

# 徹底解剖：中国AIスタートアップ「面壁智能（ModelBest）」の技術戦略とエッジLLMの市場競争力

Gemini 3.1 pro

## 1. 序論：クラウド依存からの脱却とエッジAIパラダイムの台頭

世界の人工知能(AI)研究および産業応用は現在、極めて重要なパラダイムシフトの只中にある。これまで、言語モデルの性能向上は主としてパラメータ規模の指数関数的拡大と、それを支えるクラウド上の莫大な計算資源(GPUクラスター)に依存してきた。しかし、このアプローチは、深刻なネットワーク遅延、膨大なクラウド運用コスト、およびユーザーの機密データの流出リスクという構造的な限界に直面している。こうした中、デバイス上でローカルに動作する高性能かつ軽量の「エッジAI(Edge AI)」モデルの需要が急増しており、この新たな主戦場において劇的な台頭を見せているのが、中国のAIスタートアップ「面壁智能(ModelBest)」である。

ModelBestは、中国の最高学府である清華大学の自然言語処理(NLP)ラボから2022年8月にスピンアウトした企業である<sup>1</sup>。共同創業者兼チーフサイエンティストである劉知遠(Liu Zhiyuan)氏は清華大学計算機科学系の准教授を務め、共同創業者兼CEOの李大海(Li Dahai)氏は中国有数のナレッジプラットフォーム「知乎(Zhihu)」の元CTOという強固な学術的・実業的バックグラウンドを持つ<sup>1</sup>。同社が開発した「MiniCPM」シリーズは、スマートフォン、パーソナルコンピューター、産業用ロボット、そして自動車のインテリジェント・コックピットといったエッジデバイス上で、クラウドに依存することなく完全に独立して実行可能な小規模・高性能言語モデル(SLM: Small Language Models)である。

設立からわずかな期間で、MiniCPMシリーズはGitHubやHugging Faceといった世界的なオープンソースコミュニティにおいて累計2400万回以上のダウンロードを記録し、オープンソースLLM分野におけるデファクトスタンダードの一角を占めるに至った<sup>1</sup>。さらに同社はこのほど、数億元(数十億円)規模の新規資金調達を完了し、わずか3カ月間で累計10億元(約230億円)以上の資金を集め、その企業評価額は正式にユニコーン(10億ドル以上)の水準に達している<sup>3</sup>。本レポートでは、多角的なデータ分析に基づき、ModelBestの背後にある「知識密度(Density Principle)」という独自の設計思想、アーキテクチャの革新、他社の競合モデル(Phi-3.5、Llama-3.1等)との定量的比較、そして広範な産業エコシステムへの浸透戦略を徹底的に解剖する。

## 2. 資本政策と戦略的アライアンス：インフラと実業の融合

ModelBestの近年の資金調達ラウンドを分析すると、単なる財務的支援を超えた、極めて意図的な「垂直・水平統合型」の産業ネットワーク構築の意図が読み取れる。

### 2.1 国家情報インフラとの結びつき

直近の数億元規模の調達ラウンドは、中国三大通信事業者の一角である「中国電信(China

Telecom)」が主導し、中信金石 (Citic Goldstone) や中信私募 (Citic Private Equity) といった国家資本系の大手投資機関が追随した<sup>2</sup>。この資本提携の背後にある戦略的意図は重大である。中国電信は、国内全域に広がるクラウドコンピューティング能力、5G/6Gネットワークインフラ、そして無数の末端デバイスへのアクセス経路を有している。ModelBestは、この通信インフラと自社のエッジアルゴリズムを融合させることで、司法、自動車、教育といった複雑かつセキュリティ要件の厳しい産業シナリオに対するAI実装を急速に拡大する計画である<sup>2</sup>。国家の情報インフラストラクチャープロバイダーとの提携は、中国国内市場におけるB2B(企業向け)およびB2G(政府向け)事業の展開において、圧倒的な参入障壁の構築を意味する。

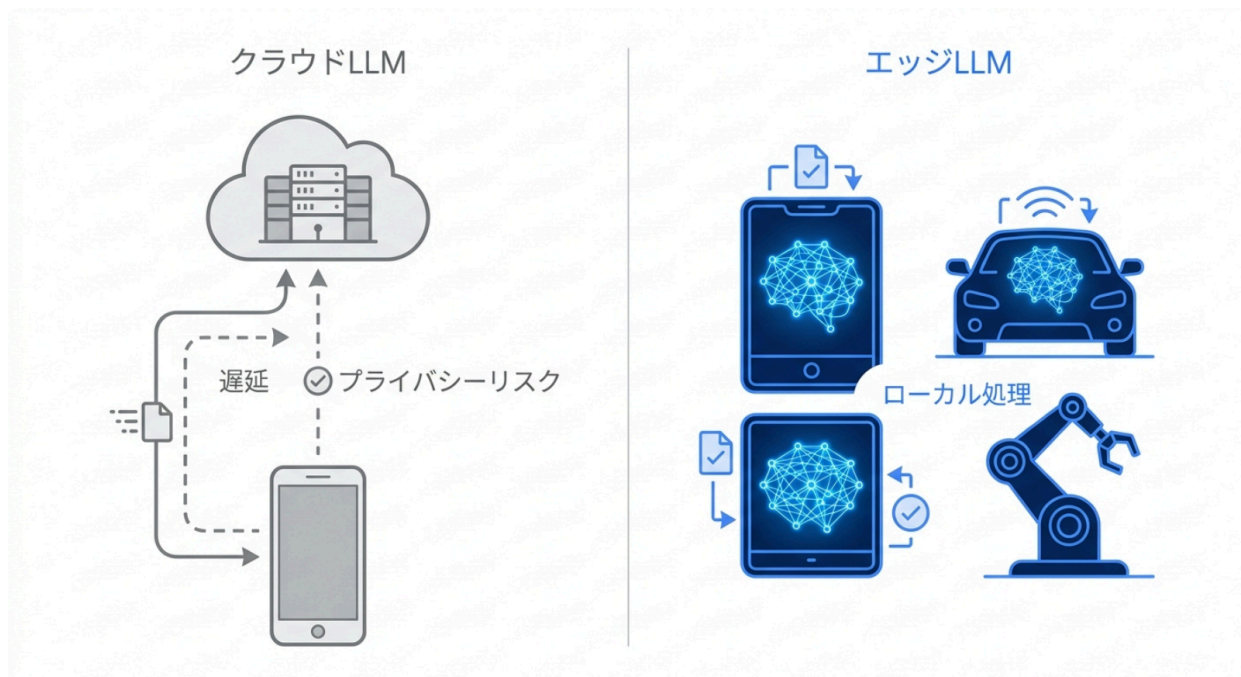
## 2.2 エンボディドAI(身体性AI)とスマートマニュファクチャリング

また、別のラウンドで主導的役割を果たした「匯川産投 (Inovance Investment)」の存在も特筆に値する<sup>3</sup>。同社の親会社であるInovance Technologyは、産業用オートメーション機器およびロボット制御コンポーネントにおける世界的リーダーである。この提携は、MiniCPMの適用領域が単なるテキスト生成やスマートフォン向け対話エージェントにとどまらず、エッジデバイスで動作する「エンボディドAI (Embodied AI: 身体性を持つAI)」や産業用ロボットの自律制御システムへと拡張されることを示唆している。データ処理の遅延が物理的な動作の遅延に直結するロボティクス分野において、ローカルで完結する高性能LLMは不可欠な技術要素である。

## 3. 核心設計思想: 「知識密度法則 (Density Principle)」

エッジデバイス(モバイル端末、車載NPU、PC等)は、クラウドサーバーと比較してメモリ容量、計算能力(TOPS)、および消費電力において極めて厳格な制約が存在する。ModelBestの技術的成功の根底にあるのは、「パラメータ数を無闇に巨大化させるのではなく、限られたパラメータ空間の中にどれだけ高密度に知識と推論能力を圧縮できるか」を追求する「Density Principle(知識密度法則)」である<sup>1</sup>。

## クラウドLLMとエッジLLM（MiniCPM）のアーキテクチャ比較



MiniCPMはクラウドサーバーへのデータ送信を必要とせず、スマートデバイスのローカルチップ（NPU/GPU）上で推論を完結させる。これにより、個人データの流出リスクを根本的に排除し、通信遅延のないリアルタイム応答を実現している。

### 3.1 「モデル風洞実験 (Model Wind Tunnel Experiments)」アプローチ

大規模言語モデルの開発において最も膨大なコストを要するのは、学習率 (Learning Rate) やバッチサイズといった最適なハイパーパラメータの探索プロセスである。巨大なパラメータサイズのモデルを直接訓練して試行錯誤することは、計算資源の甚大な浪費を招く。ModelBestは航空宇宙工学にインスピレーションを得た「モデル風洞実験 (Model Wind Tunnel Experiments)」アプローチを採用している<sup>4</sup>。

これは、Tensor Programパラメータ化スキームを活用することで、モデルのスケールが拡大しても最適なハイパーパラメータが本質的に変動しないという特性を利用した手法である<sup>4</sup>。同社の研究チームは、0.04B (4000万)、0.1B、0.2B、0.5B、0.8B、および1.2Bという6つの小規模なモデルサイズを用いて広範な学習率実験を実施した。その結果、モデルサイズが10倍に拡大しても、最適な学習率は一貫して「0.01」付近で安定することを発見した<sup>5</sup>。さらに、2.1Bスケールのモデルでの検証でも、学習率0.01が最小の損失 (Loss) をもたらすことが実証されている。

### 3.2 チンチラ最適スケージングの予測と実証

さらに、データ量とモデルサイズの関係性を定義する「Chinchilla最適データ量」に基づくスケージング則を用いることで、大規模モデルの訓練結果を極めて正確に事前予測することに成功している。

モデルサイズ	予測 / 実測	C4 データセット Loss
1.2B	実測値	2.89
2.4B	予測値(*)	2.70
7B	予測値(*)	2.45
9B	予測値(*)	2.40
<b>MiniCPM (8~9B級)</b>	実測値	<b>2.41</b>
13B	予測値(*)	2.32

上記表に示されるように、小規模モデルから導き出されたプロットに基づき、9Bモデルの最終的なC4 Lossは約2.40、7Bモデルでは約2.45になると予測された<sup>5</sup>。実際のMiniCPMの最終Lossは「2.41」であり、事前予測されたChinchilla最適モデルの数値とほぼ完全に一致した。この風洞実験的アプローチにより、ModelBestは計算資源を一切浪費することなく、目標とするLossに到達するための最適なトレーニングレシピを設計し、パラメータ単位あたりの性能(知識密度)を極限まで高めているのである。

## 4. アーキテクチャのブレイクスルー: 長文脈処理と視覚圧縮技術

エッジデバイス向けのLLMが直面する最大の技術的ハードルは、膨大なメモリ消費を伴う「長文脈の維持(Long Context)」と「高解像度画像・動画のマルチモーダル処理」である。ModelBestは、これら二つの領域において構造的な革新をもたらした。

### 4.1 MiniCPM-SALA: ハイブリッド注意機構による1Mトークン処理

一般的なTransformerベースのモデルでは、自己注意機構(Self-Attention)の計算量およびKVキャッシュ(Key-Value Cache)のメモリ使用量が、入力シーケンス長の2乗に比例して爆発的に増加

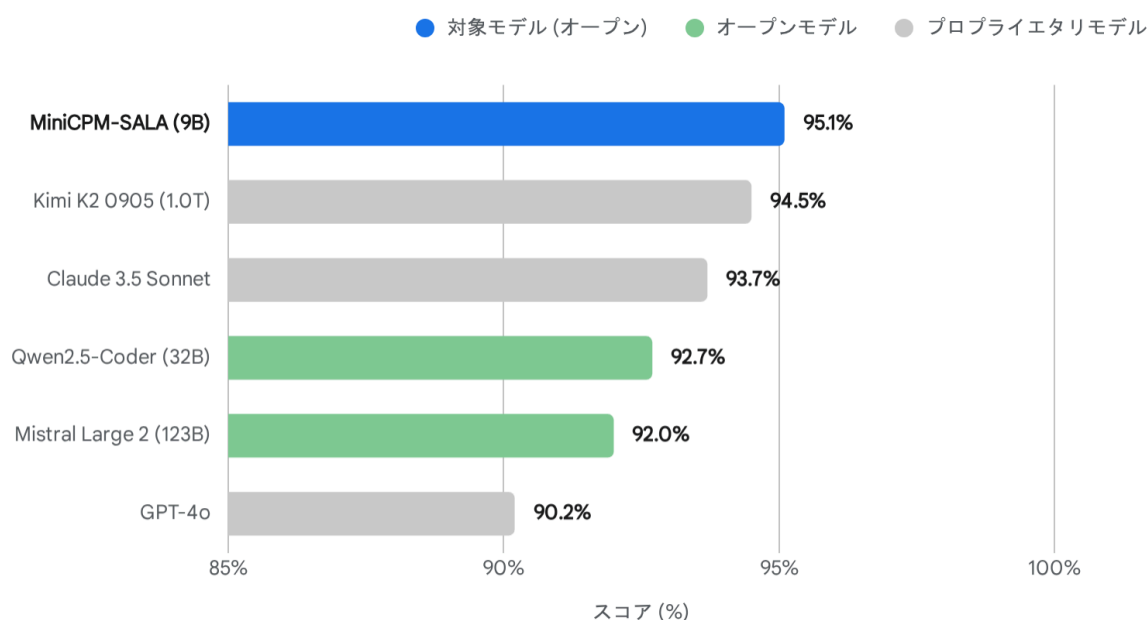
する。したがって、8Bクラスの完全自己注意(Full-Attention)モデルで数十万トークンを処理しようとする、すぐにメモリ不足(OOM: Out-Of-Memory)エラーを引き起こす<sup>6</sup>。

この物理的限界を打破するため、ModelBestは「MiniCPM-SALA(Sparse Attention and Linear Attention)」と呼ばれるハイブリッド・アーキテクチャを開発した<sup>6</sup>。このアーキテクチャは、以下の1:3の比率で層を構成する。

- **25%の層(Sparse Attention)**: 局所的な詳細情報に高い粒度でフォーカスするために、「InfLLM-V2」を用いた疎な注意機構を採用する。
- **75%の層(Linear Attention)**: 広範なコンテキスト全体を極めて高い計算効率で処理するために、「Lightning Attention」を用いた線形注意機構を採用する<sup>6</sup>。

このSALAアーキテクチャの導入により、単一のNVIDIA A6000DやRTX 5090といったコンシューマーまたはワークステーションクラスのGPU上で、なんと「最大100万(1M)トークン」のコンテキスト推論を実現した<sup>6</sup>。さらに、256Kトークンのシーケンス長における推論速度は、従来のFull-Attentionモデルと比較して最大3.5倍に達している<sup>6</sup>。

# HumanEvalベンチマークにおけるMiniCPM-SALAの性能比較



MiniCPM-SALA (9Bパラメータ) は、コーディング能力を測るHumanEvalにおいて95.1%のスコアを達成し、パラメータ数が数十倍から数百倍と推測される最先端のプロプライエタリモデル (Claude 3.5 Sonnet, GPT-4o) や、Mistral Large 2 (123B) を凌駕している。これはSALAアーキテクチャによる長文脈処理の効率化が、論理的推論タスクにおいて極めて有効であることを示している。

データソース: [llm-stats.com](https://llm-stats.com)

SALAアーキテクチャの真価は、単なるテキストの長文脈処理にとどまらない。コード生成や論理的推論タスクにおいても圧倒的な結果を残している。プログラムの機能的正しさを測定する「HumanEval」ベンチマークにおいて、MiniCPM-SALA(9Bパラメータ)は「95.1%」という驚異的なスコアを記録し、評価対象となった全66モデルの首位に立った<sup>10</sup>。これは、Anthropicの「Claude 3.5 Sonnet」(93.7%)やOpenAIの「GPT-4o」(90.2%)、さらにはMistral AIの巨大モデル「Mistral Large 2」(123B、92.0%)といった、クラウド上で稼働する世界最先端のプロプライエタリモデルをも上回る成果である<sup>10</sup>。この事実は、ハイブリッド注意機構がコード内の複雑な依存関係の追跡に極めて有効であり、パラメータ数への依存を構造的イノベーションによって凌駕できることを証明している。

## 4.2 3D-Resamplerによる動画トークンの劇的圧縮

エッジデバイスにおけるマルチモーダルLLM(MLLM)の運用において、映像データの処理はテキスト以上にリソースを消費する。通常、動画の連続するフレームをLLMに読み込ませるとトークン数が爆

発し、デバイスのメモリを即座に枯渇させてしまう。

この問題に対し、MiniCPMの最新旗艦モデル「MiniCPM-V 4.5」は、画像および動画に対する新しい統合型「3D-Resampler」を導入した<sup>11</sup>。この技術は、448x448ピクセルの動画フレームを空間的・時間的に解析し、6つのフレームをわずか「64個のビデオトークン」に共同圧縮することを可能にした。一般的なMLLMが同様の情報を処理する際に約1,536トークンを消費することを考慮すると、これは驚異の「96倍の圧縮率」に相当する<sup>11</sup>。この劇的なトークン圧縮技術の恩恵により、LLMの推論コスト（計算量とメモリ消費）を増加させることなく、最大10FPS (Frames Per Second) の高フレームレート動画や長時間の動画の空間・時間的情報を、エッジデバイス上でリアルタイムに理解できるようになった<sup>11</sup>。この技術は、常に車内外の複数のカメラ映像を処理し続けなければならない自動運転車やインテリジェント・コックピット、そして視覚フィードバックを必要とする産業用ロボットにおいて、極めて強力な競争優位性をもたらしている。

## 5. フラッグシップモデルの性能と他社競合モデルとの定量的比較

ModelBestのラインナップは、視覚理解に特化した「MiniCPM-V」シリーズと、全機能 (Omni) を統合した「MiniCPM-o」シリーズの二本柱で構成されている。これらのモデルは、同等クラスの軽量化LLMを展開するMicrosoft (Phi-3.5)、Meta (Llama-3.1-8B)、Google (Gemma-9B)、Alibaba (Qwen2.5) と直接的に市場を争っている。

### 5.1 MiniCPM-V 4.5: 視覚・言語のSOTAモデル

MiniCPM-V 4.5は、Qwen3-8BおよびSigLIP2-400Mをベースに構築された総パラメータ数8BのMLLMである<sup>11</sup>。本モデルは、8つの主要な視覚・言語ベンチマークを統合した包括的評価ツール「OpenCompass」において、平均スコア「77.0」を記録した<sup>11</sup>。このスコアは、パラメータ規模が約9倍に及ぶオープンソースの強豪「Qwen2.5-VL 72B」を上回るだけでなく、OpenAIの「GPT-4o-latest」やGoogleの「Gemini-2.0 Pro」といった最新鋭の商用モデルの性能をも凌駕している<sup>11</sup>。とりわけ、LLaVA-UHDアーキテクチャを活用した高度なOCR能力と文書解析能力に優れており、「OmniDocBench」での英語文書のエンドツーエンド解析においては、GPT-5やGemini-3 Flash、さらには専門特化型OCRツールであるDeepSeek-OCR 2よりも優れたSOTA (State-of-the-Art) 性能を実証した<sup>11</sup>。

### 5.2 ハイブリッド推論と強化学習アライメント

MiniCPM-V 4.5の機能的特長として、ユーザーの利用シナリオに応じた「Fast / Deep Thinking (速い思考 / 深い思考)」の制御可能なハイブリッド・モードの搭載が挙げられる<sup>11</sup>。日常的な単純なクエリに対しては消費電力と遅延を抑えるFast Thinkingを用い、複雑な推論タスクに対してはDeep Thinkingを動的に切り替えることができる。さらに注目すべきは、ファインチューニングにおける強化学習 (RLHF) の高度なアプローチである。事実関係の正確性が問われる視覚タスクにおいては、単純な応答には「ルールベースの報酬」を適用し、複雑な自然言語応答には「確率ベースの報酬」を適用するハイブリッドRL (強化学習) とDPO (直接選好最適化) を組み合わせることで、ハルシネーション (AIの幻覚、事実の捏造) を強力に抑制している<sup>14</sup>。このアプローチにより、自然な流暢さを損なうこ

となく、長文脈推論モードと短文脈モードの間でのクロスモード汎化能力を獲得している<sup>14</sup>。

### 5.3 競合他社SLMとのベンチマークおよびコスト比較

ここで、同クラスのオープンソースSLMであるLlama-3.1-8B-instructおよびPhi-3.5-mini-instructとの定量的な比較を行う。

指標 / モデル	MiniCPM-SALA (9B)	Llama-3.1-8B-Instruct	Phi-3.5-mini-Instruct (3.8B)
コンテキスト長 (最大)	1,000,000 (1M)	131,072 (128K)	128,000 (128K)
HumanEval (コード)	95.1%	72.6%	62.8%
ARC-C (推論)	N/A	83.4%	84.6%
MMLU (一般知識)	N/A	67.9%	N/A
RULER 128K (長文脈検索)	N/A	77.0%	63.6%
APIコスト (1M入力/出力)	ローカル運用主体	\$0.03 / \$0.03	\$0.10 / \$0.10

※一部データは複数のベンチマークからの総合値<sup>6</sup>

上記の表と分析から、以下の重要なインサイトが導き出される。

1. 長文脈における安定性の差異: コンテキスト長において、Phi-3.5は128Kをサポートするものの、実際の情報検索能力を測るRULERベンチマークでは、64Kトークンを超えると性能が急激に低下し、128Kでは63.6%に落ち込む<sup>16</sup>。一方、Llama-3.1は128Kでも77.0%の精度を保っている<sup>16</sup>。しかし、MiniCPM-SALAはこれらを遥かに凌ぐ1M(100万)トークンをサポートしており、法的文書の全件スキャンや複数冊の書籍を跨いだ分析といった領域で他を圧倒している。

2. **API処理コストの非対称性**: クラウドAPIベースの価格設定において、Llama-3.1-8Bは100万トークンあたり\$0.03という驚異的な低コストを実現しており、Phi-3.5(\$0.10)と比較して約3.3倍安価である<sup>15</sup>。しかし、MiniCPMの主たる戦略は「オンデバイスでのローカル推論」であり、デバイス購入後に追加のAPIコストは一切発生しない。推論タスクが持続的に発生するロボティクスや自動車システムにおいて、この限界費用ゼロの構造は強力な経済的優位性となる。
3. **多言語対応の設計思想**: Phi-3.5-miniはアラビア語や日本語などを含む多言語対応を大幅に強化(旧バージョン比で25~50%向上)したが、ボキャブラリサイズが32Kと制約されており、低リソース言語には追加のファインチューニングが推奨されている<sup>19</sup>。一方、MiniCPMシリーズは初期モデル(MiniCPM-Llama3-V 2.5の段階)から30言語以上をネイティブサポートし、グローバルなOCRと多言語環境におけるGPT-4Vレベルの性能をエッジ側で確保している<sup>13</sup>。

## 5.4 ハードウェア・アーキテクチャ依存の推論速度(M1 vs M3の教訓)

エッジデバイスにおけるローカル推論速度(tokens/s)は、CPUやNPUの単なる演算能力(FLOPS)ではなく、「メモリ帯域幅(Memory Bandwidth)」に強く制約されることが複数の検証で判明している<sup>20</sup>。例えば、Appleシリコンを搭載したMacBook Proの比較テストにおいて、より新しいM3 Proチップ(帯域幅150GB/s)を用いた場合、Llama 3やGemmaの推論速度は、旧型のM1 Proチップ(帯域幅200GB/s)よりもわずかに遅くなるという逆転現象が観測された<sup>20</sup>。ただし、パラメータサイズが極めて小さいPhi-3(4B)においては、キャッシュに収まりやすいためM3 Proの方がわずかに高速であった<sup>20</sup>。ModelBestのアーキテクチャ設計(特に前述の3D-Resamplerによる動画トークンの劇的圧縮や、SALAIによるKVキャッシュの効率化)は、まさにこの「エッジデバイス特有のメモリ帯域幅のボトルネック」を深いレベルで理解し、それをアルゴリズム側で回避するために最適化された結果であると分析される。

## 6. 全二重オムニモーダル通信の実装と「MiniCPM-o」

ModelBestの技術的野心の頂点に位置するのが、全機能(Omni)を統合した「MiniCPM-o」シリーズである。その最新版「MiniCPM-o 4.5(9Bパラメータ)」は、業界初となるエッジデバイス向けの「全二重・全モーダル(Full-Duplex Omni-Modal)大規模モデル」として設計された<sup>1</sup>。パラメータサイズがわずか9Bでありながら、GPT-4oやGemini 2.0 Proを凌駕し、Gemini 2.5 Flashに迫る視覚・言語能力を有する<sup>13</sup>。

### 6.1 リアルタイム全二重通信メカニズム

従来のスマートフォンやスマートスピーカーに搭載されているAIアシスタントは、ユーザーが発話を終了するのを待ってから(トリガーワードや無音区間を検知してから)クラウドへのデータ送信、推論、そして音声生成を行う「半二重(ターン制)」の通信方式であった。これに対し、MiniCPM-o 4.5は「完全二重マルチモーダル・ライブストリーミング」を実現している<sup>13</sup>。

技術的な核心は、時分割多重化(TDM: Time-Division Multiplexing)メカニズムの導入である。モデルはミリ秒単位のタイムライン上で、入カストリーム(高フレームレートの映像と音声)と出カストリーム(テキスト生成と音声合成TTS)を同期させ、互いのプロセスをブロックすることなく同時並行で処理する<sup>13</sup>。つまり、モデルはリアルタイムで「外界を見ながら、人間の指示を聞きながら、同時に話す」こ

とが可能である。

## 6.2 プロアクティブな対話と広範なデプロイメント環境

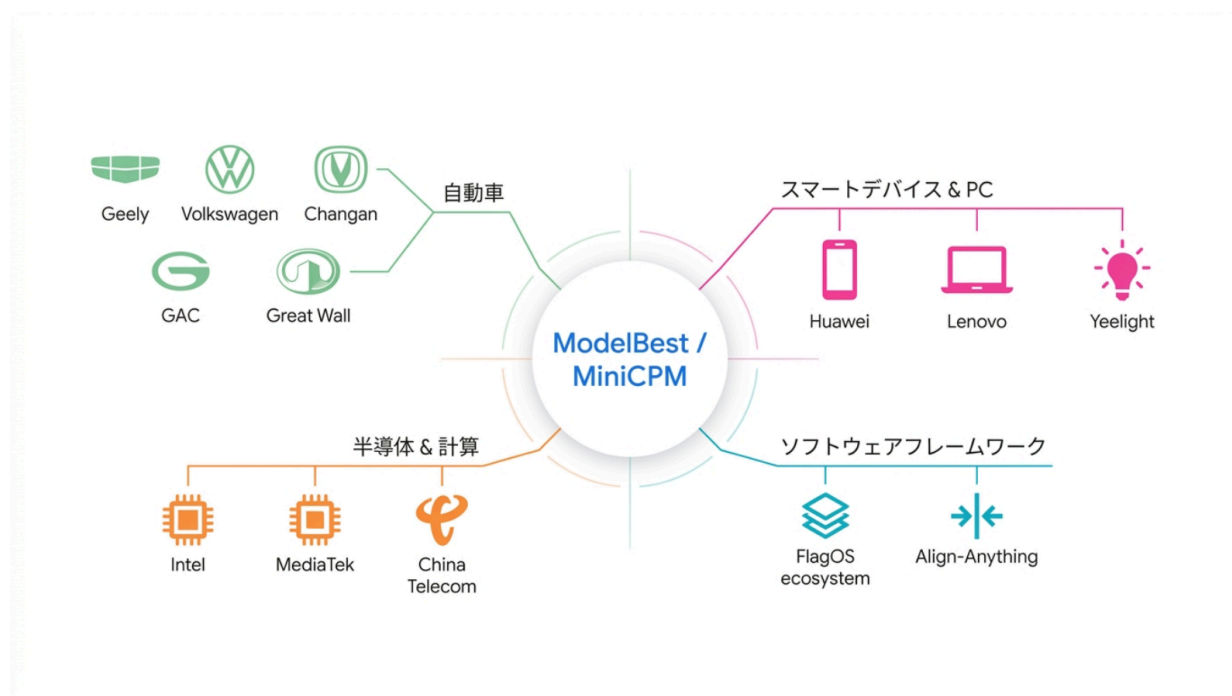
さらに、MiniCPM-o 4.5は連続的な映像・音声入力を常に監視し、1Hz(1秒間に1回)の頻度で自ら「話すべきか否か」の意思決定を行う能力を持つ<sup>13</sup>。これにより、自動車の運転手が居眠りをしているのをカメラで検知して自発的に警告を発したり、料理をしているユーザーに対して映像を基に「次に塩を加えてください」と能動的にアドバイス(プロアクティブ・インタラクション)を行ったりすることが可能となる<sup>13</sup>。

このオムニモーダル機能は、限られた研究環境だけでなく、広範なデプロイメントが公式にサポートされている。MacやNVIDIA GPU上で動作する「llama.cpp-omni」推論フレームワークとWebRTCプロトコルを統合したリアルタイムWebデモがオープンソース化されているほか、iPhoneおよびiPad上でネイティブに動作する公式のiOSアプリも提供されている<sup>13</sup>。特にiOSデバイス上でのローカル実行は、通信ネットワークへの依存を完全に排除し、ユーザーのカメラ映像や生体データといった極めてセンシティブなプライバシー情報を一切外部に送信することなく、安全かつ低遅延なAIアシスタント体験を可能にしている<sup>13</sup>。さらに、最新のMiniCPM-V 2.6等のバージョンは、iPadなどのエンドデバイス上でリアルタイムの動画理解を初めてサポートするに至っている<sup>21</sup>。

## 7. 基盤ソフトウェアの抽象化と広範なエコシステム構築

ModelBestが他のLLM開発企業と一線を画しているのは、アルゴリズムの優位性にとどまらず、最下層のAIチップ制御フレームワークから最終製品に至るまでの強固なエコシステム(垂直統合型ネットワーク)を極めて短期間で構築したことにある。

## ModelBestのエッジAI展開における戦略的パートナーシップ網



ModelBestは、単なるAIモデルの提供にとどまらず、基盤となるAIチップ制御フレームワーク（FlagOS）から、最終製品である自動車（吉利、長安）やスマートデバイス（Huawei、Lenovo）まで、ハードウェアとソフトウェアの両面で強固な産業エコシステムを構築している。

### 7.1 モバイルSoCのハードウェア限界への最適化

スマートデバイスにおけるエッジ推論の性能は、ハードウェアのSoC (System on a Chip) の特性に大きく依存する。現在、ハイエンドAndroid市場のSoCは、Qualcommの「Snapdragon 8 Gen 3」とMediaTekの「Dimensity 9300」によって二分されている<sup>23</sup>。ベンチマークテストにおいて、両者はアプリケーションのワークロード処理で拮抗しているが、ストレステストにおいては明確な差異が生じる。Snapdragon 8 Gen 3は高負荷時にも安定したパフォーマンスを維持する傾向がある一方、Dimensity 9300は冷却性能とワークロードへの依存度が高く、熱によるパフォーマンス低下（サーマルスロットリング）が顕著に現れる場面がある<sup>23</sup>。しかし同時に、Dimensity 9300は3DMark Wild Life ExtremeテストにおいてSnapdragonやAppleのA17 Proを凌駕する瞬間的なGPU性能を叩き出す能力も有している<sup>25</sup>。ModelBestはHuawei（ファーウェイ）やMediaTekといった通信・半導体プレイヤーと深い連携を図り、自社のMiniCPMモデルがこれら多様なSoC環境において、NPUやGPUのリソースを適切に配分し、サーマルスロットリングを回避しながら最適なレイテンシで推論を行えるよう、極めて深いレベルでのハードウェア最適化を推進している<sup>26</sup>。

### 7.2 ソフトウェア・スタックの抽象化：「FlagOS」と「Align-Anything」

米国による最先端GPUの輸出規制を背景に、中国市場ではHygon（海光）、Metax（沐曦）、Iluvatar

(天数智芯)、Ascend(昇騰)といった多様な国産AIチップが乱立し、ソフトウェア・スタックの断片化が深刻な課題となっている<sup>13</sup>。この問題を解決するため、ModelBestは北京智源人工知能研究院(BAAI)などが主導するオープンソースの統合システムソフトウェアスタック「FlagOS」に強力にコミットしている。FlagOSは「Model-System-Chip」の3層構造を持ち、「一度開発すれば、異なるチップ間で移行可能 (develop once, migrate across chips)」という理念を掲げている<sup>13</sup>。MiniCPM-o 4.5は、この「FlagRelease」ツールキットを介して、NVIDIA製GPUだけでなく前述の国産AIチップ6種類に最適化された形で標準提供されている。さらに、エッジ展開に不可欠な高速推論サーバーである「vLLM」や「SGLang」に対するマルチチップ推論プラグイン(vLLM-plugin-FL)も提供されており、企業はハードウェアの制約に縛られることなくMiniCPMを自在にスケールさせることが可能である<sup>13</sup>。

加えて、北京大学のアライメント・チームが開発する「Align-Anything」フレームワークのサポートも見逃せない。このフレームワークは、テキスト、画像、動画、音声といった全モダリティ(any-to-any)にわたって、人間の意図とモデルをアライメント(整合)させるためのツールである<sup>13</sup>。これにより、研究者やエンタープライズ企業はMiniCPMを用いて、SFT(教師ありファインチューニング)やDPO(直接選好最適化)、さらにはDeepSeek-R1で用いられたGRPO(Generative Real-time Preference Optimization)といった最先端の強化学習をローカル環境で独自に実行できる。エッジLLMのカスタマイズ性が極限まで高められることで、高度に専門化されたAI(医療診断、精密製造など)の社会実装が急速に進展している。

### 7.3 インテリジェント・コックピットとスマートホームへの実装

このハードウェアとソフトウェアの統合力は、最終製品レベルで既に結実している。自動車産業において、ModelBestは吉利自動車(Geely)、フォルクスワーゲン(Volkswagen)、長安自動車(Changan)、広州汽車集団(GAC)、長城汽車(Great Wall Motor)と強固な提携を結んでいる<sup>1</sup>。2024年4月に量産化された長安マツダの「MAZDA EZ-60」や、同年9月に発売された吉利のフラッグシップSUV「Galaxy M9」には、MiniCPMのマルチモーダルモデルがインテリジェント・コックピットのコアとして実装された<sup>1</sup>。これにより、通信インフラが途絶するトンネルや山間部においても、車載カメラを用いた視覚的認識や自然な音声対話が可能な次世代の「人機インターフェース(HMI)」が実現している。

また、PC領域においては「Intel AIPC Client」としてのアプリケーション展開を通じて、企業の機密データをクラウドに一切アップロードすることなくローカルで処理するセキュリティ要件を満たしたソリューションを提供している<sup>27</sup>。スマートホーム分野でも、Yeelightなどのブランドと連携し、遅延のない高度な制御インフラを構築している<sup>26</sup>。

## 8. オープンソース戦略: Apache 2.0ライセンスによる市場制圧

ModelBestの成長を強烈に後押ししているのが、そのアグレッシブなオープンソース戦略である。同社のMiniCPMシリーズ(リポジトリおよびモデルウェイト)は、商業利用に対して極めて寛容な「Apache 2.0ライセンス」のもとで公開されている<sup>11</sup>。

### 8.1 ベンダーロックインの回避と商業化の促進

一般的に、企業が製品にAIを組み込む際、OpenAIやAnthropicといったプロプライエタリ(非公開)な企業のAPIに依存することは甚大なリスクを伴う。プロバイダー側が突然の価格改定を行ったり、利用規約を変更して特定用途での使用を禁じたり、あるいは旧バージョンのモデルを非推奨(非公開)にした場合、依存企業はなす術もなくシステム移行を余儀なくされる「ベンダーロックイン」の罠に陥る<sup>28</sup>。

これに対し、Apache 2.0ライセンスで提供されるMiniCPMは、ユーザー自身のサーバーやエッジデバイスに完全にホストすることが可能であり、デプロイメントの主導権を完全にユーザー側が掌握できる<sup>28</sup>。また、コピーレフト型ライセンス(GPLなど)とは異なり、利用者がソースコードに変更を加えた場合でも、その派生物をオープンソースとして公開する義務がない<sup>29</sup>。企業は元の著作権表示や改変の通知を含めるという最小限の要件さえ満たせば、MiniCPMをコアエンジンとした独自のクローズドな商用製品を開発・販売することが可能である<sup>29</sup>。

## 8.2 デファクトスタンダードへの道程

なお、ライセンス文書内には「モデルは大量のテキストから学習してコンテンツを生成するものであり、個人的な意見や価値判断を表現する能力は持たない」「モデルの生成物は開発者の見解を代表するものではない」「生成されたコンテンツの評価と検証の責任は全的にユーザーが負う」という明確な免責条項が設けられており、法的リスクの切り離しも適切に行われている<sup>27</sup>。この寛容なライセンス形態と、前述した「1Mトークンの長文脈処理」や「96倍の動画トークン圧縮」といった他を圧倒する性能が組み合わさったことで、世界中のハードウェアメーカーやソフトウェア開発者が、法的・経済的リスクを負うことなくMiniCPMを自社製品の中核に採用しやすくなった。その結果として、GitHubやHugging FaceにおけるMiniCPMシリーズの累計ダウンロード数は2400万回という驚異的なマイルストーンを突破したのである<sup>1</sup>。これは、競合他社を抑えてエッジAI分野におけるデファクトスタンダード(事実上の標準)の地位を確立するための、極めて高度な市場制圧戦略であると評価できる。

## 9. 結論

本レポートにおける多角的な分析を通じて、中国のAIスタートアップ「面壁智能(ModelBest)」が単なる高性能LLMの開発企業にとどまらず、次世代のコンピューティング・パラダイムである「エッジAIエコシステム」の基盤を設計する中核的インフラプロバイダーとしての地位を確立しつつあることが明らかになった。

同社の競争優位性は、以下の3つの強固な柱によって支えられている。

第一に、「Density Principle(知識密度法則)」に基づく徹底的なモデルの軽量化と、アーキテクチャの構造的イノベーションである。風洞実験的アプローチによる学習コストの最小化、SALAアーキテクチャによる最大1Mトークンという前例のない長文脈処理の実現、そして3D-Resamplerによる動画トークンの96倍圧縮技術は、デバイスの物理的制約(メモリ容量、帯域幅、遅延)をアルゴリズムレベルで完全に克服した。これにより、パラメータ数において圧倒的に巨大なGPT-4oやClaude 3.5 Sonnetを一部の論理推論ベンチマーク(HumanEval 95.1%等)で凌駕するという歴史的逆転現象を引き起こしている。

第二に、基盤ソフトウェアの抽象化とオープンソース化による技術の民主化である。Apache 2.0ライセンスによるベンダーロックインの排除、そしてFlagOSやAlign-Anythingといったフレームワーク群

への適応は、企業がハードウェア(AIチップ)やクラウドAPIの制約から脱却し、完全に独自かつセキュアなAIソリューションを構築することを可能にした。

第三に、実体経済における強力なアライアンスネットワークである。中国電信、インダストリアル・ロボティクス企業(Inovance)、さらには多数の自動車メーカーやモバイルSoCベンダーとの垂直統合型の提携は、モデルの優位性を机上の空論にとどめず、自動運転車のコックピットや産業現場へ即座に量産実装する巨大なチャンネルを形成している。

今後、人工知能の主戦場は「クラウド上でいかに巨大なモデルを構築するか」という単純なスケール競争から、「限られた消費電力とローカル環境の中で、いかにリアルタイムでプロアクティブに人間と世界を理解・支援するか」というエッジ・エンボディドAI(身体性・環境適応型AI)の領域へと不可逆的にシフトしていく。MiniCPM-oが提示した「全二重オムニモーダル通信」は、その未来像を具現化した最初のマイルストーンである。設立からわずか数年でユニコーン企業へと変貌を遂げたModelBestは、MetaのLlamaシリーズやMicrosoftのPhiシリーズにとっての最大の技術的脅威であるのみならず、今後のグローバルなAIアーキテクチャの標準を再定義し、IoT空間全体を牽引していく決定的なキープレイヤーとなることは疑いようがない。

## 引用文献

1. Seeds | ModelBest Completes New Financing Round of Hundreds of Millions of Yuan, 5月 4, 2026にアクセス、  
<https://autonews.gasgoo.com/articles/news/seeds-modelbest-completes-new-financing-round-of-hundreds-of-millions-of-yuan-2041514837547327489>
2. ModelBest Secures Hundreds of Millions in Funding Led by China Telecom, Launching Deep Business Collaboration - Pandaily, 5月 4, 2026にアクセス、  
<https://pandaily.com/model-best-secures-hundreds-of-millions-in-funding-led-by-china-telecom-launching-deep-business-collaboration>
3. ModelBest secures fresh funding to accelerate on-device AI deployment - Gasgoo, 5月 4, 2026にアクセス、  
<https://autonews.gasgoo.com/articles/news/modelbest-secures-fresh-funding-to-accelerate-on-device-ai-deployment-70040208>
4. MiniCPM: Unveiling the Potential of Small Language Models with Scalable Training Strategies - arXiv, 5月 4, 2026にアクセス、<https://arxiv.org/html/2404.06395v1>
5. MiniCPM: Unveiling the Potential of End-side Large Language Models - OpenBMB Blog, 5月 4, 2026にアクセス、<https://openbmb.vercel.app/minicpm-en>
6. MiniCPM-SALA: Hybridizing Sparse and Linear Attention for Efficient Long-Context Modeling - arXiv, 5月 4, 2026にアクセス、  
<https://arxiv.org/html/2602.11761v1>
7. MiniCPM-SALA: Hybridizing Sparse and Linear Attention for Efficient Long-Context Modeling - ResearchGate, 5月 4, 2026にアクセス、  
[https://www.researchgate.net/publication/400741879\\_MiniCPM-SALA\\_Hybridizing\\_Sparse\\_and\\_Linear\\_Attention\\_for\\_Efficient\\_Long-Context\\_Modeling](https://www.researchgate.net/publication/400741879_MiniCPM-SALA_Hybridizing_Sparse_and_Linear_Attention_for_Efficient_Long-Context_Modeling)
8. openbmb/MiniCPM-SALA - Hugging Face, 5月 4, 2026にアクセス、  
<https://huggingface.co/openbmb/MiniCPM-SALA>
9. [2602.11761] MiniCPM-SALA: Hybridizing Sparse and Linear Attention for Efficient

- Long-Context Modeling - arXiv, 5月 4, 2026にアクセス、  
<https://arxiv.org/abs/2602.11761>
10. HumanEval Leaderboard - LLM Stats, 5月 4, 2026にアクセス、  
<https://llm-stats.com/benchmarks/humaneval>
  11. openbmb/MiniCPM-V-4\_5 - Hugging Face, 5月 4, 2026にアクセス、  
[https://huggingface.co/openbmb/MiniCPM-V-4\\_5](https://huggingface.co/openbmb/MiniCPM-V-4_5)
  12. MiniCPM-V 4.5 : Best LLM for Mobiles | by Mehul Gupta | Data Science in Your Pocket, 5月 4, 2026にアクセス、  
<https://medium.com/data-science-in-your-pocket/minicpm-v-4-5-best-llm-for-mobiles-94e8b91ac994>
  13. GitHub - OpenBMB/MiniCPM-o: A Gemini 2.5 Flash Level MLLM for Vision, Speech, and Full-Duplex Multimodal Live Streaming on Your Phone, 5月 4, 2026にアクセス、  
<https://github.com/OpenBMB/MiniCPM-o>
  14. MiniCPM-V 4.5: Cooking Efficient MLLMs via Architecture, Data, and Training Recipes, 5月 4, 2026にアクセス、  
<https://arxiv.org/html/2509.18154v1>
  15. Llama 3.1 8B Instruct vs Phi-3.5-mini-instruct Comparison - LLM Stats, 5月 4, 2026にアクセス、  
<https://llm-stats.com/models/compare/llama-3.1-8b-instruct-vs-phi-3.5-mini-instruct>
  16. microsoft/Phi-3.5-mini-instruct - Hugging Face, 5月 4, 2026にアクセス、  
<https://huggingface.co/microsoft/Phi-3.5-mini-instruct>
  17. Azure Llama 3.1 benchmarks : r/LocalLLaMA - Reddit, 5月 4, 2026にアクセス、  
[https://www.reddit.com/r/LocalLLaMA/comments/1e9hg7g/azure\\_llama\\_31\\_benchmarks/](https://www.reddit.com/r/LocalLLaMA/comments/1e9hg7g/azure_llama_31_benchmarks/)
  18. Interesting Model Differences Between Phi-3.5-Mini & Phi-3.5-MoE : r/LocalLLaMA - Reddit, 5月 4, 2026にアクセス、  
[https://www.reddit.com/r/LocalLLaMA/comments/1exn6wx/interesting\\_model\\_differences\\_between\\_phi35mini/](https://www.reddit.com/r/LocalLLaMA/comments/1exn6wx/interesting_model_differences_between_phi35mini/)
  19. Discover the New Multi-Lingual, High-Quality Phi-3.5 SLMs - Microsoft Community Hub, 5月 4, 2026にアクセス、  
<https://techcommunity.microsoft.com/blog/azure-ai-foundry-blog/discover-the-new-multi-lingual-high-quality-phi-3-5-slms/4225280>
  20. Comparing Llama3, Phi3 and Gemma performance on different machines - Nithin Bekal, 5月 4, 2026にアクセス、  
<https://nithinbekal.com/posts/comparing-llama3-phi3-gemma/>
  21. MiniCPM-V 2.6: A GPT-4V Level MLLM for Single Image, Multi Image and Video on Your Phone - GitHub, 5月 4, 2026にアクセス、  
<https://github.com/QAdottech/MiniCPM-V>
  22. MiniCPM-V 2.6 Training Guide - YouTube, 5月 4, 2026にアクセス、  
<https://www.youtube.com/watch?v=j7OuDEek880>
  23. Snapdragon 8 Gen 3 vs Dimensity 9300 benchmarked - Android Authority, 5月 4, 2026にアクセス、  
<https://www.androidauthority.com/snapdragon-8-gen-3-dimensity-9300-benchmarked-3395385/>
  24. Snapdragon 8 Gen 3 vs Dimensity 9300 benchmarked: There can only be one

- winner, 5月 4, 2026にアクセス、  
[https://www.reddit.com/r/Android/comments/18pmbxp/snapdragon\\_8\\_gen\\_3\\_vs\\_dimensity\\_9300\\_benchmarked/](https://www.reddit.com/r/Android/comments/18pmbxp/snapdragon_8_gen_3_vs_dimensity_9300_benchmarked/)
25. Dimensity 9300 Is 11.7 Percent Faster Than Snapdragon 8 Gen 3 In 3DMark Wild Life Extreme At The Same Power Usage; Overwhelms A17 Pro GPU - Wccfttech, 5月 4, 2026にアクセス、  
<https://wccfttech.com/3dmark-wild-life-extreme-test-dimensity-9300-beats-snapdragon-8-gen-3-and-a17-pro/>
  26. AI tech company ModelBest closes new financing round - Gasgoo, 5月 4, 2026にアクセス、  
<https://autonews.gasgoo.com/articles/icv/70035412>
  27. GitHub - OpenBMB/MiniCPM: MiniCPM4 & MiniCPM4.1: Ultra ..., 5月 4, 2026にアクセス、  
<https://github.com/OpenBMB/MiniCPM>
  28. What Is Gemma 4's Apache 2.0 License? Why It Matters More Than the Model Itself, 5月 4, 2026にアクセス、  
<https://www.mindstudio.ai/blog/gemma-4-apache-2-license-commercial-use>
  29. Open Source Licenses 101: Apache License 2.0 | FOSSA Blog, 5月 4, 2026にアクセス、  
<https://fossa.com/blog/open-source-licenses-101-apache-license-2-0/>
  30. Does it seem like MIT is not favored over Apache? : r/opensource - Reddit, 5月 4, 2026にアクセス、  
[https://www.reddit.com/r/opensource/comments/1amsq4v/does\\_it\\_seem\\_like\\_mit\\_is\\_not\\_favored\\_over\\_apache/](https://www.reddit.com/r/opensource/comments/1amsq4v/does_it_seem_like_mit_is_not_favored_over_apache/)