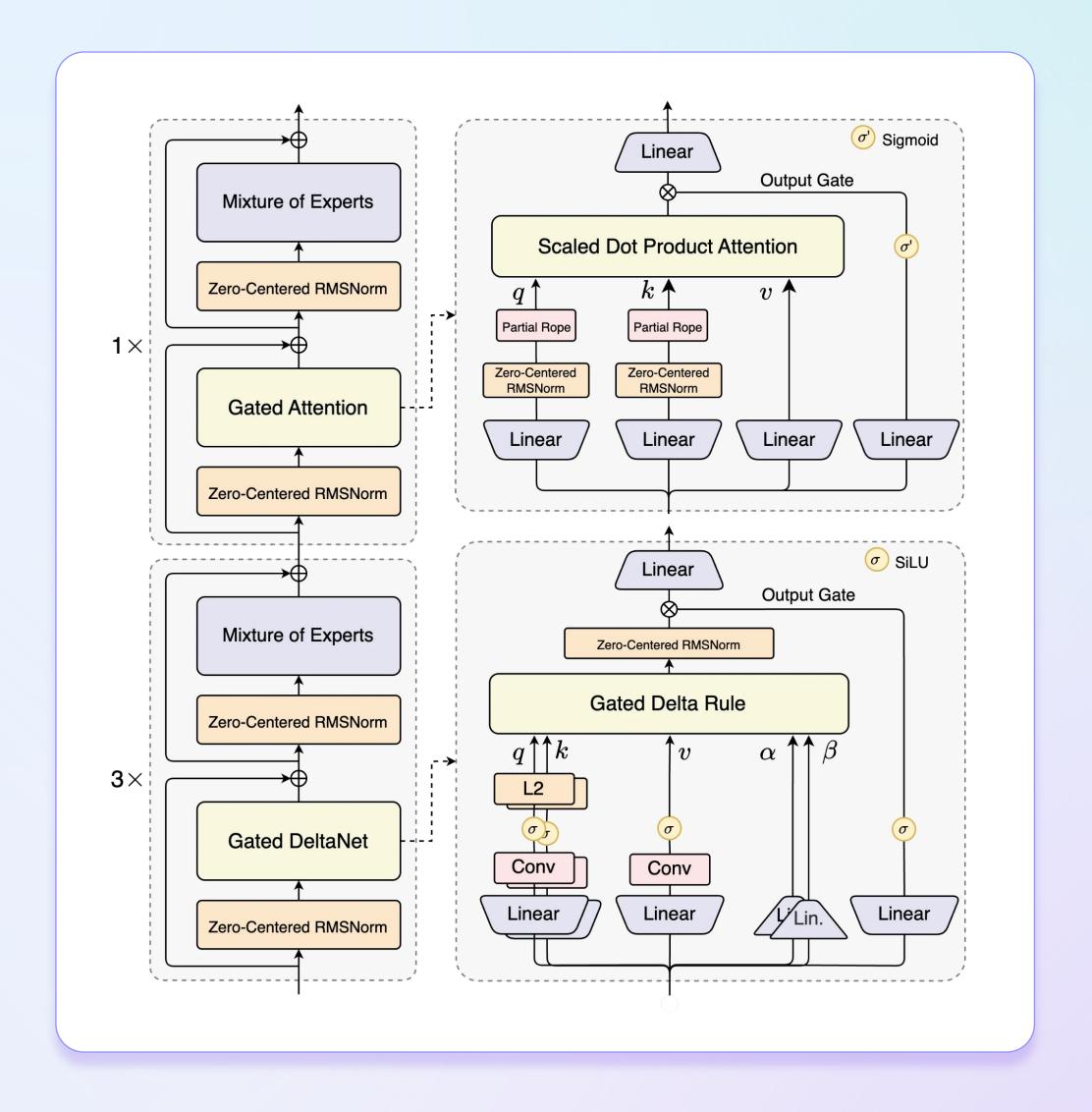
Activation Rate of Experts (Qwen)

1:50 (Qwen-next)

1:50 (Qwen-next)

マルチトークン予測メカニズム^[3] Multi- Token Prediction トレーニングの安定性に配慮した設計

- [1] Gated Delta Networks: Improving Mamba2 with Delta Rule
- [2] Gated Attention for Large Language Models: Non-linearity, Sparsity, and Attention-Sink-Free
- [3] Better & faster large language models via multi-token prediction





Qwen3 VL

Qwen VL 次世代マルチモーダルモデル

Comprehensive scale: From small to large, from Dense to MoE, from Instruct to Thinking, fully covered

包括的なスケール:小規模から大規模、高密度から低密度まで、指示から思考まで

Long sequences: Native 256K support, scalable up to 1M sequence length

長いシーケンス: ネイティブ 256K、1M まで拡張可能な長いシーケンスのサポート

Richer and more comprehensive visual world knowledge, with enhanced object recognition capabilities

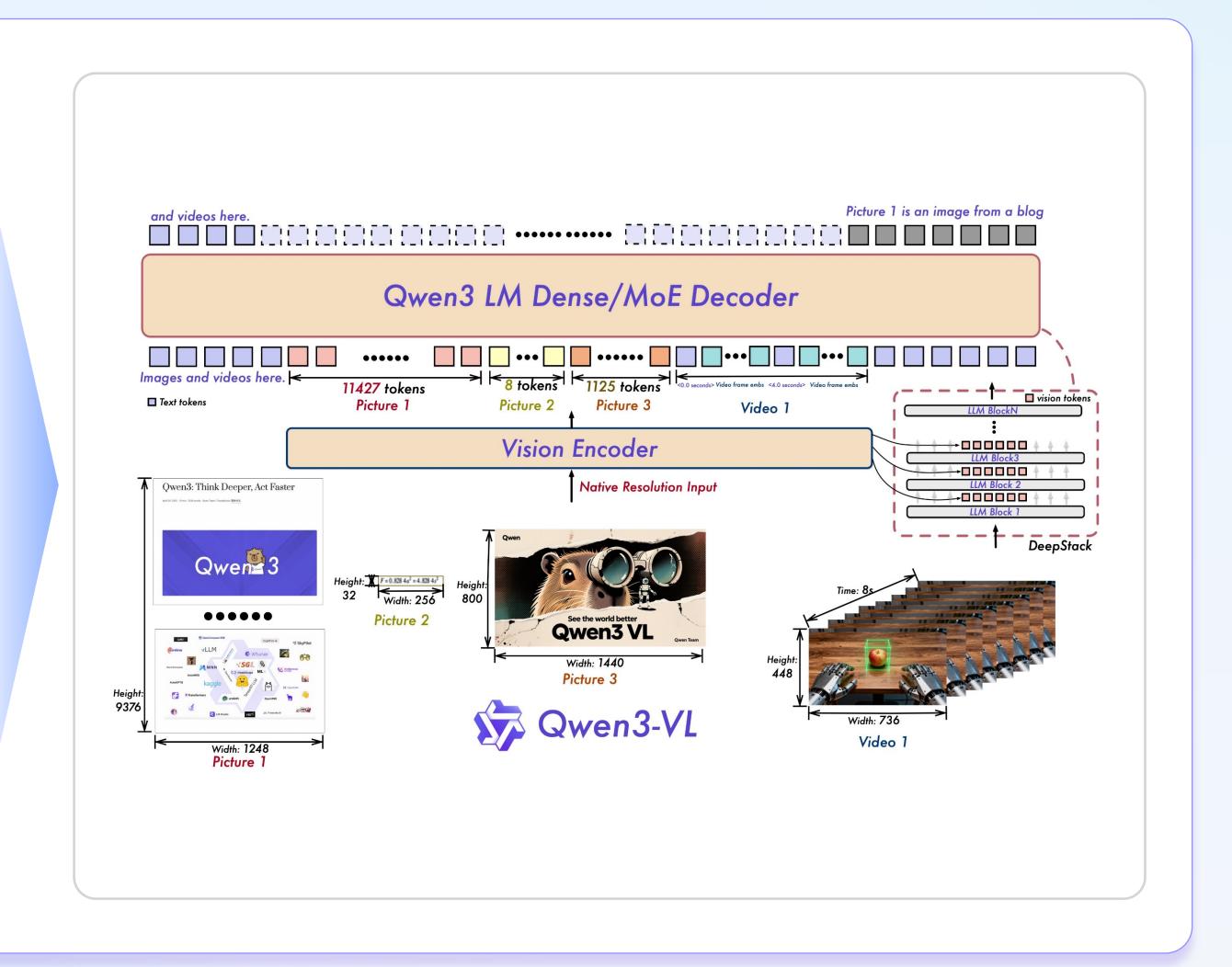
視覚世界に関するより豊かで包括的な知識と、あらゆるものを識別する能力の向上

Enhanced spatial awareness and stronger 2D/3D grounding capabilities

空間認識を改善し、2D/3Dグラウンディング機能を強化します

Superior video understanding capabilities, including long videos and spatiotemporal localization

長時間動画や時空間位置の理解を含む、より強力な動画理解機能



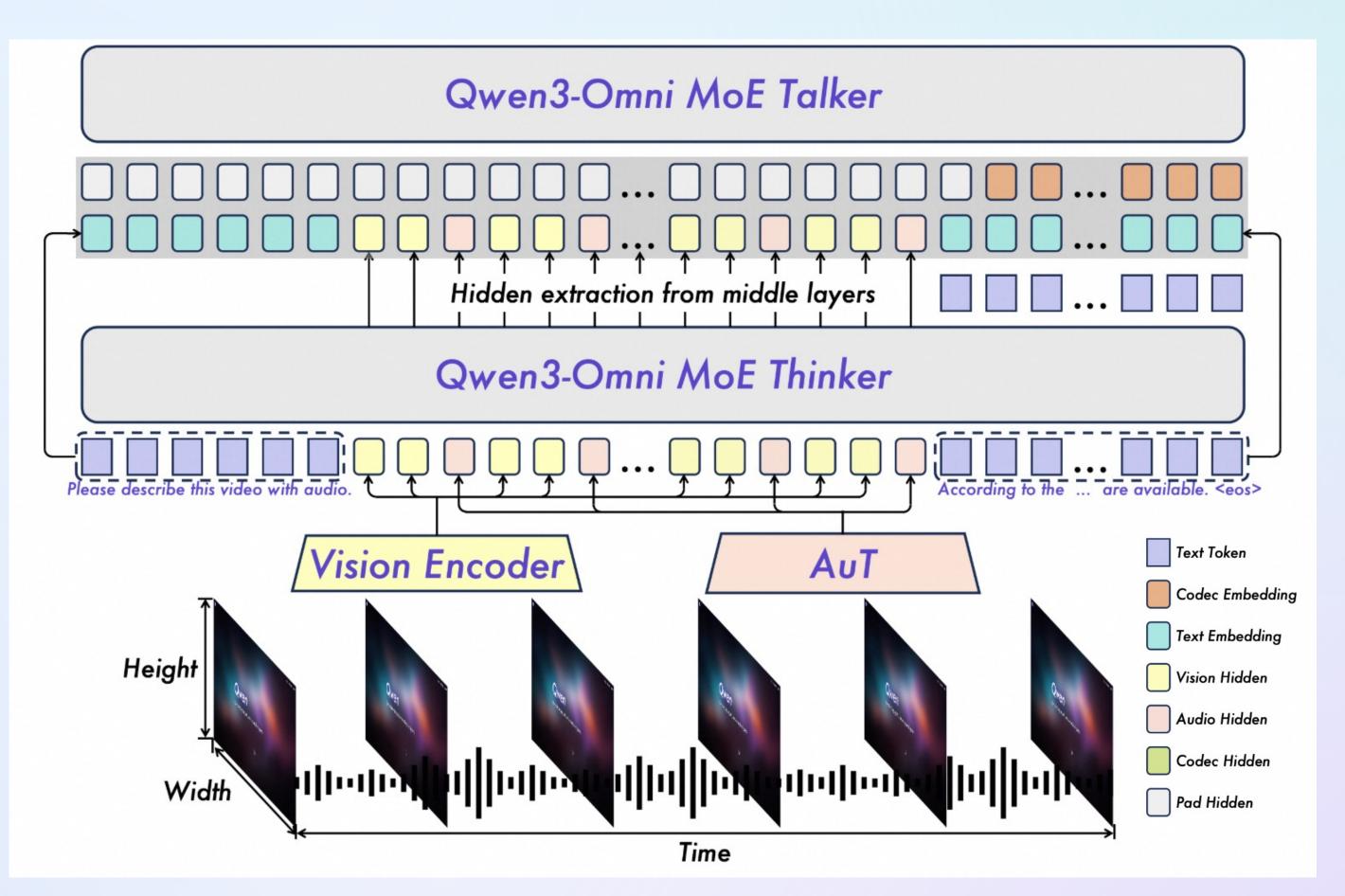


Qwen3 - Omni

エンドツーエンドのマルチモーダルモデルは、マルチモーダルインタラクションエクスペリエンスを強化します

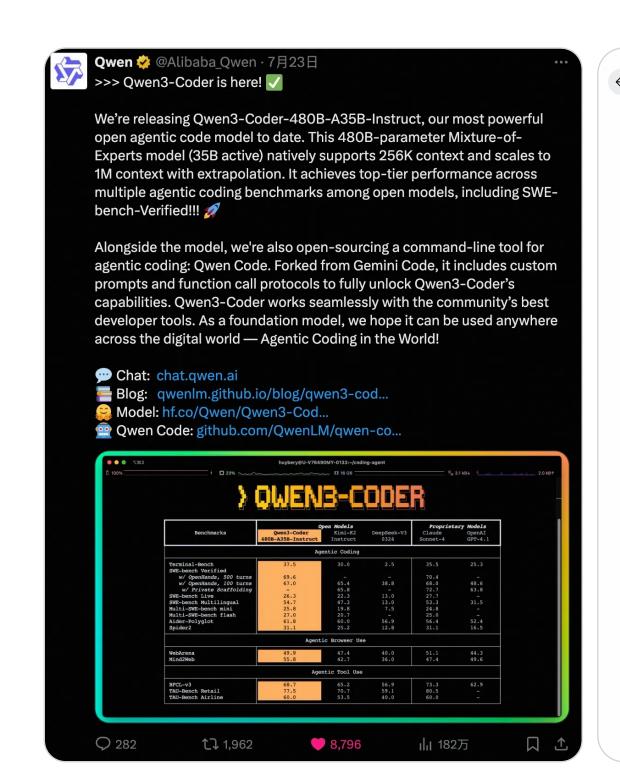
Highlight

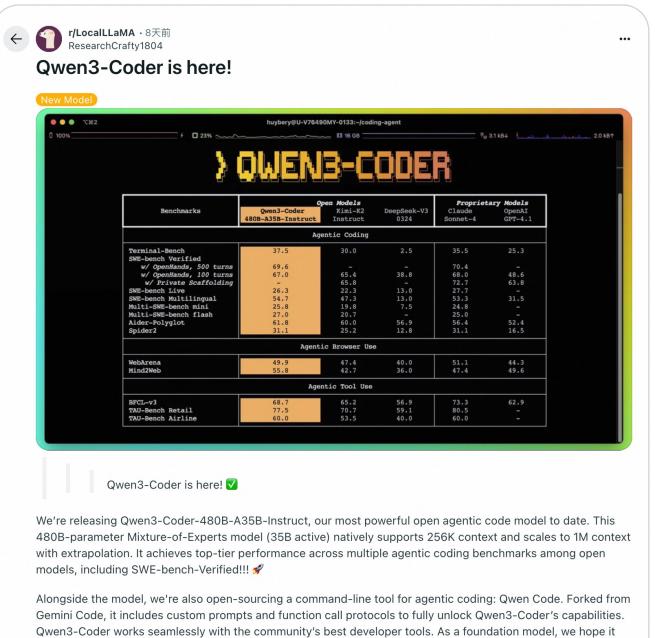
- モデルアーキテクチャ: Qwen3-OmniはThinker-Talkerアーキテクチャを採用しています。
 Thinkerはテキスト生成を担い、TalkerはThinkerから高レベルの意味表現を直接受け取り、
 ストリーミング音声トークンの生成に集中します。超低遅延ストリーミングを実現するために、
 Talkerはマルチコードブックシーケンスの自己回帰予測を使用します。各デコードステップにおいて、
 MTPモジュールは現在のフレームの残差コードブックを出力し、Code2Wavによって対応する波形に合成することで、フレームごとのストリーミングを可能にします。
- Model Architecture: Qwen3-Omni adopts the Thinker-Talker architecture.
 The Thinker is responsible for text generation, while the Talker focuses on generating streaming speech tokens by directly receiving high-level semantic representations from the Thinker.
- **AuT**: オーディオ エンコーダーは、2億時間のオーディオ データでトレーニングされた AuT モデル を使用し、非常に強力な一般的なオーディオ表現機能を備えています。
- AuT: The audio encoder adopts the AuT model trained on 200 million hours of audio data。
- MoE: Thinker と Talker はどちらも MoE アーキテクチャを採用しており、高い並行性と高速な推論をサポートします。
- MoE: Both Thinker and Talker adopt the MoE architecture, supporting high concurrency and fast inference.





Qwen3 - Coder





can be used anywhere across the digital world — Agentic Coding in the World!

♦ 1854 ♦ 🔘 262







https://chat.qwen.ai/

https://huggingface.co/Qwen

https://modelscope.cn/organization/qwen

https://github.com/QwenLM/Qwen3-VL





当社の競争優位性があなたの成功を後押しします

大規模なAIイノベー ションを加速

データインテリジェ ンスを**革**新する 効率的で信頼性の高い エンタープライズサー ビスを促進する

運用の複雑さを軽減

スタック全体にわたるテクノロジーリーダーシップ

AP最大規模の 統合サービス能力 よりコスト効率が高く、より使いやすい

ご清聴頂きありがとうございます

